

International Association of Geodesy Symposia

142

Nico Sneeuw  
Pavel Novák  
Mattia Crespi  
Fernando Sansò *Editors*

# VIII Hotine-Marussi Symposium on Mathematical Geodesy

Proceedings of the Symposium in Rome,  
17-21 June, 2013

 Springer

# International Association of Geodesy Symposia

*Chris Rizos, Series Editor*  
*Pascal Willis, Assistant Series Editor*

---

# International Association of Geodesy Symposia

*Chris Rizos, Series Editor*  
*Pascal Willis, Assistant Series Editor*

---

- Symposium 101: Global and Regional Geodynamics
- Symposium 102: Global Positioning System: An Overview
- Symposium 103: Gravity, Gradiometry, and Gravimetry
- Symposium 104: Sea Surface Topography and the Geoid
- Symposium 105: Earth Rotation and Coordinate Reference Frames
- Symposium 106: Determination of the Geoid: Present and Future
- Symposium 107: Kinematic Systems in Geodesy, Surveying, and Remote Sensing
- Symposium 108: Application of Geodesy to Engineering
- Symposium 109: Permanent Satellite Tracking Networks for Geodesy and Geodynamics
- Symposium 110: From Mars to Greenland: Charting Gravity with Space and Airborne Instruments
- Symposium 111: Recent Geodetic and Gravimetric Research in Latin America
- Symposium 112: Geodesy and Physics of the Earth: Geodetic Contributions to Geodynamics
- Symposium 113: Gravity and Geoid
- Symposium 114: Geodetic Theory Today
- Symposium 115: GPS Trends in Precise Terrestrial, Airborne, and Spaceborne Applications
- Symposium 116: Global Gravity Field and Its Temporal Variations
- Symposium 117: Gravity, Geoid and Marine Geodesy
- Symposium 118: Advances in Positioning and Reference Frames
- Symposium 119: Geodesy on the Move
- Symposium 120: Towards an Integrated Global Geodetic Observation System (IGGOS)
- Symposium 121: Geodesy Beyond 2000: The Challenges of the First Decade
- Symposium 122: IV Hotine-Marussi Symposium on Mathematical Geodesy
- Symposium 123: Gravity, Geoid and Geodynamics 2000
- Symposium 124: Vertical Reference Systems
- Symposium 125: Vistas for Geodesy in the New Millennium
- Symposium 126: Satellite Altimetry for Geodesy, Geophysics and Oceanography
- Symposium 127: V Hotine Marussi Symposium on Mathematical Geodesy
- Symposium 128: A Window on the Future of Geodesy
- Symposium 129: Gravity, Geoid and Space Missions
- Symposium 130: Dynamic Planet - Monitoring and Understanding . . .
- Symposium 131: Geodetic Deformation Monitoring: From Geophysical to Engineering Roles
- Symposium 132: VI Hotine-Marussi Symposium on Theoretical and Computational Geodesy
- Symposium 133: Observing our Changing Earth
- Symposium 134: Geodetic Reference Frames
- Symposium 135: Gravity, Geoid and Earth Observation
- Symposium 136: Geodesy for Planet Earth
- Symposium 137: VII Hotine-Marussi Symposium on Mathematical Geodesy
- Symposium 138: Reference Frames for Applications in Geosciences
- Symposium 139: Earth on the Edge: Science for a Sustainable Planet
- Symposium 140: The 1st International Workshop on the Quality of Geodetic Observation and Monitoring Systems (QuGOMS'11)
- Symposium 141: Gravity, Geoid and Height systems (GGHS2012)

# VIII Hotine-Marussi Symposium on Mathematical Geodesy

Proceedings of the Symposium in Rome,  
17–21 June, 2013

Edited by

Nico Sneeuw  
Pavel Novák  
Mattia Crespi  
Fernando Sansò

*Volume Editors*

Nico Sneeuw  
Institute of Geodesy  
University of Stuttgart  
Stuttgart  
Germany

Pavel Novák  
Department of Mathematics  
University of Western Bohemia  
Pilsen  
Czech Republic

Mattia Crespi  
Geodesy and Geomatics Division  
University of Rome "La Sapienza"  
Rome  
Italy

Fernando Sansò  
Dipartimento di Ingegneria Civile e Ambientale  
Politecnico di Milano  
Milano  
Italy

*Series Editor*

Chris Rizos  
University of New South Wales  
Sydney  
New South Wales  
Australia

*Associate Editor*

Pascal Willis  
Institut national de l'Information  
géographique et forestière  
Service de la Recherche  
et de l'Enseignement  
Saint-Mandé  
France

ISSN 0939-9585  
International Association of Geodesy Symposia  
ISBN 978-3-319-24548-5  
DOI 10.1007/978-3-319-30530-1

ISSN 2197-9359 (electronic)  
ISBN 978-3-319-30530-1 (eBook)

Library of Congress Control Number: 2016935295

© Springer International Publishing Switzerland 2016

This work is subject to copyright. All rights are reserved by the Publisher, whether the whole or part of the material is concerned, specifically the rights of translation, reprinting, reuse of illustrations, recitation, broadcasting, reproduction on microfilms or in any other physical way, and transmission or information storage and retrieval, electronic adaptation, computer software, or by similar or dissimilar methodology now known or hereafter developed. The use of general descriptive names, registered names, trademarks, service marks, etc. in this publication does not imply, even in the absence of a specific statement, that such names are exempt from the relevant protective laws and regulations and therefore free for general use.

The publisher, the authors and the editors are safe to assume that the advice and information in this book are believed to be true and accurate at the date of publication. Neither the publisher nor the authors or the editors give a warranty, express or implied, with respect to the material contained herein or for any errors or omissions that may have been made.

Printed on acid-free paper

This Springer imprint is published by Springer Nature  
The registered company is Springer International Publishing AG Switzerland

---

## Preface

This volume contains the proceedings of the VIII Hotine-Marussi Symposium on Mathematical Geodesy, which was held June 17 to 21, 2013. For a second time in row, the Symposium took place at the Faculty of Engineering of the University of Rome “La Sapienza”, Italy. Again, the symposium was hosted in the beautiful ancient *chiostro* of the Basilica of S. Pietro in Vincoli, famously known for its statue of Moses by Michelangelo.

The traditional name *mathematical geodesy* for the series of Hotine-Marussi Symposia may not fully do justice to the Symposium’s broad scope of theoretical geodesy in general. However, the name for the series has been used since 1965, i.e. the days of Antonio Marussi, which is a good reason to adhere to it. The venue of the Hotine-Marussi Symposia has traditionally been in Italy, as exemplified in the historical overview and map on the next pages.

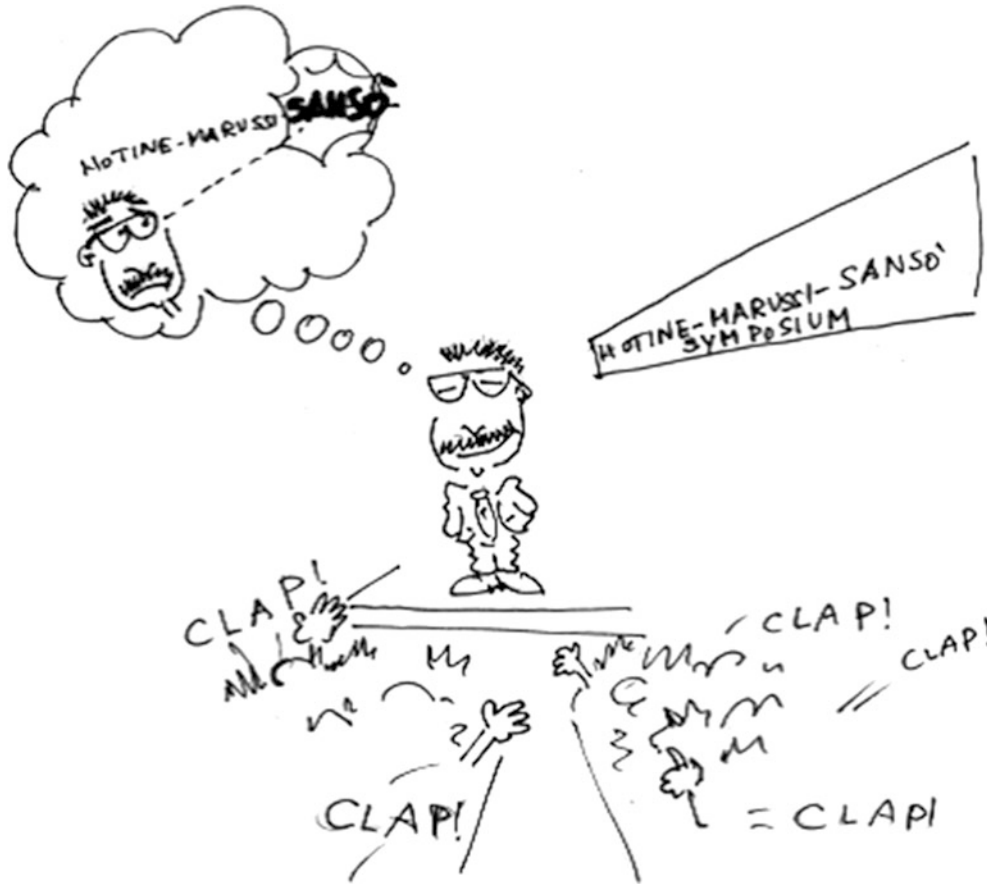
Since 2006 the series is under the responsibility of the Intercommission Committee on Theory (ICCT), a cross-commission entity within the International Association of Geodesy (IAG). The overall goal of the Hotine-Marussi Symposia has always been the advancement of theoretical geodesy. This goal is aligned with the objectives of the ICCT, which has the developments in geodetic modelling and data processing in the light of recent advances of geodetic observing systems as well as the exchange between geodesy and neighbouring Earth sciences as its central themes. Indeed, the current proceedings are testimony to the width and vibrancy of theoretical geodesy.

The Symposium attracted 90 participants who contributed 88 papers (71 oral and 17 poster), organized in eight regular sessions plus the session at the Accademia Nazionale dei Lincei. To a large extent, the sessions’ topics were modelled on the study group structure of the ICCT. The chairs of the ICCT study groups, who constituted the Symposium’s Scientific Committee, were at the same time responsible for organizing the sessions:

1. *Geodetic data analysis*  
W. Kosek, R. Gross, C. Kreemer
2. *Geopotential modeling, boundary value problems and height systems*  
P. Novák, M. Schmidt, C. Gerlach
3. *Atmospheric modeling in geodesy*  
T. Hobiger, M. Schindelegger
4. *Gravity field mapping methodology from GRACE and future gravity missions*  
M. Weigelt, A. Jäggi
5. *Computational geodesy*  
R. Čunderlík, K. Mikula
6. *Theoretical aspects of reference frames*  
A. Dermanis, T. Van Dam
7. *Digital Terrain Modeling, Synthetic Aperture Radar and new sensors: theory and methods*  
M. Crespi, E. Pottier
8. *Inverse modeling, estimation theory*  
P. Xu

A special session was organized at the Accademia Nazionale dei Lincei by Fernando Sansò, himself a member of this venerable academy, with three keynote addresses. What Fernando did

not know at this point was that the organizing committee had decided to dedicate the VIII Hotine-Marussi Symposium in his honour. To put Fernando into the spotlight four brief speeches followed by four renowned geodesists and old-time colleagues of Fernando: Michele Caputo, Sakis Dermanis, Erik Grafarend and Christian Tscherning. Please note that the latter three names represent the universities of Thessaloniki, Stuttgart and Copenhagen, respectively, from which Fernando has been awarded honorary doctorates. Each of these gentlemen reminisced about their long-term collaboration and friendship with Fernando, but they also characterized him by entertaining anecdotes. Thus the long-term commitment and dedication of Prof. Sansò was acknowledged, who has been the driving force behind the series of Hotine-Marussi Symposia over the past decades. Whether this honour might be a burden at the same time, as the cartoon seems to suggest, well: future will tell.



Credits: Riccardo Barzaghi

We want to express our gratitude to all of those who have contributed to the success of the VIII Hotine-Marussi Symposium. The aforementioned study group chairs (Scientific Committee) put much effort in organizing attractive sessions and convening them. They also took a leading role in the peer-review process, which was managed by the IAG proceedings editor Dr. Pascal Willis. We equally owe thanks to all reviewers. Although much of the review process itself remains anonymous, the complete list of the reviewers is printed in this volume as a token of our appreciation of their dedication.

Financial and promotional support was given by the Faculty of Engineering of the Sapienza University of Rome.

But most of our thanks are due to Mattia Crespi and his team of the Area of Geodesy and Geomatics (AGG), which is part of the Department of Civil, Building and Environmental Engineering (DICEA), who hosted the Symposium. It is well known that the quality of a

Local Organizing Committee (LOC) is decisive to a successful scientific meeting. Beyond responsibility for website, registration, technical support and all kinds of other arrangements, the LOC organized a visit to the Villa Farnesina and a great social event to the Capitoline Hill, including a guided museum tour and a roof-top dinner with an astonishing view over the eternal city. Through their able organization and improvisation skills Mattia Crespi and his team (Elisa Benedetti, Mara Branzanti, Paola Capaldo, Gabriele Colosimo, Nicole Dore, Francesca Fratarcangeli, Augusto Mazzoni, Andrea Nascetti, Jolanda Patruno, Francesca Pieralice and Martina Porfiri) have done more than their share in bringing the VIII Hotine-Marussi Symposium to success.

Stuttgart  
October 2014

Nico Sneeuw  
Pavel Novák  
Mattia Crespi  
Fernando Sansò





---

## Fifty Years of Hotine-Marussi Symposia

In 1959, Antonio Marussi, in cooperation with the Italian Geodetic Commission, started a series of symposia in Venice. The first three of these covered the entire theoretical definition of 3D geodesy, as delineated in discussions with renowned contemporary scientists:

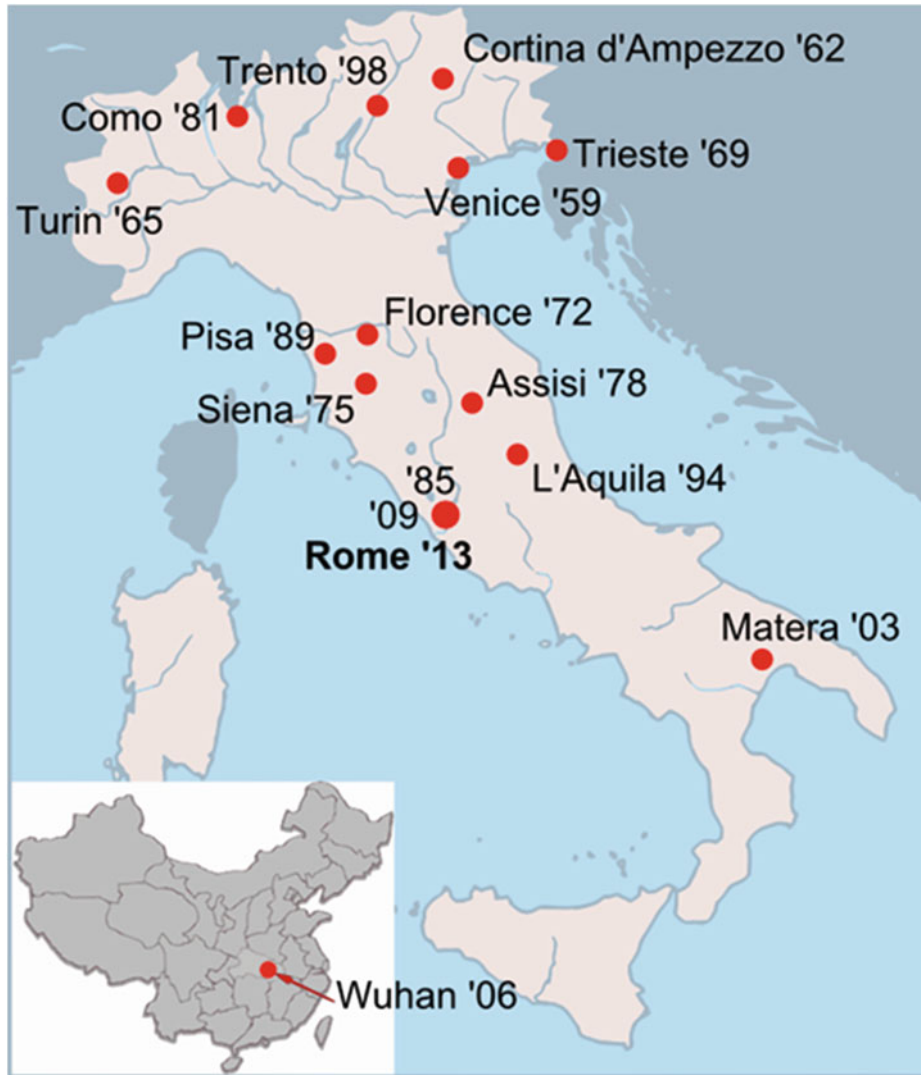
- 16–18 July 1959, Venice, 1st Symposium on Three Dimensional Geodesy, published in *Bollettino di Geodesia e Scienze Affini*, XVIII, N° 3, 1959
- 29 May–1 June 1962, Cortina d'Ampezzo, 2nd Symposium on Three Dimensional Geodesy, published in *Bollettino di Geodesia e Scienze Affini*, XXI, N° 3, 1962
- 21–22 April 1965, Turin, 3rd Symposium on Mathematical Geodesy, published by Commissione Geodetica Italiana, 1966

From the very beginning, Martin Hotine provided essential inspiration to these symposia. After his death in 1968, the following symposia bear his name:

- 28–30 May 1969, Trieste, 1st Hotine Symposium (4th Symposium on Mathematical Geodesy), published by Commissione Geodetica Italiana, 1970
- 25–26 October 1972, Florence, 2nd Hotine Symposium (5th Symposium on Mathematical Geodesy), published by Commissione Geodetica Italiana, 1973
- 2–5 April 1975, Siena, 3rd Hotine Symposium (6th Symposium on Mathematical Geodesy), published by Commissione Geodetica Italiana, 1975
- 8–10 June 1978, Assisi, 4th Hotine Symposium (7th Symposium on Mathematical Geodesy), published by Commissione Geodetica Italiana, 1978
- 7–9 September 1981, Como, 5th Hotine Symposium (8th Symposium on Mathematical Geodesy), published by Commissione Geodetica Italiana, 1981

After Marussi's death, in 1984, the symposia were finally named the Hotine-Marussi Symposia:

- 3–6 June 1985, Rome, I Hotine-Marussi Symposium (Mathematical Geodesy)
- 5–8 June 1989, Pisa, II Hotine-Marussi Symposium (Mathematical Geodesy)
- 29 May–3 June 1994, L'Aquila, III Hotine-Marussi Symposium (Mathematical Geodesy, Geodetic Theory Today), published by Springer, IAG 114
- 14–17 September 1998, Trento, IV Hotine-Marussi Symposium (Mathematical Geodesy), published by Springer, IAG 122
- 17–21 June 2003, Matera, V Hotine-Marussi Symposium (Mathematical Geodesy), published by Springer, IAG 127
- 29 May–2 June 2006, Wuhan, VI Hotine-Marussi Symposium (Theoretical and Computational Geodesy, 1st time under ICCT), published by Springer, IAG 132
- 6–10 July 2009, Rome, VII Hotine-Marussi Symposium (Mathematical Geodesy), published by Springer, IAG 137
- 17–21 June 2013, Rome, VIII Hotine-Marussi Symposium (Mathematical Geodesy), published by Springer, IAG 142



---

# Contents

## Part I Lincei Session

<b>Opening Remarks for the VIII Hotine-Marussi Symposium</b> .....	3
Michele Caputo	
<b>Fernando Sansò Laudation</b> .....	7
Michele Caputo	
<b>Global Reference Systems: Theory and Open Questions</b> .....	9
Athanasios Dermanis	

## Part II Geodetic Data Analysis

<b>Noise Analysis of Continuous GPS Time Series of Selected EPN Stations to Investigate Variations in Stability of Monument Types</b> .....	19
Anna Klos, Janusz Bogusz, Mariusz Figurski, and Wieslaw Kosek	
<b>Improvement of Least-Squares Collocation Error Estimates Using Local GOCE <math>T_{zz}</math> Signal Standard Deviations</b> .....	27
C.C. Tscherning	
<b>Multivariate Integer Cycle-Slip Resolution: A Single-Channel Analysis</b> .....	33
P.J.G. Teunissen and P.F. de Bakker	
<b>Theory of Earth Rotation Variations</b> .....	41
Richard S. Gross	
<b>Variable Seasonal and Subseasonal Oscillations in Sea Level Anomaly Data and Their Impact on Prediction Accuracy</b> .....	47
W. Kosek, T. Niedzielski, W. Popiński, M. Zbylut-Górska, and A. Wnęk	
<b>Permanent GPS Networks in Italy: Analysis of Time Series Noise</b> .....	51
R. Devoti, G. Pietrantonio, A.R. Pisani, and F. Riguzzi	
<b>VADASE: State of the Art and New Developments of a Third Way to GNSS Seismology</b> .....	59
E. Benedetti, M. Branzanti, G. Colosimo, A. Mazzoni, and M. Crespi	
<b>On the Spatial Resolution of Homogeneous Isotropic Filters on the Sphere</b> .....	67
Balaji Devaraju and Nico Sneeuw	
<b>On Time-Varying Seasonal Signals: Comparison of SSA and Kalman Filtering Based Approach</b> .....	75
Q. Chen, M. Weigelt, N. Sneeuw, and T. van Dam	
<b>Extensive Analysis of IGS REPRO1 Coordinate Time Series</b> .....	81
M. Roggero	

### Part III Geopotential Modeling, Boundary Value Problems and Height Systems

<b>Determination of <math>W_0</math> from the GOCE Measurements Using the Method of Fundamental Solutions</b> .....	91
Róbert Čunderlík	
<b>Combination of GOCE Gravity Gradients in Regional Gravity Field Modelling Using Radial Basis Functions</b> .....	101
Verena Lieb, Johannes Bouman, Denise Dettmering, Martin Fuchs, and Michael Schmidt	
<b>Rosborough Representation in Satellite Gravimetry</b> .....	109
Nico Sneeuw and Mohammad A. Sharifi	
<b>Combining Different Types of Gravity Observations in Regional Gravity Modeling in Spherical Radial Basis Functions</b> .....	115
Katrin Bentel and Michael Schmidt	
<b>Height Datum Unification by Means of the GBVP Approach Using Tide Gauges</b> ...	121
E. Rangelova, M.G. Sideris, B. Amjadiparvar, and T. Hayden	

### Part IV Atmospheric Modeling in Geodesy

<b>Computation of Zenith Total Delay Correction Fields Using Ground-Based GNSS</b> .....	131
B. Pace, R. Pacione, C. Sciarretta, and G. Bianco	
<b>Rigorous Interpolation of Atmospheric State Parameters for Ray-Traced Tropospheric Delays</b> .....	139
Camille Desjardins, Pascal Gegout, Laurent Soudarin, and Richard Biancale	
<b>Comparison of Different Techniques for Tropospheric Wet Delay Retrieval Over South America and Surrounding Oceans</b> .....	147
A. Calori, G. Colosimo, M. Crespi, and M.V. Mackern	

### Part V Gravity Field Mapping Methodology from GRACE and Future Gravity Missions

<b>The Role of Position Information for the Analysis of K-Band Data: Experiences from GRACE and GOCE for GRAIL Gravity Field Recovery</b> .....	157
A. Jäggi, G. Beutler, U. Meyer, H. Bock, and L. Mervart	
<b>Gravity Field Mapping from GRACE: Different Approaches—Same Results?</b> .....	165
Christoph Dahle, Christian Gruber, Elisa Fagiolini, and Frank Flechtner	
<b>The Effect of Pseudo-Stochastic Orbit Parameters on GRACE Monthly Gravity Fields: Insights from Lumped Coefficients</b> .....	177
U. Meyer, C. Dahle, N. Sneeuw, A. Jäggi, G. Beutler, and H. Bock	
<b>On an Iterative Approach to Solving the Nonlinear Satellite-Fixed Geodetic Boundary-Value Problem</b> .....	185
Marek Macák, Karol Mikula, Zuzana Minarechová, and Róbert Čunderlík	

### Part VI Computational Geodesy

<b>An OpenCL Implementation of Ellipsoidal Harmonics</b> .....	195
Otakar Nesvadba and Petr Holota	

<b>A Remark on the Computation of the Gravitational Potential of Masses with Linearly Varying Density</b> .....	205
Maria Grazia D'Urso	
<b>The Observation Equation of Spirit Leveling in Molodensky's Context</b> .....	213
B. Betti, D. Carrion, F. Sacerdote, and G. Venuti	
<b>Part VII Theoretical Aspects of Reference Frames</b>	
<b>Reference Station Weighting and Frame Optimality in Minimally Constrained Networks</b> .....	221
C. Kotsakis	
<b>Atmospheric Loading and Mass Variation Effects on the SLR-Defined Geocenter</b> .....	227
Rolf König, Frank Flechtner, Jean-Claude Raimondo, and Margarita Vei	
<b>Radargrammetric Digital Surface Models Generation from High Resolution Satellite SAR Imagery: Methodology and Case Studies</b> .....	233
Andrea Nascetti, Paola Capaldo, Francesca Pieralice, Martina Porfiri, Francesca Fratarcangeli, and Mattia Crespi	
<b>Part VIII Digital Terrain Modeling, Synthetic Aperture Radar and New Sensors: Theory and Methods</b>	
<b>Principles and Applications of Polarimetric SAR Tomography for the Characterization of Complex Environments</b> .....	243
Laurent Ferro-Famil, Yue Huang, and Eric Pottier	
<b>Re-gridding and Merging Overlapping DTMS: Problems and Solutions in HELI-DEM</b> .....	257
Ludovico Biagi and Laura Carcano	
<b>Single-Epoch GNSS Array Integrity: An Analytical Study</b> .....	263
A. Khodabandeh and P.J.G. Teunissen	
<b>Part IX Inverse Modeling, Estimation Theory VIII Hotine-Marussi: Geodetic Data Analysis</b>	
<b>Global to Local Moho Estimate Based on GOCE Geopotential Model and Local Gravity Data</b> .....	275
R. Barzaghi, M. Reguzzoni, A. Borghi, C. De Gaetani, D. Sampietro, and A.M. Marotta	
<b>An Overview of Adjustment Methods for Mixed Additive and Multiplicative Random Error Models</b> .....	283
Yun Shi, Peiliang Xu, and Junhuan Peng	
<b>Cycle Slip Detection and Correction for Heading Determination with Low-Cost GPS/INS Receivers</b> .....	291
Patrick Henkel and Naoya Oku	
<b>Adjusting the Errors-In-Variables Model: Linearized Least-Squares vs. Nonlinear Total Least-Squares</b> .....	301
Burkhard Schaffrin	
<b>Multivariate GNSS Attitude Integrity: The Role of Affine Constraints</b> .....	309
Gabriele Giorgi and Peter J.G. Teunissen	

---

<b>Integrating Geological Prior Information into the Inverse Gravimetric Problem: The Bayesian Approach</b> .....	317
L. Rossi, M. Reguzzoni, D. Sampietro, and F. Sansò	
<b>Effects of Different Objective Functions in Inequality Constrained and Rank-Deficient Least-Squares Problems</b> .....	325
Lutz Roese-Koerner and Wolf-Dieter Schuh	
<b>List of Reviewers</b> .....	333
<b>Author Index</b> .....	335

---

**Part I**

**Lincei Session**



---

# Opening Remarks for the VIII Hotine-Marussi Symposium

Michele Caputo

---

## Abstract

Opening remarks for the 2013 Hotine-Marussi Symposium Special Session at the Accademia Nazionale dei Lincei.

---

## Keywords

Anelasticity • Geoidal undulations • Gravity • Rheology

---

## 1 The Intrinsic Geodesy

At the previous Hotine-Marussi symposium in Rome, 2009, I told of the Marussi pendulums, of the gravity absorption experiments, of the gravitons and of the important developments of the models of the gravity field of the Earth produced and discussed a century ago by Pizzetti and Somigliana, eventually presented in the Accademia, and founders of the well known elegant modelling of the gravity field of the Earth.

Then I asked myself what comes next? To my surprise I had not mentioned Marussi's *Intrinsic Geodesy*. The reason probably is the Cayley-Darboux theorem, which was preceded by the important works of Mainardi and of Codazzi on the mapping of revolution surfaces.

Marussi knew all this so well that he assigned the theorem as subject for a thesis asking me to be *correlatore* since he was very busy and often away (Milani 1961; Sansò and

Caputo 2008). Probably because of this theorem, Marussi never encouraged me to work in this field nor went any further with the concept of intrinsic geodesy which, practically, extends to the gravity field the classic principle of embedding any problem in its most appropriate coordinate system.

In this case the possible coordinate system would have been the family of equipotential surfaces of the gravity field which, unfortunately, did not turn out to satisfy the conditions of the theorem, except locally.

In this sense it was Grafarend (1988) to see that the theory could be applied locally. In another sense I feel that I should quote Bocchio (1974) who wrote on the extension of the concept of *intrinsic* to geophysics.

Philosophically I see the word *intrinsic* as the operation of imbedding the problem in its reality. Dig into the problem, deep enough until you find the *pepita de oro*, or the truffle depending on your taste, which shows you the path. The *pepita* may be the appropriate geometric space, or the appropriate functional space, or the appropriate sensor's type and/or the best intrinsic space distribution of the data.

---

Fernando Sansò and Mattia Crespi kindly asked me to open the works of the 2013 Hotine-Marussi symposium and speak of Geodesy and of the Accademia dei Lincei, and gave me the great pleasure to be here among my friends since Geodesy and the Accademia have filled a good part of my life.

M. Caputo (✉)  
College of Geosciences, Texas A&M University, College Station  
Texas, Rome, Italy  
e-mail: [mic.caputo27@gmail.com](mailto:mic.caputo27@gmail.com)

---

## 2 Geodesy and Rheology

Let me now go back half a century when I was sharing the office, with the person who was then my mentor, his desk between the two windows in the corner, my desk in front of his but in the opposite corner. The institute had been founded by Marussi, director, and members: me and a technician; in

an old building of Trieste then still occupied by the allied troops. The library of the institute was Marussi's office. On the shelves on the side of my desk was a book with the title *Rheology*, then mysterious for me, which at that time seemed to have very little to do with geodesy or geophysics.

All I knew after 5 years of studies of the Greek language in high school was that *rheo* means movement; digging out of my memory, I arrived to *panta rei*; but that was certainly not enough to satisfy any level of curiosity.

It was a *sign* on the wall. Later I realised that rheology would be implied in the future of geophysics and geodesy to the point to make it important for some realistic results. If I may make a comparison it would be a study of macroeconomic problems without taking into account memory (Demaria 1976, Caputo 2012a, Caputo 2012b). Rheology was the last discipline to join the set forming the interdisciplinary royal crown of Geodesy and Geophysics.

For instance today we are not so concerned with elevations but with their time variations either in the case of the coast lines and of sea level as separate matters.

We measure earth tides to study, with very poor resolution as a matter of fact, some average elastic properties of the interior of the Earth, but the major expected result was the estimate (Slichter et al. 1964) of the phase lag of the Moon relative to the bulges it creates in the Earth, involving the inelastic properties of the Earth; in this case more interesting than the elastic properties, and, in turn, we found that it is rheology modelling and ruling the history of the Earth–Moon system: and its future as well as those of all planetary systems.

Always concerned with rheology, one mystery is the distribution of the continental masses on the surface of the Earth, apparently disordered, although not casually, but with little longitudinal and latitudinal symmetry, which does not seem compatible with the stability of the Earth relative to the axis of rotation (Caputo and Caputo 2013). It is the rheology which adjusts the distribution of masses in the interior in order to keep the Earth in the same position relative to the axis of rotation, at least in the recent decades. The absence of the possible and inevitable associated precessions relative to the axis of rotation proves it.

Naturally all occurs in obedience of the second principle of thermodynamics, which is embedded in the constitutive equations of rheology and makes the general equations of elasticity physically consistent when rheology is included and finally acceptable. I do not mean the elasticity as is commonly used in seismology or that existing only in elementary books but that which considers also the dissipation of elastic energy or if you like the second principle of thermodynamics.

At that time, I mean when I was in Marussi's office, the Geoid was the Tanni's (1940) Geoid or little more. Limited to several lines hundreds of kilometres long in northern Europe. Practically in the first part of the twentieth century the Geoid

was an unreachable fluid spectrum and no attention was given to the effects of rheology. You could see Geoid, since it was there on all coast lines but, since it did not have any mathematical representation, it was as for Martin Eden in Jack London's book: . . . *the instant he knew, he ceased to know*.

Today we have the Geoid and observe the tides, also for the solid Earth.

The question is then posed: if we accept that all is changing during and between repeated measurements what should we change in our strategy and techniques in observing and interpreting the results? Will Geomatics alone answer this question? Or should it be interdisciplinary?

How to go further in the concept of embedding the problem, or more generally the matter, in its reality? How about symmetries?

The *pepita* concerning the functional space of rheology seems to be in the various types of mathematical formalisms, practically representing the second principle of thermodynamics, whose use in recent decades has spread in many fields of science.

Applied in fields such as theoretical physics (e.g. Naber 2004), biology (e.g. Cesarone et al. 2004), geodesy (e.g. Baleanu et al. 2009), medicine (e.g. El Shahed 2003), diffusion (e.g. Mainardi 1996), chemistry (e.g. Martinell et al. 2006), plasmas in bounded domains (e.g. Agrawal 2002), geophysics (e.g. Iaffaldano et al. 2006), in economy and finance and in plasma turbulence (e.g. Del Castillo Negrete et al. 2004).

---

### 3 Symmetry and the Design of Networks

A particularly important step concerning intrinsic geodesy or geophysics was to find an appropriate network for the observation and the needed observables which obviously should concern Geomatics. A network intrinsic of the problem or, better, dictated by the problem and by the characteristics of the needed observable as suggested in the, perhaps only a seminal, work of Caputo (1979a, b) followed by the interesting and important meeting Optimization and Design of Geodetic Networks (Grafarend and Sansò 1985).

In accordance and with the help which may come from the principles of symmetry (Weyl 1983), I considered networks resulting from the *pflastersatz*.

For instance the selections of seismographic networks distributed among those suggested by the *pflastersatz* may improve the knowledge of the depth of earthquake which is so poorly known and also estimate the amount of lost information in the sites unreached by the network. At least such network would ensure homogeneity in the data collection.

Perhaps not surprisingly, from the distribution of the station depends the frequency of certain types of events and types of signals recorded.

One clear example is the survey of the stress tensor in the Phlaegrean Fields near Naples which showed surface deformations of the order of  $10^{-4}$  from which was inferred the possibility of an earthquake of magnitude 5 (Caputo 1979c).

I apologise for using my work for this example, not only because it was easier to dig out than someone else's work but the point is that knowing it better I can dig out its most important aspect here: my regret about it, is that I did not finish the work, the emblem of what I am trying to say. The underground is often rich, I may not say full, of cavities and the cavity effect (Neuber 1937; Harrison 1976), or the anomalies of the elastic materials may also affect the stress field which we are observing on the surface of the Earth and, therefore, we may not be extrapolated to the underground without additional information.

This in turn casts some doubts on the reliability of the estimate of the magnitude of an earthquake possibly due to the release of the estimated elastic energy stored in the local portion of the crust. In fact was also missing an accurate knowledge of the local values of the elastic parameters.

A stress field estimate based on data taken 40 years later would miss also the information on the balance between the relaxation of the elastic energy already stored and that accumulated in the additional 40 years. Obviously this is due to a lack of interdisciplinarity. Which instead, in a world of easy communication, as today, should be the emblem of most earth science research projects.

Moreover not sufficient attention was given to the large displacements of the order of half a meter over a distance of the order of a km observed in Pozzuoli and to the causes of this displacement and to the need of extending the observations at depth.

In other words since a surface network is not sufficient for a detailed and significant study of the inelastic field, additional data at depth is needed to model the important anomalies of the stress field at depth.

The same may be repeated for the networks of meteorological stations using the essential fact that temperature is not a physically meaningful parameter unless one associates to it the humidity and the pressure, that is at least the elevation of the station.

The time is ripe for a systematic, numerically quantitative, description of the state of health of regions by means of appropriate parameters which would allow comparisons between different regions and at different times. A description essential for the establishment of the priorities of the

interventions, of the detailed and of the global risk. It is the most appropriate task for Geomatics, which however should have some sort of interdisciplinarity.

The matter, to say it concisely, is to allow the description of physical objects with a limited number of parameters which, in turn, permits a direct, numerical, and therefore quantitative, comparison with its previous condition and with similar cases. The classic example is that of radioactivity where the most important property of the material is described by a single number: the average life directly related to the very simple linear differential equation of Malthus.

I am not suggesting to introduce "abstract equations" but those rightly defined as phenomenological.

These phenomenological equations, when adequately verified with experimental data, represent a step forward in respect to the usual empirical equations which are still very useful in many branches of applied science and technology. However some scientists resent the fact that they are not logically obtained from first principles ignoring that if we stick very strictly to first principles some types of progress are made very slow and difficult.

Also, when possible and when better methods are not available, we may use the method based on the numerical solution of the Cauchy problem for a stochastic differential equations (e.g. Caputo et al. 2000).

Moreover not only GPS data but also earthquake's data result from observation with many, sometimes different, instruments and, most important, at different places where the signals arrive after paths which are different in length but also travelled through different media. This, in turn, may cause remarkable effects on seismic risk and cause great difficulties in earthquake prediction. Here again interdisciplinarity is in order through the modern geodetic theories.

Rheology is one of the most efficient and less damaging mechanisms used by nature for its evolution to the inevitable attractor: the equilibrium. So memory is important for the planetary systems and down to include what is generically understood as *economy*.

However mathematical memory formalism of rheology and its irreversibility, changed all the equations of classical physics, from those of Maxwell to those of Navier-Stokes and of Fourier and general mechanics (Baleanu and Trujillo 2008).

Concerning anelasticity all attempts to use the equipartition of energy are subject to the restriction that the energy is subject to a decay with a frequency dependent mechanism. It would be surprising to find that, for all materials, the frequency dependent dissipation of energy be such to dissipate the energy of each degree of freedom in a way to preserve the equipartition. Finally most reciprocity theorems had to

be changed because of the second principle; an example at hand is that due to Graffi (1946) which was recently adjourned.

I wish to all of us good working days, pleasant weather and good stay in Roma.

## References

- Agrawal OP (2002) Solutions for a fractional diffusion wave equation defined in a bounded domain. *Nonlinear Dynam* 29(1–4):145–165. doi:[10.1023/A:1016539022492](https://doi.org/10.1023/A:1016539022492)
- Baleanu D, Trujillo JJ (2008) Fractional Euler-Lagrange and Hamilton equations within Caputo's fractional derivatives: a new perspective. In: 3rd IFAC workshop on fractional differentiation and its applications, 5–7 November. Ankara
- Baleanu D, Golmankhaneh AK, Nigmatullin RR (2009) Newtonian law with memory. *Nonlinear Dyn*. doi:[10.1007/s11071-009-9581-1](https://doi.org/10.1007/s11071-009-9581-1)
- Bocchio F (1974) From differential geodesy to differential geophysics. *Geophys J R Astron Soc* 39(1):1–10. doi:[10.1111/j.1365-246X.1974.tb05435.x](https://doi.org/10.1111/j.1365-246X.1974.tb05435.x)
- Caputo M (1979a) Topology and detection capability of seismic networks. In: Vogel A (ed) Proceedings of the international workshop on monitoring crustal dynamics in earthquake zones, Strasbourg 1978. Friedrich Vieweg & Sohn Verlag, Wiesbaden/Braunschweig, FRG
- Caputo M (1979b) Modern techniques and problems in monitoring horizontal strain. In: Vogel A (ed) Proceedings of the international workshop on monitoring crustal dynamics in earthquake zones, Strasbourg 1978. Friedrich Vieweg & Sohn Verlag, Wiesbaden/Braunschweig, FRG
- Caputo M (1979c) 2000 Years of geodetic and geophysical observations in the Phlegr. Fields near Naples. *Geophys J R Astron Soc* 56:319–328
- Caputo M (2012a) Reciprocity in elastic media with rheology. *Meccanica*. doi:[10.1007/s11012-013-9693-z](https://doi.org/10.1007/s11012-013-9693-z)
- Caputo M (2012b) The convergence of economic developments. *Non-linear Dyn Econometrics* 16(2):22. doi:[10.1515/1558-3708](https://doi.org/10.1515/1558-3708)
- Caputo M, Caputo R (2013) Mass anomalies and moments of inertia in the outer shells of the Earth. *Terranova* 1:1–10. doi:[10.1111/ter.12003](https://doi.org/10.1111/ter.12003)
- Caputo M, Ruggiero V, Sutera A, Zirilli F (2000) On the retrieval of water vapour profiles from a single GPS station. *Il Nuovo Cimento* 23(6):611–620
- Cesarone F, Caputo M, Cametti C (2004) Memory formalism in the passive diffusion across a biological membrane. *J Membr Sci* 250:79–84
- Del Castillo Negrete D, Carreras BA, Lynch VE (2004) Fractional diffusion in plasma turbulence. *Phys Plasmas* 11(8):3854–3864
- Demaria G (1976) I teoremi del punto fisso nell'analisi e nella sintesi economica. Applications of the fixed point theorem to economy, vol 43. Accademia Nazionale dei Lincei. Centro Linceo Interdisciplinare di Scienze Matematiche e loro applicazioni, pp 11–29
- El Shahed M (2003) A fractional calculus model of semilunar heart valve vibrations. In: International Mathematica symposium. Imperial College, London
- Grafarend EW (1988) The geometry of the earth's surface and the corresponding function space of the terrestrial gravitational field. *Deutsch. Geod. Komm. Bayer. Akad. Wiss., Reihe B, Heft Nr. 287*, pp 76–94
- Grafarend E, Sansò F (eds) (1985) Optimization and design of geodetic networks. Springer Verlag, Berlin
- Graffi D (1946) Sul teorema di reciprocità nella dinamica dei corpi elastici. *Memorie dell'Accademia delle scienze dell'Istituto di Bologna* 10(IV):103–109 (In Italian)
- Harrison JC (1976) Cavity and topographic effects in tilt and strain measurement. *J Geophys Res* 81(2):319–328. doi:[10.1029/JB081i002p00319](https://doi.org/10.1029/JB081i002p00319)
- Iaffaldano G, Caputo M, Martino S (2006) Experimental and theoretical memory diffusion of water in sand. *Hydrol Earth Syst Sci* 10:93–100
- Mainardi F (1996) Fractional relaxation-oscillation and fractional diffusion-wave phenomena. *Chaos Solitons Fractals* 7(9):1461–1477
- Martinell JJ, Del Castillo Negrete D, Raga AC, Williams DA (2006) Non-local diffusion and the chemical structure of molecular clouds. *Mon Not R Astron Soc* 372:213–225
- Milani N (1961) Sul significato intrinseco dell'equazione cui soddisfanno le famiglie di Lamé. In: Tesi di laurea in Scienze matematiche, relatori Marussi A. and Caputo M., Università degli Studi di Trieste, Anno Accademico 1960–61 (In Italian)
- Naber M (2004) Time fractional Schrödinger equation. *J Math Phys* 45(45):3339–3352
- Neuber H (1937) *Kerbespannungslehre*. Springer Verlag, Berlin
- Sansò F, Caputo M (2008) On the existence of Lamé orthogonal families adapted to a given potential function. *Boll Geodesia Scienze Affini LXVII(2):77–86*
- Slichter L, MacDonald GJF, Caputo M, Hager CL (1964) Report of earth tides results and of other gravity observations at Ucla. *Communications de l'Observatoire Royal de Belgique, 69, Série Geoph., VI. Marées Terrestres*, pp 124–131
- Tanni L (1940) On the continental undulations of the geoid as determined by the present gravity material. *Publ Isost Inst Int Assoc Geod*, No. 18
- Weyl H (1983) *Symmetry*. Princeton University Press, Princeton

---

# Fernando Sansò Laudation

Michele Caputo

You must know that in the last century the chairs of the Italian University *Prof* were loaded with dust because the Italian *Prof* was always away: chairing meetings, attending symposia, going to the capital city to seek subsidies for the research of his group, or attending the faculty meeting; some, as Marussi and Desio, went in scientific expeditions loaded with curiosity and of risk for their lives, I never heard that they had an insurance for this. The *Profs* were very rarely on their chair.

Fernando survived all this.

In this, as a teacher, Fernando has been very good and I like to quote an interesting result, which is not in his vita. It concerns one of his former students who, with a couple of colleagues, found that in the GPS is imbedded new type of seismograph, in other words they explored the GPS in a different frequency band and gave light to a new seismograph, they called it *VADASE*. And the credit for this is perhaps, in one of those mysterious ways of nature, to be added to the merits of Fernando.

And let me add that in these days we are gratified also by the inventions of a new shoe-box size seismograph called the *TREMINO* which may perform as those costing tens of thousand of dollars.

The invention of the *TREMINO* (S. Castellaro, M. Mucciarelli, F. Mulargia), and of the *VADASE* (G. Colosimo, M. Crespi, A. Mazzoni), let me tell you, are emblematic of the evolution in Italian Geodesy and Geophysics.

It is a cultural change as that from laser ranging for which was needed a big truck, as that of Contraves, to the GPS, from seismographs needing room size space to the *TREMINO* and of the *VADASE* which may stay in shoe size boxes.

The old generation of Italian geodesists across nineteenth and twentieth century, I am referring to that of Pizzetti-Somigliana theory, was mostly theoretical with limited

interested in the industry. In the middle of the twentieth there was some interest in the photogrammetric industry but apparently the Italian industry *missed the bus* and disappeared. Then we had the revolution of the satellites, the GPS; but the Italian contribution in this field, although of good and recognised quality, could not go very far.

The new generation of this century instead produces instruments. And part of this is also due to Fernando.

He assures us that he will continue to be busy. But I warn him from my experience, it may be because of the age, it may be because they say that the climate is changing or because of the electromagnetic pollution and that the glaciers are melting, we are bound to be more and more busy, simply because it is the only choice we have.

Simply because I do not know what else to do or I am not interested in anything but what I hope to be able to do well or I like to satisfy my curiosity. For instance asking where are those blessed roots which keep the Alps, where they are now for our games of all sorts and pleasures, or wondering if gravitons may be the cousins of neutrinos . . . .

If you believe that retirement means freedom I warn you since, beginning 30 years ago, I formally retired 3 times and the worse one was when it was not formal but practical when all sorts of *rasthaus* disappeared. I warn you because you will loose the recreation of the faculty meetings, the recreation of driving to the office and back, the recreation given by some colleagues momentarily free come to your office warning you that *it is for few minutes only*, you will loose the recreation of students asking questions, the recreation of hoping that one day we finally give some order to all the papers spread on the tables, the shelves on the walls, the chairs and floor of the office.

Often the retired person does not need to go to the office and finally may stay home to work where all escapes or recreations are gone because even the wife, who respects him so much, does not dare to interrupt his work and does not allow him to spread the papers all over the studio; because at that time the office is a *studio* to be kept in rigorous constant order.

---

M. Caputo (✉)  
Department of Geology and Geophysics, Texas A&M University,  
College Station, TX, USA  
e-mail: [mic.caputo27@gmail.com](mailto:mic.caputo27@gmail.com)

After retirement it will be endless work almost an obsession, with the fear that the curiosity will not have time to unravel all secrets of nature.

After this jocular perspective I congratulate Fernando for his life and his exceptional and innovative scientific production, for his devotion to science, for growing

and guiding such a numerous group of very good students. I also thank him for his purpose to stay around and work and I wish him long busy life in geodesy with us.

---

# Global Reference Systems: Theory and Open Questions

Athanasios Dermanis

---

## Abstract

In this review paper theoretical aspects of global reference systems are critically discussed in relation to their practical implementation through reference frames. These include the problem of the mathematical modeling of a spatiotemporal reference system for the deforming Earth, the relation of geodetic discrete-network reference systems to geophysical ones for the Earth mass continuum, the contribution of the various geodetic space techniques, the estimation issues related to the combination of the various data types, and issues relating to the compatibility of earth rotation representation. Finally issues related to future development of the International Terrestrial Reference Frame are discussed, concerning the addition of non-linear quasi-periodic terms in coordinate variation and their proper geophysical interpretation.

---

## Keywords

Earth rotation representation • Geodetic datum problem • Global geodetic networks • ITRF • Nonlinear station motions • Reference systems

---

## 1 Introduction

A reference system is merely a mathematical device within the framework of Newtonian mechanics that is conveniently used for the description of shape and its temporal deformation. It consists of a local basis  $\vec{e} = [\vec{e}_1 \ \vec{e}_2 \ \vec{e}_3]$  at a particular origin  $O$  and provides Cartesian coordinates  $\mathbf{x} = [x^1 \ x^2 \ x^3]^T$  of points  $P$  as the components of their position vectors  $\vec{x} = \overrightarrow{OP} = \vec{e} \mathbf{x}$ . In Earth related applications a (usually geocentric) global reference system separates the motion of Earth masses in space into the translational motion of the origin, the rotation of its axes around the origin (Earth rotation) and the apparent motion of Earth masses with respect to the reference system (Earth deformation).

Earth deformation forces geodesy to introduce a kinematics spatiotemporal concept of reference system, defined at every time epoch, which is much richer than the static spatial concept of classical mechanics. Truesdell and Noll (1965) who gave the axiomatic foundation of “rational” mechanics give such a limited concept:

... The position of an event can be specified only if a frame of reference, or observer, is given. Physically, a frame of reference is a set of objects whose mutual distances change comparatively little in time, like the walls of an observatory, the fixed stars, or the wooden horses on a merry-go-round. ...

The French astronomer Felix Tisserand (1845–1896) has realized the need of a spatiotemporal reference system and introduced the concept of what we now call Tisserand axes (Munk and MacDonald 1960). In his choice the temporal evolution of the orientation of the reference system axes is defined by minimizing the apparent motion of the point masses, quantified by the relative kinetic energy  $T_R = \frac{1}{2} \int_E \dot{\mathbf{x}}^T \dot{\mathbf{x}} \, dm = \min$ , which is secured by vanishing of the relative angular momentum components  $\mathbf{h}_R = \int_E [\mathbf{x} \times] \dot{\mathbf{x}} \, dm = \mathbf{0}$  (here dots denote differentiation

---

A. Dermanis (✉)  
Department of Geodesy and Surveying, Aristotle University of  
Thessaloniki, University Box 503, 54124 Thessaloniki, Greece  
e-mail: [dermanis@topo.auth.gr](mailto:dermanis@topo.auth.gr)

with respect to time,  $[\mathbf{a}\times]$  denotes the antisymmetric matrix with axial vector  $\mathbf{a}$ ,  $dm$  is the mass element and integration is carried over the whole Earth). The origin of Tisserand's reference system is the geocenter  $G$  with vanishing coordinates  $\mathbf{x}_G = \frac{1}{M} \int_E \mathbf{x} dm = \mathbf{0}$  ( $M = \text{Earth mass}$ ).

The problem of the definition of a reference system for a global geodetic network shows up in the formulation of a "reference frame" a term which in geodesy means the realization of a reference system by means of the coordinate functions  $\mathbf{x}_i(t)$  of a selected set of network points  $P_i$ . Coordinate time series provided by four fundamental space techniques VLBI, SLR, GPS and DORIS are utilized by the International Earth Rotation and Reference Systems Service (IERS) in order to realize the official International Terrestrial Reference Frame (ITRF) (Altamimi et al. 2002, 2004, 2007, 2011) and provide Earth Orientation Parameters (EOPs) describing the rotation of the Earth (Bizouard and Gambis 2009). Although the operational procedures for the ITRF formulation are now a matter of routine, there are still some recent advances as well as open problems in the theory of reference systems that may contribute to the improvement of the existing techniques. The present work is a short review of relevant results and a discussion of remaining problems for future investigations.

## 2 A Reference System Model for Geodetic Networks

The choice of reference system for a  $N$ -point three-dimensional geodetic network assigns to it a vector  $\mathbf{x} = [\dots \mathbf{x}_i^T \dots]^T$  of  $3N$  coordinates, which represents a point in  $R^{3N}$ . These coordinates are not the only ones that describe the network shape. Any arbitrary rigid transformation  $\tilde{\mathbf{x}}_i = \mathbf{R}(\boldsymbol{\theta}) \mathbf{x}_i + \mathbf{d}$  ( $\boldsymbol{\theta}$  being the rotation and  $\mathbf{d}$  the translation parameters) provides a vector  $\tilde{\mathbf{x}} \in R^{3N}$  that describes the same network shape. The submanifold  $M_{\tilde{\mathbf{x}}} = \left\{ \tilde{\mathbf{x}} \mid \tilde{\mathbf{x}}_i = \mathbf{R}(\boldsymbol{\theta}) \mathbf{x}_i + \mathbf{d}, \forall \boldsymbol{\theta}, \mathbf{d} \right\} \subset R^{3N}$  generated as  $\boldsymbol{\theta}$  and  $\mathbf{d}$  take all permissible values, is the set of all points corresponding to the same network shape. Thus  $M_{\tilde{\mathbf{x}}}$  is the shape manifold generated by  $\tilde{\mathbf{x}}$  (Dermanis 2000). Shape manifolds are naturally disjoint and through each point in  $R^{3N}$  passes only one manifold. Hence they constitute a *fibering* of an open subset of  $R^{3N}$ . The six parameters  $\boldsymbol{\theta}$ ,  $\mathbf{d}$  may serve as a set of coordinates on the six-dimensional shape manifold. For a deformable network,  $\mathbf{x}(t)$  is a curve in  $R^{3N}$  which represents the continuous time sequence of shape manifolds  $M_{\mathbf{x}(t)}$  corresponding to the shapes of the network at various time epochs  $t$ . We may define a reference system as a *section* of the shape manifold fibering, i.e. a curve

intersecting each manifold at one point. Each such curve  $\tilde{\mathbf{x}}(t)$  can be generated from the original curve  $\mathbf{x}(t)$  by means of six functions  $\boldsymbol{\theta}(t)$ ,  $\mathbf{d}(t)$  through  $\tilde{\mathbf{x}}_i(t) = \mathbf{R}(\boldsymbol{\theta}(t)) \mathbf{x}_i(t) + \mathbf{d}(t)$ . Different optimal reference systems are possible, depending on the arbitrary choice of  $\tilde{\mathbf{x}}(t_0)$ . The optimal choice  $\tilde{\mathbf{x}}(t)$ , is the shortest geodesic through  $\tilde{\mathbf{x}}(t_0)$  connecting the initial manifold  $M_{\tilde{\mathbf{x}}(t_0)}$  with the final one  $M_{\tilde{\mathbf{x}}(t_F)}$  for a time interval of interest  $t \in [t_0, t_F]$ . Such shortest geodesics are known to be perpendicular to both  $M_{\tilde{\mathbf{x}}(t_0)}$  and  $M_{\tilde{\mathbf{x}}(t_F)}$ . Since the choice of initial and final epoch is rather arbitrary,  $\tilde{\mathbf{x}}(t)$  should be perpendicular to any of the manifolds  $M_{\tilde{\mathbf{x}}(t)}$  that it crosses. Therefore the tangent vector  $\dot{\tilde{\mathbf{x}}}$  should be perpendicular to the tangent hyperplane of  $M_{\tilde{\mathbf{x}}(t)}$  at  $\tilde{\mathbf{x}}(t)$ , which is spanned by the coordinate base vectors  $\partial \tilde{\mathbf{x}}(t) / \partial q^i$ ,  $q^i$  been the elements of  $\boldsymbol{\theta}(t)$ ,  $\mathbf{d}(t)$ , i.e., by the columns of the matrix  $[\partial_{\boldsymbol{\theta}} \tilde{\mathbf{x}} \partial_{\mathbf{d}} \tilde{\mathbf{x}}]$ . Therefore the differential equations defining the optimal reference system  $\tilde{\mathbf{x}}(t)$  on the basis of a given arbitrary reference system  $\mathbf{x}(t)$  are  $[\partial_{\boldsymbol{\theta}} \tilde{\mathbf{x}} \partial_{\mathbf{d}} \tilde{\mathbf{x}}]^T \dot{\tilde{\mathbf{x}}} = \mathbf{0}$ . With  $\mathbf{x}(t)$  chosen to be barycentric ( $\frac{1}{M} \sum_i \mathbf{x}_i(t) = \mathbf{0}$ ) we arrive at the differential equations

$$\mathbf{W}(\boldsymbol{\theta}) \dot{\boldsymbol{\theta}} = \mathbf{C}^{-1} \mathbf{h}, \quad \dot{\mathbf{d}} = \mathbf{0}. \quad (1)$$

where  $\mathbf{W} = [\mathbf{w}_1 \mathbf{w}_2 \mathbf{w}_3]$  with  $\mathbf{w}_k$  being the axial vectors of the antisymmetric matrices  $[\mathbf{w}_k \times] = (\partial_{\theta_k} \mathbf{R}^T) \mathbf{R}$ ,  $\mathbf{h} = \sum_i [\mathbf{x}_i \times] \dot{\mathbf{x}}_i$  is the discrete relative angular momentum in the initial reference system and  $\mathbf{C} = -\sum_i [\mathbf{x}_i \times]^2$  is the discrete inertia matrix. The solution to the above equations is not unique but depends on integration constants  $\boldsymbol{\theta}(t_0)$ ,  $\mathbf{d}(t_0)$  or  $\tilde{\mathbf{x}}(t_0)$ . Any two solutions generate corresponding reference system curves  $\tilde{\mathbf{x}}(t)$ ,  $\tilde{\mathbf{x}}'(t)$  which are "parallel" in the sense that they are connected by a time-independent rigid transformation  $\tilde{\mathbf{x}}'_i(t) = \mathbf{Q} \tilde{\mathbf{x}}_i(t) + \mathbf{c}$ ,  $(\mathbf{Q}, \mathbf{c}) = \text{constant}$ .

Passing from an  $N$ -point discrete network to the continuous Earth body, calls for the replacement of  $R^{3N}$  with the infinite-dimensional space of the coordinates of all material points of the Earth, but the intricacies of a corresponding rigorous mathematical model are far from trivial. In any case the fibering by six-dimensional shape manifolds having the transformation parameters  $\boldsymbol{\theta}$ ,  $\mathbf{d}$  as coordinates is preserved.

It is interesting to see how a given geocentric (i.e. barycentric for all Earth points) reference system  $\mathbf{x}(t)$  can be transformed by a point-wise rigid transformation  $\tilde{\mathbf{x}}(t) = \mathbf{R}(\boldsymbol{\theta}(t)) \mathbf{x}(t) + \mathbf{d}(t)$  into a geocentric Tisserand reference system  $\tilde{\mathbf{x}}(t)$ , satisfying  $\frac{1}{M} \int_E \tilde{\mathbf{x}} dm = \mathbf{0}$  and  $\dot{\tilde{\mathbf{h}}} = \int_E [\tilde{\mathbf{x}} \times] \dot{\tilde{\mathbf{x}}} dm = \mathbf{0}$ . Carrying out the necessary computations we arrive at  $\dot{\mathbf{d}} = \mathbf{0}$  (which is one of the solutions of  $\dot{\mathbf{d}} = \mathbf{0}$ ) while  $\boldsymbol{\theta}$  satisfies again the differential equation  $\mathbf{W}(\boldsymbol{\theta}) \dot{\boldsymbol{\theta}} = \mathbf{C}^{-1} \mathbf{h}$ , with the only difference that the relative angular momentum and inertia matrix are given in this case by  $\mathbf{h} = \int_E [\mathbf{x} \times] \dot{\mathbf{x}} dm$  and  $\mathbf{C} = -\int_E [\mathbf{x} \times]^2 dm$ .



### 3 Definition of the Reference System in the ITRF Formulation

The static version of the adoption of a reference system for non-deforming networks is an old geodetic problem mostly known as the “geodetic datum problem” that emerged in the so called “free networks”, i.e. local networks which do not inherit their reference system from a pre-existing higher order network. It has also given rise to geodetic contributions to the statistical linear estimation theory without full rank, in relation to the linear(ized) model  $\mathbf{b} = \mathbf{A}\mathbf{x} + \mathbf{v}$ ,  $\mathbf{v} \sim (\mathbf{0}, \sigma^2 \mathbf{P}^{-1})$  with  $n$  observations  $\mathbf{b}$ ,  $m$  unknowns  $\mathbf{x}$  and  $\text{rank} \mathbf{A} = r < m$ . The rank deficiency and the corresponding infinity of least squares ( $\mathbf{v}^T \mathbf{P} \mathbf{v} = \min$ ) solutions for the unknown parameters is due to the use of coordinates as unknowns, while observations can only determine the geometric figure of the network. However all least-squares solutions lead to the same values for observable quantities as well as all the functions of the observables which are statistically characterized as estimable quantities. A unique solution is obtained by posing additional constraints  $\mathbf{C}^T \mathbf{x} = \mathbf{d}$  on the unknowns, which are minimal i.e. they resolve the coordinate indeterminacy without affecting the estimated geometric figure of the network. The coordinate estimates and their singular covariance matrices merely serve as a depository of information for the further computation of estimates of estimable quantities and their covariance matrices. The role of minimal constraints is that of assigning an arbitrary reference system so that coordinates can be computed. Particularly popular have been the so called inner constraints  $\mathbf{E}^T \mathbf{x} = \mathbf{0}$ , which satisfy  $\mathbf{x}^T \mathbf{x} = \min$  among all least squares solutions. The matrix  $\mathbf{E}$  results from the coordinate transformation  $\mathbf{x} \rightarrow \tilde{\mathbf{x}} = T(\mathbf{x}, \mathbf{p})$  under a change of the reference system, where  $\mathbf{p}$  are transformation parameters such that  $T(\mathbf{x}, \mathbf{0}) = \mathbf{x}$ . The linearized form of the transformation  $\tilde{\mathbf{x}} = \mathbf{x} + \mathbf{E}\mathbf{p}$  allows the determination of the desired matrix  $\mathbf{E}$ . Sometimes the total inner constraints  $\mathbf{E}^T \mathbf{x} = [\mathbf{E}_1^T \ \mathbf{E}_2^T] \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{bmatrix} = \mathbf{0}$  are replaced by partial constraints  $\mathbf{E}_1^T \mathbf{x}_1 = \mathbf{0}$  involving only a subset  $\mathbf{x}_1$  of the parameters and satisfying  $\mathbf{x}_1^T \mathbf{x}_1 = \min$ , instead.

For deformable ITRF network the choice of reference system is dominated by one extension, that of defining its temporal evolution and two restrictions. The first is the restriction to coordinate transformations “close to the identity” i.e. with very small transformation parameters  $\mathbf{p}(t)$ , which allow the replacement of  $\tilde{\mathbf{x}}(t) = (1 + s(t)) \mathbf{R}(\boldsymbol{\theta}(t)) \mathbf{x}(t) + \mathbf{d}(t)$ , with the linear approximation  $\tilde{\mathbf{x}}(t) \approx \mathbf{x}(t) + s(t)\mathbf{x}(t) + [\mathbf{x}(t) \times] \boldsymbol{\theta}(t) + \mathbf{d}(t)$  realized by  $\mathbf{R}(\boldsymbol{\theta}) \approx \mathbf{I} - [\boldsymbol{\theta} \times]$  and neglectation of second and higher order terms. More severe is the second restriction to

transformations which preserve the linear-in-time form of the ITRF coordinate model  $\mathbf{x}_i = \mathbf{x}_{i0} + (t - t_0) \mathbf{v}_i \equiv \mathbf{x}_{i0} + \Delta t \mathbf{v}_i$ , where  $\mathbf{x}_{i0}$ ,  $\mathbf{v}_i$  are the initial coordinates and constant velocities of the network point  $P_i$ . This necessitates the use of transformation parameters that are also linear in time, i.e.,  $s(t) = s_0 + \Delta t \dot{s}$ ,  $\boldsymbol{\theta}(t) = \boldsymbol{\theta}_0 + \Delta t \dot{\boldsymbol{\theta}}$ ,  $\mathbf{d}(t) = \mathbf{d}_0 + \Delta t \dot{\mathbf{d}}$ , which effectively restricts the transformation parameters to the 14 parameter set  $\mathbf{p} = [s_0 \ \boldsymbol{\theta}_0^T \ \mathbf{d}_0^T \ \dot{s} \ \dot{\boldsymbol{\theta}}^T \ \dot{\mathbf{d}}^T]^T$ . The use of minimal or inner constraints on the parameters  $\mathbf{x}$ , which due to the linearization are corrections to approximate values of the unknowns, have the disadvantage that they depend on the choice of the approximate values. An even more serious disadvantage of such “algebraic” constraints is that they have no clear physical meaning and they do not lead to a choice of reference system which is “optimal” in a physically meaningful way. A different type of physically meaningful kinematic constraints have been proposed by Altamimi and Dermanis (2009), which minimize the apparent motion of network points with respect to the reference system. They are based on a discrete version of Tisserand’s ideas where network points are treated as mass points of unit mass. The main idea is to minimize the network’s discrete relative kinetic energy  $T(t) = \frac{1}{2} \sum_i \dot{\tilde{\mathbf{x}}}_i^T(t) \dot{\tilde{\mathbf{x}}}_i(t)$  at every epoch  $t$  or equivalently to nullify the relative angular momentum  $\mathbf{h}(t) = \sum_i [\mathbf{x}_i(t)] \dot{\tilde{\mathbf{x}}}_i(t) = \mathbf{0}$ , a relation which establishes the temporal evolution of the orientation of the reference system. The origin of the system is defined by setting constant the coordinates of its barycenter  $\mathbf{x}_B(t) = \frac{1}{N} \sum_i \mathbf{x}_i(t) = \mathbf{x}_B(t_0)$ , a particular choice being  $\mathbf{x}_B(t_0) = \mathbf{0}$  (barycentric system). Scale is taken care by setting constant the mean quadratic scale of the network  $S(t) = Q(t)^{1/2}$  defined by  $Q(t) = \frac{1}{N} \sum_i [\mathbf{x}_i(t) - \mathbf{x}_B(t)]^T [\mathbf{x}_i(t) - \mathbf{x}_B(t)] = Q(t_0) = \text{const}$ . Applied to the ITRF model  $\mathbf{x}_i = \mathbf{x}_{i0} + \Delta t \mathbf{v}_i$  the kinematic constraints become

$$\sum_i [\mathbf{x}_{i0} \times] \mathbf{v}_i = \mathbf{0} \quad (2)$$

$$\frac{1}{N} \sum_i \mathbf{x}_{i0} = \mathbf{x}_B(t_0) \quad (3)$$

$$\frac{1}{N} \sum_i \mathbf{v}_i = \mathbf{0} \quad (4)$$

$$\frac{1}{N} \sum_i (\mathbf{x}_{i0} - \bar{\mathbf{x}}_0)^T (\mathbf{x}_{i0} - \bar{\mathbf{x}}_0) = Q(t_0), \quad (5)$$

$$\frac{1}{N} \sum_i (\mathbf{x}_{i0} - \bar{\mathbf{x}}_0)^T (\mathbf{v}_i - \bar{\mathbf{v}}) = \mathbf{0} \quad (6)$$

where  $\bar{\mathbf{x}}_0 = \frac{1}{N} \sum_i \mathbf{x}_{i0}$  and  $\bar{\mathbf{v}} = \frac{1}{N} \sum_i \mathbf{v}_i$ . Operational expressions in terms of corrections  $\delta \mathbf{x}_{i0} = \mathbf{x}_{i0} - \mathbf{x}_{i0}^{ap}$ ,  $\delta \mathbf{v}_i = \mathbf{v}_i - \mathbf{v}_i^{ap}$  to approximate values  $\mathbf{x}_{i0}^{ap}$ ,  $\mathbf{v}_i^{ap}$  of the initial coordinates and velocities can be found in Altamimi and Dermanis (2009).

Kinematic constraints define only the evolution of the reference system with respect to orientation (2), origin (4) and scale (6). Initial epoch constraints are either missing (for orientation) or depend on the arbitrary constants  $\mathbf{x}_B(t_0)$  for origin (3) and  $Q(t_0)$  for scale (5). The arbitrary choice of the initial epoch orientation, origin and scale leads to different but equivalent ‘‘parallel’’ reference systems, i.e., realized by coordinate functions  $\mathbf{x}_i(t)$ ,  $\tilde{\mathbf{x}}_i(t)$ , which are at any epoch related by a time-independent similarity transformation  $\tilde{\mathbf{x}}_i(t) = (1 + \lambda) \mathbf{Q} \mathbf{x}_i(t) + \mathbf{t}$  with constant  $\lambda$ ,  $\mathbf{Q}$  and  $\mathbf{t}$ . The lack of initial orientation is inherent in Tisserand reference systems; initial origin may result by selecting a barycentric  $\mathbf{x}_B(t) = \mathbf{x}_B(t_0) = \mathbf{0}$  or a geocentric one since the geocenter is an additional network point in SLR observations. The presence of the scale parameter in the coordinate transformations does not actually correspond to a deficiency but rather to the fact that different space techniques have a different unit of length. This is theoretically due to the use of different units of time realized by different sets of clocks, while additional effects come from systematic errors (tropospheric corrections, phase center corrections for satellite and ground antennas, etc.). In the past (ITRF 2005) ITRF scale was based on VLBI only, currently (ITRF 2008) a weighted combination of VLBI and SLR is used, with expected future contributions from GNSS.

#### 4 ITRF Formulation: The One-Step and the Two-Step Approach

There are two basic approaches for the formulation of the ITRF the one-step (Angermann et al. 2004; Rothacher et al. 2011; Seitz et al. 2012) and the two-step approach (Altamimi et al. 2002, 2004, 2007, 2011). They both use as pseudo-observations coordinates series estimates from space techniques

$$\mathbf{b}_T = \mathbf{A}_T \mathbf{x}_T + \mathbf{v}_T, \quad \mathbf{v}_T \sim (\mathbf{0}, \sigma^2 \mathbf{P}_T^{-1}), \quad T = V, S, G, D \quad (7)$$

(V = VLBI, S = SLR, G = GPS, D = DORIS) using the standard assumptions of the Gauss–Markov model for zero mean errors with known covariance matrices up to a scalar factor. Since the networks of different techniques are distinct, additional observations

$$\mathbf{b}_c = \mathbf{C}_V \mathbf{x}_V + \mathbf{C}_S \mathbf{x}_S + \mathbf{C}_G \mathbf{x}_G + \mathbf{C}_D \mathbf{x}_D + \mathbf{v}_c. \quad (8)$$

with weight matrix  $\mathbf{P}_c$  are utilized, which connect stations of different techniques at co-location points.

The one-step approach proceeds to the formulation of the normal equations for all data

$$\mathbf{N} \hat{\mathbf{x}} = \begin{bmatrix} \mathbf{N}_V + \mathbf{N}_{VV}^c & \mathbf{N}_{VS}^c & \mathbf{N}_{VG}^c & \mathbf{N}_{VD}^c \\ (\mathbf{N}_{VS}^c)^T & \mathbf{N}_S + \mathbf{N}_{SS}^c & \mathbf{N}_{SG}^c & \mathbf{N}_{SD}^c \\ (\mathbf{N}_{VG}^c)^T & (\mathbf{N}_{SG}^c)^T & \mathbf{N}_G + \mathbf{N}_{GG}^c & \mathbf{N}_{GD}^c \\ (\mathbf{N}_{VD}^c)^T & (\mathbf{N}_{SD}^c)^T & (\mathbf{N}_{GD}^c)^T & \mathbf{N}_D + \mathbf{N}_{DD}^c \end{bmatrix} \times \begin{bmatrix} \hat{\mathbf{x}}_V \\ \hat{\mathbf{x}}_S \\ \hat{\mathbf{x}}_G \\ \hat{\mathbf{x}}_D \end{bmatrix} = \begin{bmatrix} \mathbf{u}_V + \mathbf{u}_V^c \\ \mathbf{u}_S + \mathbf{u}_S^c \\ \mathbf{u}_G + \mathbf{u}_G^c \\ \mathbf{u}_D + \mathbf{u}_D^c \end{bmatrix} = \mathbf{u} \quad (9)$$

where  $\mathbf{N}_T = \mathbf{A}_T^T \mathbf{P}_T \mathbf{A}_T$ ,  $\mathbf{u}_T = \mathbf{A}_T^T \mathbf{P}_T \mathbf{b}_T$ ,  $\mathbf{N}_{TT'}^c = \mathbf{C}_T^T \mathbf{P}_c \mathbf{C}_{T'}$ ,  $\mathbf{u}_T^c = \mathbf{C}_T^T \mathbf{P}_c \mathbf{b}_c$ . Within the data of each technique  $\mathbf{b}_T = \mathbf{A}_T \mathbf{x}_T + \mathbf{v}_T$  there is an inherent rank deficiency due to their inability to determine a reference system to which the unknown coordinates refer. This is expressed as a deficiency in the (column) rank of the design matrix  $\mathbf{A}_T$ . If  $\tilde{\mathbf{x}}_T = \mathbf{x}_T + \mathbf{E}_T \mathbf{p}_T$  is the result of a coordinate transformation with transformation parameters  $\mathbf{p}_T$ , due to a change of the reference system, then both  $\mathbf{x}_T$  and  $\tilde{\mathbf{x}}_T$  yield the same value for the invariant observables  $\mathbf{y}_T = \mathbf{A}_T \mathbf{x}_T = \mathbf{A}_T \tilde{\mathbf{x}}_T$ . This means that  $\mathbf{A}_T \mathbf{E}_T = \mathbf{0}$  and consequently

$$\mathbf{N}_T \mathbf{E}_T = \mathbf{0}, \quad T = V, S, G, D. \quad (10)$$

No matter how the one-step approach is operationally realized, it can be shown for the sake of comparison with the two-step approach to be equivalent to a modified two-part approach (Dermanis 2011). The first part is identical to the first step of the two-step approach but it is also used in the one-step approach for preprocessing (identification of outliers and discontinuities). It involves the formulation of the separate normal equations for each technique  $\mathbf{N}_V \hat{\mathbf{x}}_V^s = \mathbf{u}_V$ ,  $\mathbf{N}_S \hat{\mathbf{x}}_S^s = \mathbf{u}_S$ ,  $\mathbf{N}_G \hat{\mathbf{x}}_G^s = \mathbf{u}_G$ ,  $\mathbf{N}_D \hat{\mathbf{x}}_D^s = \mathbf{u}_D$  and estimate computation using separate minimal constraints. The second part should replace the second step of the two-step approach in order to secure identical results with the straightforward one-step approach. It involves a least squares adjustment of the following set of uncorrelated observation equations

$$\hat{\mathbf{x}}_V^s = \mathbf{x}_V + \mathbf{e}_V, \quad \hat{\mathbf{x}}_S^s = \mathbf{x}_S + \mathbf{e}_S,$$

$$\hat{\mathbf{x}}_G^s = \mathbf{x}_G + \mathbf{e}_G, \quad \hat{\mathbf{x}}_D^s = \mathbf{x}_D + \mathbf{e}_D,$$

$$\mathbf{b}_c = \mathbf{C}_V \mathbf{x}_V + \mathbf{C}_S \mathbf{x}_S + \mathbf{C}_G \mathbf{x}_G + \mathbf{C}_D \mathbf{x}_D + \mathbf{v}_c \quad (11)$$

with corresponding weight matrices:  $\mathbf{N}_V$ ,  $\mathbf{N}_S$ ,  $\mathbf{N}_G$ ,  $\mathbf{N}_D$  and  $\mathbf{P}_c$ . The corresponding normal equations

$$\begin{bmatrix} \mathbf{N}_V + \mathbf{N}_{VV}^c & \mathbf{N}_{VS}^c & \mathbf{N}_{VG}^c & \mathbf{N}_{VD}^c \\ (\mathbf{N}_{VS}^c)^T & \mathbf{N}_S + \mathbf{N}_{SS}^c & \mathbf{N}_{SG}^c & \mathbf{N}_{SD}^c \\ (\mathbf{N}_{VG}^c)^T & (\mathbf{N}_{SG}^c)^T & \mathbf{N}_G + \mathbf{N}_{GG}^c & \mathbf{N}_{GD}^c \\ (\mathbf{N}_{VD}^c)^T & (\mathbf{N}_{SD}^c)^T & (\mathbf{N}_{GD}^c)^T & \mathbf{N}_D + \mathbf{N}_{DD}^c \end{bmatrix} \begin{bmatrix} \widehat{\mathbf{x}}_V \\ \widehat{\mathbf{x}}_S \\ \widehat{\mathbf{x}}_G \\ \widehat{\mathbf{x}}_D \end{bmatrix} = \begin{bmatrix} \mathbf{N}_V \widehat{\mathbf{x}}_V^s + \mathbf{u}_V^c \\ \mathbf{N}_S \widehat{\mathbf{x}}_S^s + \mathbf{u}_S^c \\ \mathbf{N}_G \widehat{\mathbf{x}}_G^s + \mathbf{u}_G^c \\ \mathbf{N}_D \widehat{\mathbf{x}}_D^s + \mathbf{u}_D^c \end{bmatrix} \quad (12)$$

are readily seen to be identical to those in a single step (Eq. (9)) since the separate estimates satisfy  $\mathbf{N}_V \widehat{\mathbf{x}}_V^s = \mathbf{u}_V$ ,  $\mathbf{N}_S \widehat{\mathbf{x}}_S^s = \mathbf{u}_S$ ,  $\mathbf{N}_G \widehat{\mathbf{x}}_G^s = \mathbf{u}_G$  and  $\mathbf{N}_D \widehat{\mathbf{x}}_D^s = \mathbf{u}_D$ .

In the two-step approach the first one (stacking per technique) is identical to the above first part of the one-step approach. In the second step however it is recognized that each of the separate estimates  $\widehat{\mathbf{x}}_V^s$ ,  $\widehat{\mathbf{x}}_S^s$ ,  $\widehat{\mathbf{x}}_G^s$ ,  $\widehat{\mathbf{x}}_D^s$  refers to separate reference systems which also differ from the final ITRF reference system of the sought estimates  $\widehat{\mathbf{x}}_V$ ,  $\widehat{\mathbf{x}}_S$ ,  $\widehat{\mathbf{x}}_G$ ,  $\widehat{\mathbf{x}}_D$ . For this reason transformation parameters are included in the model which becomes

$$\mathbf{b} = \begin{bmatrix} \widehat{\mathbf{x}}_V^s \\ \widehat{\mathbf{x}}_S^s \\ \widehat{\mathbf{x}}_G^s \\ \widehat{\mathbf{x}}_D^s \\ \mathbf{b}_c \end{bmatrix} = \begin{bmatrix} \mathbf{x}_V + \mathbf{E}_V \mathbf{p}_V \\ \mathbf{x}_S + \mathbf{E}_S \mathbf{p}_S \\ \mathbf{x}_G + \mathbf{E}_G \mathbf{p}_G \\ \mathbf{x}_D + \mathbf{E}_D \mathbf{p}_D \\ \mathbf{C}_V \mathbf{x}_V + \mathbf{C}_S \mathbf{x}_S + \mathbf{C}_G \mathbf{x}_G + \mathbf{C}_D \mathbf{x}_D \end{bmatrix} + \begin{bmatrix} \mathbf{e}_V \\ \mathbf{e}_S \\ \mathbf{e}_G \\ \mathbf{e}_D \\ \mathbf{v}_c \end{bmatrix} = \bar{\mathbf{A}} \bar{\mathbf{x}} + \mathbf{v} = [\mathbf{A} \ \mathbf{E}] \begin{bmatrix} \mathbf{x} \\ \mathbf{p} \end{bmatrix} + \mathbf{v} \quad (13)$$

where  $\mathbf{A}$  is the same as for (11),  $\mathbf{p} = [\mathbf{p}_V^T \ \mathbf{p}_S^T \ \mathbf{p}_G^T \ \mathbf{p}_D^T]^T$  and

$$\mathbf{E} = \begin{bmatrix} \mathbf{E}_V & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{E}_S & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{E}_G & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{E}_D \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} \end{bmatrix}. \quad (14)$$

The corresponding normal equations take in this case the extended form  $\bar{\mathbf{N}} \bar{\mathbf{x}} = \bar{\mathbf{u}}$ , where  $\bar{\mathbf{x}} = [\widehat{\mathbf{x}}^T \ \widehat{\mathbf{p}}^T]^T$ ,

$$\bar{\mathbf{N}} = \bar{\mathbf{A}}^T \mathbf{P} \bar{\mathbf{A}} = \begin{bmatrix} \mathbf{A}^T \mathbf{P} \mathbf{A} & \mathbf{A}^T \mathbf{P} \mathbf{E} \\ \mathbf{E}^T \mathbf{P} \mathbf{A} & \mathbf{E}^T \mathbf{P} \mathbf{E} \end{bmatrix} = \begin{bmatrix} \mathbf{N} & \mathbf{A}^T \mathbf{P} \mathbf{E} \\ \mathbf{E}^T \mathbf{P} \mathbf{A} & \mathbf{E}^T \mathbf{P} \mathbf{E} \end{bmatrix}$$

$$\bar{\mathbf{u}} = \bar{\mathbf{A}}^T \mathbf{P} \mathbf{b} = \begin{bmatrix} \mathbf{A}^T \mathbf{P} \mathbf{b} \\ \mathbf{E}^T \mathbf{P} \mathbf{b} \end{bmatrix} = \begin{bmatrix} \mathbf{u} \\ \mathbf{E}^T \mathbf{P} \mathbf{b} \end{bmatrix}. \quad (15)$$

However according to (10) it holds that

$$\mathbf{P} \mathbf{E} = \begin{bmatrix} \mathbf{N}_V \mathbf{E}_V & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{N}_S \mathbf{E}_S & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{N}_G \mathbf{E}_G & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{N}_D \mathbf{E}_D \end{bmatrix} = \mathbf{0} \quad (16)$$

Consequently the normal equations degenerate into

$$\bar{\mathbf{N}} \bar{\mathbf{x}} = \begin{bmatrix} \mathbf{N} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \begin{bmatrix} \widehat{\mathbf{x}} \\ \widehat{\mathbf{p}} \end{bmatrix} = \bar{\mathbf{u}} = \begin{bmatrix} \mathbf{u} \\ \mathbf{0} \end{bmatrix} \quad (17)$$

i.e. into the normal equations  $\mathbf{N} \widehat{\mathbf{x}} = \mathbf{u}$  of the second part of the one-step approach and  $\mathbf{0} \widehat{\mathbf{p}} = \mathbf{0}$ , which does not allow the determination of the transformation parameter estimate  $\widehat{\mathbf{p}}$ . In conclusion, the second step of the two-step approach needs to be replaced by a proper combination at solution level (as opposed to the combination at the normal equation level of the one-step approach), in order to secure identical results. For a comparison of these two approaches as currently applied see Appendix B of Angermann et al. (2004).

## 5 Relating the Reference System of a Global Geodetic Network to a Tisserand Reference System

From a geophysical point of view an “optimal” reference system for a global network falls short in representing the deformation of the Earth, or just the lithosphere, even when established by kinematic constraints which minimize the apparent motion of station points. For this purpose the motion of the masses of the lithosphere must be approximately inferred with the help of a geophysical model. Such a widely accepted model is that of rotating tectonic plates where within plate deformations play a minor role and may be ignored in a first approximation. If on each plate or subplate  $P_K$  lies a subnetwork  $D_K$  of the global geocentric network, information of the coordinate variation  $\mathbf{x}_i(t)$ ,  $i \in D_K$  can be used to deduce the rotation vector  $\boldsymbol{\omega}_K$  of each plate and its contribution  $\tilde{\mathbf{h}}_{P_K}$  to the relative angular momentum of the lithosphere  $\tilde{\mathbf{h}} = \sum_K \tilde{\mathbf{h}}_{P_K}$ . Setting  $\tilde{\mathbf{h}} = \mathbf{0}$  allows the determination of the rotation vector  $\boldsymbol{\omega}$  which transforms the original geocentric reference system into a Tisserand geocentric one for the lithosphere. Such a system is appropriate for comparing its observed rotation with that predicted by theories of Earth rotation. The computation algorithm consists of the following: For each subnetwork  $D_K$  the discrete matrix of inertia  $\mathbf{C}_{D_K} = -\sum_{i \in D_K} [\mathbf{x}_i \times]^2$  and the relative angular momentum  $\mathbf{h}_{D_K} = \sum_{i \in D_K} [\mathbf{x}_i \times] \dot{\mathbf{x}}_i$  with

respect to the geocenter are used to compute the rotation vector  $\boldsymbol{\omega}_K = \mathbf{C}_{D_K}^{-1} \mathbf{h}_{D_K}$  of the corresponding plate  $P_K$ . The contribution of each plate  $P_K$  to the matrix of inertia of the lithosphere is computed as  $\mathbf{C}_{P_K} = -\int_{P_K} [\mathbf{x} \times]^2 dm$  in order to compute the rotation vector

$$\boldsymbol{\omega} = \left( \sum_K \mathbf{C}_{P_K} \right)^{-1} \left( \sum_K \mathbf{C}_{P_K} \boldsymbol{\omega}_K \right) \quad (18)$$

from the original reference system  $\{\mathbf{x}_i\}$  to the Tisserand reference system of the lithosphere  $\{\tilde{\mathbf{x}}_i\}$ . Solving the generalized Euler differential equations  $[\boldsymbol{\omega} \times] = -\mathbf{R}^T \dot{\mathbf{R}} = \mathbf{R}^T \mathbf{R}$ , the parameters  $\boldsymbol{\theta}$  of the rotation matrix  $\mathbf{R}(\boldsymbol{\theta})$  are determined and finally the coordinates of the global network are converted according to  $\tilde{\mathbf{x}}_i = \mathbf{R} \mathbf{x}_i$ . The specific Tisserand system, out of infinite “parallel” ones having the same temporal evolution, depends on the chosen initial values  $\boldsymbol{\theta}(t_0)$ . It is also possible to use a model with arbitrary rigid plate motion instead of simple rotation or even to incorporate internal plate deformations deduced from the station motion of the corresponding subnetwork. It must be noted that the computation of the inertia matrices  $\mathbf{C}_{P_K}$  requires knowledge of the geometric boundaries of each plate or subplate as well as of its internal density distribution.

## 6 Compatibility of Earth Rotation Representation

The official IERS representation (Petit and Luzum 2010) of the rotation matrix  $\mathbf{R}$  converting celestial to terrestrial coordinates ( $\mathbf{x}_T = \mathbf{R} \mathbf{x}_C$ ) has the form

$$\mathbf{R} = \mathbf{W} \mathbf{D} \mathbf{Q} = \mathbf{R}_3(-F) \mathbf{R}_2(-g) \mathbf{R}_3(F + s') \mathbf{R}_3(\theta) \times \mathbf{R}_3(-E - s) \mathbf{R}_2(d) \mathbf{R}_3(E) \quad (19)$$

where  $\mathbf{Q} = \mathbf{R}_3(-E - s) \mathbf{R}_2(d) \mathbf{R}_3(E) = \mathbf{R}_3(-s) \mathbf{G}(X, Y)$  is the precession–nutation matrix,  $\mathbf{D} = \mathbf{R}_3(\theta)$  is the diurnal rotation matrix,  $\mathbf{W} = \mathbf{R}_3(-F) \mathbf{R}_2(-g) \mathbf{R}_3(F + s') = \mathbf{G}^T(\xi, \eta) \mathbf{R}_3(s') \approx \mathbf{R}_1(-y_P) \mathbf{R}_2(-x_P) \mathbf{R}_3(s')$  is the polar motion matrix, while  $x_P \approx \xi$  and  $y_P \approx -\eta$  are the coordinates of the pole. Two intermediate reference systems are the celestial intermediate system  $\vec{\mathbf{e}}^{IC} = \vec{\mathbf{e}}^C \mathbf{Q}^T$  and the intermediate terrestrial one  $\vec{\mathbf{e}}^{IT} = \vec{\mathbf{e}}^T \mathbf{W}$ , which are related by  $\vec{\mathbf{e}}^{IT} = \vec{\mathbf{e}}^{IC} \mathbf{D}^T = \vec{\mathbf{e}}^{IC} \mathbf{R}_3(-\theta)$  having a common 3rd axis along  $\vec{\mathbf{p}} = \vec{\mathbf{e}}_3^{IC} = \vec{\mathbf{e}}_3^{IT} = \vec{\mathbf{e}}^C \mathbf{p}_C = \vec{\mathbf{e}}^T \mathbf{p}_T$  the unit vector in the direction of the celestial intermediate pole (CIP), with celestial components  $\mathbf{p}_C = [X \ Y \ Z]^T$  and terrestrial ones  $\mathbf{p}_T = [\xi \ \eta \ \zeta]^T$ . The original 7 parameters (functions of time) are reduced to 5 by means of the 2 NRO (Non Rotating Origin) conditions (Capitaine et al. 1986)  $s = s(d, E) = s(X, Y)$  and  $s' = s'(g, F) = s(\xi, \eta) = s(x_P, -y_P)$  which

define the directions of the TIO (Terrestrial Intermediate Origin)  $\vec{\mathbf{e}}_1^{IT}$  and the CIO (Celestial Intermediate Origin)  $\vec{\mathbf{e}}_1^{IC}$ . For a precise definition of the NRO conditions we need the concept of the relative rotation vector  $\vec{\boldsymbol{\omega}}_{A \rightarrow B} = \vec{\mathbf{e}}^A \boldsymbol{\omega}_A = \vec{\mathbf{e}}^B \boldsymbol{\omega}_B$  between two reference systems connected by  $\vec{\mathbf{e}}^B = \vec{\mathbf{e}}^A \mathbf{R}_{A \rightarrow B}^T$  having components determined from the generalized Euler kinematic equations  $[\boldsymbol{\omega}_A \times] = \mathbf{R}_{A \rightarrow B} \dot{\mathbf{R}}_{A \rightarrow B}^T$  and  $\boldsymbol{\omega}_B = \mathbf{R}_{A \rightarrow B} \boldsymbol{\omega}_A$ . The NRO conditions are the perpendicularity conditions  $\vec{\boldsymbol{\omega}}_{T \rightarrow IT} \perp (\vec{\mathbf{e}}_3^{IT} = \vec{\mathbf{p}})$  and  $\vec{\boldsymbol{\omega}}_{C \rightarrow IC} \perp (\vec{\mathbf{e}}_3^{IC} = \vec{\mathbf{p}})$ , or in terms of the components in the intermediate systems  $(\omega_{T \rightarrow IT})_{IT}^3 = 0$ ,  $(\omega_{C \rightarrow IC})_{IC}^3 = 0$ . They produce the differential equations

$$\dot{s} = \dot{E} (\cos d - 1), \quad \dot{s}' = \dot{F} (\cos g - 1). \quad (20)$$

Any orthogonal rotation matrix  $\mathbf{R}$  depends on only three parameters, so that the original seven parameters ( $E, d, s, \theta, F, g, s'$ ) or  $(X, Y, s, \theta, \xi \approx x_P, \eta \approx -y_P, s')$  must fulfill four conditions. Such conditions cannot be established for the CIP  $\vec{\mathbf{p}} = \vec{\mathbf{e}}_3^{IC} = \vec{\mathbf{e}}_3^{IT}$ , because the latter lacks a clear and rigorous definition either physical or mathematical. Roughly speaking it is a smoothed version of the direction  $\vec{\mathbf{n}} = \omega^{-1} \vec{\boldsymbol{\omega}}$  ( $\omega = |\vec{\boldsymbol{\omega}}|$ ) of the rotation vector of the Earth  $\vec{\boldsymbol{\omega}}$  ( $\vec{\boldsymbol{\omega}}_{C \rightarrow T}$ ), where “unobserved” high frequencies of  $\vec{\boldsymbol{\omega}}$  predicted by theory have been removed from precession–nutation and included in polar motion so that the same rotation matrix is maintained. The CIP is an evolution of the CEP which replaced the instantaneous rotation vector  $\vec{\boldsymbol{\omega}}$  after Atkinson (1973) and others remarked that higher than diurnal frequencies cannot be observed as a consequence of the related temporal resolution of the then available observations. Although the idea of replacing in a model a concept with its smoothed version due to observational resolution problems may seem strange to an outsider, astronomers developed the firm belief that the CEP is “observable” while the direction of instantaneous rotation vector is not. Nevertheless, today’s observations have higher than diurnal resolution and will certainly improve in the future (see e.g. Hefty et al. 2000; Artz et al. 2012) and even astronomers recognize the possibility of observing the instantaneous rotation vector (see e.g., Bolotin et al. 1997). We will therefore seek compatibility conditions for the case where the rotation matrix has a similar to the IERS representation, with the third axes are aligned to the rotation vector direction  $\vec{\mathbf{e}}_3^{IC} = \vec{\mathbf{e}}_3^{IT} = \vec{\mathbf{n}} = \omega^{-1} \vec{\boldsymbol{\omega}}$  instead of the CIP  $\vec{\mathbf{p}}$ . This allows the use of the rigorous definition of the rotation vector  $\vec{\boldsymbol{\omega}} = \vec{\mathbf{e}}^T \boldsymbol{\omega}_T \equiv \vec{\mathbf{e}}^T \boldsymbol{\omega}$  through the generalized kinematic Euler equations  $[\boldsymbol{\omega} \times] = \mathbf{R} \dot{\mathbf{R}}^T$ . Assuming that the diurnal rotation  $\mathbf{D} = \mathbf{R}_3(\theta)$  represents the rotation of the terrestrial system, both with

respect to the direction of its rotation  $\vec{e}_3^{IC} = \vec{e}_3^{IT} = \vec{n}$  and its rotation rate  $\omega = \dot{\theta}$  yields the three conditions

$$\vec{\omega} = \dot{\theta} \vec{e}_3^{IT} = \dot{\theta} \vec{e}_3^{IC} \quad (21)$$

or in terms of components  $\omega_{IT} = \dot{\theta} \mathbf{i}_3 = \omega_{IC}$ . The explicit computation of either of the last relations is rather complicated but simple in principle. With  $\mathbf{R} = \mathbf{W}\mathbf{D}\mathbf{Q}$  in  $[\omega \times] = \mathbf{R}\dot{\mathbf{R}}^T$  it follows that  $\omega = \omega_T = \mathbf{W}\mathbf{D}\omega_Q + \mathbf{W}\omega_D + \omega_W$  where  $[\omega_Q \times] = \mathbf{Q}\dot{\mathbf{Q}}^T$ ,  $[\omega_D \times] = \mathbf{D}\dot{\mathbf{D}}^T$ ,  $[\omega_W \times] = \mathbf{W}\dot{\mathbf{W}}^T$ . Setting  $\vec{e}^T = \vec{e}^{IT} \mathbf{W}^T$  in  $\vec{\omega} = \vec{e}^T \omega = \dot{\theta} \vec{e}_3^{IT} = \vec{e}^{IT} (\dot{\theta} \mathbf{i}_3)$  gives  $\vec{e}^{IT} \mathbf{W}^T \omega = \vec{e}^{IT} (\dot{\theta} \mathbf{i}_3)$ , or  $\omega = \dot{\theta} \mathbf{W} \mathbf{i}_3 = \mathbf{W}\mathbf{D}\omega_Q + \mathbf{W}\omega_D + \omega_W$ , or  $\mathbf{W}^T \omega = \mathbf{D}\omega_Q + \omega_D + \mathbf{W}^T \omega_W = \dot{\theta} \mathbf{i}_3$ . Finally we arrive at the three conditions

$$\mathbf{R}(E + s) \begin{bmatrix} \dot{E} \sin d \\ \dot{d} \end{bmatrix} = \mathbf{R}(\theta) \mathbf{R}(F + s') \begin{bmatrix} \dot{F} \sin g \\ \dot{g} \end{bmatrix}, \quad (22)$$

$$s' - \dot{F} \cos g + \dot{F} = \dot{s} - \dot{E} \cos d + \dot{E}, \quad (23)$$

where  $\mathbf{R}(\alpha)$  denotes rotation in the plane. Instead of the four conditions required to reduce the seven rotation parameters to the three, we have only three, which fix the orientation of  $\vec{n}$  and the rotation rate  $\omega = \dot{\theta}$ , but they leave undefined the positions of  $\vec{e}_1^{IT}$  (TIO) and  $\vec{e}_1^{IC}$  (CIO), regulated by the values of  $s$  and  $s'$ . If both sides of (23) are set equal to  $\dot{f}$ , where  $f$  is an arbitrary function, then  $s$  and  $s'$  are determined from  $s' = \dot{F}(\cos g - 1) + \dot{f}$ ,  $\dot{s} = \dot{E}(\cos d - 1) + \dot{f}$  and TIO and CIO are both displaced by the same amount  $f$ , leaving the angle  $\theta$  between  $\vec{e}_1^{IT}$  and  $\vec{e}_1^{IC}$  unaltered. To resolve this indeterminacy we must resort to the NRO conditions which correspond to the particular choice  $f = 0$ ! Thus (23) splits into the two NRO conditions

$$s' - \dot{F} \cos g + \dot{F} = 0, \quad \dot{s} - \dot{E} \cos d + \dot{E} = 0, \quad (24)$$

which together with the two conditions in (22) reduce the seven parameters to the required three independent ones. The essence of conditions (22) is that when precession–nutation ( $E, d$ ) and diurnal rotation  $\theta$  are known then polar motion ( $F, g$ ) is uniquely determined! And the other way around when polar motion and diurnal rotation are known precession–nutation is uniquely determined! The above conditions should all be satisfied to assure the alignment of  $\vec{e}_3^{IT} = \vec{e}_3^{IC}$  with  $\vec{\omega}$  and the proper rate  $\omega = \dot{\theta}$ . The condition (23) or the NRO conditions (24) alone do not guarantee that  $\omega = \dot{\theta}$ . For the current IERS representation with  $\vec{e}_3^{IT} = \vec{e}_3^{IC}$  aligned to the CIP  $\omega \neq \dot{\theta}$  and specific corrections must

be applied to  $\dot{\theta}$  in order to obtain the angular velocity of Earth rotation  $\omega$  and the related correct Universal Time (UT1).

No matter what the chosen direction of  $\vec{e}_3^{IT} = \vec{e}_3^{IC}$  (axis of diurnal rotation) the resulting rotation matrix  $\mathbf{R}$  implies a mathematically compatible instantaneous rotation axis  $\vec{\omega}$  which can be computed and compared to the CIP direction  $\vec{p}$ . The separation between  $\vec{n} = |\vec{\omega}|^{-1} \vec{\omega}$  and  $\vec{p}$  are the so called Oppolzer terms (or at least one the possible definitions, see Moritz and Mueller 1987). Dermanis and Tsoulis (2007) have computed these differences by two independent methods and found that although they have the same spectral characteristics as the Oppolzer terms, their amplitudes are too large, rising to the order of tens of meters on the Earth surface.

## 7 Open Issues for Further Research

There are of course many open problems within the analysis of data in the various space techniques deserving separate reviews, here however we will concentrate to theoretical issues relating to the exploitation of data coming from these techniques.

Current approaches to the formulation of the ITRF are based on the false assumption of zero mean random errors and known covariance matrices. An analysis and comparison is needed on the effect of both biases such as quasi-periodic terms and of incorrect covariances, which, for the GPS case at least, are known to be too optimistic.

Another issue of great practical importance is the optimal merging of local or regional networks to the ITRF. This is the old network densification problem, where the main rule is that the ITRF coordinates and velocities cannot be altered. Some researchers (see e.g. Altamimi 2003; Kotsakis 2013) suggest the use of minimal constraints based on ITRF parameters. The question is which minimal constraints to use in order to best reference local networks to the “official” ITRF reference system. An answer to this problem has been given by the generalized inner constraints of Kotsakis (2013) where he minimizes the trace of the covariance matrix of the local network when the effect of the uncertainty in the ITRF parameters used in the constraints is taken into account. This however is only optimality in appearance, while optimality of the connection to the ITRF reference system is rather desired. In any case merging through minimal constraints has an important disadvantage: high quality ITRF information on the shape of the common subnetwork and its temporal variation is completely ignored.

The examination of the residuals of ITRF coordinates after the fitting of the linear-in-time model demonstrates the existence of mostly annual and also semiannual signals,

especially in the height component. If an annual signal is fitted to an annual moving window, different amplitudes and periods will result. This suggests a quasi-periodic signal with a physical significance which is best revealed by modeling as an amplitude and phase modulated signals on an annual (semiannual) carrier of the form

$$\begin{aligned} x_i(t) &= A(t) \cos(2\pi t/T_0 - \phi(t)) \\ &= a(t) \cos(2\pi t/T_0) + b(t) \sin(2\pi t/T_0), \end{aligned} \quad (25)$$

where  $T_0$  is the annual (semiannual) period of the carrier frequency. The time varying amplitude  $A(t)$  corresponds to the varying severance of weather related phenomena, while the time varying phase  $\phi(t)$  corresponds to the time shift of maxima & minima of weather related phenomena. Of course the last term of the above model is more appropriate for data analysis because of its linearity with respect to  $a(t)$  and  $b(t)$ , which in any case can be directly converted to the more physically meaningful  $A(t)$  and  $\phi(t)$ . The question how to model these functions, e.g. by piecewise linear models, can only be answered by extensive analysis of the coordinate residuals where different models may be more appropriate for different stations and coordinates. For example, some coordinates demonstrate a “saw” effect where the annual wave is steeper when descending than when ascending, a behavior which suggests a periodic phase modulation. An interesting alternative to the periodic model (25) is to use a discontinuous piecewise linear model (epoch reference frames, see e.g. Blossfeld et al. 2014).

Two main questions are open with respect to the modeling of the quasi-periodic terms: Should the quasi-periodic part of the model be included in the data analysis step for the ITRF formulation along with the linear part, or it should it be fitted a posteriori to the residuals of the linear trend? Should the definition of the kinematically optimal reference system (minimal coordinate variation) also refer to the additional quasi-periodic part (“swinging” reference system) or should it stick to the linear part related to secular plate motion?

The final challenge in any case is the efficient geophysical interpretation of the quasi-periodic variations in relation to other geophysical data. From the geodetic point of view the correlation between temporal coordinate and gravity variation, which to a great part must be have common origins, has not drawn the attention it deserves.

## References

- Altamimi Z (2003) Discussion on how to express a regional GPS solution in the ITRF. In: Proceedings of the EUREF symposium, June 4–6, Toledo, Spain
- Altamimi Z, Dermanis A (2009) The choice of reference system in ITRF formulation. *IAG Symp* 137:329–334
- Altamimi Z, Sillard P, Boucher C (2002) ITRF2000: a new release of the international terrestrial reference frame for earth science applications. *J Geophys Res* 107(B10):2214
- Altamimi Z, Sillard P, Boucher C (2004) ITRF2000: from theory to implementation. *IAG Symp* 127:157–163
- Altamimi Z, Collilieux X, LeGrand J, Garayt B, Boucher C (2007) ITRF2005: a new release of the International Terrestrial Reference Frame based on time series of station positions and Earth orientation parameters. *J Geophys Res* 112, B09401
- Altamimi Z, Collilieux X, Metivier L (2011) ITRF2008: an improved solution of the international terrestrial reference frame. *J Geod* 85:457–473
- Angermann D, Drewes H, Krügel M et al (2004) ITRS Combination Centre at DGFI: a terrestrial reference frame realization 2003. *DGK B-313, München*
- Artz T, Bernhard L, Nothnagel A, Steigenberger P, Tesmer S (2012) Methodology for the combination of sub-daily Earth rotation from GPS and VLBI observations. *J Geod* 86:221–239
- Atkinson R (1973) On the “dynamical variations” of latitude and time. *Astron J* 78:147–151
- Bizouard C, Gambis D (2009) The combined solution C04 for Earth orientation parameters consistent with International Terrestrial Reference Frame 2005. *IAG Symp* 134:265–270
- Blossfeld M, Seitz M, Angermann D (2014) Non-linear station motions in epoch and multi-year reference frames. *J Geod* 88:45–63
- Bolotin S, Bizouard C, Loyer S, Capitaine N (1997) High frequency variations of the Earth’s instantaneous angular velocity vector. Determination by VLBI data analysis. *Astron Astrophys* 317:601–609
- Capitaine N, Guinod B, Souchay J (1986) A non-rotating origin of the instantaneous equator: definition, properties and use. *Celestial Mech* 39:283–307
- Dermanis A (2000) Establishing global reference frames. Nonlinear, temporal, geophysical and stochastic aspects. *IAG Symp* 123:35–42
- Dermanis A (2011) On the alternative approaches to ITRF formulation. A theoretical comparison. *IAG Symp* 139:107–113
- Dermanis A, Tsoulis D (2007) Numerical evidence for the inconsistent separation of the ITRF-ICRF transformation into precession-nutation, diurnal rotation and polar motion. In: IERS workshop on conventions, 20–21 September 2007, Paris. <http://der.topo.auth.gr>
- Hefty J, Rothacher M, Springer T, Weber R, Beutler G (2000) Analysis of the first year of Earth rotation parameters with a sub-daily resolution gained at the CODE processing center of the IGS. *J Geod* 74:479–487
- Kotsakis C (2013) Generalized inner constraints for geodetic network densification problems. *J Geod* 87:661–673
- Moritz H, Mueller I (1987) Earth rotation. Theory and observation. Ungar, New York
- Munk WH, MacDonald GJF (1960) The rotation of the Earth. Cambridge University Press, Cambridge
- Petit G, Luzum B (eds) (2010) IERS Conventions (2010). IERS Technical Note No. 36
- Rothacher M, Angermann D, Artz T et al (2011) GGOS-D: homogeneous reprocessing and rigorous combination of space geodetic observations. *J Geod* 85:679–705
- Seitz M, Angermann D, Blossfeld M, Drewes H, Gerstl M (2012) The 2008 DGFI realization of the ITRS: DTRF2008. *J Geod* 86:1097–1123
- Truedell C, Noll W (1965) The non-linear field theories of mechanics. Springer, Berlin

---

**Part II**

**Geodetic Data Analysis**

---

# Noise Analysis of Continuous GPS Time Series of Selected EPN Stations to Investigate Variations in Stability of Monument Types

Anna Klos, Janusz Bogusz, Mariusz Figurski, and Wieslaw Kosek

---

## Abstract

The type of monument that a GPS antenna is placed on plays a significant role in noise estimation for each permanent GPS station. In this research 18 Polish permanent GPS stations that belong to the EPN (EUREF Permanent Network) were analyzed using Maximum Likelihood Estimation (MLE). The antennae of Polish EPN stations are placed on roofs of buildings or on concrete pillars. The analyzed data covers a period of 5 years from 2008 to 2013. The analysis was made on the daily topocentric coordinate changes. Firstly, the existence of the combination of white noise, flicker noise and random-walk on each of the stations was set up before the analysis, secondly – a random-walk plus white noise model was assumed, because monument instability is thought to follow random-walk. The first combination of noises did not yield any conclusions about stability of monuments, probably because of the domination of flicker noise in the time series. The second one, even if not quite correct – noises in GPS time series do not strictly reflect random-walk only-showed that concrete pillars perform better than buildings for GPS antenna locations. Unfortunately, on the basis of this it cannot be clearly stated whether they are better as monuments or not. Moreover, the stacked Power Spectral Densities (PSDs) of topocentric coordinates were obtained with Fast Fourier Transform (FFT) for each monument type. Even though stacked spectra are quite similar and do not really show any differences, PSDs made for certain station are more varied.

---

## Keywords

EPN • GPS • Maximum Likelihood Estimation • Monument instability • Noises

---

A. Klos (✉) • J. Bogusz • M. Figurski  
Faculty of Civil Engineering and Geodesy, Military University  
of Technology, Kaliskiego St., 200-908 Warsaw, Poland  
e-mail: [aklos@wat.edu.pl](mailto:aklos@wat.edu.pl)

W. Kosek  
University of Agriculture, Environmental Engineering and Land  
Surveying, Mickiewicza Av. 24/28, 30-059 Kraków, Poland

Polish Academy of Sciences, Space Research Centre, Bartycka St.  
18A, 00-716 Warsaw, Poland

---

## 1 Introduction

Each topocentric component of station coordinates (by means of North, East and Up) is considered to follow the sum of:

$$x(t) = x_0 + v_x \cdot t + \sum_{i=1}^n [A_i \cdot \sin(\omega_i \cdot t + \varphi_i)] + O_x + \sum_{j=1}^m p_j \cdot x_j^{off} + \varepsilon_x(t) \quad (1)$$

where  $x_0$  is the initial value of the coordinate component,  $v$  is



the velocity,  $A$ ,  $B$ ,  $\omega$ ,  $\phi$  are the amplitudes, angular velocity and phase shift of the  $i$ -th periodic component of a time series,  $O_x$  stands for any known outliers,  $x_{off}$  for offsets,  $p$  is the Heaviside step function that is equal to 0 or 1 depending on the position of the offset,  $\varepsilon_x$  is the noise component. The noise component  $\varepsilon_x$  is, in most cases, a combination of white noise and coloured noise with amplitudes of  $a$  and  $b_\kappa$ , respectively (Zhang et al. 1997):

$$\varepsilon_x(t) = a \cdot \alpha(t) + b_\kappa \cdot \beta(t) \quad (2)$$

As stated previously, Agnew (1992), noises in a geophysical time series are correlated in time and are well described by power-law process with a power spectrum equal to:

$$P_x(f) = P_0 \left( \frac{f}{f_0} \right)^\kappa \quad (3)$$

where  $f$  is the spatial or temporal frequency,  $P_0$  and  $f_0$  are normalising constants and  $\kappa$  is the spectral index (Mandelbrot and Van Ness 1968). Mandelbrot (1983) and Feder (1988) discussed processes with different spectral indexes and attempted to attribute causes to. Agnew (1992) proved that the power spectra of most geophysical phenomena can be described by the power-law process with spectral indexes often falling in the range of  $-3$  up to  $-1$ . The integer values of indexes indicate special types of noises. Processes with “ $\kappa = 0$ ” correspond to white noise (WH) with a flat power spectrum, “ $\kappa = -1$ ” stands for flicker noise (FL) (Mandelbrot 1983) which is commonly recognized in most GPS coordinate time series and can be present in data time series due to GNSS signal propagation errors (Wielgosz et al. 2012; Hadas et al. 2013), finally, “ $\kappa = -2$ ” is described as random-walk (RW) noise and is considered to be related to instability of monuments that GPS antennae are attached to (Johnson and Agnew 1995; Williams et al. 2004; Beavan 2005; Hill et al. 2009). In order to improve the detection of a random-walk influence from a geodetic monument is to repeat the high-precision measurements in a tectonically stable region. In addition, its appearance in a time series can be reliably detected only in time series where the appropriate length of data, sampling frequency and favourable (low) amplitudes of other noises are present in the data. The issue of noise analysis is of great importance in the determination of the reliability of the velocity field on local and regional scales (Bogusz et al. 2012, 2013). In addition, most of the permanent stations belong to active geodetic networks supporting precise positioning, hence their stability influences the user position (Grejner-Brzezinska et al. 2009).

There are a few different techniques which can be used to easily detect noise in a time series. The first one, often considered to be the most effective, accurate and precise

(Beran 1994; Williams et al. 2004), is the technique of Maximum Likelihood Estimation (MLE) (Langbein and Johnson 1997). MLE has already been used in many papers that describe noise evaluation, e.g. Beavan (2005), Bergstrand et al. (2007), Teferle et al. (2008), Bos et al. (2008). It is calculated (e.g. Williams et al. 2004) using the following:

$$lik(\hat{v}, C) = \frac{1}{(2\pi)^{N/2} (\det C)^{1/2}} \times \exp\left(-0.5 \cdot \hat{v}^T \cdot C^{-1} \cdot \hat{v}\right) \quad (4)$$

where  $lik$  is the likelihood function,  $\hat{v}$  stands for postfit residuals from linear or nonlinear models applied to data,  $N$  is the number of epochs,  $C$  is the data covariance matrix. The second method of evaluating noise, spectral analysis, is based on the evaluation of the power spectrum of the data (as in Zhang et al. 1997; Mao et al. 1999). Both King and Watson (2010) and Bogusz and Kontny (2011) used it in their analysis but Langbein and Johnson (1997) and Pilgrim and Kaplan (1998) have both stated that it is less precise than MLE. The classical definition of the periodogram is:

$$P_x(\omega) = \frac{1}{N_0} |FT_x(\omega)|^2 \quad (5)$$

where  $FT_x(\omega) = \sum_{j=1}^{N_0} x(t_j) \cdot \exp(-i \cdot \omega \cdot t_j)$  for  $j = 1, 2, \dots, N_0$

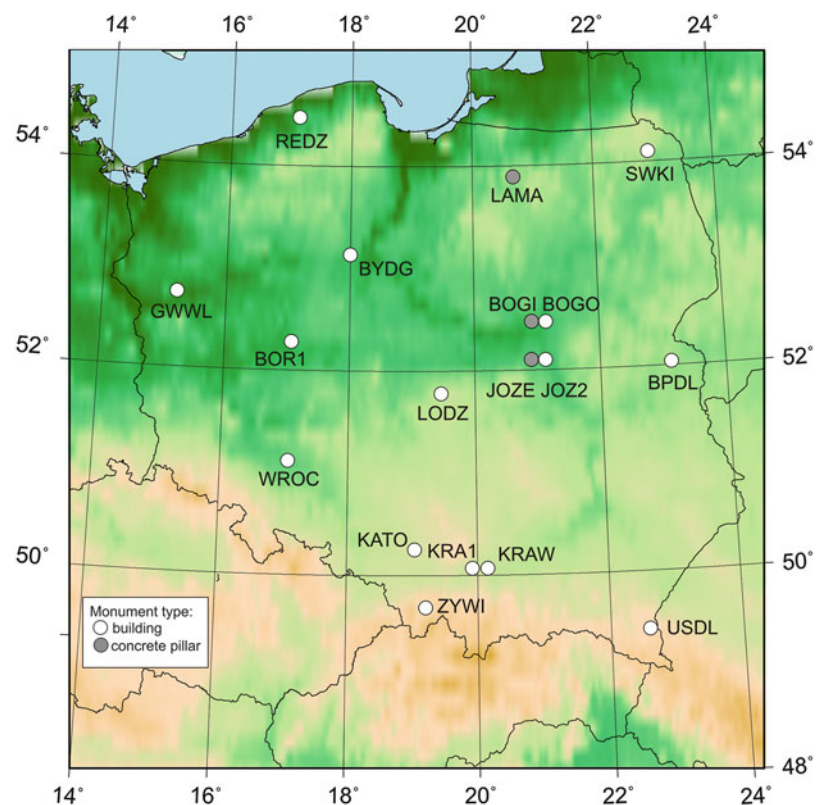
## 2 Data Analysis

Data used in this research was processed according to the EPN (Bruyninx et al. 2002) guidelines using Bernese 5.0 software (with absolute models consistent with IGS08 (Rebischung et al. 2012)) by the Centre of Applied Geomatics that cooperates at Military University of Technology as one of 18 EPN independent Local Analysis Centres. The processing strategy was performed in the Bernese 5.0 software and included parameters mentioned in Table 1. As the result, the coordinates in ITRF2008 Reference Frame were obtained (Altamimi et al. 2011). We have used 18 permanent Polish EPN stations with daily changes of topocentric coordinates (North, East, Up) (Fig. 1). One of the most important issues in studies of GPS noise is the monumentation used for instance whether the antenna is mounted on buildings or specialist concrete pillars. The main goal of this analysis is to investigate if the amplitudes of random-walk noise change with the different types of station monumentation. This investigation is a continuation of the author’s research concerning reliability in a GNSS time series (Bogusz et al. 2011).

**Table 1** Parameters used during the processing strategy in Bernese 5.0 software

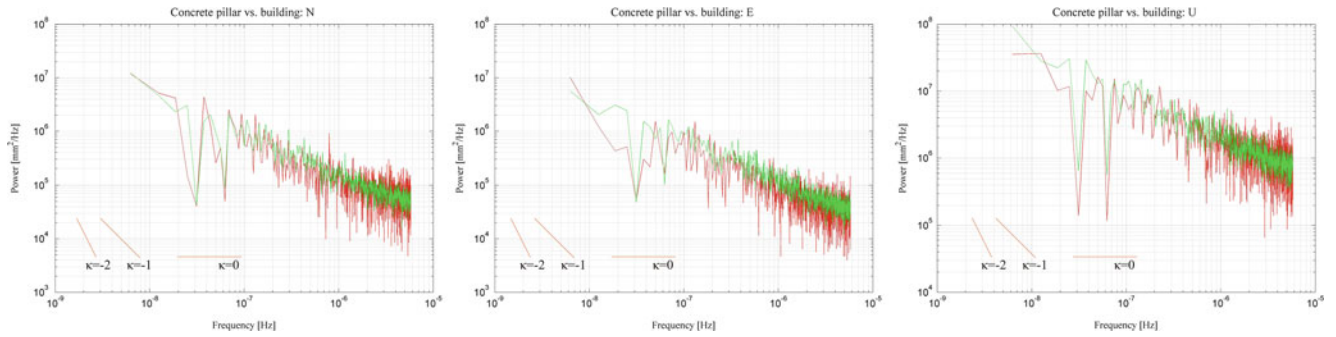
Processing strategy in Bernese 5.0 software	
Elevation angle cut-off	3 degrees, elevation dependent weighting using $\cos(z)$
Orbits and ERPs	IGS precise final orbits and ERPs
Troposphere	Saastamoinen – based dry component (Dry-Niell mapping function) as a priori model and the Wet-Niell mapping function
Ionosphere	CODE global iono models (help to increase the number of resolved ambiguities), finally cancelled out due to ionosphere-free linear combination
Ambiguity	QIF strategy – for baseline lengths shorter than 100 km – L5/L3 approach, for baselines shorter than 20 km – L1/L2 approach
Observations	Only GPS observations (RINEX format) were used with carrier phase as a basic observable (double-differences, ionosphere-free linear combination)
Planetary ephemeris model	DE405
Ocean tides model	OT_CSRC
Earth geopotential model	JGM-3
Nutation model	IERS2000
Tidal displacements	Solid Earth tides-according to IERS2003 standards
Ocean loading model	FES2004

**Fig. 1** Permanent Polish EPN stations used for the research. *Grey dots* stand for antennae placed on concrete pillars, *white ones* for antennae placed on buildings. The map was drawn in GMT software (Wessel and Smith 1998)



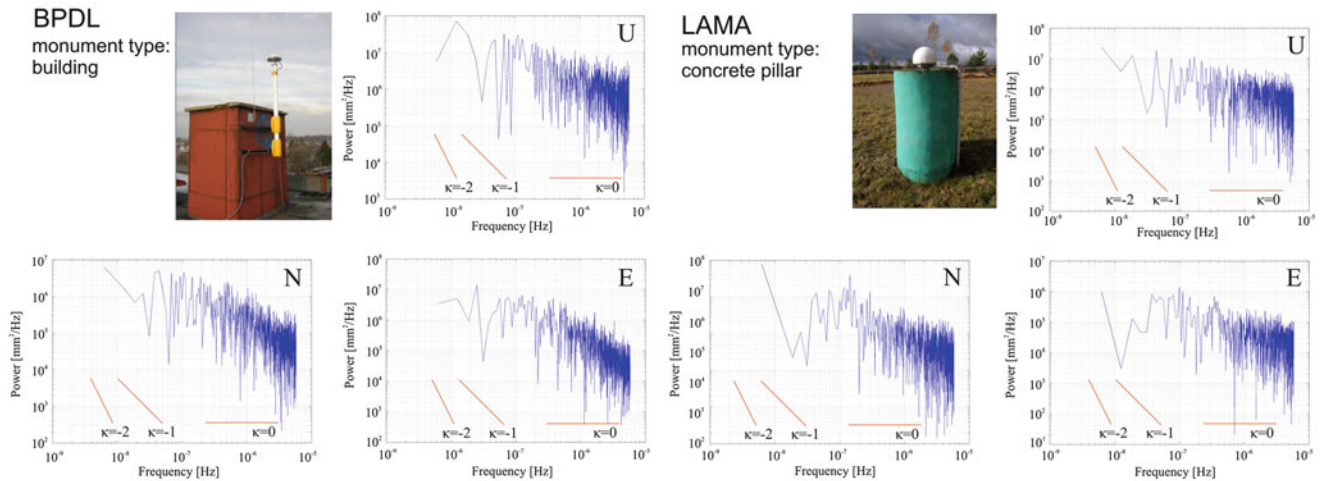
The preliminary analysis of the data involved the removal of outlying values that exceeded the criterion of three times the standard deviation ( $3\sigma$ ), the removal of seasonal components with annual and semi-annual periods as well as a linear trend with least squares. The outliers constituted about 1-3% of each time series. Any missing data was interpolated using linear interpolation before further analysis (to obtain regularly sampled data). For each topocentric component, the FFT was performed and stacked power spectral densities

were made. The scale for them was changed from linear into log-log one. The main advantage of using PSD is that the slope of the graph in log-log space corresponds to the spectral index of the dominant noise existing in the time series. Long-period components are thought to follow flicker or random-walk noise, whereas high-frequency ones follow a white noise model. For each graph, the theoretical values of “ $\kappa = -2$ ”, “ $\kappa = -1$ ” and “ $\kappa = 0$ ” were added so as to show how noise type influence a time series at low and high



**Fig. 2** Stacked Power Spectral Densities for North (*left*), East (*middle*) and Up (*right*) coordinate components. Stacked PSDs for concrete pillars (*red*) and buildings (*green*) are placed on one plot to compare

both monument types. Theoretical values of spectral indexes for white noise ( $\kappa = 0$ ), flicker noise ( $\kappa = -1$ ) and random-walk ( $\kappa = -2$ ) were also added



**Fig. 3** Exemplary PSDs for two types of antennae monumentation. *Left* – station BPD L with antennae placed on buildings. *Right* – station LAMA with antennae situated on concrete pillar. PSDs were made for

each station for N, E and U coordinate components. Theoretical values of integer spectral indexes were added to the plots

frequencies (Fig. 2). In spite of the fact that the stacked PSDs are almost the same, PSDs made for individual permanent stations show a few differences (Fig. 3). Slopes of PSDs for sites where the antennae is placed on buildings are in few cases higher than for ones mounted on concrete pillars. It indicates that monument instability is probably higher for buildings and may be caused by the settling of the building or perhaps thermal changes.

For the MLE analysis we assume that the noise present in the coordinate time series follow a combination of white noise, flicker noise and random-walk. However due to the fact that monument instability is thought to follow a random-walk noise model and it's amplitude is such (King and Williams 2009) that it may be masked by the FL and WH noise components we have analysed our time series with MLE using two different choices of noise combination: the first – WH, FL and RW and the second one – only WH and RW. All analyses were performed using the CATS software (Williams 2008).

The amplitudes of the noises in the first test show that FL dominates over WH and RW (Fig. 4). Although the typical range is between 2 and 4 mm year<sup>-0.25</sup> for the horizontal components, it is much higher for the vertical (6–12 mm year<sup>-0.25</sup>). The similarity in flicker noise amplitude over the whole region are most likely explained by the fact that FL is thought to be regionally coherent and for small areas (such as in Poland) it can reach similar values. The amplitudes of the WH component are also smaller for North and East than for Up. Unfortunately, for the first combination of noises, random-walk is close to zero for the majority of stations. Some of estimates allow for RW but its maximum amplitude reach the 1 mm year<sup>-0.5</sup> in the most extreme case (Up component, KRAW station). Such small values of RW indicate the relative stability of Polish EPN stations monuments or prove that the data is clearly not enough to detect it. The small RW noise amplitudes are consistent with the results presented in King and Williams (2009) which showed that if random-walk is considered as monument



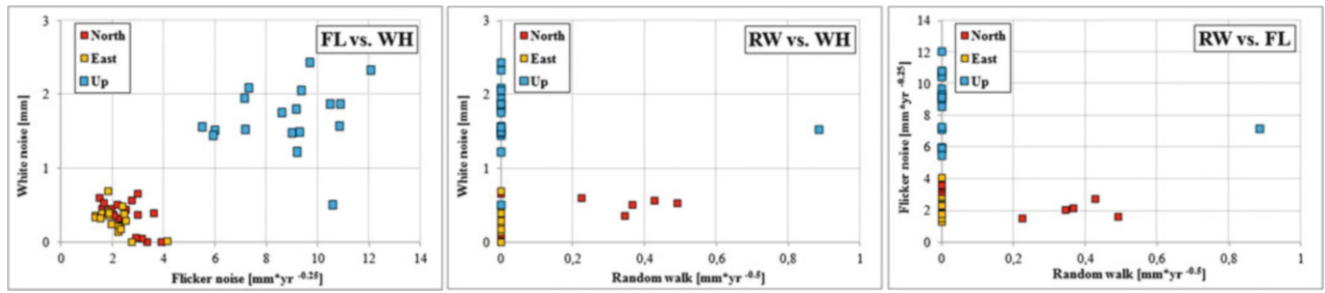
**Fig. 4** Amplitudes of white noise (grey), flicker noise (pink) and random walk (green) for North (top), East (middle) and the Up (bottom) components obtained with the MLE method. For each of these amplitudes a 1-sigma error bar was added. Different types of

monumentation (building or concrete pillar) were marked with different colours representing station abbreviations (black for buildings, grey for concrete pillars)

noise it is smaller than previously thought. They have analysed 10 short-baselines (two of them were created between stations analysed in the following research – BOGI-BOGO and JOZE-JOZ2) with the assumption of power-law or first-order Gauss Markov noise model. It was found that the amplitudes of noises for baselines are in general an order of magnitude smaller than in case of single station. They stated that random-walk is probably no higher than  $0.5 \text{ mm year}^{-0.5}$  for well monumented stations. On the other hand, we should be aware of the fact that RW amplitudes may be quite small (or estimated to be zero) because of the present domination of FL and the limited length of the data (only 5 years). On the basis of Williams et al. (2004), to detect RW with an amplitude of  $0.4 \text{ mm year}^{-0.5}$  a period of at least 30 years worth of data is needed to detect it easily. But of course, the longer the data period, the more reliable the estimations of

noises which can be obtained. Scatter plots made for the first noise combination (Fig. 5) show some dependencies between the amplitudes of WH and FL. For the Up component the amplitudes are higher and much more spread out than for the horizontal components. Unfortunately, due to the small values of RW, no dependencies between RW and WH or FL were noticed.

The noise amplitudes obtained for the combinations of WH and RW with the MLE method for horizontal components (North, East) are quite varied for random-walk noise while they remain similar for white noise at around 1 mm (Fig. 6). The highest amplitude for both North and East components were found at the station, BPD L. White noise amplitudes for the Up component are greater than for horizontal ones and are at the level of 3 mm. All of the random-walk amplitudes are higher than  $6 \text{ mm year}^{-0.5}$  while some



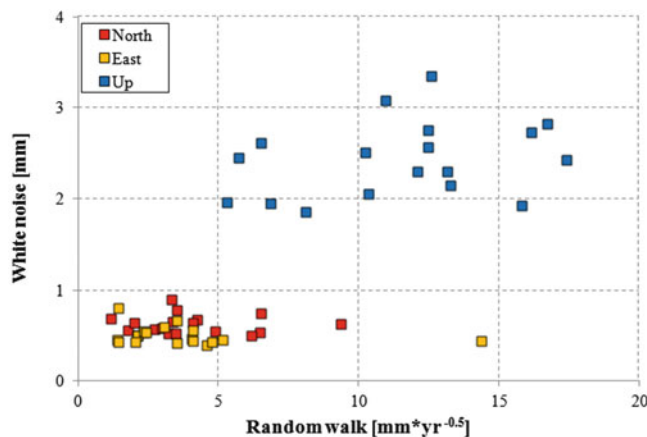
**Fig. 5** Scatter plots for different noise types: flicker noise vs. white noise (*left*), random-walk vs. white noise (*middle*), random-walk vs. flicker noise (*right*). Amplitudes for North, East and Up components are represented with different colours – *red*, *yellow* and *blue*, respectively



**Fig. 6** Amplitudes of white (*grey*) and random-walk (*green*) noise for Polish EPN stations obtained with the MLE method with the assumption of WH and RW noise only. The plots are presented for North (*top*), East (*middle*) and Up (*bottom*) components. The 1-sigma error bars were added to each of amplitude

of them even exceed the value of  $20 \text{ mm year}^{-0.5}$  (20 times greater than for the first combination). They also have greater error bars which can be the result of incorrect fit and the inappropriate assumption of WH plus RW only. This rather

unrealistic result is likely due to the domination of flicker noise in the series which, in the absence of a flicker noise component in the MLE, is misinterpreted as random-walk. Furthermore the length of data means that the MLE could



**Fig. 7** Scatter plot for amplitudes of white and random-walk noise for Polish EPN stations obtained with MLE for the assumption of WH and RW noise only. The coordinate components are represented with *red* (North), *yellow* (East) and *blue* (Up) colours

not give us significant results. The scatter plot for the WH and RW combination (Fig. 7) clearly presents dependencies between noises. It can be noted that noise amplitudes for horizontal components are placed adjacent to another while they are more scattered and not centred around one specific value for the vertical. Although the second assumption of noise combination has some drawbacks as described above, amplitudes of RW obtained with this combination appear to show that concrete pillars are better as monuments for antennae than buildings. Unfortunately, due to the fact that this assumption is not correct, it cannot be stated for sure.

### 3 Discussion

Application of the MLE algorithm to the determination of the amplitudes of white, flicker and random-walk noise showed how they influence the accuracy and reliability of the parameters that are estimated from GPS time series (e.g. velocities of permanent stations). We compared two stochastic models that were combinations of white noise, flicker noise and random-walk and alternatively white noise plus random-walk noise. We showed that the amplitudes of the assumed models influence the coordinate time series measured by Polish EPN stations. The second combination of noise unexpectedly gave values of RW up to 20 times larger than in the first combination. However if flicker noise is present in the time series and its existence is ignored, it will influence the estimated amplitudes of the random-walk component. The flicker noise was found to be at a similar amplitude over small areas, such as the size of Poland, and therefore variations in the estimated random-walk amplitudes may still be considered to be due to the influence of monument instability. The questionable point

is whether the current time series (5 years) are really long enough to detect changes related to random-walk. As showed by Williams et al. (2004) to detect easily random-walk with a magnitude of  $0.4 \text{ mm year}^{-0.5}$  the 30 years data will be needed. Of course, the longer the time series, the more reliable the estimation of random-walk.

We also applied FFT to the topocentric components to create their Power Spectral Density estimates. Their slopes in log-log space also indicate the characteristics of the noise that appears in GPS time series. We averaged or stacked the power spectra as a function of monument type. However, when plotted together we could not visibly discern any obvious differences. Individually, the PSDs made for each station did show small variations in their slopes. Taking into consideration the power spectral densities (some of the estimated spectral slopes, or indices, for buildings are closer to  $-2$  than for concrete pillars) and the second noise combination (WH plus RW) in the MLE results although somewhat lacking (time series do not simply reflect such a characteristic so it cannot be stated for sure) the results hint that antennae placed on concrete pillars are apparently more stable than those mounted on buildings.

**Acknowledgments** This research is financed by the Ministry of Science and Higher Education, grant No. 2011/01/B/ST10/05384. The authors would like to thank the reviewers for their remarks which contributed to expanding the discussion in the following research.

### References

- Agnew DC (1992) The time-domain behaviour of power-law noises. *Geophys Res Lett* 19(4):333–336
- Altamimi Z, Collilieux X, Metivier L (2011) ITRF2008: an improved solution of the International Terrestrial Reference Frame. *J Geod* 85(8):457–473. doi:10.1007/s00190-011-0444-4
- Beavan J (2005) Noise properties of continuous GPS data from concrete pillar geodetic monuments in New Zealand and comparison with data from U.S. deep drilled braced monuments. *J Geophys Res* 110:B08410. doi:10.1029/2005JB003642
- Beran J (1994) Statistics for long-memory processes. *Monogr Stat Appl Probab* 61:315
- Bergstrand S, Schnereck H-G, Lidberg M, Johansson JM (2007) BIFROST: Noise properties of GPS time series. *Dynamic Planet. Int Assoc Geodesy Symposia* 130:123–130. doi:10.1007/978-3-540-49350-1\_20
- Bogusz J, Kontny B (2011) Estimation of sub-diurnal noise level in GPS time series. *Acta Geodynamica et Geomaterialia* 8(3):273–281
- Bogusz J, Figurski M, Kroszczyński K, Szafranek K (2011) Investigation of environmental influences to the precise GNSS solutions. *Acta Geodynamica et Geomaterialia* 8(1):5–15
- Bogusz J, Figurski M, Kontny B, Grzempowski P (2012) Horizontal velocity field derived from EPN and ASG-EUPOS satellite data on the example of south-western part of Poland. *Acta Geodynamica et Geomaterialia* 9(3):349–357
- Bogusz J, Klos A, Grzempowski P, Kontny B (2013) Modelling velocity field in regular grid on the area of Poland on the basis of the velocities of European permanent stations. *Pure Appl Geophys*. doi:10.1007/s00024-013-0645-2

- Bos MS, Fernandes RMS, Williams SDP, Bastos L (2008) Fast error analysis of continuous GPS observations. *J Geod* 82:157–166. doi:[10.1007/s00190-007-0165-x](https://doi.org/10.1007/s00190-007-0165-x)
- Bruyninx C, Kenyeres A, Takacs B (2002) EPN data and product analysis for improved velocity estimation: First results. Scientific Assembly of the International-Association-of-Geodesy. *Int Assoc Geodesy Symposia* 125:47–52
- Feder JW (1988) *Fractals*. Plenum, New York
- Grejner-Brzezinska DA, Arlsan N, Wielgosz P, Hong C-K (2009) Network calibration for unfavorable reference-rover geometry in network-based RTK: Ohio CORS case study. *J Surv Eng* 135(3): 90–100
- Hadas T, Kaplon J, Bosy J, Sierny J, Wilgan K (2013) Near-real-time regional troposphere models for the GNSS precise point positioning technique. *Meas Sci Technol* 24(5), doi:[10.1088/0957-0233/24/5/055003](https://doi.org/10.1088/0957-0233/24/5/055003)
- Hill EM, Davis JL, Elosegui P, Wernicke BP, Malinkowski E, Niemi NA (2009) Characterization of site-specific GPS errors using a short-baseline network of braced monuments at Yucca Mountain, southern Nevada. *J Geophys Res* 114, B11402. doi:[10.1029/2008JB006027](https://doi.org/10.1029/2008JB006027)
- Johnson HO, Agnew DC (1995) Monument motion and measurements of crustal velocities. *Geophys Res Lett* 22(21):2905–2908. doi:[10.1029/95GL02661](https://doi.org/10.1029/95GL02661)
- King MA, Watson CS (2010) Long GPS coordinate time series: multipath and geometry effects. *J Geophys Res* 115, B04403. doi:[10.1029/2009JB006543](https://doi.org/10.1029/2009JB006543)
- King MA, Williams SDP (2009) Apparent stability of GPS monumentation from short-baseline time series. *J Geophys Res* 114, doi:[10.1029/2009JB006319](https://doi.org/10.1029/2009JB006319)
- Langbein J, Johnson H (1997) Correlated errors in geodetic time series: implications for time-dependent deformation. *J Geophys Res* 102(B1):591–603
- Mandelbrot B (1983) *The fractal geometry of nature*. W.H. Freeman, San Francisco, 466 pp
- Mandelbrot B, Van Ness J (1968) Fractional Brownian motions, fractional noises, and applications. *SIAM Rev* 10:422–439
- Mao A, Harrison CGA, Dixon TH (1999) Noise in GPS coordinate time series. *J Geophys Res* 104(B2):2797–2816
- Pilgrim B, Kaplan DT (1998) A comparison of estimators for 1/f noise. *Phys D* 114:108–122
- Reischung P, Griffiths J, Ray J et al (2012) IGS08: the IGS realization of ITRF2008. *GPS Solutions* 16(4):483–494. doi:[10.1007/s10291-011-0248-2](https://doi.org/10.1007/s10291-011-0248-2)
- Teferle FN, Williams SDP, Kierulf KP, Bingley RM, Plag HP (2008) A continuous GPS coordinate time series analysis strategy for high-accuracy vertical land movements. *Phys Chem Earth* 33:205–216. doi:[10.1016/j.pce.2006.11.002](https://doi.org/10.1016/j.pce.2006.11.002)
- Wessel P, Smith WHF (1998) New, improved version of the Generic Mapping Tools. *Released EOS Trans AGU* 79:579
- Wielgosz P, Paziewski J, Krankowski A, Kroszczyński K, Figurski M (2012) Results of the application of tropospheric corrections from different troposphere models for precise GPS rapid static positioning. *Acta Geophys* 60(4):1236–1257. doi:[10.2478/s11600-011-0078-1](https://doi.org/10.2478/s11600-011-0078-1)
- Williams SDP (2008) CATS: GPS coordinate time series analysis software. *GPS Solutions* 12:147–153. doi:[10.1007/s10291-007-0086-4](https://doi.org/10.1007/s10291-007-0086-4)
- Williams SDP, Bock Y, Fang P, Jamason P, Nikolaidis RM, Prawirodirdjo L, Miller M, Johnson D (2004) Error analysis of continuous GPS position time series. *J Geophys Res* 109, B03412. doi:[10.1029/2003JB002741](https://doi.org/10.1029/2003JB002741)
- Zhang J, Bock Y, Johnson H, Fang P, Williams S, Genrich J, Wdowinski S, Behr J (1997) Southern California permanent GPS geodetic array: error analysis of daily position estimates and site velocities. *J Geophys Res* 102(B8):18035–18055

---

# Improvement of Least-Squares Collocation Error Estimates Using Local GOCE $T_{zz}$ Signal Standard Deviations

C.C. Tscherning<sup>†</sup>

---

## Abstract

The method of Least-Squares Collocation (LSC) may be used for the modeling of the anomalous gravity potential ( $T$ ) and for the computation (prediction) of quantities related to  $T$  by a linear functional. Errors may also be estimated. However, when using an isotropic covariance function or equivalent reproducing kernel, the error estimates will be nearly constant if the used data have a good (regular) distribution. In this case the error estimate will vary only if the data distribution changes (e.g. for satellite data as a function of latitude), if data are missing in an area, or if predictions are made outside the data area.

On the other hand, a comparison of predicted quantities with observed values show that the error also varies depending on the local data standard deviation. This quantity may be (and has been) estimated using the GOCE second order vertical derivative,  $T_{zz}$ , in the area covered by the satellite.

The ratio between the nearly constant standard deviations of a predicted quantity (e.g. in a  $25^\circ \times 25^\circ$  area) and the standard deviations of  $T_{zz}$  in smaller cells (e.g.,  $1^\circ \times 1^\circ$ ) have been used as a scale factor in order to obtain more realistic error estimates. This procedure has been applied on gravity anomalies (at 10 km altitude) predicted from GOCE  $T_{zz}$ . This has given an improved agreement between errors based on the differences between values derived from EGM2008 (to degree 512) and predicted gravity anomalies.

---

## Keywords

Collocation • Error estimates • Gravity anomalies • Gravity gradients

---

## 1 Introduction

Error estimates (and error correlations) are needed for several purposes, such as: (1) an indicator of the quality of an observed or estimated quantity (e.g. Andersen and

Remmer, 1982; Balmino, 2009); (2) for the use of data in a data assimilation procedure such as estimating ocean current velocities (e.g., Bingham et al. 2011); (3) in simulation studies (Arabelos et al. 2007); (Arabelos and Tscherning, 1999, 2008); and (4) for gross-error detection (Tscherning 1991). The geodetic literature on errors, both random and systematic, is vast. However, in several cases the estimated errors are of little use due for example to missing information of the physical variation of a signal. The cause of this has often been missing knowledge of the statistical characteristics of a phenomenon. A characteristic example is the error of prediction of gravity anomalies in an inaccessible area like high mountains or large lakes. Here the situation has changed. The data collected by ESA's Gravity and Ocean Circulation Explorer satellite (GOCE) has now given us

---

<sup>†</sup>Prof. Christian Tscherning passed away October 24, 2014. Throughout his scientific career Christian has passionately and constructively contributed to discussions on theoretical geodesy. He was a regular and vociferous participant in the long series of Hotine-Marussi Symposia. The editors decided to publish his contribution to the current proceedings of the 2013 Hotine-Marussi Symposium.

C.C. Tscherning  
Niels Bohr Institute, University of Copenhagen, 2100 Copenhagen  
Oe., Denmark



an excellent picture of the gravity field variation between 83° latitude south and north. In the following it will be shown how the information can be used to improve the error estimates of quantities predicted using Least-Squares Collocation (LSC), (Moritz, 1980; Sansò and Sideris, 2013). Similar possibilities exist for the results of other gravity field modeling methods.

## 2 Error Estimates in LSC

LSC may be used for the modeling of the anomalous gravity potential ( $T$ ) and for the computation (prediction) of quantities related to  $T$  by a linear functional,  $L_P$ . The subscript  $P$  indicates a contingent point of evaluation. Error variances  $\sigma^2(L_P)$  may also be estimated. When using a covariance function  $\text{cov}(L_P, L_Q)$  or equivalent reproducing kernel, the error variance becomes,

$$\sigma^2(L_P) = \text{cov}(L_P, L_P) - \{\text{cov}(L_P, L_i)\}^T \{\text{cov}(L_i, L_j) + e_{ij}\}^{-1} \{\text{cov}(L_P, L_j)\} \quad (1)$$

where  $L_i$  and  $L_j$  are linear functional associated with the observations and  $e_{ij}$  are the variance-covariances of the noise associated with the observations. If an isotropic (rotational invariant) basic covariance function is used this equation may show where data are missing or not available or where the noise is large, see the figure in Sansò and Sideris (2013), p. 329, which is a typical example of error estimates of predicted height anomalies from gravity anomalies in the New Mexico test area.

The error, however, is strongly related to the local data variance,  $\text{cov}(L_P, L_P)$  in Eq. (1), which is constant for a specific quantity at a specific altitude due to the isotropy but which in reality will vary depending on the position.

We will show how this can be used to improve the error-estimate in situations where the observations are of the same kind, e.g. second order radial derivatives of  $T$ ,  $T_{zz}$ , as used in Arabelos et al. (2013) for the prediction of global grids of gravity anomalies at 10 km altitude. In the “ideal” situation, the data have a normal distribution. In this case the absolute value of the ratio between the error and the error estimate should follow a t-distribution, see Tscherning (1991). A possible improvement should then bring us closer to this distribution, cf. Tables 1 and 2 below.

## 3 The Trench Example

The computations were done in  $25^\circ \times 25^\circ$  blocks., and a covariance function was estimated using gravity anomalies computed from EGM2008 to maximal degree 512 (Pavlis

**Table 1** Distribution of the ratio  $\text{abs}(\text{predicted}-\text{“observed”})/\text{error estimate}$  in 1.0 intervals. 625 values used

Interval	0.0–1.0	1.0–2.0	2.0–3.0	>3.0
Number	544	66	12	3
%	87	10	2	0.48
t-distribution	68	26	6	0

**Table 2** Ratio between  $\text{abs}(\text{predicted}-\text{“observed”})$  and new error estimate grouped in 1.0 bins

Interval	0.0–1.0	1.0–2.0	2.0–3.0	>3.0
Number	453	114	35	14
%	74	19	6	1

et al. 2012) but with the contribution from ITG-Grace2010s (Mayer-Guerr et al. 2010) to degree 36 subtracted. (This has been done in order to assure a zero mean value in this block and the global set of blocks described in the next section).

The block with the “best” signal to noise ratio is found in the area bounded by  $27.5^\circ$ ,  $52.5^\circ$  in latitude and  $137.5^\circ$ ,  $162.5^\circ$  in longitude. Here the gravity field is very inhomogeneous due to a deep trench, see Figs. 1 and 2 which shows the gravity anomalies at 10 km.

The gravity anomalies were predicted using LSC from  $T_{zz}$  (minus the ITG-Grace contribution) with data spaced as close as possible to the mid-points of a  $0.1666^\circ$  grid; see Fig. 2.

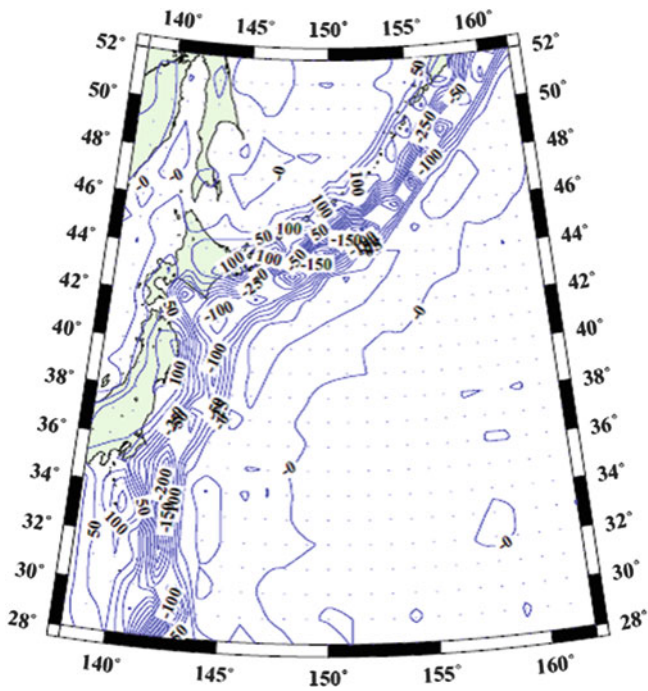
The differences between the predicted values (in a  $1^\circ$  grid) and the computed values are shown in Fig. 3, and the error estimates computed using Eq. (1) are shown in Fig. 4. The distribution of the ratio between the absolute value of the differences and the error estimate is shown in Table 1 grouped in 0.5 bins.

We can see that a very large number of values have a very small error estimate associated. This is caused by the very smooth field in the South-Eastern part of the block; see Fig. 3. We have however, from the observed gradients the actual value of the field in the area; see Fig. 5.

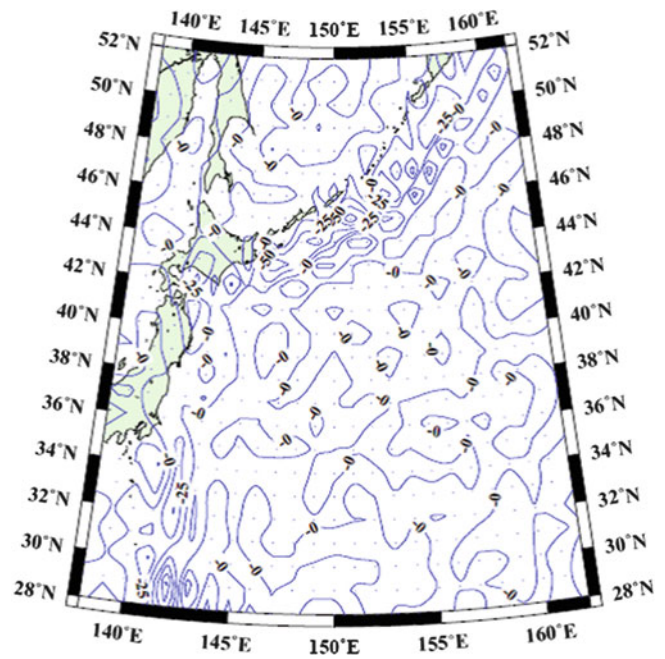
If we scale the error estimates using these local standard deviations, (see Fig. 6), we will obtain much more realistic error estimates. The mean value in the  $1^\circ$  blocks is different from zero, so alternatively we could have used the root-mean-square value.

As a scale factor we use the ratio between the error standard deviations computed from the covariance function (Fig. 8) and the local standard deviations in the  $1^\circ$  blocks (Fig. 7), which then result in a modified error-map, see Fig. 6 and a new histogram, Table 2.

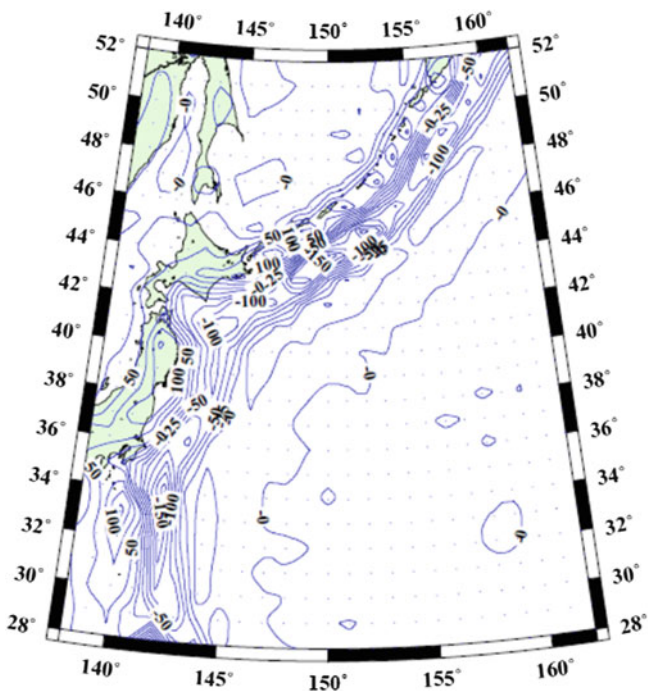
The distribution of the ratios have improved the error estimates in the sense that their distribution have become closer to the t-distribution than before see Table 1. However,



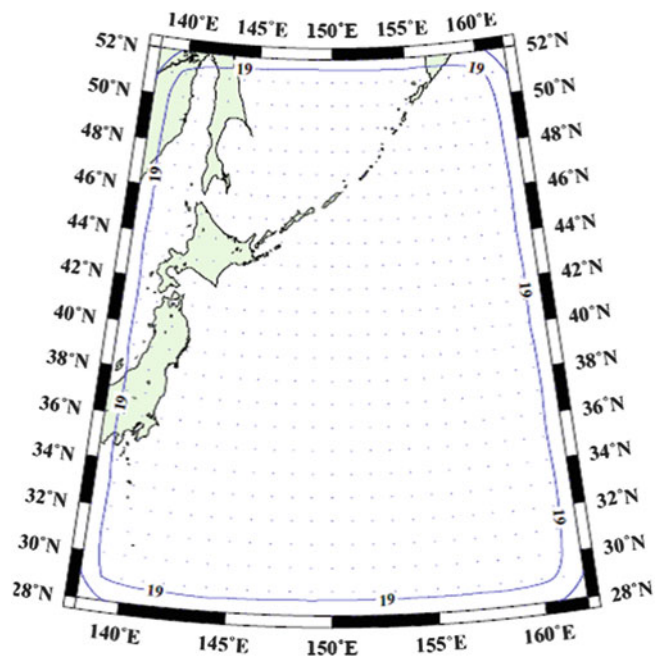
**Fig. 1** Gravity anomalies (mGal) at 10 km computed from EGM2008 minus the contribution from ITG-Grace2010s, used as equivalent to observed data



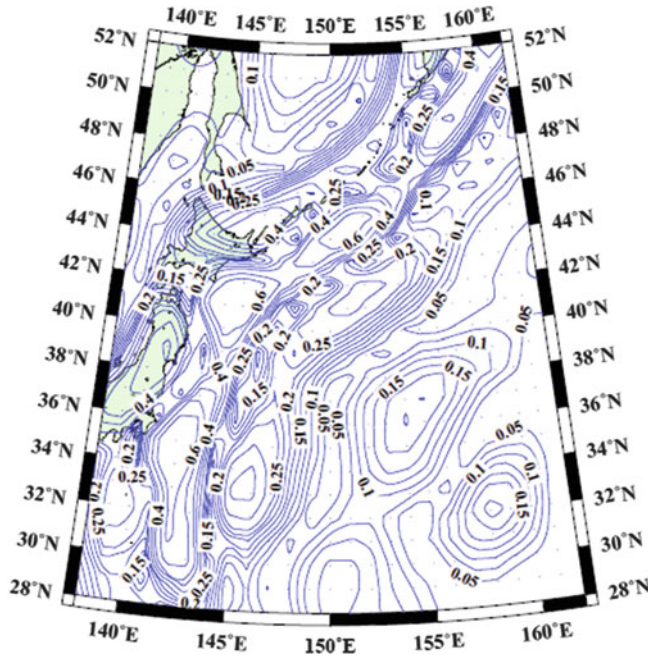
**Fig. 3** Differences between “observed” and predicted gravity anomalies at 10 km altitude, (mGal)



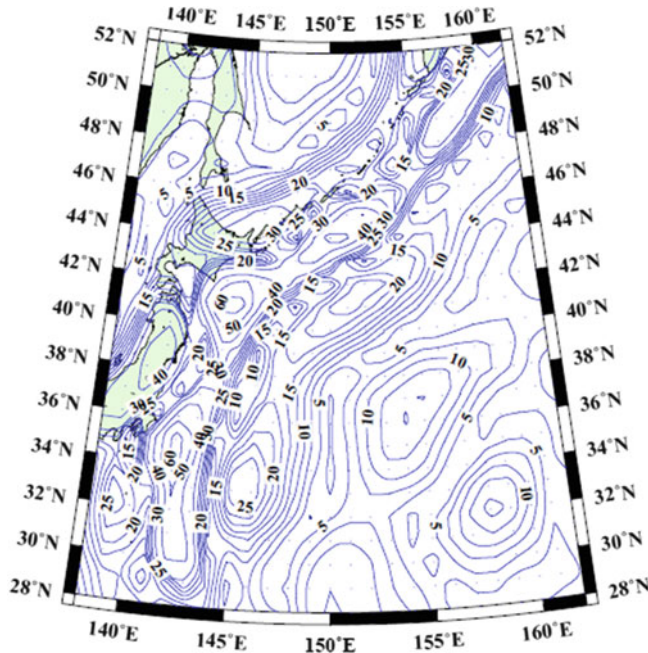
**Fig. 2** Gravity anomalies (mGal) at 10 km (minus ITG-Grace2010s) predicted from GOCE  $T_{zz}$



**Fig. 4** Error estimate of predicted gravity anomalies at 10 km using Eq. (1), (mGal)



**Fig. 5** Contour map of standard deviations of  $T_{zz}$  (minus ITG-Grace2010s to deg. 36) at satellite altitude computed for  $1^\circ$  grid blocks. Units E ( $1E = 1\text{e}6\text{tv}6s = 10^{-9}\text{s}^{-2}$ )



**Fig. 6** Scaled error estimates. Units mGal

the use of  $T_{zz}$  root-mean square values instead of standard-deviations would have given a slightly different (worse) result.

## 4 Global Scaling

The same procedure may be applied globally. The standard deviations of  $T_{zz}$  for  $1^\circ$  blocks are shown in Fig. 7. The standard deviations for the prediction of gravity anomalies in  $25^\circ \times 25^\circ$  blocks are nearly constant within a block, and are shown in Fig. 8. When these standard deviations are scaled we have a much more realistic and varying picture, see Fig. 9. The scale factor is determined based on the ratio between the local standard deviations and the standard deviation of the data in  $20^\circ \times 20^\circ$  blocks.

If  $e_c$  is the error estimate computed using collocation,  $e_{20}$ ,  $e_1$  the standard deviation of the data in the  $20^\circ$  and a  $1^\circ$  block, respectively then scaled error-estimate  $e_s$  of the  $1^\circ$  block becomes

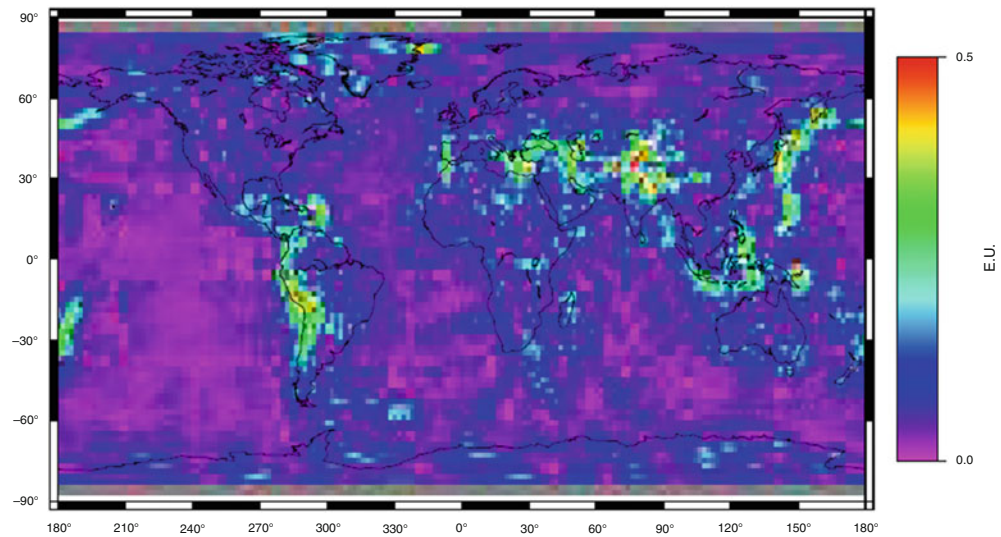
$$e_e = e_c \cdot e_1 / e_{20} \quad (2)$$

We now clearly see how the magnitude of the error is large in mountainous or trench areas and small in “smooth” areas.

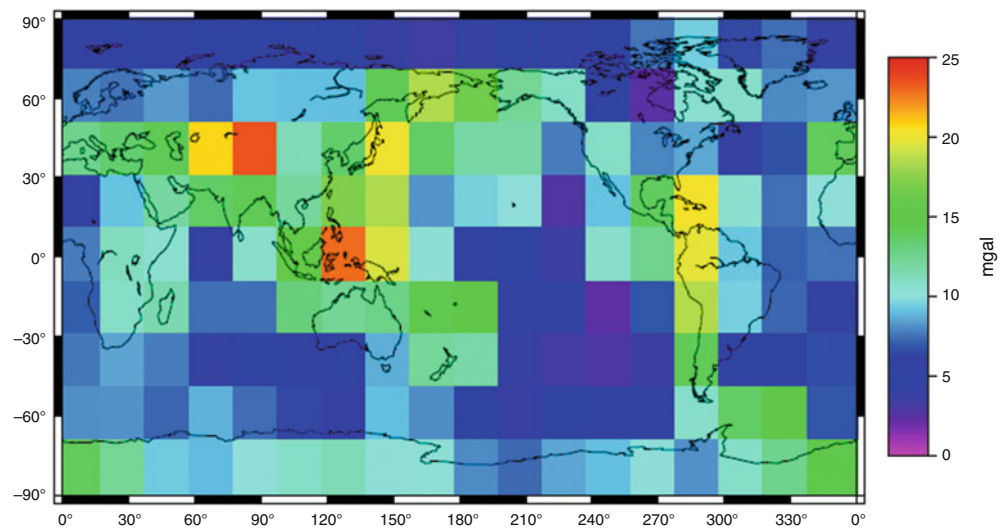
## 5 Conclusion

The scaling of LSC derived error estimates may improve the estimates, so that the variation of the error due to changing local signal standard deviation is seen. It is planned to use the procedure in order to provide error estimates of the gravity anomaly grids described in (Arabelos et al. 2013). There are however some problems remaining. How should the local signal standard deviations be computed in case of a local bias? Should the root-mean-square variation be used? Should the scaling be done regionally or by using the global root-mean-square variation? In both cases the bias is small due to the subtraction of for example the ITG-Grace 2010s field or other Earth Gravity Models to the same maximal degree.

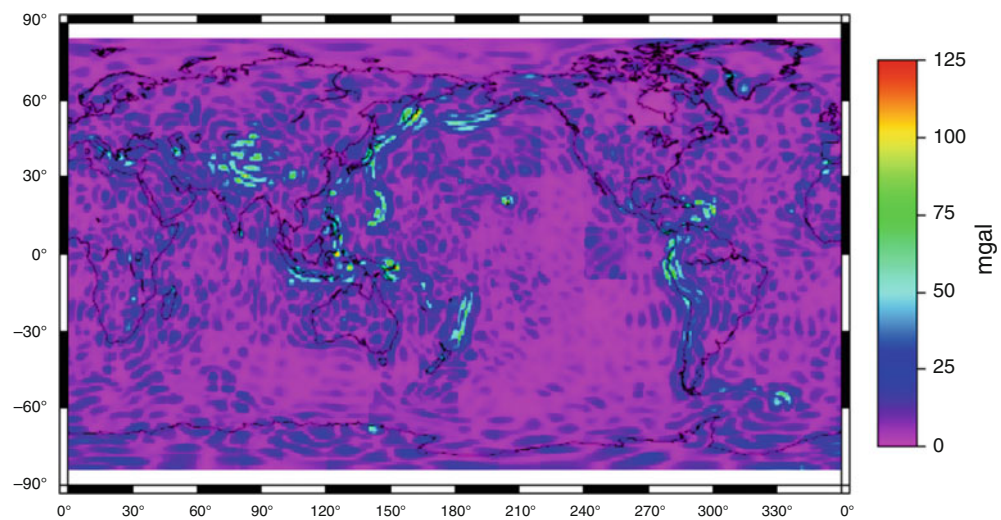
**Acknowledgment** Thanks to Prof. D.Arabelos for valuable comments. The GMT software, <http://gmt.soest.hawaii.edu/>, has been used to produce all figures.



**Fig. 7** Standard deviation of GOCE  $T_{zz}$  minus ITG-GRACE2010s (to 36) for  $1^\circ$  blocks. Units E



**Fig. 8** Mean error estimates for  $25^\circ \times 25^\circ$  blocks at 10 km gravity anomalies minus ITG-Grace2010s to deg. 36



**Fig. 9** Regionally scaled error estimates of gravity anomalies at 10 km minus ITG-Grace2010s to deg. 36. Units mGal

---

## References

- Andersen OB, Remmer O (1982) Non-random errors in the Finnish Levelling of High Precision. *Manuscripta Geodaetica* 7(4):353–373
- Arabelos D, Tscherning CC (1999) Gravity field recovery from airborne gravity gradiometer data using collocation and taking into account correlated errors. *Phys Chem Earth (A)* 24(1):19–25
- Arabelos D, Tscherning CC (2008) Error-covariances of the estimates of spherical harmonic coefficients computed by LSC, using second-order radial derivative functionals associated with realistic GOCE orbits. *J Geodesy*. doi:10.1007/s00190-008-0250-9
- Arabelos D, Forsberg R, Tscherning CC (2007) On the a-priori estimation of error-covariance functions. A feasibility study. *Geoph J Int*. doi:10.1111/j.1365-246X.2007.03460.x
- Arabelos D, Reguzzoni M, Tscherning CC (2013) Global grids of gravity anomalies and vertical gravity gradients at 10 km altitude from GOCE gradient data 2009–2011 and polar gravity. *Geophys Res Abstr* 15:EGU2013-2037
- Balmino G (2009) Efficient propagation of error covariance matrices of gravitational models: application to GRACE and GOCE. *J Geodesy* 83:989–995
- Bingham RJ, Tscherning CC, Knudsen P (2011) An initial investigation of the GOCE Error Variance-covariance matrices in the context of the GOCE user toolbox project. Proceedings 4.th Int. GOCE user workshop, Muenchen March 31 - April 1, 2011. ESA SP-696.
- Mayer-Guerr T, Kurtenbach E, Eicker A (2010) The satellite-only gravity field model ITG-Grace2010s. <http://www.igg.uni-bonn.de/apmg/index.php?id=itg-grace2010>
- Moritz H (1980) *Advanced physical geodesy*. H. Wichmann Verlag, Karlsruhe
- Pavlis NK, Holmes SA, Kenyon SC, Factor JK (2012) The development and evaluation of the Earth Gravitational Model 2008 (EGM2008). *J Geophys Res Solid Earth* (1978–2012) 117(B4)
- Sansò F, Sideris MG (eds) (2013) *Geoid Determination*. Lecture Notes in Earth System Science 110. Springer-Verlag, Berlin-Heidelberg, doi:10.1007/978-3-540-74700-0\_7
- Tscherning CC (1991) The use of optimal estimation for gross-error detection in databases of spatially correlated data. *Bulletin d'Information*, 68:79–89, Bureau Gravimétrique International

# Multivariate Integer Cycle-Slip Resolution: A Single-Channel Analysis

P.J.G. Teunissen and P.F. de Bakker

## Abstract

In this contribution we study the strength of the single-receiver, single-channel GNSS model for instantaneously resolving integer cycle-slips. This will be done for multi-frequency GPS, Galileo and BeiDou, thereby focusing on the challenging case that the slip is due to a simultaneous loss of lock on all frequencies. The analytical analysis presented is supported by means of numerical results.

## Keywords

Global navigation satellite systems (GNSS) • Integer cycle-slip • Single-channel model

## 1 Introduction

Integrity monitoring and quality control can be exercised at different stages of the GNSS data processing chain. These stages range from the single-receiver, single-channel case to the multi-receiver/antenna case, sometimes even with additional constraints included.

In the present contribution, we consider the single-receiver, single-channel model. It is a challenging model as it is the weakest of all, due to the absence of the relative receiver-satellite geometry. In Teunissen and de Bakker (2012), we studied this model's multi-frequency GNSS integrity performance against modelling errors such as code outliers, carrier-phase slips and ionospheric disturbances. By means of the minimal detectable biases (MDBs) of the uniformly most powerful invariant (UMPI) test statistics, it was shown how well these modelling errors can be found.

In (ibid.) the carrier-phase slips were allowed to be non-integer, and therefore real-valued, as well. This implied that

hypothesis testing theory with the DIA-method for the detection, identification and adaptation of the modelling error could be directly applied (Teunissen 1998a). In the present study, however, attention is restricted to *integer* slips only.

The problem of detecting and recovering from integer cycle-slips is an important one and one that has already been considered in several studies, see e.g. (Bisnath et al. 2001; Liu 2010; Carcanague 2012) for the dual-frequency case and (Dai et al. 2009; Xie et al. 2013) for the triple frequency case. It seems, however, that one is of two minds in these studies. On the one hand, namely in the detection step, one treats the slips as real-valued (i.e. integerness is not imposed), while on the other hand, after one has decided that a slip indeed occurred, one imposes the integerness by estimating it as such. This is not consistent and also not needed. Moreover, the cycle-slip detector used in these studies is often not an UMPI-test statistic.

In the present contribution we will not make the above referred to difference between *real-valued* slip-detection and *integer-valued* slip-repair. Instead we estimate the slip directly as an integer and use its probability mass function for evaluation. Consider the following slip-free and slip-biased models,

$$\begin{aligned} \mathcal{H}_0 : E(y) &= Gx, & x \in \mathbb{R}^1, D(y) &= Q_{yy} \\ \mathcal{H}_a : E(y) &= Gx + Hz, & z \in \mathbb{Z}^n, D(y) &= Q_{yy} \end{aligned} \quad (1)$$

P.J.G. Teunissen (✉) • P.F. de Bakker  
Department of Spatial Sciences, GNSS Research Centre, Curtin  
University of Technology, Perth, WA 6845, Australia

Department of Geoscience and Remote Sensing, Delft University  
of Technology, Delft, The Netherlands  
e-mail: p.debakker@curtin.edu.au; p.teunissen@curtin.edu.au

**Table 1** GPS, Galileo, BeiDou frequencies ( $f$ ) and wavelengths ( $\lambda$ ), and zenith-referenced standard deviations of undifferenced code ( $p$ ) and phase ( $\phi$ ) observables

Signal	L1	L2	L5	E1	E5a	E5b	E5	E6	B1	B3	B2
$f$ (MHz)	1575.42	1227.60	1176.45	1575.42	1176.45	1207.14	1191.795	1278.75	1561.1	1268.52	1207.14
$\lambda$ (cm)	19.0	24.4	25.5	19.0	25.5	24.8	25.2	23.4	19.2	23.63	24.83
$p$ (cm)	25	25	15	20	15	15	7	15	31	25	30
$\phi$ (mm)	1.0	1.3	1.3	1.0	1.3	1.3	1.3	1.2	1.4	1.7	1.6

These results were aggregated for GIOVE-B from Simsky et al. (2008), for GPS+GIOVE A/B from de Bakker et al. (2012) and from initial BeiDou results obtained at Curtin's GNSS Research Centre (Khodabandeh and Odolinski, 2013, BeiDou standard deviations, "Personal communication")

with  $y$  normally distributed and  $z$  denoting the integer cycle-slip. Let  $\hat{x}_0$  and  $\check{x}_a$  be the least-squares estimators of  $x$  under  $\mathcal{H}_0$  and  $\mathcal{H}_a$ , respectively. Then

$$\check{x}_a = \hat{x}_0 - G^+ H \check{z} \quad (2)$$

with  $G^+ = (G^T Q_{yy}^{-1} G)^{-1} G^T Q_{yy}^{-1}$  the least-squares inverse of  $G$  and  $\check{z}$  the integer least-squares estimator of  $z$ . In general the distribution of  $\hat{x}_0$  under  $\mathcal{H}_0$  will differ from that of  $\check{x}_a$  under  $\mathcal{H}_a$ . The latter has namely the multi-modal distribution as given in Teunissen (1999a). However, in case the probability of correct integer estimation  $P(\check{z} = z)$ , also known as success-rate, is sufficiently large, then the distribution of  $\check{x}_a$  can be approximated by a normal distribution. In that case the distribution of  $\check{x}_a$  under  $\mathcal{H}_a$  can be considered given by the distribution of  $\hat{x}_0$  under  $\mathcal{H}_0$ :

$$\check{x}_a | \mathcal{H}_a \sim \hat{x}_0 | \mathcal{H}_0 \sim N(x, Q_{\hat{x}_0 \hat{x}_0}) \quad (3)$$

Thus if the success-rate is sufficiently large, the decision whether or not a slip occurred (the so-called detection) is automatically implied in the above correction (2): if the outcome of  $\check{z}$  is zero, then  $\mathcal{H}_0$  is considered true, otherwise it is assumed that a cycle slip is detected.

In this contribution we study the single-receiver, single-channel model's ability to achieve sufficiently high success-rates for the estimated integer cycle-slip vector. This will be done for multi-frequency GPS, Galileo and BeiDou, thereby focusing on the challenging case that the slip is due to a complete loss-of-lock, i.e. a loss of lock on all frequencies. The results show that instantaneous integer cycle-slip resolution is possible for multi-frequency Galileo, but for triple-frequency GPS and BeiDou only for cut-off elevation angles larger than 25°.

## 2 The $n$ -Frequency, 1-Receiver, 2-Epoch Model

### 2.1 The Observation Equations

The carrier phase and pseudorange (code) observation equations of a single receiver that tracks a single satellite on frequency  $f_j = c/\lambda_j$  ( $c$  is speed of light,  $\lambda_j$  is  $j$ th wavelength

and  $j = 1, \dots, n$ ) at time instant  $t$  ( $t = 1, \dots, k$ ), are given as

$$\begin{aligned} \phi_j(t) &= \rho^*(t) - \mu_j i(t) + b_{\phi_j} + n_{\phi_j}(t) \\ p_j(t) &= \rho^*(t) + \mu_j i(t) + b_{p_j} + n_{p_j}(t) \end{aligned} \quad (4)$$

where  $\phi_j(t)$  and  $p_j(t)$  denote the single receiver observed carrier phase and pseudorange, respectively, with corresponding zero mean noise terms  $n_{\phi_j}(t)$  and  $n_{p_j}(t)$ . The unknown parameters are  $\rho^*(t)$ ,  $i(t)$ ,  $b_{\phi_j}$  and  $b_{p_j}$ . The lumped parameter  $\rho^*(t) = \rho(t) + c\delta t_r(t) - c\delta t^s(t) + T(t)$  is formed from the receiver-satellite range  $\rho(t)$ , the receiver and satellite clock errors,  $c\delta t_r(t)$  and  $c\delta t^s(t)$ , respectively, and the tropospheric delay  $T(t)$ . The parameter  $i(t)$  denotes the ionospheric delay expressed in units of range with respect to the *first* frequency. Thus for the  $f_j$ -frequency pseudorange observable, its coefficient is given as  $\mu_j = f_1^2/f_j^2$ . The GPS, Galileo and BeiDou frequencies and wavelengths are given in Table 1. The parameters  $b_{\phi_j}$  and  $b_{p_j}$  are the phase bias and the instrumental code delay, respectively. The phase bias is the sum of the initial phase, the phase ambiguity and the instrumental phase delay.

Both  $b_{\phi_j}$  and  $b_{p_j}$  are assumed to be time-invariant. This is allowed for relatively short time spans, in which the instrumental delays remain sufficiently constant (Liu et al. 2004). The time-invariance of  $b_{\phi_j}$  and  $b_{p_j}$  implies that only time-differences of  $\rho^*(t)$  and  $i(t)$  are estimable. We may therefore just as well formulate the observation equations in time-differenced form. Then the parameters  $b_{\phi_j}$  and  $b_{p_j}$  get eliminated and we obtain

$$\begin{aligned} \phi_j(t, s) &= \rho^*(t, s) - \mu_j i(t, s) + n_{\phi_j}(t, s) \\ p_j(t, s) &= \rho^*(t, s) + \mu_j i(t, s) + n_{p_j}(t, s) \end{aligned} \quad (5)$$

where  $\phi_j(t, s) = \phi_j(t) - \phi_j(s)$ , with a similar notation for the time-difference of the other variates.

Would we have a priori information available about the ionospheric delays, we could model this through the use of additional observation equations. In our case, we do not assume information about the *absolute* ionospheric delays, but rather on the *relative*, time-differenced, ionospheric delays. We therefore have the (pseudo) observation equation

$$i_o(t, s) = i(t, s) + n_i(t, s) \quad (6)$$

with the (pseudo) ionospheric observable  $i_o(t, s)$ . The sample value of  $i_o(t, s)$  is usually taken to be zero.

## 2.2 The Null- And Alternative Hypothesis

If we define  $\phi(t, s) = (\phi_1(t, s), \dots, \phi_n(t, s))^T$ ,  $p(t, s) = (p_1(t, s), \dots, p_n(t, s))^T$ ,  $y = (\phi(t, s)^T, p(t, s)^T, i_o(t, s))^T$ ,  $x = (\rho^*(t, s), i(t, s))^T$ , then the  $n$ -frequency, 2-epoch model can be written in compact matrix-vector form as

$$\mathcal{H}_0: \quad \mathbb{E}(y) = Gx, \quad \text{D}(y) = Q_{yy}, \quad y \in \mathbb{R}^{2n+1}, \quad x \in \mathbb{R}^2 \quad (7)$$

where

$$G = \begin{bmatrix} e_n & -\mu \\ e_n & +\mu \\ 0 & 1 \end{bmatrix}, \quad Q_{yy} = \text{blockdiag}(2Q_{\phi\phi}, 2Q_{pp}, \sigma_{di}^2) \quad (8)$$

with  $e_n$  the  $n$ -vector of ones,  $\mu = (\mu_1, \dots, \mu_n)^T$ ,  $Q_{\phi\phi}$  and  $Q_{pp}$  the  $n \times n$  variance matrices of the undifferenced phase and code observables, and scalar  $\sigma_{di}^2$  the variance of the time-differenced ionospheric delay.

In our computations we assumed the variance matrices  $Q_{\phi\phi}$  and  $Q_{pp}$  to be diagonal with its entries derived from Table 1. Since these entries are zenith-referenced, they still need to be multiplied with an elevation dependent factor to account for the elevation dependency. Based on the customary elevation-dependent models (Euler and Goad 1991), we used the following factors: 1.5 for 30°–40°, 2 for 25°–30°, and 3 for 15°–20° elevation range.

If we assume that the time series of the ionospheric delays can be modeled as a *first-order autoregressive* stochastic process, then

$$\sigma_{di}^2 = 2\sigma_i^2(1 - \beta^{|t-s|}) \quad (9)$$

For two successive epochs we have  $\sigma_{di}^2 = 2\sigma_i^2(1 - \beta)$ , while for larger time-differences the variance will tend to the white-noise value  $\sigma_{di}^2 = 2\sigma_i^2$  if  $\beta < 1$ . Thus  $\sigma_i^2$  and  $\beta$  can be used to model the level and smoothness of the noise in the ionospheric delays. We determined the approximate range of  $\sigma_{di}$ -values as given in Table 2.

Model (7) will be referred to as our null hypothesis  $\mathcal{H}_0$ . This null-hypothesis assumes that no loss-of-lock occurred between the two epochs. Would such loss-of-lock occur, however, then one or more of the  $n$  carrier-phases may become biased by an unknown number of integer cycle slips. Here we assume the worst scenario, namely that all  $n$  of the carrier-phases are affected by the loss-of-lock. Hence, instead of a one-dimensional integer cycle-slip, we consider the case of an integer cycle-slip vector  $z \in \mathbb{Z}^n$ . The model

**Table 2** Approximate range of values for  $\sigma_{di}$  (m) when sampling with intervals of 1, 10 and 30 s, respectively, using a 10 degree cut-off elevation angle

	min. $\sigma_{di}$	max. $\sigma_{di}$
1 s	$1.5 \times 10^{-3}$	$3 \times 10^{-3}$
10 s	$2 \times 10^{-3}$	$10^{-2}$
30 s	$4.5 \times 10^{-3}$	$2.5 \times 10^{-2}$

These values were obtained for mid-latitude (Delft) under moderate ionospheric conditions

for such ‘loss-of-lock’ hypothesis is given by the alternative hypothesis

$$\mathcal{H}_a: \quad \mathbb{E}(y) = [G, H] \begin{bmatrix} x \\ z \end{bmatrix}, \quad \text{D}(y) = Q_{yy}, \quad z \in \mathbb{Z}^n \quad (10)$$

with  $H = [\Lambda, 0, 0]^T$  and  $\Lambda = \text{diag}[\lambda_1, \dots, \lambda_n]$ .

## 3 Estimability of Multivariate Carrier-Phase Slip

### 3.1 Variance Matrix of Multivariate Cycle-Slip Estimator

From the structure of  $[G, H]$  in (10), it follows that the carrier phase vector  $\phi(t, s)$  will not contribute to the estimation of the parameters  $\rho^*(t, s)$  and  $i(t, s)$  under  $\mathcal{H}_a$ . These parameters are therefore solely determined by the code observables and a priori ionospheric information. As a consequence, the two-epoch cycle-slip estimator is given as the difference  $\hat{z} = \Lambda^{-1}(\phi(t, s) - \hat{\phi}(t, s))$ , where  $\hat{\phi}(t, s) = e_n \hat{\rho}^*(t, s) - \mu \hat{i}(t, s)$  is the least-squares phase estimator based solely on the code observables and a priori ionospheric information. Solving for  $\hat{\rho}^*(t, s)$  and  $\hat{i}(t, s)$ , followed by applying the variance propagation law to  $\hat{z} = \Lambda^{-1}(\phi(t, s) - e_n \hat{\rho}^*(t, s) + \mu \hat{i}(t, s))$  gives then the variance matrix of the multivariate cycle-slip estimator. The result is given in the following Lemma.

**Lemma 1 (Variance Matrix Multivariate Slip)** *The variance matrix of the least-squares estimator of  $z$  under  $\mathcal{H}_a$  is given as*

$$Q_{\hat{z}\hat{z}} = \Lambda^{-1} \left( \underbrace{2Q_{\phi\phi}}_{\text{phase}} + \underbrace{2P_{e_n} Q_{pp} P_{e_n}^T}_{\text{code; rank=1}} + \underbrace{(R_{e_n} \mu) \sigma_{di}^2 (R_{e_n} \mu)^T}_{\text{ionosphere; rank=1}} \right) \Lambda^{-1} \quad (11)$$



with  $R_{e_n} = I_n + P_{e_n}$ ,  $P_{e_n} = e_n(e_n^T Q_{pp}^{-1} e_n)^{-1} e_n^T Q_{pp}^{-1}$ ,  $P_{e_n}^\perp = I_n - P_{e_n}$ , and where

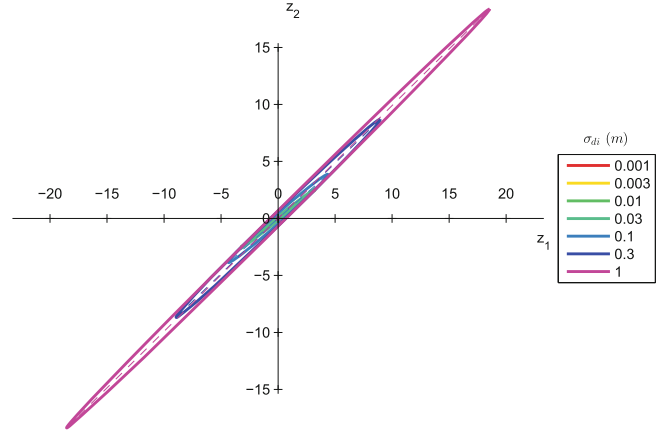
$$\sigma_{\hat{d}_i}^2 = \frac{\sigma_{d_i}^2}{1 + \frac{1}{2} \sigma_{d_i}^2 \underbrace{\|\mu - \bar{\mu}_p\|_{Q_{pp}}^2}_{\text{frequency-diversity}}}, \quad \bar{\mu}_p = P_{e_n} \mu \quad (12)$$

◇

Note that the slip variance matrix is a sum of three terms, the entries of which may differ substantially in size. The first matrix term in this sum is governed by the precision of the phase observables and will therefore have small entries. The second matrix term is governed by the precision of the code observables and will therefore have generally much larger entries than the first matrix term in the sum. The third matrix term depends, next to the precision of the code observables, also on  $\mu$  and  $\sigma_{d_i}^2$ . Its entries will become smaller, if  $\sigma_{d_i}^2$  gets smaller. This happens for smaller  $\sigma_{d_i}^2$  (smoother ionospheric delays) and/or larger  $\|\mu - \bar{\mu}_p\|_{Q_{pp}}^2$  (better code precision and/or larger frequency diversity). Thus if  $\sigma_{d_i}^2 = \infty$ , frequency diversity is needed (i.e.  $\|\mu - \bar{\mu}_p\|_{Q_{pp}}^2 \neq 0$ ) so as to avoid the entries of the third matrix term in (11) to become infinite.

Due to the relative poor code precision (as compared to phase) the confidence ellipsoid of  $\hat{z}$  is usually very elongated. This elongation gets larger if  $\sigma_{d_i}$  gets larger. Figure 1 shows the dual-frequency L1/L2 confidence ellipse of the cycle-slip estimator  $\hat{z}$  for different values of the ionospheric standard deviation. Such an elongated ellipse implies that the component of the cycle-slip in the direction of elongation is poorly estimable, whereas the component orthogonal to it, is very well estimable. Note that for larger  $\sigma_{d_i}$ , the elongation approximately points into the  $z_1 = z_2$ -direction. This explains why in those cases the wide-lane combination has good precision. But also note, that the ellipse rotates away from the  $z_1 = z_2$ -direction as  $\sigma_{d_i}$  gets smaller. This implies that for those cases other combinations than the wide-lane have better precision.

If we consider more than  $n = 2$  frequencies, it is important to point out that the second and third matrix terms of (11) are both of rank 1. This implies that the elongation of the confidence ellipsoid of  $\hat{z}$  remains restricted to two dimensions only, irrespective the value of  $n \geq 2$ . This indicates that in higher dimensions one should be able to profit from the increase in frequencies and thus better be able to successfully resolve the integer cycle-slip vector in case of a loss-of-lock. To what extent this is possible, will be investigated further in the next sections.



**Fig. 1** Dual-frequency L1/L2, 95% confidence ellipse of cycle-slip estimator  $\hat{z}$ , for different values of ionospheric standard deviation  $\sigma_{d_i}$  (cf. 11). Units along axes are cycles

#### 4 The ADOP of the Multivariate Slip

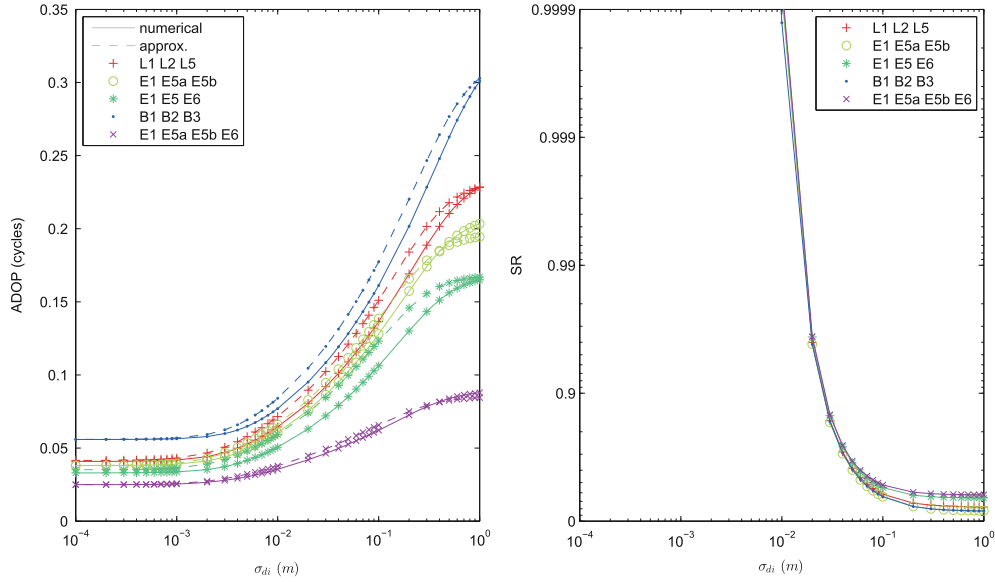
The ADOP was introduced in Teunissen (1997) as an easy-to-compute scalar diagnostic to measure the *intrinsic* model strength for successful ambiguity resolution. The ADOP is defined as the square-root of the ambiguity variance determinant taken to the power one over the number of integer ambiguities, which, in the present case, is the dimension of the integer cycle-slip vector  $z$ ,

$$\text{ADOP} = |Q_{\hat{z}\hat{z}}|^{\frac{1}{2n}} \quad (\text{cycle}) \quad (13)$$

The ADOP has the important property that it is invariant against the choice of ambiguity parametrization. Since all admissible ambiguity transformations can be shown to have a determinant of  $\pm 1$ , the ADOP does not change when one changes the definition of the ambiguities. It therefore measures the intrinsic precision of the ambiguities. As a rule-of-thumb, an ADOP smaller than about 0.10 cycle, corresponds to an ambiguity success-rate larger than 0.999.

From the variance matrix of Lemma 1, an analytical closed-form formula can be derived for the corresponding ADOP. A useful and easy-to-interpret approximation to this analytical expression is given in the following lemma.

**Lemma 2 (ADOP Rule-of-Thumb)** *If  $Q_{\phi\phi} = \sigma_\phi^2 I_n$  and  $Q_{pp} = \sigma_p^2 I_n$ , then the ambiguity dilution of precision of the multivariate cycle-slip can be approximated as*



**Fig. 2** *Left:* Multi-frequency cycle-slip ADOPs, as function of  $\sigma_{di}$ , for GPS, Galileo and BeiDou. The *dashed curves* are based on the approximation (14), while the *full curves* are based on (13); *Right:*

Multi-frequency, zenith-referenced, bootstrapped success-rates (SR) of resolving the integer cycle-slip vector  $z$ , as function of  $\sigma_{di}$ , for GPS, Galileo and BeiDou

$$\text{ADOP} \approx \left( \frac{\sqrt{2}}{\bar{\lambda}} \right) \left[ \underbrace{\alpha \left( \sigma_{\phi}^{n-1} \sigma_p \right)^2}_{\text{iono-fixed}} + (1 - \alpha) \underbrace{\left( \sigma_{\phi}^{n-2} \sigma_p^2 \right)^2}_{\text{iono-float}} \right]^{\frac{1}{2n}} \quad (14)$$

with  $\bar{\lambda} = \prod_{i=1}^n \lambda_i^{\frac{1}{n}}$  and  $\alpha = [1 + \frac{1}{2} \sum_{i=1}^n (\mu_i - \bar{\mu})^2 \sigma_{di}^2 / \sigma_p^2]^{-1}$ .  
 $\diamond$

Note that  $\bar{\lambda}$  is the geometric average of the wavelengths, whereas  $\bar{\mu}$  is the arithmetic average of the  $\mu_i$ . Also note that the term within the square brackets is a *convex* combination driven by the scalar  $\alpha$ . For  $\alpha = 1$ , the ionosphere-fixed result is obtained, while for  $\alpha = 0$  the ionosphere-float result is obtained. Thus

$$\begin{aligned} \text{ADOP}_{\alpha=1} &\approx \left( \frac{\sqrt{2}}{\bar{\lambda}} \right) \left( \sigma_{\phi}^{n-1} \sigma_p \right)^{\frac{1}{n}} && (\text{iono} - \text{fixed}) \\ \text{ADOP}_{\alpha=0} &\approx \left( \frac{\sqrt{2}}{\bar{\lambda}} \right) \left( \sigma_{\phi}^{n-2} \sigma_p^2 \right)^{\frac{1}{n}} && (\text{iono} - \text{float}) \end{aligned} \quad (15)$$

This clearly shows the roles played by the contributing factors: wavelengths ( $\bar{\lambda}$ ), phase precision ( $\sigma_{\phi}$ ), code precision ( $\sigma_p$ ) and number of frequencies ( $n$ ). It also shows the very different contributions of phase and code to either the ionosphere-fixed case or the ionosphere-float case. For instance, for the single-frequency case,  $n = 1$ , the ionosphere-fixed ADOP is driven by the code-precision only, whereas for the ionosphere-float case, the ADOP gets further magnified by  $\sigma_p / \sigma_{\phi}$ , i.e. a factor of about 100. For an  
 $\diamond$

arbitrary number of frequencies the ratio between the two ADOPs is given as

$$\frac{\text{ADOP}_{\alpha=1}}{\text{ADOP}_{\alpha=0}} \approx \left( \frac{\sigma_{\phi}}{\sigma_p} \right)^{\frac{1}{n}} \quad (16)$$

Figure 2(Left) shows, as function of  $\sigma_{di}$ , the multi-frequency, loss-of-lock cycle-slip ADOPs for GPS, Galileo and BeiDou. These results are promising as the ADOPs are all below 0.1 cycle for most of the relevant  $\sigma_{di}$ -range (cf. Table 2). In the next section we will study their success-rates.

## 5 Multi-Frequency, Cycle-Slip Resolution Success-Rates

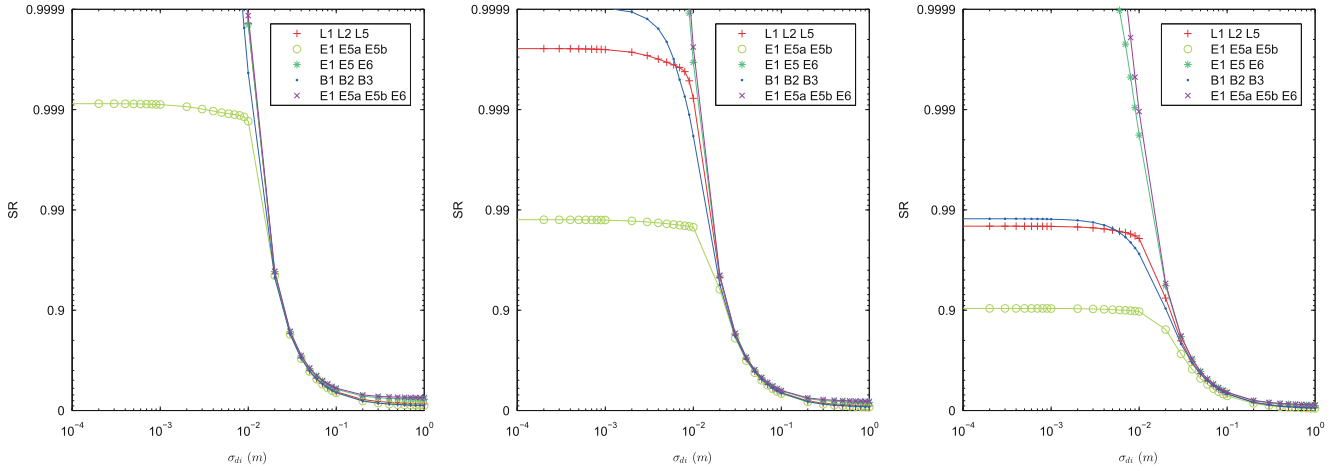
### 5.1 Bootstrapped Success-Rates

Different integer estimators can be used to solve for the integer cycle-slip vector  $z$ . Three popular integer estimators are integer rounding, integer bootstrapping and integer least-squares. As the following theorem shows, there exists a clear ordering among these three estimators.

**Theorem 1 (Teunissen 1999b)** *Let  $\hat{z} \sim N(z, Q_{\hat{z}})$  and let  $\check{z}_{\text{IR}}$ ,  $\check{z}_{\text{IB}}$ , and  $\check{z}_{\text{ILS}}$  denote the estimators of integer rounding, integer bootstrapping and integer least-squares, respectively. Then their success-rates are ordered as*

$$P(\check{z}_{\text{IR}} = z) \leq P(\check{z}_{\text{IB}} = z) \leq P(\check{z}_{\text{ILS}} = z) \quad (17)$$

$\diamond$



**Fig. 3** Multi-frequency, GPS, Galileo and BeiDou, bootstrapped success-rates (SR) of resolving the integer cycle-slip vector  $z$ , as function of  $\sigma_{di}$ , for different elevation angles (Left:  $30^\circ$ – $40^\circ$ , Middle:  $25^\circ$ – $30^\circ$ , Right:  $15^\circ$ – $20^\circ$ )

Integer rounding (IR) is the simplest, but it also has the poorest success rate. Integer least-squares (ILS) is the most complex, but also has the highest success rate of all. Integer bootstrapping (IB) sits in between. It does not need an integer search as is the case with ILS, and it does not completely neglect the information content of the ambiguity variance matrix as IR does. Moreover, bootstrapping is the only integer estimator for which an easy-to-use and exact expression can be given of its success-rate. This success-rate is given in the following theorem.

**Theorem 2 (Teunissen 1998b)** Let  $\hat{z} \sim N(z, Q_{\hat{z}\hat{z}})$ . Then the success-rate of integer bootstrapping is given as

$$P(\check{z}_{\text{IB}} = z) = \prod_{i=1}^n \left( 2\Phi\left(\frac{1}{2\sigma_{z_i|I}}\right) - 1 \right) \quad (18)$$

with  $\Phi(x) = \int_{-\infty}^x \frac{1}{\sqrt{2\pi}} \exp\{-\frac{1}{2}v^2\} dv$  and  $\sigma_{z_i|I}$  the standard deviation of the  $i$ th entry of  $\hat{z}$ , conditioned on the previous  $I = \{1, \dots, (i-1)\}$  entries.  $\diamond$

In this contribution we used the bootstrapped success-rate. It has been used after applying the decorrelating transformation of the LAMBDA method (Teunissen 1995a). For such decorrelated cases namely, the bootstrapped success-rate becomes a sharp lower bound of the ILS success-rate.

Figure 2(Right) shows the cycle-slip vector resolution success-rates for the same frequency-combinations as shown in Fig. 2(Left). As predicted by the ADOPs, the success-rates are indeed very high, 99.99% or larger for  $\sigma_{di} \leq 10^{-2}$ . However, the results of Fig. 2(Left and Right) are zenith-referenced and therefore only hold true for an elevation angle of  $90^\circ$ .

## 5.2 Success-Rates Versus Elevation

For the same frequency combinations as before, Fig. 3 now shows the success-rates for different elevation angles. The results clearly show that the success-rates get smaller as the elevation angle gets smaller. The results also show that some frequency combinations are more sensitive than others to the changes in elevation angle. For instance, the success-rate of the Galileo frequency combination E1-E5a-E5b, is the first to drop in value when the elevation gets smaller, see Fig. 3(Left). The combinations that retain a large success-rate, even at low elevation, are the triple and quadruple Galileo combinations E1-E5-E6 and E1-E5a-E5b-E6, see Fig. 3(Right). For these combinations one can expect to have at least a 99.9% success-rate up to  $\sigma_{di} = 10^{-2}$  m.

Also note, although the triple-frequency success-rates of GPS and BeiDou are too low for low elevations, their success-rates are still large for elevations in the  $25^\circ$ – $30^\circ$  range, see Fig. 3(Middle). This is a relevant finding, since as shown in Odolinski et al. (2013a), positioning with a combined GPS and BeiDou system will allow one to make use of much higher cut-off elevation angles. Similarly, it was shown in Nadarajah and Teunissen (2013) that this is true for GPS+BeiDou MC-LAMBDA attitude determination as well. Hence, for those positioning and attitude-determination applications, instantaneous loss-of-lock integer cycle slip resolution will become feasible.

## References

- Bisnath SB, Kim D, Langley RB (2001) Carrier Phase Slips: a new approach to an old problem. *GPS World*, 12(5), 46–51

- Carcanague S (2012) Real-time geometry-based cycle slip resolution technique for single-frequency PPP and RTK. In: Proceedings of the 25th ITM, pp 1136–1148
- Dai Z, Knedlik S, Loffeld O (2009) Instantaneous triple-frequency GPS cycle-slip detection and repair. *Int J Navig Obs* p 15. Doi:10.1155/2009/407231
- de Bakker PF, Tiberius CCJM, van der Marel H, van Bree RJP (2012) Short and zero baseline analysis of GPS L1 C/A, L5Q, GIOVE E1B, and E5aQ signals. *GPS Solutions*, 16, 53–64
- Euler H-J, Goad CC (1991) On optimal filtering of GPS dual frequency observations without using orbit information. *Bull Géod* 65:130–143
- Liu Z (2010) A new automated cycle-slip detection and repair method for a single dual-frequency GPS receiver. *J Geod*. Doi:10.1007/s00190-010-0426-y
- Liu X, Tiberius CCJM, de Jong K (2004) Modelling of differential single difference receiver clock bias for precise positioning. *GPS Solution* 4:209–221
- Nadarajah N, Teunissen PJG (2013) Instantaneous GPS/BeiDou/Galileo attitude determination: a single-frequency robustness analysis under constrained environments. In: Proceedings of the institute of navigation pacific PNT 2013, pp 1088–1103
- Odolinski R, Teunissen PJG, Odijk D (2013a) An analysis of combined COMPASS/BeiDou-2 and GPS single- and multiple-frequency RTK positioning. In: Proceedings of the institute of navigation pacific PNT 2013, pp 69–90
- Simsy A, Mertens D, Sleewaegen J-M, De Wilde W, Hollreiser M, Crisci M (2008) MBOC vs. BOC(1,1) multipath comparison based on GIOVE-B data. *Inside GNSS*, September/October, pp 36–39
- Teunissen PJG (1995a) The least-squares ambiguity decorrelation adjustment: a method for fast GPS integer ambiguity estimation. *J Geod* 70:65–82
- Teunissen PJG (1997) A canonical theory for short GPS baselines. Part IV: precision versus reliability. *J Geod* 71:513–525
- Teunissen PJG (1998a) Minimal detectable biases of GPS data. *J Geod* 72:236–244
- Teunissen PJG (1998b) Success probability of integer GPS ambiguity rounding and bootstrapping. *J Geod* 72:606–612
- Teunissen PJG (1999a) The probability distribution of the GPS baseline for a class of integer ambiguity estimators. *J Geod* 73:275–284
- Teunissen PJG (1999b) An optimality property of the integer least-squares estimator. *J Geod* 73:587–593
- Teunissen PJG, de Bakker PF (2012) Single-receiver single-channel multi-frequency GNSS integrity: outliers, slips, and ionospheric disturbances. *J Geod* Doi: 10.1007/s00190-012-0588-x
- Xie K, Chai H, Wang M, Pan Z (2013) Cycle slip detection and repair with different sampling interval based on compass triple-frequency. In: Sun J, et al (eds) Proceedings of China Satellite Navigation Conference (CSNC) 2013, pp 291–303

---

# Theory of Earth Rotation Variations

Richard S. Gross

---

## Abstract

The theory currently used to study small variations in the Earth's rotation that occur on time scales longer than a day is reviewed. This theory is based on the principle of the conservation of angular momentum. Using this principle, changes in the rotation of the solid Earth can be shown to be caused either by changes in the mass distribution of the solid Earth or by torques acting on the solid Earth. Such torques can be caused, for example, by the motion of the atmosphere and oceans or by the gravitational effect of the Sun, Moon, and planets. When applying this principle to the rotation of the Earth a number of simplifying assumptions are made including: (1) linearity; (2) axisymmetry; (3) equilibrium oceans; (4) Tisserand mean-mantle; (5) the core is uncoupled from the mantle; and (6) the rotational variations occur on time scales much longer than a day. While the resulting theory has been successfully used in the past to interpret the observed variations in the Earth's rotation, it is argued that the accuracy of the observations has improved to the point that the current theory is no longer adequate and that a new, more accurate theory of the Earth's rotation is needed.

---

## Keywords

Earth rotation • Length-of-day • Polar motion • Universal time

---

## 1 Introduction

The Earth's rotation changes on all observable time scales, from subdaily to decadal and longer. The wide range of time scales on which the Earth's rotation changes reflects the wide variety of processes that are causing it to change, including external tidal forces, surficial fluid processes involving the atmosphere, oceans, and hydrosphere, and internal processes acting both within the solid Earth itself and between the fluid core and the solid Earth. These changes in Earth rotation are usually studied using the principle that angular momentum is conserved as it is transferred between the solid Earth and the fluid regions with which it is in contact. Using

the principle of the conservation of angular momentum the equations governing small variations in both the rate of rotation and in the position of the rotation vector with respect to the Earth's crust are derived. As a prelude to developing an improved, more accurate theory of the Earth's rotation, particular attention is paid to the various assumptions and approximations that are made when deriving these equations. The approach taken here to derive the equations that are currently used to study small changes in the Earth's rotation follows that of Smith and Dahlen (1981) and Wahr (1982, 1983, 2005) and has been recently reviewed by Gross (2007).

---

## 2 Long-Period Equations of Motion

### 2.1 Rigid Body Rotation

The equation expressing conservation of angular momentum within a rotating terrestrial reference frame is (e.g.,

---

R.S. Gross (✉)  
Jet Propulsion Laboratory, California Institute of Technology,  
4800 Oak Grove Drive, Pasadena, CA 91109, USA  
e-mail: [Richard.Gross@jpl.nasa.gov](mailto:Richard.Gross@jpl.nasa.gov)

Goldstein 1950):

$$\frac{\partial}{\partial t} [\mathbf{h}(t) + \mathbf{I}(t) \cdot \boldsymbol{\omega}(t)] + \boldsymbol{\omega}(t) \times [\mathbf{h}(t) + \mathbf{I}(t) \cdot \boldsymbol{\omega}(t)] = \boldsymbol{\tau}(t) \quad (1)$$

where  $\boldsymbol{\tau}(t)$  represents the external torques that are acting on the Earth and where the angular momentum  $\mathbf{L}(t)$  has been separated into two parts: a part  $\mathbf{h}(t)$  due to motion relative to the rotating reference frame, and a part due to changes in the distribution of the Earth's mass and hence in its inertia tensor  $\mathbf{I}(t)$ :

$$\mathbf{L}(t) = \mathbf{h}(t) + \mathbf{I}(t) \cdot \boldsymbol{\omega}(t) \quad (2)$$

Note that strictly speaking  $\boldsymbol{\omega}(t)$  in these equations is the angular velocity of the rotating reference frame with respect to inertial space. But since the rotating frame will be later attached to the solid body of the Earth, it is also taken to be the angular velocity of the solid Earth with respect to inertial space.

The first simplifying assumption that is made when deriving the equations governing the Earth's variable rotation is to assume that the variations are 'small'. This is reasonable because observations taken over the last century show that the Earth's rotation deviates only slightly from a state of uniform rotation, being a few parts in  $10^8$  in speed, corresponding to changes of a few milliseconds in length-of-day, and being about a part in  $10^6$  in the position of the rotation pole with respect to the Earth's crust, corresponding to variations of several hundred milliarcseconds (mas) in polar motion. Assuming that the variations are 'small', Eq. (1) can be simplified by linearizing it.

Let the Earth be initially in a state of uniform rotation  $\boldsymbol{\omega}_0$  about the  $z$ -coordinate axis of the terrestrial reference frame and let the frame be oriented within the Earth in such a manner that the inertia tensor of the Earth is diagonal in that frame. In this initial state, the Earth is rotating at a constant rate  $\Omega$  about its figure axis, there are no mass displacements, and there is no relative angular momentum.

Now consider some general perturbation to this initial state that introduces both mass displacements and relative angular momentum. The perturbed relative angular momentum, instantaneous rotation vector, and inertia tensor of the Earth is:

$$\begin{aligned} \mathbf{h}(t) &= \mathbf{h}_0 + \Delta \mathbf{h}(t) \\ &= h_x(t) \hat{\mathbf{x}} + h_y(t) \hat{\mathbf{y}} + h_z(t) \hat{\mathbf{z}} \end{aligned} \quad (3)$$

$$\begin{aligned} \boldsymbol{\omega}(t) &= \boldsymbol{\omega}_0 + \Delta \boldsymbol{\omega}(t) \\ &= \Omega \hat{\mathbf{z}} + \Omega [m_x(t) \hat{\mathbf{x}} + m_y(t) \hat{\mathbf{y}} + m_z(t) \hat{\mathbf{z}}] \end{aligned} \quad (4)$$

$$\begin{aligned} \mathbf{I}(t) &= \mathbf{I}_0 + \Delta \mathbf{I}(t) \\ &= \begin{pmatrix} A & 0 & 0 \\ 0 & B & 0 \\ 0 & 0 & C \end{pmatrix} + \begin{pmatrix} \Delta I_{xx}(t) & \Delta I_{xy}(t) & \Delta I_{xz}(t) \\ \Delta I_{xy}(t) & \Delta I_{yy}(t) & \Delta I_{yz}(t) \\ \Delta I_{xz}(t) & \Delta I_{yz}(t) & \Delta I_{zz}(t) \end{pmatrix} \end{aligned} \quad (5)$$

where the hat denotes a vector of unit length, the  $h_i(t)$  are the elements of the time-dependent perturbation  $\Delta \mathbf{h}(t)$  to the initial value  $\mathbf{h}_0$  of the total relative angular momentum  $\mathbf{h}(t)$ , the  $\Omega m_i(t)$  are the elements of the time-dependent perturbation  $\Delta \boldsymbol{\omega}(t)$  to the initial value  $\boldsymbol{\omega}_0$  of the total rotation vector  $\boldsymbol{\omega}(t)$ , the  $\Delta I_{ij}(t)$  are the elements of the time-dependent perturbation  $\Delta \mathbf{I}(t)$  to the initial value  $\mathbf{I}_0$  of the total inertia tensor  $\mathbf{I}(t)$ , and  $A$ ,  $B$ , and  $C$ , the non-zero elements of the initial inertia tensor  $\mathbf{I}_0$ , are the mean principal moments of inertia of the Earth ordered such that  $A < B < C$ . Note that  $\mathbf{h}_0 = 0$  because there is no relative angular momentum in the initial state.

In Eqs. (3)–(5), the perturbations to the initial state are arbitrary. But if it is now assumed that the perturbations are small, so  $h_i(t) \ll \Omega C$ ,  $m_i(t) \ll 1$ , and  $\Delta I_{ij}(t) \ll C$ , and if just terms of first order in small quantities are kept, then the equatorial and axial components of the conservation of angular momentum equation (1) becomes:

$$\begin{aligned} \frac{1}{\sigma_r} \frac{\partial m_x(t)}{\partial t} + \left[ \frac{B(C-B)}{A(C-A)} \right]^{1/2} m_y(t) = \\ - \left( \frac{B}{A} \right)^{1/2} \left[ \frac{1}{\Omega} \frac{\partial \phi_{r,x}(t)}{\partial t} - \phi_{r,y}(t) \right] \end{aligned} \quad (6)$$

$$\begin{aligned} \frac{1}{\sigma_r} \frac{\partial m_y(t)}{\partial t} - \left[ \frac{A(C-A)}{B(C-B)} \right]^{1/2} m_x(t) = \\ - \left( \frac{A}{B} \right)^{1/2} \left[ \frac{1}{\Omega} \frac{\partial \phi_{r,y}(t)}{\partial t} + \phi_{r,x}(t) \right] \end{aligned} \quad (7)$$

$$\frac{1}{\Omega} \frac{\partial m_z(t)}{\partial t} = - \frac{1}{\Omega} \frac{\partial \phi_{r,z}(t)}{\partial t} \quad (8)$$

where the external torques  $\boldsymbol{\tau}(t)$  have been set to zero,

$$\sigma_r^2 = \left( \frac{C-A}{A} \right) \left( \frac{C-B}{B} \right) \Omega^2 \quad (9)$$

and the  $\phi_{r,i}(t)$ , known as excitation functions, are:

$$\phi_{r,x}(t) = \frac{h_x(t) + \Omega \Delta I_{xz}(t)}{\Omega \sqrt{(C-A)(C-B)}} \quad (10)$$

$$\phi_{r,y}(t) = \frac{h_y(t) + \Omega \Delta I_{yz}(t)}{\Omega \sqrt{(C-A)(C-B)}} \quad (11)$$

$$\phi_{r,z}(t) = \frac{1}{C\Omega} [h_z(t) + \Omega \Delta I_{zz}(t)] \quad (12)$$

Equations (6)–(12) describe changes in the rotation of triaxial bodies that are subject to small perturbing excitation and can be applied directly to triaxial bodies that are rigid. But the Earth is not rigid – it has a fluid atmosphere and oceans, a fluid core, and a solid crust and mantle that can deform in response not only to the applied excitation but also to changes in rotation caused by the excitation. As shown in the next section, deriving relatively simple equations for the rotational motion of the Earth that account for its deformable nature will require assuming more than what has been assumed so far, namely, that the variations in relative angular momentum, inertia tensor, and rotation are ‘small’.

## 2.2 Non-rigid Body Rotation

In general, all parts of the Earth including its fluid core, solid crust and mantle, and fluid atmosphere and oceans will respond to changes in rotation. And this response will, in general, involve both motion and mass displacements, causing changes in both relative angular momentum and in the inertia tensor. The contribution of all parts of the Earth to changes in relative angular momentum and in the inertia tensor caused by changes in rotation thus needs to be considered.

Because the crust and mantle of the Earth can deform, they can undergo motion relative to the rotating reference frame and hence can contribute to relative angular momentum. This is taken into account in Earth rotation theory by letting the rotating terrestrial reference frame in the perturbed state be oriented in such a manner that the relative angular momentum due to motion of the crust and mantle vanishes. In this frame, known as the Tisserand mean-mantle frame (Tisserand 1891), the motion of the atmosphere, oceans, and core have relative angular momentum, but the motion of the crust and mantle does not. Furthermore, if it is assumed that the oceans stay in equilibrium as the rotation of the solid Earth changes so that no oceanic currents are generated by changes in rotation, a reasonable assumption if the rotational variations occur on time scales much longer than a day, then there are also no changes in relative angular momentum due to motion of the oceans. Effects of the atmosphere can be ignored here because of its relatively small mass and inertia compared to the oceans. Thus, under the assumptions of the Tisserand mean-mantle frame and of equilibrium oceans only the core will contribute to changes in relative angular momentum caused by changes in rotation.

Hough (1895) showed for a homogeneous, incompressible, non-dissipative, fluid core with a rigid, elliptical core-mantle boundary that, in the frequency domain, the

contribution of the core to changes in relative angular momentum  $\delta h_i(\sigma)$  can be written as:

$$\begin{pmatrix} \delta h_x(\sigma) \\ \delta h_y(\sigma) \\ \delta h_z(\sigma) \end{pmatrix} = \begin{pmatrix} E & iE' & 0 \\ -iE' & E & 0 \\ 0 & 0 & \tilde{E} \end{pmatrix} \begin{pmatrix} m_x(\sigma) \\ m_y(\sigma) \\ m_z(\sigma) \end{pmatrix} \quad (13)$$

where to first order in the ellipticity  $\varepsilon_c$  of the core-mantle boundary and at frequencies  $\sigma \ll \Omega$ :

$$E = (\sigma^2/\Omega) A_c \quad (14)$$

$$E' = -\sigma(1 - \varepsilon_c) A_c \quad (15)$$

$$\tilde{E} = -\Omega C_c \quad (16)$$

where  $A_c$  and  $C_c$  are the equatorial and axial principal moments of inertia of the core. Equation (16) for  $\tilde{E}$  is obtained by assuming that there is no coupling between the core and the mantle.

Dahlen (1976) studied the passive influence of the oceans on the Earth’s rotation, including the changes in the Earth’s inertia tensor  $\delta I_{ij}$  caused by changes in the rotation of the Earth. In the absence of oceans he found:

$$\begin{pmatrix} \delta I_{xz} \\ \delta I_{yz} \\ \delta I_{zz} \end{pmatrix} = \frac{a^5 \Omega^2}{3G} \begin{pmatrix} k_2 & 0 & 0 \\ 0 & k_2 & 0 \\ 0 & 0 & n_o + \frac{4}{3}k_2 \end{pmatrix} \begin{pmatrix} m_x \\ m_y \\ m_z \end{pmatrix} \quad (17)$$

where  $a$  is the mean radius of the Earth,  $G$  is the Newtonian gravitational constant,  $k_2$  is the second-degree body tide Love number of the whole Earth, and  $n_o$  comes from the change in the mean moment of inertia of the Earth caused by the term in the centripetal potential that gives rise to purely radial deformations.

The approach followed by Dahlen (1976) neglects the effects of rotation on the Earth’s elasto-dynamics and also neglects the Earth’s ellipticity as it modifies the Earth’s elastic response. So for the wobble, Smith and Dahlen (1981) took a hybrid approach in which the body tide Love number  $k_2$  is replaced by an oceanless, wobble-effective Love number  $k_w$  that is computed from normal mode theory. This not only allows the deformation of the crust and mantle to be more accurately modeled, it also allows the effects of mantle elasticity and core structure to be included in the theory of Hough (1895) for the response of the core to changes in the rotation of the mantle.

In the presence of oceans, Dahlen (1976) found that their equilibrium influence can be written in terms of an “oceanic Love number”  $\Delta k_{ocn}$  that modifies the body tide

(or wobble-effective) Love number. Because of the non-uniform distribution of the oceans, this oceanic Love number is different for each component  $m_i$  of the Earth's rotation. However, the average of the equatorial components was taken to define a mean oceanic Love number  $\Delta k_{ocn,w}$  for the wobble, distinct from the oceanic Love number  $\Delta k_{ocn,s}$  for the spin. The non-uniform distribution of the oceans also couples the equatorial components of the Earth's rotation to each other and to the axial component by introducing off-diagonal elements in Eq. (17). However, this coupling is weak, with the off-diagonal elements being about three orders of magnitude smaller in numerical value than the diagonal elements (Gross 2007). To first order, the coupling introduced by the non-uniform distribution of the oceans can therefore be ignored, both the spin-wobble coupling and the coupling between the equatorial components.

In the above discussion about the changes in the inertia tensor caused by changes in rotation the mantle was assumed to be elastic. But the mantle is anelastic, causing a modification  $\Delta k_{an}$  to the Love number. Unfortunately, accurate models of mantle anelasticity are not available, so this modification to the Love number cannot be accurately computed. In order to account for mantle anelasticity in the rotational equations of motion at the frequencies of interest here, namely, at frequencies  $\sigma < \Omega$ , a hybrid approach was taken. The wobble-effective Love number  $k_w$  as modified by equilibrium oceans and mantle anelasticity was eliminated from the equations of motion by substituting the observed complex-valued frequency  $\sigma_o$  of the Chandler wobble for its theoretical value (which depends on the wobble-effective Love number  $k_w$  as modified by equilibrium oceans and mantle anelasticity). The resulting final equations of motion are (for details see Gross 2007):

$$\frac{1}{\sigma_o} \frac{\partial m_x(t)}{\partial t} + m_y(t) = \chi_y(t) - \frac{1}{\Omega} \frac{\partial \chi_x(t)}{\partial t} \quad (18)$$

$$\frac{1}{\sigma_o} \frac{\partial m_y(t)}{\partial t} - m_x(t) = -\chi_x(t) - \frac{1}{\Omega} \frac{\partial \chi_y(t)}{\partial t} \quad (19)$$

$$m_z(t) = -\chi_z(t) \quad (20)$$

where the excitation functions  $\chi_i(t)$  are:

$$\chi_x(t) = \frac{h_x(t) + \Omega [1 + (k'_2 + \Delta k'_{an})] \Delta I_{xz}(t)}{[C - A' + A'_m + \varepsilon_c A_c] \sigma_o} \quad (21)$$

$$\chi_y(t) = \frac{h_y(t) + \Omega [1 + (k'_2 + \Delta k'_{an})] \Delta I_{yz}(t)}{[C - A' + A'_m + \varepsilon_c A_c] \sigma_o} \quad (22)$$

$$\chi_z(t) = k_r \frac{h_z(t) + \Omega [1 + \alpha_3 (k'_2 + \Delta k'_{an})] \Delta I_{zz}(t)}{C_m \Omega} \quad (23)$$

where  $A' = (A + B)/2$  is the average equatorial principal moment of inertia of the Earth,  $C_m$  is the axial principal moment of inertia of the crust and mantle,  $A'_m = A' - A_c$  is the equatorial principal moment of inertia of the crust and mantle, and  $k_r$  is a factor, whose value is near unity, that accounts for the effects of rotational deformation on the axial component:

$$k_r = \left\{ 1 + \left[ n_o + \frac{4}{3} (k_2 + \Delta k_{ocn,s}) \right] \frac{a^5 \Omega^2}{3G} \frac{1}{C_m} \right\}^{-1} \quad (24)$$

The deformation of the Earth associated with surficial excitation processes that load the solid Earth has been taken into account in Eqs. (21)–(23) by including the second-degree load Love number  $k'_2$  where  $\Delta k'_{an}$  accounts for the effects of mantle anelasticity on the load Love number and where, because of core decoupling, the load Love number in the axial component is modified by a factor  $\alpha_3$ . Expressions for the excitation functions for processes that do not load the solid Earth can be recovered from Eqs. (21) to (23) by setting the load Love number  $k'_2 + \Delta k'_{an}$  to zero.

Numerically, the real parts of the excitation functions (21)–(23) can be written as (Gross 2007):

$$\chi_x(t) = \frac{1.608 [h_x(t) + 0.684 \Omega \Delta I_{xz}(t)]}{(C - A') \Omega} \quad (25)$$

$$\chi_y(t) = \frac{1.608 [h_y(t) + 0.684 \Omega \Delta I_{yz}(t)]}{(C - A') \Omega} \quad (26)$$

$$\chi_z(t) = \frac{0.997}{C_m \Omega} [h_z(t) + 0.750 \Omega \Delta I_{zz}(t)] \quad (27)$$

These results differ from those of Wahr (1982, 1983, 2005) by about 2%, mostly because of differences in the values of the numerical constants.

The theory of the Earth's rotation described by Eqs. (18)–(23) is a linearized theory in which it has been assumed that: (1) the perturbing excitations are small with  $h_i(t) \ll \Omega C$  and  $\Delta I_{ij}(t) \ll C$ ; (2) the rotational response of the Earth is small with  $m_i(t) \ll 1$ ; (3) the induced relative angular momentum of the core is linearly related to changes in the rotation of the solid Earth; (4) the induced deformations of the mantle, crust, and oceans are linearly related to the changes in rotation; (5) the rotating terrestrial reference frame is the Tisserand mean-mantle frame; (6) the oceans stay in equilibrium as the



rotation changes; (7) the core is uncoupled from the mantle; (8) the crust, mantle, and core are axisymmetric; (9) the rotational variations occur on time scales much longer than a day; (10) the coupling between the components of rotation introduced by a non-uniform ocean are negligibly small and hence can be ignored to first order; and (11) the difference in the oceanic Love number for the two components of polar motion is negligibly small and hence to first order can be replaced by a mean oceanic Love number for the wobble. In addition, Eqs. (14) and (15) for changes in the equatorial components of relative angular momentum caused by changes in rotation are valid only to first order in the ellipticity of the core.

The theory of the Earth's rotation described by Eqs. (18)–(23) is also a hybrid theory in which: (1) the body tide Love number  $k_2$  has been replaced with a wobble-effective Love number  $k_w$  computed from normal mode theory in order to more accurately model the structure of the core and the deformation of the crust and mantle; and (2) the theoretical Chandler wobble frequency has been replaced with its observed value  $\sigma_o$  in order to account for the effects of mantle anelasticity since no adequate theory of these effects currently exists.

### 3 Discussion and Summary

Equations (18)–(23) describe small changes in the rotation of a deformable axisymmetric body that is overlain by non-uniformly distributed equilibrium oceans and that is subject to small perturbing excitation. They are the standard equations currently used to study variations in the Earth's rotation and were developed in the late 1970s and early 1980s by Smith and Dahlen (1981) and Wahr (1982, 1983). When these equations were developed, Earth rotation observations were much less accurate than they are today. For example, in the late 1970s observations of polar motion were accurate to a few mas. Today they are accurate to better than 50 microarcseconds. This great improvement in the accuracy of the observations allows smaller signals to be studied today than could be studied when the theory was developed, such as the ellipticity of the Chandler wobble.

The Chandler wobble of a triaxial body is elliptical, as can be seen in the rigid body case from Eqs. (6) to (7), the solution of which in the absence of excitation is:

$$m_x(t) = m \cos(\sigma_r t + \alpha) \quad (28)$$

$$m_y(t) = \left[ \frac{A(C-A)}{B(C-B)} \right]^{1/2} m \sin(\sigma_r t + \alpha) \quad (29)$$

where  $m$  is the amplitude of the motion along the  $x$ -axis and  $\alpha$  is the phase of the motion. The motion described by Eqs. (28) and (29) is prograde elliptical motion of frequency  $\sigma_r$ . Using the observed values for  $A$ ,  $B$ , and  $C$  of the whole Earth (Groten 2004) and an amplitude  $m$  of 200 mas for the Chandler wobble, these equations predict that the difference between the semimajor and semiminor axes of the Chandler ellipse is 1.34 mas for a rigid Earth. Of course, the real Earth is not rigid. But the theory for the rotation of a deformable Earth, Eqs. (18)–(23), applies to an axisymmetric body, not to a triaxial body. As a result, the theory currently used to study Earth rotation variations predicts that the Chandler wobble is prograde circular, as can be seen by solving Eqs. (18) and (19) in the absence of excitation, or from Eqs. (28) to (29) by setting  $B = A$ . So the current theory of the Earth's rotation cannot be used to study the ellipticity of the Chandler wobble even though Höpfner (2003) may have observed it in modern space-geodetic polar motion observations.

The theory of the Earth's rotation that is currently used is nearly 30 years old and should be improved. It should at least be extended to describe the rotation of a triaxial body with a fluid core. It may be possible to still use the theory of Hough (1895) in this case because it was developed originally to study the motion of a fluid core caused by the rotation of a triaxial mantle. More challenging may be the extension of the theory to non-equilibrium oceans. While it may be adequate to assume that the oceans remain in equilibrium at long periods, tidal observations indicate that the oceans are strongly out of equilibrium at the fortnightly period (e.g., Gross 2009). Accounting for the dynamic nature of the oceans may be the biggest challenge to improving the theory of the Earth's rotation.

As challenging as improving the theory of the Earth's rotation is, there have already been some promising advances. Yoder and Standish (1997) and Van Hoolst and Dehant (2002) included triaxiality in their theories of the rotation of oceanless elastic bodies like Mars. More recently, Chen and Shen (2010) have developed a theory of the Earth's rotation that accounts for the triaxiality of the mantle and core, the anelasticity of the mantle, and dissipation in the oceans. And Bizouard and Zotov (2013) have developed a theory of the Earth's rotation that accounts for the triaxiality of the Earth and includes the effect of asymmetric, but still equilibrium, oceans.

In addition to these efforts, the International Astronomical Union and the International Association of Geodesy have recently established a Joint Working Group on the Theory of the Earth Rotation (Ferrándiz and Gross 2013). The purpose of the Joint Working Group is to promote the development of more accurate theories of the Earth's rotation, not of just polar motion and UT1 but also of nutation and precession.

**Acknowledgements** The work described in this paper was performed at the Jet Propulsion Laboratory, California Institute of Technology, under contract with the National Aeronautics and Space Administration. Support for that work was provided by the Earth Surface and Interior Focus Area of NASA's Science Mission Directorate.

## References

- Bizouard C, Zotov L (2013) Asymmetric effects on Earth's polar motion. *Celest Mech Dyn Astron* 116:195–212. doi:[10.1007/s10569-013-9483-x](https://doi.org/10.1007/s10569-013-9483-x)
- Chen W, Shen W (2010) New estimates of the inertia tensor and rotation of the triaxial nonrigid Earth. *J Geophys Res* 115:B12419. doi:[10.1029/2009JB007094](https://doi.org/10.1029/2009JB007094)
- Dahlen FA (1976) The passive influence of the oceans upon the rotation of the Earth. *Geophys J R Astron Soc* 46:363–406
- Ferrándiz JM, Gross RS (2013) The new joint IAG/IAU working group on theory of the Earth rotation. Presentation given at the 2013 IAG scientific assembly, Potsdam, 1–6 Sept 2013
- Goldstein H (1950) *Classical mechanics*. Addison-Wesley, Reading
- Gross RS (2007) Earth rotation variations — long period. In: Herring TA (ed) *Physical geodesy: treatise on geophysics*, vol 3. Elsevier, Oxford, pp 239–294
- Gross RS (2009) Ocean tidal effects on Earth rotation. *J Geodyn* 48:219–225. doi:[10.1016/j.jog.2009.09.016](https://doi.org/10.1016/j.jog.2009.09.016)
- Groten E (2004) Fundamental parameters and current (2004) best estimates of the parameters of common relevance to astronomy, geodesy, and geodynamics. *J Geod* 77:724–731
- Höpfner J (2003) Chandler and annual wobbles based on space-geodetic measurements. *J Geodyn* 36:369–381
- Hough SS (1895) The oscillations of a rotating ellipsoidal shell containing fluid. *Philos Trans R Soc Lond A* 186:469–506
- Smith ML, Dahlen FA (1981) The period and Q of the Chandler wobble. *Geophys J R Astron Soc* 64:223–281
- Tisserand F (1891) *Traité de Mécanique Céleste*, vol II. Gauthier-Villars, Paris (in French)
- Van Hoolst T, Dehant V (2002) Influence of triaxiality and second-order terms in flattenings on the rotation of terrestrial planets I. Formalism and rotational normal modes. *Phys Earth Planet Inter* 134:17–33
- Wahr JM (1982) The effects of the atmosphere and oceans on the Earth's wobble — I. Theory. *Geophys J R Astron Soc* 70:349–372
- Wahr JM (1983) The effects of the atmosphere and oceans on the earth's wobble and on the seasonal variations in the length of day — II. Results. *Geophys J R Astron Soc* 74:451–487
- Wahr J (2005) Polar motion models: angular momentum approach. In: Plag H-P, Chao BF, Gross RS, van Dam T (eds) *Forcing of polar motion in the chandler frequency band: a contribution to understanding interannual climate change*, vol 24. Cahiers du Centre Européen de Géodynamique et de Séismologie, Luxembourg, pp 89–102
- Yoder CF, Standish EM (1997) Martian precession and rotation from Viking lander range data. *J Geophys Res* 102(E2):4065–4080

---

# Variable Seasonal and Subseasonal Oscillations in Sea Level Anomaly Data and Their Impact on Prediction Accuracy

W. Kosek, T. Niedzielski, W. Popiński, M. Zbylut-Górska, and A. Wnęk

---

## Abstract

Weekly sea level anomaly (SLA) maps are now available courtesy of the Archiving, Validation and Interpretation of Satellite Oceanographic (AVISO) data. Using the Fourier Transform Band Pass Filter (FTBPF) variable broadband seasonal and subseasonal oscillations were computed as a function of geographic location. Irregular amplitude and phase variations in these oscillations cause the increase of prediction errors of the SLA data for a few weeks in the future. The amplitude and phase variations of the broadband annual oscillation were computed by a combination of the FTBPF and the Hilbert transform. In order to detect the impact of irregular amplitude or/and phase variations of the annual oscillation on the SLA prediction errors, standard deviations maps of amplitude time differences as well as of the products of phase time differences and amplitudes were examined. The SLA data prediction errors in certain geographic regions of the ocean seem to be caused mainly by nonlinear behaviour of the broadband annual oscillation. The nonlinearities are probably driven by mesoscale eddies, and the significant impact on SLA prediction errors was observed in the vicinity of the western boundary currents.

---

## Keywords

Fourier transform band pass filter • Prediction • Satellite altimetry • Sea level change

---

## 1 Introduction

Sea level varies across multiple spatial and temporal scales, and its changes are driven by eustatic and steric processes. A radar-based satellite altimetry technique offers absolute

observations of sea surface height (SSH) with respect to the Earth's centre of mass. The complexity of physical processes which influence these variations is considerable, and hence sea level modelling and prediction is still an ongoing challenge. Numerous researchers have scrutinised the problem of seeking models suitable for sea level prediction (Röske 1997; Gregory and Lowe 2000; Rahmstorf 2007; Niedzielski and Kosek 2009). In addition, there are dedicated systems designed for calculating and publishing the sea level prognoses on maps. Noteworthy four initiative have to be mentioned: the Ocean Prediction Center of NOAA, the HYCOM Consortium, MyOcean and Prognocean. The increased inaccuracy of sea level predictions produced by the latter system has recently been found to follow the locations where highly nonlinear ocean processes act (Chelton et al. 2011), and hence mesoscale eddies have been assumed to control departures from accurate prognoses (Niedzielski and Miziński 2013). Although the aforementioned correspondence has

---

W. Kosek • M. Zbylut-Górska • A. Wnęk  
Environmental Engineering and Land Surveying Department,  
Agriculture University of Krakow, Krakow, Poland

W. Kosek (✉) • T. Niedzielski  
Space Research Centre, Polish Academy of Sciences, Warsaw, Poland  
e-mail: [kosek@cbk.waw.pl](mailto:kosek@cbk.waw.pl)

T. Niedzielski  
Institute of Geography and Regional Development, University  
of Wrocław, Wrocław, Poland

W. Popiński  
Central Statistical Office of Poland, Warsaw, Poland

been identified a detailed diagnosis of the prediction error is needed. The objective of the work discussed in this paper is to verify the hypothesis that the reason behind inaccurate prognoses of sea level change derived in real time using a few empirical time series methods, within the system known as Prognosean (Niedzielski and Miziński 2013), is related to certain characteristics of selected annual broadband oscillation.

## 2 Prediction of Sea Level Anomaly Data

The gridded sea level anomaly (SLA) data for geographic latitudes  $\phi \in [-90; 90]$  and longitudes  $\lambda \in [0; 360]$  from AVISO (Archiving, Validation and Interpretation of Satellite Oceanographic Data) were used spatial resolution of  $1^\circ \times 1^\circ$ . The sampling interval of these data is equal to 1 week and their time span is 1992–2013. These data are computed from observations of altimetric and remote sensing satellites such as TOPEX/Poseidon, ERS 1 and 2, Jason 1 and 2, Cryosat-2 and Envisat. The aim of Prognosean is to compute real-time predictions of SLA data for lead times ranging from 1 to 14 days, making use of several methods constructed as combinations of the polynomial-harmonic (PH) model with several stochastic forecast methods, such as for instance autoregressive (AR) and threshold autoregressive (TAR) techniques. For each prediction method the maps of the root mean square (RMS) prediction errors are generated (Niedzielski and Miziński 2013). High RMS values of the SLA prediction for PH + AR and PH + TAR combinations are similar for both methods and are observed in the vicinity of the western boundary currents, independently of the prediction method applied (Fig. 1).

## 3 Analysis of Sea Level Anomaly DATA

To analyse the SLA data the Fourier transform band pass filter (FTBPF) was used (Kosek 1995; Popiński 2008), and

hence the following formula was applied:

$$u_{\phi,\lambda}(t, \omega) = FT^{-1} [FT [x_{\phi,\lambda}(t)] P(\omega, \mu)], \quad (1)$$

where  $FT$  is the Fourier transform operator,  $x_{\phi,\lambda}(t)$  is the SLA time series,  $u_{\phi,\lambda}(t, \omega)$  is the broadband oscillation with the central frequency  $\omega$ ,  $P(\omega, \mu) = \begin{cases} 1 - ((\omega - \mu)/\Lambda)^2 & \text{if } |\omega - \mu| \leq \Lambda \\ 0 & \text{if } |\omega - \mu| > \Lambda \end{cases}$  is the parabolic transmittance function,  $\mu$  is the frequency argument and  $\Lambda$  is half of the frequency bandwidth.

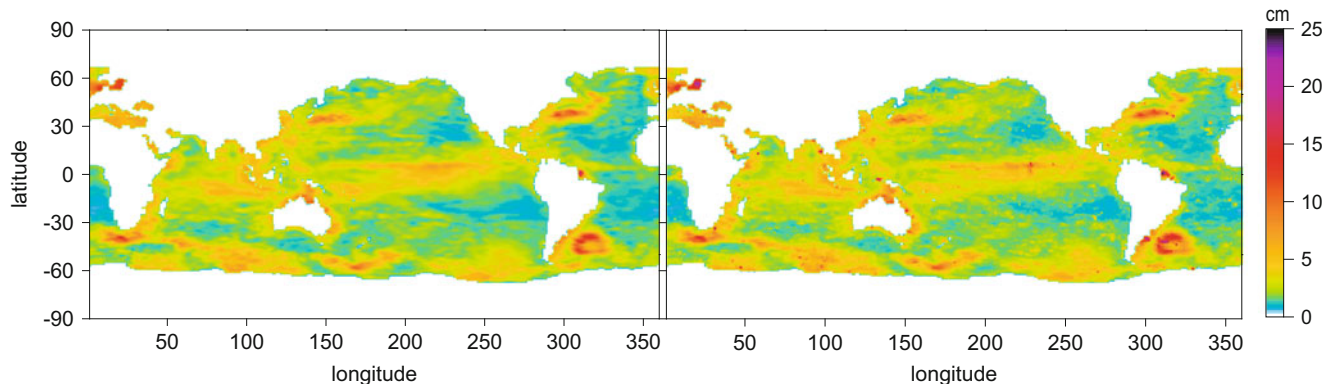
The mean amplitude spectrum as a function of geographic location and oscillation period  $T = \Delta t/\omega$  is computed as:

$$\hat{S}_{\phi,\lambda}(T) = \sqrt{\frac{2}{n-2k} \sum_{t=k+1}^{n-k} |u_{\phi,\lambda}(t, T)|^2}, \quad (2)$$

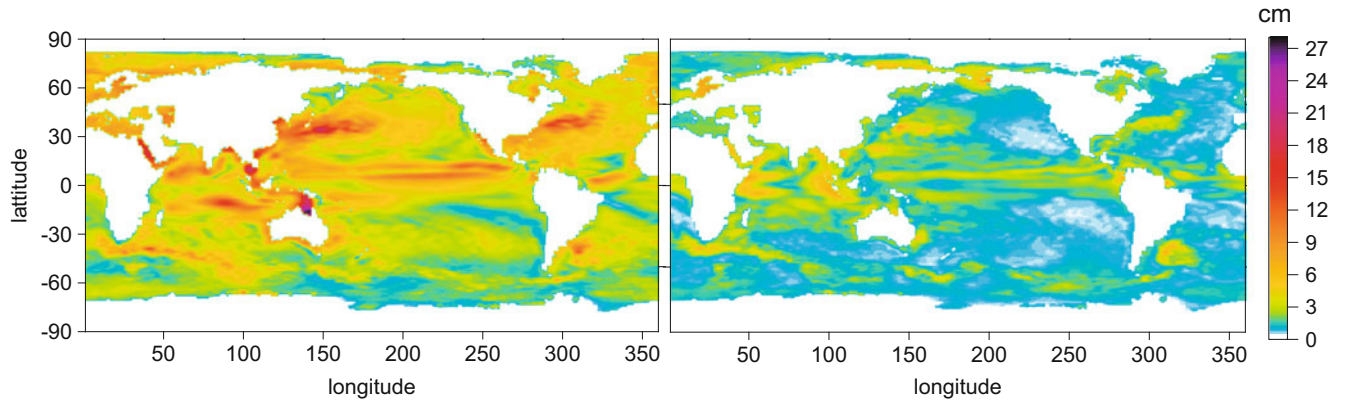
where  $\Delta t = 7$  days is the data sampling interval,  $n$  is the number of data in SLA time series,  $k = 20$  is the number of points to be dropped at the beginning and at the end of the filtered oscillations time series due to filter errors.

The annual oscillation mean amplitude computed by Eq. (2) is the greatest and reaches 25 cm at Arafura Sea, in the areas of Kuroshio, Gulf Stream and Antarctic Circumpolar currents, Red Sea, Thailand Bay and equatorial regions (Fig. 2). The amplitudes of the semi-annual oscillation are high in the areas where amplitudes of the annual oscillation attain large values. The highest values of the order of 6–9 cm are observed in the Baltic Sea, West and East Indian Ocean and Arctic regions.

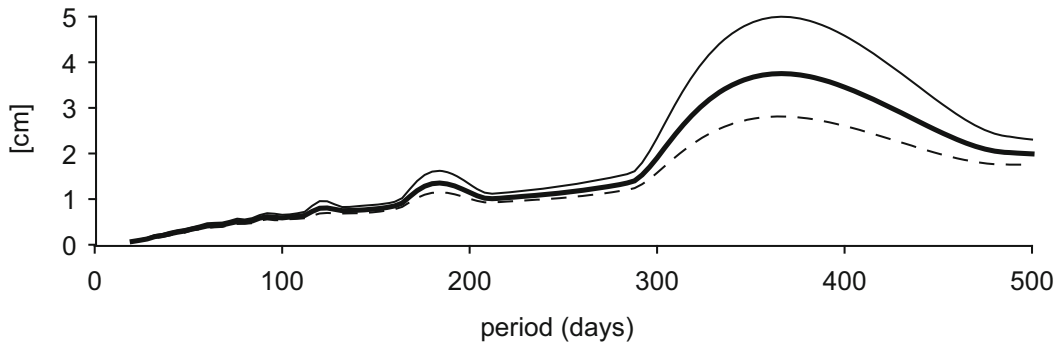
The mean amplitude of the annual oscillation for the entire ocean  $\phi \in [-90; 90]$ , the Northern  $\phi \in [0; 90]$  and Southern  $\phi \in [-90; 0]$  Hemispheres is computed by  $\hat{S}(T) = \sum_{\phi} \sum_{\lambda=1}^{360} \hat{S}_{\phi,\lambda}(T)$ . The mean amplitude of the annual oscillation in the Northern Hemisphere is twice as big as in the Southern Hemisphere (Fig. 3), as previously inferred (Kosek 2001). The peaks in these amplitude spectra correspond to



**Fig. 1** The RMS prediction error of SLA data for 14 days in the future computed with PH + AR (left) and PH + TAR (right) combinations



**Fig. 2** The mean amplitudes of the annual (*left*) and semi-annual (*right*) oscillations computed by the FTBPF for half of the bandwidth  $\Lambda = 0.01$



**Fig. 3** The mean amplitude spectrum of the entire ocean (*heavy line*), the Northern Hemisphere (*thin line*) and the Southern Hemisphere (*dashed line*), computed by the FTBPF ( $\Lambda = 0.01$ )

the integer multiplicities of the annual frequency e.g. the semi-annual, 120-day and quarter-annual oscillations, which suggests that the annual oscillation is broadband.

#### 4 Amplitude and Phase Variations of the Annual Oscillation

The broadband character of the annual oscillation suggests that its amplitude and phase are variable. To compute variations of its amplitude and phase a combination of the FTBPF and the Hilbert transform (HT) was applied. In this method the HT of the broadband oscillation  $u_{\phi,\lambda}(t, \omega)$  is used as the imaginary part to create the complex-valued time series:

$$z_{\phi,\lambda}(t, \omega) = u_{\phi,\lambda}(t, \omega) + i \cdot HT[u_{\phi,\lambda}(t, \omega)]. \quad (3)$$

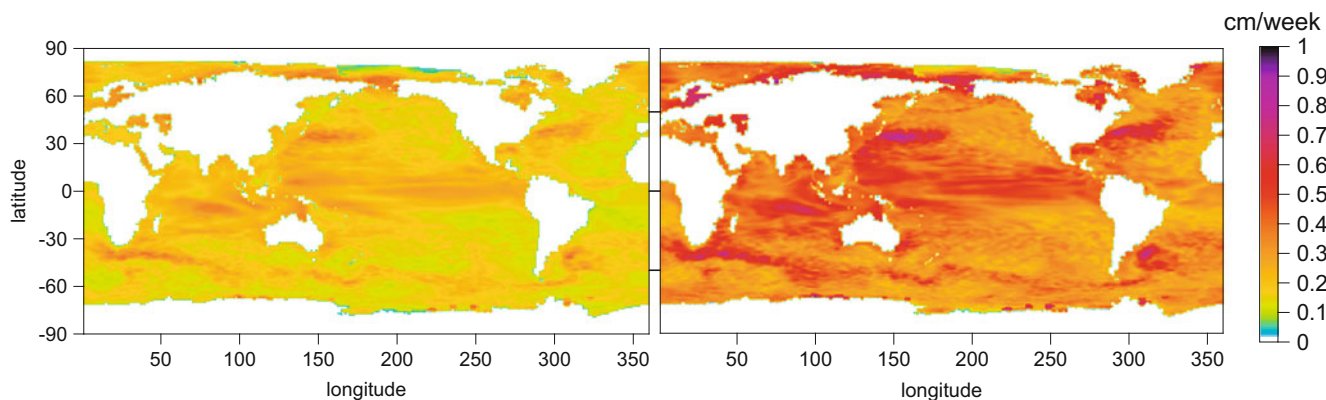
It can be proved that (Gasquet and Witomski 1999; Popiński 2008):

$$z_{\phi,\lambda}(t, \omega) = FT^{-1} [FT(x_{\phi,\lambda}(t)) \cdot P(\omega, \mu) \cdot (\text{sign}(\mu) + 1)]. \quad (4)$$

To examine their influence on prediction errors of the SLA data, variable amplitudes  $A_{\phi,\lambda}(t, \omega)$  and phases  $\phi_{\phi,\lambda}(t, \omega)$  of  $z_{\phi,\lambda}(t, \omega)$  were computed.

To detect the variable amplitudes and phases the Morlet wavelet transform (MWT) can be applied too (Kosek et al. 2006). Subsequently, the time series of the amplitude differences  $\Delta A_{\phi,\lambda}(t, \omega)$ , as well as products of amplitudes and phase differences  $A_{\phi,\lambda}(t, \omega) \cdot \Delta \phi_{\phi,\lambda}(t, \omega)$  were computed. Both, amplitude and phase differences were computed as differences between the relevant values on succeeding time moments. After subtracting the mean values of these difference based time series their standard deviations were computed (Fig. 4).

If amplitudes of oscillations are small then the phase changes do not influence oscillation variability significantly. Thus, phase differences were multiplied by amplitudes in order to estimate their impact on the irregular character of an oscillation. It can be noticed that the products of amplitudes and phase time differences have bigger standard deviations than the amplitude time differences. Therefore, phase variations of the annual oscillation have greater impact on irregular character of the annual oscillation than the



**Fig. 4** The standard deviations of the amplitude time differences (*left*) and the product of the amplitude and phase time differences (*right*) of the annual oscillation

amplitude variations. The geographic regions of high significant amplitude variations are similar to geographic regions of the phase variations of the annual oscillation.

## 5 Conclusions

The FTBPF analysis of the SLA data reveals that the annual oscillation has a broadband character. It comprises oscillations with higher frequencies being integer multiplicities of the annual frequency. The amplitude maxima of the semi-annual oscillation are located in geographic regions where the amplitude maxima of the annual oscillation occur. The mean prediction errors of the SLA data for 14 days in the future are usually considerable in geographic regions where amplitude maxima of the annual oscillation are the largest. Geographic regions of the significant phase and amplitude variations of the annual oscillation correspond to high RMS prediction errors of the SLA data. The increase of the prediction errors of sea level anomaly data is mostly caused by variable phases and amplitudes of the broadband annual oscillation. It can be argued that the irregular phase variations of the annual oscillation are the main causes of the increase of the SLA prediction errors.

**Acknowledgments** This work was supported by the Polish national science foundation NCN under grant No. 2012/05/B/ST10/02132. The research was also supported through the Homing Plus grant of the Foundation for Polish Science, contract no. Homing Plus\_2011-3/8

under leadership of Tomasz Niedzielski, founded by the European Regional Development Fund and the Innovative Economy Programme. We thank AVISO+ for providing the altimeter-derived sea level anomaly data.

## References

- Chelton DB, Schlax MG, Samelson RM (2011) Global observations of nonlinear mesoscale eddies. *Prog Oceanogr* 91:167–216
- Gasquet C, Witomski P (1999) Fourier analysis and applications – filtering, numerical computation, wavelets. Springer Verlag, New York
- Gregory JM, Lowe JA (2000) Predictions of global and regional sea-level rise using AOGCMs with and without flux adjustment. *Geophys Res Lett* 27:3069–3072
- Kosek W (1995) Time variable band pass filter spectra of real and complex-valued polar motion series. *Artif Satellites* 30:27–43
- Kosek W, Rzeszółtko A, Popiński W (2006) Phase variations of oscillations in the Earth orientation parameters detected by the wavelet technique. *Proc Journées* 2005:121–124
- Kosek W (2001) Long-term and short period global sea level changes from TOPEX/Poseidon altimetry. *Artif Satellites* 36(3):71–84
- Niedzielski T, Kosek W (2009) Forecasting sea level anomalies from TOPEX/Poseidon and Jason-1 satellite altimetry. *J Geod* 83:469–476
- Niedzielski T, Miziński B (2013) Automated system for near-real time modelling and prediction of altimeter-derived sea level anomalies. *Comput Geosci* 58:29–39
- Popiński W (2008) Insight into the fourier transform band pass filtering technique. *Artif Satellites* 43:129–141
- Rahmstorf S (2007) A semi-empirical approach to projecting future sea-level rise. *Science* 315:368–370
- Röske F (1997) Sea level forecasts using neural networks. *Ocean Dyn* 49:71–99

---

# Permanent GPS Networks in Italy: Analysis of Time Series Noise

R. Devoti, G. Pietrantonio, A.R. Pisani, and F. Riguzzi

---

## Abstract

Over the last few years numerous GPS networks in Italy have been installed and managed, mainly by local authorities and institutions. Therefore the GPS stations have been constructed with a variety of different monument types according to their needs and have been operated in fairly different environmental conditions, such as in towns or industrial regions, in the open country or mountainous regions. In this work we aim to assess the reliability and repeatability of the station positions and to study the noise property of different categories of GPS monument types. We analyze over 500 continuous GPS time series in Italy with a mean temporal length of 5.6 years. All the GPS observations were processed with the Bernese v5.0 software using a loose constraints approach. We include 45 sites in central Europe that are used as fiducial stations in the regional reference frame realization. After fitting a linear drift, offsets and annual sinusoids and after filtering a common mode movement of the whole network, the residual GPS time series represents the noise of each GPS station. We analyze the residuals using different power spectrum estimation schemes and estimate a power law noise model for each time series. The average noise characteristics are compatible with outcomes from earlier studies but we were not able to isolate distinct noise behaviors between different GPS monument types nor to ascertain a preferred monumentation, as far as noise amplitude and spectral indexes are concerned.

---

## Keywords

GPS time series • Italy • Power law noise • Spectral index

---

## 1 Introduction

The first attempt to build a nation-wide continuous GPS network was undertaken by the Italian Space Agency (ASI) in the late 1990s. Since then, ASI (<ftp://geodaf.mt.asi.it>) delivers continuous GPS data from about 30 sites and main-

tains the regional reference frame in cooperation with the European reference frame consortium (EUREF). In 2001, the Istituto Nazionale di Oceanografia e Geofisica Sperimentale (OGS) started installing a local GPS network in the Friuli region (Northeast Italy, <http://www.crs.inogs.it/frednet>) and in 2004, the Istituto Nazionale di Geofisica e Vulcanologia (INGV) started the construction of the first national GPS network (RING, <http://ring.gm.ingv.it/>). At present, the RING network consists of about 170 stations, most of them are built on short or deep-drilled braced tripods (Avallone et al. 2010). More recently an increasing number of permanent GPS sites have been installed by regional administrations and private companies, dedicated mainly to topographic applications and commercial services. These networks, although not conceived to measure long term ground deformations, have

---

R. Devoti (✉) • G. Pietrantonio • F. Riguzzi  
Istituto Nazionale di Geofisica e Vulcanologia (INGV), Centro  
Nazionale Terremoti, Via di Vigna Murata 605, 00143 Roma, Italy  
e-mail: [roberto.devoti@ingv.it](mailto:roberto.devoti@ingv.it)

A.R. Pisani  
Agenzia Spaziale Italiana (ASI), Osservazione della Terra, Via del  
Politecnico, 00133 Roma, Italy

proved to be useful in augmenting the backbone of more reliable geodynamic networks and are currently distributing their data, making it available for the scientific community. All these datasets are currently archived and processed at INGV providing over 700 RINEX files per day for a mean geometric inter-distance of about 20 km over the whole country, thus realizing an important dataset for geodynamical studies. In this perspective, the assessment of the performance of a GPS station is very important, especially in the Italian area where the crustal deformations are slow (at the few mm/year level) but gradients (strain-rates) are changing rapidly from point to point.

## 2 GPS Data and Methods

In this paper we restrict our analysis to all the permanent GPS stations located in Italy. At present over 700 stations are active and we selected 563 stations spanning a life time between 2.5 years (912 days) and 14.7 years (5,382 days). The mean length of the analyzed time series is 5.6 years (2,046 days). All the data have been processed using the *Bernese Processing Engine* (BPE) ver. 5.0 (Beutler et al. 2007). The processing strategy follows the EUREF Guidelines for EPN Analysis Centres (<http://www.epncb.oma.be/documentation/guidelines/>). The GPS orbits and the Earth's orientation parameters have been fixed to the combined IGS products and an a priori loose constraint of 10 m has been assigned to all site coordinates. The elevation-dependent phase centre corrections and absolute phase centre calibrations have been applied to the processing. The troposphere modeling consists in an a priori dry-Niell model fulfilled by the estimation of zenith delay corrections at 1-h intervals at each site using the wet-Niell mapping function; in addition one horizontal gradient parameter per day at each site is estimated. The ionosphere is not modeled a priori, it is removed by applying the ionosphere-free linear combination of L1 and L2. The ambiguity resolution is based on the QIF baseline-wise analysis. The final network solution is solved with back-substituted ambiguities, if integer; otherwise ambiguities are considered as real valued measurement biases.

The daily GPS solutions are not estimated in a given reference frame but computed in a loosely constrained reference frame. Therefore, the coordinates are randomly translated or rotated from day-to-day and their covariance matrices have large errors (on the order of meters) as a consequence of the loose constraints applied to the a priori parameters. To express the coordinate time series in a unique reference frame and to compute the real covariance matrix, we perform two main transformations. First the loose covariance matrix is projected into a well-defined reference frame imposing tight internal constraints (at the mm level), and then coor-

ordinates are transformed into the ITRF2008 by a 4-parameter Helmert transformation (translations plus scale factor); the proper set of constraints is driven by the rank deficiency of the normal matrices, a comprehensive discussion of the rank deficiency of our solutions is given in Devoti et al. (2010). The Helmert transformation uses 45 sites located in central Europe as anchor stations for the regional reference frame realization.

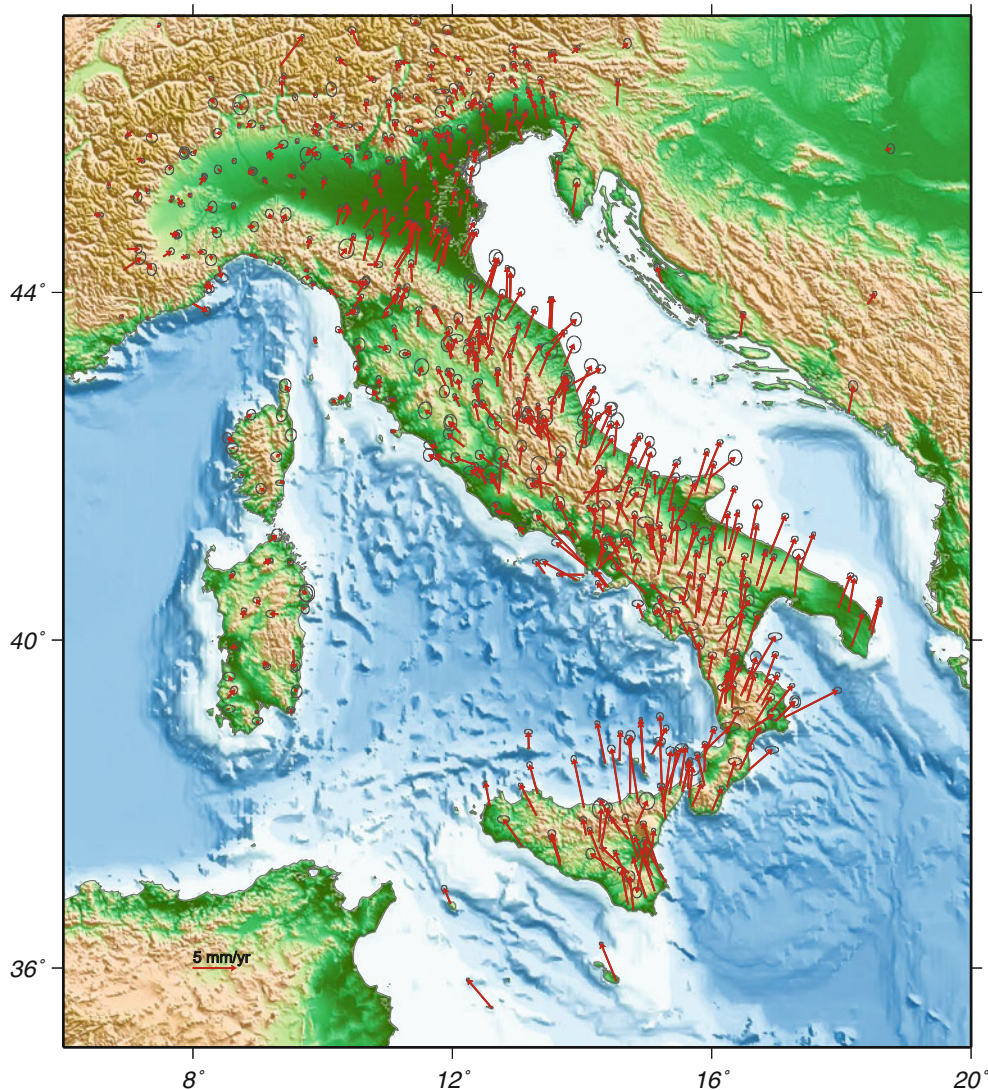
Velocities are estimated fitting simultaneously to all station coordinates a linear drift, episodic offsets and annual sinusoids. Figure 1 shows the estimated velocity field of the regional GPS network. Offsets are estimated whenever a change in the GPS equipment induces a significant step in the time series whereas seasonal oscillations are accounted for by annual sinusoids. At this stage outliers have been rejected whenever the weighted residual exceeds three times the global chi square ( $\chi^2$ ). Finally the common mode error signal has been filtered out with a procedure similar to that adopted by Wdowinski et al. (1997), we compute daily weighted mean of the residuals of a local fiducial network (57 stations) selected among the best performing stations. After the common mode filtering the time series scatter decreases by 25–35%, the final median weighted-root-mean-squared of the residuals in the Vertical, East and North components are 4.3 mm, 1.4 mm and 1.4 mm respectively.

## 3 Noise Model

GPS site monuments are very heterogeneous in type since there are many networks at regional and national scale managed by different owners. We recognize three main types (Fig. 2): shallow drilled braced tripods, formed usually by four stainless steel rods arranged in a pyramid structure and anchored to bedrock (a); a less standardized monumentation fixed with steel structures on the roof of buildings (b) and different types of concrete pillars with, usually, 1–2 m foundations (c).

One source of noise in geodetic measurements is thought to develop from random motions occurring at the connection of the geodetic instrument to the ground (Wyatt 1982, 1989). It is thought that this connection (i.e. the monument) follows an approximately random walk process, i.e. a process in which the monument position is expected to deviate from an initial position as the square-root of time (Langbein and Johnson 1997; Zhang et al. 1997). Several other papers have demonstrated that daily GPS positions are temporally correlated and not simply independent measurements using both short-baseline, regional or globally distributed GPS networks (e.g. Mao et al. 1999; Johnson and Agnew 2000; Langbein 2004; Williams et al. 2004; Beavan 2005; Langbein 2008; King and Williams 2009; King and Watson 2010). An overview of studies dealing with time-correlated noise in





**Fig. 1** Map of the estimated velocities of 563 continuous GPS stations in Italy expressed in the Eurasian reference frame. Each GPS site spans different life times, we select all those with a minimum of 2.5 years

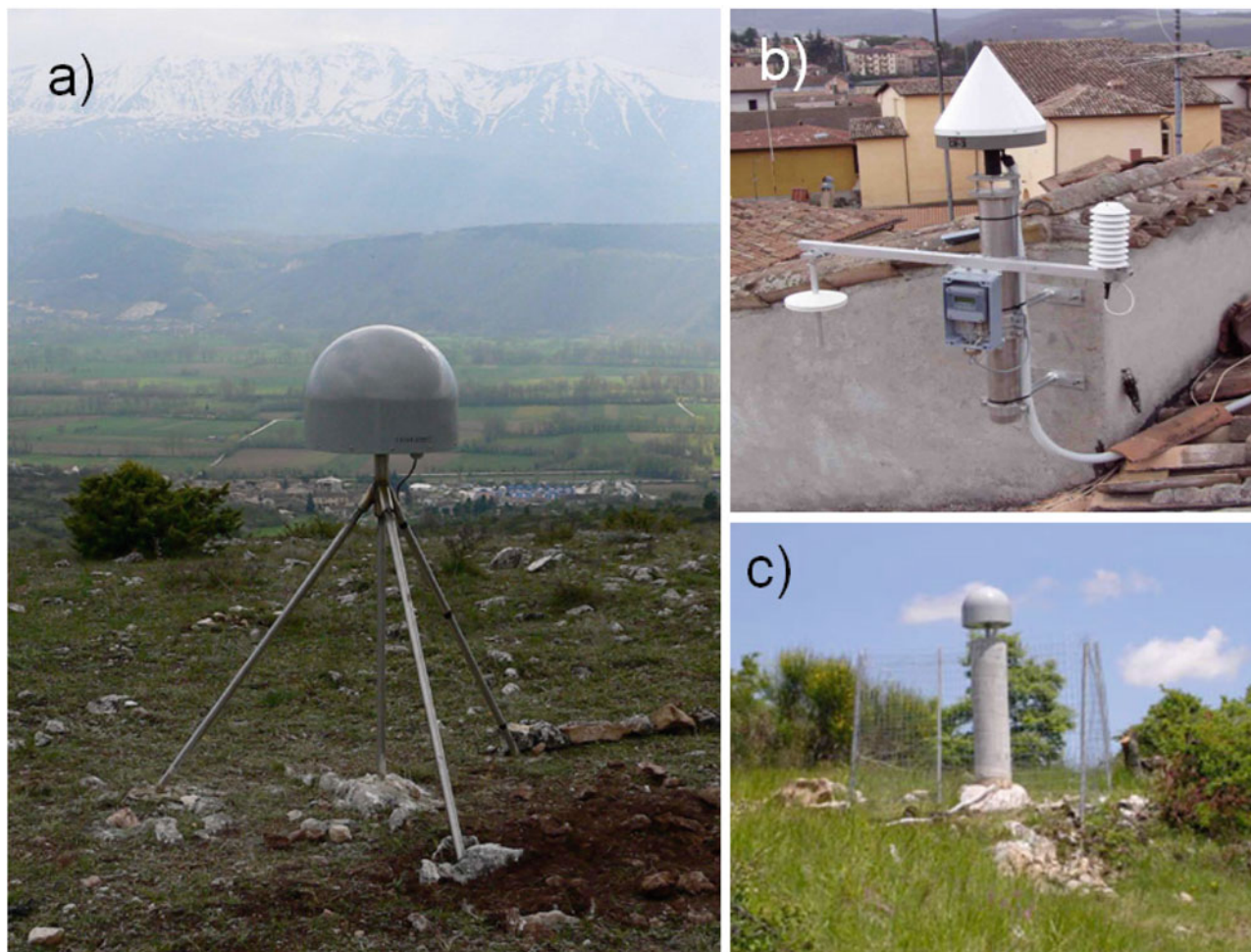
of continuous observation time, the general average time span being 5.6 years. The error ellipses represent the 95% confidence region

this field is given in Santamaria-Gómez et al. (2011). Recent analysis show that GPS time series include complex noise processes that are less correlated than simple random walk noise or in which the random walk process exhibits lower amplitudes.

Noise processes are conveniently described in the frequency domain, in a log–log plot, the noise spectrum shows a general negative slope as frequency increases. For random walk noise the slope is  $-2$ , for an intermediate noise called flicker noise, the slope is  $-1$ , whereas for uncorrelated white noise the slope is  $0$ . In order to study the noise features of our GPS time series we assume a generic power law model with amplitude  $C_0$  and spectral index  $K$  superimposed on a white noise component of amplitude  $WN$ . Thus the power spectrum can be modelled as:

$$P(f) = C_0 f^K + WN \quad (1)$$

where the unknowns are  $C_0$ ,  $K$  and  $WN$ . Using the above model we estimate the noise parameters from the power spectrum of each coordinate component (Vertical, East and North directions) using a nonlinear optimization algorithm based on the simplex search method (Lagarias et al. 1998). We test different methods to compute the power spectral density (PSD), each with positive aspects and known drawbacks. Here we restrict the analysis to the direct fast Fourier transform (FFT) and two particular algorithms that provide estimates of the power spectrum in different ways: the Lomb–Scargle (LS) periodogram and the Welch (W) power density spectrum. Data gaps in GPS time series are not interpolated nor filled with any a priori noise model,



**Fig. 2** Three typical classes of GPS monument types: shallow drilled braced tripod anchored on bedrock (a); steel structure fixed on the roof of buildings (b) and concrete pillar with usually 1–2 m foundations (c)

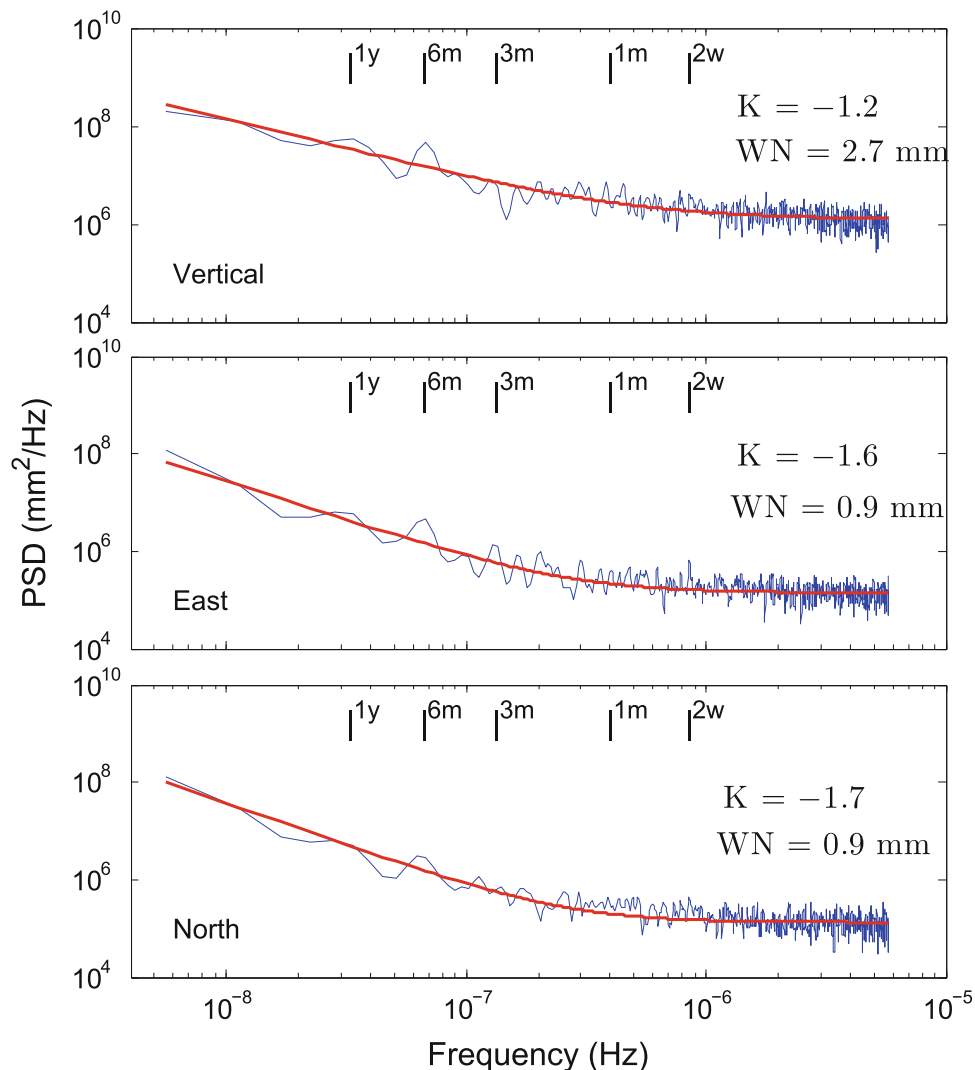
thus we prefer not to induce any artificial distortion to the data in order to minimize the effect of arbitrary a priori assumptions. The LS algorithm (Lomb 1976; Scargle 1982) was developed to deal with unevenly sampled data estimating the power spectrum using least-squares. It provides an approximation to the real spectrum and reduces to the Fourier power spectrum in the limit of equal spacing. Unfortunately the method provides a relative noisy power that does not diminish in amplitude as the sample population increase, see Scargle (1982) for a detailed discussion. Finally we test a more complex algorithm that tends to decrease the variance of the single periodogram (Welch 1967). It consists of dividing the time series data into overlapping segments, computing a modified power spectrum of each segment, and then averaging the power spectral density estimates. We apply a Hamming window to each segment and allow a 50% overlap between segments. In order to obtain homogeneous results we consider only spectral frequencies greater than  $1/(5.6 \text{ years})$  and to avoid unwanted aliasing caused by large

outliers we apply a median based filter rejecting all residuals exceeding 3.5 times the median.

## 4 Results and Discussion

Figure 3 shows the PSD of the common mode noise of the regional solution of Italy obtained with the Welch method. The spectra of the Vertical, East and North components of the network center are shown in the panels, the red line represents the fitted noise model (1). The common mode analysis covers the entire time span (15.5 years) of the GPS data, the white noise amplitude (WN) in the vertical direction is about three times larger than the horizontal. The spectral index ranges from  $-1.2$  to  $-1.7$ , being slightly lower in the horizontal rather than in the vertical components. Due to the longer time span we are probably able to detect lower frequency components of the common mode noise, i.e. shifted towards a more pronounced random walk behavior,

**Fig. 3** Power spectrum of the estimated common mode time series. The *vertical component* is shown in the *upper panel* and the *horizontal components* in the *lower two*. The *red line* represents the fitted noise model (power law + white noise). The estimated spectral index ( $K$ ) and the white noise amplitude (WN) are shown in the *insets*. Frequency tags in correspondence of 1 year, 1, 3, 6 month and 2 week periods, are also indicated in the plot



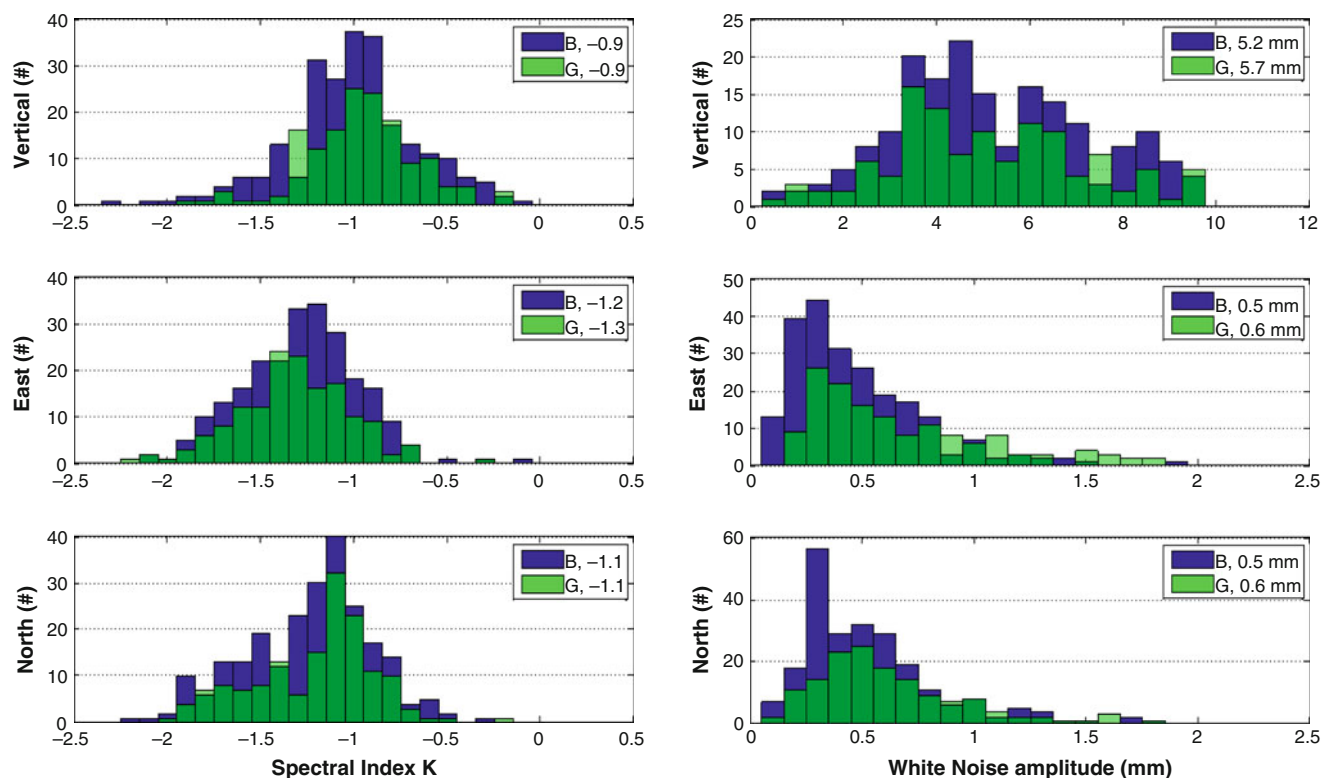
compared with the noise properties of the individual GPS time series that, on the average, have one-third of the observation time span. The common mode noise doesn't show any pronounced periodic signal in the three components, apart from a weak seasonal signature.

The GPS time series are sampled at daily frequency but the final residuals may look rather sparse and irregularly sampled after the whole analysis process. Therefore we compare the noise model parameters (1) obtained by fitting three different power spectrum estimates as described in the previous paragraph. The three classes of GPS monument types show rather similar noise behavior, Table 1 shows the median of spectral indexes ( $K$ ) and white noise amplitudes (WN) obtained from two well populated families: the roof based (B) and the ground based (G) monuments, two broad classes that aggregate respectively 238 and 151 stations. The spectral index of the vertical residuals is only slightly less than the horizontal  $K$  in all the spectral approaches, the median being around  $-1.0$ , whereas the horizontal noise

is more correlated showing a median between  $-1.1$  and  $-1.3$ . The subscript and superscript numbers in Table 1 represent the 25 and 75 percentile of the distribution. It is worth noting that the FFT and Lomb–Scargle spectra provide more scattered parameters since they do not smooth the power in any way, and furthermore the Lomb–Scargle is intrinsically a noisy estimation (see the relevant discussion in Scargle 1982). For this reason we adopt the Welch PSD estimation scheme to infer the noise parameters. Figures 4 and 5 show the histograms of estimated spectral indexes and white noise amplitudes for two relevant end members of monument types, namely the ground-based (G) against the roof-based (B) antennas (Fig. 4), and the ground-based tripods (G-tripod) against the roof-based steel antenna mast (Fig. 5). The white noise amplitudes are generally low and reflect the intrinsic uncorrelated noise of the analysis, it is on the order of the best performing SOPAC regional solutions time series (Williams et al. 2004). Both Figs. 4 and 5 provide very similar results in terms of spectral indexes,

**Table 1** Summary table of the noise amplitudes obtained with three different spectral analysis algorithms: FFT, Lomb-Scargle and Welch. On rows, K are the medians of the spectral index amplitudes, whereas WN are the medians of the white noise amplitudes, respectively of

	Vertical		North		East		
	B	G	B	G	B	G	
K	$-1.0_{-1.4}^{-0.8}$	$-1.0_{-1.2}^{-0.7}$	$-1.3_{-1.6}^{-1.0}$	$-1.1_{-1.6}^{-0.9}$	$-1.2_{-1.6}^{-0.9}$	$-1.3_{-1.7}^{-1.0}$	FFT
WN	$4.3_{3.1}^{6.4}$	$4.3_{3.1}^{6.6}$	$0.4_{0.3}^{0.6}$	$0.5_{0.4}^{0.7}$	$0.4_{0.3}^{0.6}$	$0.5_{0.4}^{0.7}$	
K	$-1.0_{-1.3}^{-0.7}$	$-0.9_{-1.2}^{-0.7}$	$-1.2_{-1.6}^{-0.9}$	$-1.2_{-1.6}^{-0.9}$	$-1.2_{-1.6}^{-0.9}$	$-1.2_{-1.6}^{-0.9}$	Lomb Scargle
WN	$4.5_{2.8}^{6.7}$	$4.9_{3.1}^{6.8}$	$0.4_{0.3}^{0.6}$	$0.5_{0.4}^{0.7}$	$0.4_{0.3}^{0.6}$	$0.5_{0.4}^{0.7}$	
K	$-0.9_{-1.2}^{-0.8}$	$-0.9_{-1.1}^{-0.7}$	$-1.1_{-1.4}^{-1.0}$	$-1.1_{-1.4}^{-1.0}$	$-1.2_{-1.5}^{-1.0}$	$-1.3_{-1.5}^{-1.1}$	Welch
WN	$5.2_{3.7}^{8.0}$	$5.7_{3.9}^{8.8}$	$0.5_{0.4}^{0.7}$	$0.6_{0.4}^{0.8}$	$0.5_{0.4}^{0.7}$	$0.6_{0.4}^{0.8}$	



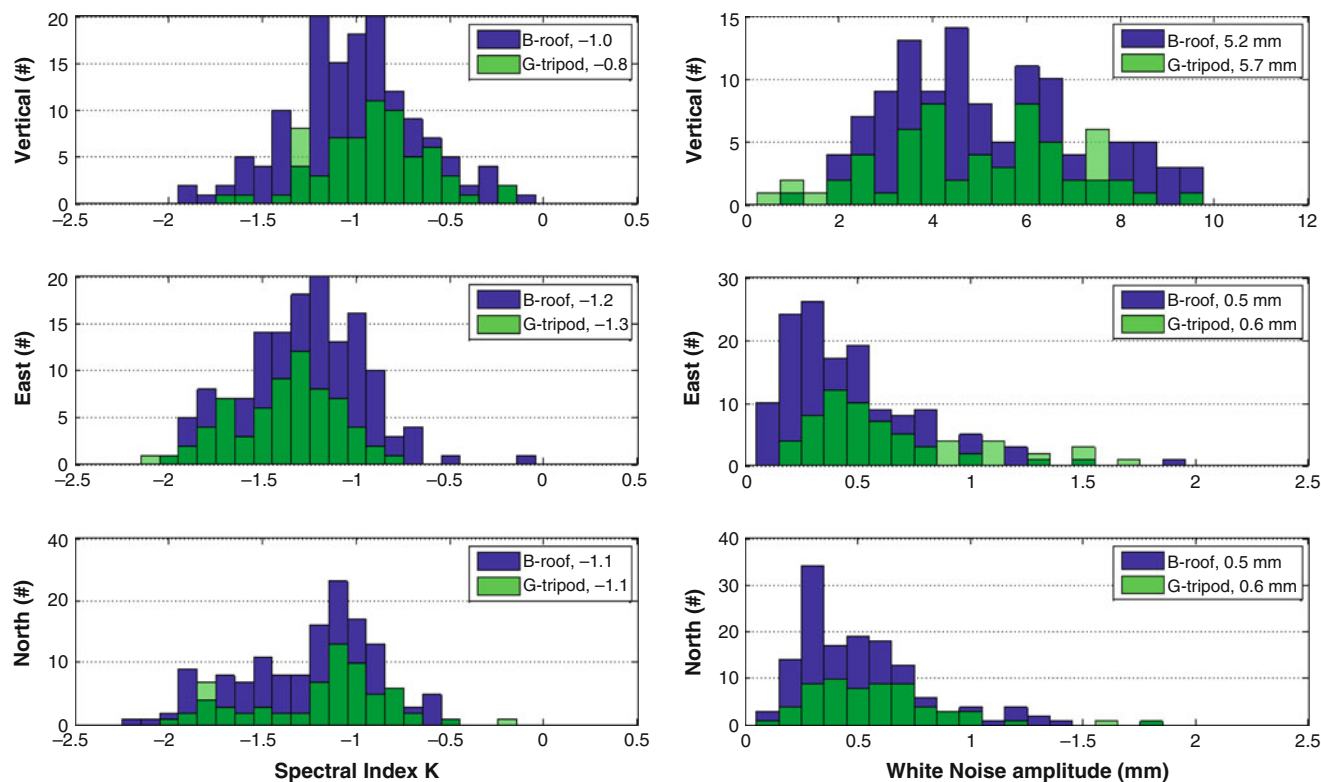
**Fig. 4** Comparison between the roof (B, blue bars) and the ground anchored (G, green bars) monuments, the two classes include a total of 238 and 151 GPS stations respectively. The K distribution of the considered stations is shown on the left panel and the

the vertical, north and east components. The subscript and superscript numbers represent the 25th and 75th percentile of the distribution. Columns labeled with B and G refer respectively to roof-based and ground-based monuments.

and no clear differences can be found between the monument types, not even between the ground-based tripods, generally considered as the best performing monuments. This partially confirms previous findings concerning regional networks mainly located in the US (e.g. Williams et al. 2004; Langbein 2008) in which the inferred spectral indexes of the power law noise range between  $-0.9$  and  $-1.4$ . Similarly a recent paper by Santamaria-Gómez et al. (2011) treating globally distributed time series concluded that 71% of the sites show a dominant flicker noise component ( $K = -1$ ). Beavan (2005) reports lower values with averages ranging between  $-0.6$

WN amplitude distribution on the right panel. The numbers in the insets refer to median values of the relevant distributions. Note that the darker green arises from the intersection of the two G and B classes

and  $-0.5$  for the continuous New Zealand GPS network. Since the latter findings are based on short time series (2–4.5 years), low K values are not surprising since the low frequency content cannot be properly resolved. A different study that assesses the stability of GPS monumentation from short-baseline time series shows spectral indexes very close to those presented in this study but with noise amplitudes at least one order of magnitude lower (King and Williams 2009). Furthermore King and Watson (2010) demonstrate that unmodeled multipath or subtle variations in the GPS constellation are able to produce several millimeters tem-



**Fig. 5** Histograms of the ground-based tripods (G-tripod, *green bars*) and the roof-based steel antenna mast (B-roof, *blue bars*), the two classes count 67 and 143 GPS stations respectively. The K distribution of the considered stations is shown on the *left panel* and the WN

amplitude distribution on the *right panel*. The numbers in the *insets* refer to median values of the relevant distributions. Note that the *darker green* arises from the intersection of the two G and B classes

poral variations in time series. Thus it is hard to believe that what we observe in our time series is a real monument motion, we instead believe that errors in geophysical models or biases in the reference frame (orbits and Earth orientation) may be the main source of the observed correlated noise. We are aware that this work is still at a preliminary stage, in future we intend to use more sophisticated approaches (e.g. Bos et al. 2013) that allow us to treat massive data and assess the noise content of extensive GPS networks in different ways. This is an important step in order to be able to resolve as close as possible the long term noise component (random-walk) and thus compute reliable uncertainties for the GPS velocities (see e.g. Langbein 2012).

## References

- Avallone A et al (2010) The RING network: improvement of a GPS velocity field in the central Mediterranean. *Ann Geophys* 53(2):39–54. doi:[10.4401/ag-4549](https://doi.org/10.4401/ag-4549)
- Beavan J (2005) Noise properties of continuous GPS data from concrete pillar geodetic monuments in New Zealand and comparison with data from U.S. deep drilled braced monuments. *J Geophys Res* 110, B08410. doi:[10.1029/2005JB003642](https://doi.org/10.1029/2005JB003642)
- Beutler G et al (2007) Bernese GPS software version 5.0. In: Dach R, Hugentobler U, Fridez P, Meindl M (eds) *Astronomical Institute, University of Bern, Bern*
- Bos MS, Fernandes RMS, Williams SDP (2013) Fast error analysis of continuous GNSS observations with missing data. *J Geod* 87:351–360. doi:[10.1007/s00190-012-0605-0](https://doi.org/10.1007/s00190-012-0605-0)
- Devoti R, Pietrantonio G, Pisani AR, Riguzzi F, Serpelloni E (2010) Present day kinematics of Italy. In: Beltrando M, Peccerillo A, Mattei M, Conticelli S, Doglioni C (eds) *J Virtual Explor* 36(2). doi:[10.3809/jvirtex.2009.00237](https://doi.org/10.3809/jvirtex.2009.00237)
- Johnson H, Agnew DC (2000) *Correlated noise in geodetic timeseries. U.S. Geological Survey Final Technical Report, FTR-1434-HQ-97-GR-03155*
- King MA, Watson CS (2010) Long GPS coordinate time series: multipath and geometry effects. *J Geophys Res* 115, B04403. doi:[10.1029/2009JB006543](https://doi.org/10.1029/2009JB006543)
- King MA, Williams SDP (2009) Apparent stability of GPS monumentation from short-baseline time series. *J Geophys Res* 114, B10403. doi:[10.1029/2009JB006319](https://doi.org/10.1029/2009JB006319)
- Lagarias JC, Reeds JA, Wright MH, Wright PE (1998) Convergence properties of the Nelder-Mead simplex method in low dimensions. *SIAM J Optim* 9(1):112–147
- Langbein J (2004) Noise in two-color electronic distance meter measurements revisited. *J Geophys Res* 109, B04406. doi:[10.1029/2003JB002819](https://doi.org/10.1029/2003JB002819)
- Langbein J (2008) Noise in GPS displacement measurements from Southern California and Southern Nevada. *J Geophys Res* 113, B05405. doi:[10.1029/2007JB005247](https://doi.org/10.1029/2007JB005247)

- Langbein J (2012) Estimating rate uncertainty with maximum likelihood: differences between power-law and flicker-random-walk models. *J Geod* 86:775–783. doi:[10.1007/s00190-012-0556-5](https://doi.org/10.1007/s00190-012-0556-5)
- Langbein J, Johnson H (1997) Correlated errors in geodetic time series: implications for time-dependent deformation. *J Geophys Res* 102:591–603
- Lomb NR (1976) Least-squares frequency analysis of unequally spaced data. *Astrophys Space Sci* 39:447–462
- Mao A, Harrison CGA, Dixon TH (1999) Noise in GPS coordinate time series. *J Geophys Res* 104:2797–2816
- Santamaria-Gómez A et al (2011) Correlated errors in GPS position time series: implications for velocity estimates. *J Geophys Res* 116, B01405. doi:[10.1029/2010JB007701](https://doi.org/10.1029/2010JB007701)
- Scargle JD (1982) Studies in astronomical time series analysis. II — statistical aspects of spectral analysis of unevenly spaced data. *Astrophys J* 263:835–853. doi:[10.1086/160554](https://doi.org/10.1086/160554)
- Wdowinski S et al (1997) Southern California permanent GPS geodetic array: spatial filtering of daily positions for estimating coseismic and postseismic displacements induced by the 1992 Landers earthquake. *J Geophys Res* 102(B8):18057–18070
- Welch PD (1967) The use of fast Fourier transform for the estimation of power spectra: a method based on time averaging over short, modified periodograms. *IEEE Trans Audio Electroacoust* AU-15:70–73
- Williams SDP et al (2004) Error analysis of continuous GPS time series. *J Geophys Res* 109, B03412. doi:[10.1029/2003JB002741](https://doi.org/10.1029/2003JB002741)
- Wyatt F (1982) Displacement of surface monuments - horizontal motion. *J Geophys Res* 87:979–989. doi:[10.1029/JB087iB02p00979](https://doi.org/10.1029/JB087iB02p00979)
- Wyatt FK (1989) Displacement of surface monuments - vertical motion. *J Geophys Res* 94:1655–1664. doi:[10.1029/JB094iB02p01655](https://doi.org/10.1029/JB094iB02p01655)
- Zhang J et al (1997) Southern California permanent GPS geodetic array: error analysis of daily position estimates and site velocities. *J Geophys Res* 102(B8):18035–18055. doi:[10.1029/97JB01380](https://doi.org/10.1029/97JB01380)

---

# VADASE: State of the Art and New Developments of a Third Way to GNSS Seismology

E. Benedetti, M. Branzanti, G. Colosimo, A. Mazzoni, and M. Crespi

---

## Abstract

In recent years, extensive work has been done to effectively exploit Global Navigation Satellite Systems (GNSS) for estimating important earthquake parameters such as the seismic moment and magnitude (i.e. GNSS Seismology). The rapid and accurate assessment of these parameters is of crucial importance to achieve reliable tsunami generation scenarios and eventually dispatch an early warning. In this framework, Geodesy and Geomatics division (AGG) of Sapienza University of Rome developed a new approach to obtain in real-time the 3D displacements of a single GNSS receiver. This solution, called VADASE (Variometric Approach for Displacement Analysis Standalone Engine), utilizes the broadcast orbits and the time differences of the high-rate (i.e. 1 Hz or more) carrier phases observations to ascertain the receiver movements over short intervals at a few centimeters accuracy level in real-time.

First we summarize the state-of-art of VADASE. Then, we illustrate the most recent developments of the algorithm, which include model refinements, single frequency (L1) capability and functionality with Galileo real data. Finally, we present the first results of an automatic procedure enabled by VADASE real-time capabilities. The epoch-by-epoch displacements (i.e. velocities) of approximately 100 stations of the IGS (International GNSS Service) high-rate (i.e. 1 Hz) network are retrieved every 15 min using VADASE, and the whole network can be characterized in terms of noise level (ranging from 1 to 5 mm/s for the horizontal and from 2 to 10 mm/s for the height); on this basis, corresponding thresholds (i.e. 3-sigma) could be set up in order to highlight significant displacements caused by an earthquake and eventually raise a tsunami alarm.

---

## Keywords

Galileo • GNSS Seismology • Real-time • Single frequency • VADASE

---

## 1 Introduction

The Global Positioning System (GPS) has been repeatedly proven to be a powerful tool to estimate coseismic dis-

placements and waveforms, with accuracies ranging from few millimeters to few centimeters. These promising results were achieved following two main strategies: the Differential Positioning (DP) and the Precise Point Positioning (PPP) (Bock et al. 1993; Kouba 2003; Larson et al. 2007; Larson 2009; Ohta et al. 2012; Xu et al. 2012; Hung and Rau 2013).

In particular, GPS-derived displacement waveforms can contribute both to the modelling of fault rupture and to the seismic moment estimation, since GPS is not affected by the saturation problems experienced by seismometers located near the epicenters of strong earthquakes. In the last years

---

E. Benedetti (✉) • M. Branzanti • G. Colosimo • A. Mazzoni • M. Crespi

Geodesy and Geomatics Division, Department of Civil, Constructional and Environmental Engineering, University of Rome La Sapienza, Via Eudossiana 18, Rome, Italy  
e-mail: [elisa.benedetti@uniroma1.it](mailto:elisa.benedetti@uniroma1.it); [gabriele.colosimo@uniroma1.it](mailto:gabriele.colosimo@uniroma1.it)

some authors (Bock et al. 2000; Langbein and Bock 2004; Blewitt et al. 2006; Bock and Genrich 2006) addressed the problem to retrieve displacement waveforms in real-time, with accuracies of few centimeters, from GPS high-rate observations (1 Hz or more).

In this context, the variometric approach VADASE (Variometric Approach for Displacements Analysis Standalone Engine) has been proposed in Colosimo et al. (2011a) and Colosimo (2013) as a third way to GNSS Seismology. The approach is based on time single differences of carrier phase observations continuously collected using a standalone GPS receiver and standard GPS broadcast products (orbits and clocks) that are available in real-time. Hence, one receiver works in standalone mode and the epoch-by-epoch displacements (equivalent to velocities) are estimated. Then, they are summed over the time interval when the earthquake occurred to retrieve coseismic displacements and waveforms. Since VADASE does not require either additional technological complexity or a centralized data analysis, in principle, it can be embedded into the GPS receiver firmware and work in real-time. Moreover, differently from DP and PPP, VADASE does not require phase ambiguity fixing and it is also able to work with single frequency data only. The effectiveness of VADASE was already proved through the application to the catastrophic Tohoku-oki earthquake (United States Geological Survey (USGS) M = 9.0, March 11, 2011, 05:46:24 UTC) (Colosimo et al. 2011b; Branzanti et al. 2013).

Here we present the state of the art of the implementation and the new developments and applications of VADASE.

In Sect. 2 a short description of the variometric approach estimation model is recalled and main developments of VADASE with respect to its first implementation are presented.

In Sect. 2.2 VADASE model refinements are discussed. In Sect. 2.3 VADASE single frequency (L1) capability is presented through its application to the Emilia earthquake (United States Geological Survey (USGS) M = 6.0, May 20, 2012, 02:03:51 UTC) Pondrelli et al. (2012). In Sect. 2.4 the first application of VADASE to Galileo real data is described. In Sect. 3 a different application of VADASE is presented: the (near) real time network monitoring. In Sect. 4 we present our conclusions and discuss future research directions for GNSS Seismology, in particular, toward the real-time application of the variometric approach, considering observations collected from geodetic (dual frequency and multi-constellation) and low-cost (single-frequency) receivers.

## 2 VADASE State of the Art and New Developments

### 2.1 VADASE Fundamentals

VADASE is an algorithm able to estimate, on the basis of carrier phase observations and broadcast orbits, the velocity of a GNSS receiver between two observations epochs. The receiver displacements waveforms, for short intervals (few minutes), can be retrieved from the estimated velocities by simple integration. The results presented in this section are all obtained using observation and broadcast orbit RINEX (Receiver Independent Exchange Format) files.

Here we recall the functional model of the least square estimation of the variometric approach in order to better assess the developments discussed in the next subsections. For a complete description of the VADASE estimation model, please refer to Colosimo et al. (2011a), Colosimo (2013), and Branzanti et al. (2013)

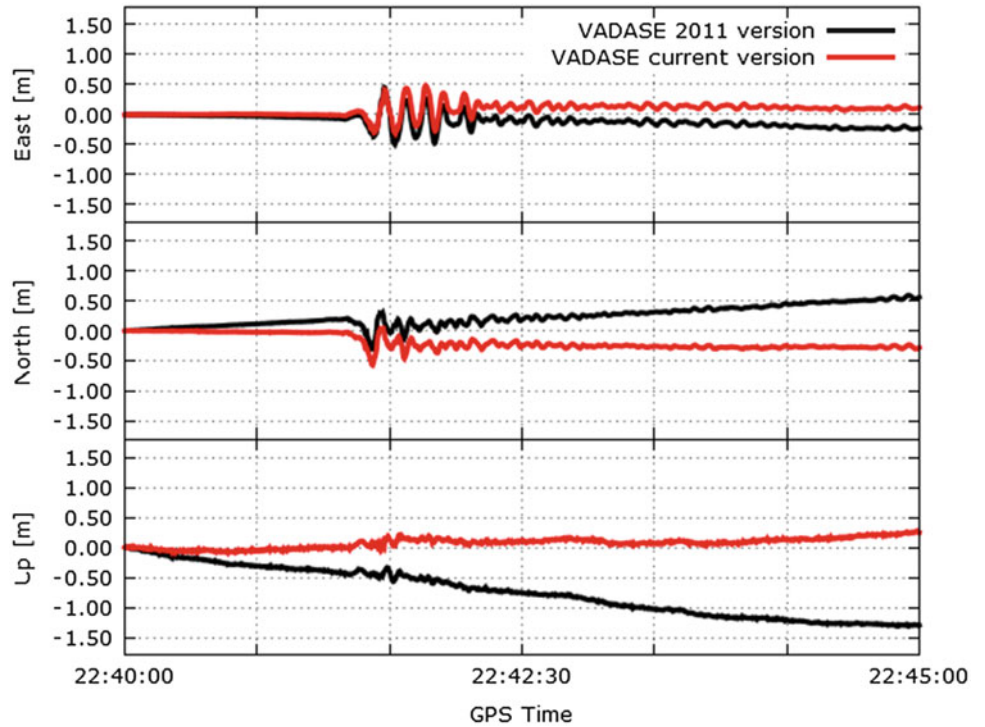
We assume that subscript  $r$  refers to a particular receiver and superscript  $s$  refers to a satellite;  $\Phi_r^s$  is the carrier phase observation of the receiver with respect to the satellite;  $\lambda$  is the carrier phase wavelength;  $\rho_r^s$  is the geometric range (i.e., the distance between the satellite and the receiver);  $c$  is the speed of light;  $\delta t_r$  and  $\delta t^s$  are the receiver and the satellite clock errors, respectively;  $T_r^s$  is the tropospheric delay along the path from the satellite to the receiver;  $p_r^s$  is the sum of the other effects (relativistic effects, phase center variations, and phase windup); and  $m_r^s$  and  $\epsilon_r^s$  represent the multipath and the noise, respectively. Equation (1) is the difference in time ( $\Delta$ ) between two consecutive epochs ( $t$  and  $t+1$ ) of carrier phase observations in the ionospheric-free combination ( $\alpha$  and  $\beta$  are the standard coefficients of L3 combination referred to the two phases L1 and L2)

$$\begin{aligned} \alpha[\lambda\Delta\Phi_r^s]_{L1} + \beta[\lambda\Delta\Phi_r^s]_{L2} = & (e_r^s \bullet \Delta\xi_r + c\Delta\delta t_r) + \\ & + ([\Delta\rho_r^s]_{OR} - c\Delta\delta t^s + \Delta T_r^s + [\Delta\rho_r^s]_{EIOI} + \Delta p_r^s) + \\ & + \Delta m_r^s + \Delta\epsilon_r^s \end{aligned} \quad (1)$$

where  $e_r^s$  is the unit vector from the satellite to the receiver,  $\Delta\xi_r$  is the (mean) velocity of the receiver in the interval  $t$  and  $t+1$ ,  $[\Delta\rho_r^s]_{OR}$  is the change of the geometric range due to the satellite's orbital motion and the Earth's rotation,  $[\Delta\rho_r^s]_{EIOI}$  is the change of the geometric range due to the variation of the solid Earth tide and ocean loading.



**Fig. 1** Comparison between VADASE 2011 undretrended solutions Colosimo et al. (2011a) vs current version. Displacements for P496 station during the Baja California, Mexico earthquake (Mw 7.2, 4 April 2010)



The term  $(e_r^s \cdot \Delta \xi_r + c \Delta \delta t_r)$  contains the four unknown parameters (the 3-D velocity  $\Delta \xi_r$  and the receiver clock error variation  $\Delta \delta t_r$ ) and  $([\Delta \rho_r^s]_{OR} - c \Delta \delta t^s + \Delta T_r^s + [\Delta \rho_r^s]_{EIOI} + \Delta p_r^s)$  is the known term that can be computed on the basis of known orbits and clocks and of proper well-known models (for a complete orbits, clocks and atmosphere error analysis, please refer to Colosimo (2013)). The least squares estimation of the 3-D velocities is based upon the entire set of variometric equation (1), which can be written for two generic consecutive epochs ( $t$  and  $t + 1$ ). The number of variometric equations depends on the number of satellites common to the two epochs, and at least four satellites are necessary in order to estimate the four unknown parameters for each consecutive epoch couple.

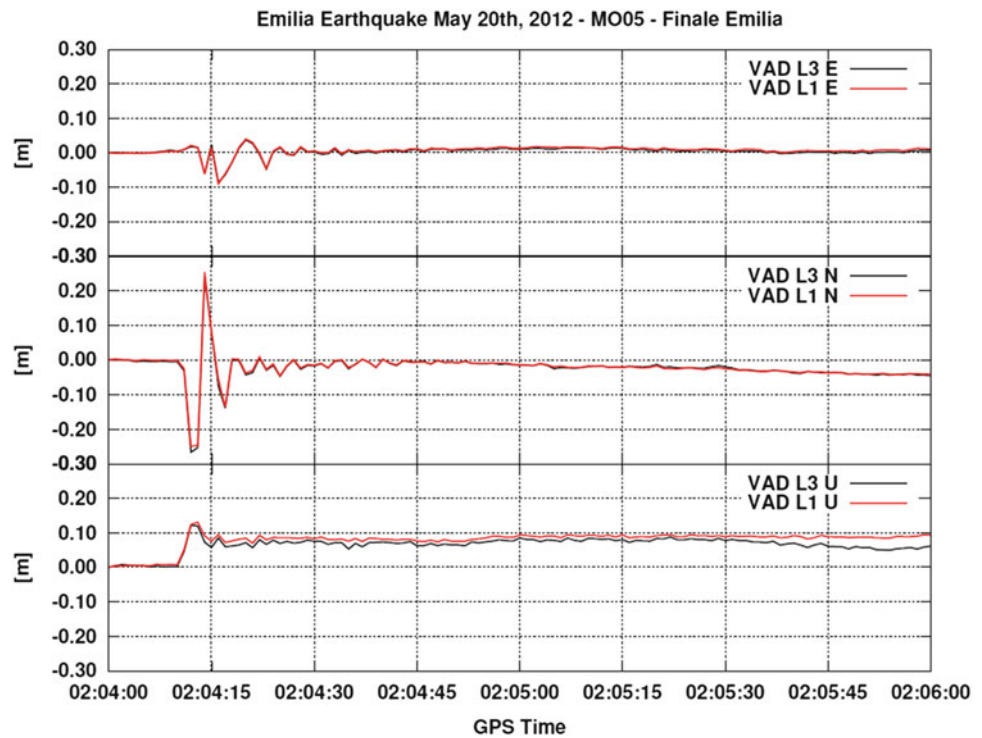
In the next subsections we present the main developments of the VADASE implementation.

## 2.2 Known Term Computation Algorithm Refinements

In Colosimo et al. (2011a), it was shown that the velocities estimated with VADASE were generally affected by bias that displayed their signature as a trend in the displacement waveforms computation obtained by simple velocities integration over time. Since the first implementation of VADASE, the software has been continuously developed

and some refinements have been carried out in the code. In particular, in the known term computation function some subroutines were added in order to improve the accuracy of the time used in the orbits computation. Specifically, an iterative routine has been added in order to perform a receiver clock estimation (in single epoch) based on code observations and a priori reference coordinates. Also a correct Sagnac effect model has been refined (Colosimo 2013). This code refinement significantly improved VADASE solutions in terms of amplitude of trends that cumulate in the displacements waveforms. In Fig. 1 is shown the comparison between the current and 2011 version of VADASE. In particular, we applied the current VADASE to the P496 5 Hz data during the Baja California, Mexico earthquake (Mw 7.2, 4 April 2010, 22:40:42 UTC) already analyzed in Colosimo et al. (2011a), where the displacement waveforms were affected by a trend. The current VADASE displacements waveforms are no longer affected by significant trends in a 5 min time interval. We want to remark that, in order to provide fast and reliable magnitude estimations (earthquake early warning), the estimation of static offset from GPS displacements waveforms is performed over short intervals. For the Tohoku-oki earthquake (USGS M = 9.0, March 11, 2011, 05:46:24 UTC), for example, reliable magnitude was estimated extracting static offsets from the Japanese GPS Earth Observation Network (GEONET) 120 s after the event time (Colosimo 2013).

**Fig. 2** Comparison between VADASE L3 and L1 solutions over the 120 s interval 02:04:00-02:06:00, May 20, 2012 (DOY 141), GPS time, 1 Hz data collected from MO05 station



### 2.3 VADASE Single Frequency (L1) Capability

Here we show the current capability of VADASE to process single frequency (L1) data applying the Klobuchar (1987) ionospheric model in order to remove the ionospheric delay. In Fig. 2 are shown displacements waveforms estimated for MO05 station 1 Hz data collected during the Emilia earthquake (United States Geological Survey (USGS)  $M = 6.0$ , May 20, 2012, 02:03:51 UTC, Pondrelli et al. 2012). For this earthquake, 7 permanent stations data were processed with VADASE and the solutions compared with the ones obtained with other well-established strategies and softwares (RMSE for L3 solutions with respect to the reference ones within 1.1 cm in horizontal and within 1.5 cm in height, L1 solutions with respect to L3 and the four reference ones within 1.7 cm in horizontal and within 1.8 cm in height, Benedetti et al. 2014) Here, in order to limit our discussion on VADASE single frequency capability, we show only VADASE L3 (ionospheric free combination) and L1 (single frequency L1) results for one station. The agreement in terms of Root Mean Square Error (RMSE) of the difference between the two solutions, over a 2 min time interval, is 0.4 cm in horizontal and 1.7 cm in height. VADASE in single frequency has been also tested with low cost receivers (U-blox). The obtained results are promising and the possibility of deploying dense networks of low cost receivers for monitoring purposes is currently under investigation.

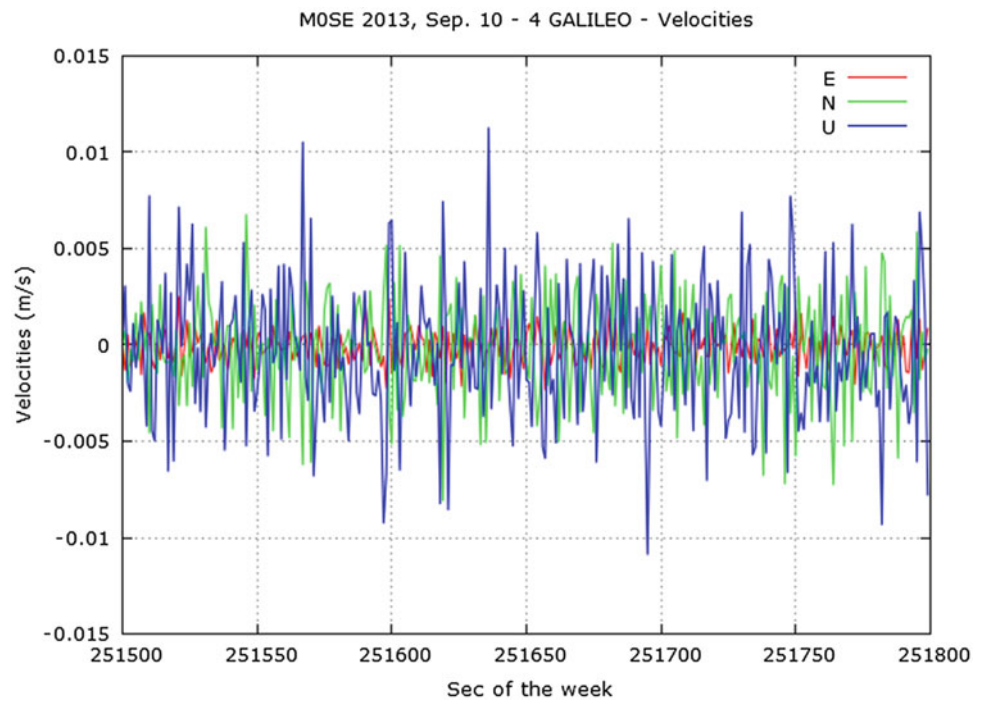
### 2.4 VADASE Galileo Capability

VADASE was already tested with Spirent simulated full Galileo constellation data (Colosimo 2013), we show now the first application of VADASE to real Galileo data.

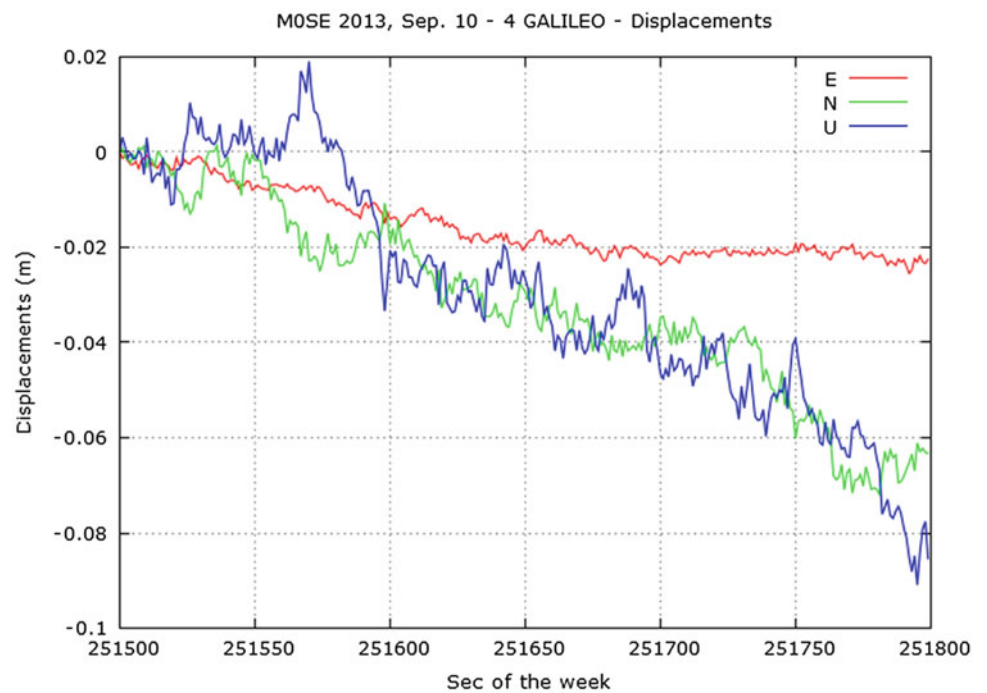
Specifically, MOSE permanent station 1 Hz RINEX 3.0 format data collected on 2013, Sep. 10 (86400 observed epochs) were processed using only Galileo observations. In standard configuration VADASE rejects non redundant solutions (i.e. less than 5 satellites observed in two consecutive epochs). For this application it was forced to run also with 4 satellites observed, since, at the moment, only Galileo E11, E12, E19 and E20 are operational.

For this first real data test, a single frequency configuration was chosen (L1C). The algorithm succeeded in all the consecutive epoch pairs with the above 4 Galileo satellites simultaneously observed (about 7,700). In Fig. 3 results for a 300 s sample are shown in terms of estimated velocities. The RMSE of the estimated velocities (compared to the null reference real velocity) are respectively at 1, 3 and 4 mm/s for East, North and Up components. In Fig. 4 the displacements obtained by simple integration of the estimated velocities are shown. Since only four Galileo satellites were used in the processing, this results are only to prove the VADASE Galileo capability and a deeper analysis of the accuracy achievable with only Galileo data will be carried out in the next future.

**Fig. 3** Sample of 300 s of VADASE estimated velocities—4 Galileo satellites data for MOSE permanent station 2013, Sep. 10



**Fig. 4** Sample of 300 s of VADASE displacement waveforms—4 Galileo satellites data for MOSE permanent station 2013, Sep. 10

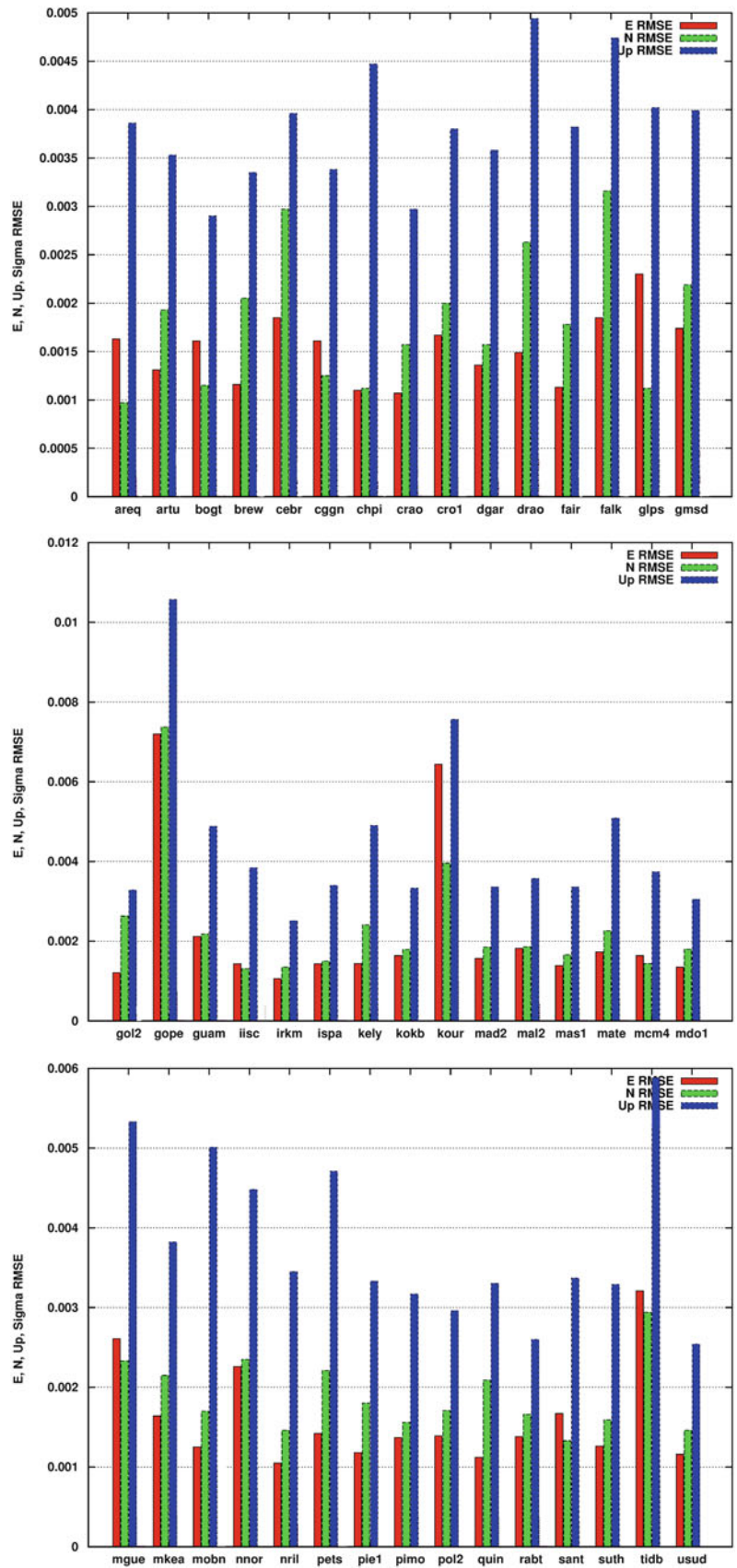


It is also important to underline that in this first only Galileo real data test, it was not possible to apply the ionospheric model correction since it was not recorded in the Galileo navigation file supplied by Leica Geosystems GR25 receiver. The VADASE ionospheric free combination for Galileo data is currently under development.

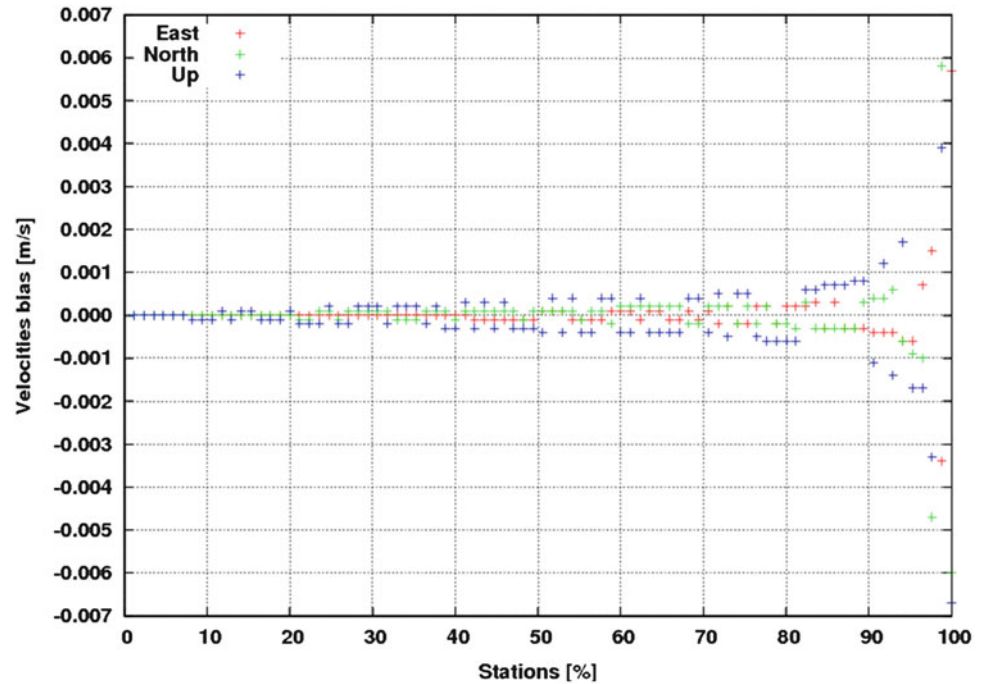
### 3 VADASE (Near) Real-Time Network Monitoring

In this section we present the application of VADASE as a tool for (near) real-time network monitoring. An automatic routine has been implemented in order to download from the

**Fig. 5** Sample statistics of VADASE application to IGS high rate sites (15-min, 1 Hz rate data). RMSE of East, North and Up velocities are in m/s



**Fig. 6** VADASE velocities bias over 15 min



high rate directory of the Crustal Dynamics Data Information System (CDDIS) the IGS (International GNSS Service) high rate stations data (Dow et al. 2009; Noll 2010). In this directory the observations (in RINEX format) of more than 100 worldwide stations are available. For each station, 15-min interval (1 Hz rate) are uploaded, as soon as they are collected. The implemented routine downloads all the available files referred to the last 15 min with respect an input time and automatically runs VADASE on the downloaded files. Then, all the VADASE output files (i.e. the epoch-by-epoch estimated parameters) are analyzed and for each station possible outliers with respect to threshold values are highlighted. Finally, statistics and plots are automatically generated.

In Fig. 5 statistics referred to one 15 min sample interval are shown for some stations. RMSE values for estimated velocities (with respect to the presumed null reference velocity) for East, North and Up components for each station are plotted. The RMSE values range (depending on permanent stations characteristics) from 1 to 5 mm/s for the horizontal and from 2 to 10 mm/s for the height. The distribution of velocities bias (over the same 15 min interval) shown in Fig. 6 supports the effectiveness of the current version of VADASE in bias mitigation (90% of stations between  $-0.001$  and  $0.001$  m/s in all components). This tool is useful to achieve two main results. On one hand it is possible to evaluate the noise level of each station in terms of velocity estimation and then monitor the stability of this level. On the other hand, it is possible to mark a threshold useful to highlight statistically significant fast movements.

## 4 Discussion and Conclusions

With respect to the first VADASE implementation some model refinements have been done. The current version is able to mitigate the biases in the velocity estimation and the detrending of the retrieved displacement waveform is no more needed when the velocity integration interval is limited to few minutes (what was proven sufficient and effective in order to estimate reliable earthquake magnitude Colombelli et al. 2013).

The VADASE single frequency capability has been fully developed. Some experiments with low cost single frequency receivers are in progress also to support the possibility of deploying dense networks of receivers for monitoring purposes.

First tests on real Galileo data have been successfully performed. The integration of the GPS and Galileo observations processing is currently under development.

An automatic routine able to continuously process 15 min length observation files of a worldwide permanent stations network as soon as they are available has been presented. Stacking the estimated solutions, it is possible to evaluate the noise level of each station in terms of velocity estimation. Once a station is characterized it is possible to mark a threshold useful to highlight statistically significant fast movements.

The software, written in C/C++ language, in its current version, runs with I/O files structure. Also a real-time data stream structure of the VADASE implementation is currently

under development in order to process real time observations broadcasted in RTCM (Radio Technical Commission for Maritime Services) format via NTRIP (Networked Transport of RTCM via Internet Protocol). Therefore it is possible to draw two different scenarios for VADASE real-time applications. In the first, VADASE runs directly on board GNSS receivers and broadcasts solutions (for example in NMEA-like (National Marine Electronics Association) format, a very light data in terms of transmission band); in this context VADASE has been already tested on board a receiver working in real-time with short interval observation files. In the second scenario VADASE runs in a dedicated server over real-time broadcasted observations; this functionality has been already tested in the frame of a cooperation with DLR Institute for Communications and Navigation at Oberpfaffenhofen. In conclusion, VADASE new developments and applications prove once more its effectiveness for GNSS Seismology.

**Acknowledgements** The authors thank the three anonymous Reviewers and the Editor in Chief for their valuable suggestions that helped improving the present work. The authors recognize the fundamental role of the International GNSS Service for delivering high-rate GNSS data in real time. The authors are indebted with Dr. Nicola Cenni, Prof. Paolo Baldi and Prof. Enzo Mantovani for providing the data of MO05 station. VADASE is subject of an international pending patent, generously supported by the University of Rome “La Sapienza”. VADASE was awarded the DLR (German Aerospace Agency) Special Topic Prize and the Audience Award at the European Satellite Navigation Competition 2010 and was partially developed thanks to 1-year cooperation with DLR Institute for Communications and Navigation at Oberpfaffenhofen (Germany).

## References

- Benedetti E, Branzanti M, Biagi L, Colosimo G, Mazzoni A, Crespi M (2014) GNSS seismology for the 2012 Mw = 6.1 Emilia Earthquake: exploiting the VADASE algorithm. *Seismol Res Lett* 85(3):649–656. doi:10.1785/0220130094
- Blewitt G, Kreemer C, Hammond WC, Plag HP, Stein S, Okal E (2006) Rapid determination of earthquake magnitude using GPS for tsunami warning systems. *Geophys Res Lett* 33(11):L11309. doi:10.1029/2006GL026145
- Bock Y, Agnew DC, Fang P et al (1993) Detection of crustal deformation from the landers earthquake sequence using continuous geodetic measurements. *Nature* 361(6410):337–340
- Bock Y, Nikolaidis RM, de Jonge PJ, Bevis M (2000) Instantaneous geodetic positioning at medium distances with the global positioning system. *J Geophys Res* 105(B12):28223–28253
- Bock Y, Genrich JF (2006) Instantaneous geodetic positioning with 10–50 Hz GPS measurements: noise characteristics and implications for monitoring networks. *J Geophys Res* 111(3):B03403. doi:10.1029/2005JB003617
- Branzanti M, Colosimo G, Crespi M, Mazzoni A (2013) GPS near-real-time coseismic displacements for the great Tohoku-oki earthquake. *IEEE Geosci Remote Sens Lett* 10(2):6265361, 372–376. doi:10.1109/LGRS.2012.2207704
- Colombelli S, Allen RM, Zollo A (2013) Application of real-time GPS to earthquake early warning in subduction and strike-slip environments. *J Geophys Res* 118(7):3448–3461. doi:10.1002/jgrb.50242
- Colosimo G, Crespi M, Mazzoni A (2011a) Real-time GPS seismology with a stand-alone receiver: a preliminary feasibility demonstration. *J Geophys Res* 116(11):B11302. doi:10.1029/2010JB007941
- Colosimo G, Crespi M, Mazzoni A, Dautermann T (2011b) Co-seismic displacement estimation: Improving tsunami early warning systems. *GIM Int* 25(5):19–23
- Colosimo G (2013) VADASE: a brand new approach to real-time GNSS seismology. Lambert Academic Publishing AG & Co KG, 180 pp, ISSN: 9783845438382 <https://www.lap-publishing.com/site/home/10>
- Dow JM, Neilan RE, Rizos C (2009) The international GNSS service in a changing landscape of Global Navigation Satellite Systems. *J Geod* 83(3–4):191–198
- Hung HK, Rau RJ (2013) Surface waves of the 2011 Tohoku earthquake: observations of Taiwans dense high-rate GPS network. *J Geophys Res* 118(1):332–345. doi:10.1029/2012JB009689
- Klobuchar JA (1987) Ionospheric time-delay algorithm for single-frequency GPS users. *IEEE Trans Aerosp Electron Syst* AES-23:325–331
- Kouba J (2003) Measuring seismic waves induced by large earthquakes with GPS. *Stu Geophys et Geod* 47(4):741–755. doi:10.1023/A:1026390618355
- Langbein J, Bock Y (2004) High-rate real-time GPS network at parkfield: Utility for detecting fault slip and seismic displacements. *Geophys Res Lett* 31(15):L15S20 1–4. doi:10.1029/2003GL019408
- Larson K, Bilich A, Axelrad P (2007) Improving the precision of high-rate GPS. *J Geophys Res* 112(5):B05422. doi:10.1029/2006JB004367
- Larson K (2009) GPS seismology. *J Geod* 83(3–4):227–233. doi:10.1007/s00190-008-0233-x
- Noll, CE (2010) The crustal dynamics data information system: A resource to support scientific analysis using space geodesy. *Adv Space Res* 45(12):1421–1440. doi:10.1016/j.asr.2010.01.018
- Ohta Y, Kobayashi T, Tsushima H, Miura S., Hino R, Takasu T, Fujimoto H, Iinuma T, Tachibana K, Demachi T, Sato T, Ohzono M, Umino N (2012) Quasi real-time fault model estimation for near-field tsunami forecasting based on RTK-GPS analysis: application to the 2011 Tohoku-oki earthquake (Mw 9.0). *J Geophys Res* 117(2):B02311. doi:10.1029/2011JB008750
- Pondrelli S, Salimbeni S, Perfetti P, Danecsek P (2012) Quick regional centroid moment tensor solutions for the Emilia 2012 (northern Italy) seismic sequence. *Ann Geophys* 55(4):615–621. doi:10.4401/ag-6159
- Xu P, Shi C, Fang R, Liu J, Niu X, Zhang Q, Yanagidani T (2012) High-rate Precise Point Positioning (PPP) to measure seismic wave motions: an experimental comparison of GPS PPP with inertial measurement units, *J Geod* 87(4):361–372. doi:10.1007/s00190-012-0606-z

---

# On the Spatial Resolution of Homogeneous Isotropic Filters on the Sphere

Balaji Devaraju and Nico Sneeuw

---

## Abstract

Interest in filtering on the sphere was rejuvenated by the necessity to filter GRACE data, which has led to the development of a variety of filters with a multitude of design methods. Nevertheless, a lacuna exists in the understanding of filters and filtered fields, especially signal leakage due to filtering and resolution of the filtered field. In this contribution, we specifically look into the latter aspect, where we take an intuitive and empirical approach instead of a rigorous mathematical approach. The empirical approach is an adaptation of the technique used in optics and photography communities for determining the resolving power of lenses. This resolution analysis is carried out for the most commonly used homogeneous isotropic filters in the GRACE community. The analysis indicates that a concrete number for the filters can only be specified as an *ideal* number. Nevertheless, resolution as a concept is described in detail by the modulation transfer function, which also provides some insight into the smoothing properties of the filter.

---

## Keywords

Empirical approach • Filters • Filtering on the sphere • Modulation transfer function • Spatial resolution

---

## 1 Introduction

The advent of the GRACE satellite mission (Tapley et al. 2004) has revived the subject of filtering on the sphere, which was long forgotten after the seminal contribution of Jekeli (1981). Filtering is and has been a central subject of GRACE data processing as the GRACE data needs to be filtered due to the presence of high-frequency noise that manifests itself as north-south stripes (Swenson and Wahr 2006). A variety of filters have been developed since the launch of the GRACE mission (e.g., Kusche 2007; Swenson and Wahr 2006), which has made it daunting to choose a filter as reflected by the some of the studies (e.g., Werth et al. 2009; Longuevergne

et al. 2010). Spatial resolution is integral to the discussion of filter choice as it determines the resolvability of features after filtering, and hence, the usability of the dataset for the given study.

In physical geodesy, spatial resolution of a gravity field, given in terms of spherical harmonic coefficients up to complete degree  $L$ , is expressed as the *half-wavelength* ( $\psi_{\frac{1}{2}}$ ) of the harmonic  $L$  at the equator.

$$\psi_{\frac{1}{2}} = \frac{\pi a_E}{L} \approx \frac{20,000}{L}, \quad (1)$$

where  $a_E$  is the semi-major axis of the ellipsoid approximating the Earth. The value  $\psi_{\frac{1}{2}}$  is the Nyquist-Shannon sampling required along the equator, and also approximately the spacing between the zeros of the Legendre polynomial of degree  $L$ . Due to the isotropic nature of spherical harmonics, this value is assumed to hold over the entire sphere. However, Laprise (1992) points out that the half-wavelength at equator

---

B. Devaraju (✉) • N. Sneeuw  
Institute of Geodesy, University of Stuttgart, Geschwister-Scholl-Str.  
24D, 70469 Stuttgart, Germany  
e-mail: devaraju@gis.uni-stuttgart.de

is one of many possible values for the resolution, and also proposes that at best such values can only be used as an upper limit. When we apply a filter to a band-limited field, we are still left with a field that is band-limited up to  $L$ , but with a different resolution. In one-dimensional Fourier analysis, the resolution of a filter is taken to be the 6 dB point, which is the filter width at half of the amplitude at the peak (Harris 1978). However, the idea of resolution has not been studied in the case of filters on the sphere. Thus, our quest in this contribution is to determine the resolution of a filtered field, albeit empirically.

We will first detail the mathematical background behind filtering in Sect. 2, following which we will illustrate the idea of spatial resolution and the methodology to determine it in Sect. 3. In the same section we will apply the method to some filters commonly used in the GRACE community to determine their effective resolution. We will then test the estimates for the resolution by filtering GOCE gravity anomalies in Sect. 4. Finally, we will summarize our findings and draw conclusions in Sect. 5.

## 2 Filters and Filtering on the Sphere

Any square integrable scalar function (e.g., gravity field)  $f(\omega)$  on a sphere ( $\Omega$ ) can be represented in terms of a spherical harmonic spectrum

$$f(\omega) = \sum_{l=0}^{\infty} \sum_{m=-l}^l F_{lm} Y_{lm}(\omega) = \sum_{l,m} F_{lm} Y_{lm}(\omega), \quad (2a)$$

$$Y_{lm}(\omega) = \begin{cases} N_{lm} P_{lm}(\cos \theta) e^{im\lambda}, & m \geq 0 \\ (-1)^m Y_{l,-m}^*(\omega), & m < 0 \end{cases}, \quad (2b)$$

$$N_{lm} = (-1)^m \sqrt{(2l+1) \frac{(l-m)!}{(l+m)!}}, \quad (2c)$$

$$F_{lm} = \int_{\Omega} f(\omega) Y_{lm}^*(\omega) d\Omega, \quad (2d)$$

$$d\Omega = \frac{1}{4\pi} \sin \theta d\theta d\lambda, \quad \omega = (\theta, \lambda). \quad (2e)$$

where  $Y_{lm}(\omega)$  are the geodetic normalized complex surface spherical harmonics of degree  $l$  and order  $m$  with  $|m| \leq l$ ;  $\theta, \lambda$  are the co-latitude and longitude of the point  $\omega$  on the sphere;  $F_{lm}$  are the geodetic normalized spherical harmonic coefficients of degree  $l$  and order  $m$ ;  $P_{lm}(\cos \theta)$  are the Associated Legendre functions of degree  $l$  and order  $m$ ; and  $N_{lm}$  is the normalization factor.

Filtering a scalar field  $f(\omega)$  given on the sphere  $\Omega$  can be performed by taking a weighted average of a region around the point of concern  $\omega$ , where the weights and the region are prescribed by a *filter function*. The filter functions on the sphere, like the covariance functions, are two-point functions in that the filter weights are specified for a pair of points. The pair constitutes the point whose filtered value is sought ( $\omega$ ), and the point, within the region specified by the filter function, to which the weight is applied ( $\omega'$ ). Thus, the filter function is denoted as  $b(\omega, \omega')$ , and the filtering operation to obtain the filtered field,  $\bar{f}(\omega)$ , is mathematically expressed as

$$\bar{f}(\omega) = \int_{\Omega'} f(\omega') b(\omega, \omega') d\Omega', \quad (3)$$

where the filter function satisfies the following condition:

$$\int_{\Omega'} b(\omega, \omega') d\Omega' = 1. \quad (4)$$

In this study, we will only be concerned about filters whose weights depend only on the spherical distance ( $\psi$ ) between the points  $\omega$  and  $\omega'$ . Such filters are rotationally symmetric and translation invariant, and hence called *homogeneous isotropic filter* functions. The homogeneity property of these filters renders this filtering operation as a convolution operation (Jekeli 1981). Further, due to their dependency only on the spherical distance, they take a special spectral form

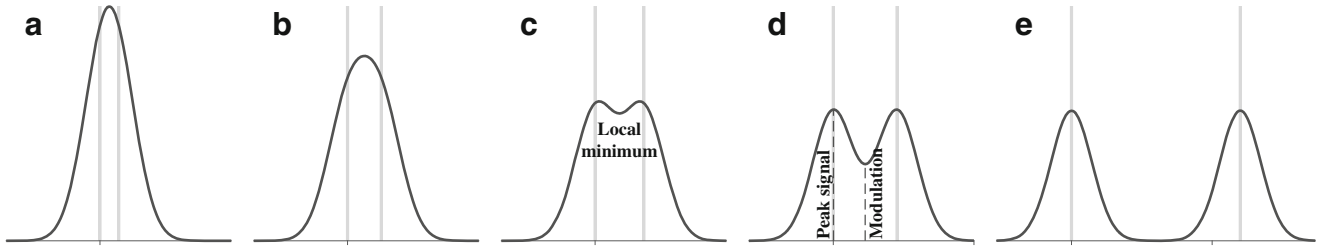
$$b(\psi) = \sum_{l=0}^{\infty} B_l \sum_{m=-l}^l Y_{lm}(\omega) Y_{lm}^*(\omega') \quad (5a)$$

$$= \sum_{l=0}^{\infty} (2l+1) B_l P_l(\cos \psi), \quad (5b)$$

where  $P_l(\cdot)$  are the unnormalized Legendre polynomials with  $\psi$  being the spherical distance between  $\omega$  and  $\omega'$ . As it can be seen from (5b) the speciality is that the spectrum of the filter depends only on the degree of the spherical harmonics, which makes them ideal for designing a variety of filters (e.g., Jekeli 1981; Sardeshmukh and Hoskins 1984). Now, combining (2a), (3) and (5a), we get the spatial and spectral forms of the filtering operation.

$$\bar{f}(\omega) = \int_{\Omega'} f(\omega') b(\psi) d\Omega' = \sum_{l,m} Y_{lm}(\omega) B_l F_{lm}. \quad (6)$$





**Fig. 1** Methodology for determining the resolution of homogeneous isotropic filters. The *light gray lines* indicate the unfiltered input signal and the *dark gray lines* indicate the filtered output. The level of

distinction between the two peaks in the filtered and resolved field can be quantified by the quantity modulation. The method is illustrated using the Gaussian filter of radius 285 km

### 3 Spatial Resolution

Spatial resolution defines the smallest possible feature that can be identified distinctly from its surroundings (Lillesand and Kiefer 1994). In signal processing, resolution is in general associated with sampling as expounded by the corresponding sampling theorem, for example, Nyquist-Shannon sampling theorem in the Euclidean space. As mentioned earlier, it becomes difficult to use such definitions once some filtering is performed on the field that changes its resolution. The optical and remote sensing communities use resolution charts to determine the resolution of lenses, which can be adapted for our problem. The resolution chart consists of a number of vertical and horizontal lines of varying thickness drawn at varying spacings. These charts are then imaged by the sensors and the least distance between the fully resolved lines is taken as the resolution of the lenses (Lillesand and Kiefer 1994). We base our method on this technique to determine the resolution of the homogeneous isotropic filters.

#### 3.1 Methodology

We define a scalar field  $g(\omega)$  on a unit sphere defined as

$$g(\omega) = \sum_{j=P,Q} \delta(\omega, \tilde{\omega}_j), \quad (7)$$

where  $\delta(\cdot, \cdot)$  is the *Dirac's pulse* on the sphere (Freedman and Schreiner 2009), located at the points P and Q, which are separated by a spherical distance of  $\psi_{PQ}$ . The Dirac's pulse is defined as a pulse with unit area at the point where it is located, which is expressed as

$$\int_{\Omega} \delta(\omega, \tilde{\omega}_j) d\Omega = \begin{cases} 1, & \omega = \tilde{\omega}_j \\ 0, & \text{elsewhere} \end{cases}. \quad (8a)$$

Its spherical harmonic spectrum is given as

$$\delta(\omega, \tilde{\omega}_j) = \begin{cases} \sum_{l,m} Y_{lm}(\omega) Y_{lm}^*(\tilde{\omega}_j), & \omega = \tilde{\omega}_j \\ 0, & \text{elsewhere} \end{cases}. \quad (8b)$$

Smoothing  $g(\omega)$  with an isotropic filter gives

$$\begin{aligned} \bar{g}(\omega) &= \int_{\Omega'} b(\psi) g(\omega') d\Omega', \\ &= \int_{\Omega'} \sum_{l,m} B_l Y_{lm}(\omega) Y_{lm}^*(\omega') \sum_{j,n,k} Y_{nk}(\omega') Y_{nk}^*(\tilde{\omega}_j) d\Omega', \\ &= \sum_j \sum_{l,m} B_l Y_{lm}(\omega) \sum_{n,k} Y_{nk}^*(\tilde{\omega}_j) \delta_{ln} \delta_{mk}, \\ &= \sum_j \sum_{l,m} B_l Y_{lm}(\omega) Y_{lm}^*(\tilde{\omega}_j). \end{aligned} \quad (9)$$

The above equation suggests that the filtered field  $\bar{g}(\omega)$  is the sum of the weights at  $\omega_P$  and  $\omega_Q$ . Since the filters that we deal with in this study are rotationally symmetric and homogeneous, (9) can be rewritten as

$$\bar{g}(\omega) = \sum_j \sum_{l,m} B_l Y_{lm}(\tilde{\omega}_j) Y_{lm}^*(\omega).$$

This suggests that the  $\bar{g}(\omega)$  is equivalent to the sum of the filter located at the points P and Q. We scrutinize the filtered field to see if the two Dirac's pulses are resolved, if not, we increase the separation  $\psi_{PQ}$  between the signals until they can be seen distinctly in the filtered field. The sequence is depicted in Fig. 1. This method was also employed by Harris (1978) to demonstrate the spectral resolution of different filter windows in the harmonic analysis of time-series. We set up the field  $g(\omega)$  for the computation by placing the Dirac's pulses at 30° (P) and 30.5° (Q) co-latitudes in the zero meridian, and increase the separation in steps of 0.1°. This

set up was chosen for ease of computation as the separation  $\psi_{PQ}$  is the difference between the two positions. All the computations were performed in the space domain in order to avoid truncation errors.

### 3.2 Modulation Transfer Functions

The concept of spatial resolution does not merely stop at the point where we are able to identify the two signals as distinct. It continues with the qualitative question of how distinct are those signals from each other before and after filtering. At the point of resolution, the two signals are distinct, but as we separate the two signals a bit further, then they are readily recognizable as two different entities (cf. Fig. 1). Therefore, there is a need for quantifying the level of distinctness of the resolved signals in comparison with the unfiltered signals. Again, this has already been treated by the remote sensing community, where they use the concept of *modulation transfer function* (MTF) to quantify the distinction between the two signals.

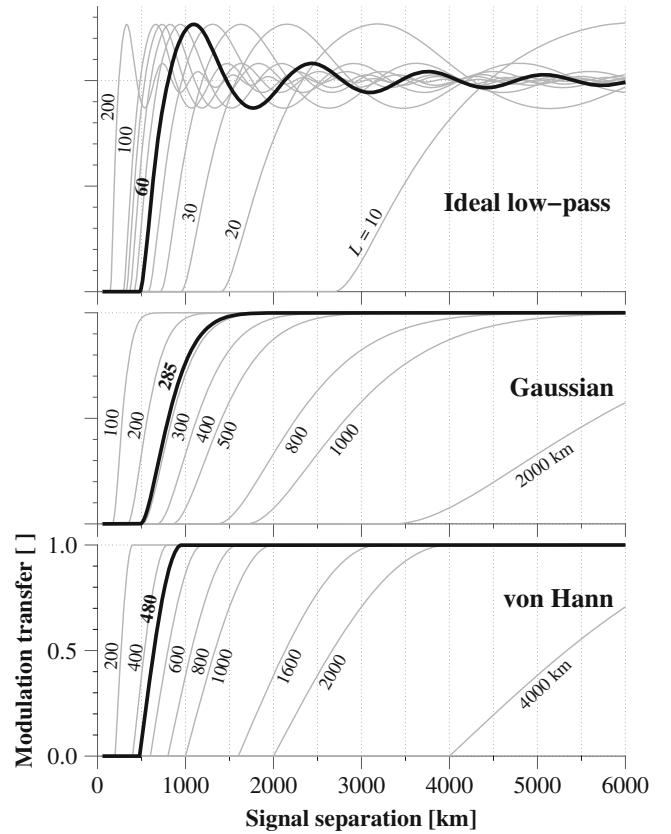
Originally, in the unfiltered field, there is no signal in the region between the Dirac pulses, but due to filtering we initially see only one ‘peak’ (Fig. 1a, b) and then a ‘valley’ between the two resolved ‘peaks’ (Fig. 1c). As we separate the signals farther apart, beyond the point of resolution, the valley deepens (Fig. 1d). We will denote the ordinate of the local minimum in the valley as *modulation*, which when zero indicates completely resolved signals (Fig. 1e). As we deal with rotationally symmetric functions, this local minimum in the valley will always occur at the mid-point of the geodesic between P and Q. Thus, modulation is

$$\text{Modulation} = 2b \left( \frac{\psi_{PQ}}{2} \right). \quad (10)$$

However, the peak value does not share such a property as it shifts from the mid-point of the geodesic between P and Q to the points P and Q themselves. Therefore, they have to be determined empirically. Now, we can set our rules for estimating resolution: When modulation is less than peak signal, we can say that the signals are resolved. Since the modulation and peak signal values are subjective to the filter function, we need to devise a relative measure, which we denote as the *modulation transfer* (MT). It is defined as

$$\text{MT} = 1 - \frac{\text{Modulation}}{\text{Peak signal}}. \quad (11)$$

As the name suggests MT indicates how the original modulation between the two signals is transferred to the filtered field. The MT takes a value zero until the signals are



**Fig. 2** Modulation transfer functions of three homogeneous isotropic filters for different filter parameters—filtering radius is the parameter for Gaussian and von Hann filters and maximum degree of spherical harmonic expansion for ideal low-pass filter. The *black lines* indicate the functions for the filters used in Fig. 3

resolved, because the peak signal always resides at the mid-point until the signals are resolved. As soon as the signals are resolved the valley starts appearing, and the value of modulation starts decreasing. Therefore, the value of MT starts to increase. Further, by plotting the MT against the signal separation for a given set of filter parameters we get a unique curve, which we denote as the *modulation transfer function* MTF (Fig. 2). The important feature of the MTF curve is the slope of the curve between 0 and 1 MT, which directly depends on the speed at which the filter function decays from peak value to zero. This is indicated by (10), where filters which decay slowly to zero (e.g., Gaussian) will not reduce quickly to zero modulation, and therefore, the corresponding MTF will have a gentle slope and *vice versa* (cf. Fig. 2).

### 3.3 Resolution of a Band-Limited Field

The term band-limited refers to the select bandwidth of frequencies that represents a given field in the spherical

**Table 1** Definitions, and spectral and spatial cross-sections of the three homogeneous isotropic filters in this study

Filter	Definition
Ideal low-pass	$B_l = \begin{cases} 1, & 0 \leq l \leq L \\ 0, & l > L \end{cases}$
Gauss	$b(\psi) = \beta \frac{e^{-\beta(1-\cos\psi)}}{1 - e^{-2\beta}}, \quad \beta = \frac{\ln(2)}{1 - \cos\psi_0}$
von Hann	$b(\psi) = \begin{cases} \frac{\eta}{2} \left(1 + \cos \frac{\pi\psi}{\psi_0}\right), & 0 \leq \psi \leq \psi_0 \\ 0, & \psi_0 \leq \psi \leq \pi \end{cases},$ $\eta = \frac{2(\pi^2 - \psi_0^2)}{\pi^2(1 - \cos\psi_0) - 2\psi_0^2}$
	Ideal low-pass    Gaussian    von Hann
$b(\psi)$	
$B_l$	

The quantity  $\psi_0$  is the smoothing radius defined for the Gaussian and the von Hann filters and  $L$  is the cut-off degree for the ideal low-pass filter

harmonic spectral domain. In spherical harmonics, the spectral frequencies are denoted by the degree ( $l$ ) and order ( $m$ ) of the spherical harmonic expansion. In general, gravity field estimates are disseminated as a spherical harmonic expansion up to a maximum spherical harmonic degree ( $L$ ), thereby making the field band-limited between the frequencies 0 and  $L$ . As per (2a), by definition the spherical harmonic expansion extends up to infinity, but in practice these are truncated up to  $L$  due to a variety of reasons: spatial sampling, measurement accuracy and computational limits. This band-limitation can be expressed as a *low-pass* filter as shown in Table 1. This filter is also referred to as the ideal low-pass filter, Shannon window and box-car filter. Due to the nature of its spectrum the filter is homogeneous and isotropic, which allows us to determine the resolution of a band-limited field using the method proposed in Sects. 3.1 and 3.2.

The MTF of the ideal low-pass filter is shown in Fig. 2 for a range of spherical harmonic degrees. The striking features of the curves are their steep slopes and the oscillation of the modulation transfer values around one. While the steep slope of the MTF indicates that the filters decay to zero very quickly, the oscillation is caused by the well-known ringing effect caused by the truncation of a harmonic series at a finite degree  $L$ . Further, the magnitude of the overshoot and its convergence to unity clearly depends on the number of spherical harmonic degrees involved in the synthesis: the more the harmonic degrees the less the overshoot and faster the convergence. This is a characteristic feature of the ideal low-pass filter.

### 3.4 Resolution of Homogeneous Isotropic Filters

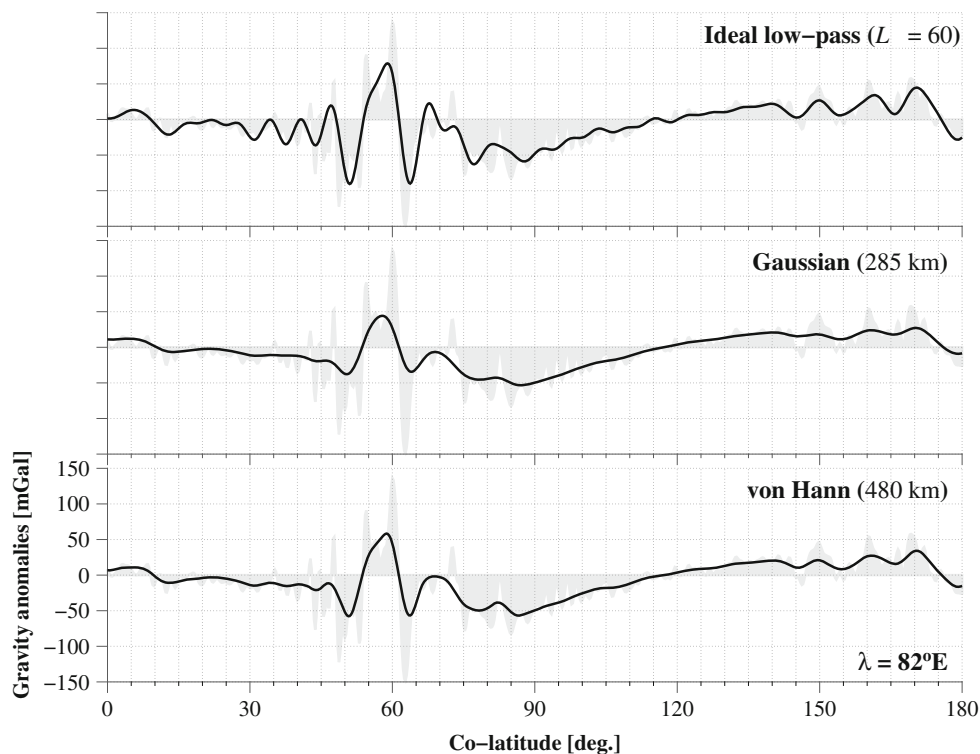
Here, we will take a look into two well-known homogeneous isotropic filters, Gaussian and von Hann filters, of which the Gaussian filter is the most widely used filter in the GRACE community due to its relative ease of implementation and also its effectiveness in smoothing out noise. These two filters are adaptations of their counterparts in one-dimensional harmonic analysis, and were adapted by Jekeli (1981) amongst a host of other filters. Their filter definitions are given in Table 1, where it is obvious that both the filters are parameterized via the filter radius. Also what is obvious is the continuous nature of the Gaussian and the piece-wise nature of the von Hann filters.

The Gaussian function is a unique function in that its spatial and spectral forms are both bell-shaped and its tail converges to zero asymptotically. Further, the filter radius of the generic form of the Gaussian filter is the spread of the filter, but for the sake of convenience the filter radius is defined at a certain fraction of the peak value—in geodesy the fraction is taken to be one-half of the peak value. It is also interesting to note that a von Hann filter, whose radius is twice that of the Gaussian filter radius, gets half of its peak value at half of the filter radius. This unusual coincidence allows us to tacitly verify whether resolution of a filter can be defined using the 6 dB point. Due to this reason, von Hann filters whose filtering radii are twice that of Gaussian filters are chosen. The MTF of Gaussian and von Hann filters are shown in Fig. 2 for a variety of filtering radii.

The MTF of the Gaussian filter shows that the resolution is much more than the filter radius, which negates the use of 6 dB point as the resolution of the filter. It rises from zero and converges to one in an asymptotic manner, while increasing along a gentle slope. Due to these characteristics, the Gaussian filter takes a lot of distance to completely resolve the two input pulses. For example, the Gaussian 500 km filter—a very widely used filter radius for the Gaussian filter in the GRACE community—needs a signal separation of at least 2,500 km to completely resolve the signals. This can be explained by the asymptotic nature of the Gaussian filter, which implies that in realistic signal scenarios the Gaussian will provide smoother signals.

The MTF curves for the von Hann filter show that the resolution is slightly more than the filter radius itself. The slope of those curves are steeper than those of the Gaussian MTF curves, but gentler than those of the ideal low-pass filter. The piece-wise nature of the filter is clearly reflected in the way the MT values increase sharply from zero. In contrast to the ideal low-pass filter, both the Gaussian and von Hann filter MTF curves do not show any overshoot beyond 1 because they do not suffer from ringing effect.

**Fig. 3** Comparison of unfiltered and filtered GOCE data along a longitudinal profile of  $82^\circ\text{E}$ . The filter parameters are given in the brackets adjacent to the label indicating the filter



#### 4 An Example: Filtering GOCE

We now put to test our estimates of the resolution of the three filters, and for this purpose we use the Gravity field and steady state Ocean Circulation Experiment (GOCE) (Rummel et al. 2011) data computed by Bruinsma et al. (2013) and apply the three filters to the GOCE gravity anomalies. For the case of the ideal low-pass filter, we chose a cut-off degree of  $L = 60$ , which provides a resolution of 490 km. The Gaussian filter of radius 285 km and von Hann filter of radius 480 km provide the same resolution as the ideal low-pass filter, and therefore, they can be compared with each other to study their resolving abilities. In Fig. 3, we show the filtered fields as a profile along the longitude  $82^\circ\text{E}$  overlaid with a  $5^\circ$  grid. The profile we have chosen here provides a lot of features that are both smaller and larger compared to the maximum achievable resolution.

The ideal low-pass filter, despite experiencing ringing effect, provides a more detailed filtered output, and the Gaussian provides the smoothest output. The Gaussian filtered profile shows very few details in addition to a strong reduction in the amplitude. In terms of detail, the von Hann filtered profile is no better, but it certainly retains a lot more amplitude than the Gaussian. Here, the level of detail and the amplitude retained is associated with the slope of the corresponding MTF curve, because the slope determines how fast the filter is able to resolve two signals completely. To

understand this let us consider the two negative peaks located at  $80^\circ$  and  $85^\circ$  which are separated by a narrow wedge. These two peaks are  $5^\circ$  ( $\approx 550$  km) apart, and therefore, we can expect them to be resolved as the separation is greater than the resolution of all the three filters. As expected the two peaks are resolved by all the three filters, but the modulation of the wedge differs according to the filter used. This is important while choosing a filter, because for changing the resolution of a given field one might be interested only in losing finer details but retaining the amplitudes as much as possible in which case a filter with a steep slope in the MTF must be preferred. The situation is opposite in the case of filtering data to reduce noise, where the details are mostly corrupted by noise, and hence a need for a smoother filtered field. This explains the success enjoyed by the Gaussian filter in negotiating the noise levels in GRACE datasets.

While the above example vindicates our method for determining resolution, there are also some examples that illustrate the ideal nature of our method. For example, the two positive peaks at  $55^\circ$  and  $60^\circ$  are again separated by  $5^\circ$ , but they are unequal signals and there is some signal between the peaks. Based on the distance between the two peaks we might expect them to be resolved, but they are unresolved mainly because they are unequal signals, with the peak at  $55^\circ$  significantly higher than the one at  $60^\circ$ . This is the reason we term the signal separation at the least non-zero modulation transfer of a given filter the *ideal resolution*.

## 5 Summary and Conclusions

In this study, we outlined a method to determine the resolution of homogeneous isotropic filters defined on the sphere. Initially, our aim was to provide a single number for the resolution, but it was shown that the idea of resolution can be qualitatively extended to describe the distinction between signals via modulation transfer. The modulation transfer when plotted against a host of signal separation distances between the distinct signals gives the modulation transfer function. An interesting feature of the modulation transfer function is that its slope describes the smoothing properties of the filter. Nevertheless, in cases where only one number needs to be specified we propose the use of the ideal resolution, which is the signal separation at the least non-zero modulation transfer. Further, it was also shown in the GOCE filtering example that care must be exercised in using these numbers as the modulation transfer function is constructed from an ideal scenario.

## References

- Bruinsma SL, Förste C, Abrikosov O, Marty JC, Rio MH, Mulet S, Bonvalot S (2013) The new ESA satellite-only gravity field model via the direct approach. *Geophys Res Lett* 40:1–6. doi:10.1002/grl.50716
- Freeden W, Schreiner M (2009) *Spherical functions of mathematical geosciences: a scalar, vectorial, and tensorial setup*. Springer, Berlin/Heidelberg
- Harris FJ (1978) On the use of windows for harmonic analysis with the discrete Fourier transform. *Proc IEEE* 66(1):51–93
- Jekeli C (1981) Alternative methods to smooth the Earth's gravity field. Tech. Rep. 327, Department of Geodetic Science and Surveying, The Ohio State University
- Kusche J (2007) Approximate decorrelation and non-isotropic smoothing of time-variable GRACE-type gravity field models. *J Geod* 81:733–749
- Laprise R (1992) The resolution of global spectral models. *Bull Am Meteorol Soc* 73(9):1453–1454
- Lillesand TM, Kiefer RW (1994) *Remote sensing and image interpretation*, 3rd edn. Wiley, New York
- Longuevergne L, Scanlon BR, Wilson CR (2010) GRACE hydrological estimates for small basins: evaluating processing approaches on the High Plains Aquifer, USA. *Water Resour Res*. doi:10.1029/2009WR008564
- Rummel R, Yi W, Stummer C (2011) GOCE gravitational gradiometry. *J Geod* 85(11):777–790
- Sardeshmukh PD, Hoskins BJ (1984) Spatial smoothing on the sphere. *Mon Weather Rev* 112:2524–2529
- Swenson S, Wahr J (2006) Post-processing removal of correlated errors in GRACE data. *Geophys Res Lett*. doi: 10.1029/2005GL025285
- Tapley BD, Bettadpur S, Watkins MM, Reigber C (2004) The gravity recovery and climate experiment: mission overview and early results. *Geophys Res Lett*. doi:10.1029/2004GL019920
- Werth S, Güntner A, Schmidt R, Kusche J (2009) Evaluation of GRACE filter tools from a hydrological perspective. *Geophys J Int* 179(3):1499–1515

---

# On Time-Variable Seasonal Signals: Comparison of SSA and Kalman Filtering Based Approach

Q. Chen, M. Weigelt, N. Sneeuw, and T. van Dam

---

## Abstract

Seasonal signals (annual and semi-annual) in GPS time series are of great importance for understanding the evolution of regional mass, e.g. ice and hydrology. Conventionally, these signals are derived by least-squares fitting of harmonic terms with a constant amplitude and phase. In reality, however, such seasonal signals are modulated, i.e., they have time-variable amplitudes and phases. Davis et al. (J Geophys Res 117(B1):B01,403, 2012) used a Kalman filtering (KF) based approach to investigate seasonal behavior of geodetic time series. Singular spectrum analysis (SSA) is a data-driven method that also allows to derive time-variable periodic signals from the GPS time series. In Chen et al. (J Geodyn 72:25–35, 2013), we compared time-varying seasonal signals obtained from SSA and KF for two GPS stations and received comparable results. In this paper, we apply SSA to a global set of 79 GPS stations and further confirm that SSA is a viable tool for deriving time variable periodic signals from the GPS time series. Moreover, we compare the SSA-derived periodic signals with the seasonal signals from KF with two different input process noise variances. Through the comparison, we find both SSA and KF obtain promising results from the stations with strong seasonal signals. While for the stations dominated by the long-term variations, SSA seems to be superior. We also find that KF with input process noise variance based on variance rates performs better than KF with the input process noise variance based on simulations.

---

## Keywords

Kalman filtering • Singular spectrum analysis • Time variable seasonal signals

---

## 1 Introduction

Over the last few decades, the Global Positioning System (GPS) has demonstrated its capability for monitoring deformations of the Earth's surface. Seasonal signals in GPS

position time series, which are well known to result from surface mass loading (e.g., Dong et al. 2002), are of great value for studying the evolution of surface mass cycles, e.g., the hydrological cycle.

Conventionally, seasonal signals are retrieved with a linear model via least squares fitting prior to or simultaneous with some noise assumption (e.g., Williams et al. 2004), which results in constant amplitudes and phases. In reality, seasonal variations are not constant from year to year, neither in amplitude nor in phase. Several studies have suggested determining the time-varying periodic signals by relying, for instance, on KF based techniques (Murray and Segall 2005; Davis et al. 2012), on piecewise continuous linear polynomials (Davis et al. 2006), on a flexible

---

Q. Chen (✉) • N. Sneeuw  
Institute of Geodesy, University of Stuttgart, Geschwister-Scholl-Str.  
24D, 70174 Stuttgart, Germany  
e-mail: [qiang.chen@gis.uni-stuttgart.de](mailto:qiang.chen@gis.uni-stuttgart.de)

M. Weigelt • T. van Dam  
Faculté des Science, de la Technologie et de la Communication,  
University of Luxembourg, 6 rue Richard Coudenhove-Kalergi,  
L-1359, Luxembourg

semi-parametric model (Bennett 2008), on non-parametric annual signal (Freymueller 2009; Tesmer et al. 2009), or on singular spectrum analysis (SSA) (Chen et al. 2013).

In line with Davis et al. (2012) and Chen et al. (2013) in this paper, we further investigate the abilities of singular spectrum analysis and the KF based approach to extract time-varying periodic signals from the GPS time series. To fulfill the goal, we first apply SSA to the time series from 79 global IGS stations with more than 11 years of weekly measurements. To obtain a better understanding of the KF approach, we adopt two scenarios to handle the input process noise variances. One is directly based on variance rates from Davis et al. (2012), the other is to simulate the random walk process variances and implement an ensemble KF approach (Weigelt et al. 2013). It should be noted that those two scenarios, in principle, follow Davis et al. (2012).

As SSA has its prominent advantages, e.g., model independence, and disadvantages, e.g., time-consumption, another purpose of this work is trying to investigate the capability of KF with different settings in separating time-variable periodic signals. However, true time variable seasonal signals buried in the GPS time series are unknown. We thus tentatively take the SSA-derived results as reference in comparison to the results from the two KF scenarios, which helps us to understand the application of KF to the GPS time series analysis.

This paper is organized as follows: in Sect. 2 we shortly outline the methodologies of SSA and KF. This is followed by a brief discussion of the strengths and weaknesses of both approaches in this section. Section 3 is the data analysis part, which demonstrates the performance of SSA. In this section, we also implement a comparison between SSA and two KF scenarios. Finally, we draw a conclusion in Sect. 4.

## 2 Methodology

### 2.1 Singular Spectrum Analysis

Following the description of SSA from Broomhead and King (1986) and Vautard and Ghil (1989), the main procedure of the technique can be summarised as follows:

1. Given a centered time series  $x_t$  ( $1 \leq t \leq N$ ), the first step is to construct a covariance matrix. We follow the VG algorithm, proposed by Vautard and Ghil (1989), to compute the covariance matrix, which is based on the lagged-covariance matrix of the process  $x_t$ . With a maximum lag (or window size),  $M$ , the matrix  $C_{VG}$  has a Toeplitz structure, i.e., constant diagonals corresponding to equal lags:

$$C_{VG} = \begin{pmatrix} c_0 & c_1 & c_2 & \cdots & c_{M-1} \\ c_1 & c_0 & c_1 & \cdots & \cdot \\ c_2 & c_1 & c_0 & \cdots & \cdot \\ \cdot & c_2 & c_1 & c_0 & \cdots & \cdot \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ \cdot & \cdot & \cdot & \cdot & \cdot & c_1 \\ c_{M-1} & \cdot & \cdot & \cdot & c_1 & c_0 \end{pmatrix}, \quad (1)$$

where entries  $c_j$ ,  $0 \leq j \leq M-1$ , are the covariance of  $x$  at lag  $j$ . Its unbiased estimates are:

$$c_j = \frac{1}{N-j} \sum_{i=1}^{N-j} x_i x_{i+j}, \quad 0 \leq j \leq M-1 \quad (2)$$

2. We apply eigenvalue decomposition to  $C$  in order to obtain the eigenvalues,  $\lambda_k$ , and eigenvectors (also called EOFs),  $E^k$ , of this matrix. These are then sorted in descending order of  $\lambda_k$ , where index  $k = 1, 2, \dots, M$ . The  $k^{\text{th}}$  principal component (PC) is

$$a_i^k = \sum_{j=1}^M x_{i+j} E_j^k, \quad 0 \leq i \leq N-M \quad (3)$$

3. We reconstruct each component of the original time series as given by Vautard et al. (1992)

$$x_i^k = \begin{cases} \frac{1}{i} \sum_{j=1}^i a_{i-j}^k E_j^k & 1 \leq i \leq M-1 \\ \frac{1}{M} \sum_{j=1}^M a_{i-j}^k E_j^k & M \leq i \leq N-M+1 \\ \frac{1}{N-i+1} \sum_{j=i-N+M}^M a_{i-j}^k E_j^k & N-M+2 \leq i \leq N. \end{cases} \quad (4)$$

4. According to Plaut and Vautard (1994), harmonic oscillations can be identified in terms of the three fundamental properties: (1) two consecutive eigenvalues are nearly equal; (2) the two corresponding time sequences described by EOFs are nearly periodic, with the same period and in quadrature; (3) the associated PCs are in quadrature.

In the practical implementation of the SSA algorithm, the choice of the lag-window size  $M$  is of great importance. Generally, the lag-window size  $M$  depends not only on the length of the data but also on the desired periodic cycles. Both empirical and mathematical rules exist in the literature

to suggest how to make a proper choice of the lag-window size  $M$ . When dealing with GPS time series, Chen et al. (2013) demonstrated that a window size of 2 or 3 years, that is, 105 or 157 weeks of weekly GPS time series, is appropriate. Attention is required since the standard SSA technique can only handle evenly sampled data. For dealing with GPS time series with gaps, we employed a modified SSA algorithm which was suggested by Schoellhamer (2001). In this modified SSA algorithm, a parameter  $f$  ( $0 \leq f \leq 1$ ) is introduced that represents a specified fraction of allowable missing data points within a given window size  $M$ . A more detailed description about the application of SSA technique to the GPS time series analysis is referred to (Chen et al. 2013).

## 2.2 Kalman Filtering

We follow the concept described in Davis et al. (2012) and implement the model (see Eq. (5)) that includes stochastic annual and semiannual terms using linear Kalman filtering.

$$\begin{aligned} x(t) = & x_0 + v(t)(t - t_0) \\ & + a'_1(t) \cos(2\pi f_0(t - t_0)) + b'_1(t) \sin(2\pi f_0(t - t_0)) \\ & + a'_2(t) \cos(4\pi f_0(t - t_0)) + b'_2(t) \sin(4\pi f_0(t - t_0)), \end{aligned} \quad (5)$$

where  $f_0 = 1$  cpy and  $t_0$  is a reference epoch,  $v(t)$  is time-variable velocity term and  $a'_1(t)$ ,  $a'_2(t)$ ,  $b'_1(t)$  and  $b'_2(t)$  are instantaneous time-variable amplitudes.

We apply the same dynamic process model (random constant for  $x_0$  and random walk for stochastic terms) but we adopt two different strategies to handle the input process variance. One scenario is the simple case which is based on the variance rates (variance rate values of  $1 \text{ mm}^2 \text{ year}^{-3}$  for the rate term and  $0.5 \text{ mm}^2 \text{ year}^{-1}$  for sinusoidal amplitudes) given by Davis et al. (2012) and multiplied by the time interval between two epochs. We call this scenario as KF 1 through this paper.

The other scenario of dealing with input process noise variance is to simulate the random walk process noise, normalize and scale them onto corresponding GPS noise level (Weigelt et al. 2013). In practice, this is done by integrating a random sequence drawn from a Gaussian distribution. In order to minimize a possible dependency on the used sequence of random numbers, an ensemble approach is employed, i.e. the estimation is repeated using different random sequences of the process noise. We obtain an ensemble of solutions which are averaged. We call the second scenario as KF 2 likewise through the paper. It should be noted that a Rauch–Tung–Striebel smoother (Rauch et al. 1965) is employed during the running of KF.

## 2.3 Strengths and Weaknesses of SSA and KF

Both SSA and KF have their own merits. One big advantage of SSA is model independence so that SSA is free of any prior model or noise assumptions. This big merit leads to its drawback at the same time: it does not produce uncertainties of the result. On the contrary, the KF approach is based not only on the functional model (the observation model (Eq. (5)) and the state transition model) but also on the stochastic model (observation noise and process noise). Together both models control the output estimates and produce uncertainties. Model dependence is the big drawback of KF because the assumed functional model and stochastic model will force the estimation process to follow prescribed behaviors. Inaccurate models will produce erroneous estimates.

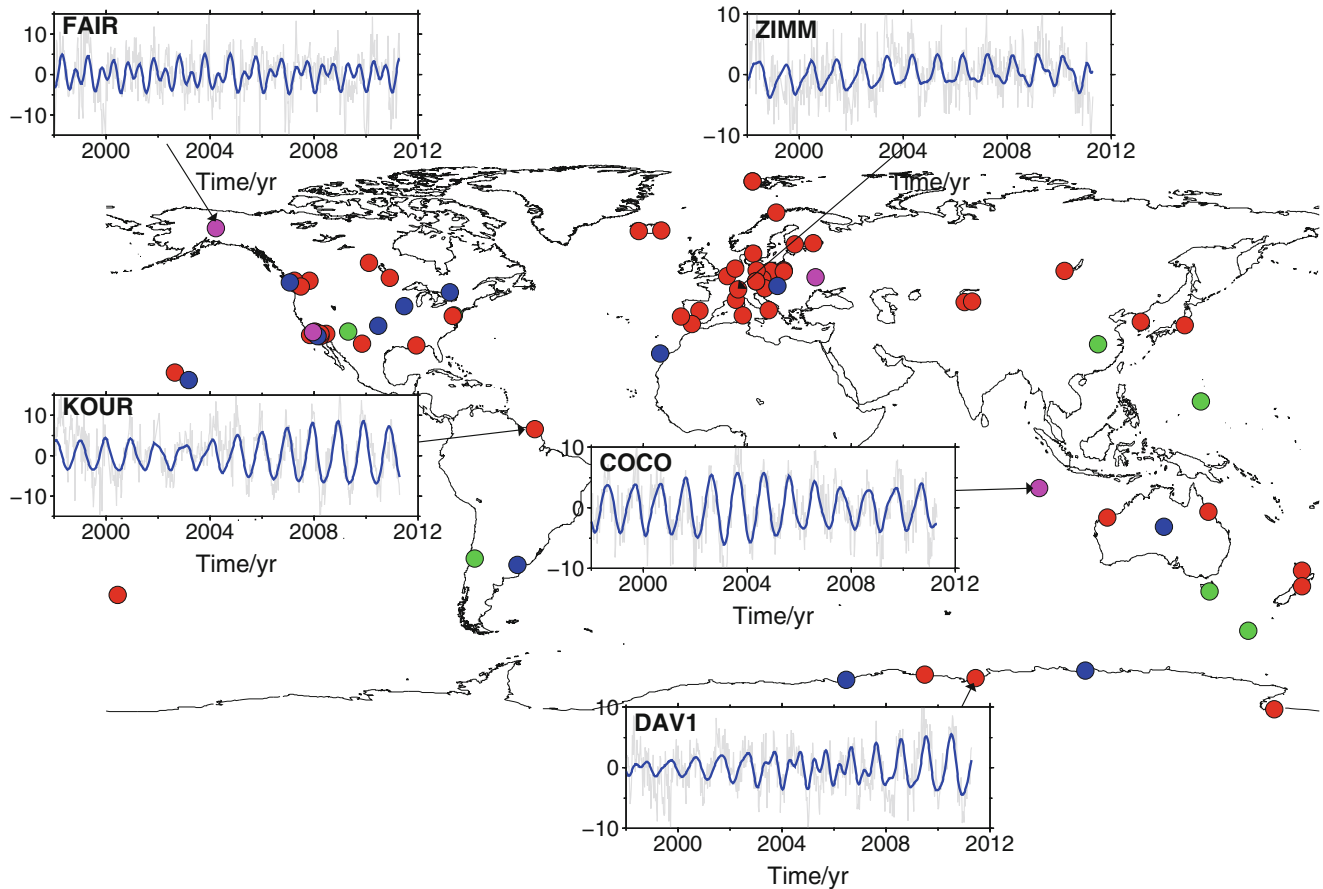
Another drawback of applying SSA in the GPS time series analysis is its time-consumption. Due to the complicated composition of the GPS time series, an empirical lag-window size  $M$  works for one GPS time series while it might not work for another. A trial and error of the window size makes the analysis process time-consuming. As a result, SSA, to a certain extent, is not a globally efficient analysis tool in comparison with the KF based approach. This motivated us to look into the model-based approaches, especially the KF based approach. On the basis of Murray and Segall (2005) and Davis et al. (2012), we search for a better understanding of functional models and stochastic models of KF. To be specific, within the framework of Davis et al. (2012), we compare two process noise models in this work.

## 3 Data Analysis

As an extension of (Chen et al. 2013), we follow its settings regarding SSA and will not repeat the details. We select 79 height time series with more than 10 years of weekly observations from IGS stations over the globe. We use the weekly height GPS time series from Collilieux et al. (2012), which preserves crustal deformation information in the coordinate time series, especially at the seasonal timescale. Figure 1 shows the station distribution.

We analyze 79 height time series with SSA and separate modulated annual and semiannual signals from original GPS time series. During the analysis, a window size of 157 weeks (3-year window) is most frequently assigned (55 stations). For the remaining stations, a 2-year window (8 stations), a 4-year (12 stations) and a 5-year window (4 stations) are applied to accommodate the specific GPS time series. To some extent, it validates the description in Chen et al. (2013) about the choices of window-size (a 2- or 3-year window should be appropriate).





**Fig. 1** Distribution of the selected 79 GPS stations and five examples of SSA analysis. *Gray lines* in each subplot are original time series and *blue curves* are seasonal signals extracted by SSA. Different colored

*dots* represent the GPS stations with different window size in the course of SSA analysis. *Green, red, blue, magenta dots* stand for the GPS stations with a window size of 105, 157, 209 and 261 weeks respectively

Annual signals appear frequently in the first and second components due to the significant annual signal in GPS height time series. We choose five examples to show the performances of SSA, see Fig. 1. The five subplots clearly demonstrate that seasonal signals derived by SSA follow the original time series quite well. KOUR and COCO contain strong annual signals. DAV1 and ZIMM show both annual and semi-annual signals. While for FAIR, the semi-annual signal dominates the whole time series.

To assess the performance of the KF based approach, we utilize the two scenarios described in Sect. 2.2. Due to the unknown true time variable seasonal signals, we tentatively take the results from SSA as the reference, i.e. ‘true’ signal, and compare with the results from KF with two process noise scenarios. Note that assuming the SSA-derived results as reference is an assumption that we can not yet fully verify. The correlation and RMS of difference between SSA and the KF based technique are employed to evaluate the performance.

Tables 1 and 2 show the statistical results of the comparison. At the annual signal level, both scenarios (KF 1 and KF 2) achieve a high correlation (more than 0.9) and

**Table 1** Comparison of correlation between SSA and two KF scenarios

	Max	Min	Mean
Only annual			
SSA vs KF 1	0.99	0.74	0.96
SSA vs KF 2	0.99	0.48	0.91
Only semi-annual			
SSA vs KF 1	0.98	0.19	0.84
SSA vs KF 2	0.91	0.19	0.67
Annual+semi-annual			
SSA vs KF 1	0.99	0.79	0.95
SSA vs KF 2	0.98	0.64	0.89

relatively good RMS values (less than 1 mm), which results from the strong annual signal buried in the GPS height time series. At this level, KF 1 performs better than KF 2. At the semi-annual signal level, correlations decrease while they are still acceptable. RMS values do not change much because of weak semi-annual signal strength. At the combination level (annual plus semi-annual), correlations are high in both scenarios. As for RMS values, with the signal strength

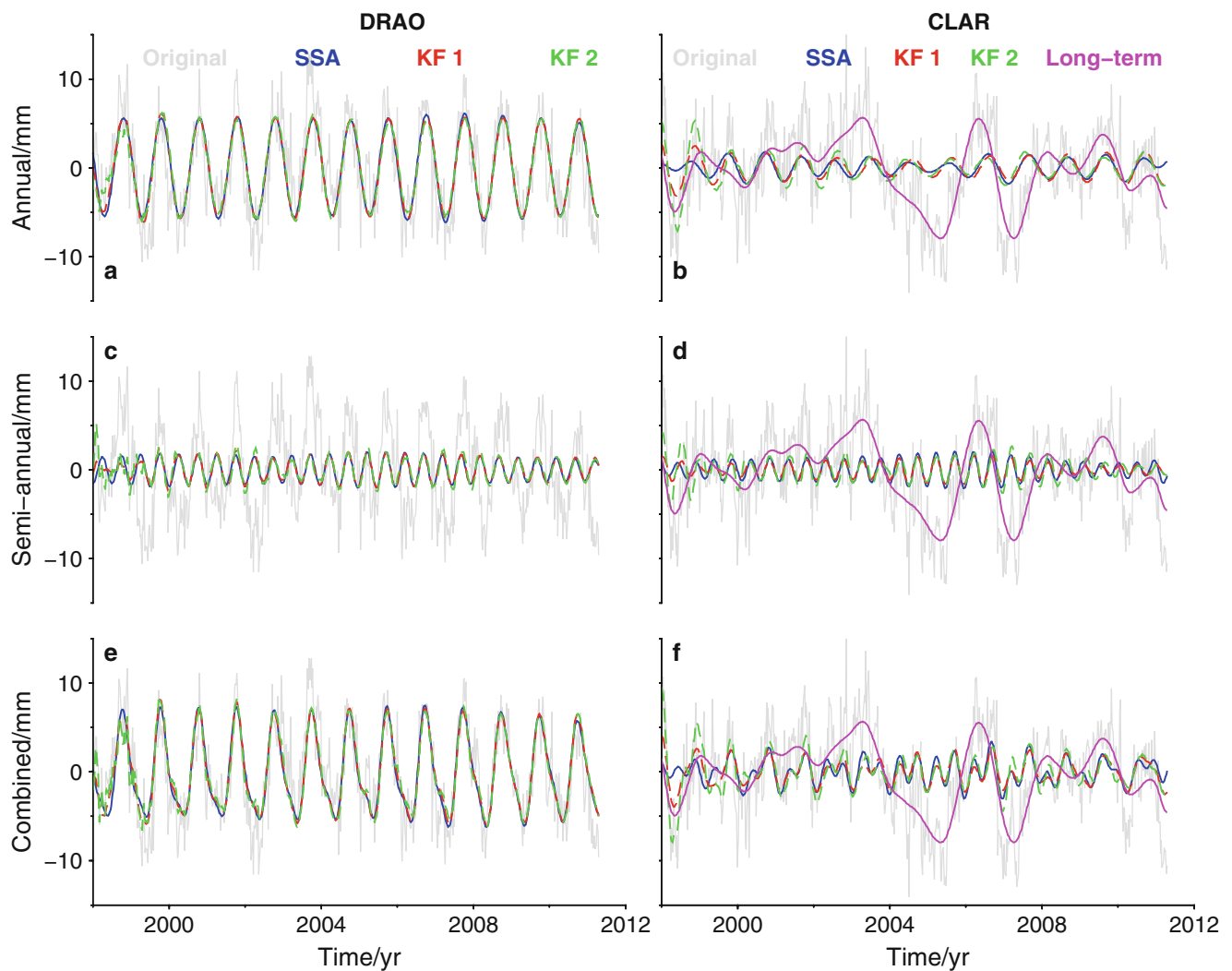
**Table 2** Comparison of RMS values (mm) between the difference of SSA and two KF scenarios

	Max	Min	Mean
<b>Only annual</b>			
SSA vs KF 1	2.0	0.3	0.6
SSA vs KF 2	3.4	0.4	0.9
<b>Only semi-annual</b>			
SSA vs KF 1	3.4	0.2	0.5
SSA vs KF 2	2.7	0.4	0.9
<b>Annual+semi-annual</b>			
SSA vs KF 1	4.3	0.4	0.8
SSA vs KF 2	8.3	1.0	2.3

increasing, KF 1 could still receive a mean RMS less than 1 mm. However, the mean RMS value of KF 2 increases to

2.3 mm. As the input process noise variances balance the predicted and observed values during the running of KF, we alter the scale factor in KF 2 by 10 times and 100 times larger or smaller, which means to enlarge or reduce the input process noise variances, to investigate its sensitivities. No matter how we alter the scale factor, the performance of KF 2 is still inferior to KF 1. It may indicate that the method of KF 2 to create input process noise variances is not suitable for this situation. One possible reason is that the KF 2 does not consider the time interval in the GPS time series when generating the noise process. In Weigelt et al. (2013), the CHAMP data they used were equally sampled while data gaps exist in the GPS time series we used.

Figure 2 show two examples of the comparison at annual, semi-annual and seasonal levels. DRAO, which is located in



**Fig. 2** Comparison of the extracted annual signals (subplot a and b), semi-annual signals (subplot c and d) and their combinations (subplot e and f) in DRAO and CLAR, respectively. The long-term variation presented in CLAR is extracted by SSA. It is interesting to find in the subplot (b) that the signal that KF obtains in an annual signature at the

beginning and the end of the CLAR coordinate time series is modeled as the long-term variation by SSA. Given the amplitude of the variabilities of the long-term and annual signals produced by SSA, this would seem to be a more reasonable partitioning of the signal, but we have no data with which to test this assumption

Penticton, Canada, has very strong seasonal signals while long-term variation dominates CLAR, which is located in Claremont, United States. In cases like DRAO, both KF scenarios obtain very promising results. Nevertheless, in cases like CLAR, the KF based technique performs poorly probably due to the insignificant annual and semiannual signals, or the inaccurate observation model we used in KF.

The solid lines in magenta in the subplots Fig. 2b, d, f are the long-term variations in CLAR which is extracted by SSA. It is interesting to find that the signal that KF obtains in an annual signature at the beginning and the end of the CLAR time series is modeled as the long-term variation by SSA. Given the amplitude of the variabilities of the long-term and annual signals produced by SSA, this would seem to be a more reasonable partitioning of the signal, but we have no data with which to test this assumption. This issue will be followed up in future work.

## 4 Conclusion

In this study, we confirm the capability of singular spectrum analysis for modeling the time variable periodic signals over 79 GPS stations. As a data-driven approach, SSA is capable of extracting amplitude and phase modulated seasonal signals from the GPS time series. In addition, we also demonstrate two ways of handling the input process noise variances in the KF process. The results show that both KF scenarios work well for the stations that have strong seasonal signals. Nevertheless, for stations which are dominated by long-term variations, SSA seems to be superior.

In terms of a comparison of two KF scenarios over 79 GPS stations, we conclude that KF 1 which is based on the variance rates (Davis et al. 2012) outperforms KF 2 that is based on the simulation of the process noise variances (Weigelt et al. 2013). It indicates that the Weigelt et al. (2013) way of dealing with input noise variances is not an optimal choice for the case of GPS time series analysis.

**Acknowledgements** We acknowledge the International GNSS Service (IGS), especially Xavier Collilieux (IGN, France), for providing the original GPS coordinate time series. We appreciate J. L. Davis and two anonymous reviewers for their valuable comments. Qiang Chen acknowledges the Chinese Scholarship Council for supporting his PhD study.

## References

Bennett RA (2008) Instantaneous deformation from continuous GPS: contributions from quasi-periodic loads. *Geophys J Int* 174(3):1052–1064. doi:10.1111/j.1365-246X.2008.03846.x

- Broomhead D, King GP (1986) Extracting qualitative dynamics from experimental data. *Physica D* 20:217–236. doi:10.1016/0167-2789(86)90031-X
- Chen Q, van Dam T, Sneeuw N, Collilieux X, Weigelt M, Rebischung P (2013) Singular spectrum analysis for modeling seasonal signals from GPS time series. *J Geodyn* 72:25–35. doi:10.1016/j.jog.2013.05.005
- Collilieux X, van Dam T, Ray J, Coulot D, Métivier L, Altamimi Z (2012) Strategies to mitigate aliasing of loading signals while estimating gps frame parameters. *J Geod* 86:1–14. 10.1007/s00190-011-0487-6
- Davis JL, Wernicke BP, Bisnath S, Niemi NA, Elosegui P (2006) Subcontinental-scale crustal velocity changes along the Pacific-North America plate boundary. *Nature* 441(7097):1131–1134. doi:10.1038/nature04781
- Davis JL, Wernicke BP, Tamisiea ME (2012) On seasonal signals in geodetic time series. *J Geophys Res* 117(B1):B01403. doi:10.1029/2011JB008690
- Dong D, Fang P, Bock Y, Cheng MK, Miyazaki S (2002) Anatomy of apparent seasonal variations from GPS-derived site position time series. *J Geophys Res* 107(B4):2075. doi:10.1029/2001JB000573
- Freyemueller J (2009) Seasonal position variations and regional reference frame realization. In: Drewes H (ed) *Geodetic reference frames, IAG symp*, vol 134, pp 191–196. doi:10.1007/978-3-642-00860-3\_30
- Murray JR, Segall P (2005) Spatiotemporal evolution of a transient slip event on the San Andreas fault near Parkfield, California. *J Geophys Res* 110(B9):B09407. doi:10.1029/2005JB003651
- Plaut G, Vautard R (1994) Spells of low-frequency oscillations and weather regimes in the northern hemisphere. *J Atmos Sci* 51(2):210–236. doi:10.1175/1520-0469(1994)051<0210:SOLFOA>2.0.CO;2
- Rauch HE, Striebel C, Tung F (1965) Maximum likelihood estimates of linear dynamic systems. *AIAA J* 3(8):1445–1450. doi:10.2514/3.3166
- Schoellhamer DH (2001) Singular spectrum analysis for time series with missing data. *Geophys Res Lett* 28(16):3187–3190. doi:10.1029/2000GL012698
- Tesmer V, Steigenberger P, Rothacher M, Boehm J, Meisel B (2009) Annual deformation signals from homogeneously reprocessed VLBI and GPS height time series. *J Geod* 83:973–988. doi:10.1007/s00190-009-0316-3
- Vautard R, Ghil M (1989) Singular spectrum analysis in nonlinear dynamics, with applications to paleoclimatic time series. *Physica D* 35(3):395–424. doi:10.1016/0167-2789(89)90077-8
- Vautard R, Yiou P, Ghil M (1992) Singular-spectrum analysis: a toolkit for short, noisy chaotic signals. *Physica D* 58(1-4):95–126. doi:10.1016/0167-2789(92)90103-T
- Weigelt M, van Dam T, Jäggi A, Prange L, Tourian MJ, Keller W, Sneeuw N (2013) Time-variable gravity signal in Greenland revealed by high-low satellite-to-satellite tracking. *J Geophys Res* 118:3848–3549. doi:10.1002/jgrb.50283
- Williams SDP, Bock Y, Fang P, Jamason P, Nikolaidis RM, Prawirodirdjo L, Miller M, Johnson DJ (2004) Error analysis of continuous GPS position time series. *J Geophys Res* 109(B3):B03412. doi:10.1029/2003JB002741

---

# Extensive Analysis of IGS REPRO1 Coordinate Time Series

M. Roggero

---

## Abstract

The work describes the analysis conducted on the IGS REPRO1 coordinate time series, in order to detect GNSS permanent stations periodic behavior. Frequency analysis requires cyclostationary time series, while observed coordinates time series are not cyclostationary because of discontinuities of different kind and origin and of long term linear or non linear trend. For this reason time series offsets and trends must be estimated and eliminated, prior to conduct the harmonic analysis.

Discontinuities are usually documented by IGS, but undocumented discontinuities also exists and need to be detected. The long term component of the signal is generally modeled as a linear trend, but the linear model is often inadequate to obtain cyclostationary residuals. An alternative model based on a discrete time Markov process will be adopted.

The study has been conducted on the up component of the REPRO1 raw coordinates time series. No correction for the atmospheric pressure loading has been applied. Harmonic analysis has been performed using the non linear least square algorithm implemented by F. Mignard in the Frequency Analysis Mapping On Unusual Sampling software (Mignard, FAMOUS, Frequency Analysis Mapping on Unusual Sampling, (OCA Cassiopee), 2003).

We obtained a complete statistic on the vertical component period, amplitude and phase. Signals at from 1 to 7 cycle per solar and draconitic year can be observed in most stations as expected, but also other signals have been detected that can be attributed to tidal model errors. Some interpretation will be given referring to recent literature.

---

## Keywords

Periodic signals • REPRO1 • Time series

---

## 1 Introduction and Motivation

Recently IGS released a first full reanalysis of all GPS data collected since 1994, the REPRO1 solution, based on the weekly SINEX solutions, from GPS week 729 (01/01/1994) to GPS week 1631 (04/16/2011). We performed a retrospective analysis of the REPRO1 coordinate time series, focusing at first on the detection, estimation, and elimination of time

series offsets, than on the long term model estimation and finally on the harmonic analysis of model residuals.

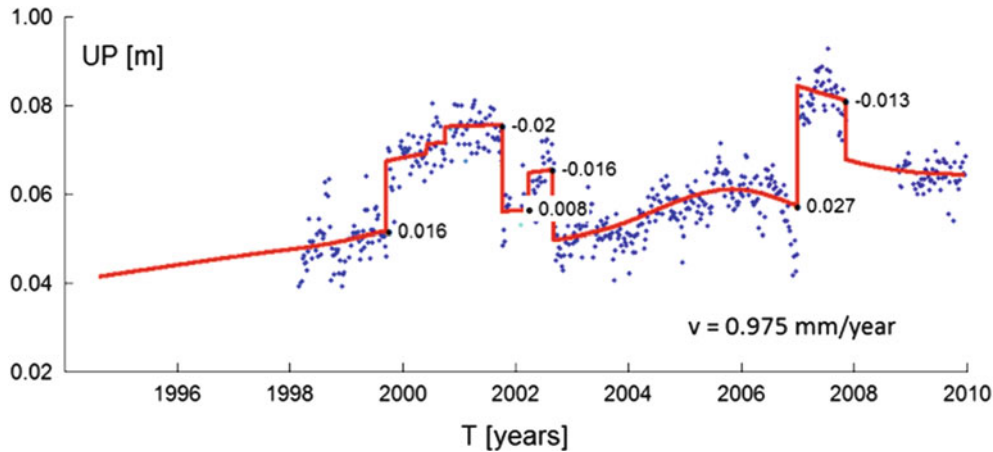
Published harmonic analysis of GPS coordinates time series have shown significant variation in the respective spectrum. The most recent studies include (Ray et al. 2008; Collilieux et al. 2007, 2010; Fritsche et al. 2009; Mtamakaya 2012). A comprehensive analysis of the whole REPRO1 data set,<sup>1</sup> providing a significant and self consistent sample, can help to understand the impact of detected signals at a global scale. Moreover regional dependent spectral signatures can be observed, even if the analysis suffers of a lack of data in

---

M. Roggero (✉)  
Politecnico di Torino, Viale Mattioli 39, 10125 Torino, Italy  
e-mail: [marco.roggero@polito.it](mailto:marco.roggero@polito.it)

---

<sup>1</sup><ftp://cddis.gsfc.nasa.gov/pub/gps/products/repro1>.



**Fig. 1** BAKO UP time series with the estimated jumps

the southern hemisphere and in the polar areas. Coordinate and formal error time series of 526 IGS stations have been extracted from REPRO1 SINEX ( $\Delta E$ ,  $\Delta N$ ,  $\Delta h$  and  $\sigma_N, \sigma_E, \sigma_h$ ) and in the present work we examined the up component ( $\Delta h$ ,  $\sigma_h$ ). Time series discontinuities are documented by IGS in the SOLN.SNX file, that reports 1206 documented discontinuities in position and velocity, over 338 stations (a mean of 3 or 4 discontinuities per station). The cause of about 25% of the discontinuities is unknown, while it is not possible to exclude the presence of other undocumented discontinuities. The reprocessing of IGS had been carried out by using fully consistent models in order to avoid model dependant discontinuities and coordinates variance anisotropy that occur in operational time series. However, the discontinuities due to hardware or monument change or to geophysical effects remain also in the reprocessed time series as noted in Steigenberger (2009). An example of discontinuity estimation is given in the BAKO station UP time series of Fig. 1, that presents eight different position discontinuities reported by IGS.

The changes in velocity, usually caused by earthquakes or by other geophysical effects, can be described by multi linear models, as shown in Perfetti (2006) and Ostini (2012). In these works the coordinate time series are represented as the sum of a long term linear or multi linear trend and a step function, taking in account for position offsets, and cyclical components. Harmonic analysis requires to be applied on cyclostationary residuals, having statistical properties that vary cyclically with time, however linear and also multi linear models seem to be often inadequate to describe the long term behavior of a station in order to obtain cyclostationarity, as it will be shown in Sect. 2.3.

To overcome the inadequacy of multi linear models applied to coordinate time series, the presented approach

is based on a discrete time Markov process modeling and focuses on three steps:

1. Detect, estimate and remove the level shifts, performing iteratively the so called detection, identification and adaptation procedure (DIA) as explained in Roggero (2012).
2. Model the long term signal constraining the system dynamic, in order to obtain cyclostationary residuals.
3. Residuals harmonic analysis by the non linear least square algorithm implemented in the FAMOUS software (Mignard 2003).

The step 1 of the proposed algorithm has been tested on synthetic data in the framework of the DOGEx project (Gazeaux et al. 2013). The signals estimated by FAMOUS will be analyzed in frequency, amplitude and phase, stacking the power spectra in order to detect the most significant effects. Some preliminary consideration will be given in Sect. 4. The software FAMOUS has been used also in Collilieux et al. (2010) to analyze the ITRF2008, that is based on entirely reprocessed GPS solutions from 1997 to 2008.

## 2 Time Series Modeling

### 2.1 Discrete Time Linear Model

GNSS time series can be modeled as discrete-time Markov process. Consider a discrete-time linear system described by a finite state vector  $x$ , evolving with known dynamics  $T$  through the epochs  $t$  ( $\in t [1, n]$ ), with system noise  $v$  (with variance-covariance matrix  $R_v$ ):

$$\begin{aligned} x_{t+1} &= T_{t+1}x_t + v_{t+1} \\ y_{t+1} &= H_{t+1}x_{t+1} + \varepsilon_{t+1} \end{aligned} \quad (1)$$

The observations  $y$  are known with observation noise  $\varepsilon$  (with variance-covariance matrix  $R_{\varepsilon\varepsilon}$ ). It has been shown by Albertella et al. (2006) that the system has the optimal solution

$$\hat{x} = (D^T W_\omega D + M^T W_\varepsilon M)^{-1} M^T W_\varepsilon y \quad (2)$$

where  $\hat{x}$  is the estimated state vector,  $D$  and  $M$  are block diagonal matrices representing respectively  $T$  and  $H$  over the considered time interval, while  $W_\omega$  and  $W_\varepsilon$  are weight matrices. The 3D state vector  $\hat{x}$  is estimated constraining the system dynamic, by setting the system noise  $\nu$  at a given value. The process is detailed in Roggero (2008, 2012).

For a system with slow dynamic as GNSS coordinate time series, the motion can be described by a constant velocity model in  $T$ , with acceleration  $\ddot{p} = 0$

$$\begin{aligned} p_{t+1} &= p_t + \dot{p}_t \cdot \delta t \\ \dot{p}_{t+1} &= \dot{p}_t \end{aligned} \quad T = \begin{bmatrix} 1 & \delta t \\ & 1 \end{bmatrix} \quad (3)$$

where the position  $p$  and the velocity  $\dot{p}$  are the two elements of the state vector  $x = [p \ \dot{p}]$ , with system noise  $\nu = [v_p \ v_{\dot{p}}]$ . The approach is equivalent to Kalman filtering and smoothing, but allows to manage the estimation of constant biases more efficiently, as will be shown in Sect. 2.2.

In state estimation the outliers are not rejected but properly weighted according to the system and observation noise.

## 2.2 Discontinuity Model

Discontinuities has been detected, estimated and removed applying the detection, identification and adaptation procedure (DIA) presented in Teunissen (1998), as applied in Perfetti (2006) and in Roggero (2012). Taking in account for discontinuities requires to modify the model (1). The bias vector  $b$  represents the time series offsets and modifies the system as follows:

$$\begin{aligned} x_{t+1} &= T_{t+1}x_t + B_{t+1}b_t + v_{t+1} \\ y_{t+1} &= H_{t+1}x_{t+1} + C_{t+1}b_t + \varepsilon_{t+1} \\ b_{t+1} &= b_t \end{aligned} \quad (4)$$

The bias vector  $b$  is constant with steps, and it is linked to the system dynamic and to the observations by the matrices  $B$  and  $C$ . The matrix  $C$ , whose elements are 0 or 1, represents the occurrence of the biases in the time series. The number of rows is equal to the number of observation epochs, while the number of columns is equal to the unknown number of jumps to be estimated. The matrix  $B$  it is assumed equal to zero if the bias affects only the observed position and not the real position. However it can be different by zero in the case of seismic displacements. These matrices can be

known a priori in the case of documented discontinuities, or determined by means of some detection criteria for undocumented discontinuities. The estimation of the extended state vector  $z = [x \ b]$  requires the inversion of a large sparse normal matrix. This matrix has a bordered block or band-diagonal structure (quasi-triangular Schur form), so it can thus be block wisely inverted by using Shur decomposition as in Roggero (2008).

The offsets detection is based on a hypothesis test which assumes as null hypothesis  $H_0$  that the time series do not have any offset. This hypothesis is tested against a certain number of alternative hypotheses  $H_A$ , with a jump in a given epoch. An alternative hypothesis can be formulated for each observation epoch or for candidate epochs only. The adequacy of the model can be verified using the ratio test, which is known to have the  $\chi^2$  distribution. After detecting the offsets, they can be estimated and removed.

Because offsets do not necessarily affect horizontal and vertical components similarly, the vertical component is studied separately using the same approach. This approach also makes it possible to consider documented and undocumented offsets, to predict the station coordinates in data gaps, and to correctly represent pre-seismic and post-seismic deformations or other non-linear behaviors.

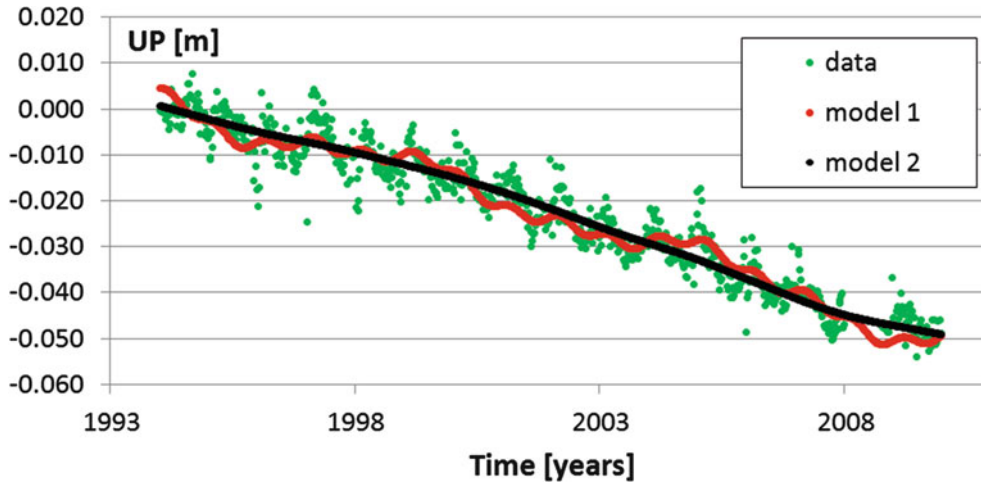
The frequency analysis shows that not only the coordinate time series present discontinuities in position and velocity, but also in their characteristic frequency signature. This last kind of discontinuity can be site dependant and in this case seems to be related to hardware change or to other site effects to be investigated. Model dependant frequency discontinuities can also exist, but have been avoided in REPRO1 solution by adopting models and methods fully consistent during the time.

## 2.3 Long Term Signals

In the long term signals we include the linear trend, the non linear and non periodic signals, and also the periodic signals with period larger than the time series length, that therefore are not estimable by harmonic analysis. From this point of view linear trend is only one component of the long term signal, often the larger one, but to obtain cyclostationary model residuals we can't neglect the non linear long term signals.

As we have seen in Sect. 2.1, the time series are modeled as a Markov process, and the system dynamic is described by a constant velocity model. A way to estimate the long term signal is by decreasing the variance  $\sigma_\nu$  of the system noise  $\nu$ ; in other words, the system noise is related to the maximum frequency of the estimated signal.

Because of the system noise depends on the system dynamic, the model can follow different dynamics by setting



**Fig. 2** Two different variance model used to estimate the long term behavior of the ALGO station in the UP component: **Model 1**: short wavelength,  $\sigma_p = 10^{-6}$  m and  $\sigma_{\dot{p}} = 10^{-5}$  m; **Model 2**: long wavelength,  $\sigma_p = 10^{-6}$  m and  $\sigma_{\dot{p}} = 10^{-6}$  m

different values for the system noise variance  $\sigma_v$ , that has been done empirically. Figure 2 shows the UP component time series of the ALGO station. Two different model have been estimated: model 1, with  $\sigma_p = 10^{-6}$  m and  $\sigma_{\dot{p}} = 10^{-5}$  m, that follows the short term signal, and model 2 with  $\sigma_p = 10^{-6}$  m and  $\sigma_{\dot{p}} = 10^{-6}$  m, that follows the long term signal.

The choice of the system variance noise is critical and depends on the sampling frequency and on the system noise. It has been fixed empirically on a subset of ten stations, randomly chosen. However the algorithm is insensitive to quite large variations of this parameter.

### 3 Harmonic Analysis

Harmonic analysis leads to the representation of the signal as a superposition of basic waves. A variety of different approaches are presently available such as Fast Fourier transformation (FFT), Frequency Analysis Mapping On Unusual Sampling (FAMOUS) and least squares spectral analysis (LSSA). LSSA software was developed in the Department of Geodesy and Geomatics Engineering at the University of New Brunswick and it is based on the developments by Vaníček (1969, 1971), Wells et al. (1985) and Pagiatakis (1999, 2000). However, all of them use a set of base functions made up of sine and cosine functions in the decomposition process, to generate a frequency spectrum. FAMOUS and LSSA have been developed as an alternative to bypass some of the limitations present in the classical Fourier methods. These limitations include the need for long time series, constant sampling rate, equally weighted data values, no presence of gaps or datum shifts all of which render the time series strongly non stationary (Fig. 3).

FAMOUS (Frequency Analysis Mapping On Unusual Sampling) makes the decomposition of a time series as

$$\psi(t) = c_0 + \sum_{i=1}^k c_i \cos(2\pi v_i^0 t) + s_i \sin(2\pi v_i^0 t) \quad (5)$$

where  $c_i$  and  $s_i$  are constant or polynomial of time:

$$\begin{aligned} c_i &= a_i^0 + a_i^1 t + a_i^2 t^2 + \dots + a_i^p t^p \\ s_i &= b_i^0 + b_i^1 t + b_i^2 t^2 + \dots + b_i^p t^p \end{aligned} \quad (6)$$

The model

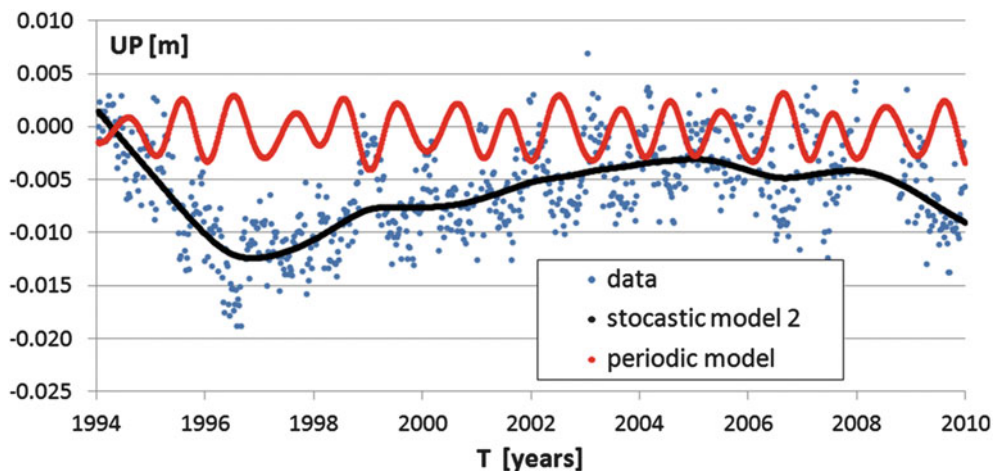
$$\min |S(t) - \psi(t)|^2 \quad (7)$$

is a non linear least square very sensitive to the starting values, solved in two steps (SVD and Levenberg-Marquardt minimization).

The solution is also given in term of  $A \cos(\omega t + \phi)$  and the signal can be reconstructed as

$$\psi(t) = c_0 + \sum_{i=1}^k A_i \cos(2\pi v_i^0 t + \phi_i) \quad (8)$$

FAMOUS allows the analysis of equally weighted data, with known or unknown a priori variance factors, assuming them to be uncorrelated; the algorithm can handle unevenly spaced time series without a pre-processing requirement. Tests of statistically significant spectral peaks are implemented, with respect to S/N ratio. FORTRAN source code is available by F. Mignard.



**Fig. 3** KOSG station 16 years time series, UP component. The process is clearly non cyclostationary, it presents a long term non linear behavior and a large data gap in the year 2008. The model 2 (*black dashed line*) with  $\sigma_p = 10^{-6}$  m and  $\sigma_{\dot{p}} = 10^{-6}$  m, follows the long term behavior

of the data. The periodic model (*red dashed line*) has been estimated by FAMOUS on the stochastic model residuals, and it is the sum of a linear component and of five different periodic signals with periods of 775.6, 527.8, 365.8, 271.0 and 41.6 days

#### 4 Analysis of IGS REPRO1 Time Series

Coordinate and formal error weekly time series of the full data set of 526 IGS permanent stations have been analyzed. The presented three steps procedure has been implemented in fortran90 integrating the FAMOUS source code and it is completely automatic. For this reason a complete reanalysis of REPRO1 data set takes only few minutes. The S/N threshold value for acceptance of the detected signals has been set equal to 3, as in Collilieux et al. (2007), Mignard (2003). For each time series we obtain:

- the estimated offsets,
- the long term signal model,
- the cyclostationary residuals,
- the frequency, amplitude and phase of the detected harmonics, with their RMS.

We must note that the analysis has been conducted on the REPRO1 raw coordinates, without taking in account for atmospheric loading correction. It was been observed in Mtamakaya (2012) that a slight improvement to coordinates repeatability may result if Atmospheric Pressure Loading were included in the processing, however this does not cause any significant reduction in spectral peaks that are still present in the REPRO1 solutions. See also (Tregoning and Watson 2009) for a quantitative analysis of atmospheric loading. The global analysis starts from the number of detected signals over the total number of stations in (%), reported in Fig. 4. We can observe three classes of signals, related to seasonal, orbital (draconitic) and tidal effects. The minimum sampling frequency (Nyquist frequency) for weekly time series is 14 days, at which a strong signal has been also detected, that will be attributed to tidal model errors.

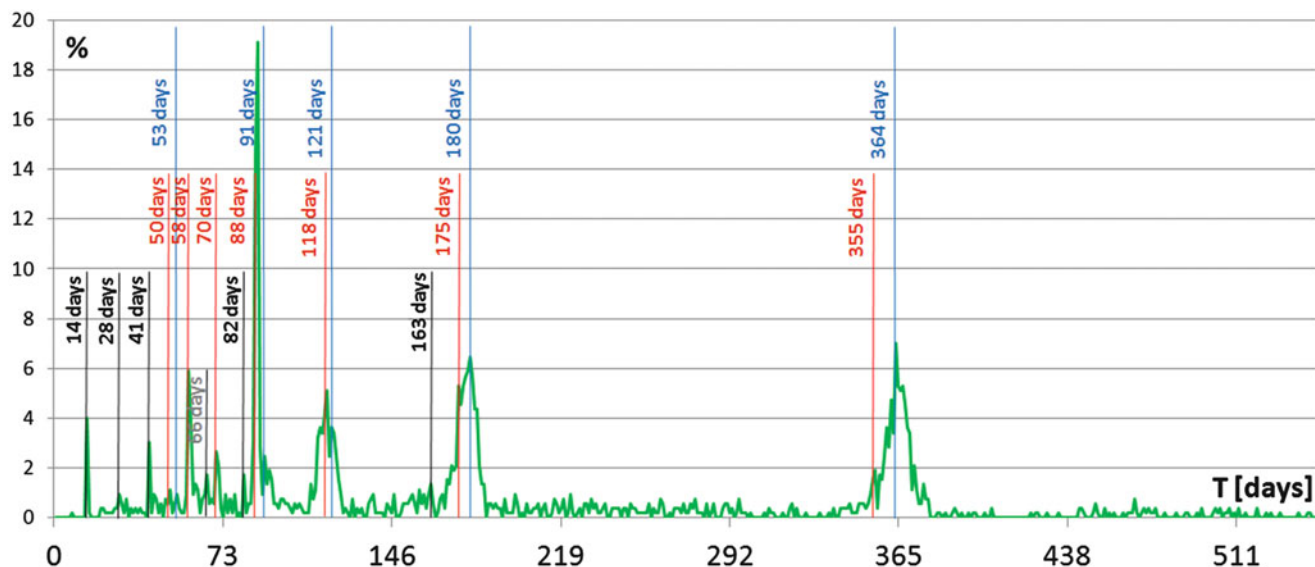
The 89% of the stations present an annual signal and signals have been detected at the 2nd, 3rd, 4th and 7th harmonics of the solar year (seasonal signals), with respective periods of 182.6, 121.8, 91.3 and 52.2 days. The annual signal amplitude, represented in Fig. 5, has a mean value of 3.1 mm and it is strongly spatially correlated. The maximum value of 1.0 cm is at the station WSLR (Whistler, Canada).

The observed signatures appear to be consistent also around the 1st, 2nd, 3rd, 4th, 6th and 7th draconitic<sup>2</sup> harmonics with respective periods of 351.2, 175.6, 117.1, 87.8, 70.2, 58.5 and 50.2 days. No signal has been detected at the higher frequencies of the draconitic harmonics, even if some signal can be aliased by the tidal harmonics. Note in the Fig. 4 that the 1st, 2nd and 3rd draconitic harmonics overlap the 1st, 2nd and 3rd seasonal harmonics. For this reason, in many time series draconitic errors cannot be distinguished by seasonal signals and contribute to them, as already noted by Rebischung et al. (2012), and beating between draconitic and seasonal harmonics can explain the annual an inter-annual amplitude variations. Draconitic and solar year are in phase every 26 years<sup>3</sup> during which the amplitude of the combined signal ranges between 0 and 6.8 mm with a beating effect. The superimposition of the seven draconitics and four solar detected harmonics results in a signal amplitude that ranges from 3 to 27 mm. The wavelength RMS are larger for the detected seasonal harmonics and smaller for the draconitic and tidal harmonics, because seasonal effects present a greater variability. The mean amplitudes are coher-

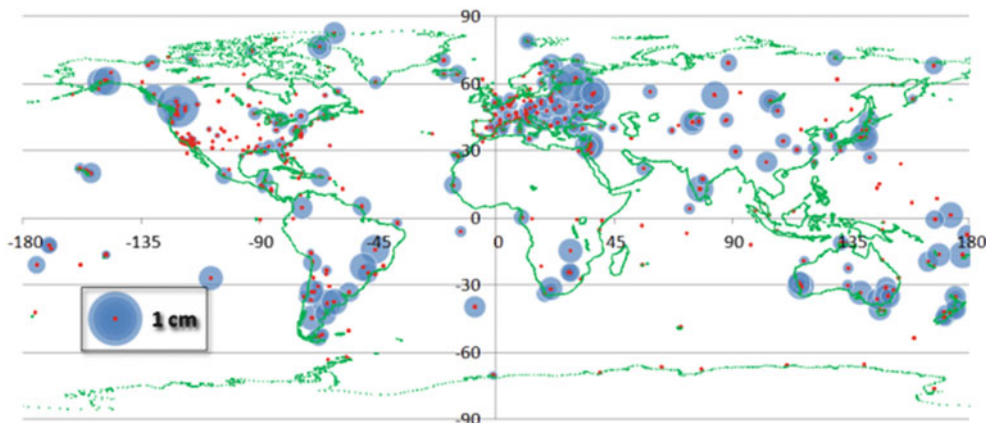
<sup>2</sup>Draconitic year is the interval of  $1.040 \pm 0.008$  cycles per year ( $351.2 \pm 2.8$  days) needed for the Sun to return to the same point in space relative to the GPS orbital nodes (as viewed from the Earth).

<sup>3</sup> $365.25 / (365.25 - 351.2) \approx 26$ .





**Fig. 4** Detected signals over the number of stations in (%). The signals evidenced in *blue* are the solar year harmonics (seasonal term), in *red* the draconitic harmonics, while in *black* are the tidal harmonics



**Fig. 5** Amplitude of the annual term

ent with the values reported by Ostini (2012), obtained by the FODITS method proposed in Ostini et al. (2008).

Two peaks has been found with period of 14 and 28 days that are doubtless related to tidal model errors. The peak at 14 days has been attributed to sub-daily EOP tidal errors by Ray et al. (2013). The relative motions of the Earth, Moon and Sun cause the tides to vary in numerous tidal cycles, the two most important ones being the spring-neap cycle and the equinoctial cycle. The spring-neap cycle is a 14.77 day cycle resulting from the tidal influence of the sun and moon either reinforcing or partially cancelling each other (neap tides). The semi-annual equinoctial cycle is caused by the tilt of the Earth, and its orbit around the Sun which leads to higher than average spring tides around the time of the equinoxes (March and September) and lower than average spring tides in June and December. Because of its seasonality this effect cannot be distinguished by other seasonal effects. The Moon crosses

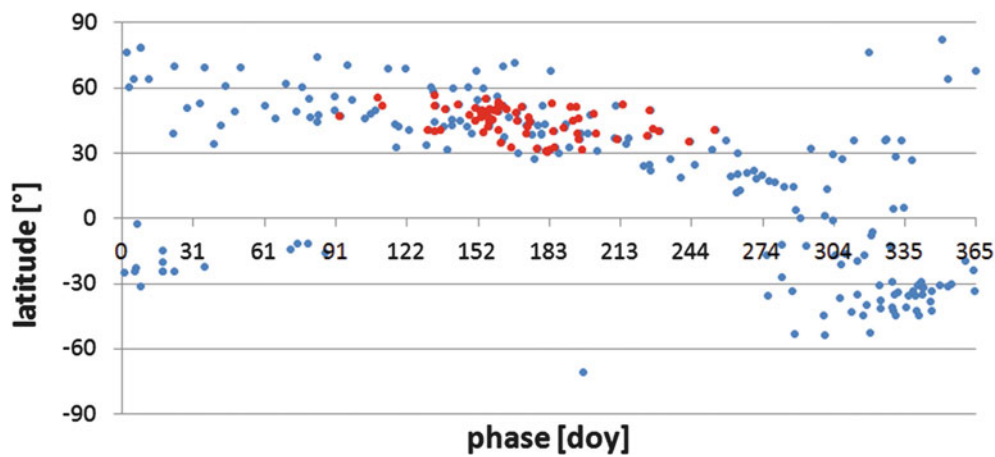
the ecliptic at the same node every  $\sim 27.3$  days, and a peak can also be observed at a frequency of 28 days. Finally three peaks are at 41 ( $=3 \cdot 14$ ), 82 ( $=6 \cdot 14$ ) and 163 ( $=12 \cdot 14$ ) days, that can also be related to tidal model errors. A synthesis of the detected signals is reported in Table 1.

The remaining signatures could be attributed to other un-modeled effects that must be investigated, such as non tidal loading displacement, high order ionosphere terms and mismodeling effect in GPS attitude models. For example, on remarkable peak is at 66 days, that cannot be related to seasonal, draconitic or tidal effects.

The phase  $\phi$  represents the signal maximum in the model (8). It seems to be spatially correlated, as can be noted in Fig. 6, where the signal phase is represented with respect to station latitudes. The European cluster is represented in red and has it maximum between doy (day of year) 120 and 220 (April–June). As consequence of the seasonal loading

**Table 1** Detected signals that are consistent with the solar, draconitic and tidal harmonics

Cpy	Harmonics (days)						Mean amplitudes (mm)		
	Expected			Estimated			Drac.	Solar	Tidal
	Drac.	Solar	Tidal	Drac.	Solar	Tidal			
1	351.2	365.3	164.0	355 ± 3	364 ± 10	163 ± 2	3.2	3.1	1.3
2	175.6	182.6	82.0	175 ± 3	180 ± 6	82 ± 1	2.0	1.7	0.8
3	117.1	121.8		118 ± 2	121 ± 5		1.3	1.1	
4	87.8	91.3	41.0	88 ± 2	91 ± 3	41 ± 1	1.1	1.5	0.7
5	70.2	73.1		70 ± 2			1.0		
6	58.5	60.9	27.3	58 ± 1		28 ± 1	1.0		1.1
7	50.2	52.2		50 ± 1	53 ± 1		1.0	0.7	
12			13.7			14 ± 1			1.1



**Fig. 6** Phase of the annual term, in term of doy of the UP maximum. It can be noted some correlation with the station latitude

effects, the phase distribution seems to be coherent with the Earth deformation global model proposed by Blewitt et al. (2001), according to which during February to March the Northern Hemisphere compresses and the Southern Hemisphere expands. The opposite pattern of deformation occurs during August to September. More uniform data from both the hemispheres are necessary to clearly identify this latitude dependant effect.

## 5 Conclusions

Spectral analysis of weekly station coordinate time series of 526 IGS sites reveals signals at the seasonal, draconitic and tidal harmonics. The analysis has been conducted on the UP component of the REPRO1 time series, while E and N are not yet analyzed and they must be considered in future works. It has been shown that the detected annual signal is spatially correlated in both amplitude and phase, and it depends on the loading changes due to the water cycle. Similar analysis must be conducted on the sub annual signals, in order to better understand their origin, that seems to be related to model errors. Some geophysical effect can also be observed at a

global level, such as the expansion of the hemispheres during the summer and their contraction during winter. Other non periodical geophysical signals can be potentially discovered in the long term signal models, not studied in the present work.

## References

Albertella A, Betti B, Sansò F, Tornatore V (2006) Real time and batch navigation solutions: alternative approaches. Bollettino SIFET 2  
 Blewitt G, Lavallée D, Clarke P, Nurutdinov K (2001) A new global mode of earth deformation: seasonal cycle detected. Science 294(5550). doi:10.1126/science.1065328  
 Collilieux X, Altamimi Z, Coulot D, Ray J, Sillard P (2007) Comparison of VLBI, GPS, SLR height residuals from ITRF2005 using spectral and correlation methods. J Geophys Res 112, B12403. doi:10.1029/2007JB004933  
 Collilieux X, Métivier L, Altamimi Z, van Dam T, Ray J (2010) Quality assessment of GPS reprocessed Terrestrial Reference Frame. GPS Solut. doi:10.1007/s10291-010-0184-6  
 Fritsche M, Dietrich R, Rülke A, Rothacher M, Steigenberger P (2009) Low-degree earth deformation from reprocessed GPS observations. GPS Solut. doi:10.1007/s10291-009-0130-7  
 Gazeaux J, Williams S, King M, Bos M, Dach R, Deo M, Moore AW, Ostini L, Petrie E, Roggero M, Teferle FN (2013) Detecting offsets

- in GPS time series: first results from the Detection of offsets in GPS experiment. *J Geophys Res Solid Earth*. doi:[10.1002/jgrb.50152](https://doi.org/10.1002/jgrb.50152)
- Mtamakaya JD (2012) Assessment of atmospheric pressure loading on the international GNSS REPRO1 solutions periodic signatures, Tech Report 282, Department of Survey Engineering, University of New Brunswick, Fredericton
- Mignard F (2003) FAMOUS, Frequency Analysis Mapping on Unusual Sampling, (OCA Cassiopee), Technical report. Obs. de la Cote d'Azur Cassiopee, Nice, <ftp://ftp.obs-nice.fr/pub/mignard/Famous/>; [http://obswww.unige.ch/~eyer/VSWG/Meeting4/MINUTES/VSWG4\\_FM2.pdf](http://obswww.unige.ch/~eyer/VSWG/Meeting4/MINUTES/VSWG4_FM2.pdf)
- Ostini L, Dach R, Meindl M, Schaer S, Hugentobler U (2008) FODITS: a new tool of the Bernese GPS Software to analyze time series. In: EUREF 2008 symposium, Brussels
- Ostini L (2012) Analysis and quality assessment of GNSS-derived parameter time series. Ph.D. dissertation, Philosophisch naturwissenschaftlichen Fakultät der Universität Bern
- Pagiatakis SD (1999) Stochastic significance of peaks in the least-squares spectrum. *J Geodesy* 73:67–78
- Pagiatakis SD (2000) Application of the least-squares spectral analysis to superconducting gravimeter data treatment and analysis; Proceedings of the workshop on High precision gravity measurements with application to geodynamics and Second GGP workshop. Cahiers Du Centre Europeen De Geodynamique et de Seismologie (ECGS) 17:103–113
- Perfetti N (2006) Detection of station coordinates discontinuities within the Italian GPS Fiducial Network. *J Geodesy*. Springer Berlin/Heidelberg ISSN 0949-7714 (Print) 1432-1394 (Online)
- Ray JR, Altamimi Z, Collilieux X, van Dam T (2008) Anomalous harmonics in the spectra of GPS position estimates. *GPS Solut.* doi:[10.1007/s10291-007-0067-7](https://doi.org/10.1007/s10291-007-0067-7)
- Ray J, Griffiths J, Collilieux X, Rebischung P (2013) Subseasonal GNSS positioning errors. *Geophys Res Lett* 40(22):5854–5860. doi:[10.1002/2013GL058160](https://doi.org/10.1002/2013GL058160)
- Rebischung P, Collilieux X, van Dam T, Ray J, Altamimi Z (2012) Analysis effects in IGS station motion time series. Presented at IGS workshop 2012, Olsztyn
- Roggero M (2008) Kinematic GPS batch processing, a source for large sparse problems. In: VI Hotine Marussi symposium on theoretical and computational geodesy. Springer-Verlag, Berlin/Heidelberg, ISBN: 3-540-74583-1
- Roggero M (2012) Discontinuity detection and removal from data time series. In: VII Hotine Marussi symposium on theoretical and computational geodesy. Springer-Verlag, Berlin/Heidelberg, ISBN:9783642220777
- Steigenberger P (2009) Reprocessing of a global GPS network, Ph.D. dissertation, Fakultät für Bauingenieur und Vermessungswesen, TU München
- Teunissen P (1998) Quality control and GPS. In: Teunissen PJG, Kleusberg A (eds) *GPS for geodesy*. Springer-Verlag, Berlin/Heidelberg/New York. ISBN 3-540-63661-7
- Tregoning P, Watson C (2009) Atmospheric effects and spurious signals in GPS analyses. *J Geophys Res* 114, B09403. doi:[10.1029/2009JB006344](https://doi.org/10.1029/2009JB006344)
- Vaniček P (1969) Approximate spectral analysis by least-squares fit. *Astrophys Space Sci* 4:387–391
- Vaniček P (1971) Further development and properties of the spectral analysis by least squares. *Astrophys Space Sci* 12:10–33
- Wells D, Vaniček P, Pagiatakis S (1985) Least squares spectral analysis revisited. Tech Report 84, Department of Survey Engineering, University of New Brunswick, Fredericton

**Geopotential Modeling, Boundary Value Problems  
and Height Systems**

---

# Determination of $W_0$ from the GOCE Measurements Using the Method of Fundamental Solutions

Róbert Čunderlík

---

## Abstract

The paper presents the method of fundamental solutions (MFS) applied for global gravity field modelling. MFS as an inherent mesh-free method is used to derive the geopotential and its first derivatives from the second derivatives observed by the GOCE satellite mission, namely from the radial components of the gravity tensor. Unknown coefficients of the approximate solution by MFS are determined at the source points located directly on the Earth's surface. Afterwards, the disturbing potential or gravity disturbance can be evaluated at any point above the Earth's surface. To get their values on the Earth's surface, singularities of the fundamental solutions need to be overcome. In this paper two strategies are used: (1) the source points are located on a fictitious boundary (FB), which is situated below the Earth's surface, or (2) ideas of the singular boundary method that isolate the singularities are implemented. The paper studies how a depth of FB influences accuracy of the MFS solutions. All particular solutions are compared with the GOCO03S satellite-only geopotential model. In all cases mean values of the residuals are smaller than  $0.04 \text{ m}^2\text{s}^{-2}$  ( $\sim 4 \text{ mm}$ ). The best agreement in terms of the standard deviation of residuals is for the FB depth of 20 km. Finally, the geopotential on the DTU10 mean sea surface is evaluated from the MFS solutions resulting in the  $W_0$  estimates. The obtained  $W_0$  values differ from ones based on GOCO03S or EGM2008 by less than  $0.1 \text{ m}^2\text{s}^{-2}$  ( $\sim 1 \text{ cm}$ ).

---

## Keywords

Geopotential on the mean sea surface • Global gravity field modelling • GOCE measurements • Method of fundamental solutions •  $W_0$  estimates

---

## 1 Introduction

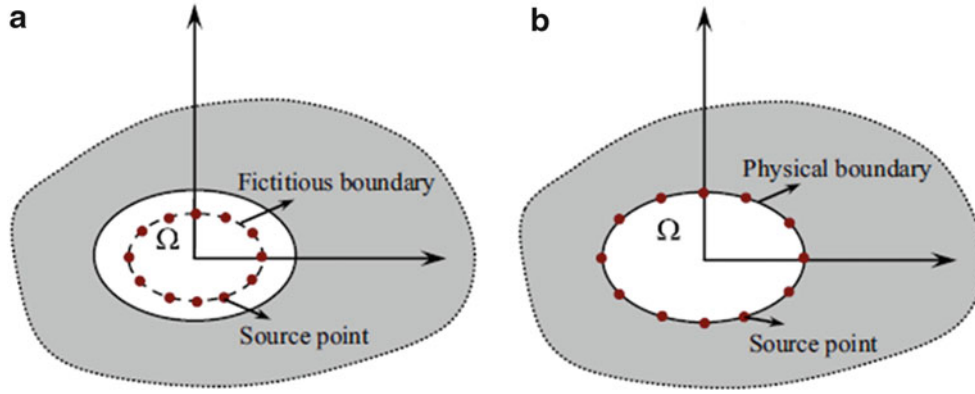
The unification of local vertical datums and establishment of the World Height System (WHS) is one of the main tasks of modern geodesy, cf. (Sideris and Fotopoulos 2012). It involves a determination of  $W_0$  as a reference value of the geopotential on the geoid. All recent  $W_0$  estimates are basically derived from global geopotential models (GGMs) that

are developed using the spherical harmonics (SH) approach, cf. (Burša et al. 2007), (Sánchez 2009), (Dayoub et al. 2012), (Čunderlík et al. 2014). However, the GOCE satellite mission, which is directly measuring the second derivatives of the geopotential, has brought new opportunities in applications of other numerical approaches for global gravity field modelling. In this paper, the method of fundamental solutions (MFS) is presented to derive the geopotential and its first derivatives on or above the Earth's surface from the second derivatives observed by GOCE.

Basic ideas for the formulation of MFS were first developed by V. D. Kupradze and M. A. Alexidze, cf. (Kupradze and Alexidze 1964). However, MFS as a computational technique was proposed by Mathon and Johnston (1977). Later,

---

R. Čunderlík (✉)  
Department of Mathematics, Slovak University of Technology  
in Bratislava, Bratislava, Slovakia  
e-mail: [cunderli@svf.stuba.sk](mailto:cunderli@svf.stuba.sk)



**Fig. 1** Distribution of the source points for the exterior potential problem using (a) the method of fundamental solution (MFS), and (b) the singular boundary method (SBM) (from the source: Gu et al. 2012)

MFS was extended to deal with inhomogeneous equations and time-dependent problems (Golberg and Chen 1994). At present, MFS has become a useful tool for solving a large variety of physical and engineering problems, cf. (Hon and Wei 2005), (Fan et al. 2009) or (Chen et al. 2011). To cure the problem of a fictitious boundary in MFS, some new techniques have recently been developed, e.g., the singular boundary method (SBM) (Chen and Wang 2010).

In this paper MFS is applied to derive the disturbing potential and its first derivatives from the radial components  $T_{rr}$  of the disturbing tensor observed by GOCE. Numerical experiments show how the depth of the fictitious boundary influences the accuracy of the obtained MFS solution on or above the Earth's surface. In case that the source points are located directly on the Earth's surface, the ideas of SBM that isolate singularities of the fundamental solution (Gu et al. 2012) are applied. Finally, the geopotential on the mean sea surface is evaluated from the different MFS solutions. It allows estimating the  $W_0$  values independently from the SH-based approach. The objective of the paper is to explore how such  $W_0$  estimates differ from ones computed from the SH-based GGMs.

## 2 MFS for the Potential Problems

MFS is a technique for the numerical solution of certain elliptic boundary value problems (BVPs) (Mathon and Johnston 1977). It belongs to the general class of the boundary collocation methods. Like the boundary element method (BEM), it is applicable when the fundamental solution of a partial differential equation (PDE) of interest is known. MFS was developed to overcome the major drawbacks of BEM, i.e., to avoid numerical integration of the singular fundamental solution by introducing a fictitious

boundary (FB) outside the physical domain. In contrast to BEM, MFS is an inherent mesh-free method and does not involve integral evaluation. Hence, it provides an efficient computational alternative for problems in higher dimensions with irregular domains.

In the following we focus on the exterior potential problem in 3D that corresponds to the geodetic BVP. Let us consider the potential field  $u$  satisfying the Laplace equation exterior a 3D domain  $\Omega$  (the Earth) (Fig. 1)

$$\nabla^2 u(\mathbf{x}) = 0, \quad \mathbf{x} \in \text{ext.}\Omega, \quad (1)$$

with the following boundary conditions (BC)

$$u(\mathbf{x}) = \bar{u}(\mathbf{x}), \quad \mathbf{x} \in \Gamma_D \quad (\text{Dirichlet BC}), \quad (2)$$

$$q(\mathbf{x}) = \frac{\partial u}{\partial \mathbf{n}}(\mathbf{x}) = \bar{q}(\mathbf{x}), \quad \mathbf{x} \in \Gamma_N \quad (\text{Neumann BC}), \quad (3)$$

where  $\Gamma_D$  and  $\Gamma_N$  construct the whole boundary of the domain  $\Omega$ , and  $\mathbf{n}$  denotes the outward normal.

An approximate solution by MFS is expressed as a linear combination of the fundamental solutions with respect to different source points

$$u(\mathbf{x}^i) = \sum_{j=1}^N \alpha_j G(\mathbf{x}^i, \mathbf{s}^j), \quad (4)$$

$$q(\mathbf{x}^i) = \frac{\partial u(\mathbf{x}^i)}{\partial \mathbf{n}_{\mathbf{x}^i}} = \sum_{j=1}^N \alpha_j \frac{\partial G(\mathbf{x}^i, \mathbf{s}^j)}{\partial \mathbf{n}_{\mathbf{x}^i}}, \quad (5)$$

where  $\mathbf{x}^i$  is the  $i$ -th observation point and  $s^j$  is the  $j$ -th source point,  $\alpha_j$  denotes the  $j$ -th unknown coefficient of the distributed source at  $s^j$ ,  $N$  represents the number of source points and

$$G(\mathbf{x}^i, s^j) = \frac{1}{4\pi |\mathbf{x}^i - s^j|}, \quad (6)$$

is the fundamental solution of the Laplace equation in 3D, which represents the basis functions of the method. For a well-posed BVP, the unknown coefficients  $\{\alpha_j\}, j = 1, \dots, N$ , can be determined by collocating  $N$  observation points with BC from Eq. (2) or (3). Once all the unknown coefficients  $\{\alpha_j\}$  are solved, physical quantities at any point inside the physical domain (i.e., the exterior of  $\Omega$  in our case) including its boundary can be easily evaluated from the field equations (4) or (5).

To avoid singularities of the fundamental solutions, the source points are located on the FB outside the computational domain. For the exterior BVP described in Eqs. (1–3), the FB is inside  $\Omega$ , i.e., below the Earth's surface (Fig. 1a). However, despite many years of great effort, the determination of the FB is largely based on experiences, especially for problems in complicated geometries and higher dimensions.

### 3 MFS for Gravity Field Modelling from the GOCE Measurements

The gravity field modelling is usually formulated in terms of the Laplace equation (1) for the disturbing potential  $T$ . The GOCE observations provide the second derivatives of the geopotential, or the disturbing potential. In this study, the radial components  $T_{rr}$  of the disturbing tensor are used to derive the unknown coefficients  $\alpha_j$  at the source points  $s^j$  using the expression

$$T_{rr}(\mathbf{x}^i) = \frac{\partial^2 T(\mathbf{x}^i)}{\partial r_{x^i}^2} = \sum_{j=1}^N \alpha_j \frac{\partial^2 G(\mathbf{x}^i, s^j)}{\partial r_{x^i}^2}, \quad (7)$$

where

$$\frac{\partial^2 G(\mathbf{x}^i, s^j)}{\partial r_{x^i}^2} = \frac{\partial^2 G_{i,j}}{\partial r_i^2} = \frac{1}{4\pi} \left[ \frac{1}{d_{ij}^3} - 3 \frac{\langle \mathbf{d}_{ij}, \mathbf{r}_i \rangle^2}{d_{ij}^5} \right], \quad (8)$$

and  $\mathbf{r}_i$  denotes the radial unit vector at  $\mathbf{x}^i$ ,  $\mathbf{d}_{ij} = \mathbf{x}^i - s^j$  and  $d_{ij} = |\mathbf{d}_{ij}|$  represents the distance between the  $i$ -th collocation point and the  $j$ -th source point. By collocating  $N$  observation

points with respect to  $N$  source points, we get the linear system of equations

$$\begin{bmatrix} \frac{\partial G_{1,1}^2}{\partial r_1^2} & \dots & \dots & \frac{\partial G_{1,N}^2}{\partial r_1^2} \\ \vdots & \ddots & & \vdots \\ \vdots & & \ddots & \vdots \\ \frac{\partial G_{N,1}^2}{\partial r_N^2} & \dots & \dots & \frac{\partial G_{N,N}^2}{\partial r_N^2} \end{bmatrix} \cdot \begin{bmatrix} \alpha_1 \\ \vdots \\ \vdots \\ \alpha_N \end{bmatrix} = \begin{bmatrix} T_{rr1} \\ \vdots \\ \vdots \\ T_{rrN} \end{bmatrix}. \quad (9)$$

Since the GOCE observations are given sufficiently far from the Earth (approximately 250 km above the Earth's surface), the source points can be located directly on the Earth's surface considering its complicated topography. Such a configuration does not generate any singularities. The unknown coefficients  $\{\alpha_j\}$  can be determined by solving the linear system of equations (9). Afterwards, the disturbing potential or its first derivatives can be easily evaluated anywhere above the Earth's surface using Eq. (4) or (5). The problem with the singularities appears when computing the gravity field quantities directly on the Earth's surface. In this case it is possible to use two different strategies: (1) to locate source points on the FB, which needs to be shifted below the Earth's surface, or (2) to apply ideas of SBM that isolate singularities of the fundamental solution at source points on the Earth's surface.

In the first approach, a main problem is to determine an optimal position of the FB. As mentioned earlier, this is largely based on experiences. Therefore, in the presented numerical experiments we will change step by step the depth of FBs, testing how it influences the resulting MFS solution on the Earth's surface.

In the second approach, the ideas of SBM (Chen and Wang 2010) are implemented to overcome singularities of the fundamental solution. Like MFS, SBM also uses the fundamental solution as the basis kernel function of its approximation. In contrast to MFS, the collocation and source points of SBM are coincident and they are all placed on the physical boundary (Fig. 1b) avoiding any FB. For the 3D exterior potential problem described in Eqs. (1–3), the SBM interpolation formulation can be expressed as

$$u(\mathbf{x}^i) = \sum_{j=1, i \neq j}^N \alpha_j G(\mathbf{x}^i, s^j) + \alpha_i u_{ii}, \quad (10)$$

$$q(\mathbf{x}^i) = \sum_{j=1, i \neq j}^N \alpha_j \frac{\partial G(\mathbf{x}^i, s^j)}{\partial \mathbf{n}_{x^i}} + \alpha_i q_{ii}, \quad (11)$$

where  $u_{ii}$  and  $q_{ii}$ , called the origin intensity factors, denote the singular terms  $G(\mathbf{x}^i, \mathbf{s}^j)$  and  $\partial G(\mathbf{x}^i, \mathbf{s}^j)/\partial \mathbf{n}$ , respectively, i.e., the diagonal elements of the SBM interpolation matrix. These singularities need to be regularized using some special treatment. Applying the regularization technique proposed in (Gu et al. 2012) and omitting details described in this paper, the original singular term  $q_{ii}$  for the Neumann boundary equation (11) can be transformed into the regular term

$$q_{ii} = \frac{1}{P_i} \left[ 1 - \sum_{j=1, i \neq j}^N P_j \frac{\partial G(\mathbf{x}^i, \mathbf{s}^j)}{\partial \mathbf{n}_{s^j}} \right], \quad (12)$$

where  $P_i$ , or  $P_j$ , is the area corresponding to surrounding of the collocation point  $\mathbf{x}^i$ , or the source point  $\mathbf{s}^j$ , respectively. In this way a distribution of the source points is taken into account. To evaluate the origin intensity factor  $u_{ii}$  for the Dirichlet boundary equation (10), an inverse interpolation technique can be used. Due to the limited extend of this paper, the readers are kindly addressed to (Gu et al. 2012) for more details.

Since the observations from GOCE are sufficiently far from the Earth's surface, the unknown coefficients  $\{\alpha_j\}$  can be determined from the linear system of equations (9). Afterwards, the origin intensity factors  $u_{ii}$  and  $q_{ii}$  need to be determined and finally Eqs. (10–11) can be used to evaluate the disturbing potential or its first derivatives at the source points directly on the Earth's surface. In this way the problem of singularities can be overcome.

## 4 Numerical Experiments

In the numerical experiments we have processed the GOCE measurements from its first 61 days period, i.e., from October 1 to December 1 2009. In particular, the radial components  $V_{rr}$  of the gravity tensor have been transformed to  $T_{rr}$  of the disturbing tensor (Fig. 2a) using parameters of the GRS-80 normal gravity field. Then the nonlinear diffusion filtering has been applied to reduce the noise included in the input data. The regularised Peron-Malik model has efficiently reduced the noise while preserving main structures (Fig. 2b), for more details see (Čunderlík et al. 2013). In the first experiment the source points have been located directly on the Earth's surface with a resolution of  $0.075^\circ$ . It has corresponded to 5,760,002 points ( $N$ ) regularly distributed over the Earth's surface. To consider the real topography, the vertical position of the source points were generated from

the SRTM30\_PLUS global topography model (Becker et al. 2009).

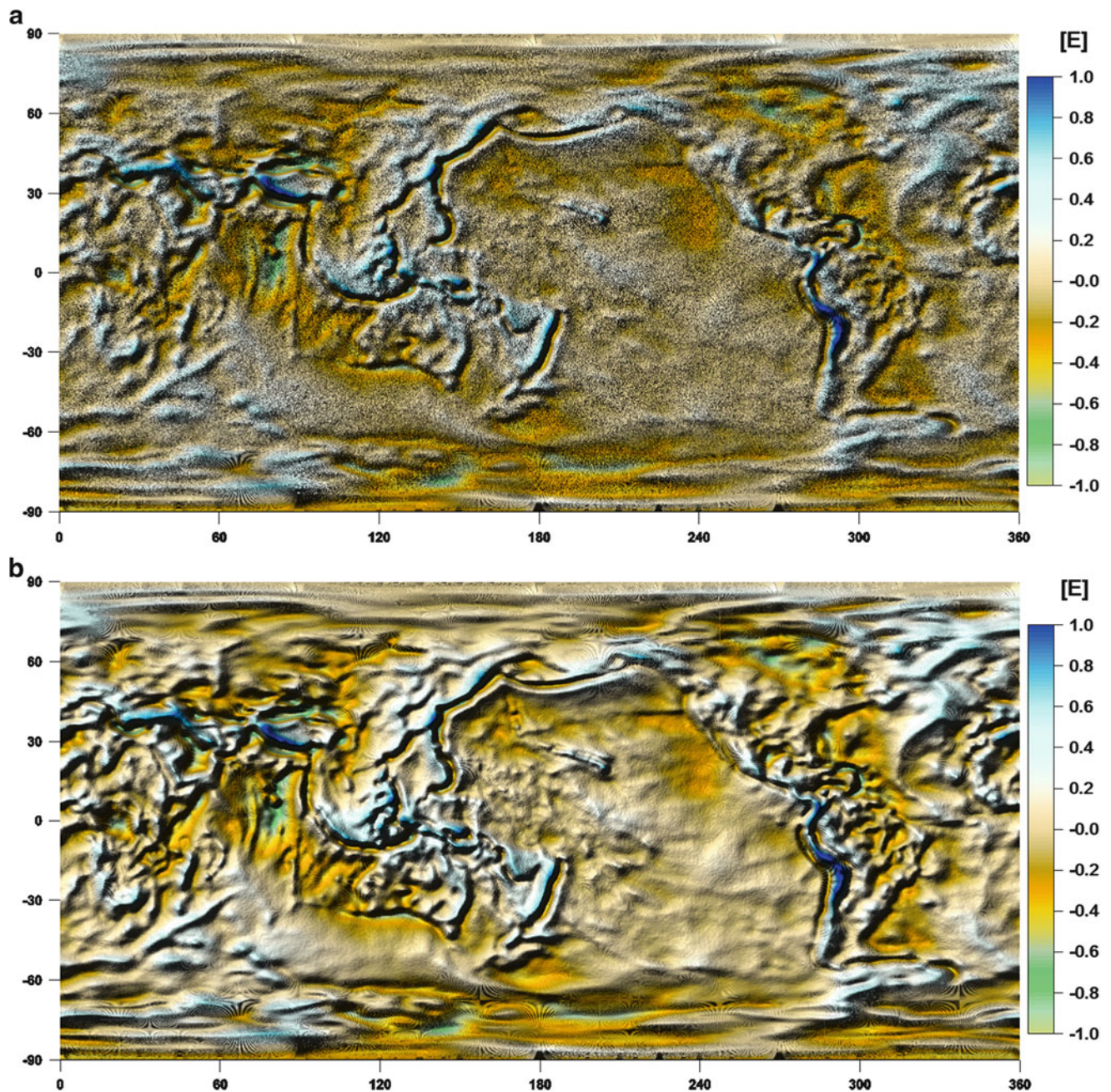
To get the linear system of equations (9), the same number of the input observations (collocation points) has been chosen. Their horizontal positions as well as ordering have been adopted from the source points. This has required an interpolation from the original GOCE measurements (firstly filtered by the nonlinear diffusion). To reduce the enormous memory requirements for the full matrix in Eq. (9), an iterative approach has been applied for the elimination of the far zones' contributions, primarily proposed for the direct BEM (Čunderlík and Mikula 2010). This approach, together with a parallel implementation using the MPI (Message Passing Interface) procedures, enables us to reach such a high level of resolution. The large-scale parallel computations were performed on the cluster with 1 TB of the distributed memory. At first, the unknown coefficients  $\{\alpha_j\}$  at the source points located directly on the Earth's surface have been determined. To obtain the disturbing potential or its derivatives at these points, the strategy of SBM has been used determining the unknown origin intensity factors  $u_{ii}$  and  $q_{ii}$  (see Sect. 3).

Afterwards, the MFS approach based on the FB has also been used. The depth of the FB has been changing step by step, namely the vertical positions of the source points, while the input GOCE observations have remained the same. For every new position of the FB, new set of the coefficients  $\{\alpha_j\}$  has been determined. From these coefficients, the disturbing potential at points on the Earth's surface (with the same positions as in the first experiment) has been evaluated avoiding the problem of singularities.

All particular solutions have been compared with two GGMs developed by the SH-based approach, namely, with the GOCO03S satellite-only model up to degree 250 (Mayer-Gürr et al. 2012) and the EGM2008 combined model up to degree 2160 (Pavlis et al. 2012). Graphs in Fig. 3 depict how the statistical characteristics of the residuals are changing depending on the FB depth. The standard deviation (STD) of the residuals is minimal for a depth of 20 km. The closer to the Earth's surface, the stronger the impact of the singularities becomes and the STD is asymptotically increasing. A special treatment of the singularities by the SBM approach (see Sect. 3) slightly improve this asymptotic worsening, however, the agreement with GOCO03S is worse than for the FB depth in the interval 5–30 km. For the FBs deeper than 30 km, the STD is increasing considerably.

On the other hand, the overall mean value of residuals is changing minimally. For FB depths in the interval



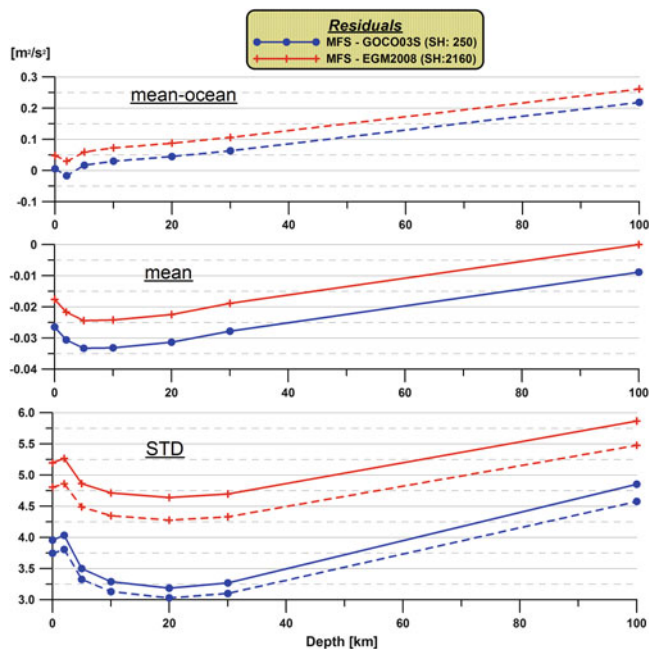


**Fig. 2** (a) The GOCE observations – the radial components  $T_{rr}$  of the disturbing tensor, and (b) after reducing noise using the nonlinear diffusion filtering

0–30 km, it changes less than  $0.01 \text{ m}^2\text{s}^{-2}$  ( $\sim 1 \text{ mm}$ ) (Fig. 3). Considering the mean value over oceans only, it varies within  $0.07 \text{ m}^2\text{s}^{-2}$  ( $\sim 7 \text{ mm}$ ). The mean values over oceans also indicate that the  $W_0$  estimates evaluated from the MFS solutions will differ from one based on GOCO03S by less than  $0.06 \text{ m}^2\text{s}^{-2}$  and from other computed from EGM2008 by less than  $0.1 \text{ m}^2\text{s}^{-2}$ . Here we remind that our solutions

are obtained by processing only 2 months of GOCE data whereas GOCO03S is based on GOCE data measured over 1 year combined with data from GRACE, CHAMP and SLR (Mayer-Gürr et al. 2012).

Figure 4 depicts the disturbing potential on the Earth's surface obtained from the MFS solution (FB depth = 20 km) and from GOCO03S up to degree 250, as well as their



**Fig. 3** An impact of the depth of the fictitious boundary on the obtained MFS solutions – statistical characteristics of the residuals between the MFS solutions and the GOCO03S model up to degree 250 and the EGM2008 model up to degree 2160 (*dashed lines* for the residuals at oceans only)

comparison. Analogously, Fig. 5 shows the first derivatives (the gravity disturbances) for both models and their comparison.

Finally, the geopotential on the DTU10 mean sea surface model (Andersen 2010) is evaluated from the MFS solutions (Fig. 6). The disturbing potential  $T$  computed from the MFS coefficients at points over oceans, whose 3D positions are interpolated from DTU10, is simply added to the normal gravity potential  $U$  evaluated at these points. Then weighted averaging of the geopotential over oceans allows us to estimate the  $W_0$  value for the selected integration area. Table 1 summarizes our  $W_0$  estimates from the MFS solutions for different FB depths. They can be considered independent from the ones obtained using the SH-based GGMs. In spite of quite large differences between the

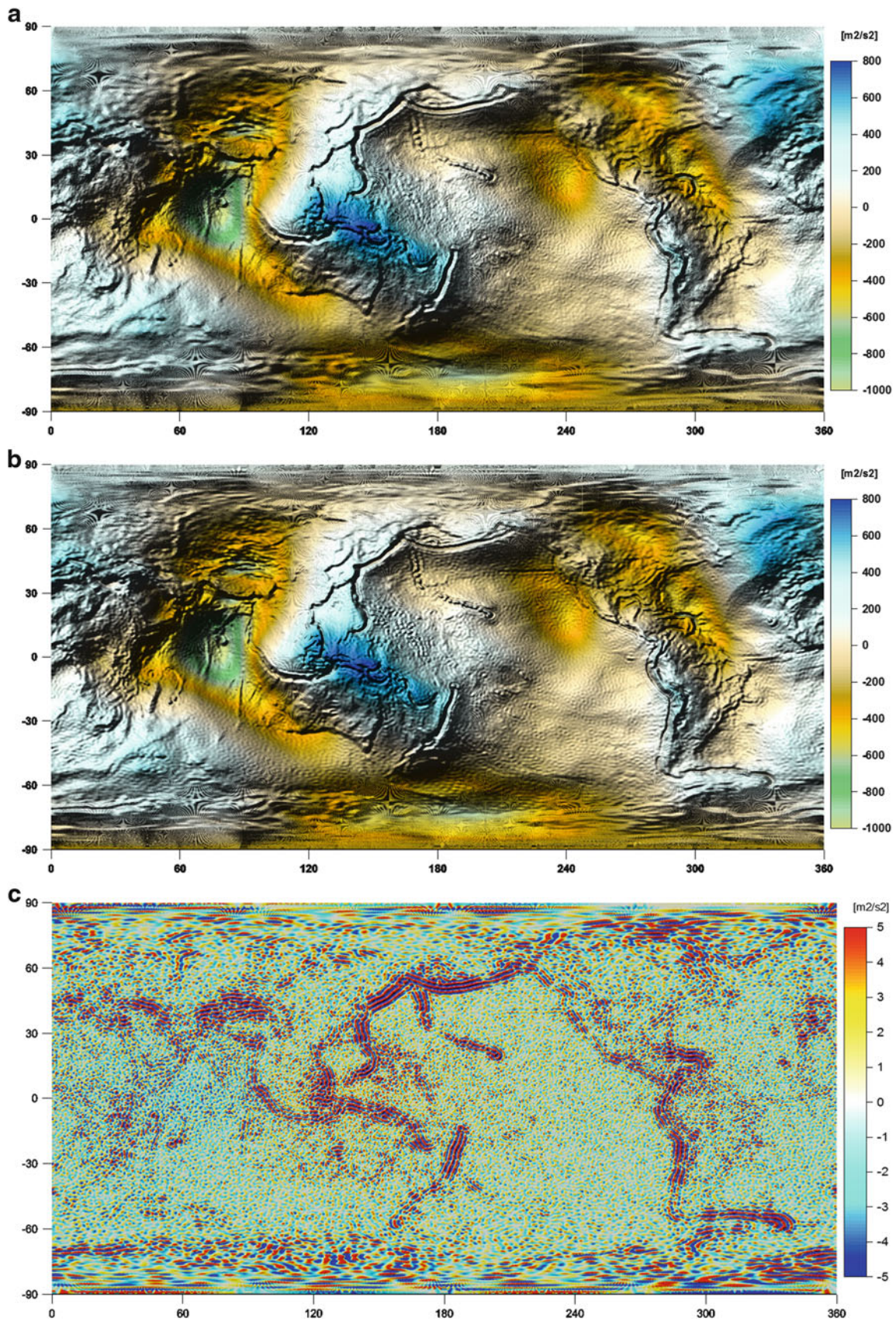
MFS solutions and GOCO03S in zones of abrupt changes of the gravity field, e.g., along edges of the lithospheric plates (Figs. 4c, 5c, and 6), the  $W_0$  estimates differ by less than  $0.1 \text{ m}^2\text{s}^{-2}$ . Such a good agreement shows that both independent approaches, i.e., MFS and SH-based methods, are providing almost the same  $W_0$  estimates.

## 5 Conclusions

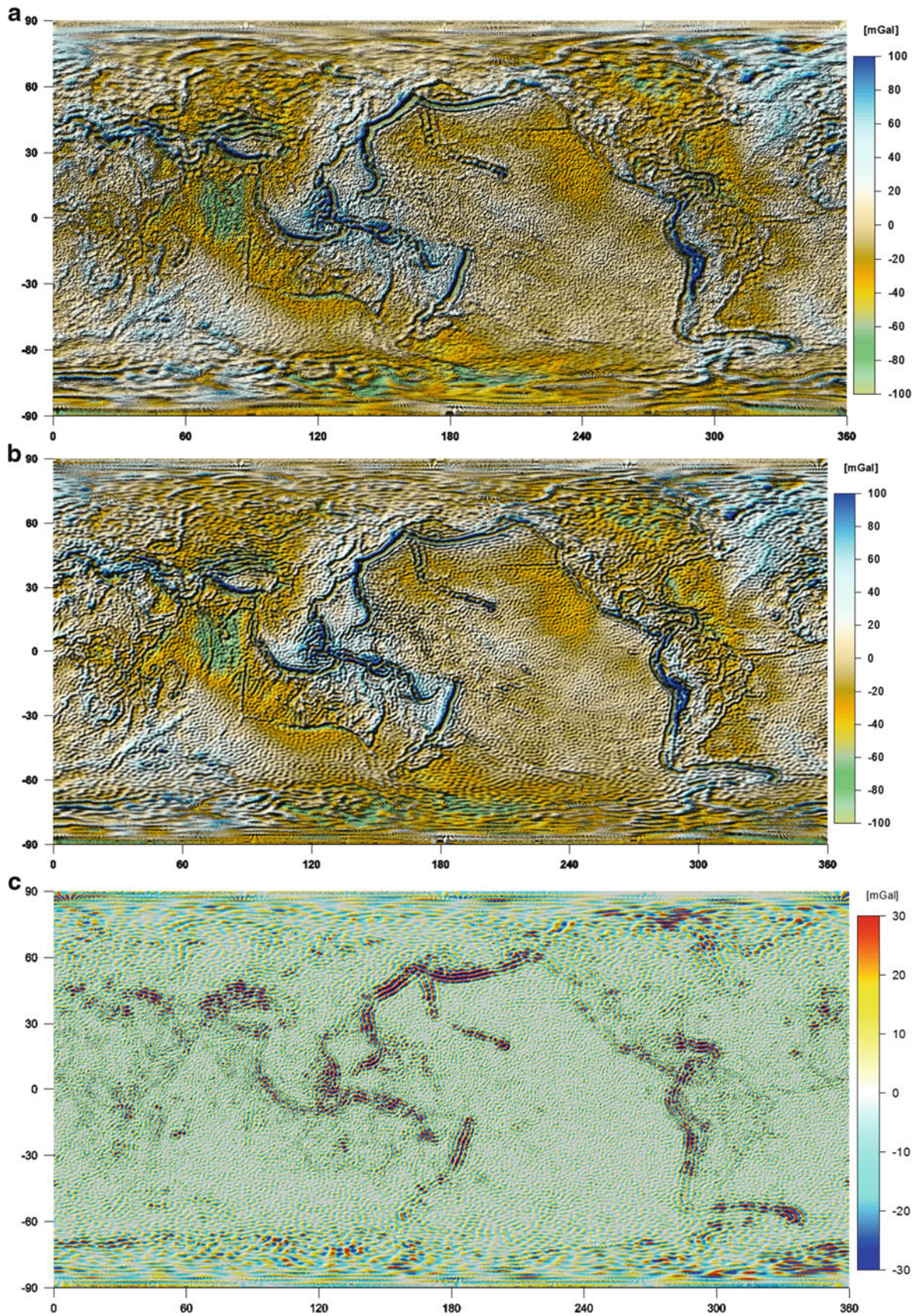
The paper demonstrates that the method of fundamental solutions is an efficient technique for global gravity field modelling. It has an advantage that the approximate solution by MFS satisfies the Laplace equation also in the computational domain with more complicated boundaries. There is no restriction to have spherical (or ellipsoidal) approximation of the Earth's surface like in the SH-based approach. On the other hand, position of the FB in the MFS approach has significant influence on the results.

In contrast to BEM, MFS as a mesh-free method does not involve integral evaluations, which make it more efficient. However, to obtain the gravity field quantities, it involves two computational steps (similarly as for SH). At first, the unknown coefficients at the source points need to be determined and then the potential or its derivatives can be evaluated.

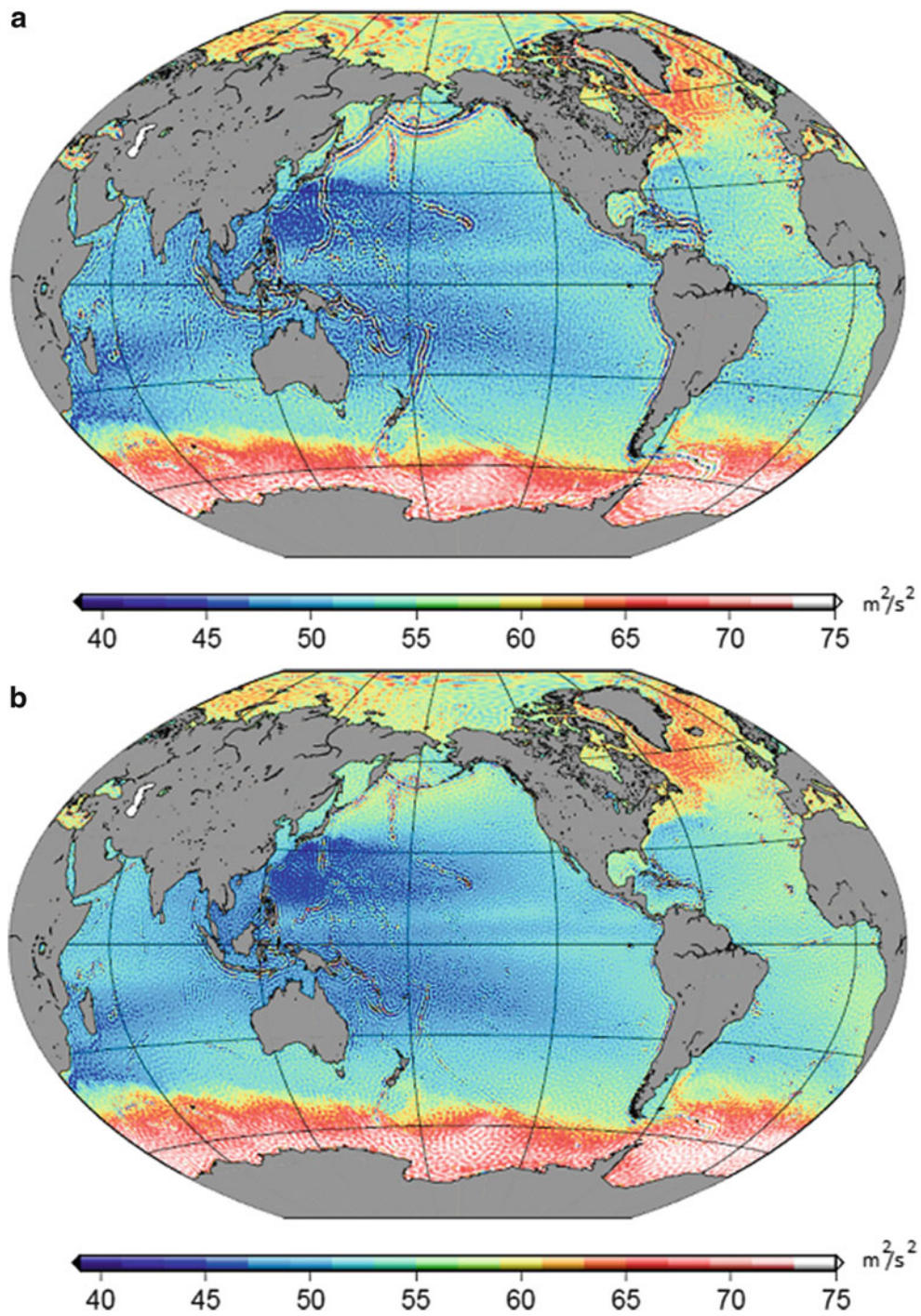
The parallel implementation of MFS and the elimination of far zones' interactions allow high-resolution modelling. In all presented numerical experiments the radial components of the gravity tensor are processed and the source points are distributed with a resolution of  $0.075^\circ$ . This yields precise global gravity field models that are in a good agreement with the SH-based GGMs, e.g., GOCO03S and EGM2008. The overall mean values of the residuals are smaller than  $0.04 \text{ m}^2\text{s}^{-2}$ . The mean values over oceans do not exceed  $0.1 \text{ m}^2\text{s}^{-2}$ . Hence, the  $W_0$  estimates evaluated from the MFS solutions differ from the ones estimated from GOCO03S or EGM2008 by less than  $0.1 \text{ m}^2\text{s}^{-2}$ . Such small differences indicate a reliability of the computed  $W_0$  estimates for the WHS realization.



**Fig. 4** The disturbing potential on the Earth's surface obtained from (a) the MFS solution with the fictitious boundary in the depth 20 km, (b) from GOCO03S model up to degree 250, and (c) the residuals between both models



**Fig. 5** The gravity disturbances on the Earth's surface obtained from (a) the MFS solution with the fictitious boundary in the depth 20 km, (b) from GOCO03S model up to degree 250, and (c) the residuals between both models



**Fig. 6** The geopotential on the DTU10 mean sea surface evaluated from (a) the MFS solution with the fictitious boundary in the depth 20 km, and (b) from GOCO03S model up to degree 250 (the constant  $62,636,800.0 \text{ m}^2/\text{s}^2$  is removed)

**Table 1**  $W_0$  estimates evaluated on the DTU10 mean sea surface model (integration area:  $82^\circ\text{S}$ – $82^\circ\text{N}$ ) from the MFS solutions with different depths of the fictitious boundaries (FB), and from the GOCO03S and EGM2008 geopotential models ( $W_0$  units:  $\text{m}^2\text{s}^{-2}$ )

FB depth (km)	MFS solution	GOCO03S (SH up to d/o 250)	EGM2008 (SH up to d/o 2160)
0	62,636,854.01	62,636,854.00	62,636,853.96
2	62,636,853.98		
5	62,636,854.02		
10	62,636,854.03		
20	62,636,854.05		
30	62,636,854.06		
100	62,636,854.22		

**Acknowledgements** The work has been supported by the grant VEGA 1/1063/11 and the project APVV-0072-11

## References

- Andersen OB (2010) The DTU10 gravity field and mean sea surface. Presented at the second international symposium of the gravity field of the Earth (IGFS2), Fairbanks, Alaska
- Becker JJ, Sandwell DT, Smith WHF, Braud J, Binder B, Depner J, Fabre D, Factor J, Ingalls S, Kim S-H, Ladner R, Marks K, Nelson S, Pharaoh A, Sharman G, Trimmer R, von Rosenberg J, Wallace G, Weatherall P (2009) Global bathymetry and elevation data at 30 arc seconds resolution: SRTM30\_PLUS. *Mar Geod* 32(4):355–371. doi:10.1080/01490410903297766
- Burša M, Kenyon S, Kouba J, Šíma Z, Vátr V, Vítek V, Vojtíšková M (2007) The geopotential value  $W_0$  for specifying the relativistic atomic time scale and a global vertical reference system. *J Geod* 81(2):103–110
- Chen W, Wang FZ (2010) A method of fundamental solutions without fictitious boundary. *Eng Anal Bound Elem* 34(5):530–32
- Chen W, Lin J, Wang F (2011) Regularized meshless method for nonhomogeneous problems. *Eng Anal Bound Elem* 35:253–257
- Čunderlík R, Minarechová Z, Mikula K (2014) Realization of WHS based on the static gravity field observed by GOCE. In: Gravity, Geoid and Height Systems, IAG Symp 141:211–220 doi:10.1007/978-3-319-10837-7\_27
- Čunderlík R, Mikula K (2010) Direct BEM for high-resolution gravity field modelling. *Stud Geophys Geod* 54(2):219–238
- Čunderlík R, Mikula K, Tunega M (2013) Nonlinear diffusion filtering of data on the earth's surface. *J Geod* 87:143–160. doi:10.1007/s00190-012-0587-y
- Dayoub N, Edwards SJ, Moore P (2012) The Gauss-Listing potential value  $W_0$  and its rate from altimetric mean sea level and GRACE. *J Geod* 86:681–694. doi:10.1007/s00190-012-1547-6
- Fan CM, Chen CS, Monroe J (2009) The method of fundamental solutions for solving convection-diffusion equations with variable coefficients. *Adv Appl Math Mech* 1:215–230
- Golberg MA, Chen CS (1994) The theory of radial basis functions applied to the BEM for inhomogeneous partial differential equations. *Bound Elem Commun* 5:57–61
- Gu Y, Chen W, Zhang J (2012) Investigation on near-boundary solutions by singular boundary method. *Eng Anal Bound Elem* 36(8):117–82
- Hon YC, Wei T (2005) The method of fundamental solution for solving multidimensional inverse heat conduction problems. *CMES Comput Model Eng Sci* 7:119–132
- Kupradze VD, Alexidze MA (1964) The method of functional equations for the approximate solution of certain boundary value problems. *USSR Comput Methods Math Phys* 4:82–126
- Mathon R, Johnston RL (1977) The approximate solution of elliptic boundary-value problems by fundamental solutions. *SIAM J Num Anal* 638–650
- Mayer-Gürr T, Rieser D, Hoeck E, Brockmann M, Schuh WD, Krasbutter I, Kusche J, Maier A, Krauss S, Hausleitner W, Baur O, Jaeggi A, Meyer U, Prange L, Pail R, Fecher T, Gruber T (2012) The new combined satellite only model GOCO03s. Presented at the GGHS-2012 in Venice, Italy, October 9–12, 2012
- Pavlis NK, Holmes SA, Kenyon SC, Factor JK (2012) The development of the earth gravitational model 2008 (EGM2008). *J Geophys Res* 117:B04406. doi:10.1029/2011JB008916
- Sánchez L (2009) Strategy to establish a global vertical reference system. In: Geodetic Reference Systems, IAG Symp 134:273–278
- Sideris M, Fotopoulos G (2012) Special issue on regional and global geoid-based vertical datums. *J Geod Sci* 2(4)

---

# Combination of GOCE Gravity Gradients in Regional Gravity Field Modelling Using Radial Basis Functions

Verena Lieb, Johannes Bouman, Denise Dettmering, Martin Fuchs, and Michael Schmidt

---

## Abstract

The satellite gravity mission GOCE measured the second-order derivatives of the Earth's gravitational potential with high accuracy. The GOCE data enrich our gravity field knowledge especially at spatial resolutions from 750 km down to 80 km. In this paper we carry out regional gravity field analysis using radial localising basis functions that permit the combination of different data types tailored to their accuracy and spectral signal content. We formulate observation equations for each individual GOCE gravity gradient as they are distinctive reflections of the gravity field and contain directional information. To optimally use the original GOCE measurements, we derive the mathematical expressions in the gradiometer reference frame. The expressions and their implementation are validated for a test area in Scandinavia by comparison with the global gravity field model GOCO03s, which yields small differences of less than  $\pm 1$  mE. The relative weighting of the observations is determined by variance component estimation. Moreover manually fixing the weights leads to smaller residuals with respect to GOCO03s, which is probably caused by systematic errors in the gradients. We demonstrate the capabilities of our method through a combination of the gradient data with terrestrial free-air anomalies. At spatial resolutions down to 40 km the terrestrial data get much larger relative weights than the GOCE data, which indicates the proper performance of the combination method.

---

## Keywords

Data combination • GOCE gravity gradients • Radial basis functions • Regional gravity field modelling • Relative weighting

---

## 1 Introduction

Equipped with a 3-axis gradiometer the satellite mission GOCE (Gravity field and steady-state Ocean Circulation Explorer) (Rummel et al. 2002) observed all second-order derivatives of the Earth's gravitational potential. The gravity gradient tensor contains the complete curvature information

of the local gravity field, with the advantage over the 1D gravity field information from GRACE (Gravity Recovery And Climate Experiment) (Tapley et al. 2004) that it can be applied in high-resolution gravity field determination (Pail et al. 2011), but also contains directional information allowing the gradients to be used for Earth interior research and for geophysical exploration (Ebbing et al. 2013).

An advantage of regional over global gravity field analysis is that one can adapt to local data availability and signal content. Well-established methods exist such as least-squares collocation (Tscherning and Arabelos 2011) or spherical splines (Eicker et al. 2007). We apply radial basis functions

---

V. Lieb (✉) • J. Bouman • D. Dettmering • M. Fuchs • M. Schmidt  
Deutsches Geodätisches Forschungsinstitut der Technischen  
Universität München (DGFI-TUM), Munich, Germany  
e-mail: [verena.lieb@tum.de](mailto:verena.lieb@tum.de)

enabling a consistent spectral combination of different observation types in order to create regional gravity fields containing maximum degree of information (Schmidt et al. 2007). The main focus of this paper lies on setting up the observation equations for the GOCE gradients in the gradiometer reference frame (GRF), which was not done so far for this method. Variance component estimation (VCE) offers the possibility of combining all six GOCE gradients in a flexible way: less accurate measurements can be down-weighted or excluded.

The GOCE measurement technique and the data set that is used are described in Sect. 2. The modelling approach itself consists of analysis and synthesis procedures which are explained in detail in Sect. 3. In Sect. 4 the results are presented and the relative weighting of the input data is discussed. We validated the regional models with a global model. Furthermore the modelling approach can be extended by combining GOCE observations with other observation types. We present an example for the combination with high resolution free-air anomalies.

## 2 Gravity Gradient Measurements from GOCE

We use the reprocessed release 2 of GOCE observations (level-2 products), available through the GOCE Virtual On-line Archive<sup>1</sup>. Three pairs of accelerometers measured the gradients in the Cartesian GRF with its  $xyz$  axes pointing approximately along-track, cross-track and in radial direction. The  $3 \times 3$  gravity gradient tensor is

$$\mathbf{V}_{ab} = \begin{bmatrix} V_{xx} & V_{xy} & V_{xz} \\ V_{yx} & V_{yy} & V_{yz} \\ V_{zx} & V_{zy} & V_{zz} \end{bmatrix} \quad (1)$$

with  $V_{xy} = V_{yx}$ ,  $V_{xz} = V_{zx}$ ,  $V_{yz} = V_{zy}$ ,  $a, b \in \{x, y, z\}$  and trace  $(\mathbf{V}_{ab}) = 0$ .  $V_{xy}$  and  $V_{yz}$  are less accurate than the other components and a rotation of the GOCE observations would reduce the accuracy in the rotated frame (Bouman 2007; Fuchs and Bouman 2011). We use observations from 02/2010 until 05/2012. The gradient errors are lowest in the measurement bandwidth (MBW) between 5 mHz and 100 mHz, above and below the MBW the errors increase rapidly. As the low part of the frequency spectrum is less accurately observed, it is removed by high-pass filtering with a cut-on frequency at the lower boundary of the MBW and filled up with model information from GOCO03s (Mayer-Gürr et al. 2012) to obtain a complete data set. Furthermore outliers and less accurate measurements have been removed.

## 3 Regional Gravity Field Modelling Approach

Our regional gravity field modelling approach uses radial basis functions that act as low-pass filters. They are related to specific frequency bands denoted as resolution levels  $j$  (Fig. 1). The basis functions can be expressed in terms of Legendre polynomials  $P_l$  (cf. Eq. (2)) developed up to a certain degree  $l = l'_j$ . This degree is related to the upper boundary of the corresponding level  $j$  with  $l'_j = 2^j - 1$ , representing the cut-off frequency of the low-pass filter. The degree is related to the spatial resolution at the Earth's surface as  $r \approx 20,000 \text{ km}/l'_j$ . Higher levels allow to model higher spatial resolutions contained in the gravity data.

In our approach we start with the choice of an appropriate level  $j = J + 1$  related to the resolution  $r$  of the input data. Next we set up the basis functions  $\phi_{J+1}$  of level  $J + 1$  which remove the high frequencies of the input data above degree  $l'_{J+1}$  (Schmidt et al. 2007). Finally we approximate gravitational potential differences  $\Delta V$  between the potential  $V$  and an appropriate global background model, in order to represent high frequency deviations for specified regions. The series expansion in terms of scaling functions  $\phi_{J+1}$  and scaling coefficients  $d_{J,q}$  reads

$$\begin{aligned} \Delta V(\mathbf{x}^p) &= \sum_{q=1}^{N_J} d_{J,q} \phi_{J+1}(\mathbf{x}^p, \mathbf{x}_q) \\ &= \sum_{q=1}^{N_J} \sum_{l=0}^{l'_{J+1}} \frac{2l+1}{4\pi} d_{J,q} \Phi_{J+1,l} \left(\frac{R}{r}\right)^{l+1} P_l(\cos \psi) \end{aligned} \quad (2)$$

for an observation point  $P(\mathbf{x}^p)$  with position vector  $\mathbf{x}^p = r \mathbf{r}^p$ . Herein  $r = |\mathbf{x}^p|$  means the radial distance and  $\mathbf{r}^p = [\cos \theta \cos \lambda, \cos \theta \sin \lambda, \sin \theta]^T$  is the unit vector depending on spherical longitude  $\lambda$  and co-latitude  $\theta$ . The number  $N_J$  of unknown scaling coefficients  $d_{J,q}$  ( $q = 1, \dots, N_J$ ) and thus the number of computation points  $Q(\mathbf{x}_q)$  on which the functions  $\phi_{J+1}$  are located depends on the level  $J + 1$ . In Eq. (2)  $\Phi_{J+1,l}$  are the Legendre coefficients,  $R$  is the mean Earth radius and  $\psi$  is the spherical distance between point  $P$  and  $Q$  (Schmidt et al. 2007). Equation (2) is given in a Terrestrial Reference Frame (TRF) in spherical coordinates, whereas the GOCE gravity gradients are measured in the Cartesian GRF. Consequently, the second-order derivatives of Eq. (2) are needed and have to be transformed into the GRF.

<sup>1</sup>eo-virtual-archive1.esa.int/Index.html



GOCE MBW												
$j$ [level]	1	2	3	4	5	6	7	8	9	10	11	12
$l$ [deg]	1	3	7	15	31	63	127	255	511	1023	2047	4095
$r$ [km]	20000	6667	2857	1333	645	317	157	78	39	20	10	5
frequency [deg]	→											

**Fig. 1** Extract of the frequency spectrum which is split into resolution levels  $j$ : upper boundary corresponds to a maximum degree  $l$ , related to the spatial resolution  $r$  at the Earth's surface. Levels where GOCE has its highest sensitivity are indicated in red (MBW)

### 3.1 Adopted Scaling Functions

The six different space dependent GOCE gravity gradients in Eq. (1) are treated as six separate measurements and thus  $K = 6$  observation equations have to be formulated. The elements  $\Delta V_{ab}$  can be expressed by

$$\Delta V_{ab} = \frac{\partial^2 \Delta V}{\partial a \partial b} = \sum_{q=1}^{N_J} d_{J,q} \phi_{J+1,ab}(\mathbf{x}^p, \mathbf{x}_q) \quad (3)$$

for level  $J + 1$  according to Eq. (2). The adopted scaling functions  $\phi_{j,ab}(\mathbf{x}^p, \mathbf{x}_q)$  read for level  $j \leq J + 1$

$$\begin{aligned} \phi_{j,xx} = & \sum_{l=0}^{l'_j} \frac{2l+1}{4\pi} \left(\frac{R}{r}\right)^{l+1} \Phi_{j,l} \\ & \cdot \left( \frac{1}{r} P_l(\cos \psi) \left(-\frac{l+1}{r}\right) + \frac{1}{r^2} \frac{\partial P_l(\cos \psi)}{\partial \theta^2} \right) \end{aligned} \quad (4)$$

$$\begin{aligned} \phi_{j,xy} = & \sum_{l=0}^{l'_j} \frac{2l+1}{4\pi} \left(\frac{R}{r}\right)^{l+1} \Phi_{j,l} \cdot \left( \frac{1}{r^2 \sin \theta} \frac{\partial P_l(\cos \psi)}{\partial \lambda \partial \theta} \right. \\ & \left. - \frac{1}{r^2 \sin^2 \theta} \frac{\partial P_l(\cos \psi)}{\partial \lambda} \right) \end{aligned} \quad (5)$$

$$\begin{aligned} \phi_{j,xz} = & \sum_{l=0}^{l'_j} \frac{2l+1}{4\pi} \left(\frac{R}{r}\right)^{l+1} \Phi_{j,l} \\ & \cdot \left( \frac{1}{r^2} \frac{\partial P_l(\cos \psi)}{\partial \theta} - \frac{1}{r} \left(-\frac{l+1}{r}\right) \frac{\partial P_l(\cos \psi)}{\partial \theta} \right) \end{aligned} \quad (6)$$

$$\begin{aligned} \phi_{j,yy} = & \sum_{l=0}^{l'_j} \frac{2l+1}{4\pi} \left(\frac{R}{r}\right)^{l+1} \Phi_{j,l} \cdot \left( \frac{1}{r} P_l(\cos \psi) \left(-\frac{l+1}{r}\right) \right. \\ & \left. + \frac{1}{r^2 \tan \theta} \frac{\partial P_l(\cos \psi)}{\partial \theta} + \frac{1}{r^2 \sin^2 \theta} \frac{\partial P_l(\cos \psi)}{\partial \lambda^2} \right) \end{aligned} \quad (7)$$

$$\begin{aligned} \phi_{j,yz} = & \sum_{l=0}^{l'_j} \frac{2l+1}{4\pi} \left(\frac{R}{r}\right)^{l+1} \Phi_{j,l} \cdot \left( \frac{1}{r^2 \sin \theta} \frac{\partial P_l(\cos \psi)}{\partial \lambda} \right. \\ & \left. - \frac{1}{r \sin \theta} \left(-\frac{l+1}{r}\right) \frac{\partial P_l(\cos \psi)}{\partial \lambda} \right) \end{aligned} \quad (8)$$

$$\begin{aligned} \phi_{j,zz} = & \sum_{l=0}^{l'_j} \frac{2l+1}{4\pi} \left(\frac{R}{r}\right)^{l+1} \Phi_{j,l} \\ & \cdot P_l(\cos \psi) \frac{(l+1)(l+2)}{r^2}. \end{aligned} \quad (9)$$

Similar expressions can be derived for other radial basis functions, as e.g. covariance functions (Tscherning 1993) or spherical splines (Eicker et al. 2007).

### 3.2 Analysis

The reduced GOCE gradients are treated as separate observations assuming that we have no error correlations. The observation equation reads for one tensor element  $\Delta V_{ab}$ , observed at the observation points  $\mathbf{x}^p$  with  $p \in \{1, \dots, P\}$  according to Eq. (3) and considering the measurement error  $e_{ab}$

$$\Delta V_{ab}(\mathbf{x}^p) + e_{ab}(\mathbf{x}^p) = \boldsymbol{\phi}_{ab}^T(\mathbf{x}^p) \mathbf{d}_J. \quad (10)$$

$\boldsymbol{\phi}_{ab}$  is the  $N_J \times 1$  vector of modified scaling functions according to Eqs. (4)–(9). In the analysis step we use the Shannon scaling function with the Legendre coefficients  $\Phi_{J+1,l}^{\text{SHA}} = 1$ , which is an ideal low-pass filter up to degree  $l'_{J+1}$  (Schmidt et al. 2007). Rotating the resulting expressions into GRF leads to the observation equations of the tensor in GRF. The  $N_J \times 1$  vector  $\mathbf{d}_J = [d_{J,1}, \dots, d_{J,N_J}]^T$  of scaling coefficients is then estimated by VCE as will be briefly explained in the following. We collect all measurements of a particular gravity gradient, so that each observation group  $\Delta \mathbf{v}_k$  with  $k \in \{1, \dots, K\}$  represents a  $P \times 1$  vector of the measurements  $\Delta V_{ab}$  and  $\boldsymbol{\phi}_k$  represents the corresponding

$P \times N_J$  matrix of scaling functions:

$$\begin{bmatrix} \Delta \mathbf{V}_{xx} \\ \Delta \mathbf{V}_{xy} \\ \Delta \mathbf{V}_{xz} \\ \Delta \mathbf{V}_{yy} \\ \Delta \mathbf{V}_{yz} \\ \Delta \mathbf{V}_{zz} \end{bmatrix} = \begin{bmatrix} \Delta \mathbf{v}_1 \\ \Delta \mathbf{v}_2 \\ \Delta \mathbf{v}_3 \\ \Delta \mathbf{v}_4 \\ \Delta \mathbf{v}_5 \\ \Delta \mathbf{v}_6 \end{bmatrix} = \Delta \mathbf{v} \quad \text{and} \quad \begin{bmatrix} \phi_1 \\ \phi_2 \\ \phi_3 \\ \phi_4 \\ \phi_5 \\ \phi_6 \end{bmatrix} = \phi. \quad (11)$$

For the combination of GOCE gradient observations with further measurement techniques the vector  $\Delta \mathbf{v}$  can be extended by other observation groups  $\Delta \mathbf{v}_k$  with  $k > K$ . The stochastic part is formulated as

$$\mathbf{D} \begin{pmatrix} \Delta \mathbf{v}_1 \\ \Delta \mathbf{v}_2 \\ \Delta \mathbf{v}_3 \\ \Delta \mathbf{v}_4 \\ \Delta \mathbf{v}_5 \\ \Delta \mathbf{v}_6 \\ \boldsymbol{\mu}_d \end{pmatrix} = \begin{bmatrix} \sigma_1^2 \mathbf{P}_1^{-1} & \mathbf{0} & \dots & \mathbf{0} \\ \mathbf{0} & \sigma_2^2 \mathbf{P}_2^{-1} & \dots & \mathbf{0} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{0} & \mathbf{0} & \dots & \sigma_d^2 \boldsymbol{\Sigma}_d \end{bmatrix}. \quad (12)$$

$\mathbf{D}$  is the covariance matrix,  $\mathbf{P}_k$  is the  $P \times P$  weighting matrix of the observation vector  $\Delta \mathbf{v}_k$ . Note, the background model is introduced as additional observation group to avoid singularity problems. Referred to Schmidt et al. (2007) the vector  $\boldsymbol{\mu}_d$  contains the expectation values of the coefficients after subtracting the background model and  $\boldsymbol{\Sigma}_d$  is the corresponding  $N_J \times N_J$  covariance matrix. The variance components (VC)  $\sigma_k^2$  and  $\sigma_d^2$  are determined iteratively according to Koch and Kusche (2002). With the estimated VCs the estimated coefficients  $\hat{\mathbf{d}}_J$  result in

$$\hat{\mathbf{d}}_J = \left( \sum_{k=1}^6 \frac{1}{\hat{\sigma}_k^2} \boldsymbol{\phi}_k^T \mathbf{P}_k \boldsymbol{\phi}_k + \frac{1}{\hat{\sigma}_d^2} \boldsymbol{\Sigma}_d^{-1} \right)^{-1} \cdot \left( \sum_{k=1}^6 \frac{1}{\hat{\sigma}_k^2} \boldsymbol{\phi}_k^T \mathbf{P}_k \Delta \mathbf{v}_k + \frac{1}{\hat{\sigma}_d^2} \boldsymbol{\mu}_d \right). \quad (13)$$

The estimated covariance matrix of the coefficients reads  $\mathbf{Q}_{dd} = \left( \sum_{k=1}^6 \frac{1}{\hat{\sigma}_k^2} \boldsymbol{\phi}_k^T \mathbf{P}_k \boldsymbol{\phi}_k + \frac{1}{\hat{\sigma}_d^2} \boldsymbol{\Sigma}_d^{-1} \right)^{-1}$ .

### 3.3 Synthesis

For the synthesis step we set up the series expansion (3) in terms of Blackman scaling functions  $\phi_{J+1,ab}$ , characterized by the Legendre coefficients  $\Phi_{J+1,l}^{\text{BLA}}$  (Schmidt et al. 2007). Compared with the Shannon kernel these functions act also band-limiting as low-pass filters up to degree  $l'_{J+1}$  according to Fig. 1, but with a smoother declining behaviour. Consequently, in the spatial domain the oscillations and sidelobes

of the Blackman functions are much smaller. Thus erroneous edge effects are significantly reduced. Inserting the estimated coefficients  $\hat{\mathbf{d}}_J$  (cf. Eq. (13)) into Eq. (3) and using Eqs. (4)–(9) with  $\Phi_{J+1,l}^{\text{BLA}}$  yield the estimated gradients of the reduced gravitational potential.

## 4 Numerical Investigations

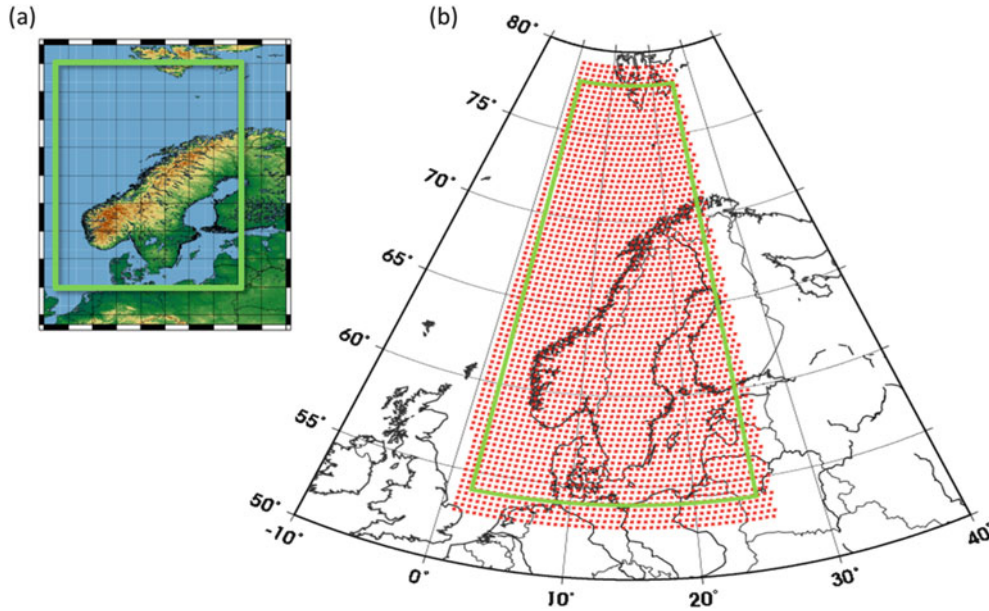
### 4.1 Study Area and Modelling Parameters

We study the Scandinavian region with an extent of  $2^\circ$  to  $25^\circ$  in longitude and  $54^\circ$  to  $78^\circ$  in latitude, see Fig. 2a. The frequency part where GOCE measures with its highest sensitivity can be seen from Fig. 1. It is highlighted in red and indicates a spatial resolution down to  $\sim 80$  km. Level  $j = 8$  is the maximum level which is completely located within the sensitivity domain of GOCE, whereas the upper part of level 9 contains a lot of noise so that only the low frequencies of  $j = 9$  deliver significant information. For our numerical investigations we consequently use the modelling approach up to level  $J + 1 = 8$ . The level depending computation points of the scaling functions can be seen in Fig. 2b (red dots). The computation area has a larger extent than the modelling area (green bordered) to diminish edge effects. The observation area containing the GOCE satellite tracks has an extension in-between both margins. Furthermore the data set is reduced by the global background model GOCO03s up to maximum degree and order 250 following Eq. (3). We used exactly the same model as for filling up the low frequencies to be consistent. The resolution of GOCO03s reaches nearly to the modelling resolution at level 8, so that most of the signal is reduced and only small deviations remain which are approximated in the estimation process.

Inserting the reduced GOCE gradients in the observation equation (10), assuming that the measurement errors are uncorrelated and have the same accuracies within an observation group  $k$ , allows us to introduce identity matrices for the weighting matrices  $\mathbf{P}_k$  in Eq. (12). As prior information we use the same model as the background model (GOCO03s). Consequently we assume that the  $N_J \times 1$  vector  $\boldsymbol{\mu}_d$  is equal to  $\mathbf{0}$  and the covariance matrix  $\boldsymbol{\Sigma}_d$  corresponds to the identity matrix.

### 4.2 Gradient Grids

As output from the synthesis procedure we obtain approximation signals which can be expressed as any functional of the gravitational potential (e.g. geoid undulations  $N$ , gravity anomalies  $\Delta g$ ). Restoring the background model, subtracting the normal potential from the reference ellipsoid



**Fig. 2** (a) Geographical location of the test area Scandinavia (green bordered) with altitude encoding topography. (b) Distribution of the grid points (red dots)

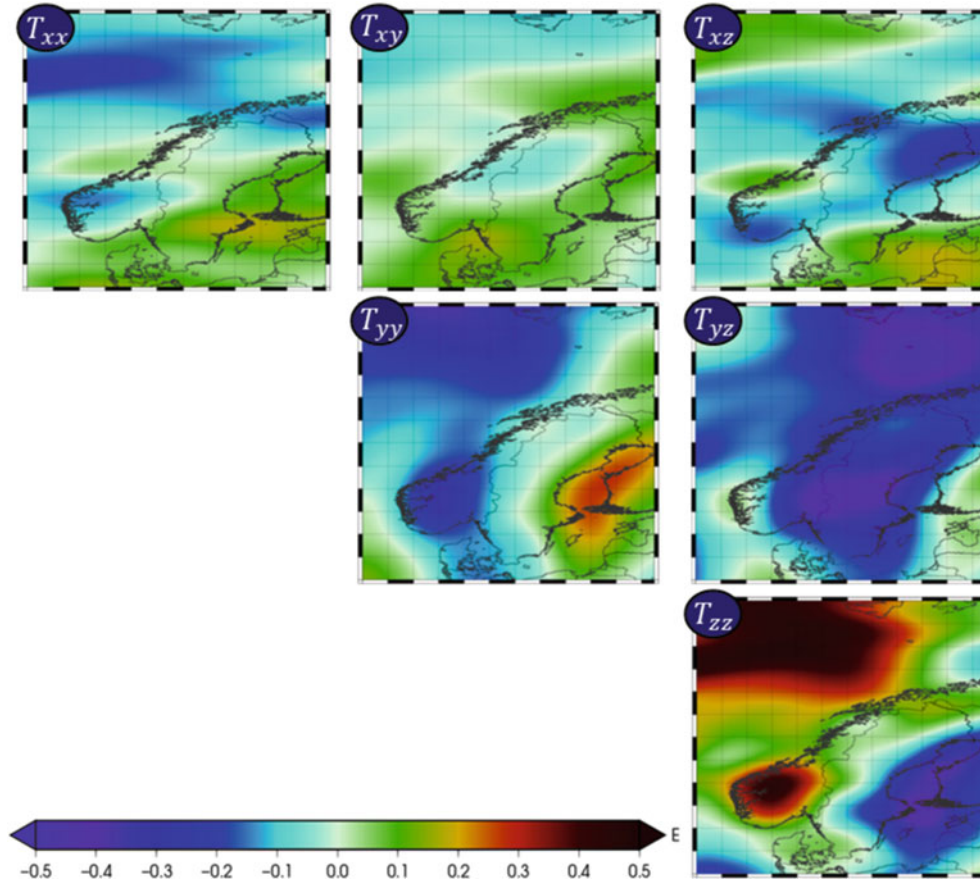
WGS84 and computing the second-order derivatives  $V_{ab}$  lead to the gradients of the disturbing potential  $T_{ab}$  for all combinations of the  $xyz$  Cartesian coordinates. Figure 3 shows the results according to the  $xyz$  tensor arrangement in a local-north-oriented frame with its axes pointing north-, west- and upwards. The modelling height corresponds to the mean GOCE orbit height of 270 km within this region. The gradients of the disturbing potential show clearly different structures depending on the different spatial directions. As expected, the radial  $zz$  component pointing directly along the field line of the Earth's gravitational potential has the largest magnitudes between  $\pm 0.5$  E. The sum of the diagonal elements should be zero according to the Laplace condition  $\text{trace}(\mathbf{T}_{ab}) = 0$ . The trace criteria gives values which are 3 orders of magnitude smaller than the signals of the single components. Considering the modelling accuracy, which depends on edge effects, oscillations of the scaling functions, smoothing and interpolation effects, the Laplace condition is therefore fulfilled. The approximation signals without restoring the background model GOCO03s vary between  $\pm 1$  mE containing additional signal to the global model but also errors from the regional approach.

### 4.3 Analysis of VCE

Variance component estimation provides a flexible tool for relative weighting of different observation groups by combining them at the level of observation equations. A large VC  $\sigma_k^2$  means hereby a low relative weight of the observation

group  $\Delta \mathbf{v}_k$  in Eq. (13). Table 1, col. (a) lists the orders of magnitude of the iteratively estimated VCs with reference to  $V_{zz}$ . The diagonal component  $V_{xx}$  obtains the same weight as the radial pointing  $V_{zz}$ . It has a smaller signal content, but also a twice smaller noise level, so that similar weighting seems to be appropriate. The less accurate components  $V_{xy}$  and  $V_{yz}$  are down-weighted by 2 and 5 orders of magnitude, respectively. Thus some information from  $V_{xy}$  is still present in the solutions while the influence of  $V_{yz}$  is negligible. The prior information is down-weighted by 2 orders of magnitude indicating that it contributes also to the output grids. The down-weighting is justified in the errors of the long wavelengths part of GOCO03s which cannot be accounted for. Against the expectation that the 4 accurate GOCE gradients should have comparable weights,  $V_{xz}$  gets a lower weight signifying that this gradient component is less accurate than the diagonal elements. We assume that this effect is specific for the Scandinavian region, as studies in other regions deliver similar VCs for the four components.

$V_{yy}$  gets the same weight as the diagonal components  $V_{xx}$  and  $V_{zz}$ , but as our test area is located near the North pole we further have to deal with systematic errors in this component (Bouman et al. 2011; Bouman and Fuchs 2012). In a second computation we thus manually fix the relative weights (Table 1, col. (b)): the VCs of  $V_{xx}$  and  $V_{zz}$  are adapted to the estimated values, but  $V_{yy}$  is down-weighted by 5 orders of magnitude. We assume a noise behaviour comparable with that of  $V_{yz}$  obtained in the estimated case.  $V_{xy}$  and  $V_{yz}$  are additionally down-weighted. Using those fixed weights we apply least squares estimation within a



**Fig. 3** Gravity gradient grids of the second-order derivatives of the disturbing potential modelled from GOCE gradient measurements at 270 km height

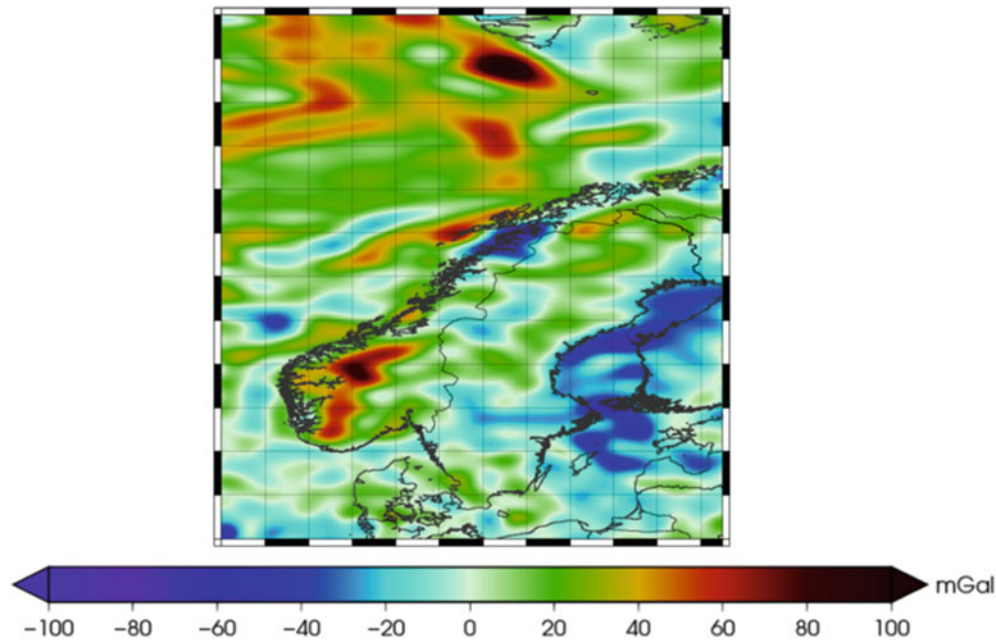
**Table 1** Relative weighting of observations. Given are the orders of magnitude of the related VCs  $\sigma_k^2$  [-]

Observation	(a) est	(b) fix	(c) fix
GOCO03s	$10^2$	$10^2$	$10^2$
$V_{xx}$	1	1	1
$V_{xy}$	$10^2$	$10^{13}$	$10^8$
$V_{xz}$	$10^3$	$10^2$	$10^2$
$V_{yy}$	1	$10^5$	$10^5$
$V_{yz}$	$10^5$	$10^{13}$	$10^8$
$V_{zz}$	1	1	1
FA			$10^{-2}$

Gauss-Markov model. We compare the results from both weighting strategies with the global GOCO03s model. The mean standard deviation of the difference grids decreases from 0.1 mE (for estimated VCs) down to 0.03 mE when setting a lower relative weight for  $V_{yy}$  component. In our study area the differences decline especially in the northern part at latitudes  $> 70^\circ$ . We conclude that this might be due to the down-weighted impact of  $V_{yy}$  and plan to investigate further studies.

#### 4.4 Combination with Free-Air Anomalies

Finer structures can be modelled by combining GOCE gravity gradients with high-resolution data sets such as free-air anomalies (FA). Figure 4 shows gravity anomalies with variations between  $\pm 100$  mGal at the Earth's surface obtained from a combination at level 9 ( $l' = 511$ ). Compared with GOCE, the FA data set (Olesen et al. 2010a,b) contains detailed information from altimetry, terrestrial and shipborne gravimetry. The FA data therefore get a higher weight than the GOCE  $V_{zz}$  gradients (2 orders of magnitude, determined with VCE). Table 1, col. (c) shows the corresponding manually fixed weights of all input data sets. The lower frequency domain of the solution is stabilised by the GOCE observations. For areas where high-frequency data are available, our regional gravity field modelling approach offers the opportunity to combine data sets which are sensitive to different frequency domains by VCs such, that the data with the highest signal content up to a specific level  $j$  contribute the most. In contrast to other gravity analysis techniques weights can be introduced individually for each resolution level.



**Fig. 4** Combination of GOCE gravity gradients and free-air anomalies at level 9 ( $j' = 511$ ) using manually fixed relative weights, Table 1, col. (c)

## 5 Summary

We derived the observation equations for GOCE gravity gradient measurements in a regional gravity field modelling approach using radial basis functions. Our aim was to use the original GOCE gravity gradients in the GRF to maintain the precision of the four accurate components. The resulting gradient grids show different structures that give information of the Earth's gravitational potential depending on different spatial directions. This advantage might further be used for research on the Earth's interior and for geophysical exploration (Ebbing et al. 2013). A validation of our regional gradients grids with GOCO03s gives differences that are smaller than  $\pm 1$  mE, which confirms that our method works properly. We also found that the use of VCs for the automatic estimation of the relative weights of the different gradient components may not be optimal. Manually down-weighting the less accurate  $V_{xy}$  and  $V_{yz}$  components as well as down-weighting the regionally less accurate  $V_{yy}$  component, significantly reduces the differences to the global GOCO03s model. Thus the optimal combination of the gradient data sets requires further study, especially in the presence of systematic errors as may be the case for  $V_{yy}$ . Moreover we demonstrated a combination of the GOCE gravity gradients with high-resolution FAs. The latter enable to model more detailed structures at higher resolution levels. With this additional information radial basis functions might offer the possibility to enrich and supplement global gravity fields in specified regions.

**Acknowledgements** The authors want to thank the European Space Agency (ESA) for funding the project GOCE+ GeoExplore. Further we are grateful to P. Willis, C.C. Tscherning and two anonymous reviewers for their valuable input.

## References

- Bouman J (2007) Alternative method for rotation to TRF. GO-TN-HPF-GS-0193. ESA-ESTEC, Noordwijk
- Bouman J, Fuchs M (2012) GOCE gravity gradients versus GOCE gravity field models. *Geophys J Int* 189(2):846–850. doi:10.1111/j.1365-246X.2012.05428.x
- Bouman J, Fiorot S, Fuchs M, Gruber T, Schrama E, Tscherning CC, Veicherts M, Visser P (2011) GOCE gravitational gradients along the orbit. *J Geodesy* 85(11):791–805. doi:10.1007/s00190-011-0464-0
- Ebbing J, Bouman J, Fuchs M, Lieb V, Haagmans R, Meekes S, Abdul Fattah R (2013) Advancements in satellite gravity gradient data for crustal studies. *Lead. Edge* 32:900–906. doi 10.1190/le32080908.1
- Eicker A, Mayer-Gürr T, Ilk KH (2007) A global CHAMP gravity field by merging of regional refinement patches. In: *Proceedings Joint CHAMP/GRACE Science Team Meeting 2004*, 0000:0001–5
- Fuchs M, Bouman J (2011) Rotation of GOCE gravity gradients to local frames. *Geophys J Int* 187(2):743–753. doi 10.1111/j.1365-246X.2011.05162.x
- Koch KR, Kusche J (2002) Regularization of geopotential determination from satellite data by variance components. *J Geodesy* 76:259–286. doi 10.1007/s00190-002-0245-x
- Mayer-Gürr T, Rieser D, Höck E, Brockmann JM, Schuh W-D, Krasbutter I, Kusche J, Maier A, Krauss S, Hausleitner W, Baur O, Jäggi A, Meyer U, Prange L, Pail R, Fecher T, Gruber T (2012) The new combined satellite only model GOCO03s. *Symposium on gravity, geoid and height systems (GGHS)*, Venice, Italy, October 9–12 2012, Gravity Observation Combination (www.goco.eu)
- Olesen O, Brünner M, Ebbing J, Gellei J, Gernigon L, Koziel J, Lauritsen T, Myklebust R, Sand M, Solheim D, Usov S (2010a) New

- aeromagnetic and gravity compilations from Norway and adjacent areas – methods and applications. In: Petroleum geology conference series, vol 7, pp 559–586
- Olesen O, Ebbing J, Gellein J, Kihle O, Myklebust R, Sand M, Skilbrei JR, Solheim D, Usov S (2010b) Gravity anomaly map, Norway and adjacent areas. Geological Survey of Norway
- Pail R, Bruinsma S, Migliaccio F, Förste C, Goiginger H, Schuh WD, Hoeck E, Reguzzoni M, Brockmann JM, Abrikosov O, Veicherts M, Fecher T, Mayrhofer R, Krasbutter I, Sansò F, Tscherning CC (2011) First GOCE gravity field models derived by three different approaches. *J Geodesy* 85:819–843. doi:10.1007/s00190-011-0467-x
- Rummel R, Balmino G, Johannessen J, Visser P, Woodworth P (2002) Dedicated gravity field missions - principle and aims. *J Geodyn* 33(1–2):3–20. doi:10.1016/S0264-3707(01)00050-3
- Schmidt M, Fengler M, Mayer-Gürr T, Eicker A, Kusche J, Sánchez L, Han S-C (2007) Regional gravity modelling in terms of spherical base functions. *J Geodesy* 81(1):17–38. doi:10.1007/s00190-006-0101-5
- Tapley BD, Bettadpur S, Watkins M, Reigber C (2004) The gravity recovery and climate experiment: mission overview and early results. *Geophys Res Lett* 31:L09607. doi:10.1029/2004GL019920
- Tscherning CC (1993) Computation of covariances of derivatives of the anomalous gravity potential in a rotated reference frame. *Manuscr Geodaet* 18(3):115–123
- Tscherning CC, Arabelos DN, Gravity anomaly and gradient recovery from GOCE gradient data using LSC and comparisons with known ground data. In: Proceedings of 4th International GOCE user workshop, ESA SP-696, 31 March–1 April 2011

---

# Rosborough Representation in Satellite Gravimetry

Nico Sneeuw and Mohammad A. Sharifi

---

## Abstract

Rosborough representations are known from and used in satellite altimetry. The radial orbit perturbation is represented in spherical coordinates and separated into a geographically mean (ascending plus descending) and a geographically variable (ascending minus descending) part. This principle is easily generalized to any functional of the gravitational field, observed along the orbit. Therefore, this type of representation can be used for gravity field recovery from missions like GRACE and GOCE.

We describe here the nature of the Rosborough formalism in terms of forward and backward rotation of spherical harmonics. A more practical derivation is subsequently given by transforming the orbital coordinates back into spherical latitude  $\phi$  and longitude  $\lambda$ , given a nominal inclination  $I$ . This transformation is greatly alleviated by making use of complex-valued functions, as opposed to the binomial series in Rosborough's original formulation.

---

## Keywords

Rosborough representation • Spaceborne gravimetry

---

## 1 Introduction

The original Rosborough problem in satellite altimetry (Rosborough 1986) deals with the question, which part of the radial orbit error  $\Delta r$  is visible in the cross-over difference at the location  $(\phi, \lambda)$ . This automatically led Rosborough to the question, how to represent (radial) orbit perturbations in spherical coordinates. Since the sphere is virtually sampled twice from satellite orbit, both from ascending and from descending tracks, a separation follows naturally:

$$\Delta r(\phi, \lambda) \leftrightarrow \Delta r^{a/d}(\phi, \lambda) \leftrightarrow \Delta r^{m/v}(\phi, \lambda).$$

---

N. Sneeuw (✉)  
Institute of Geodesy, Universität Stuttgart, Stuttgart, Germany  
e-mail: [sneeuw@gis.uni-stuttgart.de](mailto:sneeuw@gis.uni-stuttgart.de)

M.A. Sharifi  
Department of Surveying and Geomatics Engineering, University  
College of Engineering, University of Tehran, Tehran, Iran  
e-mail: [sharifi@ut.ac.ir](mailto:sharifi@ut.ac.ir)

The more general problem then consists in expressing any along-track gravitational observable in spherical coordinates, i.e. turn the time-wise approach into a space-wise approach:

$$f(t) \leftrightarrow f(u, \Lambda, I, r) \leftrightarrow f^{a/d}(\phi, \lambda) \leftrightarrow f^{m/v}(\phi, \lambda).$$

The indices a/d refer to ascending and descending orbits, whereas m/v refer to mean and variable, to be explained furtheron.

On a historical note one should mention that the question of a spatial representation of radial orbit errors was investigated around the same time by Engelis (1987). However, the name "Rosborough" seems to have stuck in literature. As such, we will talk about Rosborough problem, approach, representation and functions in the following.

Although the Rosborough representation is space-wise, it uses transfer (or sensitivity) coefficients that are rooted in the time-wise approach. Thus the advantages from both approaches can potentially be combined to yield a powerful gravity recovery method. This, however, has never been achieved beyond the description of radial orbit errors in

altimetry. One of the main hurdles for wider implementation has most likely been the use of real-valued variables, leading to quite involved algorithms, e.g. Balmino (1993). The formulation was compounded by the need of developing trigonometric functions into binomial series and to multiply such series. Moreover, the  $p$ -index, as introduced by Kaula (1966), related to the azimuthal wavenumber  $k$  by  $k = l - 2p$ , is not conducive to understanding the pertinent formulas. In terms of real-valued quantities, the algorithm was optimized in an internal technical report by Bosch (1997).

As mentioned above, the method was developed and came to fruition for the analysis of satellite altimetry. For this field of application, the method was restricted to cross-over locations – although including the concept of multiple-satellite cross-overs, e.g. Klokočník et al. (1995) – and to the gravitational functional of type *radial orbit perturbation*. Sneeuw (2003) generalized the Rosborough approach and extended it to the analysis of other along-orbit gravitational functionals. At the same time a major algorithmic improvement was presented in Sneeuw (2003) by making use of complex-valued trigonometric functions, completely eliminating the need for binomial series developments (and their products). As a result Rosborough functions could be computed in a fast and stable way up to high degree. In the era of CHAMP (Reigber et al. 2005), GRACE (Tapley et al. 2004) and GOCE (ESA 1999; Rummel 2011), however, the approach did not seem to have been used. More recently, Sharifi et al. (2013a) demonstrated by closed-loop simulation that the Rosborough approach can work with GOCE gravity gradient data. The first successful proof-of-concept, in which real GOCE data were analyzed with the Rosborough approach, was given by Sharifi et al. (2013b).

Beyond being merely an alternative formulation, the Rosborough representation separates ascending from descending arc information in a natural way, which is the reason why it was developed for cross-over analysis in satellite altimetry in the first place. Such separation is a helpful property when performing error analysis in cases where the error behaves differently between ascending and descending tracks. Examples would be (a) time tag errors in altimetry that map to different errors on ascending and descending tracks, (b) analysis of the tidal aliasing error, where the tidal phase between descending and ascending track are distinct, or (c) ionosphere-induced orbit tracking errors due to dawn-dusk orbit geometry, which may result into the geomagnetic equator to become visible in GOCE gravity field recoveries.

However, in this contribution we do not engage in discussing the merits of the Rosborough formalism relative to more standard (time-wise, space-wise) approaches. Instead we want to elucidate the nature of the Rosborough

representation and to focus on the algorithmic aspects, that were only very briefly touched upon in Sneeuw (2003).

## 2 Transforming into Orbit and Back

The nature of Rosborough functions will be revealed by the following derivation:

$$\begin{aligned}
 &V(\phi, \lambda) : \text{potential on the sphere} \\
 &\quad \Downarrow \text{transformation} \\
 &V(u, \Lambda, I, r) : \text{potential along the orbit} \\
 &\quad \text{(Kaula representation)} \\
 &\quad \Downarrow \text{transfer } H_{lmk}(r, I) \\
 &f(u, \Lambda, I, r) : \text{functional along the orbit} \\
 &\quad \Downarrow \text{reverse transformation} \\
 &f^{a/d}(\phi, \lambda, r, I) : \text{functional on 2 spheres} \\
 &\quad \text{(ascending \& descending)}
 \end{aligned}$$

The spherical latitude  $\phi$  and longitude  $\lambda$  are transformed into orbital coordinates  $u$ ,  $I$  and  $\Lambda$ , to be explained below. The above transformations refer to rotations of the coordinate system into the orbital plane (and back). As a consequence, the spherical harmonics must transform accordingly, e.g. Sneeuw (1992). Note that this is not the type of derivation followed by Rosborough (1986) or Engelis (1987), which will be the topic of Sect. 3.

### 2.1 The Geopotential on the Sphere

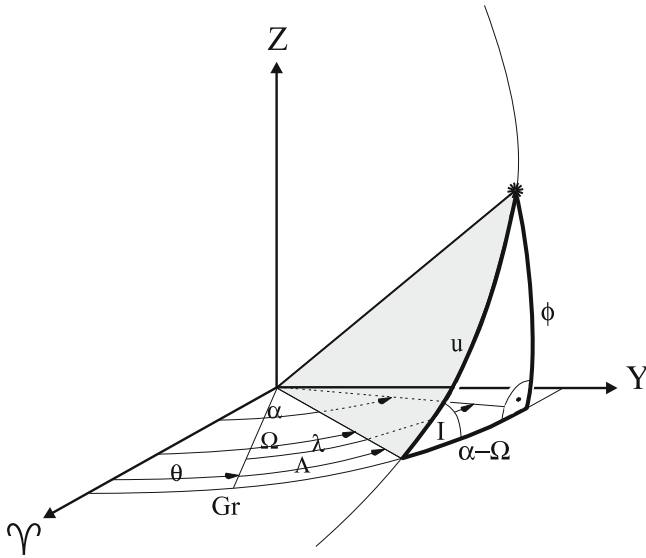
Starting point is a spherical harmonic (SH) series of the geopotential. For reasons of compactness, complex-valued quantities are employed. Moreover the dimensioning factor  $GM/r$  and the upward continuation term  $(R/r)^l$  are suppressed. It is assumed that spherical harmonics and coefficients are fully normalized, although this is not indicated notationally.

$$V(\phi, \lambda) = \sum_l \sum_m K_{lm} Y_{lm}(\phi, \lambda), \quad (1)$$

$$\text{with: } Y_{lm}(\phi, \lambda) = P_{lm}(\phi) e^{im\lambda},$$

and complex spherical harmonic coefficients  $K_{lm}$ .





**Fig. 1** Spherical triangle between the satellite, its footpoint on the equator and the instantaneous ascending node

## 2.2 The Geopotential Along the Orbit

It was shown in, e.g., Sneeuw (1992) how to arrive at a complex-valued expression in orbital variables by applying the rotation sequence  $R_z(u)R_x(I)R_z(\Lambda)$ . Use is made of the orbital coordinates  $u = \omega + \nu$  (argument of latitude),  $I$  (inclination) and  $\Lambda = \Omega - \text{GAST}$  (longitude of ascending node, with  $\Omega$  being the right ascension of the ascending node). The Greenwich actual sidereal time GAST is denoted by the variable  $\theta$  in Fig. 1. For the spherical harmonics this implies:

$$Y_{lm}(\phi, \lambda) = \sum_k D_{lmk}(\Lambda, I, u) Y_{lk}(\phi', \lambda') \quad (2)$$

$$\text{with } D_{lmk}(\Lambda, I, u) = i^{k-m} d_{lmk}(I) e^{i(ku + m\Lambda)}, \quad (3)$$

which makes use of Wigner's rotation symbols  $D_{lmk}$  and  $d_{lmk}$ . Application to (1) returns:

$$V(\phi', \lambda') = \sum_l \sum_m \sum_k K_{lm} i^{k-m} d_{lmk}(I) \cdot e^{i(ku + m\Lambda)} Y_{lk}(\phi', \lambda'). \quad (4)$$

The rotations are chosen such that the new equatorial plane ( $x'y'$ -plane) coincides with the orbital plane, hence  $\phi' = 0$ . At the same time, the last rotation  $R_3(u)$  assures that the new  $x'$ -axis will always point to the satellite, i.e.  $\lambda' = 0$ , such

that  $Y_{lk}(\phi', \lambda') = P_{lk}(0)$ . This gives rise to the inclination functions

$$F_{lmk}(I) = i^{k-m} d_{lmk}(I) P_{lk}(0). \quad (5)$$

They were introduced by Kaula (1966), although not derived along the above lines. For algorithmic aspects of efficient inclination function computation, the reader is referred to Kostecký et al. (1986) and Goad (1987) or Sneeuw (1992) among others.

With the coordinates  $\phi'$  and  $\lambda'$  now defunct it is better to express the potential in the orbital coordinates  $u, I, \Lambda$  (with  $r$  suppressed, but not defunct):  $V(u, I, \Lambda)$ .

## 2.3 Functionals Along the Orbit

The use of the inclination function is not followed here, though, since it obscures the coordinates  $\phi'$  and  $\lambda'$ . A functional of the geopotential is obtained by applying a *specific transfer*  $h_{lmk}$ . Thus from (4) one obtains the expression:

$$f(\phi', \lambda') = \sum_{l,m,k} h_{lmk} K_{lm} i^{k-m} d_{lmk}(I) \cdot e^{i(ku + m\Lambda)} Y_{lk}(\phi', \lambda'). \quad (6)$$

In this equation,  $h_{lmk}$  may contain dimensioning and upward continuation again. A collection of transfer coefficients, relevant to spaceborne gravimetric observables, is provided in Sneeuw (2000).

## 2.4 Functionals on the Sphere(s)

If one can rotate from the original sphere into the orbital system, one can surely rotate back to the sphere again. Thus, rotation (2) is reversed now:

$$Y_{lk}(\phi', \lambda') = \sum_k D_{lkp}(-u, -I, -\Lambda) Y_{lp}(\phi'', \lambda'') \quad (7)$$

$$= \sum_k i^{p-k} d_{lkp}(-I) e^{i(-ku - p\Lambda)}.$$

Naturally,  $\phi'' = \phi$  and  $\lambda'' = \lambda$ . Inserting all this in (6) yields a SH series expression of the functional:

$$f(\phi, \lambda) = \sum_{l,m,k,p} h_{lmk} K_{lm} i^{p-m} d_{lmk}(I) d_{lkp}(-I) \cdot e^{i(m-p)\Lambda} Y_{lp}(\phi, \lambda) \quad (8)$$

Now  $\Lambda$  is the longitude of the ascending node. Therefore, if  $\alpha$  denotes the right ascension of the satellite, we have  $\alpha - \Omega = \lambda - \Lambda$ , cf. Fig. 1. With  $\Lambda = \lambda - (\alpha - \Omega)$ , the latter two terms can be recast into:

$$\begin{aligned} & e^{i(m-p)\Lambda} Y_{lp}(\phi, \lambda) \\ &= e^{i(m-p)(\lambda - (\alpha - \Omega))} P_{lp}(\sin \phi) e^{ip\lambda} \\ &= e^{i(p-m)(\alpha - \Omega)} P_{lp}(\sin \phi) e^{im\lambda}. \end{aligned}$$

Performing the summations over  $p$  and  $k$  subsequently gives the expressions:

$$f(\phi, \lambda) = \sum_l \sum_m K_{lm} Q_{lm}(\phi, I) e^{im\lambda}, \quad (9)$$

$$Q_{lm}(\phi, I) = \sum_k h_{lmk} d_{lmk}(I) \Phi_{mk}(\phi, I), \quad (10)$$

$$\begin{aligned} \Phi_{mk}(\phi, I) &= \sum_p i^{p-m} d_{lkp}(-I) P_{lp}(\sin \phi) \cdot \\ &\cdot e^{i(p-m)(\alpha - \Omega)}. \end{aligned} \quad (11)$$

Expression (9) represents the functional in a series, similar to a spherical harmonic development. This type of series is referred to as *Rosborough representations*. Instead of Legendre functions, every functional  $f$  would have its own *Rosborough function*  $Q_{lm}$ . Since  $\Phi_{mk}$  contains terms in  $(\alpha - \Omega)$ , it must be separated in two functions, one for ascending tracks, the other for descending ones, cf. next section. Consequently the same holds true for  $Q_{lm}$  and the functional itself. This is treated more explicitly in the next section.

## 2.5 Isotropic Transfer

An isotropic transfer coefficient  $h_{lmk}$  does not depend on  $m$  or  $k$ , which is an azimuthal wavenumber, like  $m$ . Thus we can write  $h_{lmk} = h_l$ . For this special situation the summation over  $k$  in (8) can be performed directly. Rotating back and forth about the inclination returns a Kronecker  $\delta$ , i.e. the unit matrix:

$$\sum_k d_{lmk}(I) d_{lkp}(-I) = \delta_{mp}.$$

Thus (8) boils down to:

$$\begin{aligned} f(\phi, \lambda) &= \sum_l \sum_m \sum_p h_l K_{lm} i^{p-m} \delta_{mp} e^{i(m-p)\Lambda} Y_{lp}(\phi, \lambda) \\ &= \sum_l \sum_m h_l K_{lm} Y_{lm}(\phi, \lambda), \end{aligned}$$

which is just (1) with isotropic transfer applied. In this case, there is obviously no separation between ascending and descending contributions. Expressed in terms of the next section, the functional only contains a geographically mean contribution, no geographically variable one. Consequently, cross-over differences will vanish in case of isotropic functionals.

Inclination functions will naturally arise by rotating spherical harmonics into the orbital coordinate system. In this view, the previous derivation of generalized Rosborough functions by inverse rotation – after applying a specific transfer – reveals the nature of Rosborough functions. Nevertheless, this derivation has never been used, probably because of conceptual, but also for practical reasons. The additional summation over the  $p$ -index and the evaluation of Wigner coefficients  $d_{lkp}$  can be avoided, as shown in the next section.

Before that, two remarks are in place here. Firstly, the separation between ascending and descending track contributions, is hidden in the term  $\exp(i(m-p)\Lambda)$ . The variable  $u$  has vanished, though. Secondly, the Rosborough formalism does assume a so-called nominal orbit, in which orbital radius  $r$  and inclination  $I$  are constant. In principle, the above formulas do allow them to be variable, but it would lead to an impractical algorithm.

## 3 Rosborough Representations

Any functional of the gravitational potential along the orbit can be written by:

$$f(r, u, \Lambda, I) = \sum_l \sum_m \sum_k H_{lmk}(r, I) K_{lm} e^{i(ku + m\Lambda)}. \quad (12)$$

This equation is derived from (6) by inserting inclination function (5) and collecting  $F_{lmk}(I)$ ,  $h_{lmk}$ , dimensioning and upward continuation into the transfer coefficient  $H_{lmk}$ .

The goal now is to eliminate the orbital variables  $u$  (argument of latitude) and  $\Lambda$  (longitude of ascending node), such that (12) becomes an expression in spherical earth-fixed coordinates  $\phi$  and  $\lambda$ . The elimination is not achieved here by rotating the orbital frame back into the Earth-fixed frame, as in the previous section. We will not follow the lines of derivation as in Rosborough (1986), which led to cumbersome manipulation of real-valued trigonometric series developments with arguments of  $ku$  and  $m\Lambda$ . Instead, we manipulate the complex-valued term  $\exp(i(ku + m\Lambda))$  from (12).

For that purpose the spherical triangle between the satellite, the ascending node and the footpoint of the satellite on the equator (along its local meridian) is considered, cf. Fig. 1.

We note again that  $\lambda - \Lambda = \alpha - \Omega$ , such that

$$e^{im\Lambda} = e^{im\lambda} e^{-im(\alpha - \Omega)}. \quad (13)$$

Thus the longitude  $\lambda$  appears naturally. The sides of the spherical triangle are  $\phi$  along the satellite's meridian,  $u$  along the orbit and  $\alpha - \Omega$  along the equator. With basic spherical trigonometry, e.g. Kaula (1966), it is:

$$\cos u = \cos \phi \cos(\alpha - \Omega) \quad (14)$$

$$\sin u = \cos \phi \sin(\alpha - \Omega) / \cos I \quad (15)$$

$$\sin u = \sin \phi / \sin I \quad (16)$$

From the latter, one obtains:

$$\cos u = \pm \frac{1}{\sin I} \sqrt{\sin^2 I - \sin^2 \phi}, \quad (17)$$

where the + sign is valid for ascending tracks ( $u \in [-\pi/2, \pi/2]$ ) and the -- sign for descending tracks ( $u \in [\pi/2, 3\pi/2]$ ). This distinction is the starting point of the separation in mean and variable parts of the functional in the Rosborough representation. Moreover, (17) implies the condition  $\sin^2 \phi < \sin^2 I$ , or  $|\phi| < \pi/2 - |\pi/2 - I|$ . This means that all following formulae are only valid as long as the latitude is outside the polar gaps with spherical radius  $|\pi/2 - I|$ .

Combining (14) and (15) gives:

$$e^{i(\alpha - \Omega)} = \frac{1}{\cos \phi} (\cos u + i \sin u \cos I),$$

which still contains terms in  $u$ . Making use of (16) and (17) gives the set of equations:

$$e^{iu} = \frac{(\pm \sqrt{\sin^2 I - \sin^2 \phi} + i \sin \phi)}{\sin I}$$

$$e^{i(\alpha - \Omega)} = \frac{(\pm \sqrt{\sin^2 I - \sin^2 \phi} + i \sin \phi \cos I)}{\cos \phi \sin I}$$

Exponentiation of these with positive  $k$  and  $m$ , respectively, leads to

$$e^{iku} = \frac{(\pm \sqrt{\sin^2 I - \sin^2 \phi} + i \sin \phi)^k}{\sin^k I}$$

$$e^{-im(\alpha - \Omega)} = \frac{(\pm \sqrt{\sin^2 I - \sin^2 \phi} - i \sin \phi \cos I)^m}{\cos^m \phi \sin^m I}$$

Note that the complex conjugated was used in the latter, since  $e^{-ix} = \overline{e^{ix}}$ . This identity is also used to generalize to arbitrary  $k$  and  $m$ . Combining these results and renaming the complex-valued trigonometric function into  $\Phi$  yields:

$$\Phi_{mk}^{\pm}(\phi, I) = e^{i(ku - m(\alpha - \Omega))} \quad (18)$$

$$= \frac{1}{\cos^{|m|} \phi \sin^{|m|+|k|} I} \cdot \left( \pm \sqrt{\sin^2 I - \sin^2 \phi} + i \frac{|k|}{k} \sin \phi \right)^{|k|} \cdot \left( \pm \sqrt{\sin^2 I - \sin^2 \phi} - i \frac{|m|}{m} \sin \phi \cos I \right)^{|m|}$$

The powers are all positive. The sign of the imaginary parts depends on the sign of  $m$  and  $k$ . The following symmetry holds:  $\Phi_{-m,-k} = \overline{\Phi_{mk}}$ .

Inserting (18) into (12) results in:

$$f^{\pm}(r, \phi, \lambda, I) = \sum_l \sum_m \sum_k H_{lmk}(r, I) K_{lm} \Phi_{mk}^{\pm}(\phi, I) e^{im\lambda}. \quad (19)$$

In a next step the summation over  $k$  is performed in order to get an expression, similar to an ordinary SH series development. In this step we also revert to using the indices  $a$  and  $d$  for ascending (+) and descending (-), respectively:

$$f^{a/d}(r, \phi, \lambda) = \sum_l \sum_m K_{lm} Q_{lm}^{a/d}(\phi, I) e^{im\lambda} \quad (20)$$

$$\text{with } Q_{lm}^{a/d}(\phi, I) = \sum_k H_{lmk}(r, I) \Phi_{mk}^{a/d}(\phi, I). \quad (21)$$

Finally, the ascending and descending contributions are permuted to yield a *geographically mean* (or *geographically correlated*) part and a *geographically variable* part:

$$f^{m/v}(r, \phi, \lambda) = \sum_l \sum_m K_{lm} Q_{lm}^{m/v}(\phi, I) e^{im\lambda} \quad (22)$$

$$\text{with } Q_{lm}^m = (Q_{lm}^a + Q_{lm}^d)/2 \quad (23)$$

$$\text{and } Q_{lm}^v = (Q_{lm}^a - Q_{lm}^d)/2 \quad (24)$$

Vice versa, the ascending and descending contributions can be derived from the mean and variable parts by means of

$$Q_{lm}^a = Q_{lm}^m + Q_{lm}^v \quad (25)$$

$$Q_{lm}^d = Q_{lm}^m - Q_{lm}^v \quad (26)$$

These permutations show that representations in terms of  $m/v$  and in terms of  $a/d$  are equivalent. The geographically variable part is expressed in cross-over differences. The mean part is invisible in cross-overs.

## 4 Concluding Remarks

The representation (20) is a complex-valued counterpart of Rosborough's representation (Rosborough 1986). Depending on the transfer coefficient  $H_{lmk}$  it applies to any functional of the gravitational potential. It is valid only outside the polar gaps, i.e.  $|\phi| \leq \pi/2 - |\pi/2 - I|$ . Although the use of transfer coefficients in the framework of the time-wise approach may seem to imply validity along the satellite tracks alone, it must be pointed out that the fields themselves were rotated. Consequently the formulations is valid anywhere.

The complex-valued valued formalism and the corresponding notation avoid the necessity of the tilde-, overbar- and  $c$ - and  $s$ -variants of the functions used in Rosborough (1986). A great algorithmic advantage is that there is no need to expand (18) or its constituents further into binomial series, using functions  $Y$  and  $\Psi$ , as is done by Rosborough (1986) and in all other works that build on or improve his derivation, e.g. Bosch (1997). The complex quantities within brackets in (18) can be exponentiated and evaluated directly. Note also that from their definition (18) the  $\Phi_{mk}$ -functions obey  $||\Phi_{mk}|| = 1$ . Numerical stability is therefore guaranteed.

One could even consider to obtain  $u$  and  $(\alpha - \Omega)$  directly from  $\phi$  and  $I$  (and  $\pm$ ) by inversion of (14)–(17). The function  $\Phi_{mk}(\phi, I)$  is then evaluated directly by the very simple first line of (18).

We have provided the derivation of the Rosborough representation in complex-valued terms here in more detail than Sneeuw (2003). The only other attempt at such formulation was made in an unpublished memorandum by Balmino (1996, A note on Rosborough transformation, unpublished memorandum), who also demonstrated that binomial series expansions can be avoided when using complex notation. The complex-valued algorithm is efficient, while it avoids series expansions, and stable. By plugging in the appropriate transfer coefficient  $H_{lmk}$ , any gravity functional along the orbit can be represented. The Rosborough function  $Q_{lm}(\phi, I)$  is defined in terms of such  $H_{lmk}$ . If, for instance, one takes the representation coefficient  $H_{lmk}^{xx}$  of along-track gradients, one obtains a spatial representation of along-track gradients, i.e. in the

local orbital frame, irrespective of location on Earth. The Rosborough approach is thus a versatile methodology in spaceborne gravimetry, that deserves implementation on a broader scale. The proof-of-concept for GOCE data by Sharifi et al. (2013b) is a first step.

In this contribution we demonstrated the nature of the Rosborough representation through reverse transformation of spherical harmonics.

## References

- Balmino G (1993) Orbit choice and the theory of radial orbit error for altimetry. In: Rummel R, Sansò F (eds) Satellite altimetry in geodesy and oceanography. Lecture notes in earth sciences, vol 50, pp 244–315. Springer, Berlin
- Bosch W (1997) Geoid and orbit corrections from crossover satellite altimetry. DGFI Internal Technical Report
- Engelis T (1987) Radial orbit error reduction and sea surface topography determination using satellite altimetry. OSU Report 377, Dept. Geod. Sci. and Surv., Ohio State University
- ESA (1999) The four candidate earth explorer core missions: gravity field and steady-state ocean circulation mission. Technical Report ESA SP-1233
- Goad CC (1987) An efficient algorithm for the evaluation of inclination and eccentricity functions. Manuscr Geodaet 12:11–15
- Kaula WM (1966) Theory of satellite geodesy. Blaisdell Publishing Co, Waltham
- Klokočník J, Kostecký J, Jandová M (1995) Altimetry with dual-satellite crossovers. Manuscr Geodaet 20:82–95
- Kostecký J, Klokočník J, Kalina Z (1986) Computation of normalized inclination functions to high degree for satellite in resonances. Manuscr Geodaet 11:293–304
- Reigber C, Lühr H, Schwintzer P, Wickert J (eds) (2005) Earth observation with CHAMP: results from three years in orbit. Springer, Berlin
- Rosborough GW (1986) Satellite orbit perturbations due to the geopotential. Technical Report CSR-86-1, CSR
- Rummel R (ed) (2011) Special issue: GOCE - the gravity and steady-state ocean explorer, vol 85, number 11, pp 747–884 of Journal of Geodesy
- Sharifi MA, Safari A, Ghobadi-Far K (2013a) Rosborough formulation in satellite gravity gradiometry. Artif Satell 48(1):39–50
- Sharifi MA, Sneeuw N, Ghobadi-Far K (2013b) Analysis of GOCE data based on the Rosborough method. Poster Presentation, VIII Hotine-Marussi Symposium, Rome, Italy, 17–21 June 2013
- Sneeuw N (1992) Representation coefficients and their use in satellite geodesy. Manuscr Geodaet 17:117–123
- Sneeuw N (2000) A semi-analytical approach to gravity field analysis from satellite observations. Reihe C 527, Deutsche Geodätische Kommission
- Sneeuw N (2003) Space-wise, time-wise, torus and Rosborough representations in gravity field modelling. Space Sci Rev 108(1–2):37–46
- Tapley BD, Bettadpur S, Ries JC, Thompson PF, Watkins MM (2004) GRACE measurement of mass variability in the Earth system. Science 305:503–505

---

# Combining Different Types of Gravity Observations in Regional Gravity Modeling in Spherical Radial Basis Functions

Katrin Bentel and Michael Schmidt

---

## Abstract

With the increasing number of high-resolution gravity observations, which became available in the recent years, global Earth gravity models can be regionally refined. While global gravity models are usually represented in spherical harmonic basis functions with global support, a very promising option to model the regional refinements is the use of spherical radial basis functions with quasi-compact support. We use the approach of regional gravity modeling in spherical radial basis functions, with parameter estimation to determine the coefficients of the signal representation, on a test data set provided by the IAG-ICCT study group JSG0.3. We demonstrate on the data set for Europe that the approach is well-suited for different types of observations, such as terrestrial, aerial, and satellite-based measurements, as well as their combination. Furthermore, our results contribute to the study group's goal of inter-comparison of different modeling methodologies. Our regional modeling approach leads to relative errors of about 0.2–2% when compared to the validation data sets on the topography.

---

## Keywords

Combination of different observations • ICCT study group test data • Radial basis function • Regional gravity field modeling

---

## 1 Introduction

A study group under the umbrella of the IAG (International Association of Geodesy)—ICCT (Inter Commission Committee on Theory) between Commission 2 (Gravity Field) and Commission 3 (Earth Rotation and Geodynamics) titled as *Joint Study Group JSG0.3 Comparison of Current Methodologies in Regional Gravity Field Modeling* was established in 2011 with duration until 2015. The goal of this study group is to compare different regional

modeling methodologies and to finally outline standards and conventions for future regional gravity products. One of the activities so far was to provide synthetic test data sets which are used for inter-comparison of regional gravity modeling methodologies. Among the objectives are the choice of the type of basis function, the point grid, an appropriate methodology to solve the adjustment problem, and the consideration of errors.

Details on the study group as well as the test data can be found online at <http://jsg03.dgfi.badw.de>. Synthetic gravity observations of different types are provided for two different regions in Europe and in South America. For each region satellite-based, aerial, as well as terrestrial observations are provided, along with noise information for each observations type, and validation data sets in terms of disturbing gravity potential on the topography.

We use the test data sets in Europe for our regional gravity modeling approach in spherical radial basis functions.

---

K. Bentel (✉)

Department of Mathematical Sciences and Technology, Norwegian University of Life Sciences, IMT, Postboks 5003, 1432 Ås, Norway  
e-mail: [katrin.bentel@umb.no](mailto:katrin.bentel@umb.no)

M. Schmidt

Deutsches Geodätisches Forschungsinstitut, Munich, Germany

In Chap. 2 we explain the approach, in Chap. 3 we present our results with the individual data sets, and in Chap. 4 their combination. Finally, in Chap. 5 all modeling results are summarized and discussed. Thereby, with this article, we contribute to the goal of the study group of inter-comparison of different regional gravity modeling approaches by presenting our results with the study group's test data.

## 2 Regional Gravity Modeling in Spherical Radial Basis Functions

For regional gravity modeling, we use spherical radial basis functions, as presented in Freeden et al. (1998), Holschneider et al. (2003), or Schmidt et al. (2007) and references therein, amongst many others. We follow the approach given in Bentel et al. (2013). A regional residual gravity signal  $\Delta F$  is represented in a series expansion in spherical radial basis functions according to

$$\Delta F(\mathbf{x}) = \sum_{k=1}^K d_k B(\mathbf{x}, \mathbf{x}_k). \quad (1)$$

Thereby,  $B(\mathbf{x}, \mathbf{x}_k)$  are the radial basis functions, which depend only on the spherical distance between their location point  $\mathbf{x}_k$  and the evaluation point  $\mathbf{x}$ , and are defined as

$$B(\mathbf{x}, \mathbf{x}_k) = \sum_{n=0}^N \frac{2n+1}{4\pi R^2} \left(\frac{R}{r}\right)^{n+1} B_n P_n(\mathbf{x}, \mathbf{x}_k). \quad (2)$$

$P_n$  are the Legendre polynomials,  $R$  is the radius of a reference sphere (e.g. mean Earth radius), and  $r$  is the radius of the evaluation point  $\mathbf{x}$ . The coefficients  $B_n$  define the type of radial basis function. For the computations here, we use cubic polynomial radial basis functions, motivated by the findings in Bentel et al. (2013). They are defined by

$$B_n = \left(1 - \frac{1}{N}n\right)^2 \left(\frac{2}{N}n + 1\right), \quad (3)$$

and can be found in Freeden et al. (1998). The values for  $N$  are adjusted according to the signal which is to be modeled. With  $\mathbf{x}$ , the different types of observations are directly used at the locations at which they are obtained. The points  $\mathbf{x}_k$ , the locations for the radial basis functions, are chosen on a Reuter grid, see Freeden et al. (1998).

To determine the coefficients  $d_k$  of the regional signal representation, regularization is needed due to the downward continuation problem of gravity which is involved and due to non-uniqueness of the coefficients to be estimated. We use variance component estimation according to Koch and

Kusche (2002) to determine the variance components of the data sets and the prior information. The variance components can further be used to determine relative weighting factors between the different data sets as well as the regularization parameter with respect to the prior information. Prior information in terms of the expectation vector for the coefficients to be estimated is added. We set the vector of prior information equal to zero, because a residual signal is modeled after removing a reference field (EGM 96) up to spherical harmonic degree 60. All results presented here are obtained with sets of physically meaningful coefficients, what means they are correlated to the signal to be modeled as well as small on the margins, which are needed beyond the area of observations in order to avoid boundary effects. They are between  $2^\circ$  and  $3^\circ$  wide.

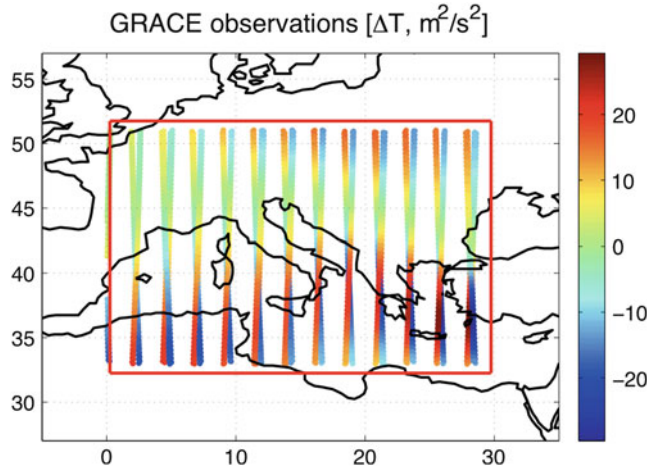
All modeling approaches are validated with the given validation data sets. From the regional gravity field representation, the disturbing potential is synthesized at the same points where validation data is available, respectively in the area where observations are available after subtracting a margin width. This is necessary to avoid boundary value effects and the margin widths are given together with the results in Table 1. Then the errors in terms of differences in each point are computed as well as a relative RMS error value in terms of percentage of the error RMS from the signal RMS of the full signal up to degree 2,190.

## 3 Different Examples of Gravity Observations in Europe

We use the test data provided by the ICCT study group ready for download at <http://jsg03.dgfi.badw.de> with the given realistic level of white noise on the observations. For the test region in Europe, two sets of synthetic observations are provided for each observation type. Satellite-based observations are available from GRACE (Tapley et al. 2004) and GOCE (Drinkwater et al. 2007), two sets of aerial observations are available for two different flight campaigns, and the terrestrial observations are available on a regular grid on the topography with two different grid spacings, one with  $30'$  and the other with  $5'$  spacing. From all different types of observations, a reference field (EGM 96) up to spherical harmonic degree 60 is removed, and restored later, when the gravity values for validation purposes are synthesized. With the regional gravity modeling approach, only a residual gravity signal is represented. Three different sets of validation data defined on the topography are provided, one in a larger area and two in a smaller area. We use the one which fits best with the area of observations for the different types of observations.

**Table 1** Summary of the modeling results

<i>Individual data sets</i>				
Type of observations	Error RMS%	$N$ kernel	Margin width [°]	
GRACE-type	2.1	300	0	
GOCE-type	1.26	350	0	
Terrestrial, 30' spacing	1.55	300	0	
Terrestrial, 5' spacing	0.21	1,200	2	
Aerial, case I	0.40	700	1	
Aerial, case II	0.37	900	1	
<i>Combination of the data sets</i>				
Type of observations	Error RMS% appr. A	Error RMS% appr. B	$N$ kernel	Margin width [°]
GRACE + terrestrial (30')	1.28	0.95	350	0
GRACE + GOCE	2.26	1.26	350	0
Aerial case I + terrestrial (5')	0.35	0.22	1,200	2
GRACE + aerial case I + aerial case II	2.3	2.7	350	2
GRACE + aerial case I + terrestrial (5')	0.65	0.17	1,000	2

**Fig. 1** GRACE observations with white noise with a standard deviation of  $0.0008 \text{ m}^2/\text{s}^2$ , together with validation area (red box)

The modeling results, in terms of RMS error after validation, for all data sets are given in the first part of Table 1. In the following, two examples are presented in more detail. The first example are GRACE-type observations. The observations, potential differences along real GRACE orbits, are given in Fig. 1. Figure 2 shows the modeling results from the GRACE observations in terms of disturbing potential on the topography on the left hand side. The plot in the center shows the validation field, and the plot on the right hand side the difference between the two previous ones, that is, the modeling errors. The second example are aerial observations for one flight campaign. Again, Fig. 3 shows the observations and Fig. 4 the modeling results.

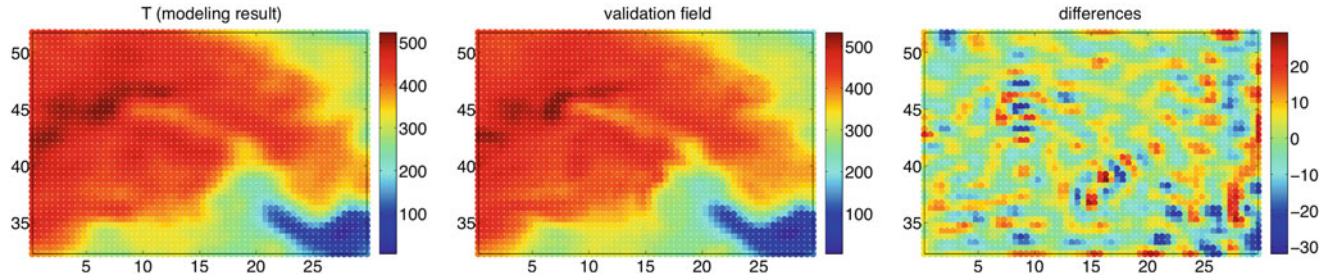
#### 4 Combination of the Different Data Sets

Investigations in combining heterogenous data sets have already been made, as for example in Panet et al. (2012) with a wavelet approach. We combine different data sets by determining a relative weighting, favourably related to the accuracies of the individual observation sets. For that purpose we use here the method of variance component estimation (VCE) as already mentioned before. To establish our linear model we first transform the observation equation as defined in Eq. (1) into the matrix equation  $\Delta\mathbf{F}_i + \mathbf{e}_i = \mathbf{A}_i \mathbf{d}$  where  $i = 1, \dots, n$  means an individual observation set. Then we combine the  $n$  single models to the combined model

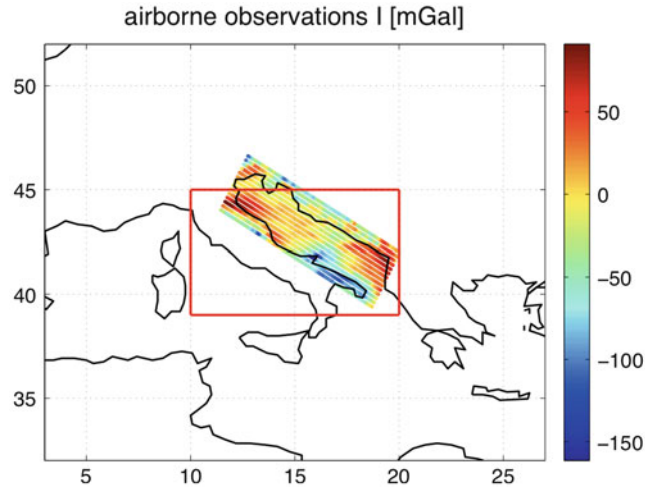
$$\begin{bmatrix} \Delta\mathbf{F}_1 \\ \Delta\mathbf{F}_2 \\ \vdots \\ \Delta\mathbf{F}_n \\ \boldsymbol{\mu} \end{bmatrix} + \begin{bmatrix} \mathbf{e}_1 \\ \mathbf{e}_2 \\ \vdots \\ \mathbf{e}_n \\ \mathbf{e}_\mu \end{bmatrix} = \begin{bmatrix} \mathbf{A}_1 \\ \mathbf{A}_2 \\ \vdots \\ \mathbf{A}_n \\ \mathbf{I} \end{bmatrix} \mathbf{d},$$

$$D\left(\begin{bmatrix} \Delta\mathbf{F}_1 \\ \vdots \\ \Delta\mathbf{F}_n \\ \boldsymbol{\mu} \end{bmatrix}\right) = \begin{bmatrix} \sigma_1^2 \mathbf{I}_1 & 0 & \cdots & 0 \\ 0 & \ddots & & \vdots \\ \vdots & & \sigma_n^2 \mathbf{I}_n & 0 \\ 0 & \cdots & 0 & \sigma_\mu^2 \mathbf{I}_\mu \end{bmatrix} \quad (4)$$

which means a *Gauss-Markov model* with unknown coefficient vector  $\mathbf{d}$  and unknown variance components  $\sigma_i^2$  with  $i = 1, \dots, n$  for the  $n$  observation sets and  $\sigma_\mu^2$  for the prior



**Fig. 2** GRACE regional modeling results for a region in Europe; all results given in disturbing potential [ $\text{m}^2/\text{s}^2$ ]; relative error RMS: 2.1%



**Fig. 3** Airborne observations (case I) with white noise with a standard deviation of 1 mGal together with validation area (red box)

information. In the following we distinguish between two approaches on variance component estimation (according to the Koch and Kusche (2002)):

- (a) We introduce the assumption  $\sigma_1^2 = \sigma_2^2 = \dots = \sigma_n^2 =: \sigma_0^2$  for the  $n$  variance components  $\sigma_i^2$ . Thus, in this approach we determine the estimations of  $\mathbf{d}$  as well as  $\sigma_0^2$  and  $\sigma_\mu^2$ . The iteratively determined variance components lead—in the point of convergence—to the solution

$$\hat{\mathbf{d}} = \left( \frac{1}{\hat{\sigma}_0^2} \sum_{i=1}^n \mathbf{A}_i^T \mathbf{A}_i + \frac{1}{\hat{\sigma}_\mu^2} \mathbf{I}_\mu \right)^{-1} \left( \frac{1}{\hat{\sigma}_0^2} \sum_{i=1}^n \mathbf{A}_i^T \Delta \mathbf{F}_i + \frac{1}{\hat{\sigma}_\mu^2} \mathbf{I}_\mu \boldsymbol{\mu} \right)$$

for the unknown coefficient vector  $\mathbf{d}$ .

- (b) Besides the unknown coefficient vector  $\mathbf{d}$  we here introduce all  $n + 1$  variance components  $\sigma_i^2$  for  $i = 1, \dots, n$  and  $\sigma_\mu^2$  defined in the model (4) as unknown parameters. With the estimation of the individual variance components the relative weighting between all observation sets and the prior information is determined. Thus, the VCE yields in the point of convergence the solution

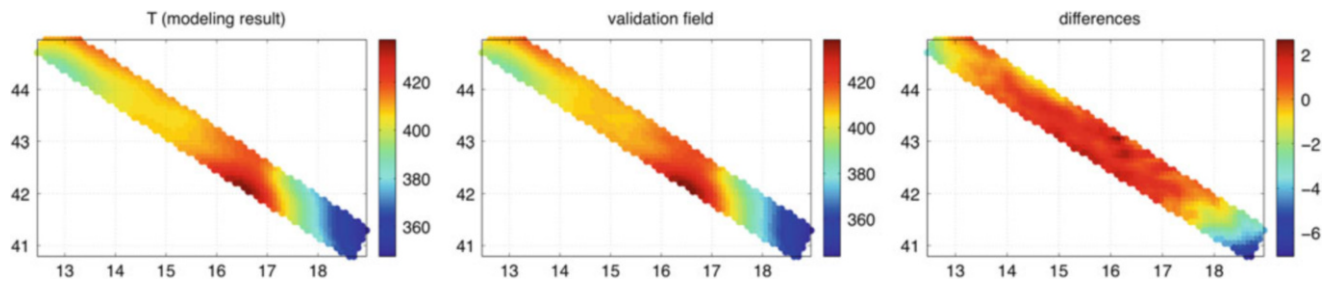
$$\hat{\mathbf{d}} = \left( \sum_{i=1}^n \frac{1}{\hat{\sigma}_i^2} \mathbf{A}_i^T \mathbf{A}_i + \frac{1}{\hat{\sigma}_\mu^2} \mathbf{I}_\mu \right)^{-1} \left( \sum_{i=1}^n \frac{1}{\hat{\sigma}_i^2} \mathbf{A}_i^T \Delta \mathbf{F}_i + \frac{1}{\hat{\sigma}_\mu^2} \mathbf{I}_\mu \boldsymbol{\mu} \right).$$

The different sets of observations are combined according to the two approaches discussed before. The modeling results from these two combinations are presented in the lower part of Table 1. In Fig. 5 the modeling results for one combination example are presented.

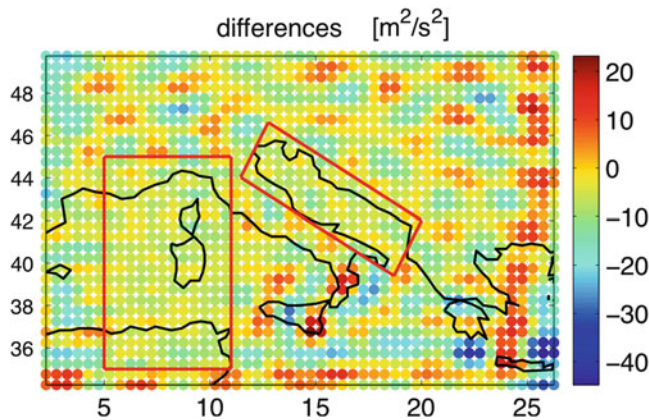
## 5 Regional Modeling Results

In Table 1, all modeling results from the ICCT study group test data for the region in Europe are summarized. For each set of observations, the relative RMS error is given in %-values together with the maximum degree in the cubic polynomial basis function and a margin width which is used in order to avoid boundary effects. The values given in [°] indicate by how much the validation area is smaller than the area of observations. For the results obtained from combination of different data sets the error RMS% values are given for both of the approaches outlined before. The results from satellite based observations lead to slightly worse results than the other observations. This is due to the downward continuation problem of gravity, which is of course included when gravity on the Earth surface is computed from observations at satellite orbit height. Downward continuation is an ill-posed problem by its physical nature, the gravity signal gets attenuated with distance from the masses. The terrestrial observations with 30' spacing lead as well to an RMS error which is not as low as from the other terrestrial data sets. This is due to the fact that also in the observations with 30' spacing, information up to spherical harmonic degree 2,190 was included. But the spacing of observations is not dense enough to sample this high frequency signal. Thus, not the full signal content can be recovered from the coarse observations. The very dense sampling of the terrestrial observations with 5' spacing as well as from the aerial observations lead





**Fig. 4** Regional modeling results for a region in Europe, from aerial observations, case I. All fields are disturbing potential, in  $[m^2/s^2]$ ; relative error RMS: 0.40%



**Fig. 5** Regional modeling results for a simultaneous analysis of GRACE and two sets of aerial observations. The plot shows the difference in the disturbing potential field on the topography synthesized from the modeling approach and the reference field. That is, the modelling errors in  $[m^2/s^2]$ . The *red boxes* indicate the area of the two flight campaigns (called case I and II). GRACE observations (see Fig. 1) are available throughout the whole area. The relative error RMS is 2.3% and the plot shows that the errors are small in the areas where aerial observations are available, but high outside

to very small errors in the recovered signal. The spacing of the sampling is dense enough to recover the maximum frequency in the signal (spherical harmonic degree 2,190). The terrestrial observations lead to even better results than the aerial observations, since in the terrestrial observations no downward continuation of gravity is included, while in the aerial observations it still plays a role.

In the data combination results in the second part of Table 1 approach B, with individual weights for the data sets, generally leads to better results than approach A. The only exception is the combination of GRACE data with the two sets of aerial results. This special case is discussed in the following in more detail.

In the combination of GRACE-type and terrestrial observations with  $30'$  spacing it can be seen how additional terrestrial observations lead to a better result. The combined error is lower than the individual errors. This also holds for the combination of GRACE and GOCE data, however, when

two satellite-based data sets are combined, the result does not improve that much. The combination of aerial and terrestrial observations leads to good results, since the error RMS of the aerial data can be significantly improved by adding terrestrial observations.

The results of combining GRACE-type observations and the two aerial data sets show that even if the two aerial data sets cover a reasonable part of the area of interest (but still less than half of the area), this is not enough to make the solution of the whole area significantly better than from GRACE observations alone. For validation, the data set with spherical harmonic degrees up to 2,190 was chosen. These high degrees can not be recovered in the areas where only GRACE data is available and with a basis function of only degree 350. Therefore, the overall RMS% error is even higher than for GRACE data alone. Furthermore, the results for approach B are even worse than for approach A, since in this non-realistic scenario, no appropriate variances can be assigned to the data sets, and the errors outside the areas of aerial observations get very high.

Finally, the combination of three different types of data sets, namely satellite-based, aerial and terrestrial observations, leads to very good modeling results. The error RMS value is amongst the lowest to be achieved.

## 6 Summary and Outlook

Different types of gravity observations can be combined in one parameter estimation step in the regional gravity field modeling approach in spherical radial basis functions. The different sets of observations can be directly used in the approach, without prior processing or gridding of the observed values. The results presented in this paper are not only useful for comparisons with other methods for the ICCT study group, but they are also a first step towards the analysis of real observations. All results shown above are obtained from simulated observations with a realistic noise level according to the ICCT study group. They lead to modeling errors between 0.2 and 2% for the different scenarios.

Due to more and more available high-resolution gravity observations, regional gravity modeling techniques play an important role, since the common approach of gravity modeling in spherical harmonic basis functions cannot accommodate regional gravity refinements appropriately. However, using real data in the regional modeling approach would be more tricky than this simulation, e.g. due to coloured noise of and correlations between the observations. In order to take the simulation study closer to processing real observations, stochastic properties of the different data types could be considered and improved in the variance component estimation step as well.

**Acknowledgements** This study was made possible through funding by a “DAAD Doktorandenstipendium” from the German Academic Exchange Service.

---

## References

- Bentel K, Schmidt M, Gerlach C (2013) Different radial basis functions and their applicability for regional gravity field representation on the sphere. *Int J Geomath* 4:67–96. doi:10.1007/s13137-012-0046-1
- Drinkwater MR, Haagmans R, Muzi D, Popescu A, Floberghagen R, Kern M, Fehringer M (2007) The GOCE gravity mission: ESA's first core Earth explorer. In: Proceedings of 3rd international GOCE user workshop, 6–8 November, 2006, Frascati, Italy, ESA SP-627, pp 1–8
- Freedon W, Gervens T, Schreiner M (1998) Constructive approximation on the sphere with applications to geoscience. Oxford Science, Oxford. ISBN 019853682
- Holschneider M, Chambodut A, Mandea M (2003) From global to regional analysis of the magnetic field on the sphere using wavelet frames. *Phys Earth Planet Inter* 135(2–3):107–124. doi:10.1016/S0031-9201(02)00210-8
- Koch KR, Kusche J (2002) Regularization of geopotential determination from satellite data by variance components. *J Geod* 76:259–268. doi:10.1007/s00190-002-0245-x
- Panet I, Kuroishi Y, Holschneider M (2012) Flexible dataset combination and modelling by domain decomposition approaches. In: Sneeuw N, Novak P, Crespi M, Sanso F (eds) VII Hotine-Marussi Symposium on Mathematical Geodesy, International Association of Geodesy Symposia, vol 137. Springer, Berlin/Heidelberg, pp 67–73. doi:10.1007/978-3-642-22078-4\_10
- Schmidt M, Fengler M, Mayer-Gürr T, Eicker A, Kusche J, Sánchez L, Han SC (2007) Regional gravity modeling in terms of spherical base functions. *J Geod* 81(1):17–38. doi:10.1007/s00190-006-0101-5
- Tapley B, Bettadpur S, Watkins M, Reigber C (2004) The gravity recovery and climate experiment: Mission overview and early results. *Geophys Res Lett* 31:L09607. doi:10.1029/2004GL019920

---

# Height Datum Unification by Means of the GBVP Approach Using Tide Gauges

E. Rangelova, M.G. Sideris, B. Amjadiparvar, and T. Hayden

---

## Abstract

In this paper, we discuss the methodology of height datum unification by means of the Geodetic Boundary Value Problem (GBVP) approach and tide gauge information. We apply the global multiple vertical datum GBVP approach with an observation equation for the datum offset written in terms of the ellipsoidal height of the mean sea level at a tide gauge, the height of mean sea level in the national vertical datum and the geoid height. An example is given for CGVD28 and NAVD88 datums in North America. A few issues related to the geoid height are studied: the magnitude of the so-called indirect bias term, as well as the GOCE global geopotential model (GGM) commission and omission errors and their effect on the accuracy of the computed datum offsets. It is shown that the indirect bias term is below 1 cm if a residual Stokes's kernel is used, which corresponds to a degree and order 180 of the GOCE GGM. The GOCE geoid commission error computed from the time-wise approach GGM of degree 180 is 2–3 cm at the North American tide gauges. The GOCE GGM omission error could affect the computed mean vertical datum offsets at the tide gauges by as much as 7 cm.

---

## Keywords

GBVP • GOCE • Height system unification • Tide gauges • Vertical datum offset

---

## 1 Introduction

The geodetic boundary value problem (GBVP) approach for height datum unification was proposed decades ago by Colombo (1980), Rummel and Teunissen (1988), Rapp and Balasubramania (1992) and others. In this study, we review the global single vertical datum GBVP (Heiskanen and Mortiz 1967) and the global multiple vertical datum GBVP (Rummel and Teunissen 1988), both in spherical approximation. For the multiple vertical datum problem in ellipsoidal approximation, Sansò and Venuti (2002) can be consulted.

---

E. Rangelova (✉) • M.G. Sideris • B. Amjadiparvar • T. Hayden  
Department of Geomatics Engineering, University of Calgary, 2500  
University Drive NW, Calgary, AB, Canada T2N 1 N4  
e-mail: [evrangel@ucalgary.ca](mailto:evrangel@ucalgary.ca)

Global or regional height system unification requires local height datum offsets that refer to one globally or regionally defined equipotential surface and are computed by means of the most accurate GOCE-based satellite global geopotential model (GGM) for long and medium gravity field wavelengths, terrestrial gravity data and ellipsoidal heights in ITRF. The local datum offsets can be computed on land using GNSS-surveyed levelling benchmarks, at the coast using GNSS-surveyed tide gauge stations, and at sea using high resolution mean dynamic topography and geoid models. Height system unification with tide gauge stations seems a natural approach as the classical height datums were typically defined by the local mean sea level computed at one or more tide gauge stations.

In this work, we compute the offsets of the Canadian and US height datums, namely CGVD28 (Cannon 1929) and NAVD88 (Zilkoski et al. 1992). Both CGVD28 and NAVD88 have the typical flaws of the continental-size

vertical datums realized through levelling data that were collected over the time span of many decades: poor absolute accuracy of heights, large coast-to-coast distortions and many local distortions, some of which are a result of significant crustal uplift or subsidence, to name a few. It should be noted that recently CGVD28 has been replaced by the new geoid-based datum CGVD2013 with which the aforementioned datum issues have been resolved. However, computing the CGVD28 offsets at the Canadian tide gauges with respect to a common equipotential surface for North America is still useful as the Canadian gravity database and high resolution digital terrain models are based on this old datum.

We apply the global multiple vertical datum GBVP approach with an observation equation for the datum offset written in terms of the ellipsoidal height of the mean sea level, the height of mean sea level in the national vertical datum and the geoid height at a tide gauge. We focus on a few issues related to the last of the three height components: the magnitude of the so-called indirect bias term in the geoid height resulting from the use of biased local gravity data in geoid computations, as well as the GOCE commission and omission geoid errors and their effect on the accuracy of the computed datum offsets. The vertical datum offsets refer to the equipotential surface defined by the value  $W_o = 62636856.0 \text{ m}^2 \text{ s}^{-2}$ , which was computed by averaging the potential of the mean water level at the North American tide gauges (Hayden et al. 2012).

## 2 Single Vertical Datum Problem

With a harmonic disturbing potential  $T = W - U$  outside the boundary surface  $\Omega$ , the global single vertical datum GBVP in spherical approximation is defined as follows:

$$\begin{cases} \Delta T = 0 \\ -\frac{\partial T}{\partial r} - \frac{2}{R}T = \Delta g - \frac{2}{R}\Delta W_o \end{cases}, \quad (1)$$

where  $R$  is the mean radius of the Earth, represented by  $\Omega$ , and  $\Delta g$  is the gravity anomaly given on the boundary surface. The datum problem can be solved for the potential upon introducing the regularity condition  $T \rightarrow 0$  when the geocentric distance  $r \rightarrow \infty$ . On the boundary surface, the  $\Delta g$  values are corrected for the unknown height datum parameter  $\Delta W_o = W_o - U_o$ , where  $W_o$  is the potential of the geoid and  $U_o$  is the normal potential of the reference ellipsoid.

When the geoid height  $N_{PGOCE}$  is computed from a GOCE-based GGM (e.g., Pail et al. 2011), the solution to the single vertical datum problem at point  $P$  is

$$N_P = N_o + N_{PGOCE} + \frac{1}{\gamma} \mathbf{S}_P \Delta g_{res}, \quad (2)$$

where  $\gamma$  is the normal gravity on the reference ellipsoid and  $N_o = \delta GM/R\gamma - \Delta W_o/\gamma$ .  $\delta GM$  is the difference in the geocentric gravitational constant  $GM$  of the geoid and  $GM^e$  of the normal ellipsoid. It is shown by Kotsakis et al. (2012) that the uncertainty of  $\delta GM$  imposes an error of 1 cm in absolute vertical positioning with respect to an arbitrary equipotential surface.

$$\mathbf{S}_P \Delta g = \frac{R}{4\pi} \iint_{\Omega} St(\psi_{PQ}) \Delta g d\Omega_Q \quad (3)$$

is Stokes's integral (Heiskanen and Mortiz 1967), where the integration is performed over the sphere  $\Omega$  with radius  $R$ . The integration kernel  $St(\psi_{PQ})$  is Stokes's function, which depends on the geocentric distance  $\psi$  between point  $P$  and the variable location  $Q$  of  $\Delta g$ .

Stokes's integral in Eq. (3) is evaluated with the residual gravity anomalies  $\Delta g_{res} = \Delta g - \Delta g_{GOCE}$ . The classical GBVP approach requires that the effect of the topographic masses on gravity is removed before the integration and is restored on the geoid afterwards.

## 3 Multiple Vertical Datum Problem

In reality, gravity anomalies do not refer to a global vertical datum defined by  $W_o$ . Instead, they refer to local vertical datums  $j$ , defined by the equipotential surfaces with potential  $W_o^j$  (the so-called local zero height level). The single vertical datum problem is transformed to a multiple vertical datum problem with as many unknown vertical datum parameters (biases)  $\delta W_o^j = W_o - W_o^j$  as vertical datums exist.

Assuming that each local zero height level is biased with respect to one reference surface (although tilts and other long-wavelength datum distortions can exist in practice), Rummel and Teunissen (1988) proposed a solution to the multiple vertical datum problem, which we adopt herein. It is assumed that the Earth is represented by the sphere  $\Omega$  that is covered by  $J+1$  non-overlapping vertical datum zones  $\Omega^j$ ,  $j=0,1,2,\dots,J$  such that  $\Omega = \Omega^0 \cup \Omega^1 \dots \cup \Omega^J$ . The zone  $\Omega^0$  defined by  $W_o$  is chosen arbitrarily as a reference datum.

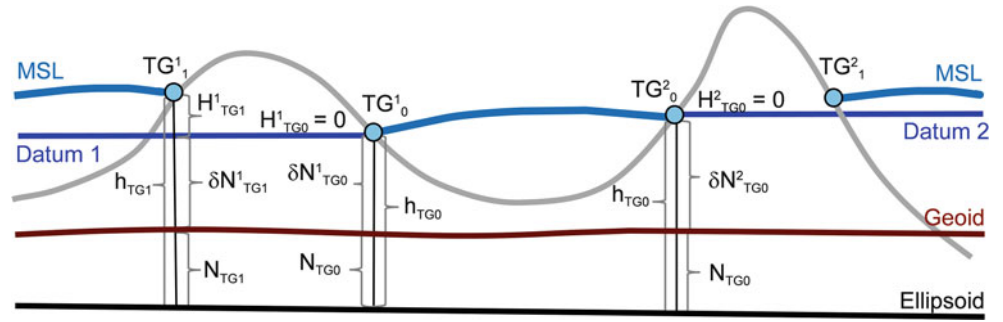
The input to Stokes's integral in Eq. (3) are the local gravity anomalies  $\Delta g^j$  corrected for the unknown bias  $\delta W_o^j$ :

$$\Delta g = \Delta g^j + \frac{2}{R} \delta W_o^j, j = 0, 1, 2, \dots, J \quad (4)$$

where the gravity anomalies  $\Delta g$  refer to  $\Omega^0$ . The term  $2\delta W_o^j/R$  is the "free-air" reduction one uses to reduce  $\Delta g^j$  from the local zero height level  $W_o^j$  to the reference level  $W_o$ .

When long and medium wavelengths of the geoid height at point  $P$  are determined from a GOCE-based GGM, i.e.,  $N_{PGOCE}$ , a residual geoid height  $N_{P_{res}}$  is computed by integrat-

**Fig. 1** Schematic representation of necessary data for computing vertical datum offsets using tide gauges and the GBVP approach



ing the residual gravity anomalies  $\Delta g_{res}^j = \Delta g^j - \Delta g_{GOCE}$  over the sphere  $\Omega$ :

$$N_{P_{res}} = \frac{R}{4\pi\gamma} \iint_{\Omega} St(\psi_{PQ}) \Delta g_{res}^j d\Omega_Q, \quad (5)$$

where, in principle, the superscript  $j$  takes on all  $J + 1$  values  $0, 1, \dots, J$ .

The full geoid height at P is represented by  $N_P = N_o + N_{P_{GOCE}} + N_{P_{res}}$ , where  $N_o = \Delta W_o / \gamma = (W_o - U_o) / \gamma$ , and is used subsequently to formulate the observation equation

$$N_P^j = N_o + N_{P_{GOCE}} + N_{P_{res}} + \delta N^j + N_P^{ind}. \quad (6)$$

$N_P^j$  on the left hand side of Eq. (6) is the GNSS-levelling geoid height computed with the ellipsoidal height  $h_P$  and the orthometric height  $H_P^j$  as  $N_P^j = h_P - H_P^j$ . It differs from the geoid height  $N_P$  on the right hand side of Eq. (6) by the direct bias term (datum offset)

$$\delta N^j = \delta W_o^j / \gamma \quad (7)$$

and the indirect bias term

$$N_P^{ind} = \frac{1}{2\pi\gamma} \delta W_o^j \iint_{\Omega^j} St(\psi_{PQ}) d\Omega_Q^j + \sum_{i=1, i \neq j}^J \frac{1}{2\pi\gamma} \delta W_o^i \iint_{\Omega^i} St(\psi_{PQ}) d\Omega_Q^i, \quad (8)$$

which is added in Eq. (6) because Stokes's integral is evaluated with the biased gravity anomalies  $\Delta g^j$  instead of  $\Delta g$ , or, to be more specific,  $\Delta g_{res}^j$  are used rather than  $\Delta g_{res}$ . The first term on the right hand side of Eq. (8) is the contribution to the geoid height at point P from the offset  $\delta W_o^j$  of the datum zone  $j$ , and the second term represents the effect on the geoid height at P from the offsets  $\delta W_o^i$  of all other datum zones  $\Omega_i$ ,  $i = 1, 2, \dots, J, i \neq j$ , excluding the reference zone  $\Omega_o$  for which the offset is zero by assumption.

Using Eq. (7), the indirect bias term in Eq. (8) can be written as

$$N_P^{ind} = \frac{1}{2\pi} \sum_{j=1}^J \delta N^j \iint_{\Omega^j} St(\psi_{PQ}) d\Omega_Q^j = \sum_{j=1}^J \delta N^j f^j. \quad (9)$$

To solve the multiple vertical datum problem, at least one point P should be given in the datum zone  $j$  with its ellipsoidal, orthometric and geoid height, and this point could be the fundamental tide gauge  $TG_o^j$  (Fig. 1). In this case, Eq. (6) can be rewritten as

$$h_{TG} - H_{TG}^j - (N_o + N_{TG_{GOCE}} + N_{TG_{res}}) = \delta N^j + \sum_{j=1}^J \delta N^j f^j, \quad (10)$$

where  $h_{TG}$  is the ellipsoidal height of the local mean sea level (MSL) at the tide gauge and  $H_{TG}^j$  is the MSL height in the datum  $j$ . The geoid height  $N_{TG} = N_o + N_{TG_{GOCE}} + N_{TG_{res}}$  is known and moved to the left hand side in the equation.

By means of Eq. (10), a linear system with  $J$  unknowns can be written with a fully populated design matrix  $A$  whose entries are the coefficients of the unknown offsets  $\delta N^j$ ,  $j = 1, \dots, J$ . The constant  $N_o$  can also be assumed unknown and estimated together with the datum offsets.

Gerlach and Rummel (2013) have shown that the indirect bias term determined by Eq. (8) can be negligibly small provided that Stokes's integral is evaluated with a residual kernel, which contains only spherical harmonic degrees  $n > n_{max}$ ; see Eq. (23) in their work. The omission of the indirect bias term results in a significantly simplified design matrix  $A$ , and the offset in each datum zone can be estimated separately from the other datum zones. One can solve a linear system of  $n$  observation equations for  $n$  tide gauges in the datum zone  $\Omega^j$  and estimate the mean offset  $\delta N^j$  by means of the least-squares adjustment model:

$$A \delta N^j = l, \quad Q_{ll} = Q_{hh} + Q_{HH} + Q_{NN_{GOCE}} + Q_{NN_{res}}, \quad (11a)$$

$$\delta\hat{N}^j = (A^T Q_{ll}^{-1} A)^{-1} A^T Q_{ll}^{-1} l, \quad Q_{\delta\hat{N}\delta\hat{N}} = (A^T Q_{ll}^{-1} A)^{-1} \quad (11b)$$

with  $l = h_{TG} - H_{TG}^j - (N_o + N_{TG_{GOCE}} + N_{TG_{res}})$  of size  $n \times 1$ , a  $n \times 1$  design matrix  $A = (1, 1, \dots, 1)^T$  and a stochastic model  $Q_{ll}$  composed of  $Q_{hh}$  of the GNSS ellipsoidal height of the local MSL,  $Q_{HH}$  of the height of MSL in the local datum,  $Q_{N_{NGOCE}}$  of the GOCE geoid height and  $Q_{N_{N_{res}}}$  of the residual geoid height.

With known error variances  $\sigma_{hh}^2$ ,  $\sigma_{HH}^2$ ,  $\sigma_{N_{NGOCE}}^2$  and  $\sigma_{N_{N_{res}}}^2$ , and assuming uncorrelated observations  $h$ ,  $H$ ,  $N_{TG_{GOCE}}$ , and  $N_{TG_{res}}$ , the mean datum offset is estimated as a weighted mean following from Eq. (11a):

$$\delta\hat{N}^j = \sum_{i=1}^n p_i l_i / \sum_{i=1}^n p_i, \quad \hat{\sigma}_{\delta\hat{N}}^2 = \hat{\sigma}_o^2 / \sum_{i=1}^n p_i \quad (12)$$

with  $p_i = \left[ (\sigma_{hh}^2)_i + (\sigma_{HH}^2)_i + (\sigma_{N_{NGOCE}}^2)_i + (\sigma_{N_{N_{res}}}^2)_i \right]^{-1}$

and the a posteriori variance factor  $\hat{\sigma}_o^2 = \sum_{i=1}^n (l_i - \delta\hat{N})^2 p_i / (n - 1)$ .

## 4 Required Data

### 4.1 Mean Sea Level in Local Height Datum and Ellipsoidal Height

Ideally, MSL should be computed from 19-year long continuous records of water levels so that nodal tides, atmospheric pressure and storm events are averaged out (Pugh 1987). For the purpose of height system unification, MSL can be computed from shorter records and/or records with data gaps in a network of tide gauge stations to possibly reduce both the effect of the GOCE-based geoid omission error on the computed mean vertical datum offsets (Gruber et al. 2012) and the effect of random data errors. In areas with significant crustal motion, heights of tide gauge benchmarks should be corrected using GNSS-derived vertical crustal velocities or geophysical models of crustal motion. In addition, water levels should be corrected for local long-term sea level changes determined at the tide gauge stations.

The ellipsoidal height of MSL at each tide gauge is computed with respect to the reference ellipsoid by reducing the ellipsoidal height of the tide gauge benchmark with the height difference between the benchmark and the chart datum obtained by precise levelling and adding the measured water level from the chart datum (Woodworth et al. 2012). The GNSS ellipsoidal heights should be given

in a common ITRF and epoch. This is not usually the case when tide gauges are surveyed by different agencies, and, in some cases, the reference frame may not be known.

Tide gauge stations used in this work are revised local reference stations from the Permanent Service for Mean Sea Level (PSMSL) with MSL data measured relative to a known benchmark. The MSL data are for the time period 1993–2002. They were corrected for the inverse barometer effect, but a correction for the nodal tide was not applied. All computations herein are performed in a conventional tide free system. The GNSS ellipsoidal heights are given in either ITRF2005 or ITRF2008. Accuracy information is not available for the GNSS ellipsoidal heights and heights of MSL in the local datum.

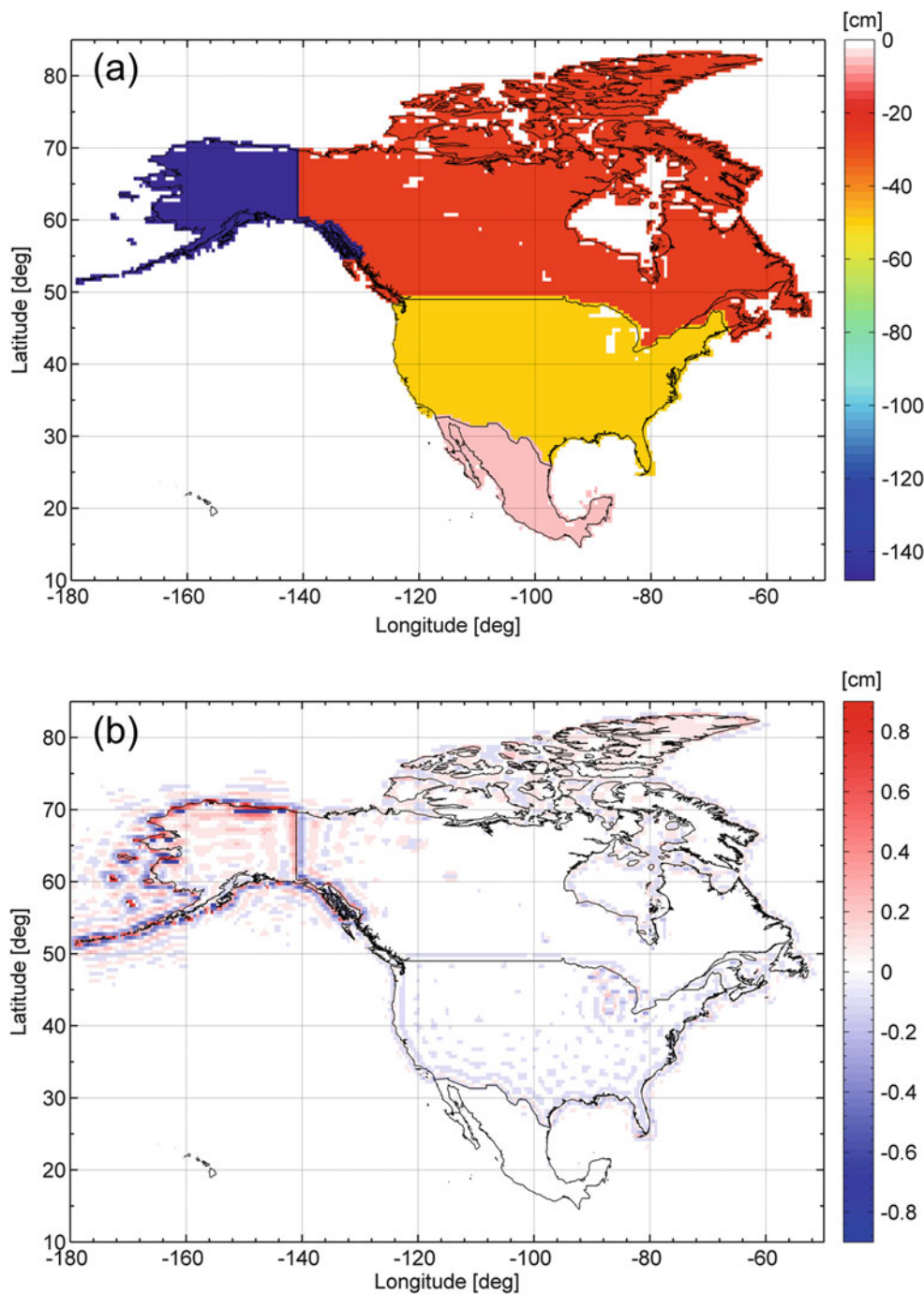
### 4.2 Geoid Height

In order to apply the GBVP approach in coastal areas with as small a geoid omission error as possible, the geoid height should be computed from a GOCE-based GGM up to the highest possible degree that keeps the high-degree commission error small. This degree is usually determined by means of evaluation of the GOCE-based GGM with local GNSS-levelling data, and typically the improvement over EGM2008 (Pavlis et al. 2012) is assessed for different spherical bands. Our evaluation for North America using GNSS-levelling geoid heights shows that the useful GOCE spherical harmonic degree (after which the GOCE geoid error becomes larger than the EGM2008 geoid error, e.g., Amjadiparvar et al. 2013) of the recent releases of the GOCE-based time-wise GGM (Pail et al. 2011) varies among Canada (degree 180), conterminous USA (degree 210), Mexico (degree 230) and Alaska (degree 230). Here, we adopt the minimum of these values, i.e., degree 180, for the whole North America and use it in the following computations of the indirect bias term and the GOCE geoid commission and omission errors.

#### 4.2.1 Indirect Bias Term

To compute the indirect bias term over North America using Eq. (8), preliminary mean offsets of NAVD88 (conterminous USA, Alaska and Mexico) and CGVD28 (Canada) with respect to the level  $W_o = 62636856.0 \text{ m}^2 \text{ s}^{-2}$  (Fig. 2a) are computed by means of GNSS ellipsoidal heights, orthometric heights and GOCE geoid heights at the first order levelling benchmarks. Figure 2b shows a map of the indirect bias term obtained by Eq. (8) using a residual kernel, where  $n_{max} = 180$ . The indirect bias term ranges between  $-0.9$  and  $0.9$  cm due to artificial effects from truncating the kernel. It can be concluded that the error in the datum offset introduced by the omission of the indirect bias term is less than 1 cm

**Fig. 2** (a) Vertical datum offsets from the level  $W_o = 62636856.0 \text{ m}^2 \text{ s}^{-2}$  computed with GNSS-levelling data. (b) Indirect bias term computed with a residual kernel with  $n_{max} = 180$



even for such an extreme offset of  $-140$  cm of NAVD88 in Alaska.

**4.2.2 GOCE Geoid Commission Error**

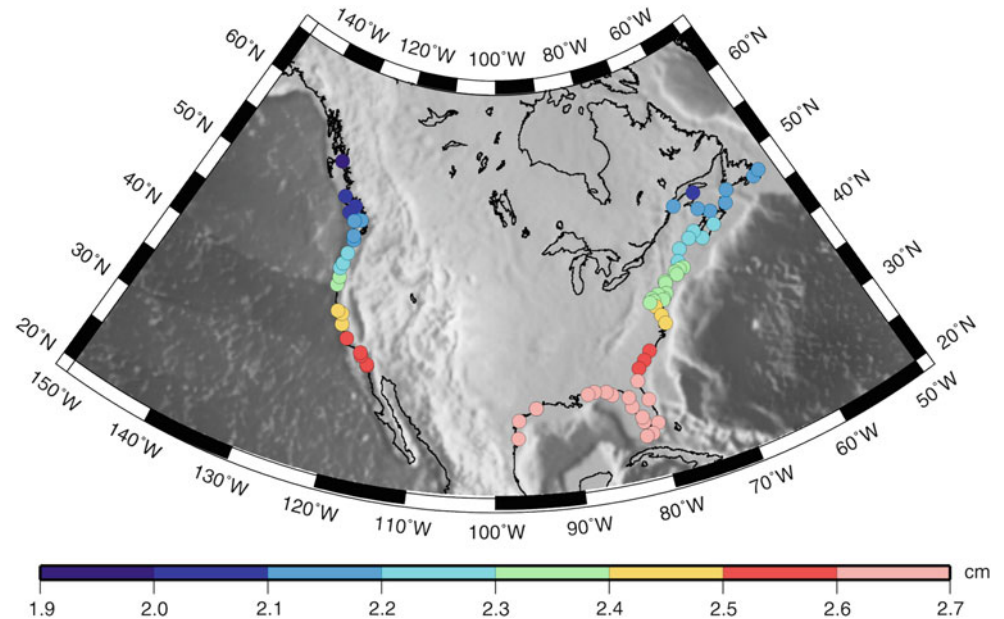
The magnitude of the GOCE geoid commission error  $\sigma_{NNGOCE}$  computed by propagating the GOCE-based time-wise GGM errors by the m-block approach (Gerlach and Fecher 2012) up to degree and order 180 has a noteworthy north-south variation. The geoid commission error in Fig. 3 varies from 1.9 cm at the northernmost North American

tide gauges to 2.7 cm at the southernmost tide gauges. A more substantial variation of 1.3 cm is produced by the propagated GOCE-based GGM errors by means of the diagonal approach (Gerlach and Fecher 2012).

**4.2.3 GOCE Omission Error**

The omission error of the GOCE-based time-wise release 4 GGM (TIM4 for short, Pail et al. 2011) of degree and order 180 is estimated by means of EGM2008 from degree and order 181 to the maximum degree 2190. Statistics of the

**Fig. 3** Geoid commission error computed at the PSMSL North American tide gauges from the time-wise GOCE GGM of a maximum degree and order 180



**Table 1** Statistics of the GOCE TIM4 geoid omission error values at the North American tide gauges computed with EGM2008

Region	Mean (cm)	Std (cm)	Min (cm)	Max (cm)
Canada Atlantic (7 TGs)	1	33	−56	38
US Atlantic (28 TGs)	−3	37	−84	72
Canada Pacific (5 TGs)	6	40	−36	64
US Pacific (17 TGs)	7	33	−49	74
Gulf of Mexico (13 TGs)	3	26	−62	28

**Table 2** Statistics of the extended GOCE TIM4 geoid omission error values at the North American tide gauges computed with USGG2012 and CGG2010 (brackets)

Region	Mean (cm)	Std (cm)	Min (cm)	Max (cm)
Canada Atlantic (7 TGs)	−3 (0)	6 (5)	−11 (−5)	8 (6)
US Atlantic (28 TGs)	−3 (2)	3 (5)	−8 (−8)	4 (15)
Canada Pacific (5 TGs)	1 (2)	6 (5)	−6 (−2)	8 (10)
US Pacific (17 TGs)	−3 (−6)	4 (6)	−10 (−16)	4 (9)
Gulf of Mexico (13 TGs)	−2 (0)	3 (5)	−7 (−5)	4 (9)

omission error at the tide gauges (12 in Canada and 58 in the USA) in five coastal regions are given in Table 1. It can be seen that the omission error has very large maxima and minima of a few dm, but it tends to cancel out at the tide gauges. The omission error of TIM4 of degree and order 180 extended with EGM2008 to the maximum degree 2190 is also estimated (Table 2) by means of comparison with the gravimetric geoids of Canada (CGG2010, Huang and Véronneau 2013) and the USA (USGG2012, <https://www.ngs.noaa.gov/GEOID/USGG2012/>). The omission error of the extended TIM4 geoid at the North American tide gauges is at the level of 3–6 cm. As the standard deviation shows, the extended TIM4 geoid agrees better with CGG2010 at the Canadian tide gauges and with USGG2012 at the US tide gauges depending on the quality of the local gravity

information used by the national geodetic agencies and some differences in the computational procedures of the two gravimetric geoids.

## 5 Vertical Datum Offsets

Examples of the mean CGVD28 and NAVD88 datum offsets with respect to the reference defined by  $W_o = 62636856.0 \text{ m}^2 \text{ s}^{-2}$  are given in Table 3. The mean datum offsets are computed by means of the least-squares adjustment model given by Eqs. (11) and (12), but an identity matrix is assumed for  $Q_{ll}$  because of the lack of information about accuracies of the ellipsoidal heights and the CGVD28 and NAVD88 heights of the local MSL. One should also take



**Table 3** Mean vertical datum offsets with respect to the level  $W_o = 62636856.0 \text{ m}^2 \text{ s}^{-2}$ 

Geoid Model	CGVD28		NAVD88 (Pacific coast)	
	$\delta N^j$ (cm)	Diff $\delta N^j$ (cm)	$\delta N^j$ (cm)	Diff $\delta N^j$ (cm)
TIM4 (degree 180)	$-28 \pm 13$	–	$-81 \pm 9$	–
TIM4 (180) + EGM2008 (2190)	$-25 \pm 10$	3	$-74 \pm 5$	7
CGG2010 ( $2' \times 2'$ )	$-24 \pm 10$	4	$-80 \pm 5$	1
USGG2012 ( $1' \times 1'$ )	$-26 \pm 10$	2	$-77 \pm 5$	4

into account that the USGG2012 gravimetric geoid is not accompanied by an error model.

The mean CGVD28 offset is computed by means of 12 Canadian tide gauges. The offset computed by means of the TIM4 geoid differs by no more than 4 cm from the offsets computed with the three high resolution geoid models. The mean NAVD88 offset is computed for the Pacific USA by means of 17 tide gauges. The large offset of 80 cm reflects the large errors of the Pacific MSL in NAVD88 as a result of the error accumulation from the datum origin in Rimouski. Differences of the mean offset computed with the TIM4 geoid and the high resolution geoid models are at the level of 1–7 cm, which shows that the GOCE geoid omission error nearly averages out at the US Pacific tide gauges.

A more realistic stochastic model is also tested, where an error of 6 cm is assumed for the difference between the GNSS ellipsoidal height and CGVD28 or NAVD88 height of the local MSL, and the geoid error is provided by the CGG2010 error model and the TIM4 geoid error model in Fig. 3. The a posteriori error, computed with Eq. (12), of the mean CGVD28 offset is 10 cm (CGG2010) and 13 cm (TIM4), similar to the a posteriori error in Table 3. The a posteriori error of the mean NAVD88 offset computed with the Pacific US tide gauges is 6 cm (CGG2010) and 11 cm (TIM4). The increase in the offset error when the TIM4 geoid is used reflects the model omission error.

## 6 Conclusions

This study demonstrated the use of Rummel and Teunissen (1988) multiple vertical datum GBVP approach in height system unification using tide gauges. For this purpose, the observation equation for the datum offset was written in terms of the ellipsoidal height of the mean sea level, the height of mean sea level in the local vertical datum and the geoid height at tide gauges. Based on the presented results, it can be concluded that

1. the indirect bias term can be omitted for North America as its value is below 1 cm if the used GOCE-based GGM is of a maximum degree and order 180;

2. with the configuration of the network of PSMSL North American tide gauges, the mean height datum offsets can be computed with an error of 10 cm or less provided that the GOCE geoid omission error is taken care of;
3. the average omission error of the geoid computed from a GOCE-based GGM extended to the full resolution of EGM2008 indicates that there is still room for improvement of the local geoid heights. Such an improvement can be achieved by designing a unified procedure for computing the residual geoid heights at the North American tide gauges with the best local gravity data and topography/bathymetry models available.

**Acknowledgements** This work is funded by ESA STSE-GOCE + Height System Unification with GOCE project and by NSERC, Canada. Phil Woodworth and Chris Hughes from National Oceanography Centre, Liverpool, UK, are acknowledged for providing the tide gauge data. Ch. Gerlach and Th. Fecher are acknowledged for the GOCE geoid commission error. The three anonymous reviewers are also acknowledged for their very useful comments and suggestions.

## References

- Amjadiparvar B, Sideris MG, Rangelova E (2013) North American height datums and their offsets: Evaluation of the GOCE-based global geopotential models in Canada and USA. *J Appl Geod* 7(3):191–203
- Cannon JB (1929) Adjustment of the Precise Level Net of Canada 1928, Publication No. 28, Geodetic Survey Division, Earth Sciences Sector, Natural Resources Canada, Ottawa, Canada
- Colombo O (1980) A World Vertical Network. Report 296, Department of Geodetic Science and Surveying, Ohio State University
- Gerlach C, Fecher F (2012) Approximations of the GOCE error variance-covariance matrix for least-squares estimation of height datum offsets. *J Geod Sci* 2(4):247–256
- Gerlach C, Rummel R (2013) Global height system unification with GOCE: a simulation study on the indirect bias term in the GBVP approach. *J Geod* 87:57–67
- Gruber T, Gerlach C, Haagmans R (2012) Intercontinental height datum connection with GOCE and GPS-levelling data. *J Geod Sci* 2(4):270–280
- Hayden T, Rangelova E, Sideris MG, Véronneau M (2012) Evaluation of  $W_0$  in Canada using tide gauges and GOCE gravity field models. *J Geod Sci* 2(4):290–301
- Heiskanen WA, Mortiz H (1967) *Physical Geodesy*. WH Freeman, San Francisco, USA, Reprint, Technical University, Graz, Austria, 1999
- Huang J, Véronneau M (2013) Canadian gravimetric geoid model 2010. *J Geod* 87:771–790

- Kotsakis C, Katsambalos K, Ampatzidis D (2012) Estimation of the zero-height geopotential level  $W_0^{LVD}$  in a local vertical datum from inversion of co-located GPS, leveling and geoid heights: a case study in the Hellenic islands. *J Geod* 86:423–439. doi:[10.1007/s00190-011-0530-7](https://doi.org/10.1007/s00190-011-0530-7)
- Pail R et al (2011) First GOCE gravity field models derived by three different approaches. *J Geod* 85:819–843
- Pavlis NK, Holmes SA, Kenyon SC, Factor JK (2012) The development and evaluation of the Earth Gravitational Model 2008 (EGM2008). *J Geophys Res* 117, B04406. doi:[10.1029/2011JB008916](https://doi.org/10.1029/2011JB008916)
- Pugh DT (1987) Tides, surges, and mean sea level. Wiley, Chichester
- Rapp RH, Balasubramania N (1992) A Conceptual Formulation of a World Height System. Report 421, Department of Geodetic Science and Surveying, Ohio State University
- Rummel R, Teunissen P (1988) Height datum definition, height datum connection and the role of the geodetic boundary value problem. *Bull Géodésique* 62:477–498
- Sansò F, Venuti G (2002) The height datum/geodetic datum problem. *Geophys J Int* 149:768–775
- Woodworth PL, Hughes CW, Bingham RJ, Gruber T (2012) Towards worldwide height system unification using ocean information. *J Geod Sci* 2(4):302–318
- Zilkoski D, Richards J, Young G (1992) Results of the General Adjustment of the North American Vertical Datum of 1988. *Survey Land Inf Syst* 52(3):133–149

---

**Part IV**

**Atmospheric Modeling in Geodesy**

---

# Computation of Zenith Total Delay Correction Fields Using Ground-Based GNSS

B. Pace, R. Pacione, C. Sciarretta, and G. Bianco

---

## Abstract

Tropospheric refraction is one of the major error sources in satellite-based positioning. The delay of radio signals caused by the troposphere ranges from 2 m at the zenith to 20 m at low elevation angles, depending on pressure, temperature and humidity along the path of the signal transmission. If the delay is not properly modelled, positioning accuracy can degrade significantly. Empirical tropospheric models, with or without meteorological observations, are used to correct these delays but they cannot model tropospheric variations exactly since they are limited in accuracy and spatial resolution resulting in up to a few decimetres error in positioning solutions. The present availability of dense ground based Global Navigation Satellite System (GNSS) networks and the state of the art GNSS processing techniques enable precise estimation of Zenith Tropospheric Delays (ZTD) with different latency ranging from Near Real-Time (NRT) to post-processing. We describe a method for computing ZTD correction fields interpolating, through Ordinary Kriging, the residuals between GNSS-derived and model-computed ZTD at continuously operating GNSS stations. At a known user location, the correction which is added to the modelled-ZTD value can be obtained through a bi-linear interpolation with the four nearest grid points surrounding it. The performance of the method has been evaluated over a 1-year period at 25 European stations belonging to the EUREF and IGS network. It is found that such an empirical tropospheric model can be improved when considering tropospheric corrections coming from ground based GNSS network.

---

## Keywords

Augmentation • GNSS positioning • Ordinary Kriging • Tropospheric model • Zenith Total Delay

---

B. Pace (✉) • R. Pacione • C. Sciarretta  
e-GEOS S.p.A. Centro di Geodesia Spaziale (CGS), Contrada  
Terlecchia, 75100 Matera, Italy  
e-mail: [bigida.pace@e-geos.it](mailto:bigida.pace@e-geos.it); [rosa.pacione@e-geos.it](mailto:rosa.pacione@e-geos.it);  
[cecilia.sciarretta@e-geos.it](mailto:cecilia.sciarretta@e-geos.it)

G. Bianco  
Agenzia Spaziale Italiana, Centro di Geodesia Spaziale (CGS),  
Contrada Terlecchia, 75100 Matera, Italy  
e-mail: [giuseppe.bianco@asi.it](mailto:giuseppe.bianco@asi.it)

---

## 1 Introduction

Tropospheric refraction is one of the major error sources in satellite-based positioning because GNSS positioning is complicated by the presence of the tropospheric propagation delay. In current positioning services, as European Geostationary Navigation Overlay System (EGNOS) in Europe, Wide Area Augmentation System (WAAS) in United States, Multi-functional Satellite Augmentation System (MSAS) in Japan, tropospheric delay corrections, unlike ionospheric corrections, are not broadcast to the user. The delays are

supposed to be corrected locally using an empirical tropospheric model adopted by the users. These models are based on the estimates of five meteorological parameters: pressure, temperature, Water Vapour (WV) pressure, temperature lapse rate and WV lapse rate which depends on user's height, latitude and day of the year (Penna et al. 2001; Collins and Langley 1997; Ueno et al. 2001). The residual delay after modelling is at a level of a few cm in the zenith direction, which may lead to a Single Point Positioning error of up to a few dm (Santerre 1991).

The present availability of dense ground based GNSS networks and the state of the art GNSS processing techniques enable precise estimation of ZTD with different latency ranging from Near Real-Time, for hourly assimilation into Numerical Weather Prediction models (Bennitt and Jupp 2012), to post-processing, useful for climate studies (Ning et al. 2012). In Europe a GNSS ground-based water vapour network has been established in the framework of the E-GVAP project (<http://egvap.dmi.dk>) set to provide its EUMETNET members with European GNSS delay and water vapour estimates for operational meteorology in Near Real-Time. The E-GVAP network consists of more than 1,800 GNSS sites, mainly in Europe.

Following the general idea outlined in Zheng et al. (2005), we propose a method for estimating ZTD corrections based on Ordinary Kriging, which takes the residuals between GNSS-derived and model-computed ZTD at continuously operating GNSS stations as input. At a known user location, the correction which is added to the modeled-ZTD value can be obtained through a bi-linear interpolation with the four nearest grid points surrounding it.

The outline of the paper is as follows: in Sect. 2 we describe the GNSS Tropo Grid Creator; we assess its performance with respect to IGS final tropospheric solution (Byun and Bar-Sever 2009), radiosonde (RS) and Very Long Baseline Interferometry (VLBI) ZTD estimates in Sect. 3. Conclusions are drawn in Sect. 4.

## 2 GNSS Tropo Grid Creator

To generate ZTD correction fields we need ZTD residuals between GNSS-derived and model-computed ZTD at continuously operating ground-based GNSS stations. Since the Zenith Hydrostatic Delay (ZHD) can be modelled and removed with an accuracy of a few millimetres (Saastamoinen 1972), the residual tropospheric delay remaining after applying a tropospheric model is mostly due to the wet component.

The GNSS-derived ZTD estimates used are the E-GVAP Italian Space Agency (ASI) NRT solutions. A detailed description of the processing strategy together with a quality

assessment of the NRT atmospheric parameters are reported in Pacione (2005) and Pacione and Vespe (2008).

The empirical tropospheric delay model used in this study is the UNB3m model, developed at the University of New Brunswick Leandro et al. (2006), which is capable of predicting ZTD with a bias value of  $-0.5$  cm and a standard deviation (STD) of 4.9 cm. The first step in the UNB3m algorithm is to obtain the meteorological parameters for a particular latitude and day of the year.

Once all meteorological parameters are determined, the ZHD and ZWD are computed according to the following equations:

$$ZHD_{MOD} = \frac{10^{-6}k_1R}{g_m}P_0\left(1 - \frac{\beta H}{T_0}\right)^{\frac{g}{R\beta}} \quad (1)$$

$$ZWD_{MOD} = \frac{10^{-6}\left(T_{m0}k'_2 + k_3\right)R}{g_m\lambda' - \beta R}\frac{e_0}{T_0}\left(1 - \frac{\beta H}{T_0}\right)^{\frac{\lambda'g}{R\beta} - 1} \quad (2)$$

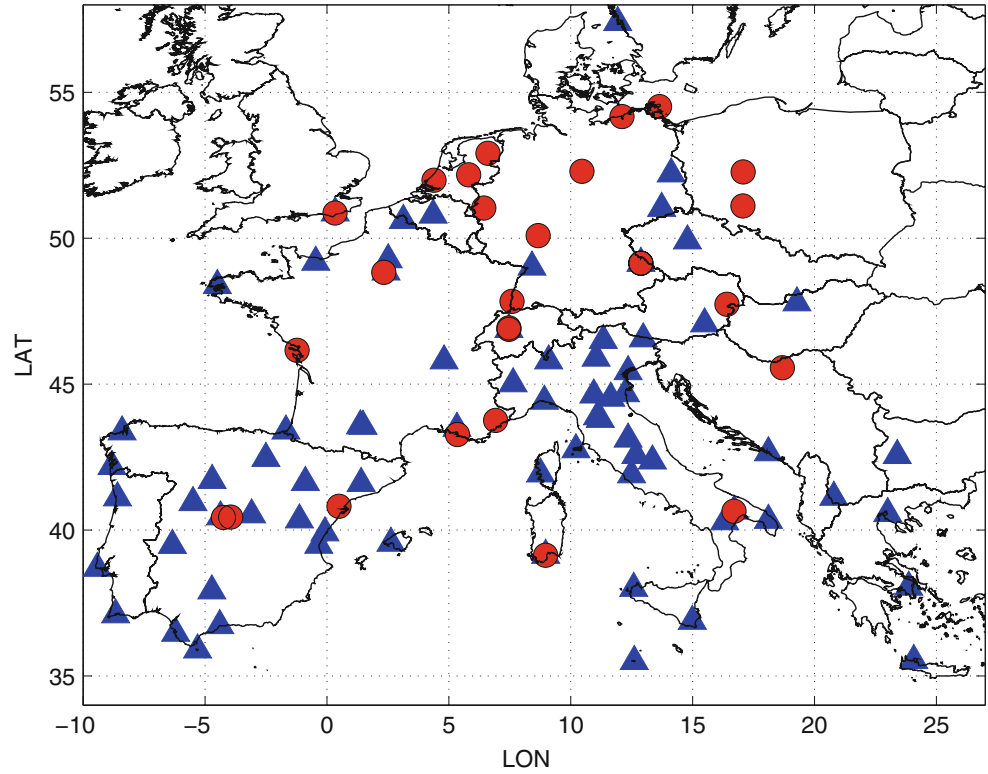
where 'MOD' stands for *modelled*,  $T_0$  (temperature),  $P_0$  (pressure),  $e_0$  (water vapour pressure),  $\lambda$  (water vapour pressure height factor) and  $\beta$  (temperature lapse rate) are the meteorological parameters,  $e_0$  is carried out following the IERS conventions (2003),  $H$  is the orthometric height in m,  $R$  is the gas constant for dry air,  $g_m = 9.784(1 - 2.66 \times 10^{-3} \cos(2\phi) - 2.8 \times 10^{-7}H)$  is the acceleration of gravity at the atmospheric column centroid,  $g$  is the surface acceleration of gravity,  $\lambda' = \lambda + 1$ ,  $T_m = T\left(1 - \frac{\beta R}{g_m\lambda'}\right)$  is the mean temperature of water vapour in Kelvin and  $k_1$ ,  $k'_2$  and  $k_3$  are refractivity constants.

Gridded corrections are obtained through Ordinary Kriging, which takes the residuals between GNSS-derived and model-computed ZTD at continuously operating GNSS stations as input. The interpolation is done over a geographical area spanning  $[35^\circ, 55^\circ]$  in latitude and  $[-10^\circ, 20^\circ]$  in longitude, both with  $0.5^\circ$  spacing. Ordinary Kriging is a powerful spatial interpolation technique, especially for irregularly spaced data points, and is widely used throughout the earth and environmental sciences. It uses known values in neighbourhood and a variogram to determine the unknown values of the location being estimated. In this study the variogram is based on a spherical correlation function between points, with a radius of 50 km, and the interpolation is based on a linear regression model.

At a user location site-ZTD correction ( $RES$ ) is obtained through a bi-linear interpolation with the four nearest grid points surrounding it:

$$RES = \sum_{i=1}^4 \omega_i RES_i \quad (3)$$

**Fig. 1** GNSS network considered for the GNSS tropo grid creator evaluation. Circle sites are the stations belonging to the European Permanent Network and International GNSS Service Network. Triangle sites are the input GNSS ZTD data



where  $RES_i$  are the gridded corrections and

$$\omega(x, y) = x^2 y^2 (9 - 6x - 6y + 4xy)$$

is the general weight function with  $x$  and  $y$ , positions of the point within the proper grid cell, calculated from:

$$x = \frac{\Delta\lambda}{\text{longitude grid interval}}, \quad y = \frac{\Delta\phi}{\text{latitude grid interval}}$$

Finally, site-ZTD is the sum of site-ZTD correction (equation (3)) and modelled-ZTD value (obtained as sum of Eqs. (1) and (2)).

### 3 GNSS Tropo Grid Creator Evaluation

The performance of the method has been evaluated over a 1-year period (January–December 2011) considering 25 European stations belonging to the European Permanent Network (EPN, <http://www.epncb.oma.be>, Bruyninx et al. 2001) and International GNSS Service (IGS, <http://igsceb.jpl.nasa.gov/>, Dow et al. 2009) Network (circle sites in Fig. 1, triangle sites are the input GNSS ZTD data). At those 25 stations we compute site-ZTD.

#### 3.1 Comparisons Against IGS ZTD Values

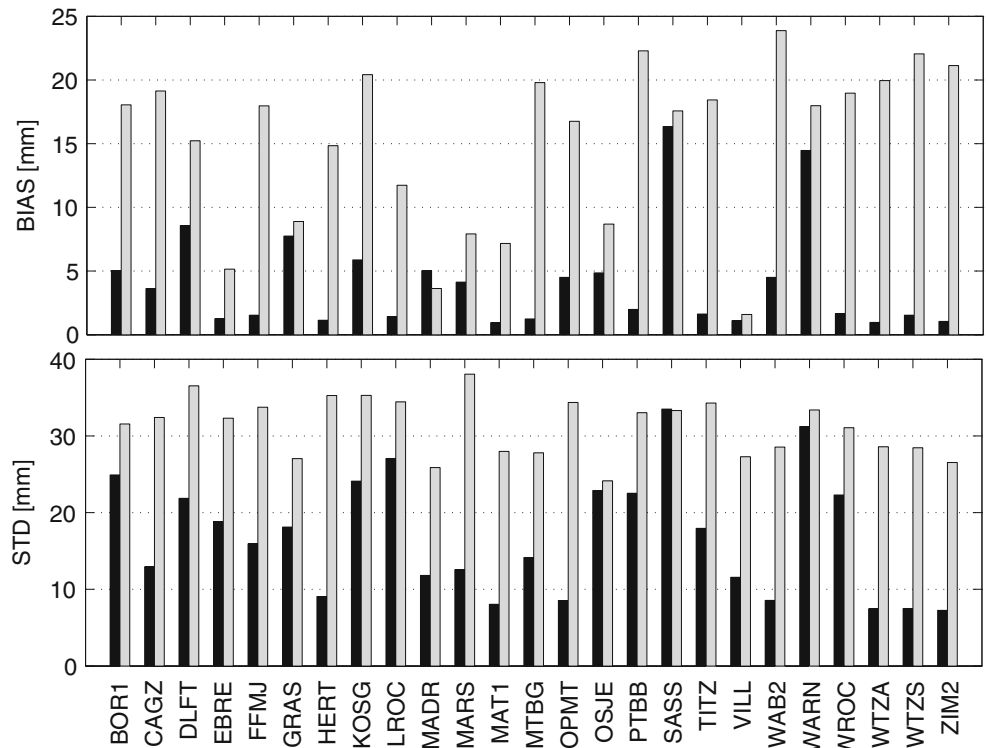
The evaluation is done via a comparison to IGS post-processing tropospheric products (Bar-Sever and Byun 2010; Byram 2011).

Figure 2 shows the statistical comparison of UNB3m-ZTD values (in grey) and site-ZTD (in black) with respect to IGS ZTD estimates for all the 25 test sites. The upper figure reports the absolute values of biases, while the bottom figure plots the standard deviation values. A decrease of about 30% for the bias and 50% for the STD is shown when site-ZTD, rather than UNB3m-ZTD values, are compared with respect to IGS estimates.

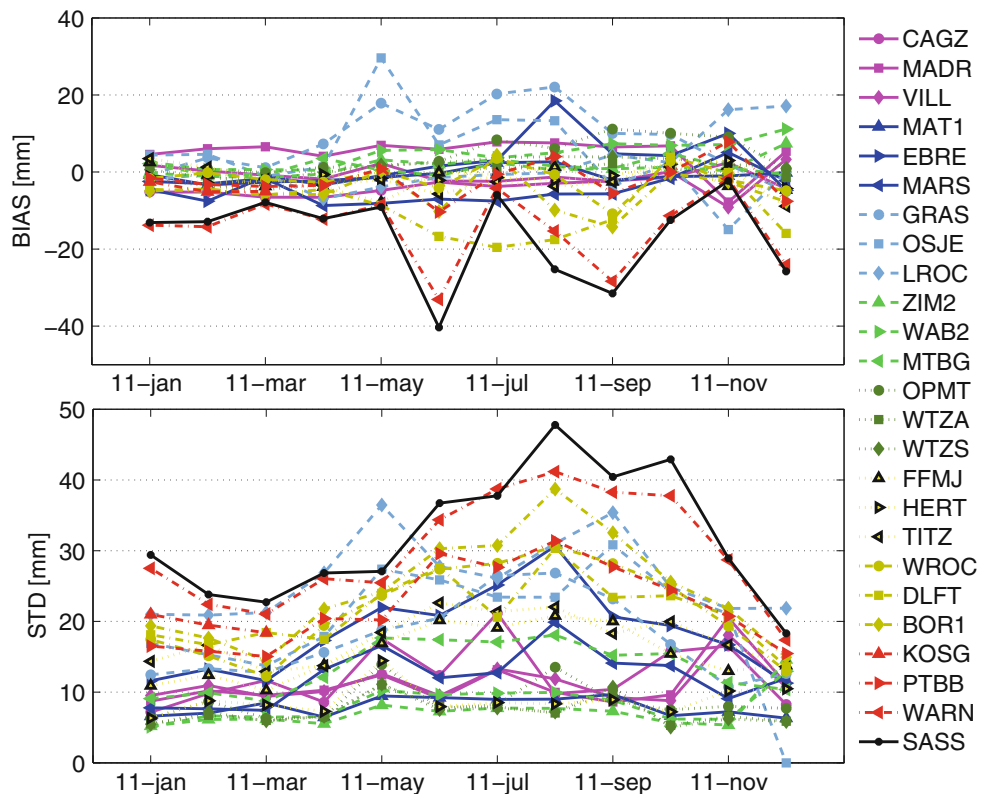
The monthly variation of the IGS ZTD values versus site-ZTDs for each test site is analyzed. Sites are sorted according to increasing latitude and increasing orthometric height to study the dependence of site-ZTD from these parameters.

Considering the latitude dependence (Fig. 3), we find that the STD increases during the summer months having the largest values for sites in the northern part of Europe that is at the boundary of the considered geographical area. As far as the height dependence is taken into account (Fig. 4), the standard deviation ranges from 5 to 50 mm with the largest values for sites having the lowest heights located in coastal areas, probably experiencing the highest humidity.

**Fig. 2** Statistical comparison of UNB3m-ZTD values (in grey) and site ZTD (in black) with respect to IGS ZTD values. Absolute values of bias (top), standard deviation (STD) (bottom)



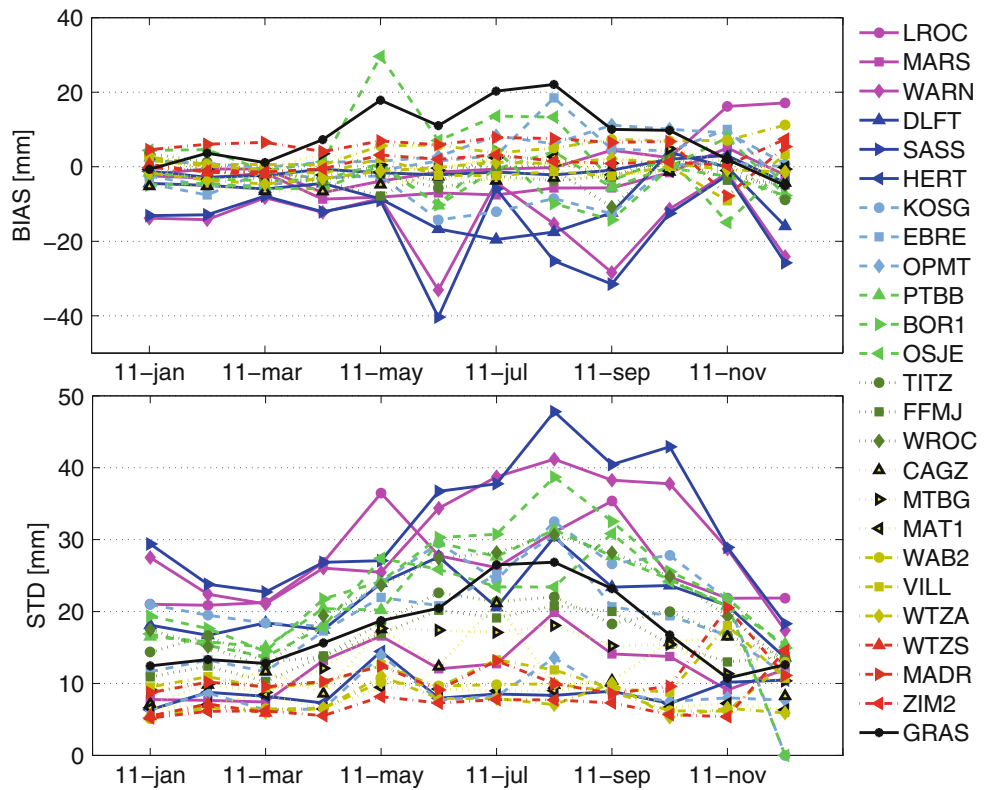
**Fig. 3** IGS ZTD vs. site-ZTD – Monthly bias (top) and STD (bottom). Sites sorted in the legend according to increasing latitude



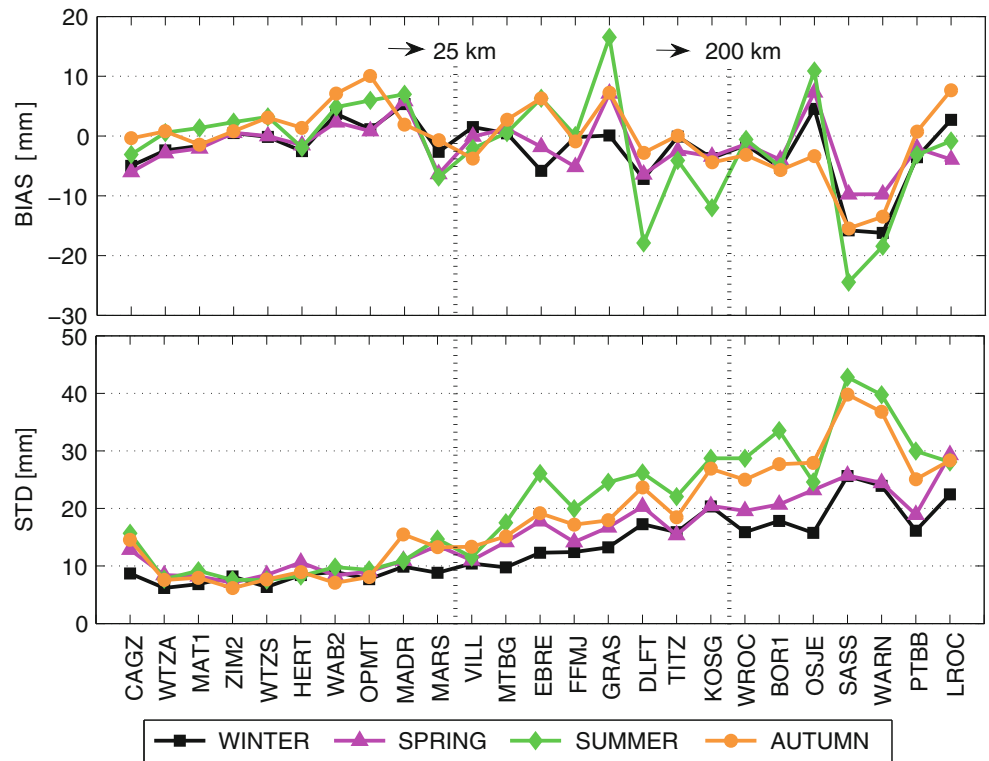
Since the GNSS sites we use as input of the GNSS tropo grid creator are not homogeneously distributed, we study how site-ZTD accuracy changes as function of the distance from the nearest GNSS input site. As can be expected the

STD increases for sites most distant from the nearest GNSS input site. In Fig. 5 sites are sorted according to increasing distances from the nearest GNSS input site and the seasonal bias and STD between IGS-ZTD and site-ZTD are plotted.

**Fig. 4** IGS ZTD vs. site-ZTD – Monthly bias (*top*) and STD (*bottom*). Sites sorted in the legend according to orthometric height



**Fig. 5** IGS ZTD versus site-ZTD – Seasonal bias (*upper*) and STD (*lower*). Sites sorted according to increasing distances w.r.t. the nearest GNSS input site. For each season the following calendar months are considered: autumn from September to November; winter from December to February; spring from March to May; summer from June to August



For each site, seasonal bias and STD are computed averaging over all the available data throughout the considered season. The bias is  $\pm 5$  mm with few exception in the summer period and for some stations at the boundary of the considered area

having a distance larger than 200 km with respect to the nearest input GNSS sites. The seasonal STD increases with the distance being in the range of [5;15] mm till 25 km, [10;30] mm till 200 km and [15;45] mm till 300 km. The



**Table 1** Radiosonde versus site-ZTD – Annual statistics. For each GPS Site is reported the corresponding radiosonde code (values in column 2) and the number of synchronized data used in comparison (values in column 5). Distance (values in column 6) is the horizontal

	RS code	BIAS (mm)	STD (mm)	# Sample	Distance from RS (km)	Distance from the GNSS (km)
HERT	3882	10.6	10.6	386	03.42	0.150
WROC	12425	1.5	23.2	578	14.49	221.738
CAGZ	16560	-0.9	17.5	666	14.49	0.002
VILL	8221	3.5	14.2	545	29.56	32.405
ZIM2	6610	5.9	11.1	573	41.02	0.019

largest STD values are found during the summer period, which can be related to the atmospheric seasonal cycle.

### 3.2 Radiosonde ZTD Versus Site-ZTD

Radiosondes measure directly humidity, temperature and pressure during their ascent, enabling determination of ZTD. They are ideal for validation of GNSS ZTDs, and have been used for that over many years (see, e.g., Vedel et al. 2001; Haase et al. 2003; Wang and Zhang 2008, 2009).

Radiosonde profiles come from World Meteorological Organization (WMO) Global Telecommunication Service (GTS) and are provided by the Danish Meteorological Institute in the framework of E-GVAP as independent data set to validate GPS ZTD data. The radiosonde profiles are passed through a program (Haase et al. 2003) that checks the quality of the profiles, converts the dew point temperatures to specific humidity, transforms the radiosonde profile to correct for the altitude offset between the GPS and the radiosonde sites and determines ZTD, ZWD and Integrated Water Vapour (IWV) compensating for the change of gravitational acceleration with height.

The annual bias and STD for 5 sites where nearby radiosonde profiles are available is reported in Table 1. Among them HERT is the closest to the radiosonde launch site (3.42 km) while ZIM2 is the most distant (41.02 km). The bias (Table 1) is positive, except for CAGZ, meaning that site-ZTD is dryer than radiosonde ZTD in agreement with what reported in Pacione et al. (2011). The STD ranges from 10 to 23 mm. It is larger than the STD reported in other studies (Pacione et al. 2011; Dousa and Bennitt 2012) where GPS and radiosonde data are compared. Moreover it is larger than the STD obtained comparing IGS and radiosonde data for the 5 sites in the same period which ranges from 6.5 mm for WROC up to 12.57 mm for CAGZ. Except CAGZ, the obtained STD values seem to be not only related to the horizontal distance between the radiosonde launch and the GPS antenna but also to the distance from the nearest GNSS input site.

distance between the radiosonde launch and the GPS antenna. Distance (values in column 7) is the horizontal distance between the nearest GNSS Tropo Grid Creator input site and the GPS antenna

**Table 2** VLBI versus site-ZTD Annual statistics (January–December 2011). For each GPS Site is reported the number of synchronized data used in comparison (values in column 5) and the correlation coefficients (CC, values in column 4)

	BIAS (mm)	STD (mm)	CC	# Sample
MAT1	0.05	9.04	0.97	875
WTZA	-0.30	9.88	0.97	2,092
WTZS	-0.14	9.88	0.97	2,092

### 3.3 VLBI ZTD Versus Site-ZTD

VLBI data provide another basis for comparison. Three test sites, namely MAT1 (Italy), WTZA (Germany) and WTRS (Germany), are co-located with VLBI radio-telescope antenna. The VLBI solutions are the ASI/CGS contribution to the International VLBI Service (IVS) tropospheric products for IVS-R1 and IVS-R4 weekly 24-h sessions.

Site-ZTD and VLBI ZTD estimates are very highly correlated, with an overall bias of the ZTD differences  $-0.13$  mm (see Table 2) and a STD of about 10 mm, which is larger than the median standard deviation of about 5 mm over all sites reported in Teke et al. 2011 where IGS and IVS ZTD estimates during the CONT08 campaign have been compared. The 5 mm standard deviation value is confirmed comparing IGS and VLBI data for the three sites in the same period.

## 4 Conclusion

Satellite-based positioning application is complicated by the presence of the tropospheric delay which is supposed to be corrected locally by empirical tropospheric models. In this paper we describe a method for estimating ZTD correction fields by using ground-based GNSS ZTD estimates. We evaluate it over a 1-year period demonstrating that such ZTD correction fields can augment empirical tropospheric models resulting in an improvement of about 30% for the bias and 50% for the standard deviation with respect to the IGS final tropospheric solutions.

The GNSS Tropo Grid Creator is running operationally using the E-GVAP ASI NRT solutions as input on hourly basis and providing ZTD correction fields over Europe in a IONEX-like format (<http://igsceb.jpl.nasa.gov/igsceb/data/format/ionex1.pdf>). Beyond the proposed method, ZTD corrections can be obtained from Numerical Weather Prediction (NWP) Model with the same level of accuracy (Teke et al. 2011) site-ZTD has.

NWP ZTD and site-ZTD can be both used in positioning service but comparison campaigns as well as field trials deserve to be carried out to highlight pros and cons of each method.

**Acknowledgements** We acknowledge the anonymous reviewers and the Editors for their comments which led to significant improvements in the manuscript.

## References

- Bar-Sever Y, Byun S (2010) Reprocessed IGS trop product now available with gradients. [IGSMail-6298]
- Bennitt GV, Jupp A (2012) Operational assimilation of GPS Zenith total delay observations into the Met Office Numerical Weather Prediction Models. *Mon Weather Rev* 140:2706–2719
- Bruyninx C, Becker M, Stangl G (2001) Regional densification of the IGS in Europe Using the EUREF Permanent GPS Network (EPN). *Phys Chem Earth* 26:531–538
- Byram S (2011) IGS Final Troposphere Product Transition to USNO. [IGSMail-6443]
- Byun SH, Bar-Sever YE (2009) A new type of troposphere Zenith Path Delay Product of the International GNSS Service. *J Geod* 83:367–373
- Collins JP, Langley RB (1997) A tropospheric delay model for the user of the wide area augmentation system. Final contract report prepared for Nav Canada, Department of Geodesy and Geomatics Engineering Technical Report No.187, University of New Brunswick, Fredericton, N.B., Canada
- Dousa J, Bennitt GV (2012) Estimation and evaluation of hourly updated global GPS Zenith Total Delays over ten months. *GPS Solutions*. doi:10.1007/s10291-012-0291-7
- Dow JM, Neilan RE, Rizos C (2009) The international GNSS service in a changing landscape of global navigation satellite systems. *J Geod* 83:191–198
- Haase J, Ge M, Vedel H, Calais E (2003) Accuracy and variability of GPS tropospheric delay measurements of water vapour in the Western Mediterranean. *J Appl Meteorol* 42:1547–1568
- Leandro RF, Santos MC, Langley RB (2006) UNB neutral atmosphere models: development and performance. In: *Proceedings of the Institute of Navigation National Technical Meeting*, 18–20 January, 2006, Monterey, CA, USA, pp 564–573
- Ning T, Elgered G, Willén U, Johansson J (2012) Evaluation of the atmospheric water vapor content in a regional climate model using ground-based GPS measurements. *J Geophys Res*. doi:10.1029/2012JD018053
- Pacione R (2005) ASI Analysis Center. COST 716 Exploitation of Ground-Based GPS for Operational Numerical Weather Prediction and Climate Applications Final Report EUR 21639, edited by G. Elgered et al., pp 35–37
- Pacione R, Vespe F (2008) Comparative studies for the assessment of the quality of near-real time GPS-derived atmospheric parameters. *J Atmos Oceanic Technol* 25:701–714
- Pacione R, Pace B, Vedel H, de Haan S, Lanotte R, Vespe F (2011) Combination methods of tropospheric time series. *Adv Space Res* 47:323–335
- Penna N, Dodson A, Chen W (2001) Assessment of EGNOS Tropospheric Correction Model. *J Navig* 54:37–55
- Saastamoinen J (1972) Atmospheric correction for the troposphere and stratosphere in radio ranging of satellites. In: Henriksen SW et al (eds) *The use of artificial satellites for Geodesy*, vol 15. *Geophysics Monograph Series*. AGU, Washington DC, pp 247–251
- Santerre R (1991) Impact of GPS satellite sky distribution. *Manuscr Geodet* 16:28–53
- Teke K, Böhm J, Nilsson T, Schuh H, Steigenberger P, Dach R, Heinkelmann R, Willis P, Haas R, García-Espada S, Hobiger T, Ichikawa R, Shimizu S (2011) Multi-technique comparison of troposphere Zenith Delays and gradients during CONT08. *J Geod* 85:395–413
- Ueno M, Hoshino K, Matsunaga K, Kawai M, Nakao H, Langley RB (2001) Assessment of atmospheric delay correction models for the Japanese MSAS. In: *Proceedings of the ION GPS 2001*, 14th International Technical Meeting of the Satellite Division of The Institute of Navigation, Salt Lake City, 11–14 September 2001, pp 2341–2350
- Vedel H, Mogensen KS, Huang XY (2001) Calculation of Zenith Delays from meteorological data comparison of NWP model, radiosonde and GPS delays. *Phys Chem Earth* 26:497–502
- Wang J, Zhang L (2008) Systematic errors in global radiosonde precipitable water data from comparisons with ground-based GPS measurements. *J Clim* 21:2218–2238
- Wang J, Zhang L (2009) Climate applications of a global, 2-hourly atmospheric precipitable water dataset derived from IGS tropospheric products. *J Geod* 83:209–217
- Zheng Y, Feng Y, Bai Z (2005) Grid residual tropospheric corrections for improved differential GPS positioning over the Victoria GPS Network (GPSnet). *J Global Positioning Syst* 4:284–290

---

# Rigorous Interpolation of Atmospheric State Parameters for Ray-Traced Tropospheric Delays

Camille Desjardins, Pascal Gegout, Laurent Soudarin, and Richard Biancale

---

## Abstract

The transformation between European Center for Medium-range Weather Forecast (ECMWF) model level assimilations and the refractivity at any given point of the neutral atmosphere has been investigated. We first present the IFS interpolations and extrapolations of each physical parameter done in operations at ECMWF. These formulae are used to compute, for example, pressure levels from model levels at ECMWF. We use this formulation to compute the pressure levels, the large majority of which are found similar to the pressure levels provided by ECMWF with an appropriate accuracy for ray-tracing. The IFS-based scheme (IFS-BS) is then presented. It is an adaptation of the interpolations and extrapolations done at ECMWF for troposphere delay computation by ray-tracing. This scheme ensures the coherence with the ECMWF meteorological model and is used in our software Horizon designed to compute the Adaptive Mapping Functions (AMF). In the IFS-BS, vertical interpolations are adapted for each thermodynamic parameter necessary to precisely rebuild the refractivity along the ray path according to the physical laws. In order to take into account the atmospheric part between the lowest model level and the Earth's topography during the ray-tracing, extrapolation of physical parameters below the lowest model level are included. The proposed scheme is expected to be relevant for applications where accuracy of refractivity is important as troposphere delay modelling for high-accuracy geodesy.

---

## Keywords

Adaptive mapping function • AMF • Atmospheric modelling • Numerical weather model • Ray tracing • Troposphere delay

---

C. Desjardins (✉)  
Centre National d'Etudes Spatiales (DCT/SI/MO - BPI 811) Collecte  
Localisation Satellites 18, avenue Edouard Belin 31401 Toulouse  
Cedex 9, France  
e-mail: [camille.desjardins@cnes.fr](mailto:camille.desjardins@cnes.fr)

P. Gegout  
Observatoire Midi-Pyrénées (CNRS/GET/GS), 14 avenue Edouard  
Belin, 31400 Toulouse, France  
e-mail: [pascal.gegout@get.obs-mip.fr](mailto:pascal.gegout@get.obs-mip.fr)

L. Soudarin  
Collecte Localisation Satellites (CLS/DOS/OBS/DOG), 11 rue  
Hermès, 31520 Ramonville St Agne, France  
e-mail: [laurent.soudarin@cls.fr](mailto:laurent.soudarin@cls.fr)

---

## 1 Introduction

For space geodetic measurements based on radio ranges such as GNSS or DORIS, the tropospheric delays induced by the neutral part of the atmosphere are still an important source of error. Indeed, estimates of neutral atmosphere are highly correlated with site displacements and receiver clock biases.

---

R. Biancale  
Centre National d'Etudes Spatiales (CNES/SI/GS), 18, avenue  
Edouard Belin, 31401 Toulouse Cedex 9, France  
e-mail: [richard.biancale@cnes.fr](mailto:richard.biancale@cnes.fr)

Accurate models of atmospheric delays have to be used to mitigate such effects. It is commonly accepted to model tropospheric delays by calculating the zenith tropospheric delay and obtaining the slant tropospheric delays with a mapping function. New mapping functions have been developed in the 2000s (Boehm et al. 2006a,b; Niell 2001) and significantly improve the geodetic positioning. Although modern mapping functions are derived from numerical weather models (NWM), most of these mapping functions ignore the azimuth dependency which is usually introduced by two horizontal gradient parameters—in north-south and east-west directions—estimated directly from observations (Chen and Herring 1997). More recently, the use of ray-traced delays through NWM directly at observation level has shown an improvement on geodetic results (Hobiger et al. 2008, 2010; Nafisi et al. 2012; Zus et al. 2012).

During the design of the Adaptive Mapping Functions (AMF) detailed by Gegout et al. (2011), our goal was to use the most information available in NWM—especially the azimuth dependency—with the aim to preserve the classical mapping function strategy. AMF are thus used to approximate thousands of atmospheric ray-traced delays using a few tens of coefficients with millimetre accuracy at low elevation. AMF have a classical form with terms which are function of the elevation  $\epsilon$  (Eq. 1). But, they also include coefficients which depend on the azimuth  $\alpha$  to represent the azimuthal dependency of ray-traced delays (Eq. 2). In addition, AMF are suitable to adapt to complex weather by changing the truncation of the successive fractions.

$$\text{AMF}(\alpha, \epsilon) = S_f \times \frac{1 + \frac{a_1}{1 + \frac{a_2}{1 + \dots}}}{\sin \epsilon + \frac{a_1}{\sin \epsilon + \frac{a_2}{\sin \epsilon + \dots}}} \quad (1)$$

where  $\epsilon$  is the elevation angle and  $S_f$  is the scale factor which can be empirically adjusted to observations.

$$\forall i \geq 1, a_i = a_{i0} + \sum_{j=1}^n C_{ij} \cos j\alpha + S_{ij} \sin j\alpha \quad (2)$$

We discuss here the vertical interpolation and extrapolation of each tropospheric propagation depending parameters.

## 2 Description of the ECMWF Model Level Data

### 2.1 The ECMWF Integrated Forecasting System (IFS)

The Integrated Forecasting System (IFS) is the ECMWF global meteorological forecasting model. The IFS coordinates are geographic latitude and longitude for horizontal

and the hybrid coordinate  $\eta$  for the vertical. In the following, these coordinates are called meteorological coordinates. The vertical hybrid coordinate  $\eta(P, P_s)$  introduced by Simmons and Burridge (1981), is a terrain-following monotonic function of the pressure  $P$  and also depends on the surface pressure  $P_s$  such that  $\eta(P_s, P_s) = 1$  and  $\eta(0, P_s) = 0$ . Pressure as a function of  $\eta$  is given by

$$P(\eta) = A(\eta) + B(\eta) \times P_s \quad (3)$$

$$\text{where } A(1) = 0, B(1) = 1 \text{ and } A(0) = B(0) = 0. \quad (4)$$

The ECMWF model uses a spectral method with spherical harmonics basis functions and triangular truncation for horizontal discretization. For the vertical, the model divides the atmosphere into  $NLEV$  layers from the model surface to  $P = 0$ . When this study was carried out,  $NLEV$  was equal to 91 in operational data. The vertical discretization is currently a finite-element scheme with cubic B-spline expansion based on (3) (Untch and Hortal 2004). Until January 2002, the finite-difference scheme defined by Simmons and Burridge (1981) was used operationally at ECMWF. The set of fixed constant coefficients  $A_{k+1/2}$  and  $B_{k+1/2}$  with  $0 \leq k \leq NLEV$  is the finite-difference discretization of  $A(\eta)$  and  $B(\eta)$ . Because only this set is provided by ECMWF, we consider here the finite-difference scheme.

## 2.2 The Model Levels

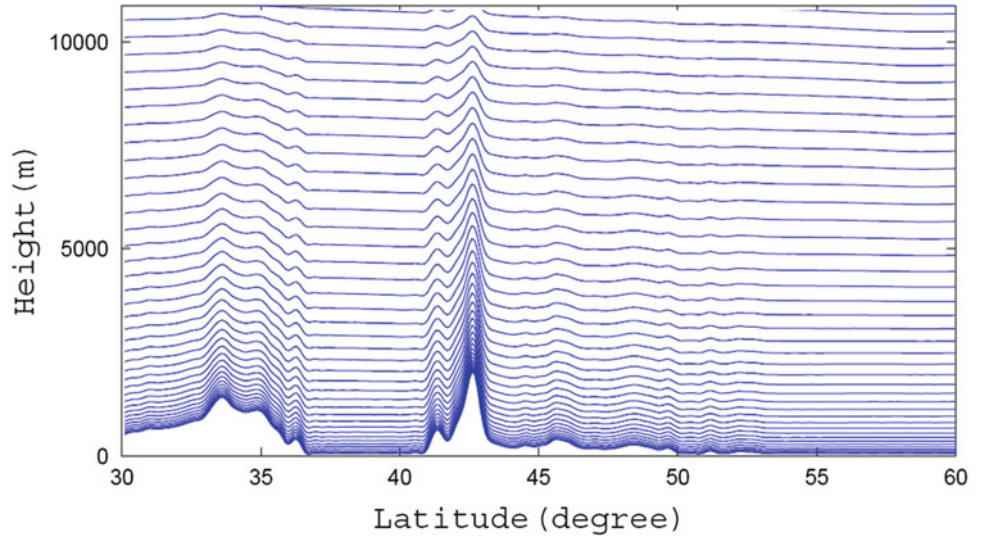
The model levels are the native levels of the IFS. They follow the model surface in the lower atmosphere and are isobars in the upper part (Fig. 1). They are not constant in height or pressure. The model level data are archived in spherical harmonics or reduced Gaussian grids. The fields are for one epoch: the pressure  $P_s$  and the geopotential  $\Phi_s$  at the model surface, the temperature  $T_k$  and the specific humidity  $q_k$  for each level  $k \in \{1, 2, \dots, NLEV\}$ . The pressures  $P_{k+1/2}$  (Eq. 5) at the interface between layers—called half-model levels—and the pressure  $P_k$  (Eq. 6) at the middle of each layer—called full-model levels or simply model levels—are recovered using  $A_{k+1/2}$  and  $B_{k+1/2}$  values. Geopotentials  $\Phi_{k+1/2}$  at half-model levels are rebuilt using the discrete analogue of hydrostatic equilibrium (Eq. 7) (ECMWF 2012).

$$0 \leq k \leq NLEV, P_{k+1/2} = A_{k+1/2} + B_{k+1/2} \times P_s \quad (5)$$

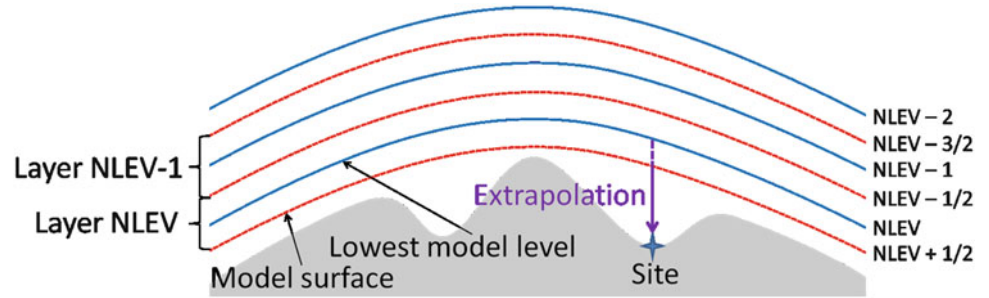
$$1 \leq k \leq NLEV, P_k = \frac{P_{k+1/2} + P_{k-1/2}}{2} \quad (6)$$

$$0 \leq k \leq NLEV, \Phi_{k+1/2} = \Phi_s + \sum_{j=k+1}^{NLEV} R_d (T_v)_j \times \ln \left( \frac{P_{j+1/2}}{P_{j-1/2}} \right) \quad (7)$$

**Fig. 1** Model levels in function of latitude at longitude 1.48°



**Fig. 2** Difference between half-model levels (red line) and model levels (blue line). Illustration of extrapolation



where  $(T_v)_j$  is the virtual temperature at layer  $j$

$$(T_v)_j = \left[ 1 + \left( \frac{R_w}{R_d} - 1 \right) q_j \right] T_j \quad (8)$$

and  $R_d$  and  $R_w$  denote respectively the specific gas constants of the dry air and the water vapour. Geopotentials  $\Phi_k$  at each model level are given by

$$1 \leq k \leq NLEV, \Phi_k = \Phi_{k+1/2} + \alpha_k R_d (T_v)_k \quad (9)$$

$$\text{with } \begin{cases} \alpha_1 = \ln 2 \\ \forall k > 1, \alpha_k = 1 - \frac{P_{k-1/2}}{P_{k+1/2} - P_{k-1/2}} \ln \left( \frac{P_{k+1/2}}{P_{k-1/2}} \right). \end{cases} \quad (10)$$

The difference between half-model levels and model levels are illustrated in Fig. 2. The model surface, also called orography by the meteorologists, defines the envelope of the real topography. The conditions (4) ensure that  $A_{NLEV+1/2} = 0$  and  $B_{NLEV+1/2} = 1$ . So using (5), it comes  $P_{NLEV+1/2} = P_s$ . The orography is thus the lowest half-model level.

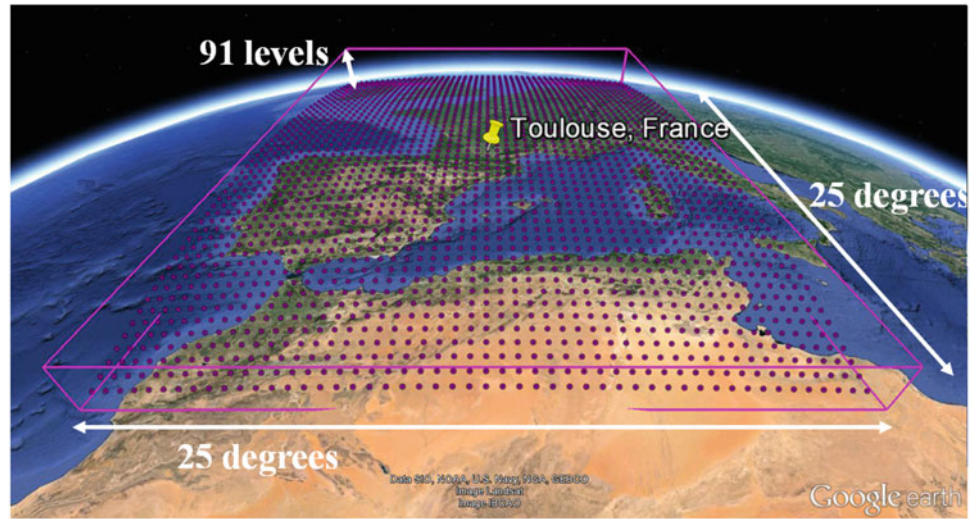
### 3 The Horizon Software and Ray-Tracing Technique

#### 3.1 The Horizon Software

The Horizon software has been developed to compute AMF for space geodetic sites from ECMWF model level data. For each site and epoch when the ECMWF data are available, one AMF is computed following the next steps:

1. *Read ECMWF model levels*: Model level fields are read and interpolated on a regular  $0.125^\circ \times 0.125^\circ$  grid centred on the site (Fig. 3). The sub-grid size is  $25^\circ$  for all sites. It is an empirical value to ensure that rays stay inside the volume even at low elevation.
2. *Rebuild the pressures and geopotentials*: The pressure and geopotential at each half- and full-model level are recovered using Eqs. (5)–(10).
3. *Convert meteorological into geodetic coordinates*: We need to convert the meteorological coordinates used in the IFS into geodetic coordinates used in the Horizon

**Fig. 3** The sub-grid of Toulouse (France) site used for ray-tracing



software. We precisely rebuild the shape of the atmosphere from geopotentials using a realistic geodetic model described by Gegout et al. (2011). Moreover, we participate in discussions on this topic under the umbrella of the International Association of Geodesy—Intercommission Committee on Theory Joint Studying Group 0.4 “Coordinate systems in numerical weather models” (see <http://hobiger.org/blog/iag-ic-ssg12/>).

4. *Compute tropospheric delays by ray-tracing technique:* Using the data pre-processed by the three first steps, ray-traced tropospheric delays of the required site are computed with  $4^\circ$  azimuth steps and non-regular elevation steps. The ray-tracing principle done here for each azimuth and elevation is: define initial conditions—site position, elevation and azimuth, compute next point using the Eikonal equation and so on until the top of the atmosphere empirically defined as a sphere of radius 6,450 km. The total tropospheric delay is relative to the distance from the site to the top end of the ray. The delay is reassessed when the parallax problem is solved for, when the precise target’s position is adjusted. At the end of this step, ray-traced delays are available for all azimuths and elevations. Hardware accelerations to reduce time consumption are discussed by Gegout et al. (2013).
5. *Determine AMF coefficients from ray-traced delays:* The ray-traced delays are used to determine AMF coefficients:  $a_{i0}$ ,  $C_{ij}$  and  $S_{ij}$  (Eq. 2). The problem is non-linear and over-determined. The Levenberg-Maquardt algorithm is used to solve it. The fit residuals defined here as the absolute difference between AMF values and ray-traced delays are typically less than 2 mm (Gegout et al. 2011).

### 3.2 Focus on Ray-Tracing Computation

During ray-tracing, the Eikonal equation is solved at each step. The Eikonal equation depends on the refractivity at the current point which is needed at any point along the ray path. The refractivity  $N$  can be expressed as a function of the pressure  $P$ , specific humidity  $q$  and temperature  $T$ :

$$N = k_1 \frac{P}{T_v} + k'_2 \frac{P_w}{T} + k_3 \frac{P_w}{T^2} \quad (11)$$

where  $P_w$  is the partial pressure of water vapour

$$P_w = \frac{q}{\frac{R_d}{R_w} + \left(1 - \frac{R_d}{R_w}\right) q} \times P \quad (12)$$

$T_v$  is the virtual temperature (Eq. 8) and  $k_1$ ,  $k'_2$  and  $k_3$  are empirical coefficients. Following Cucurull (2010), we have chosen to use the  $k_1$ ,  $k'_2$  and  $k_3$  values determined by Bevis et al. (1994). Because the refractivity depends on meteorological parameters available only on sub-grid discrete points, the refractivity computation is based on the physical values of the eight neighbour points of the current point. Vertical and horizontal interpolations are required to transform meteorological parameters of the neighbour points into the refractivity at the current point. Horizontal gradients of refractivity are hundred times smaller than vertical gradients. Bilinear interpolations may be sufficient for horizontal interpolations but not for vertical interpolations discussed later.

### 3.3 On the Use of ECMWF Model Levels

Because the model levels are not constant in height or pressure (Sect. 2.2), some precautions have to be taken to avoid vertical interpolation errors. Moreover, using model levels addresses the problem of extrapolation because the orography differs from topography (Fig. 2). The difference between these two layers can be several hundred of meters in mountainous regions. The layer between orography and topography is a dense part of atmosphere which may contain a lot of humidity. Due to this fact, this atmospheric layer can have an important impact on tropospheric delays. We so need to take into account this atmospheric layer in the best possible way on the modelling of tropospheric delays. Using model levels to compute refractivity raises some questions: How to interpolate within model levels? How to extrapolate below the orography?

### 3.4 Vertical Interpolation and Extrapolation Strategies

To provide the refractivity at the current point, our first approach was to interpolate or extrapolate exponentially in height the refractivities of the neighbour points and then interpolate horizontally. Directly interpolating the refractivity in height supposes that the refractivity varies exponentially in height between model levels. Hobiger et al. (2008) suggested that this assumption might not be reasonable because the refractivity depends on three physical parameters (Eq. 11) which have different characteristics with height and the refractivity is not computed by a linear interpolation of these parameters. These led us to develop a new vertical interpolation strategy, preserving the physical laws and IFS native discretizations. Our new approach is based on the vertical interpolations and extrapolations defined in the IFS.

## 4 Interpolations and Extrapolations Done in the IFS

In this section, we describe how the physical parameters—the geopotential  $\Phi$ , the specific humidity  $q$  and the temperature  $T$ —useful to compute the refractivity are interpolated in IFS (ECMWF 2012) in function of the pressure  $P$  and the meteorological values of the two closest model levels.

### 4.1 Geopotential

#### 4.1.1 Interpolation

The geopotential  $\Phi$  at a given pressure  $P$  is computed from model level data using the International Civil

Aviation Organization (ICAO) temperature profile (ICAO 1993; ECMWF 2012). First, the ICAO temperature  $T_k^{ICAO}$  and the ICAO geopotential  $\Phi_k^{ICAO}$  at model levels are computed using the standard temperature profile. Integrating the hydrostatic equation provides  $\Delta\Phi_k$  (Eq. 13) which is the difference between model level geopotential and the ICAO standard atmosphere at each model level.

$$\Delta\Phi_k = \sum_{j=NLEV}^{k+1} R_d \left( (T_v)_j - T_j^{ICAO} \right) \ln \frac{P_{j+1/2}}{P_{j-1/2}} + \alpha_k R_d \left( (T_v)_k - T_k^{ICAO} \right) \quad (13)$$

$(T_v)_k$  and  $\alpha_k$  are respectively defined in Eqs. (8) and (10). Then, the difference at the required pressure  $\Delta\Phi_P$  is obtained by vertical interpolation from  $\Delta\Phi_k$ . The interpolation is linear in  $\ln P$  between model levels. After computing the geopotential  $\Phi^{ICAO}$  at the required pressure and the geopotential  $\Phi_s^{ICAO}$  at the orographic pressure, the geopotential at the required pressure is obtained by

$$\Phi = \Phi_s + \Delta\Phi_P + \Phi^{ICAO} - \Phi_s^{ICAO}. \quad (14)$$

#### 4.1.2 Extrapolation Above the Highest Model Level

The geopotential is computed in the same way as interpolation assuming a constant  $\Delta\Phi_P = \Delta\Phi_1$ .

#### 4.1.3 Extrapolation Below the Orography

The geopotential is

$$\Phi = \Phi_s - \frac{R_d T_s}{\Gamma} \left[ \left( \frac{P}{P_s} \right)^\Gamma - 1 \right] \text{ where } \Gamma = \frac{\Lambda R_d}{g} \quad (15)$$

and  $T_s$  is the temperature at the orography

$$T_s = \left[ 1 - \frac{\Lambda R_d}{g} \left( \frac{P_s}{P_{NLEV-1}} - 1 \right) \right] T_{NLEV-1}. \quad (16)$$

$\Lambda = -0.0065 \text{ K} \cdot \text{m}^{-1}$  is the constant ICAO temperature gradient and  $g = 9.80665 \text{ m} \cdot \text{s}^{-2}$  is the standard gravity adopted by the World Meteorological Organization (WMO 2008).

### 4.2 Temperature

#### 4.2.1 Interpolation

The temperature  $T$  is linearly interpolated in pressure to the required pressure  $P$  between the two closest model levels, here  $k$  and  $k-1$ .

$$T = T_{k-1} + \frac{T_k - T_{k-1}}{P_k - P_{k-1}} (P - P_{k-1}) \text{ with } P_{k-1} \leq P \leq P_k \quad (17)$$

**Table 1** Number of points in percent by intervals of absolute difference  $\Delta q$  of specific humidity in  $\text{kg}\cdot\text{kg}^{-1}$  (left) and of absolute difference  $\Delta T$  of temperature in K (right)

Difference of specific humidity ( $\text{kg}\cdot\text{kg}^{-1}$ )	Number of points (%)	Difference of temperature (K)	Number of points (%)
$\Delta q < 10^{-6}$	99.97	$\Delta T < 0.1$	98.87
$10^{-6} < \Delta q < 10^{-5}$	0.03	$0.1 < \Delta T < 0.5$	1.13
$\Delta q > 10^{-5}$	< 0.01	$\Delta T > 0.5$	< 0.01

#### 4.2.2 Extrapolation Above the Highest Model Level

The temperature  $T$  is assumed to be constant and equal to the value of the highest model level  $T_1$ .

#### 4.2.3 Extrapolation Below the Lowest Model Level

Between the lowest model level and the orography, the temperature is linearly interpolated between  $T_s$  and  $T_{NLEV}$ . Below the orography, the temperature is extrapolated by a third-order polynomial in the logarithm of pressure (Eq. 18).

$$T = \left[ 1 + \Gamma \ln \frac{P}{P_s} + \frac{1}{2} \left( \Gamma \ln \frac{P}{P_s} \right)^2 + \frac{1}{6} \left( \Gamma \ln \frac{P}{P_s} \right)^3 \right] T_s \quad (18)$$

### 4.3 Specific Humidity

#### 4.3.1 Interpolation

The specific humidity  $q$  is linearly interpolated in pressure to the required pressure  $P$  between the two closest model levels, here  $k$  and  $k - 1$ .

$$q = q_{k-1} + \frac{q_k - q_{k-1}}{P_k - P_{k-1}} (P - P_{k-1}) \text{ with } P_{k-1} \leq P \leq P_k \quad (19)$$

#### 4.3.2 Extrapolations

Below the lowest and above the highest model level, the specific humidity  $q$  is assumed to be constant and equal to  $q_{NLEV}$  and  $q_1$  respectively.

### 4.4 Conversion from Model to Pressure Levels

ECMWF also produces operational data on other vertical discretizations than model levels, for example in pressure levels. The pressure levels are computed from model levels at ECMWF by post-processing. So converting model levels into pressure levels leads to a loss of vertical resolution: when the operational model included 91 model levels, there were only 25 pressure levels. To validate our implementation of the vertical interpolations and extrapolations used by the IFS

post-processing and to point out numerical and modelling errors, we test our ability to retrieve pressure levels from model levels. We first compute values at each pressure level from model levels using interpolations and extrapolations above-mentioned for one epoch (August 2nd, 2009 at 9 a.m.) and for each point of the reduced Gaussian grid. Then, we compare obtained values with ECMWF pressure level data. The results (Table 1) show that the retrieval is done with an appropriate accuracy for the ray-tracing for more than 98% of the points. However, it is important to note that this validation only shows the way to retrieve pressure levels from model levels. The difference between tropospheric delays computed from model levels and those derived from pressure levels is not considered here. For this issue, further investigations have to be done.

## 5 IFS Formulation Adapted for Ray-Tracing: The IFS-Based Scheme (IFS-BS)

### 5.1 Adaptation of the IFS Formulation for Ray-Tracing

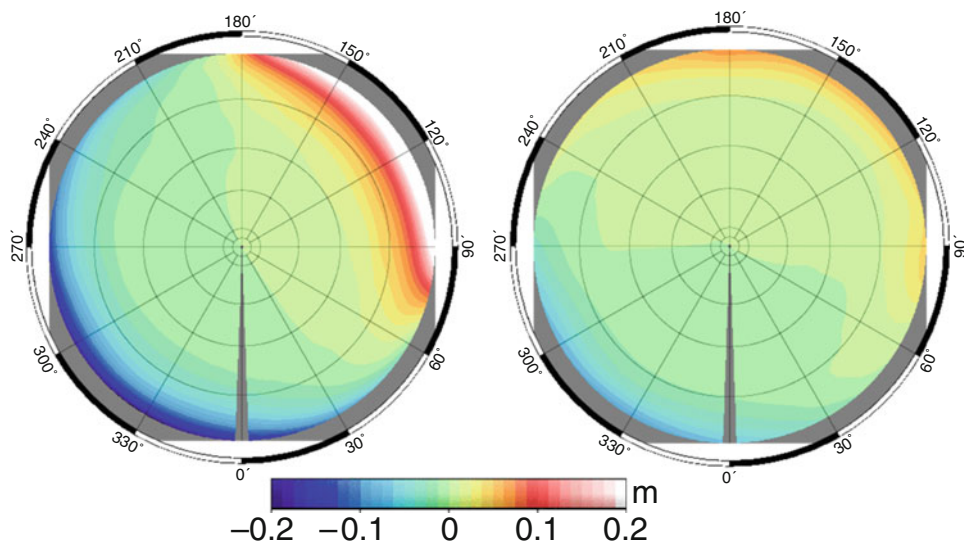
The vertical coordinate is the pressure in the scheme described in Sect. 4 whereas it is the geodetic height in the Horizon software. We have to adapt this scheme to use it in the ray-tracing step of our software. We choose to convert geodetic height into pressure via geopotential and then vertically interpolate or extrapolate the temperature and the specific humidity in pressure as described in Sects. 4.2 and 4.3. First, at the required point, the geopotential is computed from the geodetic height with the reciprocal function of the conversion which permits to have geodetic height from geopotential in the step 3 of the Horizon software (Sect. 3.1). Then, pressure is computed from geopotential at the required point. The applied method is based on the integration of the hydrostatic equilibrium (Eq. 20).

$$P = P_{k+1/2} \exp \left( \frac{\Phi_{k+1/2} - \Phi}{R_d (T_v)_k} \right) \quad (20)$$

Equation (20) is still valid above the highest model level but not below the orography. In this case, the transformation is deduced from the geopotential extrapolation done in the IFS



**Fig. 4** Azimuthal variability of the ray-traced tropospheric delay for the GNSS station located in Arequipa, Peru on May 1st, 2013 using the exponential scheme (left) and the IFS-based scheme (right)



rearranging Eq. (15) to get  $P$  (Eq. 21).

$$P = \left(1 + \frac{\Phi_s - \Phi}{R_d T_s} \Gamma\right)^{\frac{1}{\Gamma}} P_s \quad (21)$$

## 5.2 Refractivity Retrieved from the IFS-Based Scheme

The basic idea of the IFS-based scheme (IFS-BS) is to interpolate or extrapolate vertically and then horizontally each physical parameter separately at the required height from the values of the neighbour points and then compute the refractivity at the required point using interpolated values of pressure, specific humidity and temperature. The vertical interpolations and extrapolations used here are described in Sects. 4.2, 4.3 and 5.1. The IFS-BS has the advantage to adapt vertical interpolation and extrapolation strategies for each physical parameter according to the physical laws and be consistent with the vertical discretization done in the IFS.

## 5.3 Assessment of the IFS-Based Scheme

To investigate the impact of the vertical interpolations and extrapolations on tropospheric delays, the IFS-BS (Sect. 5.2) and the exponential scheme (Sect. 3.4) are compared. The difference between the two schemes can be significant especially in mountainous regions. For example, at Arequipa, Peru, where the height difference between

orography and topography is 520 m, the difference on zenith tropospheric delay between the two schemes is 5 mm. In addition, the 5° slant tropospheric delays obtained with the exponential scheme have a larger azimuthal anisotropy than when the IFS-BS is used (Fig. 4). On the contrary, at Toulouse, France, where the height difference between orography and topography is 30 m, the difference on zenith tropospheric delay is less than 0.1 mm and there is no difference on the azimuthal variability of the 5° slant tropospheric delays. The IFS-BS with its physically-based extrapolation is more consistent, especially in mountainous regions where the height difference between orography and topography is large. So, using the IFS-BS can be an improvement on the modelling of tropospheric delays.

## 6 Conclusion and Perspectives

The IFS-BS described here is based on the interpolation and extrapolation scheme adopted by ECMWF for its global meteorological forecasting model called IFS. Using this formulation ensures the coherence with model level data and permits to adapt vertical interpolations and extrapolations for each physical parameter according to its own physical law. Practical investigations of a large number of situations and sites are undergoing to provide a realistic accuracy versus GNSS measurements of the IFS-BS and the AMF. The liquid and ice water contents are not considered in this formulation although integrated in the IFS. We plan to include these parameters in our modelling in order to continue step-by-step improvements in the Horizon software.

**Acknowledgements** The ECMWF is acknowledged for providing the atmospheric three-dimensional operational data used in this study.

The Fig. 3 was performed using the Google™ Earth software.

The Horizon software and Adaptive Mapping Functions are developed by the project “Surcharges & Propagations” with the financial support of the TOSCA/CNES program. The research and technology work of Camille Desjardins was supported by CNES and CLS grants.

The authors thank three anonymous reviewers for their valuable and helpful comments.

## References

- Bevis M, Businger S, Chiswell S, Herring T, Anthes R, Rocken C, Ware R (1994) GPS meteorology: mapping zenith wet delays onto precipitable water. *J Appl Meteorol* 33(3):379–386
- Boehm J, Niell A, Tregoning P, Schuh H (2006a) Global mapping function (GMF): a new empirical mapping function based on numerical weather model data. *Geophys Res Lett* 33(7):L07304
- Boehm J, Werl B, Schuh H (2006b) Troposphere mapping functions for GPS and very long baseline interferometry from European Centre for Medium-Range Weather Forecasts operational analysis data. *J Geophys Res* 111:B02406
- Chen G, Herring TA (1997) Effects of atmospheric azimuthal asymmetry on the analysis of space geodetic data. *J Geophys Res* 102(B9):20,489–20,502
- Cucurull L (2010) Improvement in the use of an operational constellation of GPS radio occultation receivers in weather forecasting. *Weather Forecast* 25(2):749–767
- ECMWF (2012) IFS Documentation Cy38r1. European Centre for Medium-Range Weather Forecasts, Reading, UK. <http://www.ecmwf.int/research/ifsdocs/CY38r1/>
- Gegout P, Biancale R, Soudarin L (2011) Adaptive mapping functions to the azimuthal anisotropy of the neutral atmosphere. *J Geod* 85(10):661–677
- Gegout P, Oberlé P, Desjardins C, Moyard J, Brunet PM (2013) Ray-tracing of GNSS signal through the atmosphere powered by CUDA, HMPP and GPU's technologies. *IEEE J Sel Top Appl Earth Observations Remote Sens.* doi:10.1109/JSTARS.2013.2272600
- Hobiger T, Ichikawa R, Koyama Y, Kondo T (2008) Fast and accurate ray-tracing algorithms for real-time space geodetic applications using numerical weather models. *J Geophys Res* 113(D20):D20302
- Hobiger T, Shimada S, Shimizu S, Ichikawa R, Koyama Y, Kondo T (2010) Improving GPS positioning estimates during extreme weather situations by the help of fine-mesh numerical weather models. *J Atmos Sol Terr Phys* 72(2–3):262–270
- ICAO (1993) Manual of the ICAO standard atmosphere: extended to 80 km (262500 ft). International Civil Aviation Organization, Montreal, QC, Canada
- Nafisi V, Madzak M, Boehm J, Ardalan A, Schuh H (2012) Ray-traced tropospheric delays in VLBI analysis. *Radio Sci* 47(2):17
- Niell A (2001) Preliminary evaluation of atmospheric mapping functions based on numerical weather models. *Phys Chem Earth Part A* 26(6–8):475–480
- Simmons AJ, Burridge DM (1981) An energy and angular-momentum conserving vertical finite-difference scheme and hybrid vertical coordinates. *Mon Weather Rev* 109(4):758–766
- Untch A, Hortal M (2004) A finite-element scheme for the vertical discretization of the semi-lagrangian version of the ECMWF forecast model. *Q J R Meteorol Soc* 130(599):1505–1530
- WMO (2008) Guide to meteorological instruments and methods of observation, 7th edn. WMO-No. 8, World Meteorological Organization, Geneva, Switzerland
- Zus F, Bender M, Deng Z, Dick G, Heise S, Shang-Guan M, Wickert J (2012) A methodology to compute GPS slant total delays in a numerical weather model. *Radio Sci* 47(2):RS2018

---

# Comparison of Different Techniques for Tropospheric Wet Delay Retrieval Over South America and Surrounding Oceans

A. Calori, G. Colosimo, M. Crespi, and M.V. Mackern

---

## Abstract

Water vapour (WV) plays a fundamental role in several weather processes that deeply influence human activities. Satellite based radiometers, Ground based Global Navigation Satellite Systems (GNSS) and Numerical Weather Models (NWM) permit to obtain either measurements or estimates or forecasts of WV. This work presents a 2 years systematic comparison to address the agreement on the tropospheric wet delay retrieved by the three mentioned independent techniques over permanent stations belonging to SIRGAS (Sistema de Referencia para las Américas) GNSS network. SIRGAS tropospheric total delay estimations are compared with the official International GNSS Service (IGS) ones, with the measurements from the Jason-1 satellite radiometer (JMR) in terms of Zenith Wet Delays (ZWD) and, finally, with the ZWD computed from ERA Interim, the last reanalysis dataset from the European Center for Medium-Range Weather Forecasts (ECMWF). All the differences between the techniques, which were considered in order to yield a reliable comparison, are discussed. The statistical results of mean ( $\mu$ ), standard deviation ( $\sigma$ ) and correlation ( $\rho$ ), show that the highest agreement is reached between SIRGAS and IGS products ( $\mu = -0.5$  mm,  $\sigma = 5.6$  mm,  $\rho = 0.98$ ), whereas slightly worse values are obtained in the comparisons with the JMR measurements ( $\mu = -7.4$  mm,  $\sigma = 15.4$  mm,  $\rho = 0.91$ ), and the ERA Interim data ( $\mu = -1.5$  mm,  $\sigma = 16.6$  mm,  $\rho = 0.91$ ).

---

## Keywords

GNSS • Jason-1 radiometer • Numerical Weather Model • SIRGAS • ZWD retrieval

---

## 1 Introduction

Water Vapour (WV) plays a fundamental role in several weather processes that deeply influence human activities and it has been recognised as the most important among

the greenhouse gases (Mitchell 1989). Several studies have confirmed how deeply water vapour is bound to climate changes, for instance by showing the high correlation between the yearly temperature variation and the WV content in the atmosphere (Wentz and Schabel 2000). It has been clearly understood that the knowledge of high accurate WV content and its distribution in the atmosphere improves short term weather forecasts significantly. At the same time, WV reveals very rapid changes both in the temporal and in the spatial domains such that, at present, there are no theoretical models that can reliably predict its behaviour.

Retrieving WV content in the atmosphere can be performed in different ways using independent techniques: starting from the more traditional and established ones, such as radiosondes and ground-based microwave radiometers, up

---

A. Calori (✉) • M.V. Mackern  
Facultad de Ingeniería, Universidad Nacional de Cuyo, Ciudad de  
Mendoza, Mendoza, Argentina  
e-mail: [acalori@mendoza-conicet.gob.ar](mailto:acalori@mendoza-conicet.gob.ar)

G. Colosimo • M. Crespi  
DICEA-Area di Geodesia e Geomatica, University of Rome  
“La Sapienza”, Roma, Italy  
e-mail: [gabriele.colosimo@uniroma1.it](mailto:gabriele.colosimo@uniroma1.it)

to the more recent ones, such as satellite based techniques like satellite radiometers (Christensen et al. 1994), Global Navigation Satellite Systems (GNSS) (Bevis et al. 1992), Radio Occultation (Kursinski et al. 1997) and Numerical Weather Models (NWM). Since each of these techniques presents advantages and limitations, researchers' efforts have been recently focused on comparing the different approaches with the aim of combining them to retrieve WV content with the highest possible accuracy. The issues addressed in this work are related to the research activities promoted by the International GNSS Service (IGS) Troposphere Working Group.

Satellite based radiometers can provide integrated water vapor (IWV) measurements at different epochs. However, their application is limited over sea and ocean surfaces, the revisit time on the same location is rather low and reliable measurements are obtained only at certain weather conditions (e.g., no rain). Ground based GNSS stations can be used to estimate the signal delay caused by tropospheric refraction (Hogg et al. 1981). This delay, which is referred to as Zenith Total Delay (ZTD), can be unfolded into two components: the Zenith Hydrostatic Delay (ZHD) and the Zenith Wet Delay (ZWD), which are due to the contribution of the hydrostatic gases and to the water vapour, respectively. The GNSS technique has been proven capable of estimating the ZTD and then using these estimates to infer the IWV with accuracies of few millimetres (e.g., Rocken et al. 1997). Moreover, thanks to its dense station networks and to the very high temporal resolution of the estimates (up to few minutes) the interest in GNSS as IWV data source is continuously increasing. NWM, such as the European Center for Medium-Range Weather Forecasts (ECMWF), exploit data from many different sources and can be used to compute and forecast IWV all over the world with a medium-high temporal resolution (i.e., a few hours).

This work presents the results of a 2 years (i.e., June 2008–2010) comparison of the three described techniques for the determination of the tropospheric wet delay ZTD and the IWV over the South and Central American region. Initially, in order to assess the performances of SIRGAS estimations, the results were compared with the official ZTD distributed by the IGS. Then, the consistency of the products was evaluated with respect to: (1) ZWD measured by Jason-1 satellite mission; (2) ZWD computed from data of the ECMWF ERA-Interim reanalysis model. Following this introductory section, Sect. 2 describes the main features of the used techniques. The results from the comparison are discussed in Sect. 3. Finally, conclusions and future research prospects are outlined in Sect. 4.

## 2 Data Processing: Retrieving the ZTD from the Different Techniques

### 2.1 Ground Based GNSS Stations

#### 2.1.1 ZTD from the SIRGAS Network

SIRGAS-CON is the regional densification of the International Terrestrial Reference Frame (ITRF) over Latin America and Caribbean, it spans a huge extension  $-71^\circ < \phi < 20^\circ, -109^\circ < \lambda < -2^\circ$ , with altitudes up to 3.770 m and, at present, it encompasses about 250 continuously operating GNSS reference stations, 48 of them belonging to the global IGS network (Brunini et al. 2012).

Within this research work, the site-specific  $ZTD_{SIR}$  were estimated for approximately 100 GNSS SIRGAS stations (SIRGAS-CON-D-SUR) (Mackern et al. 2009) with a global formal precision of few millimetres, as described in detail in Calori et al. (2013). Figure 1 shows the overall distribution of the used GNSS stations over the South American region.

The main processing features which are relevant for the next sections are as follows:

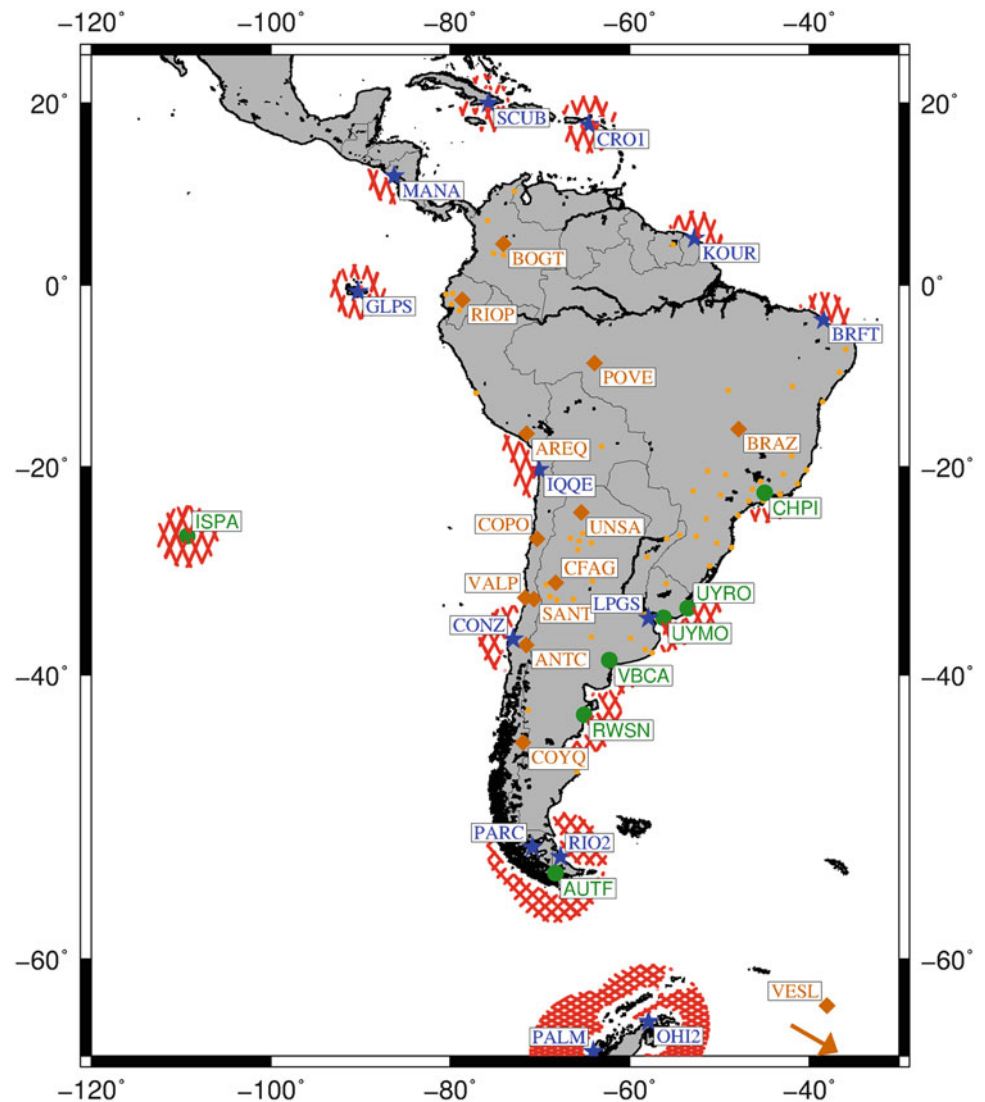
- Software: Bernese GPS Software 5.0 (Dach et al. 2009)—Differential positioning
- Elevation angle cutoff:  $3^\circ$
- Mapping function: Niell (1996) for hydrostatic and wet component
- A priori values: (Berg 1948; Saastamoinen 1972)
- Temporal resolution of ZTD estimates: 15 min (tropospheric gradients not estimated)

#### 2.1.2 IGS Tropospheric Products

Ever since 2003, a precise point positioning (PPP) approach is used within IGS to estimate ZTD values using raw GPS range measurements and the IGS Final Orbits and Clocks. This process produces one file per site per day containing a time series of ZTD with temporal resolution of 5 min and a formal precision of few millimetres (Buyn and Bar-Sever 2009). The main processing features relevant for the comparison are detailed:

- Software: GIPSY—Precise Point Positioning
- Elevation angle cutoff:  $7^\circ$
- Mapping function: Niell (1996) and GMF (Böhm et al. 2006)
- A priori values: Hydrostatic delay based on altitude (2.3 m at sea level), and 0.1 m for the wet delay
- Temporal resolution of ZTD estimates: 5 min (tropospheric gradients estimated)

**Fig. 1** Overall distribution of stations of SIRGAS-CON-D-SUR network processed within this work (*small orange circles*). Stations used in the comparison with only Jason-1 radiometer, are displayed as *green circles*. Stations used in the comparison with only IGS official products are displayed as *brown diamonds*. *Blue stars* represent stations that were used in the comparison with IGS and Jason-1 products. *Red pattern* displays the results of the automatic procedure implemented to select Jason-1 measurements that lie inside a circular area of a certain radius ( $1^\circ$ ) centred on the station



The comparison between the  $ZTD_{SIR}$  and  $ZTD_{IGS}$  was carried out over 27 sites common to both networks (Fig. 1, blue stars and brown diamonds).

## 2.2 Jason-1 and Satellite Radiometry

Jason-1 has been an altimetry satellite mission jointly operated by the French aerospace agency—Centre National d’Etudes Spatiales (CNES) and the United States National Aeronautics and Space Administration (NASA).

To retrieve the ocean topography with an accuracy of a few centimeters, Jason-1 was equipped with a Microwave Radiometer (JMR) used to measure the delay caused by the water vapor along the altimeter beam. JMR measures the brightness temperatures in the nadir direction over a circular footprint approximately between 20 and 30 km (Picot et al. 2003). Using a combination algorithm (described in p. 155 Keihm et al. 1995), the brightness temperatures can be

coupled to retrieve the delay caused by the water vapor in the atmosphere (i.e., the ZWD) with a Root Mean Square Error (RMSE) of 1.2 cm that is, however, limited to open ocean areas (Ruf et al. 1994).

For the present work, we have chosen to utilize the Geophysical Data Records (GDR) version c. Besides the altimeter measurements, GDR contain as ancillary information the hydrostatic meteorological correction (ZHD) provided by the ECMWF, which are then used to obtain the ZTD according to the standard equation

$$ZTD = ZHD + ZWD \quad (1)$$

For this research, Jason-1 GDR version c binary data corresponding to the period from June 2008 to June 2010 were downloaded from the web. Then, a tuned software was implemented to filter out only those measurements which are close to the GNSS sites. Since JMR provides reliable measurements only over open ocean areas, 20 stations were

**Table 1** Statistical results (compared product, number of stations, size of the dataset, average bias and standard deviation of the differences, correlation coefficient between the products) of the comparisons performed

Product	Stations (#)	Samples (#)	$\mu$ (mm)	$\sigma$ (mm)	$\rho$
IGS–SIR	27	1,309,868	−0.5	6.9	0.98
ERA–SIR	30	65,534	−1.5	16.6	0.91
JMR–SIR	14	1,052	−7.4	15.4	0.93
ERA–IGS	27	67,638	−2.4	14.8	0.92
JMR–IGS	11	983	−5.6	14.9	0.93
JMR–ERA	14	958	−8.5	15.5	0.94

Different techniques have been compared in terms of ZWD whereas GNSS intra-comparison refers to ZTD

selected which fulfil the geographical criteria being located within a limited distance from coastline and the height of the station. Figure 1 shows the distribution of the selected stations (green dots and blue stars) subset and the Jason-1 ground tracks (red lines) for orbit cycle 275. As described in detail in Calori et al. (2013), the differences between the radiometer measurements and the GNSS estimates were addressed to yield a reliable comparison between the techniques.

### 2.3 ERA Interim

ERA-Interim is the latest global atmospheric reanalysis product computed by the ECMWF. This contains gridded data that describe the weather as well as ocean-wave and land-surface conditions together with upper-air parameters covering the troposphere and stratosphere (Dee et al. 2011). With the purpose of retrieving the ZWD at the GNSS sites from meteorological information, the binary data in *grib* format of 3 meteorologic parameters (i.e., the mean sea level pressure ( $P_{atm}$ ), the total column water vapour (TCWV) and the 2 m temperature (2T)) were downloaded from the ECMWF web site for the time frame of the present analysis. These grids have a spatial resolution of  $0.75^\circ \times 0.75^\circ$  and a temporal resolution of 6 h (i.e., at 0, 6, 12 and 18 UTC). Here, it is also worth noting that both the TCWV and 2T are referred to the orography height ( $h_o$ ), so that some height corrections were needed to retrieve the tropospheric delays (ZHD and ZWD) at the GNSS station height. As first step, the atmospheric pressure was computed at the GNSS station height ( $h$ ) according to the standard pressure model of Berg (1948)

$$P_h = P_{atm}(1 - d \cdot h)^{5.225} \quad (2)$$

where  $d = 0.0000226$ . Then, the ZHD at the GNSS station height (i.e.,  $ZHD_{ERA,h}$ ) was retrieved following Davis et al. (1985)

$$ZHD_{ERA,h} = a \frac{P_h}{(1 - b \cdot \cos(2\phi) - c \cdot h)} \quad (3)$$

where  $a = 0.0022768$ ,  $b = 0.00266$ ,  $c = 0.28 \cdot 10^{-6}$ ,  $\phi$  is the station latitude. The mean temperature of the troposphere ( $T_m$ ) was modelled using the 2T according to Mendes et al. (2000, Eq. 17), model *UNB98Tm1*. This step was necessary in order to retrieve the ZWD at the orography height  $ZWD_{ERA,h_o}$  using the relation between the TCWV and the ZWD introduced by Askne and Nordius (1987, Eq. 25). To refer the ZWD retrieved by the ECMWF to the GNSS station height, the empirical relation proposed by Kouba (2008) was applied

$$ZWD_{ERA,h} = ZWD_{ERA,h_o} \cdot e^{-(h-h_o)/2000} \quad (4)$$

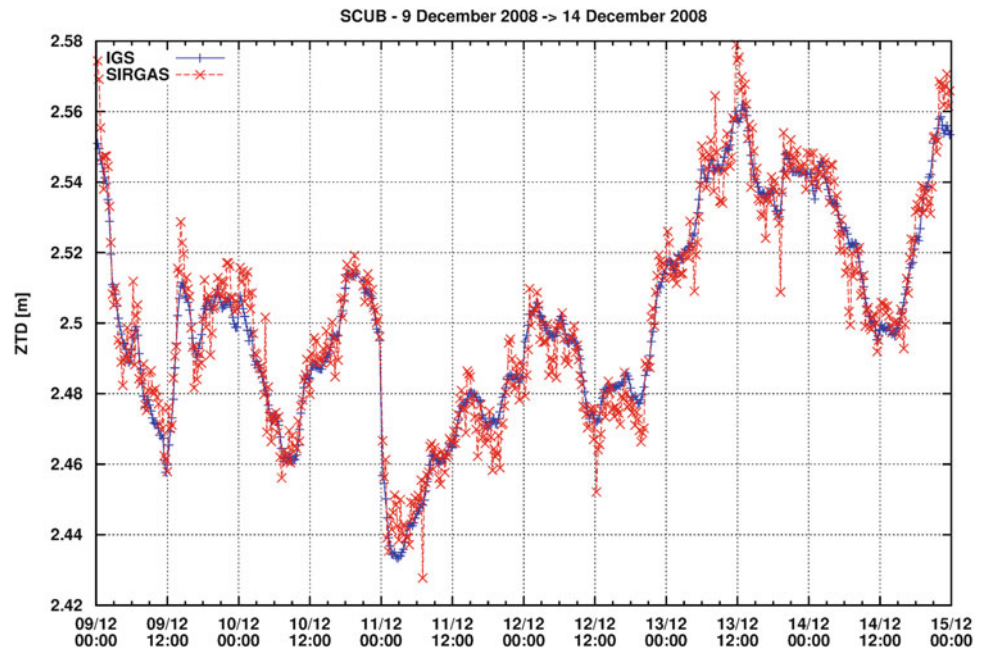
Finally, according to Eq.(1), we computed  $ZWD_{SIR}$  using  $ZHD_{ERA}$ . The comparison with  $ZWD_{SIR}$  was performed for all sites at which a comparison with either IGS or JMR values was already available (displayed as blue stars and brown diamonds in Fig. 1). At this stage of the research, no temporal interpolation was introduced so that GNSS and ECMWF were analysed only at identical times (i.e., 4 times per day).

## 3 Results and Discussion

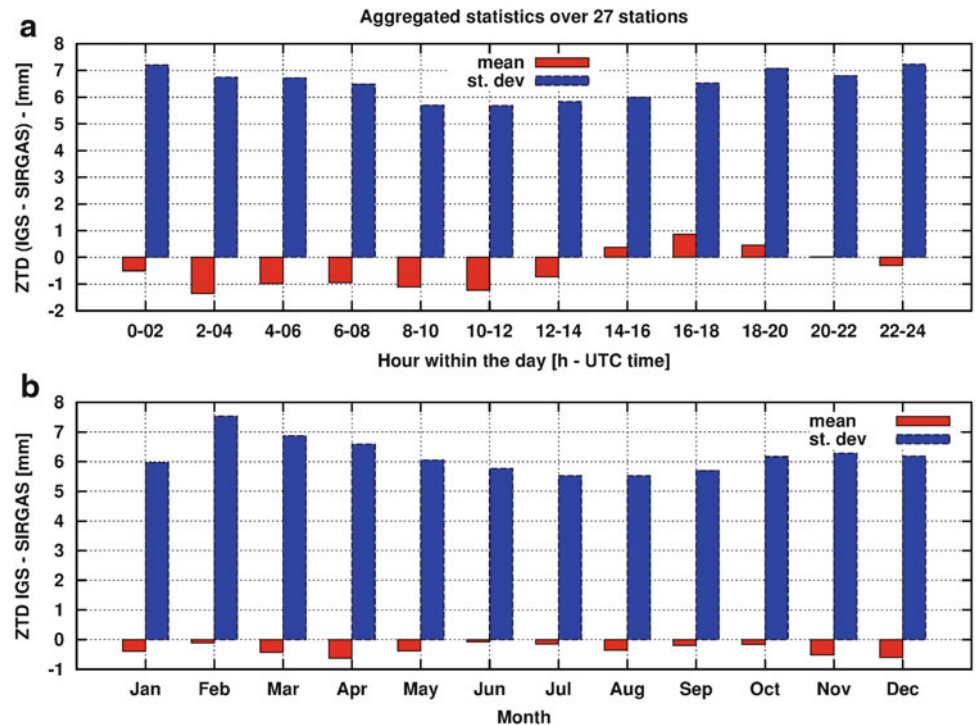
The accuracy of the tropospheric estimations retrieved from SIRGAS network was assessed in terms of consistency with three different products: (1) the official ZTD generated by IGS; (2) the ZWD computed using meteorological information provided by ERA-Interim, ECMWF; (3) the ZWD measured by the JMR aboard the Jason-1 satellite altimetry mission. The comparison was carried out from June 2008 to 2010 and, because of the inter-techniques differences, it involved separate clusters of SIRGAS stations: 27 sites for comparison 1, 30 sites for comparison 2 and 14 sites located along the coastline for comparison 3. In each comparison, the agreement between the techniques was evaluated using the bias ( $\mu$ ), the standard deviation ( $\sigma$ ) of the differences and the correlation coefficient ( $\rho$ ) of the time series as statistical indexes.

Table 1, which reports the statistical inter-techniques indicators averaged over the whole set of stations, reveals

**Fig. 2** 5 days (9–14 December 2008) time series of IGS (*blue*) and SIRGAS (*red*) ZTD estimations

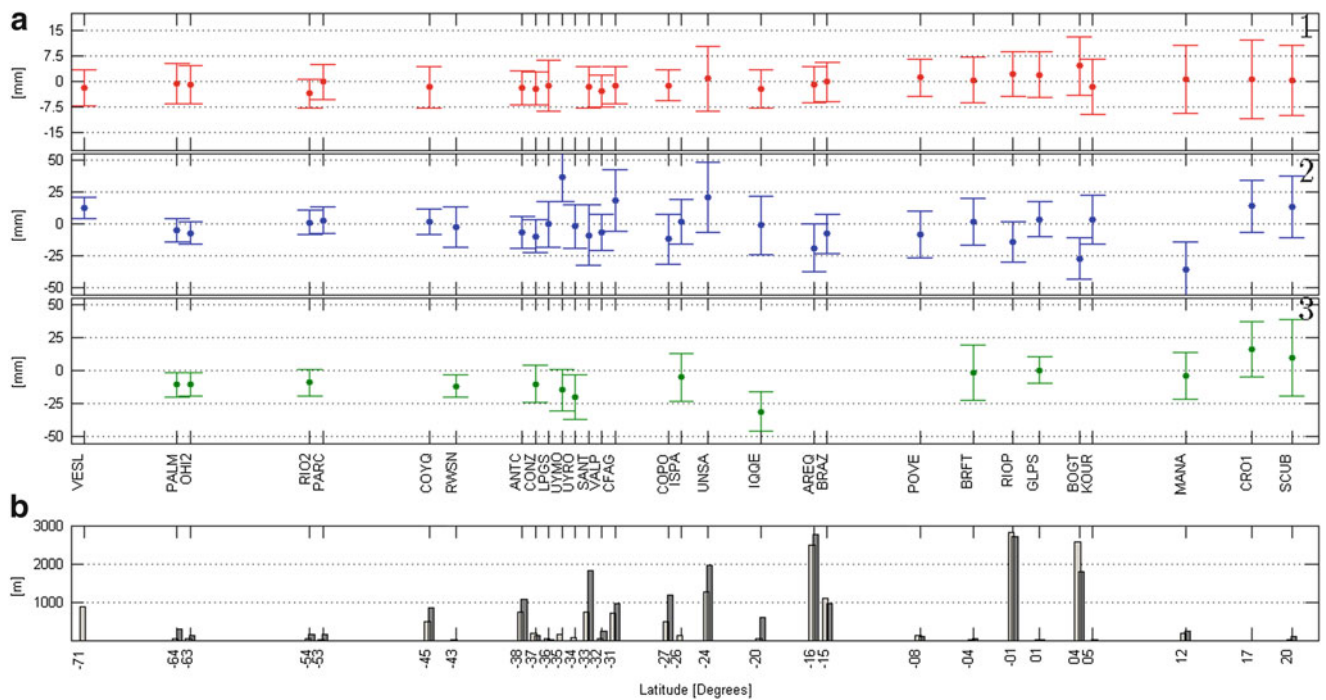


**Fig. 3** Mean and standard deviation of the differences between IGS and SIRGAS ZTD grouped every 2 h (**a**) and every month (**b**). Here, it has to be underlined that hours of the day refer to UTC time



that SIRGAS GNSS tropospheric delays agree with those obtained from the different techniques (i.e., with a bias in the difference varying from a minimum of 0.5 mm for the IGS values up to a maximum of 7.4 mm for the JMR measurements). As expected from using the same technique, the best agreement is found between SIRGAS and IGS ZTD. Nonetheless, the different strategies used to process the GPS observations (Sect. 2.1) influence the ZTD estimations: in particular, a refined analysis showed that  $ZTD_{SIR}$  are

characterized by a higher estimation noise, as it is shown in Fig. 2. Further, to investigate possible dependencies either on the epoch of the day or on the month of the year, the differences for the whole period have been grouped every 2 h and every month, respectively. Figure 3, which displays the results of the hourly and monthly comparison, does not highlight any degradation of the  $ZTD_{SIR}$  neither with the epoch of the day nor within the whole year.



**Fig. 4** (a) Bias and standard deviation over the 2 years of analysis on each SIRGAS site used in the 3 comparisons: (1)  $ZTD_{IGS} - ZTD_{SIR}$  in red; (2)  $ZWD_{ERA} - ZWD_{SIR}$  in blue; (3)  $ZWD_{JMR} - ZWD_{SIR}$  in green.

(b) The height in meters of each SIRGAS site: ellipsoidal height in light grey, orography height from ERA-Interim in dark grey

Table 1 shows that SIRGAS estimates agree to a high extent both with JMR measurements and with ERA-Interim weather data and the achieved results are fully consistent with Fernandes et al. (2013), Edwards et al. (2004), Bock et al. (2010). For each station, the results of the three comparisons are summarized in Fig. 4a; here, to investigate any latitudinal dependency of the results, the stations are sorted from north to south. Importantly, Fig. 4b displays the difference between the ellipsoidal height of the GNSS stations and the orography height of the ERA-Interim grid (i.e.,  $h$  and  $h_o$ ).

Although no clear latitudinal dependency in terms of bias is visible, Fig. 4a shows a slow decrease of standard deviation in the southern regions. The same situation is described by Teke et al. (2011) in their multi-technique comparison of ZTD and is most probably related to the lower content of WV in the colder regions as compared to the hotter regions, where the evaporation is dominating. Such effect is clearly visible in the inter-technique comparison 2 and 3. From the results of comparison 2 it is important to notice that large biases are obtained both for large and for little height differences (e.g., BOGT, UNSA and UYMO, MANA, respectively); therefore it appears difficult to infer a clear dependency between the results and the height differences.

## 4 Conclusions and Perspectives

In the period from June 2008 to June 2010, the ZTD of approximately 100 permanent stations belonging to SIRGAS network were estimated and then compared with the official ZTD distributed by the IGS. Then, the accuracy of the products was assessed in terms of consistency with respect to 2 independent techniques: (1) ZWD measured by Jason-1 satellite mission; (2) ZWD retrieved from observations data of the ECMWF ERA-Interim reanalysis model. The best agreement is reached between SIRGAS and IGS products ( $\mu = -0.5$  mm,  $\sigma = 6.5$  mm), whereas slightly worse statistical values are obtained in the comparisons with the JMR measurements ( $\mu = -7.4$  mm,  $\sigma = 15.4$  mm), and the ERA Interim data ( $\mu = -1.5$  mm,  $\sigma = 16.6$  mm).

A more detailed comparison undertaken with the IGS products confirmed that SIRGAS estimations quality is constant, independent from the local time or season of the year; at the same time, to mitigate the higher estimation noise in the final ZTD a further refinement of the processing parameters is required (e.g., using tighter constraints for parameters estimation). Overall, the achieved results are in accordance with previous researches. On one hand this testifies that



the inter-technique differences were correctly accounted for, on the other, it further confirms SIRGAS capabilities to contribute to short and long term meteorological studies.

Future investigations are oriented to evaluate the impact of including other GNSS (GLONASS, Compass, Galileo) constellation in ZTD estimation to derive reliable near real-time short weather forecast over the whole South and Central American region.

**Acknowledgements** Authors thank the three anonymous Reviewers and the Chief Editor for the valuable suggestions that thoroughly helped improving the present work. The authors recognize the fundamental role of the IGS for delivering GNSS data and products to the user community (Dow et al. 2009). ECMWF ERA-Interim data used in this study have been obtained from the ECMWF Data Server. This work was partially supported by *Progetto di cooperazione Scientifica e Tecnologica Italia-Argentina 2011–2013*.

## References

- Askne J, Nordius H (1987) Estimation of tropospheric delay for microwave from surface weather data. *Rad Sci* 22:379–386
- Berg H (1948) *Allgemeine meteorologie*. Duemmler, Bonn
- Bevis M, Businger S, Herring TA et al (1992) GPS meteorology: remote sensing of the atmospheric water vapor using the global positioning system. *J Geophys Res* 97:15787–15801
- Bock O, Willis P, Lacarra M, Bosser P. (2010) An inter-comparison of zenith tropospheric delays derived from DORIS and GPS data. *Adv Space Res* 46(10):1648–1660
- Böhm J, Niell A, Tregoning P, Schuh H (2006) Global mapping function (GMF): a new empirical mapping function based on numerical weather model data. *Geophys Res Lett* 33:L07304
- Brunini C, Sánchez L, Drewes H et al (2012) Improved analysis strategy and accessibility of the SIRGAS reference frame. *IAG Symp* 136:3–10
- Buyn S, Bar-Sever Y (2009) A new type of troposphere zenith path delay product of the International GNSS service. *J Geod* 83:367–373. doi:10.1007/s00190-008-0288-8
- Calori A, Colosimo G, Crespi M et al (2013) Zenith wet delay retrieval using two different techniques for the South American region and their comparison. *IAG Symp* 139:59–65. doi:10.1007/978-3-642-37222-3\_8
- Christensen EJ, Haines BJ, Keihm SJ et al (1994) Calibration of TOPEX/POSEIDON at platform harvest. *J Geophys Res* 99:24465–24485
- Dach R, Brockmann E, Schaer S et al (2009) GNSS processing at CODE: status report. *J Geod* 83:353–365. doi:10.1007/s00190-008-0281-2
- Davis JL, Herring TA, Shapiro II et al (1985) Geodesy by radio interferometry: effects of atmospheric modeling errors on estimates of baseline length. *Radio Sci* 20(6):1593–1607
- Dee DP, Uppala SM, Simmons AJ et al (2011) The ERA-Interim reanalysis: configuration and performance of the data assimilation system. *Q J R Meteorol Soc* 137:553–597. doi:10.1002/qj.828
- Dow JM, Neilan RE, Rizos C. (2009) The International GNSS Service in a changing landscape of Global Navigation Satellite Systems. *J Geod* 83:191–198. doi:10.1007/s00190-008-0300-3
- Edwards S, Moore P, King M (2004) Assessment of the Jason-1 and TOPEX/Poseidon microwave radiometer performance using GPS from offshore sites in the North sea. *Marine Geod* 27:717–727
- Fernandes MJ, Pires N, Lazáro C et al (2013) Tropospheric delays from GNSS for application in coastal altimetry. *Adv Space Res* 51:1352–1368. doi:<http://dx.doi.org/10.1016/j.asr.2012.04.025>
- Hogg D, Guiraud F, Decker M (1981) Measurement of excess radio-transmission length on earth-space paths. *Astron Astrophys* 95:304–307
- Keihm SJ, Janssen MA, Ruf CS et al (1995) TOPEX/Poseidon microwave radiometer (TMR). III. Wet troposphere range correction algorithm and pre-launch error budget. *IEEE Trans Geosci Remote Sensing* 33(1):147–161
- Kouba J (2008) Implementing and testing of the gridded Vienna Mapping Function 1 (VMF1). *J Geod* 82:193–205
- Kursinski E, Hajj G, Schofield J et al (1997). Observing the Earth's atmosphere with radio occultation measurements using the global positioning system. *J Geophys Res* 102:23429–23465
- Mackern MV, Mateo ML, Robin AM et al (2009) A terrestrial reference frame, coordinates and velocities for South American stations: contributions to Central Andes geodynamics. *Adv Geosci* 22:181–184
- Mendes VB, Prates G, Santos L et al (2000) An evaluation of the accuracy of models for the determination of the weighted mean temperature of the atmosphere. In: *Proceedings of ION*, pp 433–438
- Mitchell JFB (1989) The greenhouse effect and climate change. *Rev Geophys* 27:115–139
- Niell AE (1996) Global mapping functions for the atmosphere delay at radio wavelengths. *J Geophys Res* 101:3227–3246
- Picot N, Case K, Desai S et al (2003) AVISO and PODAAC User Handbook. IGDR and GDR Jason Products, SMM-MU-M5-OP-13184-CN (AVISO), JPL D-21352 (PODAAC)
- Rocken C, Van Hove T, Ware R (1997) Near real-time GPS sensing of atmospheric water vapor. *Geophys Res Lett* 24:3221–3224
- Ruf CS, Keihm SJ, Subramanya B et al (1994) TOPEX/POSEIDON microwave radiometer performance and in-flight calibration. *J Geophys Res* 99:24915–24926
- Saastamoinen J (1972) Contributions to the theory of atmospheric refraction. *Bull Gèod* 105(1):279–298
- Teke K, Böhm J, Nilsson T et al (2011) Multi-technique comparison of troposphere zenith delays and gradients during CONT08. *J Geod* 85:395–413
- Wentz J, Schabel M (2000) Precise climate monitoring using complementary satellite data sets. *Nature* 403:414–416. doi:10.1038/35000184

---

**Part V**

**Gravity Field Mapping Methodology from GRACE  
and Future Gravity Missions**

---

# The Role of Position Information for the Analysis of K-Band Data: Experiences from GRACE and GOCE for GRAIL Gravity Field Recovery

A. Jäggi, G. Beutler, U. Meyer, H. Bock, and L. Mervart

---

## Abstract

The Gravity Recovery And Interior Laboratory (GRAIL) mission orbiting the Moon and the Gravity Recovery And Climate Experiment (GRACE) mission orbiting the Earth share many conceptual commonalities. Major differences reside, however, in the absolute positioning of the spacecraft, which is accomplished by Doppler tracking from NASA's Deep Space Network (DSN) for GRAIL and by the Global Positioning System (GPS) for GRACE. Data from GRACE and from the Gravity and steady-state Ocean Circulation Explorer (GOCE) are used to investigate the role of position information. Artificially degrading either the geographical coverage or the accuracy of kinematic positions serving as input data together with continuously available K-Band inter-satellite data is shown not to be a limiting factor for gravity field recovery using the Celestial Mechanics Approach (CMA). Eventually, the CMA is applied to Level-1B data of the GRAIL mission to derive first Bernese lunar gravity field solutions.

---

## Keywords

Celestial mechanics approach • Farside of the Moon • GOCE • GRACE • GRAIL • Gravity field determination • Position information

---

## 1 Introduction

The Gravity Recovery And Interior Laboratory (GRAIL) mission, launched on September 10, 2011, has mapped the lunar gravity field with unprecedented accuracy and spatial resolution and will substantially contribute to a significantly improved understanding of the Moon's internal structure (Zuber et al. 2013b). The GRAIL mission concept was derived from the Gravity Recovery And Climate Experiment (GRACE) Earth mission (Tapley et al. 2004) and utilized a

modified GRACE inter-satellite link called the Lunar Gravity Ranging System (LGRS; Klipstein et al. 2013). Despite this important heritage, there are major differences between the GRAIL and the GRACE science payloads, such as the instrumentation used for absolute spacecraft positioning (Asmar et al. 2013). Whereas the GRACE spacecraft are equipped with Global Positioning System (GPS) receivers allowing the geolocation of the satellites with cm-accuracy at any time, e.g. Jäggi et al. (2009), the orbits of the GRAIL satellites are primarily determined by the Doppler tracking from NASA's Deep Space Network (DSN). As a consequence, the GRAIL orbits may only be constrained by Doppler measurements on the nearside of the Moon, inevitably resulting in a degradation of the reconstructed spacecraft trajectories over the Moon's farside. The trajectories reconstructed by the GRAIL science team, relying on both the DSN data and the (continuously available) inter-satellite Ka-Band data, are publicly available at NASA's Planetary Data System (PDS). They are, e.g., provided in the lunar-centred solar

---

A. Jäggi (✉) • G. Beutler • U. Meyer • H. Bock  
Astronomical Institute of the University of Bern, Sidlerstrasse 5,  
3012 Bern, Switzerland  
e-mail: [adrian.jaeggi@aiub.unibe.ch](mailto:adrian.jaeggi@aiub.unibe.ch)

L. Mervart  
Institute of Advanced Geodesy, Czech Technical University,  
Thakurova 7, 16629 Prague, Czech Republic

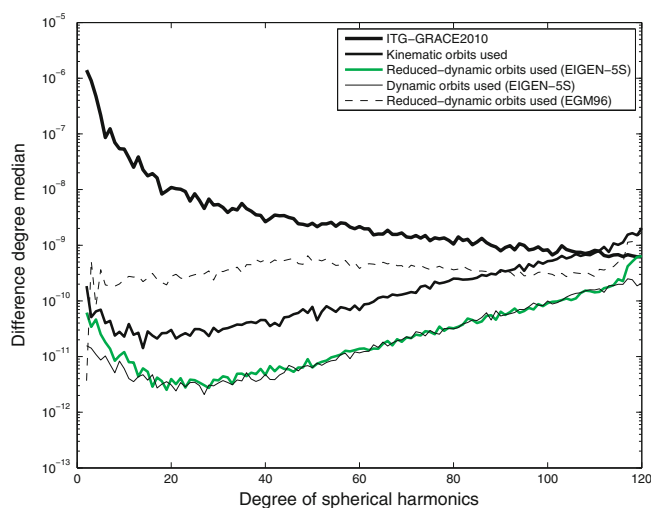
system barycentric frame (GNI1B products; Kahan 2012). Meanwhile, first GRAIL gravity fields have been released to the public (Konopliv et al. 2013; Lemoine et al. 2013).

GRACE gravity field recovery is performed at the Astronomical Institute of the University of Bern (AIUB) using the so-called Celestial Mechanics Approach (CMA; Beutler et al. 2010) for analyzing GPS, accelerometer, and inter-satellite K-band observations in a rigorously combined orbit and gravity field determination procedure. Gravity field recovery experiments based on data of GRACE and the Gravity and steady-state Ocean Circulation Explorer (GOCE) are performed here to investigate the role of position information for GRAIL gravity field determination. GOCE data are used in Sect. 2 to illustrate the consequences when using dynamic or reduced-dynamic positions as pseudo-observations for gravity field determination. In Sect. 3 the implications of using an artificially reduced geographical coverage (Sect. 3.1) and a reduced accuracy (Sect. 3.2) of position information are studied for combined orbit and gravity field recovery from kinematic GRACE positions and continuously available K-Band data by adopting the GRAIL observation scenario. Eventually, the CMA is applied to the GNI1B and KBR1B data of the GRAIL mission in Sect. 4 to derive first Bernese lunar gravity fields.

## 2 Role of Position Information for GOCE GPS-Only Gravity Field Recovery

Gerlach et al. (2003) already stated for CHAMP gravity field recovery that unbiased solutions may only be achieved when using kinematic orbit positions as pseudo-observations, but not when using reduced-dynamic orbit positions. Despite that finding reduced-dynamic orbit positions were occasionally used for gravity field recovery, e.g., together with gradiometer data, for the determination of the first GOCE gravity field model using the direct approach (Pail et al. 2011). Therefore, GOCE data are used in this section to illustrate possible consequences when using the position information from the GNI1B products for GRAIL position-only gravity field recovery, as it will also be performed in Sect. 4.

Figure 1 shows the difference degree median of GOCE gravity field solutions based on different types of orbit positions covering the period of Nov-Dec 2009 with respect to the GRACE gravity field model ITG-GRACE2010 (Mayer-Gürr et al. 2010). Apart from the independent gravity field solution based on the GOCE kinematic positions, which are derived in the frame of the GOCE High-level Processing Facility (HPF; Koop et al. 2006) as part of the GOCE Precise Science Orbit product (PSO; Bock et al. 2011), GOCE reduced-dynamic positions, being as well part of the PSO product, are also used for gravity field determination.

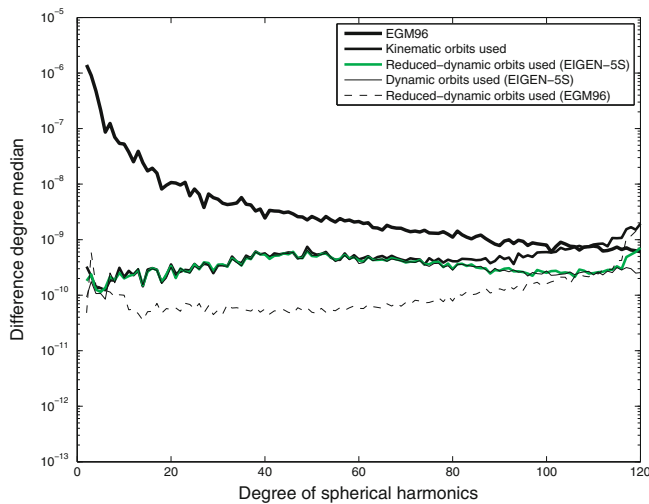


**Fig. 1** Difference degree median of GOCE GPS-only gravity field recoveries with respect to ITG-GRACE2010 when using different orbit input data

Figure 1 shows that the gravity field solution based on the reduced-dynamic positions agrees significantly better with the ITG-GRACE2010 model, because the reduced-dynamic GOCE PSO is based on the gravity field model EIGEN-5S (ESA 2010), i.e., on a gravity field model also derived from GRACE K-Band observations (Förste et al. 2008).

Figure 1 further illustrates that an even slightly better agreement with ITG-GRACE2010 may be achieved when using purely dynamic orbit positions based on EIGEN-5S, i.e., when only solving for initial conditions in the orbit determination process, but not for 6-min piecewise constant acceleration parameters as done for the reduced-dynamic orbit determination. Such a solution does, however, not provide independent gravity field information. This fact is also supported by the reduced-dynamic solution based on the gravity field model EGM96 (Lemoine et al. 1997) shown in Fig. 1, which yields a considerably inferior agreement with ITG-GRACE2010 due to the use of EGM96. Figure 2 shows the same gravity field solutions with respect to EGM96 and underlines the strong dependency on the gravity field model used for orbit determination – the solution agreeing least with ITG-GRACE2010 shows the best agreement with EGM96.

GRAIL gravity field recovery using GNI1B products as pseudo-observations instead of original DSN tracking data is thus expected to be biased towards the gravity field model underlying the GNI1B products. The situation is only mitigated insofar as the inter-satellite Ka-Band data are heavily dominating the solution for most degrees, as it is known from the experience with GRACE (Beutler et al. 2010). Nevertheless the recovery of the very long wavelength part of the gravity field has to be interpreted with care when using reduced-dynamic positions as pseudo-observations.



**Fig. 2** Difference degree median of GOCE GPS-only gravity field recoveries with respect to EGM96 when using different orbit input data

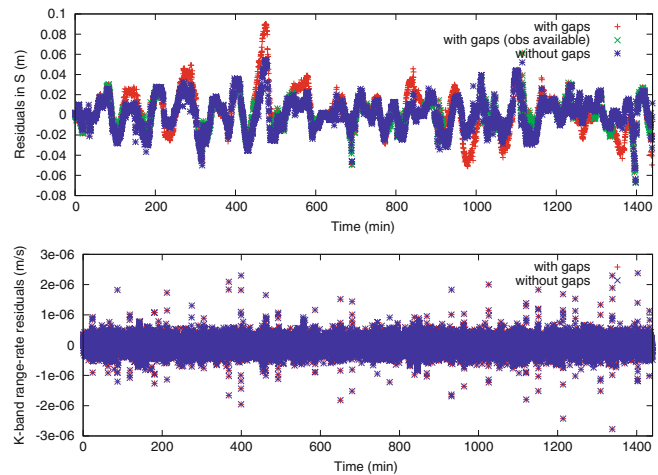
### 3 Role of Position Information for GRACE K-Band Gravity Field Recovery

GRACE gravity field determination using the CMA is based on the analysis of inter-satellite K-band measurements and GPS-derived kinematic positions (Jäggi et al. 2009). Using a priori orbits derived from the kinematic positions of both GRACE satellites and from the inter-satellite measurements, normal equations (NEQs) for both types of observations are set up on a daily basis for the unknown gravity field coefficients and for additional arc-specific parameters. The resulting technique-specific daily NEQs are then combined into one system for each daily arc. Eventually, arc-specific parameters are pre-eliminated and the combined daily NEQs are accumulated into monthly, annual or multi-annual systems. Details about the general procedure may be found in Jäggi et al. (2010) or Meyer et al. (2012).

#### 3.1 Impact of Position Coverage

In order to study the impact of missing position information over the “farside” of a planet, kinematic GRACE positions covering the period of January 2008 were artificially reduced by discarding all positions falling within the range of  $-90^\circ \leq \lambda \leq 90^\circ$  geographical longitude. Orbit and gravity field determination from the reduced/full set of kinematic positions and the K-Band range-rate data (assumed to be continuously available also over the farside) are performed within three experiments to simulate the GRAIL situation.

In a first experiment the static part of the AIUB-GRACE03S (Jäggi et al. 2011) gravity field model up to



**Fig. 3** Position (*top*) and K-Band (*bottom*) residuals of the a priori orbits of experiment 1 (see text)

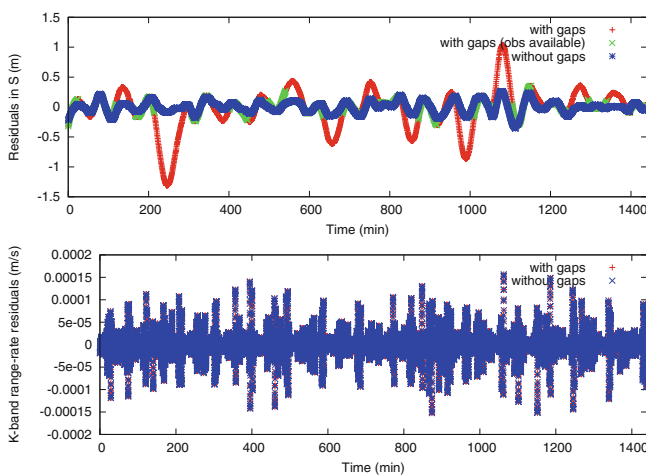
degree 160 was used to generate a priori orbits based on the reduced/full set of kinematic positions and K-Band data. Non-gravitational accelerations were modeled by taking into account the GRACE accelerometer data. Figure 3 (top) shows the orbit residuals for one particular day of the generated a priori orbits in the along-track direction for all kinematic positions, irrespective whether they were used for the orbit computation (without gaps) or not (with gaps). The kinematic positions actually used from the artificially reduced data set are marked by green points.

Figure 3 (top) shows that the absolute orbit quality is only marginally degraded when position information is lacking due to a farside effect, because high quality background models (GRACE gravity field model, accelerometer data) are used for the experiment.

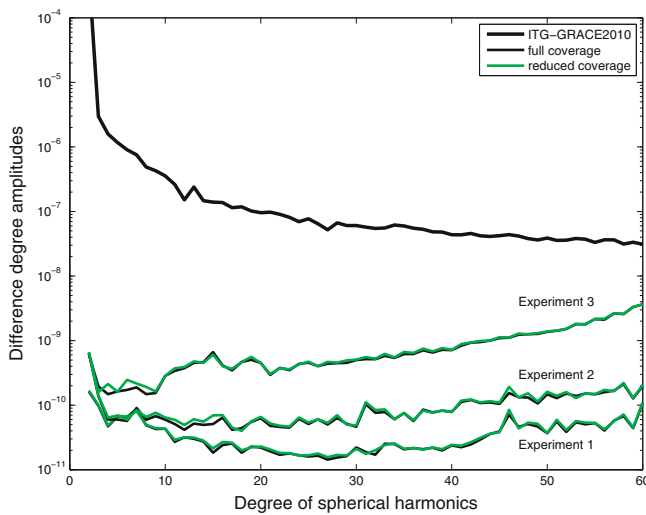
Figure 3 (bottom) shows that no differences can be seen in the K-Band residuals. When further degrading the absolute quality of the a priori orbits, e.g., by not taking accelerometer data into account (experiment 2), similar results are obtained without almost any visible impact on the K-Band residuals caused by a farside effect (not shown).

In a third experiment the EGM96 gravity field model up to degree 160 was used to generate a priori orbits based on the reduced/full set of kinematic positions and K-Band data without taking accelerometer data into account. Figure 4 (top) shows the orbit residuals of the generated a priori orbits for the same day in the along-track direction. A severe degradation of the absolute orbits with deviations up to the meter-level is observed when position information is lacking, in the K-Band residuals again no differences are seen irrespective of using the reduced or the full set of position information (see Fig. 4, bottom). A significant degradation of the general RMS-level with respect to Fig. 3 (bottom) is obvious due to the use of the EGM96 gravity field model.

GRACE monthly gravity field solutions up to degree and order 60 were generated for all experiments with the reduced/full set of position information. Figure 5 shows the difference degree amplitudes of all solutions and confirms the expectations from Figs. 3 and 4 that only marginal degradations are induced by the reduced coverage. Irrespective of the quality of the background modeling used for the three experiments, all monthly gravity field solutions are dominated by the ultra-precise K-Band data and only suffer marginally from a reduced coverage with position information. The same behavior is observed when not using a static gravity field model to reduce omission errors from degrees 61–160 (not shown), which is a more realistic scenario for GRAIL.



**Fig. 4** Position (*top*) and K-Band (*bottom*) residuals of the a priori orbits of experiment 3 (see text)



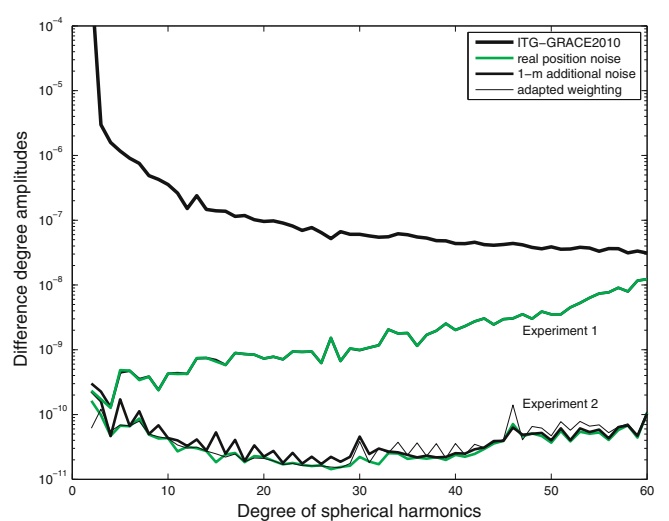
**Fig. 5** Difference degree amplitudes of GRACE monthly fields with respect to ITG-GRACE2010 when using different position information with different geographical coverage

### 3.2 Impact of Position Noise

In order to provoke an impact of degraded position information on the gravity fields of a planet, kinematic GRACE positions covering the period of January 2008 were artificially degraded in this section by adding an additional white noise of 1 m RMS. Orbit and gravity field determination from the degraded/original set of kinematic positions and the (original) K-Band range-rate data are performed within two experiments.

In view of GRAIL, where a high quality static gravity field model could not yet be used at the time of writing to reduce omission errors, the static part of the AIUB-GRACE03S gravity field model was used in a first experiment only up to degree 60 to generate a priori orbits based on the degraded/original set of kinematic positions and K-Band data, and monthly gravity field solutions up to degree 60. In a second experiment the AIUB-GRACE03S gravity field model was used in analogy to Sect. 3.1 up to degree 160 to avoid omission errors when determining the monthly field up to the same maximum degree of 60.

Figure 6 shows that, apart from degree 2, the difference degree amplitudes of the first experiment are almost identical when using the degraded/original set of kinematic positions. Only for the second experiment, where omission errors are not dominating the solution, a negative impact of the degraded positions on the recovery of the monthly solution can be seen. Figure 6 also shows, however, that the impact of the degraded positions can be “cured” for the low degrees to a large extent by adapting the weighting between the NEQ contribution of the positions and the NEQ contribution of the K-Band data according to the additionally imposed noise



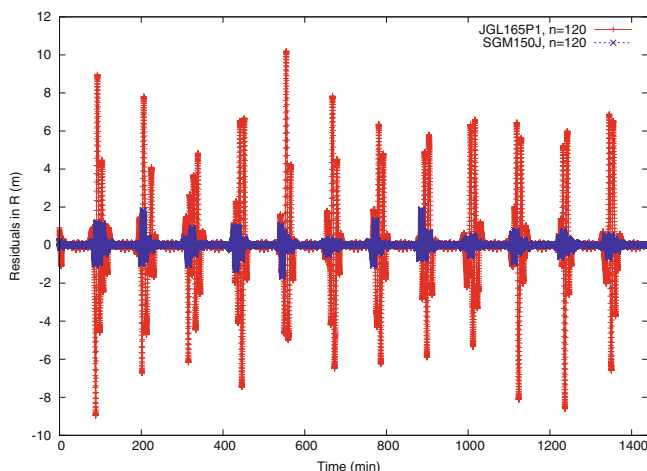
**Fig. 6** Difference degree amplitudes of GRACE monthly fields with respect to ITG-GRACE2010 when using different accuracies of the position information

level. Inferior accuracy of GRAIL position information (an accuracy of 20 cm is claimed (Konopliv et al. 2013)) might thus be indeed an issue for lunar gravity field determination based on Ka-Band data, but it should not be a limiting factor.

#### 4 First Results from GRAIL Gravity Field Recovery

GRAIL gravity field determination using the CMA based on the methodology described in Sect. 3 is now applied to the GRAIL inter-satellite Ka-band measurements and the GRAIL orbital positions stemming from the GNI1B products. The pre-GRAIL lunar gravity field models based on Lunar Prospector data (JGL165P1, Konopliv et al. (2001)) and SELENE data (SGM150J, Goossens et al. (2011)), respectively, serve as a priori gravity field models to generate a priori orbits based on the GNI1B positions and the Ka-Band data, and to estimate lunar gravity field models up to degree 120 using data collected during the GRAIL primary mission phase covering the period of Mar–May 2012. During that period the GRAIL satellites were orbiting the Moon on polar orbits at mean altitudes of about 55 km above the lunar surface. Complete coverage of the lunar surface was achieved within 27.3 days (the Moon’s rotation period), resulting in a total of three mapping cycles during the GRAIL primary mission phase (Zuber et al. 2013a).

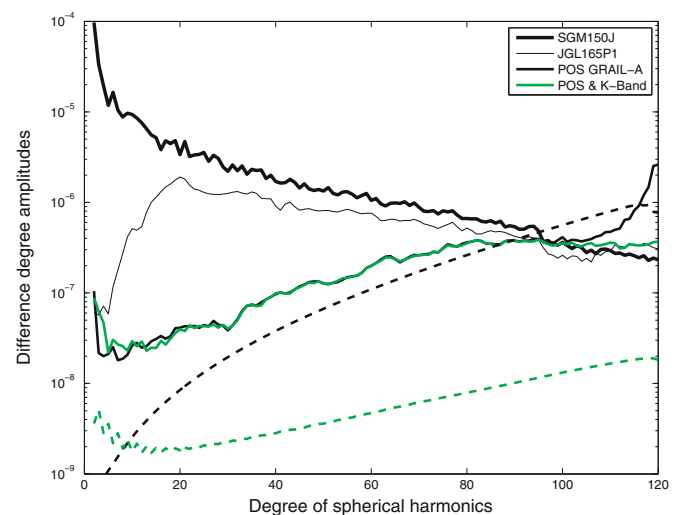
Figure 7 shows for one particular day the orbit residuals of the generated a priori orbits in the radial direction for the GNI1B positions when using the a priori gravity field model JGL165P1 and SGM150J up to degree 120, respectively. Empirical pulses were estimated every 20 min to compensate for imperfect modelings of non-gravitational accelerations,



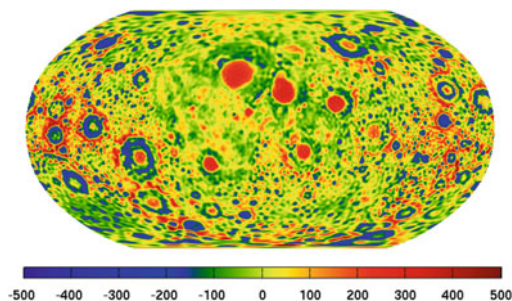
**Fig. 7** Position residuals of the a priori GRAIL orbits when using different gravity field models

e.g., caused by solar radiation pressure. Figure 7 shows that orbital fits at the level of few centimeters may be achieved for both models on the nearside of the Moon, but larger residuals are occurring on the farside. Deviations up to 10 m occur for JGL165P1, which is a model uniquely based on data covering the nearside of the Moon. SGM150J is, however, based on data partly covering also the farside of the Moon by 4-way Doppler tracking (Namiki et al. 2009), which significantly reduces the maximum residuals on the farside to the one meter-level.

Figure 8 shows the difference degree amplitudes with respect to the SELENE gravity field model SGM150J for recoveries from positions of the GRAIL-A satellite alone and for recoveries from positions of both GRAIL satellites and the inter-satellite Ka-Band range-rate data. In analogy to GRACE (Jäggi et al. 2012) essentially the same results are obtained for the recoveries irrespective of whether SGM150J or JGL165P1 are used as a priori gravity field models (not shown), which largely differ in their quality as shown in Fig. 8. Figure 8 also shows that the SELENE field only validates the quality of the recovered fields up to about degree 35. Due to the limited quality of the SELENE gravity field model (in comparison to the GRAIL data), no differences are seen for higher degrees when taking the Ka-Band data for gravity field determination into account. The error degree amplitudes reveal, however, that the enormous potential of the GRAIL Ka-Band data is by no means exploited by a gravity field solution up to degree 120. Even at this early stage recoveries up to degree 180 or higher seem to be possible with the CMA, despite the relatively large number of shortcomings in our current background modeling. The resulting free-air gravity anomalies are shown in Fig. 9,



**Fig. 8** Difference degree (*solid lines*) and error degree (*dotted lines*) amplitudes of GRAIL position-only and combined gravity field recoveries



**Fig. 9** Lunar gravity anomalies (in mgal) up to degree 120 of the combined GRAIL solution using the CMA

already revealing the most pronounced signatures of the lunar gravity field.

## 5 Summary

We studied possible consequences of using reduced-dynamic positions as pseudo-observations for Earth-based gravity field determination using GOCE data. The resulting gravity field solutions are heavily biased towards the gravity field model underlying the reduced-dynamic orbit determination, which renders the use of the GNI1B products questionable from a methodological point of view for independent GRAIL gravity field recovery. The problem is mitigated only insofar as combined GRAIL solutions are dominated by the Ka-Band data.

We used GRACE data to investigate the implications of a reduced coverage of a planetary surface with position information. Even for a GRAIL-like situation with inferior a priori gravity field information and lacking accelerometer data, almost no degradation was found due to the farside effect. Only a severely degraded accuracy of the position information led to a small degradation of the monthly solutions.

We used the CMA for lunar gravity field determination with the GRAIL positions from the GNI1B products and the Level-1B Ka-band range-rate data from the primary mission phase to generate first Bernese lunar gravity field models up to degree 120. Even our first GRAIL gravity field solutions reveal the enormous potential of the Ka-Band data. Further refinements of the CMA will primarily focus on the modeling of non-gravitational accelerations and the implementation of the DSN analysis capability.

## References

- Asmar SW, Konopliv AS, Watkins MM, Williams JG, Park RS, Kruizinga G, Paik M, Yuan DN, Fahnstock E, Strelakov D, Harvey N, Lu W, Kahan D, Oudrhiri K, Smith DE, Zuber MT (2013) The scientific measurement system of the gravity recovery and interior laboratory mission. *Space Sci Rev* 178(1):25–55
- Beutler G, Jäggi A, Mervart L, Meyer U (2010) The celestial mechanics approach: application to data of the GRACE mission. *J Geod* 84:661–681
- Bock H, Jäggi A, Meyer U, Visser P, van den IJssel J, van Helleputte T, Heinze M, Hugentobler U (2011) GPS-derived orbits for the GOCE satellite. *J Geod* 85:807–818
- ESA (2010) GOCE Standards. [http://earth.esa.int/pub/ESA\\_DOC/GOCE/GOCE\\_Standards\\_3.1.pdf](http://earth.esa.int/pub/ESA_DOC/GOCE/GOCE_Standards_3.1.pdf)
- Förste C, Flechtner F, Schmidt R, Stubenvoll R, Rothacher M, Kusche J, Neumayer KH, Biancale R, Lemoine JM, Barthelmes F, Bruinsma S, König R, Meyer U (2008) EIGEN-GL05C – A new global combined high-resolution GRACE-based gravity field model of the GFZ-GRGS cooperation. *Geophys Res Abstr* 10, EGU2008-A-03426
- Gerlach C, Földváry L, Švehla D, Gruber T, Wermuth M, Sneeuw N, Frommknecht B, Oberndorfer H, Peters T, Rothacher M, Rummel R, Steigenberger P (2003) A CHAMP-only gravity field model from kinematic orbits using the energy integral. *Geophys Res Lett* 30(20):2037
- Goossens S, Matsumoto K, Ishihara Y (2011) Improved high-resolution lunar gravity field model from SELENE and historical tracking data. AGU Fall Meeting, Abstract P44B-05
- Jäggi A, Dach R, Montenbruck O, Hugentobler U, Bock H, Beutler G (2009) Phase center modeling for LEO GPS receiver antennas and its impact on precise orbit determination. *J Geod* 83:1145–1162
- Jäggi A, Beutler G, Mervart L (2010) GRACE gravity field determination using the celestial mechanics approach – first results. In: Mertikas S (ed) Gravity, geoid and earth observation. Springer, Berlin, Heidelberg, pp 177–184
- Jäggi A, Meyer U, Beutler G, Prange L, Dach R, Mervart L (2011) AIUB-GRACE03S. <http://icgem.gfz-potsdam.de/ICGEM/>
- Jäggi A, Beutler G, Meyer U, Prange L, Dach R, Mervart L (2012) AIUB-GRACE02S – status of GRACE gravity field recovery using the celestial mechanics approach. In: Kenyon S, Pacino MC, Marti U (eds) Geodesy for planet earth. Springer, Heidelberg, pp 161–170
- Kahan D (2012) GRAIL data product software interface specification. JPL D-76383, Version 1.1, California Institute of Technology
- Klipstein WM, Arnold BW, Enzer DG, Ruiz AA, Tien JY, Wang RT, Dunn CE (2013) The lunar gravity ranging system for the gravity recovery and interior laboratory (GRAIL) mission. *Space Sci Rev* 178(1):57–76
- Konopliv AS, Asmar SW, Carranza E, Sjogren WL, Yuan DN (2001) Recent gravity models as a result of the lunar prospector mission. *Icarus* 150:1–18
- Konopliv AS, Park RS, Yuan DH, Asmar SW, Watkins MM, Williams JG, Fahnstock E, Kruizinga G, Paik M, Strelakov D, Harvey N, Smith DE, Zuber MT (2013) The JPL lunar gravity field to spherical harmonic degree 660 from the GRAIL Primary Mission. *J Geophys Res Planets* 118:1415–1434
- Koop R, Gruber T, Rummel R (2006) The status of the GOCE high-level processing facility. In: 3rd GOCE User Workshop, 6–8 November 2006. Frascati, Italy, pp 195–205, ESA SP-627
- Lemoine FG, Smith DE, Kunz L, Smith R, Pavlis EC, Pavlis NK, Klosko SM, Chinn DS, Torrence MH, Williamson RG, Cox CM, Rachlin KE, Wang YM, Kenyon SC, Salman R, Trimmer R, Rapp RH, Nerem RS (1997) The development of the NASA GSFC and NIMA Joint Geopotential Model. In: Segawa J, Fujimoto H, Okubo S (eds) Gravity, geoid and marine geodesy. Springer, Heidelberg, pp 461–469
- Lemoine FG, Goossens S, Sabaka TJ, Nicholas JB, Mazarico E, Rowlands DD, Loomis BD, Chinn DS, Caprette DS, Neumann GA, Smith DE, Zuber MT (2013) High-degree gravity models from GRAIL primary mission data. *J Geophys Res Planets* 118:1–23
- Mayer-Gürr, T, Eicker A, Kurtenbach E, Ilk K-H (2010) ITG-GRACE: global static and temporal gravity field models from GRACE data. In: Flechtner F, Gruber T, Güntner A, Manda M, Rothacher M,



- Schöne T, Wickert J (eds) System earth via geodetic-geophysical space techniques. Springer, Heidelberg, pp 159–168
- Meyer U, Jäggi A, Beutler G (2012) Monthly gravity field solutions based on GRACE observations generated with the celestial mechanics approach. *Earth Planet Sci Lett* 345–348:72–80
- Namiki N, Iwata T, Matsumoto K, Hanada H, Noda H, Goossens S, Ogawa M, Kawano N, Asari K, Tsuruta S, Ishihara Y, Liu Q, Kikuchi F, Ishikawa T, Sasaki S, Aoshima C, Kurosawa K, Sugita S, Takano T (2009) Farside gravity field of the moon from four-way doppler measurements of SELENE (Kaguya). *Science* 323(5916):900–905
- Pail R, Bruinsma S, Migliaccio F, Förste C, Goiginger H, Schug W-D, Höck E, Reguzzoni M, Brockmann JM, Abrikosov O, Veicherts M, Fecher T, Mayrhofer R, Krasbutter I, Sansó F, Tscherning CC (2011) First GOCE gravity field models derived by three different approaches. *J Geod* 85:819–843
- Tapley BD, Bettadpur S, Ries JC, Thompson PF, Watkins M (2004) GRACE measurements of mass variability in the earth system. *Science* 305(5683):503–505
- Zuber MT, Smith DE, Lehman DH, Hoffman TL, Asmar SW, Watkins MM (2013a) Gravity recovery and interior laboratory (GRAIL): mapping the lunar interior from crust to core. *Space Sci Rev* 178(1):3–24
- Zuber MT, Smith DE, Watkins MM, Asmar SW, Konopliv AS, Lemoine FG, Melosh HJ, Neumann GA, Phillips RJ, Solomon SC, Wieczorek MA, Williams JG, Goossens SJ, Kruizinga G, Mazarico E, Park RS, Yuan DN (2013b) Gravity field of the moon from the gravity recovery and interior laboratory (GRAIL) mission. *Science* 339(6120):668–671

---

# Gravity Field Mapping from GRACE: Different Approaches—Same Results?

Christoph Dahle, Christian Gruber, Elisa Fagiolini, and Frank Flechtner

---

## Abstract

GFZ as part of the GRACE Science Data System (SDS) is routinely processing time-variable global gravity field models on monthly and weekly basis throughout the whole GRACE mission period. These operational products consist of spherical harmonic coefficients which are calculated based on the so-called dynamic method, i.e. integration of variational equations. As a matter of fact, these coefficients are imperfect due to different error sources such as inaccurate background models, instrument noise and inhomogeneous sampling and thus have to be filtered during post-processing in an appropriate way. Nevertheless, the current release named GFZ RL05 shows significant improvements compared to its precursors with an average error level of only about a factor of 6 above the pre-launch estimated baseline accuracy.

Additionally, an alternative approach using radial basis functions is developed at GFZ. This approach is based on the inversion of integral equations using gradient differences as in-situ observations. The resulting gravity field products can be directly derived as gridded data making this approach also suitable for regional applications. No post-filtering is necessary, as regularization is already applied during system inversion. Additionally applying a Kalman filter, higher temporal resolution can be achieved.

This paper gives a brief overview of the methodology of both approaches and their particular strengths and weaknesses are discussed. Results from GFZ RL05 and the latest results of the radial basis function approach are compared and also validated against independent data sources.

---

## Keywords

Dynamic method • GRACE • Kalman filter • Radial basis functions • Time-variable gravity field

---

## 1 Introduction

The main objective of the GRACE mission (Tapley et al. 2004) consists of monitoring the temporal variations of the Earth's gravity field. During the past decade, an increasing

number of different GRACE releases generated by different groups have become available. The temporal resolution of these releases varies from monthly over 10-day and weekly to even daily gravity field solutions, whereas their spatial resolution naturally increases with lower temporal resolution or by applying any type of regularization but is generally limited to a few 100 km. Moreover, different approaches of gravity field recovery are applied by the processing centers. An overview of the most important approaches is given in Table 1. Although differences between different solutions have become smaller with every new release, they are still

---

C. Dahle (✉) • C. Gruber • E. Fagiolini • F. Flechtner  
GFZ German Research Centre for Geosciences, Telegrafenberg,  
14473 Potsdam, Germany  
e-mail: [dahle@gfz-potsdam.de](mailto:dahle@gfz-potsdam.de)

**Table 1** Overview of available global time-variable GRACE gravity field models

Approach	Processing center	Temporal resolution	Reference
Dynamic method	GFZ	Monthly & weekly	Dahle et al. (2012)
	CSR JPL	Monthly Monthly	Bettadpur (2012) Watkins and Yuan (2012)
	GRGS	10-day	Bruinsma et al. (2010)
Short-arc method	IGG/Bonn	Monthly & daily	Mayer-Gürr (2006), Kurtenbach et al. (2009)
Acceleration approach	TU Delft	Monthly	Liu et al. (2010)
Celestial mechanics approach	AIUB/Bern	Monthly	Meyer et al. (2012)

present. However, they cannot be simply explained by the different approaches alone, as background models and standards are also not consistent.

In this paper, two different approaches based on the same background models are compared. The first is the dynamic method, which has been used at GFZ for the operational GRACE processing since many years (Sect. 2). The second is an alternative approach based on radial basis functions (Sect. 3) which has recently been developed at GFZ (Gruber et al. 2014). Using the same background models, one can expect very similar results in case that both approaches overcome known shortcomings concerning the availability of satellite gravity data in view of the subtle temporal gravity field variations. This can either be achieved by a posteriori destriping and smoothing or by introducing spatial and temporal constraints during the solution process beforehand. On the other hand, possible differences can be related to the approaches and thus indicate their potential strengths and weaknesses. This is investigated in Sect. 4, where the results of both methods are compared and validated against independent data.

## 2 Dynamic Method

At GFZ, both monthly and weekly global GRACE gravity field models are operationally processed within the GRACE Science Data System (SDS) using the dynamic method (see Table 1). Its application to GRACE data is described e.g. in Reigber et al. (2005) or Schmidt (2007). Briefly summarized, this approach is based on numerical integration

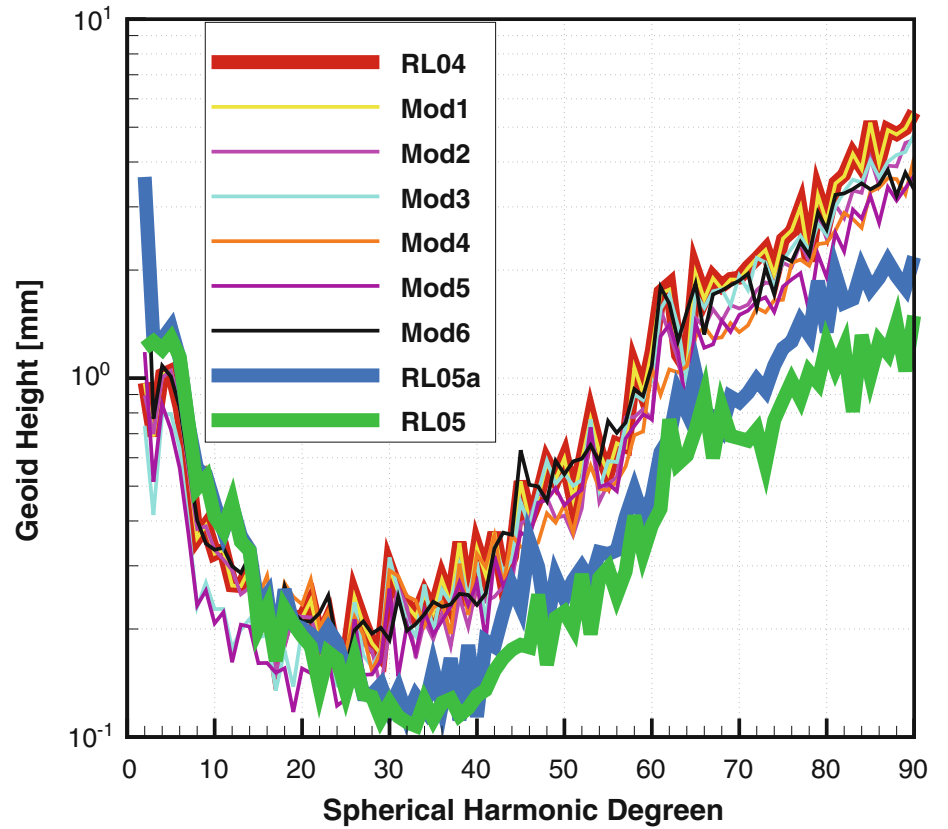
of the satellites' equations of motion and the corresponding variational equations. Then, by setting up the linearized observation equations, the unknowns, i.e. gravity field, orbit and instrument specific parameters, can be solved in a least squares adjustment. The resulting gravity field solutions are obtained in form of spherical harmonic coefficients which can be transferred to gridded values of the desired gravity field functional such as geoid undulations or gravity anomalies. However, an appropriate smoothing of the solutions has to be applied, as the derived grids are degraded by meridional striping and artefacts mainly caused by the anisotropic observation geometry of GRACE.

The current release of monthly and weekly gravity field models labelled as GFZ RL05 shows clear improvements compared to its precursor GFZ RL04, in particular when looking at the noise level which has been significantly reduced. The average RL05 error level has dropped to only about a factor of 6 above the pre-launch estimated baseline accuracy (Kim 2000), whereas this factor has been around 15 for RL04. The improvements are the results of several modifications w.r.t. RL04 described in Dahle et al. (2014) and can be summarized in the following six groups of changes, denoted as Mod1, Mod2, . . . , Mod6, resp.:

- Update of GFZ's Earth Parameter and Orbit System software (EPOS-OC), mainly due to implementation of IERS2010 standards (Mod1)
- Use of reprocessed GRACE Level-1B data (L1B RL02, provided by JPL) (Mod2)
- Improved GPS processing (Mod3)
- Updated background models (Mod4)
- Modified relative weighting of GPS and K-band range rate observations (Mod5)
- Modified parameterization of accelerometer biases (Mod6)

The individual contributions of these modifications are quantified in Fig. 1. In addition, further notable noise reduction has been achieved by keeping all orbit parameters fixed in the final run and solving only the gravity field parameters. This strategy requires an already good, i.e. sufficiently close to the geophysical truth, a priori information for the background gravity field model. For RL05, this is the case as EIGEN-6C (Shako et al. 2014) including its time-variable part is used. However, fixing the orbit parameters also regularizes the solutions towards the a priori model (Meyer et al. 2015). This effect has become visible especially in the most recent years (2009 and later), where the linearly modelled trend of EIGEN-6C is deviating from the geophysical truth in certain areas. As a consequence, GFZ has decided to provide an alternative RL05a time-series, where orbit and gravity field parameters are estimated together.

**Fig. 1** Degree variances for the April 2008 monthly solution in terms of geoid height [mm]; Mod1, Mod2, . . . , Mod6 represent the individual contribution of the corresponding modification with GFZ RL04 as reference, RL05a is the sum of all six modifications, and RL05 is the latter with fixed orbit parameters



### 3 Radial Basis Function Method

The basic idea behind the Radial Basis Function (RBF) method is the transformation of satellite instrument data to in-situ observables and subsequent inversion of gravity functionals, i.e. their corresponding integral equations, defined by reproducing kernels. Novák (2007) has already introduced a kernel function for the GRACE-type observation low-low satellite-to-satellite tracking (SST) based on inter-satellite range-accelerations (Rummel 1975). The application of the RBF method has been described in Gruber et al. (2014); in the following, only the most essential formulas are briefly introduced.

First, inter-satellite ranges  $\rho$ , range-rates  $\dot{\rho}$  and range-accelerations  $\ddot{\rho}$  taken from GRACE KBR1B-products as well as inter-satellite velocity differences  $|\delta\dot{\mathbf{r}}|$  are transformed to in-situ observables representing gradient differences of the gravity potential  $\phi$  between the two GRACE spacecrafts:

$$f(\rho, \dot{\rho}, \ddot{\rho}, |\delta\dot{\mathbf{r}}|) = \langle \nabla\phi(\mathbf{r}_A) - \nabla\phi(\mathbf{r}_B), \mathbf{e}_{\text{LOS}} \rangle \quad (1)$$

with  $\mathbf{r}_A$ ,  $\mathbf{r}_B$  the geocentric position vectors of GRACE-A and -B, resp., and  $\mathbf{e}_{\text{LOS}}$  the unit vector in line-of-sight (LOS).

The integral equation

$$f(P) = \iint_S \phi(Q) \langle \delta\nabla K(P, Q), \mathbf{e}_{\text{LOS}} \rangle dS \quad (2)$$

describes the relation between the potential  $\phi$  at grid points  $Q$  located at the Earth's surface  $S$ , approximated by the bounding sphere  $R$  (assuming that all relevant gravitational masses are embedded inside this surface) and the gravity functional  $f$  at the evaluation points  $P$  with the corresponding radius  $r_P$ . The chosen grid is equiareal to reduce the number of grid points by about 30% compared to an equiangular version. The reproducing kernel  $K(P, Q)$  reads

$$K(P, Q) = \sum_{n=2}^{\infty} (2n+1) \left(\frac{R}{r_P}\right)^{n+1} P_n(\cos\psi) \quad (3)$$

with the Legendre polynomial  $P_n(\cos\psi)$ , depending on degree  $n$  and the spherical distance  $\psi$  between  $P$  and  $Q$ .

As the potential values at the grid points  $Q$  are the unknowns to be solved for, Eq. (2) has to be inverted. This is achieved by least squares adjustment. Because the normal equation system is ill-posed, regularization is required.

This can either be achieved by a Tikhonov regularization or by adopting a Kalman filter system evolving the a priori values of the normal equation system. Whereas the former is a suitable choice for monthly solutions, using a Kalman filter enables a higher temporal resolution.

For this work, daily global Kalman filtered RBF solutions have been generated first. The feasibility of solutions with daily resolution has been shown by Kurtenbach et al. (2009). For the prediction step, covariance information for the relevant sources of mass change is required. The spatial correlations between the grid points are directly present in the kernel matrix. Estimates for the temporal correlations are derived from a hydrology model (WGHM, Döll et al. 2003), from the atmosphere and ocean de-aliasing model (AOD1B, Flechtner et al. 2014) and from other available GRACE time-series (here: GFZ RL05). Dynamic GRACE orbits, needed in Eq. (1), and corresponding background forces to reduce the in-situ observations are also taken from GFZ RL05 processing. No time-variable gravity field background model has been reduced, as the Kalman filter stochastically predicts the a priori state using the most recent estimate of the state vector.

In a second step, monthly global RBF solutions have been inverted. The daily RBF solutions are removed as additional de-aliasing product and the variances obtained from the Kalman filter serve as input for a Tikhonov regularization. The appropriate signal amplitudes are found in an optimal sense by determining a global regularization parameter using an empirical L-curve criterion to minimize the norm of the residuals and unknowns. It is worth mentioning that the generation of monthly RBF solutions is generally independent, i.e. it is not necessary to generate daily RBF solutions first, as the information for the Tikhonov regularization can also be taken from other sources, e.g. geophysical models or GRACE.

In general, the RBF method turns out to be stable even in case of poor ground track and sample coverage. Moreover, several observation types (gravity, geometry) from multiple, both space-born and terrestrial, sensor systems can be combined and commonly integrated. Global as well as regional gravity field solutions can be computed.

## 4 Comparison and Validation

For comparison and validation of results from GFZ RL05 and RBF, corresponding time-series spanning ten years (2003–2012) are analyzed. The monthly and weekly RL05 solutions have been smoothed with the DDK2 decorrelation filter (Kusche et al. 2009). The DDK2 filter has been chosen as

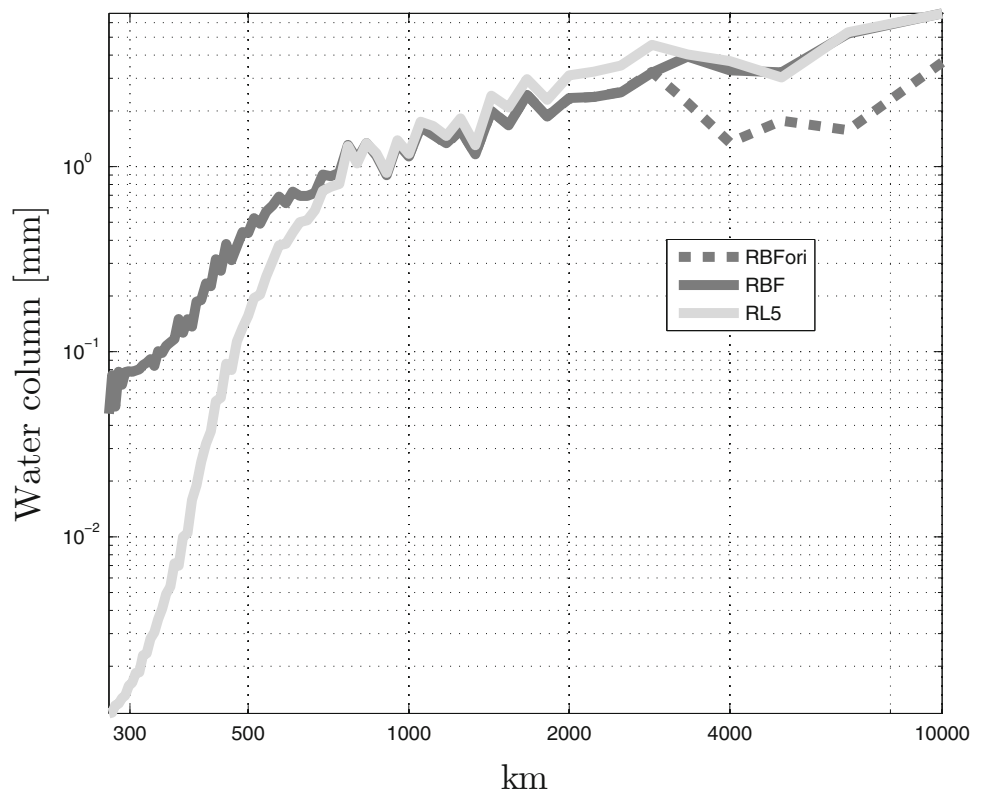
a good compromise between signal preservation and noise suppression. Its corresponding filter radius of approx. 340 km matches well the spatial resolution where the level of mm-geoid accuracy is reached for RL05 (Dahle et al. 2014). As mentioned in Sect. 3, the RBF solutions are already regularized and therefore no additional smoothing has to be applied.

First, the average monthly degree variances of both approaches are compared (Fig. 2). In order to do so, the RBF solutions, obtained as grids, have been converted into the spectral domain by spherical harmonic analysis. It becomes obvious that the RBF solutions have less power in the long wavelengths. A possible reason might be that in the RBF processing only GRACE K-band observations are used whereas GPS observations, which are essential for solving for the low degrees in the dynamic method, are omitted so far. Expanding the RBF method by GPS observations, i.e. adding orbital information in form of the equation of motion, could help in future to overcome this issue. To avoid systematic effects caused by these low degree deficiencies, the potential coefficients from degree 2 till 6 of the RBF solutions have been replaced by corresponding coefficients taken from EIGEN-6C for all further comparisons presented in this work. In the medium wavelengths (around 1,000 km spatial resolution), both methods deliver almost identical results indicating that GRACE has generally its highest sensitivity in this part of the gravity field spectrum. Looking at the shorter wavelengths, the RBF solutions show a potentially higher spatial resolution, as the smoothing applied to the RL05 solutions is damping the power in the higher degrees not only suppressing noise but likely also signal.

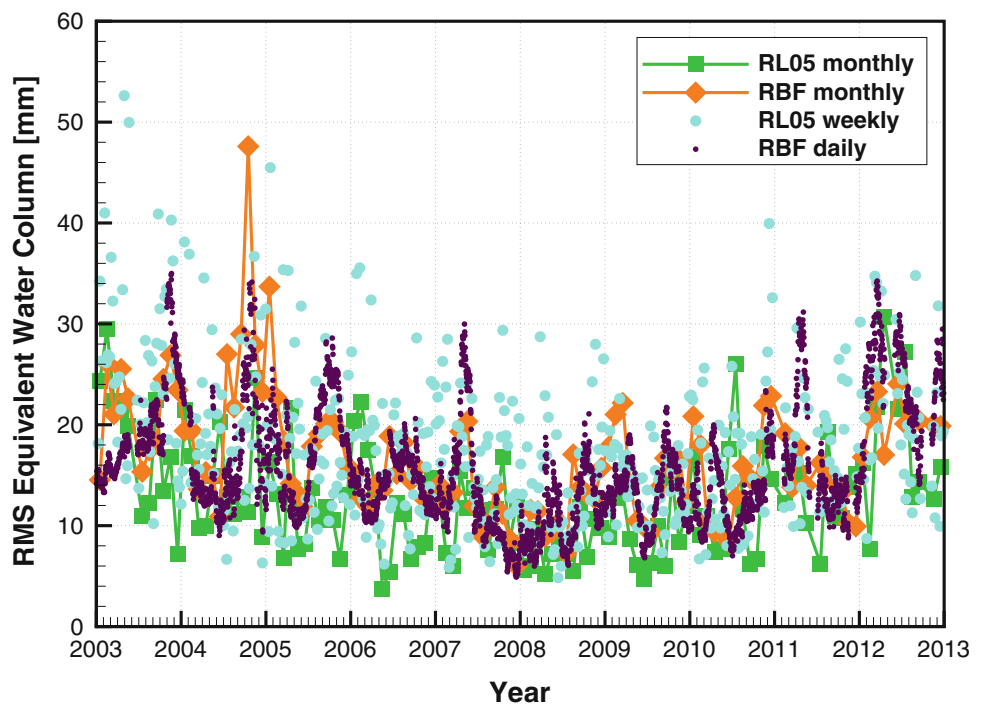
In order to compare the noise level of the different time-series, the root mean square (RMS) of equivalent water height (EWH) values located within central Sahara (where almost no mass variability is expected) have been computed (Fig. 3). The RL05 monthly time-series has the lowest noise level, but the monthly RBF solutions are very comparable for most epochs. The largest RMS values can be seen for the RL05 weekly solutions whose noise level is rather randomly distributed. The daily RBF solutions are strongly correlated with their monthly counterpart, although with sometimes larger peaks. Both RBF time-series show systematic variations possibly caused by seasonal changes of the variance of the WGHM model entering the prediction step of the daily solutions and subsequently also affecting the monthly ones. The mean noise level of all four time-series lies within 1–2 cm of EWH.

In Figs. 4 and 5, resp., basin averages of monthly RL05 and RBF solutions for Amazon and Bangladesh are plotted. Correlations between RL05 and RBF are very high

**Fig. 2** Average monthly degree variances over the period 2003–2012 in terms of water column [mm] for RBF solutions (dashed dark grey), RBF solutions with degrees 2 till 6 from EIGEN-6C (solid dark grey) and RL05 solutions (solid light grey)



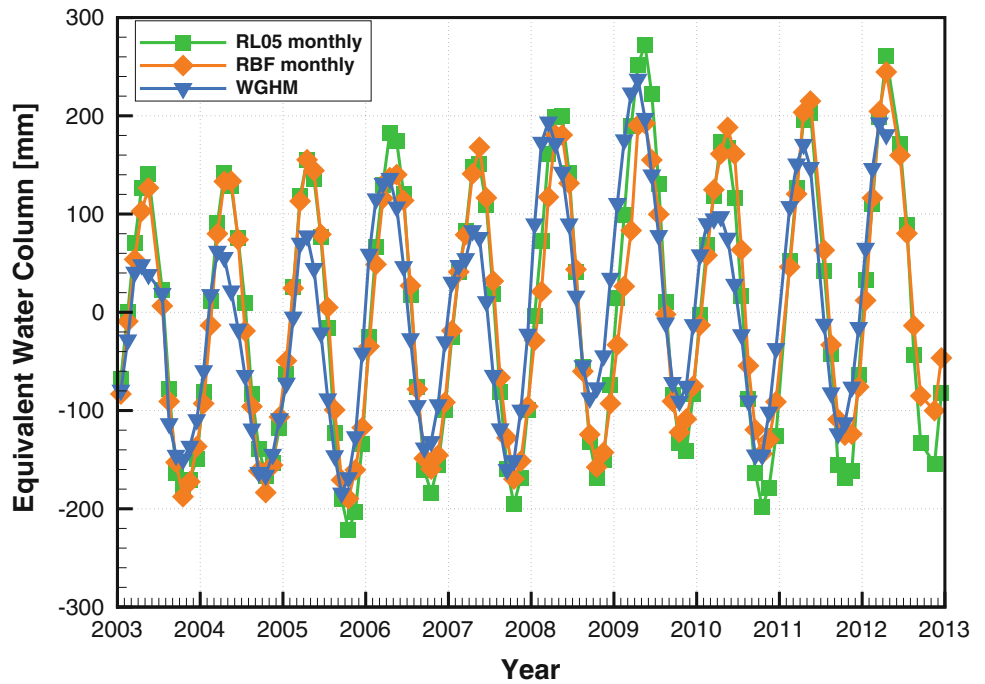
**Fig. 3** RMS of EWH values [mm] in central Sahara for RL05 monthly (green) and weekly (blue) solutions and RBF monthly (orange) and daily (purple) solutions



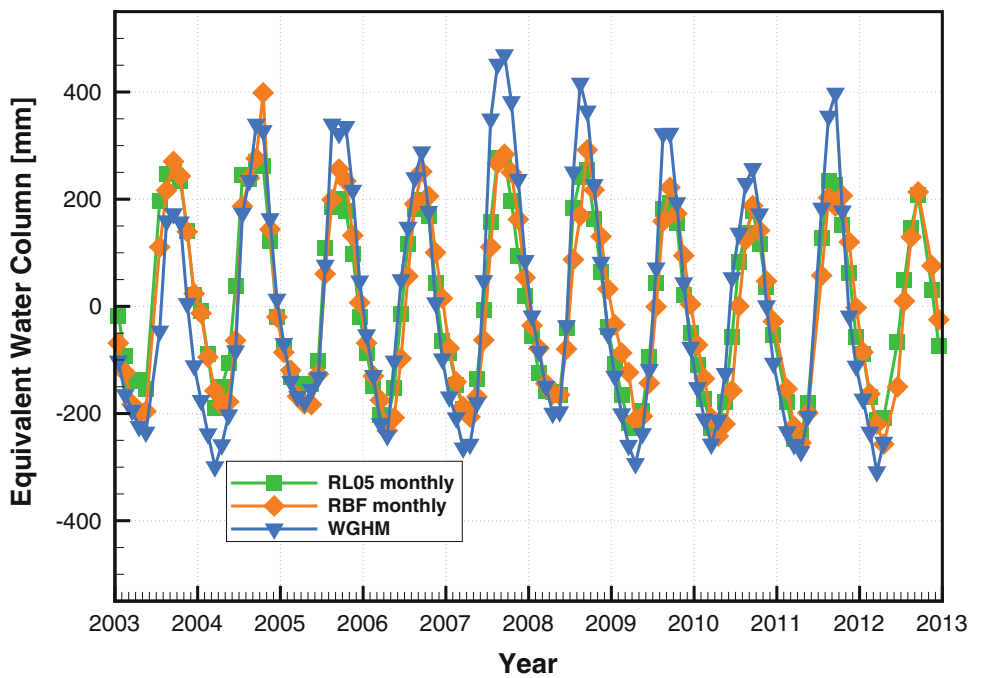
(Amazon: 98%, Bangladesh: 95%), and the correlations with (unsmoothed) monthly WGHM model output are also relatively good (see Table 2) with RL05 performing slightly

better. Amplitude differences between RL05 and RBF are generally rather small. WGHM obviously shows more pronounced year to year variations of the amplitudes. In case of

**Fig. 4** Basin averages in terms of EWH [mm] in the Amazon basin for RL05 (green) and RBF (orange) monthly solutions and WGHM (blue)



**Fig. 5** Basin averages in terms of EWH [mm] in Bangladesh for RL05 (green) and RBF (orange) monthly solutions and WGHM (blue)



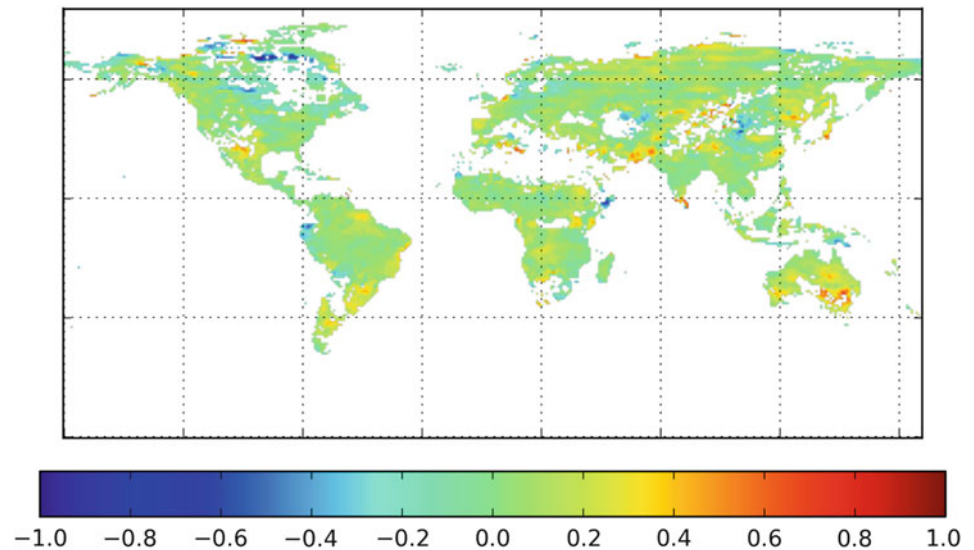
the Amazon basin, RL05 seems to better reflect these variations, whereas in Bangladesh RL05 and RBF amplitudes are almost similar. The basin averages of both submonthly time-series fit very smoothly to the corresponding monthly solutions (not shown).

The comparable correlations of RL05 and RBF w.r.t. WGHM become also visible in Fig. 6, where the differences

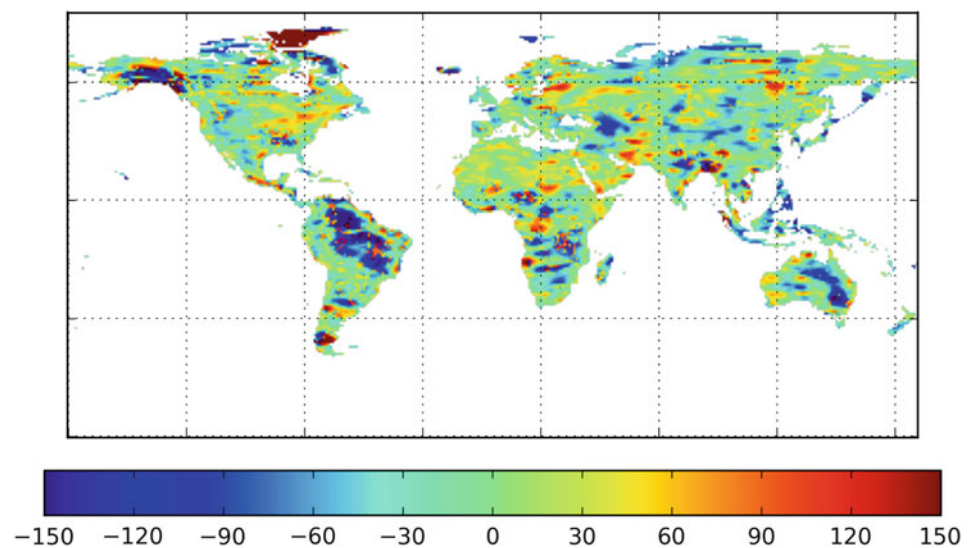
**Table 2** Correlations between monthly EWH time-series from GRACE and WGHM over the period 2003–2012

Basin	GRACE	Correlation (%)
Amazon	RL05	89
	RBF	86
Bangladesh	RL05	93
	RBF	90

**Fig. 6** Correlations between RBF and WGHM minus correlations between RL05 and WGHM; *blue areas* indicate higher correlations for RBF, *red areas* higher correlations for RL05



**Fig. 7** Amplitude differences between RBF and WGHM minus amplitude differences between RL05 and WGHM in terms of EWH [mm]; *blue areas* indicate smaller differences for RBF, *red areas* smaller differences for RL05



of these correlations are spatially plotted. Accordingly, these differences are close to zero in most regions and none of the two approaches outperforms the other in a global sense. More interesting conclusions can be drawn from Fig. 7 showing the spatial distribution of the difference of amplitude differences at each  $1^\circ \times 1^\circ$  grid point between RL05/RBF and WGHM. Amplitudes of the RBF solutions are much closer to WGHM in many regions, most prominently in the Amazon basin. From this, it can be concluded that the RBF approach is capable to better localize hydrological mass variations and suffers less from leakage effects and smearing of signal caused by smoothing.

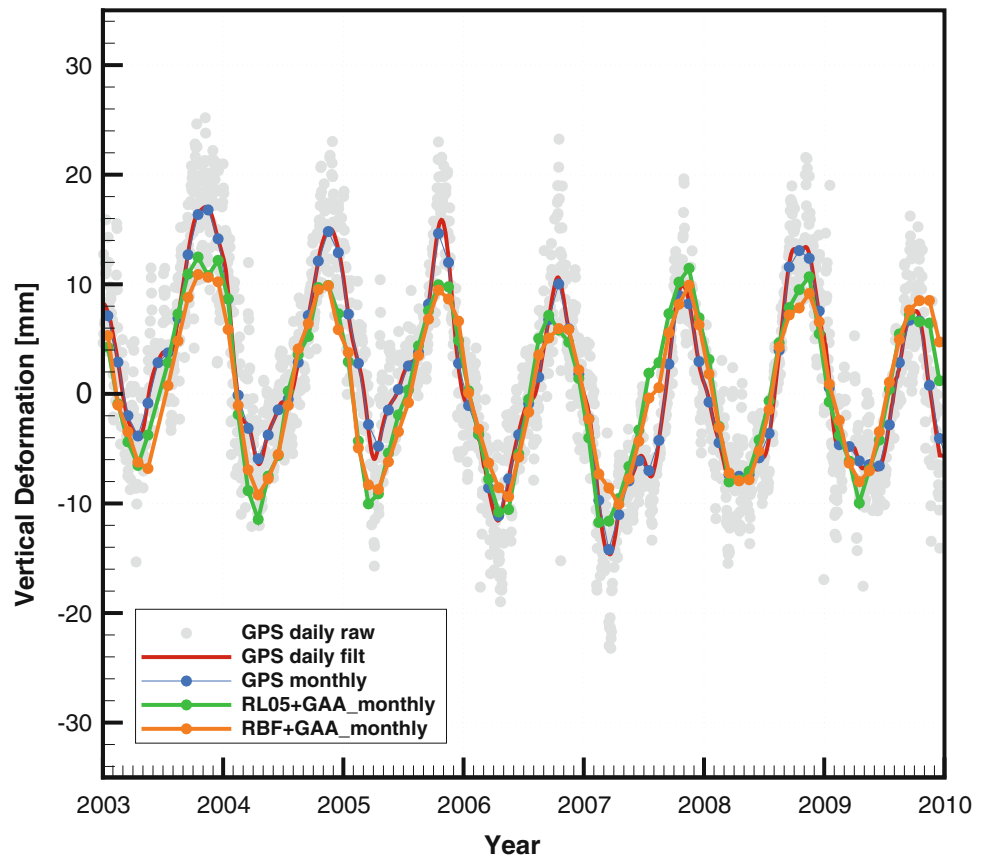
It has to be stated that the WGHM model, like other hydrology models, cannot be considered as absolute truth and in several other publications GRACE data has been

used to validate or even calibrate hydrology models (e.g. Werth et al. 2009). However, the focus of this work is on a relative comparison of different GRACE time-series rather than comparing GRACE with hydrology models. Furthermore, the discussion of Fig. 7 aims at the capability of the GRACE solutions to localize hydrology signals and in this context WGHM should perform better than GRACE and thus represent a suitable validation. The fact that stochastic a priori information from WGHM has entered the RBF solutions as described in Sect. 3 also does not affect the localization of hydrology signals in the solutions, i.e. better results for RBF can be considered as reasonable.

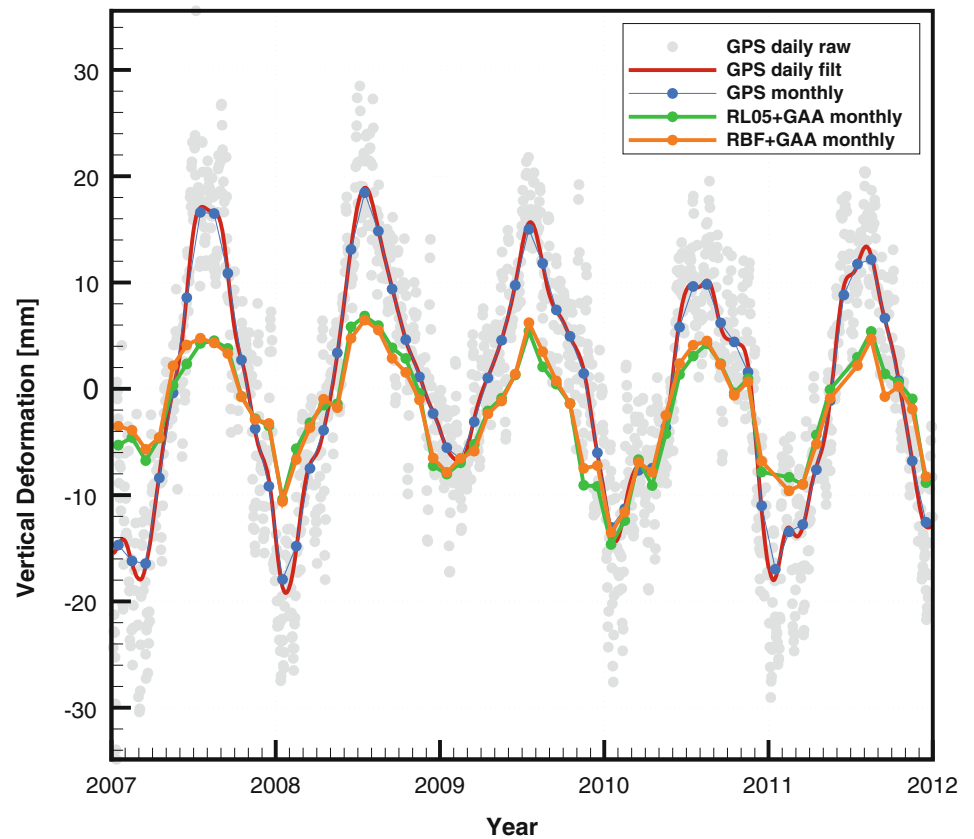
Further validation against independent data is done by comparing vertical deformations from GPS and the GRACE time-series. The GPS data used is obtained from CODE



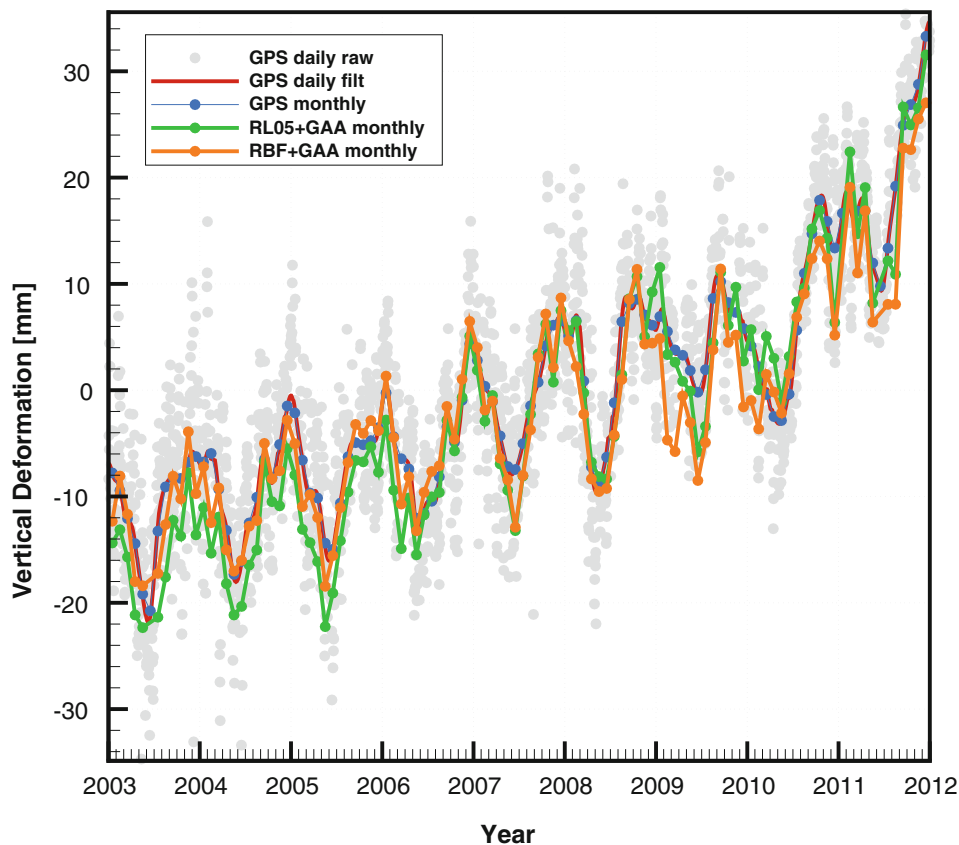
**Fig. 8** Vertical deformations [mm] at IGS Station BRAZ (15.85°S, 47.88°W) from GPS and monthly GRACE solutions



**Fig. 9** Vertical deformations [mm] at IGS Station NOVM (54.85°N, 82.91°E) from GPS and monthly GRACE solutions



**Fig. 10** Vertical deformations [mm] at IGS Station KELY (66.85°N, 50.94°W) from GPS and monthly GRACE solutions



**Table 3** Correlations between monthly vertical deformation time-series from GRACE and GPS over the period  $\Delta T$

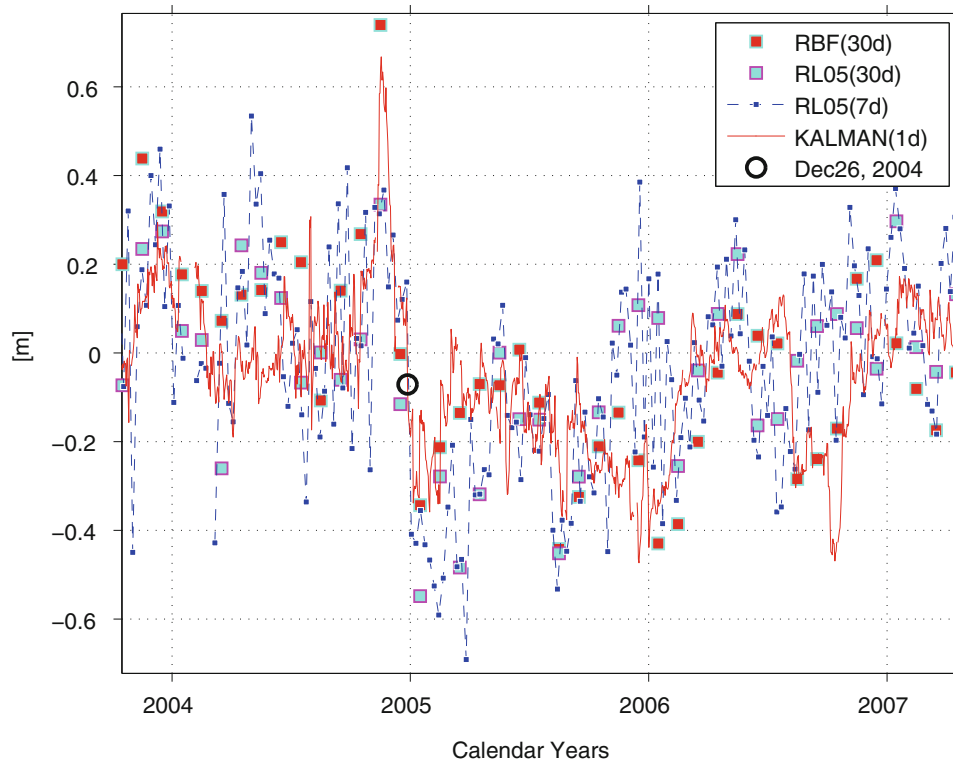
IGS station	$\Delta T$	GRACE	Correlation (%)
BRAZ	2003–2009	RL05	90
		RBF	90
NOVM	2007–2011	RL05	89
		RBF	90
KELY	2003–2011	RL05	96
		RBF	94

Reprocessing (Steigenberger et al. 2011) and consists of daily coordinate time-series of IGS stations. This original data is transformed to vertical deformations, smoothed and averaged to monthly values. The RL05 and RBF models are also converted to vertical deformations according to Tesmer et al. (2011). As variations in the GPS time-series reflect also atmospheric effects, the GAA-product (Flechtner et al. 2014) has been re-added to the GRACE solutions. Results are exemplarily shown for the GPS stations BRAZ (Brasilia,

Brazil), NOVM (Novosibirsk, Russia) and KELY (Kangerlussuaq, Greenland) in Figs. 8, 9 and 10. The corresponding correlations are listed in Table 3. Again, RL05 and RBF solutions give very similar results and both are well in phase with mass variations derived from GPS.

Finally, the benefit of GRACE time-series with higher temporal resolution is illustrated in Fig. 11. It shows EWH variations in an area of approx.  $(500 \times 500) \text{ km}^2$  size around the epicenter of the December 26th, 2004 Sumatra-Andaman earthquake. The event by itself is visible in all time-series, but the actual epoch of the earthquake cannot be captured exactly by the monthly solutions as these are 30-day averages. However, this becomes somewhat possible by analyzing the weekly and particularly the daily solutions. The latter capture the sudden mass change in the vicinity of the earthquake exactly at the actual epoch of this event and might allow for an improved description of the characteristics of such events. Although this is still work in progress, the daily Kalman solutions impart a viable enhancement, e.g. for applications which demand a higher temporal resolution.

**Fig. 11** Averages in terms of EWH [m] in an area around the epicenter of the 2004 Sumatra-Andaman earthquake ( $2^{\circ} - 6^{\circ}\text{N}$ ,  $94^{\circ} - 98^{\circ}\text{E}$ ) for monthly RBF and RL05 solutions, weekly RL05 solutions and daily RBF solutions; the actual day of the earthquake is denoted by the *black circle*



## 5 Summary and Conclusions

Two methods of gravity field recovery from GRACE have been compared: (1) the dynamic method, with GFZ RL05 used here as example, and (2) the RBF method.

The former is a standard method for gravity field processing, well-established over years, and therefore can be regarded as benchmark in case of global gravity field solutions. Its drawbacks are a relatively high computational effort, a limited spatial resolution as smoothing is required and the fact that only global solutions can be obtained.

The latter method is a new alternative approach based on inversion of integral equations with a comparatively low computational effort. The conversion of GRACE K-band observations to in-situ gravitational observations works well, as the results from the RBF method are in generally very good agreement with results from GFZ RL05. Furthermore, this approach can be used for both global and regional applications. A higher temporal resolution can be achieved by employing a Kalman filter. Considering the fact that the RBF solutions are directly given in form of gridded gravity functionals with no further smoothing needed, they appear to be more user-friendly than the standard GRACE SDS solutions. The constraining by WGHM of course adds information to the spectral density of the RBF solutions but does not predefine their spatial distribution. They are thus forced to comply with those phenomena that most likely occur on

short time scales where dominant mass redistributions of the Earth system obviously stem from continental, oceanic and atmospheric hydrology. It is then arguable that the RBF method is less suitable to derive secular trends or large-scale mass anomalies, but for short-term or even singular events such as megathrust Earthquakes it provides reasonable results. On the contrary, the post-filtering of the standard spherical harmonic solutions, which is effectively a bandpass filter, directly affects the signals' localization by annihilating higher resolution. Dealing with this leakage is a challenge with equal incertitude.

Results from both approaches have been validated against independent data. According to amplitude differences w.r.t. the WGHM model, the RBF method has a potentially higher spatial resolution and shows a better localization of mass change signals. Comparisons with in-situ GPS vertical deformations show slightly better results for GFZ RL05 than for RBF. The RBF daily solutions clearly depict the event of the 2004 Sumatra-Andaman earthquake proving that they are capable to provide additional information compared to the standard monthly or weekly solutions.

To answer the question posed in the title of this manuscript, it can be concluded that the results are not the same, but very similar. As they are based on the same background models, the appearing differences should be caused by the different approaches themselves indicating that each approach has particular strengths and weaknesses.

**Acknowledgements** This work has been funded by the German Federal Ministry of Education and Research (BMBF) with support code 03F0654A.

We would like to thank the German Space Operations Center (GSOC) of the German Aerospace Center (DLR) for providing continuously and nearly 100% of the raw telemetry data of the twin GRACE satellites.

We would also like to thank the editor, M. Weigelt, as well as U. Meyer and two anonymous reviewers for their helpful comments improving this manuscript.

## References

- Bettadpur S (2012) UTCSR level-2 processing standards document (for level-2 product release 0005). GRACE document 327–742. <ftp://podaac.jpl.nasa.gov/allData/grace/docs/>. Accessed 31 July 2013
- Bruinsma SL, Lemoine J-M, Biancale R, Valès N (2010) CNES/GRGS 10-day gravity field models (release 02) and their evaluation. *Adv Space Res* 45:587–601. doi:10.1016/j.asr.2009.10.012
- Dahle C, Flechtner F, Gruber C, König D, König R, Michalak G, Neumayer KH (2012) GFZ GRACE level-2 processing standards document for level-2 product release 0005. Scientific Technical Report STR12/02 - Data, Revised Edition, January 2013, Potsdam, 21 p. doi:10.2312/GFZ.b103-1202-25
- Dahle C, Flechtner F, Gruber C, König D, König R, Michalak G, Neumayer KH (2014) GFZ RL05: An Improved Time-Series of Monthly GRACE Gravity Field Solutions. In: Flechtner F et al (eds) Observation of the system earth from space – CHAMP, GRACE, GOCE and Future Missions, Advanced Technologies in Earth Sciences. Springer, Berlin/Heidelberg. doi:10.1007/978-3-642-32135-1\_4
- Döll P, Kaspar F, Lehner B (2003) A global hydrological model for deriving water availability indicators: model tuning and validation. *J Hydrol* 270:105–134. doi:10.1016/S0022-1694(02)00283-4
- Flechtner F, Dobslaw H, Fagiolini E (2014) AOD1B Product Description Document for Product Releases 05 (Rev. 4.2). GRACE Document 327–750
- Gruber C, Moon Y, Flechtner F, Dahle C, Novák P, König R, Neumayer KH (2014) Submonthly GRACE Solutions from Localizing Integral Equations and Kalman Filtering. In: Rizos C and Willis P (eds) Earth on the edge: science for a sustainable planet, international association of geodesy symposia, vol 139. Springer, Berlin/Heidelberg. doi:10.1007/978-3-642-37222-3\_51
- Kim J (2000) Simulation study of a low-low satellite-to-satellite tracking mission. Dissertation, University of Texas, Austin
- Kurtenbach E, Mayer-Gürr T, Eicker A (2009) Deriving daily snapshots of the Earth's gravity field from GRACE L1B data using Kalman filtering. *Geophys Res Lett* 36:L17102. doi:10.1029/2009GL039564
- Kusche J, Schmidt R, Petrovic S, Rietbroek R (2009) Decorrelated GRACE time-variable gravity solutions by GFZ, and their validation using a hydrological model. *J Geodesy* 83:903–913. doi:10.1007/s00190-009-0308-3
- Liu X, Ditmar P, Siemes C, Slobbe DC, Revtova E, Klees R, Riva R, Zhao Q (2010) DEOS Mass Transport model (DMT-1) based on GRACE satellite data: methodology and validation. *Geophys J Int* 181:769–788. doi:10.1111/j.1365-246X.2010.04533.x
- Mayer-Gürr T (2006) Gravitationsfeldbestimmung aus der Analyse kurzer Bahnbögen am Beispiel der Satellitenmissionen CHAMP und GRACE. Dissertation, University of Bonn
- Meyer U, Jäggi A, Beutler G (2012) Monthly gravity field solutions based on GRACE observations generated with the celestial mechanics approach. *Earth Planet Sci Lett* 345:72–80. doi:10.1016/j.epsl.2012.06.026
- Meyer U, Dahle C, Sneeuw N, Jäggi A, Beutler G, Bock H (2015) The effect of pseudo-stochastic orbit parameters on GRACE monthly gravity fields – insights from lumped coefficients. Submitted to international association of geodesy symposia, VIII Hotine-Marussi Symposium, Rome, 2013
- Novák P (2007) Integral inversion of SST data of type GRACE. *Stud Geophys Geod* 51:351–367. doi:10.1007/s11200-007-0020-9
- Reigber C, Schmidt R, Flechtner F, König R, Meyer U, Neumayer KH, Schwintzer P, Zhu SY (2005) An earth gravity field model complete to degree and order 150 from GRACE: EIGEN-GRACE02S. *J Geodyn* 39:1–10. doi:10.1016/j.jog.2004.07.001
- Rummel R (1975) Downward continuation of gravity information from satellite to satellite tracking or satellite gradiometry in local areas. Reports of the Department of Geodetic Science, vol 221. Ohio State University, Columbus, 50 pp
- Schmidt R (2007) Zur Bestimmung des cm-Geoids und dessen zeitlicher Variationen mit GRACE. Dissertation, Scientific Technical Report STR07/04, April 2007, Potsdam, 141 p. doi:10.2312/GFZ.b103-07042
- Shako R, Förste C, Abrikosov O, Bruinsma SL, Marty J-C, Lemoine J-M, Flechtner F, Neumayer KH, Dahle C (2014) EIGEN-6C: A High-Resolution Global Gravity Combination Model Including GOCE Data. In: Flechtner F et al (eds) Observation of the system earth from space – CHAMP, GRACE, GOCE and Future Missions, Advanced Technologies in Earth Sciences. Springer, Berlin/Heidelberg. doi:10.1007/978-3-642-32135-1\_20
- Steigenberger P, Hugentobler U, Lutz S, Dach R (2011) CODE contribution to the first IGS reprocessing campaign. Technical Report 1/2011, Institute for Astronomical and Physical Geodesy (IAPG), Technical University of Munich
- Tapley BD, Bettadpur S, Watkins M, Reigber C (2004) The gravity recovery and climate experiment: mission overview and early results. *Geophys Res Lett* 31:L09607. doi:10.1029/2004GL019920
- Tesmer V, Steigenberger P, van Dam T, Mayer-Gürr T (2011) Vertical deformations from homogeneously processed GRACE and global GPS long-term series. *J Geodesy* 85:291–310. doi:10.1007/s00190-010-0437-8
- Watkins M, Yuan DN (2012) JPL level-2 processing standards document for level-2 product release 05. GRACE document 327–744. <ftp://podaac.jpl.nasa.gov/allData/grace/docs/>. Accessed 31 July 2013
- Werth S, Güntner A, Petrovic S, Schmidt R (2009) Integration of GRACE mass variations into a global hydrological model. *Earth Planet Sci Lett* 277:166–173. doi:10.1016/j.epsl.2008.10.021

---

# The Effect of Pseudo-Stochastic Orbit Parameters on GRACE Monthly Gravity Fields: Insights from Lumped Coefficients

U. Meyer, C. Dahle, N. Sneeuw, A. Jäggi, G. Beutler, and H. Bock

---

## Abstract

The official GFZ RL05 monthly GRACE gravity models were processed in a two-step approach. In the first step the orbits were determined. In the second step corrections to the gravity field parameters were estimated, while the orbits were kept fixed. This led to a significant de-noising of the resulting monthly models, but accidentally also to a regularization, i.e., the estimated gravity field coefficients were biased towards the a priori model. We compare the GFZ RL05 models to a revised version RL05a that was determined in a common estimation of orbit and force model parameters. A large number of gravity field coefficients is significantly affected. We relate this effect to the one-hourly stochastic accelerations estimated for orbit determination, and to ignoring the correlations.

In the main part of this paper we study the interaction between pseudo-stochastic orbit parameters and gravity field coefficients. To explain this interaction we make use of a time-wise approach to gravity field determination. We apply the linear perturbation theory developed by Hill for circular orbits to compute lumped coefficients of the inter-satellite range-rate observations. We illustrate that the pseudo-stochastic orbit parameters act as a high-pass filter on the lumped coefficients spectra of the range-rates. Because the lumped coefficients are related to the spherical harmonics coefficients via a summation over all degrees, the whole range of gravity field coefficients is affected.

This result is of relevance for all approaches to gravity field estimation from orbit observations, where dynamic orbits are introduced a priori and the arc-specific parameters are kept fixed.

---

## Keywords

GRACE • Satellite gravimetry • Time-wise approach

---

U. Meyer (✉) • A. Jäggi • G. Beutler • H. Bock  
Astronomical Institute, University of Bern, Sidlerstrasse 5, 3012 Bern,  
Switzerland  
e-mail: [ulrich.meyer@aiub.unibe.ch](mailto:ulrich.meyer@aiub.unibe.ch)

C. Dahle  
GFZ German Research Centre for Geosciences, Potsdam, Germany

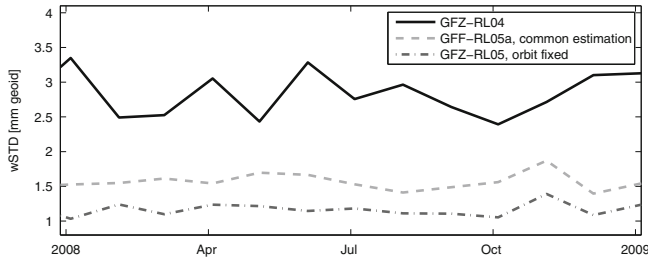
N. Sneeuw  
Institute of Geodesy, University of Stuttgart, Stuttgart, Germany

---

## 1 Introduction

Gravity field estimation using observations of, e.g., the GRACE satellite mission is a non-linear parameter estimation process. It is common practice to solve it as a generalized orbit determination problem. The ‘true’ solution is found, possibly iteratively, provided that linearization errors are small and all parameters of the physical model are determined in one common parameter estimation procedure.

A priori information of the physical models (in particular of the a priori gravity field) may affect the resulting gravity



**Fig. 1** Weighted standard deviations over the oceans for different sets of monthly solutions

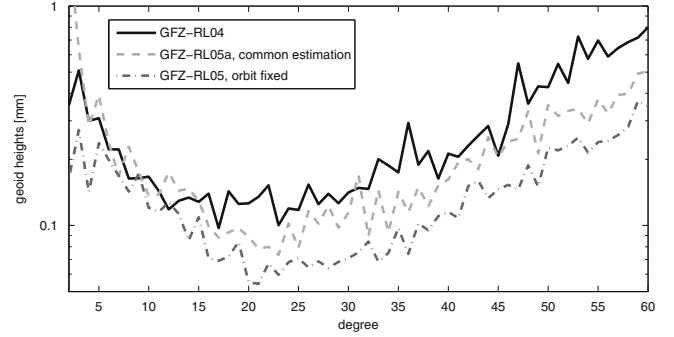
field, if the above principles are violated, e.g., when orbit determination and gravity field estimation are artificially separated. In this case dynamic orbit parameters absorb residual gravity signals and the solution is biased towards the a priori model. This happened to the official GFZ RL05 monthly gravity fields (Dahle et al. 2012) and was first reported by Meyer et al. (2013). Meanwhile the findings were confirmed independently, and, as a consequence, GFZ has produced a new release RL05a, where orbit and gravity field parameters are estimated together (Flechtner et al. 2013).

The processing strategy of GFZ RL05 was motivated by a significant de-stripping effect when the orbits of the GRACE satellites and the spherical harmonics coefficients (SHC) were estimated separately. The noise, e.g., evaluated by the standard deviation of the gridded time variable signal over the oceans (weighted by the cosine of the latitude), drops dramatically in comparison to monthly solutions, where the SHC are estimated together with the orbits (RL05a, see Fig. 1). The idea of separating orbit and gravity field estimation is not new. It has been adopted, e.g., by Lutcke et al. (2006) who reported a comparable de-stripping effect.

A number of experiments were performed at the Astronomical Institute of the University of Bern (AIUB) in order to understand this apparent gain in quality, where the processing strategy of GFZ was emulated using the Celestial Mechanics Approach (CMA; Beutler et al. 2010a,b). We found that the gain is related to the following two-step mechanism:

- solving in the first step for hourly pseudo-stochastic orbit parameters (e.g., empirical piecewise constant accelerations) while using a good a priori gravity model, i.e., a model already including trends and periodic seasonal variations up to a degree of, e.g., 30, and
- introducing in the second step the orbits from the first step for the estimation of the SHC of the gravity field. The orbit parameters are kept fixed and the correlations with the force model parameters are thereby lost.

This procedure obviously regularizes the monthly solutions, where the quality of the a priori gravity model and the spacing of the pseudo-stochastic orbit parameters are the relevant style elements.



**Fig. 2** Difference degree variances between monthly solutions (March 2008) and the time variable a priori field (EIGEN-6C), evaluated at the mid-epoch of the monthly models

We found furthermore that pseudo-stochastic orbit parameters, e.g., equally spaced by 1 h, greatly improve the consistency between the a priori gravity model and a monthly solution. The whole spectrum of degree variances is affected (Fig. 2). Subsequently, an explanation for this effect is presented. We make use of a representation of the inter-satellite range-rates in the time domain and interpret the use of pseudo-stochastic orbit parameters as a rather special high-pass filtering.

The Fourier coefficients of the range-rates are linked to the SHC of the gravity field (Sneeuw 2000) via a summation over degree  $l$ . This ‘lumping’ process explains the impact of the low frequency Fourier coefficients, dampened by signal absorption through the pseudo-stochastic orbit parameters, on the whole spectrum of SHC.

## 2 The Time-Wise Approach: Gravitational Potential Along the Orbit

A representation of the disturbing gravity field potential in osculating elements was already presented by Kaula (1966), the idea of time-wise orbit analysis for gravity field determination originates from Wagner (1983). We make use of an alternative formulation proposed by Sneeuw (2000), where the gravitational potential is written as a function of the orbital elements radius  $r$ , the inclination  $I$ , the argument of latitude  $u$ , and the argument of longitude  $\Lambda = \Omega - \theta$  (with  $\Omega$  the longitude of the ascending node and  $\theta$  the Greenwich sidereal time):

$$V(r, I, u, \Lambda) = \frac{GM}{r} \sum_{m=0}^L \sum_{k=-L}^L \sum_{l=\max(m,|k|)}^L \left(\frac{a_E}{r}\right)^l \cdot \bar{F}_{lmk}(I) \begin{cases} \bar{C}_{lm} \cos \psi_{mk} + \bar{S}_{lm} \sin \psi_{mk} \\ -\bar{S}_{lm} \cos \psi_{mk} + \bar{C}_{lm} \sin \psi_{mk} \end{cases} \begin{matrix} l-m \text{ even} \\ l-m \text{ odd} \end{matrix} \quad (1)$$

$a_E$  is the mean equatorial radius of Earth and  $GM$  the product of constant of gravity and Earth's mass. The summations run over degree  $l$  and order  $m$  of the normalized spherical harmonics coefficients  $\bar{C}_{lm}$  and  $\bar{S}_{lm}$  (the maximum degree and order being  $L$ ), as well as over the index  $k$  stemming from the inclination functions  $\bar{F}_{lmk}(I)$ , see Sneeuw (1992). The angles  $\psi_{mk} = ku + m\Lambda$  comprise the position of the satellite.

If we confine ourselves to circular orbits (close to reality for the GRACE satellites), the radius  $r$  and the inclination functions  $\bar{F}_{lmk}(I)$  are constant. In this case Eq. (1) takes the form of a Fourier series

$$V(\psi_{mk}) = \sum_{m=0}^L \sum_{k=-L}^L A_{mk}^V \cos \psi_{mk} + B_{mk}^V \sin \psi_{mk} \quad (2)$$

with constant Fourier coefficients  $A_{mk}^V$  and  $B_{mk}^V$ . These coefficients are called lumped coefficients<sup>1</sup> because they can be derived from the SHC via a summation over degree  $l$ :

$$A_{mk}^V = \sum_{l=\max(m,|k|)}^L \bar{H}_{lmk}^V \begin{cases} \bar{C}_{lm} & l-m \text{ even} \\ -\bar{S}_{lm} & l-m \text{ odd} \end{cases} \quad (3)$$

$$B_{mk}^V = \sum_{l=\max(m,|k|)}^L \bar{H}_{lmk}^V \begin{cases} \bar{S}_{lm} & l-m \text{ even} \\ \bar{C}_{lm} & l-m \text{ odd} \end{cases}. \quad (4)$$

The transfer between the two spectral representations of the gravity potential is accomplished via the transfer coefficients (Sneeuw 2000)

$$\bar{H}_{lmk}^V = \frac{GM}{r} \left(\frac{a_E}{r}\right)^l \bar{F}_{lmk}(I). \quad (5)$$

We can assign frequencies  $\dot{\psi}_{mk} = k\dot{u} + m\dot{\Lambda}$  to the lumped coefficients, and if we write the repeating frequency of the orbits  $\dot{u}/|\dot{\Lambda}|$  as integer ratio  $\beta/\alpha$ , with  $\beta$  orbital revolutions within  $\alpha$  sidereal days, we find the corresponding Fourier frequencies  $\omega_j = \dot{u}/\beta \cdot j = \dot{u}/\beta(k\beta - m\alpha)$ .

### 3 Spectral Representations of Orbit Perturbations and Inter-Satellite Range-Rates

The derivation of a spectral representation of inter-satellite range-rates, the main observation type of the GRACE mission, is detailed in Sneeuw (2000). We will only outline the general procedure and quote ready to use formulas for orbit perturbations and range-rates, because all expressions

<sup>1</sup>Not to be confused with the lumped coefficients introduced by Gooding (1971) for resonant analyses.

in Sneeuw (2000) are derived in complex notation and the transformation to real-valued expressions for computer-coding is not trivial.

First we define an orbital frame centered at and co-rotating with the satellite. Its axes are always pointing in along-track ( $x$ ), cross-track ( $y$ ) and radial ( $z$ ) direction of the satellite. We then apply the gradient operator  $\nabla$  to Eq. (1) and derive gravitational accelerations in the orbital frame.

To relate the gravitational accelerations to orbit perturbations a perturbation theory is needed. For the case of circular orbits the linear system of differential equations with constant coefficients originally developed by Hill (1878) for the description of the lunar orbit around the Earth is a convenient choice. With these equations orbit perturbations  $\Delta x$ ,  $\Delta y$  and  $\Delta z$  are determined relative to a reference point that moves with constant angular velocity  $n = \sqrt{GM/r^3}$  along a circular orbit. The linear differential equations hold only approximately for the GRACE orbits, but they have the big advantage of being solvable analytically.

Their general solution consists of the solution of the homogeneous equations and a particular solution of the inhomogeneous equations. From the latter one the spectral transfer of orbit perturbations in the orbital frame is derived. We quote the expressions for the along-track and radial perturbations, because they are needed to compute range-rates:

$$\bar{H}_{lmk}^{\Delta x} = \frac{(3n^2 + \dot{\psi}_{mk}^2)k - 2n\dot{\psi}_{mk}(l+1)}{\dot{\psi}_{mk}^2(n^2 - \dot{\psi}_{mk}^2)} \cdot \frac{GM}{a_E^2} \left(\frac{a_E}{r}\right)^{l+2} \bar{F}_{lmk}(I) \quad (6)$$

$$\bar{H}_{lmk}^{\Delta z} = \frac{2nk - (l+1)\dot{\psi}_{mk}}{\dot{\psi}_{mk}(n^2 - \dot{\psi}_{mk}^2)} \frac{GM}{a_E^2} \left(\frac{a_E}{r}\right)^{l+2} \bar{F}_{lmk}(I). \quad (7)$$

The spectral transfer becomes singular in case of resonance. Resonance occurs whenever  $\dot{\psi}_{mk} = 0$ , or  $\dot{\psi}_{mk} = \pm n$ , i.e.,  $m/k = \beta/\alpha$ , or  $m/(k \pm 1) = \beta/\alpha$ .

The lumped coefficients  $A_{mk}^{\Delta z}$  and  $B_{mk}^{\Delta z}$  of radial orbit perturbations  $\Delta z$  can be computed from the SHC by multiplication with the spectral transfer  $\bar{H}_{lmk}^{\Delta z}$  and summation over degree  $l$  according to Eqs. (3) and (4). For along-track orbit perturbations  $\Delta x$  the SHC are arranged as follows:

$$A_{mk}^{\Delta x} = \sum_{l=\max(m,|k|)}^L \bar{H}_{lmk}^{\Delta x} \begin{cases} \bar{S}_{lm} & l-m \text{ even} \\ \bar{C}_{lm} & l-m \text{ odd} \end{cases} \quad (8)$$

$$B_{mk}^{\Delta x} = \sum_{l=\max(m,|k|)}^L \bar{H}_{lmk}^{\Delta x} \begin{cases} -\bar{C}_{lm} & l-m \text{ even} \\ \bar{S}_{lm} & l-m \text{ odd} \end{cases}. \quad (9)$$

Finally, the derivation of transfer coefficients of inter-satellite ranges  $\Delta\rho$  is straightforward under the assumption

that both satellites follow each other with a certain time lag  $\tau$  (corresponding to a separation angle  $\eta = \tau n$ ) on the same circular orbit. The perturbation in the range is calculated by projecting both radial ( $\Delta x$ ) and along-track ( $\Delta z$ ) orbit perturbations onto the line of sight between the satellites. A derivation can be found in Sneeuw (2000), an alternative formulation in Visser (2005). We only quote the resulting transfer coefficients

$$\begin{aligned} \bar{H}_{lmk}^{\Delta\rho} &= 2 \cos \eta \sin\left(\eta \frac{\dot{\psi}_{mk}}{n}\right) \bar{H}_{lmk}^{\Delta x} \\ &+ 2 \sin \eta \cos\left(\eta \frac{\dot{\psi}_{mk}}{n}\right) \bar{H}_{lmk}^{\Delta z}. \end{aligned} \quad (10)$$

The summation of lumped coefficients  $A_{mk}^{\Delta\rho}$  and  $B_{mk}^{\Delta\rho}$  corresponds to Eqs. (3) and (4).

To derive the transfer coefficients for range-rates  $\Delta\dot{\rho}$  the time derivative has to be computed. Since only  $\psi_{mk}$  is time dependent ( $\dot{\psi}_{mk}$  is considered to be constant along the circular orbit), this eventually results in

$$\begin{aligned} \bar{H}_{lmk}^{\Delta\dot{\rho}} &= 2\dot{\psi}_{mk} \sin \eta \cos\left(\eta \frac{\dot{\psi}_{mk}}{n}\right) \bar{H}_{lmk}^{\Delta z} \\ &- 2\dot{\psi}_{mk} \cos \eta \sin\left(\eta \frac{\dot{\psi}_{mk}}{n}\right) \bar{H}_{lmk}^{\Delta x}. \end{aligned} \quad (11)$$

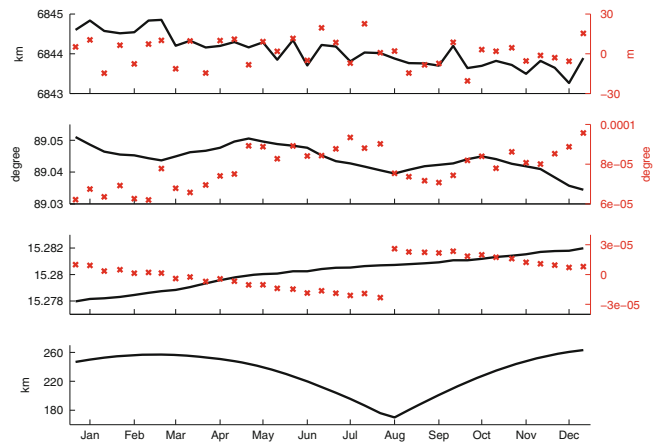
In this case the summation of lumped coefficients  $A_{mk}^{\Delta\dot{\rho}}$  and  $B_{mk}^{\Delta\dot{\rho}}$  corresponds to Eqs. (8) and (9).

## 4 Application to GRACE

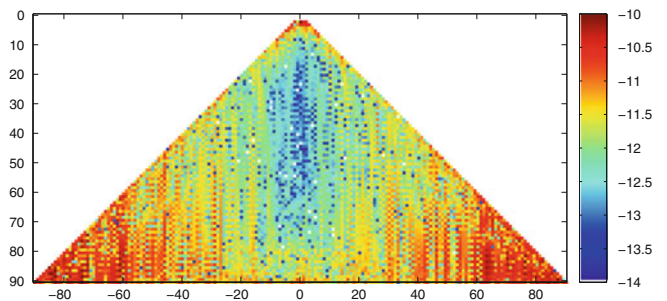
Now, we have the tools at hand to calculate the lumped coefficients of the range-rates from SHC of the gravity field. We will study the impact of the pseudo-stochastic orbit parameters on the spectra of the lumped coefficients on the basis of the monthly gravity fields of the year 2008. Therefore we have to determine mean orbital elements  $r$ ,  $I$  and orbital frequencies  $\dot{u}$  and  $\dot{\Lambda}$  from the corresponding GRACE orbits. As we only want to present a proof of concept, we ignore the subtle differences between the fixed orbits and the orbits resulting from a common estimation of arc specific and force model parameters.

In Fig. 3 mean orbital elements for GRACE A and the differences to GRACE B are presented throughout the whole year 2008. Instead of orbital frequencies the repeat ratio  $\dot{u}/\dot{\Lambda}$  is shown. The separation angle  $\eta$  can be calculated from the distance between both satellites. In August an orbit maneuver caused a jump in the orbital elements of GRACE B.

Let us first have a look at an example month (March 2008), and take the differences between the SHC of the GFZ RL05a solution (where the arc-specific parameters were pre-



**Fig. 3** Mean orbital elements of GRACE A (black line, top: radius, second: inclination, third: revolutions per sidereal day), differences between GRACE A and B (red crosses referring to right y-axis) and separation between both satellites (bottom)

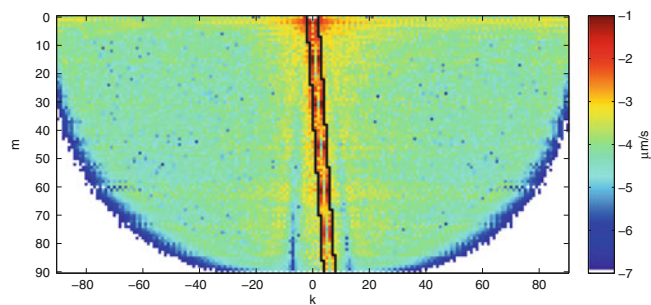


**Fig. 4** Differences between the SHC (dimensionless, color scale logarithmic) of GFZ RL05 and RL05a for March 2008

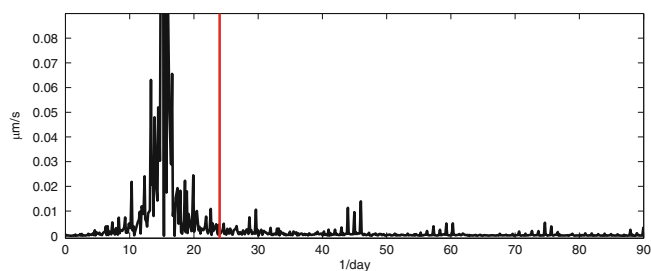
eliminated, and the correlations with gravity field parameters are correctly modeled), and the corresponding solution GFZ RL05 (where they were deleted from the normal equations, and the correlations are lost). In the triangle plot of SHC (Fig. 4) one sees differences in the very low degrees, which are intuitively correlated with the low frequency stochastic accelerations, and at high orders, where the SHC are probably dominated by noise. But one also sees vertical stripes throughout the whole range of orders, slightly more pronounced near resonant orders (15, 31, 46, 61). Even the alternation between even and odd degrees elaborated in the formulas of the lumped coefficients approach is clearly visible along these stripes.

Along the mean, circular orbit of GRACE A at the 15th of March we compute the transfer coefficients [Eqs. (6), (7), and (11)], and consequently the lumped coefficients of the inter-satellite range-rates from the differences between the two sets of SHC (Fig. 4). The difference for  $C_{20}$  is ignored, because this coefficient is not well determined from GRACE range-rate observations. Figure 5 shows the amplitudes  $\sqrt{(A_{mk}^{\Delta\dot{\rho}})^2 + (B_{mk}^{\Delta\dot{\rho}})^2}$  of the lumped coefficients,





**Fig. 5** Amplitudes of lumped coefficients of range-rates (color scale logarithmic). Frequencies below 24/day are located between the two black lines



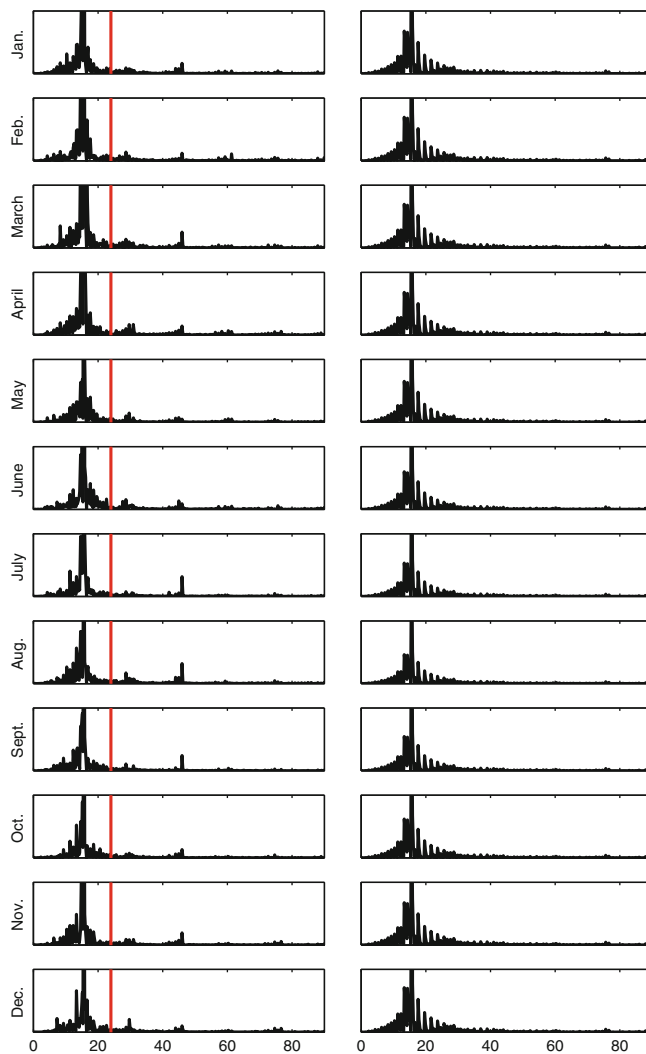
**Fig. 6** Amplitude spectrum of range-rates. The frequency of stochastic parameters 24/day is marked by a red line

ordered by index  $k$  and order  $m$ . Most of the signal is concentrated in a small band of coefficients.

Now, we compute the corresponding frequencies  $\dot{\psi}_{mk} = k\dot{u} + m\dot{\Lambda}$  for each of the lumped coefficients and mark all coefficients with frequencies below 24/day (frequency of the stochastic accelerations) in Fig. 5. The dominance of the low frequency lumped coefficients is striking. Arranging the lumped coefficients by frequency, the corresponding amplitude spectrum is derived and shown in Fig. 6. The significant part of the signal is concentrated at frequencies below 24/day. Comparable results are obtained for all the other sets of RL05 and RL05a monthly gravity fields from 2008 (not shown).

To be more flexible concerning the orbit parametrization we compute monthly models for the whole year 2008 with the CMA with processing strategies corresponding to GFZ RL05 and RL05a. Gravity field coefficients are only set up to a maximum degree and order of 60 to minimize the influence of noisy higher order SHC on the resulting amplitude spectra, presented in Fig. 7 (left). To also visualize the characteristics of the transfer coefficients  $\bar{H}_{lmk}^{\Delta\hat{\rho}}$  they are applied to a set of SHC all equal to one, and the resulting spectra are provided in Fig. 7 (right).

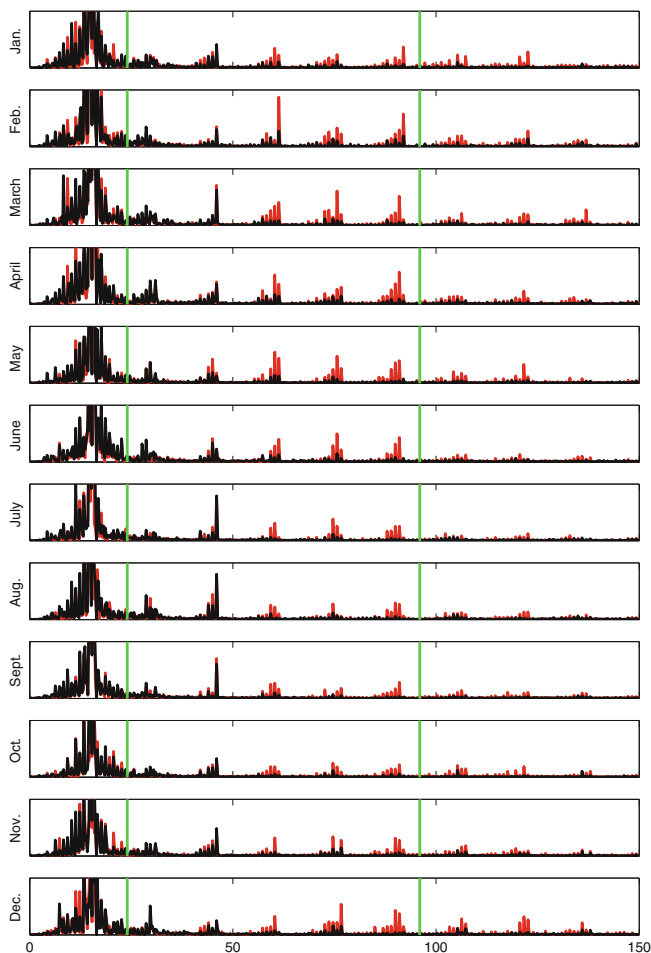
As the spectra (Figs. 6 and 7, left) were computed from the differences in the SHC between monthly solutions that were estimated, either separately, or together with the arc specific parameters, they show the signal absorbed by the



**Fig. 7** Effect of estimating and fixing stochastic accelerations at 60 min intervals on range-rates (left) and the characteristics of the corresponding transfer coefficients (right). The y-axis in the left hand plots corresponds to Fig. 6, the y-axis in the right hand plots is scaled by a factor of  $5 \cdot 10^{11}$ . Note that the amplitude at orbital resonance is far beyond the scale in all of the figures

orbit parameters, when the correlations with the force model parameters are ignored. One could conclude, that these spectra support the idea of a high-pass filter of the range-rates by the one-hourly stochastic accelerations. But the characteristics of the transfer coefficients (Fig. 7, right) reveal that most of the signal has to be expected below the 24/day frequency anyway.

We therefore performed one more experiment, where the stochastic accelerations were estimated at 15 min instead of 1-h intervals. A zoom into the resulting amplitude spectra that focuses on the frequency range beyond 24/day is presented in Fig. 8. Now we observe considerably higher amplitudes at frequencies up to 96/day. We take this as a strong indication in favor of the predicted high-pass filtering



**Fig. 8** Comparison of the effects of stochastic accelerations estimated at either 15 min (red) or 60 min (black) intervals. The two frequencies are marked by green vertical lines. In comparison to Fig. 7 the y-range was zoomed in by a factor of 4

effect. The signal visible beyond the frequency of 96/day must be attributed to the ever present noise in the estimated SHC which is amplified at multiples of the orbit frequency (and is of course also part of the signal visible at lower frequencies). Note as well, that the SHC, from which we started, were estimated along the true, slightly elliptical GRACE orbits and that any deviation from the circular orbits of the linear perturbation theory, applied for the derivation of the spectral transfer of the range-rates, tends to broaden the peaks and to fill up the computed spectra.

## 5 Summary and Conclusions

It was shown how monthly gravity field solutions are regularized by the introduction of satellite orbits fixed at some point of the estimation process of the SHC. The key parameters are the pseudo-stochastic parameters adjusted during the deter-

mination of the dynamic orbits and the quality of the a priori gravity model used for the orbits. The mechanism governing the process is a high-pass filter of the orbit perturbations represented by the pseudo-stochastic orbit parameters, intensified by breaking the correlations between orbit and gravity field parameters. This correlation could be illustrated by the introduction of a spectral presentation of the inter-satellite range-rates via lumped coefficients. Our theory explains how the signal absorbed by the pseudo-stochastic orbit parameters and visible in the spectra of the lumped coefficients of the range-rates influences the degree variances of the entire SH spectrum via the inversion of the lumping-effect (summation over degree  $l$ ).

As pseudo-stochastic parameters are common in precise orbit determination, every approach to gravity field estimation, where dynamic orbits are introduced as observations and the orbit parameters are kept fixed, will suffer from the described mechanism.

**Acknowledgements** We would like to thank P. Visser and two anonymous reviewers for their valuable comments which helped to improve the manuscript.

## References

- Beutler G, Jäggi A, Mervart L, Meyer U (2010a) The celestial mechanics approach: theoretical foundations. *J Geodesy* 84:605–624. doi:10.1007/s00190-010-0401-7
- Beutler G, Jäggi A, Mervart L, Meyer U (2010b) The celestial mechanics approach: application to data of the GRACE mission. *J Geodesy* 84:661–681. doi:10.1007/s00190-010-0402-6
- Dahle C, Flechtner F, Gruber C, König D, König R, Michalak G, Neumayer KH (2012) GFZ GRACE level-2 processing standards document for level-2 product release 0005. Scientific Technical Report STR12/02. doi:10.2312/GFZ.b103-1202-25
- Flechtner F, Dahle C, Gruber C, Sasgen I, König R, Michalak G, Neumayer KH (2013) Status GFZ RL05 and RL05a GRACE L2 products. In: Proceedings of the GRACE science team meeting 2013, Austin, Texas, 23–25 Oct 2013. Available at <http://www.csr.utexas.edu/grace/GSTM/>
- Gooding RH (1971) Lumped fifteenth-order harmonics in the geopotential. *Nat Phys Sci* 231:168–169. doi:10.1038/physci231168a0
- Hill GW (1878) Researches in the lunar theory. *Am J Math* I:5–26, 129–147, 245–260
- Kaula WM (1966) Theory of satellite geodesy. Blaisdell Publishing Company, Waltham, pp 30–37
- Luthcke SB, Rowlands DD, Lemoine FG, Klosko SM, Chinn D, McCarthy JJ (2006) Monthly spherical harmonic gravity field solutions determined from GRACE inter-satellite range-rate data alone. *Geophys Res Lett* 33:L02402. doi:10.1029/2005GL024846
- Meyer U, Jäggi A, Beutler G, Bock H (2013) The role of a priori information in gravity field determination. *Geophysical Research Abstracts* 15, EGU2013-9008. Presentation available at [http://www.berne.unibe.ch/publist/2013/pres/EGU2013\\_UM.pdf](http://www.berne.unibe.ch/publist/2013/pres/EGU2013_UM.pdf)
- Sneeuw N (1992) Representation coefficients and their use in satellite geodesy. *Manusc Geodaet* 17:117–123
- Sneeuw N (2000) A semi-analytical approach to gravity field analysis from satellite observations. Deutsche Geodätische Kommission, Reihe C 527:22–34

- 
- Visser PNAM (2005) Low-low satellite-to-satellite tracking: a comparison between analytical linear orbit perturbation theory and numerical integration. *J Geodesy* 79:160–166. doi:10.1007/s00190-005-0455-0
- Wagner CA (1983) Direct determination of gravitational harmonics from low-low GRAVSAT data. *J Geophys Res* 88:10309–10321. doi:10.1029/JB088iB12p10309

---

# On an Iterative Approach to Solving the Nonlinear Satellite-Fixed Geodetic Boundary-Value Problem

Marek Macák, Karol Mikula, Zuzana Minarechová, and Róbert Čunderlík

---

## Abstract

The paper deals with an iterative treatment of solving the nonlinear satellite-fixed geodetic boundary-value problem (NSFGBVP). To that goal we formulate the NSFGBVP consisting of the Laplace equation in 3D bounded domain outside the Earth. The computational domain is bounded by the approximation of the Earth's surface where the nonlinear boundary condition (BC) with prescribed magnitude of the gravity vector is given and by a spherical boundary placed approximately at the altitude of chosen satellite mission on which the Dirichlet BC for disturbing potential obtained from the satellite only geopotential model is applied. In case of local gravity field modelling, we add another four side boundaries where the Dirichlet BC is prescribed as well. The concept of our iterative approach is based on determining the direction of actual gravity vector together with the value of the disturbing potential. Such an iterative approach leads to the first iteration where the classical fixed gravimetric boundary-value problem with the oblique derivative BC is solved and the last iteration represents the approximation of the actual disturbing potential and the direction of gravity vector. As a numerical method for our approach, the finite volume method has been implemented. The practical numerical experiments deal with the local and global gravity field modelling. In case of local gravity field modelling, namely in the domain above Slovakia, the disturbing potential as a direct numerical result is transformed to the quasigeoidal heights and tested by the GPS-levelling. Results show an improvement in the standard deviation for subsequent iterations in solving NSFGBVP as well as the convergence to EGM2008. The differences between the last and the first iteration, which represent the numerically obtained linearization error, reach up to 10 cm. In case of global gravity field modelling, our solution is compared with the disturbing potential generated from EGM2008. The obtained numerical results show that the error of the linearization can exceed several centimeters, mainly in high mountainous areas (e.g. in Himalaya region they reach 20 cm) as well as in areas along the ocean trenches (varying from  $-2.5$  to  $2.5$  cm).

---

## Keywords

Finite volume method • Iterative approach • Nonlinear boundary value problem

---

This work was supported by grant APVV-0072-11 and VEGA 1/1063/11.

M. Macák (✉) • K. Mikula • Z. Minarechová • R. Čunderlík  
Faculty of Civil Engineering, Department of Mathematics  
and Descriptive Geometry, Slovak University of Technology,

---

Bratislava, Slovakia  
e-mail: [macak@math.sk](mailto:macak@math.sk); [mikula@math.sk](mailto:mikula@math.sk); [minarechova@math.sk](mailto:minarechova@math.sk);  
[cunderli@svf.stuba.sk](mailto:cunderli@svf.stuba.sk)

## 1 Formulation of the Nonlinear Satellite-Fixed Geodetic Boundary-Value Problem

The nonlinear geodetic boundary-value problem (BVP) has been of interest of many scientists and researchers. A uniqueness theorem for the fixed gravimetric BVP (FGBVP) was first given by Backus (1968). Later Koch and Pope (1972) presented a uniqueness proof for the nonlinear geodetic BVP. The free nonlinear BVP exactly solved by metric continuation was discussed by Grafarend and Niemeier (1971) as well as by Grafarend et al. (1989). Then Bjerhammar and Svensson (1983) used the general implicit function theorem and gave a solution of the existence and uniqueness problem in the nonlinear case. Expanding the nonlinear boundary condition into a Taylor series, based upon some reference potential field approximating the geopotential, was shown by Heck (1989). Sacerdote and Sansó (1989) further developed the idea used by Bjerhammar and Svensson for an iterative solution and they found explicit convergence conditions. They calculated the respective constant governing the convergence in the ideal case of a spherical boundary. Finally, we should mention authors Georgio Díaz, Jesús Díaz and Otero who showed the existence and uniqueness of a viscosity solution for the Backus problem (Díaz et al. 2006, 2011).

Let us consider the non-homogeneous elliptic equation of second order outside the Earth

$$\Delta W(\mathbf{x}) = 2\omega^2, \quad (1)$$

where  $W(\mathbf{x})$  is the actual gravity potential and  $\omega$  is the spin velocity of the Earth. The norm of gradient of the gravity potential  $W(\mathbf{x})$  is

$$|\nabla W(\mathbf{x})| = g(\mathbf{x}), \quad (2)$$

where  $g(\mathbf{x})$  denotes the magnitude of so-called total gravity vector. When  $g(\mathbf{x})$  is prescribed on the Earth's surface, Eq. (1) with BC (2) represents the nonlinear geodetic BVP for the actual gravity potential  $W(\mathbf{x})$ .

The actual gravity field can be expressed as a sum of the selected model field and the remainder of the actual field (Hofmann-Wellenhof and Moritz 2005), for corresponding potentials we can write

$$W(\mathbf{x}) = U(\mathbf{x}) + T(\mathbf{x}), \quad (3)$$

where  $U(\mathbf{x})$  is the normal gravity potential and  $T(\mathbf{x})$  the disturbing potential. When the model field is generated by a massive ellipsoid rotating with the Earth with the same

spin velocity  $\omega$ , its constant surface potential is equal to geopotential  $W_0$  and its mass is the same as the mass of the Earth, then the disturbing potential  $T(\mathbf{x})$  outside the Earth will satisfy the Laplace equation  $\Delta T(\mathbf{x}) = 0$ . It follows from the fact that  $T(\mathbf{x})$  does not have any centrifugal component since the centrifugal component of the Earth is the same as the centrifugal component of the chosen model.

Now let us consider the bounded domain  $\Omega$  depicted in Fig. 1. Such a domain is set in the external space above the Earth where the bottom surface  $\Gamma \subset \partial\Omega$ , where  $\partial\Omega$  denotes a boundary of  $\Omega$ , represents a part of the Earth's surface and the upper part of the boundary is at altitude of the chosen satellite mission. On the lower part of the boundary the nonlinear BC coming from (2) is given. On the upper spherical part of the domain as well as on the side boundaries, the Dirichlet-type BC (Eymard et al. 2001) obtained from satellite gravity missions is prescribed. That allows us to fix our solution to the satellite data. It is worth noting that another BC (Neumann or Newton BC) derived from satellite gravity missions suitable for the elliptic equation of second order might be taken into account as well.

Then our nonlinear satellite-fixed geodetic BVP (NSFGBVP) for the disturbing potential  $T(\mathbf{x})$  is formulated in the following form

$$\Delta T(\mathbf{x}) = 0 \quad \mathbf{x} \in \Omega, \quad (4)$$

$$|\nabla(T(\mathbf{x}) + U(\mathbf{x}))| = g(\mathbf{x}) \quad \mathbf{x} \in \Gamma, \quad (5)$$

$$T(\mathbf{x}) = T_{SAT}(\mathbf{x}) \quad \mathbf{x} \in \partial\Omega - \Gamma. \quad (6)$$

where  $T_{SAT}$  is the disturbing potential generated from a chosen satellite only model based on the spherical harmonics. It is worth to note that we are looking for a solution in a bounded domain  $\Omega$ , so we do not deal with its regularity at infinity. The influence of BC applied on side boundaries has been studied by Fašková et al. (2010).

In general, one can write the norm of the gradient of the gravity potential in the form

$$|\nabla W(\mathbf{x})| = \frac{\nabla W(\mathbf{x})}{|\nabla W(\mathbf{x})|} \cdot \nabla W(\mathbf{x}). \quad (7)$$

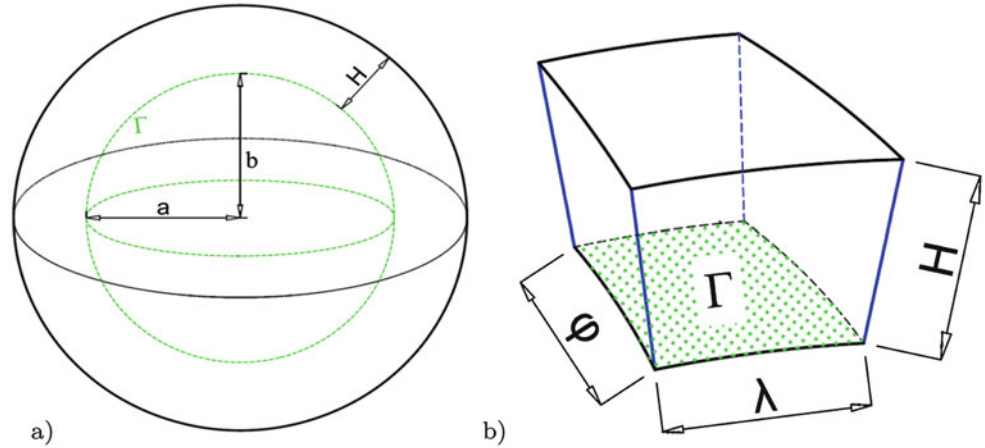
By inserting (7) in Eq. (5), we obtain

$$\frac{\nabla(T(\mathbf{x}) + U(\mathbf{x}))}{|\nabla(T(\mathbf{x}) + U(\mathbf{x}))|} \cdot \nabla(T(\mathbf{x}) + U(\mathbf{x})) = g(\mathbf{x}) \quad (8)$$

and if we denote

$$\mathbf{v}(\mathbf{x}) = \frac{\nabla(T(\mathbf{x}) + U(\mathbf{x}))}{|\nabla(T(\mathbf{x}) + U(\mathbf{x}))|}, \quad (9)$$

**Fig. 1** Sketch of the computational domain  $\Omega$  for (a) global numerical experiment, (b) local numerical experiment. The dotted boundary  $\Gamma$  represents the part of the Earth's surface,  $\varphi$  and  $\lambda$  denote latitude and longitude and  $H$  denotes the height above WGS84



we can rewrite the BC (5) as

$$\mathbf{v}(\mathbf{x}) \cdot \nabla(T(\mathbf{x})) = g(\mathbf{x}) - \mathbf{v}(\mathbf{x}) \cdot \nabla(U(\mathbf{x})) \quad \mathbf{x} \in \Gamma. \quad (10)$$

Since the unit vector  $\mathbf{v}(\mathbf{x})$ , defining the direction of the actual gravity vector, is unknown and depends on  $T(\mathbf{x})$ , BC (10) is still nonlinear, but its form allows to use an iterative approach for determining  $\mathbf{v}(\mathbf{x})$  and  $T(\mathbf{x})$  such that (4)–(5) is fulfilled. The iterative procedure for solving NSFGBVP will be defined as follows

$$\Delta T^{n+1}(\mathbf{x}) = 0 \quad \mathbf{x} \in \Omega, \quad (11)$$

$$\mathbf{v}^n(\mathbf{x}) \cdot \nabla(T^{n+1}(\mathbf{x})) = g(\mathbf{x}) - \mathbf{v}^n(\mathbf{x}) \cdot \nabla(U(\mathbf{x})) \quad \mathbf{x} \in \Gamma, \quad (12)$$

$$T^{n+1}(\mathbf{x}) = T_{SAT}(\mathbf{x}) \quad \mathbf{x} \in \partial\Omega - \Gamma, \quad (13)$$

for  $n = 0, 1, 2, \dots$ , where

$$\mathbf{v}^n(\mathbf{x}) = \frac{\nabla(T^n(\mathbf{x}) + U(\mathbf{x}))}{|\nabla(T^n(\mathbf{x}) + U(\mathbf{x}))|}, \quad (14)$$

and we start the iterations by choosing  $T^0(\mathbf{x}) = 0$ , i.e.  $W^0(\mathbf{x}) = U(\mathbf{x})$  and correspondingly for  $\mathbf{v}^0(\mathbf{x})$  we get

$$\mathbf{v}^0(\mathbf{x}) = \frac{\nabla(U(\mathbf{x}))}{|\nabla(U(\mathbf{x}))|} = \mathbf{s}(\mathbf{x}), \quad (15)$$

where  $\mathbf{s}(\mathbf{x})$  represents the direction of the normal gravity vector. One can see that in every iteration we solve the geodetic BVP for  $T^{n+1}(\mathbf{x})$  with prescribed oblique derivative vector  $\mathbf{v}^n(\mathbf{x})$ . In the first step we solve the linearized fixed gravimetric BVP (FGBVP) (Koch and Pope 1972; Holota 1997, 2005; Čunderlík et al. 2008; Fašková et al. 2010) with the oblique derivative given by

$$\mathbf{s}(\mathbf{x}) \cdot \nabla(T^1(\mathbf{x})) = g(\mathbf{x}) - \gamma(\mathbf{x}) = \delta g(\mathbf{x}), \quad (16)$$

where  $\gamma(\mathbf{x}) = |\nabla(U(\mathbf{x}))|$  and denotes a magnitude of the normal gravity vector and  $\delta g(\mathbf{x})$  denotes the gravity disturbance. In further iterations we improve the direction of the unit vector  $\mathbf{v}(\mathbf{x})$ . Such a process reduces the linearization error. Since we solve the problem iteratively, we need a stopping criterion. To that goal we use a difference of two successive iterations and stop the procedure, if in each point the inequality

$$|T^n(\mathbf{x}) - T^{n+1}(\mathbf{x})| < \varepsilon, \quad (17)$$

holds, where  $\varepsilon$  means a user-specified small real number. The last iteration represents our approximation of the disturbing potential  $T(\mathbf{x})$  and direction of gravity vector  $\mathbf{v}(\mathbf{x})$  in (4)–(5), and the sum  $T^{n+1}(\mathbf{x}) + U(\mathbf{x})$  represents the approximation of actual gravity potential  $W^{n+1}(\mathbf{x})$  in every point of the computational domain  $\Omega$ .

## 2 Numerical Solution of the Nonlinear Satellite-Fixed Geodetic Boundary-Value Problem

We can see that in each step of our iterative process (11)–(13) we deal with the oblique derivative BVP defined as

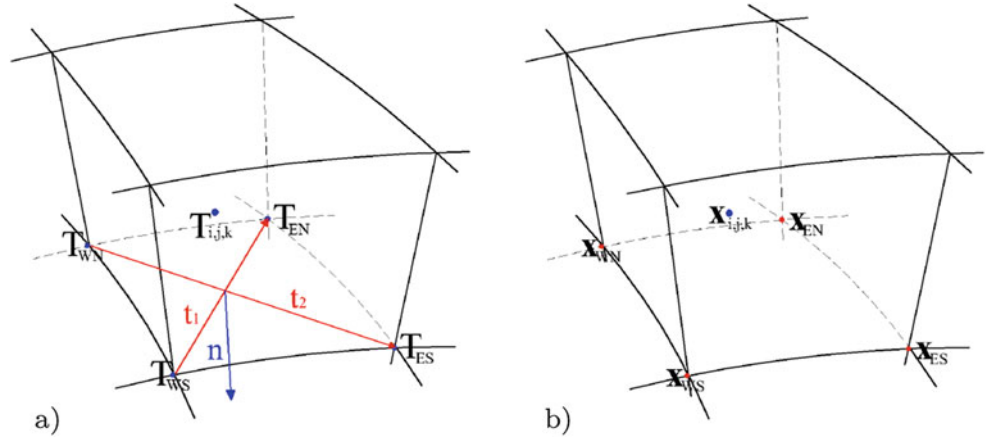
$$\Delta T(\mathbf{x}) = 0 \quad \mathbf{x} \in \Omega, \quad (18)$$

$$\mathbf{v}(\mathbf{x}) \cdot \nabla(T(\mathbf{x})) = g(\mathbf{x}) - \mathbf{v}(\mathbf{x}) \cdot \nabla(U(\mathbf{x})) = \alpha(\mathbf{x}), \quad \mathbf{x} \in \Gamma, \quad (19)$$

$$T(\mathbf{x}) = T_{SAT}(\mathbf{x}) \quad \mathbf{x} \in \partial\Omega - \Gamma. \quad (20)$$

To solve (18)–(20), we have chosen the finite volume method (FVM), (Eymard et al. 2001). In FVM we divide the computational domain  $\Omega$  into finite volumes  $p$ , multiply the Laplace equation by minus one and integrate the resulting

**Fig. 2** Brief illustration of the computational grid for an approximation of the oblique derivative. (a)  $T_{ijk}$  denotes the value of the disturbing potential in the center of volume.  $T_{WS}, T_{ES}, T_{EN}, T_{WN}$  are values of the disturbing potential in the vertices. Vectors  $\mathbf{t}_1$  and  $\mathbf{t}_2$  denote independent tangent vectors to  $\Gamma$  and  $\mathbf{n}$  the normal vector to  $\Gamma$ . (b)  $\mathbf{x}_{ijk}$  denotes position vector of the center of volume and  $\mathbf{x}_{WS}, \mathbf{x}_{ES}, \mathbf{x}_{EN}, \mathbf{x}_{WN}$  are values of the position vectors of the vertices



equation over each finite volume with a use of the divergence theorem that turns the volume integral into the surface integral,

$$-\int_p \Delta T \, dx dy dz = -\int_{\partial p} \nabla T \cdot \mathbf{n} \, d\sigma, \quad (21)$$

from where we get the *weak formulation* of the Eq. (18) in the finite volume  $p$

$$-\int_{\partial p} \frac{\partial T}{\partial n} d\sigma = 0. \quad (22)$$

Let  $q \in N(p)$  be a neighbour of the finite volume  $p$ , where  $N(p)$  denotes all neighbours of  $p$ . Let  $T_p$  and  $T_q$  be approximate values of  $T$  in  $p$  and  $q$ ,  $e_{pq}$  be a boundary of the finite volume  $p$  common with  $q$ ,  $\mathbf{n}_{pq}$  be its unit normal vector oriented from  $p$  to  $q$ ,  $m(e_{pq})$  is the area of  $e_{pq}$ . Let  $x_p$  and  $x_q$  be representative points of  $p$  and  $q$  respectively and  $d_{pq}$  their distance. If we approximate the normal derivative along the boundary of the finite volume  $p$  by

$$\frac{\partial T}{\partial n_{pq}} \approx \frac{T_q - T_p}{d_{pq}}, \quad (23)$$

we obtain from (22) and (23) the following equation for every finite volume  $p$

$$\sum_{q \in N(p)} \frac{m(e_{pq})}{d_{pq}} (T_p - T_q) = 0, \quad (24)$$

which forms together the linear system of algebraic equations. The term  $\frac{m(e_{pq})}{d_{pq}}$  defined on sides of the finite volume  $p$  is referred to as the transmissivity coefficient (Eymard et al.

2001). Then we define indices  $i = 1, \dots, n_1$ ,  $j = 1, \dots, n_2$  and  $k = 1, \dots, n_3$  in the direction of the longitude  $\lambda$ , latitude  $\varphi$  and height  $h$ , where  $n_1$ ,  $n_2$  and  $n_3$  denote the numbers of discretisation intervals in zonal, meridional and height's direction, respectively. In this way we obtain the linear system of equations that can be written in the form

$$P_{i,j,k} T_{i,j,k} - W_{i,j,k} T_{i-1,j,k} - E_{i,j,k} T_{i+1,j,k} - N_{i,j,k} T_{i,j+1,k} - S_{i,j,k} T_{i,j-1,k} - U_{i,j,k} T_{i,j,k+1} - D_{i,j,k} T_{i,j,k-1} = 0, \quad (25)$$

where  $P_{i,j,k}$ ,  $W_{i,j,k}$ ,  $E_{i,j,k}$ ,  $N_{i,j,k}$ ,  $S_{i,j,k}$ ,  $U_{i,j,k}$  and  $D_{i,j,k}$  are transmissivity coefficients and their derivation can be found in Macák et al. (2012).

The system (25) must be accompanied by the boundary conditions. In case of the Dirichlet BC, we prescribe the value of  $T_q$  on the boundary, while in case of the oblique derivative BC, a special treatment is needed. For the bottom boundary, when  $k = 1$ , we add new finite volumes  $p$  signed by index  $k = 0$ . Then we split the gradient of  $T(\mathbf{x})$  in (19) into one normal and two tangential directions

$$\nabla T = (\nabla T \cdot \mathbf{n}) \mathbf{n} + (\nabla T \cdot \mathbf{t}_1) \mathbf{t}_1 + (\nabla T \cdot \mathbf{t}_2) \mathbf{t}_2 = \frac{\partial T}{\partial n} \mathbf{n} + \frac{\partial T}{\partial t_1} \mathbf{t}_1 + \frac{\partial T}{\partial t_2} \mathbf{t}_2, \quad (26)$$

where  $\mathbf{n}$  is the unit normal vector and  $\mathbf{t}_1$ ,  $\mathbf{t}_2$  are linearly independent unit tangent vectors to  $\Gamma \subset \partial\Omega \subset R^3$ . So the BC (19) is transformed into the form

$$\frac{\partial T}{\partial n} (\mathbf{n} \cdot \mathbf{v}) + \frac{\partial T}{\partial t_1} (\mathbf{t}_1 \cdot \mathbf{v}) + \frac{\partial T}{\partial t_2} (\mathbf{t}_2 \cdot \mathbf{v}) = \alpha. \quad (27)$$

Then we approximate the normal and tangential derivatives according to notations depicted in Fig. 2

$$\frac{\partial T}{\partial n} \approx \frac{T_D - T_P}{|\mathbf{x}_D - \mathbf{x}_P|}, \quad \frac{\partial T}{\partial t_1} \approx \frac{T_{EN} - T_{WS}}{|\mathbf{x}_{EN} - \mathbf{x}_{WS}|},$$

$$\frac{\partial T}{\partial t_2} \approx \frac{T_{WN} - T_{ES}}{|\mathbf{x}_{WN} - \mathbf{x}_{ES}|},$$

where we have denoted values  $T_{i,j,k-1}$  and  $T_{i,j,k}$  by  $T_D$  and  $T_P$ , respectively. Values  $T_{EN}$ ,  $T_{WS}$ ,  $T_{WN}$ ,  $T_{ES}$  are obtained as follows

$$T_{WN} = \frac{T_{i,j,k} + T_{i,j-1,k} + T_{i,j,k-1} + T_{i,j-1,k-1} + T_{i-1,j,k} + T_{i-1,j-1,k} + T_{i-1,j,k-1} + T_{i-1,j-1,k-1}}{8},$$

$$T_{EN} = \frac{T_{i,j,k} + T_{i,j-1,k} + T_{i,j,k+1} + T_{i,j-1,k+1} + T_{i-1,j,k} + T_{i-1,j-1,k} + T_{i-1,j,k+1} + T_{i-1,j-1,k+1}}{8},$$

$$T_{WS} = \frac{T_{i,j,k} + T_{i,j+1,k} + T_{i,j,k-1} + T_{i,j+1,k-1} + T_{i-1,j,k} + T_{i-1,j+1,k} + T_{i-1,j,k-1} + T_{i-1,j+1,k-1}}{8},$$

$$T_{ES} = \frac{T_{i,j,k} + T_{i,j+1,k} + T_{i,j,k+1} + T_{i,j+1,k+1} + T_{i-1,j,k} + T_{i-1,j+1,k} + T_{i-1,j,k+1} + T_{i-1,j+1,k+1}}{8},$$

and  $\mathbf{x}_D, \mathbf{x}_P, \mathbf{x}_{EN}, \mathbf{x}_{WS}, \mathbf{x}_{WN}, \mathbf{x}_{ES}$  are their corresponding position vectors, see Fig. 2. More details can be found in Macák et al. (2012). Then the final discrete form of the oblique derivative BC is given by

$$\mathbf{v} \cdot \nabla(T(\mathbf{x})) \approx \frac{T_D - T_P}{|\mathbf{x}_D - \mathbf{x}_P|} (\mathbf{n} \cdot \mathbf{v}) + \frac{T_{EN} - T_{WS}}{|\mathbf{x}_{EN} - \mathbf{x}_{WS}|} (\mathbf{t}_1 \cdot \mathbf{v}) + \frac{T_{WN} - T_{ES}}{|\mathbf{x}_{WN} - \mathbf{x}_{ES}|} (\mathbf{t}_2 \cdot \mathbf{v}) = \alpha.$$

### 3 Numerical Experiments

The local numerical experiment was performed in the domain above Slovakia bounded by  $\varphi \in (47.0^\circ, 50.5^\circ)$  and  $\lambda \in (16.0^\circ, 23.0^\circ)$ . The bottom boundary was created using heights generated from SRTM30 PLUS (Becker et al. 2009) and the upper boundary was at the height of 240 km above WGS84, corresponding to an average altitude of the satellite orbit. The number of finite volumes was 1,000 in height, 630 in meridional and 840 in zonal directions, i.e. the resolution with respect to latitude and longitude was  $30'' \times 20''$ . We started our computations by solving the linearized FGBVP where the surface gravity disturbances were applied on the bottom boundary  $\Gamma$ . They were generated from an available dataset of terrestrial gravity data in Slovakia (Grand et al. 2001) while ellipsoidal heights of gravimetric measurements were computed from levelling heights using EGM2008 (Pavlis et al. 2012). On the upper and side boundaries, the disturbing potential generated from the GOCO03s satellite-only model (Mayer-Gürr et al. 2012) was prescribed. Computations were performed on 30 processors

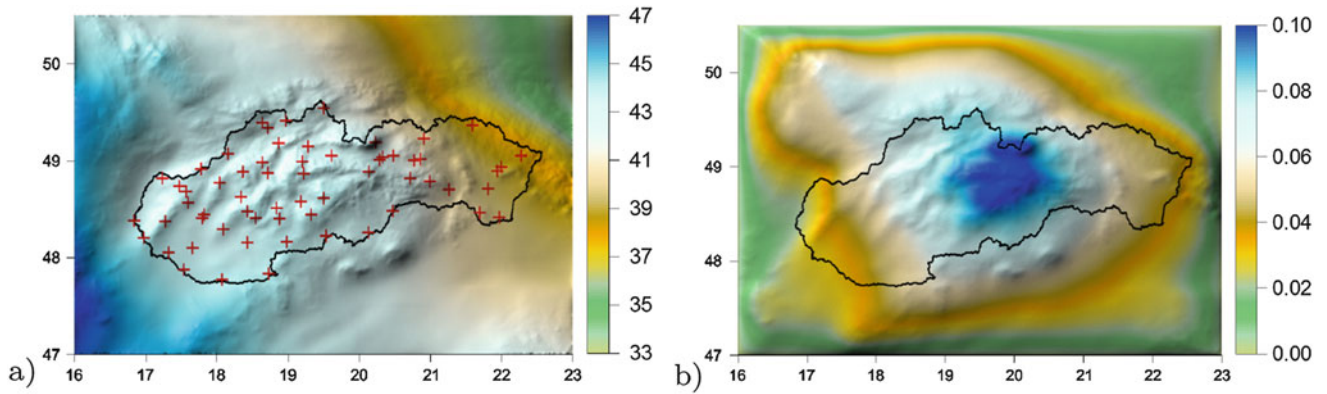
**Table 1** Statistics of residuals [m] between our NSFGBVP solution and quasigeoidal heights obtained by GPS/levelling at 61 points in the area of Slovakia

	1st iter.	5th iter.	8th iter.	10th iter.	EGM2008
Min. value	0.151	0.209	0.229	0.248	0.301
Mean value	0.284	0.325	0.348	0.352	0.437
Max. value	0.422	0.459	0.476	0.493	0.584
St. deviation	0.055	0.049	0.047	0.046	0.043

using 78 GB of distributed memory taking approximately 5 h of total CPU time per processor. To reach the prescribed stopping criterium  $\varepsilon = 10^{-3} [\text{m}^2 \text{s}^{-2}]$ , ten iterations were needed. Results are presented in Table 1 and Fig. 3. One can observe an improvement in the standard deviation for subsequent iterations in solving NSFGBVP (Table 1) as well as the convergence to EGM2008. The differences between the 10th and 1st iteration, which represent the numerically obtained linearization error, reach up to 10 cm.

The global numerical experiment dealt with the high-resolution global gravity field modelling in the computational domain  $\Omega$  bounded by the bottom boundary approximating the real Earth's surface created by using heights generated from SRTM30 PLUS and by a surface at height of 240 km above WGS84 corresponding to the average altitude of satellite orbit. The number of divisions was  $4320 \times 2160 \times 600$  leading to the resolution  $5' \times 5' \times 400 \text{ m}$ . Again we start with the linearized FGBVP consisting of gravity disturbances interpolated from the DTU10-GRAV gravity field model (Andersen 2010) and applied on the bottom boundary. On the upper boundary the disturbing potential generated from GOCO03s was prescribed. The stopping criterium was  $\epsilon = 10^{-3} [\text{m}^2 \text{s}^{-2}]$  and again, ten iterations were needed. The FVM solutions obtained in each iteration are compared with EGM2008. Statistical characteristics of residuals are presented in Table 2. Figure 4 depicts differences between the 10th and 1st iteration. They represent the numerically obtained linearization error in the linearized FGBVP. One can observe that our iteration approach improves solution mainly in areas of high mountains (e.g. in Himalaya region they reach 20 cm) as well as in areas along the ocean trenches (varying from  $-2.5$  to  $2.5$  cm).

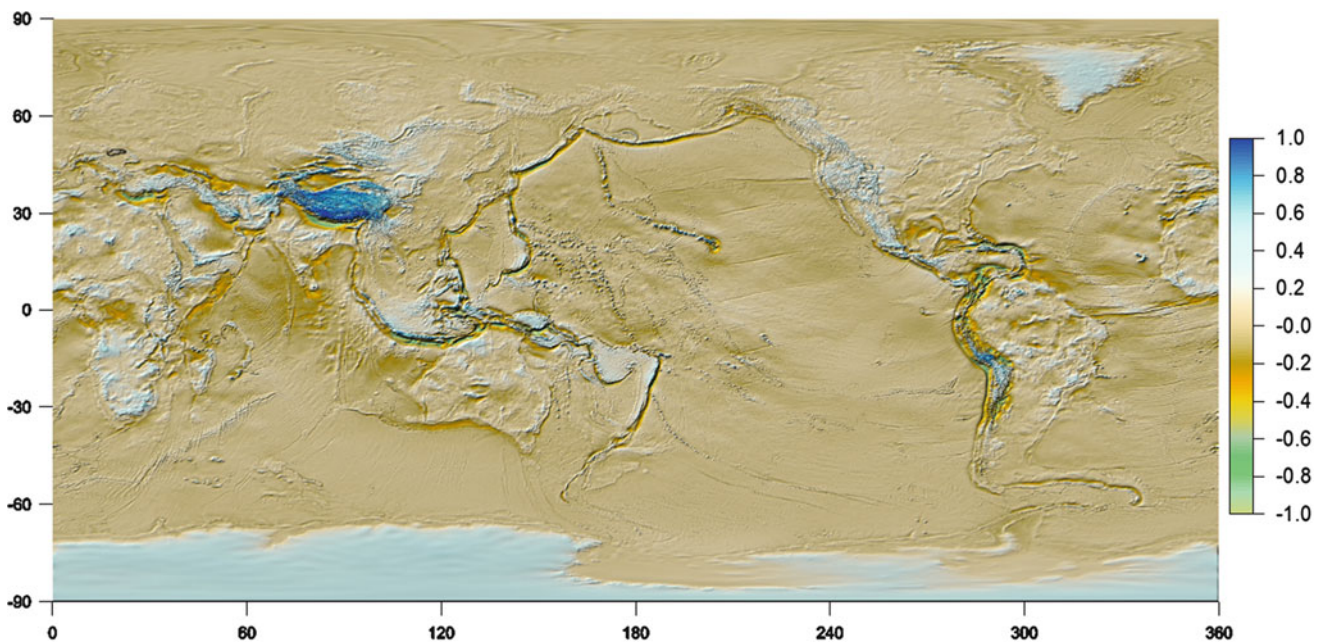




**Fig. 3** (a) Quasigeoidal heights  $\zeta$  [m] obtained by solving the NSFGBVP, (b) Differences in  $\zeta$  [m] between 10th and 1st iteration. Red crosses denote the distribution of 61 GPS/leveling points

**Table 2** Statistics of residuals [ $\text{m}^2 \text{s}^{-2}$ ] between the disturbing potential obtained by solving NSFGBVP and the disturbing potential generated from EGM2008 in the global experiment

Iter.	Min. value		Mean value		Max. value		St. dev.	
	1st	10th	1st	10th	1st	10th	1st	10th
Total	-2.150	-1.985	0.004	0.001	6.143	4.158	0.501	0.419
Sea	-0.705	-0.632	-0.021	-0.011	1.131	1.019	0.206	0.199
Land	-2.150	-1.985	0.035	0.029	6.143	4.158	0.855	0.768



**Fig. 4** Differences in  $T$  [ $\text{m}^2 \text{s}^{-2}$ ] between 10th and 1st iteration, representing the numerically obtained linearization error

## 4 Summary and Conclusions

We have presented an iterative approach to solving the nonlinear satellite-fixed geodetic boundary-value problem (NSFGBVP) defined in this paper. The NSFGBVP has been solved by the finite volume method, where the direction of

the actual gravity vector as well as the disturbing potential are updated in each iteration. In the first iteration, the linearized FGBVP is solved together with the oblique derivative problem. Next iterations treat its numerically obtained linearization error. The obtained numerical results show that the error of the linearization can exceed several centimeters, mainly in high mountainous areas and along ocean trenches.

This indicates that for precise gravity field modeling it is necessary to deal with the nonlinear geodetic BVPs avoiding the linearization error. Presented numerical experiments show that the proposed iterative approach converges while the study of its convergence from theoretical point of view will be a task of our future research.

## References

- Andersen OB (2010) The DTU10 Gravity field and mean sea surface. In: Second international symposium of the gravity field of the Earth (IGFS2), Fairbanks, Alaska
- Backus GE (1968) Application of a non-linear boundary-value problem for Laplace's equation to gravity and geomagnetic intensity surveys. *Q J Mech Appl Math* 2:195–221
- Becker JJ et al (2009) Global bathymetry and elevation data at 30 arc seconds resolution: SRTM30 PLUS. *Mar Geod* 4:355–371
- Bjerhammar A, Svensson L (1983) On the geodetic boundary value problem for a fixed boundary surface – a satellite approach. *Bull Geod* 57(1–4):382–393
- Čunderlík R, Mikula K, Mojzeš M (2008) Numerical solution of the linearized fixed gravimetric boundary-value problem. *J Geod* 82(1):15–29
- Díaz G, Díaz JI, Otero J (2006) On an oblique boundary value problem related to the Backus problem in geodesy. *Nonlinear Anal Real World Appl* 7:147–166
- Díaz G, Díaz JI, Otero J (2011) Construction of the maximal solution of Backus' problem in geodesy and geomagnetism. *Stud Geophys Geod* 55(3):415–440
- Eymard R, Gallouet T, Herbin R (2001) Finite volume approximation of elliptic problems and convergence of an approximate gradient. *Appl Num Math* 37(1–2):31–53
- Fašková Z, Čunderlík R, Mikula K (2010) Finite element method for solving geodetic boundary value problems. *J Geod* 84:135–144
- Grafarend E, Niemeier W (1971) The free nonlinear boundary value problem of physical geodesy. *Bull Geod* 101:243–261
- Grafarend E (1989) The geoid and the gravimetric boundary value problem, Report No 18. The Royal Institute of Technology (Dep of Geod), Stockholm
- Grand T, Šefara J, Pašteka R, Bielik M, Daniel S (2001) Atlas of geophysical maps and profiles. State geological institute, Bratislava, MS Geofond (in Slovak)
- Heck B (1989) On the non-linear geodetic boundary value problem for a fixed boundary surface. *Bull Geod* 63(1):57–67
- Hofmann-Wellenhof B, Moritz H (2005) *Physical geodesy*. Springer, Wien New York
- Holota P (1997) Coerciveness of the linear gravimetric boundary-value problem and a geometrical interpretation. *J Geod* 71(10):640–651
- Holota P (2005) Neumann's boundary-value problem in studies on earth gravity field: weak solution. 50 years of Research Institute of Geodesy, Topography and Cartography, Prague, pp 34, 49–69
- Koch KR, Pope AJ (1972) Uniqueness and existence for the geodetic boundary value problem using the known surface of the earth. *Bull Geod* 46:467–476
- Macák M, Mikula K, Minarechová Z (2012) Solving the oblique derivative boundary-value problem by the finite volume method. In: *ALGORITMY 2012*, 19th conference on scientific computing, pp 75–84
- Mayer-Gürr T et al (2012) The new combined satellite only model GOCO03s. Presented at the GGHS-2012 in Venice, Italy
- Pavlis NK, Holmes SA, Kenyon SC, Factor JK (2012) The development and evaluation of the earth gravitational model 2008 (EGM2008). *J Geophys Res* 117:B04406. Doi:10.1029/2011JB008916
- Sacerdote F, Sansó F (1989) On the analysis of the fixed-boundary gravimetric boundary-value problem. In: Sacerdote F, Sansó F (eds) *Proceedings of the 2nd Hotine-Marussi symposium on mathematical geodesy*, Pisa, Politecnico di Milano, pp 507–516

---

**Part VI**

**Computational Geodesy**

---

# An OpenCL Implementation of Ellipsoidal Harmonics

Otakar Nesvadba and Petr Holota

---

## Abstract

The technology progress today makes it possible to treat most of the problems of physical geodesy by means of numerical arrangements hardly imaginable earlier. Nevertheless, considering an evaluation of spheroidal (spherical and ellipsoidal) harmonic functions in our typical tasks, we still observe a huge performance gap between our demands and capabilities of common CPUs. Methods used for calculating associated Legendre functions are mostly recursive and thus sequential. Therefore, it is challenging, but feasible, to arrange the processing of Legendre functions in a way that reduces memory utilisation and admits massive parallelism. Following this aim, we developed a streaming-parallel algorithm for computing oblate spheroidal harmonic functions and their derivatives. The algorithm is free of assumptions concerning the function arguments, maximal degree/order or number of computation points and can be utilised on any data type, like a vector or scalar float, double or even integer numbers. Besides, it solves floating-point issues in the numerical treatment of Legendre functions. We demonstrate its Open Computing Language (OpenCL) implementation on a general-purpose graphics processing unit (GPGPU), which is ideal for its inexpensive computational power of some TFlops. Added performance benchmarks lead to the conclusion that our implementation on a single GPGPU device substantially outperforms recent multi-core CPUs, free of any precision penalty. Furthermore, thanks to the OpenCL standard, we can benefit from an excellent portability and scalability over heterogeneous parallel platforms. Let us note finally, that the topic presented is a matter of importance in many other application fields, not only in physical geodesy.

---

## Keywords

Associated Legendre functions • Clenshaw summation • Heterogeneous parallel computing • Oblate spheroidal harmonics • OpenCL

---

O. Nesvadba (✉)  
Land Survey Office, Pod Sídlištěm 9, 182 00 Prague, Czech Republic  
e-mail: [nesvadba@sky.cz](mailto:nesvadba@sky.cz)

P. Holota  
Research Institute of Geodesy, Topography and Cartography,  
Prague-East, 250 66 Zdice 98, Czech Republic  
e-mail: [petr.holota@pecny.cz](mailto:petr.holota@pecny.cz)

---

## 1 Introduction

In the past decade of the computer industry, we can observe a departure from serial (single-core) execution in favour of parallel processors. Such a transition, that affects almost all categories of computers, recently resulted in the emergence of *heterogeneous computing*. Heterogeneous computers can be characterised as a combination of a complex CPU (central processing unit) with a computational accelerator, usually

GPGPU (general-purpose graphics processing unit), recognised as a simple but highly-parallel processor.

Our primary motivation is to demonstrate an evaluation of truncated series of spherical and spheroidal (ellipsoidal) harmonics as an ideal application of heterogeneous computing. We believe, this application can be very useful for a broad geoscience community.

In this paper we treat the spheroidal harmonics from a numerical point of view, focusing on the optimisation of recursive algorithms and the adoption of the streaming processing schema in particular. In Sect. 3 we demonstrate an OpenCL implementation of the varying degree recursive formula algorithm for ALFs (associated Legendre functions) of the first and second kind, gradually advancing our accomplishments. Consequently, in Sect. 4, we introduce the most efficient implementation for spheroidal harmonic synthesis based on the Clenshaw summation algorithm. Finally, in Sect. 5, we make an insight into its overall performance on GPGPU and CPU devices.

## 1.1 Oblate Spheroidal Coordinates

Let us start with oblate spheroidal coordinates  $(u, \beta, \lambda)$  of parameter  $E \geq 0$ , well-known in geodesy as *ellipsoidal coordinates* (polar radius  $u$ , reduced latitude  $\beta$  and geodetic longitude  $\lambda$ ), cf. Torge (2001, Sect. 4.2.2). Nevertheless, for practical reasons, in the following we will use a slightly different coordinate notation  $(\rho, \tau, \lambda)$ , where

$$\tau = \sin \beta, \quad \beta \in \left(-\frac{\pi}{2}, \frac{\pi}{2}\right) \Rightarrow \tau \in \langle -1, 1 \rangle \quad (1)$$

and  $\rho$  is related to  $u$  by the following equation

$$\rho = \frac{a}{\sqrt{E^2 + u^2}}, \quad u \in \langle 0, \infty \rangle \Rightarrow \rho \in \langle 0, e^{-1} \rangle. \quad (2)$$

Here we have introduced a particular oblate spheroid of a linear eccentricity  $E$  and polar radius  $b > 0$ .<sup>1</sup> Its equatorial radius is therefore  $a = \sqrt{E^2 + b^2}$  and eccentricity  $e = \frac{E}{a}$ .

Relation of  $\rho, \tau, \lambda$  to Cartesian coordinates is given by

$$(x_1, x_2, x_3) = \frac{a}{\rho} (\sqrt{1 - \tau^2} \cos \lambda, \sqrt{1 - \tau^2} \sin \lambda, \tau \sqrt{1 - \rho^2 e^2}).$$

Note explicitly that we admit a special case  $E = 0$ , where spheroidal coordinates become spherical, so that  $\rho = \frac{a}{u}$ .

<sup>1</sup>In this study we are considering GRS 80 (or WGS 84) parameters  $E = 521\,854.0097$  m,  $a = 6378137$  m,  $b^2 = a^2 - E^2$  as in Moritz (1984).

## 1.2 Oblate Spheroidal Harmonic Function

Let us define a spheroidal harmonic function

$$Y_{n,m}(\rho, \tau, \lambda) = Q_{n,m}\left(i \frac{u}{E}\right) P_{n,m}(\tau) \exp(im\lambda) \quad (3)$$

of degree  $n \in \mathbb{N}_0$  and order  $m \in \mathbb{N}_0$ ,  $m \leq n$ , where  $\exp$  is the exponential function,  $P_{n,m}$ ,  $Q_{n,m}$  are associated Legendre functions (ALFs) of the first and second kind and  $i = \sqrt{-1}$ . Recall that from Eq. (2) it follows that  $\frac{u}{E} = \sqrt{e^{-2}\rho^{-2} - 1}$ .

*ALFs of the First Kind* The  $P_{n,m}$  values can be obtained from

$$P_{n,m}(\tau) = \frac{2^{-n}(2n)!}{n!(n-m)!} \frac{(1-\tau^2)^{\frac{m}{2}}}{\tau^{m-n}} {}_2F_1\left(\frac{m-n}{2}, \frac{m-n+1}{2}; \frac{1-2n}{2}; \tau^{-2}\right), \quad (4)$$

where  ${}_2F_1$  represents Gauss' hypergeometric function, see e.g. Bateman and Erdélyi (1953, Chap. II) or Abramowitz and Stegun (1964, Sect. 15). In the special case  $n = m$  Eq. (4) becomes

$$P_{m,m}(\tau) = \frac{(2m)!}{2^m m!} (1 - \tau^2)^{\frac{m}{2}}. \quad (5)$$

*ALFs of the Second Kind* Similarly, referring to Hobson (1931, Sect. V), Bateman and Erdélyi (1953, Sect. 3.2) or to Holota (2001) we can found that

$$Q_{n,m}(z) = (-1)^m 2^n \frac{n!(n+m)!}{(2n+1)!} (z^2 - 1)^{-\frac{n+1}{2}} \times \\ \times {}_2F_1\left(\frac{n+m+1}{2}, \frac{n-m+1}{2}; n + \frac{3}{2}; \frac{1}{1-z^2}\right). \quad (6)$$

Note that  $n + \frac{3}{2} - \frac{n+m+1}{2} - \frac{n-m+1}{2} = \frac{1}{2} > 0$ , so the series

$${}_2F_1\left(\frac{n+m+1}{2}, \frac{n-m+1}{2}; n + \frac{3}{2}; \frac{1}{1-z^2}\right) = \sum_{j=0}^{\infty} k_j \left(\frac{1}{1-z^2}\right)^j, \quad (7)$$

where  $k_j = k_{j-1} \frac{(n+m-1+2j)(n-m-1+2j)}{4j(n+\frac{1}{2}+j)}$ ,  $k_0 = 1$ , converges absolutely for any  $|\frac{1}{1-z^2}| < 1$ .

## 1.3 Normalised Spheroidal Harmonic Function

Just for technical reasons, let us replace  $P_{n,m}(\tau)$  by

$$p_{n,m}(\tau, \rho) = \rho^{n+1} H_{n,m} P_{n,m}(\tau), \quad (8)$$

where the constants  $H_{n,m} \in \mathbb{R} - \{0\}$  represent particular normalisation factors for the given  $n$  and  $m$ .<sup>2</sup>

<sup>2</sup>Usually we suppose  $H_{n,m} = \sqrt{(2 - \delta_{0,m})(2n+1)} \frac{(n-m)!}{(n+m)!}$ .

Similarly, instead of  $Q_{n,m}(i\frac{u}{E})$  we will use

$$q_{n,m}(\rho) = \rho^{-n-1} \frac{Q_{n,m}(i\frac{u}{E})}{Q_{n,m}(i\frac{b}{E})} \Rightarrow q_{n,m} \in (0, \infty). \quad (9)$$

Thus, in our notation we introduce a *normalised oblate spheroidal harmonic function*

$$\bar{Y}_{n,m}(\rho, \tau, \lambda) = q_{n,m}(\rho) p_{n,m}(\tau, \rho) \exp(im\lambda), \quad (10)$$

which differs from  $Y_{n,m}(\rho, \tau, \lambda)$  in Eq. (3) just by the constant factor  $\frac{H_{n,m}}{Q_{n,m}(i\frac{b}{E})}$ .

**Remark** In the special case  $E = 0$  we have  $q_{n,m}(\rho) = 1$  for any  $n, m$  and  $\rho$ , hence a normalised spheroidal harmonic function becomes a *normalised spherical harmonic function*

$$\bar{Y}_{n,m}(\rho, \tau, \lambda) = p_{n,m}(\tau, \rho) \exp(im\lambda). \quad (11)$$

## 2 Recursive Formula for Legendre Functions

### 2.1 Associated Legendre Functions of the First Kind

When treating ALFs of the first kind numerically, the most convenient way is to employ LFRF (varying-degree recursive formula for ALFs), see e.g. Hobson (1931), Abramowitz and Stegun (1964, Sect. 8.5), Holmes and Featherstone (2002). Therefore, in the  $p_{n,m}(\tau, \rho)$  notation one can write

$$p_{n,m}(\tau, \rho) = \rho \tau \frac{2n-1}{n-m} \frac{H_{n,m}}{H_{n-1,m}} p_{n-1,m}(\tau, \rho) - \rho^2 \frac{n+m-1}{n-m} \frac{H_{n,m}}{H_{n-2,m}} p_{n-2,m}(\tau, \rho). \quad (12)$$

The LFRF algorithm given by Eq. (12) with the seed values

$$p_{m-1,m}(\tau, \rho) = 0, \quad (13)$$

$$p_{m,m}(\tau, \rho) = \rho^{m+1} (1 - \tau^2)^{\frac{m}{2}} H_{m,m} (2m-1)!!, \quad (14)$$

enables us to get any  $p_{n,m}$  for  $n \geq m$ .

Although a simple recursive formula is available for calculation of Eq. (14), in this study we would prefer the explicit expression based on Eq. (5), since it allows parallel LFRF processing independent of  $m$ .

### 2.2 Associated Legendre Functions of the Second Kind

The same recurrent relations as for  $P_{n,m}$  are also valid for  $Q_{n,m}$ , cf. Hobson (1931) and Abramowitz and Stegun (1964). Nevertheless, as it has been already discussed in Sona (1995), the LFRF is useless in this case, because of exponentially growing errors. On the other hand, from the same discussion it is clear that the reverse evaluation of the LFRF should work. This ‘‘backward’’ LFRF approach has been mentioned for instance in Abramowitz and Stegun (1964, note in 8.15) and further discussed in Gil and Segura (1998), Nesvadba, 2009, *Numerical problems in evaluating high degree and order associated Legendre functions*, EGU, Vienna, Nesvadba (2011) or Fukushima (2013).

Following the references above we can write

$$q_{n-2,m}(\rho) = \frac{\sqrt{1-\rho^2 e^2}}{\sqrt{1-e^2}} \frac{2n-1}{2n-1+h_{n,m}} q_{n-1,m}(\rho) + \rho^2 \frac{h_{n,m}}{2n-1+h_{n,m}} q_{n,m}(\rho), \quad (15)$$

where  $h_{n,m}$  constants are defined as

$$h_{n,m} = i \frac{E}{b} (n-m) \frac{Q_{n,m}(i\frac{b}{E})}{Q_{n-1,m}(i\frac{b}{E})} \Rightarrow h_{n,m} \in (0, \infty). \quad (16)$$

By applying the LFRF to  $Q_{n,m}(i\frac{b}{E})$  in Eq. (16), one can show that the following recursive formula holds for any  $n \geq m$

$$h_{n,m} = \frac{e^2}{1-e^2} \frac{(n+m)(n-m)}{2n+1+h_{n+1,m}}. \quad (17)$$

Therefore, Eqs (15) and (17) represent a recursive computation scheme for any  $q_{n,m}(\rho)$ ,  $n \geq m$ . Recall that the equations are valid also in the case of  $e = 0$ , where  $h_{n,m} = 0$  and  $q_{n,m}(\rho) = 1$  for any  $n, m$  and  $\rho$ .

Moreover, if  $n = m - 1$  is allowed, Eq. (15) can serve as a ‘‘check-out’’ formula, since  $q_{m-1,m}(\rho) = 1$  holds for any  $m, \rho$  and  $e$ , cf. Nesvadba (2011).

*Hypergeometric Series for  $q_{n,m}$*  In order to get the seed values of Eq. (15), we have to return to Eq. (9) and after some algebra involving Eq. (6) we obtain

$$q_{n,m}(\rho) = \frac{{}_2F_1(\frac{n+m+1}{2}, \frac{n-m+1}{2}; n + \frac{3}{2}; e^2 \rho^2)}{{}_2F_1(\frac{n+m+1}{2}, \frac{n-m+1}{2}; n + \frac{3}{2}; e^2)}. \quad (18)$$

Similarly, Eq. (16) transforms to

$$h_{n,m} = \frac{e^2(n+m)(n-m)}{(1-e^2)(2n+1)} \frac{{}_2F_1\left(\frac{n+m+1}{2}, \frac{n-m+1}{2}; \frac{2n+3}{2}; e^2\right)}{{}_2F_1\left(\frac{n+m+1}{2}, \frac{n-m+1}{2}; \frac{2n+1}{2}; e^2\right)}. \quad (19)$$

The numerical evaluation of Eqs (18) and (19) follows directly from Eq. (7). In the summation it is sufficient just to replace the infinity with some moderate  $j_{\max}$ .<sup>3</sup>

### 3 LFRF Implementation in OpenCL

Among many programming frameworks which are suitable for LFRF implementation we have chosen OpenCL (Open Computing Language). Nevertheless, in the following we will abstract from OpenCL technical details as much as possible, in order to allow an easy re-implementation in other frameworks like CUDA (compute unified device architecture) or C++ AMP (accelerated massive parallelism).

#### 3.1 Open Computing Language

OpenCL, Khronos Group (2012), is recognised as the industry standard for heterogeneous parallel programming and computations. Since the standard has been adopted by most of the hardware vendors, the OpenCL framework is truly cross-platform. OpenCL hardware abstraction layer is able to utilise any OpenCL-enabled hardware (*device*), including a broad range of CPUs, GPGPUs or even DSPs (digital signal processors) or FPGAs (field-programmable gate arrays).

*OpenCL Execution Model* OpenCL relies on a host–device execution model. It means, there are always two kinds of programs in OpenCL. A host-side program and device programs called *compute kernels*. The compute kernels are written in a device-independent programming language based on ISO/IEC 9899:1999 standard augmented with additional keywords, data types and functions to support memory address spaces, SIMD (single instruction multiple data) vectors, atomic operations, synchronisation primitives, etc.

The OpenCL kernels are naturally parallel, i.e. they usually run in many simultaneous threads (*work items*) on a device. A global worksize, i.e. the total number of work items, is organised into smaller units, the so-called *work-groups*, which allow utilisation of a *local memory* and easy work-item synchronisation. We will refer the reader to the Khronos Group (2012) OpenCL documentation for further information since the details are often complicated.

The host layer of OpenCL provides an application interface (API) for accessing devices and coordination of the

computation across devices. Within the host program one can compile OpenCL kernels for particular OpenCL device, pass the machine code and data to the device and execute it (asynchronously with the host program) by means of device execution queues. An implementation of the host program is not a subject of this text.

#### 3.2 LFRF Algorithm for $p_{n,m}$

The LFRF algorithm given by Eq. (12) seems to be very efficient. If the coefficients  $\frac{2n-1}{n-m} \frac{H_{n,m}}{H_{n-1,m}}$  and  $\frac{n+m-1}{n-m} \frac{H_{n,m}}{H_{n-2,m}}$  are precalculated, one elementary LFRF step – defined as an evaluation of a single  $p_{n,m}$  in Eq. (12) – costs just 5 operations (four multiplications and one addition). Nevertheless, at the same time it might be strongly inefficient in a memory access,<sup>4</sup> because each LFRF step costs two memory reads and possibly one memory write (of the result).

In order to increase the memory efficiency of the LFRF algorithm it is vital to rearrange it as follows. Let us define coefficients  $I_{n,m}$

$$I_{m,m} = H_{m,m} \prod_{l=1}^m (2l-1), \quad I_{n,m} = \frac{1}{n-m} \frac{H_{n,m}}{H_{n-1,m}}, \quad n > m \quad (20)$$

and the integers  $\alpha_n, \beta_{n,m}$  (please, do not confuse it with  $\beta$ )

$$\alpha_n = 2n + 1, \quad \beta_{n,m} = (n+m)(n-m). \quad (21)$$

In the following relations, for simplicity, we omit the second index  $m$  and the arguments  $\tau, \rho$ . From Eqs (12–14) we get

$$p_m = \rho^{m+1} (1 - \tau^2)^{\frac{m}{2}} I_m, \quad p_{m-1} = 0, \quad (22)$$

$$p_n = \rho \tau \alpha_{n-1} I_n p_{n-1} - \rho^2 \beta_{n-1} I_n I_{n-1} p_{n-2}, \quad (23)$$

$$\alpha_n = \alpha_{n-1} + 2, \quad \beta_n = \beta_{n-1} + \alpha_{n-1}. \quad (24)$$

The algorithm represented by Eqs (22)–(24) utilises less memory at the expense of a slightly increased computation overhead. One LFRF step now involves 6 multiplications and 3 additions, but only one set of the coefficients  $I_{n,m}$  is required.

*Enhanced LFRF Algorithm with Extended Exponent* Since we are working mostly with ISO/IEC 60559 floating-point

<sup>3</sup>For instance  $j_{\max} < 100$  for  $\rho^2 e^2 < 0.007$  and  $n < 22,000$ .

<sup>4</sup>Typical cost of a global memory access is equal to the cost of 200–600 Flops (Floating-point operations) on our GPGPU device, cf. Advanced Micro Devices (2011).

(FP) numbers in our LFRF implementation, a FP underflow and FTZ (flush to zero) might occur with an attenuation of  $p_{m,m}$  for  $m > 0$  and  $\tau^2 \rightarrow 1$ , as can be clearly seen from Eq. (14), cf. Jekeli et al (2007).

The underflow is dangerous, because the seed values of Eq. (23) are affected, resulting in  $p_n = 0$  for any  $n \geq m$ . The underflow problem can be overcome neither by a different  $H_{n,m}$  nor by a normalisation related to  $p_{m,m}$  as suggested in Holmes and Featherstone (2002). Our approach does not rely on the extended-range FP arithmetics like in Lozier and Smith (1981) or Wittwer et al (2008), nor on the quadruple precision FP numbers. Instead we use a FP arithmetic with a dynamic exponent extension discussed in Nesvadba, 2008, *Towards the numerical evaluation of high degree and order associated Legendre functions as in EGM08*, GGEO, Chania, similar to the approach by Fukushima (2012).

According to ISO/IEC 60559, the FP number is defined as  $(-1)^s \mu 2^\xi$ , where  $\mu$  is the mantissa ( $f$  bits wide binary fraction),  $\xi$  is the exponent ( $g$  bits wide integer in two's complement) and  $s \in \{0, 1\}$  is the sign bit. We extend the FP number with an exponent extension  $\zeta$  ( $z$  bits wide integer in two's complement) that allows interpretation of the extended number as  $(-1)^s \mu 2^{\xi+\chi\zeta}$ . A constant parameter  $\chi$  is chosen to set up a reasonable computational range (band), as we will keep the exponent  $\xi$  loosely within the limits  $|\xi| \leq \chi$ .<sup>5</sup> Any FP number with  $\zeta < 0$  is treated as *subnormal* (please, do not confuse it with a synonym to “denormal” in ISO/IEC 60559). It means, it will not appear in the evaluation of numerical quantities of Eq. (10), nevertheless, if  $\zeta$  is shared among  $p_{n-1}$ ,  $p_{n-2}$ , it still may act in Eqs (22)–(23).

*Implementation in OpenCL* Because the LFRF steps are strictly sequential, we will assume that one usually needs to compute  $p_{n,m}$  values for (many) different arguments  $\tau$ ,  $\rho$ . In this case the implementation of the streaming LFRF algorithm in OpenCL is pretty straightforward, see Listing 1.

When the `pnm_lfrf_sto` kernel is submitted to a device, a global worksize  $L$  must be specified. The `get_global_id` function at line 3 in Listing 1 maps each work item to the unique index  $i \in \{0, 1, \dots, L-1\}$ . One-dimensional arrays (vectors) `Rho` and `Tau` serve as an input of  $\rho_i$  and  $\tau_i$  parameters. The `Pnm` array of a size  $L \times (N - m + 1)$  acts as an output buffer for  $p_{n,m}(\rho_i, \tau_i)$  organised row-wise, so that each row addresses  $p_{n,m}$  of particular  $n, m$ .<sup>6</sup> The work items fill up the `Pnm` array synchronously from  $p_{m+1,m}$  for  $p_{N,m}$ , where  $N$  and  $m$  are

<sup>5</sup>For instance, double precision FP numbers with  $\chi = 256$  provide the “operational band”  $2^{\pm 256} \approx 10^{\pm 77}$ , still allowing to track the numbers as small as  $2^{-549755814000} \approx 10^{-165492990300}$  (for  $z = 32$ ).

<sup>6</sup>Array addressing becomes clear from the source code. Note that the index formula at lines 14 and 25 is used just for a better code readability; the production version of the kernel utilises memory pointers.

```

1  __kernel pnm_lfrf_sto(__global const FPv *Tau,
2  __global const FPv *Rho, __global FPv *Pnm,
3  __global Iv *Zeta, const long m, const long N)
4  {
5  int i=get_global_id(0); int L=get_global_size(0);
6  const FPv tau = Tau[i];
7  const FPv rho = Rho[i];
8  FPv pn = Pnm[i]; // initial Pmm value
9  FPv pn1 = FPv(0.0); long l;
10 FPT alpha = 2*m+1; // alpha_l == 2l+1
11 FPT beta = 0; //beta_l,m = (1-m)*(1+m)
12 for(l=m+1; l<=N && (Zeta[i]); l++) {
13   FPv In = rho*sqrt((1+2.0/alpha)/(beta+alpha));
14   FPv pnt = tau*alpha*(pn*=In) - pn1*In;
15   beta += alpha; alpha += 2;
16   Pnm[i+(1-m)*L] = (FPv) (0.0);
17   Iv sowf = ((fabs(pnt)>ldexp(1.0,chi))? 2 : 0);
18   pn1 = ldexp(pn*beta, -chi*sowf);
19   pn = ldexp(pnt, -chi*sowf); // pnt*2^(-2*chi)
20   Zeta[i] += sowf;
21 }
22 for( ; l<=N; l++) {
23   FPv In = rho*sqrt((1+2.0/alpha)/(beta+alpha));
24   FPv pnt = tau*alpha*(pn*=In) - pn1*In;
25   beta += alpha; alpha += 2;
26   pn1 = pn*beta; // modif.p_l-2,m of the next step
27   Pnm[i+(1-m)*L] = pn = pnt;
28 }
29 }
30 __kernel void pmm_sto(__global const FPv *Tau,
31 __global const FPv *Rho, __global FPv *Pnm,
32 __global Iv *Zeta, const long m, const FPT Im)
33 {
34   const int i=get_global_id(0);
35   // NOT allowed fabs(Tau[i])>=1 || Rho[i]<=0 here
36   FPv pm = log2(Rho[i])*(m+1)+m*0.5*(log2(1-Tau[i])
37   +log2(1+Tau[i]));
38   Zeta[i] = 2*(convert_Iv_rtz(pm)/(2*chi));
39   Pnm[i] = Im*exp2(pm-(FPv)(chi*Zeta[i]));
40 }

```

**Listing 1** OpenCL implementation of the LFRF algorithm for  $p_{n,m}$ . Arguments  $\rho_i, \tau_i$  are provided in the input vectors `Tau` and `Rho`, resulting  $p_{n,m}(\rho_i, \tau_i)$  are saved to the `Pnm` array. The vector `Zeta` keeps exponent extensions  $\zeta_i$ , which can be initialised together with the first row of `Pnm` array by `pmm_sto` kernel.

constant parameters of the kernel. The first row of `Pnm` should be initialised with  $p_{m,m}(\rho_i, \tau_i)$  values before, that could be done by the `pmm_sto` kernel.

The LFRF loop at lines 10–19 in Listing 1 is dedicated to processing of subnormal  $p_{n,m}$ , which are presented always as zero at the output, see line 14. Within the evaluation of Eq. (23) the subnormality of  $p_{n,m}$  is tested. In the event of a *soft overflow*, i.e. in case when  $\xi > \chi$  occurs at line 15, the context variables `pn` and `pn1` are rescaled by  $2^{-2\chi}$  in combination with  $\zeta$  increment, see the lines 16–18 in Listing 1. The consequent loop at lines 20–26 processes ordinary LFRF just in the case of  $\zeta = 0$ .

Presented LFRF implementation enables us to employ SIMD execution paths by means of abstract data types `FPv`, `FPT` and `Iv`. For instance, to apply 4-way SIMD scheme we simply define `FPT` as `double`, `FPv` as `double4` and `Iv` as `int4` at the kernel compile time. Each work item processes a vector data now, multiplying achieved parallelism by the factor of four. For maximal streaming efficiency, however, the `pn`, `pn1` vectors should share  $\zeta$  in all the elements.



The LFRF coefficients  $I_{n,m}$  (for the given  $H_{n,m}$ ) are calculated “on the fly”, that is still faster, despite the consumption of 6 Flops including expensive fdiv and sqrt instructions, than a memory read. Nevertheless, further optimisation is possible by way of a local memory, see Sect. 4.2.

### 3.3 LFRF Algorithm for $q_{n,m}$

Referring to Eqs (15)–(17), we can implement the LFRF algorithm almost immediately. Nevertheless, for technical reasons it is more convenient to rearrange it as

$$q_{n-1} = \sqrt{1 - \rho^2 e^2} \alpha_n J_n q_n + \rho^2 e^2 \beta_{n+1} J_n J_{n+1} q_{n+1}, \quad (25)$$

where  $J_{n,m}$  coefficients are defined by

$$J_{n,m} = \frac{1}{\sqrt{1 - e^2}} \frac{1}{\alpha_n + h_{n+1,m}}. \quad (26)$$

Consequently we have  $J_{n-1} = (\sqrt{1 - e^2} \alpha_{n-1} + e^2 \beta_n J_n)^{-1}$ .

*Implementation in OpenCL* The implementation of a backward LFRF algorithm for  $q_{n,m}$  in Listing 2 is very similar to that of the forward LFRF. Therefore, we enjoy it to demonstrate an adoption of parallelism extended to parameter  $m$ .

Within the `qnm_lfrf_sto` kernel execution, each work item is parameterised with  $\rho_i$  and  $m_j$ , see the code at lines 3–6 and 22–24. Stress that the parallelism in  $m$  is done on a workgroup level now, so that all the work items in a workgroup share the same  $m_j$ ,  $j \in \{0, \dots, W - 1\}$ , where  $W \leq L$  is the total number of scheduled workgroups. The differences in  $m_j$  result in varying execution times between the workgroups, however, there is no need to take care about it.<sup>7</sup>

One can easily imagine, there are plenty of processing schemes possible on our kernels.

## 4 Spheroidal Harmonic Synthesis

All the tools for the numerical treatment of spheroidal harmonics have been provided in previous paragraphs. In order to demonstrate them in practice we consider one of the most essential applications: the evaluation of a linear combination of  $\bar{Y}_{n,m}$ , or the so-called *spheroidal harmonic synthesis*.

<sup>7</sup>An inactive workgroup (all the work items in the workgroup do not issue instructions) releases allocated compute resources, enabling an immediate execution of another job waiting in a device queue.

```

1  __kernel void qnm_lfrf_sto(__global const FPv *Rho,
2  __global FPv *Qnm, __global FPv *Qnm1,
3  __global const long *M, const long n,
4  __global const FPT *Hnm)
5  {
6  int i=get_global_id(0); int L=get_global_size(0);
7  const int j = get_group_id(0);
8  const FPv rho2e2 = e2*Rho[i]*Rho[i];
9  const long m = M[j];
10 FPT alpha = 2*n+1; // alpha_l== 2l+1
11 FPT beta = (n-m)*(n+m); //beta_l,m
12 FPT hn = beta*e2/(1-e2)/Hnm[j] - alpha;
13 FPv qn = Qnm[i], qn1 = Qnm1[i];
14 for(long l=n-1; l>=m; l--) {
15   FPT Jn = 1.0/sqrt(1-e2)/(alpha+hn);
16   hn = beta*Jn*e2/sqrt(1-e2);
17   FPv qnt = sqrt(1-rho2e2)*alpha*(qn*=-Jn) + qn1*Jn;
18   qn1 = rho2e2*beta*qn;
19   alpha -= 2; beta -= alpha;
20   Qnm[i+(n-l)*L] = qn = qnt;
21 }
22 }
23 __kernel void qnm_sto(__global const FPv *Rho,
24 __global FPv *Qnm, __global const long *M,
25 const long n)
26 {
27 const int i = get_global_id(0);
28 const long m = M[get_group_id(0)];
29 const FPv rho2 = Rho[i]*Rho[i];
30 FPT aj=0.5*(n+m+1), bj=0.5*(n-m+1), cj=n+1.5;
31 FPv kjr=1.0, nom=1.0, kj=1.0, denom=1.0;
32 for(int j = 1; j < jmax; j++) {
33   FPT rk = e2*aj*bj/(cj*j);
34   nom += kjr * rk*rho2;
35   denom += kj * rk;
36   aj += 1.0; bj += 1.0; cj += 1.0;
37   //if((kjr<=tol*nom)&&(kj<=tol*denom)) break;
38 }
39 Qnm[i] = nom/denom;
40 }

```

**Listing 2** OpenCL implementation of the backward LFRF given by Eq. (25). Kernel `qnm_lfrf_sto` computes  $q_{n,m}(\rho_i)$  from degree  $n$  for degree  $m$  ( $m$  may differ between workgroups) and write them to the `Qnm` array. Seed vectors `Qnm`, `Qnm1` can be initialised with the kernel `qnm_sto`. For the demonstration purposes, the coefficients  $J_{n,m}$  and  $h_{n,m}$  respectively are, again, calculated “on the fly” (at lines 12–13) with the aid of the seed constants  $h_{n,m}$  (line 9) provided in input vector `Hnm`.

Our intention is to calculate a sum

$$S_m(\tau, \rho) = \sum_{n=m}^N C_{n,m} p_{n,m}(\tau, \rho) q_{n,m}(\rho), \quad (27)$$

where  $C_{n,m} \in \mathbb{R}$  represent a set of the given coefficients.

If we recall Eqs (23) and (25), it is obvious that the directions of the LFRF processing for  $p_{n,m}$  and  $q_{n,m}$  are opposite. This is a substantial disadvantage, as it forces us to save one result, e.g.  $p_{n,m}$ , temporarily to the memory, to be read back later in the  $q_{n,m}$  multiplication. Such a constraint would break the advantage of parallel processing down to a limit given by the memory bandwidth.

### 4.1 Clenshaw Summation

To overcome the problem, let us recall a Clenshaw summation. The Clenshaw summation, originally proposed in

Clenshaw (1955), is well-known in applications of spherical harmonics, see e.g. Tscherning and Pöder (1982) and Holmes and Featherstone (2002). Using these references we can write

$$\sum_{n=m}^N C_n p_n = s_m p_m, \quad (28)$$

where  $p_m$  is given by Eq. (14) and  $s_m$  is obtained from

$$s_{n-1} = \tau \rho \alpha_n I_n s_n - \rho^2 \beta_{n+1} I_n I_{n+1} s_{n+1} + C_{n-1}, \quad (29)$$

starting with  $s_N = C_N$ ,  $s_{N+1} = 0$ . Note again, we are omitting the second index  $m$  and the arguments  $\tau, \rho$ .

As it can be seen from Eqs (25) and (29), the Clenshaw summation and backward LFRF computation schemes are similar. Therefore, we can easily compose the summation of  $C_n p_n$  with the  $q_n$  multiplication as follows

$$\beta_{n+1} = \beta_{n+2} - \alpha_{n+1}, \quad \alpha_n = \alpha_{n+1} - 2, \quad (30)$$

$$q_{n-1} = \sqrt{1 - \rho^2 e^2} \alpha_n J_n q_n + \rho^2 e^2 \beta_{n+1} J_n J_{n+1} q_{n+1}, \quad (31)$$

$$s_{n-1} = \tau \rho \alpha_n I_n s_n - \rho^2 \beta_{n+1} I_n I_{n+1} s_{n+1} + C_{n-1} q_{n-1}. \quad (32)$$

The seed values  $q_N$  and  $q_{N+1}$  are computed from Eq. (18), enabling us to put  $s_N = q_N C_N$ ,  $s_{N+1} = 0$ . Finally, we get

$$S_m = \sum_{n=m}^N C_n p_n q_n = s_m p_m, \quad (33)$$

where  $p_m$  can be determined from Eq. (22).

*Remark on the Accuracy of the Sum* In case of  $C_n$  coefficients decreasing in magnitude with increasing  $n$ , e.g. in Kaula's rule for Earth's gravity potential, the backward summation induces lower cancellation errors, see e.g. Goldberg (1991). Thus the Clenshaw summation provides a more accurate result in comparison with the forward sum of  $C_n p_n q_n$ . It can be shown, cf. Nesvadba (2011), that with a proper computation arrangement an overall achieved relative accuracy of  $S_m$  is better than  $10^{-13}$  for double or  $10^{-6}$  for single precision FP, even in the case of  $N = 21,600$ .

## 4.2 OpenCL Implementation of the Clenshaw Summation

The advantage of the Clenshaw summation is that the LFRF computation performs in a work item private memory space. Moreover, as we share  $\alpha_n$  and  $\beta_n$ , one step of the algorithm involves 18 Flops only. However, the global memory would be heavily accessed still, in order to get  $C_n$ ,  $I_n$  and  $J_n$  every step. Considering that these coefficients are shared by all the concurrent work items in a workgroup, the most convenient way in OpenCL is to copy the coefficients to a local memory.

```

1  __kernel void cspqnm_block(__global FPv *CBuff,
2  __global const FPv *Tau, __global const FPv *
3  Rho, const long m, const long n,
4  __global const FPT *Cnm, __global const FPT *Jnm)
5  {
6  int i=get_global_id(0); int L=get_global_size(0);
7  int wi=get_local_id(0); int wL=get_local_size(0);
8  const FPv tau = -Tau[i]; // private argument
9  const FPv rho = Rho[i]; // private argument
10 const FPv r2e2 = rho*rho*e2;
11 const FPv sigma = sqrt(1.0 - r2e2);
12 __local FPT In[BLOCK], Jn[BLOCK], Cn[BLOCK];
13 long dim = (n-m > BLOCK) ? BLOCK : n-m;
14 for(long j = wi; j < dim; j+=wL) {
15   In[j] = sqrt(((2*n-2*j+1))/((2*n-2*j-1)*(n-j+m)
16   *(n-j-m))); // I_n,m for std.norm
17   Jn[j] = Jnm[n-j-m]; // local copy of J_n,m
18   Cn[j] = Cnm[n-j-m]; // local copy of C_n,m
19 }
20 barrier(CLK_LOCAL_MEM_FENCE); // synchronise WG
21 FPv sn, sn1, qn, qn1;
22 load_context(CBuff, i, L, sn,sn1,qn,qn1);
23 FPT alpha = 2*n-1;
24 FPT beta = (FPT)(n-m-1)*(n+m-1);
25 #pragma unroll 8
26 for(ulong k = 0; k < dim; k++) {
27   FPv in = rho*In[k];
28   FPv snt = tau*alpha*(sn*=in) + sn1*in;
29   FPv qnt = sigma*alpha*(qn*=Jn[k]) + qn1*Jn[k];
30   sn1 = beta*sn; // s_l+1 * beta_l
31   qn1 = r2e2*beta*qn; // q_l+1*beta_l*rho2*e2
32   alpha -= 2; beta -= alpha;
33   sn = -snt + Cn[k]*(qn=qnt); //S+=C_l-1,m*q_l-1,m
34 }
35 save_context(CBuff, i, L, sn,sn1,qn,qn1);
36 }

```

**Listing 3** OpenCL kernel `cspqnm_block` for processing blocks of the Clenshaw summation. Arguments are provided in vectors  $Rho$  and  $Tau$ , context vector  $CBuff$  serves as the input of the LFRF seed values as well as the output of the results. Coefficients  $J_{n,m}$ ,  $C_{n,m}$  are provided on input,  $I_{n,m}$  are calculated once during the initialisation phase (line13).

Our implementation of the Clenshaw summation is illustrated in Listing 3. First of all, the sequence of  $s_n$  is aggregated to short blocks, e.g. 512 steps each. All the work items in the workgroup collaboratively load the local memory with the coefficients  $C_n$ ,  $I_n$ ,  $J_n$  of the respective block, see the code at lines 10–16. The `barrier` command at line 17 synchronises the work items and thus warrants the local memory consistency within the workgroup. The so-called context of the Clenshaw summation loop, i.e. a work item private variables  $sn$ ,  $sn1$ ,  $qn$ ,  $qn1$ , is temporarily saved at line 32 and restored at line 19. The  $s_n$  are, therefore, computed sequentially block by block with the `cspqnm_block` kernel. The computation of the partial sum is finalised after the last block according to Eq. (33) (though the respective implementation is omitted from Listing 3).

*FP Overflow in the Clenshaw Summation* The problem of underflow in  $p_n$  is now turned into the overflow in  $s_n$ . Nevertheless, the solution based on the FP numbers with extended exponent is still the same. Moreover, our implementation benefits from the block aggregation, where the exponent extension is processed outside of the `cspqnm_block` ker-

nel by examination (and possibly modification) of the context after each block. The respective exponent extension  $\zeta$ , necessary for the correct scale of  $C_n$  in Eq. (32), is propagated to the kernel indirectly, through the rescaled  $q_n$ ,  $q_{n1}$  at line 30. Finally, the exponent extension of  $s_m$  is adjusted with  $\zeta$  of the computed  $p_m$ .

## 5 Results and Conclusions

### 5.1 Performance

In order to demonstrate our implementation we conduct three different approaches to the calculation of Eq. (27):

- the LFRF approach to  $p_n$  (`pnm_lfrf_sto` kernel) combined with calculations of  $q_n$  based on Eq. (18),
- opposite LFRFs for  $p_n$  and  $q_n$  (`qnm_lfrf_sto` kernel) multiplied and summed up in the memory,
- the Clenshaw summation implemented by means of the local memory blocks (`cspqnm_block` kernel).
- Besides, there are examples of a spherical harmonic synthesis (`cspqnm_block` with  $e^2 = 0$  setting).

As a measure of the processing speed in Eq. (27) we will use “sn” units per second. The “sn” unit is defined as an evaluation of one elementary synthesis step, i.e. the calculation of  $s_n$  or  $C_n p_n q_n$ . Stress that the performance within our implementation is not sensitive to  $N$ , thus we can see the sn/s measure as an invariant, independent on  $n, m, \rho$  and (almost) on  $\tau$  parameters.

Computations were performed on the CPU and on the GPGPU device, both using the same code on the platform OpenCL 1.2 AMD-APP (923.1), Catalyst 13.4, GNU/Linux. As we can see from Table 1, all the tested kernels performs very well on the GPGPU. It becomes clear that the most efficient in Eq. (27) evaluation is the Clenshaw summation algorithm, mainly due to successfully masked memory latencies. Referring to the hardware documentation, Advanced Micro Devices (2011), one ideal “sn” step of Eqs (30)–(32) in double precision consume at least 28 clock cycles. A confrontation with Table 1 thus reveals, that our implementation of the `cspqnm_block` kernel gained about 75% of the peak machine performance. The device becomes saturated with the kernel at  $L = 6,144$  for FPv as double2 and FPv as float4, respectively. Therefore, one can deduce that the GPGPU keeps the kernel active in 16 simultaneous wavefronts per compute unit (CU).<sup>8</sup>

<sup>8</sup>The `cspqnm_block` kernel occupancy factor reaches 69% on our GPGPU, constrained mainly with the local memory size per CU.

**Table 1**  $S_m$  performance comparison for various computing kernels

OpenCL device	GPU AMD Radeon 6970 HD “Cayman”	CPU Intel Core i5 “Sandy Bridge”
Clock frequency	880 MHz	3 200 MHz
<b>Double precision</b>	24 CU: 768 PE	2 CU: 8 PE (AVX)
double4 FPv	675 GFlops, 180 GB/s	52 GFlops, 18 GB/s
pnm_lfrf_sto, qnm_sto, $j_{\max}=100$	<b>0.4 Gsn/s</b>	<b>4.4 Msn/s</b>
pnm_lfrf_sto, qnm_lfrf_sto	<b>1.9 Gsn/s</b>	27 Msn/s
cspqnm_block block=512	<b>18.0 Gsn/s</b>	0.47 Gsn/s
cspqnm_block block=512, e2=0	<b>29.3 Gsn/s</b>	0.70 Gsn/s
<b>Single precision</b>	24 CU: 1536 PE	2 CU: 16 PE (AVX)
float4 FPv	2.7 TFlops, 180 GB/s	102 GFlops, 18 GB/s
cspqnm_block block=64	<b>44.9 Gsn/s</b>	1.3 Gsn/s
cspqnm_block block=64, e2=0	<b>67.5 Gsn/s</b>	1.5 Gsn/s
<b>Energy efficiency, double4 FPv</b>		
OpenCL device	GPU AMD Radeon	CPU Intel Core i5
Input power	220 W	85 W
cspqnm_block block=512	<b>320 Msn/J</b>	20 Msn/J

*Spheroidal Harmonic Synthesis* Truncated series of  $C_{n,m} \bar{Y}_{n,m}$  can be evaluated with an important computational aid of Eq. (27).<sup>9</sup> A total number of sn units required for a full spheroidal harmonic synthesis for given  $(\rho, \tau, \lambda)$  in dependence on maximal degree  $N$  and  $\tau$  is estimated by

$$T(N, \tau) = (N + 1)^2 (2\sqrt{1 - \tau^2} - 1 + \tau^2) + 2(N + 1), \quad (34)$$

where the approximate relation  $m = n \cos \beta$  for a sufficient maximal order in the expansion was borrowed from Jekeli et al (2007). By applying the mean value theorem to  $\tau$  in Eq. (34) we can see that the enhanced LFRF run in a mean by 10% faster.<sup>10</sup> Dependence on  $\rho$  is not taken into account in Eq. (34), however, one can decrease  $N$  accordingly to the particular  $\rho$ , the respective  $C_{n,m}$  and the target accuracy.

As an example of the capabilities of our implementation performing on the single GPGPU device, we run a synthetic benchmark of the spheroidal harmonics synthesis

<sup>9</sup>As Eq. (27) should be applied on  $C_{n,m} \in \mathbb{C}$  in the spheroidal harmonics applications, we substitute it by two separate  $S_m$  calculations.

<sup>10</sup>Run-time criterion  $\zeta_{N,m} \neq 0$  (or resulting  $\zeta_{m,m} \neq 0$  in Clenshaw summation) based on the extended exponent approach is used for a decision about the maximal  $m$  relevant in particular synthesis. A little slowdown in enhanced LFRF caused by the exponent extension processing completely diminish with the higher orders cancellation.

for  $N = 21,600$ . According to Eq. (34), one full synthesis needs to perform at maximum 467 Msn, that agrees well with observed rates: 38 Hz in double and 95 Hz in single precision FP. Regular grid computations show even better results. For instance, a spheroidal harmonics series of  $N = 21,600$  in double precision can be expanded to a global spheroidal grid with a spatial resolution of 10 arcsecond within 30 minutes.

## 5.2 Summary

- The Clenshaw summation implemented in the OpenCL framework proved to be a very efficient way to an oblate spheroidal harmonic synthesis. The synthesis of  $\bar{Y}_{n,m}$  derivatives has been developed too. It relies on `sn1` and `qn1` variables, unfortunately, we are not able to provide more details here due to the given page limit.
- Our implementation is platform and device independent, however on the GPGPU we gained about **40 times faster** evaluation in comparison with a common CPU device. As a bonus, the computations offloaded to the GPGPU are about **16 times more energy efficient**.
- The hypergeometric series approach to  $q_{n,m}$ , e.g. Thong and Grafarend (1989) and Sebera et al (2012), seems to be **20–100 times slower** (in dependence on  $n, m$ , see also Fukushima 2013). Our implementation on GPGPU in comparison with the “traditional” approach on CPU is, therefore, accelerated by a factor of thousands.
- Spherical harmonics are obtained from the same code, just with  $e^2 = 0$  setting. Synthesis then performs about 40% faster, nevertheless, we do not recognise it as a significant handicap of the oblate spheroidal ( $e > 0$ ) case.
- In our experience, the OpenCL framework proves to be convenient and very useful for numerical tasks containing computationally intensive parallel algorithms, like in linear algebra, numerical integration, etc.

**Acknowledgements** The work on this paper was partly supported by the European Regional Development Fund (ERDF), project “NTIS – New Technologies for Information Society”, European Centre of Excellence, CZ.1.05/1.1.00/02.0090 and also by the Czech Science Foundation through Project No. 14-34595S. All this support is gratefully acknowledged.

## References

Abramowitz M, Stegun IA (1964) Handbook of mathematical functions with formulas, graphs, and mathematical tables. Dover, New York  
 Advanced Micro Devices I (2011) HD 6900 series instruction set architecture reference guide, revision 1.1  
 Bateman H, Erdélyi A (1953) Higher transcendental functions, vol 1. McGraw-Hill Book Company, New York/Toronto/London

Clenshaw CW (1955) A note on the summation of the Chebyshev series. *Math Tab Autom Comp* 9:118–120  
 Fukushima T (2012) Numerical computation of spherical harmonics of arbitrary degree and order by extending exponent of floating point numbers. *J Geod* 86:271–285  
 Fukushima T (2013) Recursive computation of oblate spheroidal harmonics of the second kind and their first-, second-, and third-order derivatives. *J Geod* 87:303–309  
 Gil A, Segura J (1998) A code to evaluate prolate and oblate spheroidal harmonics. *Comput Phys Commun* 108(2–3):267–278  
 Goldberg D (1991) What every computer scientist should know about floating-point arithmetic. *Assoc Comput Mach Comput Surv* 23(1):5–48  
 Hobson EW (1931) The theory of spherical and ellipsoidal harmonics, vol 1. Cambridge University Press, Cambridge  
 Holmes SA, Featherstone WE (2002) A unified approach to the Clenshaw summation and the recursive computation of very high degree and order normalised associated Legendre functions. *J Geod* 76:279–299  
 Holota P (2001) Variational methods in the recovery of the gravity field – Galerkin’s matrix for an ellipsoidal domain. *IAG Symposia* 123:277–283  
 Jekeli C, Lee JK, Kwon JH (2007) On the computation and approximation of ultra-high-degree spherical harmonic series. *J Geod* 81:603–615  
 Khronos Group KOWG (2012) The OpenCL specification version: 1.2, revision: 19  
 Lozier DW, Smith JM (1981) Algorithm 567 extended-range arithmetic and normalized Legendre polynomials. *ACM Trans Math Softw* 7:141–146  
 Moritz H (1984) Geodetic reference system 1980. *Bull Geod* 58(10):388–398  
 Nesvadba O (2011) Nothing to fear from ellipsoidal harmonics. EGU General Assembly, Copernicus.org. [http://presentations.copernicus.org/EGU2011-13386\\_presentation.pdf](http://presentations.copernicus.org/EGU2011-13386_presentation.pdf)  
 Sebera J, Bouman J, Bosch W (2012) On computing ellipsoidal harmonics using Jekeli’s renormalization. *J Geod* 86(9):713–726  
 Sona G (1995) Numerical problems in the computation of ellipsoidal harmonics. *J Geod* 70:117–126  
 Thong NC, Grafarend EW (1989) A spheroidal harmonic model of the terrestrial gravitational field. *Manuscr Geodaet* 14:285–304  
 Torge W (2001) Geodesy, 3rd edn. de Gruyter, Berlin  
 Tscherning CC, Poder K (1982) Some geodetic applications of Clenshaw summation. *Bollettino di geodesia e scienze affini* XLI(4):350–375  
 Wittwer T, Klees R, Seitz K, Heck B (2008) Ultra-high degree spherical harmonic analysis and synthesis using extended-range arithmetic. *J Geod* 82:223–229

---

# A Remark on the Computation of the Gravitational Potential of Masses with Linearly Varying Density

Maria Grazia D'Urso

---

## Abstract

The potential of a polyhedral body with linearly varying density has been given two different expressions in Holstein (Geophysics 68:157–167, 2003) and Hamayun et al. (J Geodesy 83:1163–1170, 2009) although in both papers the derivation is started from the same surface integral obtained by transforming the original volume integral via the Gauss theorem. Conversely, we prove that a suitable modification of the approach exploited by Hamayun et al. (J Geodesy 83:1163–1170, 2009) yields the formula derived by Holstein (Geophysics 68:157–167, 2003). Furthermore, an additional expression of the surface integral, which is also proved in this paper, allows us to derive a variant of the linear part of the potential, i.e. the integral multiplying the gradient of the density contrast, which filters the null contribution of faces containing the observation point. The new formula is specialized to the case of a prism.

---

## Keywords

Gravitational potential • Linear density variation • Polyhedron • Singularities

---

## 1 Introduction

The computation of the gravity effects (potential, gravity and tensor gradient fields) is usually carried out for mass distributions in which a constant density is assumed. This considerably simplifies the computation of the integrals, extended to the domain occupied by the given mass, which are intrinsic to the definition of the gravity effects.

However, the constant density assumption is not always realistic in geological structures, e.g. for modeling the compaction in sedimentary basins.

Preliminary studies on this issue (Chai and Hinze 1988; García-Abdeslem 1992, 2005; Gallardo-Delgado et al. 2003) have concerned the simple case of the right rectangular

parallelepiped (prism) since it represents a versatile tool for modeling complex mass distributions.

The case of polyhedral bodies has been addressed more recently since this simplifies geometric modelling of complex bodies (Pohánka 1988; Hansen 1999; Holstein 2003; Hamayun et al. 2009; Zhou 2009; D'Urso 2015a).

In this last paper the author has applied to the case of polyhedral bodies with linearly varying density a recent approach (D'Urso 2012, 2013, 2014, 2015b; D'Urso and Trotta 2015) for computing the gravity effects of bodies with uniform density and for consistently taking in account the relevant singularities.

Significantly, the same approach has been successfully applied to the solution of problems in geophysics (D'Urso and Marmo 2013), in geomechanics (Sessa and D'Urso 2013; D'Urso and Marmo 2015), in heat transfer (Rosati and Marmo 2014) and in elasticity (Marmo and Rosati 2015).

The aim of this paper is to illustrate the advantages of the formula presented in D'Urso (2015a) for the evaluation of the linear part of the potential, i.e. the integral multiplying the

---

M.G. D'Urso (✉)

DICeM - Department of Civil and Mechanical Engineering, University of Cassino and of Lazio Meridionale, Via G. Di Biasio 43, 03043 Cassino (FR), Italy  
e-mail: [URSO@unicas.it](mailto:URSO@unicas.it)

gradient of the density contrast, with respect to those derived by Holstein (2003) and Hamayun et al. (2009).

In particular the formula derived by Holstein (2003) is originated by a transformation of the volume integral expressing the linear part of the potential to a surface integral by means of Gauss theorem.

Though starting from the same surface integral, the formula derived later by Hamayun et al. (2009) appears to be quite different and the authors state in the introduction that *in comparison to the formula for the potential given by Holstein (2003), the main advantage of adopting the computational strategy developed by Pohánka (1988, 1998) is that our expression does not have the singular terms.*

Thus Hamayun et al. (2009) supplement formula (53) of their paper by a small positive number in order to avoid undefined operations occurring when an edge of a face does contain the origin of the reference frame local to the face.

We prove that a suitable modification of the approach exploited by Hamayun et al. (2009), based upon a further application of Gauss theorem in the plane to transform the surface integral to a line integral, yields the formula derived in Holstein (2003) by transforming the surface integral via Stokes theorem.

Thus we agree with Holstein (2003) that also for the linear part of the potential, the Pohánka-type exclusion zone is unnecessary, as already shown in Holstein and Ketteridge (1996) and Holstein et al. (1999), for the constant density case.

It is also shown that a suitable application of the Gauss theorem allows us to consistently take into account the singularity of the field appearing in the volume integral representing the linear part of the potential and to derive a novel expression of the related surface integral.

This last one exhibits the remarkable property of filtering the null contribution of the faces containing the observation point. A similar idea was presented in Holstein (2002b) for the constant density case.

Finally, the novel part of the formula proved in this paper is specialized to the case of a prism with a vertex coincident with the observation point.

## 2 Gravitational Potential of a Polyhedral Body with Linearly Varying Density

Let us consider an arbitrary bounded domain  $\Omega$  whose continuous mass distribution has a density  $\delta(\mathbf{s})$  varying linearly as function of a gradient  $\mathbf{g}$  and of the position vector  $\mathbf{s}$  of an arbitrary point belonging to it. Hence

$$\delta(\mathbf{s}) = \delta_o + \mathbf{g} \cdot \mathbf{s} \quad (1)$$

where  $\delta_o$  is a constant reference density evaluated at the origin  $O$  of a three-dimensional (3D) cartesian reference frame  $(O, x, y, z)$  in which the coordinates of  $\mathbf{s}$  are assigned.

Denoting by  $\mathbf{p}$  the position vector of an arbitrary point  $P$ , the gravitational potential  $U$  induced at  $P$  by the mass of  $\Omega$  is defined by the Newton integral:

$$U(P) = U(\mathbf{p}) = G \int_{\Omega} \frac{\delta(\mathbf{s})}{[(\mathbf{p} - \mathbf{s}) \cdot (\mathbf{p} - \mathbf{s})]^{1/2}} dV(\mathbf{s}) \quad (2)$$

where  $G$  is the gravitational constant.

Substituting in the previous formula the expression (1) for  $\delta(\mathbf{s})$ , adding and subtracting the quantity  $\mathbf{g} \cdot \mathbf{p}$ , setting  $\mathbf{r} = \mathbf{s} - \mathbf{p}$  and  $r = (\mathbf{r} \cdot \mathbf{r})^{1/2}$  one has

$$\begin{aligned} U(P) &= G(\delta_o + \mathbf{g} \cdot \mathbf{p}) \int_{\Omega} \frac{dV}{r} + G\mathbf{g} \cdot \int_{\Omega} \frac{\mathbf{r}}{r} dV = \\ &= G(\delta_o + \mathbf{g} \cdot \mathbf{p})U_c(P) + G\mathbf{g} \cdot \mathbf{U}_l(P) \end{aligned} \quad (3)$$

where the suffixes  $(\cdot)_c$  and  $(\cdot)_l$  have been used to denote the constant and linear contributions to  $U(P)$ , respectively.

For polyhedral bodies the computation of  $U_c(P)$  has been addressed by several authors see e.g. D'Urso (2012, 2013, 2014) for a comprehensive list of references.

In these papers the author has exploited a novel approach for computing the gravity effects of bodies with uniform density and for consistently taking into account the relevant singularities. The approach has been subsequently extended in D'Urso (2015a) to the case of polyhedral bodies with linearly varying density.

Thus, it is interesting to compare the expressions contributed in this last paper with the already available ones. Due to space limitations we shall limit ourselves to compare the expressions of  $\mathbf{U}_l$  in D'Urso (2015a) with those due to Holstein (2003) and Hamayun et al. (2009). In the sequel they will be denoted as  $\mathbf{U}_l^{DUR}$  and  $\mathbf{U}_l^{HH}$  respectively.

According to Pohánka (1998), Holstein (2003), and Hamayun et al. (2009) the integral  $\mathbf{U}_l$  in (3) can be expressed as

$$\mathbf{U}_l^{HH}(P) = \int_{\Omega} \frac{\mathbf{r}}{r} dV = \int_{\Omega} \text{grad } r dV = \int_{Fr(\Omega)} r \mathbf{n} dA \quad (4)$$

where the last equality follows from Gauss theorem. For a polyhedral body the previous expression specializes to

$$\mathbf{U}_l^{HH}(P) = \sum_{i=1}^{N_F} \mathbf{n}_i \int_{F_i} (\mathbf{r}_i \cdot \mathbf{r}_i)^{1/2} dA_i = \sum_{i=1}^{N_F} \mathbf{n}_i \int_{F_i} r_i dA_i \quad (5)$$

where  $N_F$  denotes the number of faces of the polyhedral body, the vector  $\mathbf{r}_i$  spans the  $i$ th face  $F_i$  of the boundary  $Fr(\Omega)$  of  $\Omega$  and  $\mathbf{n}_i$  is the outward unit normal to  $F_i$ .

To derive an alternative expression for  $U_l$  which takes into account the singularity at  $\mathbf{r} = \mathbf{o}$  intrinsic to the integrand in (4) we briefly report the approach contributed in D'Urso (2014); it stems from the identity

$$\operatorname{div} \frac{\mathbf{r} \otimes \mathbf{r}}{r} = \operatorname{div}[r^2 \mathbf{r}] \otimes \frac{\mathbf{r}}{r^3} = \frac{3\mathbf{r}}{r} + [r^2 \mathbf{r}] \operatorname{div} \frac{\mathbf{r}}{r^3} \quad (6)$$

and from the properties  $\operatorname{div}[\mathbf{r}/r^3] = 0$  if  $\mathbf{r} \neq \mathbf{o}$  and

$$\int_{\Omega} \varphi(\mathbf{r}) \operatorname{div} \frac{\mathbf{r}}{r^3} dV = \begin{cases} 0 & \text{if } \mathbf{o} \notin \Omega \\ m(\mathbf{o})\varphi(\mathbf{o}) & \text{if } \mathbf{o} \in \Omega \end{cases} \quad (7)$$

where  $\varphi$  is a continuous scalar field and  $m$  represents the measure, expressed in radians, of the solid angle of the intersection between  $\Omega$  and a spherical neighborhood of the singularity point  $\mathbf{o}$ . To sum up,  $\operatorname{div}(\mathbf{r}/r^3)$  represents the Dirac delta function  $\Delta$  at  $\mathbf{o}$  D'Urso and Russo (2002).

Thus, integrating (6) over  $\Omega$  one has

$$\begin{aligned} U_l^{Dur}(P) &= \int_{\Omega} \frac{\mathbf{r} dV}{r} = \frac{1}{3} \int_{\Omega} \operatorname{div} \frac{\mathbf{r} \otimes \mathbf{r}}{r} dV - \\ &\quad - \frac{1}{3} \int_{\Omega} r^2 \mathbf{r} \Delta(\mathbf{o}) dV = \frac{1}{3} \int_{Fr(\Omega)} \frac{(\mathbf{r} \cdot \mathbf{n}) \mathbf{r}}{r} dA \end{aligned} \quad (8)$$

where the first volume integrand has a removable singularity, allowing the integral to be defined even if  $\mathbf{r} = \mathbf{o}$  occurs inside or on the boundary of  $\Omega$ . Similarly, the second volume integral also has a removable singularity in the integrand, making it safe to apply Gauss' theorem. Together with result (7), we arrive at the final integral in Eq. (8).

Its expression specializes to polyhedral bodies as

$$U_l^{Dur}(P) = \frac{1}{3} \sum_{i=1}^{N_F} d_i \int_{F_i} \frac{\mathbf{r}_i}{(\mathbf{r}_i \cdot \mathbf{r}_i)^{1/2}} dA_i = \frac{1}{3} \sum_{i=1}^{N_F} d_i \int_{F_i} \frac{\mathbf{r}_i}{r_i} dA_i \quad (9)$$

$d_i$  being the distance between the point  $P$  and the face  $F_i$ .

The main advantage of the previous expression with respect to (5) is that the integral pertaining to each face is scaled by the factor  $d_i$  so that a face containing the observation point does provide a null contribution to  $U_l$ ; hence, it is easy to test whether to omit all calculations for a face.

Nevertheless, it will be shown that formulas obtained by the specialization of (5) and (9) as function of the vertices of each face are very closely related.

To this end we denote by  $P_i$  the orthogonal projection of the observation point  $P$  on  $F_i$  and assume  $P_i$  as origin of a 2D reference frame local to the face, see D'Urso (2014).

Furthermore, we decompose the vector  $\mathbf{r}_i$  as sum of a vector  $\mathbf{r}_i^\perp$  orthogonal to the face  $F_i$  and a vector  $\mathbf{r}_i^\parallel$  parallel to it by setting

$$\mathbf{r}_i = \mathbf{r}_i^\perp + \mathbf{r}_i^\parallel = (\mathbf{r}_i \cdot \mathbf{n}_i) \mathbf{n}_i + \mathbf{r}_i^\parallel = d_i \mathbf{n}_i + \mathbf{T}_{F_i} \boldsymbol{\rho}_i \quad (10)$$

where the vector  $\boldsymbol{\rho}_i = (\xi_i, \eta_i)$  represents the position vector of a generic point of the  $i$ th face with respect to  $P_i$ .

The linear operator  $\mathbf{T}_{F_i}$  maps the 2D vector  $\boldsymbol{\rho}_i$  to the 3D one  $\mathbf{r}_i^\parallel$  (D'Urso 2013, 2014) and fulfills the property  $\mathbf{T}_{F_i}^t \mathbf{T}_{F_i} = I_{2D}$  where  $(\cdot)^t$  denotes transpose and  $I_{2D}$  is the 2D identity operator.

### 3 Two Alternative Approaches for Computing the Linear Contribution $U_l$ to the Gravitational Potential

The surface integral (5) is transformed to a line integral by Hamayun et al. (2009) by using the Gauss theorem in the plane of the generic face. In this way the authors derive a formula which appears to be considerably different from the one previously derived in Holstein (2003) in spite of the fact that also this last author uses formula (5) as starting point.

Closer inspection shows that the differences are only notational. The same underlying functions are used in both cases, in similar combinations.

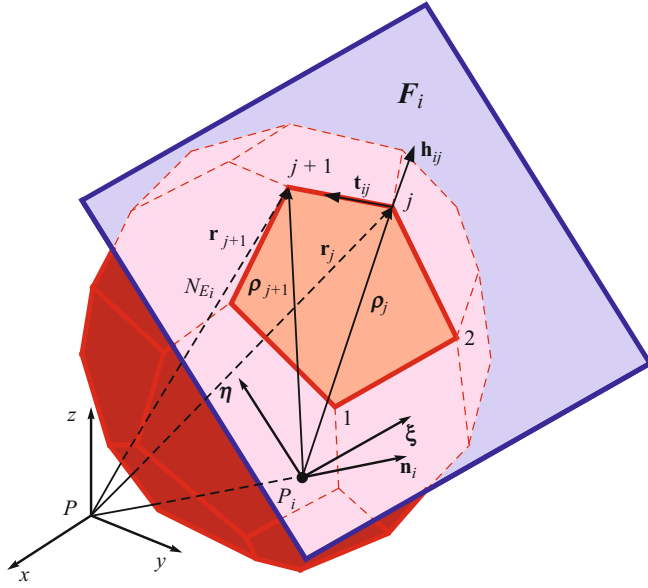
Furthermore, Hamayun et al. (2009) supplement formula (53) of their paper by a small positive number, originally introduced by Pohánka (1988, 1998), in order to avoid undefined operations when the observation does belong to an edge of a face.

In this respect they seem to be unaware of the fact that Holstein (2003) had already proved that the potential and the first gradient were free from singularities and that the gravity gradient tensor was undefined on facet surfaces, though  $U_l$  has a continuous second gradient

In order to reconcile the two approaches we substitute (10) in (5); setting  $f(\boldsymbol{\rho}_i, d_i) = (\boldsymbol{\rho}_i \cdot \boldsymbol{\rho}_i + d_i^2)^{1/2}$  we get

$$\mathbf{U}_l^{Ham} = \sum_{i=1}^{N_F} \mathbf{n}_i \int_{F_i} f(\boldsymbol{\rho}_i, d_i) dA_i \quad (11)$$

where the suffix  $(\cdot)^{Ham}$  has been used to remind of the fact that we are using the same approach exploited by Hamayun et al. (2009), i.e. the Gauss theorem, to transform the previous surface integral to a line integral.



**Fig. 1** 3D and 2D notation for vertices and edges on the generic face

To this end we apply the differential identity (Tang 2006)

$$\operatorname{div}\left[f(\boldsymbol{\rho}_i, d_i)\boldsymbol{\rho}_i\right] = 3f(\boldsymbol{\rho}_i, d_i) - \frac{d_i^2}{f(\boldsymbol{\rho}_i, d_i)} \quad (12)$$

which provides, upon application of Gauss theorem

$$\begin{aligned} \mathbf{U}_l^{Ham} &= \frac{1}{3} \sum_{i=1}^{N_F} \mathbf{n}_i \left[ d_i^2 \int_{F_i} \frac{dA_i}{f(\boldsymbol{\rho}_i, d_i)} + \int_{Fr(F_i)} f(\boldsymbol{\rho}_i, d_i) \beta_i ds_i \right] \\ &= \frac{1}{3} \sum_{i=1}^{N_F} [d_i^2 \mathbf{n}_i I_{c, F_i} + \mathbf{n}_i \Upsilon_{l, F_i}] \end{aligned} \quad (13)$$

where  $\beta_i = \boldsymbol{\rho}_i \cdot \mathbf{v}(s_i)$  and  $\mathbf{v}(s_i)$  is the 2D unit vector which is normal to  $Fr(F_i)$  at the abscissa  $s_i$  and points outside  $F_i$ .

The integral  $I_{c, F_i}$  is required as well for computing  $U_c$  in (3), see, e.g., Holstein et al. (1999) for a comparison of different expressions and D'Urso (2012, 2013).

For a polyhedral body the 2D unit normal  $\mathbf{v}(s_i)$  in (13) is constant on each side of  $F_i$  so that

$$\Upsilon_{l, F_i} = \sum_{j=1}^{N_{E_i}} [\boldsymbol{\rho}_{ij} \cdot \mathbf{v}_{ij}] \int_{E_j} (\mathbf{r}_{ij} \cdot \mathbf{r}_{ij})^{1/2} ds_{ij} \quad (14)$$

where  $E_j$  denotes the  $j$ th edge of  $F_i$ ,  $N_{E_i}$  the number of edges defining  $Fr(F_i)$ ,  $\mathbf{r}_{ij}$  ( $\boldsymbol{\rho}_{ij}$ ) the 3D (2D) vector spanning  $E_j$  and  $\mathbf{v}_{ij}$  the unit vector orthogonal to the  $j$ th edge.

Denoting by  $\mathbf{t}_{ij}$  the unit vector directed along the  $j$ th edge, oriented counter-clockwise around the normal  $\mathbf{n}_i$ , see

e.g. Fig. 1, and setting  $\mathbf{h}_{ij} = \mathbf{t}_{ij} \times \mathbf{n}_i$  we observe that

$$\boldsymbol{\rho}_{ij} \cdot \mathbf{v}_{ij} = \mathbf{T}_{F_i} \boldsymbol{\rho}_{ij} \cdot \mathbf{T}_{F_i} \mathbf{v}_{ij} = \mathbf{r}_{ij}^{\parallel} \cdot \mathbf{h}_{ij} = \mathbf{r}_{ij} \cdot \mathbf{h}_{ij} \quad (15)$$

since  $\mathbf{r}_{ij}^{\perp} \cdot \mathbf{h}_{ij} = 0$ .

Hence  $\mathbf{U}_l^{Ham}$  in (13) does coincide with the expression  $V_2$  that Holstein (2003) obtained in formula (A-6) of his paper by using Stokes theorem instead of Gauss one.

For this reason we use the apex  $(\cdot)^{HH}$  in (13) and write

$$\mathbf{U}_l^{HH} = \frac{1}{3} \sum_{i=1}^{N_F} [d_i^2 \mathbf{n}_i I_{c, F_i} + \mathbf{n}_i \Upsilon_{l, F_i}^{HH}] = \frac{1}{3} \sum_{i=1}^{N_F} [\mathbf{U}_{l_1} + \mathbf{U}_{l_2}^{HH}] \quad (16)$$

where

$$\Upsilon_{l, F_i}^{HH} = \sum_{j=1}^{N_{E_i}} \mathbf{r}_{ij} \cdot \mathbf{h}_{ij} \int_{E_j} (\mathbf{r}_{ij} \cdot \mathbf{r}_{ij})^{1/2} ds_{ij} = \sum_{j=1}^{N_{E_i}} (\mathbf{r}_{ij} \cdot \mathbf{h}_{ij}) I_{(l, F_i)_j} \quad (17)$$

In order to facilitate the comparison with the previous formula we specialize formula (9) by using the approach deducible from Holstein (2003). It is based on the identity

$$\operatorname{curl}(\varphi \mathbf{u}) = \operatorname{grad} \varphi \times \mathbf{u} \quad (18)$$

where  $\varphi(\mathbf{u})$  is a scalar (constant) vector field, (Tang 2006). Actually, by virtue of Stokes theorem one has:

$$\int_{F_i} \operatorname{grad} \varphi \times \mathbf{n}_i dA_i = - \int_{Fr(F_i)} \varphi \mathbf{t} ds_i \quad (19)$$

where  $\mathbf{t}$  is the unit vector tangent to the boundary of  $F_i$ .

On the other hand we infer from (9) and (10)

$$\mathbf{U}_l^{Dur} = \frac{1}{3} \sum_{i=1}^{N_F} \left[ d_i^2 \mathbf{n}_i I_{c, F_i} + d_i \int_{F_i} \frac{\mathbf{r}_i^{\parallel}}{(\mathbf{r}_i \cdot \mathbf{r}_i)^{1/2}} dA_i \right] \quad (20)$$

and the last integral becomes, on account of (19),

$$\begin{aligned} \mathbf{r}_{l, F_i}^{Dur} &= \int_{F_i} \frac{\mathbf{r}_i^{\parallel}}{r_i} dA_i = \int_{F_i} \frac{\mathbf{n}_i \times (\mathbf{r}_i \times \mathbf{n}_i)}{r_i} dA_i = \\ &= \int_{F_i} \mathbf{n}_i \times [\operatorname{grad} r_i \times \mathbf{n}_i] dA_i = -\mathbf{n}_i \times \int_{Fr(F_i)} r_i \mathbf{t} ds_i \end{aligned} \quad (21)$$

Thus for a polyhedral body, we finally have

$$\mathbf{r}_{l, F_i}^{Dur} = \sum_{j=1}^{N_{E_i}} \mathbf{h}_{ij} \int_{E_j} (\mathbf{r}_{ij} \cdot \mathbf{r}_{ij})^{1/2} ds_{ij} = \sum_{j=1}^{N_{E_i}} \mathbf{h}_{ij} I_{(l, F_i)_j} \quad (22)$$



and (20) becomes

$$\mathbf{U}_i^{Dur} = \frac{1}{3} \sum_{i=1}^{N_F} [d_i^2 \mathbf{n}_i I_{c,F_i} + d_i \boldsymbol{\Upsilon}_{l,F_i}^{Dur}] = \frac{1}{3} \sum_{i=1}^{N_F} [\mathbf{U}_{l_1} + \mathbf{U}_{l_2}^{Dur}] \quad (23)$$

Numerical experiments have proved that the previous formula and (16), which is quite similar, yield the same result.

#### 4 Comparative Assessment of the Two Formulas for Computing $\mathbf{U}_{l,2}^{HH}$ and $\mathbf{U}_{l,2}^{Dur}$

A comparison of (16) and (23) shows that  $\mathbf{U}_l^{HH}$  and  $\mathbf{U}_l^{Dur}$  differ only for the second addend. However, as anticipated soon after formula (9),  $\mathbf{U}_{l_2}^{Dur}$  leads to a control structure of the calculations in which an entire face can be omitted, analogously to  $\mathbf{U}_{l_1}$ , whenever it contains the observation point.

In the homogeneous density case this property had already been noted in Holstein (2002b), see section ‘‘Formula Alternatives’’ of his paper.

Conversely, this has not yet been established in the linearly variable case although formulas (19)–(21) in Holstein (2003) had been organized to exhibit the  $d_i$  factor.

It is also worth noting that the integral  $I_{(l,F_i)_j}$  in (17) and (22) hides a singularity, holding when the observation point is aligned with an edge, although it will be shown in the next section that the singularity is removable.

This is not surprising since we know from the very beginning that the integral  $\mathbf{U}_l$  is well defined. This can be proved either by a limiting process of the integral extended to a domain with an exclusion zone embodying the singularity, as outlined in Holstein (2003), or, more elegantly, by distribution theory as in (6) and (7).

Nevertheless, the well posedness of the integral defining  $\mathbf{U}_l$  does not exclude that some harmless indeterminate subexpressions have to be properly coped with in order to actually evaluate it.

This can happen in the first reduction, from a 3D to a 2D integral, or in the second one, from a 2D to a line integral. The former case characterizes the approach by Holstein (2003) and Hamayun et al. (2009) since they ignore the singularity at  $\mathbf{r} = \mathbf{o}$  in the transformation (4) leading to the surface integral.

Conversely, a removable singularity affects the line integral stemming from the second reduction carried out in this paper since Stokes theorem has been applied in (19) to the singular scalar field  $\varphi = 1/(\mathbf{r}_i \cdot \mathbf{r}_i)^{1/2}$ .

In order to express  $I_{(l,F_i)_j}$  as function of the basic input data of the polyhedral body, i.e. the coordinates of the

vertices defining the generic face  $F_i$ , we write

$$\mathbf{r}_{ij}(\lambda_j) = \mathbf{r}_j + \lambda_j(\mathbf{r}_{j+1} - \mathbf{r}_j) \quad \lambda_j \in [0, 1] \quad (24)$$

so that, being  $\lambda_j = s_{ij}/l_j$  and  $l_j$  the edge length, it turns out

$$I_{(l,F_i)_j} = l_j \int_0^1 (p_j \lambda_j^2 + 2q_j \lambda_j + t_j)^{1/2} d\lambda_j = l_j I_{l,E_j} \quad (25)$$

where it has been set  $p_j = (\mathbf{r}_{j+1} - \mathbf{r}_j) \cdot (\mathbf{r}_{j+1} - \mathbf{r}_j)$ ,  $q_j = \mathbf{r}_j \cdot (\mathbf{r}_{j+1} - \mathbf{r}_j)$ ,  $t_j = \mathbf{r}_j \cdot \mathbf{r}_j = |\mathbf{r}_j|^2$ .

Denoting by  $u_j = p_j + 2q_j + t_j$  and

$$LN_j = \ln k_j = \ln \left[ \frac{p_j + q_j + \sqrt{p_j} \sqrt{p_j + 2q_j + t_j}}{q_j + \sqrt{p_j t_j}} \right] \quad (26)$$

the integral in (25) is provided by

$$I_{l,E_j} = \frac{p_j t_j - q_j^2}{2p_j} LN_j + \frac{(p_j + q_j) \sqrt{u_j} - q_j \sqrt{t_j}}{2\sqrt{p_j}} \quad (27)$$

an expression whose numerical properties will be discussed in the next section. It can be usefully compared with formulas (17) and (30) in Holstein (2003).

The well-posedness of  $I_{(l,F_i)_j}$  depends upon the radicands and the arguments of the logarithm in (26). To this end we first observe that

$$p_j = l_j^2 > 0 \quad t_j \geq 0 \quad u_j = \mathbf{r}_{j+1} \cdot \mathbf{r}_{j+1} = |\mathbf{r}_{j+1}|^2 \geq 0 \quad (28)$$

so that the radicands in (26) and (27) turn out to be positive unless  $|\mathbf{r}_{j+1}| = \mathbf{o}$  or  $|\mathbf{r}_j| = \mathbf{o}$ . In this case, however,  $P$  belongs to  $F_i$  and, as detailed below, the contribution of the  $j$ th edge to  $\Upsilon_{l,F_i}^{HH}$  and  $\boldsymbol{\Upsilon}_{l,F_i}^{Dur}$  is null.

Moreover, both the numerator and the denominator of the logarithm in (26) are positive when the observation point does not belong to the  $j$ th edge. In the opposite case at least one of them can vanish but, as stated above, the computation of the logarithm is not required. Actually, invoking the definition of  $p_j, q_j, t_j$  in (25) one has

$$k_j = \frac{\mathbf{r}_{j+1} \cdot (\mathbf{r}_{j+1} - \mathbf{r}_j) + l_j |\mathbf{r}_{j+1}|}{\mathbf{r}_j \cdot (\mathbf{r}_{j+1} - \mathbf{r}_j) + l_j |\mathbf{r}_j|} = \frac{\mathbf{r}_{j+1} \cdot \mathbf{t}_{ij} + |\mathbf{r}_{j+1}|}{\mathbf{r}_j \cdot \mathbf{t}_{ij} + |\mathbf{r}_j|} \quad (29)$$

Should  $\mathbf{r}_{j+1} \cdot \mathbf{t}_{ij}$  or  $\mathbf{r}_j \cdot \mathbf{t}_{ij}$  be both negative, the numerically more stable expression in Holstein and Ketteridge (1996) can be usefully resorted to.

By definition of scalar product one infers that both the numerator and the denominator are non-negative unless  $d_i = 0$  and the observation point is aligned with the  $j$ th edge. Actually, under these assumptions, either  $|\mathbf{r}_{j+1}|$

or  $|\mathbf{r}_j|$  can vanish or, if both of them are non-zero, at least one between the numerator and the denominator can vanish.

All the previous cases imply that either  $\mathbf{r}_j$  or  $\mathbf{r}_{j+1}$  is the null vector or that  $\mathbf{r}_j$  and  $\mathbf{r}_{j+1}$  are parallel. Accordingly,  $LN_j$  in (26), and hence  $(I_{l,F_i})_j$  in (25), tends to  $+\infty$  or  $-\infty$  with an infinitesimally low rate.

This has no effect on  $\Upsilon_{l,F_i}^{HH}$  since  $\mathbf{r}_{ij} \cdot \mathbf{h}_{ij} = 0$  in this case and the product  $(\mathbf{r}_{ij} \cdot \mathbf{h}_{ij})I_{l,E_j} = 0$  in (17) tends to zero, what excludes the necessity of computing  $LN_j$  in (26).

The role played by  $\mathbf{r}_{ij} \cdot \mathbf{h}_{ij}$  for ensuring the well-posedness of  $\Upsilon_{l,F_i}^{HH}$  is played by  $d_i$  for  $\Upsilon_{l,F_i}^{Dur}$ . Actually,  $k_j \rightarrow 0$  or  $k_j \rightarrow +\infty$  depending on the fact that the numerator or the denominator in (26) vanish; in both cases  $LN_j = \ln k_j$  is an infinite of arbitrarily low degree. This means that

$$\lim_{d_i \rightarrow 0} d_i LN_j = \lim_{d_i \rightarrow 0} d_i \ln k_j = 0 \quad (30)$$

stating that the  $j$ th edge of  $F_i$  gives a null contribution to the quantity  $\Upsilon_{l,F_i}^{Dur}$  in (22).

Thus, irrespective of whether  $k_j$  is zero or undefined through a zero divisor, one can skip the evaluation of the contribution of the edges belonging to a face characterized by  $d_i = 0$ .

We remark that the parameter  $k_j$  is reminiscent of the parameter  $\Lambda = l_j / (|\mathbf{r}_{j+1}| + |\mathbf{r}_j|)$  introduced by Strakhov et al. (1986) and extensively used in Holstein (2002a, 2003) although  $\Lambda \in [0, 1]$  while  $k_j \in [0, +\infty)$ .

## 5 Specialization of $U_{l_2}^{Dur}$ to a Prism

Let us now specialize the general formula (23) for  $U_{l_2}^{Dur}$  to the case of a right rectangular parallelepiped (prism) whose sides are parallel to the axes of a Cartesian reference frame.

We have selected  $U_{l_2}^{Dur}$  instead of  $U_{l_2}^{HH}$  because this last one requires longer calculations in this case. Actually, for an observation point coincident with a vertex of the parallelepiped, which is the case addressed in this section,  $U_{l_2}^{HH}$  has to be computed for all the six faces of the prism while just the three faces which do not contain the origin are needed for  $U_{l_2}^{Dur}$ .

In any case the computation of  $U_{l_2}^{Dur}$  and  $U_{l_2}^{HH}$  depends essentially upon the evaluation of the integral (25). Due to space limitations the evaluation of  $I_{(l,F_i)j}$  is detailed only for the face  $F_1$ , i.e. the one orthogonal to the  $x$  axis and not containing the origin. For this reason we have shown in Fig. 2 the vertices which define the first face.

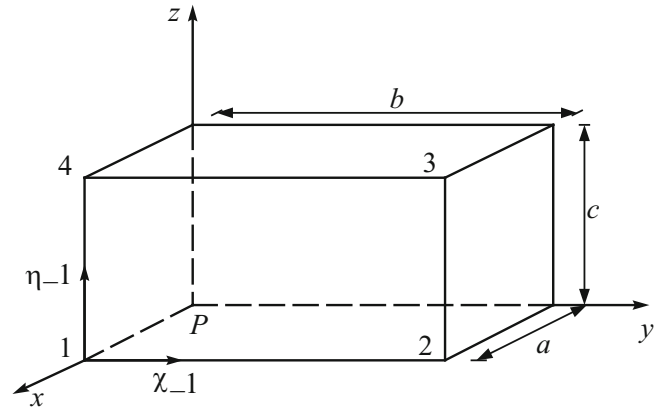


Fig. 2 Geometry of the prism and vertices of the first face

Thus, the position vectors of the vertices are  $\mathbf{r}_1 = (a, 0, 0)$ ,  $\mathbf{r}_2 = (a, b, 0)$ ,  $\mathbf{r}_3 = (a, b, c)$ ,  $\mathbf{r}_4 = (a, 0, c)$  so that

$$\begin{aligned} \mathbf{r}_2 - \mathbf{r}_1 &= (0, b, 0) & \mathbf{r}_3 - \mathbf{r}_2 &= (0, 0, c) \\ \mathbf{r}_4 - \mathbf{r}_3 &= (0, -b, 0) & \mathbf{r}_1 - \mathbf{r}_4 &= (0, 0, -c) \end{aligned} \quad (31)$$

Accordingly, one infers from (25)–(27)

$$\begin{aligned} p_1 &= b^2 & q_1 &= 0 & u_1 &= 0 & t_1 &= a^2 \\ p_1 + 2q_1 + t_1 &= a^2 + b^2 & (p_1 t_1 - q_1^2) / (2p_1) &= a^2 / 2 \end{aligned} \quad (32)$$

for the first edge of the first face

$$\begin{aligned} p_2 &= c^2 & q_2 &= 0 & u_2 &= b^2 \\ t_2 &= a^2 + b^2 & p_2 + 2q_2 + t_2 &= a^2 + b^2 + c^2 \\ (p_2 t_2 - q_2^2) / (2p_2) &= (a^2 + b^2) / 2 \end{aligned} \quad (33)$$

for the second edge of the first face

$$\begin{aligned} p_3 &= b^2 & q_3 &= -b^2 & u_3 &= b^2 + c^2 \\ t_3 &= a^2 + b^2 + c^2 & p_3 + 2q_3 + t_3 &= a^2 + c^2 \\ (p_3 t_3 - q_3^2) / (2p_3) &= (a^2 + c^2) / 2 \end{aligned} \quad (34)$$

for the third edge of the first face

$$\begin{aligned} p_4 &= c^2 & q_4 &= -c^2 & u_4 &= c^2 \\ t_4 &= a^2 + c^2 & p_4 + 2q_4 + t_4 &= a^2 \\ (p_4 t_4 - q_4^2) / (2p_4) &= a^2 / 2 \end{aligned} \quad (35)$$

for the fourth edge of the first face.

Hence, setting  $l_d = \sqrt{a^2 + b^2 + c^2}$ , formula (27) yields

$$\begin{aligned}(I_{l,F_1})_1 &= \frac{a^2}{2} \ln \frac{b + \sqrt{a^2 + b^2}}{a} + \frac{b\sqrt{a^2 + b^2}}{2} \\(I_{l,F_1})_2 &= \frac{a^2 + b^2}{2} \ln \frac{c + l_d}{\sqrt{a^2 + b^2}} + \frac{cl_d}{2} \\(I_{l,F_1})_3 &= \frac{a^2 + c^2}{2} \ln \frac{\sqrt{a^2 + c^2}}{-b + l_d} + \frac{bl_d}{2} \\(I_{l,F_1})_4 &= \frac{a^2}{2} \ln \frac{a}{-c + \sqrt{a^2 + c^2}} + \frac{c\sqrt{a^2 + c^2}}{2}\end{aligned}\quad (36)$$

Due to the symmetry of the problem at hand  $(I_{l,F_2})_j$  is obtained from  $(I_{l,F_1})_j$ , ( $j = 1 \dots 4$ ) by setting in the formula above the ordered triple  $\{b, c, a\}$  in place of  $\{a, b, c\}$ . In turn  $(I_{l,F_2})_j$  transforms to  $(I_{l,F_3})_j$  by means of the substitution  $\{b, c, a\} \rightarrow \{c, a, b\}$ .

Being also  $\mathbf{h}_{11} = (0, 0, -1)$ ,  $\mathbf{h}_{12} = (0, 1, 0)$ ,  $\mathbf{h}_{13} = (0, 0, 1)$ ,  $\mathbf{h}_{14} = (0, -1, 0)$  and  $d_1 = a$ ,  $d_2 = b$ ,  $d_3 = c$ , the expression of  $U_{l,2}^{Dur}$  in (23) can be finally computed.

## 6 Conclusions

The expression of the potential contributed by Hamayun et al. (2009) has been suitably reformulated in order to prove its equivalence with that derived in Holstein (2003).

Furthermore, the volume integral representing the linear part of the potential has been transformed by Gauss theorem to a surface integral which is different from the expression jointly used by Holstein (2003) and Hamayun et al. (2009).

When specialized to polyhedral bodies, one obtains the linear part of the potential as sum of two quantities one of which, namely  $U_{l_1}$  in (23), coincides with existing solutions.

However, differently from the solution in Holstein (2003), also  $U_{l_2}^{Dur}$  is expressed as the product of two terms, one of which is the distance  $d_i$  existing between the generic face and the observation point  $P$ . Hence, the expression (23) extends to the linear part of the potential a property already known in the constant density case (Holstein 2002b).

Since the same property is already known to hold for  $I_{c,F_i}$ , which is required for computing  $U_c$  in (3), we can conclude that the potential  $U(P)$  can be given an expression that factors the vanishing pre-multiplier for faces containing  $P$ .

**Acknowledgements** The author is particularly indebted to prof. Horst Holstein, who acted as one of the reviewers, for his continuous support, very useful suggestions and constructive criticism of the earlier versions of the manuscript. The contribution of the Editor-in-Chief, Ph.D. Pascal Willis, of the Associate Editor, Ph.D. Robert Cunderlik, and of two further anonymous reviewers is also gratefully acknowledged.

## References

- Chai Y, Hinze WJ (1988) Gravity inversion of an interface above which the density contrast varies exponentially with depth. *Geophysics* 53:837–845
- D’Urso MG (2012) New expressions of the gravitational potential and its derivatives for the prism. In: Sneeuw N, Novák P, Crespi M, Sansò F (eds) VII Hotine-Marussi international symposium on mathematical geodesy. Springer, Berlin/Heidelberg
- D’Urso MG (2013) On the evaluation of the gravity effects of polyhedral bodies and a consistent treatment of related singularities. *J Geodesy* 87:239–252
- D’Urso MG (2014) Analytical computation of gravity effects for polyhedral bodies. *J Geodesy* 88:13–29
- D’Urso MG (2015a) Gravity effects of polyhedral bodies with linearly varying density. *Celest Mech Dyn Astron* 120:349–372
- D’Urso MG (2015b) The gravity anomaly of a 2D polygonal body having density contrast given by polynomial functions. *Surv Geophys* 36:391–425
- D’Urso MG, Marmo F (2013) On a generalized Love’s problem. *Comput Geosci* 61:144–151
- D’Urso MG, Marmo F (2015) Vertical stress distribution in isotropic half-spaces due to surface vertical loadings acting over polygonal domains. *Zeit Angew Math Mech* 95:91–110
- D’Urso MG, Russo P (2002) A new algorithm for point-in polygon test. *Surv Rev* 284:410–422
- D’Urso MG, Trotta S (2015) Comparative assessment of linear and bilinear prism-based strategies for terrain correction computations. *J Geodesy* 89:199–215
- Gallardo-Delgado LA, Perez-Flores MA, Gomez-Trevino E (2003) A versatile algorithm for joint inversion of gravity and magnetic data. *Geophysics* 68:949–959
- García-Abdeslem J (1992) Gravitational attraction of a rectangular prism with depth dependent density. *Geophysics* 57:470–473
- García-Abdeslem J (2005) Gravitational attraction of a rectangular prism with density varying with depth following a cubic polynomial. *Geophysics* 70:J39–J42
- Hamayun P, Prutkin I, Tenzer R (2009) The optimum expression for the gravitational potential of polyhedral bodies having a linearly varying density distribution. *J Geodesy* 83:1163–1170
- Hansen RO (1999) An analytical expression for the gravity field of a polyhedral body with linearly varying density. *Geophysics* 64:75–77
- Holstein H (2002a) Gravimagnetic similarity in anomaly formulas for uniform polyhedra. *Geophysics* 67:1126–1133
- Holstein H (2002b) Invariance in gravimagnetic anomaly formulas for uniform polyhedra. *Geophysics* 67:1134–1137
- Holstein H (2003) Gravimagnetic anomaly formulas for polyhedra of spatially linear media. *Geophysics* 68:157–167
- Holstein H, Ketteridge B (1996) Gravimagnetic similarity in anomaly formulas for uniform polyhedra. *Geophysics* 61:1126–1133
- Holstein H, Schürholz P, Starr A, Chakraborty M (1999) Comparison of gravimetric formulas for uniform polyhedra. *Geophysics* 64:1438–1446
- Marmo F, Rosati L (2015) A general approach to the solution of Boussinesq’s problem for polynomial pressures acting over polygonal domains. *J Elast*. doi:10.1007/s10659-015-9534-5
- Pohánka V (1988) Optimum expression for computation of the gravity field of a homogeneous polyhedral body. *Geophys Prospect* 36:733–751
- Pohánka V (1998) Optimum expression for computation of the gravity field of a polyhedral body with linearly increasing density. *Geophys Prospect* 46:391–404
- Rosati L, Marmo F (2014) Closed-form expressions of the thermo-mechanical fields induced by a uniform heat source acting over an isotropic half-space. *Int J Heat Mass Transfer* 75:272–283

- Sessa S, D'Urso MG, (2013) Employment of Bayesian networks for risk assessment of excavation processes in dense urban areas. In: Proceedings of 11th International Conference on ICOSAR 2013, pp 30163–30169
- Strakhov VN, Lapina MI, Yefimov AB (1986) A solution to forward problems in gravity and magnetism with new analytical expression for the field elements of standard approximating body. *Izv Earth Sci* 22:471–482
- Tang KT (2006) *Mathematical methods for engineers and scientists*. Springer, Berlin/Heidelberg/New York
- Zhou X (2009) 3D vector gravity potential and line integrals for the gravity anomaly of a rectangular prism with 3D variable density contrast. *Geophysics* 74:143–153

---

# The Observation Equation of Spirit Leveling in Molodensky's Context

B. Betti, D. Carrion, F. Sacerdote, and G. Venuti

---

## Abstract

Spirit leveling and surface gravity observations can be expressed as orthometric height differences plus corrections which require the knowledge of the Earth crust density. For leveling increments we can write observation equations in a linearized form, according to the standard Molodensky approach, i.e., intrinsic geodesy, depending on normal height differences plus a correction term. This term is a function of the gravity anomaly at the surface level, thus not requiring any assumption on the crust density, and of the curvature of the normal field force lines, which can not be neglected for leveling profiles directed along meridians. The present work shows how to derive these observation equations. An example is presented to verify the effectiveness of the new approach.

---

## Keywords

Normal correction • Normal heights • Spirit leveling observation equation

---

## 1 Introduction

National height systems are referred to an official vertical datum and adopt different coordinate types, such as orthometric, dynamic or normal heights. A complete review of all the systems in use can be found, for instance, in Jekeli (2000) and Meyer et al. (2006). When derived from leveling measurements, all these heights need to be properly corrected using gravity information. Orthometric heights have a precise geometric meaning, being defined as the lengths of the plumb lines from the Earth's surface to the geoid, which are both physical surfaces. Yet, the computation of orthometric corrections requires the knowledge of the Earth's crust density, which can be established only approximately. Their derivation is illustrated in Heiskanen and Moritz (1967), in

Chapter 1 of Sansò and Sideris (2013) and discussed in detail in Sansò and Vaníček (2006). Dynamic heights are simply proportional to geopotential numbers and, contrary to orthometric heights, have no direct physical meaning; the computation of the corresponding corrections is quite simple and requires only gravity data on the Earth's surface. Normal heights are defined as the separation between the ellipsoid and the telluroid, whose points satisfy the relation  $U(P_T) = W(P_S)$ , where  $W$  is the geopotential,  $U$  is the normal potential,  $P_S$  is on the Earth's surface,  $P_T$  is on the telluroid and lies on the ellipsoidal normal through  $P_S$ . Therefore, for a given latitude and longitude  $(\varphi, \lambda)$ , the normal height  $h^*$  is implicitly related to the ellipsoidal height  $h$  by the equation  $U(\varphi, \lambda, h^*) = W(\varphi, \lambda, h)$ . Owing to this definition, normal heights have a geometric meaning, but the surfaces involved are not defined in terms of physical quantities. The most direct formula for their computation is  $h^* = \frac{C}{\bar{\gamma}}$ , where  $C$  is the geopotential number and  $\bar{\gamma}$  is the average of the normal gravity, along the ellipsoidal normal, between the reference ellipsoid and the telluroid. The computation, therefore, requires gravity values on the Earth's surface only. Even in the case of normal heights it may be

---

B. Betti • D. Carrion • G. Venuti (✉)  
DICEA, Politecnico di Milano, P.zza Leonardo da Vinci 32, Milan, Italy  
e-mail: giovanna.venuti@polimi.it

F. Sacerdote  
DICEA, Università degli Studi, Firenze, Via Santa Marta 3, Florence, Italy

interesting to look for the expression of a correction term, enabling us to compare directly normal height differences to leveling measurements. Such a term, called normal correction, is already presented in Heiskanen and Moritz (1967). The formula is derived in analogy to that of orthometric height differences by substituting the normal gravity field to the actual one and approximating lengths along the line of force of the normal potential with those along the normal to the reference ellipsoid. Our approach is different. We start from leveling observations and derive the related observation equations by linearizing the actual gravity field according to Molodensky's approach. Normal height differences come into the formulas as a result and prove to be the natural height coordinates to be used in connection with spirit leveling and surface gravity data. Moreover, the normal correction term, which is different from the one in Heiskanen and Moritz (1967), comes out to depend on the gravity anomaly on the ground and on the curvature of the normal field force lines, which cannot be neglected especially for leveling profiles directed along meridians.

We derive the new observation equations and review the ones presented in Heiskanen and Moritz (1967) highlighting the differences in Sect. 2.

A first check of the new equations has been performed along a leveling line in the western Alps already used for altimetry tests by Gentile et al. (2011) and Barzaghi et al. (2014). The results are reported in Sect. 3, followed by some concluding remarks in Sect. 4.

## 2 Spirit Leveling and Surface Gravity Observation Equation

The introduction of normal heights has to be framed within Molodensky's approach in the linearization of the gravity field boundary value problem. The same coordinates come out when modeling the spirit leveling increments, by expressing the actual potential as the sum of normal plus anomalous one, and the gravity acceleration as the normal counterpart plus the gravitational disturbance, as we show hereafter in more detail. We start from the infinitesimal leveling increment (cf. Sansò and Sideris 2013, Chapter 1):

$$\delta L = \mathbf{n} \cdot d\mathbf{r} = -\frac{\mathbf{g}}{g} \cdot d\mathbf{r} = -\frac{dW}{g} \quad (1)$$

where  $\mathbf{n}$  is the unit vector tangent to the gravity field plumb lines,  $d\mathbf{r}$  is the infinitesimal vector between two points in space along the leveling line,  $\mathbf{g}$  is the gravity vector and  $g$  is its modulus. We set

$$dW = dU + dT \quad (2)$$

where  $W, U, T$  are the actual, the normal and the anomalous gravity potential, respectively. Furthermore,

$$\frac{1}{g} = \frac{1}{\gamma + \delta g} \approx \frac{1}{\gamma} - \frac{\delta g}{\gamma^2} \quad (3)$$

where  $\gamma$  is the modulus of the normal gravity vector, and  $\delta g = g - \gamma$  is the gravity disturbance. Hence, neglecting higher order terms

$$-\frac{dW}{g} = -\frac{dU}{\gamma} + \frac{\delta g}{\gamma^2} dU - \frac{dT}{\gamma}. \quad (4)$$

The contribution of these three terms in the integral along the leveling line  $l_{AB}$  will be separately investigated.

As for the first term, we have:

$$-\int_{l_{AB}} \frac{dU}{\gamma} = -\int_{l_{AB}} \frac{\boldsymbol{\gamma}}{\gamma} \cdot d\mathbf{r} = \int_{l_{AB}} \mathbf{e}_\gamma \cdot d\mathbf{r} \quad (5)$$

where  $\mathbf{e}_\gamma$  does not coincide with the ellipsoidal normal  $\mathbf{v}$ , due to the curvature of the normal gravity force lines, whose expression in geodetic coordinates is (cf. Sansò and Sideris 2013, Chapter 1)

$$k = \frac{1}{\gamma(M+h)} \frac{\partial \gamma}{\partial \varphi} \quad (6)$$

where  $h$  is the ellipsoidal height,  $M = \frac{a(1-e^2)}{(1-e^2 \sin^2 \varphi)^{3/2}}$  is the reference ellipsoid meridian curvature radius,  $a$  is the ellipsoid semi-major axis and  $e$  its eccentricity. We write the unit vector  $\mathbf{e}_\gamma$  at a point P as the sum of the constant unit vector  $\mathbf{v}$ , orthogonal to the reference ellipsoid, plus an additive vector  $\boldsymbol{\delta}$  depending on  $h$  and practically directed along the meridian:

$$\mathbf{e}_\gamma(h) = \mathbf{v} + \boldsymbol{\delta} \approx \mathbf{v} + \delta \mathbf{e}_\varphi \quad (7)$$

where  $\mathbf{e}_\varphi$  is the unit vector tangent to the meridian and

$$\boldsymbol{\delta} = kh = \frac{1}{\gamma(M+h)} \frac{\partial \gamma}{\partial \varphi} h \quad (8)$$

By substituting Eq. (7) in Eq. (5) and using Eq. (8) we obtain:

$$\begin{aligned} \int_{l_{AB}} -\frac{dU}{\gamma} &= \int_{l_{AB}} \mathbf{v} \cdot d\mathbf{r} + \int_{l_{AB}} \boldsymbol{\delta} \cdot d\mathbf{r} = \\ &= h_B - h_A + \int_{l_{AB}} \frac{1}{\gamma} \frac{\partial \gamma}{\partial \varphi} h d\varphi \end{aligned} \quad (9)$$

where we have used the relation

$$\mathbf{e}_\varphi \cdot d\mathbf{r} = (M + h)d\varphi. \quad (10)$$

As for the second term in Eq. (4) we have:

$$\begin{aligned} \int_{l_{AB}} \frac{\delta g}{\gamma^2} dU &= \int_{l_{AB}} \frac{\delta g}{\gamma} \frac{\boldsymbol{\nu}}{\gamma} \cdot d\mathbf{r} = \\ &= - \int_{l_{AB}} \frac{\delta g}{\gamma} (\boldsymbol{\nu} + \boldsymbol{\delta}) \cdot d\mathbf{r} \cong - \int_{l_{AB}} \frac{1}{\gamma} \frac{\partial T}{\partial h} dh. \end{aligned} \quad (11)$$

As in fact  $\frac{\delta g}{\gamma}$  is at the most of the order of  $10^{-4}$ , the term

$\int_{l_{AB}} \frac{\delta g}{\gamma} \boldsymbol{\delta} \cdot d\mathbf{r}$  is clearly negligible while  $\frac{\delta g}{\gamma} \boldsymbol{\nu} \cdot d\mathbf{r}$  is equal to  $\frac{\partial T}{\partial h} dh$ .

As for the third term in Eq. (4), an integration by parts yields:

$$\begin{aligned} - \int_{l_{AB}} \frac{dT}{\gamma} &= - \left[ \frac{T}{\gamma} \right]_A^B - \int_{l_{AB}} \frac{T}{\gamma^2} d\gamma = \\ &= \zeta_A - \zeta_B - \int_{l_{AB}} \frac{T}{\gamma} \frac{1}{\gamma} \frac{\partial \gamma}{\partial \varphi} d\varphi - \int_{l_{AB}} \frac{T}{\gamma} \frac{1}{\gamma} \frac{\partial \gamma}{\partial h} dh = \\ &= \zeta_A - \zeta_B - \int_{l_{AB}} \zeta \frac{1}{\gamma} \frac{\partial \gamma}{\partial \varphi} d\varphi - \int_{l_{AB}} \frac{T}{\gamma} \frac{1}{\gamma} \frac{\partial \gamma}{\partial h} dh \end{aligned} \quad (12)$$

where  $\zeta = \frac{T}{\gamma}$  is the height anomaly, that is the difference between the ellipsoidal and the normal height of a point.

Adding up Eqs. (9), (11) and (12), one obtains the spirit leveling observation equation between the two extreme points  $A$  and  $B$ , along the line  $l_{AB}$ :

$$\begin{aligned} \Delta_{AB}L &= (h_B - \zeta_B) - (h_A - \zeta_A) + \\ &+ \int_{l_{AB}} (h - \zeta) \frac{1}{\gamma} \frac{\partial \gamma}{\partial \varphi} d\varphi - \int_{l_{AB}} \frac{1}{\gamma} \left( \frac{\partial T}{\partial h} - \frac{1}{\gamma} \frac{\partial \gamma}{\partial h} T \right) dh = \\ &= h_B^* - h_A^* + \int_{l_{AB}} h^* \frac{\partial \gamma}{\partial \varphi} \frac{1}{\gamma} d\varphi - \int_{l_{AB}} \frac{\Delta g}{\gamma} dh = \\ &= h_B^* - h_A^* + NC_1 + NC_2 \end{aligned} \quad (13)$$

where  $h^* = h - \zeta$  is the normal height, and the fundamental equation of the physical geodesy

$$\Delta g = - \frac{\partial T}{\partial h} + \frac{1}{\gamma} \frac{\partial \gamma}{\partial h} T \quad (14)$$

has been used.

Equation (13) gives the correction term that must be applied to the observed leveling differences to obtain normal height differences.

The contribution of the first integral can add up to a value of 5 cm for a 30 km long path in the direction of the meridian, at a height  $h = 2,000$  m, as can be derived by the following numerical expression for  $\delta$  [cf. Eq. (8)]:

$$\frac{h}{M+h} \frac{1}{\gamma} \frac{\partial \gamma}{\partial \varphi} \cong 5 \cdot 10^{-3} \frac{h}{M+h} \sin 2\varphi \quad (15)$$

obtained from the approximation of Cassinis' formula (cf. Sansò and Sideris 2013, Chapter 1) for the normal gravity vector modulus:

$$\gamma(\varphi, h) \cong 978.0327715(1 + 5.27448 \cdot 10^{-3} \sin^2 \varphi) - 0.30877h \quad (16)$$

The second integral, as  $\frac{\Delta g}{\gamma}$  can reach a maximum value of about  $10^{-4}$ , can give a contribution of about 20 cm for a height difference of 2,000 m.

Similar considerations show that even a large uncertainty, say 10 m, in the height along the leveling line, has a negligible effect on the numerical evaluation of the correction terms.

It is worth making a comparison with the model for normal height differences (Heiskanen and Moritz 1967):

$$\begin{aligned} h_B^* - h_A^* &= \Delta_{AB}L + \\ &+ \int_{l_{AB}} \frac{g - \gamma_0}{\gamma_0} \delta L + \frac{\bar{\gamma}_A - \gamma_0}{\gamma_0} h_A^* - \frac{\bar{\gamma}_B - \gamma_0}{\gamma_0} h_B^* \end{aligned} \quad (17)$$

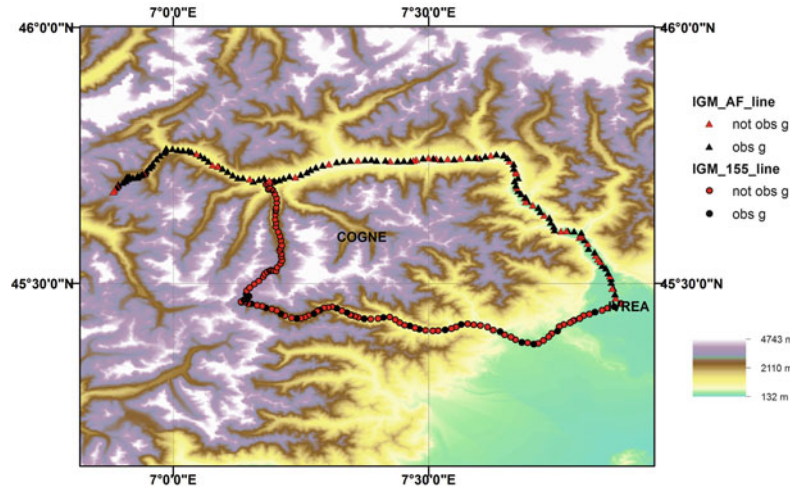
derived in analogy with the formula for orthometric height differences:

$$\begin{aligned} H_B - H_A &= \Delta_{AB}L + \\ &+ \int_{l_{AB}} \frac{g - \gamma_0}{\gamma_0} \delta L + \frac{\bar{g}_A - \gamma_0}{\gamma_0} H_A - \frac{\bar{g}_B - \gamma_0}{\gamma_0} H_B \end{aligned} \quad (18)$$

where  $\gamma_0$  is a conventional normal gravity at a given latitude, for instance  $45^\circ$ .

We start defining the orthometric height  $H$  of a point  $P$  in terms of the corresponding geopotential number. This is obtained by integrating the actual potential differential along the plumb line through the point  $P$ , from the geoid  $P_G$  to the Earth's surface itself:

$$\begin{aligned} C_P &= W(P_G) - W(P) = - \int_{P_G}^P dW = \int_{P_G}^P g dH = \\ &= H \frac{1}{H} \int_{P_G}^P g dH = H \bar{g} \end{aligned} \quad (19)$$



**Fig. 1** IGM leveling lines used for a first check of the new observation equation

and therefore

$$H = \frac{C}{g} \quad (20)$$

In the same manner, by referring to the normal potential and disregarding the normal potential plumb lines curvature, that is by considering

$$-\frac{\gamma}{\gamma} \cong \nu, \quad (21)$$

one can define the normal height of the point P. We have:

$$\begin{aligned} C_P &= W(P_G) - W(P) = U(P_E) - U(P_T) = \\ &= - \int_{P_E}^{P_T} dU \cong \int_{P_E}^{P_T} \gamma dh = h^* \frac{1}{h^*} \int_{P_E}^{P_T} \gamma dh = h^* \bar{\gamma} \end{aligned} \quad (22)$$

and therefore:

$$h^* = \frac{C_P}{\bar{\gamma}}. \quad (23)$$

While Eq. (20) requires the knowledge of the gravity along the plumb line from the geoid to the Earth's surface, which in turn depends on the unknown crust density, Eq. (23) does not. On the other hand, although the normal height has a geometric meaning, it does not refer to physical surfaces. As it will be used in the following, we recall also the dynamic height:

$$h_P^{dyn} = \frac{C_P}{\gamma_0} \quad (24)$$

which is the easiest to compute, but has no geometrical meaning.

Equations (17) and (18) can be easily obtained in the following way. Let us consider two points  $A$  and  $B$  on the Earth's surface. The difference of their dynamic heights can be written as the difference of orthometric or normal heights plus a correction term.

We have:

$$\begin{aligned} h_B^{dyn} - h_A^{dyn} &= \frac{C_B - C_A}{\gamma_0} = \frac{W_A - W_B}{\gamma_0} = \\ &= - \int_A^B \frac{dW}{\gamma_0} = \int_A^B \frac{g}{\gamma_0} \delta L = \int_A^B \frac{g - \gamma_0 + \gamma_0}{\gamma_0} \delta L = \\ &= \Delta_{AB} L + \int_A^B \frac{g - \gamma_0}{\gamma_0} \delta L. \end{aligned} \quad (25)$$

On the other hand,

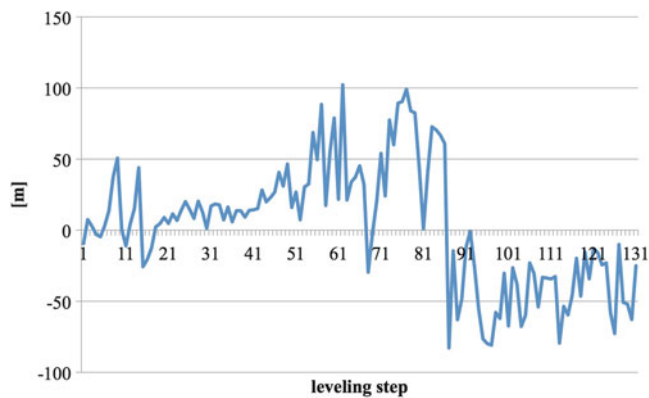
$$\begin{aligned} h_B^{dyn} - h_A^{dyn} &= \frac{C_B}{\gamma_0} - \frac{C_A}{\gamma_0} \\ &= \int_{B_G}^B \frac{g - \gamma_0 + \gamma_0}{\gamma_0} dH - \int_{A_G}^A \frac{g - \gamma_0 + \gamma_0}{\gamma_0} dH = \\ &= H_B - H_A + \int_{B_G}^B \frac{g - \gamma_0}{\gamma_0} dH - \int_{A_G}^A \frac{g - \gamma_0}{\gamma_0} dH = \\ &= H_B - H_A + \frac{\bar{g}_B - \gamma_0}{\gamma_0} H_B - \frac{\bar{g}_A - \gamma_0}{\gamma_0} H_A \end{aligned} \quad (26)$$

where  $A_G$  and  $B_G$  are on the geoid and the integrals are computed along the plumb lines between  $A_G$  and  $A$  and  $B_G$  and  $B$ . Finally, by comparing Eqs. (25) and (26) one finds Eq. (18). With the same approximation of Eq. (21), substituting in Eq. (26) the normal potential to the actual one as in Eq. (22) one obtains Eq. (17).

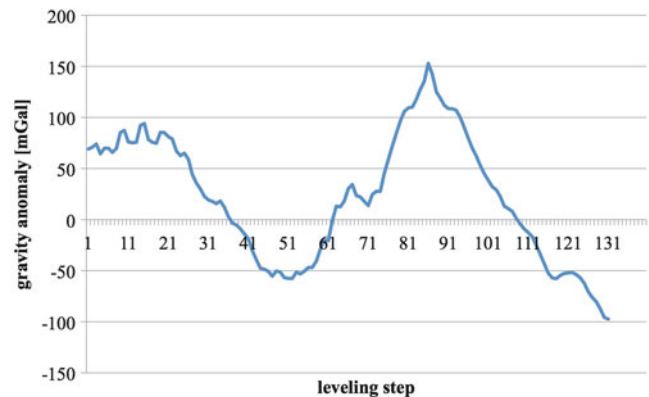
### 3 Tests on the Western Alps

A first assessment of the new observation equation was performed on the leveling line displayed in Fig. 1, which is located in the Italian western Alps. This line is divided in two



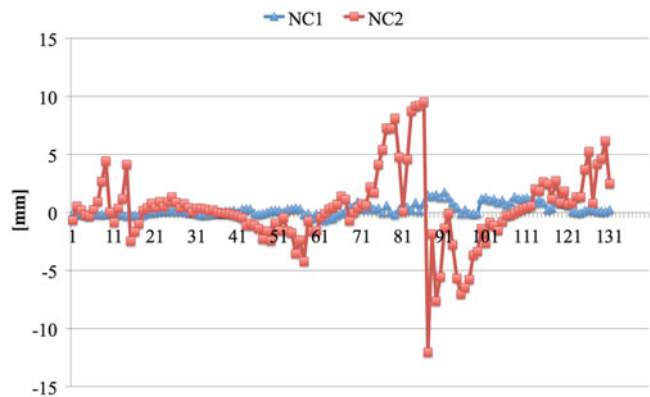


**Fig. 2** Leveling differences along line 155



**Fig. 3** Gravity anomalies along line 155

branches: lines 155 and AF, respectively marked with circles and triangles. In the following, for the sake of brevity, we will show only the results related to branch 155. This line, starting at about 250 m, reaches about 2,600 m. Its length is about 120 km and it contains about 130 benchmarks. A small number of benchmarks, all in the part of the line directed from East to West, are coupled with observed gravity data (cf. Fig. 1). The missing values were predicted from the Italian gravity model (cf. Barzaghi et al. 2014). Leveling differences, reported in Fig. 2, and surface gravity values (observed and predicted) were used to compute the normal height of each benchmark by discretizing Molodensky's observation equation (13). Gravity anomalies along the line are plotted in Fig. 3. In the test area, due to the presence of the Ivrea body,  $\Delta g$  undergoes high variations, much higher than the global standard deviation of 30 mGal (cf. Pavlis N., Global Gravitational Models, in Sansò and Sideris 2013), as well as than the Italian one of about 64 mGal (cf. Barzaghi et al. 2007). The correction terms behavior is plotted in Fig. 4. The term accounting for the normal field force line curvature,  $NC_1$ , as reported in Table 1, has a range of about 3 mm along the line and a total value of 40 mm; the main part of it, almost 37 mm out of 40 mm, is related to the part of the leveling line lying along the meridian. The term



**Fig. 4** Normal correction terms in Molodensky's observation equation along line 155.  $NC_1$  is the term accounting for the normal vertical curvature,  $NC_2$  is the term depending on the gravity anomaly

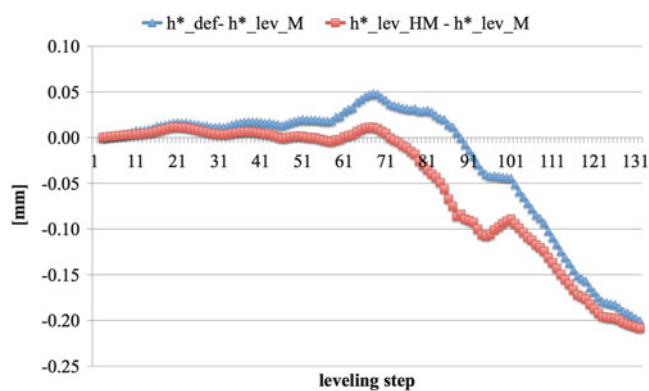
**Table 1** Statistics of the new equation correcting terms:  $NC_1$  is the term depending on the normal vertical lines curvature,  $NC_2$  is the term depending on the gravity anomaly.  $NC = -NC_1 + NC_2$  is the total correcting term

	Mean [mm]	Std [mm]	Total amount [mm]
$NC_1$	0.3	0.5	40.5
$NC_2$	0.4	3.2	46.6
$NC$	0.05	3.3	6.0

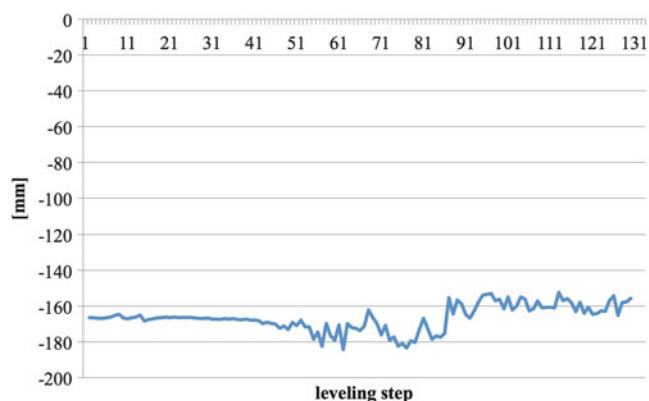
depending on the gravity anomaly,  $NC_2$ , has a much higher variation than the first one: its range is almost 20 mm and its total amount is of 46.6 mm. Normal heights computed with the new equation were compared against those derived by using the definition of Eq. (23) and the classical formula in Eq. (17). The differences, whose maximum value is around 0.2 mm, are reported in Fig. 5. Although not significantly different from zero, they have an interesting behavior: they increase more along that portion of line 155 which is directed along the meridian, showing the effect of the neglected curvature of the normal field force lines, both in the definition and in the classical observation equation. A final analysis was performed on the difference between orthometric heights, computed by Eq. (18) and normal heights computed by Eq. (13). The difference should be approximately expressed in terms of Bouguer anomaly as follows:

$$h^* - H \cong \frac{\Delta g_{Bouguer}}{\bar{\gamma}} H \quad (27)$$

The residual between the left and the right-hand side of Eq. (27) are plotted in Fig. 6. Apart from a bias of almost 17 cm, the differences standard deviation is equal to 7 mm showing a good agreement despite large variations of Bouguer anomalies in the test area. We have not investigated the reason of this bias, which is probably contained in the geopotential number attributed to the starting point of the leveling line, but this is not central to the present work. A last



**Fig. 5** Differences between normal heights computed with the new equation ( $h^*_lev\_M$ ) and those computed using the definition ( $h^*\_def$ ) or the classical equations ( $h^*\_lev\_HM$ )



**Fig. 6** Line 155: residual between the left and the right-hand side of Eq. (27)

remark is worth about the use of normal heights. As already shown in Barzaghi et al. (2014), the misclosure error of the leveling loop of Fig. 1, which amounts to 7 cm without any correction, can be reduced to 3.7 cm by applying orthometric corrections and to 2.7 cm by applying normal corrections either the classical or the new ones. Beyond the effectiveness of corrections, this result shows also the advantage of using observation equations not requiring any assumption on the Earth's crust density.

## 4 Conclusions

A new formulation of spirit leveling observation equation in the framework of Molodensky's approach has been presented. Normal height differences come naturally into the linearized equation, proving to be the proper coordinates to be used in connection with spirit leveling and surface gravity

data when no assumptions on crust density are required. The formula involves two correction terms, one depending on the curvature of the normal gravity field force lines, the other on the gravity anomaly. The contribution of the first term is not negligible for leveling lines directed along meridians, while the second term is particularly relevant for lines at high altitude, especially in zones with large gravity anomaly values, like in the western Alps. The first term highlights the dependence on latitude of the normal correction as already observed by other authors, such as Marti (2002) and Filmer et al. (2010). Normal heights derived by using the new formula are practically equivalent to those derived by using the definition or the classical normal correction; the systematic behavior of the differences is due to the approximation of lengths of arcs of the normal gravity field force lines with segments along the ellipsoidal normal, both in the normal height definition and in the classical normal correction. Moreover, despite the large variation of Bouguer anomalies in the test region, the approximate relation in Eq. (27) is essentially satisfied, but for a bias, whose origin should be better investigated.

## References

- Barzaghi R, Borghi A, Carrion D, Sona G (2007) Refining the estimate of the Italian quasi-geoid, *Bollettino di Geodesia e Scienze Affini*, Firenze, Anno LXXVI 3:145–160
- Barzaghi R, Betti B, Carrion D, Gentile G, Maseroli R, Sacerdote F (2014) Orthometric correction and normal heights for Italian leveling network: a case study. *Appl Geomatics* 6:17–25. doi:10.1007/s12518-013-0121-9
- Filmer MS, Featherstone WE, Kuhn M (2010) The effect of EGM2008-based normal, normal-orthometric and Helmert orthometric height systems on the Australian levelling network. *J Geodesy* 84:501–513
- Gentile G, Maseroli R, Sacerdote F (2011), Studio dell'effetto della gravità su circuiti chiusi della livellazione di alta precisione in presenza di dislivelli molto elevati. In: *Proceedings of ASITA 15–18 Nov 2011, Parma*, pp 1151–1158, ISBN 978-88-903132-6-4
- Heiskanen WA, Moritz H (1967) *Physical geodesy*. W.H. Freeman and Co., San Francisco
- Jekeli C (2000) *Heights, the geopotential and vertical datums*. Ohio State University Report No. 459
- Marti U (2002) Modeling of differences of height systems in Switzerland. In: Tziavos I (ed), *Gravity and geoid 2002*. Proceedings of the 3rd meeting of the international gravity and geoid commission, Thessaloniki, Greece
- Meyer TH, Roman DR, Zilkoski DB (2006) What does height really mean? Part III: height systems. *Surveying Land Inf Sci* 66(2):149–160
- Sansò F, Sideris M (2013) *Geoid determination: theory and methods*. Lecture notes in earth system sciences, vol 110. Springer, New York
- Sansò F, Vaníček P (2006) The orthometric height and the holonomy problem. *J Geodesy* 80(5):225–232

**Theoretical Aspects of Reference Frames**

---

# Reference Station Weighting and Frame Optimality in Minimally Constrained Networks

C. Kotsakis

---

## Abstract

The aim of this paper is to present a general solution of the weight choice problem for the reference stations in minimally constrained network adjustment. Our treatment is based on the optimization of the accuracy of the estimated network coordinates over all possible choices of minimum constraints on the reference stations. The optimal criterion considers the joint effect of the data and datum noise on the estimated coordinates and it is implemented over an arbitrary subset of the network stations. The final solution leads to a flexible treatment of the datum choice problem by allowing the weight matrix of the reference stations to be tuned to various options regarding the frame quality in the adjusted network.

---

## Keywords

Data noise effect • Datum choice problem • Datum noise effect • Frame optimization • Minimum constraints • Network adjustment

---

## 1 Introduction

The datum choice problem (DCP) is a fundamental issue in geodetic network adjustment with coordinate-based models. It is linked to the optimal estimation of a set of coordinates from a singular system of normal equations that are obtained by the network data analysis within a linearized least squares (LS) setting (e.g. Dermanis 1985; Schaffrin 1985; Teunissen 1985). The usual treatment of the DCP requires that a set of external constraints is used to complement the missing datum information in the available data. Several options exist for the selection and the implementation of those constraints into the network adjustment procedure, each of which has its own merits for the final estimated solution. In this study we concentrate on the so-called minimum constraints (MCs) which treat the datum defect of the geodetic network without

interfering with its estimable characteristics from the available data (Sillard and Boucher 2001).

These datum constraints are typically applied over a number of reference stations that are included in the network analysis and have a priori known coordinates with respect to the user's desired reference frame. Their general expression is given in terms of the linear system

$$\mathbf{E}(\mathbf{x} - \mathbf{x}^{\text{ext}}) = \mathbf{0} \quad (1)$$

which corresponds to the well known inner constraints or, more precisely, partial inner constraints since only a part of the network stations is involved in the datum definition process (Meissl 1969; Blaha 1971). The above constraints enforce the harmonization of the (non-estimable) frame parameters of the geodetic network with the respective frame parameters implied by the known coordinates  $\mathbf{x}^{\text{ext}}$  of the reference stations. The matrix  $\mathbf{E}$  stems from the usual Helmert transformation model using only the rows that correspond to the datum defect of the underlying network (Sillard and Boucher 2001).

---

C. Kotsakis (✉)  
Department of Geodesy and Surveying, Aristotle University  
of Thessaloniki, Thessaloniki 54124, Greece  
e-mail: [kotsaki@topo.auth.gr](mailto:kotsaki@topo.auth.gr)

The use of Eq. (1) offers a restrictive optimality for the estimated coordinates in the desired frame. In fact, the fundamental property of the network solution under these constraints is the minimization of the propagated data noise on the estimated coordinates of the reference stations. This suggests that there are two limitations which Eq. (1) is not able to handle in network adjustment problems, namely:

1. the optimal control of the propagated data noise on other network stations (apart from the reference stations), and
2. the optimal control of the random errors in the a priori reference coordinates  $\mathbf{x}^{\text{ext}}$  and their propagated effect (hereafter called *datum noise effect*) on the network solution.

Both of these issues are crucial in the context of the optimal datum choice for minimally constrained networks and they have to be treated in a more general setting than the one provided by the classic inner constraints in Eq. (1). Specifically, the use of a weight matrix for the reference stations according to the extended form of inner constraints

$$\mathbf{E}\mathbf{P}(\mathbf{x} - \mathbf{x}^{\text{ext}}) = \mathbf{0} \quad (2)$$

enables us to overcome the aforementioned limitations within a zero-order optimization scheme for geodetic networks. This result has been established in Kotsakis (2013) where the choice problem for the weight matrix  $\mathbf{P}$  was tackled on the basis of the joint minimization of the data and datum noise effects over all network stations.

The aim of this paper is to present a useful extension of the previous result by considering the minimization of the data/datum noise effects over *an arbitrary subset* of the network stations. This generalization provides a flexible treatment of the DCP by allowing the weight matrix of the reference stations to be tuned to various options regarding the frame quality of the adjusted network. Essentially, we formulate herein an MC-based scheme for geodetic network adjustment under an optimality principle for the estimated coordinates of any desired group of the network stations.

## 2 Problem Formulation

The general problem that is treated herein can be briefly described as follows. Our starting point is a singular system of normal equations (NEQ)

$$\mathbf{N}(\mathbf{X} - \mathbf{X}^0) = \mathbf{u} \quad (3)$$

which is obtained from the linearized LS adjustment of a geodetic network. It is considered that the rank defect of the above system is solely caused by the datum deficiency in the used data. Without loss of generality, we assume that any nuisance parameters have been eliminated beforehand from

the NEQ system, so that the term  $\mathbf{X} - \mathbf{X}^0$  contains only the unknown corrections to the approximate coordinates of the network stations.

The total coordinate vector in Eq. (3) is partitioned as

$$\mathbf{X} = \begin{bmatrix} \mathbf{x} \\ \mathbf{x}' \end{bmatrix} \quad (4)$$

where  $\mathbf{x}$  refers to the reference stations that are included in the network and  $\mathbf{x}'$  corresponds to the new stations whose coordinates represent the primary unknowns of the estimation problem at hand.

The rationale of our study relies on the exact inversion of the NEQ system using the datum information that is contained in the known coordinates of the reference stations. The usual datum choice is provided by the system of (partial) inner constraints in Eq. (1). The corresponding solution is obtained through the general formula (Koch 1999)

$$\hat{\mathbf{X}} = \mathbf{X}^0 + (\mathbf{N} + \mathbf{H}^T\mathbf{H})^{-1}(\mathbf{u} + \mathbf{H}^T\mathbf{c}) \quad (5)$$

where the constraint matrix  $\mathbf{H}$  and the vector  $\mathbf{c}$  are given by the following expressions

$$\mathbf{H} = [\mathbf{E} \quad \mathbf{0}] \quad (6)$$

$$\mathbf{c} = \mathbf{E}(\mathbf{x}^{\text{ext}} - \mathbf{x}^0) \quad (7)$$

The covariance (CV) matrix of the above solution has the form (Koch 1999)

$$\begin{aligned} \Sigma_{\hat{\mathbf{X}}} &= (\mathbf{N} + \mathbf{H}^T\mathbf{H})^{-1}\mathbf{N}(\mathbf{N} + \mathbf{H}^T\mathbf{H})^{-1} \\ &= \begin{bmatrix} \Sigma_{\hat{\mathbf{x}}} & \Sigma_{\hat{\mathbf{x}}\hat{\mathbf{x}'}} \\ \Sigma_{\hat{\mathbf{x}'\hat{\mathbf{x}}} & \Sigma_{\hat{\mathbf{x}'}} \end{bmatrix} = \mathbf{N}^- \end{aligned} \quad (8)$$

and it corresponds to a generalized inverse of the normal matrix with minimum trace over the reference stations.<sup>1</sup> This is a well known result in network optimization theory upon which the use of inner constraints was introduced in geodetic practice (Blaha 1971); see also Grafarend (1974) and Schmitt (1982).

The CV matrix from Eq. (8) reflects only the data noise effect in the estimated coordinates, thus ignoring the influence of random errors in the known coordinates  $\mathbf{x}^{\text{ext}}$  of the

<sup>1</sup>The minimum-trace property of the error CV submatrix  $\Sigma_{\hat{\mathbf{x}}}$  in Eq. (8) is equivalent to a simple (unweighted) LS fit of the adjusted network to the known coordinates  $\mathbf{x}^{\text{ext}}$  of the reference stations using the Helmert transformation model that involves only the non-estimable frame parameters of the underlying network.

reference stations. The latter introduce a datum-related noise in the estimated solution which reflects the uncertainty of the (non-estimable part of the) coordinate system itself for the adjusted network. The optimal control of this datum noise effect, concurrently with the data noise effect, is an important issue for the minimally constrained network adjustment, however it cannot be handled through the classic inner constraints.

An additional concern stems from the fact that the minimum-trace property of the previous CV matrix refers only to its part related to the reference stations. The coordinates of the new stations are not estimated in an optimal way under the choice of Eq. (1). A worthy enhancement, therefore, is to look for an MC matrix  $\mathbf{Q}$  to replace the classic inner-constraint matrix  $\mathbf{E}$ , so that the revised datum constraints

$$\mathbf{Q}(\mathbf{x} - \mathbf{x}^{\text{ext}}) = \mathbf{0} \quad (9)$$

yield the best accuracy for the estimated coordinates at the new stations (or any selected subset of network stations) with respect to the desired frame which is realized by the reference stations.

A similar version of the above problem was presented in Kotsakis (2013) for the case of the joint minimization of the data/datum noise effects over all network stations. The theoretical investigation in that study showed that the MC matrix should have the factorized form  $\mathbf{Q} = \mathbf{E} \mathbf{P}$ , with  $\mathbf{P}$  being a suitable weight matrix for the reference stations.

Herein we treat the case of minimizing the data/datum noise effects over an arbitrary subset of the network stations. The optimal MC matrix should again have the same factorized form while the weight matrix of the reference stations will have a more general structure than the one given in Kotsakis (2013).

### 3 General Expressions for the CV Matrix of a MC Solution

Before we proceed with the optimal datum choice in minimally constrained networks (more specifically, the optimal choice of the weight matrix for the reference stations), it is instructive to review the various CV matrices involved in the accuracy assessment of the estimated network coordinates.

In general, the total CV matrix of a MC solution can be expressed as a sum of two components

$$\Sigma_{\hat{\mathbf{X}}}^{\text{total}} = \Sigma_{\hat{\mathbf{X}}}^{\text{obs}} + \Sigma_{\hat{\mathbf{X}}}^{\text{mc}} \quad (10)$$

which contain the contributions from separate error sources, that is the data and datum noise effects, respectively. Their

analytic forms are derived through straightforward covariance propagation to Eq. (5) and they are given by the general expressions

$$\Sigma_{\hat{\mathbf{X}}}^{\text{obs}} = (\mathbf{N} + \mathbf{H}^T \mathbf{H})^{-1} \mathbf{N} (\mathbf{N} + \mathbf{H}^T \mathbf{H})^{-1} \quad (11)$$

and

$$\Sigma_{\hat{\mathbf{X}}}^{\text{mc}} = (\mathbf{N} + \mathbf{H}^T \mathbf{H})^{-1} \mathbf{H}^T \Sigma_{\mathbf{c}} \mathbf{H} (\mathbf{N} + \mathbf{H}^T \mathbf{H})^{-1} \quad (12)$$

The previous equations are valid for any network solution that is determined by an arbitrary set of MCs

$$\mathbf{H}(\mathbf{X} - \mathbf{X}^0) = \mathbf{c} \quad (13)$$

where the pseudo-observation vector  $\mathbf{c}$  is associated with a prior CV matrix  $\Sigma_{\mathbf{c}}$ .

Taking into account well known algebraic identities from the MC theory in singular NEQ systems (e.g. Kotsakis 2012), the following equivalent expressions can be also used

$$\Sigma_{\hat{\mathbf{X}}}^{\text{obs}} = (\mathbf{N} + \mathbf{H}^T \mathbf{H})^{-1} - \tilde{\mathbf{E}}^T (\mathbf{H} \tilde{\mathbf{E}}^T)^{-1} (\tilde{\mathbf{E}} \mathbf{H}^T)^{-1} \tilde{\mathbf{E}} \quad (14)$$

and

$$\Sigma_{\hat{\mathbf{X}}}^{\text{mc}} = \tilde{\mathbf{E}}^T (\mathbf{H} \tilde{\mathbf{E}}^T)^{-1} \Sigma_{\mathbf{c}} (\tilde{\mathbf{E}} \mathbf{H}^T)^{-1} \tilde{\mathbf{E}} \quad (15)$$

where  $\tilde{\mathbf{E}}$  denotes the *inner-constraint matrix for the entire network* which, in accordance to the partition of Eq. (4), is expressed as

$$\tilde{\mathbf{E}} = [\mathbf{E} \quad \mathbf{E}'] \quad \text{and} \quad \mathbf{N} \tilde{\mathbf{E}}^T = \mathbf{0} \quad (16)$$

For more details and the mathematical proofs of the preceding equations see Kotsakis (2012, 2013). Note that, for simplicity, the a priori variance factor is assumed to be equal to one.

The CV matrix  $\Sigma_{\hat{\mathbf{X}}}^{\text{obs}}$  is always singular and it contains the effect of the data noise on the MC solution. It corresponds to a reflexive generalized inverse of the normal matrix  $\mathbf{N}$  and its rank defect is equal to the datum defect of the network observational model. The trace minimization of this matrix was used as a criterion for solving the DCP in the context of network optimization theory (e.g. Blaha 1971; Schmitt 1982), thus leading to the classic type of inner constraints for geodetic network adjustment.

The CV matrix  $\Sigma_{\hat{\mathbf{X}}}^{\text{mc}}$  is also singular and it reflects the datum noise effect in the MC solution. In fact, it quantifies the accuracy of the estimated coordinates due to random errors in the pseudo-observation vector  $\mathbf{c}$ . From a geodetic viewpoint, it is a necessary component for the realistic accuracy assessment of frame realizations obtained via minimally constrained networks on a number of reference stations. In

such cases the general MCs of Eq. (13) should be formulated in the partitioned form

$$\underbrace{\begin{bmatrix} \mathbf{Q} & \mathbf{0} \end{bmatrix}}_{\mathbf{H}} \underbrace{\begin{bmatrix} \mathbf{x} - \mathbf{x}^o \\ \mathbf{x}' - \mathbf{x}'^o \end{bmatrix}}_{\mathbf{x} - \mathbf{x}^o} = \underbrace{\mathbf{Q}(\mathbf{x}^{\text{ext}} - \mathbf{x}^o)}_{\mathbf{c}} \quad (17a)$$

or equivalently

$$\mathbf{Q}(\mathbf{x} - \mathbf{x}^{\text{ext}}) = \mathbf{0} \quad (17b)$$

where  $\mathbf{Q}$  is an arbitrary MC matrix to be applied to the reference stations, which will be optimally determined in the next section.

Note that the vector  $\mathbf{c}$  is often set to zero by selecting the approximate coordinates of the reference stations to be equal to their a priori known values in the desired frame. This does not eliminate the datum noise effect which should be always accounted in terms of the matrix  $\Sigma_{\hat{\mathbf{x}}}^{\text{mc}}$  (see Eq. 12 or 15) using the auxiliary covariance expression

$$\Sigma_{\mathbf{c}} = \mathbf{Q} \Sigma_{\mathbf{x}^{\text{ext}}} \mathbf{Q}^T \quad (18)$$

where  $\Sigma_{\mathbf{x}^{\text{ext}}}$  corresponds to the prior CV matrix of the reference station coordinates  $\mathbf{x}^{\text{ext}}$ .

## 4 Optimal Datum Choice in MC Networks: A General Formulation

For the purpose of this study, the optimal datum choice is linked to the minimization of an objective functional that quantifies the accuracy of the estimated coordinates at (all or part of) the network stations. A standard option for this functional is the trace of the total CV matrix  $\Sigma_{\hat{\mathbf{x}}}^{\text{total}}$  which was analytically described in the previous section. Hence, the DCP is formulated in terms of the optimization problem

$$\min_{\mathbf{Q}} \text{tr} \left( \mathbf{S} \Sigma_{\hat{\mathbf{x}}}^{\text{total}} \mathbf{S}^T \right) \quad (19)$$

or, more explicitly

$$\min_{\mathbf{Q}} \text{tr} \left( \mathbf{S} \Sigma_{\hat{\mathbf{x}}}^{\text{obs}} \mathbf{S}^T + \mathbf{S} \Sigma_{\hat{\mathbf{x}}}^{\text{mc}} \mathbf{S}^T \right) \quad (20)$$

where  $\mathbf{Q}$  is the sought MC matrix and  $\mathbf{S}$  corresponds to a “selection matrix” for the participating stations in the optimality principle. Note that the MCs are applied only to the reference stations (see Eq. 17) while the optimality principle may refer to any subset of network stations.

Considering the partition scheme in Eq. (4), some examples of the selection matrix are  $\mathbf{S} = [\mathbf{I} \ \mathbf{0}]$ ,  $\mathbf{S} = [\mathbf{0} \ \mathbf{I}]$ ,  $\mathbf{S} = \mathbf{I}$  which can be used for the accuracy optimization of the estimated coordinates at the reference stations, at the new stations, or at all network stations, respectively.

The MC matrix that satisfies the optimality principle in Eq. (20) can be derived from the equation

$$\frac{\partial \text{tr} \left( \mathbf{S} \Sigma_{\hat{\mathbf{x}}}^{\text{obs}} \mathbf{S}^T \right)}{\partial \mathbf{Q}} + \frac{\partial \text{tr} \left( \mathbf{S} \Sigma_{\hat{\mathbf{x}}}^{\text{mc}} \mathbf{S}^T \right)}{\partial \mathbf{Q}} = \mathbf{0} \quad (21)$$

The dependence of  $\Sigma_{\hat{\mathbf{x}}}^{\text{obs}}$  and  $\Sigma_{\hat{\mathbf{x}}}^{\text{mc}}$  on the matrix  $\mathbf{Q}$  stems from Eqs. (14) and (15) taking also into account the relationships in Eqs. (17) and (18). After some lengthy derivations the solution of the last equation is obtained as

$$\mathbf{Q} = \mathbf{E}(\Sigma + \Sigma_{\mathbf{x}^{\text{ext}}})^{-1} \quad (22)$$

where the matrix  $\Sigma$  is defined by the formula

$$(\mathbf{N} + \mathbf{S}^T \mathbf{S} \tilde{\mathbf{E}}^T \tilde{\mathbf{E}} \mathbf{S}^T \mathbf{S})^{-1} = \left[ \begin{array}{c|c} \Sigma & \mathbf{L} \\ \hline \mathbf{L}^T & \Sigma' \end{array} \right] \quad (23)$$

(the above partitioning is compatible with the one introduced in Eq. 4). The proof of the above result for the case  $\mathbf{S} = \mathbf{I}$  is given in Kotsakis (2013) whereas the proof for an arbitrary selection matrix  $\mathbf{S} \neq \mathbf{I}$  can easily be obtained as a straightforward extension of the derivations given in that paper.

### 4.1 General Remarks

The optimal weight matrix for the reference stations in MC network adjustment has the general form

$$\mathbf{P} = (\Sigma + \Sigma_{\mathbf{x}^{\text{ext}}})^{-1} \quad (24)$$

where  $\Sigma_{\mathbf{x}^{\text{ext}}}$  is the CV matrix of their prior coordinates and  $\Sigma$  is an auxiliary matrix obtained by Eq. (23). We underline that the last expression stems from a formal optimization scheme which has led to the factorized form of Eq. (22), thus proving that the weighted inner constraints is indeed the appropriate tool to ensure special optimal properties for the realized frame in a minimally constrained network.

The weight matrix  $\mathbf{P}$  depends on two components each of which has a distinct role in the MC network adjustment. The first component is responsible for minimizing the propagated data noise on the estimated coordinates of a group of network

stations that is specified by the selection matrix  $\mathbf{S}$ . The second component, on the other hand, is related to the filtering of the random errors in the known reference coordinates from the final network solution. This dual role of the weight matrix is dictated by the joint presence of the data and datum noise effects, both of which influence in their own way the frame quality in the adjusted network.

## 4.2 Minimization of the Data Noise Effect

If we ignore the random errors in the known reference coordinates (i.e.  $\Sigma_{\mathbf{x}}^{\text{ext}} = \mathbf{0}$ ) then the weighted MCs take the form

$$\mathbf{E} \Sigma^{-1} (\mathbf{x} - \mathbf{x}^{\text{ext}}) = \mathbf{0} \quad (25)$$

and the resulting network solution will have the minimum data noise effect at the stations specified by the selection matrix  $\mathbf{S}$  which is hidden in the weight matrix  $\Sigma^{-1}$ ; see Eq. (23).

If the selection matrix involves only the reference stations of the underlying network, that is  $\mathbf{S} = [\mathbf{I} \ \mathbf{0}]$ , then it can be shown that Eq. (25) is reduced to the form

$$\mathbf{E} \mathbf{E}^T \mathbf{E} (\mathbf{x} - \mathbf{x}^{\text{ext}}) = \mathbf{0} \quad (26)$$

which, due to the invertibility of the matrix  $\mathbf{E} \mathbf{E}^T$ , is equivalent to

$$\mathbf{E} (\mathbf{x} - \mathbf{x}^{\text{ext}}) = \mathbf{0} \quad (27)$$

In this special case, therefore, we reproduce the classic (unweighted) inner constraints whose optimality is solely related to the minimization of the data noise effect at the reference stations.

## 4.3 Minimization of the Datum Noise Effect

If we consider the datum noise minimization in the MC solution without accounting for the data noise effect, that is

$$\min_{\mathbf{Q}} \text{tr} \left( \mathbf{S} \Sigma_{\widehat{\mathbf{x}}}^{\text{mc}} \mathbf{S}^T \right) \quad (28)$$

then the optimal datum choice will be provided by the weighted MCs

$$\mathbf{E} (\Sigma_{\mathbf{x}}^{\text{ext}})^{-1} (\mathbf{x} - \mathbf{x}^{\text{ext}}) = \mathbf{0} \quad (29)$$

In this case the weight matrix of the reference stations is independent of the selected stations that participate in the

optimality principle. In contrast to the data noise effect, the minimization of the datum noise effect, over all or part of the network, requires a fixed weighting of the reference stations in terms of their prior CV matrix. In fact, if we take into account Eqs. (15)–(18) then we obtain the covariance decomposition formulae

$$\Sigma_{\widehat{\mathbf{x}}}^{\text{mc}} = \tilde{\mathbf{E}}^T \Sigma_{\theta} \tilde{\mathbf{E}} \quad (30)$$

and

$$\Sigma_{\theta} = (\mathbf{Q} \mathbf{E}^T)^{-1} \mathbf{Q} \Sigma_{\mathbf{x}}^{\text{ext}} \mathbf{Q}^T (\mathbf{E} \mathbf{Q}^T)^{-1} \quad (31)$$

where the matrix  $\Sigma_{\theta}$  describes the accuracy with which the (non-estimable) frame parameters are defined in the minimally constrained network. The datum noise minimization is equivalent to the trace minimization of  $\Sigma_{\theta}$  (see Kotsakis 2013) and it will not be affected by the selection matrix  $\mathbf{S}$  that appears in Eq. (28).

## 5 Conclusions

The weight choice problem of the reference stations in minimally constrained networks has been investigated in this paper. Our treatment is based on the optimization of the total accuracy (considering both the data and datum noise effects) of the estimated coordinates over all possible choices of MCs on the reference stations. As a result of this procedure, we obtained a weighted type of MCs which encompasses the classic (unweighted) inner constraints as a special option within a more general setting for the datum choice problem.

In contrast to Kotsakis (2013) the current treatment allows the accuracy optimization of the minimally constrained solution over an arbitrary subset of network stations and not necessarily over the entire network. Hence, the weight matrix of the reference stations is not generally unique since it depends on the network stations that participate in the optimality principle of Eq. (20). It is noted that the use of (the inverse of) the prior CV matrix of the known reference coordinates as a weight matrix is warranted only for minimizing the datum noise effect in the estimated network coordinates – it does not contribute to the optimal control of the data noise effect over all or part of the network.

A useful extension of the present study is the treatment of the weight choice problem for the reference stations in the case of non-minimal datum constraints. In our current approach the matrix  $\mathbf{E}$  refers only to the non-estimable frame parameters and the implementation of the datum constraints  $\mathbf{Q} (\mathbf{x} - \mathbf{x}^{\text{ext}}) = \mathbf{0}$  does not affect any estimable frame characteristics in the underlying network. The case where the constraint matrix  $\mathbf{Q}$  refers also to estimable frame



parameters (e.g. scale in the case of GNSS networks) is more complicated since several properties that have been used in this paper's algebraic derivations will simply not hold true.

---

## References

- Blaaha G (1971) Inner adjustment constraints with emphasis on range observations. Department of Geodetic Science, The Ohio State University, OSU Report No. 148, Columbus, Ohio
- Dermanis A (1985) Optimization problems in geodetic networks with signals. In: Grafarend EW, Sanso F (eds) Optimization and design of geodetic networks. Springer, Berlin, pp 221–256
- Grafarend EW (1974) Optimization of geodetic networks. *Boll Geod Sci Affi XXXIII*:351–406
- Koch K-R (1999) Parameter estimation and hypothesis testing in linear models, 2nd edn. Springer, Berlin
- Kotsakis C (2012) Reference frame stability and nonlinear distortion in minimum-constrained network adjustment. *J Geod* 86(9):755–774
- Kotsakis C (2013) Generalized inner constraints for geodetic network densification problems. *J Geod* 87(7):661–673
- Meissl P (1969) Zusammenfassung und Ausbau der inneren Fehlertheorie eines Punkthaufens. Deutsche Geodätische Kommission Reihe A 61:8–21
- Schaffrin B (1985) Aspects of network design. In: Grafarend EW, Sanso F (eds) Optimization and design of geodetic networks. Springer, Berlin, pp 549–597
- Schmitt G (1982) Optimization of geodetic networks. *Rev Geophys Space Phys* 20(4):877–884
- Sillard P, Boucher C (2001) A review of algebraic constraints in terrestrial reference frame datum definition. *J Geod* 75:63–73
- Teunissen P (1985) Zero order design: generalized inverse, adjustment, the datum problem and S-transformations. In: Grafarend EW, Sanso F (eds) Optimization and design of geodetic networks. Springer, Berlin, pp 11–55

---

# Atmospheric Loading and Mass Variation Effects on the SLR-Defined Geocenter

Rolf König, Frank Flechtner, Jean-Claude Raimondo, and Margarita Vei

---

## Abstract

The geocenter time series can conveniently be inferred by evaluating satellite orbit perturbations in a dynamical model. We base our approach on nearly 30 years of Satellite Laser Ranging (SLR) observations to the LAGEOS satellites and recent model standards. We model station deformations by atmospheric pressure loading and we model orbit perturbations by mass variations in the atmosphere. Both effects on the geocenter realization are analyzed.

---

## Keywords

Atmospheric mass variations • Atmospheric pressure loading • Geocenter • SLR

---

## 1 Introduction

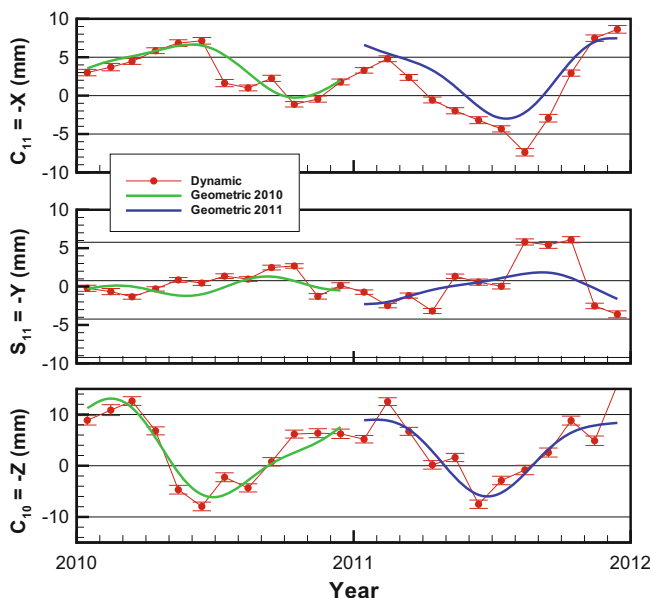
Satellite Laser Ranging (SLR) contributes to the determination of the International Terrestrial Reference Frame (ITRF) besides the other space geodetic techniques as Global Positioning Service (GPS), Doppler Orbitography and Radiopositioning Integrated by Satellite (DORIS), and Very Long Baseline Interferometry (VLBI). The SLR data are gathered within the International Laser Ranging Service (ILRS, Pearlman et al. (2002)). In the most recent ITRF realization, the ITRF2008 (Altamimi et al. 2011), SLR is the only one to define the origin of the reference frame which by definition coincides with the center of mass of the Earth system or the geocenter. The sensitivity of SLR to the geocenter can also be exploited to monitor variations of the geocenter. First results are reported by Watkins and Eanes (1997), a recent series is published by Cheng et al. (2013). A general observation is that the geocenter shows seasonal variations with amplitudes in the millimeter range.

Geocenter time series show amongst other model deficiencies mass redistributions in the Earth system, mass redistribution phenomena include among others atmospheric pressure field variations and non-tidal mass transport in the oceans and atmosphere. The atmosphere acts in two ways. Atmospheric pressure loading (APL) changes the coordinates of the stations that realize the Terrestrial Reference Frame (TRF). And the atmospheric mass redistributions contribute to the changing gravity field. The atmospheric mass variations act also on the oceans which in turn lead to mass variations within Earth's water and land mass. Those effects are available in form of the so-called GRACE Atmosphere and Ocean Dealiasing (AOD) products (Flechtner and Dobslaw 2013). It should be noted that these models are forward models in GRACE data evaluation, not the result of GRACE observations, and not perfect models.

Depending on the approach of how the geocenter time series are being determined the variations are given either with respect to the center of the TRF or with respect to the center of the gravity field model. We evaluate SLR data by Precise Orbit Determination (POD) by our software system EPOS-OC (see Zhu et al. (2004)). In a differential orbit and parameter improvement process all dynamic, geometric and measurement models are set up to estimate the parameters describing the orbit, parameters for mitigating the influence of model errors, and those parameters we are interested in.

---

R. König (✉) • F. Flechtner • J.-C. Raimondo • M. Vei  
GFZ German Research Centre for Geosciences, c/o DLR  
Oberpfaffenhofen, 82234 Wessling, Germany  
e-mail: koenigr@gfz-potsdam.de



**Fig. 1** Comparison of dynamic and geometric geocenter solutions

For the geocenter the degree one terms of the gravity field are solved for, thus the results are relative to the center of the TRF.

Kar (1997) recommends to infer geocenter motion rather from a geometric solution than using degree one terms because non-zero degree one terms result in a rotating reference system where corrections are needed to account for the centrifugal, Euler, and Coriolis accelerations. For the following analysis these accelerations are neglected. In order to assess the possible error contribution from that both the dynamic and the geometric approaches are compared in a short test period of 2 years, the results are displayed in Fig. 1. For the geometric approach EPOS-OC estimates bias, annual and semi-annual periodic motions of the TRF per year. So a one-to-one comparison of both approaches is not possible because of the different number of solved for parameters. Also Kar (1997) points out that the one-to-one relationship of the geometric and dynamic approach is difficult to achieve in real world applications. Figure 1 reveals similarities in the seasonal behaviour but also discrepancies when looking at monthly differences which show standard deviations of 1.8, 2.0, and 2.9 mm in X, Y, and Z respectively. There are no significant differences between the annual amplitudes of the dynamic and the geometric solutions.

For the solution a datum has to be chosen. Cheng et al. (2013) fix the coordinates of the SLR ground stations to a certain reference frame. We are going to apply Helmert conditions that care for the datum defect. The choice of a consistent datum in SLR applications over long time periods is practically impossible. Therefore we put some efforts in qualifying our datum choice for determining geocenter time series.

Once the proper datum is found the geocenter time series is compared to cases when APL is applied to the TRF and when AOD mass variations are applied to the gravity field.

## 2 Data

LAGEOS launched in 1976 (ILRS 2013) is a cannon ball satellite in an orbit with a semi-major axis of 12,270 km, an inclination of  $109.8^\circ$ , and an eccentricity of 0.004. Its twin LAGEOS-2 launched in 1992 orbits with a semi-major axis of 12,160 km, an inclination of  $52.6^\circ$ , and an eccentricity of 0.014. We use SLR data to LAGEOS from 1983 to 2011 for this analysis, and SLR data to LAGEOS-2 from 1992 to 2011. The orbital fits of the operational analysis is shown in Fig. 2. It becomes clear that the quality of the data steadily improves from the early years until about the advent of LAGEOS-2. From then on the RMS of orbital fits size around the centimeter. For LAGEOS 1,750,000 Normal Points are evaluated, for LAGEOS-2 1,165,000. For the LAGEOS data evaluation before the advent of LAGEOS-2, 15-day arcs are used, from then on 7-day arcs. So in the early years the mean number of NPs per arc lies around 2,200, later on it lies around 1,400 per satellite with a light but steady increase over time.

The APL concerns deformations of the solid Earth due to atmospheric pressure. It is computed in-house for all the sites of this study based on the theory of Farrel (1972). Input surface pressure fields come from the European Centre for Medium-Range Weather Forecasts (ECMWF 2013). For actual periods the ECMWF operational data are used, thus those data that are input for the generation of the GRACE AOD products. As these data have only been archived since the beginning of the GRACE mission, for our purpose the archive needs to be expanded back in time. This is done via the re-analysis data ERA-40 (ECMWF 2013) covering mid 1957 to mid 2002 and alternatively ERA-Interim (ECMWF 2013) covering 1979 up to date.

Site displacements due to APL exhibit sub-daily to seasonal signals with amplitudes up to 20 mm. An example is shown in Fig. 3 for station Maidanak. Maidanak is located in the middle of Asia. Therefore the effect of APL is quite pronounced. Figure 3 also shows a comparison of the in-house derived displacements with those generated independently by Petrov and Boy (2004) with a good agreement between the two.

The GRACE AOD products provide short term mass variations from atmosphere and ocean. Hydrological mass redistributions are not included on purpose because hydrological models are of lower quality. Within the GRACE project the AOD products thus serve to take off the signals from atmosphere and ocean that otherwise corrupt the signal from hydrology. GRACE gravity field time series

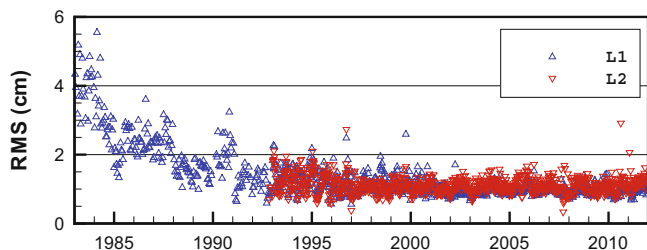


Fig. 2 LAGEOS-1 and -2 orbital fits

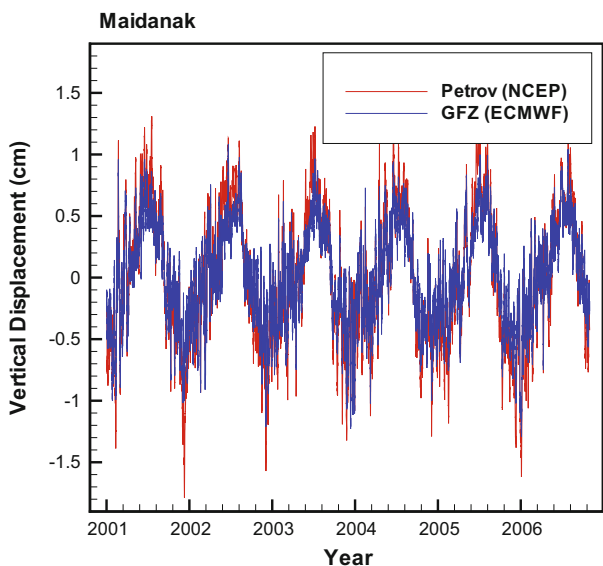


Fig. 3 Atmospheric loading site displacements for station Maidanak

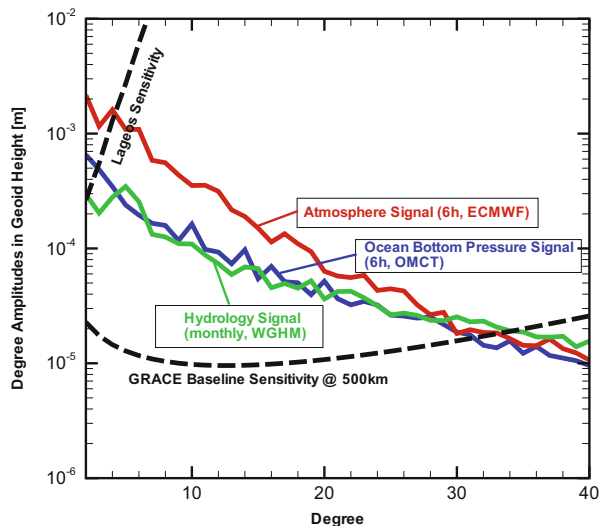


Fig. 4 Signal level and sensitivities

show therefore mainly variations due to the hydrological cycle (Schmidt et al. 2008). This is demonstrated in Fig. 4 where the sensitivity level of GRACE is displayed versus the various signals in the frequency domain. Figure 4 also shows

the sensitivity level of LAGEOS. In comparison to the signal curves it becomes clear that LAGEOS is also sensitive to the atmospheric and oceanic signals in the very low degrees. Thus again if atmospheric and oceanic signals are taken off LAGEOS may also sense hydrology in the very long wavelengths.

### 3 Approach

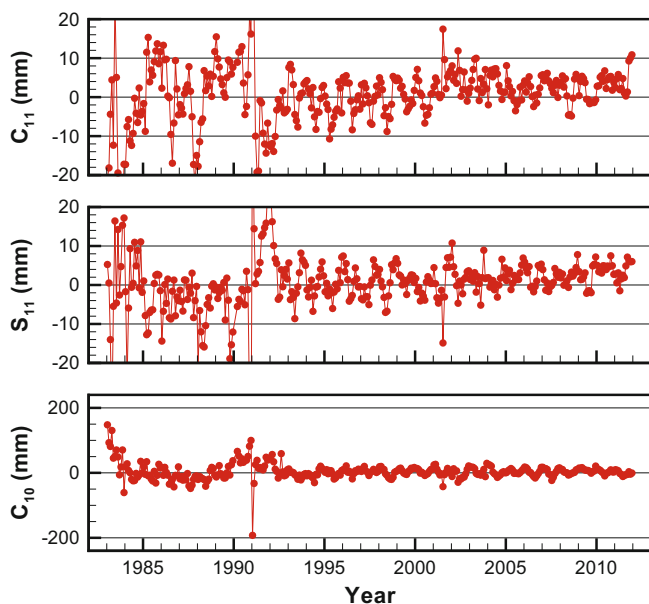
The geocenter is estimated via the dynamic approach by evaluating SLR data to LAGEOS and LAGEOS-2 with EPOS-OC. The standards mostly follow the Conventions 2010 (Petit and Luzum 2010) of the International Earth Rotation and Reference Systems Service (IERS 2013). In particular we take for the gravity field model the static part of EIGEN-6C (Förste et al. 2011), for the a priori station coordinates the IERS adoption SLRF2008 (IERS 2013) of the ITRF2008, and for the Earth orientation parameters the “EOP 08 C04 (IAU2000)” series (IERS 2013).

The solved for parameters of interest are the station coordinates and the degree one to five spherical harmonic coefficients of the gravity field. The degree one and two spherical harmonics are stacked to monthly solutions, the degree three to five spherical harmonics and the station coordinates are stacked to global solutions. The degree two coefficients are also solved for monthly because over these time spans considerable variations are taking place (see Cheng and Tapley (2004) that, if not modelled, could contaminate the degree one results. The C(1,1), S(1,1), and C(1,0) parameters represent the dynamic geocenter solution and correspond to X-axis, Y-axis, and Z-axis offsets respectively of a geometric solution.

The datum defect of this approach is six, three translations and three rotations are free. There is no defect in scale for it is provided by the highly precise SLR range observations. Therefore we apply three no-net-translation and three no-net-rotation conditions (NNRT), i.e. the Helmert conditions, for the station coordinates to the solution. At this point it must be noted that a common datum for all the 15-day and 7-day solutions can practically not be found as the composition of the station network changes considerably over time.

### 4 Results

A first time series of geocenter estimates is given in Fig. 5. For this no APL and no AOD is applied. For the datum the NNRT conditions are applied to 8 stations out of all which show a good performance in terms of number of observations and of quality. Figure 5 shows some evident systematics for the early years that is the period before the advent of LAGEOS-2.

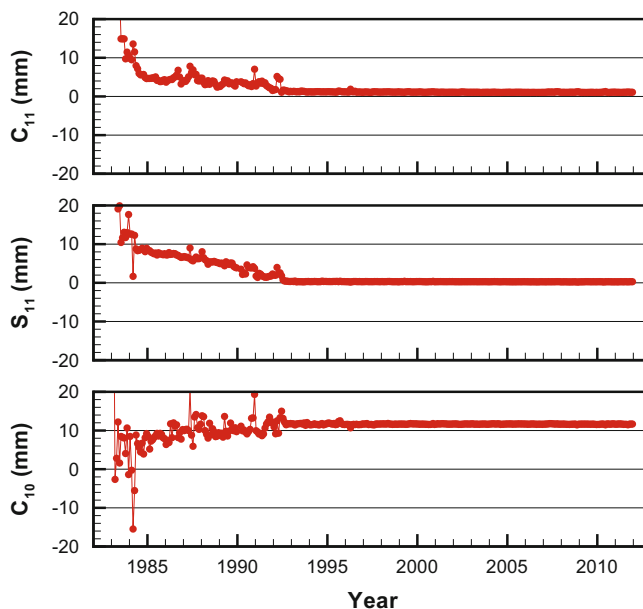


**Fig. 5** Basic geocenter time series

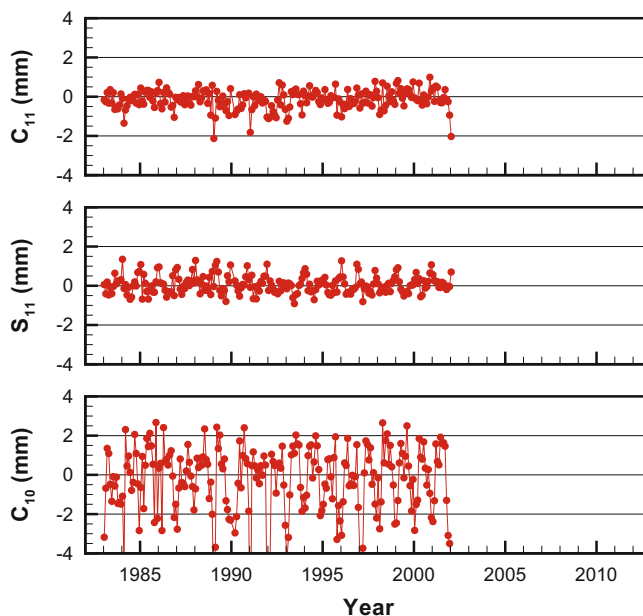
In order to check whether these systematics depend on the choice of the datum, a second solution is generated where the NNRT conditions are applied not only to 8 but to all 122 stations that showed up in the analysis period of 29 years. Figure 6 displays the differences of this new geocenter time series to the previous one. Indeed the significant impact of the choice of the datum becomes evident. The early years of the analysis period are characterized by varying biases, later on when both satellites are present the bias gets stable. Therefore a reliable solution becomes available from 1993 onwards.

Next the APL derived from ERA-40 data is applied to the solutions, the difference with respect to the solution without APL is given in Fig. 7. The analysis period is restricted to the years 1983–2001 (the ERA-40 data end shortly after) in order to have a common datum for the following comparisons. In this way the results will be based on the same datum and therefore differences between different results will not suffer from a datum-induced bias. Figure 7 shows seasonal signals also in the differences, and significant biases or significant trends can not be observed.

Alternatively also APL derived from ERA-Interim data is applied to the solutions. Table 1 compiles the statistics for 109 monthly values spanning the years 1993–2001 for all cases when APL is applied from the two data bases and when no APL is applied. In summary the application of APL does not introduce any biases or trends to the geocenter time series. There is hardly any difference between the application of APL based on ERA-40 or on ERA-Interim data. It is obvious that the application of APL removes some signal, in particular the yearly amplitudes of the C(1,0) series



**Fig. 6** Differences of geocenter time series with different datums



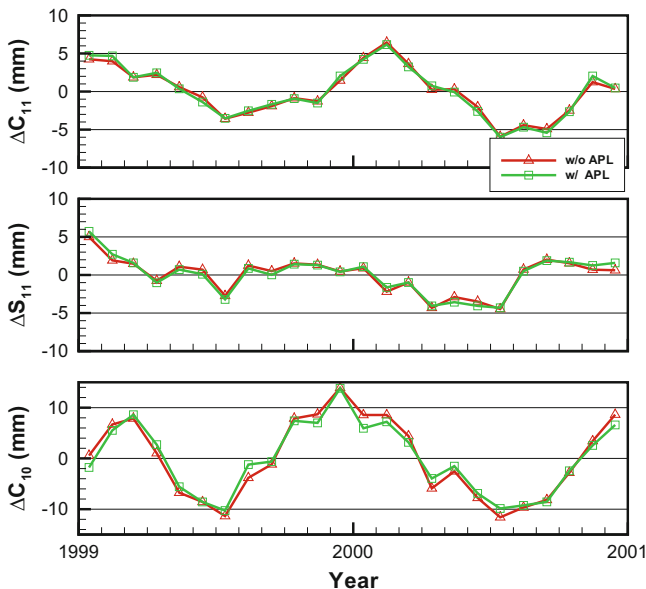
**Fig. 7** Differences of geocenter time series with APL from ERA-40 and without APL

get damped by approximately 1.7 mm (though statistically hardly at the significance level). Sošnica et al. (2013) also report a reduction in the Z amplitude by 1.1 mm. Figure 8 makes the signal reduction obvious by zooming into two geocenter time series when APL is applied or not. For these 2 years standard deviations of the differences are 0.4/0.5/1.4 mm for C(1,1)/S(1,1)/C(1,0) respectively.

To see the impact of AOD on geocenter determination, the same period as before, namely 1993–2001, is analyzed where

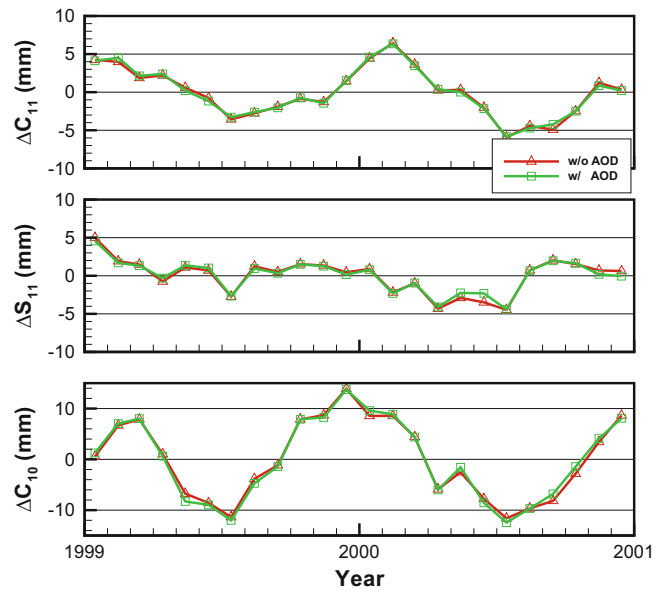
**Table 1** Comparison of geocenter series with APL from ERA-40 and from ERA-Interim and with AOD from AOD1B RL04 to a series without APL and without AOD for the years 1993–2001

	No APL No AOD	APL ERA-40	APL ERA-Interim	AOD AOD1B RL04
Mean [mm]				
C(1,1)	5.3±0.4	5.3±0.4	5.3±0.4	5.1±0.4
S(1,1)	-0.4±0.4	-0.4±0.4	-0.4±0.4	-0.2±0.3
C(1,0)	3.6±1.0	3.3±1.0	3.4±1.0	3.0±1.0
Trend [mm/a]				
C(1,1)	0.3±0.2	0.3±0.2	0.3±0.2	0.4±0.2
S(1,1)	-0.3±0.1	-0.2±0.1	-0.2±0.1	-0.3±0.1
C(1,0)	0.7±0.4	0.7±0.4	0.7±0.4	0.7±0.4
Annual amplitude [mm]				
C(1,1)	2.2±0.6	2.2±0.6	2.2±0.6	2.3±0.6
S(1,1)	3.8±0.4	4.1±0.4	4.1±0.4	3.4±0.4
C(1,0)	10.6±1.2	8.9±1.2	8.9±1.2	11.1±1.2
Semi-annual amplitude [mm]				
C(1,1)	0.3±0.6	0.3±0.6	0.3±0.6	0.4±0.6
S(1,1)	0.9±0.4	0.8±0.4	0.8±0.4	0.9±0.4
C(1,0)	3.5±1.2	4.1±1.2	4.1±1.2	3.5±1.1



**Fig. 8** Blow up of geocenter time series with and without APL

the AOD1B RL04 products are adopted. The AOD1B RL04 products are taken from Flechtner et al. (2008), a series that extends back in time to serve the pre-GRACE era. Recent GRACE AOD1B RL05 products (Flechtner and Dobsław 2013) start in 2001 only. Table 1 is augmented to give also the statistics for the AOD case for comparison. It seems that now the amplitudes get slightly amplified when AOD is considered. However this statement is below the significance level. Figure 9 provides a close up of the same 2 years as in Fig. 8 of geocenter estimates where AOD is or is not applied.



**Fig. 9** Blow up of geocenter time series with and without AOD

For these 2 years standard deviations of the differences are 0.3/0.4/0.7 mm for C(1,1)/S(1,1)/C(1,0) respectively. So the variations are a factor of up to 2 smaller compared to the APL case.

An overlap between AOD1B RL04 and RL05 products is given for the year 2001 only. To check for a bias between the two products, geocenter solutions for the year 2001 only are generated with either AOD1B RL04 or RL05 applied. As the same data base is used, the datum choice does not introduce any bias. For 12 monthly differences non-significant biases result with values of 0.0±0.1/-0.4±0.1/0.2±0.3 for C(1,1)/S(1,1)/C(1,0) respectively. So there seems to be no bias between the solutions when the AOD1B RL05 or AOD1B RL04 products are applied.

## 5 Conclusions

29 years of LAGEOS and 20 years of LAGEOS-2 data up to 2011 are processed in order to derive a geocenter time series via the dynamic method in form of monthly C(1,1), S(1,1) and C(1,0) coefficients. The fixing of the datum introduces a bias, from 1993 onwards with the advent of LAGEOS-2 a stable solution is achieved. The geocenter time series shows annual amplitudes of 2–4 mm in C(1,1) and S(1,1) and 9–11 mm in C(1,0). The results in C(1,1) and S(1,1) are well in agreement e.g. with the geometric solutions by Cheng et al. (2013) and Sošnica et al. (2013) but about 50% larger in C(1,0). In order to investigate the impact of atmospheric loading, site displacements derived from ECMWF pressure fields are adopted for the solution. This reduces the amplitudes of the geocenter time series in the

millimeter range. Biases by using APL based either on ERA-40 or on ERA-Interim data are not observed. Also the impact of atmospheric mass variations on the geocenter solution is examined. The impact of AOD is some factors smaller than that of APL but with a tendency to increase the amplitudes. Biases between the use of AOD1B RL04 and RL05 products are not observed.

**Acknowledgements** SLR data, station coordinates and Earth orientation parameters are provided by ILRS and IERS of the International Association of Geodesy (IAG).

## References

- Altamimi Z, Collilieux X, Metivier L (2011) ITRF 2008: an improved solution of the international terrestrial reference frame. *J Geod* 85(8):457–473
- Cheng MK, Tapley BD (2004) Variations in the Earth's oblateness during the past 28 years. *J Geophys Res* 109(B9):B09402
- Cheng MK, Ries JC, Tapley BD (2013) Geocenter variations from analysis of SLR data. In: Altamimi Z, Collilieux X (eds) Reference frames for applications in geosciences. International Association of Geodesy Symposia, vol 138. Springer, Berlin, Heidelberg, pp 19–26
- ECMWF (2013) European Centre for Medium-Range Weather Forecasts, Data Archive Services web site. <http://www.ecmwf.int/products/data/archive/index.html>. Accessed 16 September 2013
- Farrel WE (1972) Deformation of the earth by surface loads. *Rev Geophys* 10(3):761–797
- Flechtner F, Thomas M, König R (2008) A long-term model for nontidal atmospheric and oceanic mass redistributions and its implications on LAGEOS-derived solutions of Earth's oblateness. Scientific Technical Report STR08/12, GFZ German Research Centre for Geosciences, Potsdam, ISSN 1610-0956
- Flechtner F, Dobszlaw H (2013) AOD1B product description document for product release 05. GRACE document GRACE 327–750 (GR-GFZ-AOD-0001), <http://isdc.gfz-potsdam.de/index.php?name=UpDownload&req=getit&lid=619>. Accessed 16 September 2013
- Förste C, Bruinsma S, Shako R, Marty J-C, Flechtner F, Abrikosov O, Dahle C, Lemoine J-M, Neumayer H, Biancale R, Barthelmes F, König R, Balmino G (2011) EIGEN-6 - A new combined global gravity field model including GOCE data from the collaboration of GFZ-Potsdam and GRGS-Toulouse. EGU General Assembly 2011. Geophysical Research Abstracts, 13:EGU2011-3242-2
- IERS (2013) International Earth rotation and reference systems service web site. <http://www.iers.org>. Accessed 16 September 2013
- ILRS (2013) International laser ranging service web site. <http://ilrs.gsfc.nasa.gov>. Accessed 16 September 2013
- Kar S (1997) Long-period variations in the geocenter observed from laser tracking of multiple Earth satellites. The University of Texas at Austin, 1997. Publication Number: AAI9802916; ISBN: 9780591529647
- Pearlman MR, Degnan JJ, Bosworth JM (2002) The international laser ranging service. *Adv Space Res* 30(2):135–143
- Petit G, Luzum B (2010) IERS conventions (2010). Bundesamt für Kartographie und Geodäsie, Frankfurt am Main. <http://www.iers.org/TN36>. Accessed 16 September 2013
- Petrov L, Boy J-P (2004) Study of the atmospheric pressure loading signal in VLBI observations. *J Geophys Res* 109(B3):B03405
- Schmidt R, Flechtner F, Meyer U, Neumayer K-H, Dahle C, König R, Kusche J (2008) Hydrological signals observed by the GRACE satellites. *Surv Geophys* 29(4–5):319–334
- Sośnica K, Thaller D, Dach R, Jäggi A, Beutler G (2013) Impact of loading displacements on SLR-derived parameters and on the consistency between GNSS and SLR results. *J Geod* 87: 751–769
- Watkins MM, Eanes RJ (1997) Observations of tidally coherent diurnal and semidiurnal variations in the geocenter. *Geophys Res Lett* 24(17):2231–2234
- Zhu S, Reigber C, König R (2004) Integrated adjustment of CHAMP, GRACE and GPS data. *J Geod* 78(1–2):103–108

---

# Radargrammetric Digital Surface Models Generation from High Resolution Satellite SAR Imagery: Methodology and Case Studies

Andrea Nascetti, Paola Capaldo, Francesca Pieralice, Martina Porfiri, Francesca Fratarcangeli, and Mattia Crespi

---

## Abstract

The goal of this paper is to investigate the potential of high resolution SAR satellite imagery for DSMs generation using the radargrammetric technique. This study is methodological, devoted to illustrate both the fundamental advantages of this approach and also its drawbacks. As for photogrammetry, the achievable accuracy level of a radargrammetric generated DSM is strictly related both to the image orientation and to the image matching procedure. A rigorous orientation model based only on metadata information and an innovative matching strategy have been developed, so that a complete suite for the DSMs generation through radargrammetry has been embedded in SISAR, a scientific software developed at the Geodesy and Geomatic Division of the University of Rome “La Sapienza”. Here we discuss the results coming from two COSMO-SkyMed SpotLight stereo pairs (ascending and descending) acquired over the area of Como (Northern Italy), characterized by a mixed land cover (flat urban area, steep forested mountain slopes). Three DSMs (ascending, descending and merged) have been generated and compared with a LiDAR DSM; the accuracy of the merged product is around 7 m, better than the accuracy of the ascending and descending DSMs (around 8–10 m).

---

## Keywords

Digital surface models • High resolution SAR • Methodology • Radargrammetry • Stereo

---

## 1 Introduction

Synthetic Aperture Radar (SAR) satellite systems may give important contributions in terms of Digital Surface Models (DSMs) generation, considering their complete independence from logistic constraints on the ground and weather conditions. Starting from the SAR data, two different methods may be used to generate DSMs: the well-known interferometric and the radargrammetric one. In particular, the radargrammetric approach was first used in the 1950s; then

it was less and less used, due to the quite low amplitude resolution of SAR imagery, if compared to their high phase resolution (Leberl 1990). Only in the last years the importance of the radargrammetric approach is rapidly growing due to the new high-resolution imagery (up to 1 m GSD) (Raggam et al. 2010; Balz et al. 2013; Gutjahr et al. 2014). Therefore, here we focus on the present potential of high resolution SAR satellite imagery for DSMs generation with a radargrammetric stereo-mapping approach, so that the goal of this paper is methodological, devoted to illustrate both the fundamental advantages of this approach and also its drawbacks. As regards the pros, it is worth to mention the independence from image coherence, unlike interferometric approach, which can be guaranteed only with an extremely short (tens of seconds) revisit time, the parsimony (it can work with just a couple of images), and therefore the short time required for imagery collection (from tens of minutes

---

A. Nascetti (✉) • P. Capaldo • F. Pieralice • M. Porfiri • F. Fratarcangeli • M. Crespi  
Geodesy and Geomatic Division - DICEA - University of Rome “La Sapienza”, Roma, Italy  
e-mail: [andrea.nascetti@uniroma1.it](mailto:andrea.nascetti@uniroma1.it)



to few hours), also thanks to the mentioned independence from illumination and weather. Concerning the cons, the well known deformations of SAR imagery may cause remarkable difficulties with complex morphologies and have to be duly accounted for the acquisition planning. The achievable accuracy level of a DSM generated through radargrammetry is strictly related both to the image orientation and matching processes. Generally, SAR stereo restitution has been modelled utilizing two Doppler equations and two range equations (Chen and Dowman 1996, 2001) and the optimum geometric configuration is when the target is observed in opposite-side view; however it causes large geometric and radiometric disparities hindering the image matching procedure. A good compromise is to use a same-side configuration stereo pair with a base to height ratio ranging from 0.25 to 2 (commonly between 0.35 and 0.70), to increase the efficiency in the correlation image process.

A rigorous orientation model based only on metadata information and an innovative matching strategy have been developed. A complete suite for the DSMs generation through the radargrammetric approach has been embedded in SISAR (Software per Immagini Satellitari ad Alta Risoluzione), a scientific software developed at the Geodesy and Geomatic Division of the University of Rome “La Sapienza” in the IDL (Interactive Data Language, [www.exelisvis.com](http://www.exelisvis.com)) environment. In order to demonstrate the radargrammetric mapping potential of high resolution SAR satellite imagery, several tests were carried out using data with different acquisition modes (SpotLight, StripMap) and coming from different platforms (COSMO-SkyMed, TerraSAR-X, RADARSAT-2). Here, for the sake of brevity, we can discuss only the results related to the COSMO-SkyMed imagery acquired on the Como area (Northern Italy), characterized by mixed morphology and land cover. Sections 2 and 3 summarize the fundamentals of radargrammetric orientation and matching, also underlying the crucial role of the orientation model within matching, whereas Sect. 4 discusses the results regarding Como area; in the end, some conclusions and possible future prospects are outlined.

## 2 Radargrammetric Orientation Model

The radargrammetry technique performs a 3D reconstruction based on the determination of the sensor-object stereo model, in which the position of each point on the object is computed as intersection of two radar rays coming from different positions and therefore with two different incidence angles. Actually, these radar rays can be simply modelled as two segments of measured lengths centered in two different positions (along two different satellite orbits), so that the intersection generating each object point is one of the two

possible intersections between two circumferences centered in the two different positions and laying into two planes orthogonal to the two different satellite orbits, whose radii are equal to the segment measured lengths. Consequently, the model is based on two standard equations. The first equation of (1) represents the general case of zero-Doppler projection: in zero-Doppler geometry the target is acquired on a heading that is perpendicular to the flying direction of satellite; the second equation of (1) is the slant range constrain. The couple of equations in a local Cartesian system reads (Capaldo et al. 2011):

$$\begin{cases} V_{XS} \times (X_S - X_P) + V_{YS} \times (Y_S - Y_P) + \\ V_{ZS} \times (Z_S - Z_P) = 0 \\ \sqrt{(X_S - X_P)^2 + (Y_S - Y_P)^2 + (Z_S - Z_P)^2} \\ - (D_S + CS \times I) = 0 \end{cases} \quad (1)$$

where

$X_P, Y_P, Z_P$	are the coordinates of the ground point $P$ (time independent)
$X_S, Y_S, Z_S$	are the coordinates of the satellite sensor (time dependent)
$v_{SX}, v_{SY}, v_{SZ}$	are the components of the satellite sensor velocity (time dependent)
$D_S$	is the so-called near range
$CS$	is the column spacing
$I$	is the column position of point $P$ on the image

The relationship between image coordinate  $J$  and the time  $t$ , can be expressed by a linear relation

$$t = t_0 + \frac{1}{PRF} J \quad (2)$$

in which the start time of the acquisition ( $t_0$ ) and the Pulse Repetition Frequency ( $PRF$ ), the sampling frequency in azimuth direction, are involved.

Starting from the few state vector available in metadata, the orbit segment must be reconstructed using some kind of interpolation in order to compute the satellite position for each line. Here, it was adopted the Lagrange polynomials: these interpolation is sufficiently accurate to model the short orbital segment and its well-known problems at the edges do not affect the modelling since the images are acquired in the central part of the orbital segment. Additionally, using a standard divide and conquer algorithm it was possible to find in a rapid and accurate way the epoch when satellite orbit is perpendicular to the line of sight between the sensor and the ground point. In recent studies has been proven that, given the high precision of the electronic metadata parameters and the high accuracy of position and velocity state vectors, the orientation can be performed with high accuracy without the refinement of Ground Control Points (GCPs) (Capaldo et al.

2011). This fact represents a significant advantage of high resolution SAR imagery with respect to the optical ones, whose precise orientation requires at least a few GCPs.

Furthermore, the presented radargrammetric orientation model can be conveniently parameterized using the Rational Polynomial Coefficients (RPCs), recognized as a suitable model for high resolution SAR imagery (Zhang and Zhu 2008; Zhang et al. 2010). In particular, a tool for RPCs generation, based on a so-called terrain independent scenario, is implemented in SISAR (Capaldo et al. 2012) and it has been used to reduce the computational effort of matching process, indeed the Rational Polynomial Functions (RPFs) provide a direct mathematical relationship (rapid and easy to compute) to convert ground to image coordinates (see Sect. 3.2).

### 3 SAR Image Matching

The development of a fully automatic, accurate and reliable image matching method that adapts to different images and scene contents is a challenging problem. Dissimilarities between SAR images due to occlusion, geometric distortions, radiometric differences and speckle noise must be taken into account and this is one of the reasons why many different image matching approaches have been developed in recent years. Hereafter the basic features of our original matching procedure, presently under patenting by the University of Rome “La Sapienza”, are outlined.

#### 3.1 Area Selection and Filtering

At the beginning of the image matching procedure, it is mandatory to select an area of interest and a coarse height range (approximate maximum and minimum terrain ellipsoidal heights), in order to reduce the object space and to remarkably decrease the processing time.

Moreover, SAR imagery are affected by a particular form of noise called speckle and that they exhibit a grainy appearance (*salt-and-pepper*) hindering target recognition and correct matching. In order to reduce speckle, three different kind of well-known (but never applied to high-resolution SAR imagery) adaptive spatial filters (Lee, Kuan, GammaMap) have been considered for a preprocessing enhancement. Nevertheless, thanks to a number of tests, it was highlighted that these spatial filters significantly increase the number of points at the expense of vertical accuracy, since they mitigate the speckle but smooth the image features.

Starting from this experimental awareness, an original filtering procedure *dynamic filtering* has been developed, in order to maximize not only the number of points, but also their quality. Unlike the traditional preprocessing techniques, the image filtering is done directly during the matching

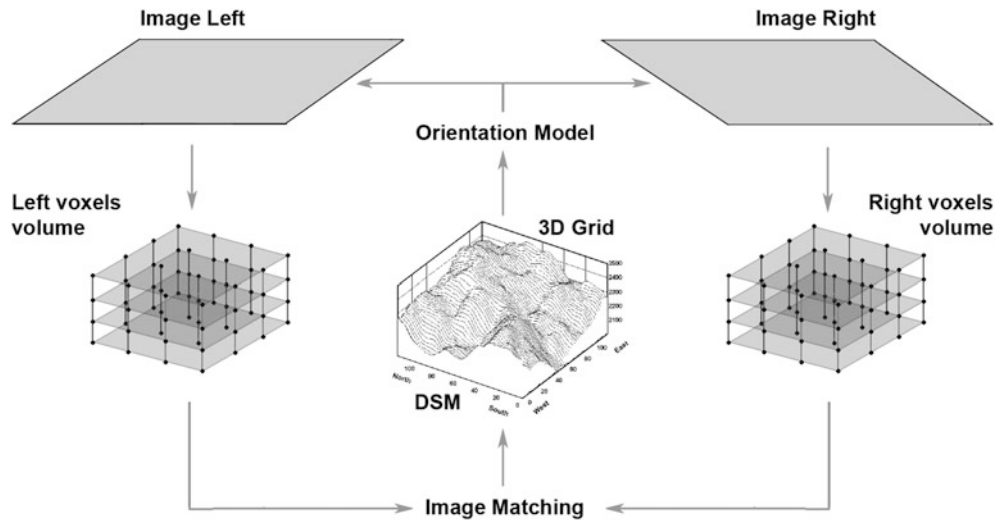
procedure; the leading idea is to find all the possible matched point using the raw imagery and, only after, apply speckle spatial filters (i.e. Lee, Kuan, GammaMap) to search points in areas where the previous search failed. This allows to operate at several pyramidal levels (with different resolution) independently and in different ways, e.g. making one or more filtering cycles.

#### 3.2 Image Matching Strategy

The image matching strategy is based on an hierarchical solution with a geometrical constrain, and the corresponding points (actually, so-called primitives) are searched using an area based matching criterion and analysing the signal-to-noise ratio (SNR) (Ma et al. 2004). In this sense, the peculiarity of the proposed algorithm is to use the image orientation model (re-parameterized in form of RPCs) to limit the search area of the corresponding primitives, allowing a fast and robust matching. Primitives are searched directly in the object space re-projecting and re-sampling the stereo images on a regular grid in the ground geometry, that is in the ground reference system. In fact starting from a ground point with a selected height, the orientation model provides point image coordinates and thus it is possible to back-transfer the SAR radiometric information from slant-range to ground geometry. Therefore it is important to underline that the effectiveness of the matching algorithm is strictly related to the orientation model accuracy, since the strategy is based on geometrical constrains.

From the practical point of view, after images preprocessing and area selection, a 3D grid is generated in ground geometry, with several layers slicing the entire height range. Starting from this 3D grid, by means of the orientation model, the two images are re-projected on each layer creating two voxel sets (one for left and one for right image). Through this process (see Fig. 1), the two generated voxel sets contain the geometrically corrected radiometric information in the same ground reference system. At this point, for each horizontal position (X,Y) of the 3D grid, the main objective is to identify the correct height comparing the two voxel sets. This correct height corresponds to the best matching of the two voxels (for left and right image) at the same height; therefore, to this aim, the search can be conveniently carried out along vertical paths.

During the algorithm development different primitive models have been considered (i.e. Area Based Matching or Feature Based Matching), and the experimental results have highlighted that a normalized cross-correlation (NCC) linked with a signal-to-noise ratio analysis is the more efficient and accurate method. Overall, for each horizontal position (X,Y),



**Fig. 1** Geometrical constrain and voxel generation

the search of the corresponding primitives consists of the following steps:

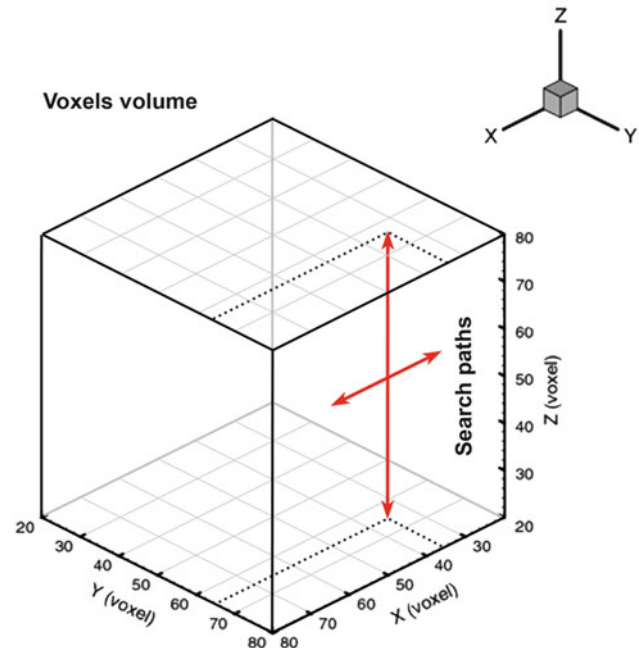
- compute the NCC values along the vertical search path
- find the maximum NCC value along the vertical search path (see Fig. 2)
- analyse the NCC profile and compute the vertical Signal-to-Noise Ratio (SNR) according to the formula:

$$SNR_v = \frac{1 + \rho_{max}}{1 + \bar{\rho}} \quad (3)$$

where  $\rho_{max}$  and  $\bar{\rho}$  are respectively the maximum and mean value of NCC along the vertical search path; note that this search mainly examines the correspondence of primitives in West-East direction for each horizontal layer, that is orthogonally to the direction of the orbits of the considered SAR satellite

- to strength the matching a second search is performed moving the correlation windows in North-South direction in the selected horizontal layer (see Fig. 2), starting from the height corresponding to the found NCC maximum value; accordingly to the same formulation [Eq. (3)] a second value  $SNR_p$  is computed
- if  $\rho_{max}$  and both  $SNR_v$  and  $SNR_p$  are higher than the respectively chosen thresholds, the primitives are considered matched and the height value for the horizontal position (X,Y) is finally determined

At the end of this process, after investigating and finding all the corresponding primitives for each (X,Y) position, an irregular DSM (point cloud) in (X,Y,Z) coordinates is obtained.



**Fig. 2** Search paths

### 3.3 Pyramidal Approach

The described matching strategy is used in a coarse-to-fine hierarchical solution, following a standard pyramidal scheme based on a multi-resolution imagery approach. The well known advantage of this technique is that at lower resolution it is possible to detect larger structures whereas

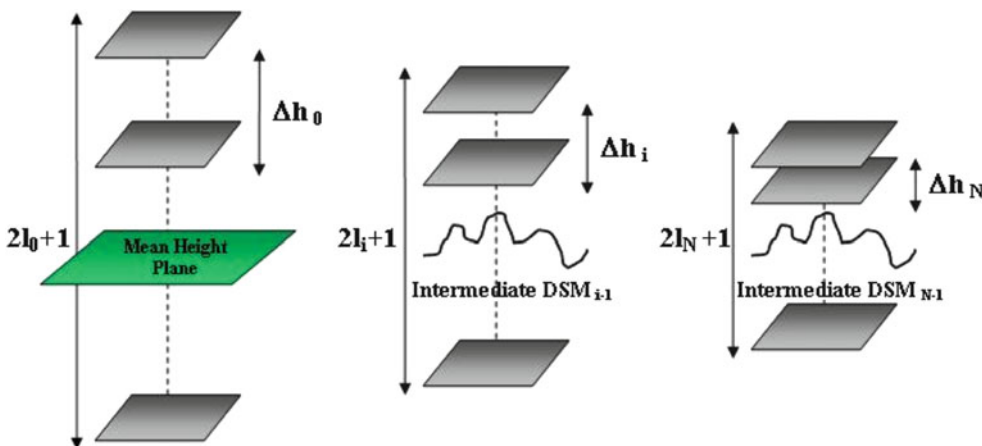


Fig. 3 Coarse-to-fine approach

Table 1 Como stereopairs characteristics

	ASC stereopair		DESC stereopair	
Date	7/8/11	17/6/11	24/6/11	28/6/11
Area (km)	10 × 10	10 × 10	10 × 10	10 × 10
Inc. Ang. (deg)	28.9	50.8	27.8	55.4
Orbit	Asc	Asc	Desc	Desc
Look side	Right	Right	Right	Right
B/H	0.6		0.8	

at higher resolutions small details are progressively added to the already obtained coarser DSM. The procedure is started choosing a suitable image multi-looking considering the original image resolution.

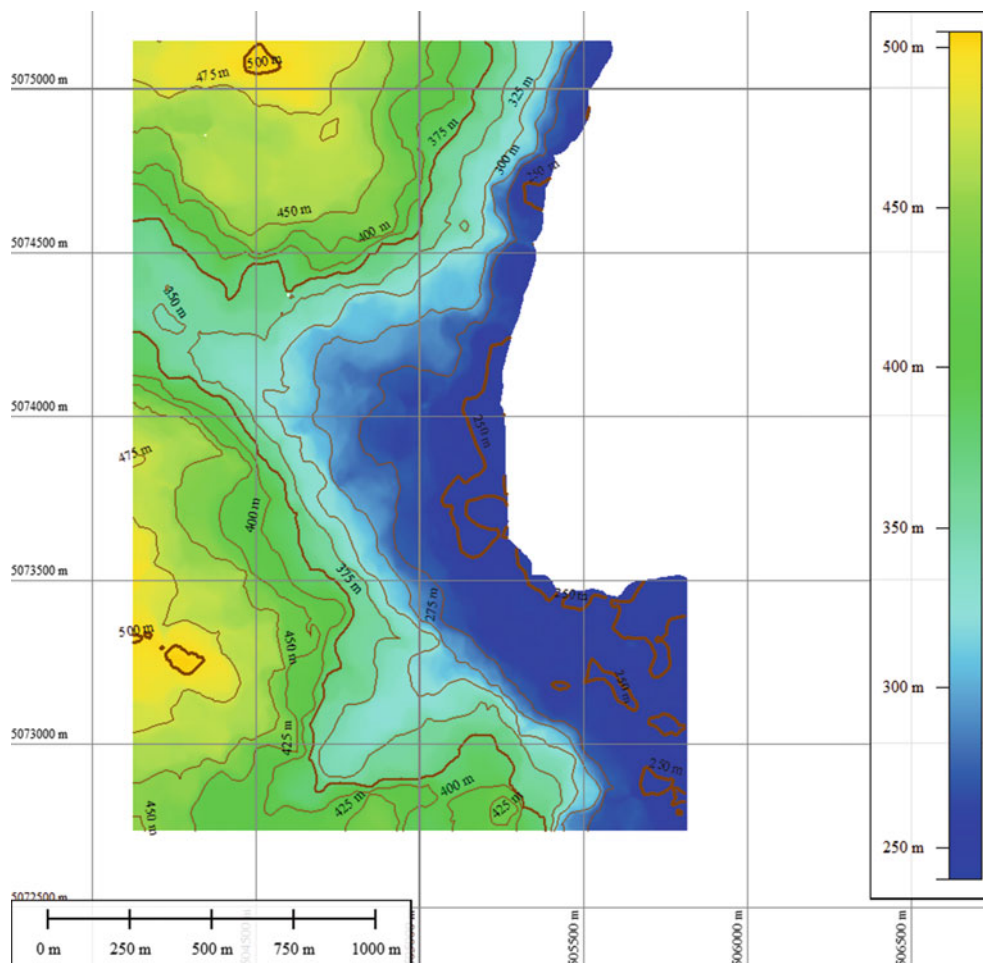
In this way, at each pyramid step, an intermediate DSM is extracted and it is modeled by the triangular irregular network (TIN) using a 2D Delauney triangulation method. Further, DSM is interpolated on a regular grid in the ground reference system, becoming the input for the next pyramid level. Correspondingly, for each horizontal position (X,Y), the height coming from the DSM obtained in the previous pyramid step is selected as starting point for the vertical search path, whereas at the first iteration just a plane with a mean elevation is set as reference DSM. In this respect it is worth to underline that, differently from other approaches Raggam et al. (2010), no external DSMs (for example SRTM DEM or ASTER DEM) are needed to guide matching.

As long as the resolution and pyramid level increase, and the DSM approaches the final solution, the mentioned discretization of the entire height range is correspondingly refined, so that the height step between the layers of the 3D grid and also the number of the considered layers decrease (Fig. 3).

## 4 Experimental Results

Several stereo pairs, with different image types (SpotLight, StripMap) coming from different SAR sensors (COSMO-SkyMed, TerraSAR-X, RADARSAT-2), have been processed and accuracy evaluations have been performed, comparing the generated DSMs with more accurate reference DSMs usually coming from airborne LiDAR. These comparisons and quality assessment have been carried out using the scientific software DEMANAL (K. Jacobsen, Leibniz University of Hannover, Germany). Here, for the sake of brevity, we can discuss only the results related to the COSMO-SkyMed SpotLight imagery (a descending and an ascending same-side stereo pair), with a GSD of 1 m, acquired over the area of Como (Northern Italy), characterized by mixed land cover (flat urban area, steep forested slopes) (Table 1).

Three DSMs have been generated (ascending, descending and merged), estimating the heights on a 3 × 3 m grid by a standard kriging interpolation and they have been compared with LiDAR DSM with a posting of 1 × 1 m and a height accuracy of 0.25 m. In particular, the merged DSM (Fig. 4) has been generated using a combination of the point clouds that have been previously filtered, removing the matched points with lower correlation. The accuracy (RMSE) ranges from 8–10 m for the ascending and the descending DSMs. The accuracy of descending DSM is worse than the accuracy of ascending one and it is due to the lower quality of radiometric information in one of the descending image. In particular, one image presents some small zones affected by SAR artefact that hinder the matching procedure. The accuracy decreases to 7 m in the merged product (Table 2), with remarkably leptokurtic error distribution. The improvement



**Fig. 4** Como merged DSM

**Table 2** Como DSMs accuracy assessment [m]

DSM	BIAS	ST. DEV.	RMSE	LE95
Asc.	-2.07	7.69	7.96	22.32
Desc.	-1.98	10.76	10.94	34.18
Merged	-1.80	7.30	7.52	19.71

of accuracy in the merged DSM is around 0.5 m, with respect to the accuracy of the ascending DSM since the accuracy of the descending one is corrupted by the images radiometric low quality. It is evident both the impacts of acquisition configuration (ascending vs. descending, causing different foreshortening and layover) and the benefit of considering different acquisition configuration over areas characterized by steep slope. Even if in this case, due to the mentioned poor radiometry of one descending image, the merged DSM accuracy improvement is limited the use of at least two stereo pairs acquired under different look sides can be an effective strategy to mitigate the SAR distortion problems, obtaining a more dense point cloud.

## 5 Conclusions and Future Prospects

The paper discusses the present potential of high resolution satellite SAR imagery for DSMs generation with a radargrammetric stereo-mapping approach and presents the main features of the radargrammetric procedure defined and implemented in SISAR package. It outlines the orientation model and, with more attention, it focuses on the original matching strategy, presently patent pending by the University of Rome “La Sapienza”. It is based on area based primitive model and on an hierarchical solution with geometrical constrain. The leading idea has been to search the corresponding primitives directly in the object space, re-projecting and re-sampling the stereo images into a 3D ground grid. The correspondences are looked analysing the signal-to-noise ratio (SNR) along two perpendicular search paths. Moreover, a specific speckle dynamic filtering technique has been designed and embedded into the radargrammetric procedure, based on three standard speckle filters (Lee, Kuan, GammaMap).

The complete radargrammetric processing chain has been developed and implemented using the IDL development environment. In order to demonstrate its mapping potential, several tests were carried out using high resolution SAR satellite imagery with different acquisition mode (Spotlight, Stripmap) and coming from different platforms (COSMO-SkyMed, TerraSAR-X, RADARSAT-2). A homogeneous DSM assessment procedure has been considered in the different tests, based on the comparison with a reference ground truth using the scientific software DEMANAL. Summarizing the main results of other tests, the DSMs vertical accuracy is strictly related to the terrain morphology and land cover. In case of limited SAR distortions (layover and foreshortening) the observed RMSE values range from 3 to 4 m over bare soil and forest to 6–7 m in more complex urban areas. Otherwise, with strong SAR distortions, as the considered area of Como, the accuracy becomes worse and the terrain morphology may be conveniently reconstructed using at least two same-side stereo pairs with different acquisitions (ascending and descending).

Finally, it was demonstrated that radargrammetric stereo-mapping approach appears a valuable tool to supply topographic information and it is likely to become an effective complement/alternative to InSAR technique, since it may work using a couple of images with a good performance even over areas (forested or vegetated areas) characterized by low coherence values.

Although the experimental results have demonstrated that StereoSAR approach has the capability to give good and encouraging results, there are still a lot of challenging issues which need to be considered for further improvements. Hereafter we propose some ideas for the future:

- efficiency improvement in urban areas: to model the complicated urban morphology, specific algorithms must be investigated, also accounting for remarkable features as double bounces or building shadows; in this respect, some preliminary investigations applying semiglobal matching (Hirschmuller 2008)
- accounting for polarimetric information: studying algorithms and techniques for optimizing DSMs generation from full SAR polarimetric data through radargrammetry; in particular, the potential of polarimetric imagery and their derived products (i.e. span, entropy, H-A classification maps) should be investigated in order to enhance the image matching
- interferometry and radargrammetry tight integration: the two techniques should be considered to exploit the 3D mapping potential of high resolution satellite SAR data:

radargrammetric DSMs can be used within the InSAR processing chain to simplify the unwrapping process in order to avoid areas affected by phase jumps

**Acknowledgements** The Authors are indebted with:

- Prof. K. Jacobsen for the DEMANAL software
- Dr. R. Lanari, PI of the Italian Space Agency Announcement of Opportunity for COSMO-SkyMed project “*Exploitation and Validation of COSMO-SkyMed Interferometric SAR data for Digital Terrain Modelling and Surface Deformation Analysis in Extensive Urban Areas*”, within which the Como imagery were made available
- Regione Lombardia, for making available the LiDAR DSM

## References

- Balz T, Zhang L, Liao M (2013) Direct stereo radargrammetric processing using massively parallel processing. *ISPRS J Photogramm Remote Sens* 79:137–146. ISSN:09242716. doi:10.1016/j.isprsjprs.2013.02.014
- Capaldo P, Crespi M, Fratarcangeli F, Nascetti A, Peralice F (2011) High-resolution SAR radargrammetry: a first application with COSMO-SkyMed SpotLight imagery. *IEEE Geosci Remote Sens Lett* 8(6):1100–1104, 5936096. ISSN:1545598X. doi:10.1109/LGRS.2011.2157803
- Capaldo P, Crespi M, Fratarcangeli F, Nascetti A, Peralice F (2012) A radargrammetric orientation model and a RPCs generation tool for COSMO-SkyMed and TerraSAR-X high resolution SAR. *Ital J Remote Sens* 44(1):55–67. ISSN:2279–725. doi:10.5721/ItJRS20124415
- Chen PH, Dowman IJ (1996) Space intersection from ERS-1 synthetic aperture radar images. *Photogramm Rec* 15(88):561–573. ISSN:0031868X. doi:10.1111/0031-868X.00064
- Chen PH, Dowman IJ (2001) A weighted least squares solution for space intersection of spaceborne stereo SAR data. *IEEE Trans Geosci Remote Sens* 39(2):233–240. ISSN:01962892. doi:10.1109/36.905231
- Gutjahr K, Perko R, Raggam H, Schardt M (2014) The epipolarity constraint in stereo-radargrammetric DEM generation. *IEEE Trans Geosci Remote Sens* 52:5014–5022. ISSN:01962892. doi:10.1109/TGRS.2013.2286409
- Hirschmuller H (2008) Stereo processing by semiglobal matching and mutual information. *IEEE Trans Pattern Anal Mach Intell* 30(2):328–341. ISSN:01628828. doi:10.1109/TPAMI.2007.1166
- Leberl F (1990) Radargrammetric image processing. Artech House, Norwood
- Ma Y, Soatto S, Kosecka J, Shankar Sastry S (2004) An invitation to 3D vision: from images to geometric models. Springer, New York
- Raggam H, Gutjahr K, Perko R, Schardt M (2010) Assessment of the stereo-radargrammetric mapping potential of TerraSAR-X multi-beam Spotlight data. *IEEE Trans Geosci Remote Sens* 48(2):971–977. ISSN:01962892. doi:10.1109/TGRS.2009.2037315
- Zhang G, Zhu XY (2008). A study of the RPC model of TerraSAR-X and COSMO-SKYMED SAR imagery. *Int Arch Photogramm Remote Sens Spat Inf Sci XXXVII(B1):321–324*
- Zhang G, Fei WB, Li Z, Zhu XY, Li DR (2010) Evaluation of the RPC model for spaceborne SAR imagery. *Photogramm Eng Remote Sens* 76(6):727–733

**Digital Terrain Modeling, Synthetic Aperture Radar  
and New Sensors: Theory and Methods**

---

# Principles and Applications of Polarimetric SAR Tomography for the Characterization of Complex Environments

Laurent Ferro-Famil, Yue Huang, and Eric Pottier

---

## Abstract

Despite its widely recognized capabilities for mapping and characterizing large areas, 2-D Synthetic Aperture Radar (SAR) imaging meets serious limitations over volumetric media, due to its incapacity to discriminate scattering contributions in the elevation direction. This paper proposes some methods for characterizing complex volumetric environments using polarimetric SAR tomography, a 3-D imaging technique based on the use of diversely polarized electromagnetic waves acquired from different trajectories. The use of polarimetric diversity permits to both improve the tomographic separation between different components of complex volumetric media and to characterize the EM behavior of the observed environments. A set of spectral estimation techniques, adapted to tomographic focusing, are tested against signal models accounting for the statistical complexity of hybrid volumetric environments. Due to their statistical adaptivity, their robustness to mismodeling and their accuracy, spectral estimators based on weighted subspace fitting criteria are selected and extended to the polarimetric case. The effectiveness of the proposed approaches is assessed over real data and for three different applications, related to urban area 3-D mapping using a minimal set of images, tropical forest structure characterization using low frequency waves, and under-foliage concealed vehicle imaging.

---

## Keywords

3-D imaging • Polarimetry • Spectral estimation • Synthetic Aperture Radar • Tomography

---

## 1 Introduction

A SAR is an active device used to form 2-D maps of the electromagnetic (EM) reflectivity of environments, with a resolution generally ranging from some decimeters to some meters, which are well adapted to environmental cartography, physical parameter estimation and change monitoring.

---

L. Ferro-Famil (✉) • E. Pottier  
University of Rennes 1, IETR, Rennes, France  
e-mail: [Laurent.Ferro-Famil@univ-rennes1.fr](mailto:Laurent.Ferro-Famil@univ-rennes1.fr);  
[Eric.Pottier@univ-rennes1.fr](mailto:Eric.Pottier@univ-rennes1.fr)

Y. Huang  
Intermap Technologies, Inc., Calgary, AB, Canada  
e-mail: [yhuang@intermap.com](mailto:yhuang@intermap.com)

Due to its intrinsic 2-D nature, single-channel SAR imaging meets strong limitations over volumes, as it cannot separate and characterize contributions from scatterers located at different elevations. Volumetric environments are generally considered as media distributed in the elevation direction and penetrated by EM waves over a vertical path larger than a wavelength, and may be modeled as a superposition of several layers, characterized by continuous or discrete densities of reflectivity and by their extinction properties, which all depend on the incident wave frequency (Treuhaff and Siqueira 2000). Data measured by SAR systems operated with one or several modes of diversity, like frequency, polarization, space, may be used to emphasize the response of specific components of a volume (Mougin et al. 1999) or characterize scattering phenomena which can be related to some of the physical properties of the observed environment



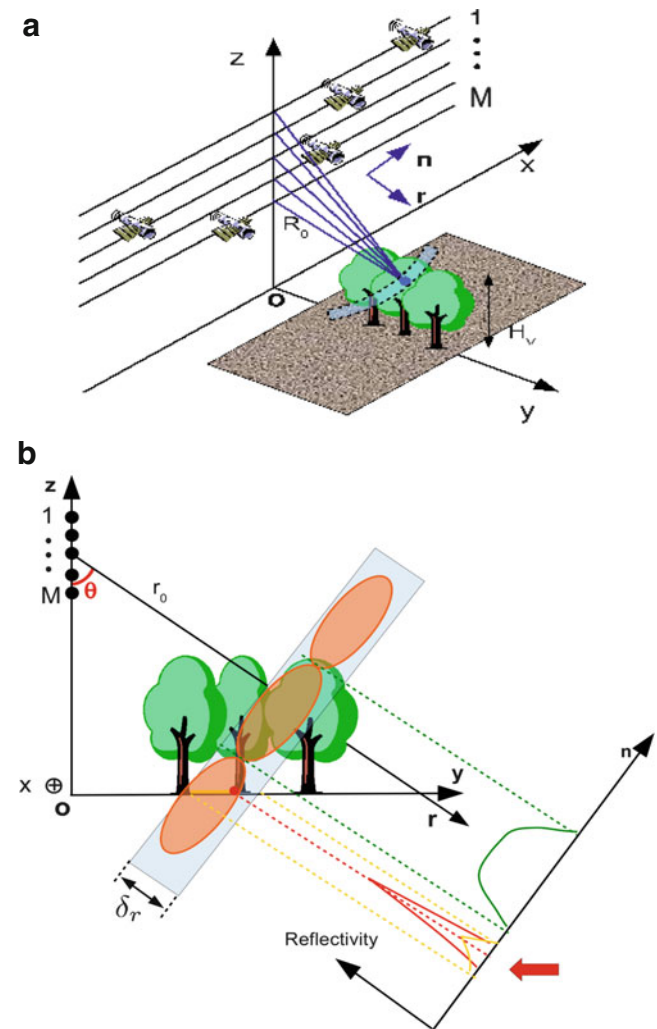
(Askne et al. 1997). In particular, the polarimetric response of a medium or object is known to be tightly related to some key physical properties, like its structure, orientation, dielectric properties (Cloude 2009) and has been widely used for the analysis of vegetated and built-up areas (Freeman and Durden 1998; Papathanassiou and Cloude 2001; Cloude 2009). Nevertheless such approaches generally rely on strong assumptions and have limited domains of validity and precisions. SAR tomography (TomSAR) is a natural solution to this problem as it is based on the use of more than two SAR acquisitions which are combined in order to form an additional synthetic aperture in the elevation direction. This 3-D SAR imaging mode has been successfully applied to the characterization of natural (Neumann et al. 2009; Tebaldini 2009; Tebaldini and Rocca 2012; Huang et al. 2011) and artificial (Fornaro and Serafino 2006; Budillon et al. 2011; Zhu and Bamler 2012; Sauer et al. 2011; Huang and Ferro-Famil 2009) environments, using ground-, air- and space-borne sensors, operated at various frequency bands (Tebaldini and Ferro-Famil 2013). The geometrical configuration of a TomSAR measurement may be characterized in terms of its height ambiguity, related to the spacing between acquisition tracks, and its Fourier vertical resolution, which depends on the extent of the aperture in elevation. Hence, for a regularly spaced array with a given ambiguous height, i.e. baseline, the vertical Fourier resolution improves with the number of images, although it may be affected by range decorrelation over natural environments (Gatelli et al. 1994). Realistic and unambiguous TomSAR data sets being generally composed of a moderate, or small, number of acquisitions, 3-D focusing using classical Fourier imaging techniques may lead to a poor vertical resolution. As it is shown in this paper, SAR tomography may be considered as a spectral estimation problem, and a wide variety of techniques may be used to focus data in the vertical direction with a substantially improved resolution (Fornaro and Serafino 2006; Budillon et al. 2011; Zhu and Bamler 2012; Sauer et al. 2011; Huang et al. 2012). The combination of SAR tomography with polarimetric diversity permits to both improve the performance of SAR tomography in terms of separation between different components of complex volumetric media and to further characterize the EM behavior of the observed environments through the estimation of scattering mechanisms (Sauer et al. 2011; Huang and Ferro-Famil 2009; Huang et al. 2012; Tebaldini 2009; Tebaldini and Rocca 2012; Aguilera et al. 2013; Minh et al. 2014; Frey and Meier 2011; Neumann et al. 2009). Part 2 presents the basic principles of SAR tomography and introduces polarimetric TomSAR (PolTomSAR) signal models corresponding to different kinds of observed media. Polarimetric spectral estimation techniques approaches are then introduced and compared in the frame of 3-D High Resolution (HR) focusing. In part 3 are shown some application results, obtained

with PolTomSAR data sets acquired in various configurations and dealing with urban area mapping with a minimal data set, tropical forest structure characterization and under-foliage concealed vehicle imaging.

## 2 Basics of Polarimetric SAR Tomographic Imaging

### 2.1 Geometrical Configuration

SAR tomographic (TomSAR) imaging is based on the acquisition of  $M$  SAR signals along slightly shifted trajectories  $t_i$ , as illustrated on Fig. 1. After focusing, compensating and co-registering the acquired signals,  $M$  2-D SAR images are obtained,  $s_i(x, r)$ , where  $x$  and  $r$  represent azimuth and slant range coordinates, respectively. As shown on Fig. 1, due to its intrinsic cylindrical ambiguity, classical 2-D SAR imaging is



**Fig. 1** Geometry of a TomSAR measurement. (a) Tomographic acquisition, (b) tomographic vs 2-D SAR resolution cells

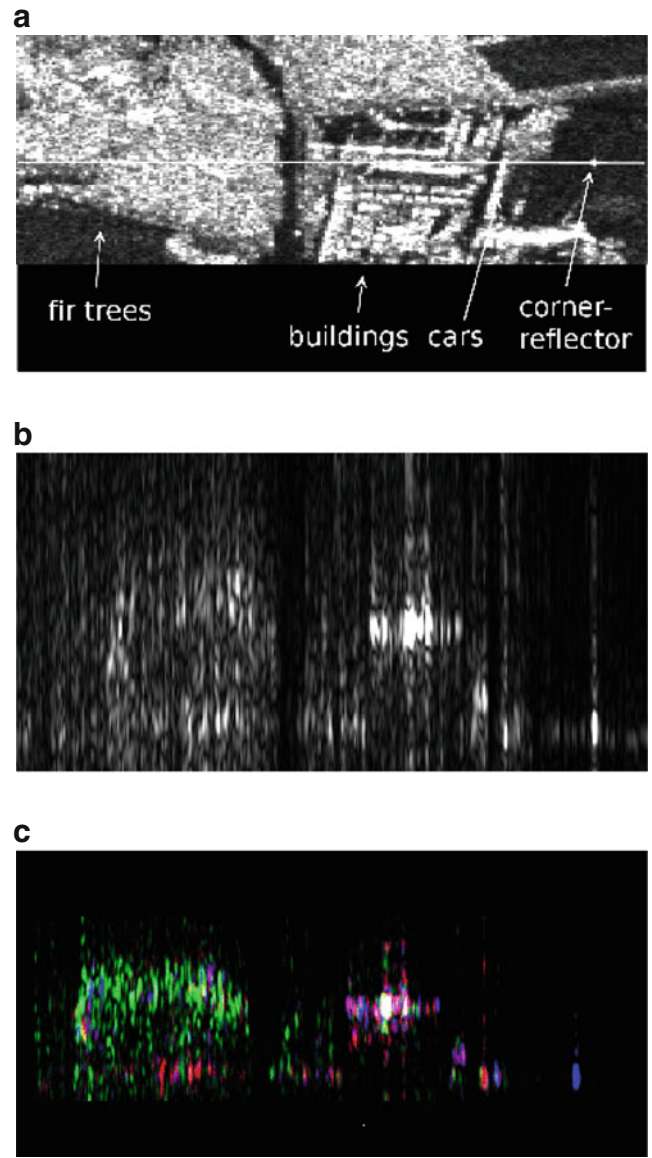
not well adapted to the characterization of volumetric media. For a given focusing position  $(x_0, r_0)$ , and considering a simplified rectangular azimuth-range impulse response, the focused SAR results from the coherent integration of the reflectivity density  $a(\mathbf{r})$ , with  $\mathbf{r} = (x, y, z)$ , of the measured medium, as:

$$s_i(x_0, r_0) = \int_{\mathcal{C}(x_0, r_0)} a(\mathbf{r}) e^{j\mathbf{k}_i \cdot \mathbf{r}} d\mathbf{r} \quad i = 1 \dots M \quad (1)$$

where  $\mathbf{k}_i$  stands for the two-way wavenumber of the  $i^{\text{th}}$  acquisition and  $\mathcal{C}(x_0, r_0)$  represents one resolution cell with a cylindrical slice shape, whose limits are given by  $|x - x_0| < \frac{\delta x}{2}$ ,  $|r - r_0| < \frac{\delta r}{2}$ , and  $z_{v_{\min}} < z < z_{v_{\max}}$ , with  $\delta x$  and  $\delta r$ , the SAR resolution in azimuth and range, respectively, and  $z_{v_{\min}}, z_{v_{\max}}$ , the upper and lower limits of the observed volumetric medium. Depending on the nature of the environment under observation, the reflectivity density of an environment may be described by a sum of discrete contributions generated by a set of point-like scatterers,  $a(\mathbf{r}) = \sum_i a_i \delta(\mathbf{r} - \mathbf{r}_i)$ , or under the form of a continuous function, having a stochastic behavior conferred by the speckle phenomenon, i.e. that may be assimilated to a spatially modulated white noise, with  $E(a(\mathbf{r})) = 0$  and  $E(a(\mathbf{r})a^*(\mathbf{r} + d\mathbf{r})) = \sigma(\mathbf{r})\delta(d\mathbf{r})$  (Bamler and Hartl 1998). Assuming that  $a(\mathbf{r})$  represents a continuous reflectivity density which is a function of the elevation position only, i.e. that  $\delta x$  and  $\delta y$  have sufficiently low values and variations in these directions within a resolution cell may be neglected, the cross-correlation between two images,  $\rho_{ij} = E(s_i s_j^*)$  may be expressed as:

$$\rho_{ij}(x_0, r_0) = \int_{\mathcal{C}(x_0, r_0)} \sigma_{ij}(z) e^{j(k_{z_i} - k_{z_j})z} dz \quad (2)$$

where the geometry-dependent vertical wavenumber,  $k_{z_i}$ , accounts for variations of the interferometric phase from one image to the other. The derivation of (2) from (1) is based on the fact that  $a_{i,j}(\mathbf{r})$  are spatially white random processes, i.e.  $E(a_i(\mathbf{r})a_j^*(\mathbf{r} + d\mathbf{r})) = \sigma_{ij}(\mathbf{r})\delta(d\mathbf{r})$ , where  $|\sigma_{ij}(\mathbf{r})| \leq \sqrt{\sigma_i(\mathbf{r})\sigma_j(\mathbf{r})}$  is the temporarily stable backscatter coefficient common to both measurements (Askne et al. 1997). The purpose of TomSAR imaging, illustrated on Fig. 1, is to improve the vertical resolution of the SAR measurement in order to characterize the volumetric backscatter coefficient  $\sigma(z)$  in a more accurate way. To do so,  $\rho_{ij}$  is computed for all the available image pairs, and the resulting set of coefficients are processed using (2) in the frame of a classical problem of spectral estimation, for which a series of solutions can be applied (Gini and Lombardini 2005; Stoica and Moses



**Fig. 2** First airborne tomographic SAR experiment. Courtesy of Dr. Andreas Reigber, DLR, Germany. (a) HH SAR image, (b) single-polarization tomogram along a path, (c) polarimetric tomogram  $\frac{|hh-vv|}{|hv|}$

2005; Fornaro and Serafino 2006; Budillon et al. 2011; Zhu and Bamler 2012; Sauer et al. 2011; Huang et al. 2012). The first airborne TOMSAR experiment was conducted by the DLR using their ESAR sensor at L band over the test site of Oberpfaffenhofen, Germany (Reigber and Moreira 2000). The data set consists of 14 SAR images acquired over quasi-parallel tracks within a short period of time. Some results of this pioneer work, presented in Fig. 2, show that tomography can be used to reliably locate scatterers in elevation, estimate building heights and image forest canopies.

## 2.2 Polarimetric SAR Tomography

Polarimetric SAR measurements are generally performed using an orthogonal polarization basis, whose elements may be used at both the emission and reception of SAR signals. In the case of horizontally or vertically polarized antennas, the fully polarimetric response for a given resolution cell may be represented using a scattering matrix, given, in the  $(h, v)$  basis, by:

$$\mathbf{S} = \begin{bmatrix} S_{hh} & S_{hv} \\ S_{vh} & S_{vv} \end{bmatrix} \in \mathbb{C}^{2 \times 2} \quad (3)$$

where  $S_{pq}$ , with  $p, q = h$  or  $v$  represents the scattering coefficient when the EM signal is emitted through the polarization channel  $q$  and received on channel  $p$ . Polarimetric SAR systems generally have collocated or quasi-collocated antenna systems and in this case the scattering matrix is symmetric, i.e.  $S_{hv} = S_{vh}$  (Lee and Pottier 2008). Tomograms computed over different polarimetric channels may be combined to appreciate the 3-D polarimetric behavior of the media under observation, as schematically represented in Fig. 1 and shown in Fig. 2. Specific types of scattering mechanism are associated to colors: blue represents single-bounce reflection over rough surface, i.e. ground or building roofs, red can be associated to double-bounce reflections, characteristic of ground-tree trunk or dihedral-like objects, whereas green indicates scattering by anisotropic particles and can be associated to forest canopy in this case. The study presented in Reigber and Moreira (2000) demonstrated for the first time the potential of polarimetric diversity for the 3-D characterization of volumetric media. As this is shown in the following, PolTomSAR may be used not only to identify basic scattering mechanism, but also to further discriminate different components of complex media and estimate some of their geo-physical parameters.

## 2.3 Tomographic Signal Models

In order to be estimated with numerical techniques, the density of reflectivity,  $a(x, r, z)$ , used in the SAR signal formulation of (1), is generally associated with a set of discrete contributions, called sources. Discrete media, like urban environments where a resolution cell contains a low number of dominant contributions, are well modeled using a small number of sources. Oppositely, natural media having a continuous density of reflectivity, like forests, require an infinite number of sources to be modeled adequately.

### 2.3.1 Single-Polarization TomSAR Signal Models

Considering an azimuth-range resolution cell that contains  $n_s$  backscattering contributions from scatterers located at

different heights and assuming no decorrelation between the different acquisitions, the data vector measured by  $M$  SAR acquisitions,  $\mathbf{y} \in \mathbb{C}^{M \times 1}$ , can be formulated as follows:

$$\mathbf{y}(l) = \sum_{i=1}^{n_s} s_i \mathbf{a}(z_i) + \mathbf{n}(l) = \mathbf{A}(\mathbf{z})\mathbf{s} + \mathbf{n}(l) \quad (4)$$

where  $l = 1, \dots, L$  indicates one of the  $L$  independent realizations of the signal acquisition, also called looks. The source signal vector,  $\mathbf{s} = [s_1, \dots, s_{n_s}]^T$ , contains the unknown complex reflection coefficient of the  $n_s$  scatterers, and  $\mathbf{n} \in \mathbb{C}^{M \times 1}$  represents the complex additive noise, assumed to be Gaussianly distributed and to be white in time and space, i.e.  $\mathbf{n} \sim \mathcal{N}(\mathbf{0}, \sigma_n^2 \mathbf{I}_{(M \times M)})$  and  $E(\mathbf{n}(l)\mathbf{n}^\dagger(k)) = \sigma_n^2 \mathbf{I}_{(M \times M)} \delta_{l,k}$ . The interferometric phase information associated to a source located at the elevation position  $z$  above the reference focusing plane is given by the steering vector  $\mathbf{a}(z) = [1, \exp(jk_{z_2}z), \dots, \exp(jk_{z_M}z)]^T$ , where  $k_{z_j}$  is the two-way vertical wavenumber between the master and the  $j^{\text{th}}$  acquisition tracks. The steering matrix  $\mathbf{A}(\mathbf{z})$  consists of  $n_s$  steering vectors corresponding each to a backscattering source  $\mathbf{A}(\mathbf{z}) = [\mathbf{a}(z_1), \dots, \mathbf{a}(z_{n_s})]$  with  $\mathbf{z} = [z_1, \dots, z_{n_s}]^T$ , the vector of unknown source heights. Considering now interferometric decorrelation between different acquisitions, the initial model in (4) may be reformulated as a sum of contributions from random sources (Gini and Lombardini 2005):

$$\mathbf{y}(l) = \sum_{i=1}^{n_s} \mathbf{x}_i(l) \odot \mathbf{a}(z_i) + \mathbf{n}(l) \quad (5)$$

where  $\odot$  represents the element-wise product between two vectors and  $\mathbf{x}_i \in \mathbb{C}^{M \times 1}$  accounts for both the reflection coefficient of the  $i$ th source,  $s_i$  and its potential variations between the  $M$  acquisitions or over the  $L$  realizations. Depending on the type of scatterer under observation, the composite signal  $\mathbf{y}(l)$  may follow different behaviors, that are linked to the statistical properties of the source signal  $\mathbf{x}_i$  (Huang et al. 2012).

For *distributed scatterers*, characterized by a scattering response having a random behavior conferred by the speckle effect, the received signal component  $\mathbf{y}_{u_i}(l)$  may be represented using an *Unconditional Model* (UM) (Stoica and Nehorai 1990), which accounts for the random nature of the source signal using a multiplicative term given by  $\mathbf{x}_i(l) = s_i \mathbf{x}_{u_i}(l) \in \mathbb{C}^{M \times 1}$ , with  $\mathbf{x}_{u_i}(l) \sim \mathcal{N}(\mathbf{0}, \mathbf{C}_i)$  (Gini and Lombardini 2005). The reflectivity of the source is given by  $\sigma_i = |s_i|^2$ , whereas the  $(M \times M)$  covariance matrix  $\mathbf{C}_i$  describes the interferometric coherence for a source. This kind of source signal is well adapted to the modeling of scattering over natural environments, like rough surfaces, ground and natural volumes.

*Deterministic scatterers* have a highly coherent EM response and their response,  $\mathbf{y}_{c_i}(l)$ , may be represented using a *Conditional Model* (CM) (Stoica and Nehorai 1990), which considers a frozen source signal over all the observations,  $\mathbf{x}_i(l) = s_i \mathbf{1}_{(M \times 1)}$ . This behavior is generally related to specular scattering mechanisms and can be observed over coherent scatterers like calibrators, facets facing the radar, double-bounce reflections over dihedral-like objects having smooth surfaces (like ground-trunk interactions and double bounce reflections between an object and the ground), or may be linked to resonant behaviors over quasi-periodic media (Ferro-Famil et al. 2005; Ferro-Famil and Pottier 2007).

*Hybrid scatterers* consist of a mixture of coherent and distributed scatterers (Ferro-Famil et al. 2005; Ferro-Famil and Pottier 2007) and may be modeled as  $\mathbf{y}(l) = \mathbf{y}_c(l) + \mathbf{y}_u(l)$  (Sauer et al. 2011), where each component, conditional or unconditional, is composed of several contributions. This type of signals can be frequently encountered when dealing with intermediate-resolution SAR images.

### 2.3.2 PolTomSAR Signal Models

The polarimetric scattering matrix, introduced in (3), can be vectorized using, for instance, the Pauli basis matrix set (Lee and Pottier 2008), in order to build a strictly equivalent target vector,  $\mathbf{v} = \frac{1}{\sqrt{2}}[S_{hh} + S_{vv}, S_{hh} - S_{vv}, 2S_{hv}] = s \mathbf{k}$  where  $\sigma = |s|^2 = \mathbf{v}^\dagger \mathbf{v}$  represents the polarimetric span (Lee and Pottier 2008), or reflectivity, of the scatterer response, and  $\mathbf{k} = [k_1, k_2, k_3]^T \in \mathbb{C}^{3 \times 1}$  represents a unitary polarimetric target vector, i.e.  $\mathbf{k}^\dagger \mathbf{k} = 1$ . In a PolTomSAR configurations,  $M$  polarimetric SAR signals are acquired and denoted  $\mathbf{v}_j$ , with  $j = 1, \dots, M$  the track index. A  $3M$  element PolTom received signal,  $\mathbf{y}_P$ , is then formed by stacking the TomSAR responses for each polarization channel  $\mathbf{y}_P = [\mathbf{y}_1^T, \mathbf{y}_2^T, \mathbf{y}_3^T]^T \in \mathbb{C}^{3M \times 1}$  where  $\mathbf{y}_p \in \mathbb{C}^{M \times 1}$ , with  $p = 1, 2$  or  $3$ , represents the TomSAR response for the  $p$ th polarimetric channel, i.e.  $[\mathbf{y}_p]_j = [\mathbf{v}_j]_p$ . Using this convention of representation, the polarimetric steering vector of the  $i$ th source is given by  $\mathbf{a}(z_i, \mathbf{k}_i) = \mathbf{k}_i \otimes \mathbf{a}(z_i)$  and may be parametrized using 5 real coefficients. The corresponding steering matrix writes  $\mathbf{A}(\mathbf{z}, \mathbf{K}) = [\mathbf{a}(z_1, \mathbf{k}_1), \dots, \mathbf{a}(z_{n_s}, \mathbf{k}_{n_s})] \in \mathbb{C}^{M \times n_s}$  with  $\mathbf{K} = [\mathbf{k}_1, \dots, \mathbf{k}_{n_s}]$ . Similarly to the single polarization expression given in (4), the received PolTomSAR signal  $\mathbf{y}_P(l)$  may be formulated as

$$\mathbf{y}_P(l) = \mathbf{A}(\mathbf{z}, \mathbf{K})\mathbf{s}(l) + \mathbf{n}(l) \in \mathbb{C}^{3M \times 1} \quad (6)$$

with  $\mathbf{s}(l) \in \mathbb{C}^{n_s \times 1}$ , a realization of the complex response of the  $i$ th source. Diverse model assumptions given for single-polarization signals, can be similarly used for in the PolTomSAR case.

## 2.4 Tomographic Focusing Techniques

In order to preserve the compactness of this paper, the presentation of the different tomographic focusing techniques is restricted in the following to the general description of their principles and performance. More details and references to specific studies may be found in Gini and Lombardini (2005), Huang et al. (2012), Stoica and Moses (2005) and Sauer et al. (2011). The objective of tomographic focusing is to estimate the reflectivity, scattering vector and height  $\sigma, \mathbf{k}_i, z_i$  of each source using from the covariance matrix of the received signal  $\mathbf{R} = \mathbb{E}(\mathbf{y}\mathbf{y}^\dagger) = \mathbf{A}\mathbf{R}_{xx}\mathbf{A}^\dagger + \sigma_n^2 \mathbf{I}_{M \times M}$ . In practice, for a locally Gaussian statistical behavior, a maximum likelihood estimate of  $\mathbf{R}$  may be computed from  $L$  independent locations surrounding the pixel under analysis, as  $\hat{\mathbf{R}} = \frac{1}{L} \sum_{l=1}^L \mathbf{y}(l)\mathbf{y}^\dagger(l)$  and may be used instead of  $\mathbf{R}$  to perform tomography. The number of sources,  $n_s$ , is in general unknown and needs to be estimated from the measured data. Some commonly used Model Order selection techniques based on statistical approaches, e.g. ITC, MDL, AIC (Wax and Ziskind 1991), may be used to determine  $\hat{n}_s$ . Once  $n_s$  is determined, eigenstructure matrices can be estimated from the sample covariance matrix  $\hat{\mathbf{R}} = \hat{\mathbf{E}}_s \hat{\mathbf{\Lambda}}_s \hat{\mathbf{E}}_s^\dagger + \hat{\mathbf{E}}_n \hat{\mathbf{\Lambda}}_n \hat{\mathbf{E}}_n^\dagger$  where  $\hat{\mathbf{E}}_s$  and  $\hat{\mathbf{E}}_n$  are represent the estimated signal and noise subspaces, respectively.

### 2.4.1 Single Polarization Tomography

Single Polarization (SP) tomography can be derived from the data covariance matrix  $\hat{\mathbf{R}}$  using spectral estimation techniques. Some examples of mono- and multi-dimensional estimators are given in the following.

#### Mono-dimensional Estimators

These approaches determine  $\hat{\mathbf{z}}$ , an estimate of the elevation of the scatterers under observation, as the coordinates of the  $\hat{n}_s$  largest local maxima of an continuous objective function  $P(\mathbf{z})$ ,  $\hat{\mathbf{z}} = \arg \max_{\mathbf{z}, loc} P(\mathbf{z})$ . The Fourier Beamformer (FB) and Capon's method, also called the adaptive beamformer, are non-parametric spectral estimation techniques, i.e. they do not require to estimate the number of sources, and their objective function is given by the continuous estimate of the reflectivity, obtained by height-varying linear filtering as  $P_{B,C}(\mathbf{z}) = \hat{\sigma}_{B,C}(\mathbf{z}) = \mathbf{h}_{B,C}^\dagger(\mathbf{z}) \hat{\mathbf{R}} \mathbf{h}_{B,C}(\mathbf{z})$ , with  $\mathbf{h}_B^\dagger(\mathbf{z}) = \mathbf{a}^\dagger(\mathbf{z})/\sqrt{M}$  and  $\mathbf{h}_C^\dagger(\mathbf{z}) = \mathbf{R}^{-1} \mathbf{a}/(\mathbf{a}^\dagger \mathbf{R}^{-1} \mathbf{a})$  (Stoica and Moses 2005). The selection of discrete sources from peaks of the reflectivity spectrum confers to the FB and Capon estimation techniques an important sensitivity to the acquisition configuration, and in particular to the presence of spurious sidelobes related to an irregular baseline sampling. The FB is known to show a low resolution and may then overlook

some closely spaced scatterers, whereas Capon's technique possesses an improved resolution but a reduced radiometric accuracy. MUSIC is a subspace-based mono-dimensional technique (Stoica and Moses 2005), whose objective function is a measure of the orthogonality between a steering vector  $\mathbf{a}(z)$  and the estimated noise subspace  $\hat{\mathbf{E}}_n$  and is given by:  $P_M(z) = 1/||\mathbf{a}^\dagger(z)\hat{\mathbf{E}}_n||^2$ . An estimate of the reflectivity vector  $\hat{\mathbf{s}}$  can be obtained from  $\hat{\mathbf{z}}$  using a Least Squares (LS) approach (Stoica and Moses 2005; Gini and Lombardini 2005). Nonparametric approaches like FB and Capon, are generally used to globally appreciate the structure of a volumetric medium and the main trends of the continuous reflectivity distribution in elevation. For the analysis of discrete spectral components, they may fail to discriminate closely spaced scatterers due either to their limited resolution, or to the presence of side lobes that may induce an erroneous estimation of the source location. MUSIC generally presents better resolution and performance for the analysis of discrete sources, but, like all parametric methods, MUSIC is sensitive to data modeling errors, and in particular those related to the estimated number of sources  $\hat{n}_s$ . Moreover, MUSIC is known to perform badly in the presence of correlated scatterers, due to the singularity of the the source signal covariance matrix (Stoica and Nehorai 1990). One of the main advantages of such techniques resides in the low numerical complexity of the mono-dimensional optimization they are based on.

### Multi-dimensional Estimators

Maximum Likelihood (ML) techniques aim to estimate  $(\hat{\mathbf{z}}, \hat{\boldsymbol{\sigma}}, \hat{\sigma}_n^2)$ , the set of parameters maximizing  $\mathcal{L}(\mathbf{z}, \boldsymbol{\sigma}, \sigma_n^2)$ , the data likelihood function, where the index  $u, c$  indicate the statistical model, UM or CM, under consideration. The optimization being separable in  $\hat{\mathbf{z}}$  and  $\hat{\boldsymbol{\sigma}}$ , highly concentrated expressions have been derived in both cases (Ottersten et al. 1993; Stoica and Nehorai 1990). Weighted Subspace Fitting (WSF) techniques are based on the comparison of the eigenstructure of the data covariance matrix with the model accounting for a number of steered sources plus noise: for a correct height vector guess,  $\mathbf{A}(\mathbf{z})$  should be orthogonal to the noise space and parallel to the signal one. Two weighted LS fitting cost functions have been presented in the literature Viberg et al. (1995),  $Q_{NSF}(\mathbf{z}) = ||\hat{\mathbf{E}}_n^\dagger \mathbf{A}(\mathbf{z})||_{\mathbf{W}}^2$  and  $Q_{SSF}(\mathbf{z}) = ||\hat{\mathbf{E}}_s - \mathbf{A}(\mathbf{z})\mathbf{T}||_{\mathbf{W}}^2$  where NSF and SSF respectively stand for Noise and Signal Subspace Fitting. The linear transformation matrix  $\mathbf{T}$  can be replaced by its LS estimate and the source heights are found using  $\hat{\mathbf{z}}_{WSF} = \arg \min_{\mathbf{z}} Q_{WSF}(\mathbf{z})$  where the suffix WSF indicates ones of the methods, NSF or SSF, mentioned above. The weighting matrix,  $\mathbf{W}$ , aims to correct for potential discrepancies between the considered model and real data as well as for errors occurring during the estimation of the data covariance matrix. It has been shown in Viberg et al.

(1995), that any hermitian positive semi definite weighting matrix  $\mathbf{W}$  yields consistent parameter estimates, and that a consistent estimate of  $\mathbf{W}$  permits to obtain minimum variance estimates, which asymptotically reach the Cramer-Rao lower bound (Ottersten et al. 1993). The parametric methods presented in this section require a  $n_s$ -dimensional minimization. Compared with mono-dimensional ones, these methods are more robust and generally lead to global optima, but at an expensive computational cost. WSF techniques may reach the same level of accuracy than ML techniques at a reduced computational cost and are susceptible to better adapt to complex scattering configurations.

### 2.4.2 Fully Polarimetric Tomography

Similarly to the SP case, Fully Polarimetric (FP) tomography can be performed from the data covariance matrix  $\hat{\mathbf{R}}_P = \frac{1}{L} \sum_{l=1}^L (\mathbf{y}_P(l)\mathbf{y}_P^\dagger(l))$  using mono- and multi-dimensional FP spectral estimators, whose objectives are to estimate the source locations in elevation,  $\hat{\mathbf{z}}$ , their polarimetric target vectors,  $\hat{\mathbf{K}}$ , and their reflectivities  $\boldsymbol{\sigma}$ .

#### Mono-dimensional Estimators

The elevation and scattering mechanisms of the different scatterers are estimated as the coordinates of the  $\hat{n}_s$  largest local maxima of a polarimetric objective function  $P(\mathbf{z}, \mathbf{k})$ , obtained by replacing the steering vector  $\mathbf{a}$  by its FP expression  $\mathbf{a}(\mathbf{z}, \mathbf{k})$  in the FB, Capon and MUSIC criteria mentioned in the preceding section. The optimization of the resulting cost functions with respect to the polarimetric target vector can be performed efficiently using an eigendecomposition. This important aspect maintains the complexity of the focusing a value close to the one obtained in the SP case. These methods are computationally efficient but may reach some of the limitations mentioned in the SP case.

#### Multi-dimensional Estimators

The FP-ML and FP-WSF estimators may be formulated from the SP expressions by replacing the SP steering matrix  $\mathbf{A}(\mathbf{z})$  by  $\mathbf{A}(\mathbf{z}, \mathbf{K})$ . The optimization of the corresponding criteria requires, for  $\hat{n}_s$  sources, a search of the optimal parameters over a  $5\hat{n}_s$ -dimensional space, and general implies an excessive computational burden. Some computationally optimized methods have been proposed for the ML techniques in Wax and Ziskind (1991); Ottersten et al. (1993), whereas an analytical solution proposed in Huang et al. (2012) permits to maintain the computational cost of the FP-NSF to the one of the SP case. Compared with SP tomography, polarimetric tomography can localize scatterers more accurately due to the additional power of discrimination brought by polarization diversity, and the physical properties of the observed media can be further characterized from their scattering mechanism.

### 2.4.3 Performance Assessment

The performance of the aforementioned tomographic estimators is investigated with a data set of  $M = 5$  SAR images, simulated with a regular spacing in elevation  $\Delta k_z = 0.1$ , with  $SNR = 20\text{dB}$  and  $L = 256$  independent realizations. The simulation considers two sources with height difference  $\Delta h$  and whose statistical behavior is steered by the correlation coefficient  $\rho$ . Coherent scatterers are obtained for  $\rho = 1$ , whereas distributed ones correspond to  $\rho = 0$  and hybrid scatterers have  $0 < \rho < 1$ . Figure 3a illustrates the interesting resolution properties of the WSF estimators for uncorrelated scatterers, as they work well for  $\Delta h \geq 0.4m$ , while the performance of MUSIC and Capon's techniques degrade significantly for  $\Delta h < 2m$  and  $\Delta h < 4m$ , respectively. Figure 3b shows that the NSF estimator provides the most accurate estimate for uncorrelated or partially correlated signals ( $\rho < 0.95$ ), whereas the SSF estimator copes well with highly correlated signals ( $\rho \geq 0.95$ ). MUSIC cannot deal with highly correlated sources due to its extreme sensitivity to the quasi-singularity of  $\mathbf{R}_{xx}$ . Similar simulations are run in the FP case, with 3-element scattering vectors,  $\mathbf{k}_1$  and  $\mathbf{k}_2$ , separated by an angular distance,  $\zeta$ , defined as  $\cos \zeta = \frac{|\mathbf{k}_1^\dagger \mathbf{k}_2|}{\|\mathbf{k}_1\| \|\mathbf{k}_2\|}$ . Figure 3c shows that polarization diversity between two scatterers improves height resolution, especially for the FP-Capon estimator which reaches the same resolution than MUSIC when  $\zeta > 60^\circ$ . The FP-NSF estimator performs best for any polarization similarity value.

## 3 Applications

The potential of PolTomSAR techniques for the 3-D characterization of complex environments is illustrated in the following with three very different applications.

### 3.1 Urban Remote Sensing Using a Minimal TomSAR Configuration

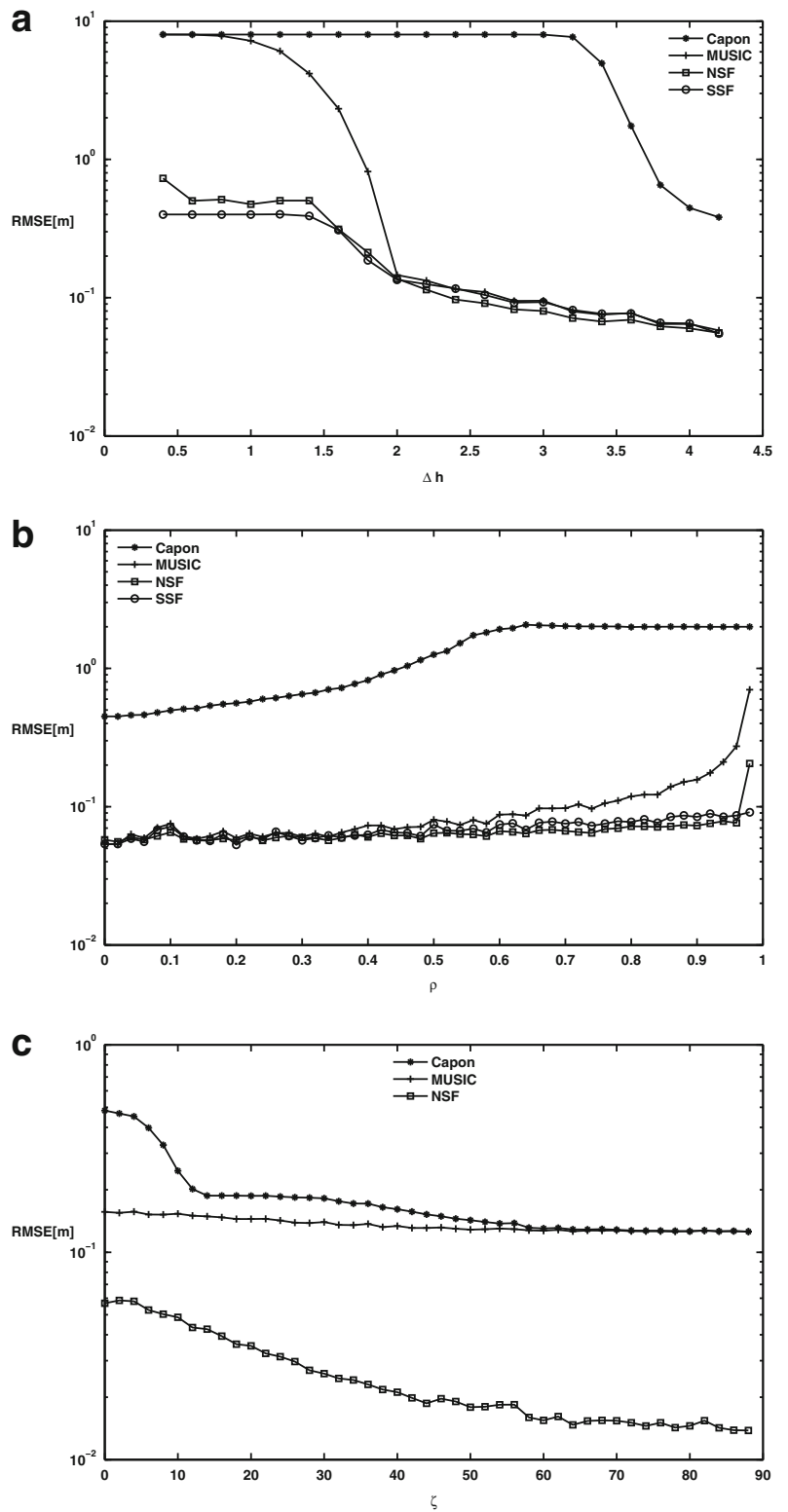
Dense urban environments may generate complex SAR responses. As depicted in Fig. 4a, diverse scattering patterns, like double bounce reflection due to wall-ground reflections, surface scattering from roofs and from the ground, volumetric scattering due to potential vegetation, can be encountered over such volumetric media, and may fall into a single SAR resolution cell. The test data set for this application was acquired over Dresden, Germany, by the DLR ESAR sensor at L band. It consists of three fully polarimetric SAR images with an intermediate resolution of  $1.5m \times 3m$ , measured over a short period of time. Unlike other approaches, based on the use of time series of numerous SAR acquisitions (Fornaro and Serafino 2006;

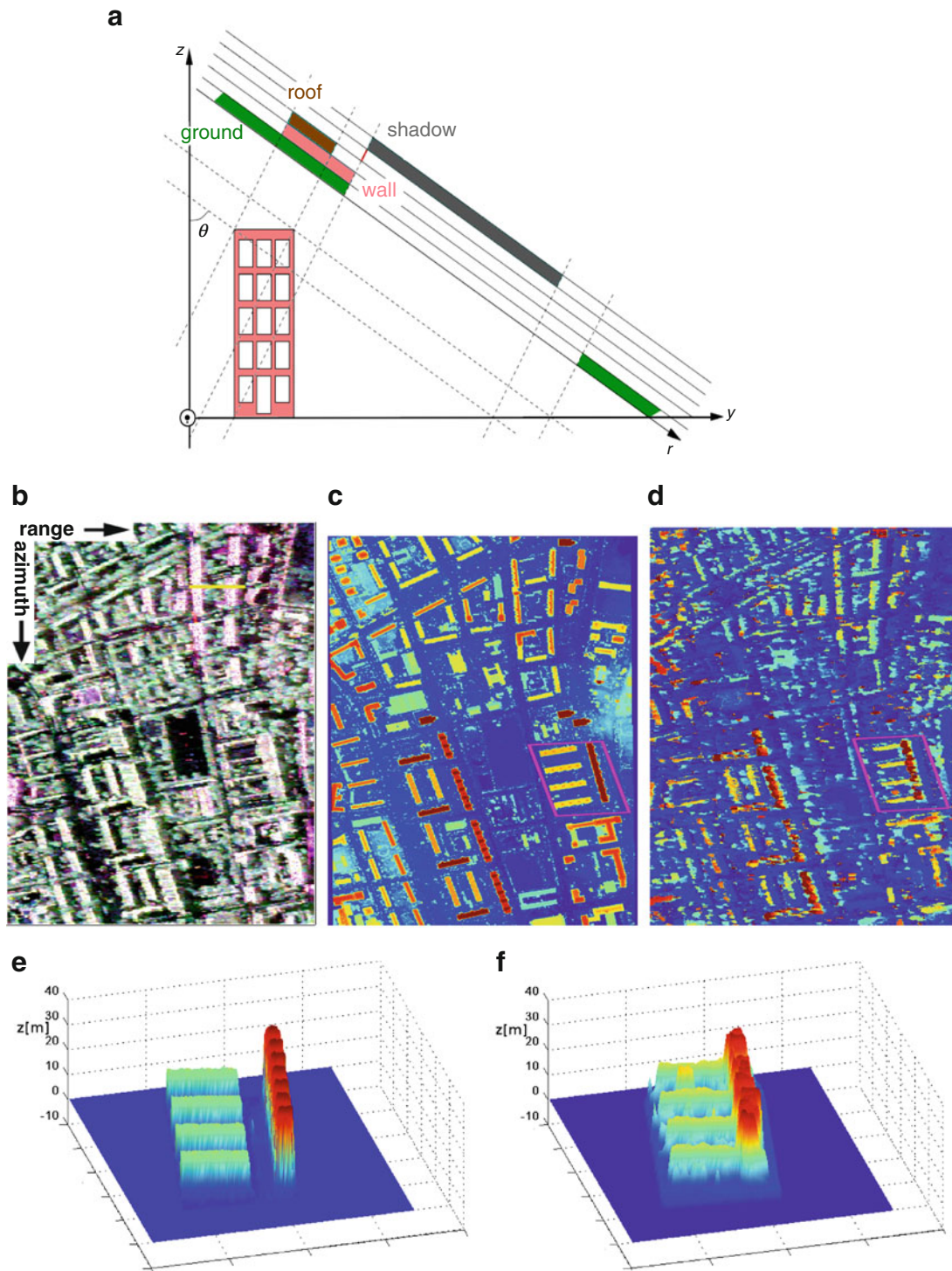
Budillon et al. 2011; Zhu and Bamler 2012), this applications deals with the minimum number of images required to perform tomography and may be considered as an extreme case of 3-D imaging using coherent SAR data processing. The vertical Fourier resolution being particularly coarse, around 10m, the accurate mapping of building shapes requires the use of HR estimation techniques. Here again, this application differentiates itself from those based on time series analysis, for which sparse signal reconstruction techniques are well adapted (Budillon et al. 2011; Zhu and Bamler 2012), since the number of available images, equal to three, is comparable to the number of dominant responses that can be expected within a resolution cell. Over the different spectral estimators presented above, the model-adaptive FP-NSF technique proposed in Huang and Ferro-Famil (2009) and Huang et al. (2012) was found to be the most suitable for such an extreme configuration as it fully exploits polarimetric diversity in order to further separate closely spaced scatters and can benefit from the high level of coherence over natural environments in order to reduce potential biases induced by the observed vegetation covers. The results depicted in Fig. 4 indicate that over buildings and bare ground, PolTomSAR gives results that are comparable to Lidar measurements. Over parts covered by vegetation, associated to green colors in Fig. 4b, significant differences may be found due to the penetration of EM waves through the canopy of trees. The 3-D view of a specific set of buildings is given in Fig. 4f. One may note the very good matching between Lidar and PolTomSAR profiles, excepted for a gap between the four oriented buildings and the main one, which appears filled in the PolTomSAR reconstruction. This effect is due to the side looking geometry of a SAR acquisition, which unlike Lidar measurements is performed at incidence angles far away from nadir. The resulting projection of the scattering contributions in slant range occupies the space between the buildings.

### 3.2 Tropical Forest Characterization at P Band

This application deals with a six-image P-band data set, acquired over the tropical forest test area of Paracou, French Guiana, by the ONERA SETHI sensor, in the frame of the TropiSAR campaign whose objectives concerned the estimation of tropical forest biomass using SAR Minh et al. (2014). A tropical forests is considered here as a widespread canopy lying over a potential vegetation layer on the ground and the ground itself. From a spectral estimation point of view, the canopy is modeled as a continuous spectrum, whereas the ground is an impenetrable medium that possesses an isolated localized phase center and is characterized by a discrete spectrum. It has been verified in Huang et al. (2011) that, in

**Fig. 3** Tomographic performance simulation results. (a) TomSAR height rmse vs  $\Delta h$ ,  $\rho = 0$ , (b) TomSAR height rmse vs  $\rho$ ,  $\Delta h = 4m$ , (c) PolTomSAR height rmse vs  $\zeta(^{\circ})$ ,  $\rho = 0$

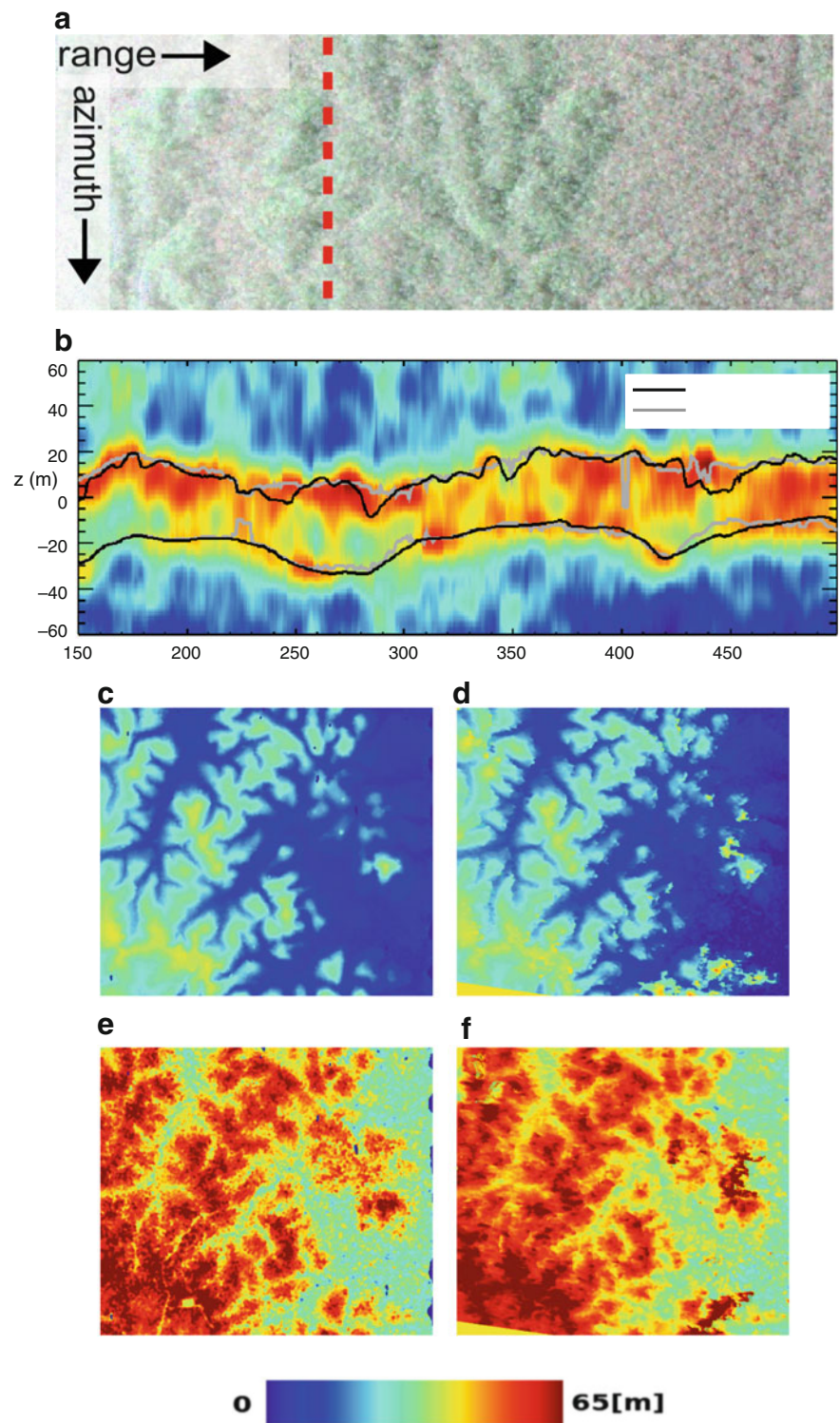


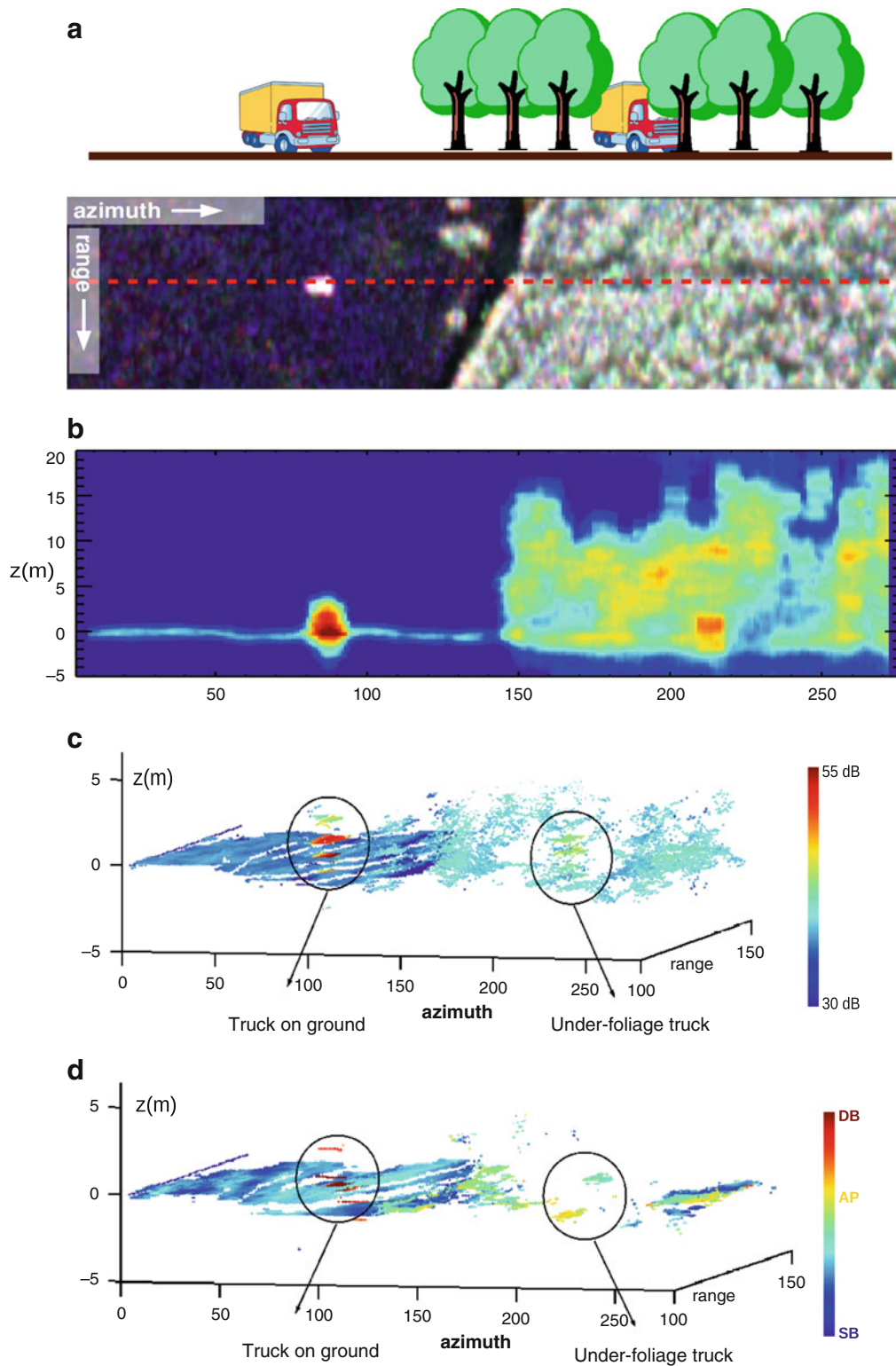


**Fig. 4** 3-D characterization of buildings in a dense urban environment. (a) SAR scattering patterns of a building, (b) PolSAR image, (c) Lidar h (0–30 m), (d) PolTomSAR h, (e) Lidar 3-D view (detail), (f) PolTomSAR 3-D view



**Fig. 5** Forest structure characterization using PolTomSAR and comparison with Lidar measurements. **(a)** Color-coded PolSAR image and path, **(b)** FP span tomogram along the path, **(c)** Lidar  $z_g$ , **(d)** estimated  $z_g$ , **(e)** Lidar  $z_{top}$ , **(f)** estimated  $z_{top}$





**Fig. 6** PolTomSAR imaging of an under-foilage object. (a) Color-coded PolSAR image, (b) FP-Capon span, (c) SSF intensity tomogram, (d) FP-NSF tomogram of scattering mechanisms

the presence of both CM and UM responses, the performance of tomographic estimators optimized for one of the models degrade due to their lack of adaptation. A novel hybrid tomographic approach proposed in Huang et al. (2011) in order to deal with mixed spectra, which provides robust estimates for both tree top heights,  $z_{top}$ , and the underlying ground topography,  $z_g$ , is presented here. This technique estimates the ground component using a high-performance parametric approach, e.g. the WSF estimator, and localizes the tree top by applying an energetic criterion on the response of the Capon focusing technique, well adapted to the analysis of continuous media, i.e. here the canopy (Tebaldini 2009). The results shown in Fig. 5b indicate that this hybrid PolTomSAR approach provides tree top and the ground elevation estimates that are similar to the LiDAR ones. In Fig. 5b, a comparison led at a local scale and over a larger area, shows that the estimated  $z_{top}$  and  $z_g$  match well with the Lidar ones in terms of texture and mean values, excepted for some overestimated areas, where topographic effects strongly affect the measured signals.

### 3.3 Under-Foliage Concealed Object Imaging

This application considers the detection of vehicles located under a temperate forest cover and observed by an airborne sensor at L band. The conditions of the experiment are sketched in Fig. 6a which indicates a selected profile containing the responses of a truck outside the forest cover and another hidden below the forest canopy. As it can be seen from the PolSAR image, the under-foliage truck response cannot be directly discriminated from the surrounding forest one using 2-D imaging. The data set for this study consists of 27 PolSAR images acquired by the DLR's ESAR sensor at L band. As shown in Fig. 6b, the reflectivity estimated by Capon's method permits to appreciate both the shape of the hidden vehicle and the canopy profile, as previously mentioned in Nannini et al. (2008), but considering the complex structure of such objects, HR approaches are required to both discriminate their closely spaced scattering features in the vertical direction and extract their response. The TomSAR response of an under-foliage object is considered as a deterministic component embedded within a speckle affected environment and hence is associated to a series of complex scattering centers. Such complex features require the use of statistically adaptive spectral analysis techniques, as it has been shown in Huang and Ferro-Famil (2009) and Huang et al. (2012). A 3-D reflectivity tomogram of the area, performed over azimuth bins and for varying range positions using the SP-SSF approach is displayed in Fig. 6c. Limiting the reconstruction height to 4 m above the terrain permits to isolate the under-foliage truck response,

with a reconstructed shape similar to the uncovered one, and with a slightly higher reflectivity than the one of the surrounding environment. A 3-D reconstruction of scattering mechanisms, i.e. Single Bounce reflection, Anisotropic Scattering and Double Bounce reflection, is performed using the FP-NSF estimator and common polarimetric analysis techniques (Lee and Pottier 2008). This approach permits to visualize the corresponding scattering pattern and shows that the truck outside the forest generates a strong double bounce reflection due to the ground-truck interaction. The shape of the truck beneath the canopy is precisely reconstructed with scattering mechanisms oscillating between SB and AS due to sidelobe effects from the canopy repines. Using PolTomSAR techniques, the under-foliage truck can be well described both in terms of shape and scattering patterns.

## 4 Conclusion

This paper has presented some principles of polarimetric SAR tomography, a technique able to characterize 3-D environment using electromagnetic waves. Some high resolution tomographic focusing techniques have been introduced and their advantages and drawbacks have been discussed. Some extensions handling polarimetric diversity have been given and the performance of different estimators have been compared. A set of studies dealing with real data, related to urban area and forest remote sensing or to the imaging of under-foliage objects, have been presented in order to illustrate the potential of this imaging mode for a wide variety of applications.

## References

- Aguilera E, Nannini M, Reigber A (2013) Wavelet-based compressed sensing for sar tomography of forested areas. *IEEE Trans Geosci Remote Sens* PP(99):1–13
- Askne J, Dammert P, Ulander LMH, Smith G (1997) C-band repeat-pass interferometric sar observations of the forest. *IEEE Trans Geosci Remote Sens* 35(1):25–35
- Bamler R, Hartl P (1998) Synthetic aperture radar interferometry. *Inverse Prob* 14(4):R1–R54
- Budillon A, Evangelista A, Schirinzi G (2011) Three-dimensional sar focusing from multipass signals using compressive sampling. *IEEE Trans Geosci Remote Sens* 49(1):488–499
- Cloude SR (2009) Polarisation applications in remote sensing. Oxford University Press, Oxford
- Ferro-Famil L, Pottier E (2007) Urban area remote sensing from L-band PolSAR data using time-frequency techniques. In: *Urban remote sensing joint event* 11–13 April 2007, p 1, 6
- Ferro-Famil L, Reigber A, Pottier E (2005) Nonstationary natural media analysis from polarimetric sar data using a two-dimensional time-frequency decomposition approach. *Can J Remote Sens* 31(1):20–29
- Fornaro G, Serafino F (2006) Imaging of single and double scatterers in urban areas via sar tomography. *IEEE Trans Geosci Remote Sens* 44(12):3497–3505

- Freeman A, Durden S (1998) A three-component scattering model for polarimetric sar data. *IEEE Trans Geosci Remote Sens* 36(3):963–973
- Frey O, Meier E (2011) 3-d time-domain sar imaging of a forest using airborne multibaseline data at l- and p-bands. *IEEE Trans Geosci Remote Sens* 49(10):3660–3664
- Gatelli F, Guamieri A, Parizzi F, Pasquali P, Prati C, Rocca F (1994) The wavenumber shift in sar interferometry. *IEEE Trans Geosci Remote Sens* 32(4):855–865
- Gini F, Lombardini F (2005) Multibaseline cross-track sar interferometry: a signal processing perspective. *IEEE Aerosp Electron Syst Mag* 20(8):71–93
- Huang Y, Ferro-Famil L (2009) 3-D characterization of buildings in a dense urban environment using L-band Pol-InSAR data with irregular baselines. In: *IEEE IGARSS 2009*, vol. 3, 12–17 July 2009, p III–29, III–32
- Huang Y, Ferro-Famil L, Lardeux C (2011) Polarimetric SAR tomography of tropical forests at P-Band. In: *IEEE IGARSS 2011*, 24–29 July 2011, p 1373, 1376
- Huang Y, Ferro-Famil L, Reigber A (2012) Under-foliage object imaging using sar tomography and polarimetric spectral estimators. *IEEE Trans Geosci Remote Sens* 50(6):2213–2225
- Lee JS, Pottier E (2009) Polarimetric radar imaging from basics to applications, 1st edn., February 2009, p 422
- Minh DHT, Toan TL, Rocca F, Tebaldini S, d’Alessandro M, Villard L (2014) Relating p-band synthetic aperture radar tomography to tropical forest biomass. *IEEE Trans Geosci Remote Sens* 52(2):967–979
- Nannini M, Scheiber R, Horn R, Moreira A (2012) First 3-D reconstructions of targets hidden beneath foliage by means of polarimetric SAR tomography. *IEEE Geosci Remote Sens Lett* 9(1):60, 64
- Mougin E, Proisy C, Marty G, Fromard F, Puig H, Betoulle JL, Rudant JP (1999) Multifrequency and multipolarization radar backscattering from mangrove forests. *IEEE Trans Geosci Remote Sens* 37(1):94–102
- Neumann M, Ferro-Famil L, Reigber A (2009) Estimation of forest structure, ground, and canopy layer characteristics from multibaseline polarimetric interferometric sar data. *IEEE Trans Geosci Remote Sens* 48(3):1086–1104
- Ottersten B, Viberg M, Stoica P, Nehorai A (1993) Radar array processing, chap. Exact and large sample approximations of maximum likelihood techniques for parameter estimation and detection in array processing. Springer, New York, pp. 99–151
- Papathanassiou K, Cloude S (2001) Single-baseline polarimetric sar interferometry. *IEEE Trans Geosci Remote Sens* 39(11):2352–2363
- Reigber A, Moreira A (2000) First demonstration of airborne sar tomography using multibaseline l-band data. *IEEE Trans Geosci Remote Sens* 38(5):2142–2152
- Sauer S, Ferro-Famil L, Reigber A, Pottier E (2011) Three-dimensional imaging and scattering mechanism estimation over urban scenes using dual-baseline polarimetric insar observations at l-band. *IEEE Trans Geosci Remote Sens* 49(11):4616–4629
- Stoica P, Moses R (2005) Spectral analysis of signals. Prentice Hall, Saddle River
- Stoica P, Nehorai A (1990) Performance study of conditional and unconditional direction-of-arrival estimation. *IEEE Trans Acoust Speech Signal Process* 38(10):1783–1795
- Tebaldini S (2009) Algebraic synthesis of forest scenarios from multibaseline polinsar data. *IEEE Trans Geosci Remote Sens* 47(12):4132–4142
- Tebaldini S, Ferro-Famil L (2013) High resolution three-dimensional imaging of a snowpack from ground-based sar data acquired at x and ku bands. In: *Proc. IGARSS*
- Tebaldini S, Rocca F (2012) Multibaseline polarimetric sar tomography of a boreal forest at p- and l-bands. *IEEE Trans Geosci Remote Sens* 50(1):232–246
- Treuhaft R, Siqueira P (2000) Vertical structure of vegetated land surfaces from interferometric and polarimetric radar. *Radio Sci* 35(1):141–178
- Viberg M, Stoica P, Ottersten B (1995) Array processing in correlated noise fields based on instrumental variables and subspace fitting. *IEEE Trans Signal Process* 43(5):1187–1199
- Wax M, Ziskind I (1991) Detection and localization of multiple sources via the stochastic signal model. *IEEE Trans Signal Process* 39(11):2450–2456
- Zhu XX, Bamler R (2012) Demonstration of super-resolution for tomographic sar imaging in urban environment. *IEEE Trans Geosci Remote Sens* 50(8):3150–3157

---

# Re-gridding and Merging Overlapping DTMS: Problems and Solutions in HELI-DEM

Ludovico Biagi and Laura Carcano

---

## Abstract

Neighboring or partly overlapping Digital Terrain Models (DTMs) that are stored as grids with similar resolutions and accuracies must be merged to produce a unified model: a reference frame transformation and a re-gridding on a common output grid have to be applied. A direct approach foresees firstly the reference frame transformation of the input DTMs, then their re-gridding by interpolation on the final grid. In the individual re-gridding of a DTM, no smoothing is required, because the input data are not raw observations but the nodes of an already smoothed model. Furthermore, re-gridding of huge amount of data could be required: therefore an efficient numerical method should be applied. A iso-determined interpolation approach based on local bicubic surfaces is investigated in this paper. In particular, the problem of an ill conditioned interpolation due to critical spatial distributions of the input elevations is discussed and different techniques to overcome the problem are compared. Finally, different approaches are discussed to average partly overlapping DTMs in their overlaps.

---

## Keywords

Bicubic interpolation • DTM • Ill conditioning • Re-gridding • Regularization

---

## 1 Introduction

Elevation data are numerically stored in Digital Elevation Models (DEM, EI Sheimy et al. 2005; Li et al. 2005), that can be realized by sampling elevations for a certain number of significant points and by storing the sample of 3D dimensional coordinates. Different data models can be adopted, such as contour lines, grids (or elevation matrices) and Triangular Irregular Networks (TIN): this paper is focused on grid models.

Gridded DEMs are georeferenced regular matrices of  $(x_i, y_i)$  nodes, whose elevations  $H_i$  are stored. The horizontal coordinates of the nodes can be either in a cartographic

projection ( $x$ : East,  $y$ : North) or geographic ( $x$ :  $\lambda$ ,  $y$ :  $\varphi$ ). Typically, the horizontal spacing between nodes (the grid resolution  $\Delta x$  and  $\Delta y$ ) is equal in  $x$  and  $y$  directions: this is not a strict requirement but quite a standard. The correct georeferencing of a grid requires the knowledge of the reference frame and coordinates system; then, additional metadata are needed: typically the grid origin (i.e. the coordinates of the lower left node),  $\Delta x$ ,  $\Delta y$  and the total number of rows and columns are stored.

DEM refers to the generic family of elevation models, that are distinguished in Digital Surface Models (DSMs), which represent the actual surface (including buildings, woods, etc.) and Digital Terrain Models (DTMs) which represent the elevation of bare soil: our focus is on DTMs. Stored elevations are typically orthometric heights: when the acquisition technique produces ellipsoidal heights, these are converted into orthometric.

Section 1.1 poses the general problems of re-gridding and merging of DTMs. In Sect. 2, HELI-DEM specific needs are

---

L. Biagi (✉) • L. Carcano  
Politecnico di Milano, DICA, Geomatics Laboratory at Como Campus,  
Via Valleggio 11, IT-22100 Como, Italy  
e-mail: ludovico.biagi@polimi.it

discussed. In the following sections, the methods and the results will be discussed.

## 1.1 The Addressed Problems

DTMs have to be merged that are georeferenced in different reference frames and originally gridded in different coordinates systems, with different resolutions and accuracies: they have to be transformed to a common reference frame and coordinates system and must be re-gridded on the output grid. Moreover, a proper merging technique has to be applied where they overlap.

A reference frame transformation is a three dimensional transformation of the X, Y, Z (or  $\varphi$ ,  $\lambda$ , h) coordinates of a point. However, a point of a DTM is a node of an horizontal grid for which the relevant elevation is stored: therefore, the reference frame transformation is applied only to the horizontal coordinates of the node. The output of the transformation of the whole DTM is a list of transformed horizontal nodes and their terrain elevations. The nodes are again almost regularly distributed but no more on an exactly oriented grid. Moreover, a re-gridding could be needed with a different resolution or in a different coordinates system. Therefore, an interpolation of the elevations is needed from the input to the output nodes. To transform and re-grid a DTM, two opposite approaches are possible.

1. A transformation of reference frame and coordinates system is applied to the input nodes, that are then used to interpolate elevations on the output grid. This approach will be called Direct Transformation.
2. The nodes of the output grid are back transformed to the input reference frame and coordinates system. The elevations of the input grid are interpolated on these horizontal coordinates, then the interpolated elevations are assigned to the relevant nodes of the output grid. This approach will be called Inverse Transformation.

In merging overlapping input DTMs, their preliminary cross-validation (Buckley and Mitchell 2004) is needed to check biases and anomalies. Even in most favorable cases, local differences exist and should be filtered. Two opposite approaches are possible.

1. The input DTMs are transformed to the output reference frame and coordinate system: a unified transformed dataset is created and is interpolated on the output grid. This approach will be called Merging by Interpolating a Unified Database (Merg-IUD).
2. Each input DTM is individually re-gridded on the output grid: the individual interpolations are averaged in the overlapping nodes. This approach will be called Merging by Averaging Individual Interpolations (Merg-AII).

Note that a Direct Transformation can be combined with both the merging approaches. On the contrary, the inverse transformation necessarily implies Merg-AII.

This paper studies the Direct Transformation approach and compares the two merging approaches, in order to better understand their advantages and disadvantages for HELI-DEM project goals.

## 1.2 HELI-DEM Project: Data Acquisition and Preliminary Analyses

HELI-DEM project (HELvetia-Italy Digital Elevation Model, Biagi et al. 2011, [www.helidem.eu](http://www.helidem.eu)) is funded by “PO Italia-Svizzera 2007–2013 Fondo Europeo di Sviluppo Regionale”: its main aim is the estimation of a unified DTM for the alpine and subalpine area between Italy (Piedmont, Lombardy) and Switzerland (Ticino and Grisons Cantons), with a resolution of about 20 m, by properly merging the available DTMs. In particular, four neighboring DTMs have been acquired for the project.

1. The regional Lombardy DTM: it is in Roma40 reference frame (Donatelli et al. 2002), Gauss Boaga projection. It has an horizontal resolution of 20 m, vertical standard deviation (ZS) of about 10 m (LE95 = 20 m) in the mountains.
2. Two DTMs are available for Piedmont: both are in ETRF89-IGM95, UTM projection. They have respectively: the former horizontal resolution of 50 m, ZS  $\approx$  2.5 m (LE95 = 5 m), the latter horizontal resolution of 5 m, ZS  $\approx$  1 m (LE95 = 2 m).
3. The SwissTopo DTMCH25 is available for Switzerland, gridded in geographic ETRF89 coordinates. It has resolution of 1'' (about 28 m in  $\varphi$  and 20 m in  $\lambda$ ), ZS  $\approx$  3 m (LE95 = 6 m) in mountain areas.

The input DTMs will be used to produce a unified output DTM that will be georeferenced in ETRF2000 (Boucher and Altamimi 2011) and will cover the alpine area contained in the geographic rectangle with boundaries  $\lambda = 7.80^\circ$  East and  $\lambda = 10.70^\circ$  East in longitude,  $\phi = 45.10^\circ$  North e  $\phi = 46.70^\circ$  North in latitude. The output DTM will be gridded in geographical coordinates with a spatial resolution of  $2 \times 10^{-4}$  degrees, about 15 m in longitude and 22 m in latitude. The output grid will be composed of 8,000 rows and 14,500 columns, for a total of 116 M of nodes.

The four input DTMs that will be merged have been already cross-validated. In particular, shifts and biases between them have been investigated in the overlap borders: the results are satisfactory and have been extensively reported in Biagi et al. (2012), Biagi et al. (2013), Carcano (2014).

Input DTMs have been already transformed to ETRF2000 and are now stored as lists of geographic horizontal coordinates and terrain elevations, to be interpolated on the nodes of the final grid: in total, about  $350 \times 10^6$  elevations are available.

## 2 Re-gridding by Bicubic Surfaces

To grid observations, firstly a parametric models has to be chosen: the parameters are then estimated by the input set. Many approaches and models exist: a first classification can be in interpolation and approximation (Davis 1975; Christakos 1992). In interpolation, a model is estimated that passes through all the observations. In approximation, statistical methods are applied to estimate a smoother model from the observations.

When raw observations are used to estimate a digital model, approximation must be adopted because this allows the filtering of observations errors and outliers: one example is discussed in Biagi and Negretti (2004). On the contrary, in re-gridding, the input data are the nodes of a model that has been already filtered and checked against observation errors. In this case, actual details of the model can be lost in the approximation smoothing: therefore, interpolation should provide better results. Furthermore, often the re-gridding of huge amount of data is required: therefore an efficient numerical method should be applied.

Typically, GIS software (O’Sullivan and Unwin 2003) implements either splines or polynomial surfaces (Kidner 2003) to interpolate grid nodes: these standard approaches will be investigated for a re-gridding application.

Particularly, we will focus on the use of the well known and classical bilinear and bicubic surfaces. Previous works (see for example Rees 2000) state the better accuracy of bicubic surfaces. In any case preliminary tests have been performed also on our case study to compare bilinear and bicubic surfaces.

Several grids (case studies corresponding to different type of terrain) have been extracted from Lombardia DTM. Each one of them is re-gridded on nodes that are exactly in the middle of the original nodes. The output grid is back interpolated on the nodes of the original grid. Re-gridding is performed by adopting both bilinear and bicubic interpolations. The original and the final grids are compared (Table 1): synthetically, in flat areas no significant differences exist, while in rough terrain bicubic provides better results. Our project involves mainly mountainous areas: therefore bicubic polynomial surface will be adopted.

In the following, a small case study will be used for the tests. An output grid of  $400 \times 500$  ( $2 \times 10^5$ ) nodes has to be interpolated: the grid is centered on  $\lambda = 9.22^\circ$  East,

**Table 1** Internal checks

	Bias (m)	Std (m)	Max (m)	Min (m)
Bilinear	0.1	2	52	-56
Bicubic	0.0	1	36	-35

A gridded DTM is interpolated on new nodes in the middle of the original nodes. The output grid is back interpolated on the nodes of the original grid. Comparisons between bilinear and bicubic interpolations. Statistics of the differences between original and twice interpolated data

$\phi = 45.86^\circ$  North (Triangolo Lariano area, near Como Lake, Fig. 1) and has the afore mentioned resolution of  $2 \times 10^{-4}$  degrees. Terrain elevations have mean equal to 660 m, standard deviation of 345 m, minimum and maximum respectively of 197 and 1,679 m. Both Lombardia and SwissTopo DTMs are available in this test area: their differences have zero mean, standard deviation of 20 m and maximum of 100 m.

A bicubic surface is given by

$$H(x, y) = a_{00} + a_{10}x + a_{01}y + a_{20}x^2 + a_{11}xy + a_{02}y^2 + a_{30}x^3 + a_{21}x^2y + a_{12}xy^2 + a_{03}y^3 + a_{31}x^3y + a_{22}x^2y^2 + a_{13}xy^3 + a_{32}x^3y^2 + a_{23}x^2y^3 + a_{33}x^3y^3 \quad (1)$$

At least  $(x_i, y_i, H_i)$  sixteen observations are needed to estimate  $a_{ij}$  parameters: the system can be written as follows

$$\mathbf{z} = \begin{bmatrix} H_1 \\ H_2 \\ \dots \\ H_n \end{bmatrix} = \mathbf{A}\xi = \begin{bmatrix} 1 & x_1 & y_1 & \dots & x_1^3 & y_1^3 \\ 1 & x_2 & y_2 & \dots & x_2^3 & y_2^3 \\ \dots & \dots & \dots & \dots & \dots & \dots \\ 1 & x_n & y_n & \dots & x_n^3 & y_n^3 \end{bmatrix} \begin{bmatrix} a_{00} \\ a_{10} \\ a_{01} \\ \dots \\ a_{33} \end{bmatrix} \quad (2)$$

If 16 observations are used, the parameters are estimated by the solution of the iso-determined system

$$\xi = \mathbf{A}^{-1}\mathbf{z} \quad (3)$$

If more observations are used, a redundant system is built and can be solved by Least Squares (LS, Koch 1987) approach.

$$\xi = \mathbf{N}^{-1}\mathbf{A}^T\mathbf{Q}^{-1}\mathbf{z}, \quad \mathbf{N} = \mathbf{A}^T\mathbf{Q}^{-1}\mathbf{A} \quad (4)$$

The estimated parameters are used to estimate elevations in any inner point  $P = (x_p, y_p)$ . By applying (3) bicubic surface acts as interpolator, while (4) produces an approximation. Let consider the bicubic interpolation of a gridded DTM: to estimate the height in  $P$ , the 16 nodes of the 4 rows  $\times$  4 columns around it are used. Before estimation, a normalization of the coordinates is needed, like for example

$\theta_i = (q_i - q_p)/(q_{\max} - q_{\min})$ , where  $q_i$  stands either for  $x$  or  $y$ : this is needed for a decent conditioning of the system. Moreover, with the coordinates normalization

$$\widehat{H}_P(0, 0) = \widehat{a}_{00} \quad (5)$$

In our re-gridding case, the input points are almost regularly placed, but are not exactly gridded. To interpolate, the 16 nearest points of the input dataset are selected around each output node: considering their regular placement, system (3) is generally solved without any instability problem. However, in few cases, the system is ill conditioned and badly invertible. In these cases the polynomials may lead to extreme and unstable results. The problem can be easily checked by computing the condition number of  $\mathbf{A}$  (Press et al. 1992). According to the so called Singular Value Decomposition (SVD), the following holds

$$\mathbf{A} = \mathbf{U}\mathbf{\Lambda}\mathbf{V}^T \quad (6)$$

where  $\mathbf{U}$  and  $\mathbf{V}$  are orthonormal,  $\mathbf{\Lambda}$  is diagonal: its elements  $\lambda_i$  are called singular values and are real numbers. In case  $\mathbf{A}$  has full rank, all  $\lambda_i$  are positive. In case  $\mathbf{A}$  is singular or ill conditioned, some  $\lambda_i$  are zero or almost zero. The Condition Number (in the following  $CN$ ) of  $\mathbf{A}$  is defined as the ratio between the maximum and the minimum  $\lambda_i$ : in particular, on our case study,  $CN$  exceeds  $10^2$ ,  $10^3$  and  $10^5$  respectively in 82.5%, 0.025% and 0.01% (20) of the nodes.

When  $\mathbf{A}$  is ill conditioned, the interpolation can fail, leading to anomalous estimates. In particular  $CNs$  exceeding  $10^5$  provide very critical interpolations, that can easily identified by a simple visual inspection. No outliers appear for smaller  $CNs$ : however, interpolation problems could remain that are not evident.

## 2.1 Regularization of the System and Tests

To regularize the system, a threshold on  $CN$  can be set and different regularization approaches can be adopted: two of them have been implemented and will be compared.

More observations can be added to the nearest 16: a redundant system is built and is solved by LS. In this case, an iterative approach is implemented: at each step a new observation (the nearest that has not yet been included) is added to the observations vector,  $\mathbf{N}$  matrix is recomputed and its condition number is checked against a threshold. Iterations stop when the resulting system is well conditioned. This approach will be called Regularization by Redundant Observations (Reg-RO).

A different solution is based on the SVD of  $\mathbf{A}$ : in case the matrix is ill conditioned, some  $(\lambda_i/\lambda_{\max})^{-1}$  tend to infinite. To solve the system, the relevant  $(\lambda_i)^{-1}$  are set to zero: in

**Table 2** Internal checks on interpolation approaches

	Bias (m)	Std (m)	Max (m)	Min (m)
BI – Obs	0.3	6.4	92.6	−1,031
Reg-RO – Obs	0.3	5.3	39.3	−56.8
Reg: BI-RO-1	0.0	0.5	6.3	−3.7
Reg: BI-RO-2	148.9	581.3	2,714	−127.0
Reg: RO-AE	0.1	0.5	4.3	−2.5
(SwissTopo) EB-Obs	0.3	5.0	−25.2	37.2

Input data are used to produce output grid by three approaches: iso-determined bicubic interpolation (BI), regularized bicubic by Redundant Observations (Reg-RO), regularized bicubic by Annihilating Eigenvectors (Reg-AE). The output grid is back-interpolated on the input nodes: the original and the final elevations are compared. First five lines: Lombardy tests. BI-Obs: BI results, differences in the input points between original and back interpolated elevations. Reg-RO – Obs: Reg-RO results, same differences. Reg: BI-RO-1: differences between BI and Reg-RO in the (858) regularized nodes with  $10^{2.5} \leq CN < 10^5$ . Reg: BI-RO-2: differences between BI and Reg-RO in the (20) regularized nodes with  $10^5 \leq CN$ . Reg: RO-AE: differences between Reg-RO and RegAE interpolations in all the regularized nodes. SwissTopo tests: only BI-Obs is reported

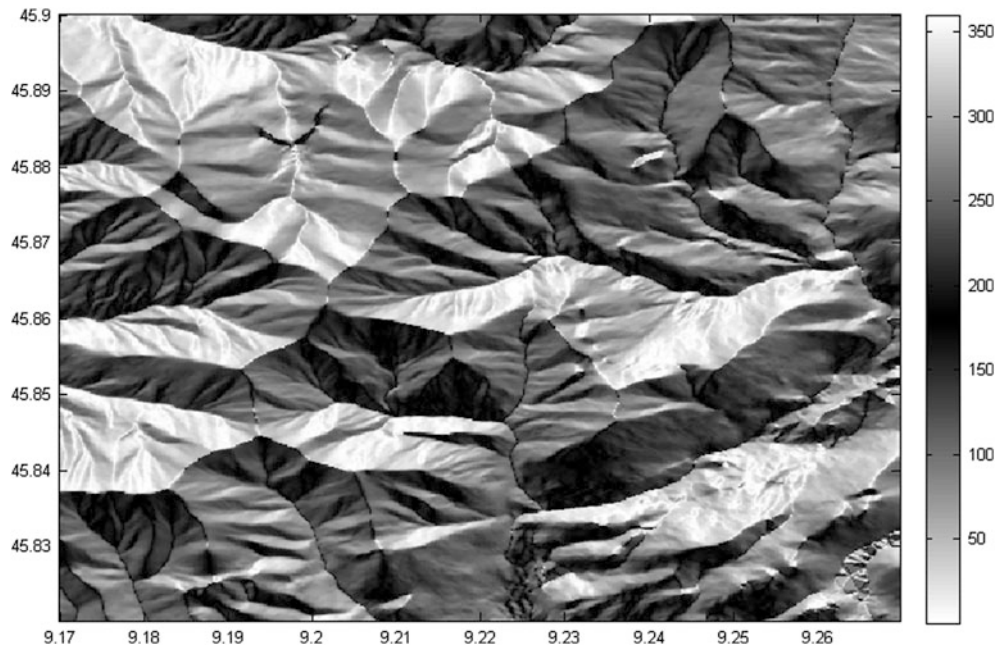
this way the inverse of  $\mathbf{A}$  is filled, and  $\mathbf{A}^{-1}$  is computed. This approach will be called Regularization by Annihilating Eigenvectors (Reg-AE).

Note that both Reg-RO and Reg-AE lead to approximation because the resulting surface no longer passes for the input observations. Both of them have been implemented and tested on the case study. A threshold  $CN = 10^{2.5}$  has been set on  $\mathbf{A}$ : the threshold is clearly over-conservative but allows to test the regularization effects in an acceptable number of nodes (878, 0.45% of the total). Reg-RO threshold ( $CN$  of  $\mathbf{N}$ ) is set to  $10^5$ . In Reg-RO, one, two and three redundant observations are added respectively in 170, 608, 100 cases, while in Reg-AE one, two and three eigenvectors are annihilated respectively in 856, 19, 3 cases. This was in some way expected, because annihilation is selective on the less significant eigenvalues, while redundant observations are added simply on a distance criterion.

To verify the effect of regularization, an internal check is possible. The output grid is back-interpolated on the input nodes: the original and the final elevations can be compared. Moreover, a morphometric analysis is possible (Li et al. 2005). Without regularization, roughness parameters increase for nodes with high  $CN$ : on the contrary, after the regularization, they are homogeneous between regularized and not regularized output nodes. Comparisons (Table 2) clearly show that the regularization improves critical interpolations, without over-smoothing the input dataset. Reg-RO and Reg-AE provide consistent results and only ten regularized interpolations deviate more than 5 m from the corresponding iso-determined interpolations.

In case of SwissTopo input data, the maximum  $CN$  is smaller than  $10^{2.5}$ : indeed, SwissTopo nodes never present a





**Fig. 1** Case study: merging of two datasets to produce a unified DTM. 3-D shaded model obtained by averaging the individual interpolations of Lombardy and SwissTopo dataset. Angles in degrees, clockwise from North direction

critical spatial configuration with respect to the output nodes. The statistics of interpolation are similar to those provided by the regularization on Lombardy dataset.

### 3 Comparisons Between Different Methods to Merge Overlapping DTMS

Input DTMs should be merged: in this section, Merg-IUD and Merg-AII are compared. To test Merg-IUD, one dataset has been formed and used by unifying the Lombardy and SwissTopo datasets. At first, the conditioning of the interpolation is analyzed.  $CN$  exceeds  $10^{2.5}$ ,  $10^5$ ,  $10^7$ ,  $10^9$  respectively in 99.5%, 83.25%, 24.75% and 0.62% (1,256) of the nodes. The problem is clearly due to the spatial pattern of the unified Lombardy and SwissTopo datasets. Moreover, local differences between the heights of the two DTMs exist and in this case must be filtered. Therefore, at least 32 observations are used to estimate each node by LS. About 12,000 (6%) nodes require a further regularization (threshold  $CN \geq 10^5$ ). The output grid contains significant artifacts and presents artificial roughness, that are clearly due to the local differences between the two input datasets, that are not efficiently smoothed. A smoother solution could be obtained by increasing the minimum number of input points: however, as an alternative, we would test the results provided by Merg-AII.

Two DTMs have been re-gridded (Reg-RO) from Lombardy and SwissTopo input datasets: their differences are spatially correlated and rather smooth: they have mean equal to 2.1 m and standard deviation equal to 20 m. 21 differences exceed 100 m and the maximum is 123 m. The two DTMs are averaged to produce a final grid (Fig. 1): the result does not present artifacts and its roughness is consistent with those of the individual interpolations. Therefore, in our case study, Merg-AII provide less instabilities and smoother results.

### 4 Conclusions

To apply a reference frame transformation on a DTM, two opposite approaches can be adopted. In the Direct Transformation, the input DTM is transformed and the transformed nodes are used to interpolate the output grid. In the Inverse transformation, the horizontal positions of the output nodes are back transformed on the input grid, and elevations are interpolated in their horizontal positions. In both the approaches, re-gridding is required. In this paper, iso-determined interpolation based on local polynomial surfaces has been investigated. Indeed, in our particular application, no smoothing is required because the input data are the nodes of a model and not raw observations.

To produce a unified DTM, neighboring and partly overlapping DTMs are transformed to a common grid and merged: to merge, two different approaches are possible. In Merg-IUD, the input datasets are unified and then are used to grid the output DTM while in Merg-AII, an individual interpolation of each dataset is followed by the average of the results.

The first part of the paper compares bilinear and bicubic interpolations: on our case study, the latter provides better results. Then, the problems of bicubic are discussed in the Direct Transformation. In particular, critical spatial horizontal distributions of the input nodes can cause an ill conditioning of the interpolation systems and, consequently, sparse outliers and irregularities. To regularize, redundant observations can be added and a least squares system is solved; alternatively, the not significant eigenvectors of the system can be identified and annihilated by SVD analysis. On a test case study, both the approaches provide satisfactory and consistent results and do not over-smooth the input dataset.

In the second part, Merg-IUD and Merg-AII are compared: the former approach provides unsatisfactory results and no further investigations seem to be needed on it. More analyses will be performed on Merg-AII: a particular focus will be given to its combination with the Inverse Transformation, that in a future work will be implemented and compared with the direct one.

**Acknowledgments** The paper has been supported by HELI-DEM project. Suggestions of F. Sansò, M. Crespi and the two anonymous reviewers have been really fruitful and have given space to further improvements.

## References

- Biagi L, Negretti M (2004) Environmental thematic maps prediction and easy probabilistic classification. *Trans GIS* 8(2):263–274, Plate 21
- Biagi L, Brovelli MA, Campi A, Cannata M, Carcano L, Credali M, De Agostino M, Manzano A, Sansò F, Siletto G (2011) Il progetto HELI-DEM (Helvetia-Italy digital elevation model): scopi e stato di attuazione. *Bollettino SIFET* 1/2011
- Biagi L, Carcano L, De Agostino M (2012) DTM cross validation and merging: problems and solutions for a case study within the HELI-DEM project. *Int Arch Photogramm Remote Sens Spatial Inf Sci* XXXIX-B4
- Biagi L, Carcano L, Lucchese A, Negretti M (2013) Creation of a multiresolution and multiaccuracy DTM: problems and solutions for a case study. *ISPRS Archives*, vol XL-5/W3, WG V/3. The role of geomatics in hydrogeological risk, 27–28 February 2013, Padua, pp 63–71
- Boucher C, Altamimi Z (2011) Memo: specifications for reference frame fixing in the analysis of a EUREF GPS campaign, Version 8. <http://etrs89.ensg.ign.fr/memo-V8.pdf>
- Buckley S, Mitchell HL (2004) Integration, validation and point spacing optimization of digital elevation models. *Photogramm Record* 19(108):277–295
- Carcano L (2014) Merging local DTMs: HELI-DEM project, problems and solutions, PhD Thesis, Politecnico di Milano, DICA, Geomatics Laboratory of Como Campus. <http://geomatica.como.polimi.it>
- Christakos G (1992) Random field models for earth sciences. Academic, New York
- Davis PJ (1975) Interpolation and approximation. Dover Publication Inc, New York
- Donatelli D, Maseroli R, Pierozzi M (2002) La trasformazione tra sistemi di riferimento utilizzati in Italia. *Boll Geod Sci Affini Anno LXI(4):247–310*
- El-Sheimy N, Valeo C, Habib A (2005) Digital terrain modeling – acquisition, manipulation, and applications. Artech House, Boston
- Kidner D (2003) Higher-order interpolation of regular grid digital elevation models. *Int J Remote Sens* 24(14):2981–2987
- Koch KR (1987) Parameter estimation and hypothesis testing in linear models. Springer, Berlin
- Li Z, Zhu Q, Gold C (2005) Digital terrain modeling – principles and methodology. CRC, Boca Raton
- O’Sullivan D, Unwin D (2003) Geographic information analysis. Wiley, Hoboken
- Press WH, Teukolsky SA, Vetterling WT, Brian PF (1992) Numerical recipes in C: the art of scientific computing, 2nd edn. Cambridge University Press, Cambridge
- Rees WG (2000) The accuracy of digital elevation models interpolated to higher resolutions. *Int J Remote Sens* 21(1):7–20

---

# Single-Epoch GNSS Array Integrity: An Analytical Study

A. Khodabandeh and P.J.G. Teunissen

---

## Abstract

In this contribution we analyze the integrity of the GNSS array model through the so-called uniformly most powerful invariant (UMPI) test-statistics and their corresponding minimal detectable biases (MDBs). The model considered is characterized by multiple receivers/satellites with known coordinates where the multi-frequency carrier-phase and pseudo-range observables are subject to atmospheric (ionospheric and tropospheric) delays, receiver and satellite clock biases, as well as instrumental delays. Highlighting the role played by the model's misclosures, analytical multivariate expressions of a few leading test-statistics together with their MDBs are studied that are further accompanied by numerical results of the three GNSSs GPS, Galileo and BeiDou.

---

## Keywords

Array model • GNSS misclosures • Integrity • Minimal detectable bias (MDB) • Uniformly most powerful invariant (UMPI) test-statistic

---

## 1 Introduction

The notion of the GNSS array model, here, refers to an array of antennas tracking the multi-frequency carrier-phase/pseudo-range observables in the presence of atmospheric effects. The coordinates of antennas and satellites are assumed to be known. The definition presented is rather general in the sense that even the medium-scale control networks can also be considered as an array. Examples of such are the continuously operating reference

station (CORS) networks sending corrections to the RTK-and/or PPP-RTK users (de Jonge 1998; Odijk et al. 2014), a set of antennas mounted on rigid platforms improving the position/attitude of points in its vicinity (Teunissen 2010, 2012), and the ground based augmentation systems (GBASs) supporting safe flight procedures such as landing, departure and surface operations at an airport (Khanafseh et al. 2012; Giorgi et al. 2012). Despite their different applications, all of the aforementioned arrays are, however, utilized for the purpose of the *same* functionality, that is, providing accurate corrections for the users. Ensuring the integrity and reliability of the corrections, even at the pre-analysis level, is therefore of great importance, see e.g., Teunissen (1998); Teunissen and de Bakker (2012).

Integrity monitoring and quality control of the GNSS array model is the topic of this contribution. We confine our study to the *single-epoch* scenario as it is indeed the ultimate goal of the near real-time applications and, at the same time, brings us conservative thresholds of the reliability measures of the corresponding multi-epoch scenario. Our strategy commences with the model's misclosures. Although the GNSS misclosures can be treated as diagnostic tools in

---

A. Khodabandeh (✉)  
Department of Spatial Sciences, GNSS Research Centre, Curtin  
University of Technology, Perth, WA 6845, Australia  
e-mail: [amir.khodabandeh@curtin.edu.au](mailto:amir.khodabandeh@curtin.edu.au)

P.J.G. Teunissen  
Department of Spatial Sciences, GNSS Research Centre, Curtin  
University of Technology, Perth, WA 6845, Australia

Department of Geoscience and Remote Sensing, Delft University  
of Technology, Delft, The Netherlands  
e-mail: [p.teunissen@curtin.edu.au](mailto:p.teunissen@curtin.edu.au)

their own right, we make use of certain linear functions of them to formulate the UMPI test-statistics which give rise to the highest probability of the detection for a class of critical regions. For an overview of the underlying principles of the UMPI test, see Arnold (1981), and for its applications to hypothesis testing in linear models, see e.g. Teunissen (2000).

The test-statistics to be studied are (1) the array-, antenna- and satellite-detectors in which the overall/local validity of the model is tested, (2) the celebrated  $w$ -test-statistic for the purpose of outlier identification and (3) the atmospheric detectors well suited to the small-scale arrays. The detectability of the tests is formulated via the corresponding MDBs where the associated numerical illustrations, emphasized on the three GNSSs GPS, Galileo and BeiDou, are also given.

## 2 Array Model and the GNSS Misclosures

Consider a single antenna, say antenna  $r$  ( $r = 1, \dots, n$ ), that tracks  $s$  number of commonly-viewed satellites on frequency  $j$  ( $j = 1, \dots, f$ ). One can then put the corresponding *undifferenced* carrier-phase observations on each frequency, as the  $s$ -vectors  $\phi_{r,j}$  ( $j = 1, \dots, f$ ), into a higher-dimensioned vector  $\phi_r = [\phi_{r,1}^T, \dots, \phi_{r,f}^T]^T$ . Doing the same to the pseudo-range observations  $p_r$  and collecting observations of all  $n$  antennas, the final  $sf \times n$  matrices of carrier-phase and pseudo-range data of the array can be, respectively, formulated as

$$\Phi = [\phi_1, \dots, \phi_n], \quad P = [p_1, \dots, p_n]$$

The satellite/receiver-dependent biases are, respectively, canceled out by applying the between-receiver single-differenced (SD) operator  $D_n$  and the between-satellite SD operator  $D_s$  (Teunissen 1997). The multivariate representation of the double-differenced (DD) observation equations of the array model, under the null hypothesis  $H_0$ , reads then

$$\begin{aligned} \mathbf{E}\{(I_f \otimes D_s^T) \Phi D_n\} &= (e_f \otimes D_s^T g) \tau^T D_n - (\mu \otimes I_{s-1}) \\ &\quad \times D_s^T \iota D_n + (\Lambda \otimes I_{s-1}) Z \\ \mathbf{E}\{(I_f \otimes D_s^T) P D_n\} &= (e_f \otimes D_s^T g) \tau^T D_n + (\mu \otimes I_{s-1}) \\ &\quad \times D_s^T \iota D_n \end{aligned} \quad (1)$$

$$\begin{aligned} \mathbf{D}\{\text{vec}[(I_f \otimes D_s^T) \Phi D_n]\} &= D_n^T D_n \otimes Q_\phi \otimes D_s^T W_s^{-1} D_s \\ \mathbf{D}\{\text{vec}[(I_f \otimes D_s^T) P D_n]\} &= D_n^T D_n \otimes Q_p \otimes D_s^T W_s^{-1} D_s \end{aligned} \quad (2)$$

where the  $s$ -vector  $g$  contains functions mapping the slant tropospheric delays (STDs) onto the zenith tropospheric delays (ZTDs)  $\tau = [\tau_1, \dots, \tau_n]^T$ . The  $s \times n$  matrix  $\iota$  is introduced as  $\iota = [\iota_1, \dots, \iota_n]$ , with  $\iota_r$  being the  $s$ -vector of the (first-order) slant ionospheric delays of antenna  $r$ . The  $f$ -vector  $\mu$  contains the ionospheric coefficients  $\mu_j = \lambda_j^2 / \lambda_1^2$ , with  $\lambda_j$  being the wavelengths positioned on the  $f \times f$  diagonal matrix  $\Lambda$ . The matrix  $Z$  contains the integer-valued DD ambiguities. The  $f \times f$  positive-definite matrices  $Q_\phi$  and  $Q_p$  are the cofactor matrices of the phase and pseudo-range observable-type. The  $s \times s$  diagonal matrix  $W_s$  captures the satellite elevation dependency of the observations.  $I$  and  $e$ , respectively, denote the identity matrix and the vector of ones, where the subscripts indicate their size. The operator  $\otimes$  denotes the Kronecker product.  $\mathbf{E}\{\cdot\}$  and  $\mathbf{D}\{\cdot\}$  are the mathematical expectation and dispersion operators, respectively. The operator  $\text{vec}[\cdot]$  vectorizes the associated matrix.

Using model (1) and (2), we are interested to check the validity of the model against unaccounted effects. To do so, we therefore work with the conditioned equations of (1) and the corresponding misclosures. The idea to employ the conditioned equations rather than the commonly-used observation equations is motivated by the desire to characterize the intrinsic behavior of the array model in relation to the possible misspecifications. This is indeed realized by forming the GNSS misclosures showing the contribution of the observations to the redundancy of the model.

### 2.1 GNSS-Based Decoupled Misclosures

Although the misclosures of (1) can be formed in many different ways, we form those that are group-wise *uncorrelated* and at the same time have easy interpretations. The GNSS-based decoupled misclosures, in case of the ambiguity-float scenario, are introduced as follows (cf. Appendix)

*i: Frequency-differenced misclosures:*

$$M_1 = [(D_f^T \mu)^\perp D_f^T \otimes c_{d|\tau}^2 \bar{g}^T W_s] P D_n \quad (3)$$

*ii: Atmosphere-free misclosures:*

$$M_2 = [\mu^\perp \otimes (D_s^T g)^\perp D_s^T] P D_n$$

where  $(\cdot)^\perp$  denotes the orthogonal complement basis matrix. We introduce the  $s$ -vector  $\bar{g} = g + (c_{d\tau}/c_\tau^2)e_s$ , in which the satellite-domain (co)variance-type scalars  $c_{d|\tau}^2$ ,  $c_{d\tau}$  and  $c_\tau^2$  are computed as

$$\begin{aligned} c_\tau^2 &= \frac{e_s^T W_s e_s}{[e_s^T W_s e_s][g^T W_s g] - [g^T W_s e_s]^2} \\ c_{d\tau} &= \frac{-g^T W_s e_s}{[e_s^T W_s e_s][g^T W_s g] - [g^T W_s e_s]^2}, \quad c_{d|\tau}^2 = \frac{1}{e_s^T W_s e_s} \end{aligned} \quad (4)$$

With regard to (3),  $M_1$  and  $M_2$ , respectively, contribute to the model's redundancy of size  $(f - 2)$  and  $(f - 1)(s - 2)$  per baseline. After fixing ambiguities, similar expressions can be obtained for the phase and phase-and-code misclosures.

## 2.2 Atmosphere-Aided Decoupled Misclosures

The GNSS-based misclosures, presented in (3), contain the *complete* information needed to check and to study the quality of the observation matrices  $\Phi$  and  $P$  in (1). One may, however, strengthen the model by using a-priori atmospheric information, i.e. the spatial dependency of the atmospheric delays. In case of not-too-large arrays, the differential atmospheric delays, with amount of uncertainty, would thus play the role of pseudo-observables as

$$\begin{aligned} E\{D_s^T \iota D_n\} &= D_s^T \iota D_n, & \text{with } D\{D_s^T \iota D_n\} &= \sigma_\iota^2 D_n^T D_n \otimes D_s^T W_s^{-1} D_s \\ E\{D_n^T \tau\} &= D_n^T \tau, & \text{with } D\{D_n^T \tau\} &= \sigma_\tau^2 D_n^T D_n \end{aligned} \quad (5)$$

with  $\sigma_\iota^2$  and  $\sigma_\tau^2$  being the a-priori ionospheric and tropospheric variances, respectively.

Appending the preceding equations to (1) does increase the redundancy of the model by comparing the GNSS-based estimators of the differential atmospheric delays with their pseudo-observable ones ( $s$  redundant observations per baseline). This, in a similar way to (3), provides us with the atmosphere-aided misclosures.

## 3 UMPI Test-Statistics and Their MDBs

Given the GNSS decoupled misclosures introduced in the previous section, we are now in a position to form various test-statistics.

**Theorem 1 (UMPI Test-Statistic and Its MDB)** *Let the alternative hypothesis  $H_a$  be related to the null hypothesis  $H_o$  as  $E\{\text{vec}[Y]|H_a\} = E\{\text{vec}[Y]|H_o\} + C_Y \nabla$ , where the  $q$ -vector of misspecifications  $\nabla$  is linked to the observations by the full-rank design matrix  $C_Y$ . Given a representation for the model's misclosures as  $M = B^T \text{vec}[Y]$  under  $H_o$ , the UMPI test-statistic  $T_q$  and its MDB are respectively given by*

$$T_q = \frac{\text{tr}\{Q_{MM}^{-1} P_{C_M} M M^T\}}{\text{tr}\{P_{C_M}\}} \quad (6)$$

$$\|\nabla\| = \sqrt{\frac{v_{q,\alpha,\gamma}}{d_Y^T C_M^T Q_{MM}^{-1} C_M d_Y}}, \quad \nabla = \|\nabla\| d_Y \quad (7)$$

where  $Q_{MM} = D\{M\}$  and  $C_M = B^T C_Y$ , with the projector  $P_{C_M} = C_M (C_M^T Q_{MM}^{-1} C_M)^{-1} C_M^T Q_{MM}^{-1}$ . The scalar  $v_{q,\alpha,\gamma}$  is the  $\chi^2$ -noncentrality parameter to be determined by the power of the test  $\gamma$  and the probability of false alarm  $\alpha$ . The operator  $\text{tr}\{\cdot\}$  denotes the trace of a matrix, whereas  $\|\cdot\|^2 = (\cdot)^T (\cdot)$  is the squared-norm of a vector.

*Proof* see Appendix.  $\square$

The above theorem shows how the multivariate representation of the UMPI test-statistic is realized through the model's misclosures and the type of misspecifications, i.e.  $C_Y$ . We remark that the test-statistic  $T_q$  follows central and noncentral  $F$ -distribution under  $H_o$  and  $H_a$ , respectively, that is,  $T_q|H_o \sim F(q, \infty, 0)$  and  $T_q|H_a \sim F(q, \infty, v)$ , with the first two arguments  $q, \infty$  being the degrees of freedom and  $v$  the noncentrality parameter.

As to the GNSS array model, one may formulate a rather general structure for the misspecification design matrix  $C_Y$  at the *undifferenced* level. The following structure has been adopted in this study

$$E\{\tilde{Y}|H_a\} = E\{\tilde{Y}|H_o\} + (C_f \otimes C_s) \nabla C_n^T \quad (8)$$

where the full-rank matrices  $C_f, C_s$  and  $C_n$  specify the type of misspecification  $\nabla$  in the frequency-, satellite- and antenna-domain, respectively. The role of the *atmosphere-corrected* observation matrix  $\tilde{Y}$  can be taken by  $\tilde{P}$  and  $\tilde{\Phi}$  or both of them, in which we define

$$\begin{aligned} \tilde{P} &= P - (e_f \otimes g)\tau - (\mu \otimes I_s)\iota \\ \tilde{\Phi} &= \Phi - (e_f \otimes g)\tau + (\mu \otimes I_s)\iota \end{aligned} \quad (9)$$

**MDB-Parametrization** In general there is no unique solution for the MDB of the misspecifications of a *multi dimensional* type. This issue can be properly circumvented through an MDB-parametrization as follows

$$\nabla = \|\text{vec}[\nabla]\| (d_f \otimes d_s) d_n^T \quad (10)$$

with  $d_f, d_s$  and  $d_n$  being, respectively, the frequency-, satellite- and antenna-domain vectors such that their resultant vector  $d_n \otimes d_f \otimes d_s$  is of length 1, i.e. a direction vector.

In the following, a few important test-statistics, together with their MDBs, will be specialized by setting  $C_f, C_s$  and  $C_n$  in (8) to certain structures.

### 3.1 Array-, Antenna- and Satellite-Detectors

First one needs to check the validity of the array model against any type of misspecification that might potentially occur. The stated validity can be either of an *overall* type

or of a *local* type. The overall validity of the model is tested through the array-detector characterized by the following setting

$$\text{array-detector: } C_f \mapsto I_f, \quad C_s \mapsto D_s, \quad C_n \mapsto D_n \quad (11)$$

Depending on the a-priori atmospheric variances  $\sigma_t^2$  and  $\sigma_r^2$ , several expressions can be formulated. In case of atmosphere-fixed scenario, i.e.  $\sigma_t^2 = 0$  and  $\sigma_r^2 = 0$ , the array-detector can be shown to take the following form (cf. Appendix)

$$T_q = \frac{1}{q} \text{tr}\{[Q_p^{-1} \otimes W_s P_{e_s}^\perp] \tilde{P} P_{D_n} \tilde{P}^T\} + \frac{1}{q} \text{tr}\{[Q_\phi^{-1} \otimes W_s P_{e_s}^\perp] \tilde{\Phi} P_{D_n} \tilde{\Phi}^T\} \quad (12)$$

with the projectors  $P_{e_s}^\perp = I_s - c_{d|t}^2 e_s e_s^T W_s$  and  $P_{D_n} = I_n - (1/n) e_n e_n^T$ . The degrees of freedom  $q$  is determined upon choosing the following scenarios

$$\begin{aligned} \text{codeless data } (Q_p^{-1} = 0) &: \Rightarrow q = f(s-1)(n-1) \\ \text{phaseless data } (Q_\phi^{-1} = 0) &: \Rightarrow q = f(s-1)(n-1) \\ \text{code+phase data} &: \Rightarrow q = 2f(s-1)(n-1) \end{aligned} \quad (13)$$

The corresponding MDB, in accordance with (7), reads

$$\begin{aligned} & \|\text{vec}[\nabla]\| \\ &= \frac{v_{q,\alpha,\gamma}^{\frac{1}{2}}}{\sqrt{[d_f^T (Q_p^{-1} + Q_\phi^{-1}) d_f][d_s^T D_s^T W_s P_{e_s}^\perp D_s d_s][d_n^T D_n^T P_{D_n} D_n d_n]}} \end{aligned} \quad (14)$$

Clearly a judgment on the size of the MDB cannot be easily made since it depends on the three vectors  $d_f$ ,  $d_s$  and  $d_n$ . Keeping fixed two vectors out of which however, one can still gain information on the sensitivity of the MDB to the contributing factors like the number of frequencies/satellites and the quality of the observables. This idea leads to locally validate the model by testing observations of a particular antenna and/or those of a particular satellite. We can therefore characterize the antenna-/satellite-detector upon the following setting

$$\begin{aligned} \text{antenna-detector: } & C_f \mapsto I_f, \quad C_s \mapsto D_s, \quad C_n \mapsto u_r^n \\ \text{satellite-detector: } & C_f \mapsto I_f, \quad C_s \mapsto u_i^s, \quad C_n \mapsto D_n \end{aligned} \quad (15)$$

where  $u_r^n$  denotes the canonical  $n$ -vector containing zeros except the  $r^{\text{th}}$  element equal to one. The canonical vector

$u_i^s$  is defined similarly. With this setting, in an analogous way to (11) and (12), expressions of the stated test-statistics as well as their MDBs can be obtained. For the atmosphere-fixed case, the degrees of freedoms of the antenna-/satellite-detector are  $q = f(s-1)$  and  $q = f(n-1)$ , respectively.

### 3.2 $w$ -Test-Statistic and the MDB of Single Outliers

We now focus our attention to the well-known  $w$ -test-statistic employed for the purpose of identification of a single erroneous observation (Baarda 1968). The structure of  $C_Y$  is then set to

$$w\text{-test-statistic: } C_f \mapsto u_j^f, \quad C_s \mapsto u_i^s, \quad C_n \mapsto u_r^n \quad (16)$$

Similar to the array-detector, depending on the scenarios considered, several expressions can be given to the  $w$ -test-statistic. The structures of the corresponding MDB do however follow the same pattern. Let us, for the moment, consider STDs rather than ZTDs in the model. The code-outlier MDB can be shown to read as

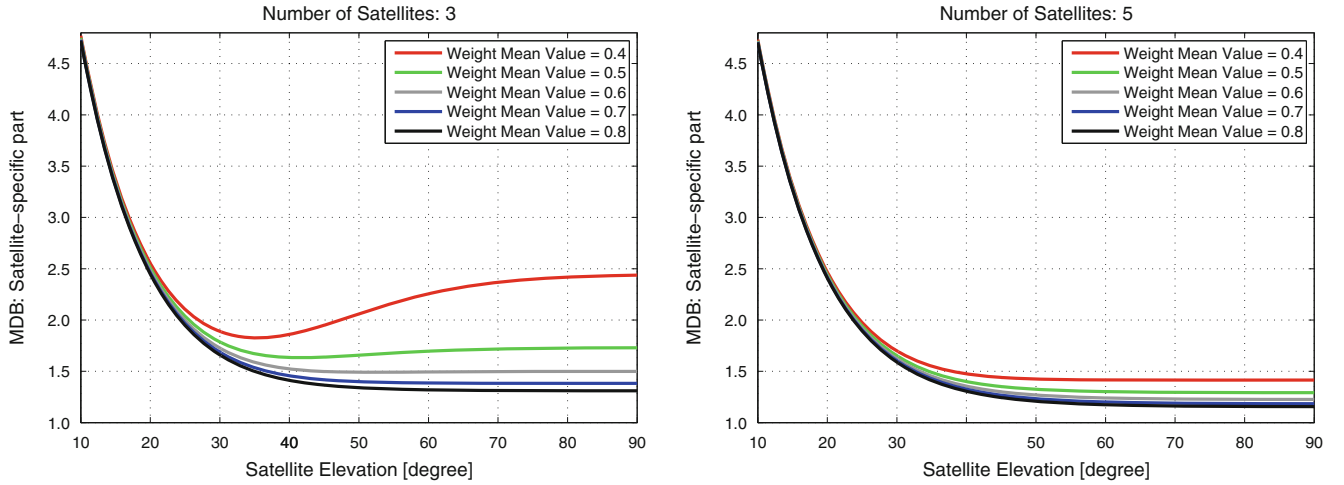
$$\begin{aligned} \|\text{vec}[\nabla]\|_{STD} &= v_{1,\alpha,\gamma}^{\frac{1}{2}} \times \left[\frac{n}{n-1}\right]^{\frac{1}{2}} \times [w^i (1 - [\frac{w^i}{\bar{w}}]_s^{\frac{1}{s}})]^{-\frac{1}{2}} \\ &\quad \times \left[\frac{\sigma_{p_j}^2}{1 - (\sigma_{p_j}^2 / \sigma_{p_j}^2)}\right]^{\frac{1}{2}} \end{aligned} \quad (17)$$

Four contributing elements show up themselves in the above MDB that are described in the following:

**Noncentrality Parameter  $v_{1,\alpha,\gamma}$**  Given the fixed probability of false alarm  $\alpha$ , the  $\chi^2$ -noncentrality parameter increases as the power of the test  $\gamma$  increases. In other words, for a given fixed model and a fixed  $\alpha$ , the higher the power of the test is sought, the larger the MDB becomes.

**Antenna-Specific Part  $[n/(n-1)]^{\frac{1}{2}}$**  This clearly shows that an increase in the number of antennas  $n$  could decrease the size of the outlier MDB considerably, would one, in the beginning, consider a *limited* number of antennas (e.g.  $n = 2$  or  $n = 3$ ). However, the stated MDB *does not* significantly decrease in size by adding an extra antenna when dealing with an array of a *large* number of antennas.

**Satellite-Specific Part  $[w^i (1 - [\frac{w^i}{\bar{w}}]_s^{\frac{1}{s}})]^{-\frac{1}{2}}$**  This term depends on three factors, namely, 1) the elevation-dependent weight  $w^i$  of an individual satellite, say  $i$ , 2) the mean value of  $w^i$  ( $i = 1, \dots, s$ ) denoted by  $\bar{w}$ , and 3) the number of common satellites  $s$ . In this study, we make use of the exponential elevation weighting strategy to form the diagonal



**Fig. 1** Satellite-specific part of the outlier MDBs as function of the satellite elevation for different satellite configuration. The overall satellite configuration has been characterized by the weight mean value ‘ $\bar{w}$ ’

elements of matrix  $W_s$ , i.e.  $w^i$  (Euler and Goad 1991)

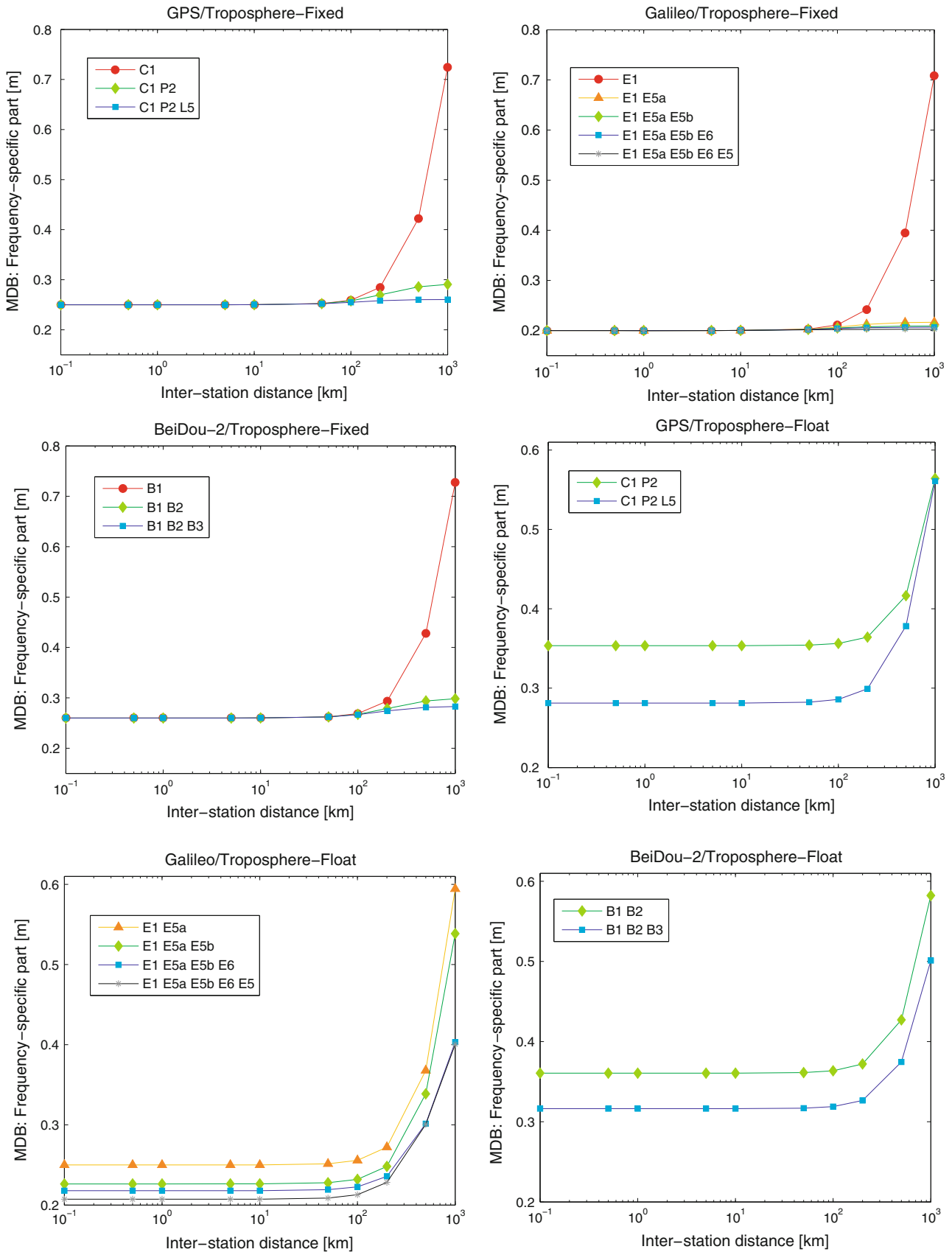
$$w^i = [1 + 10 \exp(-\frac{\epsilon^i}{10^\circ})]^{-2}, \quad i = 1, \dots, s \quad (18)$$

where  $\epsilon^i$  is the elevation of satellite  $i$  [degree] with respect to the reference antenna. Note that the elevation-dependent weight  $w^i$  should not be confused with the  $w$ -test-statistic.

Figure 1 depicts the satellite-specific part as function of the elevation of an individual satellite. The graphs have been presented for different values of  $\bar{w}$  reflecting the overall configuration of the satellites with respect to the array ( $0.4 \leq \bar{w} \leq 0.8$ ). This has been done for two cases, the case where the number of satellites is  $s = 3$  (left-panel) and the other one with  $s = 5$  (right-panel). As shown, the size of the MDB of an outlier, occurred in a single observation of satellites of *low elevation* (e.g.  $10^\circ \leq \epsilon^i \leq 20^\circ$ ), is governed by the elevation of the corresponding satellite only, irrespective of the number/configuration of the satellites. In case of satellites of a higher elevation, the scenario would change as the number of satellites starts taking an active role as well. Considering a limited number of common satellites, it is interestingly observed that the MDB does not *generally* decrease as the elevation of the corresponding satellite increases (see the red thick line in Fig. 1, left-panel). In this case, in addition to the satellite elevation, the overall satellite configuration would also contribute to the size of the MDB. The stated contribution does however get insignificant once the number of satellite increases (see Fig. 1, right-panel). In the situations where the number of satellites is large enough (e.g. more than 5), one can therefore simply consider the elevation of each satellite *individually* to analyze the corresponding outlier MDB.

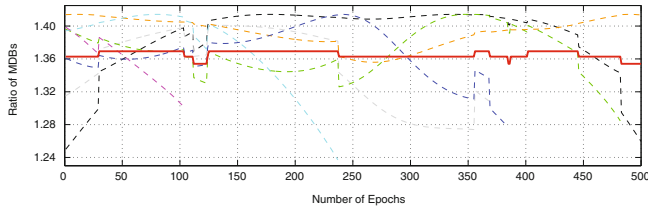
**Frequency-Specific Part**  $[\sigma_{p_j}^2 / (1 - (\sigma_{\hat{p}_j} / \sigma_{p_j})^2)]^{\frac{1}{2}}$  In addition to the variance of an individual pseudo-range observable-type  $\sigma_{p_j}^2$ , this term is also dependent on the variance of the *adjusted* observable-type denoted by  $\sigma_{\hat{p}_j}^2$ . This quantity in turn is a function of the a-priori atmospheric variances, the quality of the other pseudo-range observable-types through  $Q_p$  and the ionospheric vector  $\mu$ .

Figure 2 shows the frequency-specific part as function of the inter-station distance for the three GNSSs GPS, Galileo and BeiDou. In order to link the inter-station distance to the ionospheric variance  $\sigma_\tau^2$ , use has been made of that given in Schaffrin and Bock (1988). The graphs have been plotted for the troposphere-fixed case  $\sigma_\tau^2 = 0$  (top-panel) as well as the troposphere-float case  $\sigma_\tau^2 \rightarrow \infty$  (bottom-panel). As shown, the frequency-specific part behaves almost unchanged up to a certain inter-station distance (in this study around 100 [km]). In contrast to the single-frequency data (red dots), the MDB associated with the multi-frequency data does not significantly change as the inter-station distance increases (troposphere-fixed case). Because of the dispersive nature of the ionospheric effects (i.e. dependency on the frequencies), the *GNSS-based* misclosures would, in addition to the atmosphere-aided ones, also contribute to the  $w$ -test-statistic, whereas they do vanish in case of single-frequency data (cf. (3)). It is also important to note that there is no redundancy in the single-frequency troposphere-float scenario, thus giving rise to infinite MDBs. The single-frequency case is therefore excluded from the graphs of the bottom-panel. We remark, due to a generally better precision of the Galileo’s signals, that the associated results illustrate a superior performance to those of GPS and BeiDou. The pseudo-range zenith-referenced standard deviation is taken as 25 [cm] for GPS/BeiDou, and as 20 [cm] for Galileo .



**Fig. 2** Frequency-specific part of the code-outlier MDBs [m] as function of the inter-station distance [km] for three GNSSs GPS, Galileo, and BeiDou





**Fig. 3** Reduction-factors (ambiguity-float) of the code-outlier MDB (dashed lines) due to mapping the slant tropospheric delays (STDs) to their zenith counterparts (ZTDs) compared to the *rule-of-thumb* formula (red thick line) over time (a GPS data-set). Different colors have been used for different satellites

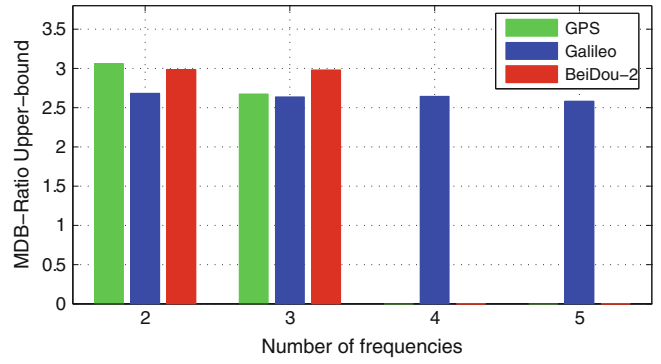
Similar to the satellite-specific part, as the model gets stronger (i.e.  $\sigma_{p_j}^2 \approx 0$ ), one can only consider the quality of that individual observable-type on frequency  $j$  (i.e.  $\sigma_{p_j}^2$ ). As an example, for an array of 4 antennas with the probability of false alarm  $\alpha = 0.01$ , the code-outlier MDB is about 79 [cm] ( $\gamma = 0.8$ ) and 89 [cm] ( $\gamma = 0.9$ ). In case of phase-slip MDB of two successive epochs, the MDB is about 3.9 [mm] ( $\gamma = 0.8$ ) and 4.4 [mm] ( $\gamma = 0.9$ ). The zenith-referenced standard deviations of the pseudo-range and carrier-phase observables are, respectively, set to  $\sigma_{p_j} = 20$  [cm] and  $\sigma_{\phi_j} = 1$  [mm].

### 3.2.1 MDB Reduction-Factor: From the STD-Based Model to the ZTD-Based Model

As stated so far, the MDB given in (17) refers to the STD-based model. One can now ask to what extent the MDB decreases by mapping the STDs to their ZTDs. Following the same procedure as before, the MDB of the ZTD-based model can be formulated that reveals the gain in terms of the reduction of the MDB. Although in addition to the number of frequencies/satellites, the stated reduction-factor does also depend on the tropospheric mapping functions  $g$  and the elevation-dependent weight matrix  $W_s$  (see Fig. 3), one can however present a *rule-of-thumb* expression as its rough value, namely (ionosphere-fixed scenario)

$$\begin{aligned} & \text{Before ambiguity-fixing} \\ & \frac{\|\text{vec}[\nabla]\|_{STD}}{\|\text{vec}[\nabla]\|_{ZTD}} \approx \left[1 + \frac{1}{f-1} \left(\frac{s-2}{s-1}\right)\right]^{\frac{1}{2}} \\ & \text{After ambiguity-fixing} \\ & \frac{\|\text{vec}[\nabla]\|_{STD}}{\|\text{vec}[\nabla]\|_{ZTD}} \approx 1 \end{aligned} \quad (19)$$

According to the above equations, before fixing ambiguities the reduction-factor is mostly governed by the number of frequencies  $f$  but not too much by the number satellites  $s$ . As the number frequency increases, the gain in terms of MDB reduction gets less. After fixing ambiguities, the reduction-factor becomes almost 1 meaning that the code-outlier MDB remains almost unchanged by even strengthening the model through mapping the tropospheric delays to their ZTDs.



**Fig. 4** Upper bounds of the reduction-factor of the ionospheric MDB due to excluding zenith tropospheric delays (ZTDs) from the underlying model for three GNSSs GPS (green bars), Galileo (blue bars), and BeiDou (red bars)

### 3.3 Atmospheric Detectors and Their MDBs

In most applications dealing with the small-scale arrays, one needs to check as to whether there are significant dispersive/nondispersive effects or not. Taking the atmosphere-fixed scenario as the null hypothesis, the atmospheric detectors are defined as

$$\begin{aligned} & \text{tropospheric-detector: } C_f \mapsto e_f, C_s \mapsto g, C_n \mapsto D_n \\ & \text{ionospheric-detector: } C_f \mapsto \mu, C_s \mapsto D_s, C_n \mapsto D_n \end{aligned} \quad (20)$$

Despite the complexity of the atmospheric MDBs, we can evaluate them in a *relative* sense. For instance, one can analyze the reduction of the ionospheric MDB when the differential ZTDs are assumed to be a-priori known via the following bounds (codeless data)

$$1 \leq \frac{\|\text{vec}[\nabla]\|_{\tau}}{\|\text{vec}[\nabla]\|} \leq \left[1 + \frac{\bar{\mu}^2}{\sigma_{\mu}^2}\right]^{\frac{1}{2}} \quad (21)$$

with  $\bar{\mu} = (1/f) \sum_{j=1}^f \mu_j$  and  $\sigma_{\mu}^2 = (1/f) \sum_{j=1}^f (\mu_j - \bar{\mu})^2$ . The ionospheric MDBs with and without ZTDs are denoted by  $\|\text{vec}[\nabla]\|_{\tau}$  and  $\|\text{vec}[\nabla]\|$ , respectively.

According to (21), the detectability of the differential ionosphere can get better at most  $\left[1 + (\bar{\mu}/\sigma_{\mu})^2\right]^{\frac{1}{2}}$  times, if one excludes the differential ZTDs from the model. For the current systems, the stated value is around 3 (cf. Fig. 4).

## 4 Concluding Remarks

In this contribution, the UMPI test-statistics as well as their MDBs, associated with the array model, were studied. With the aid of the GNSS decoupled misclosures, a few impor-

tant examples such as the array-detector,  $w$ -test-statistic and the ionospheric-detector were discussed. In particular, we showed that as the model gets stronger, one can simply, in case of outlier's MDB, analyze the single-channel/frequency scenario instead.

**Acknowledgements** P.J.G. Teunissen is the recipient of an Australian Research Council Federation Fellowship (project number FF0883188).

## Appendix

*Proof of (3)* The model's misclosures, forming the condition equations, can be formulated through pre-multiplying the corresponding observation vector by an orthogonal complement basis matrix of the design matrix (Teunissen 2000). In case of the single-epoch ambiguity-float scenario, the carrier-phase observations are all reserved to determine the DD ambiguities, thus leaving the code observations to contribute to the redundancy of the model. Given the observations Eq. (1), the code-only design matrix  $A$ , together with its orthogonal complement basis matrix  $B$ , can therefore be expressed as (per baseline)

$$A \mapsto [e_f \otimes D_s^T g, \mu \otimes I_{s-1}] \Rightarrow$$

$$B^T \mapsto \begin{bmatrix} (D_f^T \mu)^{\perp T} D_f^T \otimes c_{d\tau}^2 g^T D_s (D_s^T W_s^{-1} D_s)^{-1} \\ \mu^{\perp T} \otimes (D_s^T g)^{\perp T} \end{bmatrix} \quad (22)$$

from which (3) follows. That the misclosures  $M_1$  and  $M_2$  are mutually uncorrelated follows from the identities  $D_s^T \bar{g} = D_s^T g$ , and  $(D_s^T g)^{\perp T} D_s^T g = 0$ .  $\square$

*Proof of Theorem 1* Equation (6) is indeed another expression of the UMPI test-statistic  $T_q$  presented in Teunissen (2000). In terms of the model's misclosures  $M$ ,  $T_q$  and its MDB-squared  $\|\nabla\|^2$  read

$$T_q = \frac{1}{q} M^T Q_{MM}^{-1} P_{C_M} M \quad (23)$$

$$\|\nabla\|^2 = \frac{v_{q,\alpha,\gamma}}{d_Y^T C_M^T Q_{MM}^{-1} C_M d_Y} \quad (24)$$

To complete the proof, we thus need to show

$$\begin{aligned} \text{tr}(Q_{MM}^{-1} P_{C_M} M M^T) &= M^T Q_{MM}^{-1} P_{C_M} M, \\ \text{tr}(P_{C_M}) &= q \end{aligned} \quad (25)$$

The first expression follows from the trace-property  $\text{tr}(UV) = \text{tr}(VU)$  for any matrices  $U$  and  $V$  of an appropriate size, and the fact that the trace of a scalar is equal to the scalar itself. The second expression follows from the equality between the trace of a projector and its rank, that is

$$\text{tr}(P_{C_M}) = \text{rank}(P_{C_M}) = q, \quad (26)$$

since  $\text{rank}(C_M) = q$ .  $\square$

*Proof of (12)* In case of the atmosphere-fixed scenario, no differential atmospheric delays are to be estimated, i.e.  $\mu = 0$  and  $g = 0$ . This yields  $\mu^{\perp} = I_f$  and  $(D_s^T g)^{\perp} = I_{s-1}$ . According to (3), the frequency-difference misclosures  $M_1$  vanishes, and the vectorized version of the atmosphere-free misclosures  $M_2$  takes the following form

$$M_{\tilde{P}} = [D_n^T \otimes I_f \otimes D_s^T] \text{vec}[\tilde{P}] \quad (27)$$

with the variance matrix (cf. (2))

$$Q_{M_{\tilde{P}} M_{\tilde{P}}} = D_n^T D_n \otimes Q_p \otimes D_s^T W_s^{-1} D_s \quad (28)$$

Upon choosing the array-detector structure (11), matrix  $C_M$  of  $M_{\tilde{P}}$ , introduced in Theorem 1, reads then

$$C_{M_{\tilde{P}}} = D_n^T D_n \otimes I_f \otimes D_s^T D_s \quad (29)$$

Similar expressions are formulated for the carrier-phase observations  $\tilde{\Phi}$ , in case the ambiguities are fixed to their integers. The structures of  $M_{\tilde{\Phi}}$ ,  $Q_{M_{\tilde{\Phi}} M_{\tilde{\Phi}}}$  and  $C_{M_{\tilde{\Phi}}}$  are thus identical to those of  $\tilde{P}$ . Substituting  $M = [M_{\tilde{P}}^T, M_{\tilde{\Phi}}^T]^T$ ,

$$Q_{MM} = \begin{bmatrix} Q_{M_{\tilde{P}} M_{\tilde{P}}} & 0 \\ 0 & Q_{M_{\tilde{\Phi}} M_{\tilde{\Phi}}} \end{bmatrix}, \quad C_M = \begin{bmatrix} C_{M_{\tilde{P}}} & 0 \\ 0 & C_{M_{\tilde{\Phi}}} \end{bmatrix}, \quad (30)$$

an application of Theorem 1 gives (cf. (6))

$$T_q = \frac{\text{tr}\{Q_{M_{\tilde{P}} M_{\tilde{P}}}^{-1} P_{C_{M_{\tilde{P}}}} M_{\tilde{P}} M_{\tilde{P}}^T\} + \text{tr}\{Q_{M_{\tilde{\Phi}} M_{\tilde{\Phi}}}^{-1} P_{C_{M_{\tilde{\Phi}}}} M_{\tilde{\Phi}} M_{\tilde{\Phi}}^T\}}{\text{tr}\{P_{C_{M_{\tilde{P}}}}\} + \text{tr}\{P_{C_{M_{\tilde{\Phi}}}}\}} \quad (31)$$

The proof follows then from

$$P_{C_{M_{\tilde{\Phi}}}} = P_{C_{M_{\tilde{P}}}} = I_{n-1} \otimes I_f \otimes I_{s-1}, \quad (32)$$

and

$$\begin{aligned} \text{tr}\{Q_{M_{\tilde{P}} M_{\tilde{P}}}^{-1} M_{\tilde{P}} M_{\tilde{P}}^T\} &= \text{tr}\{[Q_p^{-1} \otimes W_s P_{e_s}^{\perp}] \tilde{P} P_{D_n} \tilde{P}^T\}, \\ \text{tr}\{Q_{M_{\tilde{\Phi}} M_{\tilde{\Phi}}}^{-1} M_{\tilde{\Phi}} M_{\tilde{\Phi}}^T\} &= \text{tr}\{[Q_{\phi}^{-1} \otimes W_s P_{e_s}^{\perp}] \tilde{\Phi} P_{D_n} \tilde{\Phi}^T\} \end{aligned} \quad (33)$$

with the projectors  $P_{e_s}^{\perp} = W_s^{-1} D_s (D_s^T W_s^{-1} D_s)^{-1} D_s^T$ , and  $P_{D_n} = D_n (D_n^T D_n)^{-1} D_n^T$ .

The proof of (14), (17), (19) and (21) goes along the same lines as the proof of (12).  $\square$

---

## References

- Arnold SF (1981) *The theory of linear models and multivariate analysis*, vol 2. Wiley, New York
- Baarda W (1968) A testing procedure for use in geodetic networks. Technical Report, Publications on Geodesy, New Series, vol 2, no. 5. Netherlands Geodetic Commission, Delft
- Euler HJ, Goad CC (1991) On optimal filtering of GPS dual frequency observations without using orbit information. *J Geod* 65(2):130–143
- Giorgi G, Henkel P, Gunther C (2012) Testing of a statistical approach for local ionospheric disturbances detection. In: *Proceedings of the IEEE-ION Position Location and Navigation Symp (PLANS)*, 2012. IEEE, USA, pp 167–173
- de Jonge PJ (1998) A processing strategy for the application of the GPS in networks. PhD Thesis, Publication on Geodesy, 46, Delft University of Technology, Netherlands Geodetic Commission, Delft
- Khanafseh S, Pullen S, Warburton J (2012) Carrier phase ionospheric gradient ground monitor for GBAS with experimental validation. *Navigation* 59(1):51–60
- Odiijk D, Teunissen PJG, Khodabandeh A (2014) Single-frequency PPP-RTK: theory and experimental results. *IAG Symp* 139:167–173
- Schaffrin B, Bock Y (1988) A unified scheme for processing GPS dual-band phase observations. *Bull Géod* 62(2):142–160
- Teunissen PJG (1997) GPS double difference statistics: with and without using satellite geometry. *J Geod* 71(3):137–148
- Teunissen PJG (1998) Minimal detectable biases of GPS data. *J Geod* 72(4):236–244
- Teunissen PJG (2000) *Testing theory: an introduction*. Series on Mathematical Geodesy and Positioning. Delft University Press, Delft
- Teunissen PJG (2010) Integer least-squares theory for the GNSS compass. *J Geod* 84(7):433–447
- Teunissen PJG (2012) A-PPP: Array-aided precise point positioning with global navigation satellite systems. *IEEE Trans Signal Process* 60(6):2870–2881
- Teunissen PJG, de Bakker PF (2012) Single-receiver single-channel multi-frequency GNSS integrity: outliers, slips, and ionospheric disturbances. *J Geod* 87(2):161–177

---

**Part IX**

**Inverse Modeling, Estimation Theory VIII  
Hotine-Marussi: Geodetic Data Analysis**

---

# Global to Local Moho Estimate Based on GOCE Geopotential Model and Local Gravity Data

R. Barzaghi, M. Reguzzoni, A. Borghi, C. De Gaetani, D. Sampietro, and A.M. Marotta

---

## Abstract

Collocation approach has been applied to get a global Moho model in spherical approximation based on a GOCE geopotential model. A simple single layer model, with known density contrast, has been considered and a linearized relationship between the spherical harmonic coefficients of the anomalous potential and those of the Moho depth has been derived. This allows the covariance propagation from gravity to Moho depth. The derived covariance functions are then used in the collocation estimate of the global Moho depth. In order to be as close as possible to the considered model, reductions for the gravity signal related to topography/bathymetry have been applied. Simulated and real data tests have been performed and the obtained global solution has been compared with Moho estimates available in literature.

The obtained global Moho has been then used as a starting solution for a regional refinement assuming planar approximation. In this second step the computation has been performed in the Central Mediterranean area, based on collocation, local gravity and topography/bathymetry data.

---

## Keywords

Moho • Gravity data • GOCE model • Collocation

---

## 1 Introduction

The Mohorovičić discontinuity (Moho) separates the lower crust from the upper mantle (Turcotte and Schubert 1982). This discontinuity causes seismic wave refraction and reflection and can be thus investigated via seismic methods (Parker 1973; Lebedev et al. 2013). Seismic methods can give accurate estimates that are, however, obtained along sections

which are usually sparse and not homogeneously distributed on the Earth.

The mean density variation between the lower crust and the upper mantle is usually set at  $0.4 \text{ g/cm}^3$  (Anderson 1989), even though this value has been recently revised and raised to values ranging from  $0.448$  to  $0.485 \text{ g/cm}^3$  (Tenzer et al. 2012; Sjöberg and Bagherbandi 2011). The gravimetric signal that is related to the Moho has, at the Earth surface, a strong signature and a standard deviation that can range from 50 to 100 mGal. So, also gravity can be used to estimate this surface provided that proper constraints are used to overcome the intrinsic non-uniqueness of the related inverse gravimetric problem (Tarantola 2005; Sampietro and Sansó 2012). Gravity data are densely and homogeneously distributed on the Earth and can thus give valuable information on the Moho structure. Global geopotential models can be profitably used to this purpose. Particularly, satellite dedicated gravity missions (Reigber et al. 1999; Tapley et al.

---

R. Barzaghi (✉) • M. Reguzzoni • A. Borghi • C. De Gaetani  
DICA, Politecnico di Milano, Milano, Italy  
e-mail: [riccardo.barzaghi@polimi.it](mailto:riccardo.barzaghi@polimi.it)

D. Sampietro  
GReD, Politecnico di Milano, Milano, Italy

A.M. Marotta  
Dip. "Ardito Desio", Università degli Studi di Milano, Milano, Italy

2004; Albertella et al. 2002) allowed the estimate of reliable satellite only models that can be used in combination with ground based gravity data (Shin et al. 2007) for estimating the Moho depths.

Note that, Moho estimates can be also obtained using both gravity data and seismic derived Moho depths (Eshagh et al. 2011).

In this work, a method for the inversion of gravity to get the Moho depth based on collocation is presented (Krarup 1969; Tscherning 1985; Moritz 1989). To this aim, the gravity signal related to the Moho discontinuity is expressed, at global scale, as a linearized functional of the Moho depth variation with respect to a suitable mean depth sphere. Furthermore, in order to guarantee uniqueness, it is assumed to have a single layer model with a known constant density contrast (Barzaghi and Sansò 1988). This global estimate, which can be also integrated with known Moho depths coming from seismic analyses, can be then used to refine the estimated Moho depths at local scale, following a similar scheme (Barzaghi et al. 1992; Knudsen 1993; Arabelos et al. 2007) and taking into account detailed density and local gravity information.

## 2 The Methodology

Assume to have a single layer model in spherical approximation (see Fig. 1).

The sphere having radius  $R_2$  is considered to be the mean Earth sphere while the Moho depth is ranging around an internal sphere of radius  $R_1$ . Thus the Moho topography is described as  $R_1 + \delta R_1(\theta, \lambda)$ . We further assume that

$$\int_{\sigma} \delta R_1(\theta, \lambda) d\sigma = 0 \quad (1)$$

where  $\sigma$  is the unit sphere and  $(\theta, \lambda)$  are, respectively, colatitude and longitude.

Assume also that the anomalous potential  $T(P)$  outside the external sphere is given by the mass contained in the volume bounded by  $R_1 + \delta R_1(\theta, \lambda)$  and  $R_2$ , having density contrast  $\delta\rho_{12} = \rho_1 - \rho_2$ , being  $\rho_1$  the crust density and  $\rho_2$  the mantle one. In these hypotheses and referring to Fig. 1, we can write

$$\begin{aligned} T(P) &= G \int_V \frac{\delta\rho_{12}(Q)}{r_{PQ}} dv_Q = \\ &= G \int_{\sigma} d\sigma_Q \int_{R_1+\delta R_1}^{R_2} dr_Q r_Q^2 \frac{\delta\rho_{12}(Q)}{r_{PQ}} = \\ &= G \int_{\sigma} d\sigma_Q \int_{R_1}^{R_2} dr_Q r_Q^2 \frac{\delta\rho_{12}(Q)}{r_{PQ}} + \end{aligned}$$

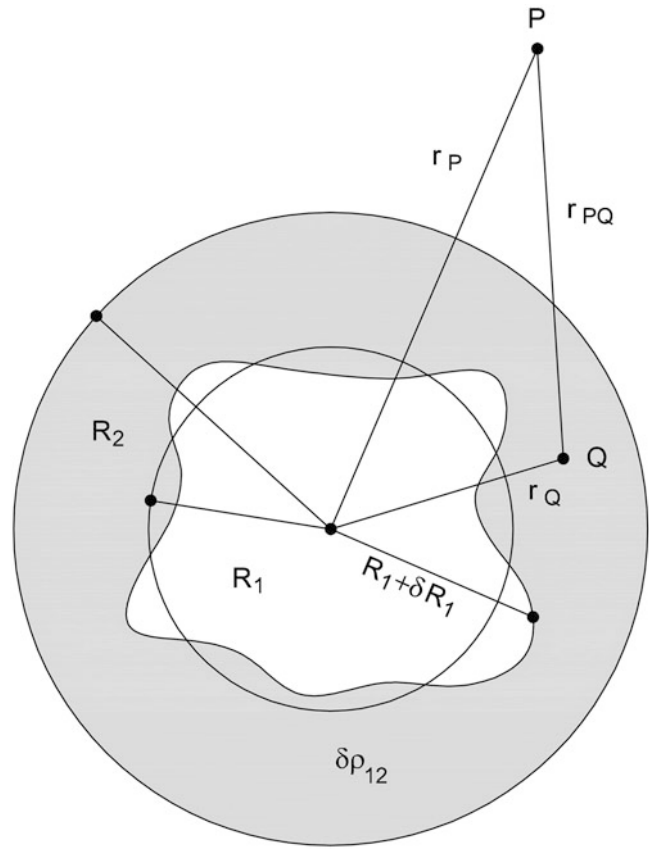


Fig. 1 The single layer model describing the Moho

$$\begin{aligned} &+ G \int_{\sigma} d\sigma_Q \int_{R_1+\delta R_1}^{R_1} dr_Q r_Q^2 \frac{\delta\rho_{12}(Q)}{r_{PQ}} = \\ &= \bar{T}(P) + \delta T(P) \end{aligned} \quad (2)$$

where  $G$  is the gravitational constant.

In our analysis we only consider the second term which allows the estimation of  $\delta R_1(\theta, \lambda)$ . This can be written in the following form

$$\begin{aligned} \delta T(P) &= \\ &= G \int_{\sigma} d\sigma_Q \int_{R_1+\delta R_1}^{R_1} dr_Q r_Q^2 \frac{\delta\rho_{12}(Q)}{r_{PQ}} \cong \\ &\cong \frac{G\delta\bar{\rho}_{12}}{r_P} \int_{\sigma} d\sigma_Q \int_{R_1+\delta R_1}^{R_1} dr_Q r_Q^2 \sum_{n=0}^{+\infty} \frac{r_Q^n}{r_P^n} P_n(\cos \psi_{PQ}) \\ &= \frac{G\delta\bar{\rho}_{12}}{r_P} \int_{\sigma} d\sigma_Q \sum_{n=0}^{+\infty} \frac{1}{r_P^n} \frac{[R_1^{n+3} - (R_1 + \delta R_1)^{n+3}]}{(n+3)} \\ &\quad \times P_n(\cos \psi_{PQ}) \end{aligned} \quad (3)$$

using the standard expansion of  $1/r_{PQ}$  in terms of Legendre polynomials (Mac Millan 1958). In this equation, we made a further simplification by considering a mean constant density contrast  $\delta\bar{\rho}_{12}$  instead of the varying density contrast  $\delta\rho_{12}(Q)$ . If we now linearize the term

$$(R_1 + \delta R_1)^{n+3} \cong R_1^{n+3} \left[ 1 + (n+3) \frac{\delta R_1}{R_1} \right] \quad (4)$$

we obtain

$$\delta T(P) = -G\delta\bar{\rho}_{12}R_1 \int_{\sigma} d\sigma_Q \delta R_1(Q) \sum_{n=1}^{+\infty} \left( \frac{R_1}{r_P} \right)^{n+1} \times P_n(\cos \psi_{PQ}) \quad (5)$$

where now the series starts from  $n=1$  by virtue of (1).

Using the decomposition formula which allows expressing the Legendre polynomials in terms of spherical harmonics (Mac Millan 1958), we obtain

$$\begin{aligned} \delta T(P) &= -G\delta\bar{\rho}_{12}R_1 \int_{\sigma} d\sigma_Q \delta R_1(Q) \sum_{n=1}^{+\infty} \left( \frac{R_1}{r_P} \right)^{n+1} \\ &\quad \times P_n(\cos \psi_{PQ}) = \\ &= -G\delta\bar{\rho}_{12}R_1 \int_{\sigma} d\sigma_Q \delta R_1(Q) \sum_{n=1}^{+\infty} \left( \frac{R_1}{r_P} \right)^{n+1} \\ &\quad \times \frac{1}{2n+1} \sum_{m=-n}^n Y_{nm}(P) Y_{nm}(Q) = \\ &= -4\pi G\delta\bar{\rho}_{12}R_1 \sum_{n=1}^{+\infty} \sum_{m=-n}^n \frac{1}{2n+1} \left( \frac{R_1}{r_P} \right)^{n+1} \\ &\quad \times \delta R_{nm}^{(1)} Y_{nm}(P) \end{aligned} \quad (6)$$

where

$$\delta R_{nm}^{(1)} = \frac{1}{4\pi} \int_{\sigma} d\sigma_Q \delta R_1(Q) Y_{nm}(Q) \quad (7)$$

By evaluating this function on any external sphere of radius  $R_e > R_2$ , we finally have

$$\begin{aligned} \delta T^0(P) &= \sum_{n=1}^{+\infty} \sum_{m=-n}^n \delta T_{nm}^0 Y_{nm}(P) = \delta T(P)|_{P \in R_e} = \\ &= -4\pi G\delta\bar{\rho}_{12}R_1 \sum_{n=1}^{+\infty} \sum_{m=-n}^n \frac{1}{2n+1} \left( \frac{R_1}{R_e} \right)^{n+1} \\ &\quad \times \delta R_{nm}^{(1)} Y_{nm}(P) \end{aligned} \quad (8)$$

which implies

$$\delta T_{nm}^0 = -\frac{4\pi G\delta\bar{\rho}_{12}R_1}{2n+1} \left( \frac{R_1}{R_e} \right)^{n+1} \delta R_{nm}^{(1)} \quad (9)$$

### 3 The Covariance Propagation Law for the Collocation Estimation of the Moho

The geodetic applications of collocation formula are based on the covariance propagation law to the linear functional of the anomalous potential (Moritz 1989). This rule allows defining the covariance function of any geodetic observation which can be expressed as a linear(ized) functional of the anomalous potential. In the same way, the propagation law can be applied to (8) and (9) which express the observed  $\delta T^0(P)$  values as a linear functional of  $\delta R_1(P)$ .

Since, due to Eq. (8),  $\delta T^0(P)$  is given by

$$\delta T^0(P) = \sum_{n=1}^{+\infty} \sum_{m=-n}^n \delta T_{nm}^0 Y_{nm}(P) \quad (10)$$

it follows (Moritz 1989) that its auto-covariance function  $C_{\delta T^0 \delta T^0}(P, Q)$  can be defined as

$$\begin{aligned} C_{\delta T^0 \delta T^0}(P, Q) &= \sum_n \sigma_n^2 (\delta T^0) P_n(\cos \psi_{PQ}) \\ \sigma_n^2 (\delta T^0) &= \sum_{m=-n}^n (\delta T_{nm}^0)^2 \end{aligned} \quad (11)$$

Based on that and applying covariance propagation law, one can get the cross-covariance  $C_{\delta R_1 \delta T^0}(P, Q)$  as

$$\begin{aligned} C_{\delta R_1 \delta T^0}(P, Q) &= \sum_{n=1}^{+\infty} \left( -\frac{2n+1}{4\pi G\delta\bar{\rho}_{12}R_1} \right) \left( \frac{R_e}{R_1} \right)^{n+1} \\ &\quad \times \sigma_n^2 (\delta T^0) P_n(\cos \psi_{PQ}) \end{aligned} \quad (12)$$

Both covariances are needed to compute the collocation estimate of  $\delta R_1(P)$  which is expressed by

$$\delta \widehat{R}_1(P) = C_{\delta R_1 \delta T^0} [C_{\delta T^0 \delta T^0} + \sigma_v^2 I]^{-1} \delta T^0 \quad (13)$$

The  $\delta \widehat{R}_1$  estimate can be then computed following the standard collocation procedure. Given the observed  $\delta T^0(P)$  values on any sphere  $R_e > R_2$ , the empirical covariance is computed and then fitted using the model (11). In this way,  $\sigma_n^2(\delta T^0)$  and  $\sigma_v^2$  (the variance of the noise associated to the observed  $\delta T^0$ ) are estimated. Having defined the  $\sigma_n^2(\delta T^0)$  values, the cross-covariance between  $\delta T^0(P)$  and  $\delta R_1(P)$  can be derived too, taking (12) into account. This allows the

**Table 1** The statistics of the differences between predicted and known Moho depths

	Mean (km)	$\sigma$ (km)	Min (km)	Max (km)
SIMULATION A	0.005	0.057	-0.528	0.541
SIMULATION B	0.004	1.403	-13.570	22.742

computation of formula (13) that gives the wanted estimate of the Moho depth, i.e.  $\widehat{R}_{Moho} = R_1 + \delta \widehat{R}_1$ .

## 4 Input Data and Results

Two simulations have been set up to prove the method feasibility. The estimated GEMMA Moho depth (Reguzzoni et al. 2013; Reguzzoni and Sampietro 2015) has been assumed as the known surface for computing a synthetic potential signal on the external sphere  $R_e = R_2 + 50 \text{ km} = (6,371 + 50) \text{ km}$ , assuming a constant density contrast  $\delta \bar{\rho}_{12} = -0.4 \text{ g/cm}^3$ . The simulated potential has been computed either using Eq. (8) (SIMULATION A), with  $\delta R_{nm}^{(1)}$  coming from harmonic analysis of the known  $\delta R_1(\theta, \lambda)$  function (up to degree  $n = 180$ ), and a quadrature formula of the integral over the given layer (Reguzzoni et al. 2013) (SIMULATION B). The two data sets have been used for computing the empirical covariances that are needed to set the auto- and cross-model covariances to be used in Eq. (13). The estimated depths have then been compared with the known Moho surface and the statistics are shown in Table 1.

From these values, it can be concluded that the method is able to retrieve the given Moho depth in a quite accurate and precise way. SIMULATION A is clearly better than SIMULATION B because in that case only linear terms were considered in generating the potential signal. However, it must be underlined that in SIMULATION B absolute differences higher than  $3\sigma$  are around 2%. Also, they are in areas, such as the Andes, where the Moho structure is complex and sharp variations occur. Since the devised model is based on a linearization, such high frequencies cannot be retrieved properly. In these areas, local investigations are then needed and second order refinements are to be estimated starting from the global smooth solution and local data. Naturally, this is a close loop simulation that only proves that the mathematical and the numerical formulations are sound, provided that the hypotheses on the model are fulfilled.

Real data have then been analysed taking into account the GOCE model obtained via direct approach (Pail et al. 2011-release R4). The anomalous potential signal implied by this model has been evaluated at 50 km altitude and up to degree 180, according to what has been done in the simulations. Corrections have then been computed to come as close as possible to the model given in Sect. 2. It must be underlined that only major corrections have been taken

into account (we recall that, in our mind, this is only the first step in the Moho estimating procedure which can be then refined using detailed density and gravity information at local level). As for the topography/bathymetry compensation, the ETOPO1 DTM (Amante and Eakins 2009) has been considered at a grid step of  $30'$ . The potential effect of the masses above sea level has been computed at 50 km altitude assuming a constant  $2.67 \text{ g/cm}^3$  density. This has been done following the quadrature formula method used in the simulation. In the same way, the bathymetry effect has been evaluated using a  $1.9 \text{ g/cm}^3$  seawater density contrast. The potential signal coming from the positive topography has been then removed from the GOCE model signal while the bathymetry effect has been added. Ice sheets from the ETOPO1 model have been also compensated, removing the potential signal of ice masses above sea level with a standard  $0.98 \text{ g/cm}^3$  water density and adding the signal of those below sea level (mainly in Antarctica) with a  $1.92 \text{ g/cm}^3$  density contrast. The derived data set has been then used in the devised collocation procedure assuming that, after the above-mentioned reductions, these data can be considered as the  $\delta T^0(\theta, \lambda)$  signal in Eq. (8). The empirical covariance function of this signal has been evaluated and fitted with the model covariance coming from the harmonic analysis of  $\delta T^0(\theta, \lambda)$  itself. The plot of the empirical and the model covariance is given in Fig. 2

Collocation formula (13) has been then applied with  $R_e = (6,371 + 50) \text{ km}$ ,  $\delta \bar{\rho}_{12} = (2.9 - 3.3) \text{ g/cm}^3 = -0.4 \text{ g/cm}^3$  and  $R_1 = (6,371 - 23) \text{ km}$  (the mean depth at  $-23 \text{ km}$  for the Moho has been assumed based on literature; Bassin et al. 2000).

The estimated Moho depths have been computed on an equal-area  $1^\circ \times 1^\circ$  grid and are displayed in Fig. 3.

Furthermore, the estimated depths have been interpolated on two sections in areas of expected strong Moho variations.

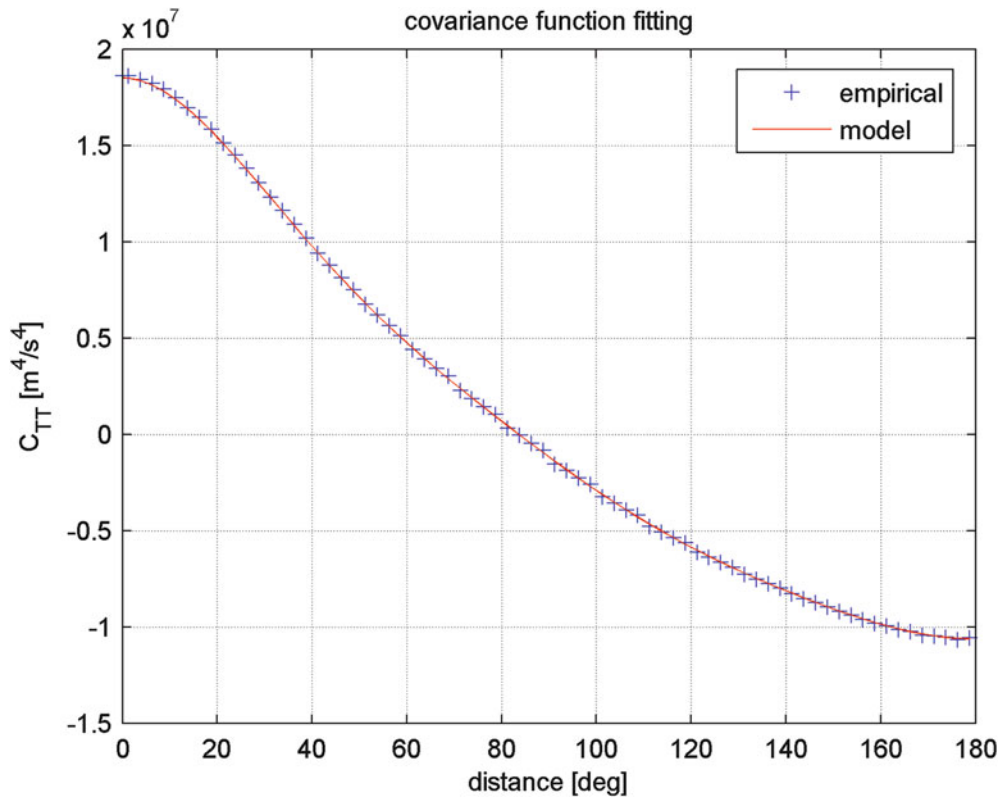
These two sections are across the Andes (section at constant latitude  $\varphi = -23^\circ$ ; Fig. 4a) and the Himalaya region (section at constant latitude  $\varphi = 33^\circ$ ; Fig. 4b). The obtained values (labelled as *Pred*) have been plotted together with three other existing Moho models, i.e. the CRUST2 (Bassin et al. 2000), the model by Meier et al. (2007) and the GEMMA estimate (Reguzzoni and Sampietro 2015).

The new estimated Moho is coherent with these models along these sections, particularly with the CRUST2 model. In Table 2, the statistics over the whole global grid of the residuals among all these models are listed.

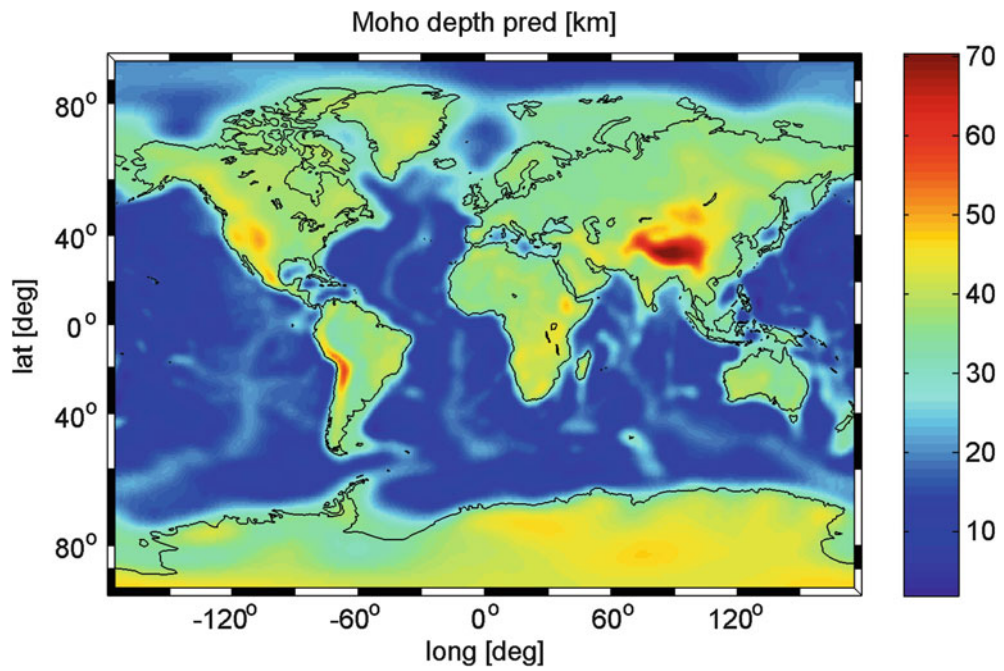
As one can see from Table 2, the differences between the Moho model based on collocation and the other considered models are statistically equivalent with those between them.

This predicted global Moho depth has been then considered as a starting point for a local estimate in the Central Mediterranean. In the area  $35^\circ < \varphi < 49^\circ$ ,  $4^\circ < \lambda < 20^\circ$ , local gravity data (Barzaghi et al. 2007) have been selected

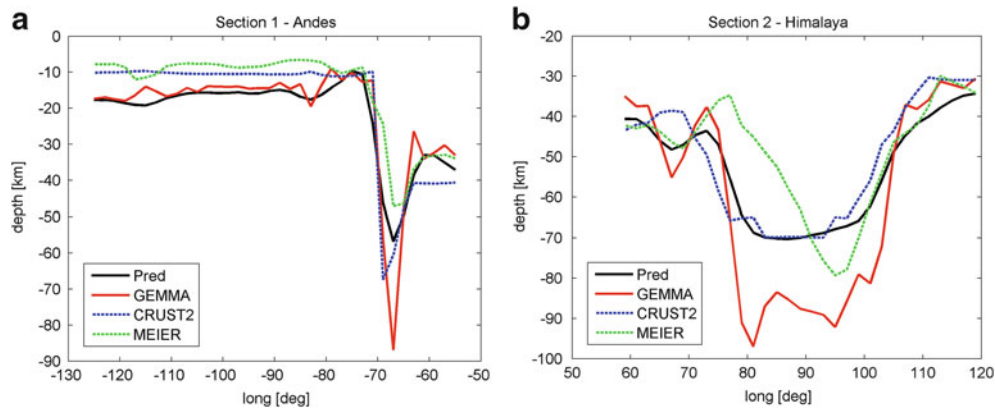




**Fig. 2** The empirical and the model covariance of the reduced GOCE data



**Fig. 3** The predicted Moho depth



**Fig. 4** (a) The Andes section (b) The Himalaya section

**Table 2** The statistics of the differences among the considered Moho models

Model difference	Mean (km)	$\sigma$ (km)	Min (km)	Max (km)
Pred-GEMMA	2.723	4.563	-35.833	20.775
Pred-CRUST2	2.116	5.392	-31.257	27.125
Pred-MEIER	4.307	5.917	-19.871	31.676
GEMMA-CRUST2	-0.607	5.937	-30.662	38.920
GEMMA-MEIER	1.584	7.000	-30.338	61.857
CRUST2-MEIER	2.191	5.938	-21.233	49.712

**Table 3** The statistics of the differences between the collocation based and the ESC Moho

Model difference	Mean (km)	$\sigma$ (km)	Min (km)	Max (km)
Local solution-ESC	-2.7	10.7	-25.7	21.2

and the Bouguer topographic/bathymetric reduction has been applied based on SRTM/DTM ([www2.jpl.nasa.gov/srtm/](http://www2.jpl.nasa.gov/srtm/)) and NOAA bathymetry ([www.ngdc.noaa.gov/mgg/bathymetry](http://www.ngdc.noaa.gov/mgg/bathymetry)). The Bouguer gravity anomalies have been then used in a local collocation approach in planar approximation, as described in Barzaghi et al. (1992).

The use of planar approximation for Moho estimation is justified by the dimension of the area. At regional scale (regions of the order of  $10^\circ \times 10^\circ$ ) the spherical and planar solution are known to be practically the same, showing differences smaller than 0.5 km (Sampietro 2011).

In this preliminary test, the global solution has been only used to define the mean depth  $H$  and for mapping the estimated  $\varepsilon$  value, i.e. the depth anomaly with respect to  $H$  (the density contrast has been set to  $\delta\bar{\rho}_{12} = -0.4 \text{ g/cm}^3$ ). The obtained estimate is plotted in Fig. 5.

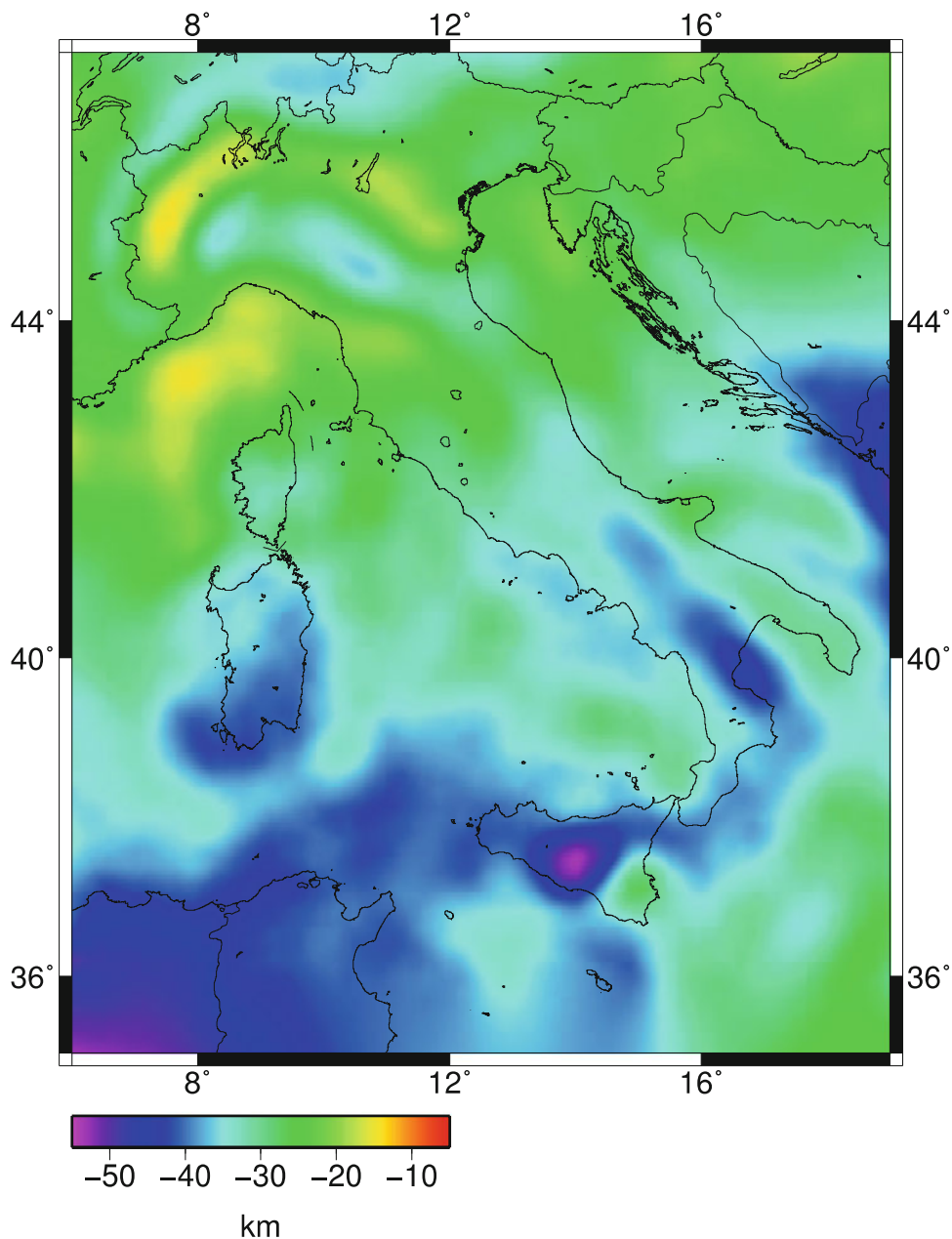
Although this solution has, as expected, a higher frequency pattern with respect to the global solution, it is still quite poor if compared with e.g. the ESC estimate (Grad et al. 2009) as it is clearly seen in Table 3.

Thus, it seems that a more refined data reduction for local crustal density anomalies is required in order to get a more detailed and reliable solution.

## 5 Remarks and Conclusions

The global Moho estimate derived from a properly reduced GOCE potential model proved to be effective. Based on a simple single layer model in spherical approximation, collocation method can provide a feasible Moho surface that is, from a first preliminary comparison, statistically equivalent to other Moho estimates from literature. This estimate can be thus assumed as a reference surface for local refinements. In this paper this has been done partially, i.e. using the global solution for only defining the mean depth  $H$  and as a reference surface for mapping the local estimate. In the future, a deeper interaction between the global and the local procedure will be studied and devised. In estimating the global solution, a potential signal that is coherent with a known DTM/bathymetry is defined. This gravity signal and the related DTM/bathymetry can be used in reducing the data in the local area. Furthermore, a residual terrain correction can be computed with respect to the used DTM/bathymetry (i.e. the local gravity observations are reduced for the global Moho effect and a coherent residual terrain effect). In this way, reduced gravity data can be obtained that reflect high frequency structures of the Moho. Further refinements, always at local scale, can be obtained by using density anomaly information thus hopefully improving the final estimate. Also, since collocation allows the joint use of gravity and depth estimates, refinements could be obtained in a joint estimate based on gravity/Moho depth data.

**Fig. 5** The final Moho estimated depth based on local data



## References

- Albertella A, Migliaccio F, Sansò F (2002) GOCE: the Earth gravity field by space gradiometry. *Celest Mech Dyn Astron* 83(1–4):1–15
- Amante C, Eakins BW (2009) ETOPO1 1 arc-minute global relief model: procedures, data sources and analysis. NOAA Technical Memorandum NESDIS NGDC-24
- Anderson DL (1989) *Theory of the Earth*. Blackwell, Boston/Oxford/London/Edinburgh/Melbourne
- Arabelos DN, Mantzios G, Tsoulis D (2007) Moho depths in the Indian ocean based on the inversion of satellite gravity data. In: Ip W-H, Chen Y-T (eds) *Advances in geosciences*, vol 9, solid Earth, ocean science & atmospheric sciences. World Scientific, Singapore, pp 41–52
- Barzaghi R, Sansò F (1988) Remarks on the inverse gravimetric problem. *Geophys J R Astron Soc* 92(3):505–511, The Universities Press, Belfast
- Barzaghi R, Gandino A, Sansò F, Zenucchini C (1992) The collocation approach to the inversion of gravity data. *Geophys Prospect* 40:429–451
- Barzaghi R, Borghi A, Carrion D, Sona G (2007) Refining the estimate of the Italian quasi-geoid. *Bollettino di Geodesia e Scienze Affini*, n 3
- Bassin C, Laske G, Masters G (2000) The current limits of resolution for surface wave tomography in North America. *EOS Trans AGU* 81:F897
- Eshagh M, Bagherbandi M, Sjöberg L (2011) A combined global Moho model based on seismic and gravimetric data. *Acta Geodaetica et Geophysica Hungarica* 46(1):25–38

- Grad M, Tira T, ESC Working Group (2009) The Moho depth map of the European plate. *Geophys J Int* 176:279–292
- Knudsen P (1993) Integrated inversion of gravity data. Final report of Norsk hydro R&D project, Kort og Matrikelstyrelsen, Copenhagen. ISBN 8774500872
- Krarpur T (1969) A contribution to the mathematical foundation of physical geodesy. Meddelelse No. 44, Geodætisk Institut, København
- Lebedev S, Adam JMC, Meier T (2013) Mapping the Moho with seismic surface waves: a review, resolution analysis, and recommended inversion strategies. doi:[10.1016/j.tecto.2012.12.030](https://doi.org/10.1016/j.tecto.2012.12.030)
- Mac Millan WD (1958) The theory of the potential. Dover Publications, New York
- Meier U, Curtis A, Trampert J (2007) Global crustal thickness from neural network inversion of surface wave data. *Geophys J Int* 169:706–722
- Moritz H (1989) Advanced physical geodesy. Wichmann, Karlsruhe
- Pail R, Bruinsma SL, Migliaccio F, Förste C, Goiginger H, Schuh WD, Höck E, Reguzzoni M, Brockmann JM, Abrikosov O, Veicherts M, Fecher T, Mayrhofer R, Krasbutter I, Sansò F, Tscherning CC (2011) First GOCE gravity field models derived by three different approaches. *J Geod* 85(11):819–843. doi:[10.1007/s00190-011-0467](https://doi.org/10.1007/s00190-011-0467)
- Parker RL (1973) The rapid calculation of potential anomalies. *Geophys J R Astron Soc* 31:447–455
- Reguzzoni M, Sampietro D (2015) GEMMA: An Earth crustal model based on GOCE satellite data. *International Journal of Applied Earth Observation and Geoinformation*, vol 35, Part A, pp 31–43. doi:[10.1016/j.jag.2014.04.002](https://doi.org/10.1016/j.jag.2014.04.002)
- Reguzzoni M, Sampietro D, Sansò F (2013) Global Moho from the combination of the CRUST2.0 model and GOCE data. *Geophys J Int* 195(1):222–237. doi:[10.1093/gji/ggt247](https://doi.org/10.1093/gji/ggt247)
- Reigber C, Casper R, Pätzold W (1999) The CHAMP geopotential mission. In: IAA 2nd international symposium on small satellites for Earth observation, Berlin, pp 25–28
- Sampietro D (2011) GOCE exploitation for Moho modeling and applications. In: Proceedings of the 4th international GOCE user workshop, vol 31, Munich, ESA Publication, SP-696
- Sampietro D, Sansò F (2012) Uniqueness theorems for inverse gravimetric problems. In: VII Hotine-Marussi symposium on mathematical geodesy. Springer, Berlin/Heidelberg, pp 111–115
- Shin YM, Xu H, Braitenberg C, Fang J, Wang Y (2007) Moho undulations beneath Tibet from GRACE-integrated gravity data. *Geophys J Int* 170(3):971–985
- Sjöberg LE, Bagherbandi M (2011) A method of estimating the Moho density contrast with a tentative application by EGM2008 and CRUST2.0. *Acta Geophys* 58:1–24
- Tapley BD, Bettapur S, Watkins M, Reigber C (2004) The gravity recovery and climate experiment: mission overview and early results. *Geophys Res Lett* 31(9):L09607
- Tarantola A (2005) Inverse problem theory and methods for model parameter estimation. SIAM, Philadelphia
- Tenzer R, Hamayun, Novák P, Gladkikh V, Vajda P (2012) Global crust-mantle density contrast estimated from EGM2008, DTM2008, CRUST2.0, and ICE-5G. *Pure Appl Geophys* 169(9):1663–1678. doi:[10.1007/s00024-011-0410-3](https://doi.org/10.1007/s00024-011-0410-3)
- Tscherning CC (1985) Local approximation of the gravity potential by least squares collocation. In: Schwarz KP (ed) Proceedings of the international summer school on local gravity field approximation, Beijing, Publ. 60003, University of Calgary, Calgary
- Turcotte DL, Schubert G (1982) Geodynamics: applications of continuum physics to geological problems. Wiley, New York

---

# An Overview of Adjustment Methods for Mixed Additive and Multiplicative Random Error Models

Yun Shi, Peiliang Xu, and Junhuan Peng

---

## Abstract

Geodetic adjustment theory has been developed on the basis of a linear or nonlinear Gauss-Markov model, in which the random errors of measurements are always assumed to be independent of the true values of measurements themselves and naturally added to the functional model. However, modern geodetic instruments and geodetic imaging systems have clearly shown that the random errors of such measurements consist of two parts: one is of local nature and has nothing to do with the quantity under observation, and the other is proportional to the true value of measurement. From the statistical point of view, these two types of errors are called *additive* and *multiplicative* errors, respectively. Obviously, the conventional geodetic adjustment theory and methods for Gauss-Markov models with additive errors cannot theoretically meet the need of processing measurements contaminated by mixed additive and multiplicative random errors. This paper presents an overview of parameter estimation methods for processing mixed additive and multiplicative random errors. More specifically, we discuss two types of methods to estimate parameters in a mixed additive and multiplicative error model, namely, quasi-likelihood and least-squares-based methods. From this point of view, we extend the conventional adjustment theory and methods and give a solid theoretical foundation to process geodetic measurements contaminated by mixed additive and multiplicative random errors. Finally, we further discuss parameter estimation with prior information.

---

## Keywords

Generalized estimating equation • Least squares • Multiplicative and additive errors • Quasi-likelihood

---

Y. Shi (✉)

School of Geomatics, Xi'an University of Science and Technology,  
Xi'an 710054, PR China

e-mail: [shiyun0908@hotmail.com](mailto:shiyun0908@hotmail.com)

P. Xu

Disaster Prevention Research Institute, Kyoto University, Uji, Kyoto  
611-0011, Japan

e-mail: [pxu@rcep.dpri.kyoto-u.ac.jp](mailto:pxu@rcep.dpri.kyoto-u.ac.jp)

J. Peng

School of Land Science and Geomatics, China University of  
Geosciences, Beijing, PR China

---

## 1 Introduction

Geodetic adjustment theory has been developed by assuming the following model of measurements:

$$\left. \begin{aligned} \mathbf{y} &= \mathbf{f}(\boldsymbol{\beta}) + \boldsymbol{\epsilon} \\ E(\mathbf{y}) &= \mathbf{f}(\boldsymbol{\beta}) \\ D(\mathbf{y}) &= \mathbf{W}^{-1}\sigma^2 \end{aligned} \right\}, \quad (1)$$

where  $\mathbf{y}$  is a vector of measurements,  $\mathbf{f}(\boldsymbol{\beta})$  is the mathematical or functional model which describes the physical or geometrical relationships between the measurements,  $\boldsymbol{\beta}$

is the real-valued vector of unknown parameters to be estimated,  $\epsilon$  is the random error vector of the measurements. Very often, we also assume that  $\epsilon$  is of zero mean and variance-covariance matrix  $\mathbf{W}^{-1}\sigma^2$ , with  $\mathbf{W}$  being a given weight matrix of measurements and  $\sigma^2$  an unknown positive scalar (the variance of unit weight),  $E(\cdot)$  and  $D(\cdot)$  stand for the expectation and variance-covariance matrix of the measurements, respectively. The most important feature of adjustment model (1) is that the random errors  $\epsilon$  are added to the functional model  $\mathbf{f}(\boldsymbol{\beta})$ . In other words, the sizes or magnitudes of random errors are independent of the true values of measured quantities.

However, in geodetic practice, we know that this assumption is not necessarily always true. For example, we know that the accuracy of an EDM, VLBI and/or GPS baseline is proportional to the length of the baseline itself, namely,

$$\sigma_L^2 = a^2 + b^2 L^2, \quad (2)$$

(see e.g., Ewing and Mitchell 1970; MacDoran 1979; Seeber 2003; Petrov et al. 2010), where both  $a$  and  $b$  are constants. Physically, the constant  $a$  may be more specific to the local environment of stations and  $b$  more to the path of propagation of light/electronic waves (see e.g., Xu et al. 2013). From the statistical point of view, the accuracy formula (2) is equivalent to the following representation of random errors:

$$\epsilon_L = \epsilon_a + L \epsilon_b, \quad (3)$$

where  $\epsilon_L$  is the random error of  $L$ , and  $\epsilon_a$  and  $\epsilon_b$  stand for the random errors of mean zero and variances  $a^2$  and  $b^2$ , respectively, if  $\epsilon_a$  and  $\epsilon_b$  are assumed to be statistically independent. The error representation (3) clearly indicates that the random error  $\epsilon_L$  is proportional to the measured baseline. In geodetic practice, both  $\epsilon_a$  and  $\epsilon_b$  are generally assumed to be normally distributed. For other modern space observation technology such as SLR (see e.g. Pearlman et al. 2002; Seeber 2003) and DORIS (see e.g. Willis et al. 2010), since they essentially utilize electromagnetic waves for observation and go through the same physical media as VLBI and GPS, we conjecture that errors of SLR and DORIS baselines should also show multiplicative error behavior, which will be a topic of research in the future.

Modern geodetic technology also fully utilizes coherent imaging systems such as Synthetic Aperture Radar (SAR) images and Light Detection And Ranging (LiDAR). As is well known, SAR images are contaminated by speckle noise (see e.g. Goodman 1976; Ulaby et al. 1986; Oliver 1991; López-Martínez et al. 2011) and the corresponding observational equation can be represented as follows:

$$y_{ij} = s_{ij}(1 + \epsilon_{ij}), \quad (4)$$

where  $y_{ij}$  is the measurement,  $s_{ij}$  the true (or noiseless) value of the signal and  $\epsilon_{ij}$  the random error with zero mean and variance  $\sigma^2$ . Intensity measurements of SAR type are usually assumed to have a gamma-distribution. Other imaging systems would also produce Gaussian multiplicative random errors (see e.g., Tian et al. 2001). Range measurements of LiDAR are also shown to be contaminated by multiplicative speckle errors (see e.g., Flamant et al. 1984; Wang and Pruitt 1992; Hill et al. 2003).

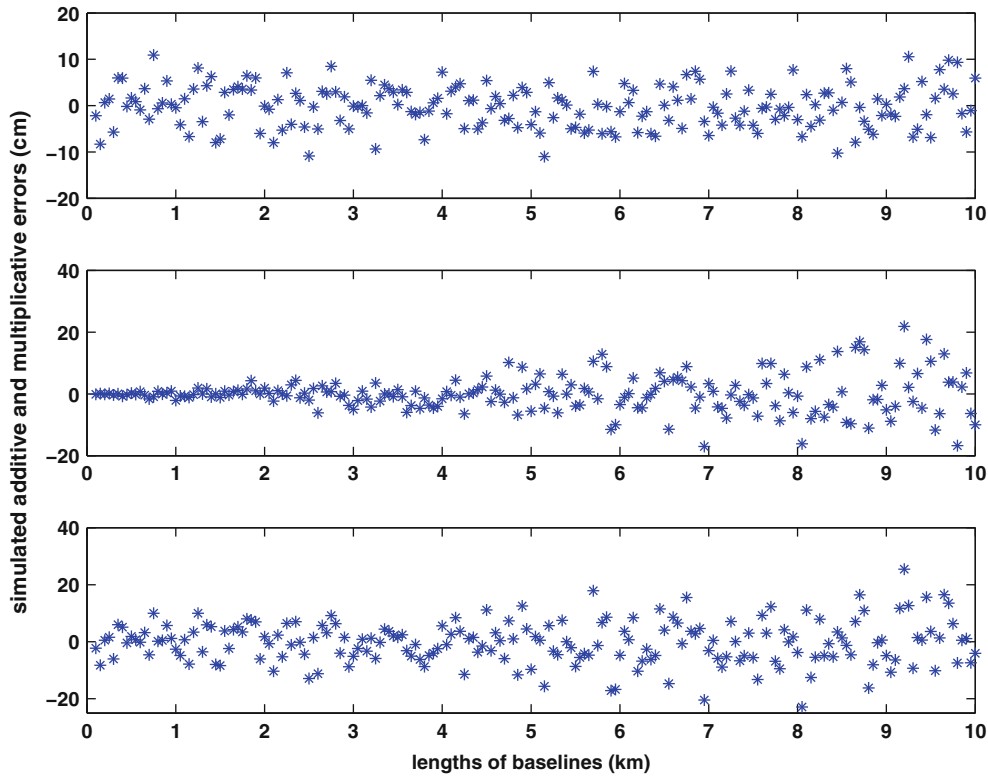
The paper is organized as follows. Section 2 will first define mixed additive and multiplicative error models. In Sects. 3 and 4, we will discuss two important classes of methods for parameter estimation in mixed additive and multiplicative error models, namely, quasi-likelihood and least squares (LS). Computational algorithms will be briefly given. If the reader is interested in other methods such as cumulant moment and variational methods, he may refer to, Swami (1994) and Aubert and Aujol (2008), for example. Actually, variational methods assume gamma-distributions for intensity measurements and then add an extra smoothness or regularized term to the log-likelihood of the gamma-distributions for de-speckling or de-noising multiplicative random errors, as can be seen in Xu (1999) and Aubert and Aujol (2008). We will then extend the bias-corrected LS method to the case with prior information in Sect. 5. Finally, we will then finish our paper with some concluding remarks in Sect. 6.

## 2 Mixed Additive and Multiplicative Error Models

We will now extend the conventional Gauss-Markov adjustment model (1) to account for both additive and multiplicative errors. The new starting model of adjustment becomes

$$\left. \begin{aligned} \mathbf{y} &= \mathbf{f}(\boldsymbol{\beta}) \odot (\mathbf{1} + \boldsymbol{\epsilon}_m) + \boldsymbol{\epsilon}_a \\ E(\mathbf{y}) &= \mathbf{f}(\boldsymbol{\beta}) \\ E(\boldsymbol{\epsilon}_m) &= \mathbf{0}, \quad D(\boldsymbol{\epsilon}_m) = \boldsymbol{\Sigma}_m \\ E(\boldsymbol{\epsilon}_a) &= \mathbf{0}, \quad D(\boldsymbol{\epsilon}_a) = \boldsymbol{\Sigma}_a \end{aligned} \right\}, \quad (5)$$

where  $\mathbf{y}$  and  $\mathbf{f}(\boldsymbol{\beta})$  have been defined in (1),  $\odot$  stands for the Hadamard product of matrices and/or vectors,  $\mathbf{1}$  for the vector with all its elements being equal to unity, both of the random errors  $\boldsymbol{\epsilon}_m$  and  $\boldsymbol{\epsilon}_a$  are of mean zero and variance-covariance matrices  $\boldsymbol{\Sigma}_m$  and  $\boldsymbol{\Sigma}_a$ , respectively. Since the random vector  $\boldsymbol{\epsilon}_m$  is multiplied to the true value of measurements  $\mathbf{f}(\boldsymbol{\beta})$ ,  $\boldsymbol{\epsilon}_m$  has been naturally called *multiplicative errors*. As in the case of the additive error model (1),  $\boldsymbol{\epsilon}_a$  in (5) is additive. Accordingly,  $\boldsymbol{\epsilon}_a$  is called *additive errors*. To illustrate additive and multiplicative random errors, we simulate the random errors of baselines, with the constants  $a$  and  $b$  set to 0.05 m and 10 ppm, respectively. The generated



**Fig. 1** Illustration of the additive and multiplicative random errors of baselines: the *upper panel* – the additive errors; the *middle panel* – the multiplicative errors; and the *lower panel* – the mixed additive and multiplicative errors

errors with the lengths of baselines are illustrated in Fig. 1. It is obvious from the simulated random errors in the upper panel of Fig. 1 that the additive random errors uniformly scatter over different lengths of baselines. The multiplicative errors in the middle panel of the same figure show a clear trend of fan shape, with the amplitudes of errors increasing with the increase of lengths of baselines.

Assuming that  $\epsilon_m$  and  $\epsilon_a$  are statistically independent and applying the error propagation law to each of the measurements  $\mathbf{y}$ , we have

$$\sigma_{y_i}^2 = f_i^2(\boldsymbol{\beta})\sigma_{m_i}^2 + \sigma_{a_i}^2, \quad (6)$$

where  $\sigma_{y_i}^2$  is the variance of the  $i$ th measurement of  $\mathbf{y}$ , and  $\sigma_{m_i}^2$  and  $\sigma_{a_i}^2$  are the  $i$ th diagonal elements of  $\boldsymbol{\Sigma}_m$  and  $\boldsymbol{\Sigma}_a$ , respectively. It is obvious from (6) that the larger the true value of measurement  $f_i(\boldsymbol{\beta})$ , the noisier the corresponding measurement  $y_i$ . When applying the same error propagation law to the measurement vector  $\mathbf{y}$ , we can obtain the variance-covariance matrix of the measurements  $\mathbf{y}$  as follows:

$$\boldsymbol{\Sigma}_y(\boldsymbol{\beta}) = \mathbf{D}_{f\boldsymbol{\beta}}\boldsymbol{\Sigma}_m\mathbf{D}_{f\boldsymbol{\beta}} + \boldsymbol{\Sigma}_a, \quad (7)$$

where  $\mathbf{D}_{f\boldsymbol{\beta}}$  is a diagonal matrix with its  $i$ th diagonal element being equal to  $f_i(\boldsymbol{\beta})$ . The elements of  $\boldsymbol{\Sigma}_y(\boldsymbol{\beta})$  are obviously the functions of the parameters  $\boldsymbol{\beta}$ . For simplicity, we will use

$\boldsymbol{\Sigma}_y$  to denote the variance-covariance of  $\mathbf{y}$ . If necessary, we can also readily take the correlation between  $\epsilon_m$  and  $\epsilon_a$  into account, which will not be discussed in this paper, however.

If  $\mathbf{f}(\boldsymbol{\beta})$  is linear, then the error model (5) becomes

$$\left. \begin{aligned} \mathbf{y} &= (\mathbf{A}\boldsymbol{\beta}) \odot (\mathbf{1} + \boldsymbol{\epsilon}_m) + \boldsymbol{\epsilon}_a \\ E(\mathbf{y}) &= \mathbf{A}\boldsymbol{\beta} \\ E(\boldsymbol{\epsilon}_m) &= \mathbf{0}, \quad D(\boldsymbol{\epsilon}_m) = \boldsymbol{\Sigma}_m \\ E(\boldsymbol{\epsilon}_a) &= \mathbf{0}, \quad D(\boldsymbol{\epsilon}_a) = \boldsymbol{\Sigma}_a \end{aligned} \right\}, \quad (8)$$

where  $\mathbf{A}$  is a given design matrix, which will be assumed to be of full rank. The model (5) will be called mixed additive and multiplicative error models in the remainder of this paper. In accordance with (8), we can rewrite  $\mathbf{D}_{f\boldsymbol{\beta}}$  as  $\mathbf{D}_{a\boldsymbol{\beta}}$ , whose diagonal elements are equal to  $\mathbf{a}_i\boldsymbol{\beta}$ , where  $\mathbf{a}_i$  is the  $i$ th row of the matrix  $\mathbf{A}$ .

### 3 The Quasi-Likelihood Method

The quasi-likelihood method was first proposed by Wedderburn (1974). It has since become a statistical method to estimate the parameters in the model of type (4) and widely applied in many areas of science and engineering. Actually, the multiplicative error model (4) has been better known in statistics as the generalized linear model and

well documented in statistical books (see e.g. McCullagh and Nelder 1989; Heyde 1997, chapter 5.3), if the function of signal  $s_{ij}$  can be represented linearly by a number of unknown parameters  $\beta$ .

Wedderburn (1974) started with a set of independent measurements  $y_i$  ( $i = 1, 2, \dots, n$ ), with expectations  $\bar{y}_i$  and variances  $\sigma_i^2(\bar{y}_i)$ , and then defined the quasi-likelihood function  $\text{QLF}(y_i, \bar{y}_i)$  as follows:

$$\frac{\partial \text{QLF}(y_i, \bar{y}_i)}{\partial \bar{y}_i} = \frac{y_i - \bar{y}_i}{\sigma_i^2(\bar{y}_i)}, \quad (9)$$

where the variance of each  $y_i$  is assumed to be the function of  $\bar{y}_i$ . By letting the expression (9) equal zero over all the measurements  $y_i$ , Wedderburn (1974) was then able to estimate the unknown parameters from the measurements  $\mathbf{y}$ . If  $\bar{y}_i$  can further be represented linearly by a number of unknown parameters  $\beta$  and if the measurements  $\mathbf{y}$  are assumed to be correlated, then (9) can be rewritten as follows:

$$\frac{\partial \text{QLF}(\beta)}{\partial \beta} = \mathbf{A}^T \Sigma_y^{-1}(\beta)(\mathbf{y} - \mathbf{A}\beta), \quad (10)$$

where  $\Sigma_y(\beta)$  is the variance-covariance matrix of  $\mathbf{y}$  whose elements are all the functions of the unknown parameters  $\beta$ . The quasi-likelihood function is proved to be equal to the maximum likelihood function, if the distribution of  $y_i$  is exponential. In general, quasi-likelihood is different from maximum likelihood, however.

Although Wedderburn (1974) defined the quasi-likelihood function  $\text{QLF}(y_i, \bar{y}_i)$  through the differential equation (9),  $\text{QLF}(y_i, \bar{y}_i)$  is really not required for parameter estimation. Actually, all what we need for parameter estimation is the expression on the right hand side of (9), which is completely defined by  $y_i$ , its expectation  $\bar{y}_i$  and its variance  $\sigma_i^2(\bar{y}_i)$ . When the quasi-likelihood method is applied to the mixed additive and multiplicative error model (8), we have the system of normal equations:

$$\mathbf{A}^T \Sigma_y^{-1}(\hat{\beta}_{ql})(\mathbf{y} - \mathbf{A}\hat{\beta}_{ql}) = \mathbf{0}, \quad (11)$$

where  $\hat{\beta}_{ql}$  stands for the quasi-likelihood estimate of  $\beta$ . Obviously, the system of normal equations (11) is nonlinear and can generally be solved by using numerical methods. Very often, one can use the Gauss-Newton method to find the solution to (11). The quasi-likelihood estimator  $\hat{\beta}_{ql}$  is asymptotically unbiased (see e.g., McCullagh 1983) and its variance-covariance matrix, denoted by  $D(\hat{\beta}_{ql})$ , is then given approximately by

$$D(\hat{\beta}_{ql}) = (\mathbf{A}^T \Sigma_y^{-1} \mathbf{A})^{-1}. \quad (12)$$

It is seen from (11) that the equation system (11) has completely defined the quasi-likelihood estimator  $\hat{\beta}_{ql}$ , no matter whether we can or cannot solve for the quasi-likelihood function  $\text{QLF}(\mathbf{y}, \beta)$  through the differential equation of type (9). The system of normal equations clearly indicates that an estimator can simply be constructed through a system of equations. As a result of this, the system of equations like (11) has been called *generalized estimating equations* (see e.g., Crowder 1995; Desmond 1997; Heyde 1997; Kukusha et al. 2010; Fitzmauric 1995).

## 4 Least-Squares-Based Methods

Although quasi-likelihood has become a standard method for parameter estimation in multiplicative error models, its associated quasi-likelihood function may hardly be derived for a general nonlinear function  $\mathbf{f}(\beta)$ . Even if such a function can indeed be found by solving the corresponding differential equations, it may not be connected with any physically meaningful distribution function. As a result, Xu and Shimada (2000) alternatively proposed LS-based methods to estimate the unknown parameters in the mixed additive and multiplicative error model (8). In this section, we will briefly discuss the ordinary LS, the weighted LS and bias-corrected weighted LS methods. If the reader is interested in the error analysis of adjusted measurements and the corrections of measurements and/or the estimation of the variance of unit weight in multiplicative error models, she/he is referred to Shi et al. (2014).

### 4.1 The Ordinary LS Method

When applying the ordinary LS method to estimate the unknown parameters  $\beta$  in the model (8), we have the following optimization objective function:

$$\min : F_1(\beta) = (\mathbf{y} - \mathbf{A}\beta)^T (\mathbf{y} - \mathbf{A}\beta). \quad (13)$$

The solution to (13) is the ordinary LS estimate of  $\beta$ , which is denoted by  $\hat{\beta}_{LS}$  and given by

$$\hat{\beta}_{LS} = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{y}. \quad (14)$$

The variance-covariance matrix of  $\hat{\beta}_{LS}$  is then given as follows:

$$D(\hat{\beta}_{LS}) = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \Sigma_y \mathbf{A} (\mathbf{A}^T \mathbf{A})^{-1}, \quad (15)$$

where  $\Sigma_y$  is the variance-covariance matrix of the measurements  $\mathbf{y}$ .



If we assume that the signal within a small area of pixels in a coherent image with multiplicative noises is identical, namely,  $s_{ij}$  remains unchanged in such a small area, then all the corresponding measurements  $y_{ij}$  are of the same variances. In other words, the weights of measurements  $y_{ij}$  are all equal to each other. As a result, the estimate of  $s_{ij}$  is simply equal to the mean value of all  $y_{ij}$  in the area. Actually, it is exactly the local mean filter for de-noising images contaminated by multiplicative noises.

## 4.2 The Weighted LS Method

When applying the weighted LS method to the model (8), we have the following minimization problem:

$$\min : F_2(\boldsymbol{\beta}) = (\mathbf{y} - \mathbf{A}\boldsymbol{\beta})^T \boldsymbol{\Sigma}_y^{-1} (\mathbf{y} - \mathbf{A}\boldsymbol{\beta}). \quad (16)$$

To derive the weighted LS estimate of  $\boldsymbol{\beta}$ , we can differentiate  $F_2(\boldsymbol{\beta})$  of (16) with respect to  $\boldsymbol{\beta}$  and let it be equal to zero. After some lengthy derivations, we can finally obtain the system of normal equations as follows:

$$(\mathbf{A}^T \hat{\boldsymbol{\Sigma}}_y^{-1} \mathbf{A}) \hat{\boldsymbol{\beta}} - \mathbf{A}^T \hat{\boldsymbol{\Sigma}}_y^{-1} \mathbf{y} - \mathbf{G}_1 (\mathbf{A} \hat{\boldsymbol{\beta}} - \mathbf{y}) = \mathbf{0}, \quad (17)$$

where

$$\mathbf{G}_1 = \begin{bmatrix} (\mathbf{A} \hat{\boldsymbol{\beta}} - \mathbf{y})^T \hat{\boldsymbol{\Sigma}}_y^{-1} \mathbf{D}_{ae_1} \boldsymbol{\Sigma}_m \hat{\mathbf{D}}_{a\hat{\boldsymbol{\beta}}} \hat{\boldsymbol{\Sigma}}_y^{-1} \\ (\mathbf{A} \hat{\boldsymbol{\beta}} - \mathbf{y})^T \hat{\boldsymbol{\Sigma}}_y^{-1} \mathbf{D}_{ae_2} \boldsymbol{\Sigma}_m \hat{\mathbf{D}}_{a\hat{\boldsymbol{\beta}}} \hat{\boldsymbol{\Sigma}}_y^{-1} \\ \vdots \\ (\mathbf{A} \hat{\boldsymbol{\beta}} - \mathbf{y})^T \hat{\boldsymbol{\Sigma}}_y^{-1} \mathbf{D}_{ae_t} \boldsymbol{\Sigma}_m \hat{\mathbf{D}}_{a\hat{\boldsymbol{\beta}}} \hat{\boldsymbol{\Sigma}}_y^{-1} \end{bmatrix},$$

$\mathbf{D}_{ae_i}$  is a diagonal matrix with its  $k$ th diagonal element being equal to  $(\mathbf{a}_k \mathbf{e}_i)$ ,  $\mathbf{e}_i$  is the  $i$ th natural basis vector of dimension  $t$ ,  $\hat{\mathbf{D}}_{a\hat{\boldsymbol{\beta}}}$  is the estimate of  $\mathbf{D}_{a\boldsymbol{\beta}}$  by replacing  $\boldsymbol{\beta}$  with its corresponding weighted LS estimate  $\hat{\boldsymbol{\beta}}$ . Following Xu et al. (2013), we can solve for the weighted LS estimate of  $\boldsymbol{\beta}$  through the following iteration procedures:

$$\hat{\boldsymbol{\beta}}_{k+1} = (\mathbf{A}^T \hat{\boldsymbol{\Sigma}}_{y_k}^{-1} \mathbf{A})^{-1} \{ \mathbf{A}^T \hat{\boldsymbol{\Sigma}}_{y_k}^{-1} \mathbf{y} + \mathbf{G}_{1k} (\mathbf{A} \hat{\boldsymbol{\beta}}_k - \mathbf{y}) \}, \quad k = 0, 1, \dots \quad (18)$$

where  $\hat{\boldsymbol{\Sigma}}_{y_k}$  and  $\mathbf{G}_{1k}$  stand for computing  $\hat{\boldsymbol{\Sigma}}_y$  and  $\mathbf{G}_1$  at the point of  $\hat{\boldsymbol{\beta}}_k$ .

It is obvious from (17) that the weighted LS estimate  $\hat{\boldsymbol{\beta}}$  is nonlinear and is expected to be biased. Xu et al. (2013) derived the bias of  $\hat{\boldsymbol{\beta}}$  in the mixed additive and multiplicative error model (8), which is denoted by  $\mathbf{b}(\hat{\boldsymbol{\beta}})$  and is simply listed as follows:

$$\mathbf{b}(\hat{\boldsymbol{\beta}}) = E(\mathbf{b}(\boldsymbol{\beta})) = \mathbf{N}^{-1} \mathbf{g}_2, \quad (19)$$

where  $\mathbf{N} = \mathbf{A}^T \boldsymbol{\Sigma}_y^{-1} \mathbf{A}$  and  $\mathbf{g}_2$  is given by

$$\mathbf{g}_2 = \begin{bmatrix} \text{tr}\{\mathbf{D}_{ae_1} \boldsymbol{\Sigma}_m \mathbf{D}_{a\boldsymbol{\beta}} \boldsymbol{\Sigma}_y^{-1}\} \\ \text{tr}\{\mathbf{D}_{ae_2} \boldsymbol{\Sigma}_m \mathbf{D}_{a\boldsymbol{\beta}} \boldsymbol{\Sigma}_y^{-1}\} \\ \vdots \\ \text{tr}\{\mathbf{D}_{ae_t} \boldsymbol{\Sigma}_m \mathbf{D}_{a\boldsymbol{\beta}} \boldsymbol{\Sigma}_y^{-1}\} \end{bmatrix}.$$

By limiting themselves to the linear term of  $\hat{\boldsymbol{\beta}}$  with respect to the random errors  $\boldsymbol{\epsilon}_m$  and  $\boldsymbol{\epsilon}_a$ , Xu et al. (2013) also derived the first order approximation of the variance-covariance matrix of the weighted LS estimate  $\hat{\boldsymbol{\beta}}$ , which is denoted by  $D_1(\hat{\boldsymbol{\beta}})$  and given as follows:

$$D_1(\hat{\boldsymbol{\beta}}) = (\mathbf{A}^T \boldsymbol{\Sigma}_y^{-1} \mathbf{A})^{-1}. \quad (20)$$

After taking the bias (19) into account, we obtain the approximate mean squared error (MSE) matrix of  $\hat{\boldsymbol{\beta}}$  as follows:

$$\begin{aligned} \mathbf{M}(\hat{\boldsymbol{\beta}}) &= D_1(\hat{\boldsymbol{\beta}}) + \mathbf{b}(\hat{\boldsymbol{\beta}}) [\mathbf{b}(\hat{\boldsymbol{\beta}})]^T \\ &= (\mathbf{A}^T \boldsymbol{\Sigma}_y^{-1} \mathbf{A})^{-1} + (\mathbf{A}^T \boldsymbol{\Sigma}_y^{-1} \mathbf{A})^{-1} \mathbf{g}_2 \mathbf{g}_2^T (\mathbf{A}^T \boldsymbol{\Sigma}_y^{-1} \mathbf{A})^{-1}, \end{aligned} \quad (21)$$

where  $\mathbf{M}(\hat{\boldsymbol{\beta}})$  stands for the MSE matrix of  $\hat{\boldsymbol{\beta}}$ .

## 4.3 The Bias-Corrected Weighted LS Method

Bias analysis in Xu and Shimada (2000) and Xu et al. (2013) clearly indicates that the bias of the weighted LS estimate is solely caused by the non-zero term of derivatives of the variance-covariance matrix  $\boldsymbol{\Sigma}_y$  with respect to  $\boldsymbol{\beta}$ . Thus both works propose deleting the corresponding term in the normal equations, namely the third term on the right hand side of the normal equations (17), to remove the bias from the weighted LS estimate  $\hat{\boldsymbol{\beta}}$ . As a result, they are able to construct the bias-corrected weighted LS estimate of  $\boldsymbol{\beta}$ .

When the same idea is applied to the mixed additive and multiplicative error model (8), they derive the bias-corrected weighted LS estimate of  $\boldsymbol{\beta}$ , which is denoted by  $\hat{\boldsymbol{\beta}}_{bc}$  and solved through the following system of normal equations:

$$(\mathbf{A}^T \hat{\boldsymbol{\Sigma}}_y^{-1} \mathbf{A}) \hat{\boldsymbol{\beta}}_{bc} - \mathbf{A}^T \hat{\boldsymbol{\Sigma}}_y^{-1} \mathbf{y} = \mathbf{0}, \quad (22)$$

where  $\hat{\boldsymbol{\beta}}_{bc}$  is the bias-corrected WLS estimate of  $\boldsymbol{\beta}$ . Equivalently,  $\hat{\boldsymbol{\beta}}_{bc}$  can be formally rewritten as follows:

$$\hat{\boldsymbol{\beta}}_{bc} = (\mathbf{A}^T \hat{\boldsymbol{\Sigma}}_y^{-1} \mathbf{A})^{-1} \mathbf{A}^T \hat{\boldsymbol{\Sigma}}_y^{-1} \mathbf{y}, \quad (23)$$

which is unbiased up to the second order approximation (Xu et al. 2013). The variance-covariance matrix of  $\hat{\beta}_{bc}$  is denoted by  $D(\hat{\beta}_{bc})$  and given by

$$D(\hat{\beta}_{bc}) = (\mathbf{A}^T \hat{\Sigma}_y^{-1} \mathbf{A})^{-1}. \quad (24)$$

Because the matrix  $\hat{\Sigma}_y$  depends on  $\hat{\beta}_{bc}$ , (23) is actually a nonlinear system of equations and can, in general, be solved numerically. If the Gauss-Newton method is used to solve for the bias-corrected weighted LS estimate, we have the following iterative formula:

$$\hat{\beta}_{bc}^{k+1} = \hat{\beta}_{bc}^k - (\mathbf{A}^T \hat{\Sigma}_{yk}^{-1} \mathbf{A})^{-1} \mathbf{A}^T \hat{\Sigma}_{yk}^{-1} (\mathbf{A} \hat{\beta}_{bc}^k - \mathbf{y}), \quad k = 0, 1, \dots \quad (25)$$

(see also McCullagh 1983; McCullagh and Nelder 1989; Dennis and Schnabel 1996; Xu et al. 2013).

We should like to note that given some approximate values, say  $\beta_0$ ,  $\Sigma_{m0}$  and  $\Sigma_{a0}$ , we can then replace  $\hat{\Sigma}_y$  of (22) with  $\Sigma_{y0}$  (computed at  $\beta_0$ ,  $\Sigma_{m0}$  and  $\Sigma_{a0}$ ), which is exactly the conventional practice of adjustment of geodetic networks such as EDM, VLBI and GPS baselines. In other words, the conventional weighted LS adjustment of baseline networks can be interpreted as a special case of the bias-corrected weighted LS method with given approximate values. Nevertheless, the effectiveness of using approximate values would depend on how far away these approximate values deviate from their true values, as also pointed out and demonstrated in Xu et al. (2013).

## 5 Mixed Additive and Multiplicative Random Error Models with Prior Information

In this section, we will extend the parameter estimation in mixed additive and multiplicative random error models to the case with prior information. Prior information will only be assumed in the form of the first two moments on the unknown parameters  $\beta$ , i.e. its prior mean  $\mu$  and prior variance-covariance matrix. Bearing the concept of additive and multiplicative random errors in mind, we will accordingly assume two types of prior variance-covariance matrices. As in the case of Gauss-Markov models with additive random errors, the first type of prior variance-covariance matrices is assumed to be independent of  $\beta$  and symbolically denoted by  $\Sigma_0$ . If prior information on  $\beta$  is obtained from measurements contaminated by multiplicative errors other than the measurements  $\mathbf{y}$  in the mixed additive and multiplicative error model (8), then the prior variance-covariance matrix will surely be dependent on  $\beta$ , as can be readily seen in (24), for example. Thus, the second type of prior variance-covariance matrices is assumed to be the functions of  $\beta$  and

denoted by  $\Sigma_\beta$ . Of course, prior information can also be presented in the form of prior distributions. In this case, one can use Bayesian inference to estimate the unknown parameters. For more information, the reader is referred to Xu (1999).

If the model (8) is combined with the first type of prior variance-covariance matrices, namely,  $\Sigma_0$ , then the corresponding generalized (weighted) LS objective function will become

$$\min : F_3(\beta) = (\mathbf{y} - \mathbf{A}\beta)^T \Sigma_y^{-1} (\mathbf{y} - \mathbf{A}\beta) + (\beta - \mu)^T \Sigma_0^{-1} (\beta - \mu). \quad (26)$$

According to the bias analysis in Xu and Shimada (2000) and Xu et al. (2013), we know that  $\Sigma_y$  will create a bias in the solution to (26). Since the prior variance-covariance matrix  $\Sigma_0$  is independent of  $\beta$ , it will not contribute extra terms to the bias of the solution. By following the same rationale as in Sect. 4.3, we can ignore the dependence of  $\Sigma_y$  on  $\beta$  as if it were independent of  $\beta$  and, as a result, readily construct the bias-corrected estimator of  $\beta$  with prior information, denoted by  $\hat{\beta}_{bc}^{p0}$ , as follows:

$$\hat{\beta}_{bc}^{p0} = (\mathbf{A}^T \hat{\Sigma}_y^{-1} \mathbf{A} + \Sigma_0^{-1})^{-1} (\mathbf{A}^T \hat{\Sigma}_y^{-1} \mathbf{y} + \Sigma_0^{-1} \mu). \quad (27)$$

The first order accuracy of  $\hat{\beta}_{bc}^{p0}$  is then given by

$$D(\hat{\beta}_{bc}^{p0}) = (\mathbf{A}^T \Sigma_y^{-1} \mathbf{A} + \Sigma_0^{-1})^{-1}. \quad (28)$$

If the second type of prior information is combined with the measurements from the model (8), we should then have the following optimization problem:

$$\min : F_4(\beta) = (\mathbf{y} - \mathbf{A}\beta)^T \Sigma_y^{-1} (\mathbf{y} - \mathbf{A}\beta) + (\beta - \mu)^T \Sigma_\beta^{-1} (\beta - \mu). \quad (29)$$

Obviously, both of  $\Sigma_y$  and  $\Sigma_\beta$  will now directly contribute terms to the bias of the optimal solution to the optimization problem (29), according to Xu and Shimada (2000) and Xu et al. (2013). In the similar manner to (27), we can construct the bias-corrected estimator of  $\beta$  with prior mean  $\mu$  and prior variance-covariance matrix  $\Sigma_\beta$ , which is denoted by  $\hat{\beta}_{bc}^{p1}$  and simply listed as follows:

$$\hat{\beta}_{bc}^{p1} = (\mathbf{A}^T \hat{\Sigma}_y^{-1} \mathbf{A} + \hat{\Sigma}_\beta^{-1})^{-1} (\mathbf{A}^T \hat{\Sigma}_y^{-1} \mathbf{y} + \hat{\Sigma}_\beta^{-1} \mu), \quad (30)$$

which is unbiased up to the second order approximation, and its first order accuracy is given by

$$D(\hat{\beta}_{bc}^{p1}) = (\mathbf{A}^T \Sigma_y^{-1} \mathbf{A} + \Sigma_\beta^{-1})^{-1}. \quad (31)$$

As in the case of the bias-corrected weighted LS estimation, the bias-corrected stochastic inference (or LS

collocation) with prior information in mixed multiplicative and additive error models is essentially of the same form as in the case of purely additive error models, as correctly pointed out by one of the reviewers. However, unlike stochastic inference in additive error models, the bias-corrected stochastic inference with prior information in mixed multiplicative and additive error models requires that  $\hat{\beta}_{bc}^{pl}$  of (30) be computed iteratively, since the right hand side of (30) contains the unknowns  $\hat{\beta}_{bc}^{pl}$  as well. Nevertheless, if one would simply apply the conventional principle of stochastic inference to mixed multiplicative and additive error models, one would end up with a biased estimator, which would not be of the same form as in the case of purely additive error models.

In the one-dimensional case, namely,

$$y_{ij} = (1 + \varepsilon_{mij})s_{ij} + \varepsilon_{aij},$$

then the bias-corrected LS estimate of  $s_{ij}$  with prior information can be rewritten as follows:

$$\hat{s}_{ij} = \mu_{ij} + \frac{\sigma_{\mu}^2}{\sigma_m^2 s_{0ij}^2 + \sigma_a^2 + \sigma_{\mu}^2} (y_{ij} - \mu_{ij}), \quad (32)$$

where  $\mu_{ij}$  is the prior mean of  $s_{ij}$ ,  $\sigma_{\mu}^2$  is the prior variance of  $\mu_{ij}$ ,  $s_{0ij}$  is some approximate value of  $s_{ij}$ , and  $\sigma_m^2$  and  $\sigma_a^2$  are the variances of the multiplicative and additive errors  $\varepsilon_{mij}$  and  $\varepsilon_{aij}$ , respectively. By properly choosing the values of  $\mu_{ij}$ ,  $\sigma_{\mu}^2$ ,  $s_{0ij}$ ,  $\sigma_m^2$  and  $\sigma_a^2$ , one can then construct the filter by Kuan et al. (1985) for image de-noising.

## 6 Concluding Remarks

Geodetic adjustment has been developed on the basis of Gauss-Markov models with additive random errors. The most important feature of such a Gauss-Markov model with additive random errors is that the accuracy of a measurement has nothing to do with the true value of the measurement. However, geodetic practice has clearly demonstrated that random errors of EDM, VLBI and GPS baselines indeed change with the length of a baseline. In other words, random errors of such types usually consist of two parts: one behaves more or less constant and may reflect only random effects of local nature, while the other is proportional to the length of the baseline and could, very likely, reflect total effect of the propagation path between the two stations. Such error characteristics are part of modern geodetic coherent imaging systems such as SAR and LiDAR. Obviously, the conventional adjustment theory that has been developed on

the assumption of additive random errors cannot theoretically meet the need to process geodetic measurements that are contaminated by multiplicative and/or mixed additive and multiplicative random errors.

In this paper, we have briefly reviewed two types of methods for parameter estimation in mixed additive and multiplicative error models, namely, quasi-likelihood and least-squares-based methods, with or without prior information. Quasi-likelihood, though theoretically connected with distributions, can be used directly for parameter estimation without any assumption on distributions. If there exist multiple solutions to the generalized estimating equations, no criterion is available for quasi-likelihood to pick up the right solution, however. The LS-based methods have a clearly defined objective function. Thus the sense of optimality of LS-based estimates is well defined. For the linear model (8), quasi-likelihood, the ordinary and bias-corrected weighted LS methods can all warrant an unbiased estimate of the unknown parameters, while the weighted LS method will generally lead to a biased estimate. Quasi-likelihood and the bias-corrected weighted LS method are more efficient than the ordinary LS method. We have also extended the bias-corrected LS estimate to the case with prior information, which can either be given in the form of prior mean and a parameter-free or a parameter-dependent prior variance-covariance matrix.

**Acknowledgements** This work is partially supported by the National Foundation of Natural Science of China (Nos.41204006, 41374016) and the project SKLGED2013-4-8-E, and the Grant-in-Aid for Scientific Research (C25400449). The authors thank the reviewers and the editor very much for their constructive comments, which help clarify some points of the paper.

## References

- Aubert G, Aujol J-F (2008) A variational approach to removing multiplicative noise. *SIAM J Appl Math* 68:925–946
- Crowder M (1995) On the use of a working correlation matrix in using generalised linear models for repeated measures. *Biometrika* 82:407–410
- Dennis Jr. JE, Schnabel RB (1996) Numerical methods for unconstrained optimization and nonlinear equations. *SIAM classics in applied mathematics*. SIAM, Philadelphia
- Desmond AF (1997) Optimal estimating functions, quasi-likelihood and statistical modelling. *J Stat Plan Inference* 60:77–123
- Ewing CE, Mitchell MM (1970) Introduction to geodesy. Elsevier, New York
- Fitzmauric GM (1995) A caveat concerning independence estimating equations with multivariate binary data. *Biometrics* 51:309–317
- Flamant PH, Menzies RT, Kavaya MJ (1984) Evidence for speckle effects on pulsed CO<sub>2</sub> lidar signal returns from remote targets. *Appl Optics* 23:1412–1417
- Goodman JW (1976) Some fundamental properties of speckle. *J Opt Soc Am* 66:1145–1150
- Heyde CC (1997) Quasi-likelihood and its applications. Springer, New York

- Hill AC, Harris M, Ridley KC, Jakeman E, Lutzmann P (2003) Lidar frequency modulation vibrometry in the presence of speckle. *Appl Optics* 42:1091–1100
- Kuan DT, Sawchuk AA, Strand TC, Chavel P (1985) Adaptive noise smoothing filter for images with signal-dependent noise. *IEEE Trans Pattern Anal Mach Intell PAMI-7*:165–177
- Kukusha A, Malenkova A, Schneeweiss H (2010) Optimality of quasi-score in the multivariate mean-variance model with an application to the zero-inflated Poisson model with measurement errors. *Statistics* 44:381–396
- López-Martínez C, Fàbregas X, Pipia L (2011) Forest parameter estimation in the Pol-InSAR context employing the multiplicative-additive speckle noise model. *ISPRS J Photogram Remote Sens* 66:597–607
- MacDoran PF (1979) Satellite emission radio interferometric earth surveying series — GPS geodetic system. *Bull Géod* 53:117–138
- McCullagh P (1983) Quasi-likelihood functions. *Ann Stat* 11:59–67
- McCullagh P, Nelder J (1989) *Generalized linear models*, 2nd edn. Chapman and Hall, London
- Oliver CJ (1991) Information from SAR images. *J Phys D Appl Phys* 24:1493–1514
- Pearlman MR, Degnan JJ, Bosworth JM (2002) The international laser ranging service. *Adv Space Res* 30:135–143
- Petrov L, Gordon D, Gipson J, MacMillan D, Ma C, Fomalont E, Walker R, Carabajal C (2010) Precise geodesy with the very long baseline array. *J Geodesy* 83:859–876
- Seeber G (2003) *Satellite geodesy*, 2nd edn. de Gruyter, Berlin
- Shi Y, Xu PL, Peng JH, Shi C, Liu JN (2014) Adjustment of measurements with multiplicative errors: error analysis, estimates of the variance of unit weight, and effect on volume estimation from LiDAR-type digital elevation models. *Sensors* 14:1249–1266. doi:10.3390/s140101249
- Swami A (1994) Multiplicative noise models: parameter estimation using cumulants. *Signal Process* 36:355–373
- Tian H, Fowler B, Gamal El (2001) Analysis of temporal noise in CMOS photodiode active pixel sensor. *IEEE J Solid State Circuits* 36:92–101
- Ulaby F, Kouyate F, Brisco B, Williams T (1986) Textural information in SAR images. *IEEE Trans Geosci Remote Sens* 24:235–245
- Wang JY, Pruitt PA (1992) Effects of speckle on the range precision of a scanning lidar. *Appl Opt* 31:801–808
- Wedderburn R (1974) Quasi-likelihood functions, generalized linear models, and the Gauss-Newton method. *Biometrika* 61:439–447
- Willis P, Fagard H, Ferrage P, Lemoine FG, Noll CE, Noomen R, Otten M, Ries JC, Rothacher M, Soudarin L, Tavernier G, Valette JJ (2010) The international DORIS service, toward maturity. In: Willis P (ed) *DORIS: scientific applications in geodesy and geodynamics*. *Adv Space Res* 45(12):1408–1420
- Xu PL (1999) Despeckling SAR-type multiplicative noise. *Int J Remote Sens* 20:2577–2596
- Xu PL, Shimada S (2000) Least squares estimation in multiplicative noise models. *Commun Stat B* 29:83–96
- Xu PL, Shi Y, Peng JH, Liu JN, Shi C (2013) Adjustment of geodetic measurements with mixed multiplicative and additive random errors. *J Geodesy* 87:629–643

---

# Cycle Slip Detection and Correction for Heading Determination with Low-Cost GPS/INS Receivers

Patrick Henkel and Naoya Oku

---

## Abstract

Precise attitude determination with low-cost GPS receivers requires integer ambiguity resolution and reliable cycle slip correction. In this paper, a tree search of cycle slips is proposed, which combines double difference GPS carrier phases from all visible satellites, gyroscope and acceleration measurements, and a priori information on the baseline length between both GPS receivers. The proposed method was verified in both a slalom drive with high vehicle dynamics and a drive below trees with shadowed GPS signals: The residuals of the fixed phase measurements were reduced to less than 15 cm throughout the measurement period.

---

## Keywords

Attitude determination • Cycle slip correction • Integer ambiguity resolution

---

## 1 Introduction

The availability of mass-market GPS receivers with carrier phase tracking has led to a wide range of new applications of RTK and attitude determination. However, the measurements of low-cost GPS receivers differ in three aspects from the measurements of geodetic receivers: First, code multipath is much larger due to the small size of the patch antennas. It can be 10m even in open-sky conditions, which is a challenge for ambiguity resolution. Secondly, half cycle slips occur much more frequently and at multiple satellites simultaneously. Today, cycle slip detection and correction can be performed reliably for *geodetic* receivers with inertial sensors: Du and Gao (2012) differenced the carrier phase measurements between two satellites and two subsequent epochs such that clock offsets, ambiguities, biases and

atmospheric delays are eliminated (except for the drift which is in general negligible for periods of less than 1 s). This leaves the change in position and the cycle slips as unknowns. The change in position is predicted by an inertial sensor. Thus, cycle slips can be determined on a satellite by satellite basis by simple rounding. Du and Gao (2012) applied a cascaded approach to two dual-frequency Novatel OEM4 receivers, i.e. the widelane cycle slips were determined first. Subsequently, the extra-widelane cycle slips were resolved and, finally, the L1 and L2 cycle slips were derived from the widelane and extra-widelane cycle slips.

Dai et al (2009) proposed a cycle slip detection and correction method for triple frequency GNSS receivers. They used two triple frequency geometry-free phase combinations and performed an integer least-squares estimation using the LAMBDA method of Teunissen (1995).

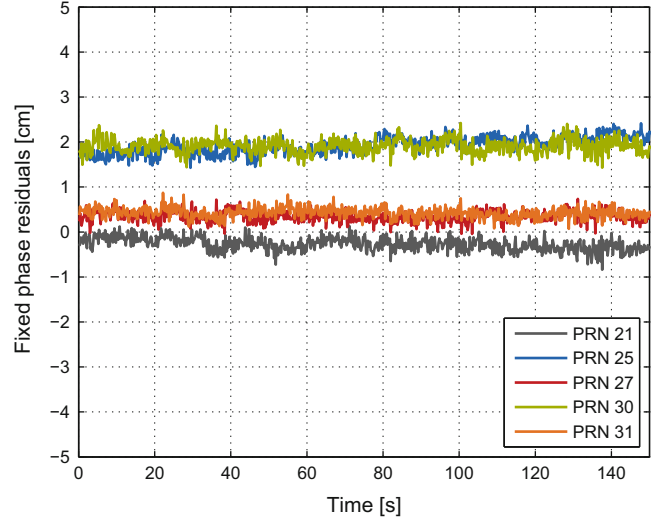
Colombo et al (1999) also proposed a cascaded cycle slip detection for dual frequency receivers. The widelane cycle slips were fixed as described by Du and Gao (2012). The individual L1 and L2 cycle slips were then directly fixed with the help of the widelane cycle slips. They used an Ashtech receiver and high-grade IMU with a gyroscope drift of only 3°/h, and were able to correct 99.1% of cycle slips for data gaps of 5 s.

---

P. Henkel (✉)  
Technische Universität München, Theresienstrasse 90, 80333  
München, Germany  
e-mail: [patrick.henkel@tum.de](mailto:patrick.henkel@tum.de)

N. Oku  
Advanced Navigation Solutions - ANAVS, Friedrichshafener Str. 1,  
82205 Gilching, Germany

**Fig. 1** Residuals of fixed phase solution during initialization: the residuals are clearly smaller than the wavelength and are drift-free, which indicates a correct integer ambiguity resolution. The biases of up to 2 cm arise from a small bias in the baseline length constraint due to antenna phase center offsets



These methods can not be used for low-cost single-frequency mass-market GNSS receivers and inertial sensors: The receivers show half cycle slips, i.e. jumps of integer multiples of  $\lambda/2 = 9.5$  cm, and substantial phase and code multipath. The measurement model must be strengthened, e.g. by including baseline a priori information as described by Henkel and Kiam (2013).

A third difference between low-cost and geodetic GNSS receivers is that the receiver clock offsets are in the order of milliseconds instead of nanoseconds (see Henkel and Gunther 2013), which implies that the satellite movement within the time span of the differential receiver clock offset is no longer negligible. A synchronization correction is required to preserve the integer property of the double difference (DD) ambiguities.

Therefore, we include the receiver clock offset explicitly in our model for the DD carrier phase measurements. As the baseline is short for attitude determination, atmospheric and orbital errors can be neglected and we obtain

$$\begin{aligned}
& (\lambda\varphi_1^k(t + \delta\tau_1) - \lambda\varphi_1^l(t + \delta\tau_1)) \\
& - (\lambda\varphi_2^k(t + \delta\tau_2) - \lambda\varphi_2^l(t + \delta\tau_2)) \\
& = \|x_1(t + \delta\tau_1) - x^k(t + \delta\tau_1 - \Delta\tau_1^k)\| \\
& - \|x_1(t + \delta\tau_1) - x^l(t + \delta\tau_1 - \Delta\tau_1^l)\| \\
& - \|x_2(t + \delta\tau_2) - x^k(t + \delta\tau_2 - \Delta\tau_2^k)\| \\
& + \|x_2(t + \delta\tau_2) - x^l(t + \delta\tau_2 - \Delta\tau_2^l)\| \\
& + \lambda N_{12}^{kl} + \lambda/2 \Delta N_{12}^{kl}(t) + m_{\varphi_{12}^{kl}}(t + \delta\tau_1, t + \delta\tau_2) \\
& + \varepsilon_{\varphi_{12}^{kl}}(t + \delta\tau_1, t + \delta\tau_2), \tag{1}
\end{aligned}$$

with the carrier wavelength  $\lambda$ , the undifferenced phase measurement  $\varphi_r^k$  of receiver  $r$  and satellite  $k$ , the receiver clock

offset  $\delta\tau_r$ , the receiver position  $x_r$ , the satellite position  $x^k$ , the DD integer ambiguity  $N_{12}^{kl}$ , the DD half cycle slip  $\Delta N_{12}^{kl}$ , the DD phase multipath  $m_{\varphi_{12}^{kl}}$  and the DD phase noise  $\varepsilon_{\varphi_{12}^{kl}}$ . The synchronization correction is given by

$$\begin{aligned}
& \lambda \Delta\varphi_{12}^{kl}(t + \delta\tau_1, t + \delta\tau_2) \tag{2} \\
& = (e_1^k(t + \delta\tau_1))^T (x_1(t + \delta\tau_1) - x^k(t + \delta\tau_1 - \Delta\tau_1^k)) \\
& - (e_1^l(t + \delta\tau_1))^T (x_1(t + \delta\tau_1) - x^l(t + \delta\tau_1 - \Delta\tau_1^l)) \\
& - (e_1^k(t + \delta\tau_2))^T (x_1(t + \delta\tau_2) - x^k(t + \delta\tau_2 - \Delta\tau_2^k)) \\
& + (e_1^l(t + \delta\tau_2))^T (x_1(t + \delta\tau_2) - x^l(t + \delta\tau_2 - \Delta\tau_2^l)).
\end{aligned}$$

We linearize the norms in (1), apply the synchronization correction of (2) and a cycle slip correction to write the linearized corrected DD phase and code measurements in matrix-vector notation as

$$\begin{pmatrix} \lambda\varphi_{12} \\ \rho_{12} \end{pmatrix} = H b_{12} + A N_{12} + \begin{pmatrix} m_{\varphi_{12}} \\ m_{\rho_{12}} \end{pmatrix} + \begin{pmatrix} \varepsilon_{\varphi_{12}} \\ \varepsilon_{\rho_{12}} \end{pmatrix} \tag{3}$$

with  $H$  and  $A$  being the implicitly defined DD geometry and ambiguity coefficient matrices, and  $b_{12}$  being the baseline between both receivers. The integer least-squares estimation of the DD integer ambiguities and baseline coordinates is improved by some a priori information  $l$  on the baseline length, which leads to the constrained integer least-squares estimation problem:

$$\begin{aligned}
& \min_{b_{12} \in \mathbb{R}^{3 \times 1}, N_{12} \in \mathbb{Z}^{K \times 1}} \left\| \begin{pmatrix} \lambda\varphi_{12} \\ \rho_{12} \end{pmatrix} - H b_{12} - A N_{12} \right\|_{\Sigma_{\varphi_{12}}^{-1}}^2 \\
& \text{s. t. } \|b_{12}\| \stackrel{!}{=} l. \tag{4}
\end{aligned}$$

A solution to this problem was developed in Teunissen (2006, 2010). Figure 1 shows the fixed phase residuals of this

estimator for DD phase measurements from two u-blox LEA 6T GPS receivers mounted on the roof of a car. The residuals are clearly smaller than the wavelength and are drift-free, which indicates a correct integer ambiguity resolution. The biases of up to 2 cm arise most likely from a bias in the baseline length constraint due to antenna phase center offsets.

## 2 Initial Calibration of Inertial Sensors

We start with a rough alignment by transforming the measured acceleration  $a^s$  and angular rotation rate  $\omega^s$  from the sensor-fixed (s-) frame to the body-fixed (b-) frame (aligned with longitudinal and transversal axis of car), i.e.

$$\begin{aligned} a_{\text{rough}}^b &= C_s^b a^s = R_1(\phi_s^b) R_2(\theta_s^b) R_3(\psi_s^b) a^s \\ \omega_{\text{rough}}^b &= C_s^b \omega^s = R_1(\phi_s^b) R_2(\theta_s^b) R_3(\psi_s^b) \omega^s, \end{aligned} \quad (5)$$

where the roll angle  $\phi_s^b$ , the pitch angle  $\theta_s^b$  and the yaw angle  $\psi_s^b$  are approximated from the mounting of the sensor on the body and  $R_i$  is the rotation around the  $i$ -th axis. Subsequently, we average  $\omega_{\text{rough}}^b$  in static conditions to determine the biases  $b_{\omega}^b$ , which are then subtracted from the measurements:

$$\omega^b = \omega_{\text{rough}}^b - b_{\omega}^b. \quad (6)$$

The acceleration measurements are also averaged over time in static conditions to reduce the noise. The obtained  $\bar{a}^b$  is expressed in terms of the Euler angles  $\phi$ ,  $\theta$ ,  $\psi$  and the gravitational acceleration  $g$ , i.e.

$$\begin{aligned} \bar{a}^b &= C_n^b \bar{a}^n \approx R_1(\phi) R_2(\theta) R_3(\psi) (0, 0, g)^T + b_{\bar{a}}^b + \eta_{\bar{a}}^b \\ &= g \cdot (-\sin \theta, \cos \theta \sin \phi, \cos \theta \cos \phi)^T + b_{\bar{a}}^b + \eta_{\bar{a}}^b, \end{aligned} \quad (7)$$

with the bias  $b_{\bar{a}}^b$  and noise  $\eta_{\bar{a}}^b$ . The biases and misalignment errors  $\Delta\theta$  and  $\Delta\phi$  between GPS and INS sensors are determined by least-squares estimation once the car is moving, i.e.

$$\min_{\Delta\theta, \Delta\phi, b_{\bar{a}}^b} \left\| \bar{a}^b - g \cdot \begin{pmatrix} -\sin(\theta + \Delta\theta) \\ \cos(\theta + \Delta\theta) \sin(\phi + \Delta\phi) \\ \cos(\theta + \Delta\theta) \cos(\phi + \Delta\phi) \end{pmatrix} - b_{\bar{a}}^b \right\|^2 \quad (8)$$

The roll and pitch angles can then be derived from (7) without the need of knowing  $g$  as

$$\begin{aligned} \phi &= \text{atan} \left( (\bar{a}_y^b - b_{\bar{a}_y}^b) / (\bar{a}_z^b - b_{\bar{a}_z}^b) \right) \\ \theta &= \text{atan} \left( -(\bar{a}_x^b - b_{\bar{a}_x}^b) / \sqrt{(\bar{a}_y^b - b_{\bar{a}_y}^b)^2 + (\bar{a}_z^b - b_{\bar{a}_z}^b)^2} \right). \end{aligned} \quad (9)$$

We initialize the heading from the GPS fixed solution. Once the Euler angles and a rough estimate of the absolute position (longitude  $\lambda$ , latitude  $\varphi$ ) is available from GPS, the coordinate frame transformation from the b-frame to the ECEF e-frame is determined as

$$C_b^e = C_n^e C_b^n \quad (10)$$

with

$$C_b^n = \begin{pmatrix} -\sin \varphi \cos \lambda & -\sin \lambda & -\cos \varphi \cos \lambda \\ -\sin \varphi \sin \lambda & \cos \lambda & -\cos \varphi \sin \lambda \\ \cos \varphi & 0 & -\sin \varphi \end{pmatrix} \quad (11)$$

and  $C_b^n = (C_n^b)^{-1} = (R_1(\phi) R_2(\theta) R_3(\psi))^{-1}$ . The rotation matrix  $C_e^b = (C_b^e)^{-1}$  is then transformed to a Quaternion as described in Jekeli (2001):

$$q = \frac{1}{\|[q_a, q_b, q_c, q_d]\|} \cdot [q_a, q_b, q_c, q_d]^T \quad (12)$$

with the four quaternion elements  $q_a$ ,  $q_b$ ,  $q_c$  and  $q_d$ .

## 3 Integration of Orientation with Quaternions

Jekeli (2001) derived the time-derivative of  $C_e^b$  as

$$\dot{C}_e^b = C_e^b \Omega_{be}^e, \quad (13)$$

which represents a differential equation with unknown  $C_e^b$ . The skew-symmetric matrix  $\Omega_{be}^e$  is given by

$$\Omega_{be}^e = \begin{pmatrix} 0 & -\omega_3 & \omega_2 \\ \omega_3 & 0 & -\omega_1 \\ -\omega_2 & \omega_1 & 0 \end{pmatrix}, \quad (14)$$

where the angular rotation rates  $\omega_i$  of the e-frame w.r.t. the b-frame are obtained from (6) by subtracting the earth rotation rate, i.e.

$$(\omega_1, \omega_2, \omega_3)^T = \omega^b - C_e^b \cdot (0, 0, \omega_E)^T =: \omega_{be}^b. \quad (15)$$

The differential equation of (13) shall be solved with Quaternions. Jekeli (2001) transformed the  $3 \times 3$  matrix equation of (13) to the  $4 \times 1$  vector equation

$$\dot{q} = \frac{1}{2} A_q q, \quad (16)$$

with the quaternion  $q$  and the matrix of angular velocities  $A_q$ . The latter one is given by

$$A_q = \begin{pmatrix} 0 & \omega_1 & \omega_2 & \omega_3 \\ -\omega_1 & 0 & \omega_3 & -\omega_2 \\ -\omega_2 & -\omega_3 & 0 & \omega_1 \\ -\omega_3 & \omega_2 & -\omega_1 & 0 \end{pmatrix}. \quad (17)$$

We performed the integration of the Quaternion with the third order Runge-Kutta method (see Jekeli 2001), i.e. the Quaternion at time  $t_{n+1}$  is given by

$$\begin{aligned} q(t_{n+1}) &= q(t_n + h) \\ &= q(t_n) + h \cdot \left( \frac{1}{6}\delta q_0 + \frac{2}{3}\delta q_1 + \frac{1}{6}\delta q_2 \right), \end{aligned} \quad (18)$$

where  $h = 2\delta t$  denotes the integration time and  $\delta q_0$ ,  $\delta q_1$  and  $\delta q_2$  denote the coefficients given by

$$\begin{aligned} \delta q_0 &= \frac{1}{2}A_q(t_n)q(t_n) \\ \delta q_1 &= \frac{1}{2}A_q(t_n + \frac{h}{2})(q(t_n) + h/2\delta q_0) \\ \delta q_2 &= \frac{1}{2}A_q(t_n + h)(q(t_n) - h\delta q_0 + 2h\delta q_1). \end{aligned} \quad (19)$$

Once the integrated quaternion is determined, it is transformed back to a rotation matrix.

## 4 Cycle Slip Detection and Correction

The fixed double DD carrier phase measurement for satellite pair  $\{k, l\}$  is modeled as

$$\lambda(\varphi_{12}^{kl} - \check{N}_{12}^{kl}) = e^{kl}b_{12} + \lambda/2\Delta N_{12}^{kl} + \varepsilon_{12}^{kl}, \quad (20)$$

with the single difference  $e^{kl}$  between the unit vectors pointing from the satellites to the receiver, the baseline vector  $b_{12}$  between both receivers, the carrier wavelength  $\lambda$ , the fixed DD integer ambiguity  $N_{12}^{kl}$ , the unknown DD cycle slip (CS)  $\Delta N_{12}^{kl}$  and the DD phase noise  $\varepsilon_{12}^{kl}$ . Solving (20) for  $\Delta N_{12}^{kl}$  yields

$$\Delta \check{N}_{12}^{kl} = \left\lceil \frac{1}{\lambda/2} \left( \lambda(\varphi_{12}^{kl} - \check{N}_{12}^{kl}) - e^{kl}\hat{b}_{12}^{\text{IMU}} \right) \right\rceil, \quad (21)$$

where  $\lceil \cdot \rceil$  denotes the rounding operator and  $\hat{b}_{12}^{\text{IMU}}$  denotes the baseline estimate of the IMU that was obtained in four steps: First, the IMU was initialized with the validated GPS attitude solution of the last GPS measurement epoch. Subsequently, the orientation was integrated using IMU measurements between the epochs of the last and current GPS measurement. The integration was performed with quaternions as described

in the previous section. In a third step, the quaternions were transformed to Euler angles and, finally, the horizontal baseline estimate was obtained from the heading  $\psi_{\text{IMU}}$  and baseline length a priori information  $l$  as

$$\hat{b}_{12}^{\text{IMU}} = l \cdot (\sin(\psi_{\text{IMU}}), \cos(\psi_{\text{IMU}}))^T. \quad (22)$$

The cycle slip correction of (21) might be erroneous if the initial calibration of the IMU was erroneous.

In this case, the cycle slip correction (CSC) can be improved by jointly estimating the CSC and baseline coordinates using fixed DD phase measurements from all available satellites, the IMU baseline estimate of (22) and the baseline length a priori information, i.e.

$$\min_{b_{12}, \Delta N_{12}} \|z_{12} - Hb_{12} - A\Delta N_{12}\|_{\Sigma_{\varphi}^{-1}}^2 \quad \text{s. t.} \quad \|b_{12}\| = l, \quad (23)$$

with the combined GPS/INS measurement vector

$$z_{12} = \begin{pmatrix} \lambda(\varphi_{12} - \check{N}_{12}) \\ \hat{b}_{12}^{\text{IMU}} \end{pmatrix}, \quad (24)$$

and  $H$  being the redefined  $(K+3) \times 3$  DD geometry matrix,  $A(\lambda/2)$  is the extended  $(K+3) \times q$  CSC coefficient matrix for  $q$  DD cycle slips, and  $l$  describes the a priori information on the baseline length. The minimization of (23) includes a search of  $\Delta N_{12}$  inside a predefined search space volume  $\chi^2$  and an iterative computation of  $b_{12}$  for each integer candidate vector  $\Delta N_{12}$ . We use the orthogonal decomposition of Teunissen (1995) to rewrite the sum of squared errors of (23) as

$$\begin{aligned} &\|z_{12} - Hb_{12} - A\Delta N_{12}\|_{\Sigma_{\varepsilon_{12}}^{-1}}^2 \\ &= \|\Delta \hat{N}_{12} - \Delta N_{12}\|_{\Sigma_{\Delta \hat{N}_{12}}^{-1}}^2 + \|\check{b}_{12}(\Delta N_{12}) - b_{12}\|_{\Sigma_{\check{b}_{12}}^{-1}}^2 \\ &\quad + \|P_{\check{A}}^{\perp} P_H^{\perp} z_{12}\|_{\Sigma_{z_{12}}^{-1}}^2, \end{aligned} \quad (25)$$

with  $P_H^{\perp}$  being the projector on the orthogonal complement of the range space of  $H$  and  $\check{A} = P_H^{\perp} A$ . The first term on the right side was further developed by Teunissen (1995):

$$\|\Delta \hat{N}_{12} - \Delta N_{12}\|_{\Sigma_{\Delta \hat{N}_{12}}^{-1}}^2 = \sum_{l=1}^k \frac{(\Delta N_{12}^l - \Delta \hat{N}_{12}^{l1, \dots, l-1})^2}{(\sigma_{\Delta \hat{N}_{12}^{l1, \dots, l-1}})^2}, \quad (26)$$

with  $\Delta \hat{N}_{12}^{l1, \dots, l-1}$  being the  $l$ -th conditional cycle slip estimate. Setting (26) into (25), defining the search space volume  $\chi^2$  as an upper bound on (25), adding the baseline length



constraint as a zero term, and solving for the  $k$ -th ambiguity yields:

$$\begin{aligned} \frac{(\Delta N_{12}^k - \Delta \hat{N}_{12}^{k|1, \dots, k-1})^2}{\sigma_{\Delta \hat{N}_{12}^{k|1, \dots, k-1}}^2} &\leq \chi^2 - \|P_A^{\frac{1}{2}} P_H^{\frac{1}{2}} z_{12}\|_{\Sigma_{z_{12}}^{-1}}^2 \quad (27) \\ &- \sum_{l=1}^{k-1} \frac{(\Delta N_{12}^l - \Delta \hat{N}_{12}^{l|1, \dots, l-1})^2}{(\sigma_{\Delta \hat{N}_{12}^{l|1, \dots, l-1}})^2} \\ &- \min_{b_{12}, \mu} \left( \|\check{b}_{12}(N_{12}) - b_{12}\|_{\Sigma_{b_{12}}^{-1}}^2 + \mu \cdot (\|b_{12}\|^2 - l^2) \right) \end{aligned}$$

with the Lagrange multiplier  $\mu$ .

We set the partial derivative of the last term w.r.t.  $b_{12}$  to zero and solve it for  $b_{12}$  to obtain

$$\hat{b}_{12}(\mu) = (\Sigma_{b_{12}}^{-1} + \mu \mathbf{1})^{-1} \Sigma_{b_{12}}^{-1} \check{b}_{12}(N_{12}) \quad (28)$$

Setting  $\hat{b}_{12}(\mu)$  into the length constraint finally results in a root finding problem:

$$f(\mu) = \|\hat{b}_{12}(\mu)\|^2 - l^2 \stackrel{!}{=} 0. \quad (29)$$

As the roots of  $f(\mu)$  can not be found in closed form, we use the iterative Newton method. The Lagrange parameter  $\mu$  is given at the  $(n+1)$ -th iteration by

$$\mu^{(n+1)} = \mu^{(n)} - f(\mu) / \left. \frac{\partial}{\partial \mu} f(\mu) \right|_{\mu=\mu^{(n)}} \quad (30)$$

with

$$\frac{\partial}{\partial \mu} f(\mu) = 2(\hat{b}_{12}(\mu))^T \frac{\partial}{\partial \mu} (\hat{b}_{12}(\mu)). \quad (31)$$

Let  $\tilde{\Lambda}(\mu) = \Sigma_{z_{12}}^{-1} + \mu \cdot \mathbf{1}$ , then

$$\hat{b}_{12}(\mu) = (\tilde{\Lambda}(\mu))^{-1} \Sigma_{z_{12}}^{-1} z_{12}, \quad (32)$$

and

$$\begin{aligned} \frac{\partial}{\partial \mu} (\hat{b}_{12}(\mu)) &= \frac{\partial}{\partial \mu} (\tilde{\Lambda}^{-1}(\mu)) \Sigma_{z_{12}}^{-1} z_{12} \\ &= -\tilde{\Lambda}^{-1}(\mu) \frac{\partial}{\partial \mu} (\tilde{\Lambda}(\mu)) \tilde{\Lambda}^{-1}(\mu) \Sigma_{z_{12}}^{-1} z_{12} \\ &= -\tilde{\Lambda}^{-1}(\mu) \tilde{\Lambda}^{-1}(\mu) \Sigma_{z_{12}}^{-1} z_{12}. \quad (33) \end{aligned}$$

## 5 Measurement Analysis

In this section, the proposed cycle slip detection and correction method is verified with real data. We used the following measurement setup:

- measurement period:
  - week number: 1738, TOW  $\in \{159,156 \text{ s}, 159,812 \text{ s}\}$
- measurement location (area of test drive):
  - longitude  $\in \{11.41869^\circ, 11.42815^\circ\}$
  - latitude  $\in \{47.99410^\circ, 47.99723^\circ\}$
- measurement equipment:
  - 2 LEA 6T GPS receivers, 5 Hz, u-blox
  - 2 single frequency patch antennas, u-blox
  - 1 MPU 9150 inertial sensor, 100 Hz, Invensense
- installation:
  - mounting of GPS antennas on roof of car
  - alignment of baseline between antennas with longitudinal axis of car
  - baseline length:  $\bar{l} = 1.32 \text{ m}$ ,  $\sigma_{\bar{l}} = 1 \text{ cm}$

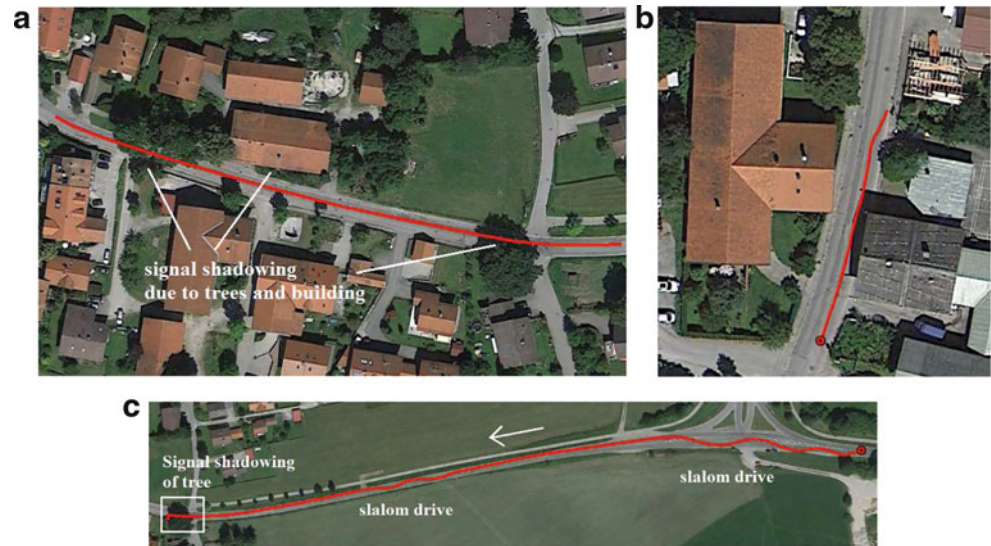
Figure 2 shows selected sections of the test drive. In the first two sections, the car was passing below trees and between two buildings. The code and carrier phase signals of all visible satellites are temporarily affected by significant multipath. The last subfigure shows a slalom drive with high receiver dynamics.

Figure 3 shows the fixed double difference phase residuals *before* cycle slip detection and correction over time. In the upper left corner, a skyplot shows the satellite geometry. We applied an elevation mask of  $20^\circ$ . At 176, 241 and 401 s, the car was passing below trees or between two buildings, which resulted in phase jumps of more than 10 half cycles. The largest residuals correspond to the satellites of lowest elevation. Numerous additional slips of  $\{-2, -1, 0, +1, +2\}$  cycles can be observed in between the major jumps.

Figure 4 shows the fixed double difference phase residuals *after* cycle slip detection and correction (CSC) using only GPS measurements in (23). The residuals are substantially lower than in Fig. 3 but there remain numerous undetected cycle slips.

Figure 5 shows the fixed double difference phase residuals *after* GPS/INS-based cycle slip detection and correction. The residuals of all satellites are substantially reduced to less than 15 cm. Variations of the residuals of more than 1 cm correspond mainly to the double difference phase multipath. The subplot in the lower right corner shows the fixed phase residuals only for the three satellites of highest elevation. One can observe three sections with severe multipath.

**Fig. 2** Selected sections of the test drive: the first two sections are in a high multipath environment. The last section is characterized by high receiver dynamics. (a) 136–156 s: passing below trees, (b) 172–177 s: passing between two buildings, (c) 330–404 s: slalom drive



**Fig. 3** Fixed phase residuals before cycle slip detection and correction: The residuals are jumping by more than 10 half cycles at 176, 241 and 401 s, where the car was passing below trees or between two buildings. The largest residuals correspond to the satellites of lowest elevation. The satellite geometry is shown in the skyplot (20° elevation mask) in the upper left corner. Numerous additional slips of  $\{-2, -1, 0, +1, +2\}$  cycles can be observed in between the major jumps

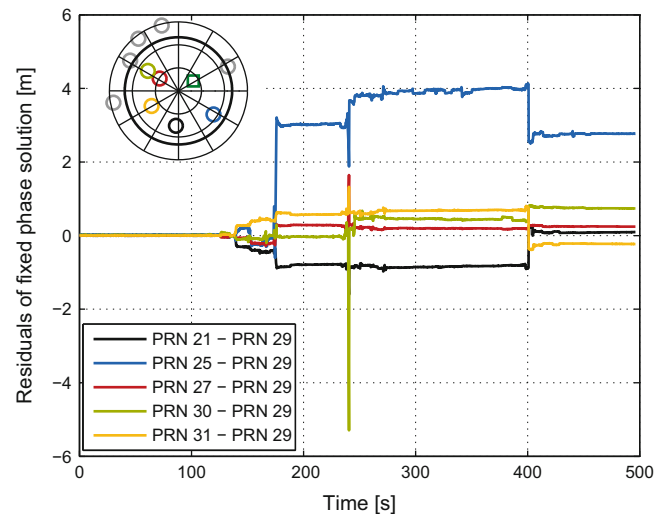


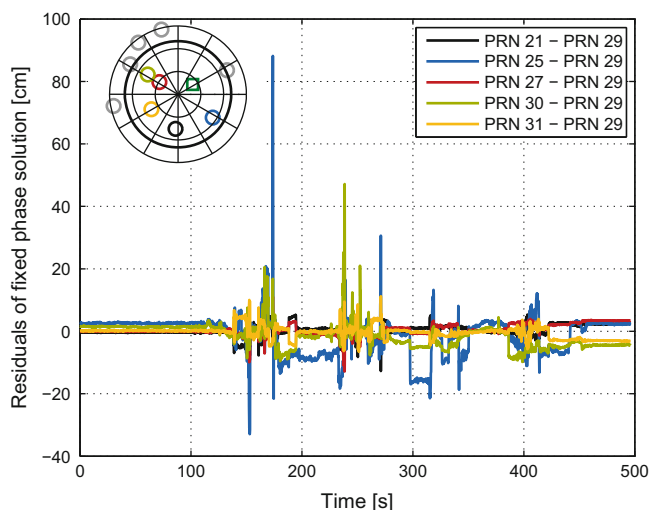
Figure 6 shows a comparison of the heading estimates with GPS-only CSC and with GPS/INS combined CSC. In the first 100 s, the car was standing. The noise of the heading estimates is in the order of  $0.1^\circ$  only. In the test drive, the heading with GPS-only CSC differs by up to  $30^\circ$  from the heading with GPS/INS combined CSC, as the GPS-only CSC can not correct for all cycle slips in heavy multipath environments. The enlarged periodic heading variations between 330 and 380 s indicate high receiver dynamics and correspond to the slalom drive. The GPS-only based CSC corrects at 341 s its erroneous ambiguities and, thus, follows the combined GPS/INS solution. The enlarged heading between 390 and 405 s shows some ripples in the heading estimate for the GPS/INS CSC. In this section, the car was passing below a tree and all code and carrier phases were affected by substantial multipath.

The reliability of the GPS/INS combined CSC depends on the drift of the IMU. We re-initialize the IMU at every

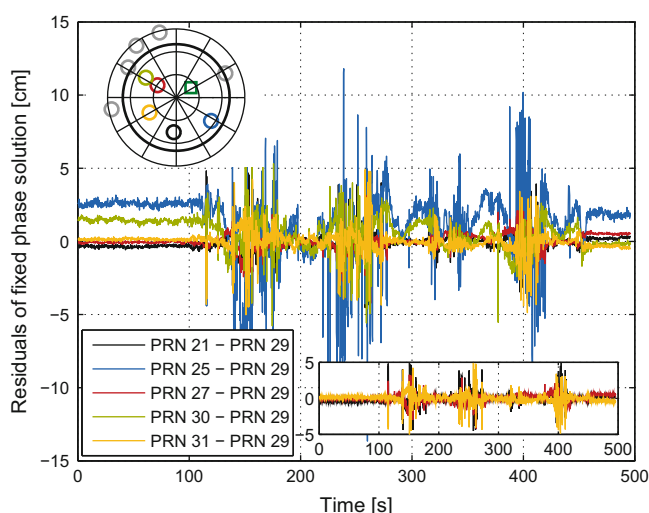
GPS measurement epoch after cycle slip correction with the GPS solution. Consequently, the drift of the IMU only between two subsequent GPS measurement epochs (0.2 s) is relevant. Figure 7 shows the difference between the heading of the IMU without continuous GPS-based calibration and the heading of the GPS/INS-combined solution. We can observe a slight continuous drift of less than  $0.5^\circ$  with temporarily increased variations. The increased variations of up to  $1^\circ/0.2$  s are most likely caused by heading errors of the GPS/INS-combined solution (in multipath environments) and not by changes of the IMU's drift. A cycle slip correction can still be performed reliably.

The proposed CSC of (27) jointly determines the baseline coordinates and integer cycle slips using both GPS DD carrier phases and the IMU-predicted baseline. It finds the optimized trade-off between minimizing the GPS measurement residuals and minimizing the IMU-predicted baseline residuals. We compare the performance of this optimized

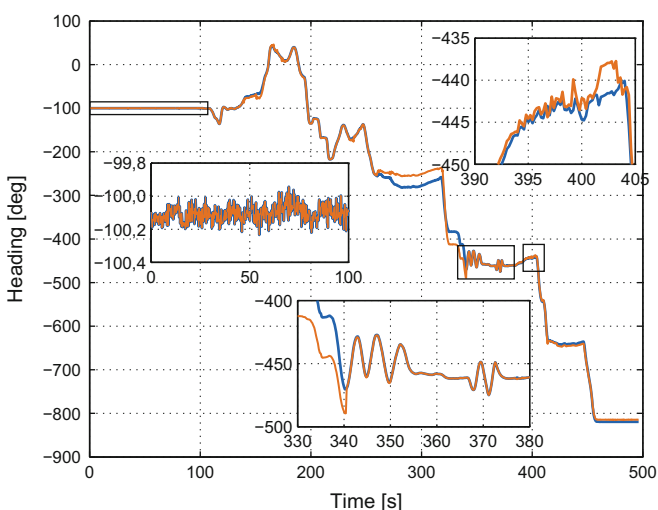
**Fig. 4** Fixed phase residuals after GPS-based cycle slip detection and correction: the residuals are substantially reduced but there remain numerous undetected cycle slips



**Fig. 5** Fixed phase residuals after GPS/IMU-based cycle slip detection and correction: the residuals of all satellites are reduced to less than 15 cm. The “noise” in the residuals corresponds to the double difference phase multipath. The subplot in the lower right corner shows the fixed phase residuals only for the three satellites of highest elevation. One can observe three sections with severe multipath



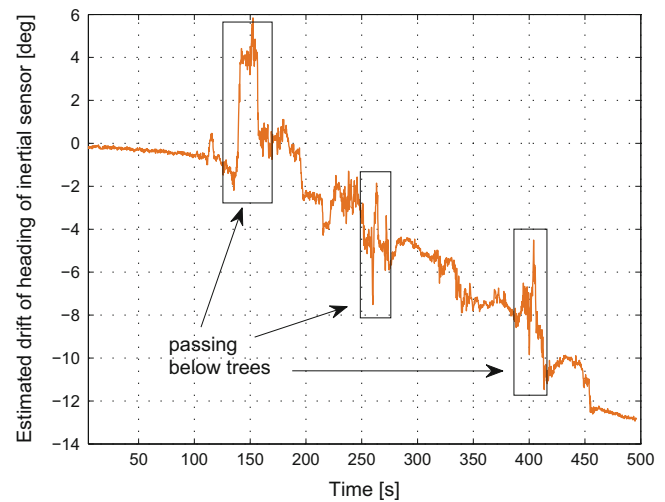
**Fig. 6** Heading determination with GPS/INS combined cycle slip correction (marked in blue) and GPS-only (marked in orange) cycle slip correction: the noise level of both heading estimates is in the order of 0.1° in static conditions and increases to 1° in high multipath environments (e.g. passing below a tree between 395 and 405 s). The periodic variations between 330 and 380 s indicate high receiver dynamics. The GPS-based heading temporarily differs from the GPS/INS-based heading due to some uncorrected half cycle slips



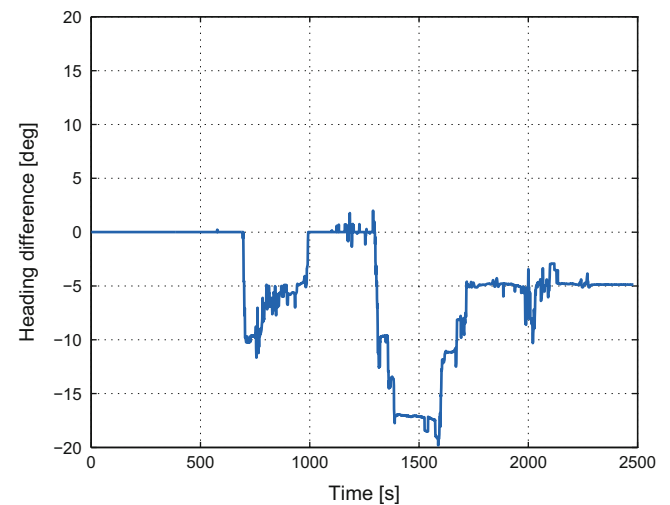
approach to the performance of an integer least-squares CSC estimation with the baseline vector considered known and fixed from the IMU. This alternative approach also combines GPS and INS measurements. However, the integer least-

squares estimation only of the CSC does not guarantee a final baseline estimate to be close to the one from the IMU. Figure 8 shows that the heading of the integer least-squares estimation can differ by up to 20° from the combined

**Fig. 7** Drift of heading of inertial sensor without continuous GPS-based calibration: the change of the heading between two subsequent epochs is in general less than  $0.5^\circ$ . The largest variation of  $1^\circ/0.2$  s occurs during the passing below trees, where the coupled GPS/INS was itself more noisy. A cycle slip correction can still be performed reliably



**Fig. 8** Cycle slip correction with GPS/INS: the heading of the integer least-squares cycle slip estimation with fixed IMU baseline differs by up to  $20^\circ$  from the heading of the combined cycle slip and baseline estimation using both GPS and IMU measurements. This indicates that a pure integer least-squares estimation is not sufficient. As the combined estimator finds an optimized trade-off between minimizing the squared measurement residuals and minimizing the baseline vector residuals, it finds the correct cycle slip correction with highest success rate



solution. This shows again the need for the proposed cycle slip detection and correction method, which combines all available measurement and a priori information and minimizes the sum of the measurement residuals and a priori information residuals.

The reliability of cycle slip detection and correction can also be verified by adding artificial cycle slips to the raw phase measurements. We have tested all  $\{-2, -1, 0, +1, 2\}$  cycle slip combinations for  $\{1, 2, 3\}$  simultaneous cycle slips at a moment of high receiver dynamics ( $t = 342$  s), and observed that all cycle slip combinations were correctly found.

## 6 Conclusion

In this paper, a method for reliable cycle slip detection and correction was proposed for *low-cost* GPS receivers with *high receiver dynamics* in *challenging* environments. It combines double difference GPS carrier phases from all visible

satellites, gyroscope and acceleration measurements, and a priori information on the baseline length. It performs a tree search for finding the cycle slip corrections. The estimator was first tested by a simulation of artificial cycle slips. It turned out to be extremely powerful if the IMU was properly calibrated, i.e. it was able to correct simultaneous cycle slips at all visible satellites. Subsequently, the proposed method was verified in a test drive with two low-cost GPS receivers and a low-cost 9-axes INS. The residuals of the fixed phase measurements after cycle slip correction remained less than 15 cm in both a slalom drive with high receiver dynamics and passages below trees.

## References

- Colombo OL, Bhapkar UV, Evans AG (1999) Inertial-aided cycle-slip detection/correction for precise, long-baseline kinematic GPS. In: Proceedings of the 12th ION GNSS, Nashville, pp 1915–1922
- Dai Z, Knedlik S, Loffeld O (2009) Instantaneous triple-frequency GPS cycle-slip detection and repair. *Intl J Navig Obs* 2009:15. Article ID 407231

- Du S, Gao Y (2012) Inertial aided cycle slip detection and identification for integrated PPP GPS and INS. *Sensors* 12(11):14,344–14,362
- Henkel P, Gunther C (2013) Attitude determination with low-cost GPS/INS. In: Proceedings of the 26-th ION GNSS+, Nashville, pp 2015–2023
- Henkel P, Kiam JJ (2013) Maximum a posteriori probability estimation of integer ambiguities and baseline. In: IEEE proceedings of the 55th symposium ELMAR, Zadar, pp 353–356
- Jekeli C (2001) Inertial navigation systems with geodetic applications. Walter de Gruyter
- Teunissen P (1995) The least-squares ambiguity decorrelation adjustment: a method for fast GPS ambiguity estimation. *J Geod* 70:65–82
- Teunissen P (2006) The LAMBDA method for the GNSS compass. *Artif Satell* 41(3):89–103
- Teunissen P (2010) Integer least-squares theory for the GNSS compass. *J Geod* 84:433–447

---

# Adjusting the Errors-In-Variables Model: Linearized Least-Squares vs. Nonlinear Total Least-Squares

Burkhard Schaffrin

---

## Abstract

It has long been known that the Errors-In-Variables (EIV) Model is a special case of the nonlinear Gauss–Helmert Model (GHM) and can, therefore, be adjusted by standard least-squares techniques in iteratively linearized GH-Models, which is the approach by Helmert (Adjustment Computations Based on the Least-Squares Principle (in German), 1907) and – later – by Deming (Phil Mag 11:146–158, 1931; Phil Mag 17:804–829, 1934).

Apart from the fact that there are, at least, two other nonlinear models that are equivalent to the above GH-Model, thus allowing two more classical least-squares approaches based on iterative linearization, it was the seminal paper by Golub and van Loan (SIAM J Numer Anal 17:883–893, 1980) in which they proved that a purely nonlinear approach can be followed as well, thereby avoiding any model linearization. They called such an approach “Total Least-Squares adjustment” by which any normal equations may be replaced by a simple eigenvalue problem, as long as only diagonal dispersion matrices are involved.

Here, an attempt will be made to show the differences and parallels in various algorithms, even in the fully weighted case, which obviously all generate the same results, but without necessarily showing equal efficiency in doing so, as is well known since the publications by Schaffrin and Wieser (J Geodesy 82:415–421, 2008), Fang (Weighted Total Least-Squares solutions with applications in geodesy, 2011), and Mahboub (J Geodesy 86:359–367, 2012).

---

## Keywords

Errors-In-Variables Models • Total Least-Squares • Equivalent nonlinear models • Linearized Least-Squares

---

## 1 Introduction

The Errors-In-Variables (EIV) Model has recently seen a lot of attention since, in accordance with Golub and van Loan (1980), it can be treated in its *nonlinear* form by a least-squares approach that they coined “*Total Least-Squares adjustment*”. It eventually leads to a (generalized) eigenvalue problem that needs to be solved in lieu of the sequence of

normal equations that would result from a traditional “*Least-Squares adjustment*” within *iteratively linearized* models. The latter approach dates, at least, back to Helmert (1907), but has as well been used by Deming (1931, 1934) for the approximation of curves and, more recently, by Neitzel (2010) to determine the parameters of a similarity transformation.

In contrast, the nonlinear Total Least-Squares (TLS) approach which, in its original formulation, could tolerate only “*element-wise weighting*” and thus only *diagonal* weight matrices, has since been generalized in several steps by Schaffrin and Wieser (2008), Fang (2011), and Mahboub (2012) to now accept any *positive-definite weight matrices*.

---

B. Schaffrin (✉)  
Division of Geodetic Science, School of Earth Sciences, The Ohio  
State University, Columbus, OH, USA  
e-mail: [schaffrin.1@osu.edu](mailto:schaffrin.1@osu.edu)

This development will be presented in the following Sect. 2, thereby showing how the more specialized algorithms can be derived from the more general ones by simplification.

Moreover, it should be noted that progress has also been made towards the use of *positive-semidefinite dispersion matrices* in TLS adjustment, which may be handled as described by Schaffrin et al. (2014). These cases are quite relevant whenever the random error matrix needs to show a certain pattern or structure after the adjustment. Due to the limited space, these advanced methods will not be discussed below.

Instead, attention will be paid to a *triplet of classical nonlinear models* that all can be constructed to be *equivalent* to the EIV-Model and, furthermore, may undergo a sequence of Least-Squares adjustments via iterative linearization which, in the end, converge to the very same TLS solution. This will be the theme in Sect. 3 although many details have to be left out; for those, see Schaffrin (2015).

## 2 Nonlinear TLS Adjustment in an EIV-Model

### 2.1 Fang's Algorithm

Let the EIV-Model be defined by

$$y = \underbrace{(A - E_A)}_{n \times m} \xi + e_y, \quad rkA = m < n, \quad (1a)$$

$$e := \begin{bmatrix} e_y \\ e_A := \text{vec} E_A \end{bmatrix} \sim \left( \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \sigma_o^2 \begin{bmatrix} P_y^{-1} & 0 \\ 0 & P_A^{-1} \end{bmatrix} \right) =: \sigma_o^2 P^{-1} \quad (1b)$$

where

$y$  is the  $n \times 1$  observation vector;

$A$  is the  $n \times m$  (random) coefficient matrix with full column rank (aka "data matrix");

$E_A$  is the  $n \times m$  (unknown) random error matrix associated with  $A$ ;

$\xi$  is the  $m \times 1$  (unknown) parameter vector;

$e_y$  is the  $n \times 1$  (unknown) random error vector associated with  $y$ ;

$e_A$  is the  $nm \times 1$  vectorial form of the matrix  $E_A$ ;

$Q$  is the  $n(m+1) \times n(m+1)$  block-diagonal pos.- def. cofactor matrix;

$P := Q^{-1}$  is the corresponding block-diagonal pos.- def. weight matrix;

$\sigma_o^2$  is the (unknown) variance component (unit- free);

$\text{Cov}\{e_y, \text{vec} E_A\} = 0$  for the sake of simplicity.

The model generalizes the one used by Schaffrin and Wieser (2008) where a Kronecker product structure for

$$Q_A = P_A^{-1} = Q_o \otimes Q_x \quad (2)$$

was assumed, as well as the one used by Golub and von Loan (1980) who only allowed *diagonal* cofactor matrices with

$$Q_o := I_m, \quad Q_x := Q_y = \text{Diag}(p_1^{-1}, \dots, p_n^{-1}) = P_y^{-1}. \quad (3)$$

The objectives of a *nonlinear Total Least-Squares (TLS) adjustment* are now based on the principle

$$e_y^T P_y e_y + e_A^T P_A e_A = \min. \quad \text{s.t.} \quad (1a), \quad (4)$$

which can be given the equivalent form of a *Lagrange target function*, namely:

$$\phi(e_y, e_A, \xi, \lambda) := e_y^T P_y e_y + e_A^T P_A e_A + 2\lambda^T [y - A\xi - e_y + (\xi^T \otimes I_n) e_A] = \text{stationary}. \quad (5)$$

Consequently, the Euler-Lagrange necessary conditions result in the following system of *nonlinear "normal equations"*:

$$\frac{1}{2} \frac{\partial \phi}{\partial e_y} = P_y \tilde{e}_y - \hat{\lambda} \doteq 0 \quad (6a)$$

$$\frac{1}{2} \frac{\partial \phi}{\partial e_A} = P_A \tilde{e}_A + (\hat{\xi}^T \otimes I_n) \hat{\lambda} \doteq 0, \quad (6b)$$

$$\frac{1}{2} \frac{\partial \phi}{\partial \xi} = -(A - \tilde{E}_A)^T \hat{\lambda} \doteq 0, \quad (6c)$$

$$\frac{1}{2} \frac{\partial \phi}{\partial \lambda} = y - A\hat{\xi} - \tilde{e}_y + (\hat{\xi}^T \otimes I_n) \tilde{e}_A \doteq 0, \quad (6d)$$

which still needs to be reduced by partial elimination since the sufficient condition is fulfilled as

$$\frac{1}{2} \frac{\partial^2 \phi}{\partial \begin{bmatrix} e_y \\ e_A \end{bmatrix} \partial \begin{bmatrix} e_y^T \\ e_A^T \end{bmatrix}} = \begin{bmatrix} P_y & 0 \\ 0 & P_A \end{bmatrix} \text{ is pos.-def.} \quad (7)$$

Now, (6a, b) are transformed to provide the *residual vectors* through

$$\tilde{e}_y = Q_y \hat{\lambda} \quad \text{and} \quad \tilde{e}_A = -Q_A (\hat{\xi} \otimes I_n) \hat{\lambda} \quad (8a)$$

so that (6d) can be rewritten as

$$y - A\hat{\xi} = \left[ Q_y + (\hat{\xi} \otimes I_n)^T Q_A (\hat{\xi} \otimes I_n) \right] \cdot \hat{\lambda} =: Q_1 \cdot \hat{\lambda}, \tag{8b}$$

with  $Q_1 = Q_1(\hat{\xi})$  being *nonsingular*, thus leading to

$$\hat{\lambda} = Q_1^{-1} (y - A\hat{\xi}) \tag{9}$$

and, together with (6c), to the system

$$\begin{bmatrix} Q_1 & (A - \tilde{E}_A) \\ (A - \tilde{E}_A)^T & 0 \end{bmatrix} \begin{bmatrix} \hat{\lambda} \\ \hat{\xi} \end{bmatrix} = \begin{bmatrix} y - \tilde{E}_A \hat{\xi} \\ 0 \end{bmatrix} \tag{10}$$

Obviously, the estimated parameter vector is now obtained as in Fang (2011, p.27) via

$$\hat{\xi} = \left[ (A - \tilde{E}_A)^T Q_1^{-1} (A - \tilde{E}_A) \right]^{-1} (A - \tilde{E}_A)^T Q_1^{-1} \cdot (y - \tilde{E}_A \hat{\xi}) \tag{11}$$

and allows updates for  $Q_1$ ,  $\hat{\lambda}$ , and  $\tilde{e}_A$ , from which a new estimate  $\hat{\xi}$  results.

The Total Sum of weighted Squared Residuals (TSSR) may now readily be computed from

$$\begin{aligned} \tilde{e}_y^T P_y \tilde{e}_y + \tilde{e}_A^T P_A \tilde{e}_A &= \hat{\lambda}^T \left[ Q_y + (\hat{\xi} \otimes I_n)^T Q_A (\hat{\xi} \otimes I_n) \right] \hat{\lambda} \\ &= \hat{\lambda}^T Q_1 \hat{\lambda} = \hat{\lambda}^T (y - A\hat{\xi}) =: \text{TSSR} \end{aligned} \tag{12}$$

so that a suitable variance component estimate may be obtained through

$$\hat{\sigma}_o^2 = \hat{\lambda}^T (y - A\hat{\xi}) / (n - m) = \text{TSSR} / (n - m) \tag{13}$$

as the redundancy in model (1a, b) is still  $n - m$ .

Alternatively, system (10) can be given the *asymmetric* form

$$\begin{bmatrix} Q_1 & A \\ (A - \tilde{E}_A)^T & 0 \end{bmatrix} \begin{bmatrix} \hat{\lambda} \\ \hat{\xi} \end{bmatrix} = \begin{bmatrix} y \\ 0 \end{bmatrix} \tag{14}$$

which would then provide the estimated parameter vector through

$$\hat{\xi} = \left[ (A - \tilde{E}_A)^T Q_1^{-1} A \right]^{-1} (A - \tilde{E}_A)^T Q_1^{-1} y \tag{15}$$

and should lead to a similar iteration as before. Note that (15) also appears as formula (21) in Xu et al. (2012), but

essentially represents a variant of Fang’s algorithm; also, cf. Fang (2013) where further alternatives are presented.

## 2.2 Mahboub’s Algorithm

On the other hand, combining (9) with (6c) leads to the following sequence of identities:

$$\begin{aligned} A^T Q_1^{-1} (y - A\hat{\xi}) &= A^T \hat{\lambda} = \tilde{E}_A^T \hat{\lambda} = (\hat{\lambda}^T \otimes I_m) \text{vec}(\tilde{E}_A^T) = \\ &= (\hat{\lambda}^T \otimes I_m) \cdot (K \tilde{e}_A) = (I_m \otimes \hat{\lambda}^T) \tilde{e}_A = \\ &= - \left[ (I_m \otimes \hat{\lambda})^T Q_A (\hat{\xi} \otimes I_n) \right] \hat{\lambda} =: -R_1 \cdot \hat{\lambda} = \\ &= -R_1 \cdot Q_1^{-1} (y - A\hat{\xi}) \end{aligned} \tag{16}$$

where  $K$  denotes a  $nm \times nm$  “commutation matrix” that is also known as “vec-permutation matrix”; for more details, see Magnus and Neudecker (2007).

Obviously, (16) translates into the estimated parameter vector

$$\hat{\xi} = [(A^T + R_1) Q_1^{-1} A]^{-1} (A^T + R_1) Q_1^{-1} y \tag{17a}$$

with  $R_1 = R_1(\hat{\xi}, \hat{\lambda})$  and, from (16), with

$$R_1 \hat{\lambda} = -\tilde{E}_A^T \hat{\lambda} \tag{17b}$$

without necessarily implying that  $R_1 = -\tilde{E}_A^T$ . Therefore, the sequence of solutions to (15) may differ from the sequence of solutions to (17a) when iteratively updating  $Q_1$ ,  $\hat{\lambda}$ , and  $R_1$ , before a new parameter vector estimate  $\hat{\xi}$  can be found; yet the ultimate convergence points will be the same.

Again, the TSSR can be computed from (12) which will lead to the variance component estimate in (13).

## 2.3 A New Variant of Mahboub’s Algorithm

After giving (16) the form

$$A^T Q_1^{-1} (y - A\hat{\xi}) = - \left[ (I_m \otimes \hat{\lambda})^T Q_A (I_m \otimes \hat{\lambda}) \right] \hat{\xi}, \tag{18a}$$

the estimated parameter vector may as well be obtained from

$$\hat{\xi} = \left[ A^T Q_1^{-1} A - \left[ (I_m \otimes \hat{\lambda})^T Q_A (I_m \otimes \hat{\lambda}) \right] \right]^{-1} A^T Q_1^{-1} y \tag{18b}$$



thus allowing updates for  $Q_1$  and  $\hat{\lambda}$ . This algorithm will be further explored in the near future.

## 2.4 The Schaffrin–Wieser Algorithm

This algorithm was designed for the somewhat more special case where the cofactor matrix  $Q_A$  can be split into a *Kronecker product*, thereby indicating that all columns have cofactor matrices proportional to each other. This implies

$$Q_A = Q_o \otimes Q_x \Rightarrow Q_1 = Q_y + (\hat{\xi}^T Q_o \hat{\xi}) \cdot Q_x \quad (19)$$

and thus

$$A^T Q_1^{-1} (y - A\hat{\xi}) = A^T [Q_y + (\hat{\xi}^T Q_o \hat{\xi}) \cdot Q_x]^{-1} \cdot (y - A\hat{\xi}) = A^T \hat{\lambda} = \tilde{E}_A^T \hat{\lambda} \quad (20)$$

with

$$\tilde{e}_A = - (Q_o \hat{\xi} \otimes Q_x) \hat{\lambda} = -\text{vec} (Q_x \hat{\lambda} \hat{\xi}^T Q_o) = \text{vec} \tilde{E}_A. \quad (21)$$

(20) and (21) together generate the identity

$$A^T [Q_y + (\hat{\xi}^T Q_o \hat{\xi}) \cdot Q_x]^{-1} (y - A\hat{\xi}) = -Q_o \hat{\xi} \cdot (\hat{\lambda}^T Q_x \hat{\lambda}) =: -Q_o \hat{\xi} \cdot \hat{v} \quad (22a)$$

suggesting the *iteration*

$$\hat{\xi} = \left( A^T [Q_y + (\hat{\xi}^T Q_o \hat{\xi}) \cdot Q_x]^{-1} A - Q_o \hat{v} \right)^{-1} \cdot A^T [Q_y + (\hat{\xi}^T Q_o \hat{\xi}) \cdot Q_x]^{-1} y \quad (22b)$$

with

$$\hat{v} := (\hat{\lambda}^T Q_x \hat{\lambda}) \quad \text{and} \quad \hat{\lambda} := [Q_y + (\hat{\xi}^T Q_o \hat{\xi}) \cdot Q_x]^{-1} \cdot (y - A\hat{\xi}) \quad (22c)$$

while (12) and (13) generate first the TSSR and then a suitable variance component estimate.

## 2.5 The Golub–van-Loan Algorithm

Now, the condition (19) is further specialized to

$$Q_x := Q_y \Rightarrow Q_1 = (1 + \hat{\xi}^T Q_o \hat{\xi}) \cdot Q_y \quad (23a)$$

and

$$\hat{\lambda} = (1 + \hat{\xi}^T Q_o \hat{\xi})^{-1} Q_y^{-1} (y - A\hat{\xi}) \quad (23b)$$

so that (22a) becomes

$$A^T Q_y^{-1} (y - A\hat{\xi}) = -Q_o \hat{\xi} \cdot \hat{v} (1 + \hat{\xi}^T Q_o \hat{\xi}) =: -Q_o \hat{\xi} \cdot \sigma_{\min}^2 \quad (24a)$$

with

$$\begin{aligned} \sigma_{\min}^2 &= (\hat{\lambda}^T Q_y \hat{\lambda}) (1 + \hat{\xi}^T Q_o \hat{\xi}) = \\ &= (y - A\hat{\xi})^T Q_y^{-1} (y - A\hat{\xi}) / (1 + \hat{\xi}^T Q_o \hat{\xi}), \end{aligned} \quad (24b)$$

and this, from (24a, b), becomes

$$\begin{aligned} \sigma_{\min}^2 \cdot (1 + \hat{\xi}^T Q_o \hat{\xi}) &= y^T Q_y^{-1} (y - A\hat{\xi}) + (\hat{\xi}^T Q_o \hat{\xi}) \cdot \sigma_{\min}^2 \Rightarrow \\ \Rightarrow \sigma_{\min}^2 &= y^T Q_y^{-1} (y - A\hat{\xi}) = \text{TSSR} \end{aligned} \quad (24c)$$

(24a) and (24c) allow the problem to be rephrased as a *generalized eigenvalue problem*, specifically as:

$$\begin{bmatrix} A^T Q_y^{-1} A & A^T Q_y^{-1} y \\ y^T Q_y^{-1} A & y^T Q_y^{-1} y \end{bmatrix} \begin{bmatrix} \hat{\xi} \\ -1 \end{bmatrix} = \begin{bmatrix} Q_o & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} \hat{\xi} \\ -1 \end{bmatrix} \cdot \sigma_{\min}^2 \quad (25)$$

with the variance component estimate

$$\hat{\sigma}_o^2 = \sigma_{\min}^2 / (n - m) \quad (26)$$

The original situation, treated by Golub and van Loan (1980), was characterized by the further specializations

$$Q_o := I_m \quad \text{and} \quad Q_y := \text{Diag} (p_1^{-1}, \dots, p_n^{-1}) = P^{-1} \quad (27)$$

which, in turn, lead to the *standard eigenvalue problem*

$$\begin{bmatrix} A^T P A & A^T P y \\ y^T P A & y^T P y \end{bmatrix} \begin{bmatrix} \hat{\xi} \\ -1 \end{bmatrix} = \begin{bmatrix} \hat{\xi} \\ -1 \end{bmatrix} \cdot \sigma_{\min}^2 \quad (28)$$

whose solution provides the *Total Least-Squares Solution* (TLSS).

In the next section, a few *equivalent models* will be presented for which, traditionally, an identical weighted LEast-Squares Solution (LESS) would have been found after *iterative linearization*.

### 3 Traditional Models, Equivalent to the EIV-Model

#### 3.1 The Nonlinear Gauss–Helmert Model

Here, the new vectors

$$Y := \text{vec}[y|A] \quad \text{and} \quad e := \text{vec}[e_y|E_A] \quad (29)$$

are introduced. Then,

$$\underline{b} \left( \mu := \underbrace{Y - e}_{n(m+1) \times 1}, \xi \right) := (y - e_y) - (A - E_A) \xi = 0, \quad (30a)$$

$$e \sim \left( 0, \sigma_o^2 \begin{bmatrix} Q_y & 0 \\ 0 & Q_A \end{bmatrix} = \sigma_o^2 Q \right), \quad (30b)$$

with the *nonlinear* vector-valued vector function

$$\underline{b} : R^{(n+1)(m+1)-1} \rightarrow R^n, \quad (30c)$$

due to the term  $E_A \cdot \xi$ , forms an equivalent *Gauss–Helmert Model* that would traditionally be linearized for an iterative Least-Squares adjustment.

The truncated Taylor series, following Pope (1972), then reads:

$$0 = \underline{b}(\mu, \xi) \approx \underline{b}(\mu_o, \xi_o) + \left. \frac{\partial \underline{b}(\mu, \xi)}{\partial [\mu^T | \xi^T]} \right|_{\mu_o, \xi_o} \cdot \begin{bmatrix} \mu - \mu_o \\ \xi - \xi_o \end{bmatrix} \quad (31)$$

with suitable approximations  $\xi_o$  and  $\mu_o := Y - \underset{\sim}{0}$  where  $\underset{\sim}{0}$  here denotes a “stochastic zero vector” of size  $n(m+1) \times 1$ . This leads first to

$$0 \approx \underline{b}(\mu_o, \xi_o) + \left. \frac{\partial \underline{b}}{\partial \mu^T} \right|_{\mu_o, \xi_o} \cdot \left( \underset{\sim}{0} - e \right) + \left. \frac{\partial \underline{b}}{\partial \xi^T} \right|_{\mu_o, \xi_o} \cdot (\xi - \xi_o), \quad (32)$$

then to

$$\underline{b}(\mu_o, \xi_o) + \left. \frac{\partial \underline{b}}{\partial \mu^T} \right|_{\mu_o, \xi_o} \cdot \underset{\sim}{0} \approx \left. \frac{\partial \underline{b}}{\partial \mu^T} \right|_{\mu_o, \xi_o} \cdot e - \left. \frac{\partial \underline{b}}{\partial \xi^T} \right|_{\mu_o, \xi_o} \cdot (\xi - \xi_o), \quad (33)$$

and eventually to the *linearized Gauss–Helmert Model*:

$$w_o := \underline{b}(Y = \mu_o + \underset{\sim}{0}, \xi_o) \approx B^{(o)} \cdot e + A^{(o)} \cdot (\xi - \xi_o), \quad (34a)$$

$$B^{(o)} := [I_n | -(\xi_o^T \otimes I_n)], \quad A^{(o)} := A - \underset{\sim}{0}, \quad (34b)$$

$$e \sim \left( 0, \sigma_o^2 Q = \sigma_o^2 \begin{bmatrix} Q_y & 0 \\ 0 & Q_A \end{bmatrix} \right). \quad (34c)$$

Note that the weighted LEast-Squares Solution (LESS) is now being formed through the *normal equations*

$$\left[ (A^{(o)})^T (Q_1^{(o)})^{-1} A^{(o)} \right] \cdot \widehat{\xi}^{(1)} = (A^{(o)})^T (Q_1^{(o)})^{-1} \cdot (y - \underset{\sim}{0} \cdot \xi_o) \quad (35a)$$

with

$$Q_1^{(o)} := B^{(o)} Q (B^{(o)})^T = Q_y + (\xi_o \otimes I_n)^T Q_A (\xi_o \otimes I_n), \quad (35b)$$

and the *residual vectors* through

$$\tilde{e}^{(1)} := \begin{bmatrix} \tilde{e}_y^{(1)} \\ \tilde{e}_A^{(1)} \end{bmatrix} = \begin{bmatrix} Q_y \\ -Q_A (\widehat{\xi}^{(1)} \otimes I_n) \end{bmatrix} (Q_1^{(o)})^{-1} \cdot (y - \underset{\sim}{0} \cdot \xi_o - A^{(o)} \cdot \widehat{\xi}^{(1)}) \quad (35c)$$

Looking at the next and all the following iteration steps, it becomes clear that this represents one specific *iterative solver* of Fang’s TLS normal equations (11).

For more details, see Fang (2011, ch. 4.4), Snow (2012, ch. 4), and the forthcoming OSU-Report by Schaffrin (2015), as well as Neitzel (2010) for a specific application.

#### 3.2 The Nonlinear Gauss–Markov Model

In this case, the expectation of the data matrix  $A$  is introduced as a new  $n \times m$  “parameter matrix”

$$\Xi_A := A - E_A \quad \text{with} \quad \xi_A := \text{vec} \Xi_A, \quad (36)$$

leading to the equivalent *Gauss–Markov Model*

$$y = (\xi \otimes I_n)^T \cdot \xi_A + e_y =: a(\xi, \xi_A) + e_y, \quad (37a)$$

$$e_y \sim (0, \sigma_o^2 Q_y), \quad (37b)$$

with the *nonlinear* vector-valued vector function

$$a : R^{(n+1)m} \rightarrow R^n \quad (37b)$$

due to the term  $\Xi_A \cdot \xi$ . The linearization of model (37a, b) with respect to the approximations  $\xi_o$  and  $\xi_A^{(o)} := \text{vec}(A^{(o)}) = \text{vec}(A - \underset{\sim}{0})$ , where  $\underset{\sim}{0}$  now denotes a “stochastic zero matrix” of size  $n \times m$ , then leads first to

$$\xi_A - \xi_A^{(o)} = \underset{\sim}{0} - e_A, \quad e_A \sim (0, \sigma_o^2 Q_A), \quad (38a)$$

$$\begin{aligned} y - e_y &\approx \underline{a}(\xi_o, \xi_A^{(o)}) + \frac{\partial \underline{a}(\xi, \xi_A)}{\partial [\xi^T, \xi_A^T]} \Big|_{\xi_o, \xi_A^{(o)}} \cdot \begin{bmatrix} \xi - \xi_o \\ \xi_A - \xi_A^{(o)} \end{bmatrix} \\ &= A^{(o)} \cdot \xi_o + A^{(o)} \cdot (\xi - \xi_o) + (\xi_o \otimes I_n)^T (\xi_A - \xi_A^{(o)}), \end{aligned} \quad (38b)$$

and finally to the *linearized Gauss–Markov Model*

$$\begin{bmatrix} y - A^{(o)} \cdot \xi_o \\ \underset{\sim}{0} \end{bmatrix} = \begin{bmatrix} A^{(o)} & | & (\xi_o \otimes I_n)^T \\ 0 & | & I_{nm} \end{bmatrix} \begin{bmatrix} \xi - \xi_o \\ \xi_A - \xi_A^{(o)} \end{bmatrix} + \begin{bmatrix} e_y \\ e_A \end{bmatrix}, \quad (39a)$$

$$\begin{bmatrix} e_y \\ e_A \end{bmatrix} \sim \left( \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \sigma_o^2 \begin{bmatrix} Q_y & 0 \\ 0 & Q_A \end{bmatrix} = \sigma_o^2 Q \right). \quad (39b)$$

After a number of further manipulations, the weighted LESS for model (39a, b) can be shown to fulfill the “normal equations”

$$\left[ (A^{(o)})^T (Q_1^{(o)})^{-1} A^{(o)} \right] \cdot \widehat{\xi}^{(1)} = (A^{(o)})^T (Q_1^{(o)})^{-1} y \quad (40a)$$

with

$$\begin{aligned} Q_1^{(o)} &:= \left[ Q_y^{-1} - Q_y^{-1} (\xi_o \otimes I_n)^T \cdot \right. \\ &\quad \left. \cdot \left[ Q_A^{-1} + (\xi_o \xi_o^T \otimes Q_y^{-1}) \right]^{-1} (\xi_o \otimes I_n) Q_y^{-1} \right]^{-1} \end{aligned} \quad (40b)$$

$$= Q_y + (\xi_o \otimes I_n)^T Q_A (\xi_o \otimes I_n) \quad (40c)$$

which nicely corresponds to (35a, b). More details can be found in the forthcoming OSU-Report by Schaffrin (2015).

### 3.3 The Model of Direct Observations with Nonlinear Constraints

Now, the expectation of the observation vector  $y$  is introduced as just another parameter vector  $\xi_y$  of size  $n \times 1$  so that the new model combines the *direct observation equations*

$$\begin{bmatrix} y \\ \text{vec} A \end{bmatrix} = \begin{bmatrix} \xi_y \\ \xi_A \end{bmatrix} + \begin{bmatrix} e_y \\ e_A \end{bmatrix}, \quad (41a)$$

$$\begin{bmatrix} e_y \\ e_A \end{bmatrix} \sim \left( \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \sigma_o^2 \begin{bmatrix} Q_y & 0 \\ 0 & Q_A \end{bmatrix} = \sigma_o^2 Q \right),$$

with the *nonlinear constraints*

$$\xi_y - \Xi_A \cdot \xi = 0 \quad (41b)$$

which might be linearized into

$$\left[ A^{(o)} \mid -I_n \mid (\xi_o \otimes I_n)^T \right] \begin{bmatrix} \xi - \xi_o \\ \xi_y - \xi_y^{(o)} \\ \xi_A - \xi_A^{(o)} \end{bmatrix} = 0. \quad (42)$$

In the already mentioned OSU-Report by Schaffrin (2015), it will be shown how the resulting iterative LESS's do converge to the Total Least-Squares Solution.

For another take on this model, refer to Donevska et al. (2011) who stress the equivalence to orthogonal regression as applied by Deming (1931, 1934).

## 4 Conclusions

It has been clarified that the TLS approach towards the EIV-Model requires a *nonlinear treatment* of the *nonlinear model*. A number of different algorithms have been presented to generate the Total Least-Squares Solution from a certain set of nonlinear normal equations. A triplet of conventional nonlinear models has also been considered, suggesting that the LEast-Squares Solutions from iterative linearization do converge to the nonlinear TLS-Solution in all three cases. Most of the details, however, will be published in a forthcoming OSU-Report, due to the space restrictions for these Proceedings.

**Acknowledgment** This author is very much indebted to the thorough reading of the text by three anonymous reviewers. Their recommendations have improved the readability quite substantially. Also appreciated are many discussions with his long-time collaborator Kyle Snow.

---

## References

- Deming WE (1931) The application of least squares. *Phil Mag* 11:146–158
- Deming WE (1934) On the application of least squares II. *Phil Mag* 17:804–829
- Donevska S, Fišerova E, Horn K (2011) On the equivalence between orthogonal regression and linear model with type II constraints. *Acta Univ Palacki Olomuc, Fac rer nat, Mathematica* 50(2):19–27
- Fang X (2011) Weighted Total Least-Squares solutions with applications in geodesy. Publ. No. 294, Department of Geodesy and Geoinformatics, Leibniz University, Hannover
- Fang X (2013) Weighted Total Least-Squares: necessary and sufficient conditions, fixed and random parameters. *J Geodesy* 87:733–749
- Golub GH, van Loan CF (1980) An analysis of the Total Least-Squares problem. *SIAM J Numer Anal* 17:883–893
- Helmert FR (1907) Adjustment computations based on the Least-Squares Principle (in German), 2nd edn. Teubner, Leipzig
- Magnus JR, Neudecker H (2007) Matrix differential calculus with applications in statistics and economics, 3rd edn. Wiley, Chichester
- Mahboub V (2012) On weighted Total Least-Squares for geodetic transformations. *J Geodesy* 86:359–367
- Neitzel F (2010) Generalization of Total Least-Squares on example of unweighted and weighted 2D similarity transformation. *J Geodesy* 84:751–762
- Pope A (1972) Some pitfalls to be avoided in the iterative adjustment of nonlinear problems. In: Proceedings of the 38th Ann. ASPRS Meeting, Amer. Soc. of Photogrammetry: Falls Church, pp 449–477
- Schaffrin B, Wieser A (2008) On weighted Total Least-Squares adjustment for linear regression. *J Geodesy* 82:415–421
- Schaffrin B (2015) The Errors-In-Variables (EIV) model. Nonlinear Total Least-Squares (TLS) adjustment or iteratively linearized Least-Squares (LS) adjustment? OSU-Report, Division of Geodetic Science, School of Earth Sciences, The Ohio State University, Columbus
- Schaffrin B, Snow K, Neitzel F (2014) On the Errors-In-Variables model with singular covariance matrices. *J Geodetic Sci* 4:28–36
- Snow K (2012) Topics in Total Least-Squares adjustment within the Errors-In-Variables model: singular covariance matrices and prior information, Report No. 502, Division of Geodetic Science, School of Earth Sciences, The Ohio State University, Columbus
- Xu P, Liu JN, Shi C (2012) Total Least-Squares adjustment in partial Errors-In-Variables models. Algorithms and statistical analysis. *J Geodesy* 86:661–675

---

# Multivariate GNSS Attitude Integrity: The Role of Affine Constraints

Gabriele Giorgi and Peter J.G. Teunissen

---

## Abstract

In this work we analyze the integrity properties of an affine-constrained estimator applied to arrays of GNSS antennas. GNSS pseudorange and carrier phase measurements from multiple antennas whose relative positions are known are cast in a linearly-constrained observation model. The linear constraints are inherent to an affine transformation that is applied to the baseline coordinates. The affine transformation yields enhanced redundancy, thus improving the model integrity properties with respect to the unconstrained model. The extent of the improvement is measured in terms of internal and external reliability.

---

## Keywords

Affine constrained attitude model • Attitude determination • Galileo • GNSS • GPS • ILS • MDB • Multivariate ILS • Reliability

---

## 1 Introduction

GNSS carrier phase attitude determination enables precise estimations of a body orientation in space, see e.g. Bar-Itzhack et al. (1998), Cohen (1992) and Giorgi (2011). In order to derive the orientation of the body with respect to a reference frame, an array of GNSS antennas is employed. The distances between antennas are surveyed a priori, such that the coordinates of each baseline are known in a local reference frame (e.g., a frame integral with the body). The estimation of the orientation of the local frame with respect to the frame in which the GNSS measurements are expressed is the aim of attitude determination.

Carrier phase-based attitude estimations are characterized by a much higher estimation precision than code-

only processing, but the carrier phase measurements are inherently ambiguous by an unknown number of integer cycles (Strang and Borre 1997; Teunissen and Kleusberg 1998). These have to be resolved to their correct integer value. Application of Integer Least-Squares (ILS) to resolve the integer ambiguities guarantees the highest possible success rate among the admissible integer estimators, see Teunissen (1993) and Teunissen (1999). ILS methods have been recently studied in the context of nonlinear constrained models such as the GNSS-based attitude estimation problem (Giorgi 2011; Giorgi et al. 2012; Teunissen 2007). The baselines formed by the GNSS antenna array can be re-parameterized in terms of an attitude matrix. An admissible attitude matrix is orthonormally-constrained (OC) (Shuster 1993). If these constraints are integrated in the estimation of the whole parameter space, i.e., the attitude matrix and the integer ambiguities, a nonlinear estimation problem has to be solved. It is shown in Giorgi et al. (2011) and Teunissen (2007) how to solve for the OC GNSS attitude model, and it is demonstrated how the rigorous inclusion of the orthonormality constraint largely enhances ambiguity resolution.

A model of intermediate strength between the unconstrained (UC) and the OC approaches is obtained by operating the re-parameterization of the baselines in terms of

---

G. Giorgi (✉)  
Institute for Communications and Navigation, Technische Universität München, München, Germany  
e-mail: [gabriele.giorgi@tum.de](mailto:gabriele.giorgi@tum.de)

P.J.G. Teunissen  
Global Navigation Satellite Systems Research Centre, Curtin University of Technology, Perth, Australia  
e-mail: [P.Teunissen@curtin.edu.au](mailto:P.Teunissen@curtin.edu.au)

the attitude matrix, but neglecting the associated nonlinear constraints. This yields an affine-constrained (AC) GNSS attitude model, as introduced in Giorgi and Teunissen (2013) and Teunissen (2012). This model has the advantage that it avoids the computational complexity of the OC GNSS attitude model, while it still has a significantly improved ambiguity resolution performance over its unconstrained counterpart.

More than improving the ambiguity resolution performances, the AC attitude model is characterized by enhanced integrity properties with respect to the UC model. Due to the affine-transformation of a parameter subset, a larger observations-to-unknowns redundancy enables smaller minimum detectable biases (MDBs), i.e., the minimum bias magnitude that can be detected, and smaller bias-to-noise ratios (BNRs), i.e., a measure of the influence of undetected errors on the parameter estimation. The MDBs and BNRs are used to measure the internal and external reliability of the observation model, respectively. It is shown in this contribution the impact that adopting the AC model over the UC model has on the model integrity properties.

This contribution is structured as follows. Section 2 introduces a theorem on the integrity properties of affine-constrained linear models.

Section 3 formulates the GNSS functional and stochastic model in multivariate form, for one-, two- and three-dimensional antenna arrays, tracking GNSS signals on an arbitrary number of frequencies with two or more antennas.

Section 4 analyzes the impact of the affine constraints on the model internal and external reliability. The MDBs and BNRs are given in analytical form, and the difference between the UC and AC approaches is analyzed from a theoretical as well as a numerical standpoint.

## 2 Reliability and Constraints: A Theorem

The Minimal Detectable Bias (MDB) and the Bias-to-Noise-Ratio (BNR) are two important diagnostic measures that describe the model reliability after statistical testing, see Baarda (1968) and Teunissen (2006). Let the  $m$ -vector of observables  $y$  be distributed under the hypothesis  $\mathcal{H}$  as  $y \sim N(E(y|\mathcal{H}), Q_{yy})$ , and consider the null-hypothesis  $\mathcal{H}_U^0$  and the alternative hypothesis  $\mathcal{H}_U^a$ ,

$$E(y|\mathcal{H}_U^0) = Ax \quad \text{and} \quad E(y|\mathcal{H}_U^a) = Ax + c_y \nabla \quad (1)$$

in which  $A \in \mathbb{R}^{m \times n}$  and  $c_y \in \mathbb{R}^m$  are given (full rank) matrices, and  $x$  and  $\nabla$  are unknown parameters. The suffix ‘U’ has been used to discriminate this pair of hypotheses from their constrained counterparts,  $\mathcal{H}_C^0$  and  $\mathcal{H}_C^a$ , respectively (see 3).

In Chapter 4 of Teunissen (2006) it is shown that the MDB and BNR of uniformly most powerful invariant (UMPI) testing  $\mathcal{H}_U^0$  against  $\mathcal{H}_U^a$  are given as

$$\begin{aligned} \text{MDB}_U &= \sqrt{\frac{\lambda_0}{\|P_A^\perp c_y\|_{Q_{yy}}^2}}, \\ \text{BNR}_U &= \text{MDB}_U \|P_{Ac_y}\|_{Q_{yy}} \end{aligned} \quad (2)$$

with noncentrality parameter  $\lambda_0$ , the weighted squared-norm  $\|\cdot\|_{Q_{yy}}^2 = (\cdot)^T Q_{yy}^{-1} (\cdot)$  and the orthogonal projector  $P_A = A(A^T Q_{yy}^{-1} A)^{-1} A^T Q_{yy}^{-1}$ .

The following theorem shows how the MDBs and BNRs change when both the null- and alternative hypothesis are strengthened with the same constraints.

**Theorem (Constraints, MDBs and BNRs)** *Let the hypotheses of (1) be constrained as*

$$\begin{aligned} E(y|\mathcal{H}_C^0) &= Ax, & K^T x &= 0 \\ E(y|\mathcal{H}_C^a) &= Ax + c_y \nabla, & K^T x &= 0 \end{aligned} \quad (3)$$

with given full rank constraint matrix  $K \in \mathbb{R}^{n \times p}$  ( $p \leq n$ ). Then the unconstrained MDBs and BNRs of (1) are related to their constrained counterparts of (3) as

$$\frac{\text{MDB}_U^2}{\text{MDB}_C^2} = 1 + \frac{\|P_L^\perp c_{\hat{x}}\|_{Q_{\hat{x}\hat{x}}}^2}{\|P_A^\perp c_y\|_{Q_{yy}}^2} \quad (4)$$

$$\frac{\text{BNR}_U^2}{\text{BNR}_C^2} = \frac{\text{MDB}_U^2}{\text{MDB}_C^2} \left[ 1 - \frac{\|P_L^\perp c_{\hat{x}}\|_{Q_{\hat{x}\hat{x}}}^2}{\|P_{Ac_y}\|_{Q_{yy}}^2} \right]^{-1} \quad (5)$$

with basis matrix  $L$  spanning the null space of  $K^T$ , and

$$\begin{aligned} P_L^\perp &= I - L(L^T Q_{\hat{x}\hat{x}}^{-1} L)^{-1} L^T Q_{\hat{x}\hat{x}}^{-1} \\ &= Q_{\hat{x}\hat{x}} K(K^T Q_{\hat{x}\hat{x}} K)^{-1} K^T \\ c_{\hat{x}} &= Q_{\hat{x}\hat{x}} A^T Q_{yy}^{-1} c_y \\ Q_{\hat{x}\hat{x}} &= (A^T Q_{yy}^{-1} A)^{-1} \end{aligned} \quad (6)$$

*Proof* From the definition of the MDB and its geometric interpretation (Teunissen 2006), we have

$$\frac{\text{MDB}_U^2}{\text{MDB}_C^2} = \frac{\|P_{AL}^\perp c_y\|_{Q_{yy}}^2}{\|P_A^\perp c_y\|_{Q_{yy}}^2} = 1 + \frac{(\|P_{Ac_y}\|_{Q_{yy}}^2 - \|P_{AL} c_y\|_{Q_{yy}}^2)}{\|P_A^\perp c_y\|_{Q_{yy}}^2} \quad (7)$$

With

$$\begin{aligned} Q_{\hat{x}\hat{x}} &= (A^T Q_{yy}^{-1} A)^{-1} \\ P_A &= A Q_{\hat{x}\hat{x}} A^T Q_{yy}^{-1} \\ P_{AL} &= AL(L^T Q_{\hat{x}\hat{x}}^{-1} L)^{-1} L^T A^T Q_{yy}^{-1} \\ P_L &= L(L^T Q_{\hat{x}\hat{x}}^{-1} L)^{-1} L^T Q_{\hat{x}\hat{x}}^{-1} \end{aligned} \quad (8)$$

it follows that

$$\begin{aligned} \|P_{Ac_y}\|_{Q_{yy}}^2 &= \|c_{\hat{x}}\|_{Q_{\hat{x}\hat{x}}}^2 \\ \|P_{AL} c_y\|_{Q_{yy}}^2 &= \|P_L c_{\hat{x}}\|_{Q_{\hat{x}\hat{x}}}^2 \end{aligned} \quad (9)$$

Substitution of (9) into (7) proves the result. The derivation of the BNR-ratio goes along similar lines.  $\square$

This theorem shows how the MDBs and BNRs improve (i.e. get smaller) when constraints are added to both the null- and alternative hypothesis. Such improvement will be absent if  $\|P_L^\perp c_{\hat{x}}\|_{Q_{\hat{x}\hat{x}}}^2 = 0$ , i.e. if  $c_{\hat{x}} = 0$  or  $P_L^\perp c_{\hat{x}} = 0$ . The first case occurs if the solution for  $x$  under  $\mathcal{H}_0^a$  is invariant for the modelling bias, i.e. when  $c_y$  is orthogonal to the range space of  $A$ . The second case occurs when  $K^T c_{\hat{x}} = 0$ , i.e. when the effect of the modeling bias on  $x$  is not felt in the constraint.

In the following we will apply the above theorem to the multivariate GNSS affine-constrained model.

### 3 Multivariate GNSS Observation Models

Consider an array of  $r + 1$  GNSS antennas forming  $r$  independent baselines. The  $2mf$  DD GNSS pseudorange and carrier phase observations obtained by simultaneously tracking  $m + 1$  satellites on  $f$  frequencies, are cast in the columns of a  $2mf \times r$  observation matrix  $Y$ . The DD multivariate UC GNSS observation model is formulated as

$$E(Y) = GB + NX \quad , \quad D(\text{vec}Y) = P \otimes Q_{yy} \quad (10)$$

with  $B \in \mathbb{R}^{3 \times r}$  ,  $X \in \mathbb{Z}^{mf \times r}$

The coordinates of the baselines forming the array are cast in the columns of matrix  $B$ , whereas  $X$  contains the integer-valued unknown ambiguities (the array is assumed small enough to neglect the atmospheric delays). The entries of the  $2mf \times 3$  matrix  $G$  are the differenced line-of-sight vectors, and the  $2mf \times mf$  matrix  $N$  contains the wavelengths.

The matrices  $P$  and  $Q_{yy}$  of the dispersion  $D(\text{vec}Y)$  are given as  $P = \frac{1}{2} D_r^T D_r = \frac{1}{2} (I_r + e_r e_r^T)$  and  $Q_{yy} = 2\Sigma \otimes D_m^T D_m$ , with cofactor matrices  $\Sigma = \text{blockdiag}[\Sigma_P, \Sigma_\Phi]$ ,  $\Sigma_P = \text{diag}(\sigma_{p_1}^2, \dots, \sigma_{p_f}^2)$ ,  $\Sigma_\Phi = \text{diag}(\sigma_{\phi_1}^2, \dots, \sigma_{\phi_f}^2)$  containing the undifferenced code and phase variances, and where  $D_t^T = [-e_t, I_t]$  is an  $t \times (t + 1)$  differencing matrix.

We define a local frame in which the baseline coordinates are invariant, and introduce a  $q \times r$  matrix  $F$  whose entries are the local baseline coordinates. Parameter  $q$  denotes the rank of matrix  $F$ :  $q = 1$  for configurations of  $r$  antennas aligned in the same direction,  $q = 2$  for configurations of  $r$  coplanar antennas and  $q = 3$  for configurations of  $r$  non-coplanar antennas. The transformation between the local coordinate system and the reference system in which the observations are collected is defined by a rotation matrix  $R$  as

$$B = RF \quad ; \quad R \in \mathbb{O}^{3 \times q} \quad , \quad F \in \mathbb{R}^{q \times r} \quad (11)$$

The attitude matrix belongs to the class of  $3 \times r$  orthonormal matrices  $\mathbb{O}$ , i.e., matrix  $R$  fulfill the nonlinear constraints defined by  $R^T R = I_q$ . Substitution of relationship (11) into model (10) gives the OC attitude model:

$$E(Y) = GRF + NX \quad , \quad D(\text{vec}Y) = P \otimes Q_{yy} \quad (12)$$

with  $R \in \mathbb{O}^{3 \times q}$  ;  $X \in \mathbb{Z}^{mf \times r}$

The solution of the OC model is inherently more complex than the UC model, due to the nonlinear constraints. The solution of model (12) has been given in Giorgi (2011), Giorgi et al. (2012), Giorgi et al. (2011), Teunissen (2007), and will not be discussed any further in this work.

A model of intermediate complexity is obtained by adopting the transformation (11), but disregarding the nonlinear constraints. The AC GNSS attitude model is then formulated as (Teunissen 2012)

$$E(Y) = GRF + NX \quad , \quad D(\text{vec}Y) = P \otimes Q_{yy} \quad (13)$$

with  $R \in \mathbb{R}^{3 \times q}$  ;  $X \in \mathbb{Z}^{mf \times r}$

Hence, the orthonormality constraint  $R \in \mathbb{O}^{3 \times q}$  has been replaced by the constraint  $R \in \mathbb{R}^{3 \times q}$ . Note that the transformation (11) allows to reduce the number of unknown parameters in (13) for those antenna configurations in which the span of the baselines is smaller than their number ( $q < r$ ).

Since model (13) is linear, a classic ILS solution can be derived, as shown in Giorgi et al. (2011), Teunissen (2007) and Teunissen (2012).

#### 3.1 Alternative Hypotheses

Models (10) and (13) will be treated as our null hypotheses  $\mathcal{H}_{uc}^0$  and  $\mathcal{H}_{ac}^0$ , respectively. In order to check for errors and/or biases in the observation vector, or model misspecifications, we introduce two alternative functional models that will be compared to the representations (10) and (13). These alternative hypotheses are

$$\begin{aligned} \mathcal{H}_{uc}^a : E(Y) &= GB + NX + C\gamma, \quad D(\text{vec}Y) = P \otimes Q_{yy} \\ \mathcal{H}_{ac}^a : E(Y) &= GRF + NX + C\gamma, \quad D(\text{vec}Y) = P \otimes Q_{yy} \end{aligned} \quad (14)$$

The scalar  $\gamma \in \mathbb{R}$  denotes the error/bias magnitude, whereas matrix  $C$  defines the observation(s) affected by the error. We assume that matrix  $C$  can be expressed as the (outer) product between two vectors:  $C = cd^T$ , in which the  $2mf$ -vector  $c$  can be used to specify which observable is biased, and the  $r$ -vector  $d$  can be used to select the antenna of the biased observable. Thus, for instance, if the  $i$ th observable, of the

$j$ th satellite, of the  $k$ th antenna is assumed biased, then  $c$  and  $d$  are chosen as  $c = (I_{2f} \otimes D_m^T)(u_i^{2f} \otimes u_j^{m+1})$  and  $d = D_r^T u_k^{r+1}$ , in which  $u_i^n$  denotes a canonical unit vector of dimension  $n$ , with the 1 in its  $i$ th slot.

In the following section we analyze the integrity properties of the UC and AC observation models.

## 4 Integrity Properties

The MDBs and BNRs of models (10)–(13) tested against the corresponding alternative hypotheses in (14) follows from (2) as

$$\begin{aligned} MDB_{uc} &= \sqrt{\frac{v_0}{\sigma_{\gamma_{uc}}^2}} = \sqrt{\frac{v_0}{d^T P^{-1} d \bar{c}_{N,\bar{G}}^T Q_{yy}^{-1} \bar{c}_{N,\bar{G}}}} \\ BNR_{uc} &= MDB_{uc} \sqrt{d^T P^{-1} d (c^T Q_{yy}^{-1} P_N c + c^T Q_{yy}^{-1} P_{\bar{G}} c)} \end{aligned} \quad (15)$$

and

$$\begin{aligned} MDB_{ac} &= \sqrt{\frac{v_0}{\sigma_{\gamma_{ac}}^2}} = \sqrt{\frac{v_0}{\sigma_{\gamma_{uc}}^2 + d^T P_S P^{-1} d \bar{c}_N^T P_{\bar{G}} \bar{c}_N}} \\ BNR_{ac} &= \frac{MDB_{ac}}{\sqrt{d^T P^{-1} d c^T Q_{yy}^{-1} P_N c + d^T P_S^\perp P^{-1} d c^T Q_{yy}^{-1} P_{\bar{G}} c}} \end{aligned} \quad (16)$$

with

$$\begin{aligned} \bar{c}_{N,\bar{G}} &= [I - P_{\bar{G}}] \bar{c}_N = [I - P_{\bar{G}}][I - P_N] c \\ \bar{G} &= [I - P_N] G \\ P_{\bar{G}} &= \bar{G} [\bar{G}^T Q_{yy}^{-1} \bar{G}]^{-1} \bar{G}^T Q_{yy}^{-1} \\ P_N &= N [N^T Q_{yy}^{-1} N]^{-1} N^T Q_{yy}^{-1} \\ P_S &= I - P_S^\perp = I - P^{-1} F^T (F P^{-1} F^T)^{-1} F \end{aligned} \quad (17)$$

Matrix  $P_S$  is the projector of rank  $r - q$  that projects onto the null space of the body frame baseline matrix  $F$ . This matrix reduces to the zero-matrix when the number of baselines  $r$  equals their span  $q$ .

Application of the Theorem given in Sect. 2 to the UC and AC observations models yields the following ratios between MDBs and between BNRs:

$$\omega = \frac{MDB_{uc}}{MDB_{ac}} = \sqrt{1 + \frac{d^T P_S P^{-1} d \bar{c}_N^T Q_{yy}^{-1} P_{\bar{G}} \bar{c}_N}{d^T P^{-1} d \bar{c}_N^T Q_{yy}^{-1} P_{\bar{G}}^\perp \bar{c}_N}} \quad (18)$$

and

$$\tau = \frac{BNR_{uc}}{BNR_{ac}} = \frac{\omega}{\sqrt{1 - \frac{d^T P_S P^{-1} d c^T Q_{yy}^{-1} P_{\bar{G}} c}{d^T P^{-1} d (c^T Q_{yy}^{-1} P_N c + c^T Q_{yy}^{-1} P_{\bar{G}} c)}}} \quad (19)$$

with  $P_{\bar{G}}^\perp = I - P_{\bar{G}}$ . For  $q = r$  the MDBs and BNRs of the two models UC and AC are equal. However, for  $q < r$  the AC model is characterized by smaller MDBs and BNRs, or equivalently, the AC model is capable of detecting the same error magnitude with a larger power of detection. The extent of the improvement, quantified by the ratios  $\omega$  and  $\tau$ , depends on the number of baselines ( $r$ ), their relative geometry ( $P_S$ ), the measurement quality  $Q_{yy}$ , and the given satellite geometry distribution ( $G$ ). Note that the ratios in (18)–(19) are independent from a geometrical scaling of the whole GNSS antenna array.

### 4.1 Numerical Example

We provide in this section a numerical example of the improvement obtained in terms of MDBs and BNRs ratios in two single-constellation case studies. We study both GPS and Galileo signals, assuming 4 to 10 satellites for each case, tracked on frequency L1, L2 and L5 (GPS), and E1-E5a-E5b-E5-E6 (Galileo), by an array of 4 coplanar antennas. The local baseline coordinates are

$$F = \begin{bmatrix} 1 & 0 & 0.5 \\ 0 & 1 & 1 \end{bmatrix} \quad (20)$$

The DD observation noise ( $Q_{yy}$ ) is composed by using the variances reported in Table 1, and the GPS and Galileo simulated satellite positions are illustrated in the skyplots of Figs. 1 and 2. The DD observations are formed by differencing with respect to satellite ‘1’.

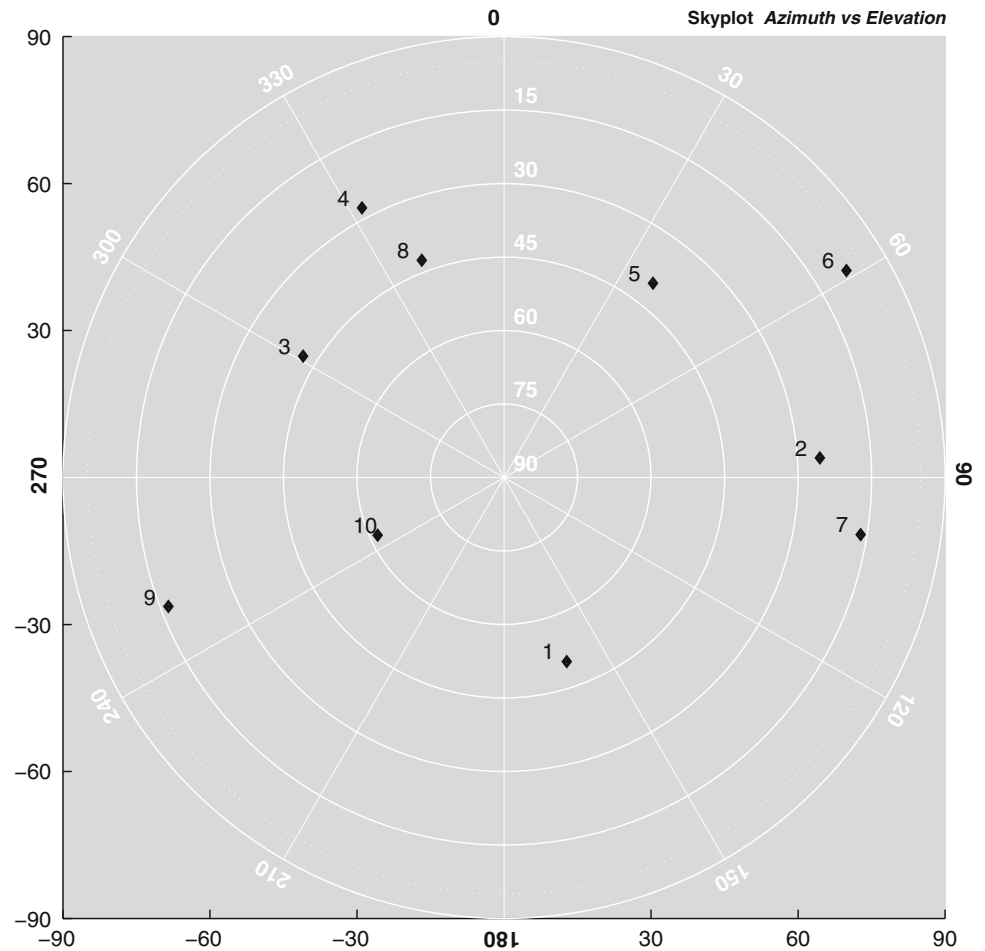
Table 2 reports several values of ratio  $\omega$  as function of the number of satellites and channels (frequencies) tracked by the antenna array. We assume an error occurring on the third receiver, relative to a single satellite (first DD observation) on frequency L5. The improvement obtained with the AC GNSS model with respect to the UC model is rather large. In the weakest measurement scenario, i.e., only 5 satellites tracked on the single-frequency L5, the minimum detectable bias is twice as large for the UC model compared to the AC model. Stronger scenarios, i.e., higher number of satellites and/or multifrequency observations, are characterized by lower values, although in many cases sensibly larger than the unit.

**Table 1** Standard deviations of GPS and Galileo undifferenced pseudorange and carrier phase ( $\Phi$ ) observables

	L1	L2	L5	E1	E5a	E5b	E5	E6
$\sigma_p$ (cm)	25	25	15	20	15	15	7	15
$\sigma_\phi$ (mm)	1.0	1.3	1.3	1.0	1.3	1.3	1.3	1.2



**Fig. 1** GPS simulated constellation, skyplot



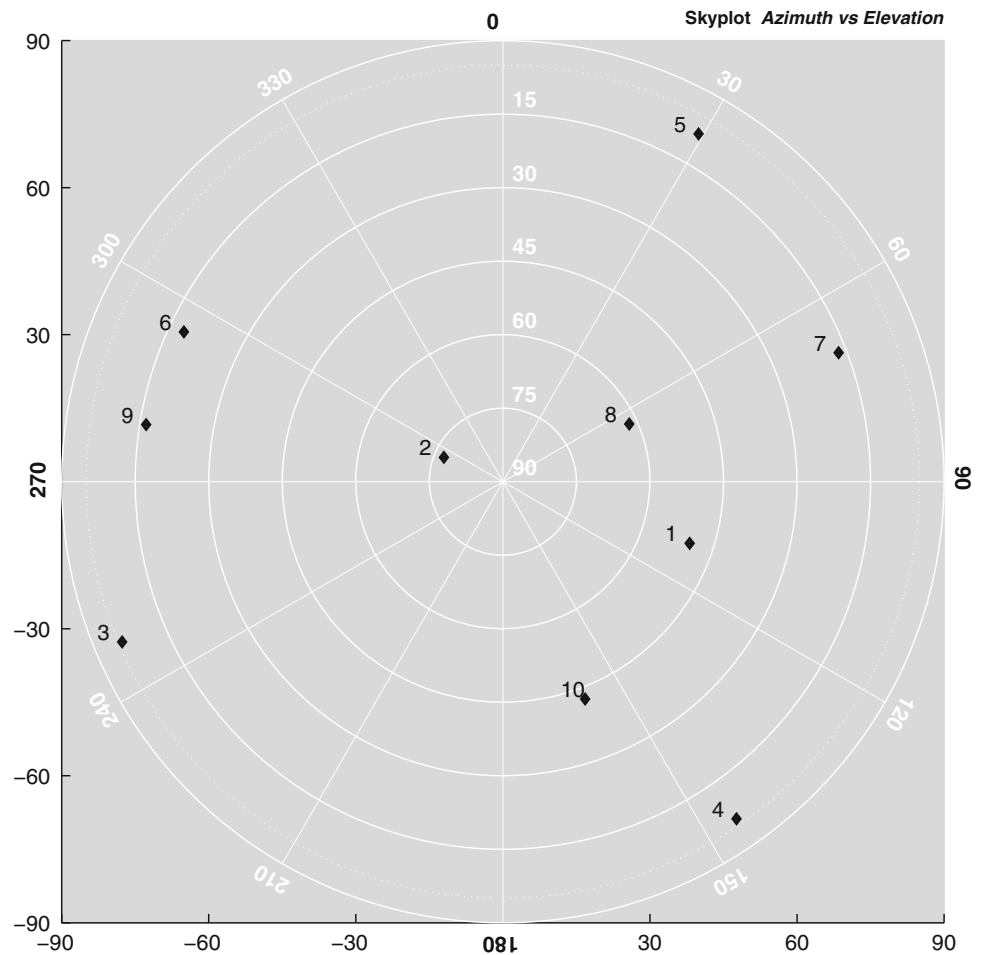
The same interpretation applies to the corresponding ratio  $\tau$  between BNRs, given in Table 3 as function of the number of satellites and channels. The improvement obtained by adopting the AC model is amplified with respect the MDBs ratio, with BNRs ratios as large as three for the weakest scenario. The impact of an undetected error in the parameter estimation is thus largely reduced when adopting the AC model.

The second case study focuses on a simulated Galileo constellation. Galileo signals are more precise than GPS signals (cfr. Table 1), and will generally provide smaller absolute values for the MDBs and BNRs. However, we analyze here the relative performance between the UC and

AC models, rather than the absolute improvement obtained by employing more precise observations.

Tables 4 and 5 reports the ratios between MDBs and BNRs for the Galileo scenario, relative to the detectability of an outlier on the third receiver, first observation at frequency L5. The larger gain associated to the AC model is again obtained in weaker scenarios, when limited number of satellites and/or single-frequency observations are available, with MDBs twice as large and BNRs three times as large when using the UC model. The extent of the improvement is then comparable to the one of the GPS case study, although with slight differences due to the satellite geometry and to the different accuracy between signals within the same GNSS.

**Fig. 2** Galileo simulated constellation, skyplot



**Table 2** Ratio  $\omega$  for GPS-only satellites, as function of number of satellites and frequencies available

# sat	L5	L1 + L5	L2 + L5	L1 + L2 + L5
4	–	1.58	1.58	1.32
5	2.06	1.38	1.38	1.24
6	1.10	1.07	1.07	1.05
7	1.06	1.04	1.04	1.03
8	1.06	1.04	1.04	1.03
9	1.06	1.04	1.04	1.03
19	1.06	1.04	1.04	1.03

The ratio  $\omega$  refers to the detectability of an outlier on the third receiver, first observation at frequency L5

**Table 3** Ratio  $\tau$  for GPS-only satellites, as function of number of satellites and frequencies available

# sat	L5	L1 + L5	L2 + L5	L1 + L2 + L5
4	–	2.31	2.31	1.93
5	3.02	2.03	2.03	1.81
6	1.61	1.56	1.56	1.54
7	1.55	1.52	1.52	1.51
8	1.55	1.53	1.53	1.51
9	1.55	1.52	1.52	1.51
10	1.55	1.52	1.52	1.51

The ratio  $\tau$  refers to the presence of an outlier on the third receiver, first observation at frequency L5

Overall, the strengthening of the observation model obtained by imposing linear constraints yields an effective reduction of the MDBs and BNRs, thus proving the

advantage of affine-constrained models in terms of reliability for those configurations of antennas whose relative positions can be modeled as spatially invariant.

**Table 4** Ratio  $\omega$  for Galileo-only satellites, as function of number of satellites and frequencies available

# sat	E1 + E5 E5a + E5b E1 + E5a					
	E5	E1 + E5	E5a + E5	+E6	+E5 + E6	+E5 + E6
4	–	2.32	1.86	1.60	1.35	1.40
5	1.44	1.34	1.28	1.24	1.17	1.18
6	1.34	1.27	1.23	1.20	1.14	1.15
7	1.35	1.28	1.24	1.20	1.14	1.16
8	1.17	1.14	1.13	1.11	1.09	1.09
9	1.16	1.14	1.12	1.11	1.08	1.09
10	1.14	1.12	1.11	1.09	1.07	1.08

The ratio  $\omega$  refers to the detectability of an outlier on the third receiver, first observation at frequency E5

**Table 5** Ratio  $\tau$  for Galileo-only satellites, as function of number of satellites and frequencies available

# sat	E1 + E5 E5a + E5b E1 + E5a					
	E5	E1 + E5	E5a + E5	+E6	+E5 + E6	+E5 + E6
4	–	3.39	2.72	2.34	1.97	2.05
5	2.11	1.95	1.88	1.81	1.71	1.73
6	1.96	1.85	1.80	1.75	1.67	1.69
7	1.98	1.87	1.81	1.76	1.68	1.69
8	1.72	1.68	1.65	1.63	1.59	1.60
9	1.70	1.67	1.64	1.62	1.58	1.59
10	1.67	1.64	1.62	1.60	1.57	1.58

The ratio  $\tau$  refers to the detectability of an outlier on the third receiver, first observation at frequency E5

## 5 Conclusions

In linear models, linear constraints on parameter subsets potentially enable improved error detection, thanks to an implicit enhanced observations-to-unknowns redundancy. A consistent improvement in terms of integrity properties can be obtained by adopting the affine-constrained model over the unconstrained formulation. The improvement was captured analytically by expressing the ratio between minimum detectable biases (and bias-to-noise ratios) obtained in the linearly-constrained and unconstrained models.

The theorem finds a direct application in arrays of GNSS antennas whose relative distances do not vary. Following a re-parameterization of the baseline coordinates in terms of an attitude matrix, the observation model becomes linearly constrained when disregarding the orthonormality of the attitude matrix. The GNSS attitude model with affine

constraints yields enhanced internal and external reliability, as demonstrated with several numerical examples that provide evidence of the improved performance of the affine-constrained model in terms of minimum detectable biases and bias-to-noise ratios.

**Acknowledgements** The research of Peter J.G. Teunissen has been supported by an Australian Research Council Federation Fellowship (project number FF0883188).

## References

- Baarda W (1968) A testing procedure for use in geodetic networks. Netherlands Geodetic Commission Publications on Geodesy, vol. 2, issue 5, 97 p
- Bar-Itzhack IY, Montgomery P, Garrick J (1998) Algorithm for attitude determination using global positioning system. *J Guid Control Dyn* 21(6):846–852
- Cohen CE (1992) Attitude determination using GPS. Ph.D. Thesis, Stanford University, Palo Alto, CA
- Giorgi G (2011) GNSS carrier phase-based attitude determination. Estimation and applications. Delft University of Technology
- Giorgi G, Teunissen PJG (2013) Low-complexity instantaneous ambiguity resolution with the affine-constrained GNSS attitude model. *IEEE Trans Aerosp Electron Syst* 49(3):1745–1759
- Giorgi G, Teunissen PJG, Verhagen S, Buist PJ (2011) Instantaneous ambiguity resolution in GNSS-based attitude determination applications: the MC-LAMBDA method. *J Guid Control Dyn* 35(1):51–67
- Giorgi G, Teunissen PJG, Verhagen S, Buist PJ (2012) Integer ambiguity resolution with nonlinear geometrical constraints. *IAG Symp* 137:39–45
- Shuster MD (1993) A survey of attitude representations. *J Astronaut Sci* 41(4):439–517
- Strang G, Borre K (1997) Linear algebra, geodesy, and GP. Wellesley-Cambridge Press, Wellesley
- Teunissen PJG (1993) Least squares estimation of the integer GPS ambiguities. Invited lecture, Section IV theory and methodology, IAG general meeting, Beijing also in: LGR series No 6, Delft Geodetic Computing Center, Delft University of Technology
- Teunissen PJG (1999) An optimality property of the integer least-squares estimator. *J Geod* 73(11):587–593
- Teunissen PJG (2006) Testing theory: an introduction. Series on mathematical geodesy and positioning, 2nd edn. Delft Academic Press, Orlando
- Teunissen PJG (2007) A general multivariate formulation of the multi-antenna GNSS attitude determination problem. *Artif Satell* 42(2):97–111
- Teunissen PJG (2012) The affine constrained GNSS attitude model and its multivariate integer least-squares solution. *J Geod* 86:547–563
- Teunissen PJG, Kleusberg A (1998) GPS for geodesy, 2nd edn. Springer/Heidelberg, Berlin/New York

---

# Integrating Geological Prior Information into the Inverse Gravimetric Problem: The Bayesian Approach

L. Rossi, M. Reguzzoni, D. Sampietro, and F. Sansò

---

## Abstract

It is well known that the inverse gravimetric problem is generally ill-posed and therefore its solution requires some restrictive hypotheses and strong numerical regularization. However, if these initial assumptions are improperly used, the final results could be theoretically and physically admissible but far from the actual mass density distribution. In this work, a Bayesian approach to estimate the mass density distribution from gravity data coupled with a-priori geological information is presented. It requires to model the masses in voxels, each of them characterized by two random variables: one is a discrete label defining the type of material (or the geological unit), the other is a continuous variable defining the mass density (considered constant inside the single voxel). The a-priori geological information is translated in terms of this model, providing for each class of material the mean density and the corresponding variability and for each voxel the a-priori most probable label. Basically the method consists in a simulated annealing aided by a Gibbs sampler with the aim to find the MAP (maximum a posteriori) of the posterior probability distribution of labels and densities given the observations and the a-priori geological model. Some proximity constraints between labels of adjacent voxels are also introduced into the solution.

The proposed Bayesian method is here tested on two simulated scenarios. In particular the first is an example of bathymetry recovering, while the second a salt dome shape estimation. These experiments show the capability of the method to correct the possible inconsistencies between the a-priori geological model and the gravity observations: 86% and 60% of wrong voxels have been corrected in the first and second test respectively.

---

## Keywords

Bayesian approach • Inverse gravimetric problem • Monte Carlo Markov Chain method

---

L. Rossi (✉) • M. Reguzzoni  
DICA, Politecnico di Milano, P.za Leonardo da Vinci 32, 20133  
Milano, Italy  
e-mail: [lorenzo1.rossi@polimi.it](mailto:lorenzo1.rossi@polimi.it)

F. Sansò  
DICA, Politecnico di Milano, Como Campus, via Valleggio 11,  
22100 Como, Italy

D. Sampietro  
GReD s.r.l., via Cavour 2, 22074 Lomazzo, Italy

---

## 1 Introduction

The intrinsic indetermination of the inverse gravimetric problem is well known and the description of the whole set of possible internal masses, given the external gravity potential, has been fully described on a purely mathematical ground (e.g. Parker 1975; Sampietro and Sansò 2012). However, in order to obtain realistic solutions, some constraints should be added in the solution of the inverse gravimetric problem. For instance the solution can be derived from the “experience” of an operator assisted by fast forward algorithms (Parker

1973; Caratori Tontini et al. 2009; Gordon et al. 2012) and from generic geological information by means of trial and error procedures.

Another possible solution could be to add severely constraints in terms of mass contrast leading to search for the geometry of discontinuity surfaces (Barbosa et al. 1997, 1999; Fedi and Rapolla 1999; Fedi 2006). This approach is commonly called non-linear inverse problem due to the non linearity of the functional relating gravity observations and geometrical parameters of the sources. On the contrary in the so-called linear inversion (Last and Kubik 1983; Guillen and Menichetti 1984; Barbosa and Silva 1994) there is a linear relation in terms of Newtonian integral between the mass density and the functional of the external gravity potential, which is usually described as a summation on volume elements (voxels). Considering the linear problem the relation between data and unknowns is univocal when the number of voxel is conveniently taken smaller than the number of observations. However this relation is highly unstable because, as the dimension of voxels decreases, we are approaching the continuous setting, where non-uniqueness is large, as recalled above. Therefore the solution is usually obtained by imposing proper constraints. This can be done either under deterministic models (Medeiros and Silva 1996) or stochastic ones (Tarantola and Valette 1982; Tarantola 2002). In any case this approach is reconducted to the optimization of some non-linear, often quadratic, functionals of the gravity observations and the unknown mass distribution. This optimization can be obtained by Monte Carlo Markov Chain methods, including simulated annealing (Nagihara and Hall 2001; Roy et al. 2005), as it is very well known in literature. Naturally the relation between sources and observations, i.e. the forward model, can be conveniently reckoned using a Fourier approach that greatly speeds up the computational time.

This paper is in the flow of the above way of reasoning, but trying to incorporate also the interactive approach mentioned at the beginning by modelling the geological information in a Bayesian mode as prior probability. This is already present in geophysical literature even coupling gravimetric and magnetic observations (e.g. Bosch 1999, 2004; Bosch and McGaughey 2001; Mosegaard and Tarantola 2002; Bosch et al. 2006; Guillen et al. 2008). In particular we propose here an approach similar to the one shown in Guillen et al. (2008) in which a field of discrete variables (namely geological units) is introduced as an additional unknown, with some prior information. As it will be explained in the following, the main differences with respect to Guillen et al. (2008) are in the way the prior information is formalized and in the algorithm used to find the solution of the inverse

problem. Note that this work represents only a preliminary study, mainly focused on the mathematical formalization of the problem and that the improvement of the method is still a matter of investigation. Wishing to estimate a MAP (Maximum A Posteriori) of our posterior distribution, we are facing an optimization problem with part of the variables which are discrete. The proposed solution resorts to an application of a Gibbs sampler combined with a simulated annealing (Smith and Roberts 1993; Sansò et al. 2011), as it can be found in a large part of literature; here the application of the method to the image analysis, with the seminal paper by Geman and Geman (1984), is worth being mentioned.

A remark however can be put forward already in this introduction, namely that while image analysis deals only with “local” observations, i.e. observations that solely depend on the pixel to be updated in the Gibbs sampler, in our case any variation of density at any point will instead affect all the observable gravity anomalies wherever they are.

## 2 Problem Formalization

Similarly to Guillen et al. (2008) the inversion algorithm is developed assuming that some geological information is available in the studied region. In details, we suppose to know a list of all the possible geological units present in the area and their approximate geometrical distribution (e.g. from geological sections). We also suppose to know for each geological unit the most probable density and its variability (e.g. from literature). However, while in Guillen et al. (2008) only the boundaries of the geological units can be modified, in case merging separated portions of features or removing isolated ones, in the proposed method the formalization of the prior probability allows a more general solution to the problem, e.g. the possibility to generate new features.

In the following we formalize these assumptions in a Bayesian scheme: we start from the Bayes theorem in the usual form (Bayes 1763; Box and Tiao 2011):

$$P(\mathbf{x}|\mathbf{y}) \propto \mathcal{L}(\mathbf{y}|\mathbf{x}) P(\mathbf{x}) \quad (1)$$

where  $\mathbf{y}$  is a vector of observable quantities, while  $\mathbf{x}$  is a vector of body parameters. The investigated volume is split into voxels,  $V_i$ , with index  $i = 1, 2, \dots, N$ ; each voxel will carry two parameters  $(\rho_i, L_i)$  where  $\rho_i$  is the voxel mass density and  $L_i$  is a “label” attributing to  $V_i$  the presence of a certain geological unit chosen from the a-priori archive (e.g. water, sediment, salt, rock of a given type, etc.). So  $\rho_i$  is a continuous variable and  $L_i$  a discrete one among the  $M$  integers denoting the various materials.

Crucial is the way in which the prior probability  $P(\mathbf{x})$  is supplied, namely the shape of the distribution  $P(\mathbf{x}) = P(L_1, \rho_1; L_2, \rho_2; \dots; L_N, \rho_N)$ . We assume that:

$$P(\mathbf{x}) = \prod_{i=1}^N P(\rho_i | L_i) \cdot P(\mathbf{L}) = \prod_{i=1}^N P(\rho_i | L_i) \cdot P(L_1, L_2, \dots, L_N) \quad (2)$$

meaning that, once a label  $L_i = \ell$  has been chosen for  $V_i$ , the corresponding density will follow the law  $P(\rho_i | L_i = \ell)$ , which in our case is a normal distribution:

$$P(\rho_i | L_i = \ell) \sim \mathcal{N}(\bar{\rho}_\ell, \sigma_\ell^2) \quad (3)$$

with the mean  $\bar{\rho}_\ell$  and the variance  $\sigma_\ell^2$  given by geological literature. In this respect a comprehensive set of rock properties can be found for instance in Christensen and Mooney (1995). As for the prior  $P(\mathbf{L}) \equiv P(L_1, L_2, \dots, L_N)$ , we assume to have a Gibbs distribution (Azencott 1988):

$$P(\mathbf{L}) \propto e^{-\mathcal{E}(\mathbf{L})} \quad (4)$$

where the energy  $\mathcal{E}(\mathbf{L})$  depends only on the values  $\ell_i^o$  of  $L_i$  provided by the geological model, as well as from cliques (Geman and Geman 1984) of order two expressing the fact that the value of  $L_i$  is more likely to be equal to the value of the labels of the nearest neighbour voxels according to the following rules:

$$P(L_i = \ell | \mathbf{L}_{\Delta_i}) \propto e^{-\gamma s^2(L_i, \ell_i^o) - \lambda \sum_{j \in \Delta_i} q^2(L_i, L_j)} \quad (5)$$

where  $\gamma, \lambda$  are parameters to be empirically tuned,

$$s^2(L_i, \ell_i^o) = s_i^2 = \begin{cases} 0 & \text{if } L_i = \ell_i^o \\ \alpha_i & \text{if } L_i \neq \ell_i^o \end{cases} \quad (6)$$

$$q^2(L_i, L_j) = q_{ij}^2 = \begin{cases} a_i & \text{if } L_i = L_j \\ a_{ij} & \text{if } L_i \neq L_j \end{cases} \quad (7)$$

with  $V_j \in \Delta_i$  and  $\Delta_i$  is the neighbourhood of the voxel  $V_i$  defined by the cliques of order two, as mentioned above.

Note that given the geological model it is possible to create a table of proximity of geological units and then, by tuning  $\alpha_i, a_i$  and  $a_{ij}$ , to create a hierarchy of the most probable values for  $L_i$ . For example supposing to have three units,  $\ell = \{1, 2, 3\}$ , and a proximity table as the

	1	2	3
1	X	X	
2	X	X	X
3		X	X

**Fig. 1** Example of proximity table. The geological unit 1 can be close to unit 2, but not to unit 3

one presented in Fig. 1, this translates into the following definition:

$$s_i^2 = \begin{cases} 0 & \text{if } L_i = \ell_i^o \\ \alpha & \text{if } L_i \text{ is a geological neighbour of } \ell_i^o \\ \beta & \text{if } L_i \text{ is not a geological neighbour of } \ell_i^o \end{cases} \quad (8)$$

$$q_{ij}^2 = \begin{cases} a & \text{if } L_i = L_j \\ b & \text{if } L_i \text{ is a geological neighbour of } L_j \\ c & \text{if } L_i \text{ is not a geological neighbour of } L_j \end{cases} \quad (9)$$

with  $\beta > \alpha > 0$  and  $c > b > a$ .

Summarizing, the geological information enters into the solution providing the set of the possible geological units (i.e. the possible labels) with their mean density and its variability, the neighborhood relationship between the different geological units and the most probable value  $\ell_i^o$  of each voxel. All these data can be derived from basin geological studies (e.g. geological sections or maps) or through geophysical techniques.

Two remarks are in order: the first is that  $\mathbf{L}$ , with prior  $P(\mathbf{L})$ , is indeed a Markov random field (MRF), see Rozanov (1982). The second is that the final result of our optimization will depend from the chosen value of all the constants, which have to be tuned on the specific example.

As always for a MRF, the characteristics, namely the conditional distributions (5), determine a joint distribution  $P(\mathbf{L})$  such that:

$$\log P(\mathbf{L}) \propto -\frac{1}{2}\gamma \sum_{i=1}^N s^2(L_i, \ell_i^o) - \frac{1}{2}\lambda \sum_{i=1}^N \sum_{j \in \Delta_i} q^2(L_i, L_j). \quad (10)$$

The logarithm of the posterior distribution (1) will be written as:

$$\log P(\mathbf{x}|\mathbf{y}) = \log P(\boldsymbol{\rho}, \mathbf{L} | \Delta \mathbf{g}^o) \propto \propto -\frac{1}{2} (\Delta \mathbf{g}^o - \mathbf{A}\boldsymbol{\rho})^T \mathbf{C}_{\Delta \mathbf{g}^o}^{-1} (\Delta \mathbf{g}^o - \mathbf{A}\boldsymbol{\rho}) +$$

$$\begin{aligned}
& -\frac{1}{2}(\boldsymbol{\rho} - \bar{\boldsymbol{\rho}})^T \mathbf{C}_\rho^{-1} (\boldsymbol{\rho} - \bar{\boldsymbol{\rho}}) - \frac{1}{2}\gamma \sum_{i=1}^N s^2(L_i, \ell_i^o) + \\
& -\frac{1}{2}\lambda \sum_{i=1}^N \sum_{j \in \Delta_i} q^2(L_i, L_j) \quad (11)
\end{aligned}$$

where we recall that  $\Delta \mathbf{g}^o$  is the vector of observed gravity anomalies,  $\mathbf{C}_{\Delta g}$  its noise covariance matrix,  $\mathbf{A}$  is the forward modelling operator from densities to gravity anomalies,  $\boldsymbol{\rho}$  and  $\bar{\boldsymbol{\rho}}$  the vectors of components  $\rho_i$  and  $\bar{\rho}_i = \bar{\rho}(\ell_i)$ ,  $\mathbf{C}_\rho$  the corresponding covariance matrix and  $s^2(L_i, \ell_i^o)$ ,  $q^2(L_i, L_j)$  given by (6) and (7). This is the target function we want to maximize with respect to  $\rho_i$  and  $L_i$ .

The maximization of (11), due to the fact that some variables are discrete, is never an easy task, as we know from other important problems in geodesy, e.g. the GNSS initial phase ambiguity fixing (De Lacy et al. 2002). The idea, mutated from image analysis, is to apply a Gibbs sampler, chained with a simulated annealing (Casella and Robert 1999). In order to apply it to both the variables ( $\rho_i, L_i$ ), which are functions of the voxel  $V_i$ , we have simplified the problem by considering  $\rho_i$  as a discrete variable too. In practice we have substituted the normal distribution (3) with a discrete distribution on  $K$  values, e.g. on five argumental values taken at the average  $\bar{\rho}_\ell$ , and at  $\bar{\rho}_\ell \pm \sigma_\ell$ ,  $\bar{\rho}_\ell \pm 2\sigma_\ell$  respectively. Of course to each argument the proper probability is assigned, according to the normal law. Once this is done, the Gibbs sampler is applied by drawing one couple ( $\rho_i, L_i$ ) at a time, holding fixed all the other values and following a simple updating routine. The probabilities of the sampling are computed from (11) letting  $\rho_i$  run over its  $K$  values and  $\ell_i$  run over  $1, 2, \dots, M$ ; in this way we have a table of  $K \times M$  knots with their probabilities.

Actually the probability of  $\mathbf{x}$  is modulated by introducing a “temperature” parameter  $T$ :

$$P_T(\mathbf{x}) \propto e^{\frac{1}{T} \log P(\mathbf{x}|\mathbf{y})} \quad (12)$$

and  $T$  is slowly reduced at each step (e.g. by 5% of its value). In this way starting from a very large  $T$ , we obtain a sequence of samples converging in probability to the point  $\bar{\mathbf{x}}$  where the maximum of  $\log P(\mathbf{x}|\mathbf{y})$  is achieved (Azencott 1988).

### 3 Numerical Experiment

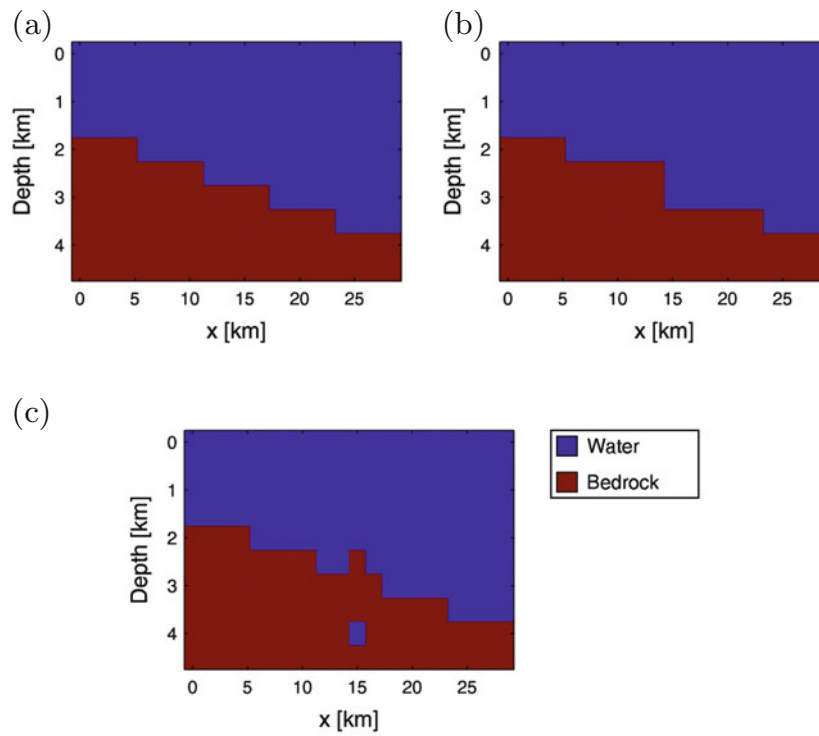
In order to assess the effectiveness of the presented Bayesian approach, which is able to consider also qualitative geological information, two simple experiments are carried out. They consist in recovering the mass density distribution of 3D synthetic models from their gravitational field. The density of each voxel is assumed to be equal to the mean

density of the associated geological unit and moreover the model is assumed constant along one planar direction, i.e. all the vertical cross sections in this direction are equal. From this reference model the two inputs of the inversion algorithm, i.e. the gravitational signal and the approximate geological model, are simulated. In particular the latter is obtained by slightly modifying the labels of the reference model. The inversion algorithm is therefore applied and the result is compared with the reference model in a closed-loop test.

In this work we will present two numerical examples: the first simulates the recovering of a bathymetry, while the second consists in recovering the shape of a salt dome.

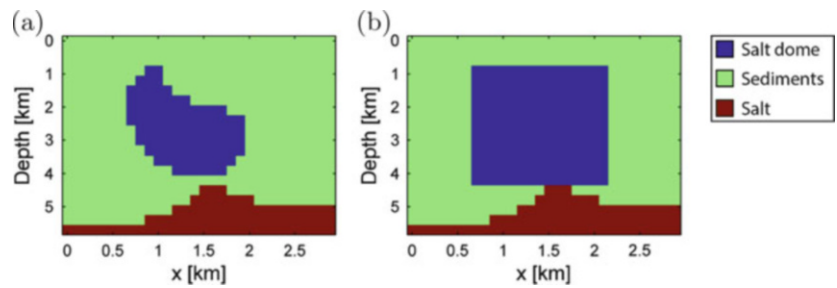
In the bathymetry model only two geological units are considered, water and bedrock, defined by  $\rho_w = 1,000 \text{ kg m}^{-3}$ ,  $\sigma_w = 5 \text{ kg m}^{-3}$  and  $\rho_b = 2,900 \text{ kg m}^{-3}$ ,  $\sigma_b = 50 \text{ kg m}^{-3}$  respectively. The investigated area is a square of 30 km side and has a depth of 5 km. A vertical cross section of the synthetic model, displayed in terms of “labels”, is represented in Fig. 2a. The volume is modelled by means of 1,200 rectangular prisms, each of them of dimensions 1.5 km (x)  $\times$  5.0 km (y)  $\times$  0.5 km (z) and its gravitational observations are simulated by means of Nagy equations (see Nagy 1966) in a noiseless scenario. In particular the observations are generated on a regular grid at an altitude of 250 m and with a spatial resolution of 1 km, thus simulating the result of an aerogravimetric flight. As explained above, the geological model is simulated by slightly modifying the reference model as shown in Fig. 2b. The two parameters  $\lambda$  and  $\gamma$  are empirically set to the values of 0.833 and 0.733 respectively and finally the values of the labels are randomly initialized from a uniform distribution (i.e. drawn with an infinite temperature in the simulating annealing). The solution is obtained in about 5,000 iterations and about 4 h on a common personal computer. A vertical cross section of the resulting synthetic model is depicted in Fig. 2c showing how the error in the geological model is properly corrected. In fact 86% of the wrong labels are corrected and the error on density has a standard deviation of  $216 \text{ kg m}^{-3}$ .

In the salt dome experiment three geological units are considered: salt dome ( $\rho_{dome} = 2,000 \text{ kg m}^{-3}$ ,  $\sigma_{dome} = 50 \text{ kg m}^{-3}$ ), salt ( $\rho_{salt} = 2,700 \text{ kg m}^{-3}$ ,  $\sigma_{salt} = 50 \text{ kg m}^{-3}$ ) and sediments ( $\rho_{sed} = 3,000 \text{ kg m}^{-3}$ ,  $\sigma_{sed} = 50 \text{ kg m}^{-3}$ ). The volume is modelled by means of 2,400 voxels, each of them with size of 0.4 km (x)  $\times$  0.1 km (y)  $\times$  0.3 km (z). The investigated area has a planar size of 3 km  $\times$  2 km and has a depth of 6 km. The geological units of a vertical cross section of the synthetic model are shown in Fig. 3a. The gravitational signal is simulated using point masses into a white noise scenario (noise standard deviation  $\sigma_{\Delta g} = 1 \text{ mGal}$ ). The simulated geological model is shown in Fig. 3b. In Fig. 4 three examples of the prior distribution are depicted, thus showing its dependence from the function  $s^2$  and  $q^2$  defined in (8)

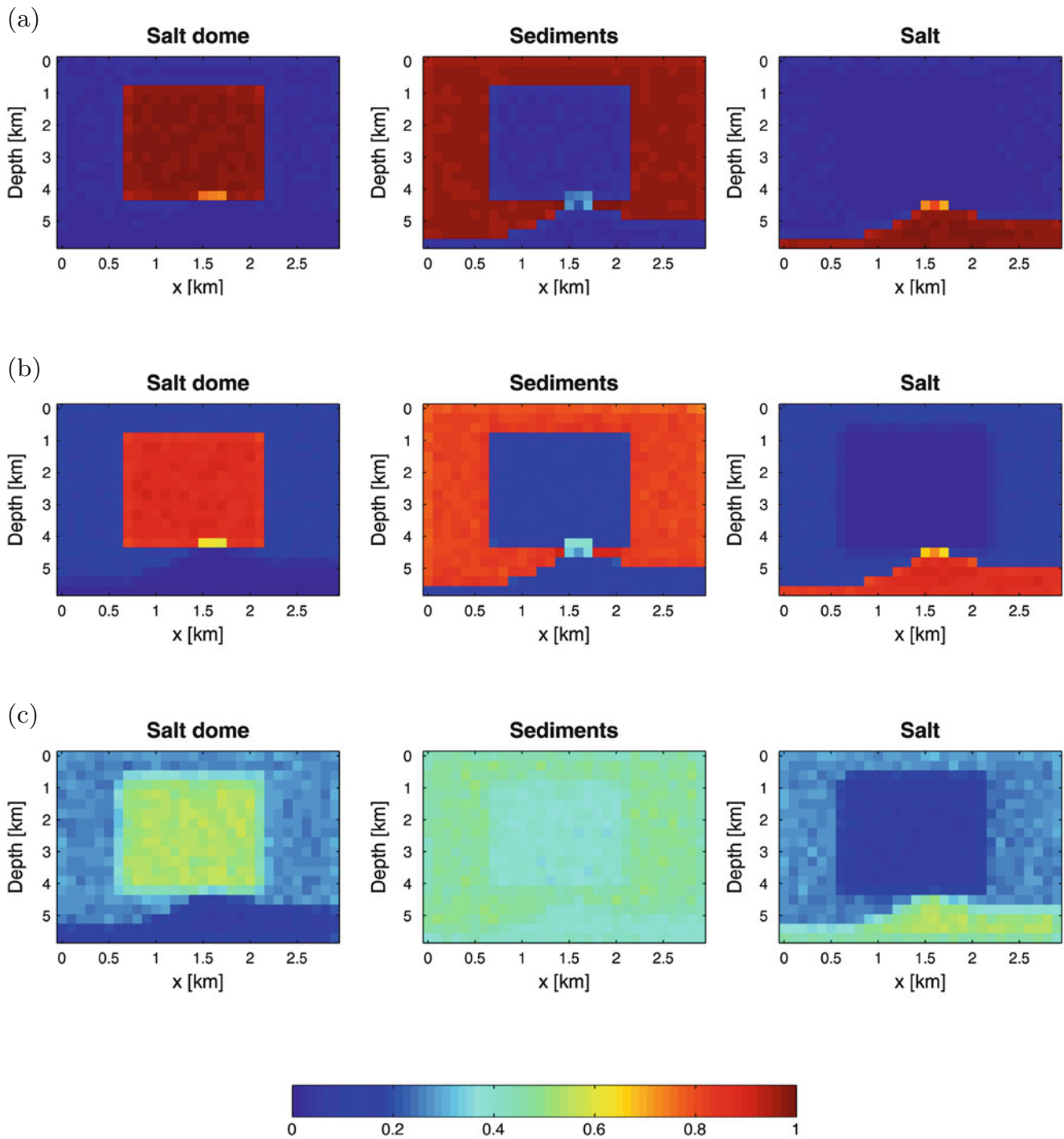


**Fig. 2** Vertical cross sections representing the geological units (“labels”) of the bathymetry test. (a) reference model; (b) geological model; (c) solution

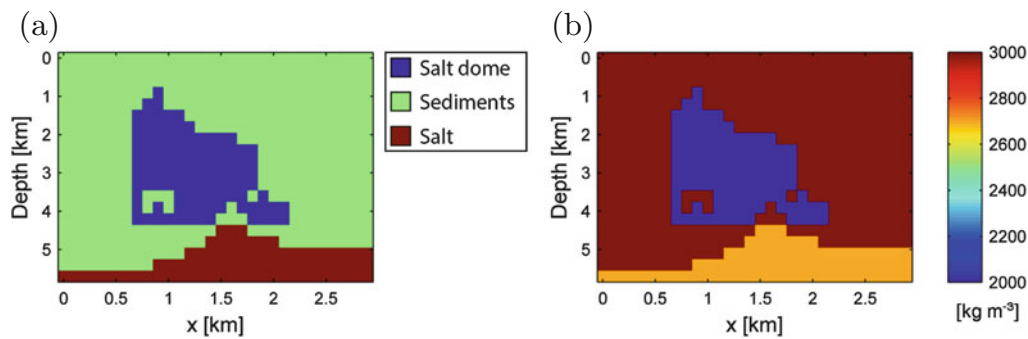
**Fig. 3** Vertical cross sections representing the geological units (“labels”) of inputs to the salt dome test. (a) reference model; (b) geological model







**Fig. 4** Vertical cross sections representing the relative frequency of each geological unit (“labels”) obtained from 2,000 realizations of the prior distribution. Each row is computed assuming different values of the prior parameters. (a)  $\gamma = 0.6$ ,  $\lambda = 0.03$ ,  $s^2 = \{0, 1, 10\} \forall i$  and  $q^2 = \{0, 1, 10\} \forall i, j$ ; (b)  $\gamma = 0.6$ ,  $\lambda = 0.03$ ,  $s^2 = \{0, 0.5, 5\} \forall i$  and  $q^2 = \{0, 0.5, 5\} \forall i, j$ ; (c)  $\gamma = 0.6$ ,  $\lambda = 0.03$ ,  $s^2 = \{0, 0.01, 0.1\} \forall i$  and  $q^2 = \{0, 0.5, 2\} \forall i, j$



**Fig. 5** Vertical cross sections representing the solution of the salt dome test. (a) geological units; (b) density

and (9). These sample distributions are obtained by counting the occurrences of each geological unit for each voxel and then computing the corresponding relative frequencies over 2,000 samples. From these three examples it can be noticed that the  $s^2$  function controls the “certainty” of the geological unit of each voxel (the closer are the numerical values of the parameters  $\alpha$  and  $\beta$  in (8), the more non-informative is the prior), while  $q^2$  is related to the “certainty” of the geological unit boundaries (the closer are the numerical values of  $a$ ,  $b$  and  $c$  in (9), the more unreliable are the boundaries). As for  $\lambda$  and  $\gamma$ , they are constants that practically controls the relative weight in the prior (2) between the density information and the geometrical one.

The solution is carried out by computing the prior fixing  $\lambda = 0.03$ ,  $\gamma = 0.6$ ,  $\alpha = 1$ ,  $\beta = 10$ ,  $a = 0$ ,  $b = 1$  and  $c = 10$ , see Fig. 4a, in about 2 h and 200 iterations and it is shown in Fig. 5. In this case the algorithm is able to recover about 60% of the wrong voxels and the error on density has a standard deviation of  $244 \text{ kg m}^{-3}$ . It can be seen from the salt dome experience that the algorithm is able to properly recover the shallowest part of the investigated volume, while the deepest one still present uncorrect features. This is probably due to the fact that the functions  $s^2$  and  $q^2$  are defined in the same way for the whole region, while a dependence at least on the vertical coordinate should be included.

## 4 Conclusions and Future Works

In the present paper a Bayesian approach to invert gravity data with the support of a given geological model has been studied. The method works properly at least in the performed preliminary test scenarios. Actually, the two main limiting factors are the choice of all the parameters playing a role in the formulation of the a-priori probability and the computational time.

In this respect it would be useful, in order to limit the impact of user decisions on the solution, to implement a semi-automatic determination of the optimal numerical val-

ues of the  $s^2$  and  $q^2$  functions and of the  $\lambda$  and  $\gamma$  parameters. These parameters in fact can modulate how close/far the final solution is from the geological model and from the gravity observations.

The order of magnitude of these parameters, as seen from the numerical experiments, is strongly linked with the extension of the investigated volume, with the total number of voxels and with the “certainty” of the geological model. A further foreseen improvement is to consider possible dependences of  $s^2$  and  $q^2$  from the voxel position, thus allowing the prior to be more informative where the geological model is considered more reliable (e.g. in presence of borehole logging).

Last but not least, the algorithm needs to be numerically optimized in order to increase the model resolution. This step will imply a relevant growth of the total number of variables, thus increasing the total computational burden.

**Acknowledgements** The authors would like to thank the anonymous reviewers for their valuable comments and suggestions that contributed to improving the manuscript.

## References

- Azencott R (1988) Simulated annealing. *Seminaire Bourbaki* 30:223–237
- Barbosa VCF, Silva JBC (1994) Generalized compact gravity inversion. *Geophysics* 59(1):57–68
- Barbosa VCF, Silva JBC, Medeiros WE (1997) Gravity inversion of basement relief using approximate equality constraints on depths. *Geophysics* 62(6):1745–1757
- Barbosa VCF, Silva JBC, Medeiros WE (1999) Gravity inversion of a discontinuous relief stabilized by weighted smoothness constraints on depth. *Geophysics* 64(5):1429–1437
- Bayes T (1763) An essay toward solving a problem in the doctrine of chances. *Philos Trans R Soc Lond* 53:370–418
- Bosch M (1999) Lithologic tomography: from plural geophysical data to lithology estimation. *J Geophys Res Solid Earth* 104(B1):749–766
- Bosch M (2004) The optimization approach to lithological tomography: combining seismic data and petrophysics for porosity prediction. *Geophysics* 69(5):1272–1282
- Bosch M, McGaughey J (2001) Joint inversion of gravity and magnetic data under lithologic constraints. *Lead Edge* 20(8):877–881

- Bosch M, Meza R, Jiménez R, Hönig A (2006) Joint gravity and magnetic inversion in 3D using Monte Carlo methods. *Geophysics* 71(4):G153–G156
- Box GEP, Tiao GC (2011) Bayesian inference in statistical analysis. Wiley, New York
- Caratori Tontini F, Cocchi L, Carmisciano C (2009) Rapid 3-D forward model of potential fields with application to the Palinuro Seamount magnetic anomaly (southern Tyrrhenian Sea, Italy). *J Geophys Res Solid Earth* 114(B2):1978–2012
- Casella G, Robert CP (1999) Monte Carlo statistical methods. Springer, New York
- Christensen NI, Mooney WD (1995) Seismic velocity structure and composition of the continental crust: a global view. *J Geophys Res Solid Earth* 100(B6):9761–9788
- De Lacy MC, Sansò F, Rodriguez-Caderot G, Gil AJ (2002) The Bayesian approach applied to GPS ambiguity resolution. A mixture model for the discrete-real ambiguities alternative. *J Geodesy* 76(2):82–94
- Fedi M (2006) DEXP: a fast method to determine the depth and the structural index of potential fields sources. *Geophysics* 72(1):I1–I11
- Fedi M, Rapolla A (1999) 3-D inversion of gravity and magnetic data with depth resolution. *Geophysics* 64(2):452–460
- Geman S, Geman D (1984) Stochastic relaxation, Gibbs distributions, and the Bayesian restoration of images. *IEEE Trans Pattern Anal Mach Intell PAMI-6*:721–741
- Gordon AC, Mohriak WU, Barbosa VCF (2012) Crustal architecture of the Almada Basin, NE Brazil: an example of a non-volcanic rift segment of the South Atlantic passive margin. *Geol Soc Lond Spec Publ* 369:215–234
- Guillen A, Menichetti V (1984) Gravity and magnetic inversion with minimization of a specific functional. *Geophysics* 49(8):1354–1360
- Guillen A, Calcagno P, Courrioux G, Joly A, Ledru P (2008) Geological modelling from field data and geological knowledge: part II. Modelling validation using gravity and magnetic data inversion. *Phys Earth Planet In* 171(1):158–169
- Last BJ, Kubik K (1983) Compact gravity inversion. *Geophysics* 48(6):713–721
- Medeiros WE, Silva JBC (1996) Geophysical inversion using approximate equality constraints. *Geophysics* 61(6):1678–1688
- Mosegaard K, Tarantola A (2002) Probabilistic approach to inverse problems. *Int Geophys* 81:237–265
- Nagihara S, Hall SA (2001) Three-dimensional gravity inversion using simulated annealing: constraints on the diapiric roots of allochthonous salt structures. *Geophysics* 66(5):1438–1449
- Nagy D (1966) The gravitational attraction of a right rectangular prism. *Geophysics* 31(2):362–371
- Parker RL (1973) The rapid calculation of potential anomalies. *Geophys J R Astron Soc* 31(4):447–455
- Parker RL (1975) The theory of ideal bodies for gravity interpretation. *Geophys J Int* 42(2):315–334
- Roy L, Sen MK, Blankenship DD, Stoffa PL, Richter TG (2005) Inversion and uncertainty estimation of gravity data using simulated annealing: an application over Lake Vostok, East Antarctica. *Geophysics* 70(1):J1–J12
- Roazanov YA (1982) Markov random fields. Springer, New York
- Sampietro D, Sansò F (2012) Uniqueness theorems for inverse gravimetric problems. *IAG Symp* 137:111–115
- Sansò F, Reguzzoni M, Triglione D (2011) Metodi Monte Carlo e delle Catene di Markov: una introduzione (in Italian). Maggioli Editore.
- Smith AFM, Roberts GO (1993) Bayesian computation via the Gibbs sampler and related Markov chain Monte Carlo methods. *J R Stat Soc Ser B (Methodological)* 55(1):3–23
- Tarantola A (2002) Inverse problem theory: methods for data fitting and model parameter estimation. Elsevier Science, Amsterdam
- Tarantola A, Valette B (1982) Inverse problems = quest for information. *J Geophys* 50(3):150–170

---

# Effects of Different Objective Functions in Inequality Constrained and Rank-Deficient Least-Squares Problems

Lutz Roese-Koerner and Wolf-Dieter Schuh

---

## Abstract

Rank-deficient estimation problems often occur in geodesy due to linear dependencies or underdetermined systems. Well-known examples are the adjustment of a free geodetic network or a finite element approximation with data gaps. If additional knowledge about the parameters is given in form of inequalities (e.g., non-negativity), a rank-deficient and inequality constrained adjustment problem has to be solved.

In Roese-Koerner and Schuh (J Geodesy, doi:10.1007/s00190-014-0692-1) we proposed a framework for the rigorous computation of a general solution for rank-deficient and inequality constrained least-squares problems. If the constraints do not resolve the manifold of solutions, a second minimization is performed in the nullspace of the design matrix. This can be thought of as a kind of pseudoinverse, which takes the inequality constraints into account.

In this contribution, the proposed framework is reviewed and the effect of different objective functions in the nullspace optimization step is examined. This enables us to aim for special properties of the solution like sparsity ( $L^1$  norm) or minimal maximal errors ( $L^\infty$  norm). In a case study our findings are applied to two applications: a simple bivariate example to gain insight into the behavior of the algorithm and an engineering problem with strict tolerances to show its potential for classic geodetic tasks.

---

## Keywords

Convex optimization • Inequality constrained least-squares •  $L^1$  norm •  $L^2$  norm •  $L^\infty$  norm • Nullspace minimization • Rank defect

---

## 1 Introduction

Geodesists often encounter rank-deficient optimization problems. This can be either due to underdetermined equation systems (e.g., resulting from a second order design of a geodetic network with more weights to be estimated than entries in the criterion matrix) or due to external parameters, which cannot be estimated from the observations (e.g., a

datum defect). These cases do not result in one unique but in a manifold of solutions. If the situation gets even more complicated and additional knowledge about the parameters in form of linear inequality constraints is given, no closed formulas exist. Inequalities for nonnegative quantities like run-time delay in GPS, SAR or VLBI or e.g., slope constraints in surface fitting are typical constraints.

In the unconstrained case, usually a second objective function is introduced to the problem to enforce a unique solution despite the rank defect. A classic choice would be to minimize the length of the solution vector with respect to the  $L^2$  norm (cf. Gill et al. 1991, pp. 230–234). While unconstrained problems with a rank defect are well studied, little is known about the inequality constrained case.

---

L. Roese-Koerner (✉) • W.-D. Schuh  
Institute of Geodesy and Geoinformation, University of Bonn, Bonn,  
Germany  
e-mail: roese-koerner@geod.uni-bonn.de; schuh@geod.uni-bonn.de

The special case of non-negative least-squares with a possible rank defect and additional inequalities was solved by Schaffrin (1981). However, this method cannot be generalized easily. Werner and Yapar (1996) proposed a method for computing a rigorous general solution of inequality constrained problems with a rank defect. Unfortunately, their method is solely suited for small-scale problems as it involves arbitrarily testing of subsets of the constraints, which becomes a limiting factor in the multivariate case. Additional work by Xu et al. (1999) is focused on the stabilization of ill-conditioned linear complementarity problems. Dantzig (1998) described a method to compute one arbitrary particular solution despite a possible rank defect for problems with a linear or quadratic objective function. However, no description of the manifold of solutions is given.

All these algorithms are either tailor-made for special problems or yield only one of an infinite number of particular solutions or are restricted to small-scale problems. Therefore, in Roesse-Koerner and Schuh (2014) we developed a framework for the rigorous computation of a general solution of inequality constraint least-squares problems. Based on this framework, here we set out to compute a solution with desirable properties like e.g., sparsity, exploring the opportunities which result from the choice of the objective function in the second minimization step.

## 2 Inequality Constrained Estimation

In the following, we focus on a linear Gauss-Markov model (GMM)

$$\boldsymbol{\ell} + \mathbf{v} = \mathbf{A}\mathbf{x}. \quad (1)$$

$n \times 1$  vector  $\boldsymbol{\ell}$  contains the observations and  $n \times 1$  vector  $\mathbf{v}$  the corresponding residuals.  $\mathbf{A}$  is the design matrix and  $\mathbf{x}$  comprises the  $m$  parameters to be estimated. Following the least-squares principle, the (weighted) sum of squared residuals shall be minimized

$$\Phi(\mathbf{v}) = \mathbf{v}^T \boldsymbol{\Sigma}^{-1} \mathbf{v} \dots \min. \quad (2)$$

$\boldsymbol{\Sigma}$  is the possibly fully populated variance-covariance matrix of the observations. The objective function shall be minimized with respect to  $p$  inequality constraints which have to be fulfilled strictly. Linear inequality constraints can be formulated as

$$\mathbf{B}^T \mathbf{x} \leq \mathbf{b}. \quad (3)$$

$m \times p$  matrix  $\mathbf{B}$  is called constraint matrix and  $p \times 1$  vector  $\mathbf{b}$  is the corresponding right-hand-side of the constraints. The

constraints can be subdivided in active constraints  $\mathbf{B}_a, \mathbf{b}_a$ , which hold as equality constraints at the optimal solution  $\bar{\mathbf{x}}$  and inactive constraints  $\mathbf{B}_i, \mathbf{b}_i$ , which hold as strict inequalities

$$\mathbf{B}_a^T \mathbf{x} = \mathbf{b}_a, \quad \mathbf{B}_i^T \mathbf{x} < \mathbf{b}_i. \quad (4)$$

Minimizing (2) subject to (3) will be referred to as the inequality constrained least-squares (ICLS) problem, which can be expressed as a quadratic program (QP) in standard form

INEQUALITY CONSTRAINED LEAST-SQUARES

**objective funct.:**  $\Phi(\mathbf{x}) = \mathbf{v}^T \boldsymbol{\Sigma}^{-1} \mathbf{v} \dots \text{Min}$

**constraints:**  $\mathbf{B}^T \mathbf{x} \leq \mathbf{b}$

**optim. variable:**  $\mathbf{x} \in \mathbb{R}^m$ .

(5)

A QP is a convex optimization problem with a quadratic objective function and linear constraints. As QPs are well studied, there is a variety of existing algorithms to solve them efficiently. Therefore, it can be beneficial to reformulate the described problem as a QP. As it is not known beforehand which constraints will be active in the optimal solution, only iterative algorithms exist. Most of them can be subdivided into two classes: simplex methods (e.g. Dantzig's simplex method for quadratic programming, Dantzig 1998, pp. 490–498) and interior-point methods (e.g., primal-dual methods, Boyd and Vandenberghe 2004, pp. 609–613).

## 3 Rank-Deficient ICLS Problems

In the following, we focus on determining a unique rigorous general solution of the ICLS problem (5) with a rank-deficient design matrix  $\mathbf{A}$ ,

$$Rg(\mathbf{A}) = r < m, \quad d = m - r. \quad (6)$$

The proposed framework consists of three major parts. First, inequality constraints are not taken into account, a general solution of the rank-deficient unconstrained problem is computed, and a transformation of parameters is performed (described in Sect. 3.1). The next step depends on whether the manifold of solutions and the feasible set intersect. In case of an intersection (*case 1*, Sect. 3.2), there is still a manifold and a second minimization procedure is carried out in the nullspace. If there is no such intersection (*case 2*, Sect. 3.3), Dantzig's method is used to compute a particular solution of the original problem. The resulting particular solution can be shown to be unique if no active constraint is parallel to the manifold. In turn, an active parallel constraint calls for a nullspace optimization.

### 3.1 Transformation of Parameters

As described in Roesse-Koerner and Schuh (2014), the introduction of linear inequality constraints can result in a shift and/or a restriction of the manifold but never in a rotation. Therefore, it is instructive to first compute a general solution of the unconstrained ordinary least-squares (OLS) problem

$$\tilde{\mathbf{x}}^{OLS}(\boldsymbol{\lambda}) = \mathbf{x}_p^{OLS} + \mathbf{X}_{hom} \boldsymbol{\lambda}. \quad (7)$$

Subsequently, the constraints (3) can be reformulated with respect to the  $d$  free parameters  $\boldsymbol{\lambda}$ , the particular solution  $\mathbf{x}_p^{OLS}$  and the homogenous solution  $\mathbf{X}_{hom}$ ,

$$\mathbf{B}^T (\mathbf{x}_p^{OLS} + \mathbf{X}_{hom} \boldsymbol{\lambda}) \leq \mathbf{b}. \quad (8)$$

With the substitutions  $\mathbf{B}_\lambda^T := \mathbf{B}^T \mathbf{X}_{hom}$  and  $\mathbf{b}_\lambda := \mathbf{b} - \mathbf{B}^T \mathbf{x}_p^{OLS}$ , (8) reads

$$\mathbf{B}_\lambda^T \boldsymbol{\lambda} \leq \mathbf{b}_\lambda. \quad (9)$$

If these constraints are contradictory (as it can be examined by solving a feasibility problem, cf. Boyd and Vandenberghe 2004, pp. 579–580), there is no intersection of manifold and feasible set and we proceed as described in Sect. 3.3. Otherwise, we proceed as described in Sect. 3.2.

### 3.2 Case 1: Intersection

In case of an intersection, we aim for the rigorous computation of a unique particular solution. Therefore, a second optimization problem is introduced (e.g., minimizing the length of the solution vector with respect to the norm  $L^p$ ). As this minimization takes place in the nullspace of design matrix  $\mathbf{A}$ , the value of the objective function (2) of the original problem does not change

<p style="text-align: center; margin: 0;">NULLSPACE OPTIMIZATION PROBLEM</p> <p><b>objective funct.:</b> <math>\Phi_{NS}(\mathbf{x}_p^{ICLS}(\boldsymbol{\lambda})) \dots \text{Min}</math></p> <p><b>constraints:</b> <math>\mathbf{B}_\lambda^T \boldsymbol{\lambda} \leq \mathbf{b}_\lambda</math></p> <p><b>optim. variable:</b> <math>\boldsymbol{\lambda} \in \mathbb{R}^d</math>.</p>	(10)
--	------

The minimization yields optimal free parameters  $\tilde{\boldsymbol{\lambda}}$ , whose insertion in (7) results in

$$\tilde{\mathbf{x}}^{ICLS} = \mathbf{x}_p^{OLS} + \mathbf{X}_{hom} \tilde{\boldsymbol{\lambda}}, \quad (11)$$

which is a unique particular solution that fulfills all constraints.

### 3.3 Case 2: No Intersection

If manifold and feasible region are disjunct and there is no active parallel constraint, the constraints have resolved the manifold. Therefore, it is sufficient to compute one particular solution of the constraint problem – e.g., with Dantzig’s simplex method for QPs. It can be shown that

$$\tilde{\mathbf{x}}^{ICLS} = \mathbf{x}_p^{ICLS} \quad (12)$$

is the unique solution. Instead, if at least one active constraint is parallel to the manifold, there is a shift of the manifold and a second objective function has to be introduced, as described in Sect. 3.2.

A more detailed description of the framework is provided in Roesse-Koerner and Schuh (2014).

## 4 Nullspace Optimization

In this section, the choice of the (second) objective function  $\Phi_{NS}$  in the nullspace optimization problem (10) is discussed. As the whole minimization takes place in the nullspace of the design matrix, the value of the original objective function (2) will not change. Nonetheless, choosing a suitable second objective function for a particular problem can be helpful to achieve properties like sparsity of the parameter vector.

Two different parts of the objective function can be distinguished, which will be examined in some detail: the *functional relationship* and the *norm*, with respect to which the minimization (or maximization) shall be performed.

### 4.1 Functional Relationship

The *functional relationship* strongly depends on the application. However, there are two main concepts which are applicable to a big variety of problems.

First, the length of the parameter vector  $\mathbf{x}$  can be minimized. This can e.g., be beneficial if not absolute coordinates but coordinate differences are estimated, which should be close to the initial coordinates.

More sophisticated approaches include a weighted minimization of the length of the parameter vector. For example if prior knowledge about the magnitude of the parameters is given. Prominent examples are Kaula’s rule of thumb in gravity field estimation or the demand for a decay of the amplitudes of higher frequencies in signal processing to achieve a square integrable function. A mathematical example for the minimization of the length of the parameter vector is provided in Sect. 5.1.

The second main concept is to maximize the distance to the constraints

$$\|\mathbf{B}^T \mathbf{x}(\boldsymbol{\lambda}) - \mathbf{b}\| \dots \max \quad (13a)$$

$$\iff \|\mathbf{B}^T \mathbf{x}_P + \mathbf{B}^T \mathbf{X}_{hom} \boldsymbol{\lambda} - \mathbf{b}\| \dots \max. \quad (13b)$$

This could be beneficial, if the constraints constitute a kind of outermost threshold. One of the main advantages of the use of inequality constraints is, that they do not influence the result, if they are not active. Therefore, it is usually not possible to provide a buffer to the boundary of the feasible set without losing estimation quality (shown by an increased value of the original objective function). However, due to the optimization in the nullspace of  $\mathbf{A}$ , we are in the unique position to apply such a buffer without deteriorating the estimate. In Sect. 5.2 an application is described, in which this type of *functional relationship* is applied.

## 4.2 Norms

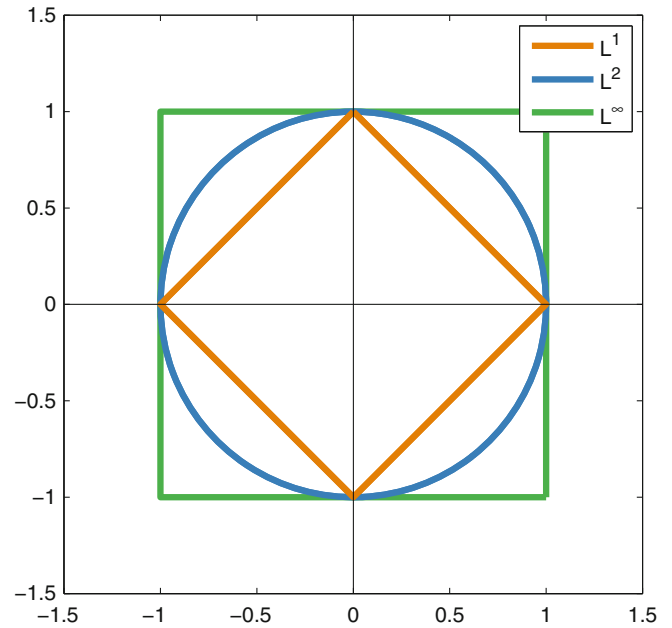
In the following, we will point out some general aspects of different *norms* and their influence on the most classical *functional relationship*: the length of the solution vector.  $L^p$  norms

$$\|\mathbf{x}\|_p := (x_1^p + x_2^p + \dots + x_m^p)^{1/p}, \quad p = 1, 2, \dots, \infty$$

used in adjustment theory include the  $L^1$  norm, the  $L^2$  norm and the  $L^\infty$  norm (cf. Jäger et al. 2005; Boyd and Vandenberghe 2004, pp. 125–128 and p. 635, respectively).

Most often, the length of a vector is minimized with respect to the  $L^2$  norm (also known as Euclidean norm). This is a quite natural choice, which can easily be visualized geometrically. However, as the influence of one element on the norm decreases if its absolute value becomes smaller, a minimization with respect to the  $L^2$  norm will seldom result in a sparse vector. This can be verified by examining the *blue*  $L^2$  norm unit sphere depicted in Fig. 1.

A naive choice to achieve maximal sparsity of a vector would be a minimization with respect to the  $L^0$  norm (e.g., the number of nonzero elements). However, a minimization with respect to the  $L^0$  norm is a combinatorial problem and computationally very demanding. Therefore, e.g., Candes et al. (2006) approximated the  $L^0$ -minimization-problem by the  $L^1$ -minimization-problem in the context of compressed sensing. They showed, that a minimization with respect to the  $L^1$  norm in most cases yields sparse results, too. This is used in many compressed sensing algorithms. Minimization with respect to the  $L^1$  norm is a convex optimization problem and can be formulated as linear program (cf. Dantzig 1998,



**Fig. 1** Two dimensional *unit spheres* of  $L^1$  (orange),  $L^2$  (blue) and  $L^\infty$  norm (green). The  $L^1$  diamond evolves from the summation of  $x$  and  $y$  value, while for the  $L^\infty$  square, only the value of the biggest quantity is decisive. Figure modified and extended from Tibshirani (1996)

pp. 60–62). As can be seen in Fig. 1, the corners of the *orange*  $L^1$  norm unit sphere coincide with the coordinate grid. This is equivalent to the statement, that one parameter is zero there, yielding a sparse solution.

Application of the  $L^\infty$  norm (also called Chebyshev norm) results in a parameter vector with minimal maximal value, that is usually not sparse. See Fig. 1 for the corresponding unit sphere. The  $L^\infty$  norm is often applied, when trying to maximize the distance to the constraints (cf. Sect. 4.1), in order to maximize the minimal buffer to the boundary.

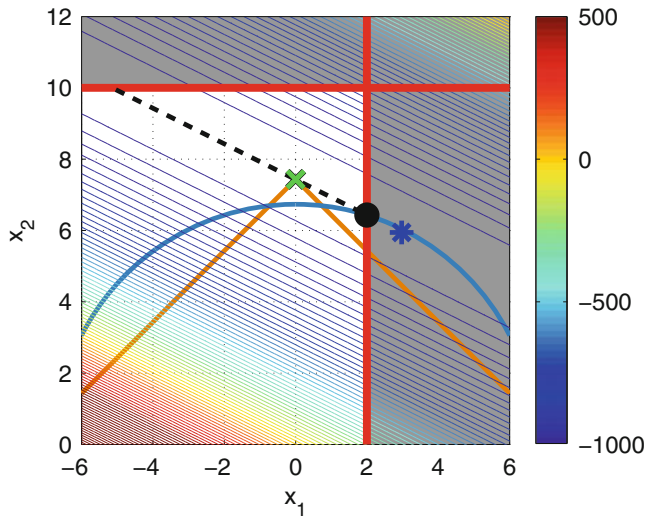
## 5 Case Studies

The effect of a nullspace optimization will be demonstrated in two examples: a very simple bivariate one and a more sophisticated network adjustment problem.

### 5.1 Case Study 1: Bivariate Example

In this case, a least-squares estimate of the two summands of a weighted sum is to be calculated

$$\ell_i + v_i = x_1 + 2x_2 \quad (14)$$



**Fig. 2** Contour lines of the objective function of example 1. Red lines represent constraints, the infeasible region is shaded. The dashed black line indicates the manifold of solutions. The green cross is the  $L^1$  norm solution, the black circle the  $L^2$  solution. For comparison the solution using a pseudoinverse (blue star) is shown, too. Appropriately scaled unit spheres are depicted in orange ( $L^1$ ) and light blue ( $L^2$ )

subject to the constraints

$$x_1 \leq 2, \quad x_2 \leq 10 \quad (15)$$

using the framework described in Sect. 3. A contour plot of the objective function and the constraints is given in Fig. 2. Red lines represent constraints, the infeasible region is shaded and the dashed black line indicates the manifold of solutions. We assume the observations

$$\boldsymbol{\ell}^T = [23.2 \ 16.4 \ 12.9 \ 8.2 \ 13.7], \quad (16)$$

to be uncorrelated. The observation equations read

$$\boldsymbol{\ell} + \mathbf{v} = \begin{bmatrix} 1 & 2 \\ 1 & 2 \\ 1 & 2 \\ 1 & 2 \\ 1 & 2 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \mathbf{A}\mathbf{x}, \quad (17)$$

with a clearly rank-deficient design matrix  $\mathbf{A}$ . Setting up the normal equations and applying the Gauss-Jordan algorithm yields a general OLS solution

$$\mathbf{x}(\lambda) = \begin{bmatrix} 14.88 \\ 0 \end{bmatrix} + \begin{bmatrix} -2 \\ 1 \end{bmatrix} \lambda = \mathbf{x}_P + \mathbf{X}_{hom}\lambda. \quad (18)$$

As there is an intersection of manifold and feasible set (cf. Fig. 2), the constraints are reformulated with respect to the free parameter  $\lambda$  and a second optimization problem in the

nullspace of the design matrix has to be solved. We chose to examine the different effects of minimizing the length of the parameter vector with respect to the  $L^1$  or  $L^2$  norm

2D EXAMPLE: NULLSPACE OPTIMIZATION

**objective funct.:**  $\|\mathbf{x}_P + \mathbf{X}_{hom}\lambda\|_p \dots \text{Min}$

**constraints:**  $\mathbf{B}_\lambda^T \lambda = \begin{bmatrix} -2 \\ 1 \end{bmatrix} \lambda \leq \begin{bmatrix} -12.88 \\ 10 \end{bmatrix} = \mathbf{b}_\lambda$

**optim. variable:**  $\lambda \in \mathbb{R}$ .

Depending on the chosen norm this results in either

$$\tilde{\mathbf{x}}_P^{ICLS,L^2} = \begin{bmatrix} 2.00 \\ 6.44 \end{bmatrix} \quad \text{or} \quad \tilde{\mathbf{x}}_P^{ICLS,L^1} = \begin{bmatrix} 0.00 \\ 7.43 \end{bmatrix}. \quad (19)$$

As expected, utilization of the  $L^1$  norm in the nullspace optimization step yields in a sparse solution without an increase in the sum of squared residuals of the original problem.

## 5.2 Case Study 2: Network Adjustment

The second case study is based on an engineering problem. We assume that some prefabricated building material shall be fitted between other elements so that the parts can be welded together. In order to make welding possible, tolerances have to be fulfilled strictly.

Figure 3 depicts the test case. Twenty-six distance measurements (black lines) are performed between the ten points  $P1$  to  $P10$  (black dots). Their 20 coordinates are the parameters to be estimated in a GMM (1). As no datum is defined, estimating absolute values of the coordinates is a rank-deficient problem.

Points  $P3$  to  $P6$  are located at the left-hand side of the gap the new part is supposed to fill, and the points  $P7$  to  $P10$  are located on its right-hand side.  $P1$  and  $P2$  are external points to stabilize the network. It shall be determined if the new part fits between both lines of points. This can be achieved by setting up the 16 linear constraints

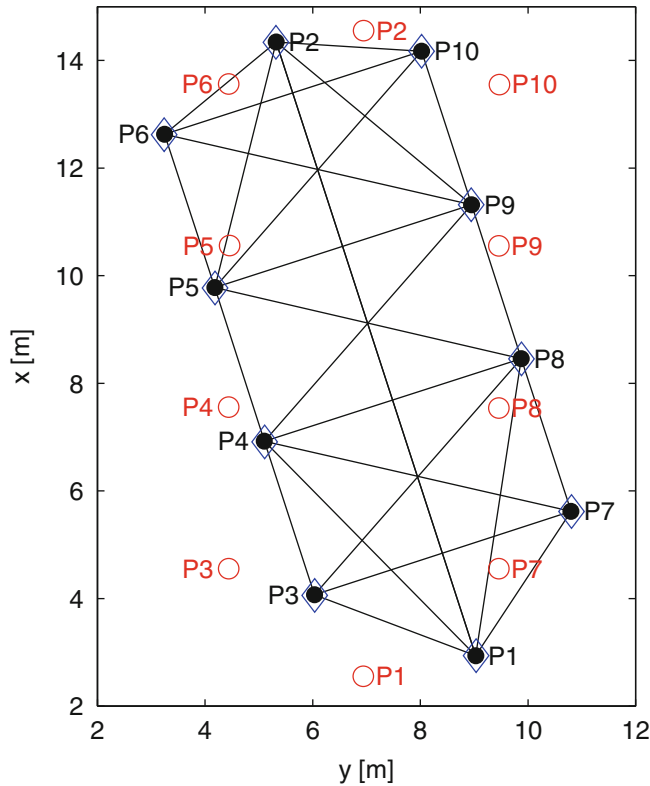
$$y_{7,8,9,10} - y_{3,4,5,6} \leq 5.03 \text{ m} \quad (20)$$

and the 16 linear constraints

$$y_{7,8,9,10} - y_{3,4,5,6} \geq 5.00 \text{ m}, \quad (21)$$

resulting in a rank-deficient ICLS problem in form of (5). While the first constraints guarantee, that the new part is not allowed to be wider than 5.03 m, the latter assure, that it is not smaller than 5.00 m (otherwise the gap would be too big for welding). The constraints force the estimated points to align





**Fig. 3** Example 2: Distance measurements (black lines) are performed between points  $P1$  to  $P10$  (black dots). A particular OLS solution (blue diamonds) and the ICLS solution (red circles) with maximal minimal distance to the constraints are shown

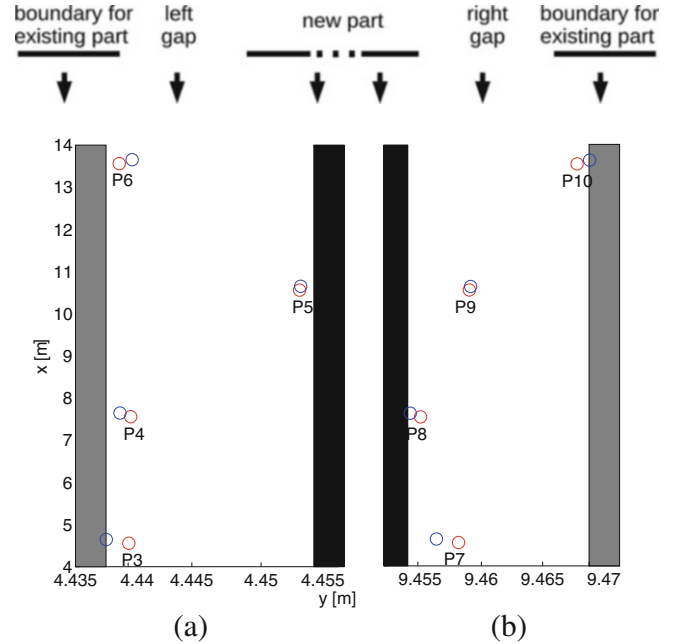
almost parallel to the  $x$  axis (cf. red circles in Fig. 3). If more than two of the 32 constraints mentioned above are active, the new part will not meet the tolerances. Incompatible elements can be detected via an analysis of the Lagrange multipliers (cf. Roesse-Koerner et al. 2012).

If the manifold is not resolved through the introduction of constraints, a nullspace optimization has to be performed. This can be used to maximize the minimal distance to the constraints

$$\Phi_{NS,\infty} = \|\mathbf{B}^T \mathbf{x}(\boldsymbol{\lambda}) - \mathbf{b}\|_{\infty} \dots \max. \quad (22)$$

Using the Chebyshev norm is always beneficial if tolerances instead of standard deviations are given. The optimization problem was solved using the CVX software (Grant and Boyd 2008). Results are shown in Fig. 3. In the chosen scenario, no constraint is active. Therefore, the new part will fit in the gap and welding is possible.

Figure 4 shows the welding boundary for the existing parts (gray area), the new part (black area) and the “gaps” at its left-hand (Fig. 4a) and right-hand (Fig. 4b) side. Please



**Fig. 4** Boundary for the existing parts (gray area), the new part (black area) and the “gaps” at the left-hand (a) and right-hand (b) side of the new part. The axes are scaled differently and there is a breach in the  $y$  axis, so most of the new part is not shown. Adjusted coordinates of the above mentioned ICLS estimate with maximal minimal distance to the constraints (red circles) are compared with those of an ICLS estimate with  $\Phi_{NS,L^2}$  as nullspace objective function (blue circles)

note the different scales in  $x$  and  $y$  direction and the breach in the  $y$  axis. Adjusted coordinates of the ICLS estimate with  $\Phi_{NS,\infty}$  (red circles) are compared with those of an ICLS estimate with

$$\Phi_{NS,L^2} = \|\mathbf{x}\|_{L^2} \dots \min \quad (23)$$

as nullspace objective function (blue circles).

While both estimates provide a decision if the new part will fit, only the adjustment with maximal minimal distance to the constraints allows to determine how well the new part will fit. This can be seen in Fig. 4, where for this estimate the minimal distance to the constraints is at least 2.5 mm at each side (namely for the points 5, 6, 8 and 10). In contrast, the blue points 3 and 10 are exactly on the boundary. So there is clearly a benefit in choosing a suited objective function for the nullspace optimization.

Due to space limitations we restricted ourselves to a brief description of the application and the figurative results presented in Figs. 3 and 4. Further information (e.g., the observations and the functional and stochastic model) as well as quantitative results can be obtained from the authors.

## 6 Conclusion

A framework for the computation of a rigorous general solution of rank-deficient ICLS problems has been reviewed. It has been shown that it is possible to obtain a solution with certain predefined optimality properties, if the manifold of solutions is not resolved by the constraints. As this results from a minimization process in the nullspace of the design matrix, the sum of squared residuals remains unchanged. Therefore, the described approach can be beneficial for practical applications as useful properties, like e.g., sparsity or maximal minimal distances, can be obtained without sacrificing estimation quality.

## References

- Boyd S, Vandenberghe L (2004) Convex optimization. Cambridge University Press, Cambridge
- Candes E, Romberg J, Tao T (2006) Robust uncertainty principles: exact signal reconstruction from highly incomplete frequency information. *IEEE Trans Inf Theory* 52(2):489–509. doi:10.1109/TIT.2005.862083
- Dantzig G (1998) Linear programming and extensions. Princeton University Press, New Jersey
- Gill P, Murray W, Wright M (1991) Numerical linear algebra and optimization, vol 1. Addison Wesley Publishing, Redwood City
- Grant M, Boyd S (2008) Graph implementations for nonsmooth convex programs. In: Blondel V, Boyd S, Kimura H (eds) Recent advances in learning and control. Lecture notes in control and information sciences. Springer, London, pp 95–110
- Jäger RR, Müller T, Saler H, Schwäble R (2005) Klassische und robuste Ausgleichungsverfahren. Wichmann-Verlag, Heidelberg
- Roese-Koerner L, Schuh WD (2014) Convex optimization under inequality constraints in rank-deficient systems. *J Geodesy*. doi:10.1007/s00190-014-0692-1
- Roese-Koerner L, Devaraju B, Sneeuw N, Schuh WD (2012) A stochastic framework for inequality constrained estimation. *J Geodesy* 86(11):1005–1018. doi:10.1007/s00190-012-0560-9
- Schaffrin B (1981) Ausgleichung mit Bedingungs-Ungleichungen. *Allgemeine Vermessungs-Nachrichten* 88. Jg.:227–238
- Tibshirani R (1996) Regression shrinkage and selection via the Lasso. *J R Stat Soc* 58(1):267–288
- Werner HJ, Yapar C (1996) On inequality constrained generalized least squares selections in the general possibly singular Gauss-Markov model: a projector theoretical approach. *Linear Algebra Appl* 237/238:359–393. doi:10.1016/0024-3795(94)00357-2
- Xu P, Cannon E, Lachapelle G (1999) Stabilizing ill-conditioned linear complementarity problems. *J Geodesy* 73:204–213. doi:10.1007/s001900050237

---

## List of Reviewers

Hussein A. Abd-Elmotaal  
Zuheir Altamimi  
Ali Reza Amiri-Simkooei  
Detlef Angermann  
Alireza A. Ardalan  
Georges Balmino  
Oliver Baur  
Battista Benciolini  
Ludovico Biagi  
Christian Bizouard  
Gerd Boedecker  
Johannes Boehm  
Jianqing Cai  
Don Chambers  
Jianli Chen  
Minkang Cheng  
Mattia Crespi  
James L. Davis  
Maria Clara de Lacy  
Demitris Delikaraoglou  
Athanasios Dermanis  
Jan Douša  
Ian Dowman  
Annette Eicker  
Artu Ellmann  
Mehdi Eshagh  
Xing Fang  
Yanming Feng  
Jeff Freymueller  
Toshio Fukushima  
Gabriele Giorgi  
Rossen Grebenitcharsky  
Richard Gross  
Christine Hackman  
Rüdiger Haas  
J. Kusche  
Shin-Chan Han  
Thomas Mejer Hansen  
Bernhard Heck  
Thomas Hobiger  
Simon Holmes  
Horst Holstein  
Jianliang Huang  
Karsten Jacobsen  
Juraj Janak  
Christopher Jekeli  
Shuanggen Jin  
Wei-Wen Kao  
Wolfgang Keller  
Steve C. Kenyon  
Jaroslav Klokocnik  
K.-R. Koch  
Wieslaw Kosek  
Christopher Kotsakis  
John Langbein  
Bofeng Li  
Zhizhao Liu  
Bryant Loomis  
Ambrogio Manzino  
Urs Marti  
Gwendoline Métivier  
Mary C. Meyer  
Ulrich Meyer  
Hiroshi Munekane  
Arthur Niell  
Felipe Nievinski  
Tobias Nilsson  
Pavel Novak  
Susan E. Owen  
Erricos C. Pavlis  
Marc Pierrot-Deseilligny  
Waldemar Popinski  
Hannes Raggam  
Paul Rebischung  
John C. Ries  
Reiner Rummel  
Fausto Sacerdote  
David Salstein

Fernando Sanso  
Alvaro Santamaria-Gomez  
Michael Schmidt  
Wolf-Dieter Schuh  
Josef Sebera  
Michael G. Sideris  
Nico Sneeuw  
Robert Tenzer

Peter J.G. Teunissen  
Carl Christian Tscherning  
Ilias Tziavos  
Giovanna Venuti  
Pieter NAM Visser  
Jinling Wang  
Simon David Paul Williams  
Hasan Yildiz

## Author Index

### A

Abramowitz, M., 196, 197  
Agnew, D.C., 20, 52  
Agrawal, O.P., 4  
Aguilera, E., 244  
Albertella, A., 83, 276  
Alexidze, M.A., 91  
Altamimi, Z., 11, 12, 15, 20, 227, 258  
Amante, C., 278  
Amjadiparvar, B., 121–127  
Andersen, O.B., 27, 96, 189  
Anderson, D.L., 275  
Angermann, D., 10  
Arabelos, D., 27, 28, 30, 101, 276  
Arnold, S.F., 264  
Artz, T., 14  
Askne, J., 150, 244, 245  
Asmar, S.W., 157  
Atkinson, R., 14  
Aubert, G., 284  
Aujol, J.-F., 284  
Avallone, A., 51  
Azencott, R., 319, 320

### B

Baarda, W., 266, 310  
Backus, G.E., 186  
Bagherbandi, M., 275  
Baleanu, D., 4, 5  
Balmino, G., 27, 110  
Balz, T., 233  
Bamler, R., 244, 245, 249  
Bar-Itzhack, I.Y., 309  
Bar-Sever, Y., 132, 133, 148  
Barzaghi, R., 214, 217, 218, 275–281  
Bassin, C., 278  
Bateman, H., 196  
Bayes, T., 318  
Beavan, J., 20, 52, 56  
Becker, J.J., 94, 189  
Benedetti, E., 59–66  
Bennett, R.A., 76  
Bennitt, G.V., 132, 136  
Bentel, K., 115–120  
Beran, J., 20  
Berg, H., 148, 150  
Bergstrand, S., 20

Bettadpur, S., 166  
Betti, B., 213–218  
Beutler, G., 52, 157–162, 177–182  
Bevis, M., 142, 148  
Biagi, L., 257–262  
Biancale, R., 139–145  
Bianco, G., 131–137  
Bingham, R.J., 27  
Bisnath, S.B., 33  
Bizouard, C., 10, 45  
Bjerhammar, A., 186  
Blaha, G., 221–223  
Blewitt, G., 60, 87  
Blossfeld, M., 16  
Bocchio, F., 3  
Bock, H., 157–162, 177–182  
Bock, O., 152  
Bock, Y., 59, 60, 267  
Boehm, J., 140  
Bogusz, J., 19–25  
Böhm, J., 148  
Bolotin, S., 14  
Borghi, A., 275–281  
Borre, K., 309  
Bos, M.S., 20, 57  
Bosch, M., 318  
Bosch, W., 110, 114  
Boucher, C., 221, 258  
Bouman, J., 101–107  
Box, G.E.P., 318  
Boy, J.-P., 228  
Boyd, S., 326–328, 330  
Branzanti, M., 59–66  
Broomhead, D., 76  
Bruinsma, S.L., 72, 166  
Brunini, C., 148  
Bruyninx, C., 20, 133  
Buckley, S., 258  
Budillon, A., 244, 245, 249  
Burrige, D.M., 140  
Burša, M., 91  
Buyn, S., 148  
Byram, S., 133  
Byun, S., 132, 133

### C

Calori, A., 147–153  
Candes, E., 328

Cannon, J.B., 121  
 Capaldo, P., 233–239  
 Capitaine, N., 14  
 Caputo, M., 3–6  
 Caputo, R., 4  
 Caratori Tontini, F., 318  
 Carcanague, S., 33  
 Carcano, L., 257–262  
 Carrion, D., 213–218  
 Casella, G., 320  
 Cesarone, F., 4  
 Chai, Y., 205  
 Chelton, D.B., 47  
 Chen, C.S., 92  
 Chen, G., 140  
 Chen, P.H., 234  
 Chen, Q., 75–80  
 Chen, W., 45, 92, 93  
 Cheng, M.K., 227–229, 231  
 Christakos, G., 259  
 Christensen, E.J., 148  
 Christensen, N.I., 319  
 Clenshaw, C.W., 201  
 Cloude, S.R., 244  
 Cohen, C.E., 309  
 Collilieux, X., 77, 81, 82, 85  
 Collins, J.P., 132  
 Colombelli, S., 65  
 Colombo, O., 121, 291  
 Colosimo, G., 59–66, 147–153  
 Crespi, M., 59–66, 147–153, 233–239  
 Crowder, M., 286  
 Cucurull, L., 142  
 Cunderlík, R., 91–100, 185–191

**D**

Dach, R., 148  
 Dahle, C., 165–174, 177–182  
 Dahlen, F.A., 41, 43, 45  
 Dai, Z., 33, 291  
 Dantzig, G., 326–328  
 Davis, J.L., 75–77, 80, 150  
 Davis, P.J., 259  
 Dayoub, N., 91  
 de Bakker, P.F., 33–38, 263  
 De Gaetani, C., 275–281  
 de Jonge, P.J., 263  
 De Lacy, M.C., 320  
 Dee, D.P., 150  
 Dehant, V., 45  
 Del Castillo Negrete, D., 4  
 Demaria, G., 4  
 Deming, W.E., 301, 306  
 Dennis, J.E., 288  
 Dermanis, A., 9–16, 221  
 Desjardins, C., 139–145  
 Desmond, A.F., 286  
 Dettmering, D., 101–107  
 Devaraju, B., 67–73  
 Devoti, R., 51–57  
 Díaz, G., 186  
 Dobslaw, H., 227, 231  
 Döll, P., 168  
 Donatelli, D., 258  
 Donevska, S., 306

Dong, D., 75  
 Dousa, J., 136  
 Dow, J.M., 65, 133, 153  
 Downman, I.J., 234  
 Drinkwater, M.R., 116  
 Du, S., 291  
 Durden, S., 244  
 D'Urso, M.G., 205–211

**E**

Eakins, B.W., 278  
 Eanes, R.J., 227  
 Ebbing, J., 101, 107  
 Edwards, S., 152  
 Eicker, A., 101, 103  
 El Shahed, M., 4  
 El-Sheimy, N., 257  
 Engelis, T., 109, 110  
 Erdélyi, A., 196  
 Eshagh, M., 276  
 Euler, H.-J., 35, 267  
 Ewing, C.E., 284  
 Eymard, R., 186–188

**F**

Fagiolini, E., 165–174  
 Fan, C.M., 92  
 Fang, X., 301, 303, 305  
 Farrel, W.E., 228  
 Fašková, Z., 186, 187  
 Featherstone, W.E., 197, 199, 201  
 Fecher, F., 125  
 Feder, J.W., 20  
 Fedi, M., 318  
 Fernandes, M.J., 152  
 Ferrándiz, J.M., 45  
 Ferro-Famil, L., 243–254  
 Figurski, M., 19–25  
 Filmer, M.S., 218  
 Fitzmauric, G.M., 286  
 Flamant, P.H., 284  
 Flechtner, F., 165–174, 178, 227–232  
 Fornaro, G., 244, 245, 249  
 Förste, C., 158, 229  
 Fotopoulos, G., 91  
 Fratarcangeli, F., 233–239  
 Freedman, W., 69, 116  
 Freeman, A., 244  
 Frey, O., 244  
 Freymueller, J., 76  
 Fritsche, M., 81  
 Fuchs, M., 101–107  
 Fukushima, T., 197, 199

**G**

Gallardo-Delgado, L.A., 205  
 Gambis, D., 10  
 Gao, Y., 291  
 García-Abdeslem, J., 205  
 Gasquet, C., 49  
 Gatelli, F., 244  
 Gazeaux, J., 82  
 Gegout, P., 139–145

Geman, D., 318, 319  
 Geman, S., 318, 319  
 Genrich, J.F., 60  
 Gentile, G., 214  
 Gerlach, C., 123, 125, 158  
 Ghil, M., 76  
 Gil, A., 197  
 Gill, P., 325  
 Gini, F., 245–248  
 Giorgi, G., 263, 309–315  
 Goad, C.C., 35, 111, 267  
 Golberg, M.A., 92  
 Goldberg, D., 201  
 Goldstein, H., 42  
 Golub, G.H., 301, 302, 304  
 Gooding, R.H., 179  
 Goodman, J.W., 284  
 Goossens, S., 161  
 Gordon, A.C., 318  
 Grad, M., 280  
 Grafarend, E., 3, 4, 186, 203, 222  
 Graffi, D., 6  
 Grand, T., 189  
 Grant, M., 330  
 Gregory, J.M., 47  
 Grejner-Brzezinska, D.A., 20  
 Gross, R.S., 41–45  
 Groten, E., 45  
 Gruber, C., 165–174  
 Gruber, T., 124  
 Gu, Y., 92, 94  
 Guillen, A., 318  
 Gunther, C., 292  
 Gutjahr, K., 233

**H**

Haase, J., 136  
 Hadas, T., 20  
 Hamayun, P., 205–207, 209, 211  
 Hansen, R.O., 205  
 Harris, F.J., 68, 69  
 Harrison, J.C., 5  
 Hartl, P., 245  
 Hayden, T., 121–127  
 Heck, B., 186  
 Hefty, J., 14  
 Heiskanen, W.A., 121, 122, 213–215  
 Helmert, F.R., 301  
 Henkel, P., 291–298  
 Herring, T.A., 140  
 Heyde, C.C., 286  
 Hill, A.C., 284  
 Hill, E.M., 20  
 Hill, G.W., 179  
 Hinze, W.J., 205  
 Hirschmuller, H., 239  
 Hobiger, T., 140, 143  
 Hobson, E.W., 196, 197  
 Hofmann-Wellenhof, B., 186  
 Hogg, D., 148  
 Holmes, S.A., 197, 199, 201  
 Holota, P., 187, 195–203  
 Holschneider, M., 116  
 Holstein, H., 205–209, 211  
 Hon, Y.C., 92

Höpfner, J., 45  
 Hortal, M., 140  
 Hoskins, B.J., 68  
 Hough, S.S., 43, 45  
 Huang, J., 126  
 Huang, Y., 243–254  
 Hung, H.K., 59

**I**

Iaffaldano, G., 4

**J**

Jäger, R.R., 328  
 Jäggi, A., 157–162, 177–182  
 Jekeli, C., 67, 68, 71, 199, 202, 213, 293, 294  
 Johnson, H., 20, 52  
 Johnson, H.O., 20  
 Johnston, R.L., 91, 92  
 Jupp, A., 132

**K**

Kahan, D., 158  
 Kaplan, D.T., 20  
 Kar, S., 228  
 Kaula, W.M., 110, 111, 113, 178  
 Keihm, S.J., 149  
 Ketteridge, B., 209  
 Khanafseh, S., 263  
 Khodabandeh, A., 263–270  
 Kiam, J.J., 292  
 Kidner, D., 259  
 Kiefer, R.W., 69  
 Kim, J., 166  
 King, M.A., 20, 22, 52, 56  
 Kleusberg, A., 309  
 Klipstein, W.M., 157  
 Klobuchar, J.A., 62  
 Klokočník, J., 110  
 Klos, A., 19–25  
 Knudsen, P., 276  
 Koch, K.-R., 104, 116, 118, 186, 187, 222, 259  
 König, R., 227–232  
 Konopliv, A.S., 158, 161  
 Kontny, B., 20  
 Koop, R., 158  
 Kosek, W., 19–25, 47–50  
 Kostelecký, J., 111  
 Kotsakis, C., 15, 122, 221–226  
 Kouba, J., 59, 150  
 Krarup, T., 276  
 Kuan, D.T., 289  
 Kubik, K., 318  
 Kukusha, A., 286  
 Kupradze, V.D., 91  
 Kursinski, E., 148  
 Kurtenbach, E., 166, 168  
 Kusche, J., 67, 104, 116, 118, 168

**L**

Lagarias, J.C., 53  
 Langbein, J., 20, 52, 56, 57, 60  
 Langley, R.B., 132

Laprise, R., 67  
 Larson, K., 59  
 Last, B.J., 318  
 Leandro, R.F., 132  
 Lebedev, S., 275  
 Leberl, F., 233  
 Lee, J.S., 235, 238, 246, 247, 254  
 Lemoine, F.G., 158  
 Li, Z., 257, 260  
 Lieb, V., 101–107  
 Lillesand, T.M., 69  
 Liu, X., 34, 166  
 Liu, Z., 33  
 Lomb, N.R., 54  
 Lombardini, F., 245–248  
 Longuevergne, L., 67  
 López-Martínez, C., 284  
 Lozier, D.W., 199  
 Luthcke, S.B., 178  
 Luzum, B., 14, 229

**M**

Ma, Y., 235  
 Mac Millan, W.D., 277  
 Macák, M., 185–191  
 MacDonald, G.J.F., 9  
 MacDoran, P.F., 284  
 Mackern, M.V., 147–153  
 Magnus, J.R., 303  
 Mahboub, V., 301  
 Mainardi, F., 4  
 Mandelbrot, B., 20  
 Mao, A., 20, 52  
 Marmo, F., 205  
 Marotta, A.M., 275–281  
 Marti, U., 218  
 Martinell, J.J., 4  
 Mathon, R., 91, 92  
 Mayer-Gürr, T., 28, 94, 95, 102, 158, 166, 189  
 Mazzoni, A., 59–66  
 McCullagh, P., 286, 288  
 McGaughey, J., 318  
 Meier, E., 244  
 Meier, U., 278  
 Meissl, P., 221  
 Mendes, V.B., 150  
 Menichetti, V., 318  
 Mervart, L., 157–162  
 Meyer, T.H., 213  
 Meyer, U., 157–162, 166, 177–182  
 Mignard, F., 82, 84, 85  
 Mikula, K., 94, 185–191  
 Milani, N., 3  
 Minarechová, Z., 185–191  
 Minh, D.H.T., 244, 249  
 Mitchell, H.L., 258  
 Mitchell, J.F.B., 147  
 Mitchell, M.M., 284  
 Miziński, B., 47, 48  
 Mooney, W.D., 319  
 Moreira, A., 245, 246  
 Moritz, H., 15, 28, 186, 196, 213–215, 276, 277  
 Mortiz, H., 121, 122  
 Mosegaard, K., 318  
 Moses, R., 245–248

Mougin, E., 243  
 Mtamakaya, J.D., 81, 85  
 Mueller, I., 15  
 Munk, W.H., 9  
 Murray, J.R., 75, 77

**N**

Naber, M., 4  
 Nadarajah, N., 38  
 Nafisi, V., 140  
 Nagy, D., 320  
 Namiki, N., 161  
 Nannini, M., 254  
 Nascetti, A., 233–239  
 Negretti, M., 259  
 Nehorai, A., 246–248  
 Neitzel, F., 301, 305  
 Nelder, J., 286, 288  
 Nesvadba, O., 195–203  
 Neuber, H., 5  
 Neudecker, H., 303  
 Neumann, M., 244  
 Niedzielski, T., 47–50  
 Niell, A., 140, 148  
 Niemeier, W., 186  
 Ning, T., 132  
 Noll, C.E., 65  
 Noll, W., 9  
 Nordius, H., 150  
 Novák, P., 167

**O**

Odiijk, D., 263  
 Odolinski, R., 38  
 Ohta, Y., 59  
 Oku, N., 291–298  
 Olesen, O., 106  
 Oliver, C.J., 284  
 Ostini, L., 82, 86  
 O'Sullivan, D., 259  
 Ottersten, B., 248

**P**

Pace, B., 131–137  
 Pacione, R., 131–137  
 Pagiatakis, S.D., 84  
 Pail, R., 101, 122, 124, 125, 158, 278  
 Panet, I., 117  
 Papatthanassiou, K., 244  
 Parker, R.L., 275, 317  
 Pavlis, N.K., 28, 94, 124, 189  
 Pearlman, M.R., 227, 284  
 Peng, J.H., 283–289  
 Penna, N., 132  
 Perfetti, N., 82, 83  
 Petit, G., 14, 229  
 Petrov, L., 228, 284  
 Picot, N., 149  
 Piallice, F., 233–239  
 Pietrantonio, G., 51–57  
 Pilgrim, B., 20  
 Pisani, A.R., 51–57  
 Plaut, G., 76



Poder, K., 201  
 Pohánka, V., 205–207  
 Pondrelli, S., 60, 62  
 Pope, A., 305  
 Pope, A.J., 186, 187  
 Popiński, W., 47–50  
 Porfiri, M., 233–239  
 Pottier, E., 243–254  
 Press, W.H., 260  
 Pugh, D.T., 124

**R**

Raggam, H., 233, 237  
 Rahmstorf, S., 47  
 Raimondo, J.-C., 227–232  
 Rangelova, E., 121–127  
 Rapolla, A., 318  
 Rapp, R.H., 121  
 Rau, R.J., 59  
 Rauch, H.E., 77  
 Ray, J.R., 81, 86  
 Rebischung, P., 20, 85  
 Rees, W.G., 259  
 Reguzzoni, M., 275–281, 317–323  
 Reigber, A., 245, 246  
 Reigber, C., 110, 166, 275  
 Remmer, O., 27  
 Riguzzi, F., 51–57  
 Robert, C.P., 320  
 Roberts, G.O., 318  
 Rocca, F., 244  
 Rocken, C., 148  
 Roese-Koerner, L., 325–331  
 Roggero, M., 81–87  
 Rosati, L., 205  
 Rosborough, G.W., 110, 112, 114  
 Röske, F., 47  
 Rossi, L., 317–323  
 Rothacher, M., 12  
 Roy, L., 318  
 Rozanov, Y.A., 319  
 Ruf, C.S., 149  
 Rummel, R., 72, 101, 110, 121–123, 127, 167  
 Russo, P., 207

**S**

Saastamoinen, J., 132, 148  
 Sacerdote, F., 186, 213–218  
 Sampietro, D., 275–281, 317–323  
 Sánchez, L., 91  
 Sansò, F., 4, 28, 121, 186, 213–215, 217, 275, 276, 317–323  
 Santamaria-Gómez, A., 53, 56  
 Santerre, R., 132  
 Sardeshmukh, P.D., 68  
 Sauer, S., 244, 245, 247  
 Scargle, J.D., 54, 55  
 Schabel, M., 147  
 Schaffrin, B., 221, 267, 301–306, 326  
 Schmidt, M., 101–107, 115–120  
 Schmidt, R., 166, 229  
 Schmitt, G., 222, 223  
 Schnabel, R.B., 288  
 Schoellhamer, D.H., 77  
 Schreiner, M., 69

Schubert, G., 275  
 Schuh, W.D., 325–331  
 Sciarretta, C., 131–137  
 Sebera, J., 203  
 Seeber, G., 284  
 Segall, P., 75, 77  
 Segura, J., 197  
 Seitz, M., 12  
 Serafino, F., 244, 245, 249  
 Sessa, S., 205  
 Shako, R., 166  
 Sharifi, M.A., 109–114  
 Shen, W., 45  
 Shi, Y., 283–289  
 Shimada, S., 286–288  
 Shin, Y.M., 276  
 Shuster, M.D., 309  
 Sideris, M., 91, 213–215, 217  
 Sideris, M.G., 121–127  
 Sillard, P., 221  
 Silva, J.B.C., 318  
 Simmons, A.J., 140  
 Simsky, A., 34  
 Siqueira, P., 243  
 Sjöberg, L., 275  
 Slichter, L., 4  
 Smith, A.F.M., 318  
 Smith, J.M., 199  
 Smith, M.L., 41, 43, 45  
 Smith, W.H.F., 21  
 Sneeuw, N., 67–73, 75–80, 109–114, 177–182  
 Snow, K., 305  
 Sošnica, K., 230, 231  
 Sona, G., 197  
 Soudarin, L., 139–145  
 Standish, E.M., 45  
 Stegun, I.A., 196, 197  
 Steigenberger, P., 82, 173  
 Stoica, P., 245, 247, 248  
 Strakhov, V.N., 210  
 Strang, G., 309  
 Svensson, L., 186  
 Swami, A., 284  
 Swenson, S., 67

**T**

Tang, K.T., 208  
 Tanni, L., 4  
 Tapley, B.D., 67, 101, 110, 116, 157, 165, 229, 275  
 Tarantola, A., 275, 318  
 Tebaldini, S., 244, 254  
 Teferle, F.N., 20  
 Teke, K., 136, 137, 152  
 Tenzer, R., 275  
 Tesmer, V., 76, 173  
 Teunissen, P., 83, 121, 122, 127, 221, 291, 292, 294  
 Teunissen, P.J.G., 33–38, 263–270, 309–315  
 Thong, N.C., 203  
 Tian, H., 284  
 Tiao, G.C., 318  
 Tibshirani, R., 328  
 Tisserand, F., 43  
 Torge, W., 196  
 Tregoning, P., 85  
 Treuhaft, R., 243

Trotta, S., 205  
Truesdell, C., 9  
Trujillo, J.J., 5  
Tscherning, C.C., 27–31, 101, 103, 201, 276  
Tsoulis, D., 15  
Turcotte, D.L., 275

**U**

Ueno, M., 132  
Ulaby, F., 284  
Untch, A., 140  
Unwin, D., 259

**V**

Valette, B., 318  
van Dam, T., 75–80  
Van Hoolst, T., 45  
van Loan, C.F., 301, 302, 304  
Van Ness, J., 20  
Vandenberghe, L., 326–328, 330  
Vaníček, P., 84, 213  
Vautard, R., 76  
Vedel, H., 136  
Vei, M., 227–232  
Venuti, G., 121, 213–218  
Véronneau, M., 126  
Vespe, F., 132  
Viberg, M., 248  
Visser, P., 180

**W**

Wagner, C.A., 178  
Wahr, J.M., 41, 44, 45, 67  
Wang, F.Z., 92, 93  
Wang, J., 136  
Wang, J.Y., 284  
Watkins, M., 166  
Watkins, M.M., 227  
Watson, C., 85  
Watson, C.S., 20, 52, 56  
Wax, M., 247, 248  
Wdowinski, S., 52

Wedderburn, R., 286  
Wei, T., 92  
Weigelt, M., 75–80  
Welch, P.D., 54  
Wells, D., 84  
Wentz, J., 147  
Werner, H.J., 326  
Werth, S., 67, 171  
Wessel, P., 21  
Weyl, H., 4  
Wielgosz, P., 20  
Wieser, A., 301, 302  
Williams, S.D.P., 20, 22, 23, 25, 52, 55, 56, 75  
Willis, P., 284  
Wittwer, T., 199  
Wnęk, A., 47–50  
Woodworth, P.L., 124  
Wyatt, F., 52

**X**

Xie, K., 33  
Xu, P.L., 59, 283–289, 303, 326

**Y**

Yapar, C., 326  
Yoder, C.F., 45

**Z**

Zbylut-Górska, M., 47–50  
Zhang, G., 235  
Zhang, J., 20, 52  
Zhang, L., 136  
Zheng, Y., 132  
Zhou, X., 205  
Zhu, S., 227  
Zhu, X.X., 244, 245, 249  
Zhu, X.Y., 235  
Zilkoski, D., 121  
Ziskind, I., 247, 248  
Zotov, L., 45  
Zuber, M.T., 157, 161  
Zus, F., 140