

Integrated approach for the development across Europe of user oriented climate indicators for GFCS high-priority sectors: Agriculture, disaster risk reduction, energy, health, water and tourism

Work Package 3

Deliverable 3.1c

INDECIS Quality Control Software and Manual:
INQC, beta version
Dr. Enric Aguilar
Center for Climate Change, C3
Universitat Rovira i Virgil (URV) de Tarragona (Spain)



This report arises from the Project INDECIS which is part of ERA4CS, an ERA-NET initiated by JPI Climate, and funded by FORMAS (SE), DLR (DE), BMWFW (AT), IFD (DK), MINECO (ES), ANR (FR), with co-funding by the European Union's Horizon 2020 research and innovation programme

| | |
|---|-------|
| 1. Authorship and licensing..... | 3 |
| 2. An overview | 3 |
| 3. Preparing your computer to run INQC..... | 4 |
| 4. Wrapper Functions and Jump-Start option | 5 |
| 5. Quality Control Tests..... | 5 |
| 6. Parametrization of wrapper functions and personalized runs of INQC | 13 |
| 7. Testing INQC and..... | 16 |
| 8. Expected evolution | 16 |
| Table 1: INQC Tests Description | 5 |
| Table 2 Default parametrization of the temperature() function, as ran by inqc() | 13 |
| Table 3 Default parametrization of the precip() function as ran by inqc() | 15 |
| Figure 1: INQC output example. Daily Maximum Temperature. Each column presents the result of one test applied to the data. The qc value is either 0 (pass) 1 (does not pass). | 4 |

1. Authorship and licensing

This code is provided free under the terms of the GNU Lesser General Public License as published by the Free Software Foundation, version 3.0 of the License. It is distributed under the terms of this license 'as-is' and has not been designed or prepared to meet any Licensee's particular requirements. The author and his institution make no warranty, either express or implied, including but not limited to, warranties of merchantability or fitness for a particular purpose. In no event will they will be liable for any indirect, special, consequential or other damages attributed to the Licensee's use of The Library. In downloading this code you understand and agree to these terms and those of the associated LGP License. See the GNU Lesser General Public License for more details (<http://www.gnu.org/licenses/lgpl.html>) or contact the Free Software Foundation, Inc., 59 Temple Place - Suite 330, Boston, MA 02111-1307, USA.

2. An overview

This document complements the documents D3.1a and D3.1c for the completion of the INDECIS' deliverable 3.1 INDECIS Quality Control Software and Manual. Here we present the software INDECIS QC, beta version (inqc_beta.R), created by Enric Aguilar, Center for Climate Change, C3, Universitat Rovira i Virgili, Tarragona (Spain), and licenced under the terms expressed in Section 1.

The software will be made available at : www.indecis.eu/software/INQC.html [pending]

Contact person: enric.aguilar@urv.cat

INDECIS QC (INQC, from now onwards) is designed to quality control European Climate Assessment and Dataset (ECA&D) daily data of maximum, minimum and average temperature, precipitation, sea level pressure, relative humidity, wind speed, snow depth, cloud coverage and sunshine duration.

INQC works applying a series of tests to the data. The result of each test (see Figure 1) is either *0 (pass)* or *1 (does not pass)*. At this point (beta version) no decision tool is provided, so users need to filter out those values which, according to the tests failed and their particular purpose, should not be considered for further climatological analyses.

| STAID | SOUID | date | value | weirddate | dupli | large | small | jump | flat | roundmax | friki | IQROUTliers | blocks | rounding | txtn |
|-------|--------|----------|-------|-----------|-------|-------|-------|------|------|----------|-------|-------------|--------|----------|------|
| 5490 | 136216 | 20080202 | 37 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 5490 | 136216 | 20080203 | 32 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 5490 | 136216 | 20080204 | 37 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 5490 | 136216 | 20080205 | 27 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 5490 | 136216 | 20080206 | 37 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 5490 | 136216 | 20080207 | 36 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 5490 | 136216 | 20080208 | 47 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 5490 | 136216 | 20080209 | 55 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 5490 | 136216 | 20080210 | 53 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 5490 | 136216 | 20080211 | 43 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 5490 | 136216 | 20080212 | 24 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 5490 | 136216 | 20080213 | 65 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 5490 | 136216 | 20080214 | 19 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 5490 | 136216 | 20080215 | -6 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 5490 | 136216 | 20080216 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 5490 | 136216 | 20080217 | 83 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 5490 | 136216 | 20080218 | 78 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 5490 | 136216 | 20080219 | 63 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 5490 | 136216 | 20080220 | 20 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 5490 | 136216 | 20080221 | 39 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 5490 | 136216 | 20080222 | 87 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 5490 | 136216 | 20080223 | 77 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 5490 | 136216 | 20080224 | 87 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 5490 | 136216 | 20080225 | 87 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 5490 | 136216 | 20080226 | 47 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 5490 | 136216 | 20080227 | 75 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 5490 | 136216 | 20080228 | 57 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |

Figure 1: INQC output example. Daily Maximum Temperature. Each column presents the result of one test applied to the data. The qc value is either 0 (pass) 1 (does not pass).

3. Preparing your computer to run INQC

INQC is designed to quality control ECA&D series. The requirements to run it are the following:

- R (developed and tested under RStudio Version 1.2.1069 and R version 3.3.2)
- An INQC folder :, e.g. `~/INQC`, This folder will store:
 - o The INQC code, stored as : `~/INQC/inqc_beta.R`
 - o ECA&D stations files (blended version), for each variable to be used, e.g. `~/INQC/ECA_blend_source_tx.txt` (these files can be downloaded from ECA&D)
- A quality control folder, named to your preference (e.g. Sweden for Swedish data): `~/Sweden`
- A raw data folder, created into your data: `~/Sweden/raw` [this folder name MUST be “raw”]. Raw data series must be non-blended ECA&D series (other formats are not supported and will not be supported, see <https://www.ecad.eu//dailydata/index.php> for information) and stored in this folder, with no sub-folders
- A qc'd data folder, where INQC will store the results: `~/Sweden/QC` [this folder name MUST be “QC”, capital letters]

NOTE: “~” stands for “any path before”

4. Wrapper Functions and Jump-Start option

After successfully completing the steps described in Section 2, INQC can be ran using the pre-set up defaults (see Table 1 in section 4 for full description) :

- Open R; set working directory to ~/INQC
- **inqc(homefolder = ~/INQC)**: Quality controlling all variables (maximum, minimum and average temperature, precipitation, sea level pressure, relative humidity, wind speed, snow depth, cloud coverage and sunshine duration)
- Quality controlling ONE variable:
 - **temperature(home=~/INQC/,element='TX')**: daily maximum temperature
 - **temperature(home=~/INQC/,element='TN')**: daily minimum temperature
 - **temperature(home=~/INQC/,element='TG')**: daily average temperature
 - **precip(home=~/INQC/)**: daily accumulated precipitation
 - **relhum(home=~/INQC/)**: relative humidity
 - **selepe(home=~/INQC/)**: sea level pressure
 - **snowdepth(home=~/INQC/)**: snow depth
 - **sundur(home=~/INQC/)**: sunshine duration
 - **windspeed(home=~/INQC/)**: windspeed

5. Quality Control Tests

Table 1: INQC Tests Description

| Test | Description and objective | Parameters |
|------|---------------------------|------------|
| | | |

| | | |
|--------------------------|---|---|
| <i>badfriki</i> | isolates extreme values which are not continuous in the distribution. If the gap is larger than a pre-set big margin, the value is flagged. | <p><i>date</i>: a vector of dates, in ECA&D YYYYMMDD format</p> <p><i>value</i>: the corresponding vector of values</p> <p><i>margin</i>: the maximum allowed difference between contiguous values in the empirical distribution</p> <p>call example:</p> <p><i>badfriki(date,value,margin=80)</i>, this call would flag values for which the difference with the preceeding value in the empirical distribution is larger than 8°C (expressed as 80 1/10ths of degree). For example, if the second largest value is 28°C (280) and the largest is 37°C (370), the later would be flagged as an outlying value</p> |
| <i>computecal</i> | produces a calendar with 3 variables: year, month, day between two given years. | <p><i>fy</i>: first year</p> <p><i>ly</i>: last year</p> <p>call example:</p> <p><i>computecal(fy=1900,ly=2018)</i>, would return a year,month,day dataframe with dates between 1900 and 2018</p> |
| <i>drywetlong</i> | detects episodes of too many consecutive wet or dry days. Uses a peak over threshold approach and a pareto distribution fit over the observed sequences | <p><i>x</i>: values</p> <p><i>ret</i>: pseudo return period for the POT-pareto, computed using the <i>parteogadget</i> auxiliary functions</p> <p>call example:</p> <p><i>drywetlong(x,ret=300)</i>, this would flag those sequences with longer length that to the 300 y pseudo-return period of with a pot-pareto approach, i.e. will flag “too long” dry or wet sequences.</p> |

| | | |
|---------------------------|--|---|
| | | |
| <i>duplas</i> | detects duplicated dates | <i>x</i> : a vector of dates in ECA&D format call example: <i>duplas(x)</i> , would flag any date appearing more than once |
| <i>flat</i> | detects consecutive equal values. Can be adapted to detect consecutive equal decimal part of the values | <i>y</i> : a data vector <i>maxseq</i> : the maximum number of contiguous repetitions of a value (e.g., if 3, sequences of 4 will be flagged) call example: <i>flat(y,maxseq=3)</i> , this would flag any streak of 4 or more consecutive values. |
| <i>IQRoutliers</i> | computes outliers centralized around a day, using a number of days around it and based on the Inter Quartile Range. Creates a tolerance interval centred around each day of the year, using all the present values in the empirical distribution for the designed window. Values outside the interval, are flagged as outliers | <i>date</i> : a vector of dates, in ECA&D YYYYMMDD format <i>value</i> : the corresponding vector of values <i>level</i> : number of IQR to be added to percentile 75 and subtracted to percentile 25 to determinate the tolerance interval. Values outside this interval, will be declared as outliers, <i>window</i> : an odd number representing the length of the window for which the outliers will be computed. Note: uses auxiliary function <i>julian</i> . call example: IQRoutliers(date,value,level=3>window=11), would flag outliers in value using a window of 11 days (e.g. for July 6 th : July, 1 st to July 11 th) |
| <i>jumps</i> | to label interdiurnal differences considered | <i>x</i> : vector of values |

| | | |
|---------------------|---|--|
| | to large | <p>maxjump: maximum difference allowed</p> <p>call example:</p> <p>jumps(x,maxjump=150) would flag all consecutive days for which the difference is 150 (e.g. 15°C for temperature, expressed as 150 1/10th s of degree)</p> |
| paretogadget | Returns the positions exceeding the value corresponding to a return period based on pareto distro and peak over threshold approach | <p>x: values</p> <p>ret: pseudo-return period for the pot-pareto distribution approach. Uses potpareto and returnpareto</p> <p>call example:</p> <p>paretogadget(x,ret=300), this would flag all values exceeding the value corresponding to the pot-pareto pseudo return period of 300 years</p> |
| physics | given a data vector, will compare the values to a specified threshold, considered to be the limit of physically possible values. In some cases. In some cases, the limitation is a consideration (e.g. 60°C), in others, it comes from the nature of the variable (e.g. 0 mm) | <p>x: vector of values</p> <p>nyu: comparison threshold, expressed in the same units of the ECA&D variable (e.g. in 1/10 of degree for temperature)</p> <p>compare: logical operation for the comparison of the vector of values to the threshold: 1 larger; 2 larger equal; 3 smaller; 4 smaller equal; 5 equal</p> <p>call example:</p> <p>physicscs(x,nyu=0,compare=3) would flag all values smaller than 0</p> |
| potpareto | Fits a pareto distribution to a series of values using as “threshold” the value representing a given | <p>y: values</p> <p>thres: quantile to compute the threshold</p> <p>call example:</p> |

| | | |
|------------------------|--|---|
| | quantile of the empirical distribution | potpareto(y,thres=0.99) , would fit a pareto distro using the quantile 0.99 of the y vector |
| putjulian | Adds julian calendar numbers, from 1 to 366 | x : a dataframe with year, month, day, value call example: putjulian(x) , will return a data frame with year, month, day, julian, value |
| repeatedvalue | This function tracks values which repeat too many times and, given the typical decaying distribution of the variable (designed for precipitation) are considered too large to repeat that many times | x : vector of values margin : the difference in frequency the nearest value friki : the minimum value to be considered call example: repeatedvalue(x,margin=20,friki=150) would flag any value larger than 15 mm (expressed as 150 1/10 th of mm) which repeats 20 times more than the previous value in the empirical distribution. For example, if 40 mm appears 25 times and the nearest value in the distribution is 38 and appears 5 times, all “40s” will be labelled. |
| returnpotpareto | For a given pareto distribution, returns the value representing a requested return period | pato : a pareto distribution fitted with potpareto ret : pseudo return period w : parameter to equate to return period to a temporal interval (recall the approach is not block maxima but peak over threshold. Typical value of w to equate the return period to years is 1.65 (See Wilks (2011), Statistical Analysis for the Atmospheric Sciences) call example: returnpotpareto(pato,ret=300,w=1.65) , would return the value associated to the return period |

| | | |
|------------------------|---|---|
| | | of 300 years. |
| rounding | splits data by month and looks if a decimal value is repeated too many times | <p>y: the vectors of values</p> <p>blocksize: the maximum number of equal decimal values allowed in a block</p> <p>call example:</p> <p>rounding(y,blocksize=20), would flag all occurrences of 20 or more values with the same decimal part in a month</p> <p>NOTE: monthly blocks are far from perfect, but they speed up the process, in comparison to sequential blocks. A fast way to do sequential blocks will be sought in future versions.</p> |
| roundprecip | splits data by month and looks if a decimal value is repeated too many times. A requested value can be excluded | <p>y: the vectors of values</p> <p>blocksize: the maximum number of equal decimal values allowed in a block</p> <p>exclude: the value to be excluded, for example in precipitation 0 should not be considered</p> <p>call example:</p> <p>rounding(y,blocksize=20, exclude=0), would flag all occurrences of 20 or more values with the same decimal part in a month, except for 0.</p> <p>NOTE: monthly blocks are far from perfect, but they speed up the process, in comparison to sequential blocks. A fast way to do sequential blocks will be sought in future</p> |
| suspectacumprec | Detects values above a threshold preceded by a given number of “no precipitation days” | <p>datos: a two columns vector, date and data, in ECA&D format</p> <p>limit: the value above which the function will search</p> <p>tolerance: the number of “non precip days”</p> |

| | | |
|----------------|---|---|
| | | <p>before the value checked that will result in flagging that value</p> <p>call example:</p> <p>suspectacumprec(datos=x[,3:4],limit=2000,tolerance=10), will flag all the values above 2000 (200 mm expressed in 1/10th of mm) which are preceded by days with no precip, either NA or 0.</p> |
| toomany | Splits data by month or year and looks if a value is repeated too many times | <p>y: two columns with date (in ECA&D format, YYYYMMDD) and vector of values</p> <p>blockmany: maximum number of values tolerated in a block</p> <p>scope: this variable controls whether the “block” are the months (1) or the years (2)</p> <p>exclude: defaulted to NULL, if specified will exclude the value or values specified. Takes a single value (e.g. 0, which should repeat many times in precipitation series) or could take a vector, expressed in the R vectorial form, e.g. exclude = c(0,0.1). Note: As an evolution, it is intended to add the possibility of excluding a range of values (e.g., smaller than 3)</p> <p>call example:</p> <p>toomany(y=x[,3:4], blockmany=15,scope = 1, exclude=0) , this call would label any value expect for 0, repeating more than 15 times in particular month.</p> |
| txtn | Compares daily maximum and daily minimum temperature and flags those values where TX is larger or | <p>y: a vector of values</p> <p>id: the file name, which is passed on to closestation auxiliary function to identify the “equivalent” tx or tn station. This is not trivial,</p> |

| | | |
|------------------|---|---|
| | equal than TN | <p>as ECA&D does not provide “direct relations”. See the auxiliary function for details</p> <p>home: home folder (this is used to locate and open the “equivalent” tx or tn station</p> <p>call example:</p> <p>txtn(y,id= TX_SOUID135829,home='./Sweden')</p> <p>This call would flag all the tx values in this series which are smaller or equal to the values of the series determined to be the corresponding TN series</p> |
| weirddate | Finds impossible dates (e.g. 19881420 or 19881131) or years out of the range of the range set by the first and the last records in the file | <p>x: vector of dates</p> <p>call example:</p> <p>weirddate(x), would return any existing “impossible” or out of range date.</p> |

6. Parametrization of wrapper functions and personalized runs of INQC

In this section we list the tests and preset values included for each wrapper function. A call to *inqc()* would run all variables with exactly these settings. For personalized settings, individual wrappers for each variable should be prepared. We provide one table for each variable.

Table 2 Default parametrization of the *temperature()* function, as ran by *inqc()*

| FUNCTION CALL | <i>temperature (home='../Sweden/',large=500,small=-500,maxjump=150,maxseq=3,margina=80,level=3>window=11,roundmax=10,blocksize=10,step=30,blockmanymonth=15,blockmanyyear=180,blocksizearound=20,element='TX')</i> | <ul style="list-style-type: none"> - The <i>home</i> parameter is superseded when this is called from <i>inqc()</i>; - The <i>element</i> parameter is altered for with “TN” and “TG” for daily maximum and daily minimum temperature respectively | |
|----------------------|---|--|---|
| Test | Parametrization | Variable in the qc'd file | Notes |
| <i>badfriki</i> | <i>margina</i> = 80 (sets the <i>margina</i> parameter) | friki | |
| <i>duplas</i> | - | duplas | |
| <i>flat</i> | <i>maxseq</i> = 3 (sets the <i>maxseq</i> parameter) | flat | |
| <i>flat</i> | <i>roundmax</i> = 15 (sets the <i>maxseq</i> parameter) | roundmax | This function is ran twice. The second call studies the decimal part (e.g. 15.0, 12.0, 10.0, 8.0 ...) are part of the same |

| | | | |
|--------------------|--|--------------|---|
| | | | “flat” sequence. For this reason the parameter is set to a larger value, 15 as default |
| IQROUTliers | level : 3 (sets the level parameter) window : 11 (sets the window parameter) | IQROUTliers | |
| jumps | maxjump = 150 (sets the maxjump parameter) | jump | |
| physics | large = 500 (sets the nyu parameter of the function) | large | Values above the parameter are flagged |
| physics | small = - 500 (sets the nyu parameter of the function) | small | Values below the parameter are flagged |
| rounding | blocksize round = 20 (sets the blocksize parameter) | rounding | |
| toomany | blockmanymonth = 15 (sets the blockmany parameter) | toomanymonth | Ran with scope =1 (not parametrized in the temperature() function), splitting the series by month |
| toomany | blockmanyyear = 180 (sets the blockmany parameter) | toomanyyear | Ran with scope =2 (not parametrized in the temperature()) |

| | | | |
|------------------|---|-----------|--|
| | | | function), splitting the series by month |
| txtn | - | txtn | Not ran for TG |
| weirddate | - | weirddate | |

Table 3 Default parametrization of the **precip()** function as ran by **inqc()**

| FUNCTION CALL | precip(home='~/INQC/',large=5000,small=0,ret=500,retornoracha=1000,margin=20,friki=150,blocksizearound=20,excludo=0,blockmanymonth=15,blockmanyyyear=180,exclude=0,limit=2000,tolerance=10,element='RR') | - The home parameter is superseded when this is called from inqc() ; | |
|------------------------|--|--|----------------------------|
| Test | Parametrization | Variable in the qc'd file | Notes |
| paretogadget | ret = 300 (sets the ret parameter) | paretogadget | |
| duplas | - | duplas | |
| suspectacumprec | limit = 2000 (sets the limit parameter) tolerance = 2000 (sets the tolerance parameter) | suspectacumprec | |
| repeatedvalue | margin = 20 (sets the margin parameter) friki = 150 (sets the friki parameter) | repeatedvalue | |
| roundprecip | blocksizearound = 20 (sets the blocksize parameter) excludo = 0 (sets the excluded parameter) | | |
| physics | large = 5000 (sets the nyu parameter of the function) | large | Values above the parameter |

| | | | |
|------------------|--|--------------|---|
| | | | are flagged |
| <i>physics</i> | <i>small</i> = 0 (sets the <i>nyu</i> parameter of the function) | small | Values below the parameter are flagged |
| <i>toomany</i> | <i>blockmanymonth</i> = 15 (sets the <i>blockmany</i> parameter) | toomanymonth | Ran with <i>scope</i> =1 (not parametrized in the temperature() function), splitting the series by month |
| <i>toomany</i> | <i>blockmanyyear</i> = 180 (sets the <i>blockmany</i> parameter) | toomanyyear | Ran with <i>scope</i> =2 (not parametrized in the temperature() function), splitting the series by month |
| <i>weirddate</i> | - | weirddate | |

TBD: tables for the other variables.

7. Testing INQC and

INQC will be tested using Baboon Benchmark

8. Expected evolution

- Additional functions

- Results interpreter and default decisions
- More comfortable parametrization of the jump-start functions
- Addition INDECIS' website and to GitHub repository, with sample data