

Best Practices in Research Data Management

Hauke Sonnenberg & Michael Rustler

2018-06-08

Contents

Introduction	5
I Best Practices	7
1 Data-Related Project Controlling	9
2 Naming Conventions	11
2.1 Acronyms	11
3 Folder Structures	13
3.1 Project Folder Structure	13
3.2 Rawdata Folder Structure	13
4 Metadata	15
4.1 Special Files	15
5 Raw Data	17
6 Data Processing	19
6.1 Automatic Data Import into R	19
II Old Chapters (German)	21
7 Einführung	23
8 Projektbeschreibungen	25
8.1 Projektgruppen	25
8.2 Projekt: Spree2011 (2007)	25
8.3 Projekt: MIACSO	26
8.4 Projekt: KURAS	26
8.5 Projekt: Flusshygiene	26
8.6 Projekt: DEMOWARE	26
8.7 Projekt:OPTIWELLS	26
8.8 Projekt: RWE	26
8.9 Projekt: AQUANES	27
9 Projektbausteine	29
9.1 Monitoring	29
9.2 Grundwassersimulation	29
9.3 Kanalnetzsimulation	30
9.4 LCA	30

10 Problembereiche	31
10.1 Datenaustausch	31
10.2 Datenablage	31
10.3 Datenimport	31
10.4 Metadaten	31
10.5 Zeitstempel	32
11 Allgemeine Methoden	35
11.1 Excel	35
11.2 Programmierung	35
11.3 Versionsverwaltung	35
11.4 Kollaboratives Schreiben	36
12 Konkrete Aufgaben und Lösungen	37
12.1 Monitoring: Erfassen von Metadaten	37
12.2 InfoWorks: Erzeugen von Regen-Inputdateien	37
12.3 InfoWorks: Erzeugen von Real Time Control (RTC)-Dateien	37
12.4 Erzeugen einer Ereignisliste {ereignisliste}	38

Introduction

This document describes best practices in research data management.

It is the outcome of the FAKIN project.

It is mainly based on the experiences of the authors.

Part I

Best Practices

Chapter 1

Data-Related Project Controlling

At the start of a research project

- Choose a Project Acronym and store it in `PROJECTS.txt`.
- Check if the organisations that you expect to get data from are listed in `ORGANISATIONS.txt` and extend the list in the file if necessary.
- Create a subfolder for your project and subfolders for the organisations in the Rawdata Folder Structure.

At the start of a project or if an employee or trainee enters the project

- Give an introduction to our research data management as described in this document.

Regularly during the project

- Check if the folder structure within your project's rawdata subfolder still complies with the Rawdata Folder Structure and clean the structure, if not.

Chapter 2

Naming Conventions

2.1 Acronyms

Acronyms are unique, clear names for objects. They should

- be short but meaningful and easy to remember,
- be all lowercase,
- consist of only alphanumeric letters (**a-z**, **0-9**) or the hyphen (**-**).

2.1.1 Acronyms for Projects

At the start of a project we define a project acronym. This acronym is intended to be used in file and folder names.

Whenever we want to indicate the relation to a certain project in a file or folder name, we use the project acronym in exactly the typing that was defined. This is important as we want to distinguish between raw data, processed data and project results in our Folder Structures.

The project acronyms are defined in a simple text file **PROJECTS.txt** in the **//server/projects\$** folder, see Project Folder Structure.

2.1.2 Acronyms for Organisations

It is very important to know the owners of data. Therefore we define unique acronyms for the owners of data that we use. The acronyms are defined in a special file **ORGANISATIONS.txt**

Chapter 3

Folder Structures

3.1 Project Folder Structure

We said that we want to concentrate on the folder structures within the project folders. Nevertheless, we would like to give a recommendation on how the project folders could be organised within its top level folder. In this structure, there are no subfolders for the different departments any more.

```
//server/projects$  
- PROJECTS.txt  
- project-1/  
- project-2/  
- project-3/
```

In the `projects$` folder

- each subfolder name should appear in the file `PROJECTS.txt`
- there should not be any folder on the top level that does not represent a project.
- there should be no other files on the top level as the files that are described in this documentation.
- there are no subfolders representing departments any more. The mapping of projects to departments is done in the file `PROJECTS.txt`

3.2 Rawdata Folder Structure

We will create a network folder `//server/rawdata$` in which all files have set the read-only property. We suggest to store raw data by project first and by the organisation that owns (i.e. generated, provided) the data second. This could look like this:

```
//server/rawdata$  
- ORGANISATIONS.txt  
- PROJECTS.lnk [Symbolic Link to PROJECTS.txt in //server/projects$]  
- flusshygiene  
  - bwB  
  - kwB  
  - uba  
  - ...  
- ogre  
  - kwB  
  - bwB
```

- uba
- ...
- ...

Restrictions/Conventions:

- Each top-level folder should represent a project, i.e. should be defined in the top level file **PROJECTS.txt**.
- Each possible owner should be defined in the top level file **ORGANISATIONS.txt**.
- The naming convention for the organisations is the same as for projects.

Chapter 4

Metadata

4.1 Special Files

We propose to define some special files that contain metadata related to files and folders. To indicate that these files have a special meaning, the file names are all uppercase.

4.1.1 File PROJECTS.txt

This file contains the project acronyms as we want to use them e.g. as top-level folder names in our project folder structure. The projects are grouped by department.

Possible content of PROJECTS.txt:

```
# Department SUW (Surface Water)
dswt: DSWT
flusshygiene: Flusshygiene
kuras: KURAS
mia-cso: MIACSO
monitor: MONITOR
ogre: OGRE
reliable-sewer: RELIABLE_SEWER
sema: SEMA
sema-berlin: SEMA Berlin
sema-berlin-2: SEMA Berlin 2
spree-2011: SPREE2011
spree-2011-2: SPREE2011 "reloaded"

# Department GRW (Groundwater)
optiwells: OPTIWELLS
optiwells-2: OPTIWELLS 2
wellma: WELLMA

# Department WWT (Wastewater Treatment)
...
```

In the file PROJECTS.txt the project acronyms appear in alphabetical order. They map the acronym to a project name or a project title and the year of the start of the project.

Question: Do we already have a place where “official” metadata about projects are stored? If yes, the acronym could be included there. But then, everybody should know about it!

4.1.2 File ORGANISATIONS.txt

Possible content of ORGANISATIONS.txt

```
bwb: Berliner Wasserbetriebe  
kwb: Kompetenzzentrum Wasser Berlin  
uba: Umweltbundesamt
```


Chapter 5

Raw Data

As raw data we define data that we receive from a device or from a project partner.

Most of our research results are based on data. We acknowledge the importance of raw data by

- storing them in a special place where it is specially secured
- describing them with metadata

Rawdata are stored in the rawdata folder structure

Chapter 6

Data Processing

6.1 Automatic Data Import into R

In the following we describe how data can be imported into the R-Programming Environment

6.1.1 Import Data From One Excel File

- Save the original file in the rawdata zone.

6.1.2 Import Data From Many Excel Files

6.1.2.1 Files Are In the Same Format

Import Excel files of the same format by

- defining a function that is able to read the data from that file
- calling this function in a loop for each file to import.

6.1.2.2 Files Are In Different Formats

We developed a general approach of importing data from many Excel files in which the formats (e.g. more than one table area within one sheet, differing numbers of header rows) differ from file to file.

Part II

Old Chapters (German)

Chapter 7

Einführung

Hier wollen wir die Ergebnisse unseres Forschungsprojekts FAKIN (Forschungsdatenmanagement an kleinen Instituten) veröffentlichen.

Da wir viele, zum Teil ganz unterschiedliche Projekte bearbeiten, sind auch unterschiedliche Bereiche des Forschungsdatenmanagements betroffen.

Wir wollen hier zunächst versuchen, Fallbeispiele herauszuarbeiten, in die sich die verschiedenen Projekte einordnen lassen. Für jedes Fallbeispiel wollen wir dann Gemeinsamkeiten in Bezug auf Anforderungen und Problemstellungen herausarbeiten und für diese dann existierende Lösungen einsammeln oder gezielt neue Lösungen entwickeln. Diese Lösungen sollen dann in zukünftigen Projekten, die sich den jeweiligen Fallbeispielen zuordnen lassen, angewendet werden.

Dazu brauchen wir Eure Mitarbeit! Wer kann Beiträge zur Beschreibung von Problemfällen und eigenen Lösungen liefern?

Chapter 8

Projektbeschreibungen

8.1 Projektgruppen

Nach welchen Kriterien können wir unsere Projekte gruppieren, wie können wir sie grob klassifizieren?

Was folgt daraus in Bezug auf das Datenmanagement?

Wenn wir Ähnlichkeiten finden, können wir in Zukunft gezielt auf die Lösungen, die in ähnlichen Projekten erzielt (und beschrieben!) wurden, zurückgreifen.

Mögliche Gruppierung nach:

- Abteilung: Grundwasser, Oberflächengewässer und Kanal, Abwasser
- Finanzierung: Sponsoring, EU-Projekt, BMBF-Projekt, Auftrag
- Anzahl Partnern: große Projekte mit vielen Partnern, kleine Projekte mit wenigen Partnern

Im folgenden fangen wir einfach mal an und versuchen, die Projekte, an denen wir mitgearbeitet haben, kurz zu beschreiben und zu charakterisieren. Wenn wir im Rahmen der Projekte bestimmte Lösungen erarbeitet haben, beschreiben wir diese kurz unter Problembereiche oder Konkrete Aufgaben und Lösungen.

Aufbauend auf dieser ersten Aufstellung von Projekten, Aufgabenbereichen und Lösungen wollen wir dann mit Euch zusammen weitere Projekte in die Betrachtung mit einbeziehen. Wir wollen gemeinsam die Charakterisierung und Gruppierung und natürlich insbesondere die Aufgabenbeschreibungen und Lösungen erweitern und verbessern.

8.2 Projekt: Spree2011 (2007)

- Meine Diplomarbeit
- eigene Messungen über einen begrenzten Zeitraum
 - Wasserstand und Durchfluss in einem Kanalschacht -> Text/CSV
 - Regen -> Text/CSV
- kontinuierliche Messungen der Berliner Wasserbetriebe
 - Fördermengen der Pumpwerke
 - Wasserstände in den Pumpwerken
 - Regen an benachbarten Regenschreibern

Aufgaben und angewandte Methoden nach Themen

- Modellierung/Kanalnetz/Kalibrierung:
 - InfoWorks: Erzeugen von Regen-Inputdateien

- InfoWorks: Erzeugen von Real Time Control (RTC)-Dateien

Aufgetretene Fragen:

- Wo werden Präsentationen abgelegt?
- Wo werden die Rohdaten abgelegt? Z.B. Persönlicher Ordner “Diplomarbeit”
- Worauf beziehen sich die InfoWorks-Zeitstempel? -> Themenbereich Metadaten

Best of aus meiner persönlichen Logdatei:

8.3 Projekt: MIACSO

8.4 Projekt: KURAS

8.5 Projekt: Flusshygiene

8.6 Projekt: DEMOWARE

Entstandene R Pakete:

- Grundwassermodellierung
 - kwb.hantush
 - kwb.vs2dh
 - kwb.demoware
- Quantitatives mikrobiologisches Risikomanagement
 - kwb.qmra: wird im Rahmen von AQUANES(#aquanes) weiter genutzt

8.7 Projekt:OPTIWELLS

Entstandene R Pakete:

- Grundwassermodellierung
 - kwb.wtaq
- (Druck)Rohrleitungsnetzsimulation (EPANET)
 - kwb.epanet

8.8 Projekt: RWE

(Semi)automatisierte Erstellung eines komplexen MODFLOW Modells (mehrere Layer, mehrere hunderte Entnahmefrühen mit zeitlich variierender Entnahmemenge sowie Hinzufügen/Entfernen von Brunnen innerhalb des Simulationszeitraums) in Python mittels flopy sowie Entwicklung der Modellszenarien auf Github (siehe hier: Maxflow).

Input Christian !?!

8.9 Projekt: AQUANES

Entstandenes R Paket:

- aquanes.report: Datenimport, zeitliche Aggregation und interaktive Visualisierung von Betriebsdaten von Pilotanlagen (u.a. Berlin-Tiefwerder, Berlin-Schönerlinde, Basel Lange-Erlen, Haridwar) sowie deren Verknüpfung mit Labordaten.
 - Herausforderungen: zeitlich hoch aufgelöste Betriebsdaten (~ 10 Mio. Datenpunkte pro Monat) erforderten massive Performance Optimierung um die Visualisierung der Rohdaten im Tool über den Testbetriebszeitraum (18 Monate) auf Rechnern mit limitierten RAM Ressourcen (~ 8 GB) zu ermöglichen.
 - Nutzung: wird von den Projektpartnern zum strukturierten Datenimport genutzt um dann darauf ggf. eigene weitere Analysen darauf aufzusetzen. Für die beiden Berliner Standorte wird zudem die Visualisierung der Betriebsdaten routinenäßig genutzt, um vom Regelbetrieb abweichende Zustände besser indentifizieren zu können.

Chapter 9

Projektbausteine

Wir können “Projektbausteine” definieren, die in verschiedenen Projekten vorkommen. Diese werden im Folgenden einzeln betrachtet.

9.1 Monitoring

In vielen Projekten führen wir eigene Messungen durch. Dazu verwenden wir eigene Messgeräte.

In Monitoringprojekten spielen Metadaten eine besonders wichtige Rolle. Wir betreiben viele Messgeräte mit Akkus. Der Zeitpunkt des Austauschs des Akkus ist eine wichtige Metainformation. Wenn Akkus zu spät gewechselt werden können wichtige Daten verloren gehen bzw. verpasst werden! Auch die Reinigung von Messgeräten ist wichtig, z.B. von Sonden mit optischen Messfenstern. Diese muss einerseits regelmäßig erfolgen. Andererseits ist es wichtig zu wissen, wann Reinigungen erfolgten, damit Messwerte entsprechend interpretiert werden können.

Im Projekt MIA-CSO haben wir erst angefangen, diese Informationen in einer Exceldatei zu erfassen. Angesichts der vielen Messgeräte, Pumpen, Kompressoren und anderer Geräte wurde das aber bald unübersichtlich. Später habe ich eine MS Access Anwendung dafür entwickelt, die seither in einigen Projekten zum Einsatz kam.

Problemstellungen:

- verschiedene Geräte verschiedener Hersteller -> verschiedene, meist nicht standardisierte Dateiformate
- Zeitstempel der Messung
 - Stellt das Gerät auf Sommerzeit um? -> Problembereich Zeitstempel
 - Geht die interne Uhr immer richtig? Wie gewährleiste ich das?

Beispielprojekte: MIACSO, KURAS, Flusshygiene

9.2 Grundwassersimulation

Beispielprojekte: OPTIWELLS, DEMOWARE, RWE

9.3 Kanalnetzsimulation

9.3.1 Datenaustausch mit BWB

Bei der Kanalnetzsimulation arbeiten wir eng mit den Berliner Wasserbetrieben zusammen. Diese haben und pflegen die Modellnetze.

Eine Schwierigkeit besteht im Austausch von Netzen und in der Versionsverwaltung. In welcher Version wurde ein Netz geschickt?

Gibt es dafür ein halbwegs vernünftiges Vorgehen? Wie ich es sehe, besteht eine große Gefahr, dass beim Hin- und Herschicken von Modellnetzen zwecks Abstimmung versehentlich wieder auf eine ältere Version zurückgegriffen wird und neueste Änderungen auf einer alten Version aufsetzen anstatt auf der aktuellen.

Neben Modellnetzen werden auch andere Bausteine, z.B. Regenreihen oder RTC-Bausteine ausgetauscht.

Fragen:

- Wie wird der Datenaustausch dokumentiert?
- Wie werden die ausgetauschten Objekte versioniert? Das müsste einheitlich auf beiden Seiten des Austauschs geschehen.

9.3.2 Modellkalibrierung

Trockenwetterkalibrierung

Die Kalibrierung der Abflüsse bei Trockenwetter wird anhand einer mittleren Trockenwetterganglinie gemacht. Dafür müssen Trockenwettertage aus den Messdaten ausgewählt werden. Wie wird das gemacht? Ist nachvollziehbar, nach welchen Kriterien die Auswahl erfolgte? Wie wird der Mittelwert über die Ganglinien berechnet? Was geschieht mit Ausreißern, was mit fehlenden Werten? Ist die Auswahl und Berechnung reproduzierbar und lässt sie sich leicht auf neue Daten anwenden? Das wäre eine typische Problemstellung, die sich mit Programmierung lösen lässt. Ggf. gibt es aber auch funktionierende Excellösungen dafür. Wie macht ihr das?

Bei der Modellkalibrierung wird versucht, die Modellparameter, die die Abflussbildung auf der Oberfläche beschreiben, so einzustellen, dass die simulierten Abflüsse und Wasserstände den gemessenen Abflüssen und Wasserständen entspricht. Dazu müssen die gemessenen Regenmengen so aufbereitet werden, dass sie in das Modell importiert werden können. Die erste Aufgabe besteht darin, die vollständige Regenmessreihe in Ereignisse zu zerlegen. Das ist wichtig, da InfoWorks ereignisbezogen z.B. die Evaporation berechnet. Nach der Zerlegung in Ereignisse müssen die Ereignisdaten noch in das richtige Format gebracht werden.

Siehe

- Erzeugen einer Ereignisliste
- InfoWorks: Erzeugen von Regen-Inputdateien

Beispielprojekte: MIACSO, KURAS, Flusshygiene

9.4 LCA

LCA steht für Life Cycle Assessment, bedeutet also soviel wie Lebenszyklusbewertung.

Chapter 10

Problembereiche

10.1 Datenaustausch

Es gibt internen Datenaustausch und Austausch mit Projektpartnern. Oft werden Dateien intern per E-Mail verschickt. Das sollten wir nicht tun, da wir unnötig Kopien von Dateien anlegen. Dateien sollten im besten Fall nur an einem Ort abgelegt sein. Dort sollten sie allerdings durch ein Backup-System gesichert sein. Natürlich muss der Zugriff auf diesen Ort gewährleistet sein. Einem Studenten mit eingeschränkten Rechten eine Datei per mail zuzuschicken, weil sie an dem Ort, an dem sie gespeichert ist, für ihn nicht zugänglich ist, ist meines Erachtens nicht die richtige Lösung. Es muss dann ein Ort geschaffen werden, an dem der gemeinsame Zugriff möglich ist. Laut Bodo ist der Ordner **Exchange** im Projektordner ein solcher Platz. Aus dem Grund der allgemeinen Zugänglichkeit wird er auch in manchen Projekten als allgemeiner Datenablageplatz “misbraucht”. Es kann aber niemandem ein Vorwurf gemacht werden, da ja nirgendwo dokumentiert ist, was in diesem Ort abgelegt werden soll.

Das können wir nun ändern:

Laut Bodo sind Dateien im Exchange-Ordner nur kurzfristig abzulegen, um z.B. Studenten Zugriff darauf zu geben. Die Dateien sollen nach dem Gebrauch wieder aus dem Ordner gelöscht werden. Demnach sollte der Ordner im Normalfall leer sein. Dies ist in einigen Projekten nicht der Fall.

Wir brauchen eine einheitliche Definition der Bedeutungen von Ordnern.

Wie die zu dokumentieren sind sollten wir in Metadaten beschreiben.

10.2 Datenablage

10.3 Datenimport

10.4 Metadaten

Beispiel:

Wenn wir Zeitreihen von den Berliner Wasserbetrieben bekommen, dann haben diese entweder eine Spalte mit Zeitstempeln oder zwei Spalten mit Zeitstempeln, wobei die erste den Anfang und die zweite das Ende des Zeitraums enthält, auf den sich die Messwerte in der entsprechenden Zeile beziehen.

Wenn nur eine Spalte Zeitstempel enthält, dann ist nicht klar, ob der Zeitstempel den Anfang oder das Ende (oder gar einen anderen Zeitpunkt innerhalb) des Zeitraums angibt.

Bei einer Zeitreihe mit fünfminütigen Messwerten, z.B.

Zeit	Wert
2018-06-05 16:20:00	0
2018-06-05 16:25:00	1
2018-06-05 16:30:00	2
2018-06-05 16:35:00	3
2018-06-05 16:40:00	4

ist nicht klar, ob sich der Wert 1 auf den Zeitraum 16:20 - 16:25 oder auf den Zeitraum 16:25 - 16:30 bezieht.

Dies ist ein Beispiel für eine sehr konkrete Metadatenanforderung. An dieser Stelle gehen wir meistens von Annahmen aus oder hat jemand von Euch schon mal eine Aussage darüber von einem BWB-Mitarbeiter erhalten?

Die Kenntnis über den Bezugspunkt des Zeitstempels ist zum Beispiel wichtig beim

- Erzeugen von Regen-Inputdateien für die Software InfoWorks

In der obigen Beispielzeitreihe ist auch nicht klar, auf welche Zeitzone sich die Zeitstempel beziehen. Die Zeitzone bzw. die Angabe des UTC-Offsets ist eine wichtige Metainformation. Am besten wäre es, wenn alle Zeitstempel diese Information enthielten, z.B.

```
2018-06-05 16:20:00+02
2018-06-05 16:25:00+02
```

Dann bräuchte man diese Information nicht extra zu führen.

Metadaten werden oft erst bei Bedarf erhoben, nämlich genau dann, wenn man sich z.B. die obigen Fragen nach Zeitstempelbezug und Zeitzone stellt. Dann wird nachgefragt und im besten Fall wird eine Antwort darauf gegeben, z.B. in Form einer E-Mail. Die Metadaten sind dann in der E-Mail enthalten.

Fragen:

- Wo und unter welchem Namen legen wir diese E-Mail ab?

Beispiel für Metadaten aus meiner persönlichen Logdatei:

Regendaten gelesen; ca. 10 Impulse gegen 14:00; Dannecker geschrubbt, neuer Akku: Batterie defect, alte Batterie drin gelassen, Daten auslesen nicht möglich; H=55cm

Das müssen wir irgendwie besser hinkriegen, aber wie?

Die grundlegenden Fragen sind:

- Wer ist für die Ablage von Metadaten verantwortlich?
- Wo werden die Metadaten abgelegt? Direkt bei den Daten oder in einer übergeordneten Datei/Datenbank?
- In welchem Format werden die Metadaten abgelegt? Einfaches Textformat vs. XML-Format, das sich wohl einfach nur mit Hilfe eines Programms editieren lässt.

10.5 Zeitstempel

Frage:

Wie stellen wir die Messgeräte ein? Sollen sie die Uhr automatisch umstellen oder nicht?

Diskussion:

Aus Sicht der Datenverarbeitung bereitet uns die Umstellung von Winterzeit (= Normalzeit) auf Sommerzeit und von Sommerzeit auf Winterzeit immer wieder Schwierigkeiten. Warum?

Umstellung von Winterzeit auf Sommerzeit

Die Uhr wird um eine Stunde von 02:00 CET (*Central European Time*) auf 03:00 CEST (*Central European Summer Time*) **vorgestellt**. Merke: Die Stühle eines Cafés werden **vor** den Laden gestellt.

In einer Zeitreihe ergibt sich daraus am Tag der Zeitumstellung (z.B. am 25.03.2018) eine Lücke in den Zeitstempeln:

```
2018-03-25 01:45:00 CET
2018-03-25 01:50:00 CET
2018-03-25 01:55:00 CET
2018-03-25 03:00:00 CEST
2018-03-25 03:05:00 CEST
2018-03-25 03:10:00 CEST
```

Wenn die als Text vorliegenden Zeitstempel in einen numerischen Zeitstempel umgewandelt werden, gibt es, auch wenn die Angabe CET oder CEST fehlt, im Grunde kein Problem, da sie eindeutig sind. Voraussetzung ist, dass bei der Umwandlung angegeben wird, dass diese Zeitstempel in der Zeitzone “Europe/Berlin” (in der es Sommer- und Winterzeit gibt) aufgenommen wurden.

In der Programmiersprache R ist das auf unseren Rechnern das Standardverhalten, da die Zeitzone vom Betriebssystem abgefragt wird:

```
as.POSIXct("2018-03-25 01:55:00")
## [1] "2018-03-25 01:55:00 UTC"
as.POSIXct("2018-03-25 03:00:00")
## [1] "2018-03-25 03:00:00 UTC"
```

Interessant an dieser Stelle ist das Verhalten, wenn ein Zeitstempel, der in der Zeitzone “Europe/Berlin” nicht existiert, angegeben wird:

```
as.POSIXct("2018-03-25 02:30:00")
## [1] "2018-03-25 02:30:00 UTC"
```

Leider tritt an dieser Stelle kein Fehler auf!

Umstellung von Sommerzeit auf Winterzeit

Die Uhr wird um 03:00 CEST wieder um eine Stunde auf die Normalzeit 02:00 CET **zurückgestellt**. Merke: Die Stühle werden wieder in den Laden **zurück** gestellt.

In einer Zeitreihe ergeben sich daraus am Tag der Zeitumstellung (z.B. am 28.10.2018) Duplikate in den Zeitstempeln:

```
2018-10-28 02:15:00 CEST
2018-10-28 02:30:00 CEST
2018-10-28 02:45:00 CEST
2018-10-28 02:00:00 CET
2018-10-28 02:15:00 CET
2018-10-28 02:30:00 CET
```

Ohne die Information CEST bzw. CET sind die Zeitstempel (für sich genommen) nicht eindeutig. In Dateien, die wir von den BWB bekommen (z.B. Regendaten) würden die Zeitstempel zum Beispiel in dieser Form vorkommen:

```
28.10.2018 02:15:00
28.10.2018 02:30:00
28.10.2018 02:45:00
28.10.2018 02:00:00
28.10.2018 02:15:00
28.10.2018 02:30:00
```

Erst aus der Reihenfolge (erstes Auftreten ist CEST, zweites Auftreten ist CET) lässt sich die Zuordnung rekonstruieren. Dies müssen wir beachten, wenn wir die Daten bearbeiten und ich denke, dass wir dies in den seltensten Fällen tun!

Wir sollten eine einheitliche, eindeutige Vorgehensweise entwickeln.

TODO: Ggf. weitere Beschreibungen in R-Workshop, den ich mal gegeben habe...

Vorschlag:

- Nach Möglichkeit die Messgeräte so einstellen, dass der ausgegebene Zeitstempel den UTC-Offset enthält. Am besten wäre ein Zeitstempel in einem Format, das die Differenz gegenüber *Greenwich Mean Time (GMT)* bzw. *Coordinated Universal Time (UTC)* enthält, wie wir es ggf. aus unseren E-Mail-Programmen oder Terminkalendern kennen: 2018-06-01 23:18 GMT+02:00.
- Manchmal gibt es auch die Möglichkeit, dass zusätzlich zum lokalen Zeitstempel der UTC-Offset (+01 im Winter und +02 im Sommer) als extra Spalte ausgegeben wird. Von dieser Möglichkeit sollte dann Gebrauch gemacht werden.
- Ansonsten empfehlen wir, die Geräte so einzustellen, dass sie die Zeitumstellung nicht mitmachen (also immer die Normalzeit, also Winterzeit, aufzeichnen). Dann sind die Zeitstempel eindeutig. Aufpassen ist dann allerdings bei der Verarbeitung und Interpretation der Daten, dazu mehr unter Datenimport

Werkzeuge:

- KWB-eigenes R-Paket `kwb.datetime`. Dieses enthält unter anderem die Funktionen
 - `date_range_CEST()` zur Ermittlung der Tage, an denen die Zeitumstellung stattfindet.

TODO:

- Wichtigste Funktionen des Pakets `kwb.datetime` (z.B. zur Ermittlung der vollständigen Zeitstempel inklusive UTC-Offset aus fortlaufenden lokalen Zeitstempeln) dokumentieren und hier referenzieren.
- Sollten wir andere Pakete benutzen und wenn ja, welche? Was bietet z.B. `lubridate`? Es scheint mir allerdings, dass unsere Problematik im WWW etwas unterbeleuchtet ist.

Chapter 11

Allgemeine Methoden

11.1 Excel

Wie arbeite ich “sauber” mit Excel?

Best Practices zum Arbeiten mit Excel:

- z.B. Zellbezüge benennen, dadurch werden Formeln besser lesbar

Ihr werdet es kaum glauben, aber am Anfang meiner Zeit am KWB habe ich noch mit Excel gearbeitet. Ich habe auch komplexe Sachen gemacht und auch Excel-Makros programmiert. Ein Beweis findet sich in meiner persönlichen Logdatei:

Fr, 14.09.07 08:15-19:30 benutzerdefinierte Excelfunktionen zur Modellgüte in Personl.xls,
Modul1; TW-Kalibrierung Wochentag fertig

Ich würde heute nicht mehr empfehlen, Excel-Makros zu programmieren. Es ist umständlich, es gibt keine Bibliotheken und der Quellcode lässt sich nicht unabhängig von der Exceldatei verwalten, so dass keine ordentliche Versionsverwaltung möglich ist. Und wir wollen nicht mehr ohne Versionsverwaltung programmieren!

Wenn wir programmieren, dann sollten wir das einheitlich in R tun. R ist frei, es gibt eine großartige Community und wir haben mittlerweile eine große Expertise erlangt.

11.2 Programmierung

Wir verwenden die freie Programmiersprache R. Aller Programmcode soll unter Versionsverwaltung stehen. Das ist angesichts unserer über 1000 Skripte und komplexer Abhängigkeiten zwischen den Paketen und Skripten unabdingbar.

11.3 Versionsverwaltung

11.3.1 Subversion

Wir verwenden einen KWB-internen Subversion Versionsverwaltungsserver, mit dem wir über die Software Tortoise SVN, die direkt in den Dateexplorer integriert ist, kommunizieren.

Im Rahmen von FAKIN legen wir fest, dass alle Programmcodes unter Versionsverwaltung gestellt werden sollen. Das gehört rein ins QMS!

11.3.2 Git

Wir haben auch einen öffentlich zugänglichen GitHub Account:

<https://github.com/KWB-R>.

Auf diesem legen wir nach und nach R-Pakete ab, deren Veröffentlichung keinen Beschränkungen unterliegt, z.B. weil der Code im Rahmen eines Auftragsprojektes erstellt wurde.

11.4 Kollaboratives Schreiben

Herausforderung:

Projektantrag mit vielen Partnern. Wie kann gewährleistet werden das immer am aktuellen Dokument gearbeitet wird?

Lösungsmöglichkeit:

Ulf hat gute Erfahrungen mit Office365 (von KWR zur Verfügung gestellt) gemacht. Auch das parallele Arbeiten von knapp 10 Personen gleichzeitig stellt kein Problem da, da von einem in Bearbeitung befindliche Abschnitte automatisch für andere blockiert werden.

Wenn größere Abschnitte von einem zu bearbeiten sind empfiehlt es sich diese kurzzeitig aus dem synchronisierten Dokument zu entfernen, die Änderungen in einer nicht synchronisierten Kopie der Word Datei zu machen und erst zum Schluss wieder im Dokument einzufügen.

Ulf strebt an, dass auch am KWB eine Office365 Lizenz erworben wird.

Chapter 12

Konkrete Aufgaben und Lösungen

12.1 Monitoring: Erfassen von Metadaten

Zur Protokollierung von Routine- und Wartungsarbeiten an den Geräten in Messkontainern habe ich eine komplexe MS Access Anwendung geschrieben. Mit dieser können Messstellen, Geräte, Aktionen (an Geräten) und Personen (die die Aktionen durchführen) verwaltet sowie Aktionen geplant und dokumentiert werden.

Probleme und Fragen

- Wo liegt sie?
- Wie wird sie bedient? Es gibt keine Beschreibung.
- Die Software wird nicht gepflegt
- Sollten wir sie weiter verwenden? Bestimmt gibt es mittlerweile Apps (womöglich kostenfrei), die auf Tablets laufen, benutzerfreundlich und gut dokumentiert sind.

12.2 InfoWorks: Erzeugen von Regen-Inputdateien

Methode: Perl-Skript `iw_event_writer.pl`

- Wo liegt dieses Skript? Erst lokal bei Hauke, mittlerweile unter Subversion, wurde glaube ich von Mathias in ein R-Skript umgewandelt...
- Wie kann es verwendet werden? Dazu benötigt man Perl, das ist standardmäßig auf keinem Rechner installiert. Auf einem Linux-Rechner wäre das der Fall... Sollten wir einen Linux-Rechner für solche Fälle vorhalten? Nun haben wir ja Rupelton...

12.3 InfoWorks: Erzeugen von Real Time Control (RTC)-Dateien

Ziel: Steuerung der Pumpwerke im Modell, wie sie in der Realität gefördert haben.

Benötigte Daten:

- Regendaten der Berliner Wasserbetriebe

Benötigte Metadaten:

- Worauf bezieht sich der Zeitstempel in einer kontinuierlichen Zeitreihe, auf den Anfang oder das Ende des Zeitintervalls mit der Länge eines Zeitschritts?

- Auf welchen UTC-Offset bezieht sich der Zeitstempel? Meiner Erfahrung nach ist bei den Daten der Berliner Wasserbetriebe immer der Lokalzeitstempel angegeben, in ihren Datenreihen gibt es also an den Tagen der Zeitumstellung Lücken bzw. Duplikate.

Manuelle Methode:

- RTC-Modul in Infoworks-Software erzeugen
- RTC-Modul als Text exportieren
- Pumpwerksdaten nach CSV konvertieren
- Pumpwerksdaten im Texteditor mit Suchen/Ersetzen auf das von InfoWorks geforderte Format bringen
- Formatierte Pumpwerksdaten in der RTC-Textdatei an den richtigen Stellen einfügen

Automatische Methode:

- Perl-Skript `rtc_generate.pl`. Dieses macht die obigen Schritte automatisch.
- Problem: lässt sich nur auf Rechner mit Perl laufen lassen.

12.4 Erzeugen einer Ereignisliste {ereignisliste}

Es können Funktionen aus dem Paket `kwb.event` verwendet werden.

TODO: Schreiben einer Vignette, die einem einen Leitfaden gibt. Es könnte ein Ausschnitt aus einem existierenden Skript verwendet werden.

Wir brauchen unbedingt Testdaten, die wir verwenden können, um bestimmte Dinge zu veranschaulichen oder als Eingangsdaten für den Test von Funktionen.

Aber wo legen wir die Daten hin und dürfen wir das eigentlich? Oder sollten wir künstliche (generierte) Daten verwenden?