



Scientists who stare at data

Dipl.-Ing. Jan Grüner | MathWorks Automotive Conference 2019



About

TU Berlin // FVB

- Chair of Naturalistic Driving Observation for Energetic Optimisation and Accident Avoidance
- Focus on Electromobility

www.fvb.tu-berlin.de

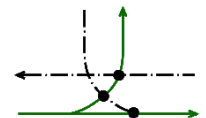
Myself

- Research Assistant
- Matlab since 2009
 - Incl. several years of teaching Matlab
- PHD Thesis: “Usage behavior of hybrid vehicles”

jan.gruener@tu-berlin.de

Matlab Projects

- IDCB (Create individual driving cycles from user data)
- Database Toolbox (wrapper functions for [MySQL Database Connector](#))
- AMPERE (Usage behavior of hybrid vehicles)
- WeatherDB (Weather DB / API to make use of data provided by the DWD)

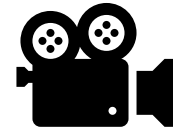




We are living in the information processing age

We record everything & everywhere

- Cameras
- Cars
- Smart-X (Phones, Homes)
- IoT



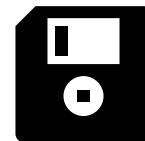
We can transmit data from A to B

- Worldwide
- Instantly / Fast
- Everything is connected



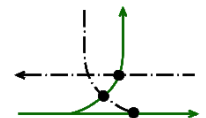
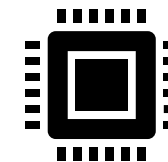
We can store huge amounts of data

- Fast storage
- Large storage



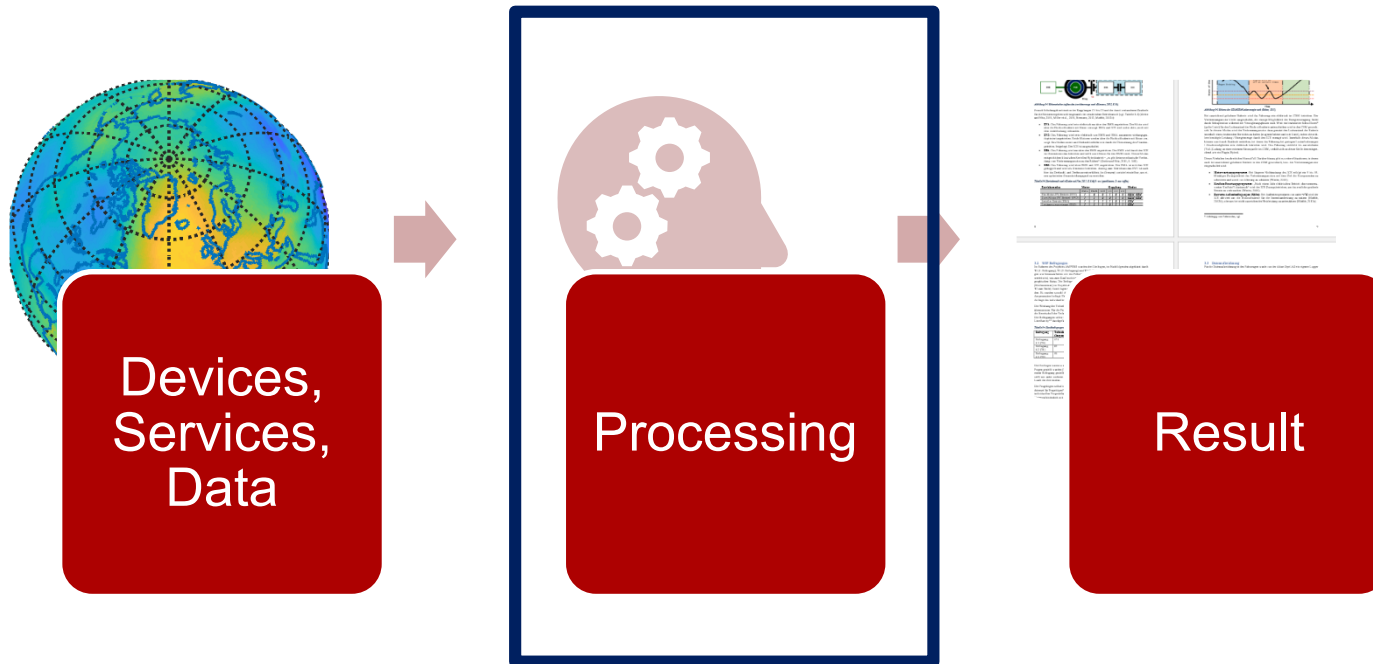
We have the computational power to analyze it

- Server
- Cluster
- Databases
- Software



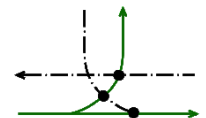


The challenge is no longer getting data, it's processing it



(and making sense of)

- With endless possibilities comes complexity
- Complexity creates chaos
 - Data formats
 - Competing standards
 - Methods
- A huge amount of information





In tech / data we (blindly) trust

tElement_noData		ukt_zehn_min_ff_19940505_19991231_00788.txt				16											
27106x16 table		STATIONS_ID;MESS_DATUM;QN;FF_10;DD_10				neStampSortHuman											
26851	15x2 string	788;199405050820;	1;	1.6;	100	2017-10-17 07:01:45.578											
26852	15x2 string	788;199405050830;	1;	2.1;	90	2017-10-17 07:01:46.866											
26853	15x2 string	788;199405050840;	1;	1.9;	90	2017-10-17 07:01:48.153											
26854	15x2 string	788;199405050900;	1;	1.5;	100	2017-10-17 07:01:49.443											
26855	15x2 string	788;199405050910;	1;	1.4;	130	2017-10-17 07:01:49.443											
26856	15x2 string	788;199405050920;	1;	1.2;	140	2017-10-17 07:01:50.748											
26857	15x2 string	788;199405050930;	1;	1.5;	170	2017-10-17 07:01:52.035											
26858	15x2 string	788;199405050940;	1;	1.8;	180	2017-10-17 07:01:53.323											
26859	15x2 string	788;199405050950;	1;	1.9;	180	2017-10-17 07:01:54.611											
26860	15x2 string	g_2014_U2_U6_1230_U01_gm1/U1s45c940_1.csv				17 07:01:55.897											
26861	15x2 string	1	timestamp;	valid;	Fahrmodus;	valid;	Energie APM aus HV	17 07:01:55.898									
26862	15x2 string	2	s;	;	Wh;	h;	m;	;	m/s^2;	km;	1/min;	1/min;	1/min;	17 07:01:57.184			
26863	15x2 string	3	13951;	1;	Normal;	1;	0.058;	1;	12;	1;	30;	1;	EV1;	1;	0.0;	-1;	17 07:01:58.472
26864	15x2 string	4	14202;	-1;	;	;	;	;	;	;	;	;	;	;	;	;	17 07:01:59.759
26865	15x2 string	5	14451;	-1;	;	;	;	;	;	;	;	;	;	;	;	;	17 07:02:11.356
26866	15x2 string	6	14702;	-1;	;	;	;	;	;	;	;	;	;	;	;	;	17 07:02:12.643
26867	15x2 string	7	14952;	1;	Normal;	1;	0.141;	-1;	;	;	;	;	;	;	;	;	17 07:02:13.931
26868	15x2 string	8	15203;	-1;	;	;	;	;	;	;	;	;	;	;	;	;	17 07:02:15.220

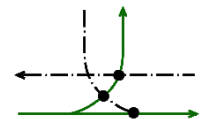
6 Columns documented ...

Consecutive datastream

Wrong unit

The pessimistic view on hardware / services

- Sensors are cheap, the device isn't
- Updates?
- Security?
- Testing?
- Documentation?
- Support?
- May change at any time

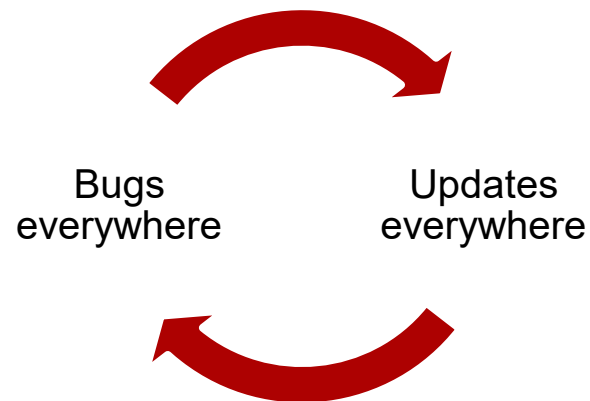




Question the data and your work (and everyone else)

The situation

- Assumptions → Can become false over time
- Documentation → Outdated (best case)
- Human error → nobody is perfect

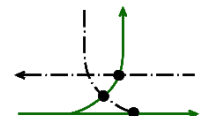


- (Variable-)Names will not change over time
- The documentation is complete

- "
- The RMS value is inserted into a circular list of 401 entries
 - No event is fired until this list is initially full
 - **TBD** should this list be persisted after each reboot?
- "
- Final document**

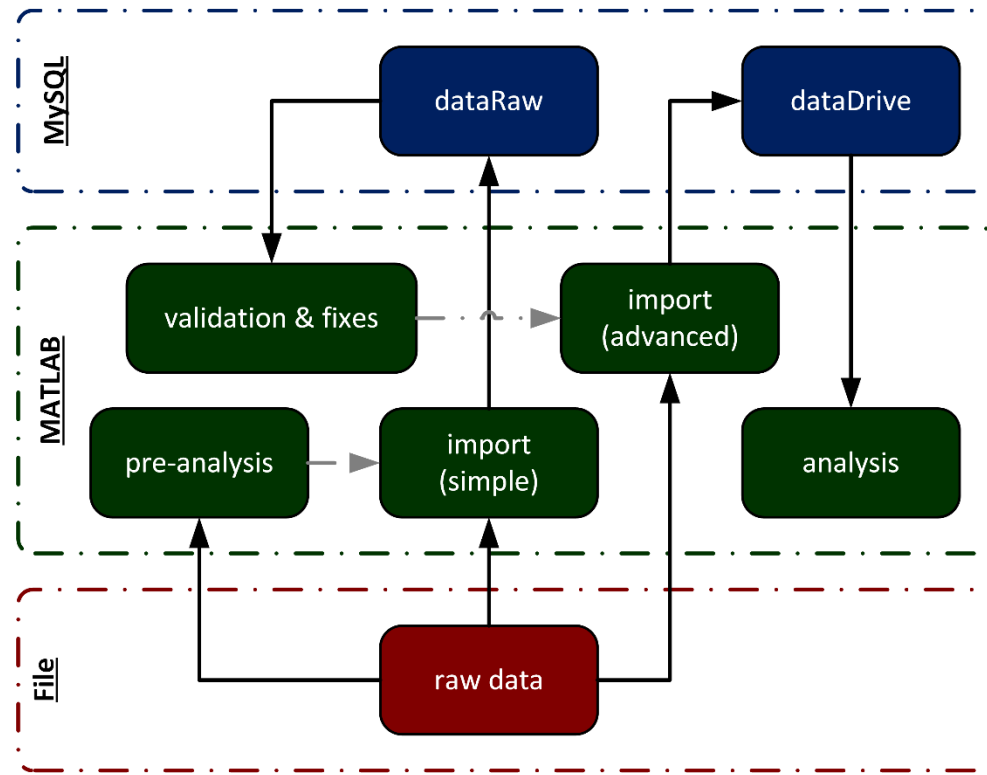
Rules of data processing

- Assume the data is broken
- Mistakes will happen
- Re-check your data / existing processes regularly



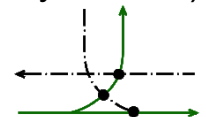


Workflow (an example)



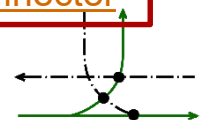
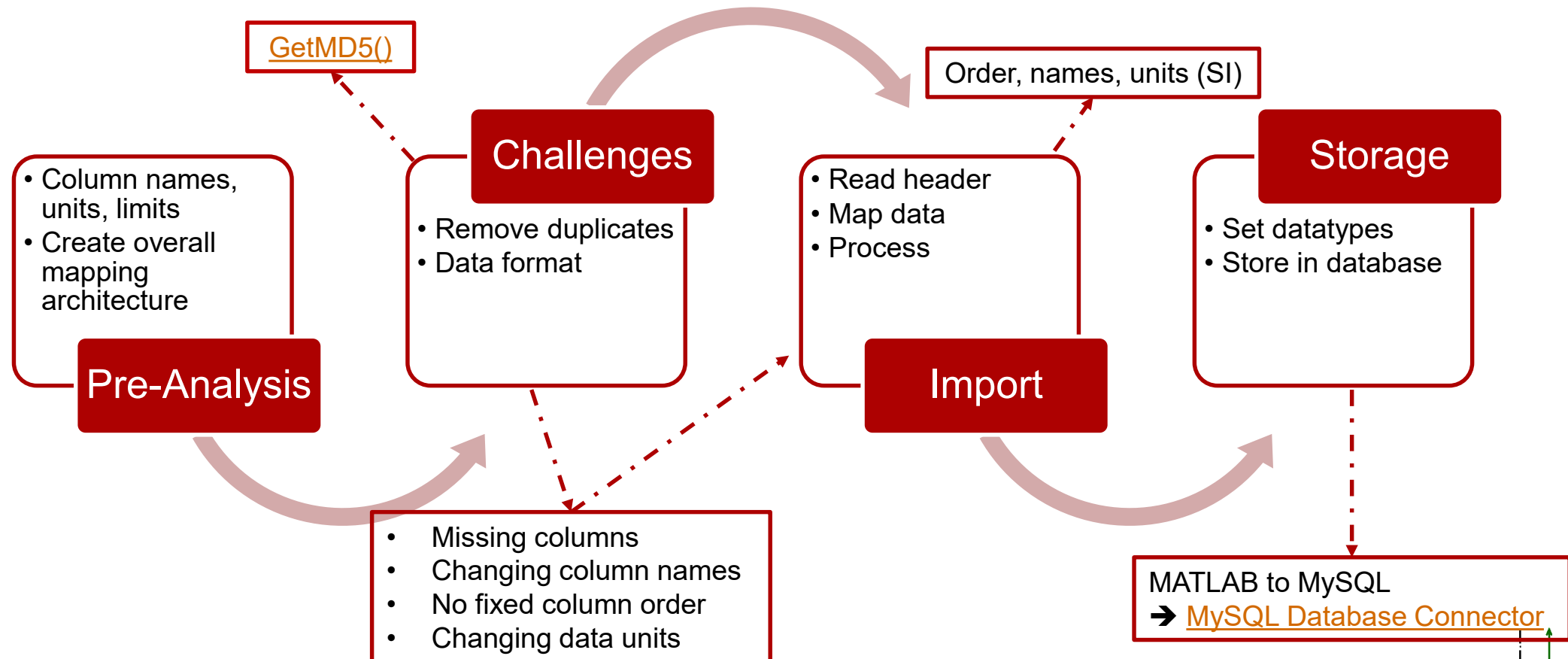
Why Matlab

- University (Everything is DIY)
 - You want it, you build it
 - You built it, you run / manage it
- Maintenance / teaching experience
 - Easy to learn / to understand (for students)
 - Avoid multiple languages
 - Complexity
 - Knowledge required
- Development
 - Easy to visualize data along the way
 - Ready to use toolboxes (or create your own)



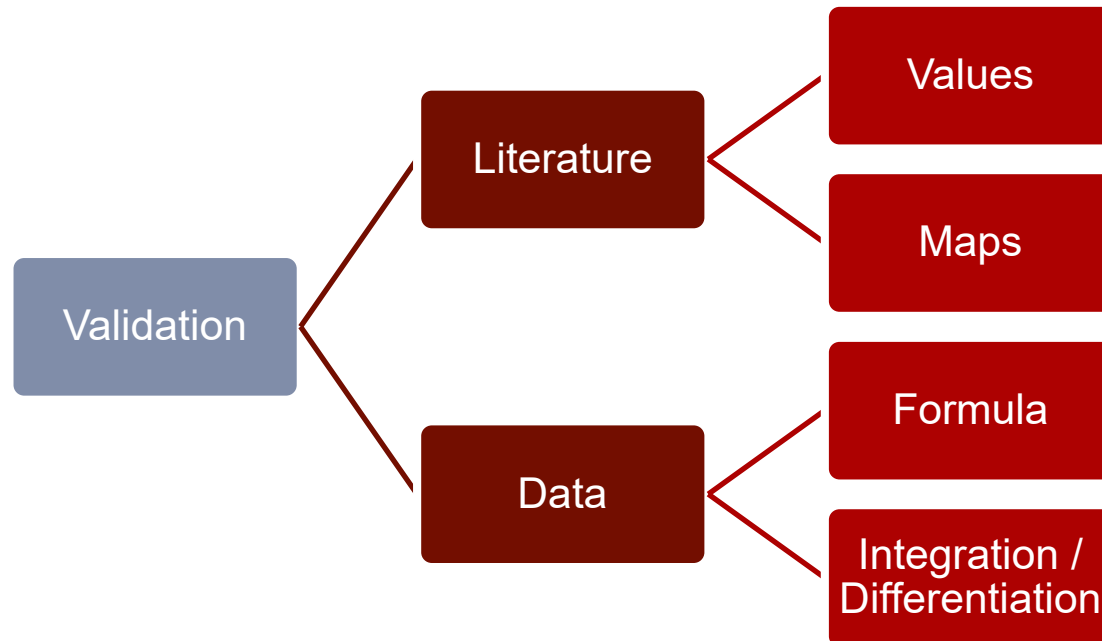


Import (a simple) (CSV)





Validation & Fixes



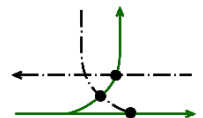
Literature

- Not always correct either
- Contradicting sources

Data

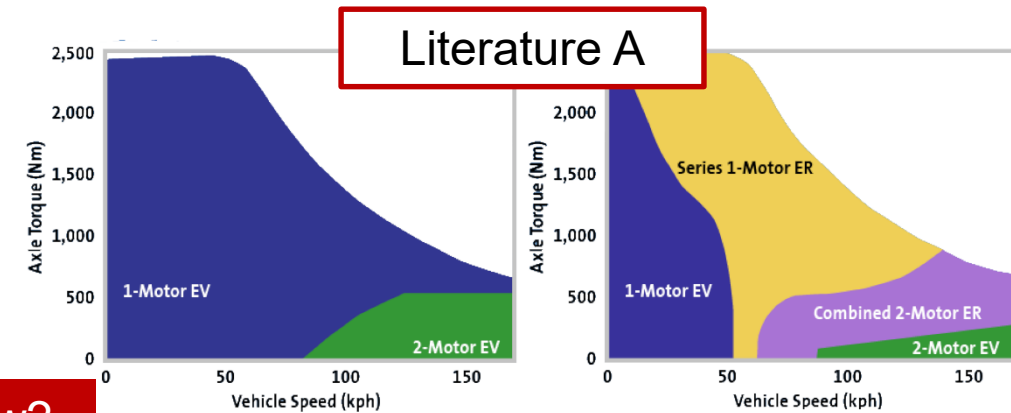
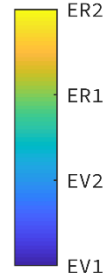
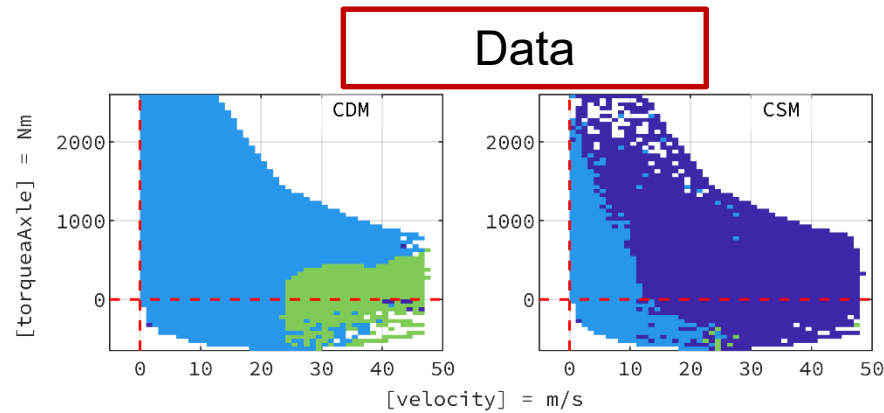
- Requires Formulas
 - Documentation?
 - Generally known?
- Signal can be estimated with “basic” math

**Knowledge is
always required!**

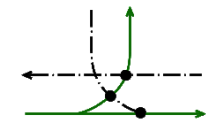
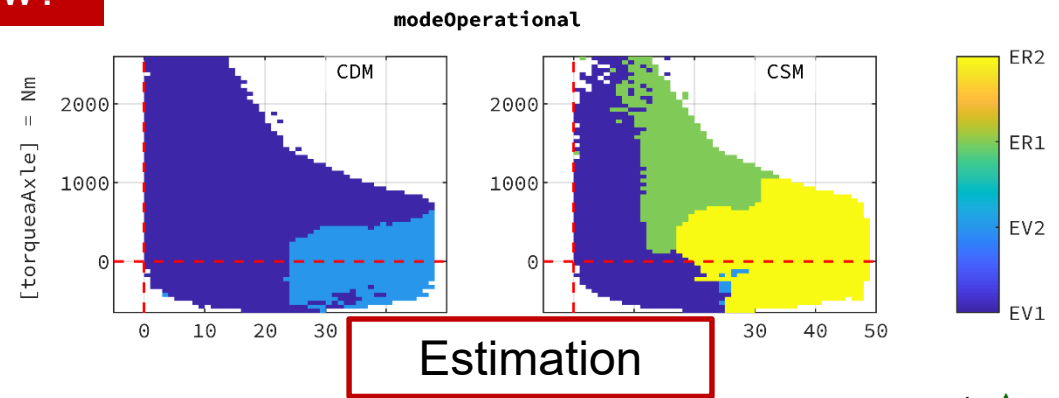
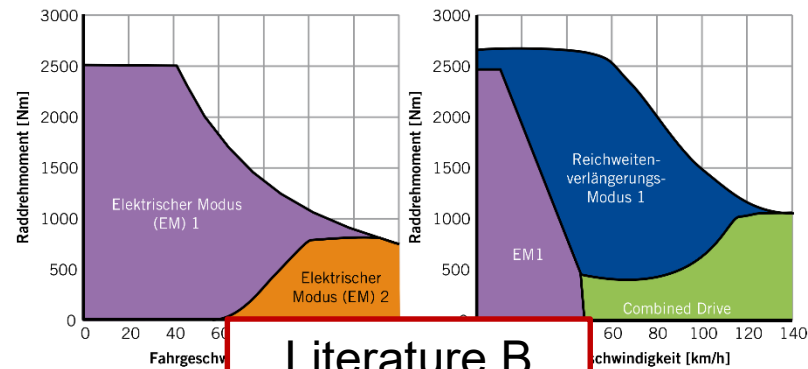




Literature



And now?





Literature

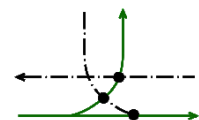
Relative Positive Acceleration (RPA)

- Descriptor for the dynamics of a (driving) cycle

#	Formula	Literature
(1)	$RPA = \frac{1}{x} \sum a_i^+ * v_i$	(Bratt und Ericsson, 2000, S. 5, The European Commission, 2016, S. 12)
(2)	$RPA = \frac{1}{x} \int v * a^+$	(Ericsson, 2000, S. 11)
(3)	$RPA = \prod_{k=1}^{end} v(k) * a(k) / mean(v)$	(Blanco-Rodriguez, Vagnoni und Holderbaum, 2016, S. 653)

- Assumptions make the difference
- Documentation

a : acceleration
v : velocity
x : distance

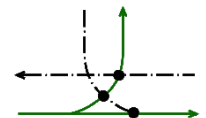
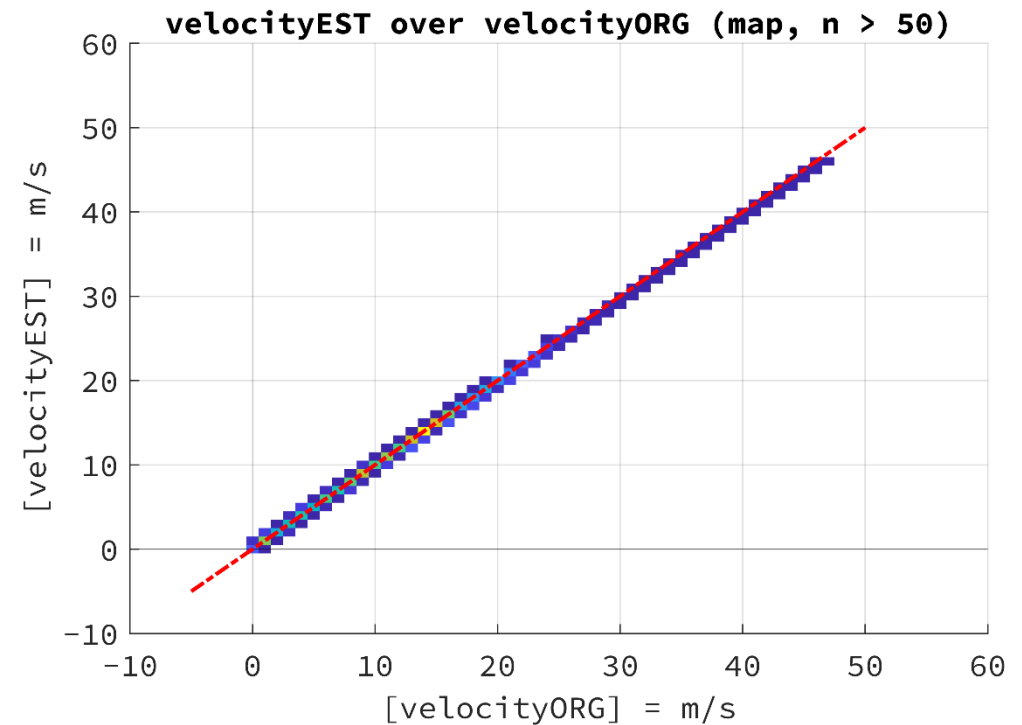
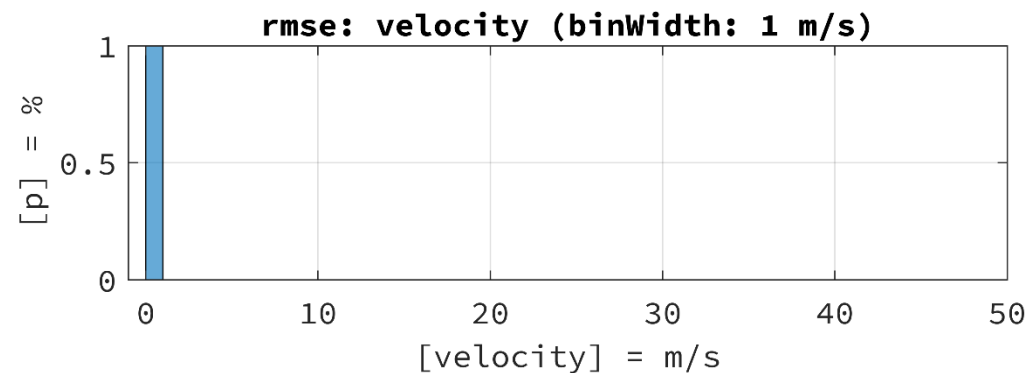




Data

Velocity

- Formula
 - $velocity_{EST} = \frac{2 * \pi * r_{wheel} * speed_{EMB}}{ratio_A * ratio_B}$
- (expected) Value range
 - 0 : 170 km/h
- Conditions
 - Validation only possible in a specific electric mode

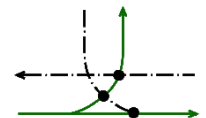
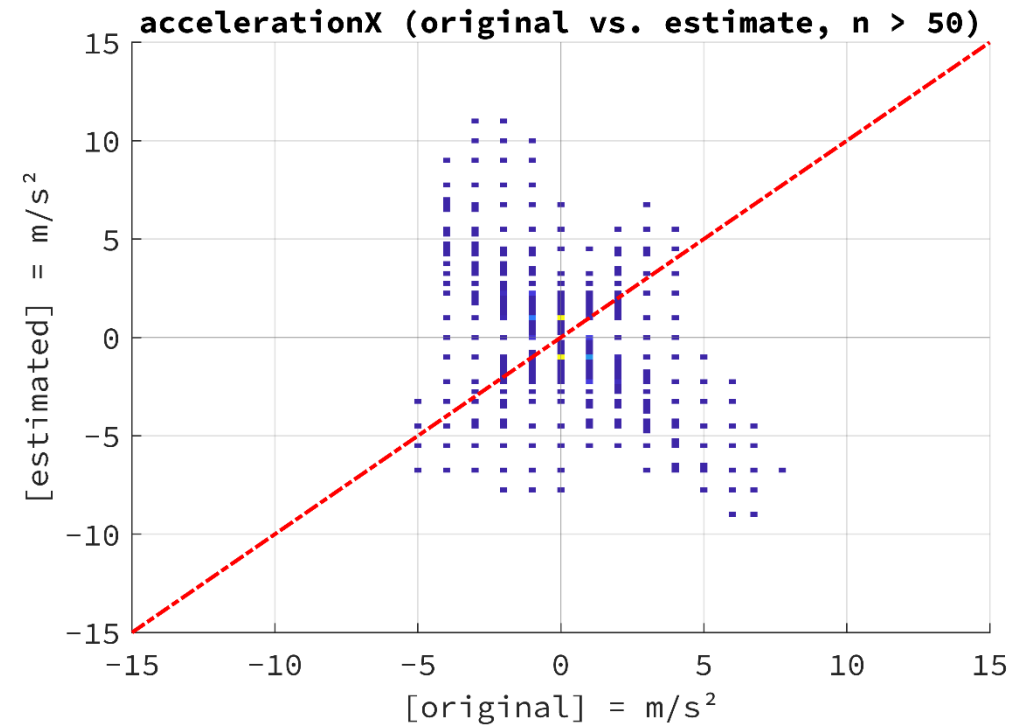
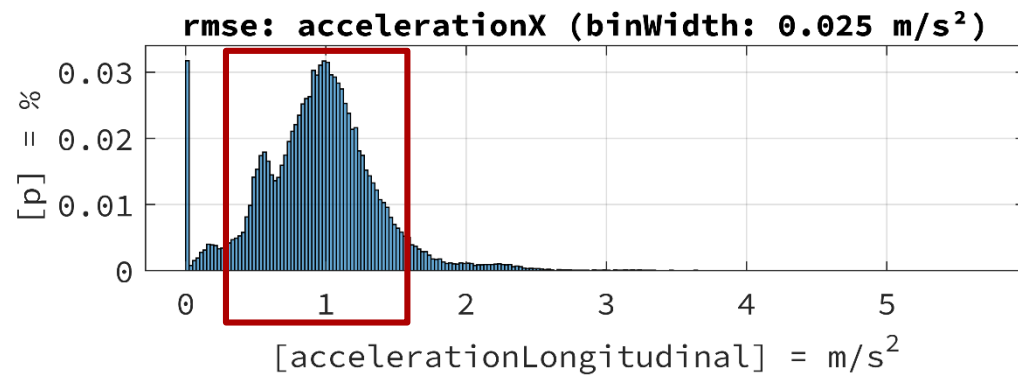




Data

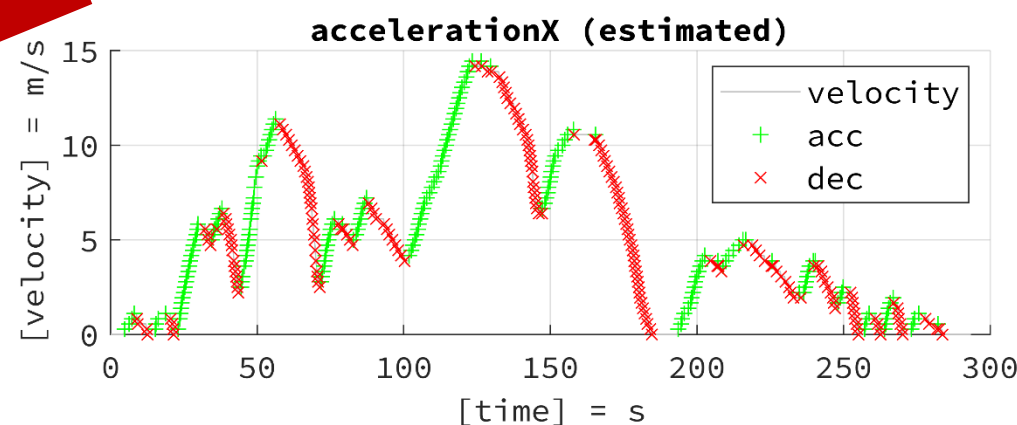
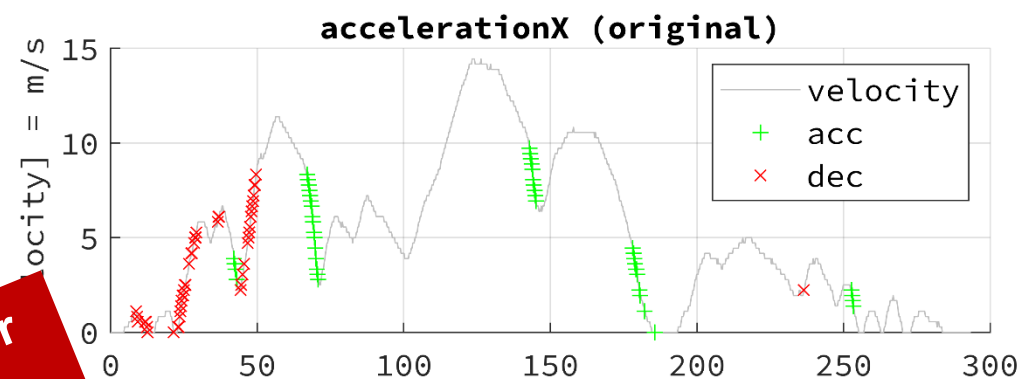
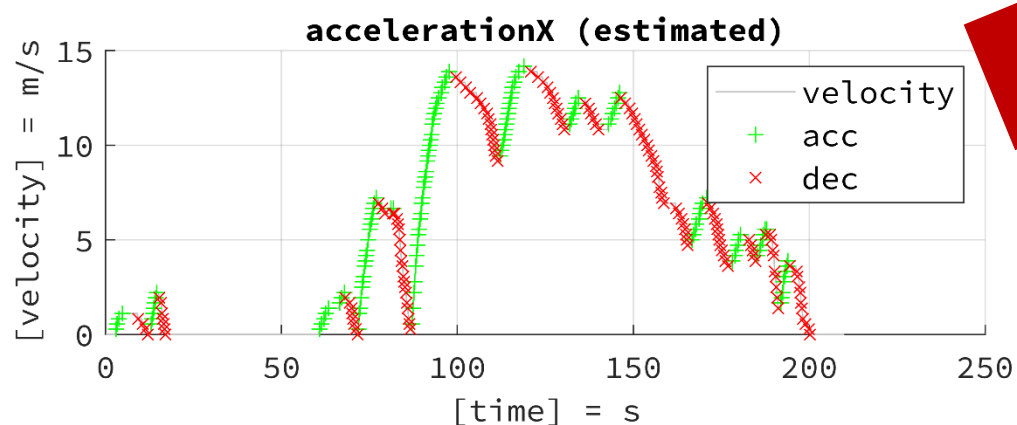
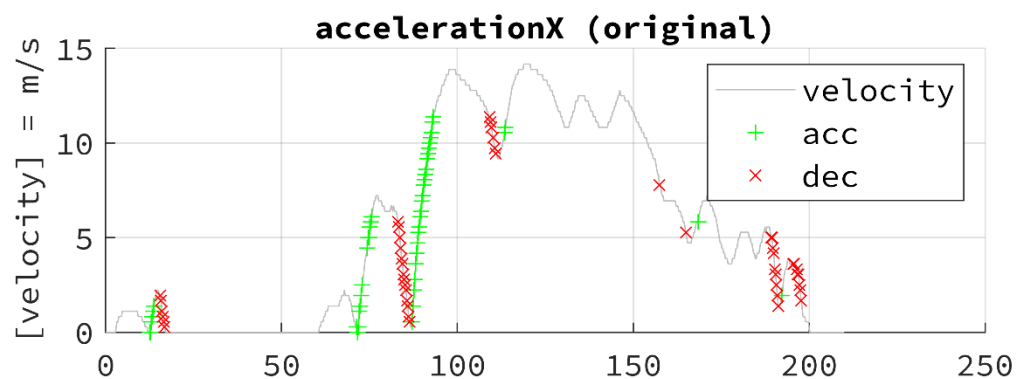
Acceleration

- Formula
 - $acceleration_{EST} = \frac{\partial velocity}{\partial time}$
- (expected) Value range
 - $-10 \text{ m/s}^2 : +10 \text{ m/s}^2$

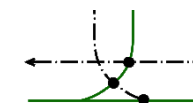




Data



Affects other values

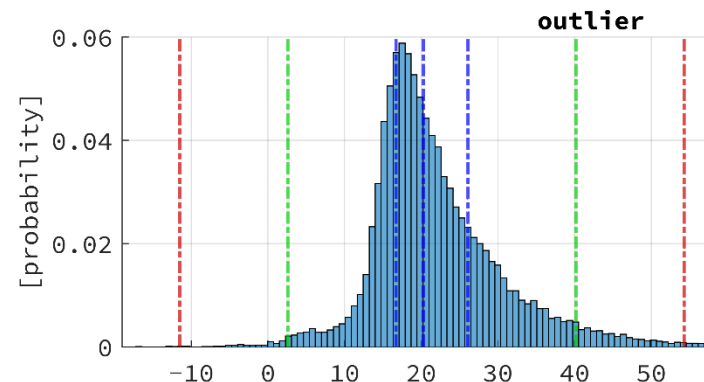




(fixed) Data

Steps (so far)

- Look at all data
 - RMSE values for identification
 - Original vs. Estimate histograms
 - Literature (if available)
- Look at individual data
 - Plot individual trips
- Correct / Fix
 - Your code
 - Assumptions



Quantile

Q_{50} : 20.3

Q_{25} : 16.7

Q_{75} : 26.1

D_{25} : 3.6

D_{75} : 5.8

Fence (Mild)

low: 2.6

high: 40.2

p: 4.2%

Fence (Extreme)

low: -11.5

high: 54.3

p: 1.0%

Info

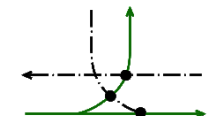
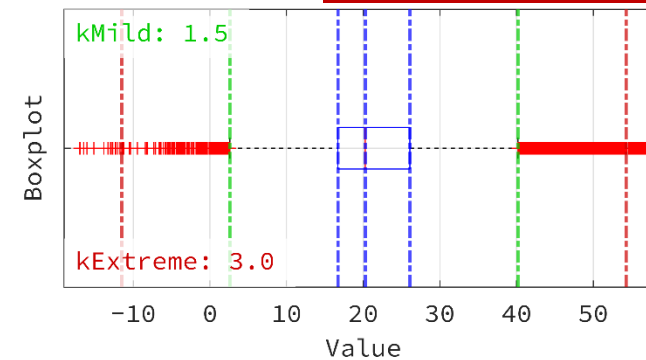
n_{DP} : 42960

x_{Min} : -174.8

x_{Max} : 149.0

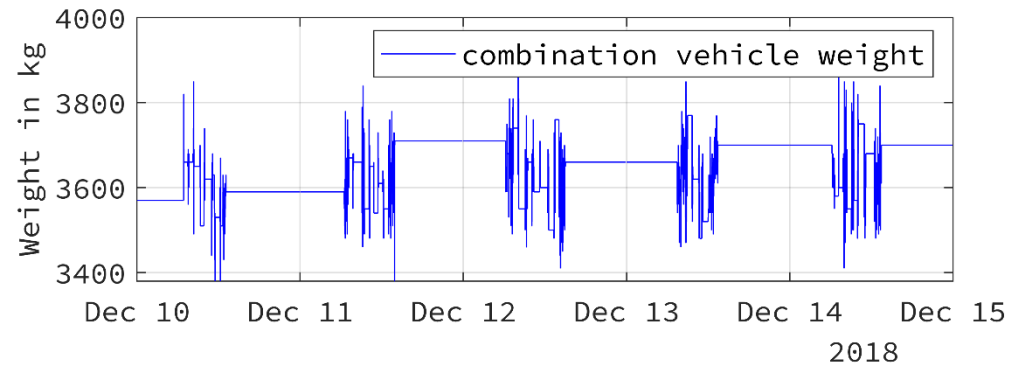
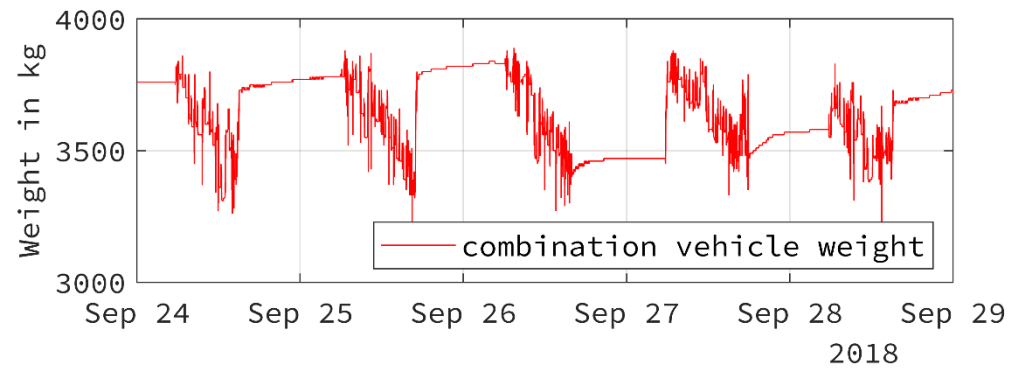
x-axis limited
to $Q50 \pm 4.0\sigma$

Handling outliers





(monitor) Data

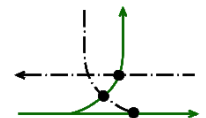


Bugs

- Sensor stopped working correctly for unknown reasons
- Bug reported in January (still not fixed)

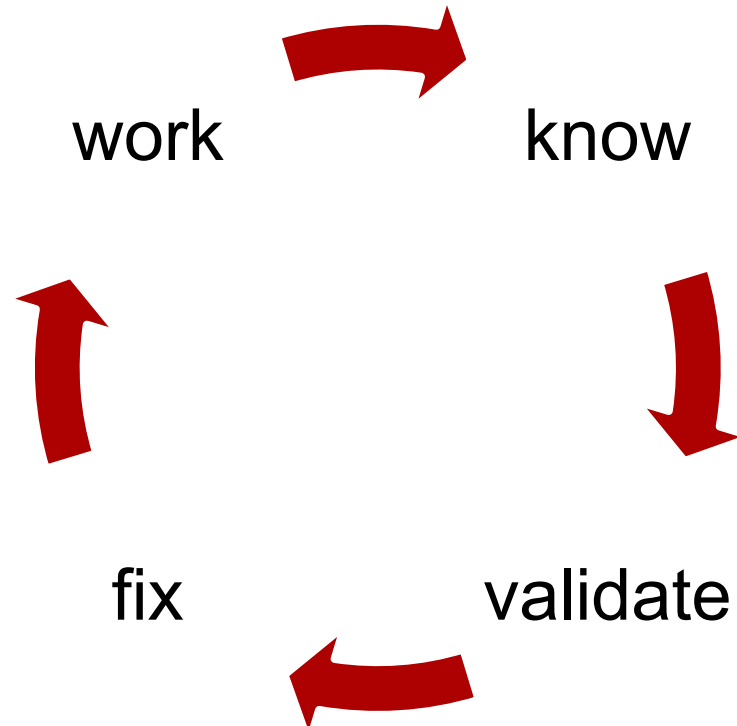
Steps

- ?

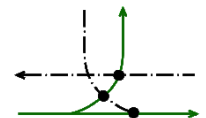




Lessons learned



- There will be new issues
 - Re-check your (raw-)data & assumptions
- Make sure you are not the problem
 - Test your code
 - Review your code
- Make proper bug reports
 - Provide examples (figures if possible!)
 - Don't expect the bugs to get fixed
- Documentation saves time
 - Use it whenever possible
 - Read it properly
 - Create your own (and keep it up to date)



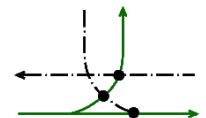


Lessons learned

- Nobody is perfect, neither is software
- Nothing is ever really finished
- Expect things to change / go wrong
- Look at your data and understand what you are looking at
- Know what your data is supposed to look like
- Never stop learning
- Decent hardware is a big help

Unexpected MATLAB lessons (over the years)

- `read(datastore)` doesn't always read the whole file
- `webread()` may not always return a variable (or produce an error)
- use `wget()` instead of `mget()`
- Different behaviour of the same function in different releases
- MATLAB and Linux work best from a terminal in software mode (best guess: buggy gfx driver)
- Never stop upgrading to the next release





End

