

Time-series

and how to cook them with ML

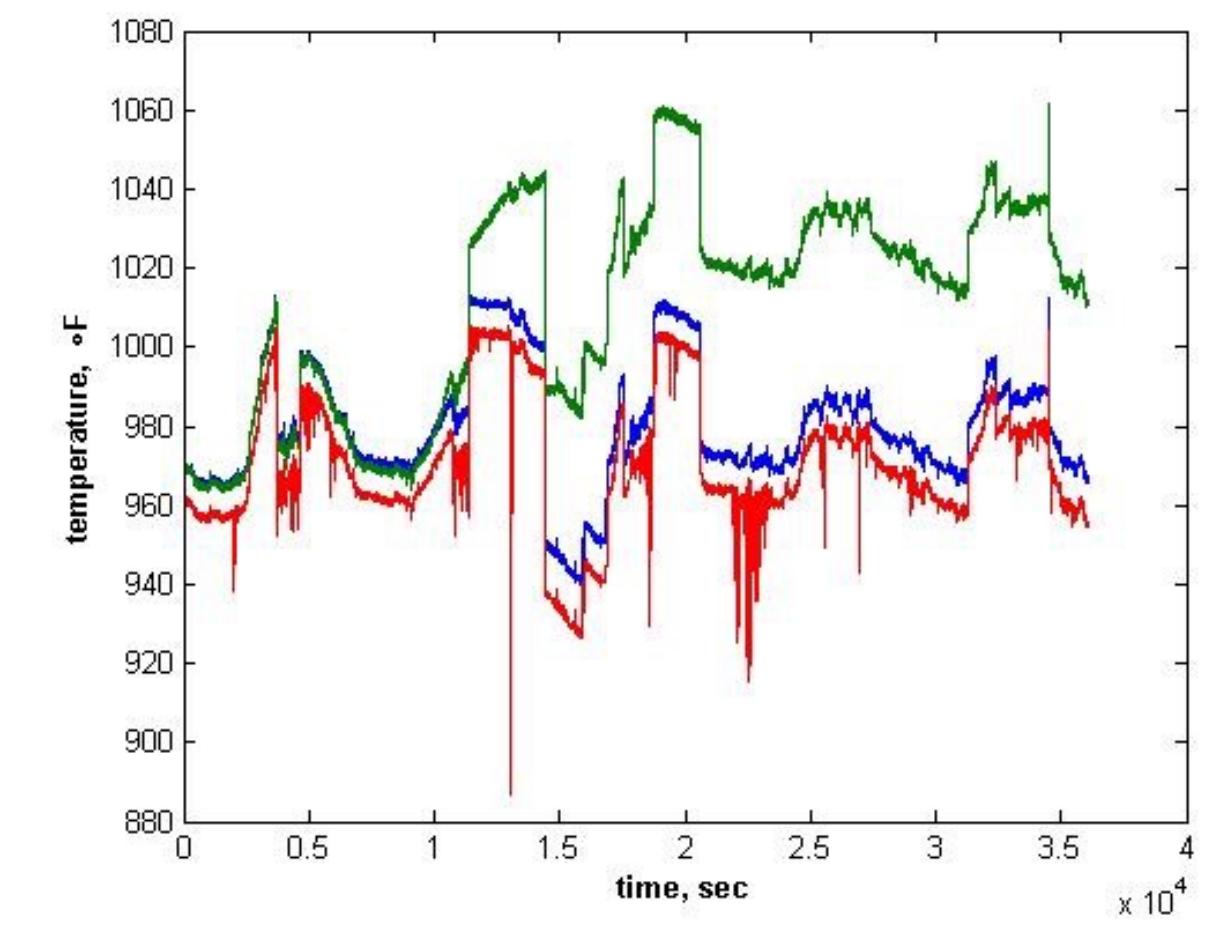
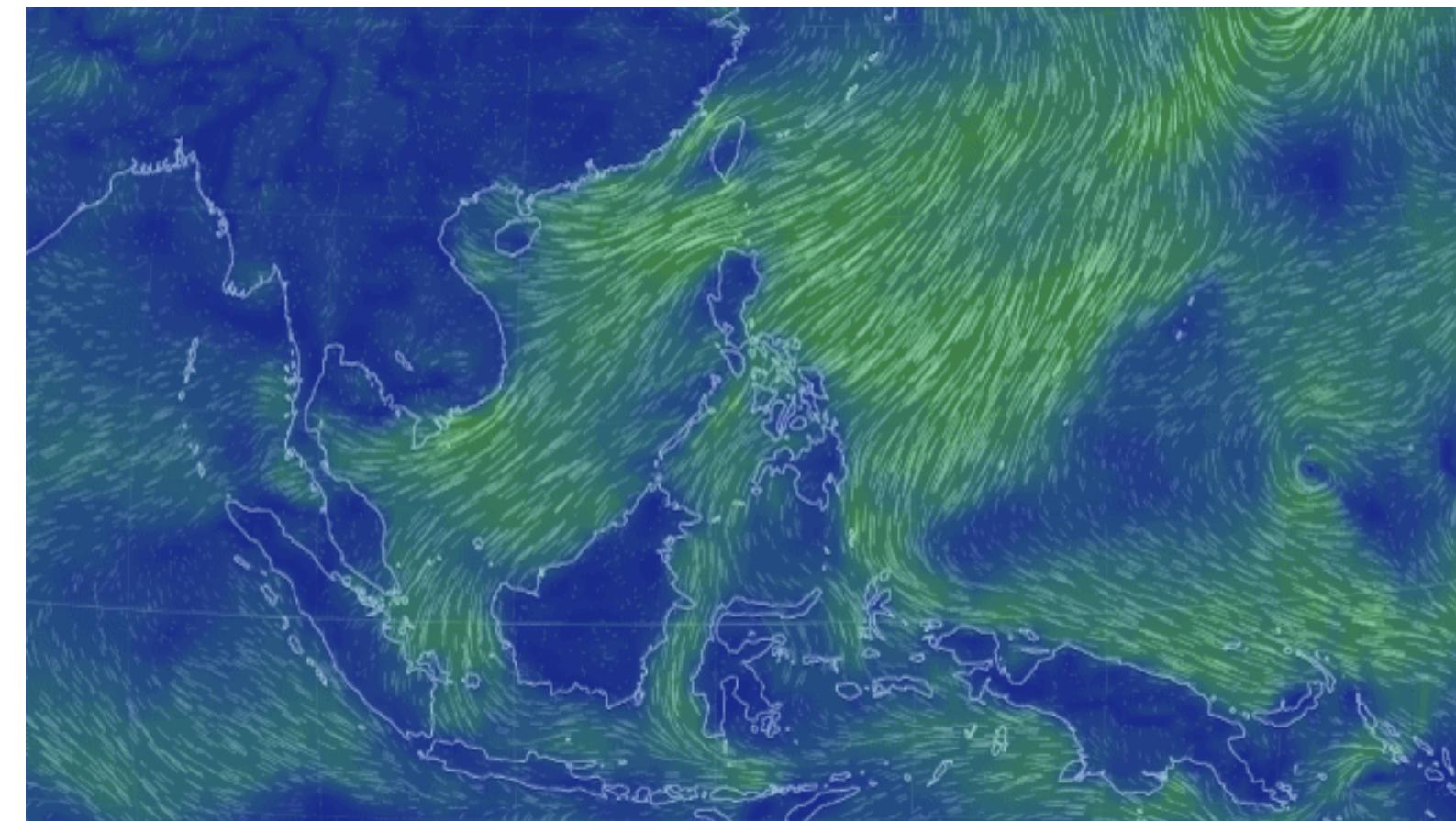
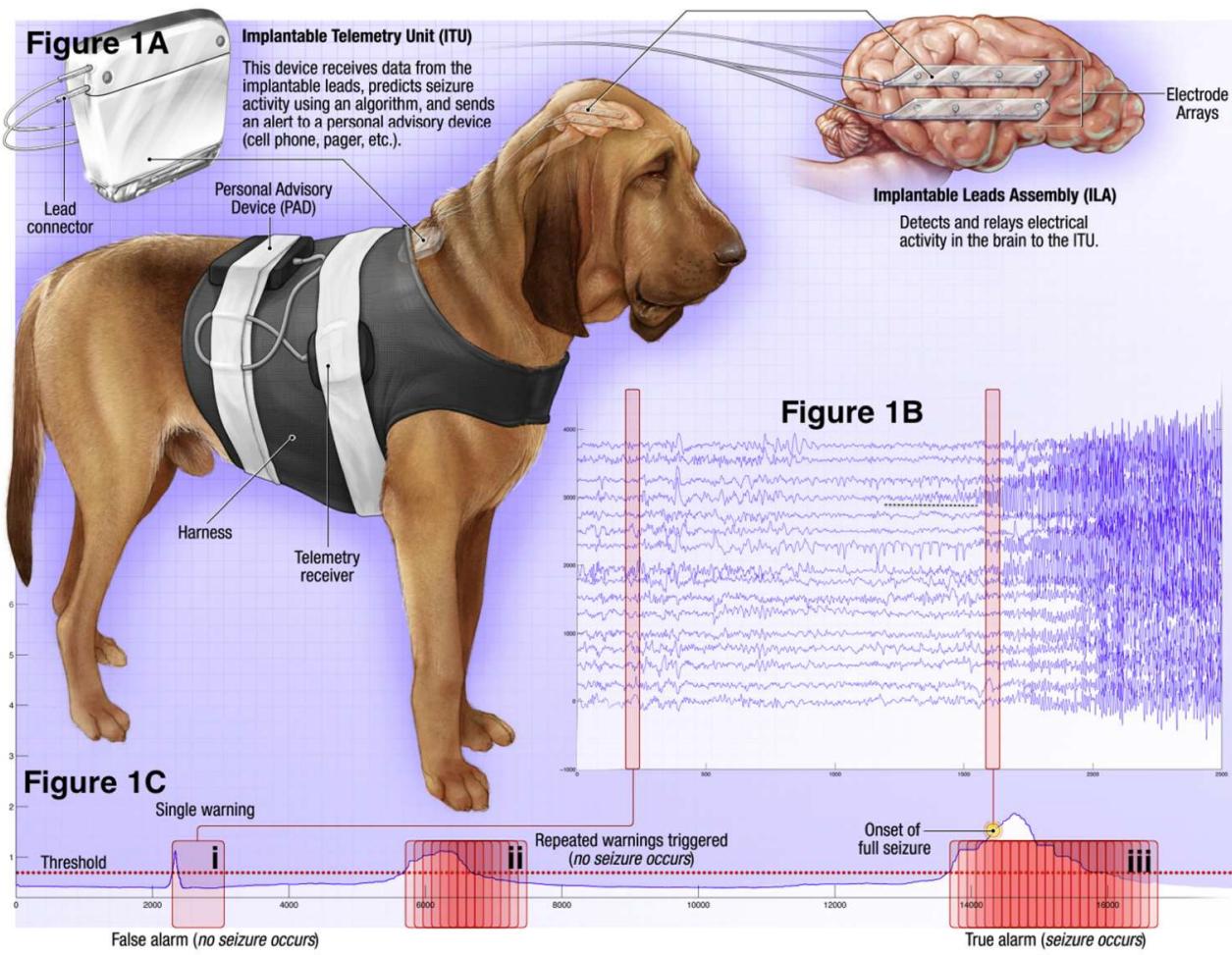
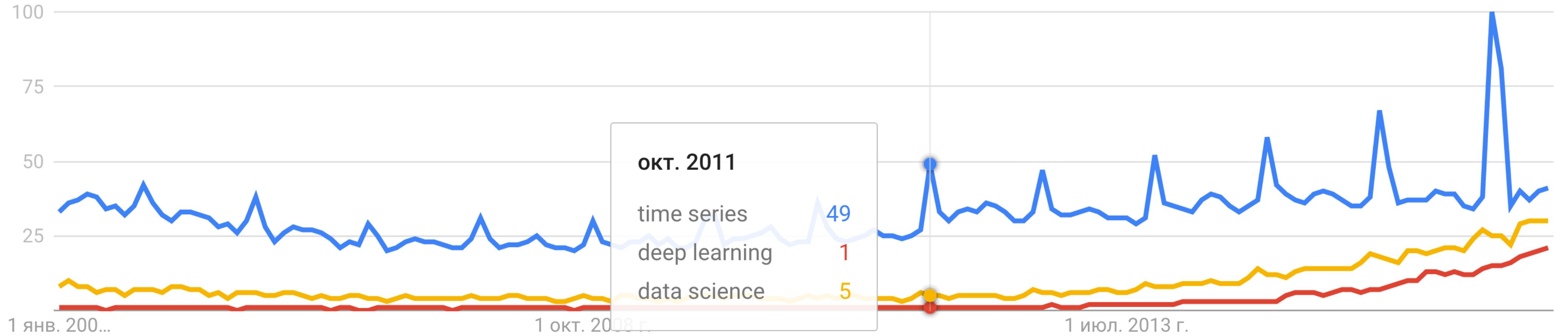
Alex Natekin

Time-series

- ▶ “Series of data points indexed in time order”
- ▶ “Collection of observations of well-defined data items obtained through repeated measurements over time”

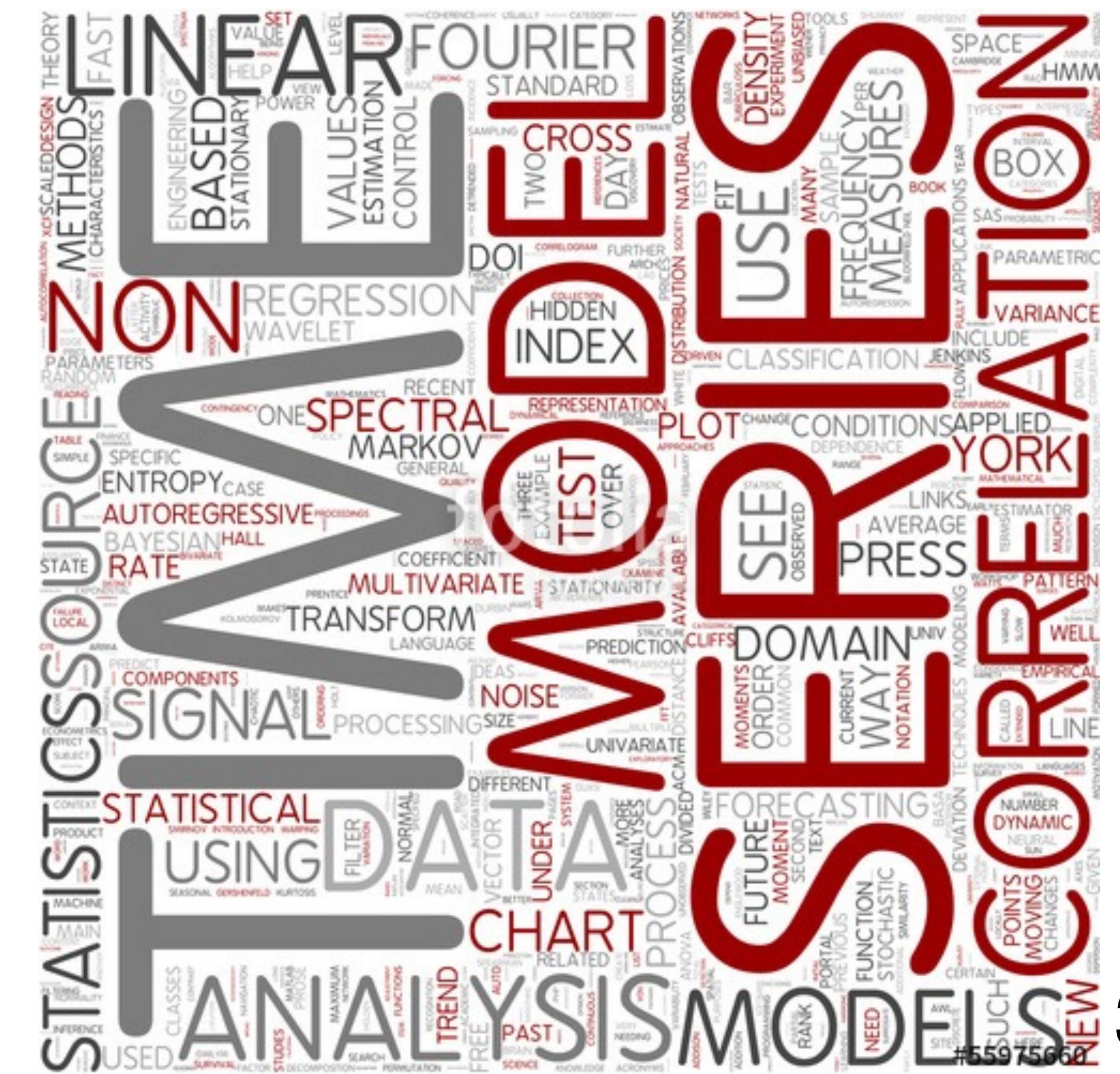
Date	UserId	Sum	Consumed (ml)	Places
2016-12-01	10780	800.00	68	“Outside”
2016-12-05	10780	120.00	60	“Inside”
2016-12-09	10780	145.00	62	“Inside”
2016-12-14	10780	1600.00	145	“Both”
2016-12-19	10780	900.00	71	“Outside”
2016-12-24	10780	600.00	65	“Inside”
2016-12-25	10780	666.00	NA	“Inside”
2016-12-26	10780	1200.00	NA	“Outside”
2016-12-27	10780	1400.00	125	“Outside”
2016-12-28	10780	2300.00	NA	“Both”

Time-series examples



Much prior work

- ARIMA, GARCH, others from statistical school
 - Wavelets and other signal processing approaches
 - RNN/LSTM get hype
 - Casual Kaggle-style ML, ...



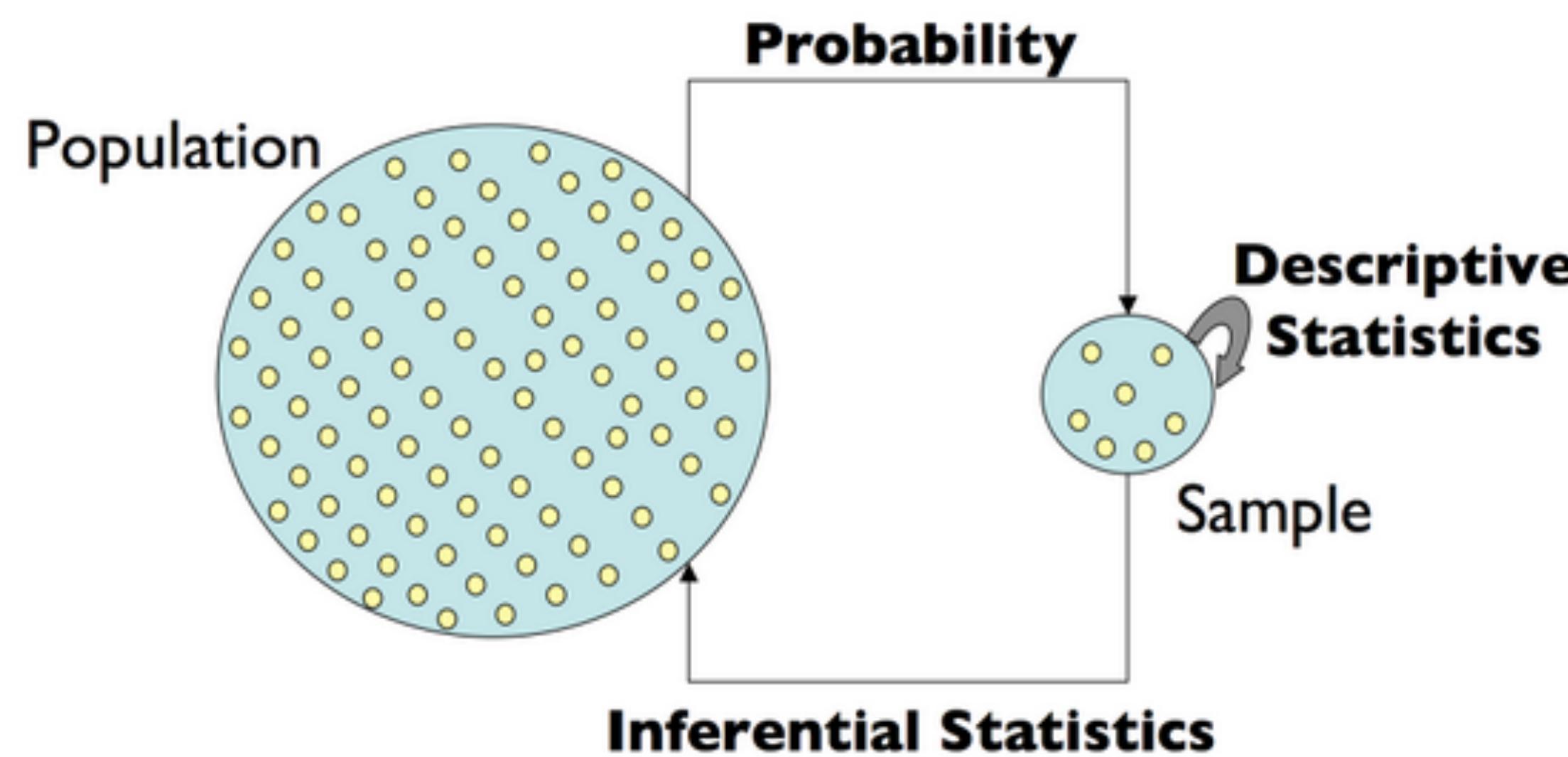
Forecasting

- ▶ Choose what to predict:
value/aggregate/class
- ▶ Guess the dependency:
univariate/naive/all/mixed
- ▶ Proceed with setup of
other ML problem
components

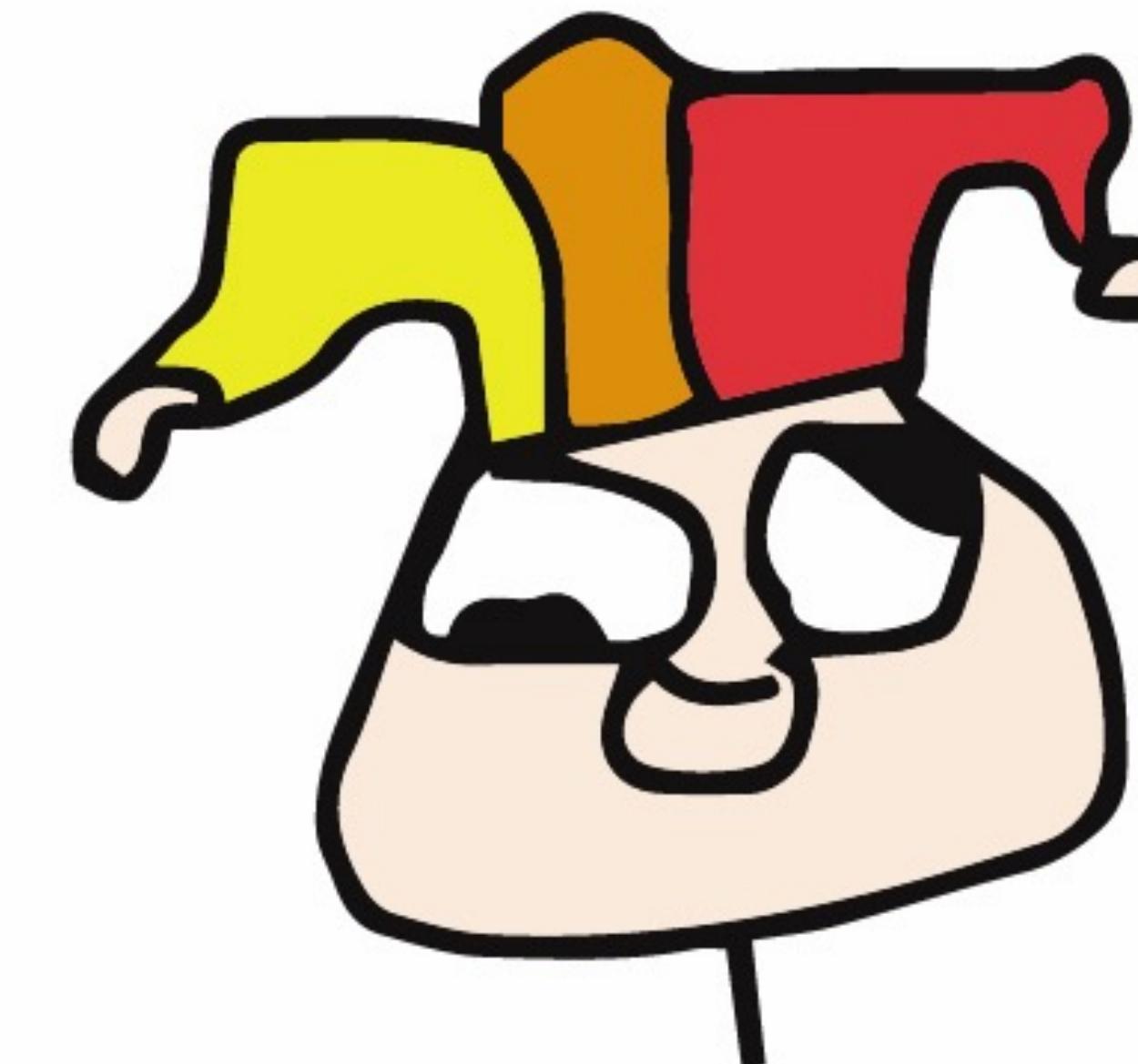
Date	UserId	Sum	Consumed (ml)	Places
2016-12-01	10780	800.00	68	"Outside"
2016-12-05	10780	120.00	60	"Inside"
2016-12-09	10780	145.00	62	"Inside"
2016-12-14	10780	1600.00	145	"Both"
2016-12-19	10780	900.00	71	"Outside"
2016-12-24	10780	600.00	65	"Inside"
2016-12-25	10780	666.00	NA	"Inside"
2016-12-26	10780	1200.00	NA	"Outside"
2016-12-27	10780	1400.00	125	"Outside"
2016-12-28	10780	2300.00	NA	"Both"

ML models for time-series

Central dogma of
Statistics:

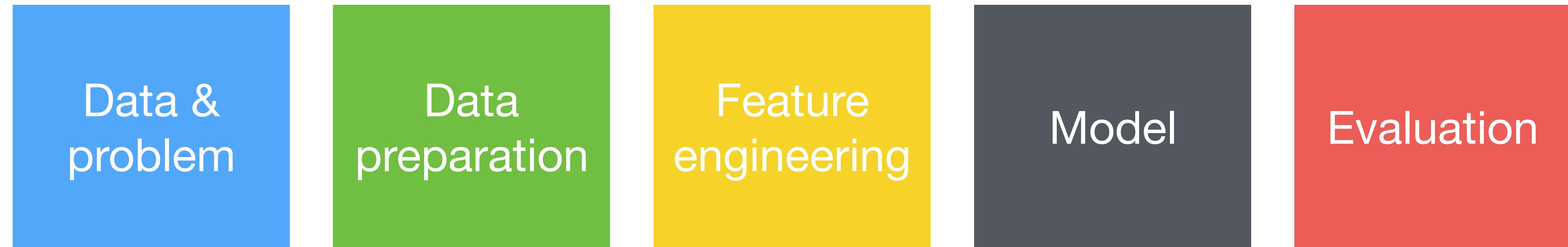


Central dogma of
Machine Learning:



**Everything
is a feature**

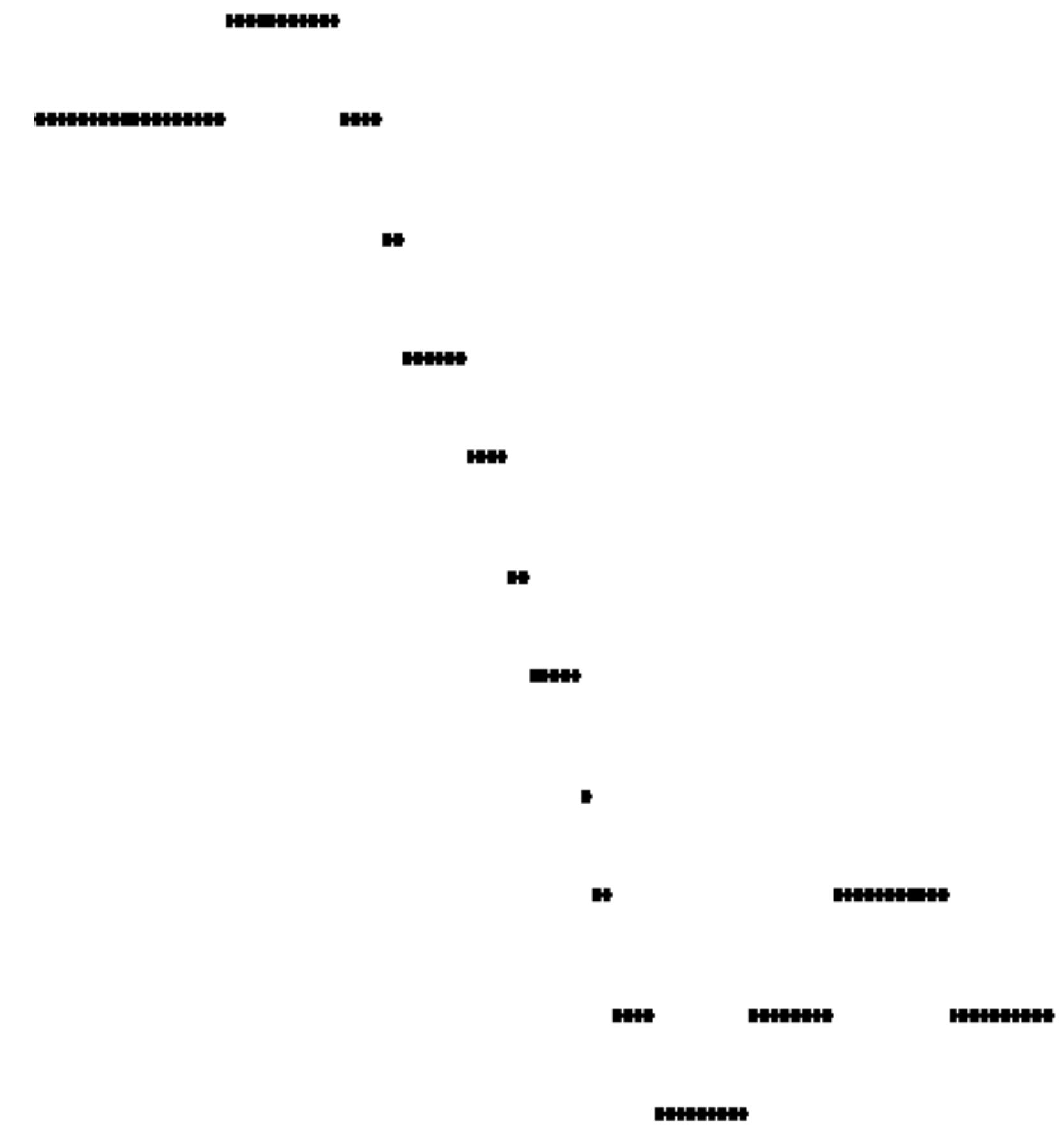
Time-series ML pipeline



- | | | | | |
|---|--|---|--|--|
| <ul style="list-style-type: none">‣ Value vs diff‣ Regression/ classification /anomalies‣ Univariate sensors vs panel data | <ul style="list-style-type: none">‣ Consistent spacing‣ Impute missing values‣ Aggregates/ pooling | <ul style="list-style-type: none">‣ Stationarised series‣ AR, MA and aggregates‣ Dense represent-ns‣ Kaggle-style | <ul style="list-style-type: none">‣ Linear‣ Casual ML & boosting‣ LSTM & RNN‣ Markov-style‣ Kaggle-style | <ul style="list-style-type: none">‣ Time-CV‣ Test-policy (rolling or not)‣ Simulations (small data) |
|---|--|---|--|--|

TS data preparation

- ▶ Check spacing. If non-uniform, choose time grid (and fill missing)
- ▶ Missing values: zero vs LOCF
(last observation carried forward)
vs interpolation
- ▶ Downsample: use minute-level aggregates instead of seconds ...

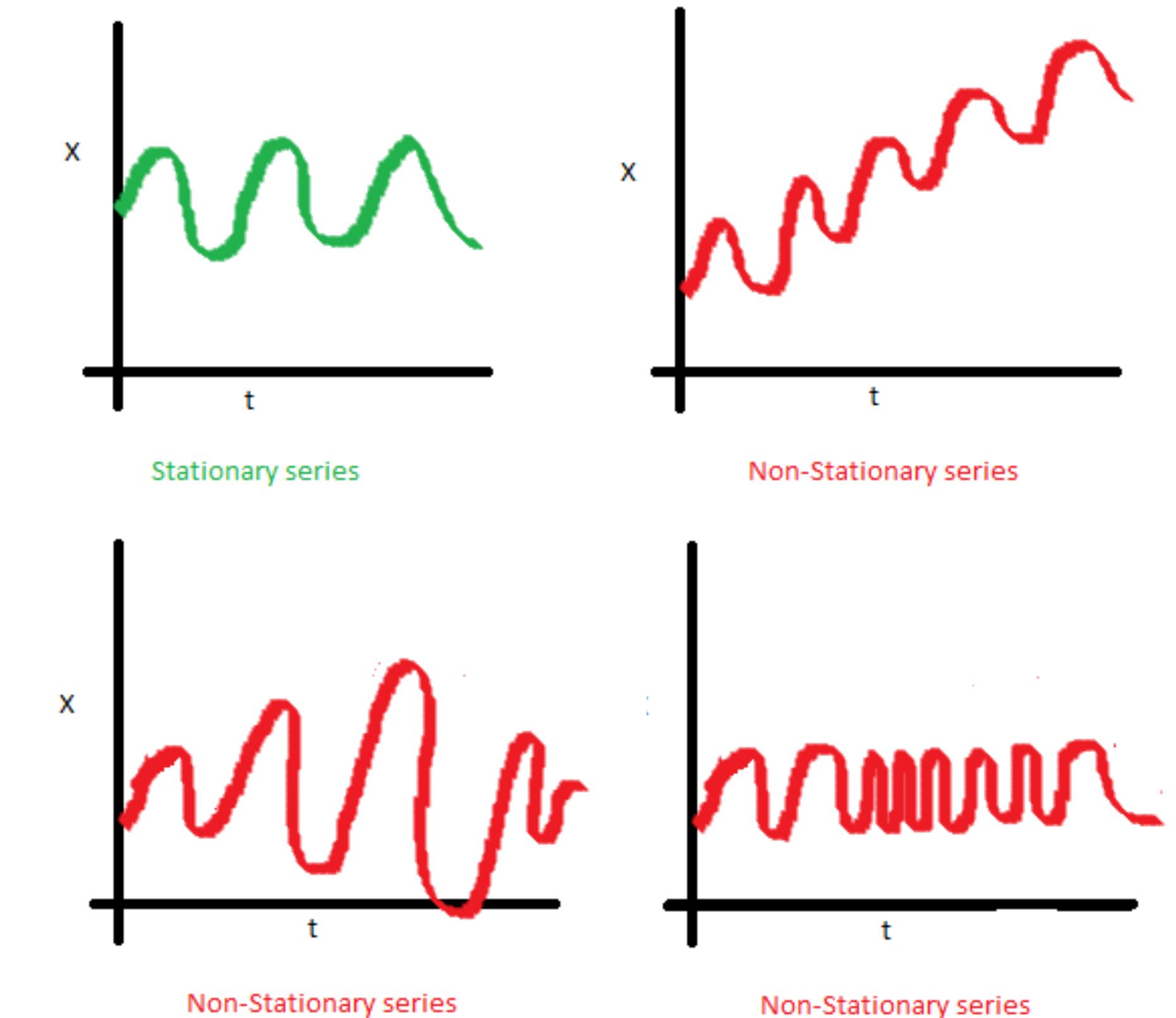


TS feature engineering (1/5)

- Stationarization: detrending & differencing (can be repeated)

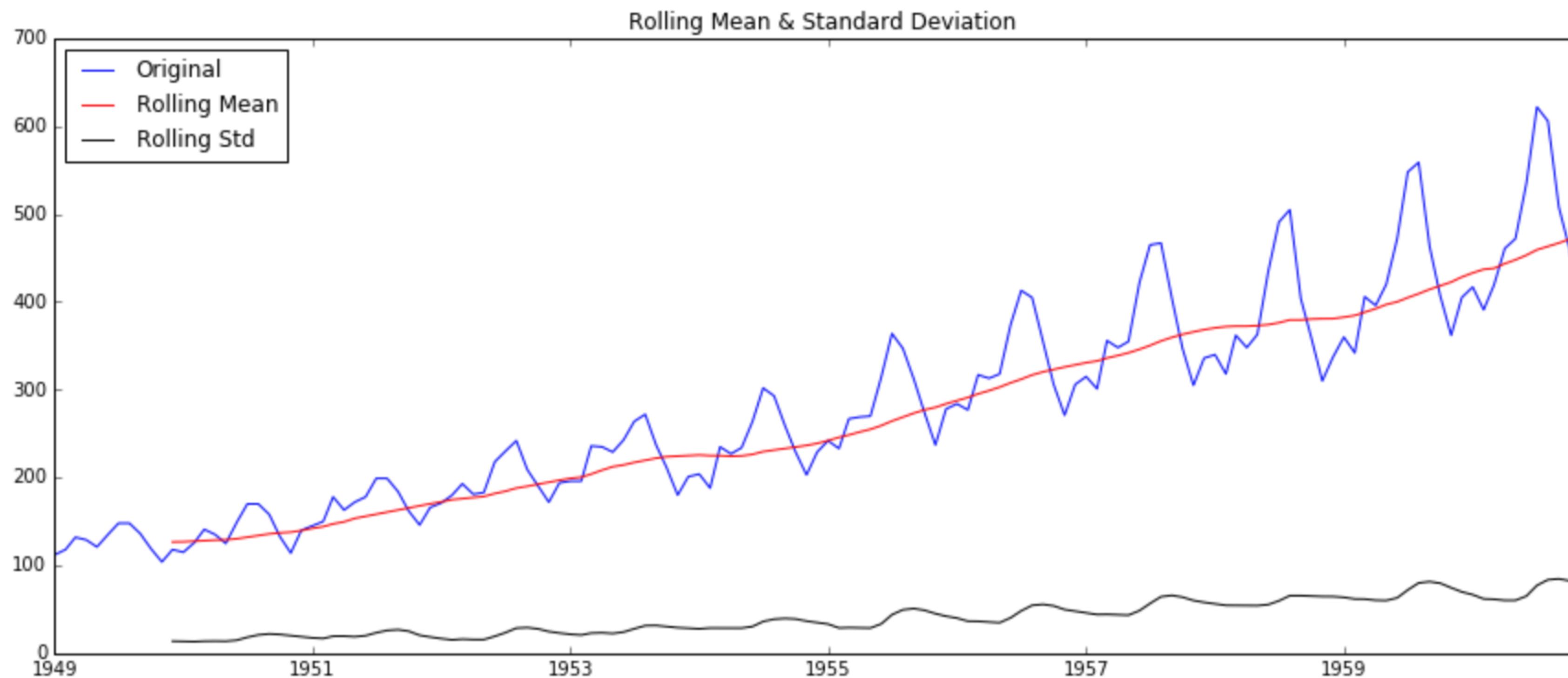
- Log-transform

- **NB:** you can keep both original and diff ones for models



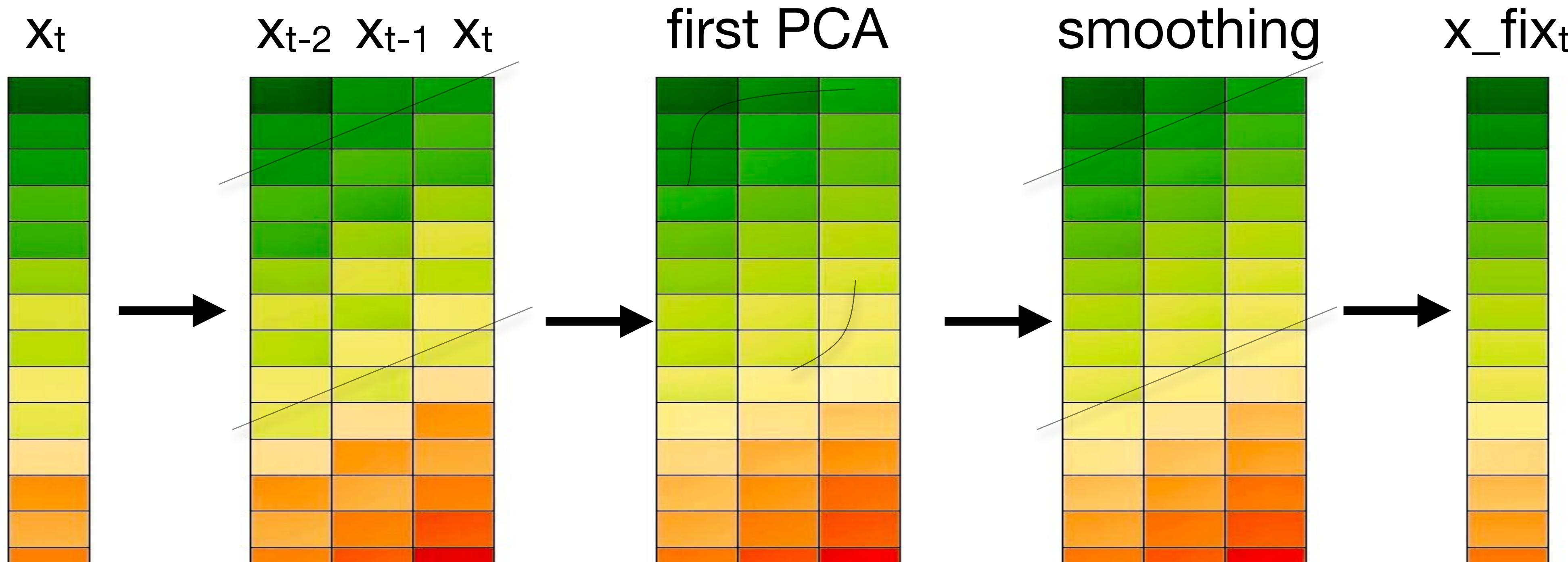
TS feature engineering (2/5)

- ▶ Rolling statistics (mean, sd, min, max, ...)
- ▶ Rolling aggregates (counts, incl. other features)



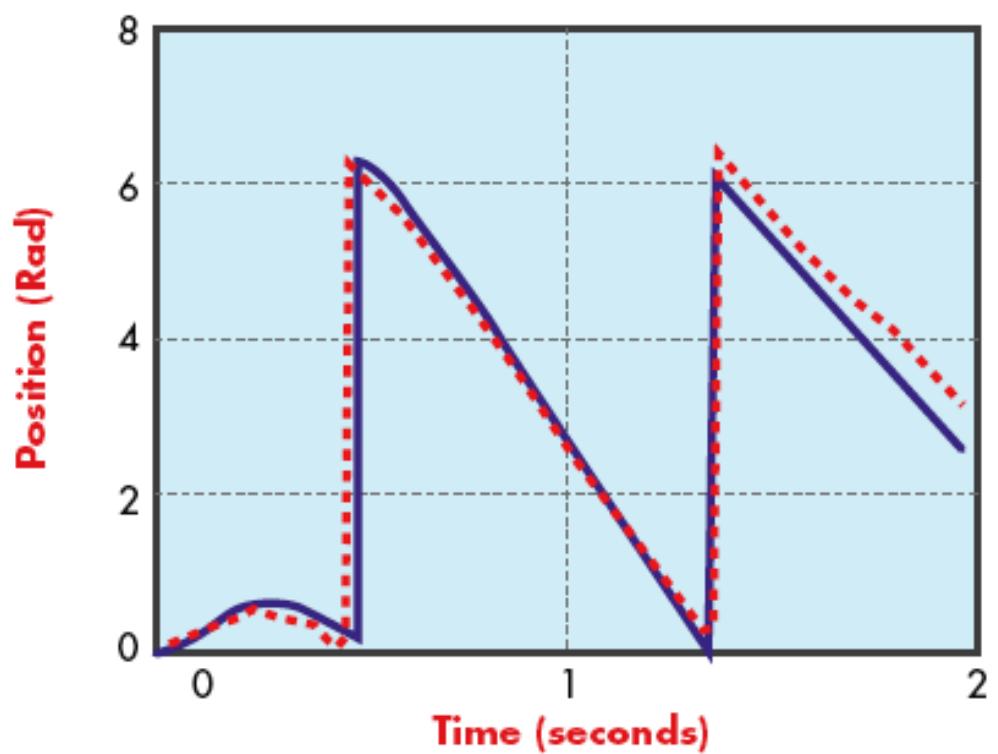
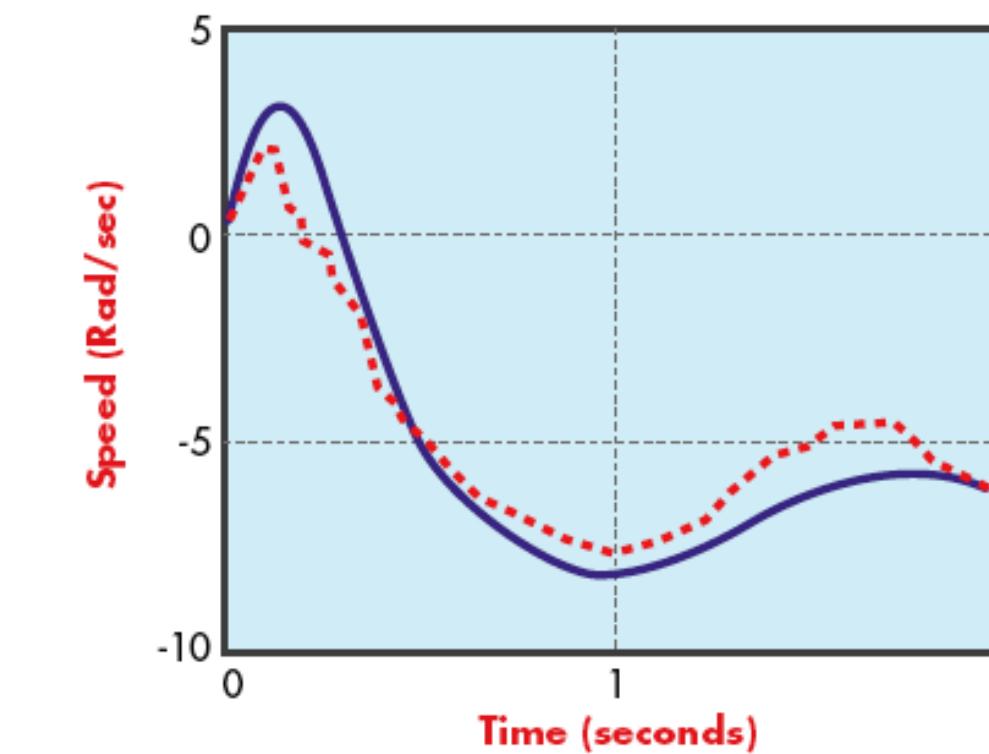
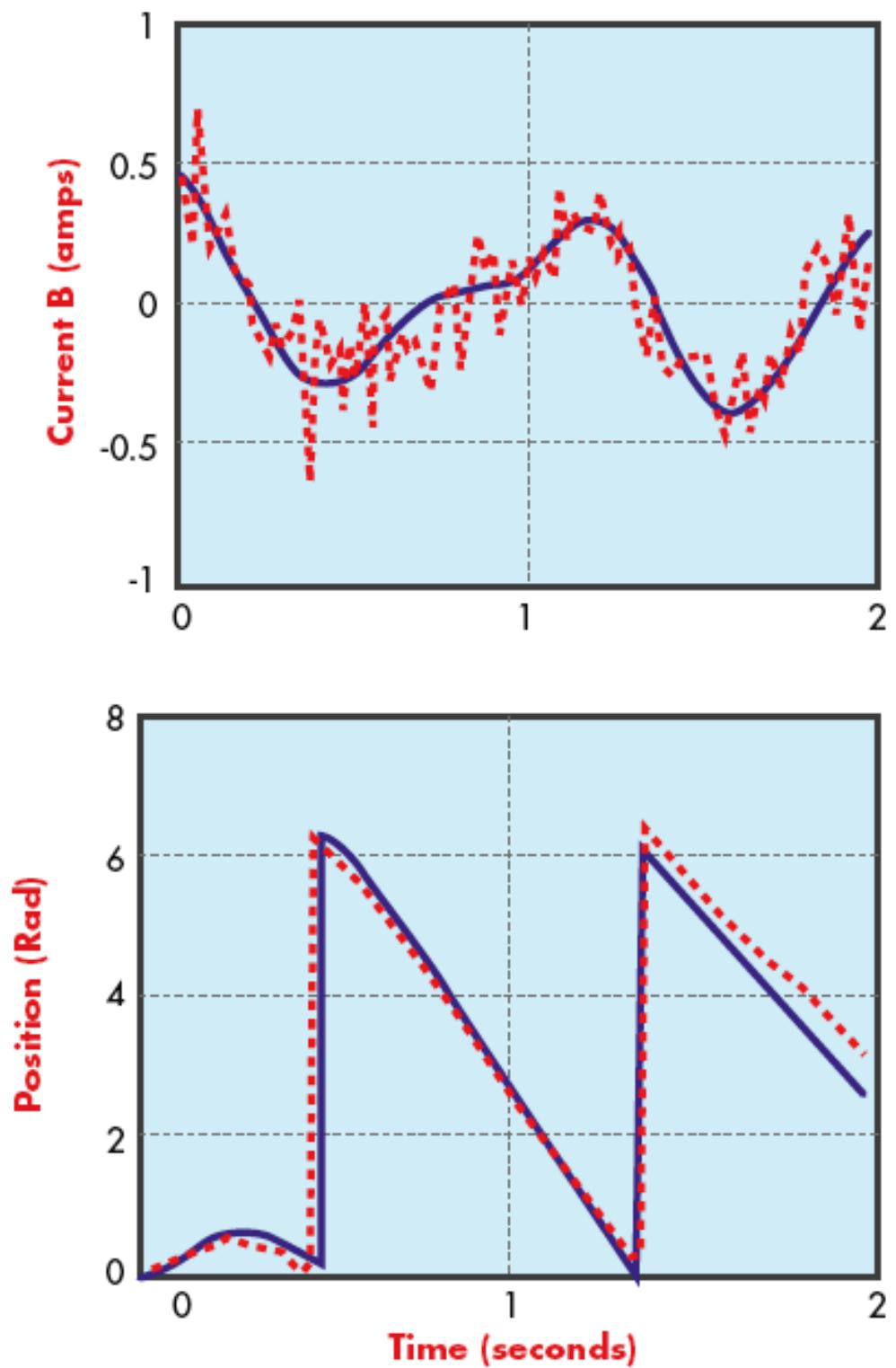
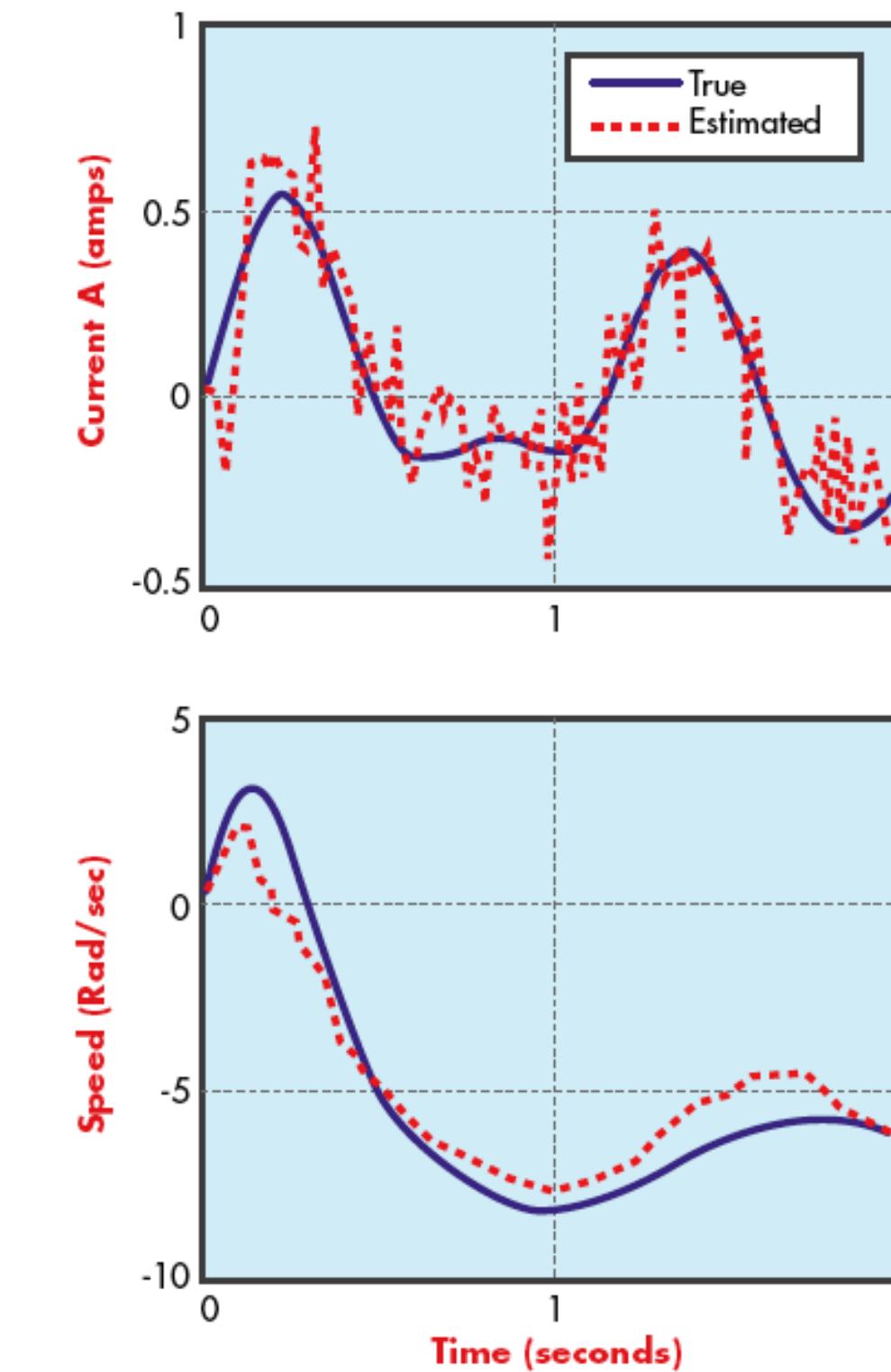
TS feature engineering (3/5)

- AR: lagged values for “auto-regression”
- SSA and respective goodies (trends, periodics, ...)



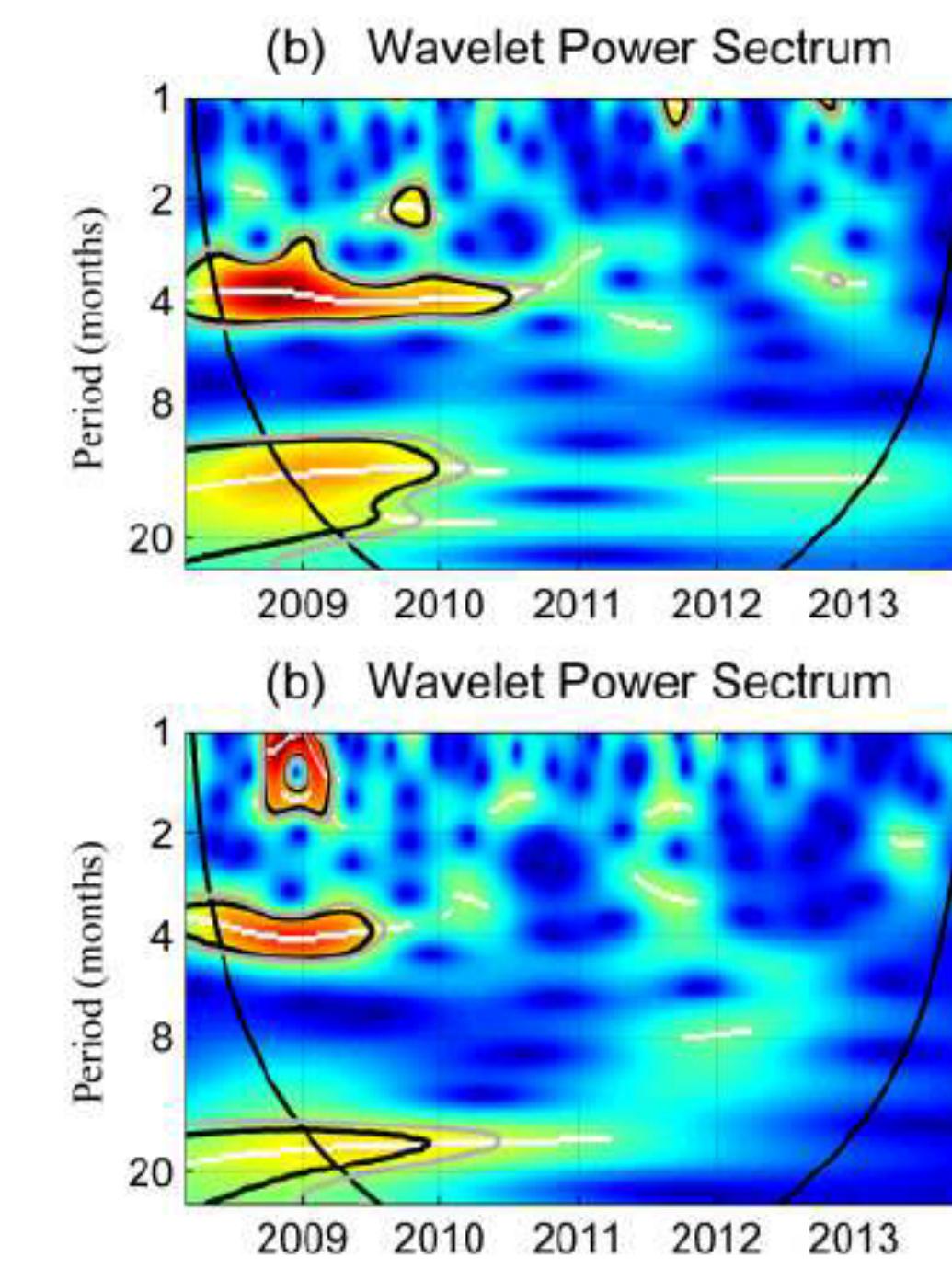
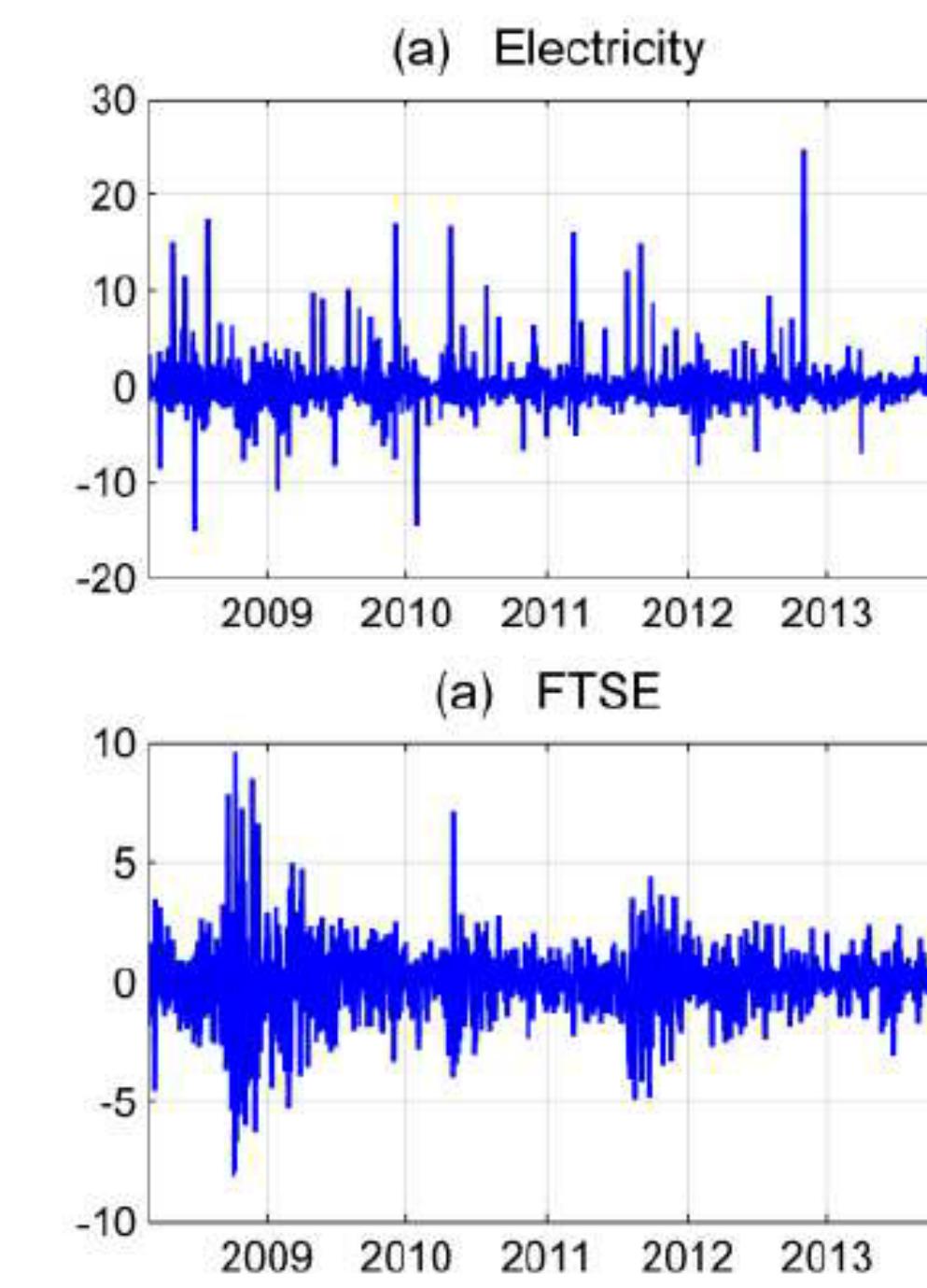
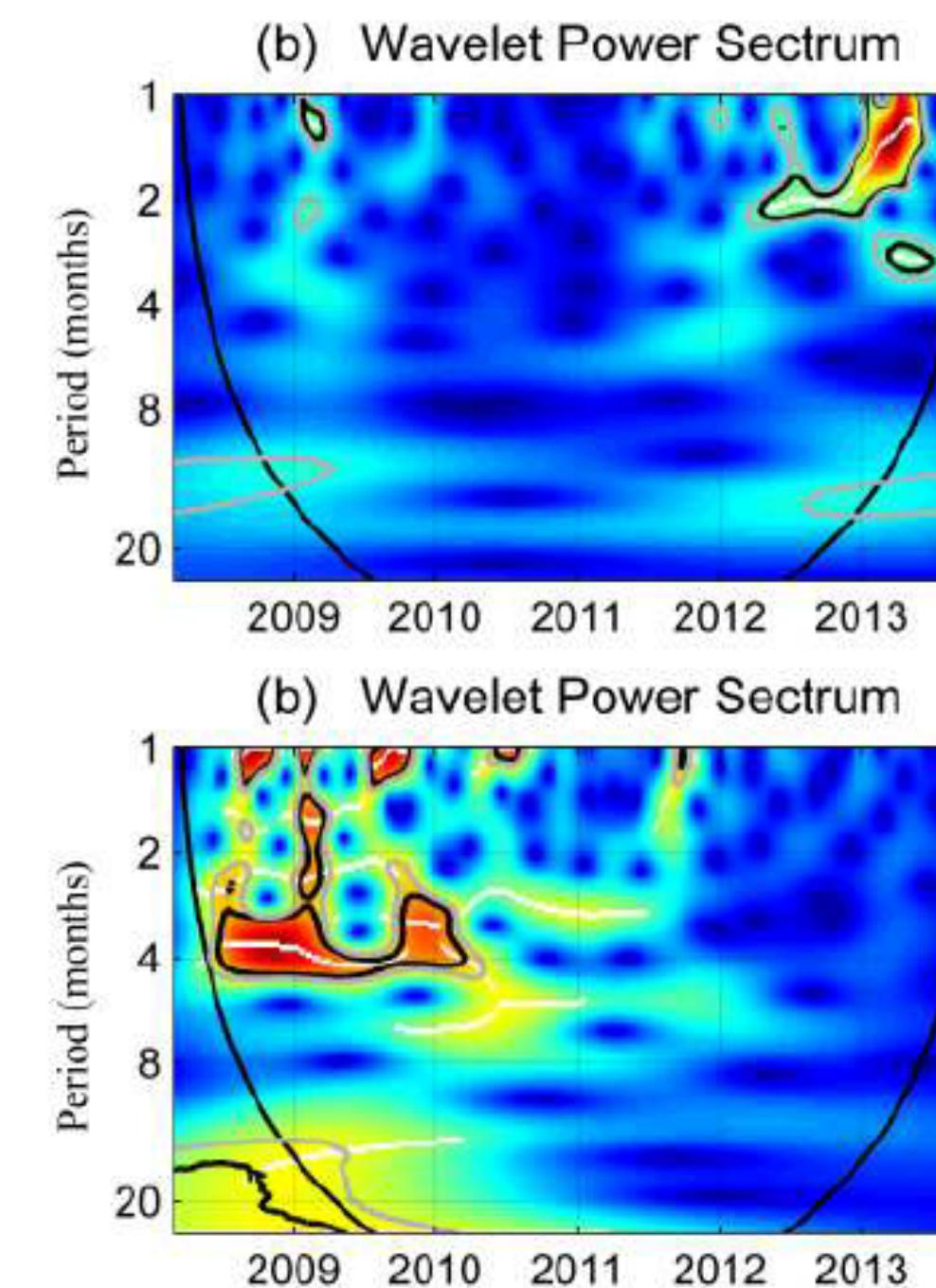
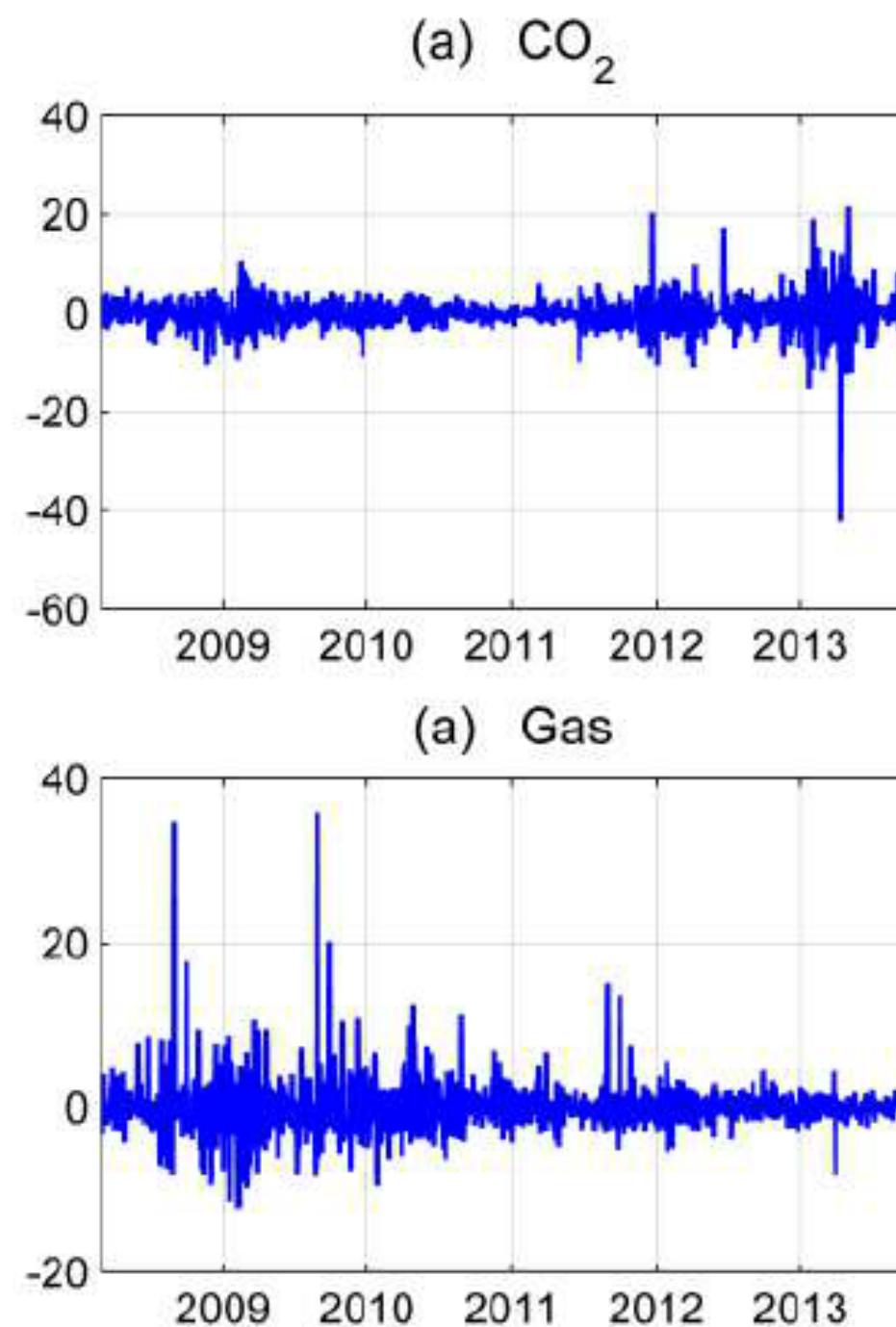
TS feature engineering (4/5)

- ▶ Exponential smoothing and other designed AR-functions
- ▶ Kalman filter family
- ▶ Robust regression filters

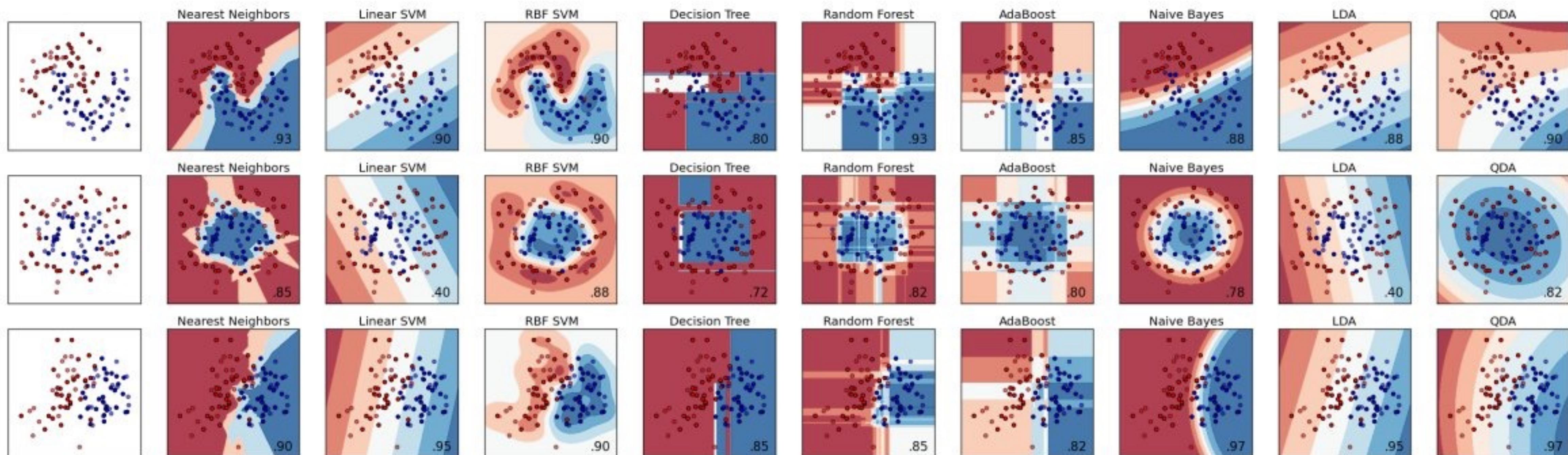


TS feature engineering (5/5)

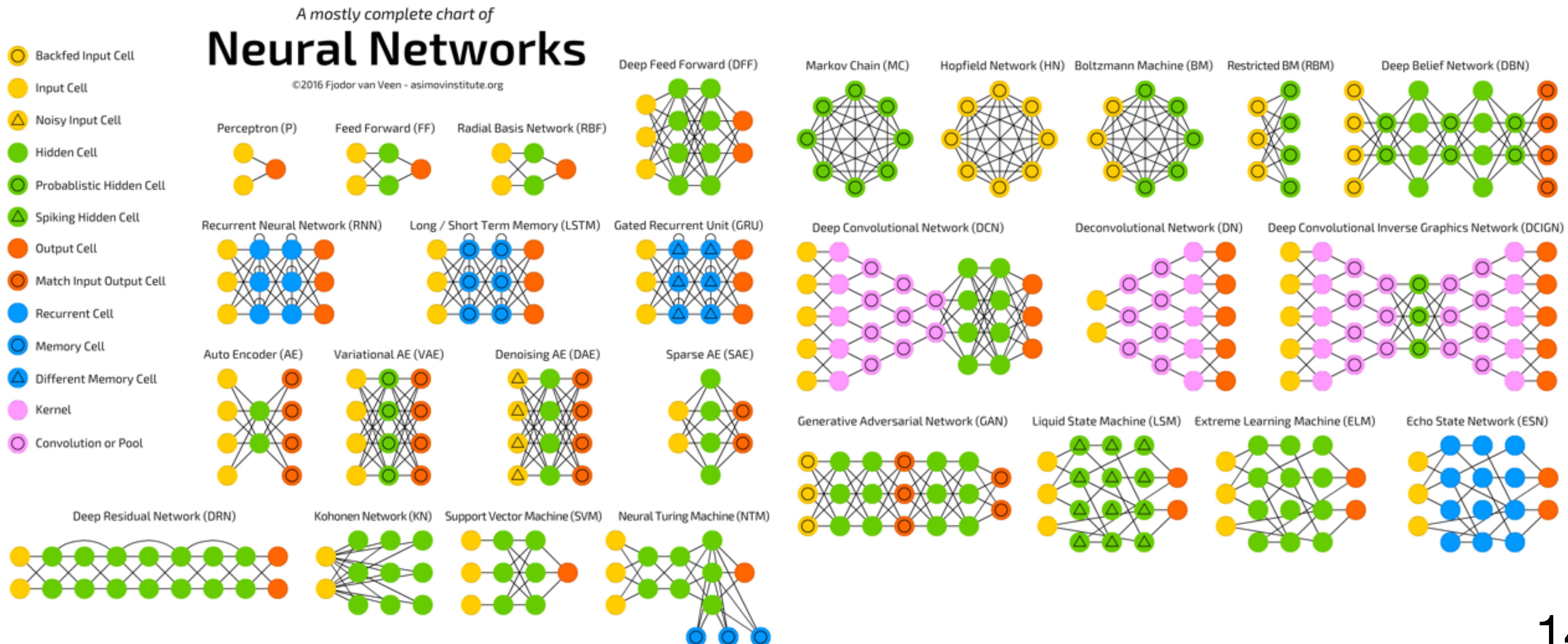
► Fourier and wavelet power spectrum



ML models vary a lot, there is no free lunch

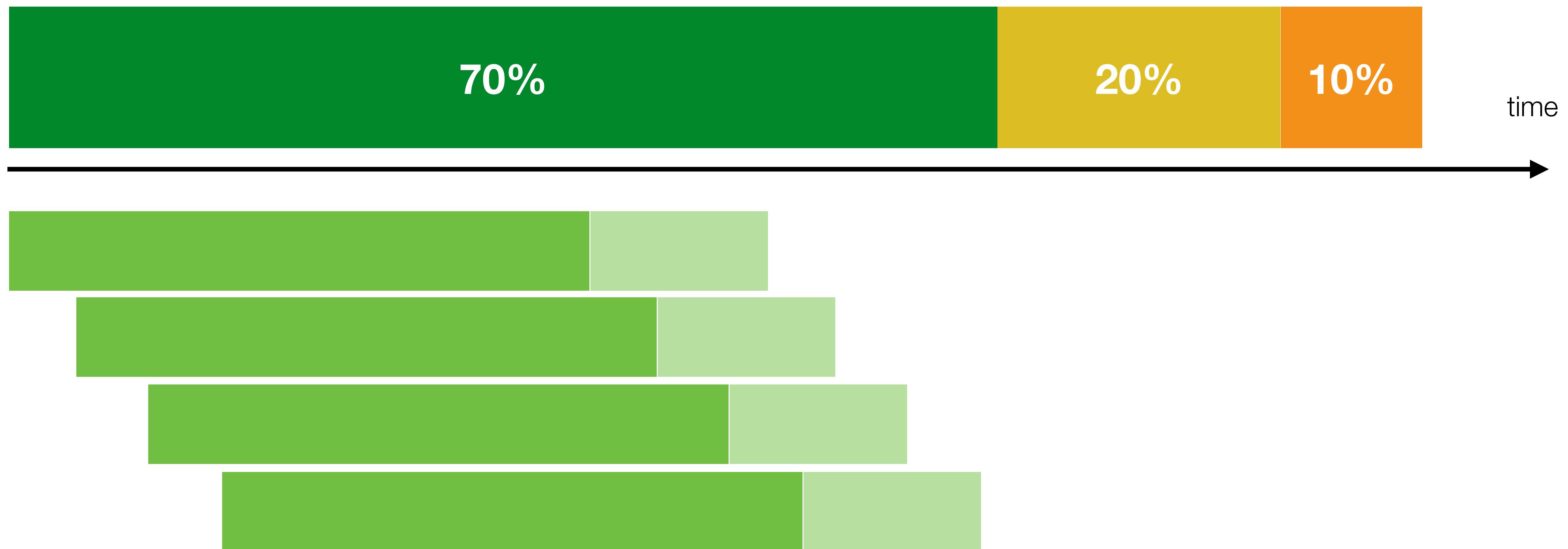


NN - one of many TS options:

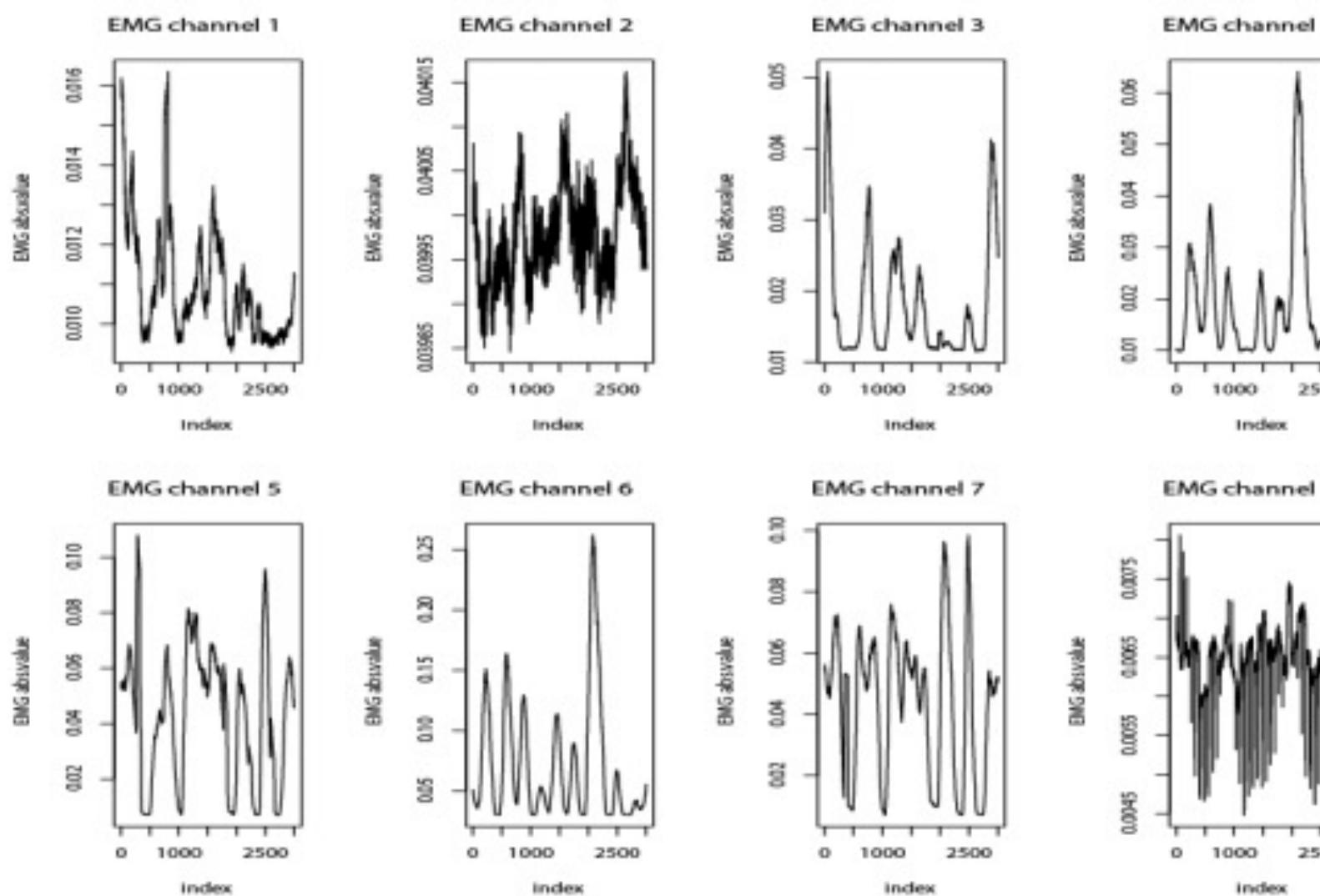
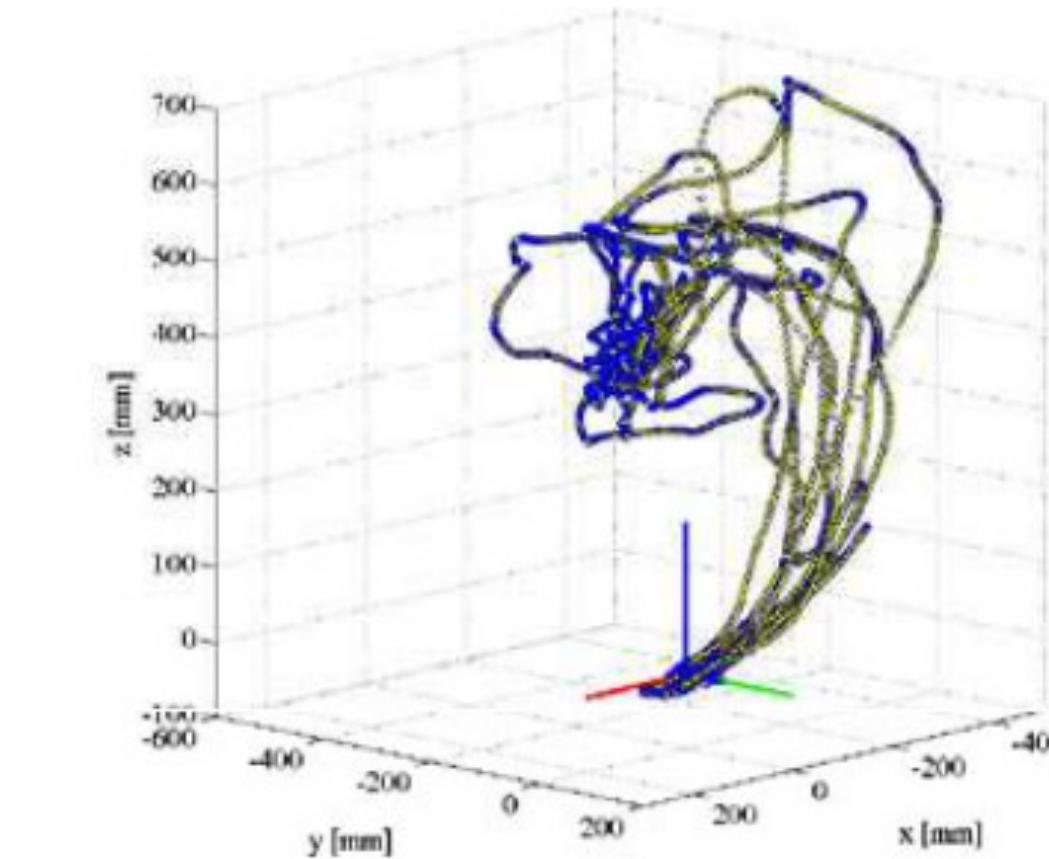
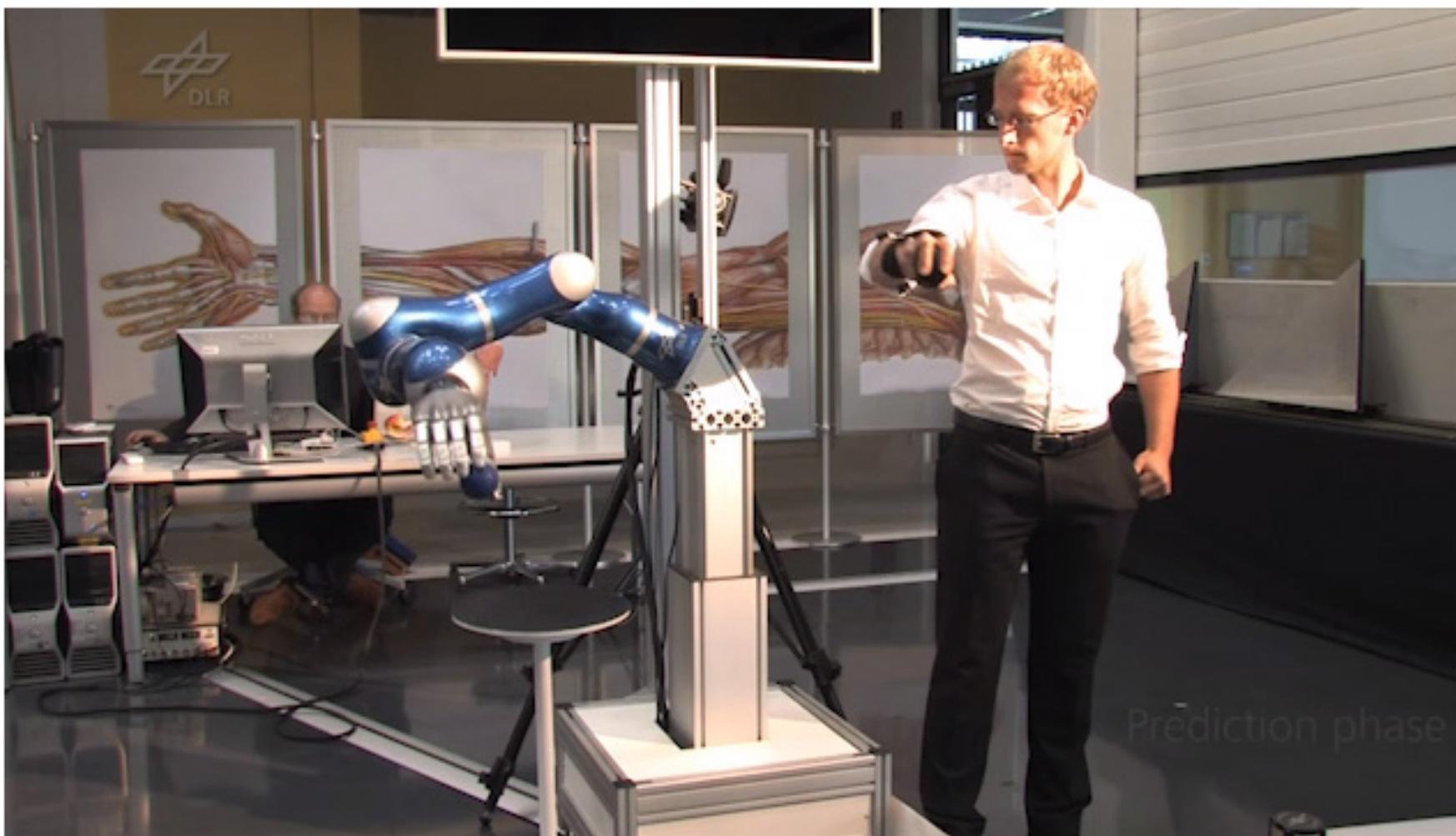




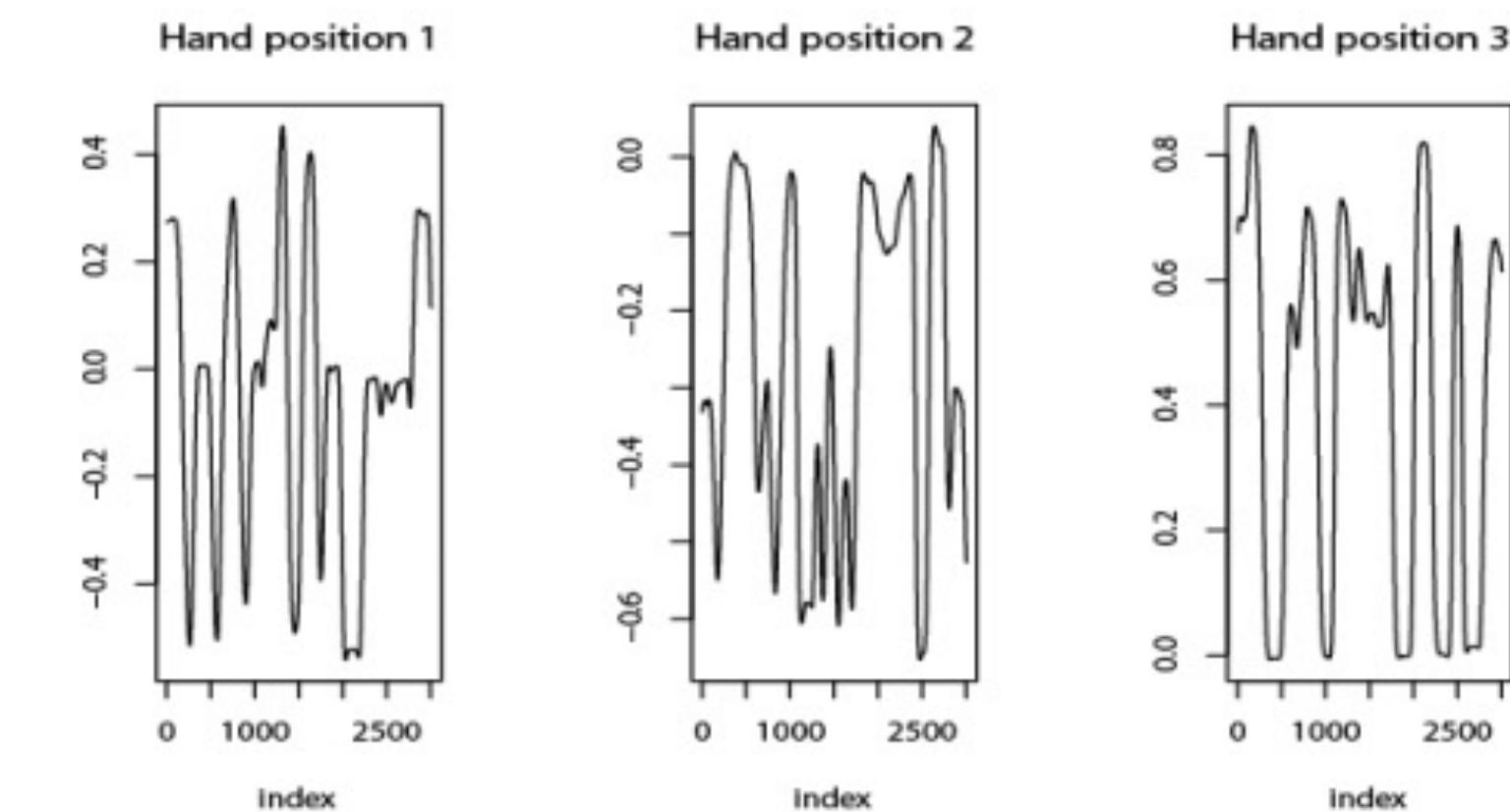
Time-series require time-CV.
K-fold leads to leaks and fails



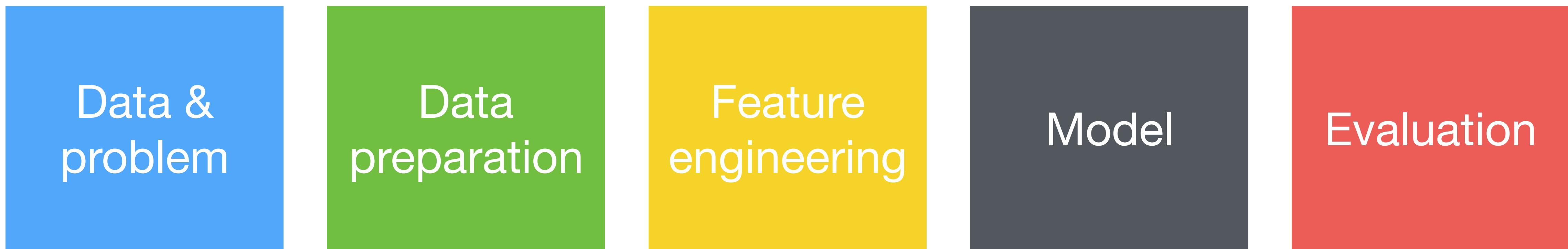
ML cases for time-series:



$$x \xrightarrow{f} y$$

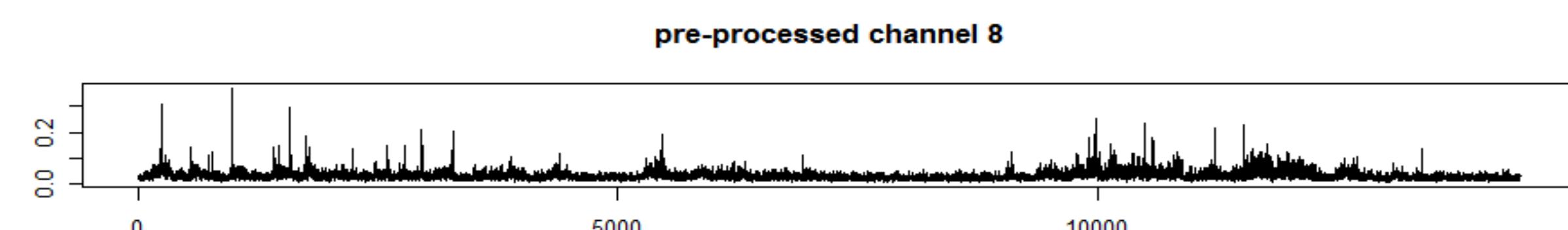
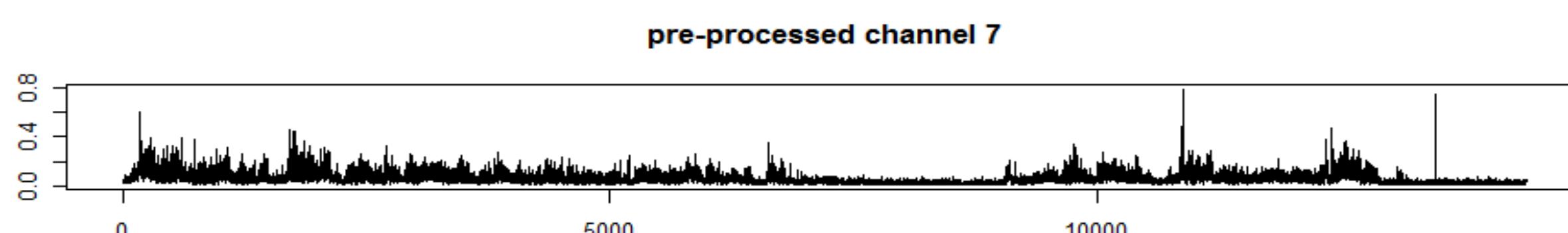
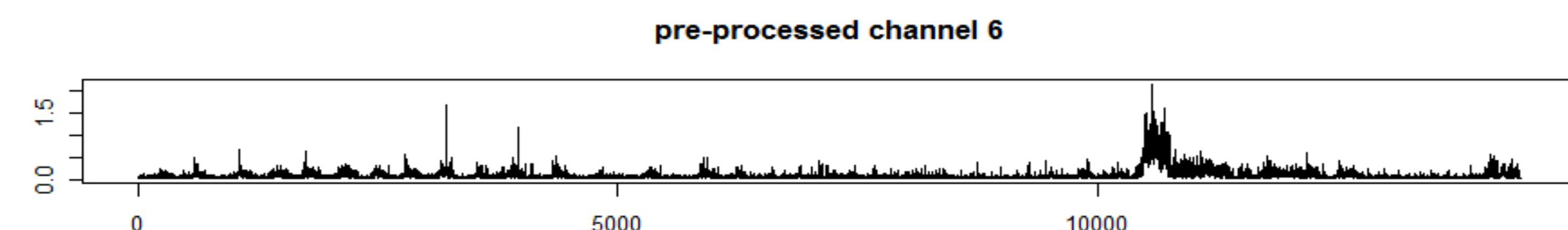
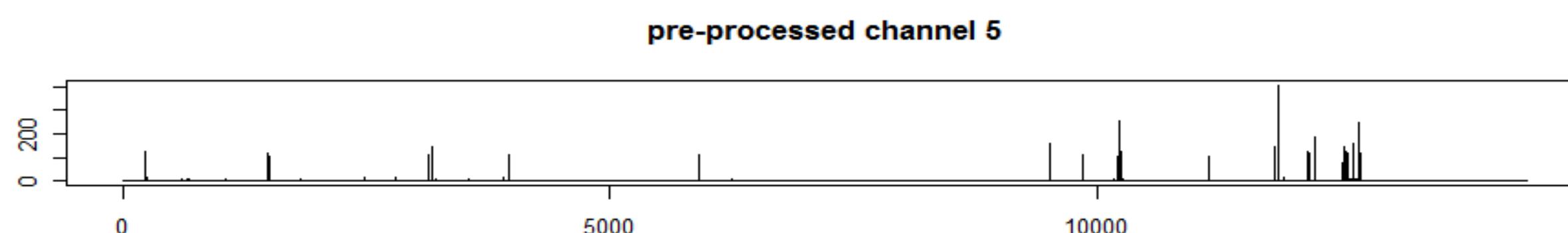
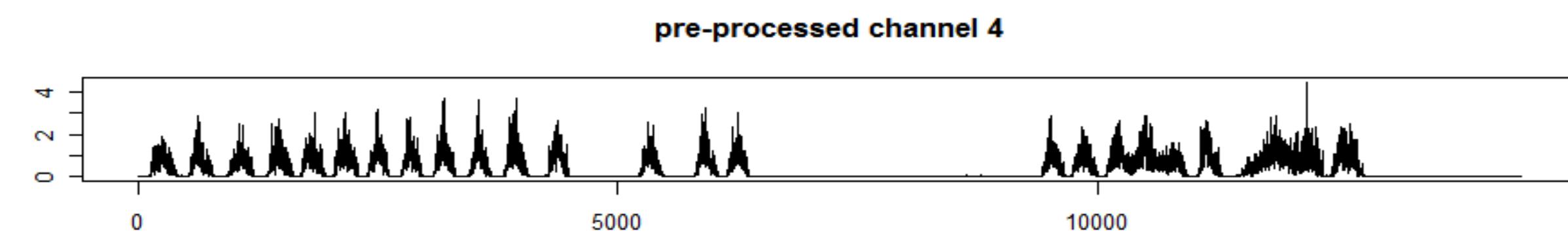
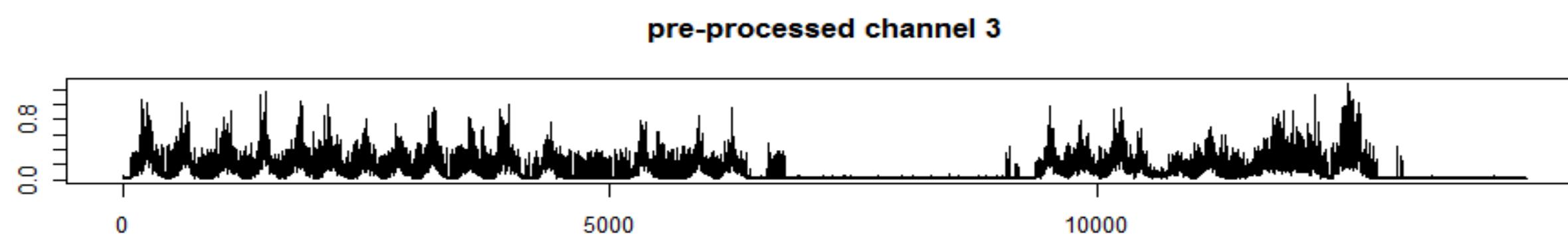
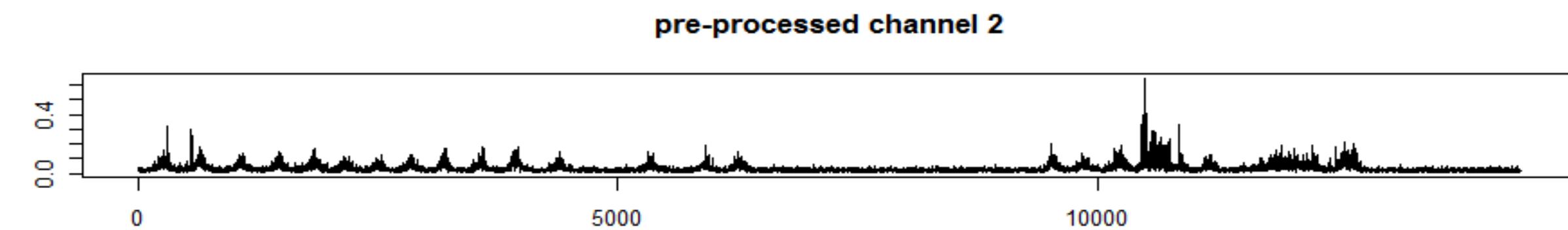
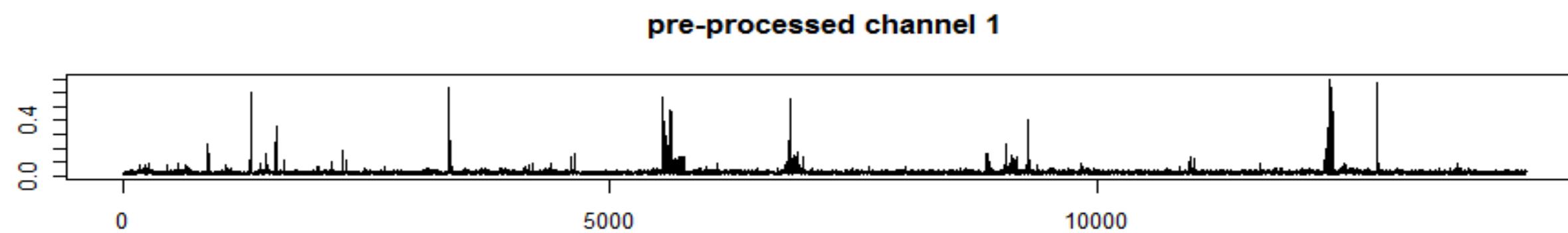


Case 1: ML pipeline



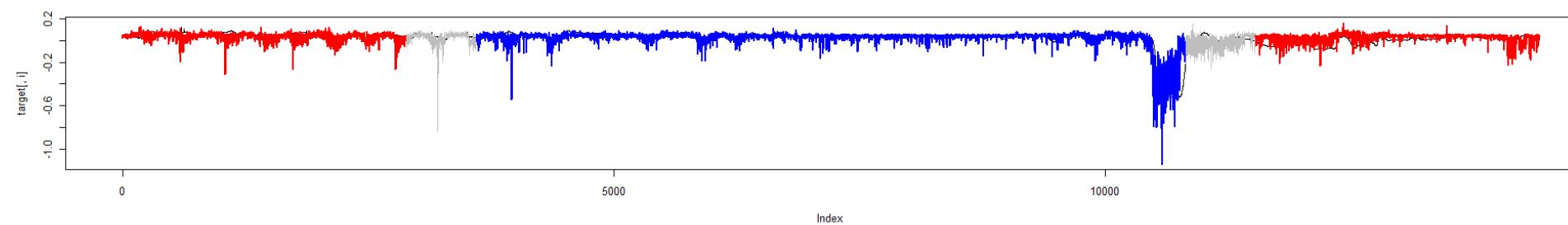
- | | | | | |
|-----------------------------------|---|--|---|--------------------------------|
| ‣ 8 noisy EMG sensors | ‣ Target and input have different frequencies | ‣ Stationary signals, but take diffs too | ‣ Classic ML does great: even linear regression | ‣ Split folds by sequences |
| ‣ Predict 3D position of the hand | ‣ Choose to aggregate inputs by max in 10 | ‣ Work on low frequencies, rolling stats are great | ‣ Cook GBM | ‣ Add post-processing (smooth) |

Case 1: input features

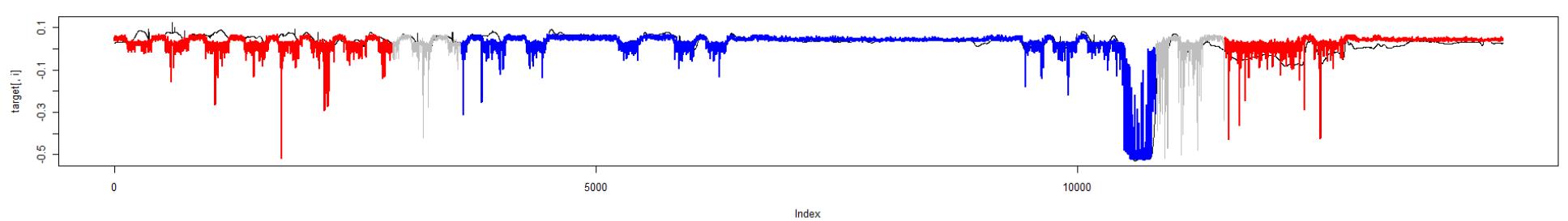


Case 1: naive baseline models

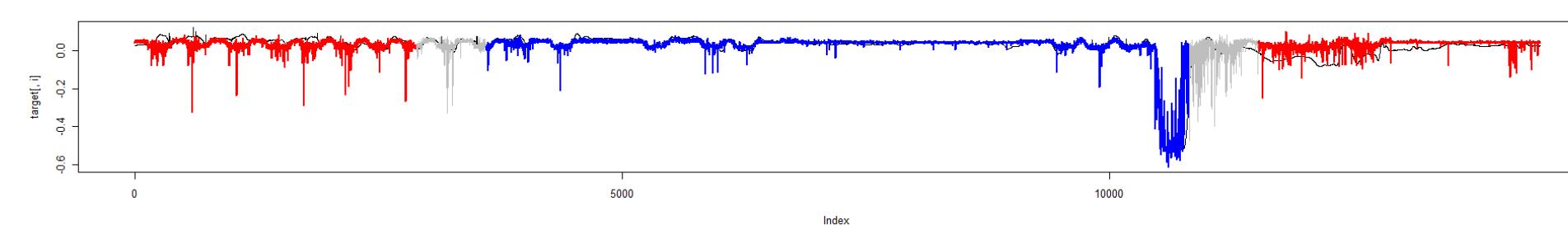
LR:



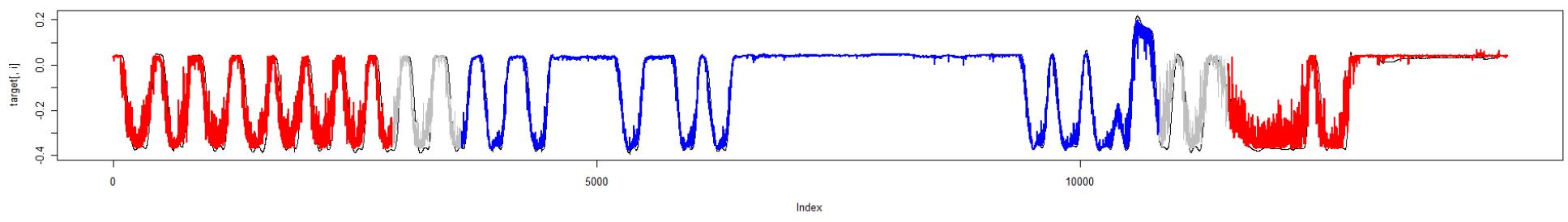
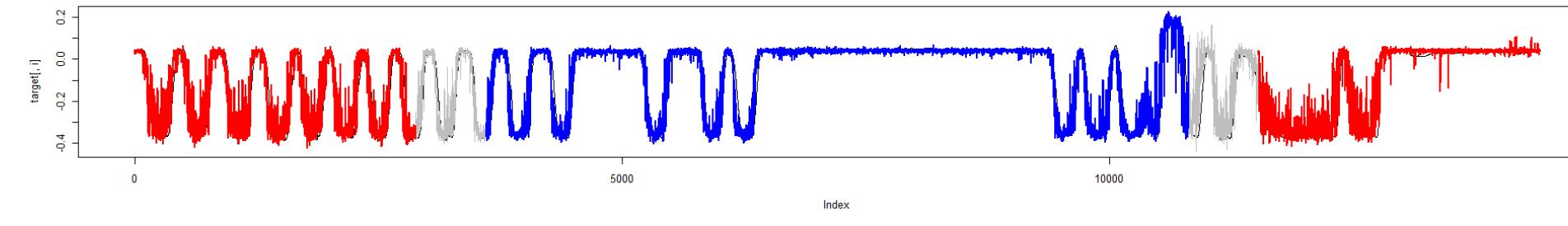
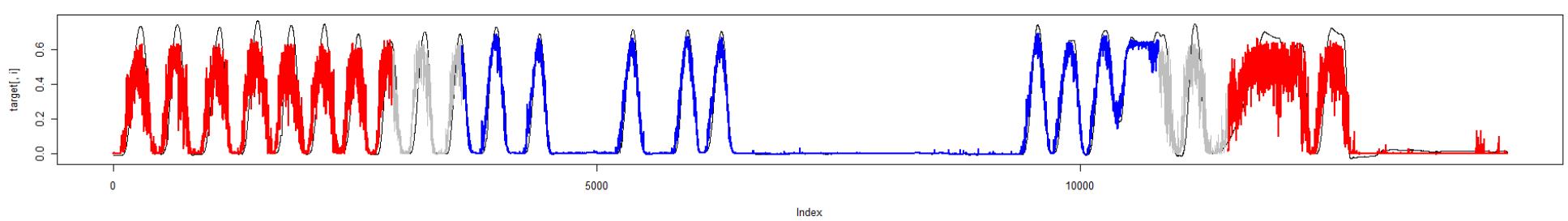
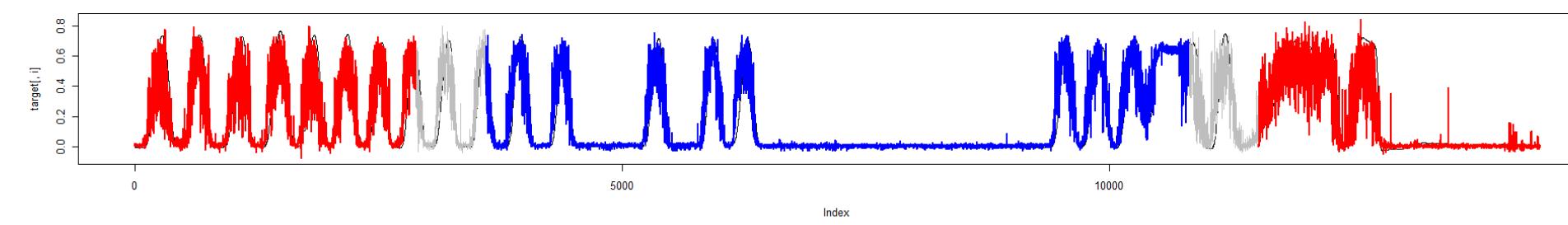
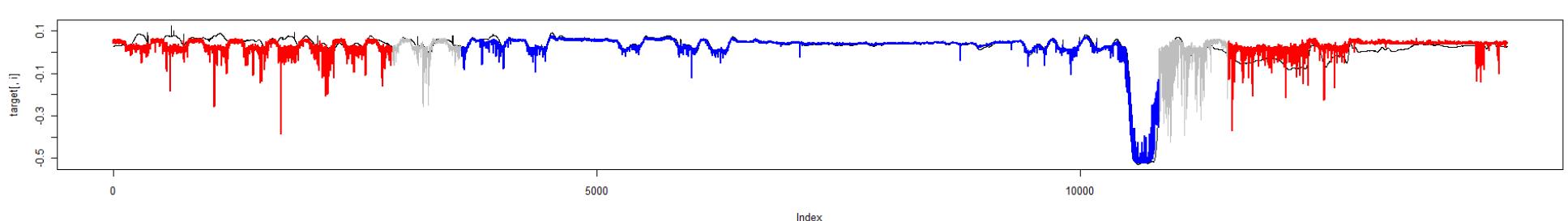
KNN:



SVM:

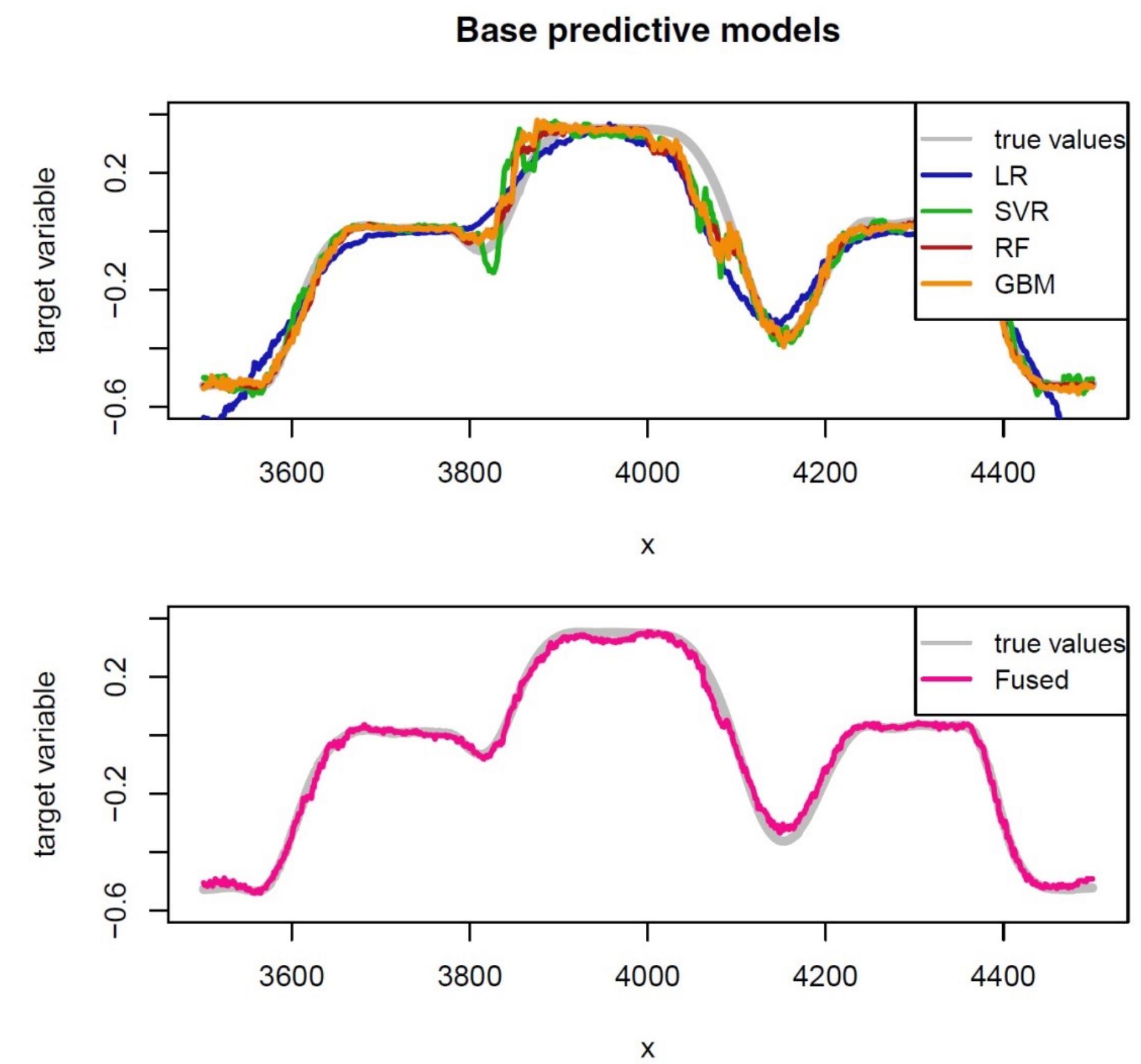
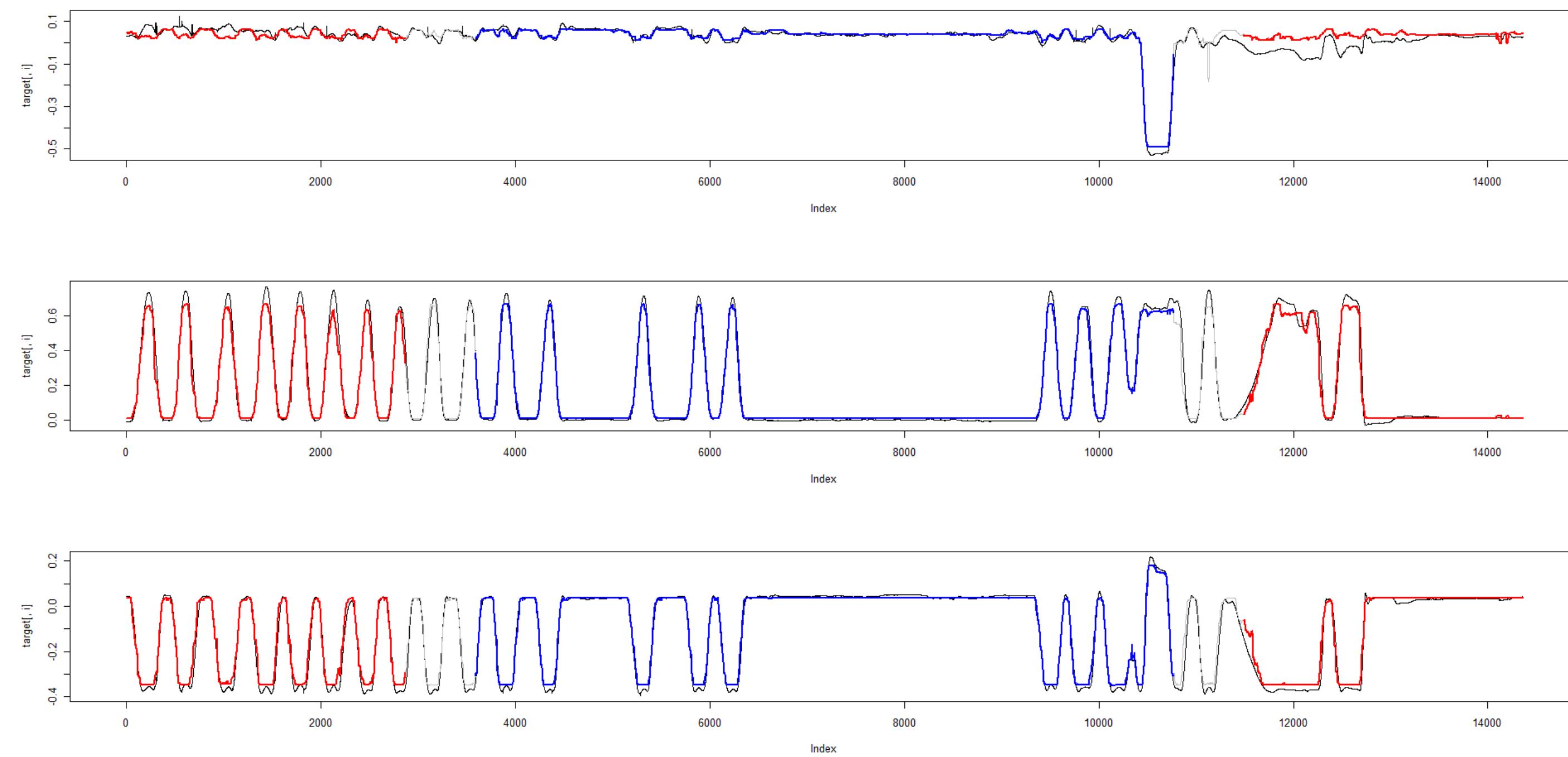


RF:

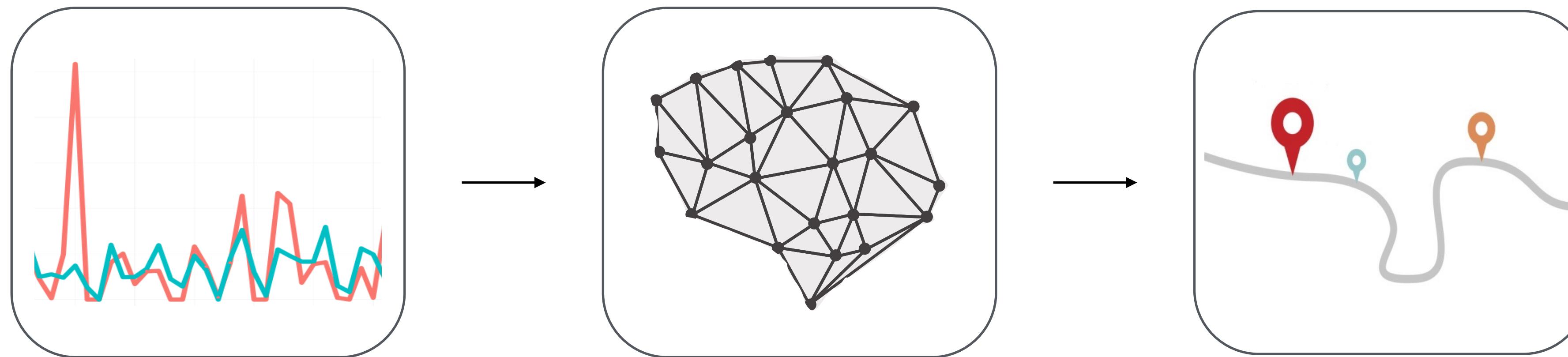


Case 1: result & post-process

- GBM with post-processing



ML cases for time-series:

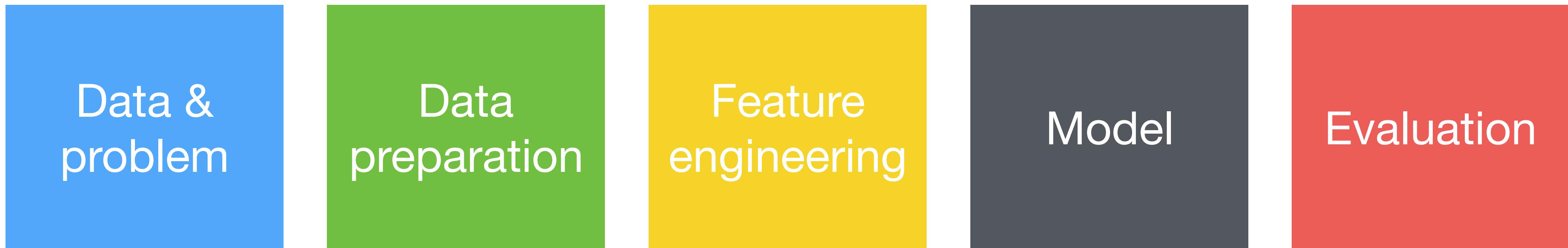


Cash demand prediction for each ATM device

Encashment candidate evaluation and simulations

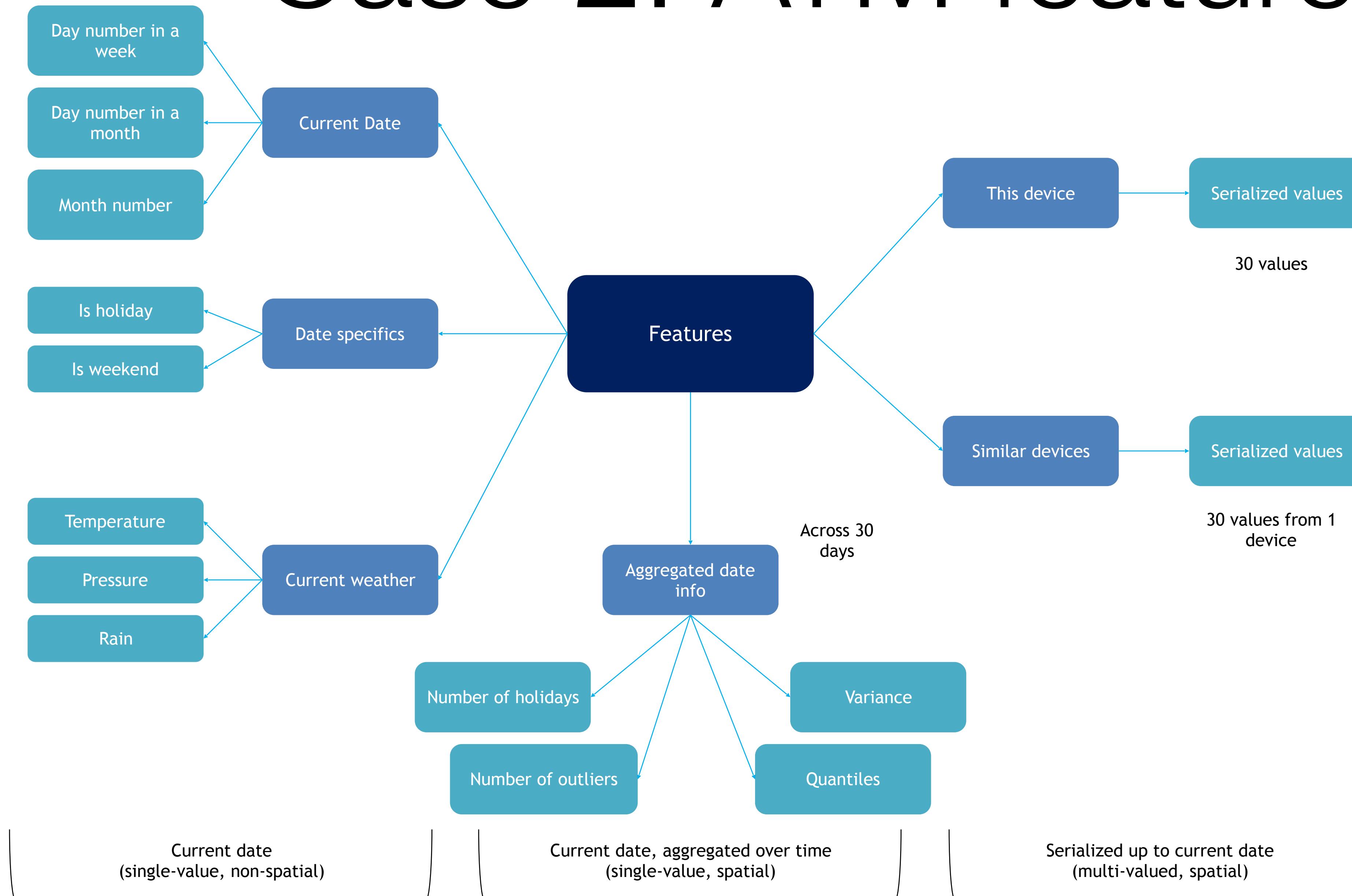
Optimize logistics for chosen candidates

Case 2: ML pipeline



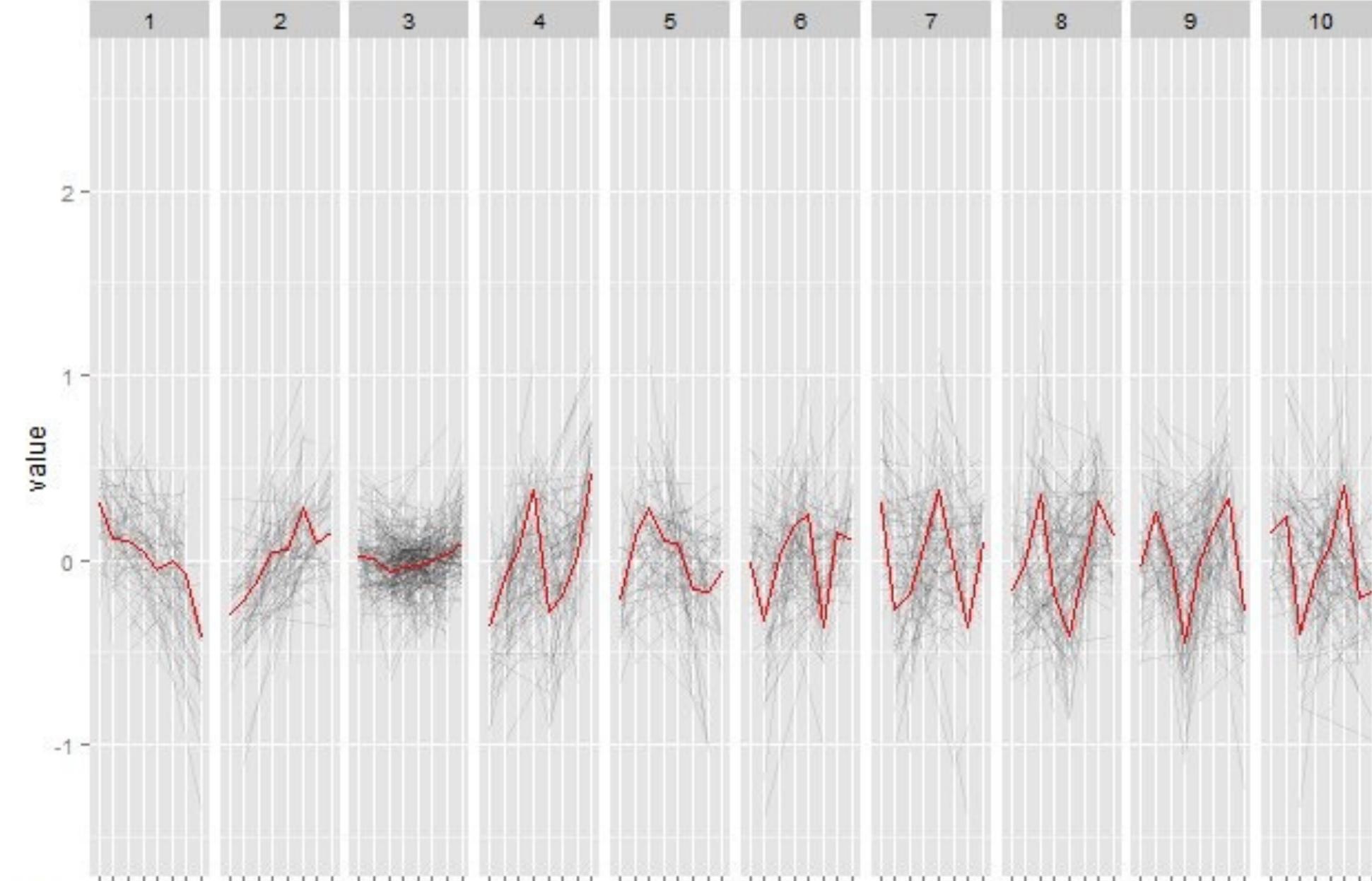
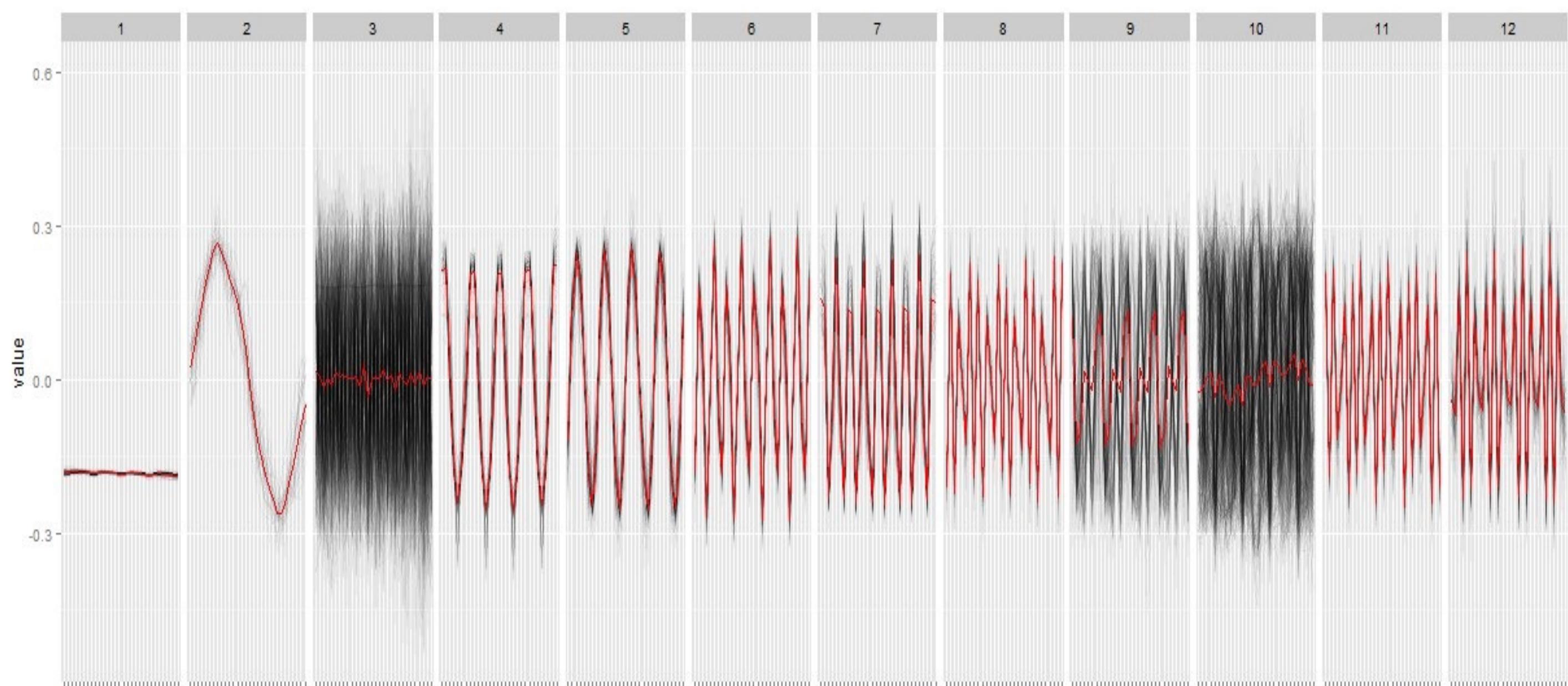
- A lot of external data sources besides cash
- Predict diffs and several total sums
- Fill NA to zero
- Adjust frequencies and time-grids of all external data
- Day-to-day
- AR and aggregates
- Aggregates on top of KNN
- Tried a lot
- GBM again was good
- Had to compete with ARIMA hard
- Time-CV
- Fair validation and test one year ahead

Case 2: ATM features

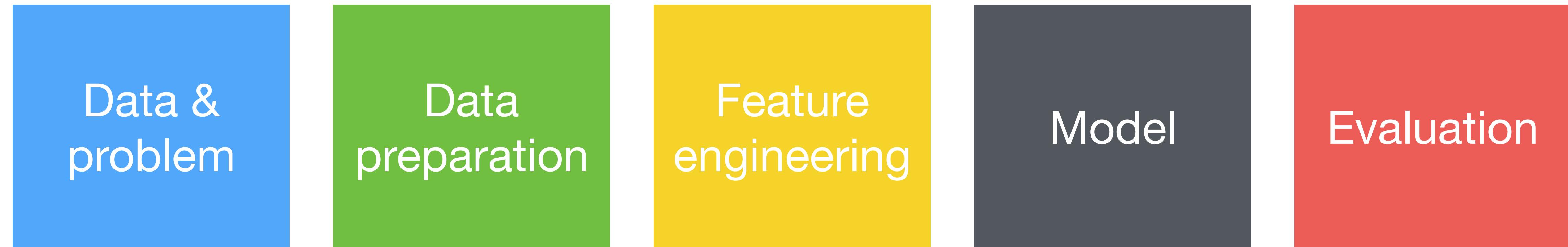


Case 2: curious stuff

- TS pattern mining was cool, but models already got all the info



Time-series ML pipeline



- | | | | | |
|---|--|---|--|--|
| <ul style="list-style-type: none">‣ Value vs diff‣ Regression/ classification /anomalies‣ Univariate sensors vs panel data | <ul style="list-style-type: none">‣ Consistent spacing‣ Impute missing values‣ Aggregates/ pooling | <ul style="list-style-type: none">‣ Stationarised series‣ AR, MA and aggregates‣ Dense represent-ns‣ Kaggle-style | <ul style="list-style-type: none">‣ Linear‣ Casual ML & boosting‣ LSTM & RNN‣ Markov-style‣ Kaggle-style | <ul style="list-style-type: none">‣ Time-CV‣ Test-policy (rolling or not)‣ Simulations (small data) |
|---|--|---|--|--|

Can you consistently forecast
values of non-stationary time series
with a tree-based model like GBM?

- A. Sure, watch me xgboost it
- B. Maybe, if we work with diffs
- C. No, trees can't extrapolate
- D. No, non-stationarity hurts

Can you consistently forecast
values of non-stationary time series
with a tree-based model like GBM?

- A. Sure, watch me xgboost it
- B. Maybe, if we work with diffs
- C. No, trees can't extrapolate
- D. No, non-stationarity hurts

Summary

- ▶ Key time-series differences from classic ML:
proper time-CV, specific preparation issues
- ▶ The big ML picture for time-series makes it easier
to solve practical problems
- ▶ On the other hand the magic disappears and you
forget names of all those methods

Thank you!

ods.ai

@natekin

alex.natekin@dmlabs.org

