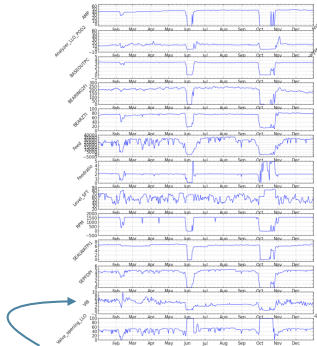


The Signals

- 25 continuous signals
 - sensors, automated: AMP, vibration, RPM, ...
 - relatively clean data: some noise but complete
- 3 discrete signals
 - log files, manual: maintenance, cleaning
 - logging is error prone: noisy / incomplete data



Maintenance when vibration crosses threshold value
Goal: Forecast Vibration Signal



Data Volume

- Timespan: 2008 - 2012
- asynchronous and different sample-rates
 - 53k to 5671k values per signal
- Original data delivered in .csv files
 - 156 files: one file per signal, per year.
 - 1.51Gb in total

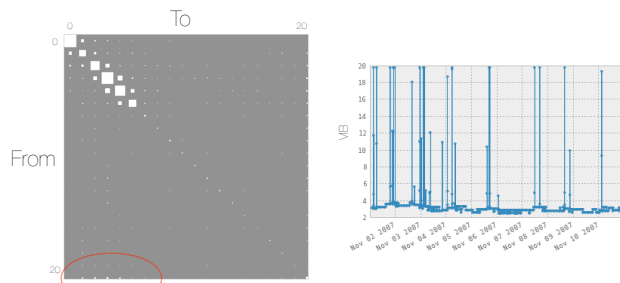
We have **“Big Data”**
but how can we extract useful knowledge?



Data Pre-processing

Peaks Analyses

After a peak, the signal drops again immediately



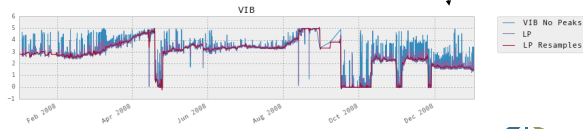
Bi-gram: visual representation of transition probabilities



Peak filtering

- We are only interested in the long term trend
- Peaks reduce model accuracy
 - Not part of the trend, irrelevant
 - introduce noise in model
- Remove peaks using a low-pass filter

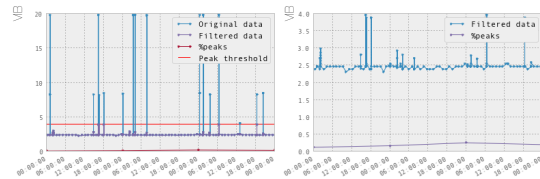
Expert knowledge



Additional Peak features

- We filter peaks out
 - long term vibration trend remains
 - vibration is a manifestation of slowly evolving degradation process
- There might still be some information in the peaks:
 - e.g. degradation causes more peaks
 - Add feature: % of peaks over a time interval (1 day)
 - Add feature: Variation of the signal over a time interval

Expert knowledge

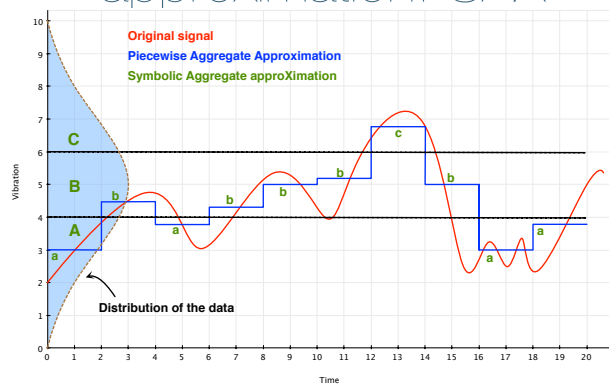


Signal synchronisation and representation

- Signals are out of sync and have different sample rates
 - Necessary to resample the signal
- High sample-rate causes huge feature vectors
 - high input dimensionality, noisy values, overfitting
 - under-sample, smooth the signal
- Signals lie in different ranges:
 - RPM [0-500], Vibration [0-10]
 - map values to same domain (SAX)

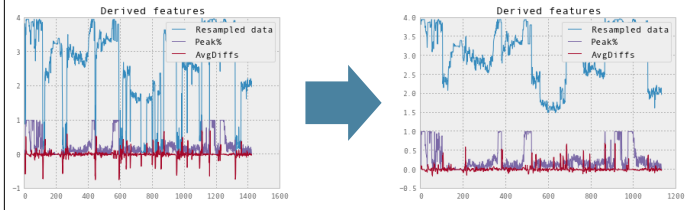


Symbolic Aggregate approximation: SAX



Ignore idle status

- Idle status not part the trend
- We do not want to model/learn when the machine is idle
- Remove part of signal where the machine is idle (low RPM values)



HDF5

- Hierarchical Data Format
 - High performance data-format
 - Bindings for C/C++/Java/Python/Matlab/Fortran...
- Parsing the csv files to time-series objects takes 10minutes = very slow
- Parsed time-series stored in an HDF store
 - Loading the data from the HDF store is nearly instantaneous (limited by disk-speed)
 - load time from 10min to <1sec
 - 60% file-size reduction from 1.51GB to 392MB

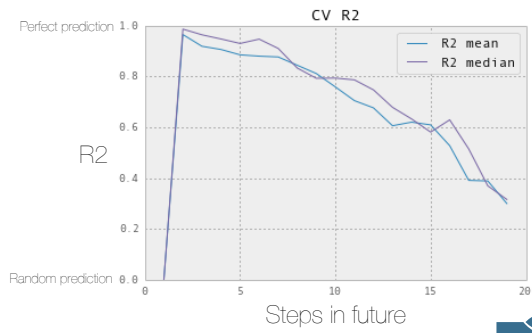
Predictive Modelling

Learning Task

- Input: Variables of the SAX signals with history of 5 days
- Output: Predicted vibration in 10 days
- Model: Random Forest
 - Can deal with multiple input streams
 - Can deal with continuous variables
 - Easy to implement in production
 - Robust against non-informative features (automatic feature selection)



Predicting VIB: future steps



Conclusions

- Lots of data doesn't equal good data
- Preprocessing is important
- Expert knowledge helps design informative features
- We can't model what we haven't seen before
 - Models need to be updated regularly

