

Automating Feature Engineering for Supervised Learning?

methods, open-source tools and prospects.

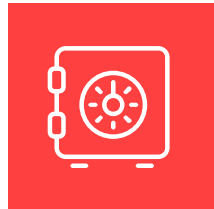
Dr. Thorben Jensen

Jonas Jostmann

INFORMATIONSFABRIK



MÜNSTER





INFORMATIONSFABRIK
DATEN VERSTEHEN. ENTSCHEIDUNGEN TREFFEN.

**Towards automating [supervised] machine learning:
Benchmarking tools for hyperparameter tuning**

Dr. Thorben Jensen
tjensen@informationsfabrik.de



#pydatabin
@pydataberlin

Content



- Feature Engineering
- Open-Source Tools
- Case Study

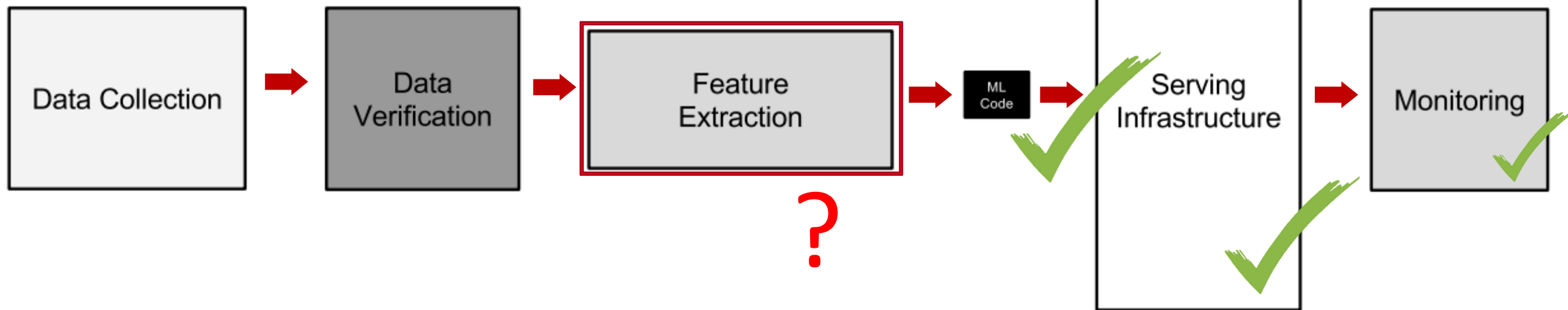
Feature Engineering

Workload in Supervised Learning



ML
Code

Supervised Learning & AutoML



Machine Learning eliminates the need for detailed programming effort.
(Arthur Samuel, 1959)

AutoML *eliminates the need for detailed Machine Learning effort.*

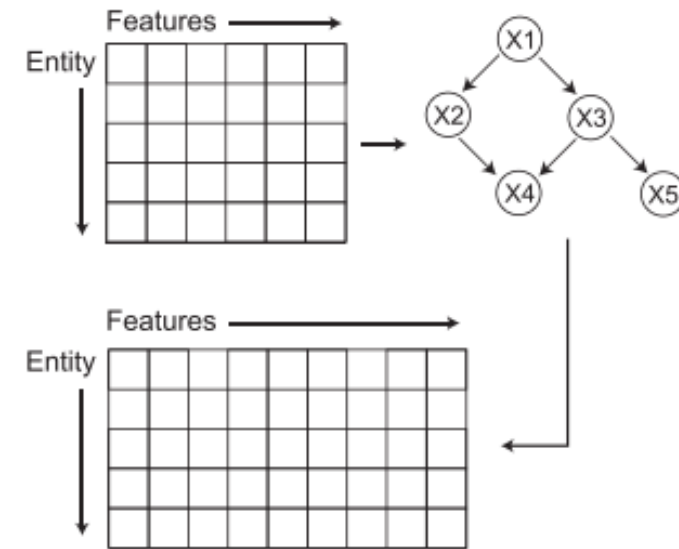
Sculley, David, et al. "Hidden technical debt in machine learning systems." *Advances in neural information processing systems*. 2015.

Samuel, Arthur L. "Some studies in machine learning using the game of checkers. II—recent progress." *Computer Games I*. Springer, New York, NY, 1988. 366-400.

Feature Engineering in Practice



Think of explanations



DATE	SALES (TARGET)	WEEKDAY	OPEN	STATE HOLIDAY	„BRIDGE DAY“	LAG(OPEN)
2014-04-29	5923	Friday	1	0	0	1
2014-04-30	5870	Saturday	1	0	0	1
2014-05-01	0	Thursday	0	1	0	1
2014-05-02	6790	Friday	1	0	1	0
2014-05-03	5498	Saturday	1	0	0	1

Current Challenge & Vision for Automated Feature Engineering



CHALLENGES

Time-consuming task

Good features require domain knowledge

Features often not transferable between data



VISION

- Easy setup (data and config)
- Handle time in data
- Features self-explanatory
- High potential complexity of features
- Minimal „garbage“ features

Open Source Tools

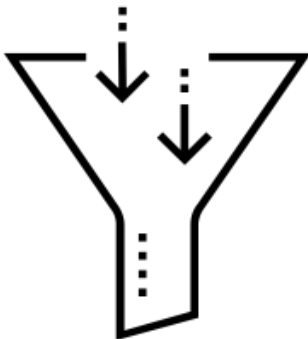
Options for Automated Feature Engineering



GENERATION



SELECTION



Options for Automated Feature Engineering



GENERATION

DIMENSIONALITY REDUCTION

Combine similar features, e.g. PCA.

TRANSFORMATION

Apply operators (+, *, Min, ...)

RELATIONAL OPERATIONS

SQL-like operations

TIME-SERIES FEATURES

Pre-defined functions (e.g. trend)

PRE-TRAINED EMBEDDINGS

Extract „features“ with pre-trained models

META-LEARNING

Predict good features based on previous datasets

SELECTION

FILTERING

Select features by testing interrelation with target variable.

WRAPPING

Search best feature set, by model performance.



Options for Automated Feature Engineering



GENERATION

DIMENSIONALITY REDUCTION

Combine similar features, e.g. PCA.

TRANSFORMATION

Apply operators (+, *, Min, ...)

RELATIONAL OPERATIONS

SQL-like operations

TIME-SERIES FEATURES

Pre-defined functions (e.g. trend)

PRE-TRAINED EMBEDDINGS

Extract „features“ with pre-trained models

META-LEARNING

Predict good features based on previous datasets

(autofeat)

TSFRESH



SELECTION

FILTERING

Select features by testing interrelation with target variable.

TSFRESH

WRAPPING

Search best feature set, by model performance.



(autofeat)

Database Tables:

OrderID	Customer ID
1	2
2	...
3	...
4	2
...	...

ID	OrderID	ProductID
1	1	3
2	1	1
3
4	4	3
...

ProductID	Price
1	\$100
2	...
3	\$200
4	...
...	...

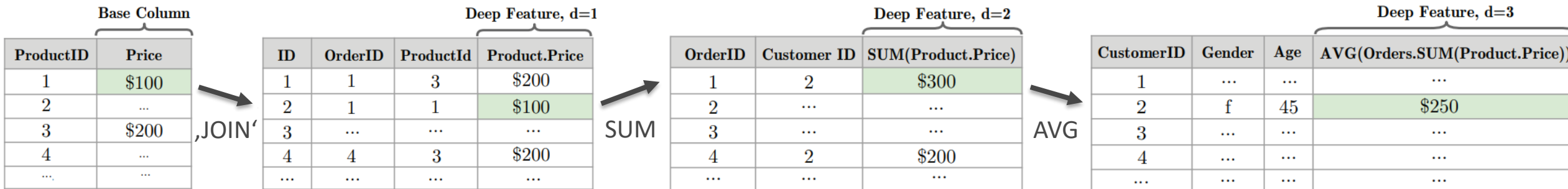
Features on Customer?

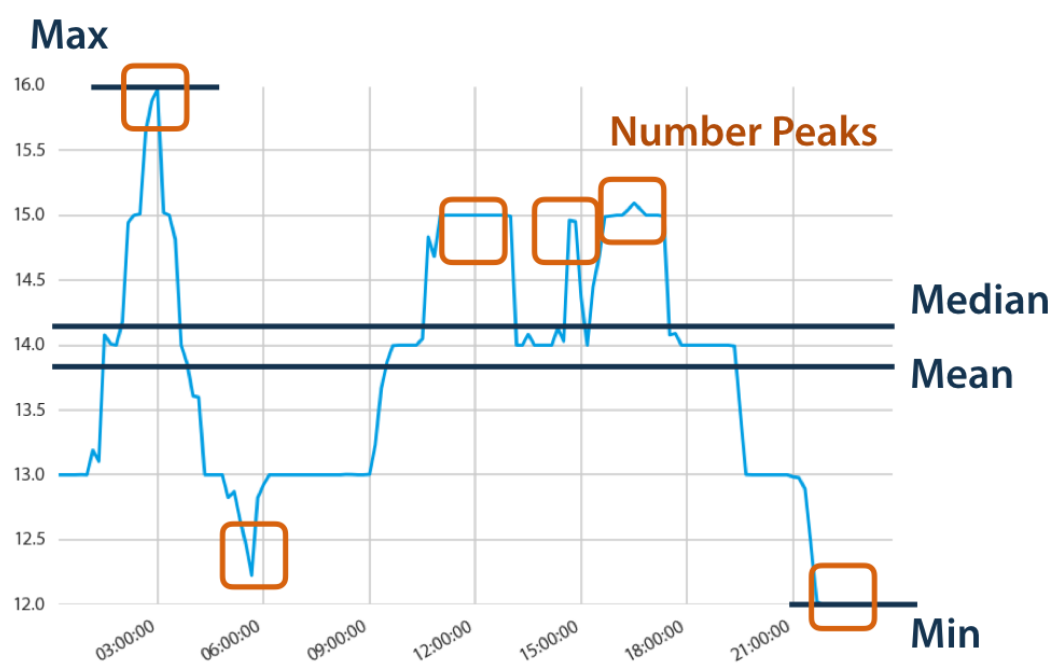
Aggregations (per value):

- sum
- trend
- time_since_first
- ...

Transformations (per entity):

- month
- weekday
- num_words
- ...



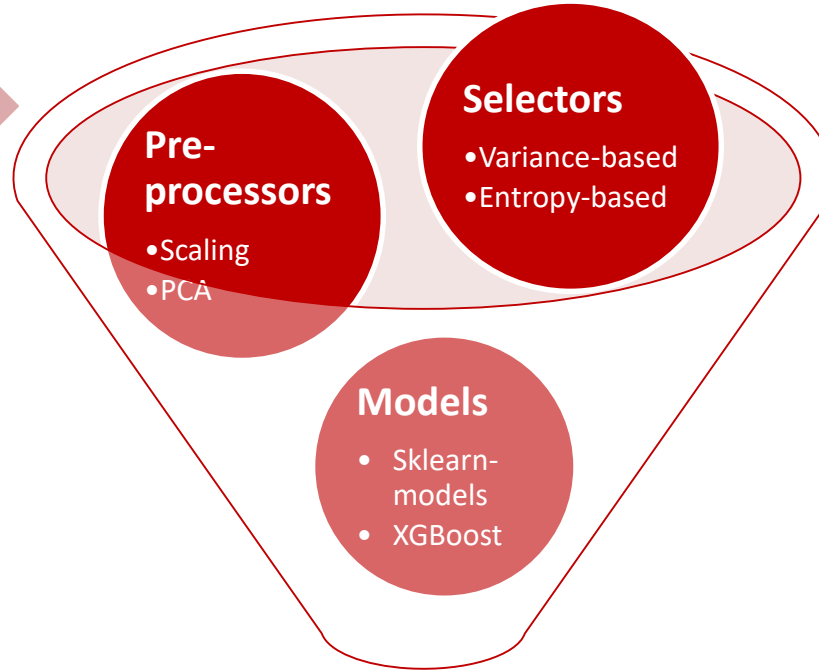


Numerous features:

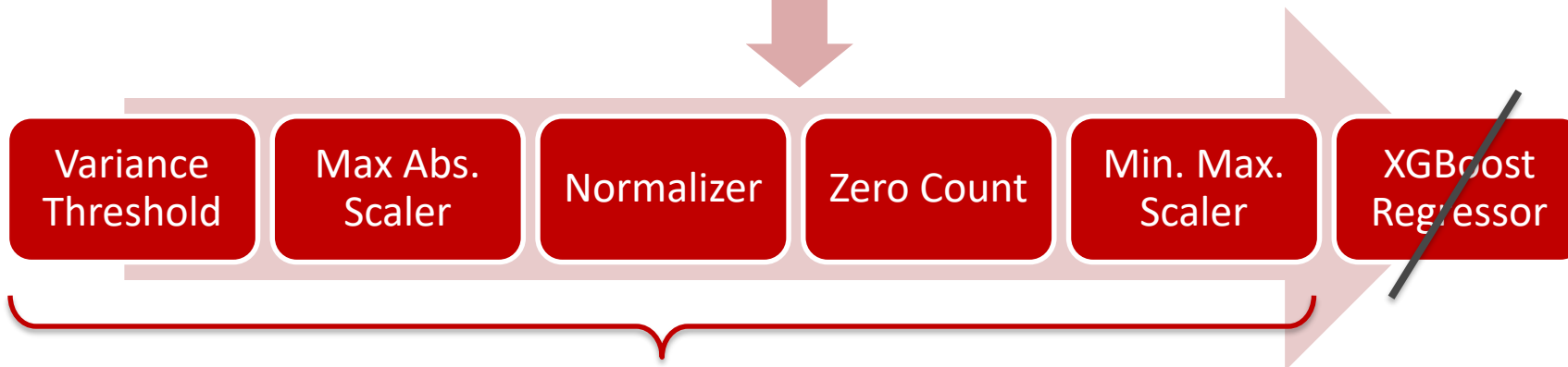
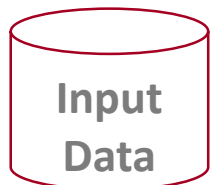
- `linear_trend(x, "slope")`
- `mean_change(x)`
- `autocorrelation(x, lag: int)`
- `spkt_welch_density(x, ...)`
- ...

Filtering most promising features:

1. Variables tested individually
2. Statistical tests on interaction with target variable
3. Discarding variables with low significance



Genetic Programming ...



„Feature Engineering“ Pipeline



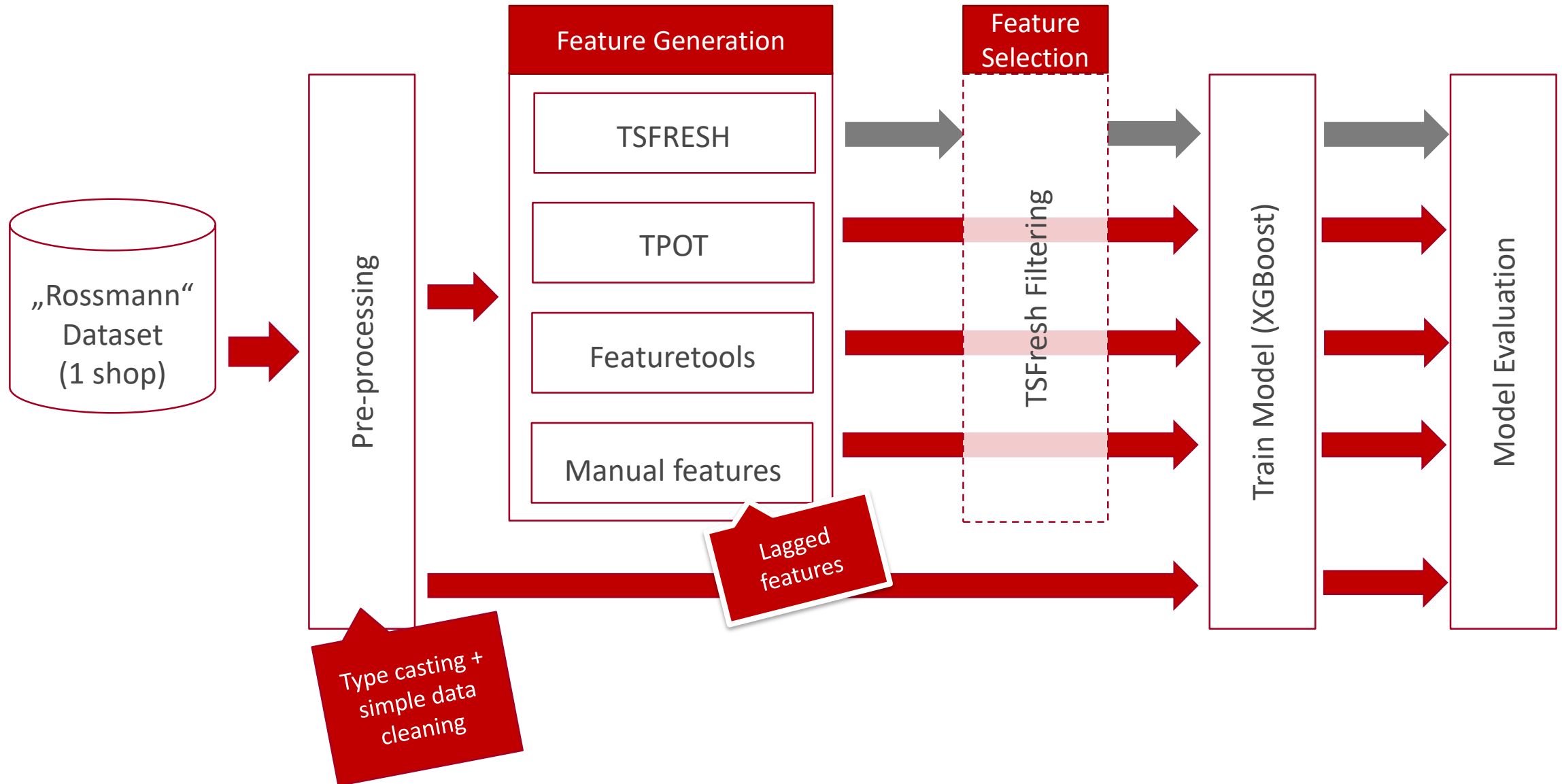
Case Study



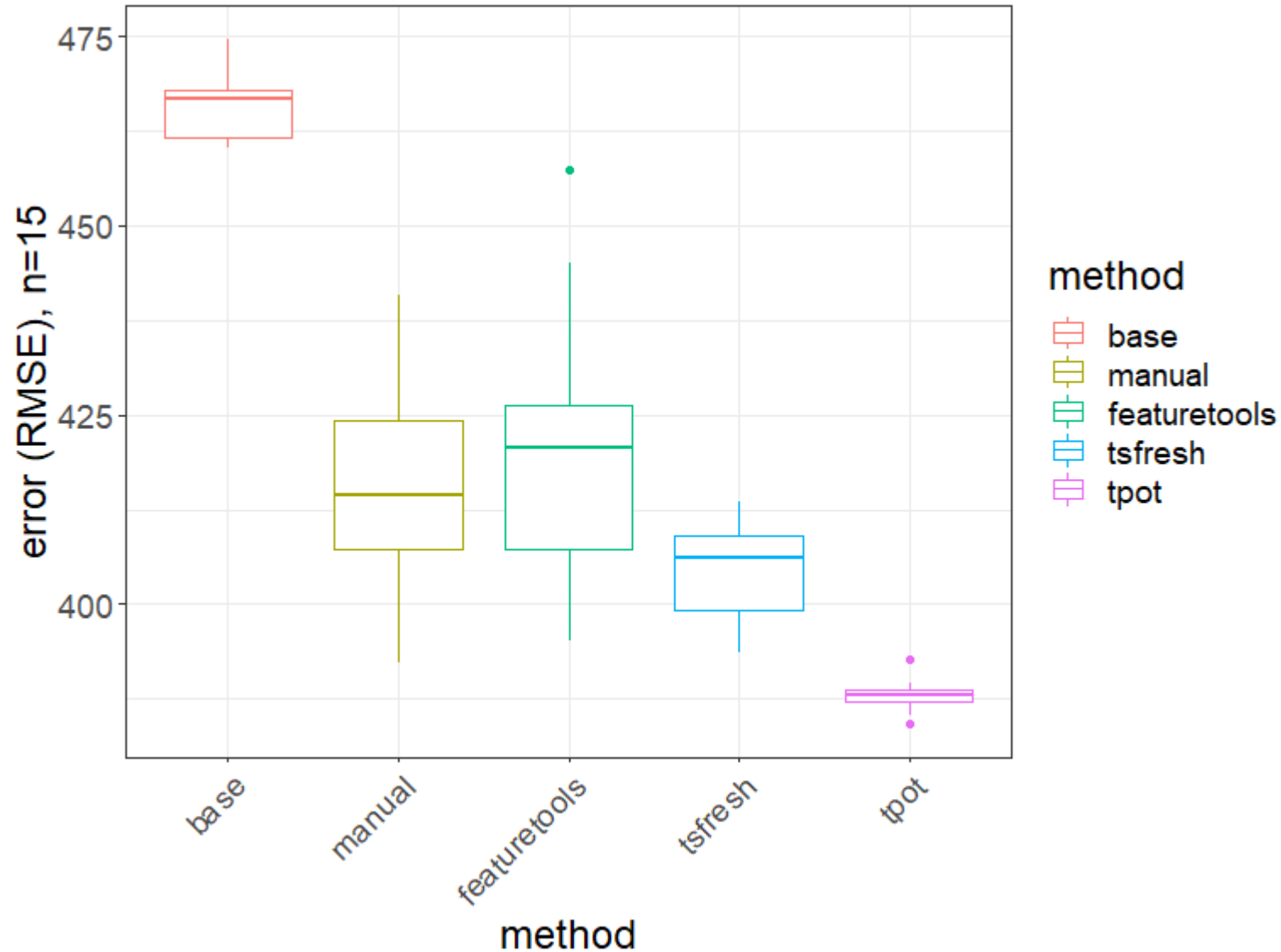
Goals and Dataset

- Goals:
 - Test libraries at default settings
 - Achieve **general understanding** of tools
 - Compare performance
- Machine-learning Task
 - Time series forecasting („Rossmann Challenge“)
 - XGBoost with fixed settings
- Out of scope
 - Not finding best library *in general*

Design of Experiment



Results Performance: methods at default usage

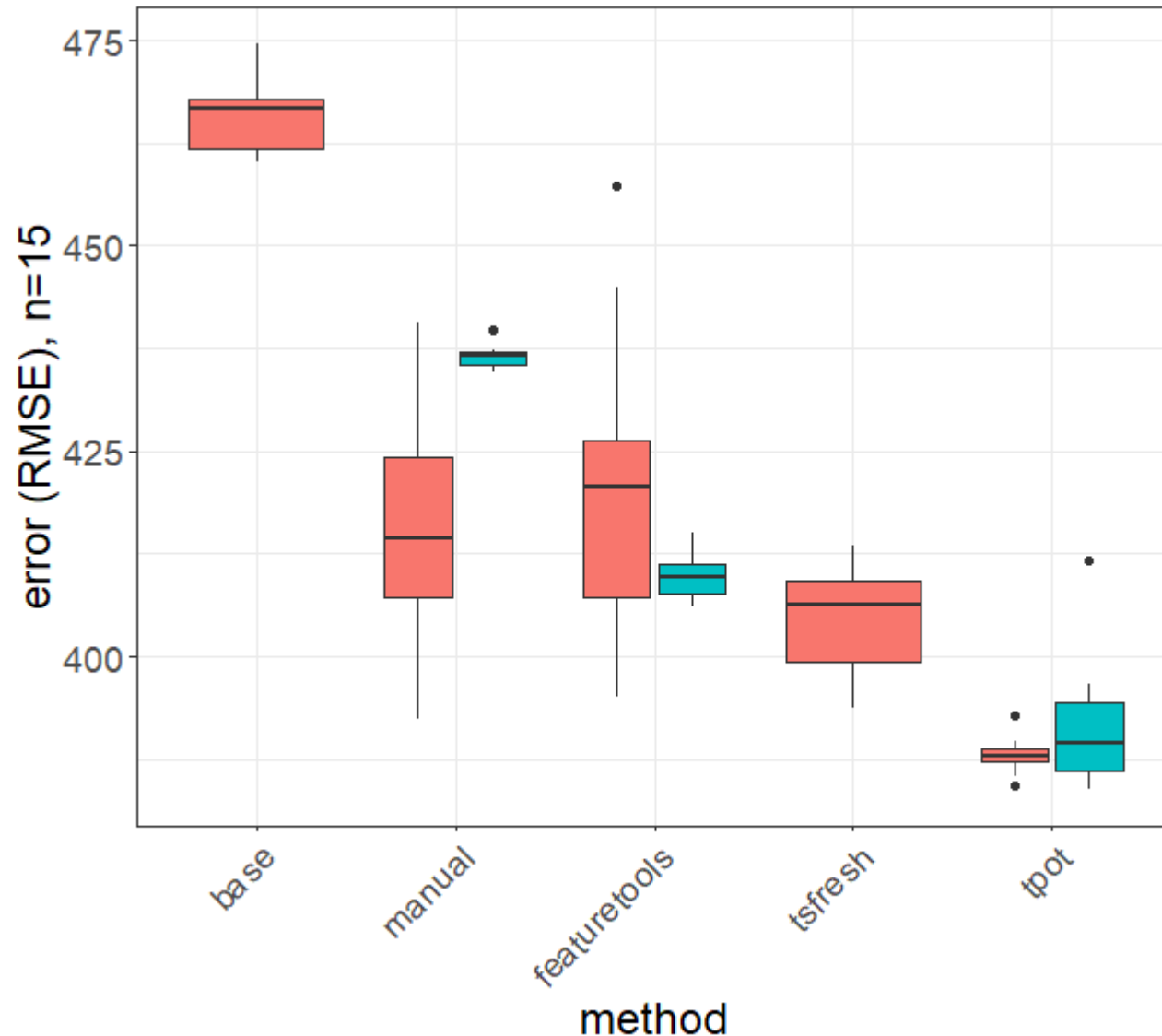


Order of performance:

1. TPOT
2. TSFRESH
3. Featuretools / Manual
5. Base

Adding features helped

Results Performance: selecting features with TSFRESH-method

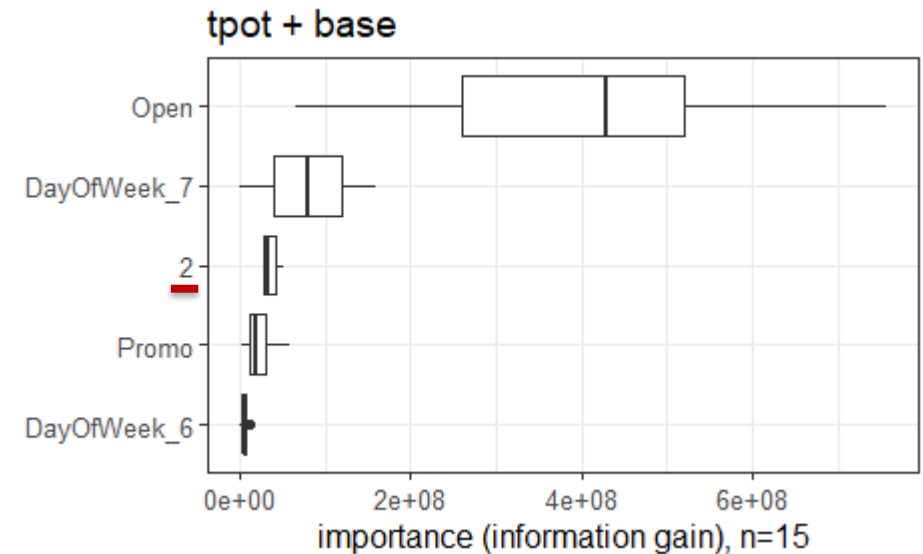
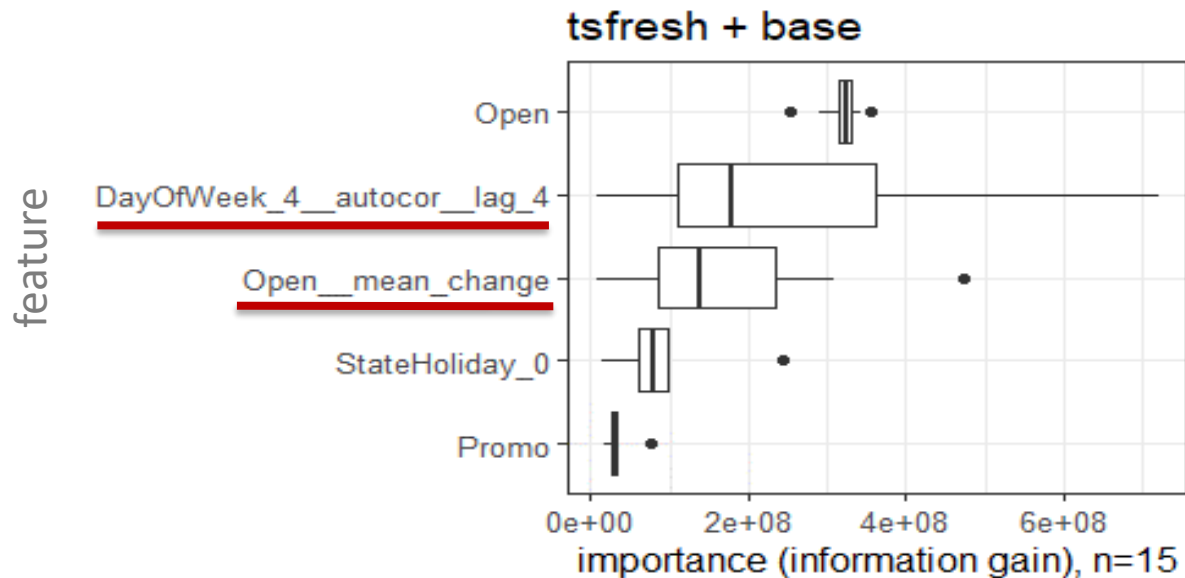
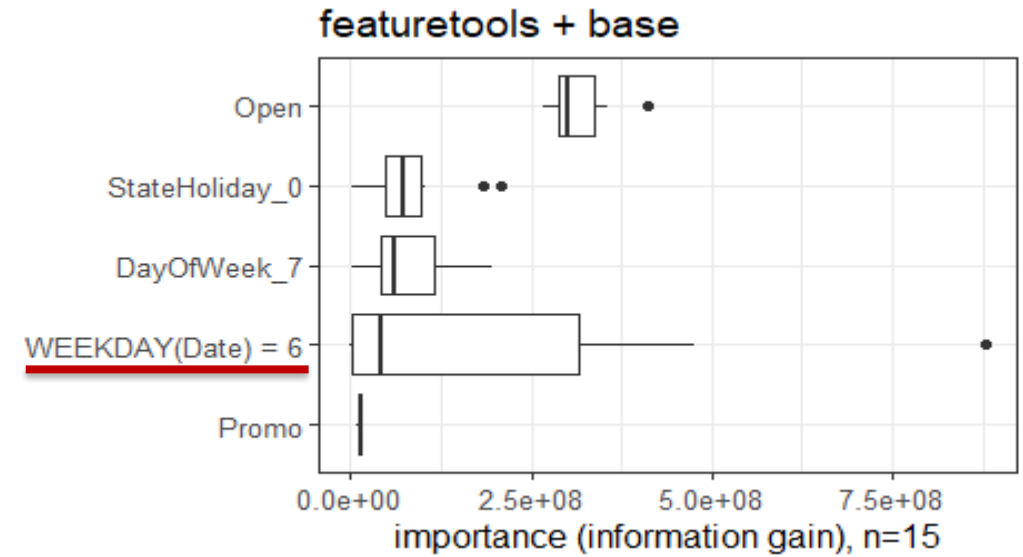
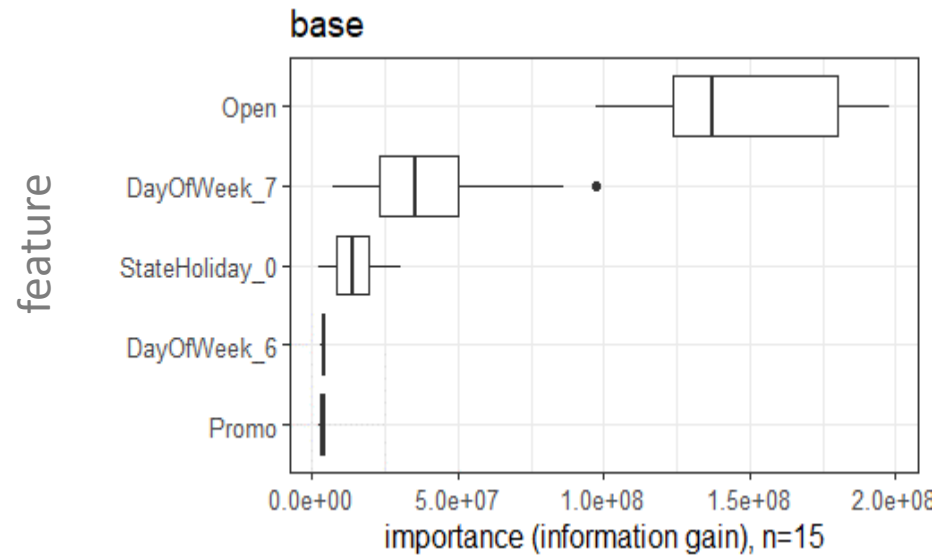


additional selection

- FALSE
- TRUE

- „Filtering“ benefited featuretools
- Filtering did not help others:
 - Manual features are often meaningful
 - TPOT's features selected by optimizer

Comparing engineered features to ,base' features



Comparison of Libraries (at default settings)



TSFRESH



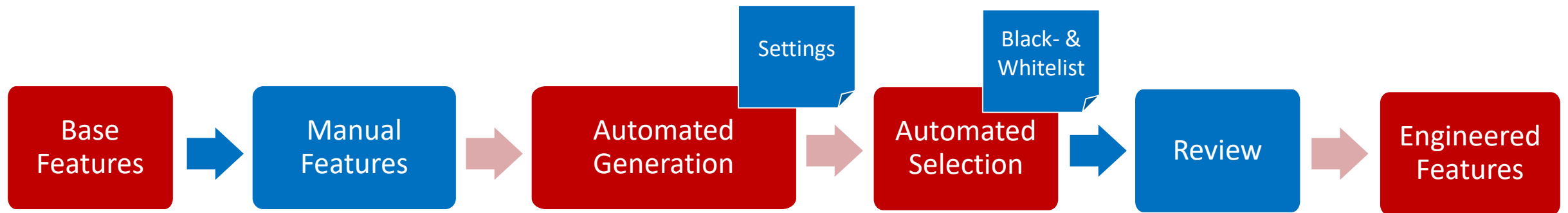
Easy setup (data and config)			
Handle time in data			
Features self-explanatory			
High potential complexity of features			
Minimal „garbage“ features			

Discussion & Conclusion



Conclusion on „Automated Feature Engineering“

- Preparation of data was still required
 - E.g. preventing time-related data leakage
- „Full automation“ possible with default settings, but stays below potential
 - Default settings might be limiting, e.g. featuretools
- Human intervention would be beneficial



Outlook



- Code open to Reviews and Pull Requests:
github.com/informationfabrik/feature-engineering
- Can we find „better“ default parameters for libraries?
- Open Source Libraries for other methods, e.g. Meta Learning?
- Best way to combine manual and automated Feature Engineering?
 - What do you do or recommend?



TAK

DANK U WEL

谢谢

GRACIAS

KÖSZÖNÖM

CHOKRANE

СПАСИБО

TERIMA KASIH

THANK YOU

VIELEN DANK

GRAZIE

DZIĘKUJĘ

MERCI

TESEKKÜR EDERIM

ขอบคุณครับ

TÄNAN

ARIGATÔ

HVALA