# Title for title page

Author 1
Author 2
Author 3

Ladislaus von Bortkiewicz Chair of Statistics
C.A.S.E. – Center for Applied Statistics
and Economics
Humboldt–Universität zu Berlin
http://lvb.wiwi.hu-berlin.de
www.case.hu-berlin.de

# Outline

# Formal Problem Setting

- *training set*: inputs $X = (x_1, \ldots, x_n) \in \mathbb{R}^{n \times d}$ and labels $Y = (y_1, \ldots, y_n) \in \mathbb{R}^n$
- *test set*: inputs $X' = (x'_1, \ldots, x'_t) \in \mathbb{R}^{t \times d}$ without labels

Find a function

$$f : X \to Y \tag{1}$$

s.t. the *test set* labels are predicted as accurately as possible, i.e.

$$f(X') \approx Y' \tag{2}$$

# Outline

1. Introduction    ✓
2. Pre-processing Steps
3. Model Selection
4. Variable Importance and Dimensionality Reduction
5. Results and Conclusion

# Pre-processing

Several transformations and cleaning steps needed before putting the data into an algorithm, e.g.



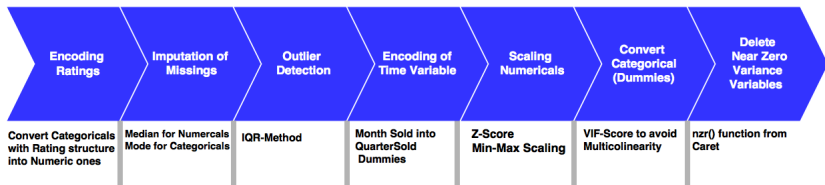| Encoding Ratings | Imputation of Missings | Outlier Detection | Encoding of Time Variable | Scaling Numericals | Convert Categorical (Dummies) | Delete Near Zero Variance Variables |
|---|---|---|---|---|---|---|
| Convert Categoricals with Rating structure into Numeric ones | Median for Numercals Mode for Categoricals | IQR-Method | Month Sold into QuarterSold Dummies | Z-Score Min-Max Scaling | VIF-Score to avoid Multicolinearity | nzr() function from Caret |

Figure 1: Workflow of Pre-Processing Steps

All transformation need to be preformed on the test set as well!

```
 1  basic_preprocessing = function(X_com, y, scaler="gaussian")
 2  {
 3      source("replace_ratings.R")
 4      source("convert_categoricals.R")
 5      source("impute_data.R")
 6      source("encode_time_variables.R")
 7      source("impute_outliers.R")
 8      source("scale_data.R")
 9      source("delete_nearzero_variables.R")
10      X_ratings = replace_ratings(X_com)
11      X_imputed = naive_imputation(X_ratings)
12      X_no_outlier = data.frame(lapply(X_imputed, iqr_outlier))
13      X_time_encoded = include_quarter_dummies(X_no_outlier)
14      X_scaled = scale_data(X_time_encoded, scale_method = scaler)
15      X_encoded = data.frame(lapply(X_scaled, cat_to_dummy))
16      X_com = delect_nz_variable(X_encoded)
17      idx_train = c(1:length(y))
18      train = cbind(X_com[idx_train, ]
19      test = X_com[-idx_train, ]
20      return(list(train = train, X_com = X_com, test = test))
21  }
```

Q dataProcessing

# Outline

1. Introduction   ✓
2. Pre-processing Steps   ✓
3. Model Selection
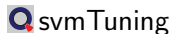4. Variable Importance and Dimensionality Reduction
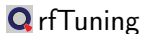5. Results and Conclusion

# Optimizing Hyper-parameters

**Algorithm 1:** t-time k-fold crossvalidation and gridSearch

1  **foreach** *i in 1:t* **do**
2  | Randomly split the data into k folds of the same size
3  | **foreach** *j in 1:k* **do**
4  | | Use *j*th fold as test set and the union of remaining folds as training set
5  | | **foreach** *p in 1:grid* **do**
6  | | | Fit model on training set using parameter set *p*
7  | | | Predict on test set and calculate RMSE
8  | | **end**
9  | **end**
10 | **foreach** *p in 1:grid* **do**
11 | | Calculate average RMSE over the $t \times k$-runs
12 | **end**
13 | choose *p* with the lowest RMSE
14 **end**

**Q** xgbTuning          **Q** rfTuning          **Q** svmTuning

# Taking on the curse of Dimensionality

Problem:

- ⊡ many variables (99 after pre-processing)
- ⊡ small training set ($n = 1460$)
- ⊡ variables are correlated with each other

Our approaches:

- ⊡ Variable selection through variable importance ranking
- ⊡ Extract a smaller set of variable using PCA

# Outline

1. Introduction ✓
2. Pre-processing Steps ✓
3. Model Selection ✓
4. Variable Importance and Dimensionality Reduction ✓
5. Results and Conclusion

## Results

- ⊡ Gaussian SVR with all variable is the single best model
- ⊡ PCA did not work well
- ⊡ Models perform best with the full set of variables as Figure **??** suggested

| Inputs | Gaussian SVR | Random Forest | GBM |
|---|---|---|---|
| All Variables | **0.1308** | 0.1484 | 0.1333 |
| Top 30 | 0.1323 | 0.1515 | 0.1436 |
| PCA | 0.1607 | 0.1657 | 0.1657 |

Table 1: RMSE of submitted predictions

Github: finalModels

# Outline

1. Introduction ✓
2. Pre-processing Steps ✓
3. Model Selection ✓
4. Variable Importance and Dimensionality Reduction ✓
5. Results and Conclusion ✓

# References

📄 Breiman, Leo
*"Random Forest." Machine learning, 45(1), 5-32, (1999)*
available on http://machinelearning202.pbworks.com

📄 Chen, Tianqi, and Carlos Guestrin
*"XGBoost: Reliable Large-scale Tree Boosting System", Proceedings of the 22nd International Conference on Knowledge Discovery and Data Mining Pages 785-794 (2015)*
available on http://learningsys.org

📄 De Cock, Dean
*"Ames, Iowa: Alternative to the Boston housing data as an end of semester regression project" Journal of Statistics Education 19.3 (2011)*
available on https://ww2.amstat.org

# References

📊 Friedman, Jerome H.
*"Greedy function approximation: a gradient boosting machine."*
*Annals of statistics 1189-1232 (2001).*
available on https://www.jstor.org/journal/annalsstatistics

📄 Kuhn, Max, and Kjell Johnson
*"Applied predictive modeling". New York: Springer (2013)*

📊 Vapnik, Vladimir, Steven E. Golowich, and Alex Smola
*"Support vector method for function approximation, regression estimation, and signal processing." Advances in neural information processing systems 281-287 (1997)*
available on https://semanticscholar.org