



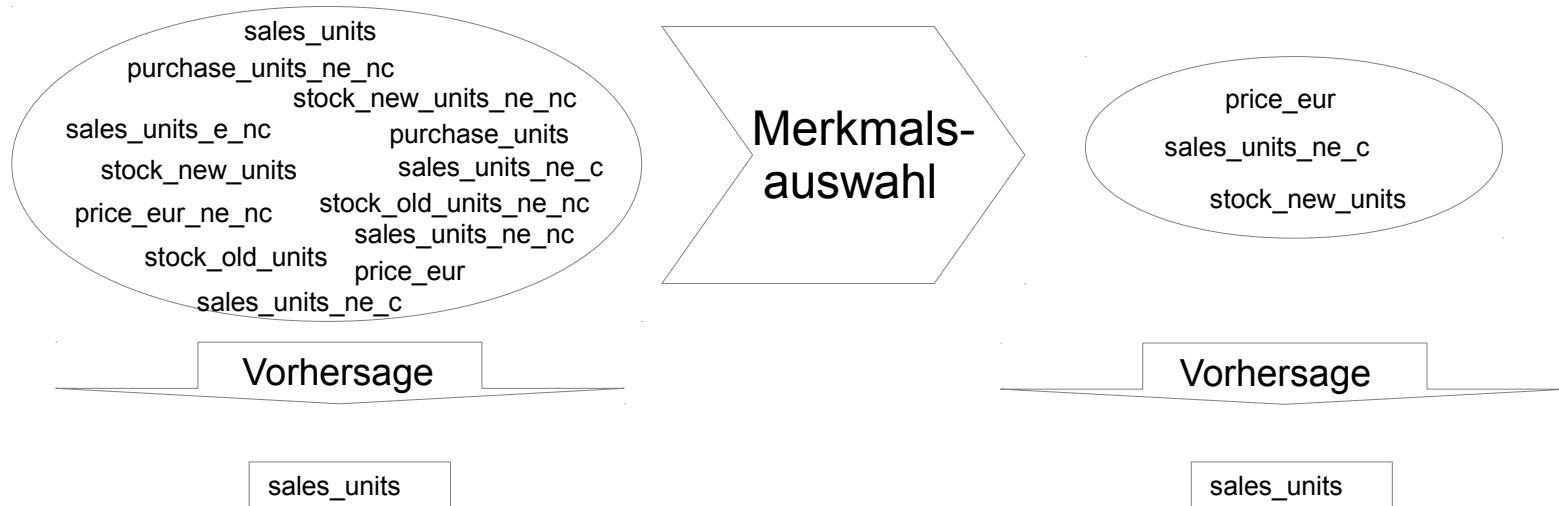
Marcel Spranger

*Merkmalsauswahl zur Optimierung von
Prognoseprozessen auf Verkaufsdaten*

15. Dezember 2014

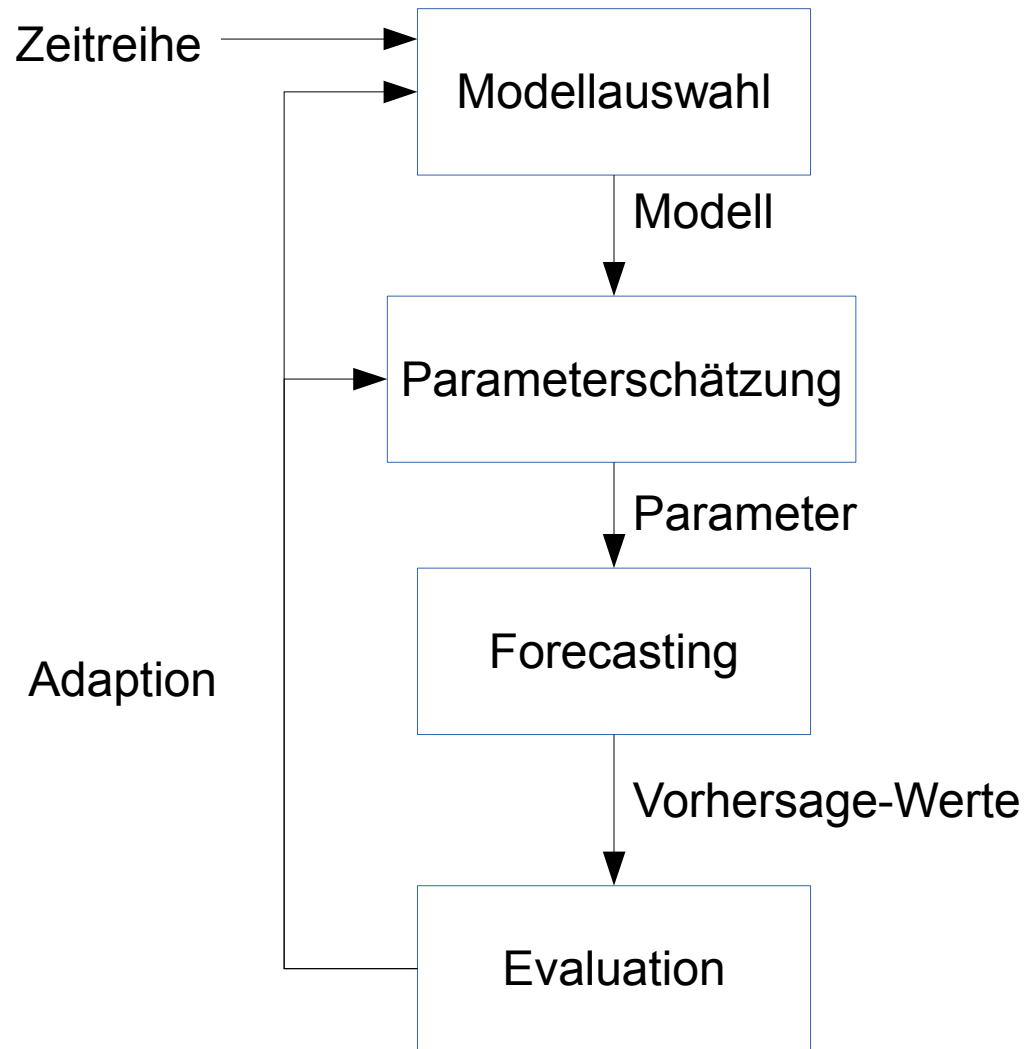


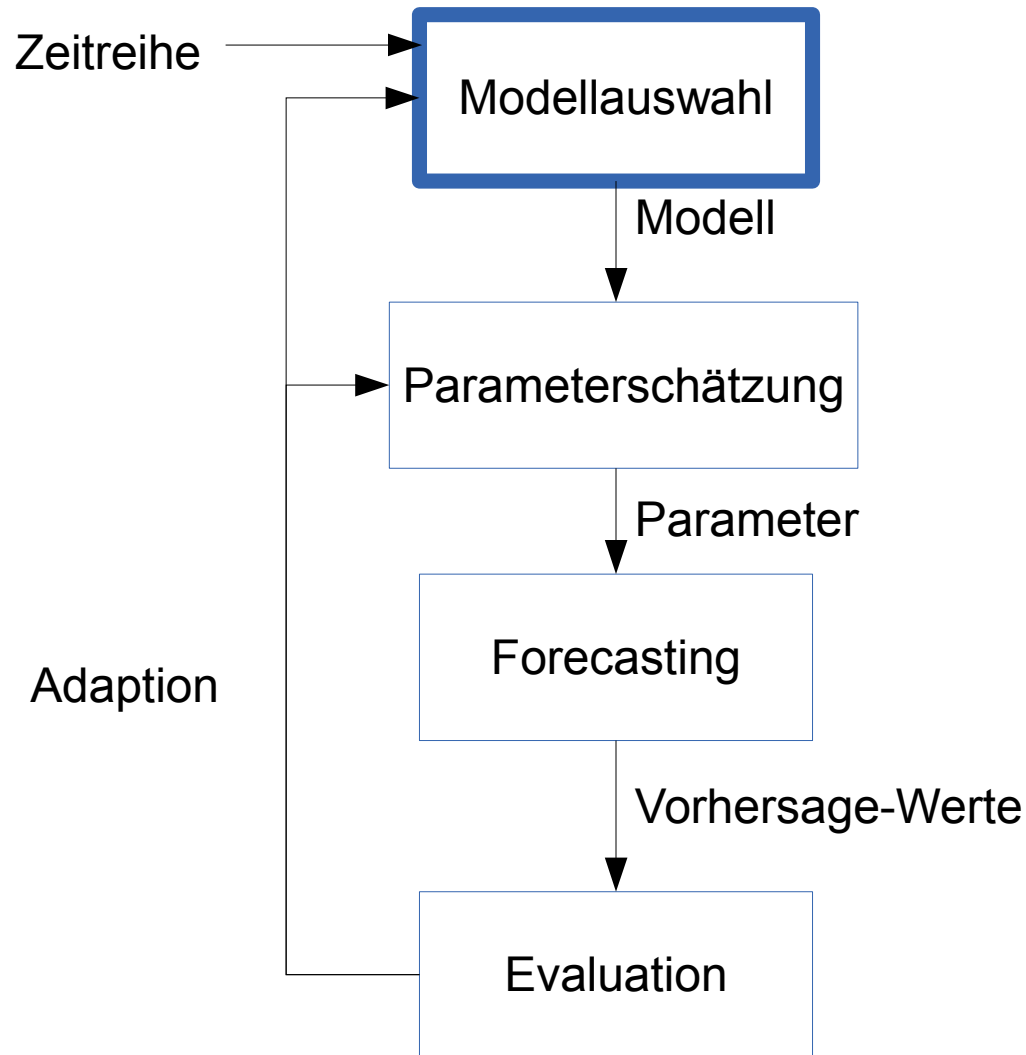
- Zeitreihenvorhersage: Prognose zukünftiger Werte einer Zeitreihe
- Cross-Sectionale Forecasting: Prognose des Zielattributs anhand weiterer Attribute
- Beispiele:
 - Vorhersage von Energieverbrauch und Energieerzeugung
 - Bruttosozialprodukt, Arbeitslosenquote
 - Verkaufszahlen

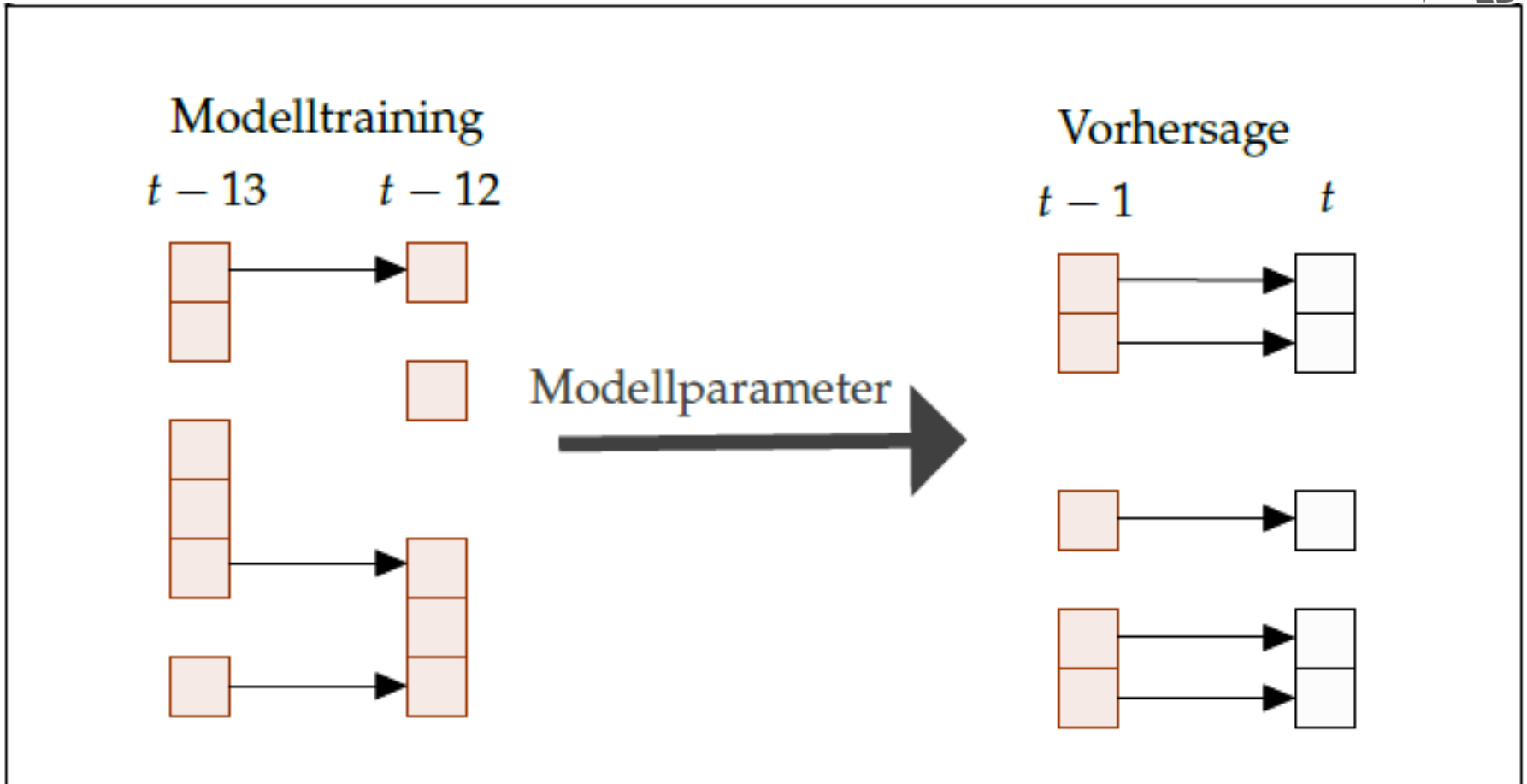




- ❖ Motivation
- ❖ Grundlagen
- ❖ Wrapper - Ansatz
- ❖ Filter - Ansatz
- ❖ Hybrider Ansatz
- ❖ Evaluation
- ❖ Zusammenfassung







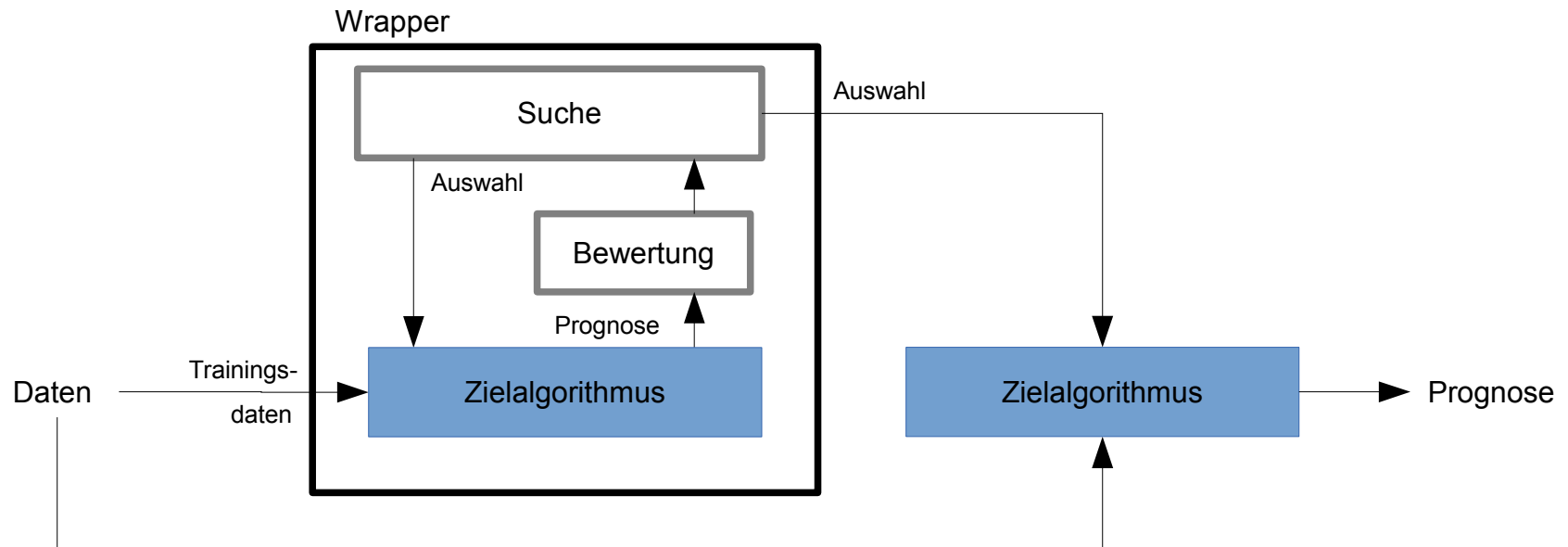
$$y = \beta_0 + \sum_{i=1}^n \beta_i x_i + u$$



- ❖ Motivation
- ❖ Grundlagen
- ❖ **Wrapper - Ansatz**
- ❖ Filter - Ansatz
- ❖ Hybrider Ansatz
- ❖ Evaluation
- ❖ Zusammenfassung

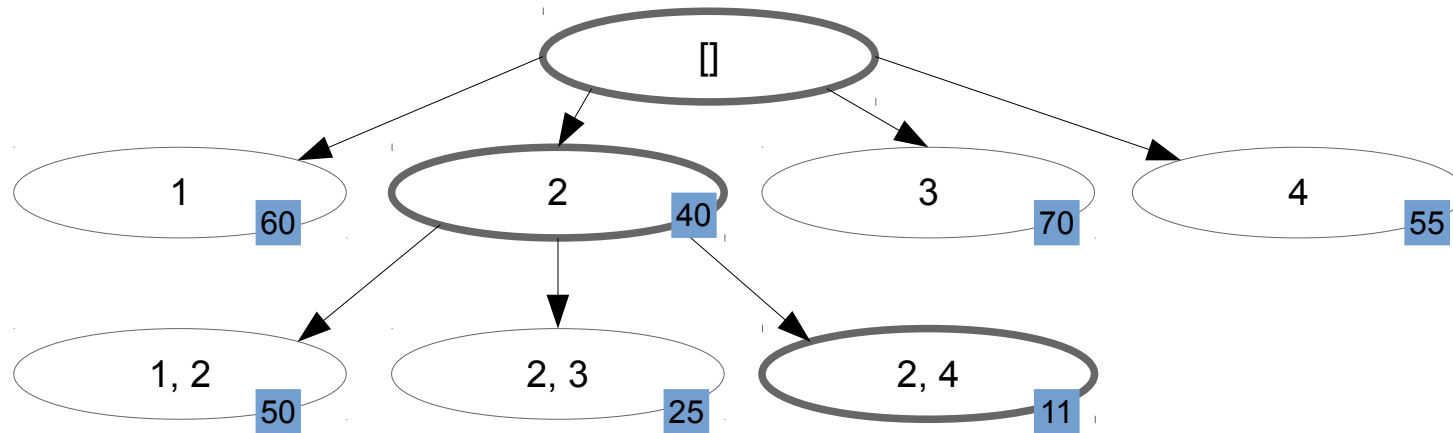


- Idee: Zielalgorithmus als Blackbox zur Bewertung der Auswahl
- Bewertung wird auf den Trainingsdaten ausgeführt
- Suche nach Attributskombination mit bester Bewertung
- Liefert im allgemeinen sehr gute Ergebnisse, da stark an den Zielalgorithmus angepasst
- Jedoch sehr aufwendig





Forward Selection & Backward Elimination



Forward Selection

- Beginne mit leerer Attributskombination
- Überprüfe, welche Attributshinzunahme zur stärksten Verbesserung führt
- Füge so lange Attribute hinzu, bis keine Verbesserung mehr stattfindet

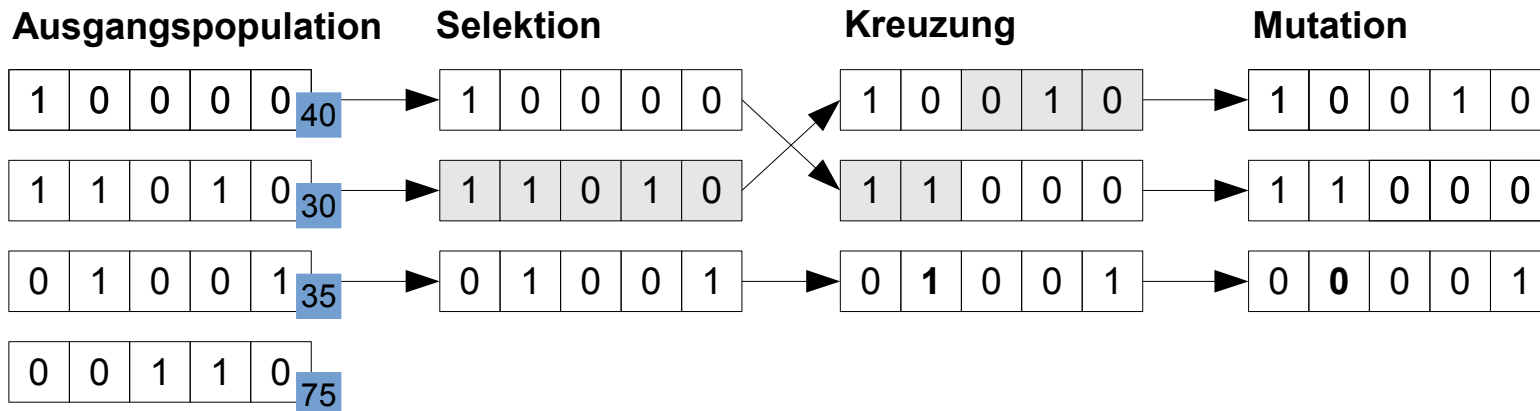
Backward Elimination

- Nehme alle Attribute in Kombination auf
- Überprüfe, welche Attributentfernung zur stärksten Verbesserung führt
- Entferne so lange, bis keine Verbesserung mehr stattfindet



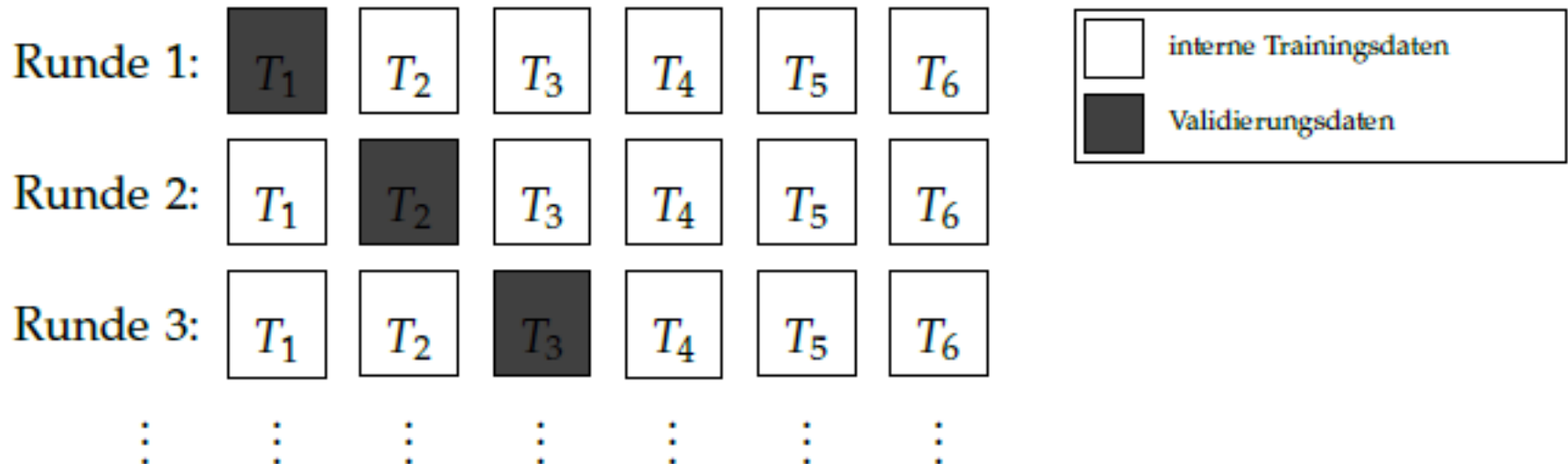
Genetischer Algorithmus

- Parameterbelegung stellt ein Individuum dar
- Durch Kreuzen zweier Individuen und Mutationen entstehen neue Individuen
- Nur Individuen mit ausreichender Bewertung dürfen in der nächsten Iteration neue Nachfahren erzeugen





- Problem: Überfittung
- Trainingsdaten werden mehrmals in interne Trainingsdaten und Validierungsdaten aufgeteilt
- Als Rückgabe dient die durchschnittliche Bewertung





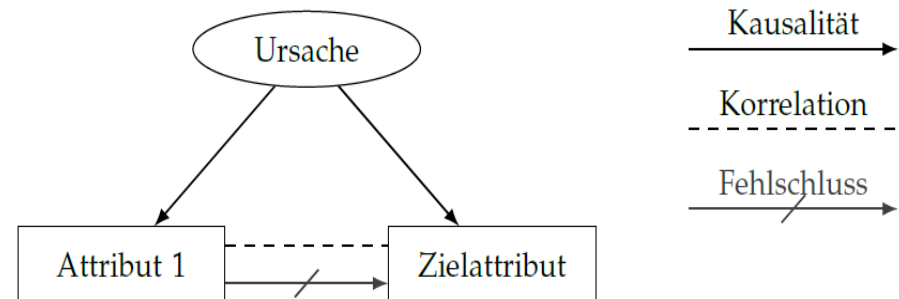
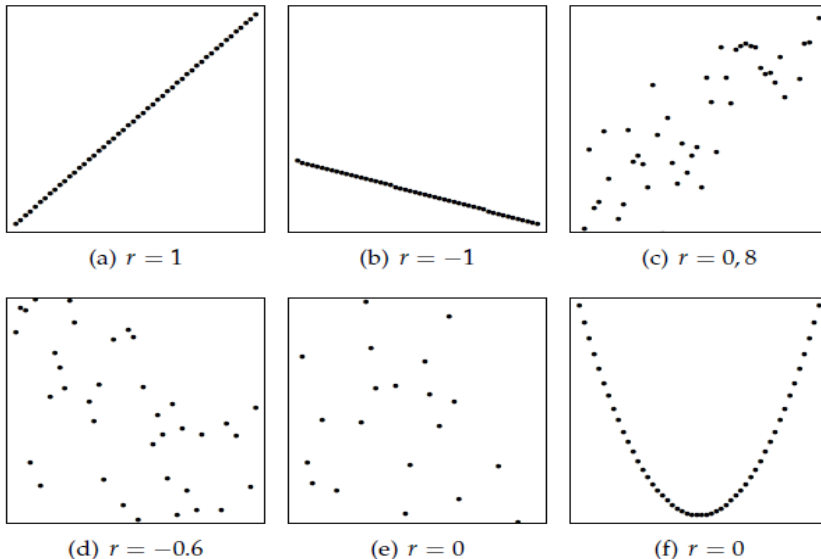
- ❖ Motivation
- ❖ Grundlagen
- ❖ Wrapper - Ansatz
- ❖ Filter - Ansatz**
- ❖ Hybrider Ansatz
- ❖ Evaluation
- ❖ Zusammenfassung



- Idee: die besten Merkmale werden anhand statistischer Eigenschaften herausgefiltert
- Unabhängig von später verwendeten Algorithmus, deswegen weniger genau als Wrapper-Methode
- Dafür weniger Aufwendig, daher schneller

- Korrelationskoeffizient nach Bravais und Pearson

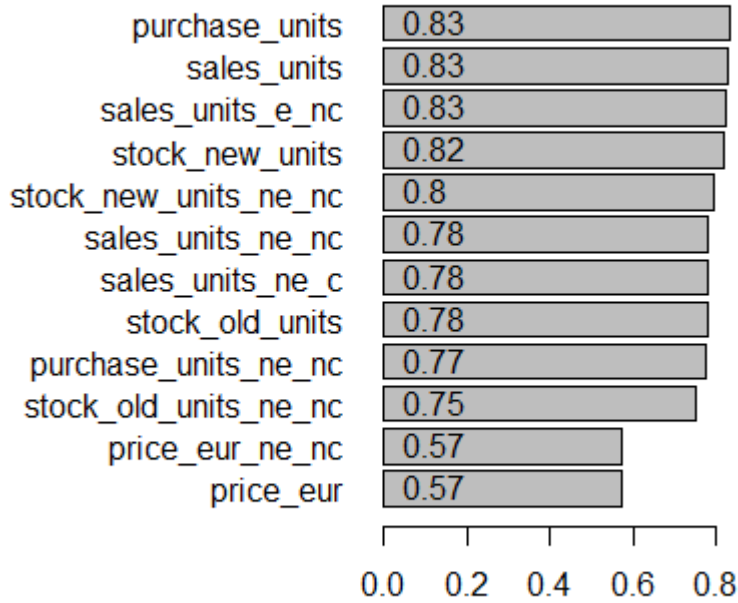
$$r = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2 \sum_{i=1}^n (y_i - \bar{y})^2}}$$



> Korrelationskoeffizient als Filter

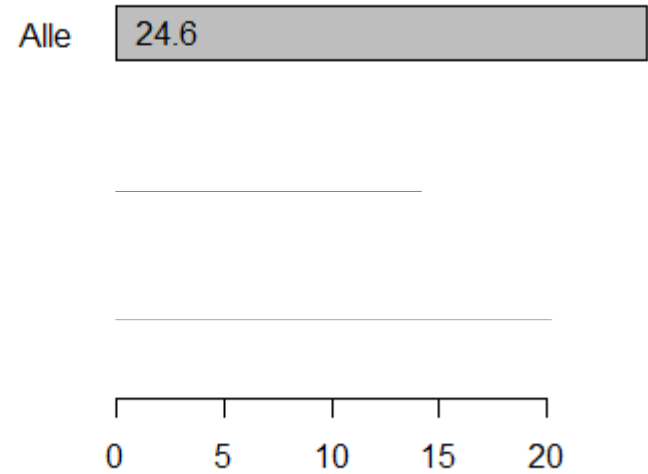


cooling

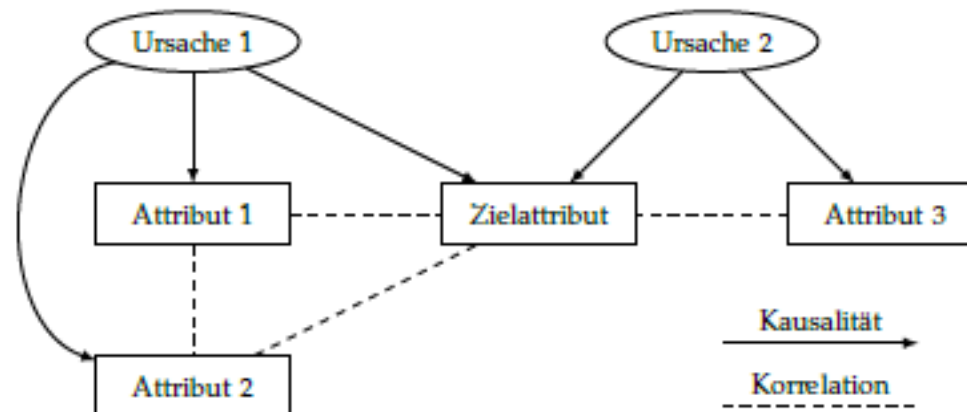


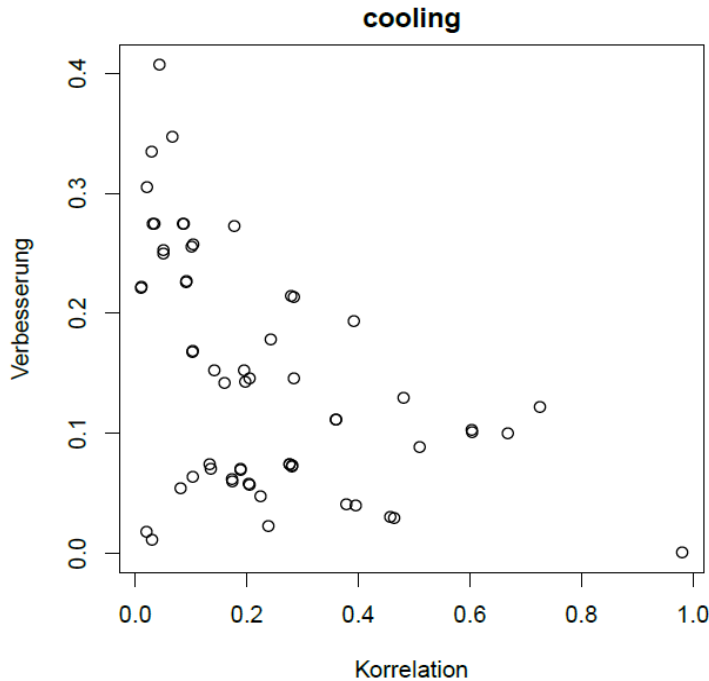
Korrelation

cooling



Fehler





- starke Korrelation → kaum Verbesserung
- Geringe Korrelation
→ starke Verbesserung möglich
- Aber: Geringe Korrelation kein Garant für Verbesserung!

- Korrelationsbasierte Filter nach Hall
- Merkmale sollen hohe Korrelation zu Zielattribut besitzen, jedoch geringe Korrelation untereinander

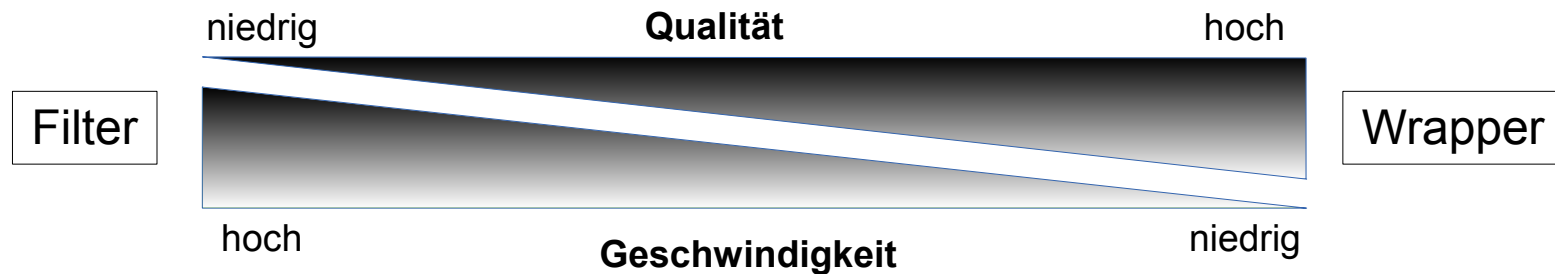
$$Merit_S = \frac{k \bar{r}_{cf}}{\sqrt{k + k(k-1) \bar{r}_{ff}}}$$



- ❖ Motivation
- ❖ Grundlagen
- ❖ Wrapper - Ansatz
- ❖ Filter - Ansatz
- ❖ **Hybrider Ansatz**
- ❖ Evaluation
- ❖ Zusammenfassung



- Filter-Ansatz schnell, aber ungenau
- Wrapper-Ansatz genau, aber langsam



→ Kombination beider Ansätze

Filter-Wrapper

Filter gibt nur bestimmte Anzahl von Kombinationen zurück, auf die der Wrapper-Ansatz durchgeführt wird.

Merit+

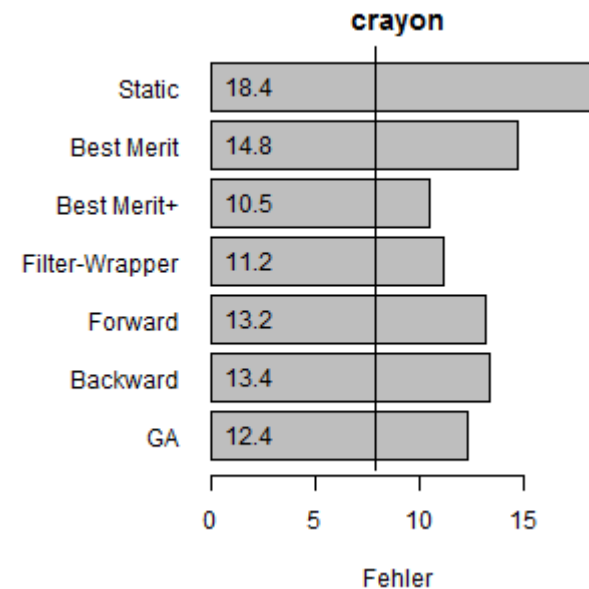
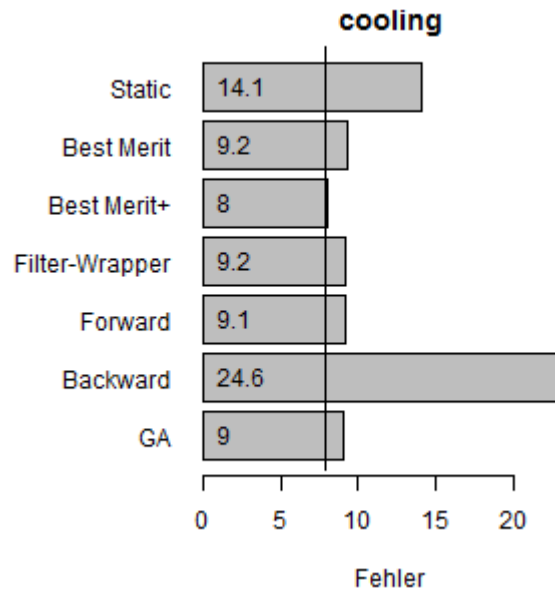
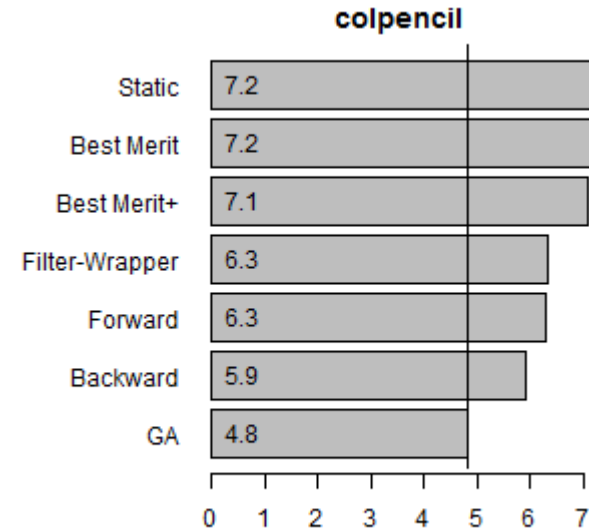
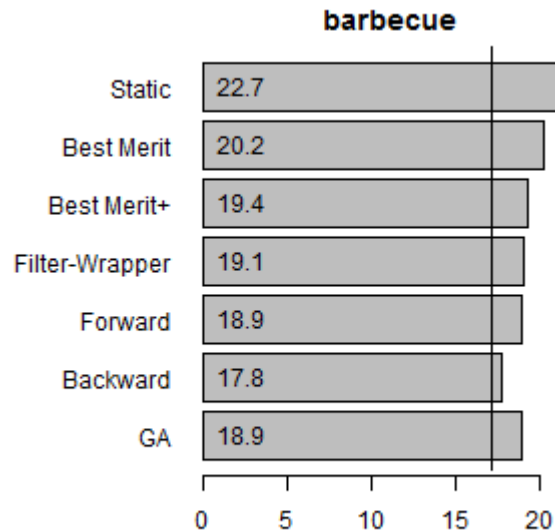
Wrapper-Ansatz zur Korrelationsbestimmung der Einzelattribute.



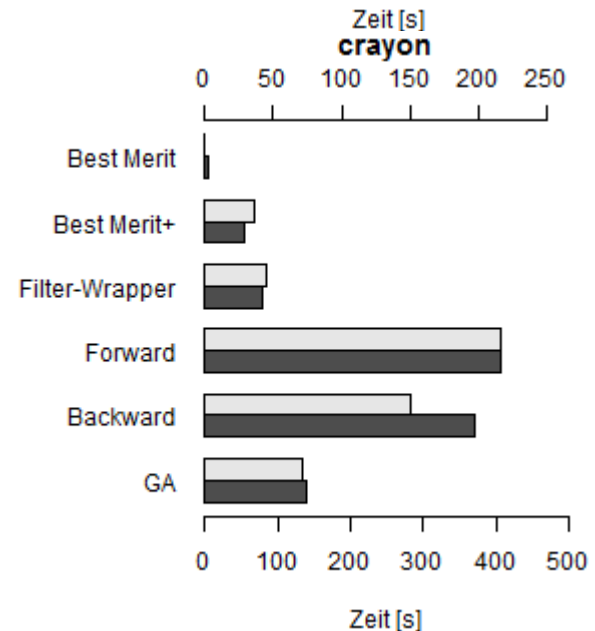
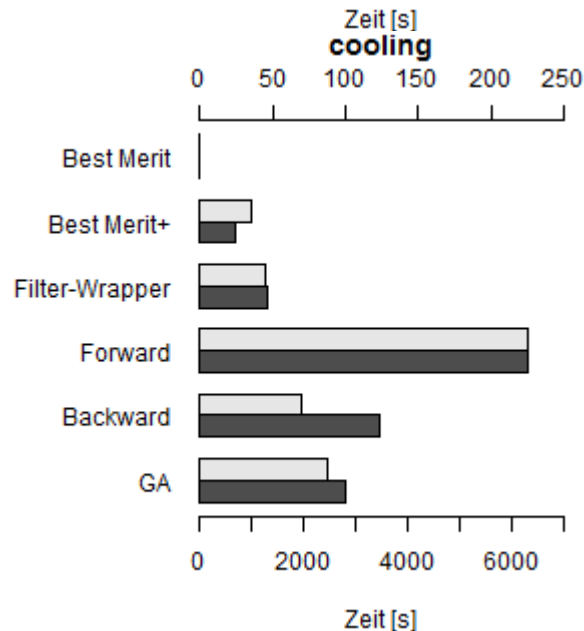
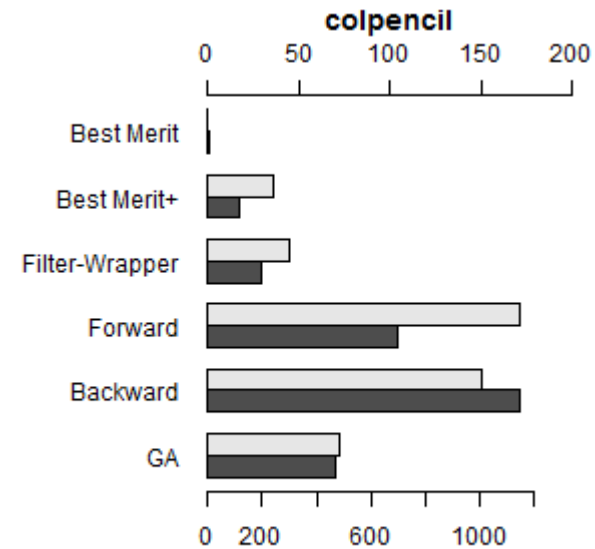
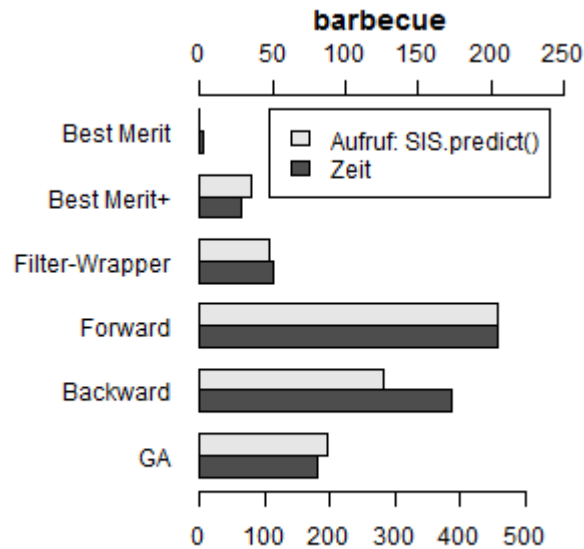
- ❖ Motivation
- ❖ Grundlagen
- ❖ Wrapper - Ansatz
- ❖ Filter - Ansatz
- ❖ Hybrider Ansatz
- ❖ **Evaluation**
- ❖ Zusammenfassung



- Bestimmung der Merkmale mithilfe Filter-, Wrapper- oder Hybriden Ansatz anhand der Trainingsdaten
- Verwendung dieser Merkmale zur Vorhersage der Evaluationsdaten
- Bestimmung der Abweichung/ des Fehlers der entstandenen Vorhersagen
- Methoden:
 - Static: Ziel-Attribut als Regressionsattribut (zum Vergleich)
 - Best Merit
 - Best Merit+
 - Filter-Wrapper: 15 besten Filter-Ergebnisse werden von Wrapper überprüft
 - Forward Selection
 - Backward Elimination
 - GA
 - Bestmögliches Ergebnis (zum Vergleich)



> Evaluation





- Wahl der Modellattribute hat entscheidenden Einfluss auf die Vorhersagequalität
- Automatisierte Merkmalsauswahl ist potentiell in der Lage, Qualität der Vorhersage erheblich zu Verbessern im Vergleich zur statischen Attributsauswahl
- Bis auf Ausnahmen lieferten sowohl Filter-, Wrapper- und kombinierte Ansätze gute bis sehr gute Ergebnisse
- Ein klarer Favorit konnte jedoch nicht bestimmt werden
- Future Work:
 - Nichtlineare Zusammenhänge
 - Hohe Attributanzahl
 - Wann ist welche Methode anzuwenden?



Marcel Spranger

*Merkmalsauswahl zur Optimierung von
Prognoseprozessen auf Verkaufsdaten*

15. Dezember 2014