



Next

- DM und KDD, Phasen
- Aufgabenstellungen des DM
- Wissensrepräsentation
- Entscheidungsbäume I – Repräsentation
- Entscheidungsbäume II – Lernen
- **Entscheidungsbäume III – Praktisch**
 - Ausblick ID3 zu C4.5
 - Ein Problem
 - Tools
 - RapidMiner
- Performance von Klassifikatoren



ID3-Algorithmus

- Wie beschrieben mit Gain als Auswahlkriterium des Attributes
 - Entwickler Ross Quinlan (Sydney), exakt das Beispiel im Paper
 - J. R. Quinlan: **Induction of decision trees**, Machine Learning, vol. 1, 1986.
-
- Kein Backtracking => nicht optimal, lokale Minima mgl.
 - ID3 sehr erfolgreich, Zerlegung anhand des Informationsgewinns einleuchtend
-
- Kleine Probleme in der Praxis, deshalb von Quinlan weiterentwickelt zu
 - **C4.5** (J48 in Weka, DecisionTree in RapidMiner) [Quinlan1993]
 - C5.0 (kommerziell, unveröff. Verfahren, ähnlich C4.5 aber schneller, C-Code single-threaded version under GPL -> R-Package)



Probleme des ID3 – Ein Beispiel

Hinter einem Vorhang verbirgt sich ein Mitarbeiter oder Professor des Fachbereiches. Sie sollen herausbekommen, wer es ist.

Welches Attribut fragen Sie als erstes ab?

Problem des Information Gain (möglichst reine Knoten)

- **Attribute mit vielen Attributwerten** werden **stark bevorzugt**, denn sie enthalten sehr viel Information – bis zur Identifikation
- Andere Maße um das beste Attribut zu bestimmen
 - **Gain ratio** (möglichst reine Knoten, aber nicht zu viele)
 - Gini Index ...



Probleme des ID3

- Attribute mit vielen Werten (Primärschlüssel, ID)
- Numerische Attribute nicht erlaubt
- Fehlende Werte
 - Bei der Baumkonstruktion
 - Beim Klassifizieren
- Overfitting: Trainingsdaten werden korrekt klassifiziert, neue Daten hingegen nicht
 - Prepruning (TDIDT rechtzeitig stoppen) vermeidet Overfitting

Diese Probleme sind gelöst im C4.5

- Operator ‚Decisiontree‘ im RapidMiner
- Details in Kap. 6.1 in [WF01]

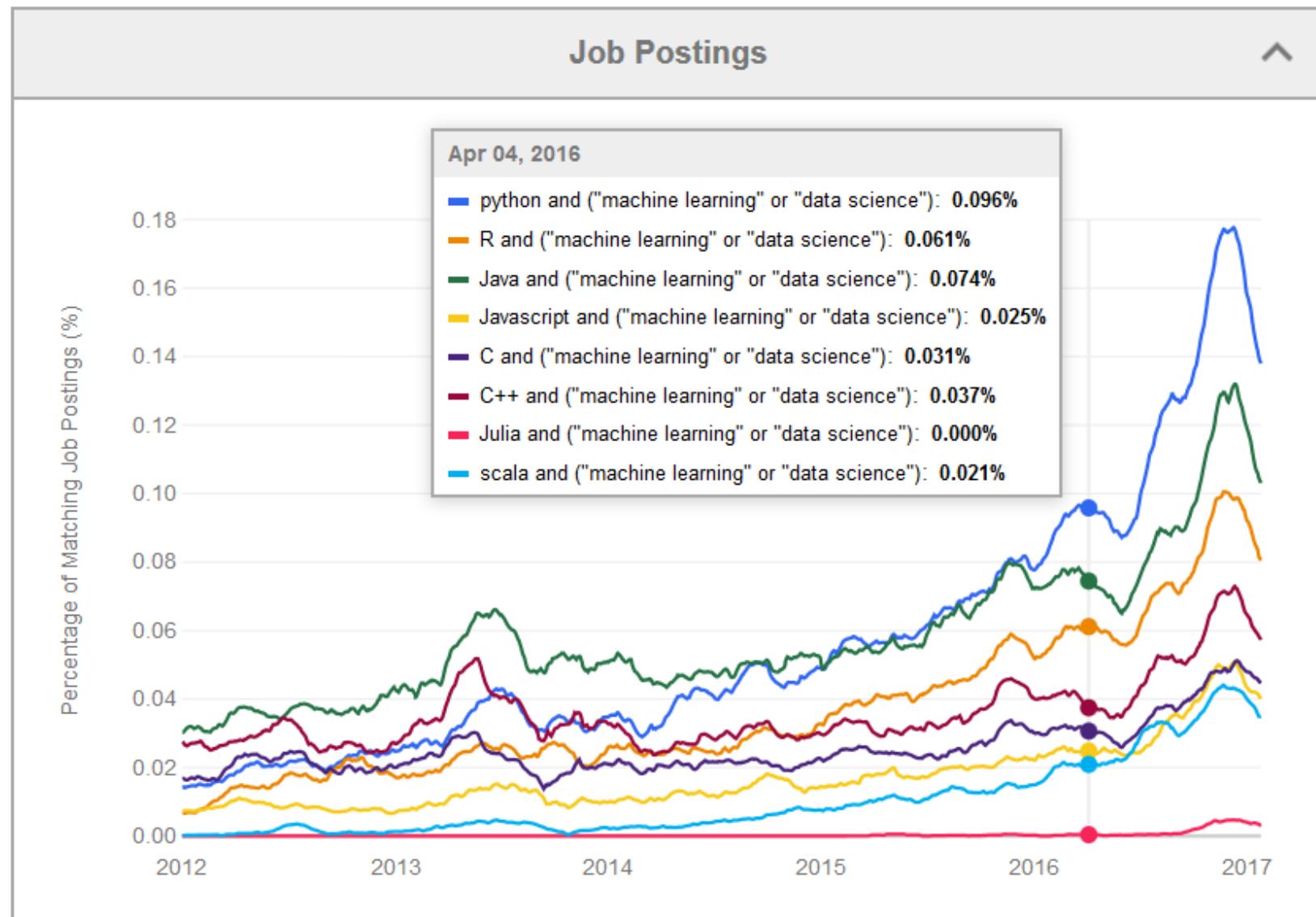


Tools und Sprachen

Kostenfreie Nutzung:

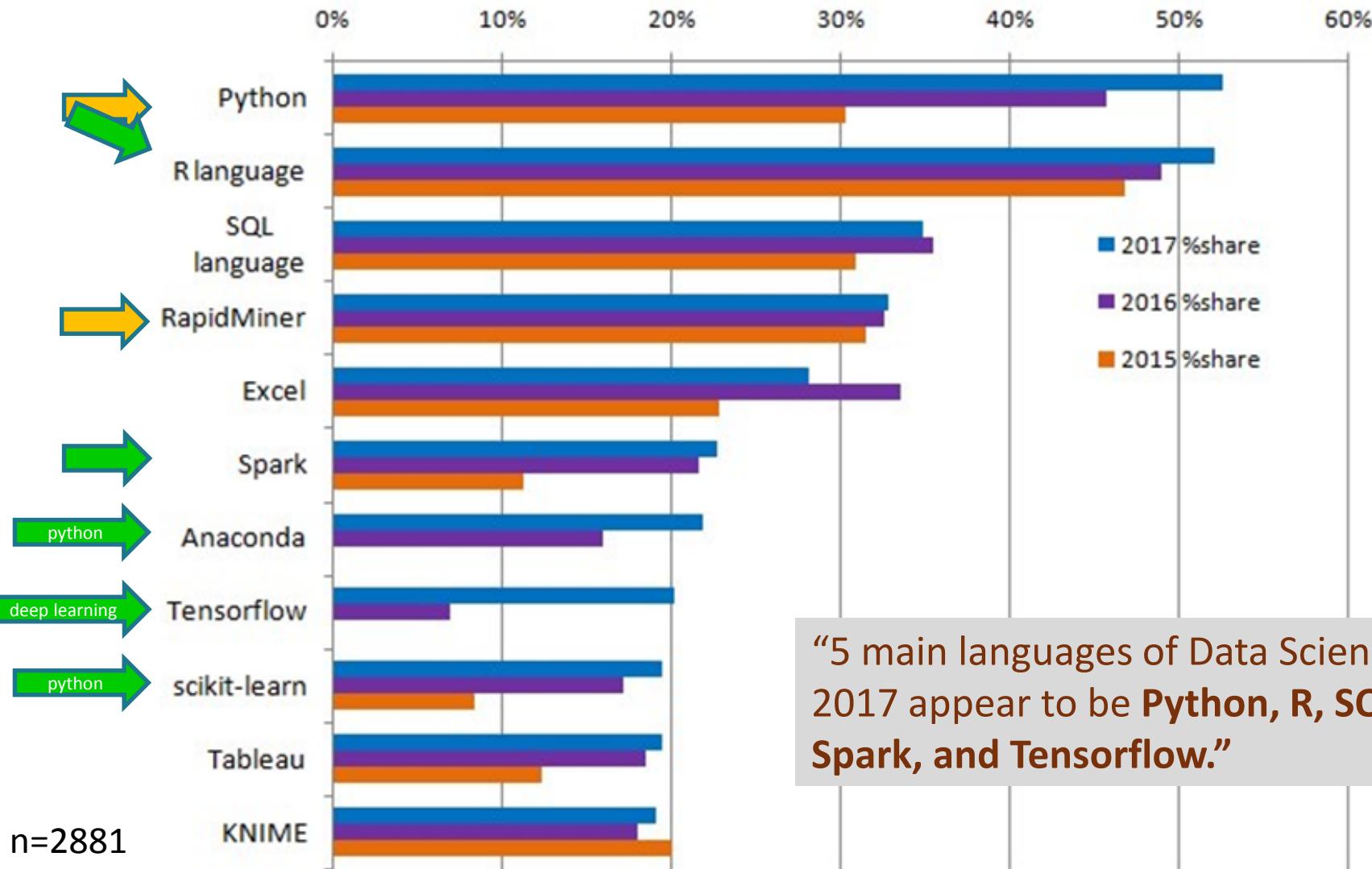
- oft Duale Lizenz: GPL + kommerzielle mit Service und Erweiterungen
- **RapidMiner** Open + Closed Source, Uni. Dortmund + RapidMiner GmbH
- **R** GPL, Programmiersprache
- **Python** GPL, Programmiersprache
- **WEKA** GPL, Universität von Waikato (Neuseeland)
- KNIME GPL + dual, Uni. Konstanz + KNIME.com GmbH (Zürich)
- Orange GPL + Qt, University of Ljubljana (Slowenien)
- Rattle in R, Open Source, Fa. Togaware (Australien)

Kommerziell: bspw. SPSS Modeller



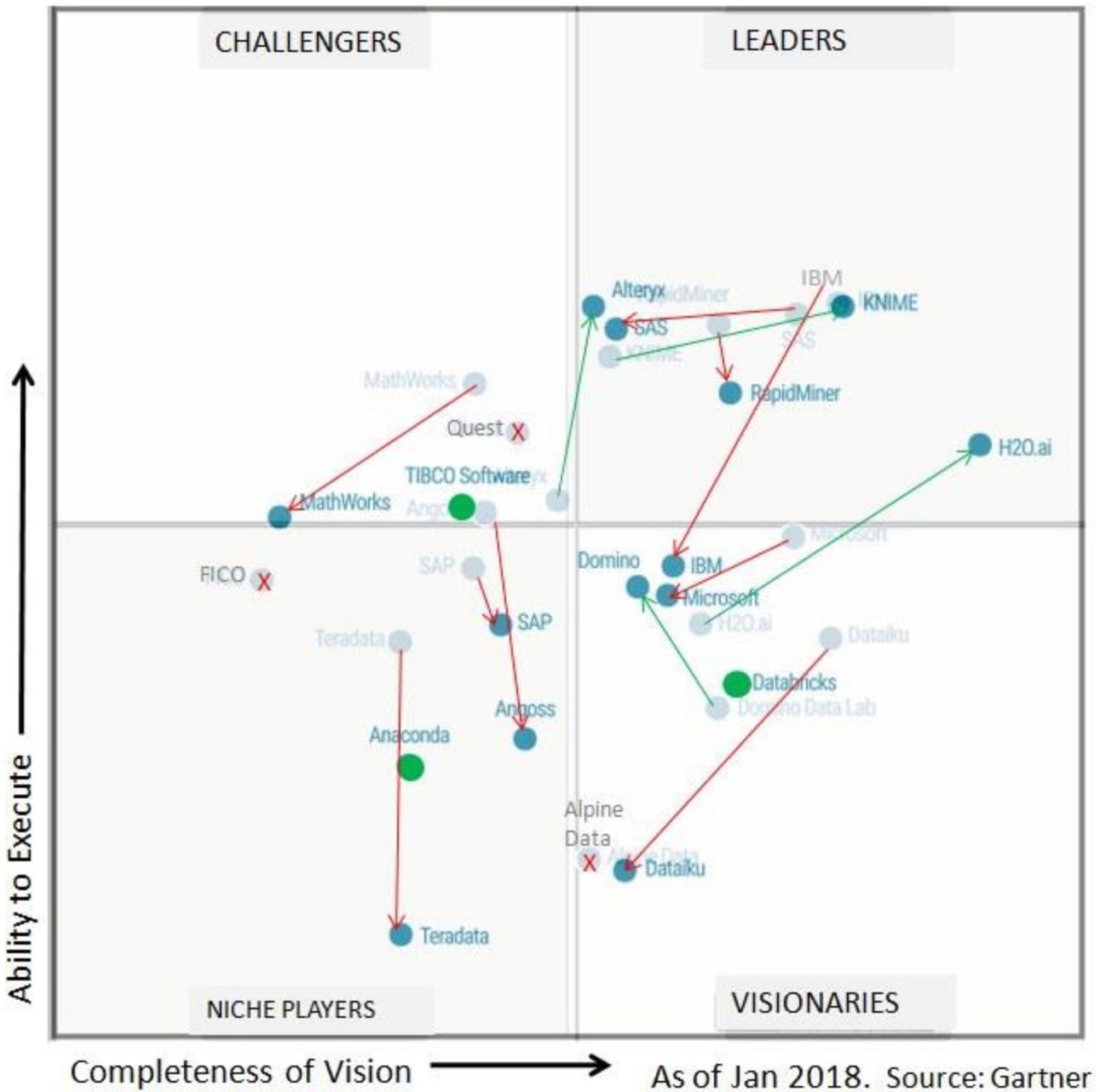
Jean-Francois Puget, <http://www.kdnuggets.com/2017/01/most-popular-language-machine-learning-data-science.html>,
Abfrage auf <https://www.indeed.com/jobtrends>

KDnuggets Analytics, Data Science, Machine Learning Software Poll, top tools share, 2015-2017



<https://www.kdnuggets.com/2017/05/poll-analytics-data-science-machine-learning-software-leaders.html>

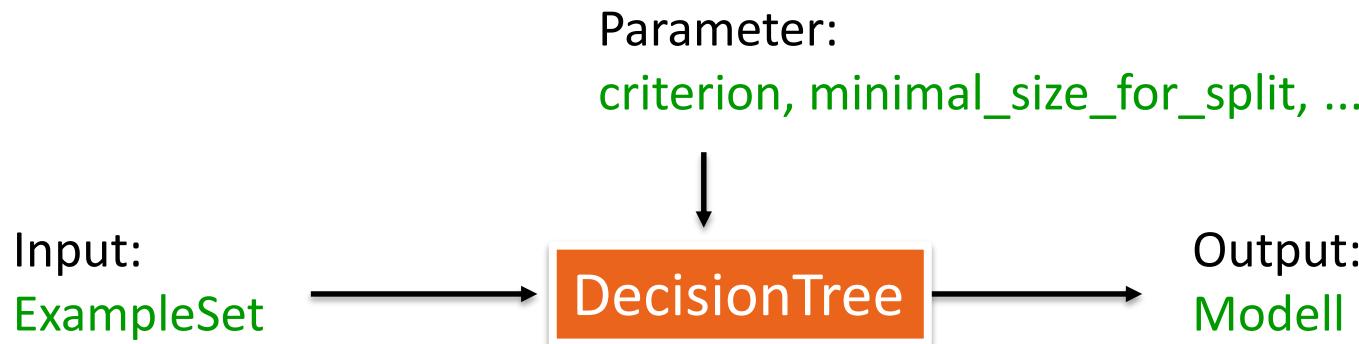






RapidMiner Studio

- Data Mining Tool
- ehemals YALE (Yet Another Learning Environment)
- Open Source (Community Edition = Vorgängerversion), JAVA
- freie akademische Lizenz
- GUI oder offline
- DM-Prozess als Operatorgraph (XML-File)
==> **Visuelles Programmieren**
- Operator: erwartet Input und Parameter, liefert Output





Seminar-Angebote zu RapidMiner, Stand 2018



2 Day Classes

This version consists of 4 hours live online training on two consecutive days (4 hrs. each). The live portion is amended by an additional two times 4 hours of self-study and solving exercises. Depending on your time zone you will receive the live portion in the morning or the after noon of Monday and Tuesday or Wednesday and Thursday, respectively, so that you are left with the afternoons or the mornings of the next day for completion of the exercises. On the following Friday you can optionally attend a live Q&A session.

4 hr live sessions start at:

10am EST/4pm CET/7:30pm IST

[Delivery Datasheet](#)

Course fee is \$1400 or 1250€ per course and student**.

Something urgent came up and you can't attend the full class? Well that will be no problem. You will be provided access to online materials covering all the course content and a Q&A online forum for 60 days starting from your first training day.

RapidMiner Analyst Bootcamp

Are you looking for the fast-track option which offers you a great deal? Well then the 'RapidMiner Analyst Bootcamp' offering is what you are looking for. In this bundle you attend the 2-day classes of 'RapidMiner & DataScience' in one week and then you are granted access to complete the 'RapidMiner Analyst Certification' exam during the following 60 days.

2hr live sessions start at:

10am EST/4pm CET/7:30pm IST

[Analyst Bootcamp Datasheet](#)

The bundle which includes both 'RapidMiner & DataScience: Foundations', 'RapidMiner & DataScience: Advanced' and the 'RapidMiner Analyst Certification' exam is \$2800 or 2500€ per student**.

VILT Classes Schedule

Delivery Style	Course	Dates	Register
2 Day Class	Text & Web Mining with RapidMiner	Apr 16 & 17	in USD / EUR
2 Day Class	RapidMiner & DataScience: Foundations	May 14 & 15	in USD / EUR
2 Day Class	RapidMiner & DataScience: Advanced	May 16 & 17	in USD / EUR
Analyst Bootcamp Bundle	RapidMiner & DataScience + Certification	May 14 - 18	in USD / EUR
2 Day Class	Text & Web Mining with RapidMiner	on request	Waiting List
2 Day Class	RapidMiner Server: Deployment and Web Apps	on request	Waiting List
2 Day Class	RapidMiner & DataScience: Foundations	Jun 18 & 19	Register
2 Day Class	RapidMiner & DataScience: Advanced	Jun 20 & 21	Register
2 Day Class	RapidMiner & DataScience: Foundations	Jul 23 & 24	Register
2 Day Class	RapidMiner & DataScience: Advanced	Jul 25 & 26	Register
2 Day Class	RapidMiner & DataScience: Foundations	Aug 20 & 21	Register
2 Day Class	RapidMiner & DataScience: Advanced	Aug 22 & 23	Register
2 Day Class	RapidMiner & DataScience: Foundations	Sep 17 & 18	Register
2 Day Class	RapidMiner & DataScience: Advanced	Sep 19 & 20	Register



** If you are working for an official RapidMiner partner, in the academic field or if you are based in an emerging markets country (GDP (PPP) per capita < 25k USD as defined by IMF) you might be entitled for a discount. Please reach out to apply: [Contact Us](#)

None of the dates fits you calendar? [Leave us your contact details](#) and we will keep you update on next available dates!



Vorgängerversion als Open Source

<https://github.com/rapidminer>

```
19 package com.rapidminer.operator.learner.tree;
20
21 import com.rapidminer.example.ExampleSet;
22 import com.rapidminer.operator.OperatorCapability;
23 import com.rapidminer.operator.OperatorDescription;
24 import com.rapidminer.operator.OperatorException;
25 import com.rapidminer.parameter.ParameterType;
26
27 import java.util.LinkedList;
28 import java.util.List;
29
30
31 /**
32  * This operator learns decision trees without pruning using nominal
33  * attributes only. Decision trees
34  * are powerful classification methods which often can also easily be
35  * understood. This decision tree
36  * learner works similar to Quinlan's ID3.
37  *
38  * @author Ingo Mierswa
39  */
40 public class ID3Learner extends AbstractTreeLearner {
41     public ID3Learner(OperatorDescription description) {
42         super(description);
43     }
44 }
```



Splitkriterien in Entscheidungsbäumen:

AbstractTreeLearner.java

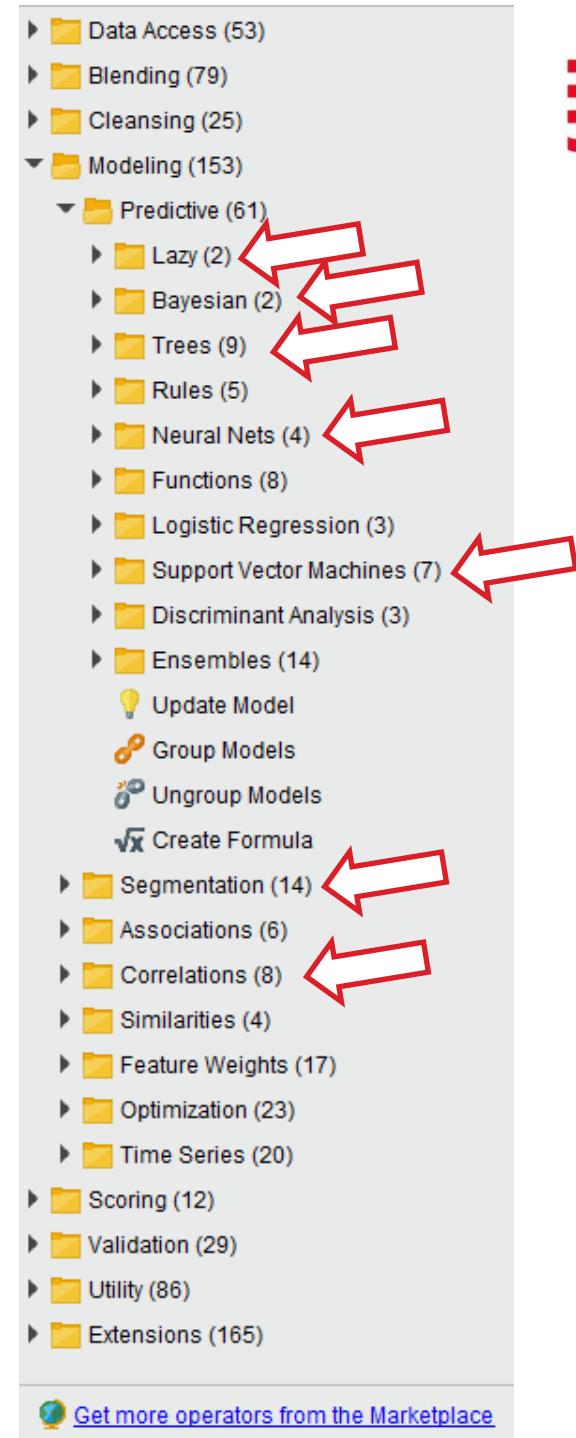
```
23 import com.rapidminer.example.Attribute;
24 import com.rapidminer.example.ExampleSet;
25 import com.rapidminer.example.Statistics;
26 import com.rapidminer.operator.Model;
27 import com.rapidminer.operator.OperatorDescription;
28 import com.rapidminer.operator.OperatorException;
29 import com.rapidminer.operator.UserError;
30 import com.rapidminer.operator.learner.AbstractLearner;
31 import com.rapidminer.operator.learner.PredictionModel;
32 import com.rapidminer.operator.learner.tree.criterions.AbstractCriterion;
33 import com.rapidminer.operator.learner.tree.criterions.AccuracyCriterion;
34 import com.rapidminer.operator.learner.tree.criterions.Criterion;
35 import com.rapidminer.operator.learner.tree.criterions.GainRatioCriterion;
36 import com.rapidminer.operator.learner.tree.criterions.GiniIndexCriterion;
37 import com.rapidminer.operator.learner.tree.criterions.InfoGainCriterion; InfoGainCriterion;
38 import com.rapidminer.parameter.ParameterType;
39 import com.rapidminer.parameter.ParameterTypeDouble;
40 import com.rapidminer.parameter.ParameterTypeInt;
41 import com.rapidminer.parameter.ParameterTypeStringCategory;
42
```



RapidMiner Operatoren

- Data preprocessing (Blending, Cleansing)
- Feature operators
- Machine learning algorithms für numerische **Prognose und Klassifikation**: **SVM, Decisiontrees, Lazy, Bayes, Rules, KNN, k-NN, ..., Metalernen**
- Clustern
- Performance evaluation
- Visualization
- Extensions: Process Mining, Text Processing ...
- In- and output

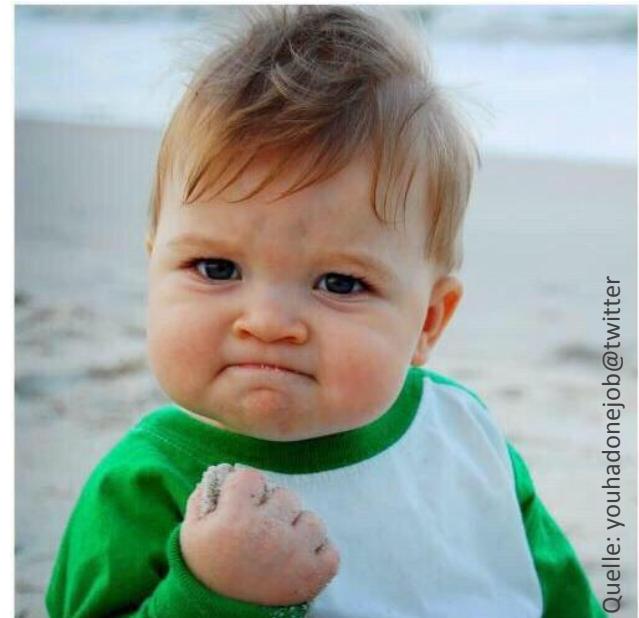
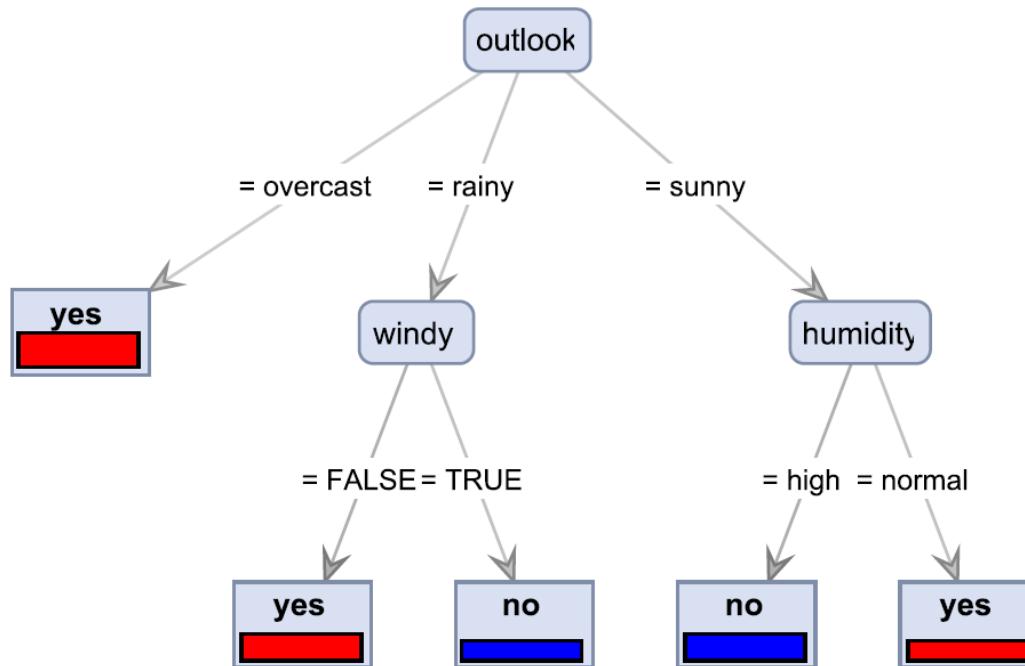
übersichtlich im PDF
„RapidMiner 8 Operator Reference Manual“





ID3 in RapidMiner ...

... erzeugt genau den uns schon bekannten Baum.



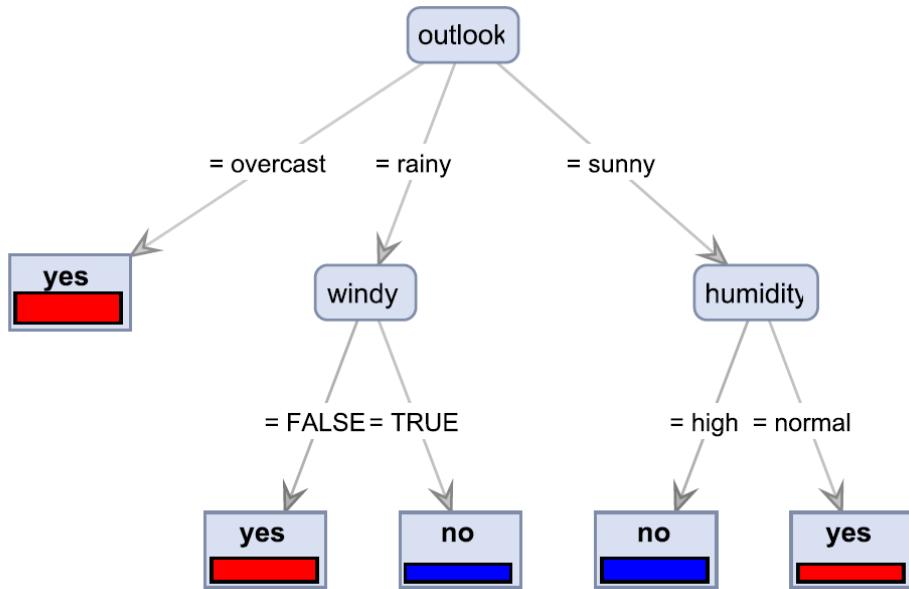
Quelle: youhaddonjob@twitter



Wir kennen nun: DT, DT
anwenden, DT erstellen, Tool
-> Applikation Spongebob: magic



Baum als Text



Tree

outlook = overcast: yes
outlook = rainy
| windy = FALSE: yes
| windy = TRUE: no
outlook = sunny
| humidity = high: no
| humidity = normal: yes



Next

- DM und KDD, Phasen
- Aufgabenstellungen des DM
- Wissensrepräsentation
- Entscheidungsbäume I – Repräsentation
- Entscheidungsbäume II – Lernen
- Entscheidungsbäume III – Praktisch
 - **Erinnerung an die Übung: RapidMiner, Scatterplot, ID3, Decisiontree**
- **Performance von Klassifikatoren**
- Ethik



In der Übung:

- Protokoll
- Datenmenge öffnen, visualisieren, Streudiagramm
- Baum mit ID3 aus Datenmenge erzeugen
- Modell anwenden auf unbekannten Datensatz
- Entscheidungsbaum der Tiere – Tiefe 20, Tiefe 4

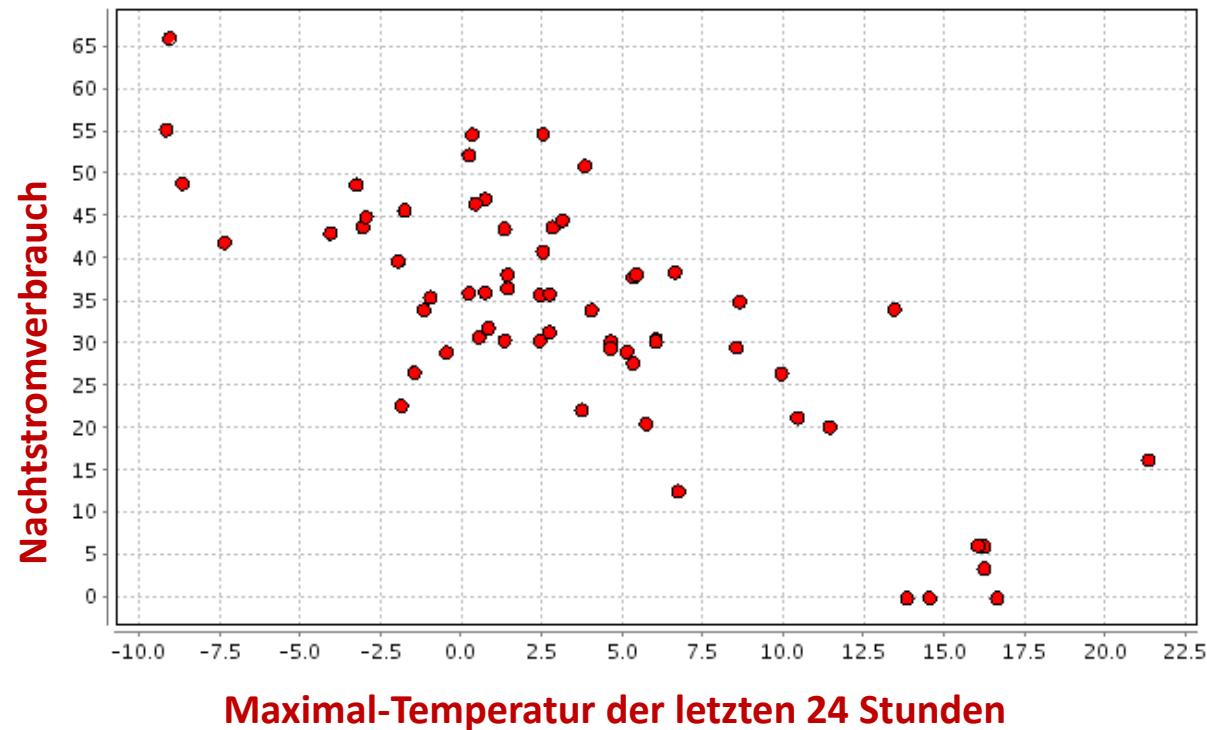
Streudiagramm (Scatterplot)

- 2 Variablen als Achsen, ein Punkt pro Datensatz
- Überdeckung von Punkten durch **Jitter** vermeidbar
- Explorative Analyse: **bivariate Korrelationen**, Cluster, Kompaktheit, Streuung

Beispiel:

Wovon hängt der
Nachtstrom-Verbrauch
ab?

Ursache oder nur
Korrelation?





Business Understanding	Data Understanding	Data Preparation	Modeling	Evaluation	Deployment
Determine Business Objectives <i>Background Business Objectives</i> <i>Business Success Criteria</i>	Collect Initial Data <i>Initial Data Collection Report</i>	Select Data <i>Rationale for Inclusion/Exclusion</i>	Select Modeling Techniques <i>Modeling Technique</i> <i>Modeling Assumptions</i>	Evaluate Results <i>Assessment of Data Mining Results w.r.t. Business Success Criteria</i> <i>Approved Models</i>	Plan Deployment <i>Deployment Plan</i>
Assess Situation <i>Inventory of Resources Requirements, Assumptions, and Constraints</i> <i>Risks and Contingencies</i> <i>Terminology</i> <i>Costs and Benefits</i>	Describe Data <i>Data Description Report</i>	Clean Data <i>Data Cleaning Report</i>	Generate Test Design <i>Test Design</i>	Review Process <i>Review of Process</i>	Plan Monitoring and Maintenance <i>Monitoring and Maintenance Plan</i>
Determine Data Mining Goals <i>Data Mining Goals</i> <i>Data Mining Success Criteria</i>	Explore Data <i>Data Exploration Report</i>	Construct Data <i>Derived Attributes</i> <i>Generated Records</i>	Build Model <i>Parameter Settings</i> <i>Models</i> <i>Model Descriptions</i>	Determine Next Steps <i>List of Possible Actions</i> <i>Decision</i>	Produce Final Report <i>Final Report</i> <i>Final Presentation</i>
Produce Project Plan <i>Project Plan</i> <i>Initial Assessment of Tools and Techniques</i>	Verify Data Quality <i>Data Quality Report</i>	Integrate Data <i>Merged Data</i>	Format Data <i>Reformatted Data</i>	Assess Model <i>Model Assessment</i> <i>Revised Parameter Settings</i>	Review Project <i>Experience Documentation</i>

Figure 3: Generic tasks (**bold**) and outputs (*italic*) of the CRISP-DM reference model

LV Data Mining 2016/Mein erster ID3* – RapidMiner Studio Academia 7.0.001 @ boersch2

Design-, Resultview

File Edit Process View Connections Cloud Settings Extensions

Run

Views: Design Results

Questions?

Repository

- Samples
- DB
- LV Data Mining 2015 (Ingo Boersch)
- LV Data Mining 2016 (Ingo Boersch) **Selected**
- Local Repository (Ingo Boersch)
- Cloud Repository (disconnected)

Operators

tree

- Modeling (10)
 - Predictive (9)
 - Trees (8)
 - Decision Tree
 - Decision Tree (Multiway)
 - Decision Tree (Weight-Based)
 - ID3 **Selected**
 - CHAID
 - Decision Stump
 - Random Tree
 - RF
- Rules
- Feature Weights (1)
- Weight by Tree Importance

+ Get More Operators

Process

Process

```

graph LR
    ReadCSV1[Read CSV] --> SetRole[Set Role]
    SetRole --> ID3[ID3]
    ID3 --> WriteExcel[Write Excel]
    ReadCSV2[Read CSV (2)] --> ApplyModel[Apply Model]
    ApplyModel --> WriteExcel
  
```

Parameters

Write Excel

excel file: Please set this parameter

file format: xlsx

sheet name: RapidMiner Data

date format: I-dd HH:mm:ss

number format: #.0

Parameter, Hilfe des ausgewählten Operators

Required parameter missing
Click on Write Excel to display its parameters; supply a value for excel_file.

Got it! Run anyway (F11)

Prozess-Design, (XML-Ansicht über Menü->View -> Show Panel -> XML)

Leverage the Wisdom of Crowds to get better results. Activate Wisdom of Crowds

Activate Wisdom of Crowds

Help

Write Excel
RapidMiner Studio Core

Synopsis
This operator writes an ExampleSet to a Excel spreadsheet file.
[Jump to Tutorial Process](#)

Einlesen des CSV-Files

File Edit Process View Connections Cloud Settings Extensions //LV Data Mining 2016/Mein erster ID3* – RapidMiner Studio Academia 7.0.001 @ boersch2

Repository Process Parameters Help

1 2 3 4 5 6

1. Open the Operators palette and expand the CSV section.

2. Double-click the Read CSV operator.

3. A process canvas appears with a single Read CSV operator.

4. Click the output port of the Read CSV operator.

5. In the Parameters palette, set the csv file to "weather.nominal.csv".

6. Click the "Activate Wisdom of Crowds" button at the bottom of the process canvas.

Leverage the Wisdom of Crowds to get operator recommendations based on your process design!

Activate Wisdom of Crowds

Read CSV
RapidMiner Studio Core

Synopsis
This operator is used to read CSV files.
[Jump to Tutorial Process](#)

File Edit Process View Connections Cloud Settings Extensions //LV Data Mining 2016/Mein erster ID3* – RapidMiner Studio Academia 7.0.001 @ boersch2

Repository Process Parameters Help

1 2 3 4 5 6

1. Open the Operators palette and expand the CSV section.

2. Double-click the Read CSV operator.

3. A process canvas appears with a single Read CSV operator.

4. Click the output port of the Read CSV operator.

5. In the Parameters palette, set the csv file to "weather.nominal.csv".

6. Click the "Activate Wisdom of Crowds" button at the bottom of the process canvas.

Leverage the Wisdom of Crowds to get operator recommendations based on your process design!

Activate Wisdom of Crowds

Read CSV
RapidMiner Studio Core

Synopsis
This operator is used to read CSV files.
[Jump to Tutorial Process](#)



Find data, operators...etc All Studio ▾

Parameters X

Loop Files

directory: n__2 Eric Bundel\tweets\train\neg

filter type glob

57% of users kept 'glob':

- glob (default)
- regex

filter by glob

recursive

enable macros

reuse results

Daten sammeln für
Entscheidungs-
unterstützung

Bspw. auch in ID
DIACOS: ICD->OPS

Ein Prozess-Ergebnis:

Daten

Miner Studio Academia 7.0.001 @ boersch2

Settings Extensions

Views:

Design

Results

Result History

ExampleSet (Read CSV)



Data



Statistics



Charts



Advanced Charts



Annotations

ExampleSet (14 examples, 0 special attributes, 5 regular attributes)

Filter (14 / 14 examples): all

Row No.	outlook	temperature	humidity	windy	play
1	sunny	hot	high	FALSE	no
2	sunny	hot	high	TRUE	no
3	overcast	hot	high	FALSE	yes
4	rainy	mild	high	FALSE	yes
5	rainy	cool	normal	FALSE	yes
6	rainy	cool	normal	TRUE	no
7	overcast	cool	normal	TRUE	yes
8	sunny	mild	high	FALSE	no
9	sunny	cool	normal	FALSE	yes
10	rainy	mild	normal	TRUE	no
11	sunny	mild	normal	TRUE	yes
12	overcast	mild	high	TRUE	yes
13	overcast	hot	normal	FALSE	no
14				TRUE	no

Nun Resultview

Ansichten der Daten

Data-View:
Tabellarische Darstellung der Daten
(das Daten-'Rechteck')

Repository

+ Add Data

- Samples
- DB
- LV Data Minin
- LV Data Minin
- Local Reposit
- Cloud Reposit

File Edit Process



Statistics-View: Beschreibung der Spalten (Attribute, Merkmale, Variablen)

Results

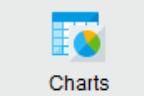
Result History



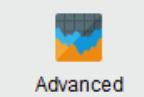
Data



Statistics



Charts



Advanced Charts



Annotations

Name

Type

Missing

Statistics

Filter (5 / 5 attributes):

Search for Attributes



outlook

Polynomial

0

overcast (4)

Most
rainy (5)Values
rainy (5)

temperature

Polynomial

0

Least
hot (4)Most
mild (6)Values
mild (6),

humidity

Polynomial

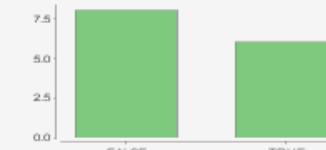
0

Least
normal (7)Most
high (7)Values
high (7),

windy

Polynomial

0



Open

play

Polynomial

0

Least
no (5)Most
yes (9)Values
yes (9),

Namenspalte zeigt die besondere Rolle: bspw. „label“ ist die Klassenspalte. „play“ ist hier also noch keine Klassenspalte, sondern ein normales reguläres Attribut.

Das ausgewählte Attribut zeigt die empirische Verteilung seiner Attributwerte

Repository

+ Add Data

- Samples
- DB
- LV Data Minin
- LV Data Minin
- Local Reposit
- Cloud Reposit



>

>

<

<

Examples: 14 Special Attributes: 0 Regular Attributes: 5



Views:

Design

Results

Result History X ExampleSet (Set Role) X Repository

Chart style: play no yes

Scatter

Scatter Multiple Scatter Matrix Scatter 3D Scatter 3D Color Bubble Parallel

Scatter Series Survey SOM BLOCK Density

Series Multiple

Histogram Histogram Color Bars Bars Stacked Pareto Andrews Curves Distribution

Annotations

Web Quartile Quartile Color Quartile Color Matrix Pie Pie 2D Ring

**Chart-View: Datenvisualisierung
- bspw. Scatter(plot) = Streudiagramm**

Übung

The screenshot shows the RapidMiner Studio interface. On the left, there's a vertical sidebar with icons for Data, Statistics, Charts, Advanced Charts, and Annotations. The 'Charts' icon is selected. In the main area, there's a 'Result History' tab and an 'ExampleSet (Set Role)' tab. Above the chart grid, there's a 'Chart style:' dropdown set to 'play' with options 'no' and 'yes'. A blue arrow points from the 'Charts' sidebar to the 'Scatter' chart thumbnail. Another blue arrow points from the 'Scatter' chart thumbnail to an orange callout box containing the text 'Chart-View: Datenvisualisierung - bspw. Scatter(plot) = Streudiagramm'. To the right of the callout is a teal button with the word 'Übung'. The main area displays a grid of 24 chart thumbnails, each with a name below it. The charts include various types such as Scatter, Scatter Multiple, Scatter Matrix, Scatter 3D, Scatter 3D Color, Bubble, Parallel, Series, Survey, SOM, BLOCK, Density, Histogram, Histogram Color, Bars, Bars Stacked, Pareto, Andrews Curves, Distribution, Web, Quartile, Quartile Color, Quartile Color Matrix, Pie, Pie 2D, and Ring.



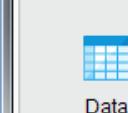
Views:

Design

Results

Result History

ExampleSet (Read CSV)



Data



Statistics



Charts



Advanced Charts



Annotations

Chart style:



Scatter

play no yes

X-Axis:

temperature

 Log scale

y-Axis:

humidity

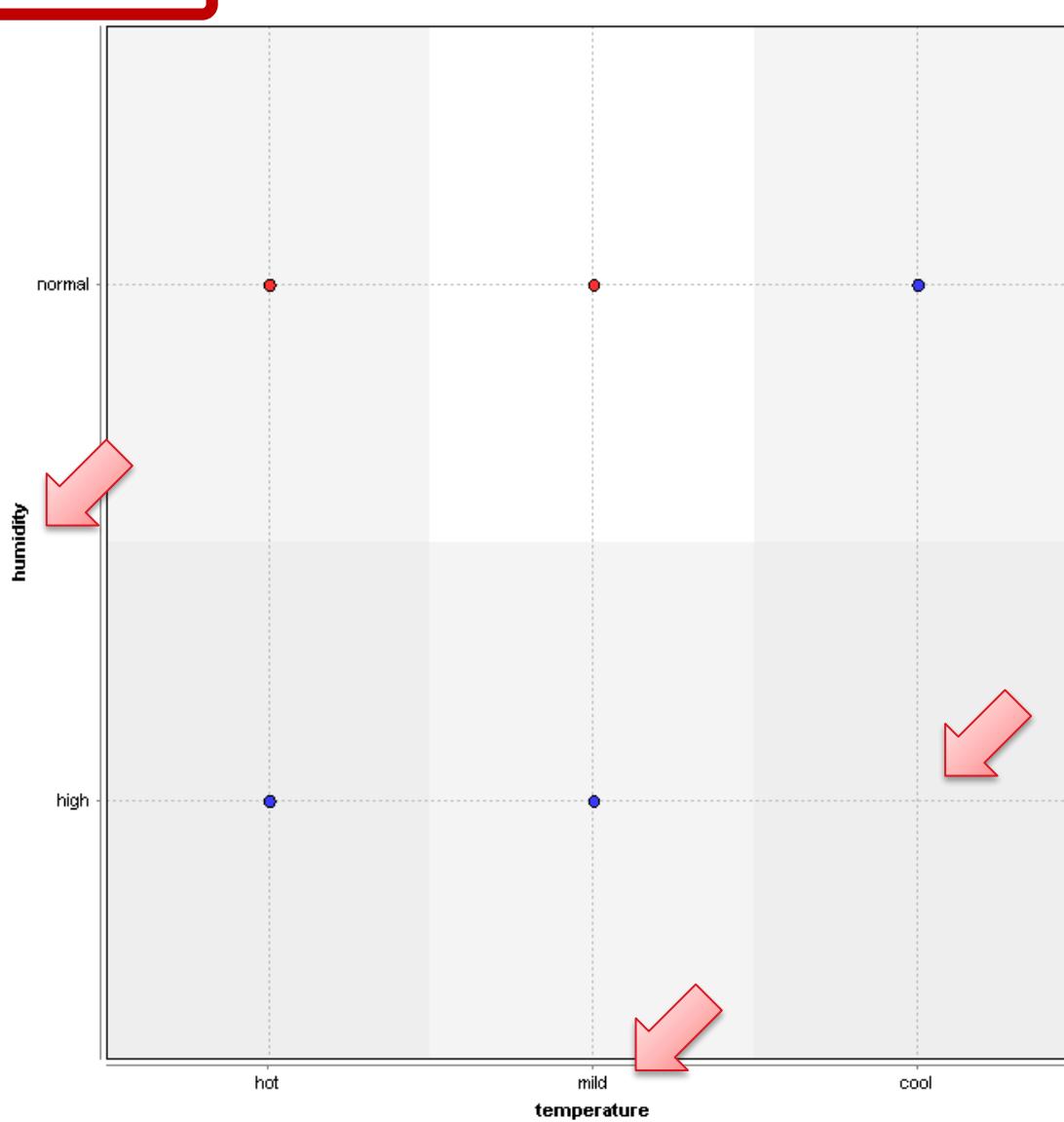
 Log scale

Color Column:

play

 Log scale

Jitter:

 Rotate labels

Repository

+ Add Data

- Samples
- DB
- LV Data Minin
- LV Data Minin
- Local Reposit
- Cloud Reposit



Views:

Design

Results

Result History

ExampleSet (Read CSV)



Data



Statistics



Charts



Advanced Charts



Annotations

Chart style:

play no yes

X-Axis:

temperature

 Log scale

y-Axis:

humidity

 Log scale

Color Column:

play

 Log scale

Jitter:

 Rotate labels

humidity

high

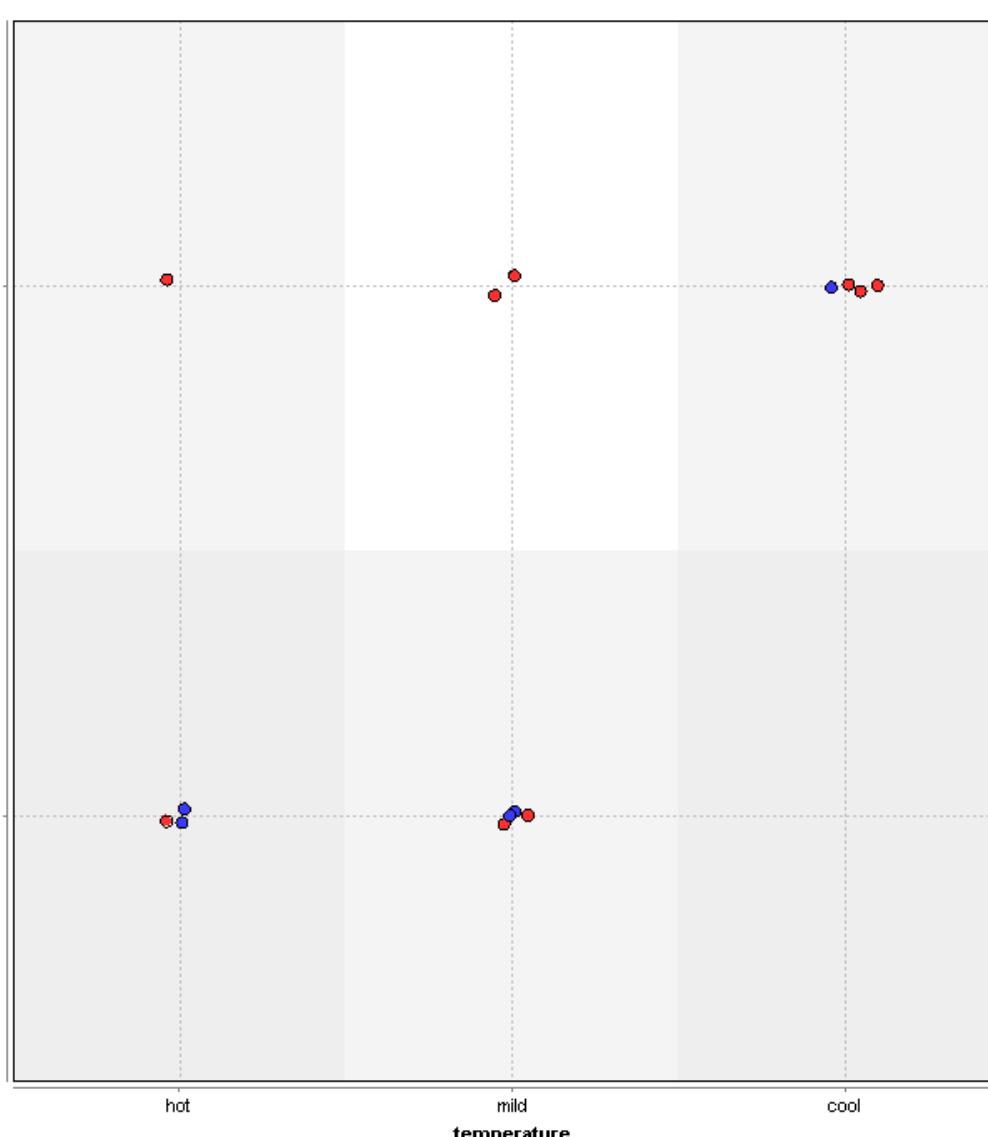
hot

temperature

normal

mild

cool



Repository

+ Add Data

- Samples
- DB
- LV Data Minin
- LV Data Minin
- Local Reposit
- Cloud Reposit

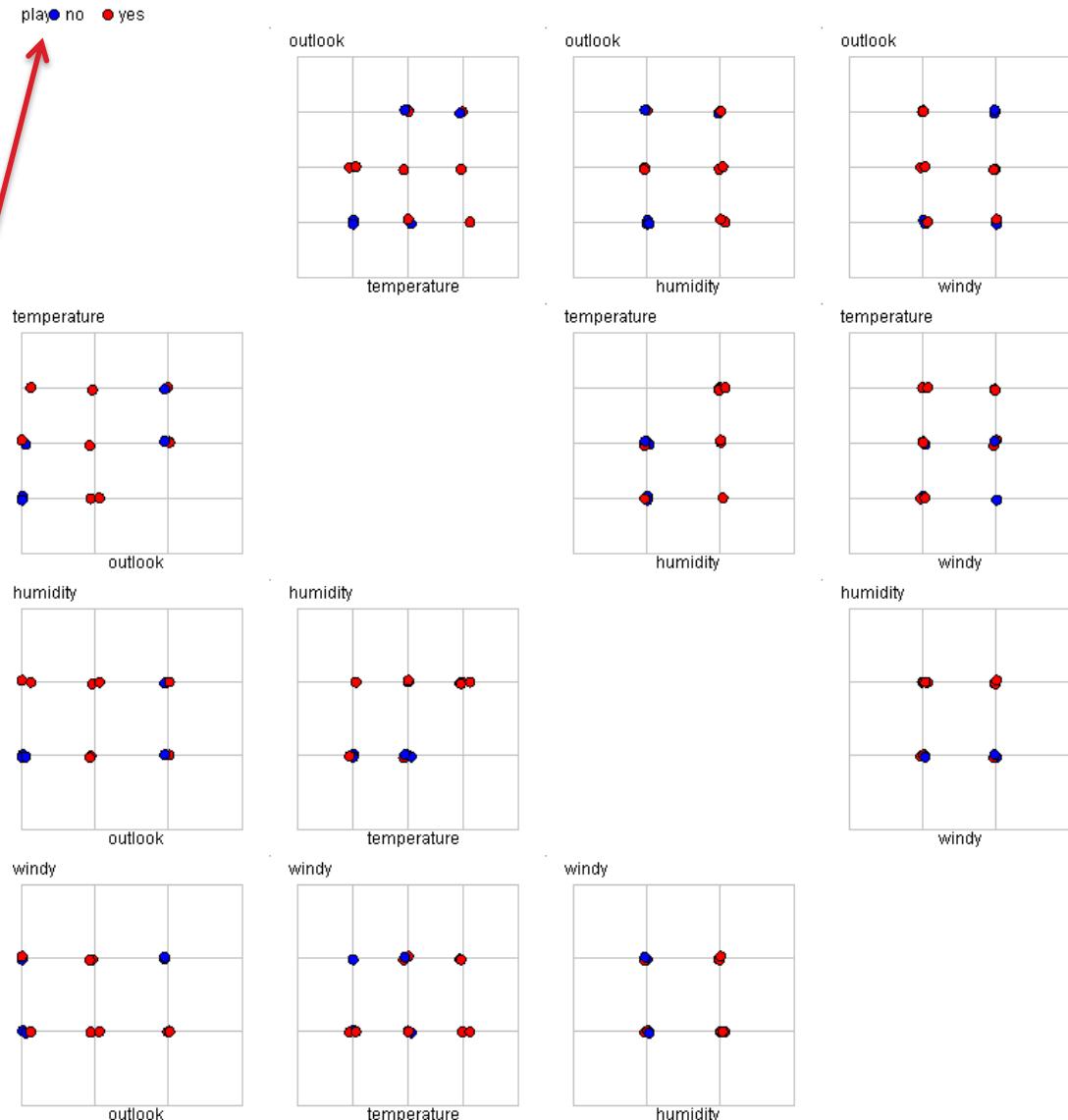


Matrix aller Streudiagramme

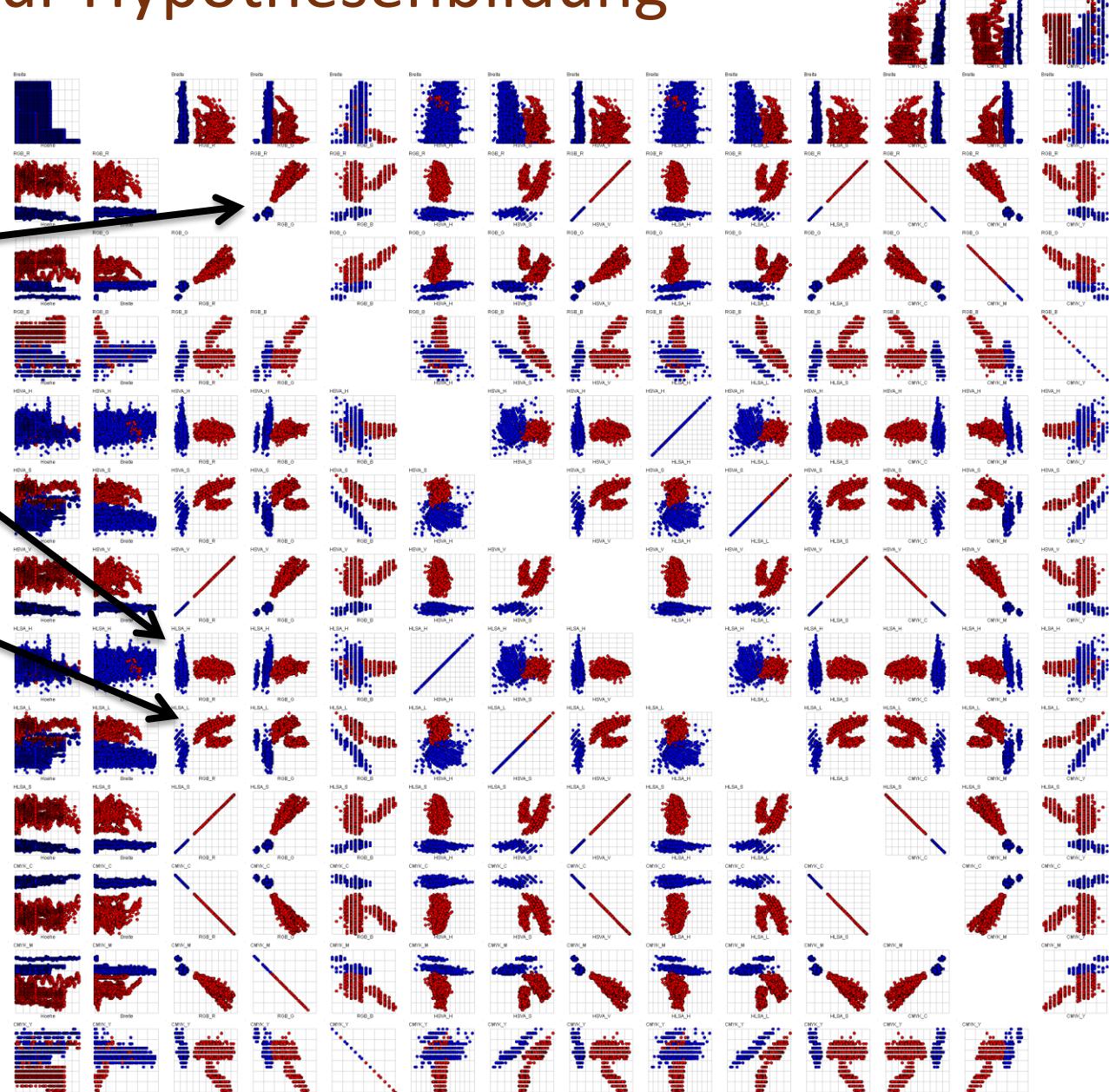
- Auch
Scatterplotmatrix
(SPLOM)

k Variablen ->
 $k \times k$ - Matrix

Matrixelemente:
Streudiagramme
Einfärbung (**blau, rot**)
durch Zielklasse



Reale SPLOM für Hypothesenbildung



Welche Merkmale sind abhangig? —

Welche Merkmale eignen sich für die Klassifizierung?

Cluster

Visual Analytics (Sonderheft Informatikspektrum , Dezember 2010)



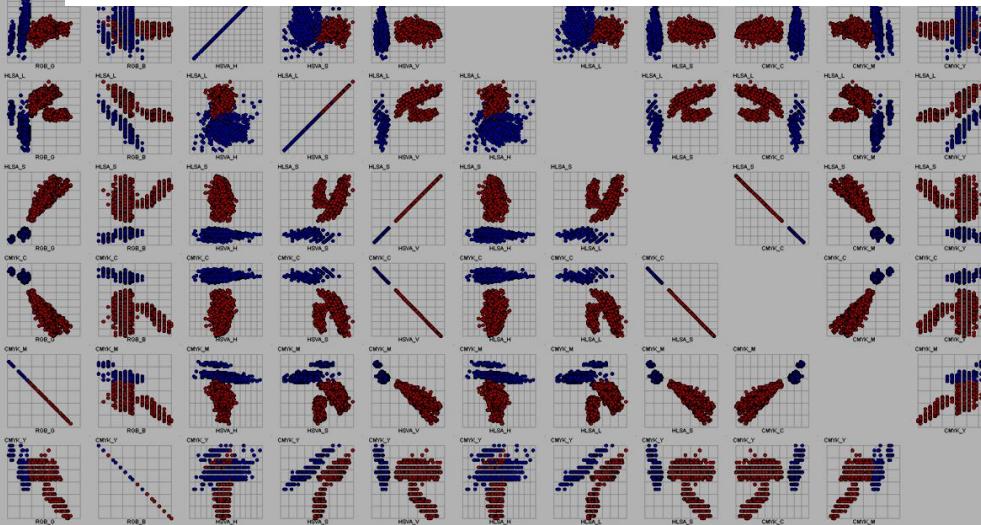
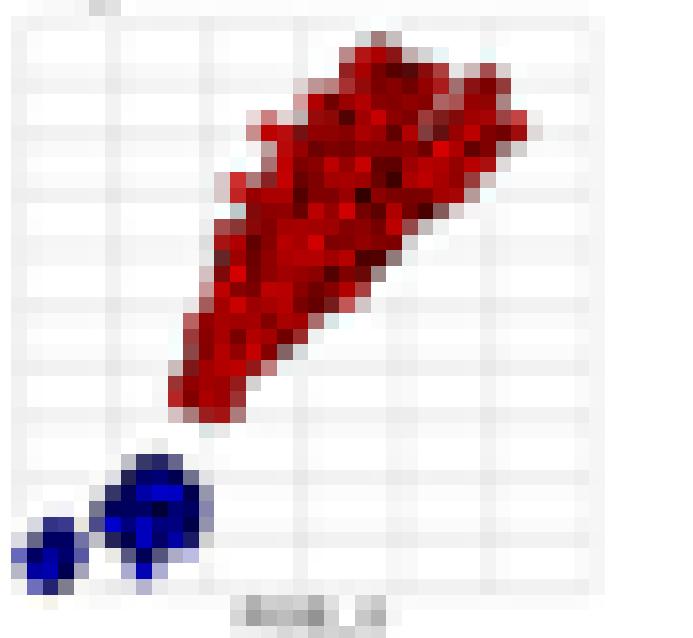
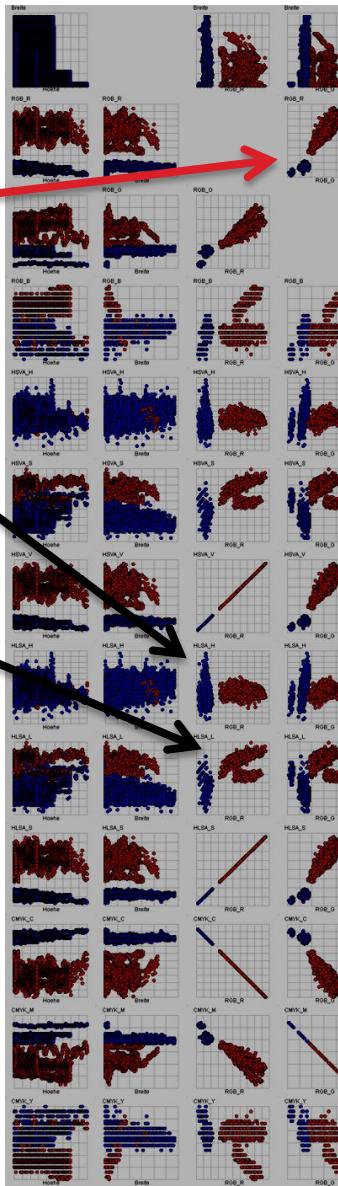
Reale SPLOM für Hypothesenbildung

Welche Merkmale sind abhängig?

Welche Merkmale eignen sich für die Klassifizierung?

Cluster

Visual Analytics
(Sonderheft
Informatikspektrum
, Dezember 2010)





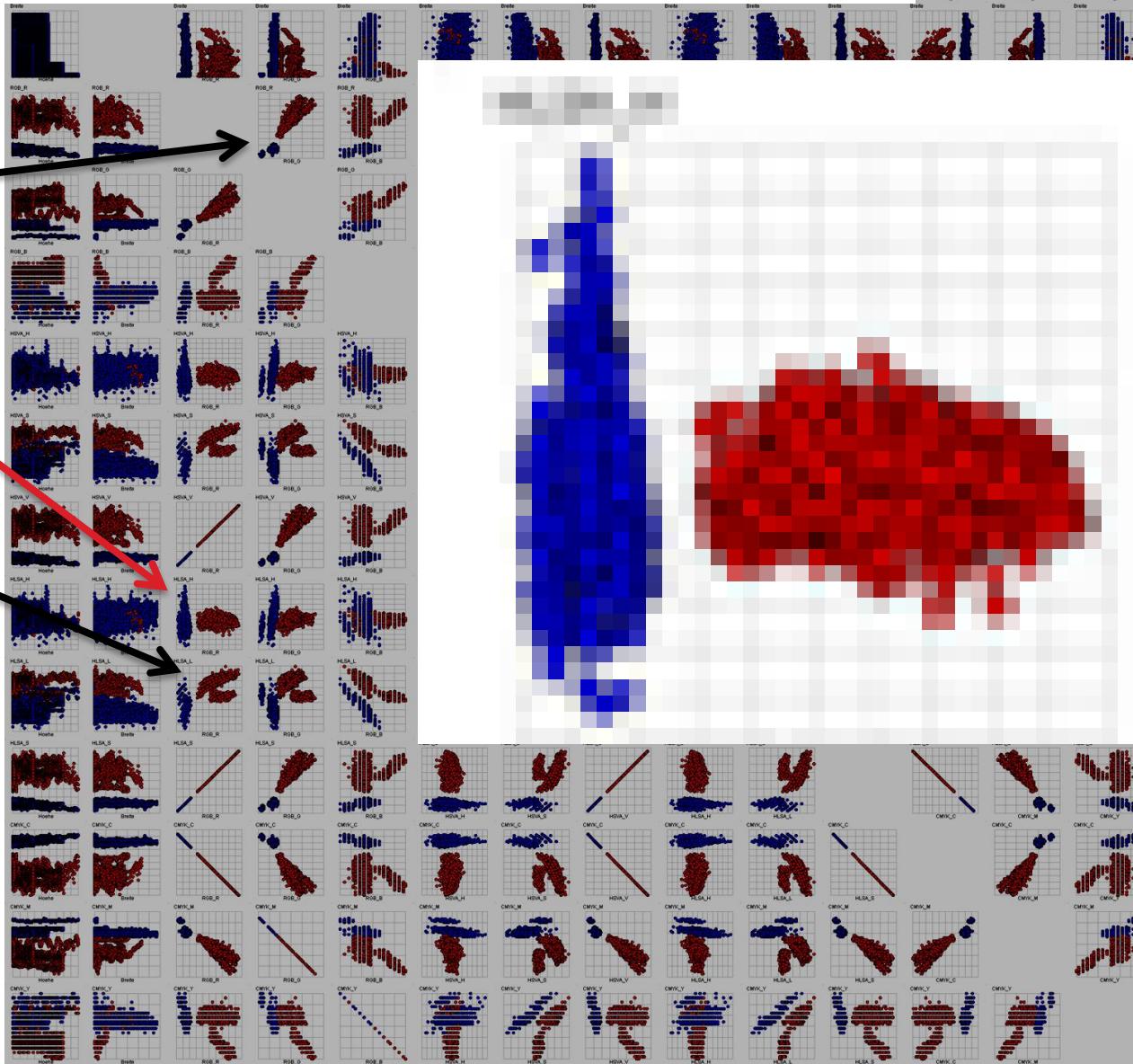
Reale SPLOM für Hypothesenbildung

Welche Merkmale sind
abhängig?

Welche Merkmale
eignen sich für die
Klassifizierung?

Cluster

Visual Analytics
(Sonderheft
Informatikspektrum
, Dezember 2010)





Reale SPLOM für Hypothesenbildung

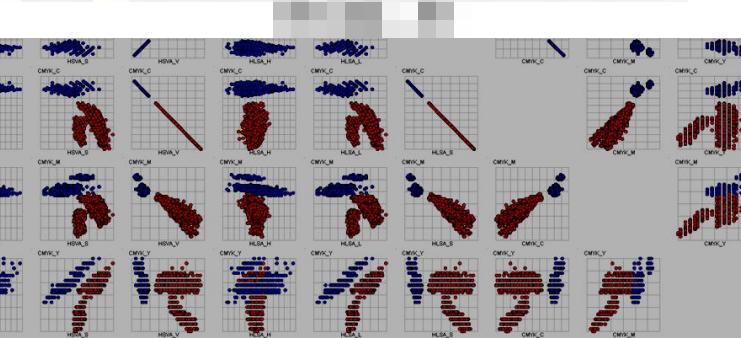
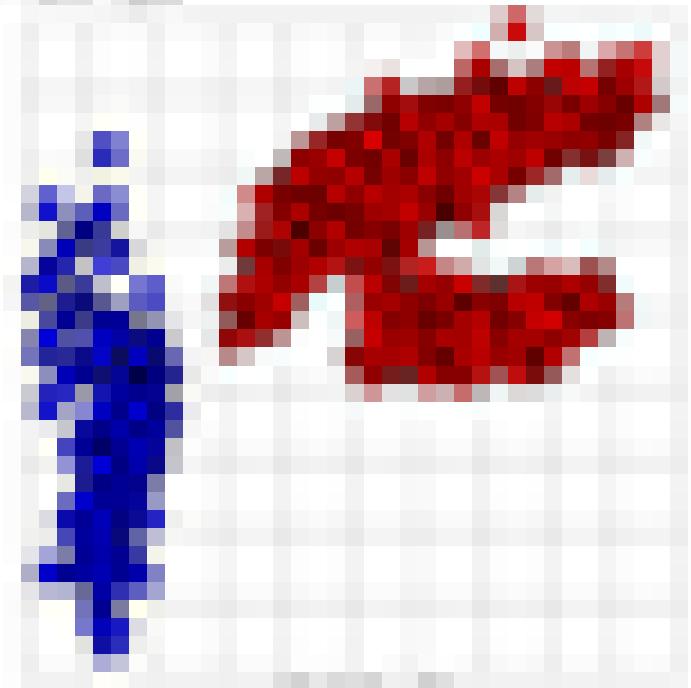
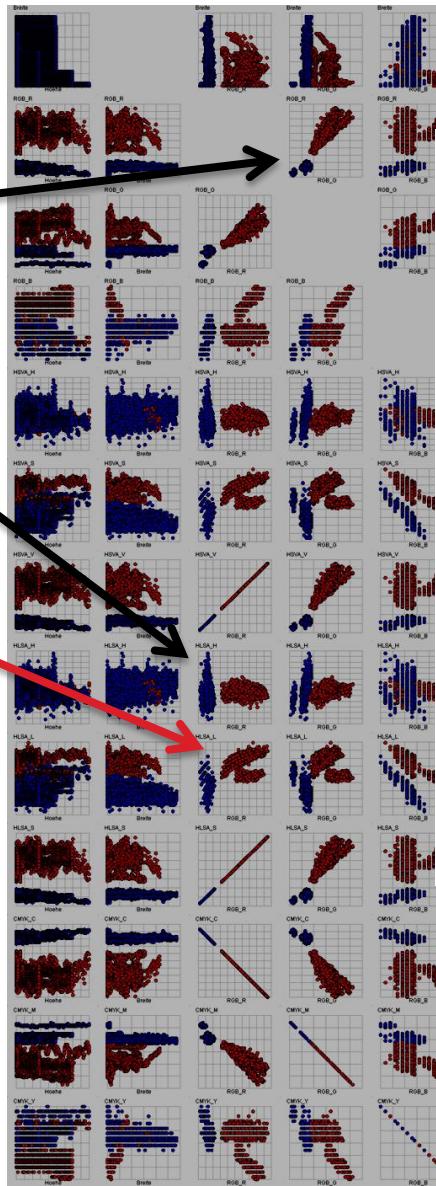
Welche Merkmale sind abhangig? ---

Welche Merkmale eignen sich für die Klassifizierung?

Cluster

Visual Analytics

(Sonderheft Informatikspektrum , Dezember 2010)





Spongebob und Crabs

- Lernen von Objektbeschreibungen durch Zeigen von Beispielen
- Diplomarbeit Benjamin Kieper:
Entwurf und Implementierung einer Anwendung zum dialogbasierten, überwachten Lernen von Objektmodellen aus Bildern



Einführung

WAS ist Data Mining?

- Data Mining (DM) und Knowledge Discovery in Databases (KDD)
- Aufgabenstellungen des DM, viele Beispiele
- Wissensrepräsentation

WIE funktioniert es?

- Entscheidungsbäume:
 - Repräsentation,
 - Lernen / Konstruieren
 - Praktisch

WIE gut sind die Modelle?

- Performance von Klassifikatoren



Einschätzung des Gelernten

Fragen in der Praxis:

1. Welche Lernmethoden sollen auf welches Problem angewendet werden?
 2. Welcher trainierte Klassifikator (Modell) soll letztendlich verwendet werden?
- ⇒ Lernmethoden und Modelle müssen **verglichen** werden.
- ⇒ Performance von Klassifikatoren



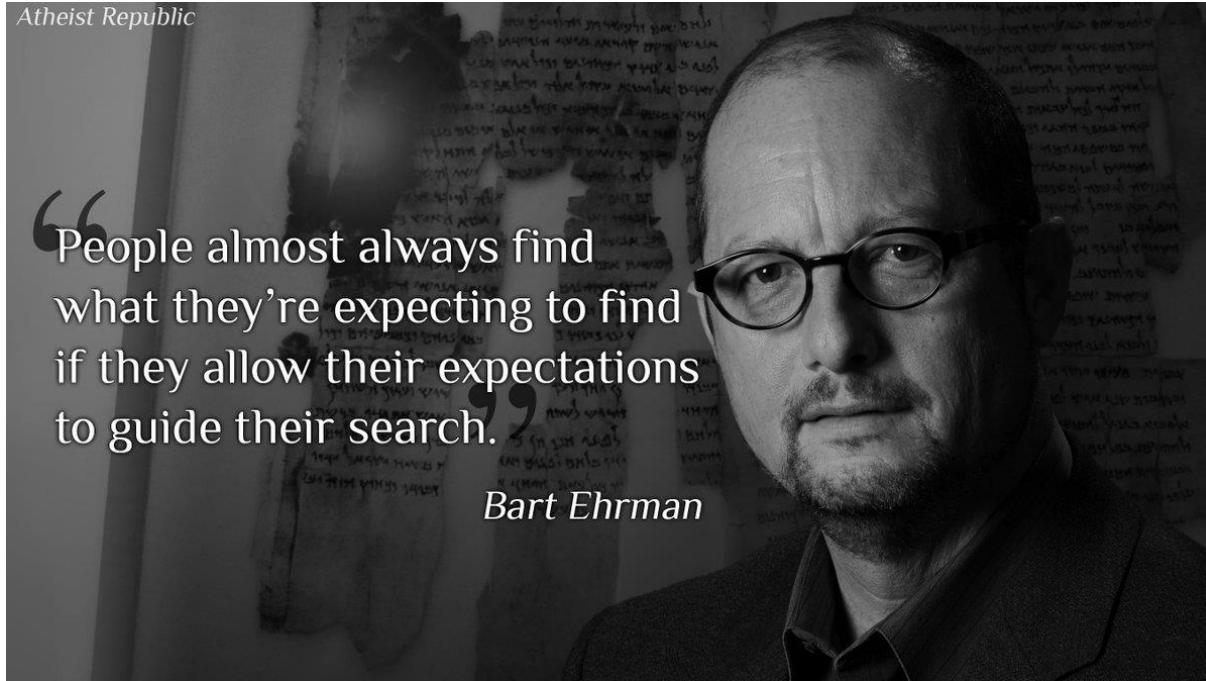
Next

- Performance von Klassifikatoren, Einschätzung des Gelernten,
 - Konfusionsmatrix
 - Fehlermaße: Fehlerrate, Erfolgsrate
 - **Trainingsfehler**
 - Bsp kNearest Neighbor -> Der Trainingsfehler hat ein Problem
 - [Overfitting]
- Fehlerschätzung – Wie sieht die Praxis aus:
 - Holdout-Methode + Stratifikation + Wiederholung -> **Kreuzvalidierung**
 - [leave-one-out, bootstrap]
 - **Praktisches Vorgehen**
 - Automated ML kann das auch



“People almost always find what they’re expecting to find if they allow their expectations to guide their search.”

Bart Ehrman



Expectation Bias

The tendency for experimenters to believe, certify, and **publish** data that **agree** with their expectations for the outcome of an experiment, and to disbelieve, discard, or downgrade the corresponding weightings for data that appear to conflict with those expectations.

https://en.wikipedia.org/wiki/List_of_cognitive_biases



Güte eines Modells

Wie kann die Prognosegüte eines Modells beurteilt werden:

Wir schätzen die Performance (bspw. Fehlerrate, Erfolgsrate), die das Modell bei seiner Anwendung in der Praxis erreichen wird.

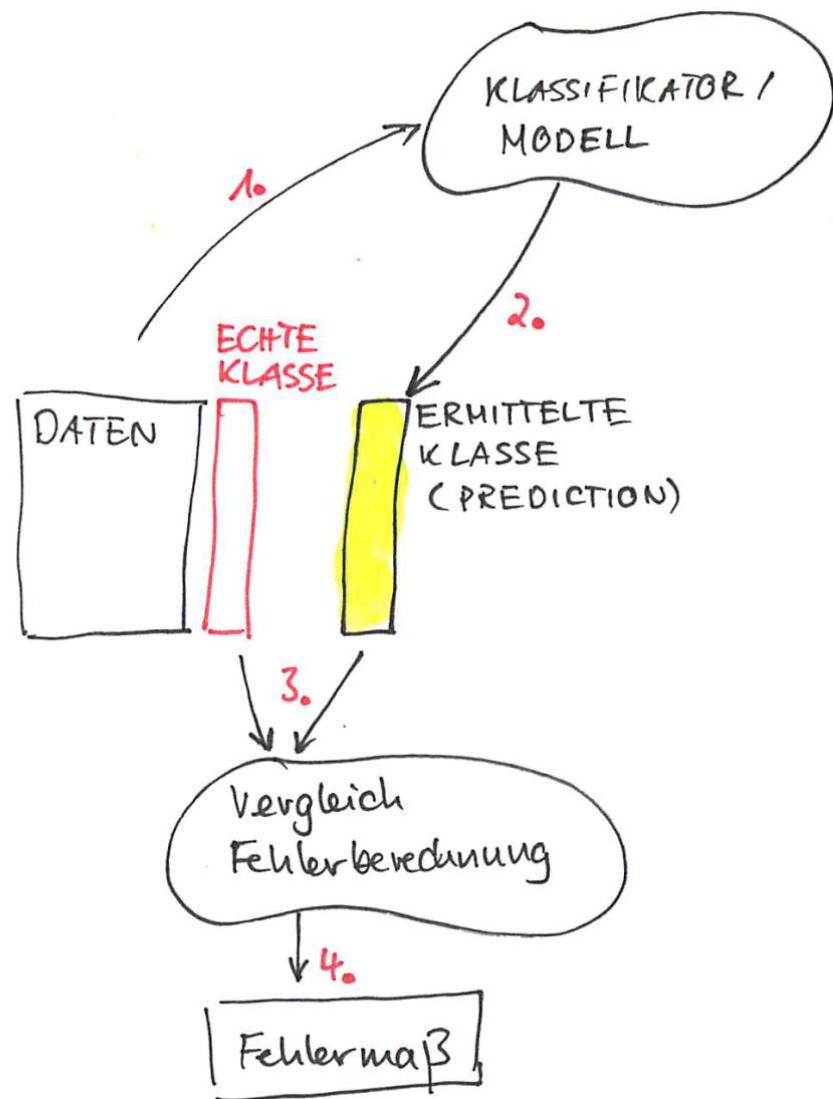
Wie? Zwei Wege (mindestens)

- Der einfache Weg
- Der bessere Weg



Modell auf Daten mit bekannter Klassifikation anwenden

1. Modell erstellen
2. Modell auf die Daten **anwenden** und Prognose der Klasse ermitteln (*prediction*)
3. Performance berechnen aus **Vergleich** der ermittelten mit der tatsächlichen Klasse
- Vielzahl von Performancemaßen





Konfusionsmatrix

Gegenüberstellung der Häufigkeiten der **ermittelten** Klassen (*predicted*) und der **tatsächlichen** Klassen (*real*) eines Testfalles.

Dimension **K x K** mit K als Anzahl der Klassen

		Tatsächliche Klasse				Σ
		K1	K2	K3	K4	
Ermittelte Klasse	K1	12	7	3	4	26
	K2	1	9	8	2	20
	K3	6	4	9	13	32
	K4	94	12	3	5	114
Σ		113	32	23	24	$\frac{192}{192}$

Hauptdiagonale =
richtig
Vorhergesagte



Aus der Konfusionsmatrix: Erfolgsrate, Fehlerrate



- **Erfolgsrate** = #richtig Vorhergesagte / #Alle
- **Fehlerrate** = #falsch Vorhergesagte / #Alle
- Engl. *accuracy, classification error*

		Tatsächliche Klasse				Σ
		K1	K2	K3	K4	
Ermittelte Klasse	K1	12	7	3	4	26
	K2	1	9	8	2	20
	K3	6	4	9	13	32
	K4	94	12	3	5	114
	Σ	113	32	23	24	$\frac{192}{192}$

Beispiel:

$$\begin{aligned}\text{Erfolgsrate} &= \\ (12+9+9+5) &/ 192 \\ &= 18.2\%\end{aligned}$$

Fehlerrate =

$$\begin{aligned}(7+3+\dots+12+3) &/ 192 \\ &= 81.8\%\end{aligned}$$

$$18.2 + 81.8 = \underline{100}$$



Spezialfall:

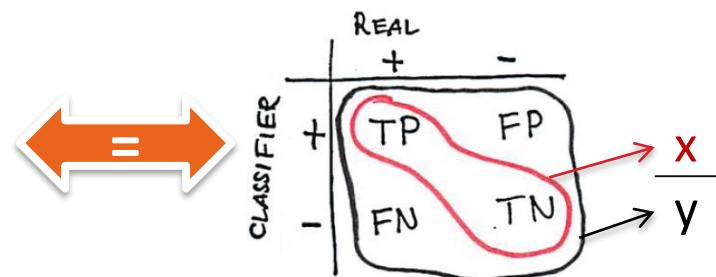
Konfusionsmatrix bei binärer Klassifikation

Genau 2 Klassen bspw: {ja, nein}, {positiv, negativ}, {0,1}, {wahr, falsch}, {krank, gesund} ...

→ Konfusionsmatrix 2 x 2:

		Tatsächliche Klasse		
		positiv	negativ	
Ermittelte Klasse	positiv	TP = True positive	FP = False positiv	Fehler 1. Art
	negativ	FN = False negative	TN = True negative	

- Erfolgsrate = $(TP+TN) / (TP + FP + FN + TN)$
- Fehlerrate = $(FP+FN) / (TP + FP + FN + TN)$

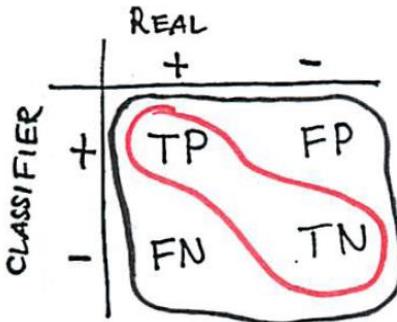




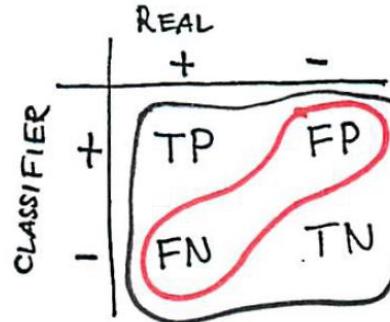
Diagnose mit 95% <Gütemaß>

krank (+, positiver Befund), gesund (-, negativer Befund)

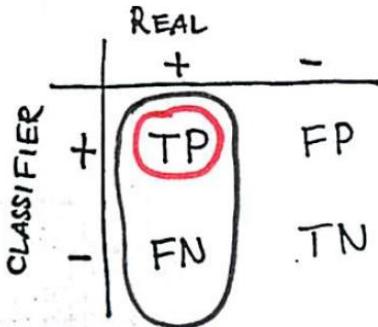
grüne Performancewerte möglichst hoch, rote niedrig



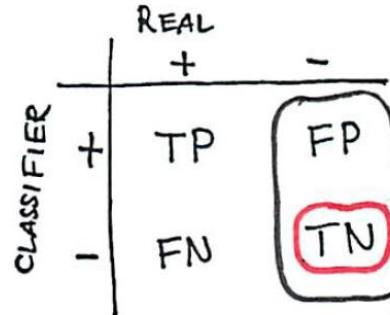
accuracy



classification error



sensitivity, recall,
true positive rate (TPR)
 $= \text{TP} / (\text{TP} + \text{FN})$



specificity,
true negative rate (TNR)
 $= \text{TN} / (\text{TN} + \text{FP})$

95% Erfolgsrate:

Die Diagnose stimmt in 95% der Fälle.

95% Fehlerrate:

95% der Diagnosen sind falsch.

95%: Sensitivität:

95% der Kranken werden als krank erkannt.

95% Spezifität:

95% der Gesunden werden als gesund erkannt.

Lies: 95% der schwarz-markierten sind rot-markiert.

Beispiel vom Anfang: Herzinfarkt: Sensitivity = 81.4% Specificity = 92.1% PPV = 72.9% Accuracy = 89.9%



Weitere Maße

		REAL	
		+	-
CLASSIFIER	+	TP	FP
	-	FN	TN

precision,
positiv predictive value (PPV)
 $= \text{TP} / (\text{TP} + \text{FP})$

		REAL	
		+	-
CLASSIFIER	+	TP	FP
	-	FN	TN

false positive rate (FPR)
 $= \text{FP} / (\text{FP} + \text{TN})$

		REAL	
		+	-
CLASSIFIER	+	TP	FP
	-	FN	TN

negative predictive
value (NPV)
 $= \text{TN} / (\text{FN} + \text{TN})$

		REAL	
		+	-
CLASSIFIER	+	TP	FP
	-	FN	TN

false discovery rate (FDR)
 $= \text{FP} / (\text{TP} + \text{FP})$

95% Positiver Vorhersagewert:

Von krank Erklärten sind 95%
krank (und 5% gesund)

95% Negativer Vorhersagewert:

95% der als gesund Erklärten
sind gesund.

95% Ausfallrate:

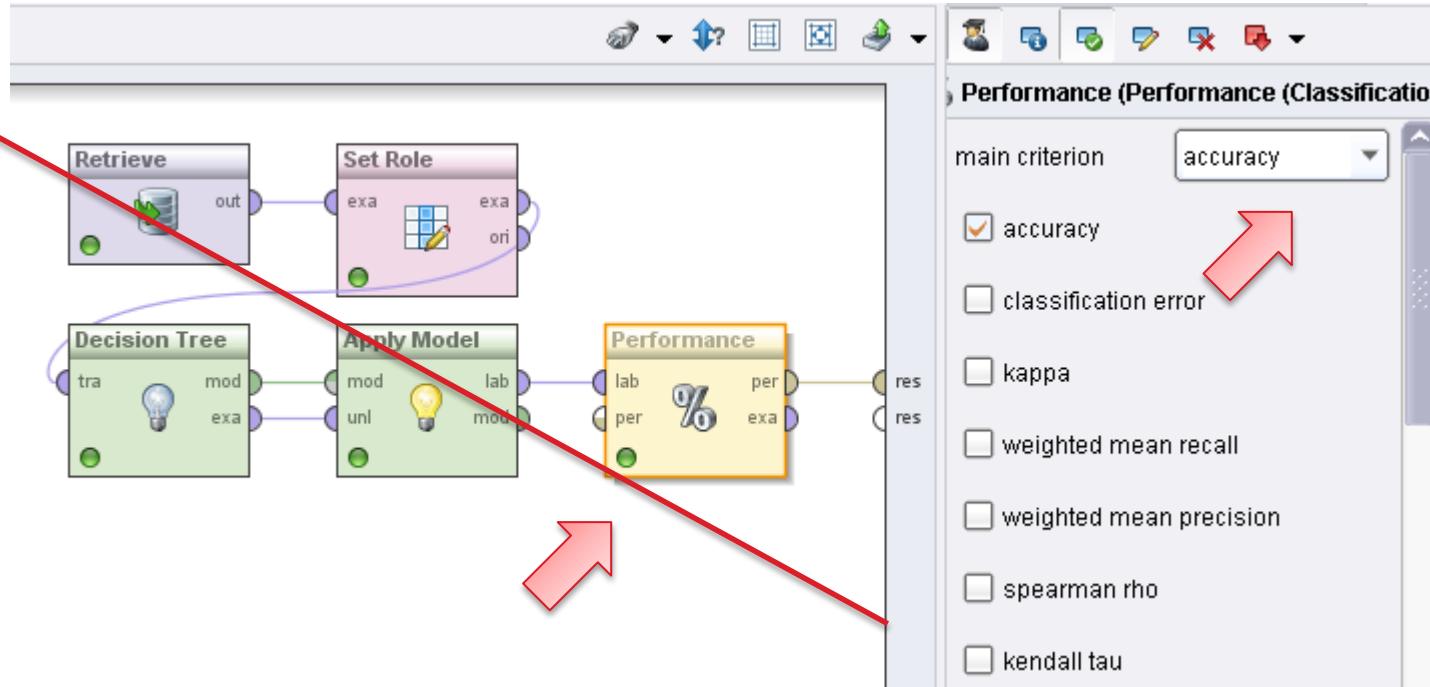
95% der Gesunden werden als
krank erklärt.

95% false discovery rate:

95% der krank Erklärten sind
gesund.



So nicht!



Wo ist das Problem?

Die Performance eines Klassifizierers auf seinen Trainingsdaten ist kein guter Indikator für seine Leistung auf anderen Daten (Testmenge.)

Die Evaluierung ist viel zu optimistisch!



Trainingsfehler

auch Resubstitutionsfehler, Resubstitution = „nochmal verwenden“

- **Trainingsfehler := Fehlerrate eines Klassifizierers auf seinen Trainingsdaten**
- Die Trainingsdaten werden also *nochmal* verwendet (Training und Test)
- Entsprechend **Trainingserfolgsrate** etc.

Warum ist der Trainingsfehler keine gute Schätzung für die spätere Leistung des Modells?

Gegenbeispiel: k-Nearest Neighbor (k-NN) - ein sog. Lazy learner („faul“)

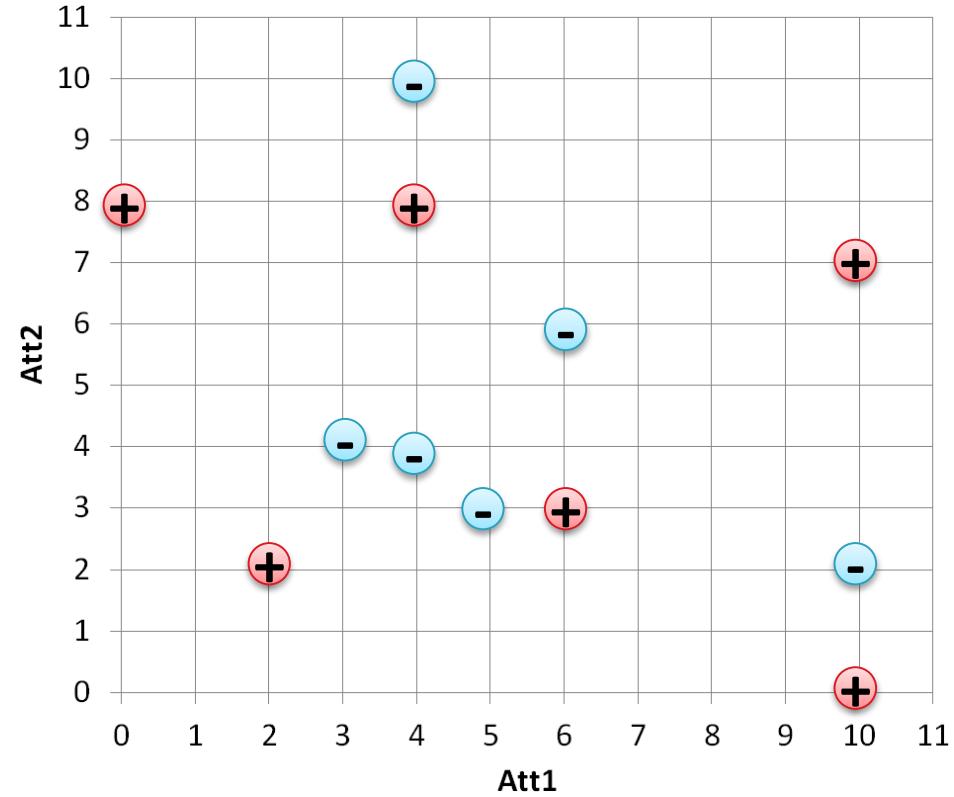
- Lernen: Speichern aller Examples (Modell == Trainingsmenge)
- Anwenden auf unbekanntes x: Finden der k ähnlichsten Beispiele in der Trainingsmenge -> deren häufigste Klasse ausgeben
- Quiz: Wie groß ist der Trainingsfehler eines k-NN mit k = 1 ?



Trainingsfehler eines k-Nearest Neighbor mit k = 1

Att1	Att2	Klasse
6	3	+
10	0	+
5	3	-
6	6	-
4	8	+
2	2	+
0	8	+
3	4	-
5	3	-
10	7	+
4	4	-
10	2	-
4	10	-

=



- Modell == Trainingsmenge
- Wie wird ein Datensatz X aus der Trainingsmenge, bspw. X = (6,6), durch den k-NN mit k=1 klassifiziert?

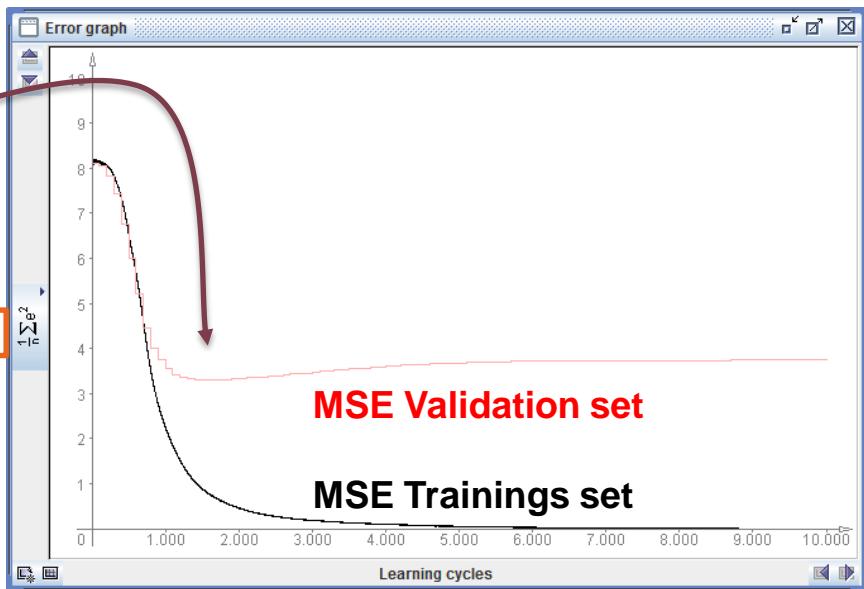
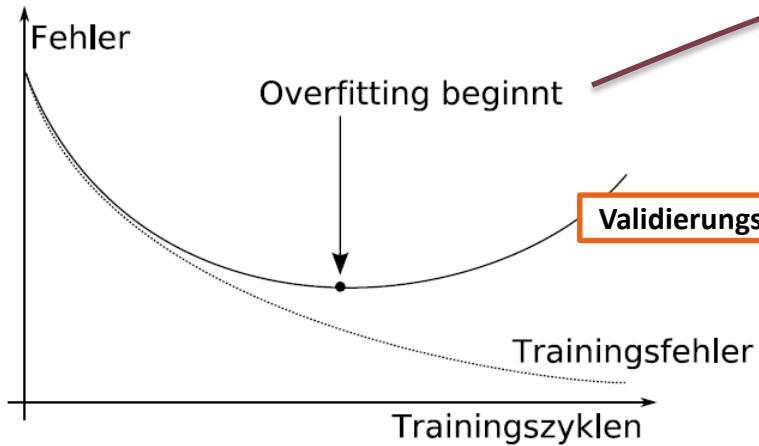
Die Trainings-Erfolgsrate eines k-NN mit k = 1 ist stets 100%.

Das ist eine **sehr** optimistische Schätzung für den zukünftigen Fehler, oder?

Overfitting

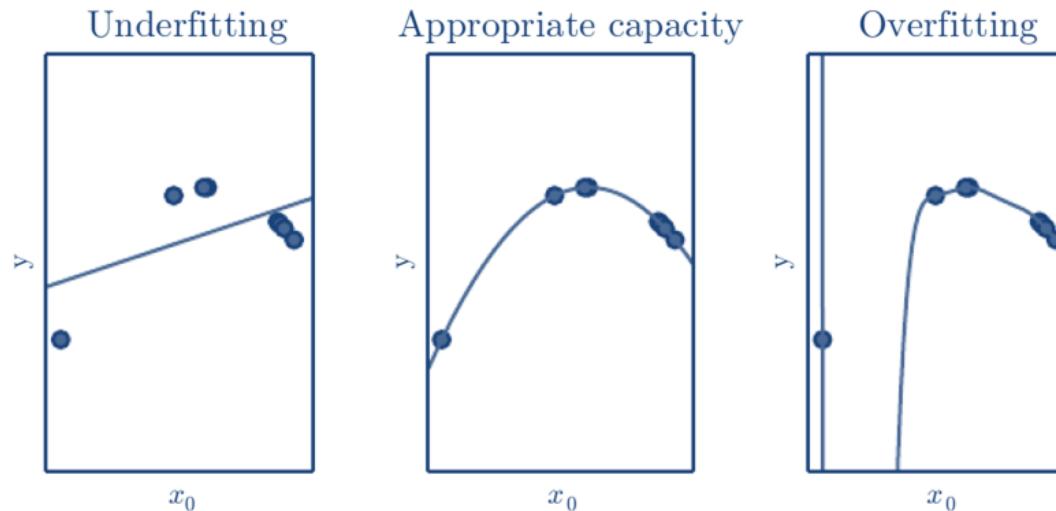
Def: Leistung auf Trainingsdaten zwar sehr gut, aber nicht im späteren Einsatz

- überangepasste Wissensrepräsentation, **Auswendiglernen** der Trainingsdaten,
- **Generalisierungsfähigkeit fehlt**
- nicht nur bei neuronalen Netzen



Kapazität, Ausdruckskraft, Parameteranzahl

(Geheime) quadratische Funktion $y = f(x_0)$ erzeugt Datenpunkte



Fit und Predict:

Links: lineare Funktion, hoher Trainingsfehler, **hoher Testfehler**

Mitte: quadratische Funktion, geringer Trainingsfehler, **geringer Testfehler**

Rechts: Polynom 9-ten Grades, geringer Trainingsfehler, **hoher Testfehler**

[GBC16] Goodfellow, Ian ; Bengio, Yoshua ; Courville, Aaron: Deep Learning. MIT Press, 2016. – <http://www.deeplearningbook.org>



Maßnahmen gegen Overfitting

- Mehr Daten,
 - gerne auch künstlich erzeugte (*real-world variation*)
 - *data set augmentation* (Bilder verzerrn, Rauschen einbringen (Daten, Gewichte))
- *early stopping*: Training bis zur Verschlechterung beim *validation set*
- Größe (Freiheitsgrade) der Modelle beschränken,
 - Pruning bei DT
 - wenige hidden neurons (aber gerade deep learning: riesige Neuronanzahl)
 - *parameter sharing*, bspw. *convolution* bei deep learning
 - *regularization* bei NN (Strafe für große Gewichte)

Quiz: Was stimmt?

- falsch** 
1. Ein Lernalgorithmus, der auf den Trainingsdaten einen hohe Erfolgsrate erreicht, liefert Modelle, die beim späteren Einsatz eine hohe Erfolgsrate zeigen.
 2. Ein Lernalgorithmus, der auf den Trainingsdaten einen niedrige Erfolgsrate erreicht, liefert Modelle, die beim späteren Einsatz eine niedrige Erfolgsrate zeigen.
- richtig** 

Wie können wir den zukünftigen Fehler besser schätzen?



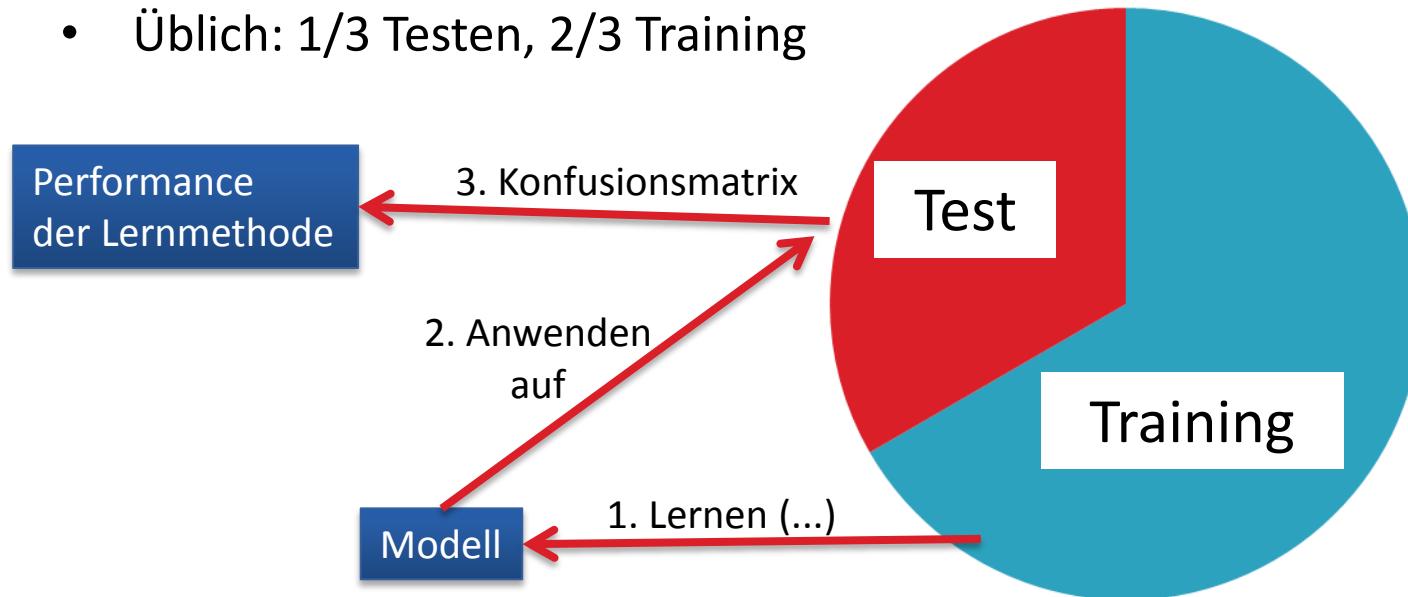
Die Trainings-Erfolgsrate einer k-NN mit $k = 1$ ist stets 100%. Ein hochangepasster Klassifikator.

Der Trainingsfehler ist nicht geeignet, den Fehler während der späteren Anwendung des Klassifikators zu schätzen.

- Trotzdem wird diese Methode oft verwendet.
 - Als Heuristik: ein großer Trainingsfehler zeigt, dass die Modellierung mit diesem Lernverfahren nicht gelingt.
- Einfacher Weg - Trainingsfehler
- Besserer Weg - kommt jetzt

Holdout-Methode

- Holdout-Methode reserviert Daten für das **Testen** und verwendet nur den Rest fürs **Training**
- 2 disjunkte Mengen
 - Trainingsdaten: Erstellung des Klassifikators, Lernen
 - Testdaten: Schätzung der Fehlerrate des Klassifikators
- Üblich: 1/3 Testen, 2/3 Training





Stratifikation

Problem: Möglicherweise sind die für die Testmenge gewählten Daten nicht repräsentativ

Stratifikation:

- Datenmenge so teilen, dass die relativen Klassenhäufigkeiten in Test- und Trainingsdaten übereinstimmen

-> **Stratifizierte Holdout Methode**

Beispiel:

Die Daten enthalten in der Klassenspalte 33% Damen und 67% Herren.

Dann sollte die Teilung in Trainings- und Testmenge so erfolgen, dass beide ebenfalls 33% Damen und 67% Herren enthalten (und nicht alle Damen zufällig in der Testmenge landen).



Wiederholtes Holdout

- Idee: Holdout-Methode **mehrmals mit unterschiedlichen Datenteilungen wiederholen** und
- den Mittelwert alle Fehlerraten bilden
- damit erfolgt eine Kompensierung der Fehler, die durch die Auswahl der Stichprobe (Teilung) entstanden sind
- Das wichtigste Verfahren des wiederholten Holdout:
Kreuzvalidierung (engl. cross validation, CV)



n-fache Kreuzvalidierung (CV)

Datensätze in n Mengen (**Partitionen, folds**) aufteilen, z.B. n=3 oder n=10.

Für i=1 bis n:

- i-te Menge ist Testmenge
 - restliche n-1 Mengen sind Trainingsmenge -> Lerner erzeugt Modell i
 - bestimme die Leistung des Modells auf der Testmenge -> Fehler i
- > n Fehlerwerte

Bestimme Mittelwert und ggf. Streuung des Fehlers -> Ergebnis der CV

- Das Ergebnis der CV ist eine Schätzung des zu erwartenden Generalisierungsfehlers auf unbekannten Daten.
- Das Modell für den späteren Einsatz wird mit allen verfügbaren Daten erstellt (sog. 100%-Modell)
- Ergebnis des Data-Mining-Prozesses = **Modell + Fehlerschätzung des Modells**

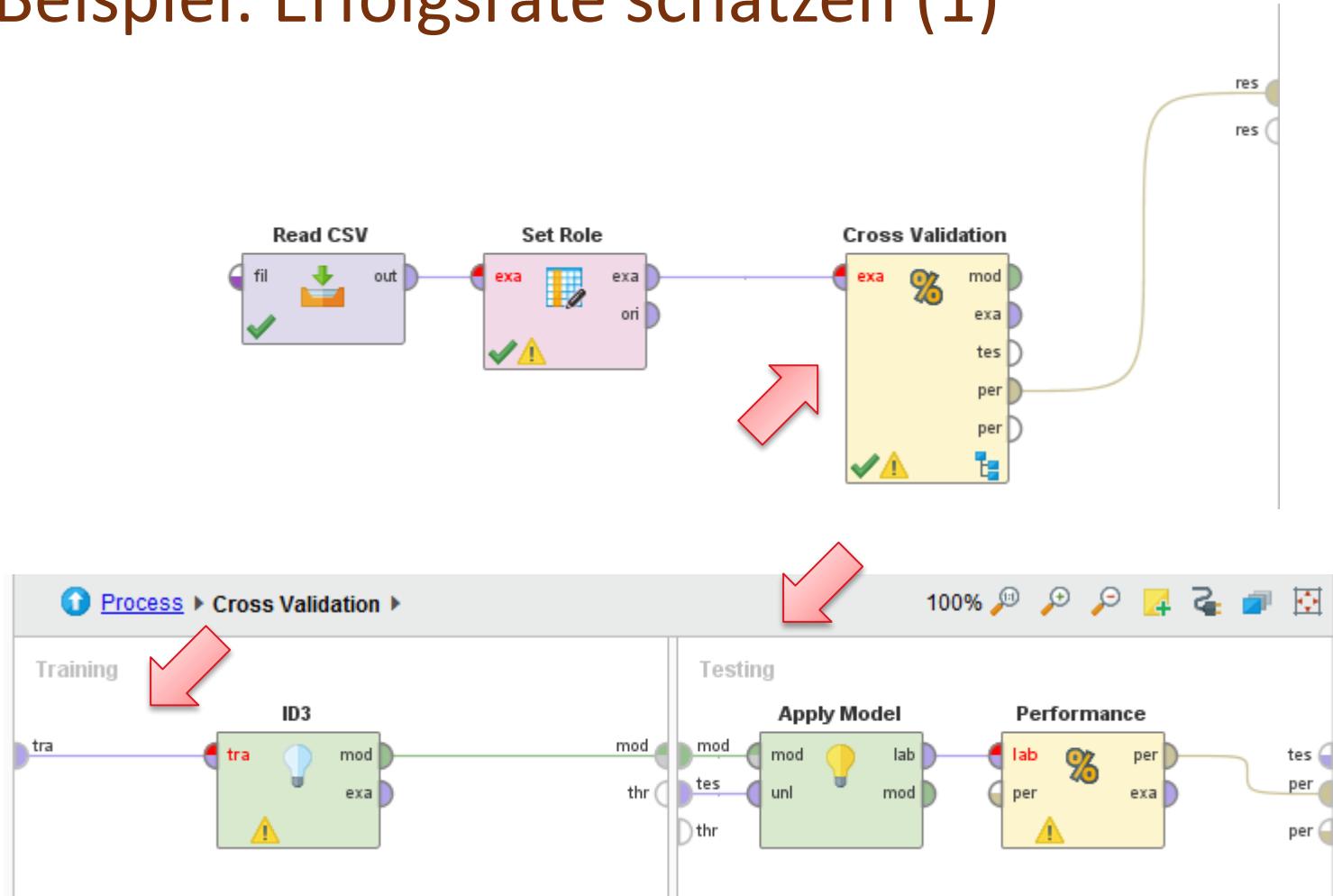


Stratifizierte Kreuzvalidierung

- Besser: Stratifikation
 - Relative Klassenhäufigkeiten in allen Partitionen annähernd gleich zur Verteilung in der Gesamtmenge
 - -> **stratifizierte Kreuzvalidierung**
- Noch besser: **wiederholte Kreuzvalidierung**
 - Bsp: 10mal wiederholte 10fache Kreuzvalidierung
 - -> 100 Modelle zu lernen und zu testen
- in RapidMiner: Cross Validation
- Beispiel: Erfolgsrate eines ID3 und NN auf der Datenmenge weather vergleichen



Beispiel: Erfolgsrate schätzen (1)

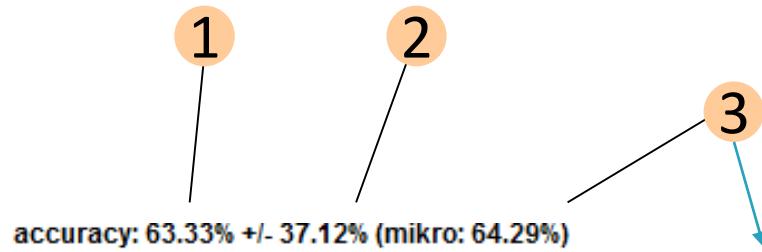


weather.nominal.csv, n=5



Erfolgsrate schätzen (2)

- Erfolgsrate (ID3) = 63.33%



	true no	true yes	class precision
pred. no	3	3	50.00%
pred. yes	2	6	75.00%
class recall	60.00%	66.67%	



Erfolgsrate schätzen (3)

5 Konfusionsmatrizen:

0	2
1	1

1	1
0	1

0	0
1	1

1	0
0	2

real

1	0
0	2

predicted

						Mittelwert
Classification_error (FP+FN)/All	3/3	1/3	1/2	0/3	0/3	36.67%
Accuracy (TP+TN)/All	0/3	2/3	1/2	3/3	3/3	63.33% 1

Standardabweichung der Erfolgsraten e_i : $s = \sqrt{\frac{1}{n} (e_i - \bar{e})^2} = 0,371184291$ 2

Kumulierte Konfusionsmatrix:

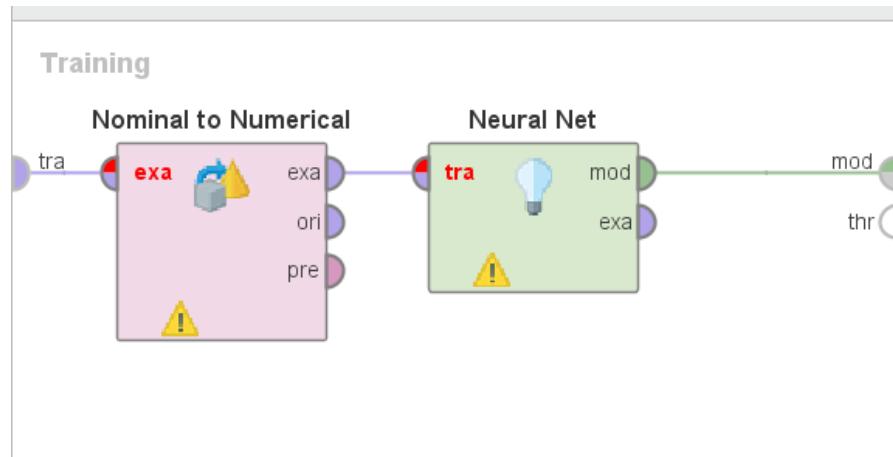
3	3
2	6

Erfolgsrate der kumulierten Konfusionsmatrix: $9/14 = 0,642857143$ 3



Erfolgsrate schätzen (4)

- Neuronales Netz
- Erfolgsrate (KNN) = 63.33 %



Sie wissen noch, wie ein Backpropagation-Netzwerk funktioniert?
Die gesamte Theorie ist hier nur ein kleiner Operator. Unter vielen.



Andere Schätzverfahren – Leave one out

- Leave one out
 - n-fache Kreuzvalidierung mit $n=$ Instanzenanzahl
 - Testmenge = eine Instanz, Rest ist Trainingsmenge
 - Keine zufällige Stichprobenwahl -> deterministisch
 - Hoher Rechenaufwand: n-mal lernen
 - Realistische Schätzung
 - keine praktische Bedeutung



Andere Schätzverfahren - Bootstrap

- Bootstrap
 - aus einer Datenmenge mit n Instanzen wird n -mal zufällig eine Instanz gewählt und in die Trainingsdaten kopiert
 - Nicht gewählte Instanzen -> Testdaten
 - $p=1/n$ für jede Instanz, dass sie ausgewählt wird
 - $p=1-1/n$, dass sie nicht ausgewählt wird
 - Potenzierung über n Auswahlvorgänge führt zur Wahrscheinlichkeit, dass eine bestimmte Instanz in der Trainingsmenge vorkommt von:

$$1 - (1 - 1/n)^n \xrightarrow{n \rightarrow \infty} 1 - e^{-1} \approx 0.632$$

- Pessimistische Schätzung, weil Trainingsmenge nur 63% der ursprünglichen Instanzen enthält
- **geeignet für kleine Datenmengen, beliebig oft wiederholbar**



Praktisches Vorgehen, bspw. beim Data Mining Cup, Data Challenge

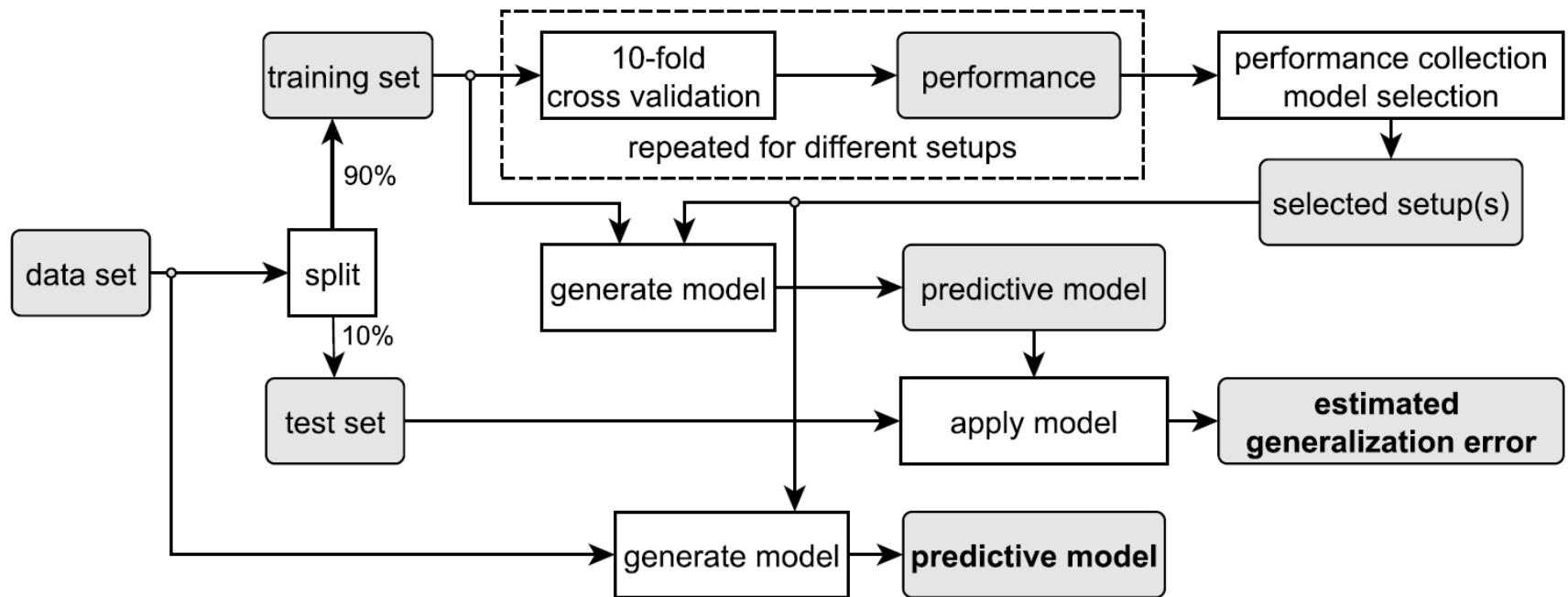


Fig. 7 The evaluation process used to estimate the performance of each setup.

Boersch, Ingo; Füssel, Uwe; Gresch, Christoph; Großmann, Christoph; Hoffmann, Benjamin:

Data mining in resistance spot welding: A non-destructive method to predict the welding spot diameter by monitoring process parameters.

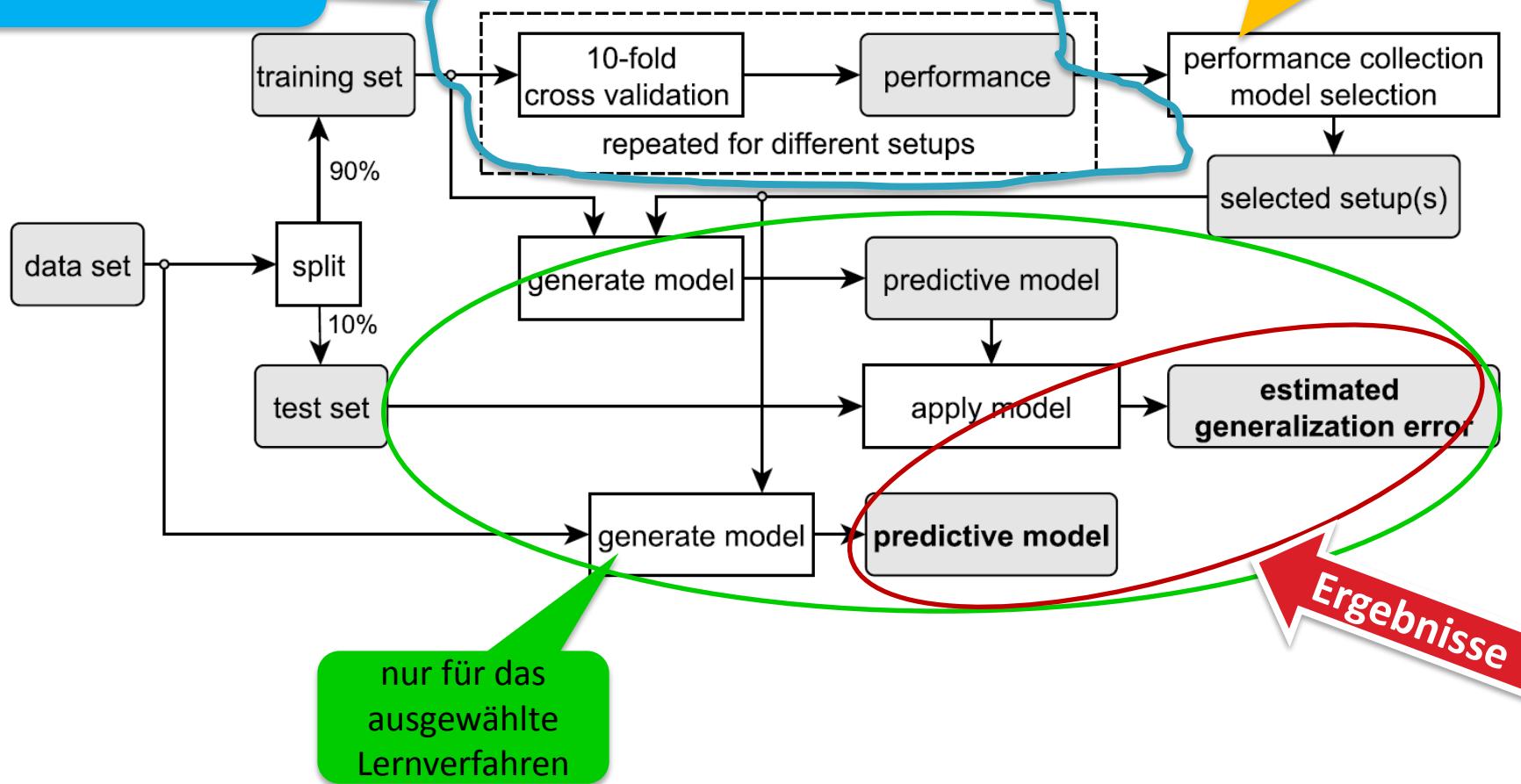
In: The International Journal of Advanced Manufacturing Technology (2016), 1--15. <http://dx.doi.org/10.1007/s00170-016-9847-y>, ISSN 1433-3015



Praktisches Vorgehen, bspw. beim Data Mining Cup, Data Challenge

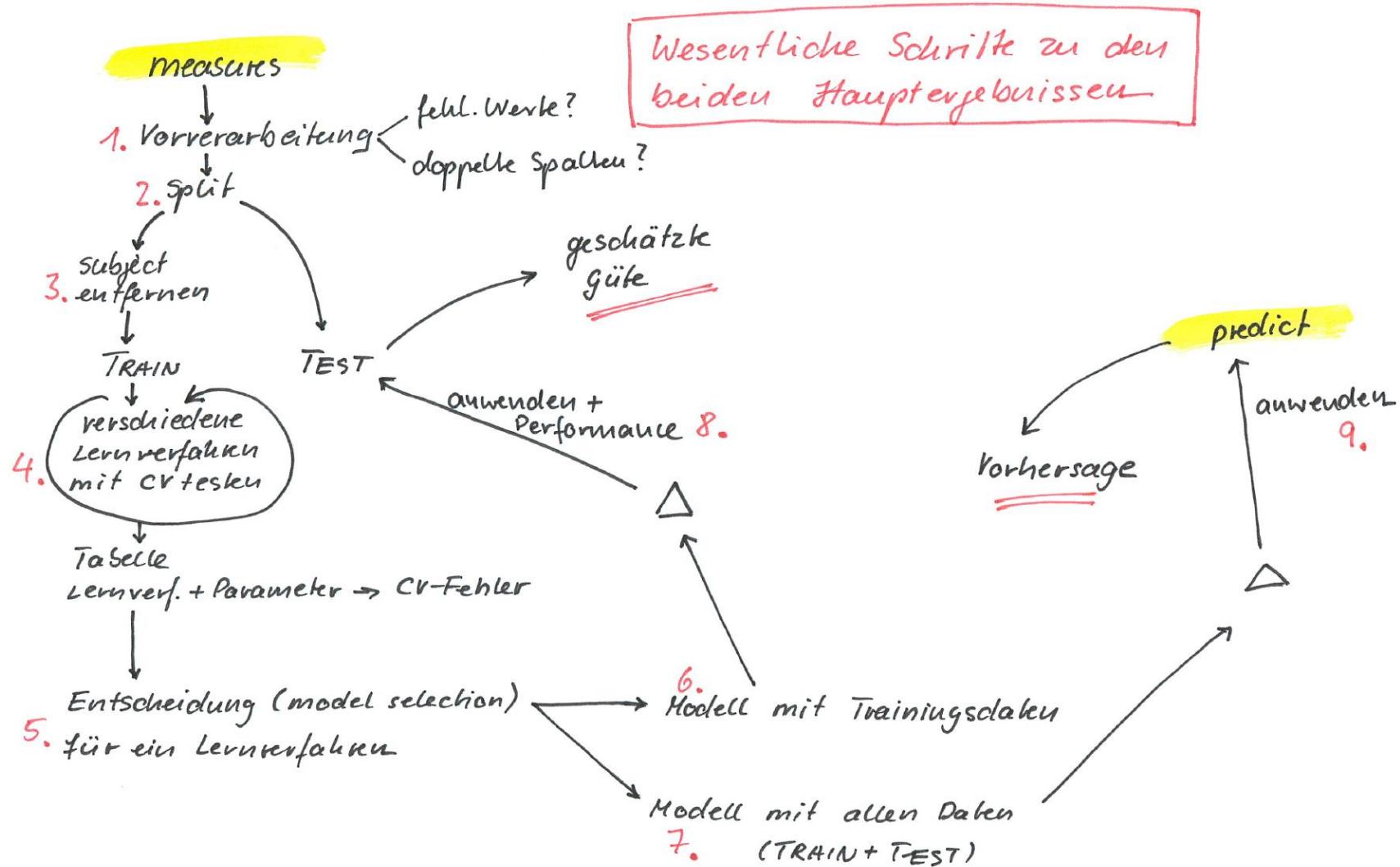
Modelle probieren,
optimieren, dabei vom
CV-Fehler leiten lassen

CV-Fehler in einer
Tabelle sammeln –
Basis für
model selection





Praktisches Vorgehen, konkreter für Data Challenge

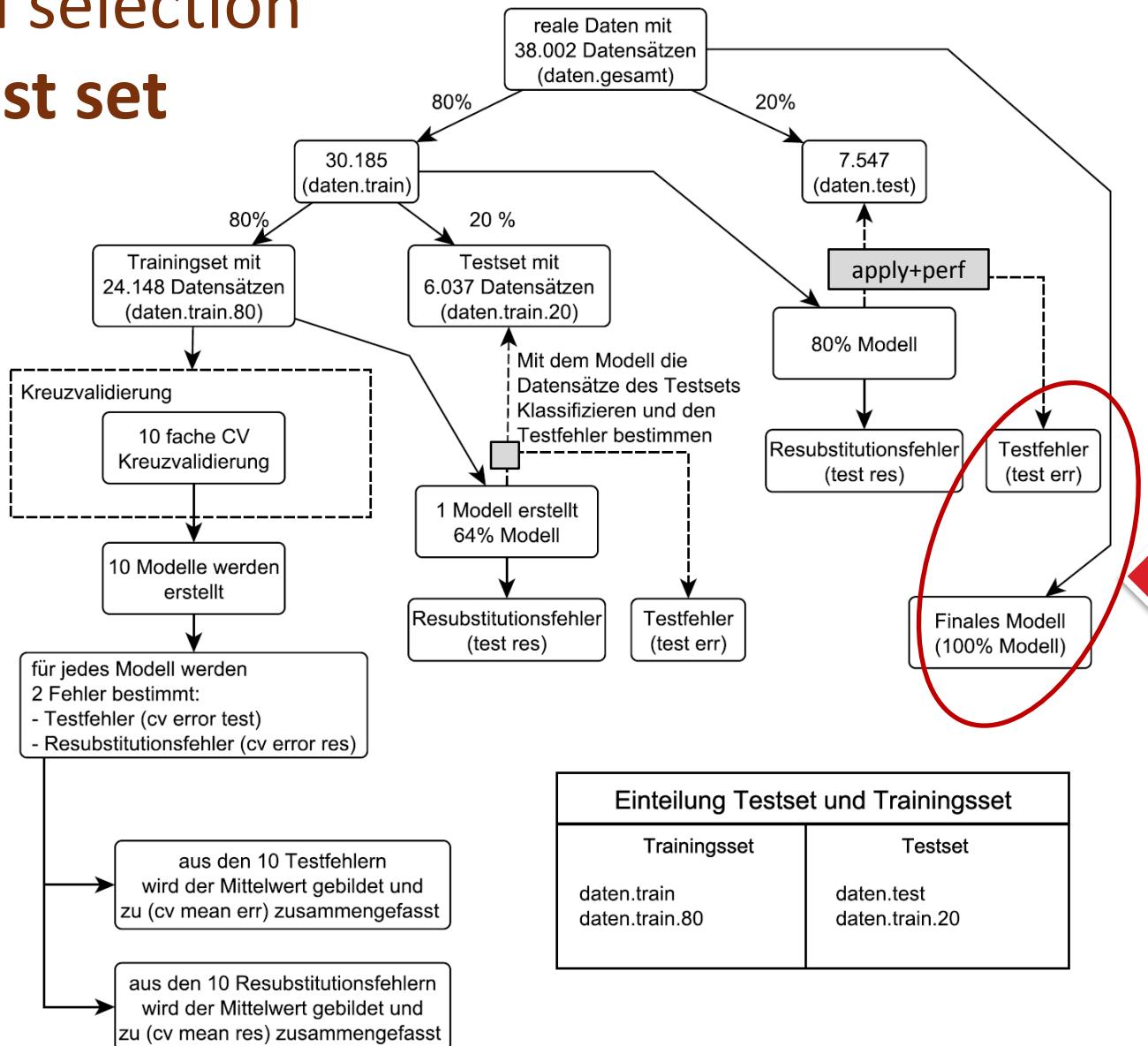




model selection mit test set

Warum?

- Suche
- Modell
- Güte



Ergebnis



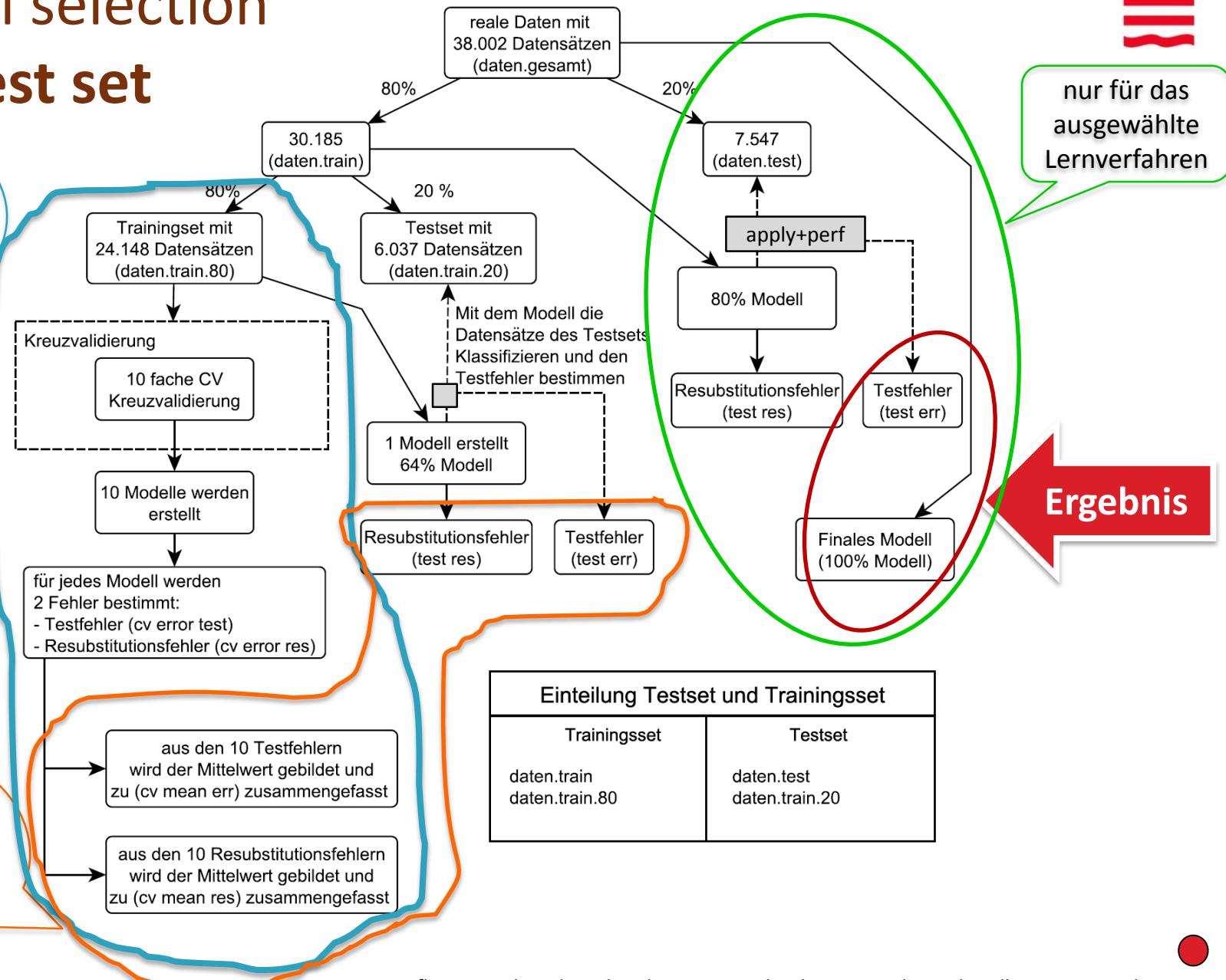
model selection mit test set

Modelle probieren, optimieren, dabei vom CV-Fehler leiten lassen

Warum?

- Suche
- Modell
- Güte

in einer Tabelle sammeln – Basis für model selection





Trend: Self-Service-Analytics

- einfach zu bedienende Analysewerkzeuge
- ohne Vorkenntnisse, ohne Programmierung
- datenbasierte Meetings

Beispiele RapidMiner

1. Welches Modell ist bei meinen Daten anwendbar?
<http://mod.RapidMiner.com/#app>
2. Welches von denen ist das beste mit den besten Parametern?
Auto Model
 - nur in Educational und Professional Licence
 - Prediction, Clustering, Outlier
 - <https://docs.RapidMiner.com/latest/studio/auto-model/>



http://mod.RapidMiner.com/#app

MOD by rapidminer About

Filters

Column types

- Numerical
- Binary
- Categorical

Target type

- No target
- Numerical
- Binary
- Categorical

Number of columns

- 10s
- 100s
- 1000s
- 10,000s

Number of rows

- 1000s
- 10,000s
- 100,000s
- 1,000,000s

Advanced options

- Updatable
- Can handle missings
- Uses row weights

Applicable Models Show all

Segmentations (10)

- k-Means ~50.0% DOCS
- Agglomerative Clustering ~9.4% DOCS
- DBSCAN ~6.9% DOCS
- X-Means ~5.8% DOCS
- k-Medoids ~5.8% DOCS
- k-Means (Kernel) ~4.7% DOCS
- k-Means (fast) MEM! ~4.5% DOCS
- Support Vector Clustering ~4.2% DOCS
- Expectation Maximization Clustering ~3.3% DOCS
- Random Clustering ~2.7% DOCS

Correlations (7)

- Correlations Matrix ~71.4% DOCS
- Covariance Matrix ~8.8% DOCS
- ANOVA Matrix ~7.7% DOCS
- Grouped ANOVA ~3.8% DOCS
- Mutual Information Matrix ~3.8% DOCS
- Rainflow Matrix ~1.3% DOCS
- Transition Matrix ~1.1% DOCS

Community



Filters

Column types

- Numerical
- Binary
- Categorical

Target type

- No target
- Numerical
- Binary
- Categorical

Number of columns

- 10s
- 100s
- 1000s
- 10,000s

Number of rows

- 1000s
- 10,000s
- 100,000s
- 1,000,000s

Advanced options

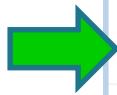
- Updatable
- Can handle missings
- Uses row weights

Applicable Models

[Show all](#)

Predictive (18)

Decision Tree	~19.3%	DOCS
Naive Bayes	~12.6%	DOCS
k-NN	MEM!	~11.2% DOCS
Neural Net	~8.2%	DOCS
Support Vector Machine (LibSVM)	~4.7%	DOCS
Rule Induction	MEM!	~2.8% DOCS
Random Forest	MEM!	~2.5% DOCS
Default Model	~1.5%	DOCS
Linear Discriminant Analysis	~1.5%	DOCS
Naive Bayes (Kernel)	~1.4%	DOCS
Random Tree	~1.1%	DOCS
CHAID	~1.0%	DOCS
Decision Stump	~1.0%	DOCS
Quadratic Discriminant Analysis	~0.3%	DOCS
Regularized Discriminant Analysis	~0.2%	DOCS
Gradient Boosted Trees	~0%	DOCS
Deep Learning	~0%	DOCS
Generalized Linear Model	~0%	DOCS





Auto Model: Prediction Task

Training geeigneter Modelle, teilweise mit Parameteroptimierung

Auto Model

Views: Design Results Auto Model

Information

Results: Classification

This is the final step of Auto Model, where you can inspect the generated models together with other results. The output depends on the data and the choices you made. For example, if you deactivated the calculation of correlations or decision trees, those results will not be displayed. Other results might be shown only for certain types of problems. Lift Charts, for example, are only available for two-class problems.

Please note that the results are calculated in the background. However, you can immediately start to inspect the results as they are completed. You can stop background execution by pressing the Stop button at the bottom. Calculations which are not completed when execution is stopped won't be available. You can go back and make changes after the execution is finished or after you stopped it.

We at RapidMiner do not believe in black boxes. This is why you can always open the process which created the model and all related results. Simply click on a model result and on Open Process at the bottom of the screen. This will show you the process which performs all necessary data preprocessing and model optimization. You can use this process for deploying the model or as a starting point for further optimizations.

We will now discuss the possible results in detail below.

Overview

Accuracy

Model	Accuracy	Run Time
Naive Bayes	69.0%	3 s
Generalized Linear Model	19.1%	18 s
Deep Learning	83.6%	10 min 50 s
Decision Tree	19.1%	2 s
Random Forest	62.0%	16 s
Gradient Boosted Trees	89.2%	27 min 21 s

Runtime (ms)

General

Back Open Process

Results

- General
- Data
- Weights
- Correlations
- Comparison
- Naive Bayes
- Generalized Linear Model
- Deep Learning
- Decision Tree
- Random Forest
- Gradient Boosted Trees

Information

This section shows generic

<new process*> – RapidMiner Studio Educational 8.1.001 @ kilab10n

File Edit Process View Connections Cloud Settings Extensions Help

Views: Design Results Auto Model 08 X All Studio Search

Auto Model

Select Data Select Task Prepare Target Select Inputs Model Types Results

Results

Optimal Parameters Number Of Trees: 100 Maximal Depth: 7

Performance for Parameters

Number of Trees	Maximal Depth	Performance
20	2	0.477
60	2	0.518
100	2	0.533
140	2	0.520
20	4	0.529
60	4	0.584

Number of Trees Maximal Depth Performance

Back Open Process

Information

Results: Classification

This is the final step of Auto Model, where you can inspect the generated models together with other results. The output depends on the data and the choices you made. For example, if you deactivated the calculation of correlations or decision trees, those results will not be displayed. Other results might be shown only for certain types of problems. Lift Charts, for example, are only available for two-class problems.

Please note that the results are calculated in the background. However, you can immediately start to inspect the results as they are completed. You can stop background execution by pressing the Stop button at the bottom. Calculations which are not completed when execution is stopped won't be available. You can go back and make changes after the execution is finished or after you stopped it.

We at RapidMiner do not believe in black boxes. This is why you can always open the process which created the model and all related results. Simply click on a model result and on Open Process at the bottom of the screen. This will show you the process which performs all necessary data preprocessing and model optimization. You can use this process for deploying the model or as a starting point for further optimizations.

We will now discuss the possible results in detail below.

General

This section shows generic

File Edit Process View Connections Cloud Settings Extensions Help

Views: Design Results Auto Model

08 All Studio Search

Auto Model

Select Data Select Task Prepare Target Select Inputs Model Types Results

Results

Deep Learning - Simulator

General

- Data
- Weights
- Correlations

Comparison

- Overview

Naive Bayes

- Model
- Simulator
- Performance

Generalized Linear Model

Deep Learn

- Model
- Simulator**
- Performance

Decision Tree

Random Forest

- Model
- Simulator
- Performance
- Optimal Parameters

Gradient Boosted Trees

Find the Optimal Inputs

Click 'Optimize' to generate optimal inputs for above.

Optimize

Information

Results: Classification

This is the final step of Auto Model, where you can inspect the generated models together with other results. The output depends on the data and the choices you made. For example, if you deactivated the calculation of correlations or decision trees, those results will not be displayed. Other results might be shown only for certain types of problems. Lift Charts, for example, are only available for two-class problems.

Please note that the results are calculated in the background. However, you can immediately start to inspect the results as they are completed. You can stop background execution by pressing the Stop button at the bottom. Calculations which are not completed when execution is stopped won't be available. You can go back and make changes after the execution is finished or after you stopped it.

We at RapidMiner do not believe in black boxes. This is why you can always open the process which created the model and all related results. Simply click on a model result and on Open Process at the bottom of the screen. This will show you the process which performs all necessary data preprocessing and model optimization. You can use this process for deploying the model or as a starting point for further optimizations.

We will now discuss the possible results in detail below.

Most Likely: LAYING

Activity	Probability (%)
WALKING...	~1%
WALKING...	~1%
WALKING	~1%
STANDING	~1%
SITTING	~1%
LAYING	100%

Important Factors for LAYING

Factor	Value	Supports LAYING?
fBodyAccMag-maxInds	-1,0	No
fBodyGyro-maxInds-Y	0,44	Yes
fBodyGyro-maxInds-Z	0,8	Yes
fBodyAcc-maxInds-X	0,41536208	Yes
tBodyAccJerk-entropy()-Z	-1,0	No
fBodyAcc-entropy()-X	-1,0	No
fBodyAccJerk-maxInds-X	0,49376381	Yes

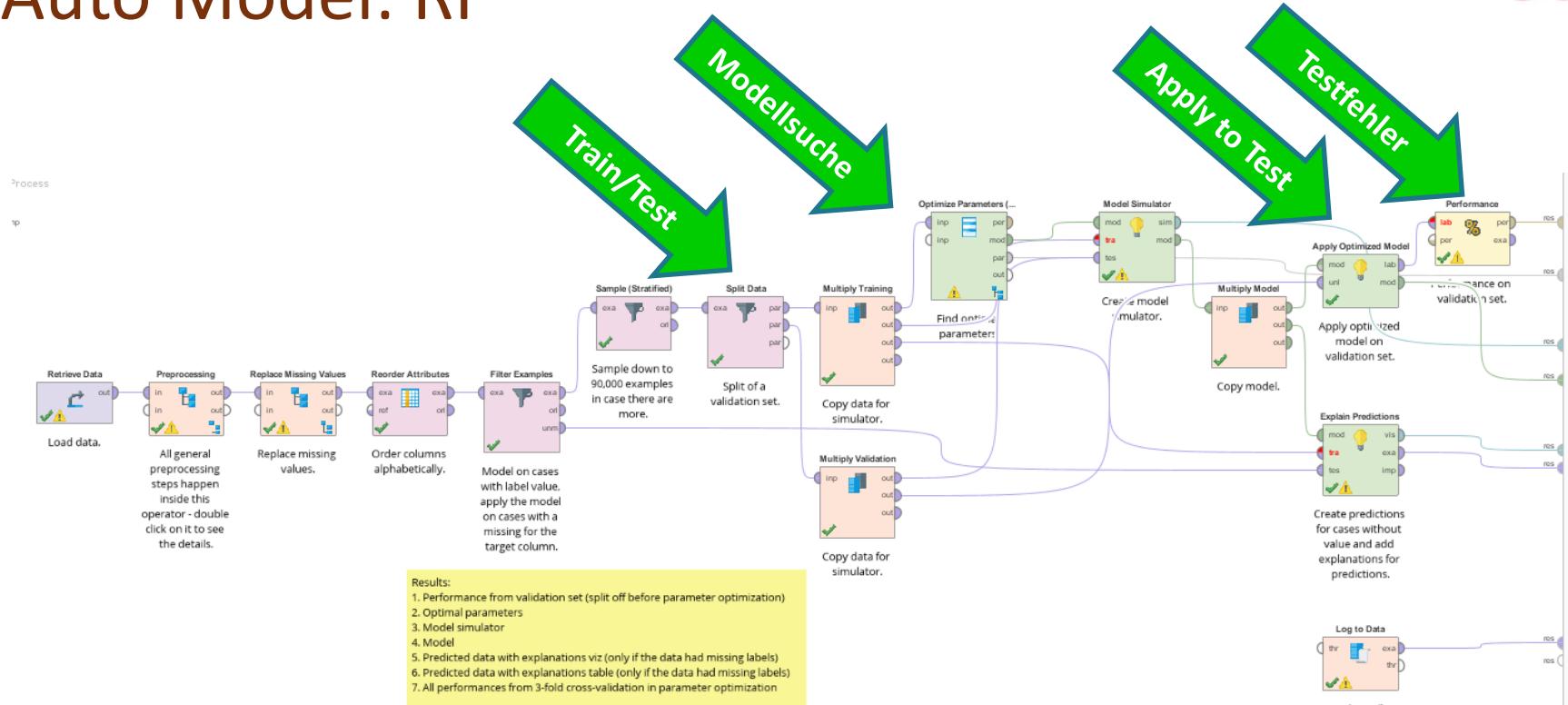
General

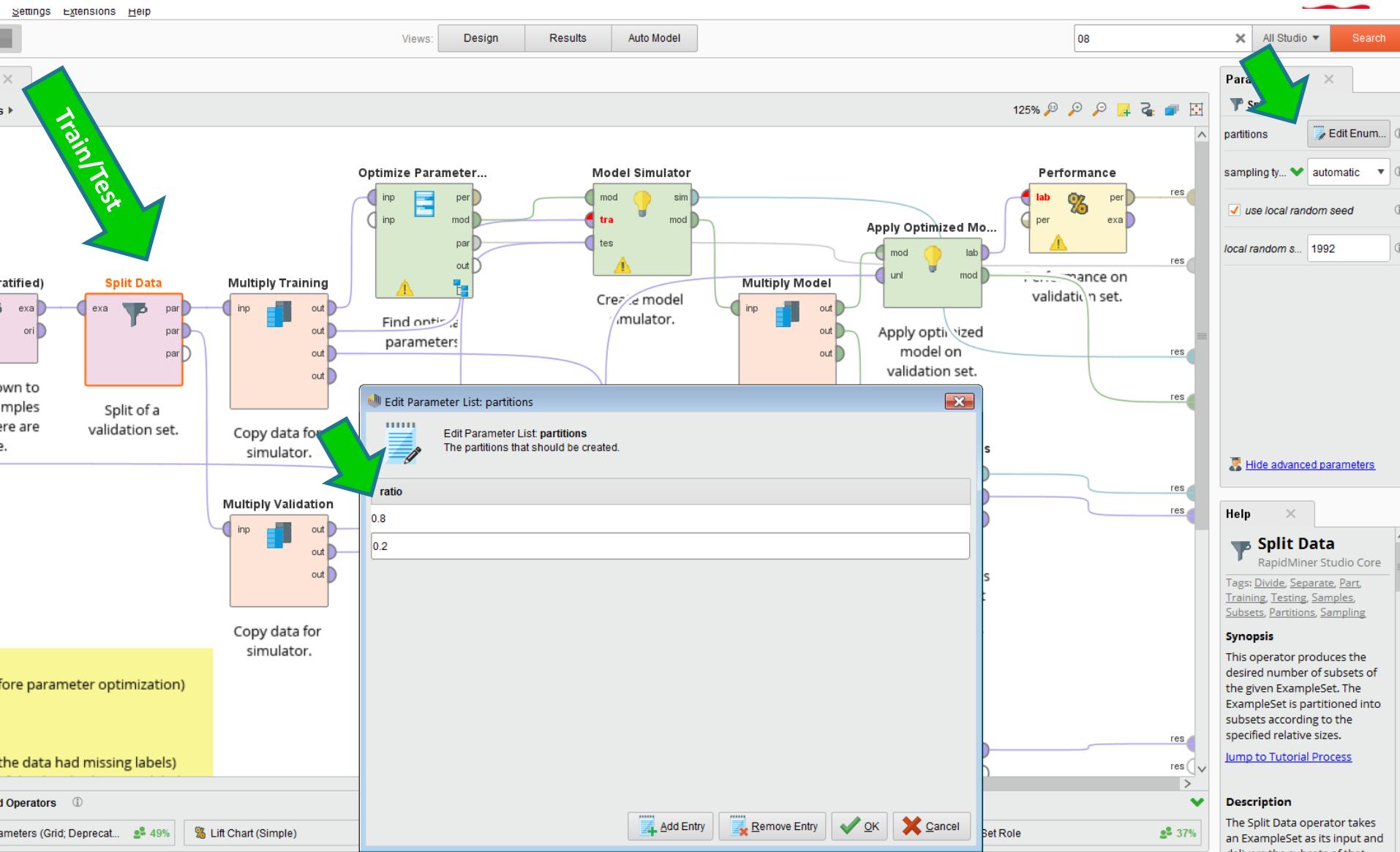
This section shows generic

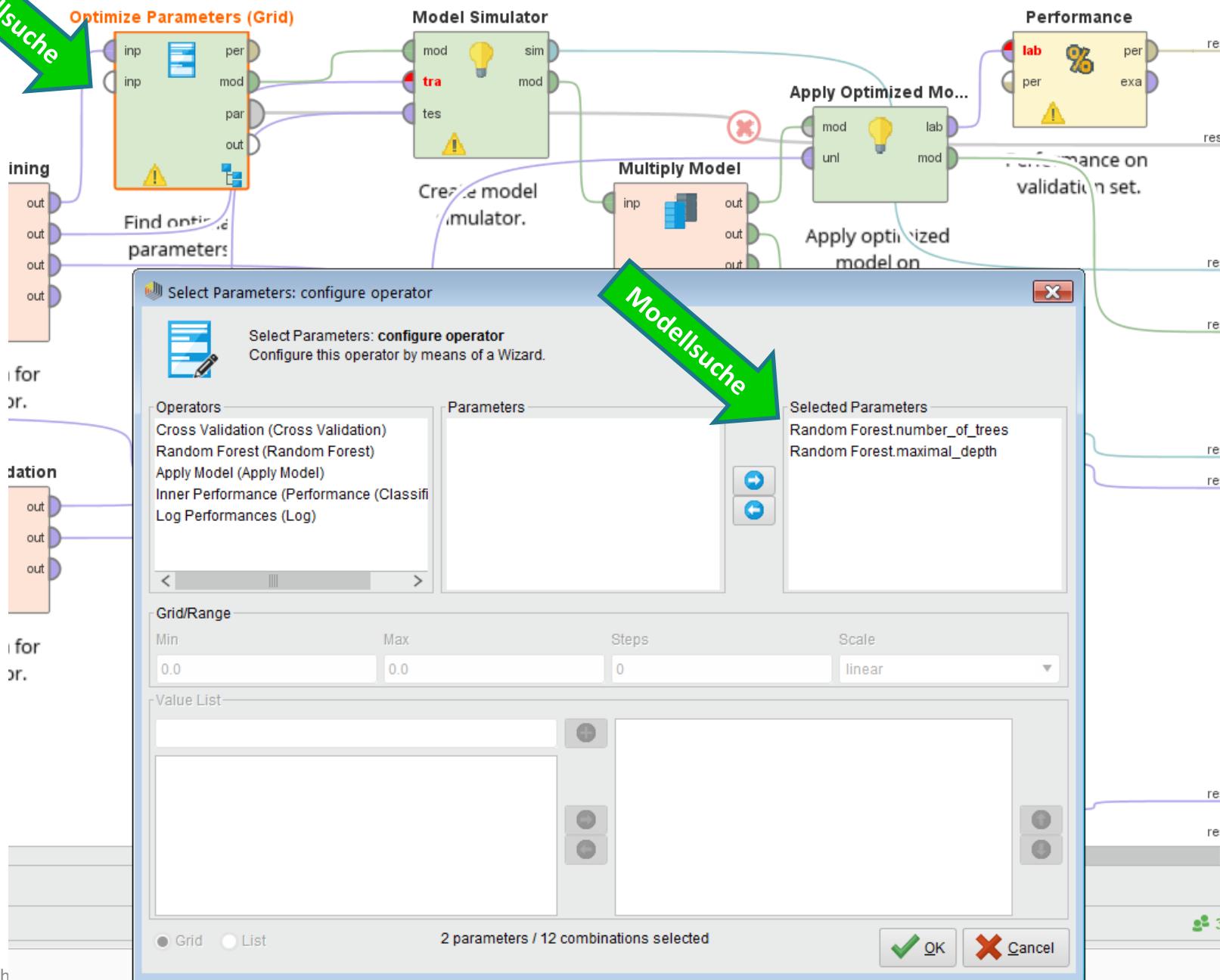




Auto Model: RF









Modellsuche

The screenshot shows a KNIME workflow titled "Optimize Parameters (Grid)". The workflow consists of two main nodes: "Cross Validation" and "Log Performances". The "Cross Validation" node has several input ports: "inp", "exa", "mod", "exa", "tes", "per", "per", and a "mod" port which is connected to the "out" port of the "Log Performances" node. The "Log Performances" node has four output ports labeled "thr", each connected to a "thr" port on the "Cross Validation" node. A large green arrow points from the "Modellsuche" label towards the "Cross Validation" node.

Parameters

Cross Validation

- split on batch attribute
- leave one out

number of folds: 3

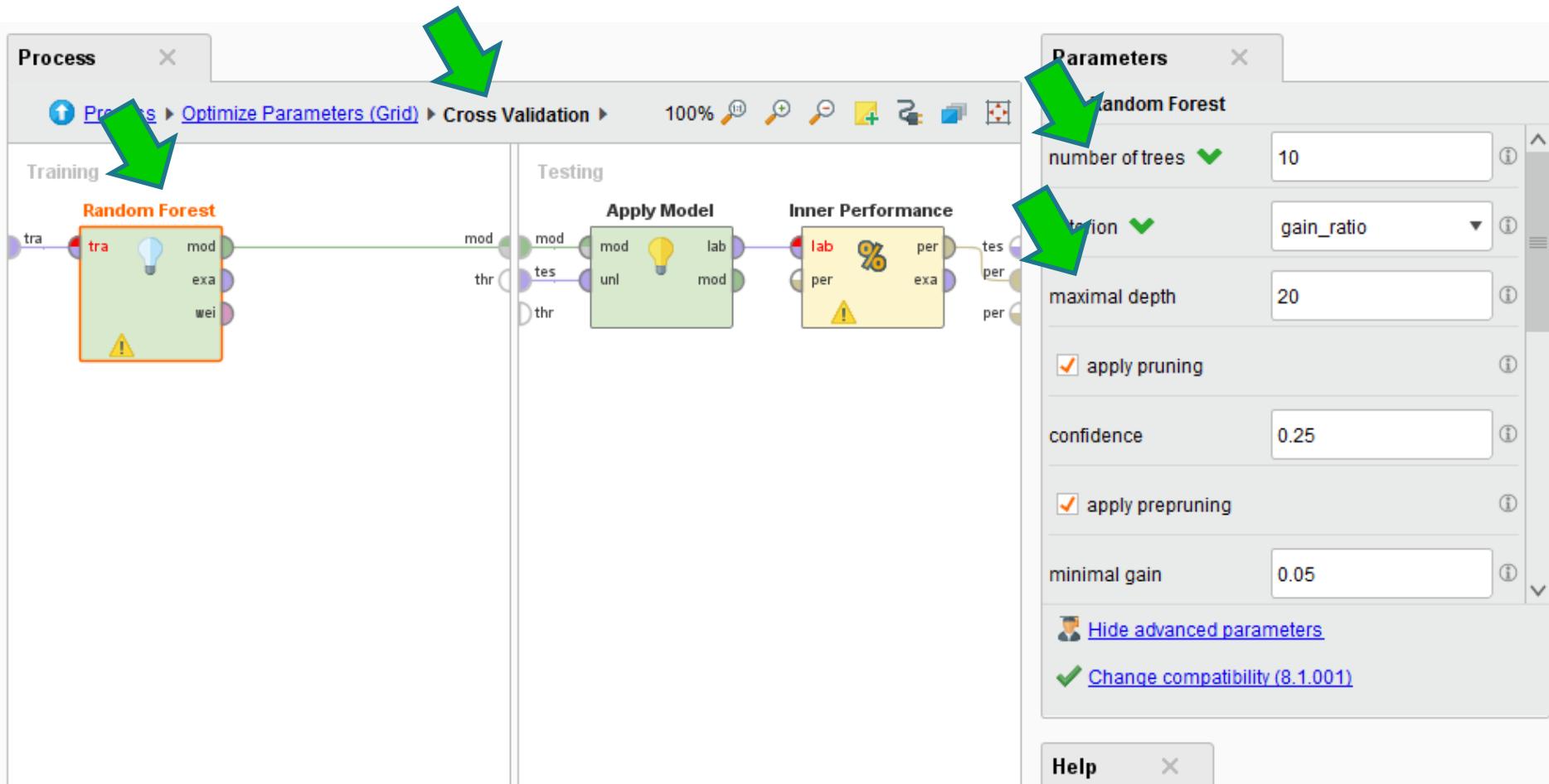
sampling type: automatic

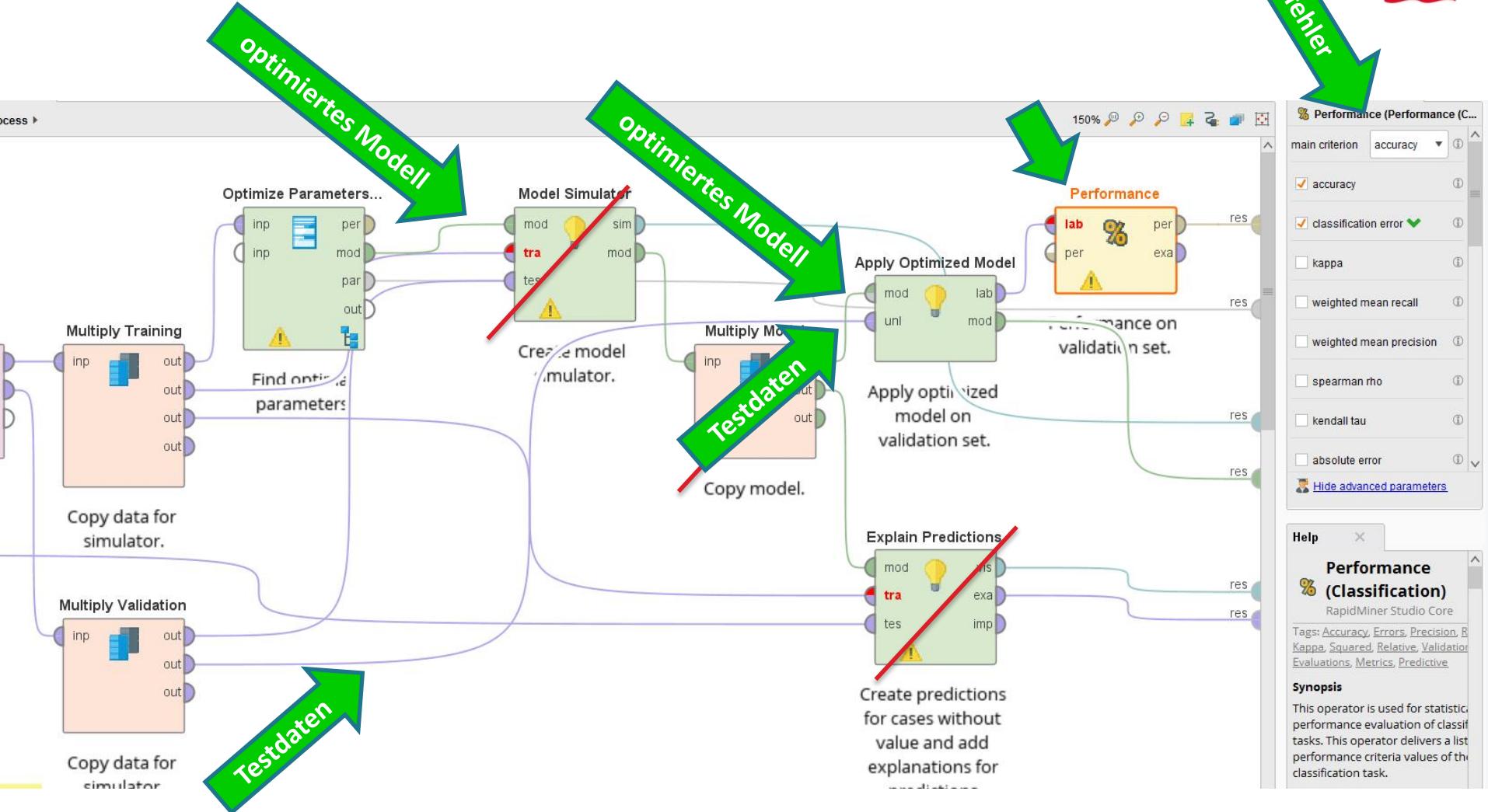
use local random seed

local random seed: 1992

enable parallel execution

[Hide advanced parameters](#)







Auto Model: Prediction Task

Training geeigneter Modelle, teilweise mit Parameteroptimierung

Auto Model

Views: Design Results Auto Model

08 All Studio Search

Information

Results: Classification

This is the final step of Auto Model, where you can inspect the generated models together with other results. The output depends on the data and the choices you made. For example, if you deactivated the calculation of correlations or decision trees, those results will not be displayed. Other results might be shown only for certain types of problems. Lift Charts, for example, are only available for two-class problems.

Please note that the results are calculated in the background. However, you can immediately start to inspect the results as they are completed. You can stop background execution by pressing the Stop button at the bottom. Calculations which are not completed when execution is stopped won't be available. You can back and make changes after the execution is finished or after you stopped it.

We at RapidMiner do not believe in black boxes. This is why you can always open the process which created the model and all related results. Simply click on a model result and on Open Process at the bottom of the screen. This will show you the process which performs all necessary data preprocessing and model optimization. You can use this process for deploying the model or as a starting point for further optimizations.

We will now discuss the possible results in detail below.

General

This section shows generic...

Auto Model

Select Data Select Task Prepare Target Select Inputs Model Types Results

Results

General Comparison Overview Naive Bayes Generalized Linear Model Deep Learning Decision Tree Random Forest Gradient Boosted Trees

Overview

Accuracy

Model	Accuracy
Naive Bayes	69%
Generalized Linear Model	19%
Deep Learning	84%
Decision Tree	19%
Random Forest	62%
Gradient Boosted Trees	89%

Runtime (ms)

Model	Run Time
Naive Bayes	3 s
Generalized Linear Model	18 s
Deep Learning	10 min 50 s
Decision Tree	2 s
Random Forest	16 s
Gradient Boosted Trees	27 min 21 s

Back Open Process



Rückblick

Performance von Klassifikatoren

- Konfusionsmatrix
- Fehlermaße
- Trainingsfehler nicht als Schätzung geeignet
- Bsp. k-Nearest Neighbor
- (Overfitting)
- Fehlerschätzung: Kreuzvalidierung
- Praktisches Vorgehen, 100%-Modell