



**Prof. Dr. Adrian Ulges**

**Empolis Workshop “Machine Learning”**

# **Machine Learning 101**

**Hochschule RheinMain**

Department DCSM (*Design, Computer Science, Media*)

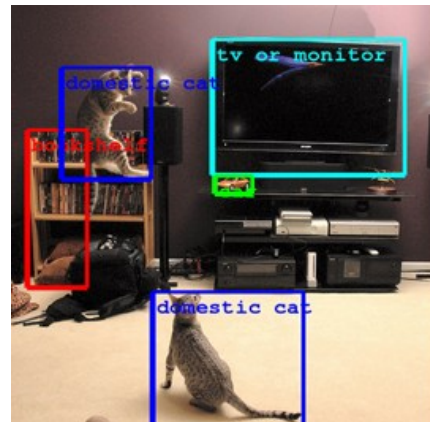
## ► Machine Learning 101

- Grundlagen, Begriffsbildung
- Design-Cycle, Benchmarking, Cross-Validation
- **Python-Beispiel: Titanic**

# ML: Definition

**Prof. Adrian Ulges**  
Fachbereich DCSM / Informatik  
Hochschule RheinMain

- **Arthur Samuel (1959):** „The field of study that gives computers the ability to learn **without being explicitly programmed**“
- **Tom Mitchell (1998):** „A computer program is said to **learn** from **experience  $E$**  with respect to some **task  $T$**  and some performance **measure  $P$** , if its performance on  **$T$** , as measured by  **$P$** , improves with Experience  **$E$** .“



Google		CTR	Search
		About 44,000,000 results (0.04 seconds)	Go to Google.com Advanced search
Everything	What is CTR in google adsense? - Web Development 101	43.2 %	
Images	Clickthrough rate - Wikipedia, the free encyclopedia		
Videos	Clickthrough rate - Wikipedia, the free encyclopedia	30.7 %	
News	CTR - Wikipedia, the free encyclopedia		
Shopping	CTR may stand for Institute, Computing Tabulating Recording Corporation ...	23.3 %	
More	Choose the right Wikipedia, the free encyclopedia		
share OLD	"Choose the right" is a saying or motto among members of the Church of Jesus ...	19.7 %	
ange location	How does search from adsense ...	15.1 %	
s web	CTR - What does CTR stand for? Acronyms and abbreviations by the ...	14.3 %	
ps from Australia	Acceptance, Defiance, CTR, Center, CTR, Center, CTR, Contractor, CTR, Click Through Rate, CTR, Corporate Threat Reduction ...	11.4 %	
y time	CTR Photos - Advance Wedding and Formal Photography	10.1 %	
est	Outstate Wedding and Portrait Photography - Affordable packages for all occasions ...	8.9 %	
at 24 hours	CTR Pacific Pro Ltd - Commercial Bricklayers	8.3 %	
at 2 days	CTR Pacific is a Canberra private commercial brick and blocklaying company ...		
at week	What is CTR in google adsense? - Web Development 101		
at month	10 May 2007 ... was asked the following question this morning by an adsense publisher and I thought some of you might also find it useful so I am sharing ...		
at year	adsense and click measurement - Google		
stom range...	Christ The Redeemer		
re search tools	Laurelton Anglican Church - Parish of Camden Haven		
nothing different	What is CTR in google adsense? - Web Development 101		
& through rate	CTR - What does CTR stand for? Acronyms and abbreviations by the ...		
ster	CTR - What does CTR stand for? Acronyms and abbreviations by the ...		
diothoracic ratio	CTR - What does CTR stand for? Acronyms and abbreviations by the ...		
	CTR - What does CTR stand for? Acronyms and abbreviations by the ...		
	CTR - What does CTR stand for? Acronyms and abbreviations by the ...		
	CTR - What does CTR stand for? Acronyms and abbreviations by the ...		
	CTR - What does CTR stand for? Acronyms and abbreviations by the ...		
	CTR - What does CTR stand for? Acronyms and abbreviations by the ...		
	CTR - What does CTR stand for? Acronyms and abbreviations by the ...		
	CTR - What does CTR stand for? Acronyms and abbreviations by the ...		
	CTR - What does CTR stand for? Acronyms and abbreviations by the ...		
	CTR - What does CTR stand for? Acronyms and abbreviations by the ...		
	CTR - What does CTR stand for? Acronyms and abbreviations by the ...		
	CTR - What does CTR stand for? Acronyms and abbreviations by the ...		
	CTR - What does CTR stand for? Acronyms and abbreviations by the ...		
	CTR - What does CTR stand for? Acronyms and abbreviations by the ...		
	CTR - What does CTR stand for? Acronyms and abbreviations by the ...		
	CTR - What does CTR stand for? Acronyms and abbreviations by the ...		
	CTR - What does CTR stand for? Acronyms and abbreviations by the ...		
	CTR - What does CTR stand for? Acronyms and abbreviations by the ...		
	CTR - What does CTR stand for? Acronyms and abbreviations by the ...		
	CTR - What does CTR stand for? Acronyms and abbreviations by the ...		
	CTR - What does CTR stand for? Acronyms and abbreviations by the ...		
	CTR - What does CTR stand for? Acronyms and abbreviations by the ...		
	CTR - What does CTR stand for? Acronyms and abbreviations by the ...		
	CTR - What does CTR stand for? Acronyms and abbreviations by the ...		
	CTR - What does CTR stand for? Acronyms and abbreviations by the ...		
	CTR - What does CTR stand for? Acronyms and abbreviations by the ...		
	CTR - What does CTR stand for? Acronyms and abbreviations by the ...		
	CTR - What does CTR stand for? Acronyms and abbreviations by the ...		
	CTR - What does CTR stand for? Acronyms and abbreviations by the ...		
	CTR - What does CTR stand for? Acronyms and abbreviations by the ...		
	CTR - What does CTR stand for? Acronyms and abbreviations by the ...		
	CTR - What does CTR stand for? Acronyms and abbreviations by the ...		
	CTR - What does CTR stand for? Acronyms and abbreviations by the ...		
	CTR - What does CTR stand for? Acronyms and abbreviations by the ...		
	CTR - What does CTR stand for? Acronyms and abbreviations by the ...		
	CTR - What does CTR stand for? Acronyms and abbreviations by the ...		
	CTR - What does CTR stand for? Acronyms and abbreviations by the ...		
	CTR - What does CTR stand for? Acronyms and abbreviations by the ...		
	CTR - What does CTR stand for? Acronyms and abbreviations by the ...		
	CTR - What does CTR stand for? Acronyms and abbreviations by the ...		
	CTR - What does CTR stand for? Acronyms and abbreviations by the ...		
	CTR - What does CTR stand for? Acronyms and abbreviations by the ...		
	CTR - What does CTR stand for? Acronyms and abbreviations by the ...		
	CTR - What does CTR stand for? Acronyms and abbreviations by the ...		
	CTR - What does CTR stand for? Acronyms and abbreviations by the ...		
	CTR - What does CTR stand for? Acronyms and abbreviations by the ...		
	CTR - What does CTR stand for? Acronyms and abbreviations by the ...		
	CTR - What does CTR stand for? Acronyms and abbreviations by the ...		
	CTR - What does CTR stand for? Acronyms and abbreviations by the ...		
	CTR - What does CTR stand for? Acronyms and abbreviations by the ...		
	CTR - What does CTR stand for? Acronyms and abbreviations by the ...		
	CTR - What does CTR stand for? Acronyms and abbreviations by the ...		
	CTR - What does CTR stand for? Acronyms and abbreviations by the ...		
	CTR - What does CTR stand for? Acronyms and abbreviations by the ...		
	CTR - What does CTR stand for? Acronyms and abbreviations by the ...		
	CTR - What does CTR stand for? Acronyms and abbreviations by the ...		
	CTR - What does CTR stand for? Acronyms and abbreviations by the ...		
	CTR - What does CTR stand for? Acronyms and abbreviations by the ...		
	CTR - What does CTR stand for? Acronyms and abbreviations by the ...		
	CTR - What does CTR stand for? Acronyms and abbreviations by the ...		
	CTR - What does CTR stand for? Acronyms and abbreviations by the ...		
	CTR - What does CTR stand for? Acronyms and abbreviations by the ...		
	CTR - What does CTR stand for? Acronyms and abbreviations by the ...		
	CTR - What does CTR stand for? Acronyms and abbreviations by the ...		
	CTR - What does CTR stand for? Acronyms and abbreviations by the ...		
	CTR - What does CTR stand for? Acronyms and abbreviations by the ...		
	CTR - What does CTR stand for? Acronyms and abbreviations by the ...		
	CTR - What does CTR stand for? Acronyms and abbreviations by the ...		
	CTR - What does CTR stand for? Acronyms and abbreviations by the ...		
	CTR - What does CTR stand for? Acronyms and abbreviations by the ...		
	CTR - What does CTR stand for? Acronyms and abbreviations by the ...		
	CTR - What does CTR stand for? Acronyms and abbreviations by the ...		
	CTR - What does CTR stand for? Acronyms and abbreviations by the ...		
	CTR - What does CTR stand for? Acronyms and abbreviations by the ...		
	CTR - What does CTR stand for? Acronyms and abbreviations by the ...		
	CTR - What does CTR stand for? Acronyms and abbreviations by the ...		
	CTR - What does CTR stand for? Acronyms and abbreviations by the ...		
	CTR - What does CTR stand for? Acronyms and abbreviations by the ...		
	CTR - What does CTR stand for? Acronyms and abbreviations by the ...		
	CTR - What does CTR stand for? Acronyms and abbreviations by the ...		
	CTR - What does CTR stand for? Acronyms and abbreviations by the ...		
	CTR - What does CTR stand for? Acronyms and abbreviations by the ...		
	CTR - What does CTR stand for? Acronyms and abbreviations by the ...		
	CTR - What does CTR stand for? Acronyms and abbreviations by the ...		
	CTR - What does CTR stand for? Acronyms and abbreviations by the ...		
	CTR - What does CTR stand for? Acronyms and abbreviations by the ...		
	CTR - What does CTR stand for? Acronyms and abbreviations by the ...		
	CTR - What does CTR stand for? Acronyms and abbreviations by the ...		
	CTR - What does CTR stand for? Acronyms and abbreviations by the ...		
	CTR - What does CTR stand for? Acronyms and abbreviations by the ...		
	CTR - What does CTR stand for? Acronyms and abbreviations by the ...		
	CTR - What does CTR stand for? Acronyms and abbreviations by the ...		
	CTR - What does CTR stand for? Acronyms and abbreviations by the ...		

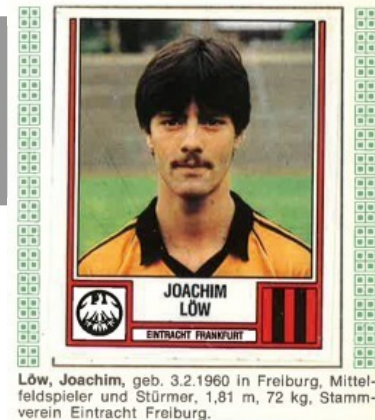
TITLE: Impact of position on clickthrough rate in search results  
Segment: Australian Retail/e-Commerce Traffic, Year: 2011  
Credits: <http://dejanseo.com.au>

# ML: Beispiele

- **Large-scale hierarchical Text Classification**  
*(categorize Wikipedia articles into one of 325,056 categories)*
- **Sentiment Analysis on Movie Reviews**
- **Job Recommendation**  
*(predict which jobs users will apply to)*
- **Axa Driver Telematics challenge**  
*(given GPS routes, identify a cars driver)*
- **Flight Quest**  
*(optimize flight routes based on wheather and traffic)*
- **TFI Restaurant Revenue Prediction**  
*(predict annual sales of restaurants)*
- **Algorithmic Trading** challenge
- **Whale Detection** Challenge
- **Merck molecular Activity** Challenge
- **Psychopathy Prediction** based on Tweets
- Forecasting **ESC Votings**

# ML: Eingabedaten

- Die **Eingabedaten** für Lernverfahren sind üblicher Weise aus **Samples** und **Labels**.
- **Samples**  $\mathbf{x}_1, \dots, \mathbf{x}_n$ : beschreiben Objekte der Welt (*Bilder / Dokumente / Personen / Queries / Situationen / ...*)
- Samples bestehen üblicher Weise aus **Merkmalen**
- Typen von Merkmalen ...
  - **nominal** (*Kategorien ohne Ordnung, z.B. Farben*)
  - **ordinal** (*Kategorien mit Ordnung, z.B. Häufigkeiten*)
  - **Intervallskaliert** (*Zahlen mit Abstandsmaß, z.B. Kalendertage*)
  - **verhältnisskaliert** (*Zahlen mit Nullpunkt, z.B. Gewicht*)
  - **Fehlend!**
- Alternativen (z.B. objekt-orientiert, Graphen): *Heute nicht*
- **Labels**  $\mathbf{y}_1, \dots, \mathbf{y}_n$ : repräsentieren eine Aussage über Objekte
  - Eine **Zuordnung** zu einer Klasse / Gruppe
  - Ein reellwertiges **Attribut** (z.B. ein Preis)

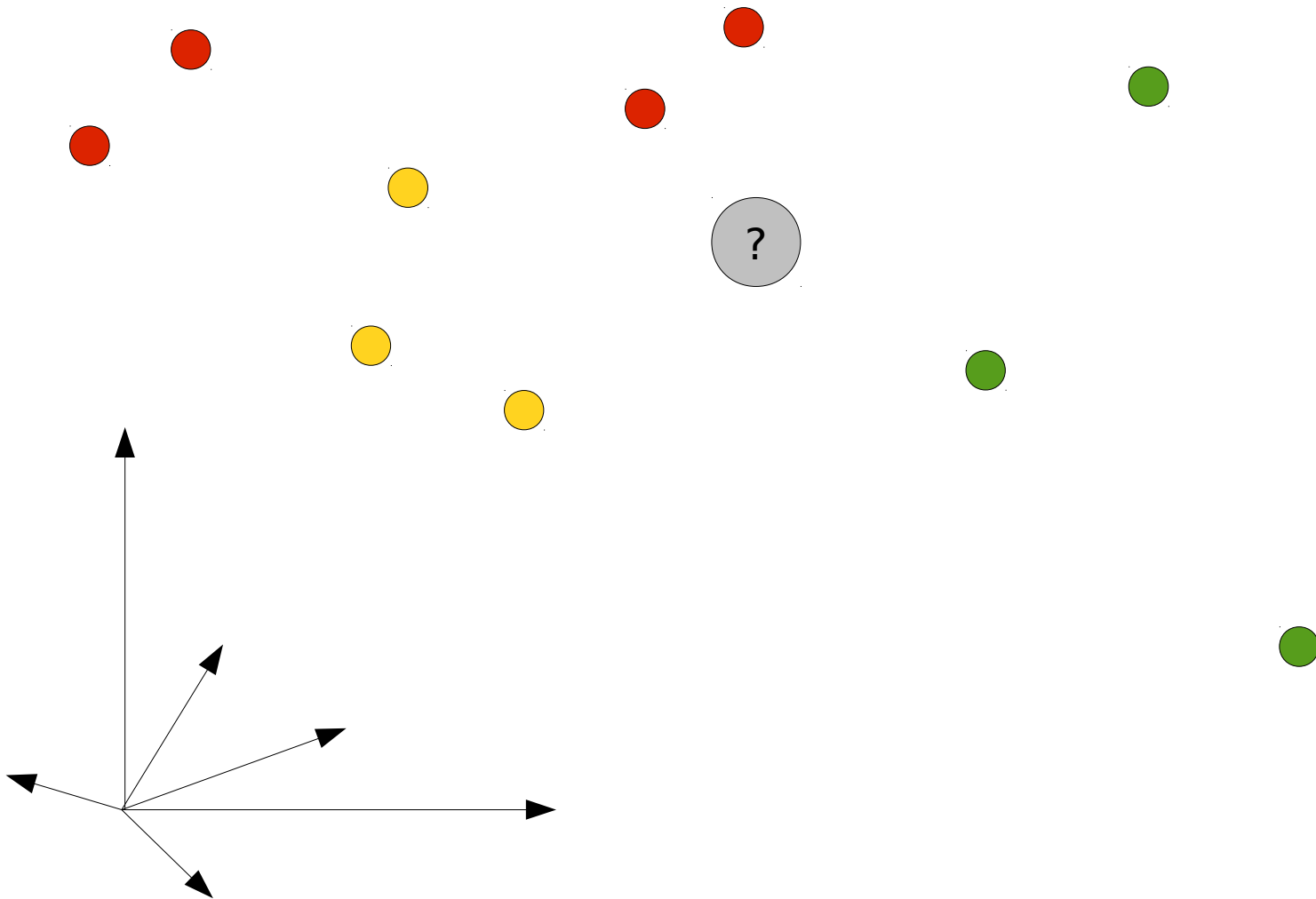


# ML: Eingabedaten

Prof. Adrian Ulges

Fachbereich DCSM / Informatik  
Hochschule RheinMain

	A	B	C	D	E	F	G	H	I	J	K	L
1	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin	Embarked
2	1	0	3	Braund, Mr. Owen Harris	male	22	1	0	A/5 21171	7.25		S
3	2	1	1	Cumings, Mrs. John Bradley (Florence Briggs Thayer)	female	38	1	0	PC 17599	71.2833	C85	C
4	3	1	3	Heikinen, Miss. Laina	female	26	0	0	STON/O2. 3101282	7925		S
5	4	1	1	Futrelle, Mrs. Jacques Heath (Lily May Peel)	female	35	1	0		113803 53.1	C123	S
6	5	0	3	Allen, Mr. William Henry	male	35	0	0		373450 8.05		S
7	6	0	3	Moran, Mr. James	male		0	0		330877 8.4583		Q
8	7	0	1	McCarthy, Mr. Timothy J	male	54	0	0		17463 51.8625	E46	S
9	8	0	3	Palsson, Master. Gosta Leonard	male	2	3	1		349909 21075		S
10	9	1	3	Johnson, Mrs. Oscar W (Elisabeth Vilhelmina Berg)	female	27	0	2		347742 11.1333		S
11	10	1	2	Nasser, Mrs. Nicholas (Adele Achem)	female	14	1	0		237736 30.0708		C
12	11	1	3	Sandstrom, Miss. Marguerite Rut	female	4	1	1	PP 9549	16.7	G6	S
13	12	1	1	Bonnell, Miss. Elizabeth	female	58	0	0		113783 26.55	C103	S
14	13	0	3	Saunderscock, Mr. William Henry	male	20	0	0	A/5. 2151	8.05		S
15	14	0	3	Andersson, Mr. Anders Johan	male	39	1	5		347082 31275		S
16	15	0	3	Vestrom, Miss. Hulda Amanda Adolfina	female	14	0	0		350406 7.8542		S
17	16	1	2	Hewlett, Mrs. (Mary D Kingcome)	female	55	0	0		248706 16		S
18	17	0	3	Rice, Master. Eugene	male	2	4	1		382652 29125		Q
19	18	1	2	Williams, Mr. Charles Eugene	male		0	0		244373 13		S
20	19	0	3	Vander Planke, Mrs. Julius (Emelia Maria Vandemoortele)	female	31	1	0		345763 18		S
21	20	1	3	Masselmani, Mrs. Fatima	female		0	0		2649 7225		C
22	21	0	2	Fynney, Mr. Joseph J	male	35	0	0		239865 26		S
23	22	1	2	Beesley, Mr. Lawrence	male	34	0	0		248698 13	D56	S
24	23	1	3	McGowan, Miss. Anna "Annie"	female	15	0	0		330923 8.0292		Q
25	24	1	1	Sloper, Mr. William Thompson	male	28	0	0		113788 35.5	A6	S
26	25	0	3	Palsson, Miss. Torborg Danira	female	8	3	1		349909 21075		S
27	26	1	3	Asplund, Mrs. Carl Oscar (Selma Augusta Emilia Johansson)	female	38	1	5		347077 31.3875		S
28	27	0	3	Emir, Mr. Farred Chehab	male		0	0		2631 7225		C
29	28	0	1	Fortune, Mr. Charles Alexander	male	19	3	2		19950 263	C23 C25 C27	S
30	29	1	3	O'Dwyer, Miss. Ellen "Nellie"	female		0	0		330959 7.8792		Q
31	30	0	3	Todoroff, Mr. Lalio	male		0	0		349216 7.8958		S
32	31	0	1	Uruchurtu, Don. Manuel E	male	40	0	0	PC 17601	27.7208		C
33	32	1	1	Spencer, Mrs. William Augustus (Marie Eugenie)	female		1	0	PC 17569	146.5208	B78	C
34	33	1	3	Glynn, Miss. Mary Agatha	female		0	0		335677 7.75		Q
35	34	0	2	Wheadon, Mr. Edward H	male	66	0	0	C.A. 24579	10.5		S
36	35	0	1	Meyer, Mr. Edgar Joseph	male	28	1	0	PC 17604	82.1708		C
37	36	0	1	Holverson, Mr. Alexander Oskar	male	42	1	0		113789 52		S
38	37	1	3	Mamee, Mr. Hanna	male		0	0		2677 7.2292		C
39	38	0	3	Cann, Mr. Ernest Charles	male	21	0	0	A/5. 2152	8.05		S
40	39	0	3	Vander Planke, Miss. Augusta Maria	female	18	2	0		345764 18		S
41	40	1	3	Nicola-Yarred, Miss. Jamila	female	14	1	0		2651 11.2417		C
42	41	0	3	Ahlin, Mrs. Johan (Johanna Persdotter Larsson)	female	40	1	0		7546 9475		S
43	42	0	2	Turpin, Mrs. William John Robert (Dorothy Ann Wonnacott)	female	27	1	0		11668 21		S
44	43	0	3	Kraeff, Mr. Theodor	male		0	0		349253 7.8958		C



Überwachtes Lernen

Unüberwachtes Lernen

Halbüberwachtes Lernen

Transduktives vs.  
Induktives Lernen

Reinforcement  
Learning

Active  
Learning

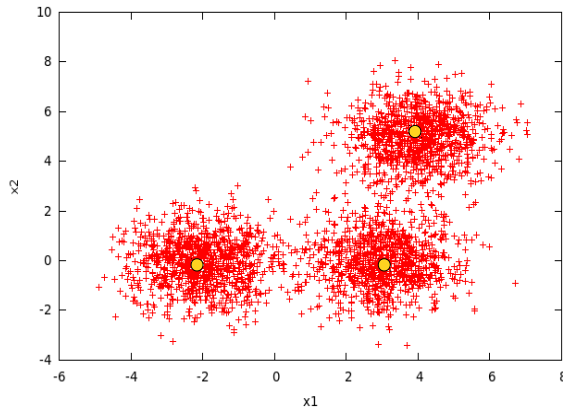
- **Überwachtes Lernen**
  - **Gegeben:** Trainingssamples  $\mathbf{x}_1, \dots, \mathbf{x}_n$  und Labels  $\mathbf{y}_1, \dots, \mathbf{y}_n$
  - **Aufgabe:** Ordne Samples  $\mathbf{x}$  die richtige Klasse  $\mathbf{y}$  zu
- **Unüberwachtes Lernen**
  - **Gegeben:** Trainingssamples  $\mathbf{x}_1, \dots, \mathbf{x}_n$  **ohne Labels**
  - **Wichtig**, denn wir besitzen oft viele Daten aber keine Labels!
  - **Aufgabe:** Inferenz der **Struktur der Daten**
- **Halbüberwachtes Lernen (Mischform)**
  - **Gegeben:** Trainingssamples  $\mathbf{x}_1, \dots, \mathbf{x}_n$ , **manche** mit Labels
- **Weitere Unterscheidungen:**
  - Transduktives vs. Induktives Lernen
  - Reinforcement learning, **active learning**



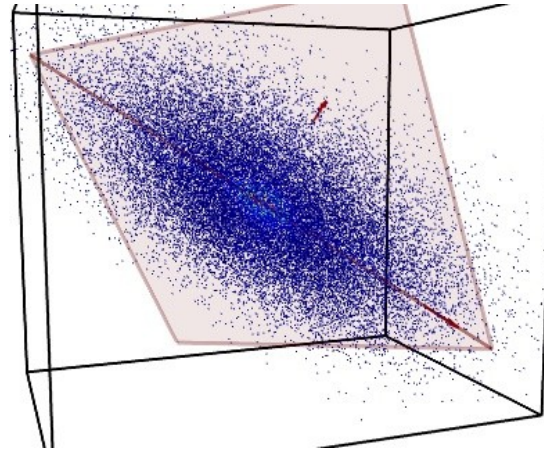
# Unüberwachtes Lernen

Prof. Adrian Ulges

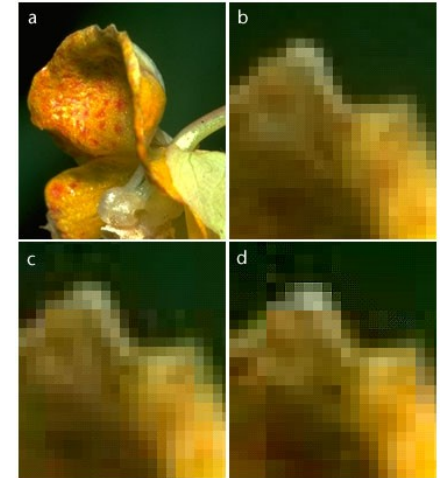
Fachbereich DCSM / Informatik  
Hochschule RheinMain



**Clustering:** Teile die Daten in kohärente Gruppen ein



**Dimensionality Reduction:**  
Komprimiere die Daten



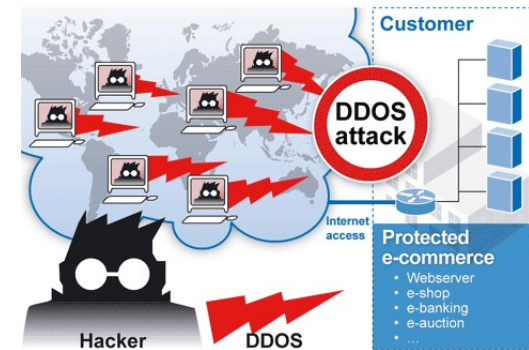
amazon.com

Recommended for You

Amazon.com has new recommendations for you based on [items](#) you purchased or told us you own.



**Itemset Mining:** Finde häufige Substrukturen in den Daten



**Anomaly Detection:** Finde Outlier / Ausreißer in den Daten

- **Zum überwachtem Lernen:** Hier unterscheiden wir zwei Fragestellungen: **Regression** und **Klassifikation**
- **Klassifikation**
  - Die Labels  $y_1, \dots, y_n$  sind **nominale Variablen** (d.h., **Klassen**)
  - Aufgabe: Ordne Sample  $x$  einer Klasse  $y$  zu
  - **Beispiel:** SPAM-Filter, Optical Character Recognition (OCR)
- **Regression**
  - Die Labels  $y_1, \dots, y_n$  sind **reellwertige Variablen**
  - Aufgabe: Ordne einem Sample  $x$  einen Zahlenwert zu
  - **Beispiel:** Fahrtdauer-Prognose im Routenplaner

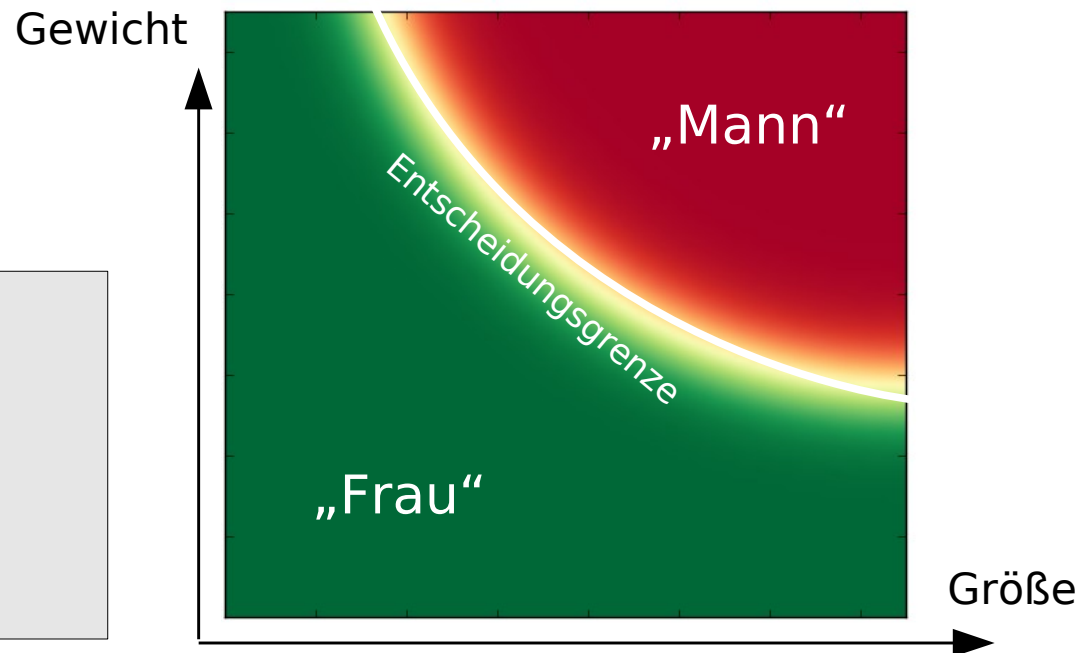
- **Szenario:** *Klassifikation* (ordne ein **Sample x** einer **Klasse y** zu)
- **Ziel:** Berechne  $P(y|x)$  für **jede Klasse y**. Hieraus ergeben sich *Entscheidungsregionen* und *-Grenzen*.

## Dummy-Beispiel

$\mathbf{x} = (\text{Größe}, \text{Gewicht})$

$P(C=\text{„Frau“}|\mathbf{x})$  ist niedrig

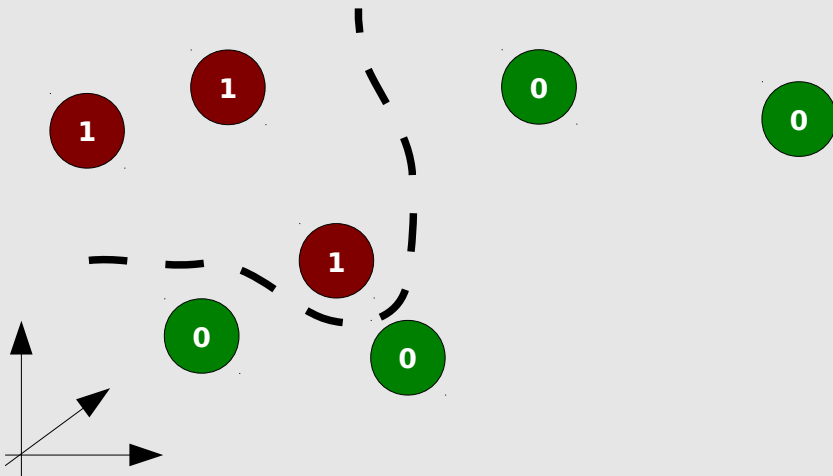
$P(C=\text{„Frau“}|\mathbf{x})$  ist hoch



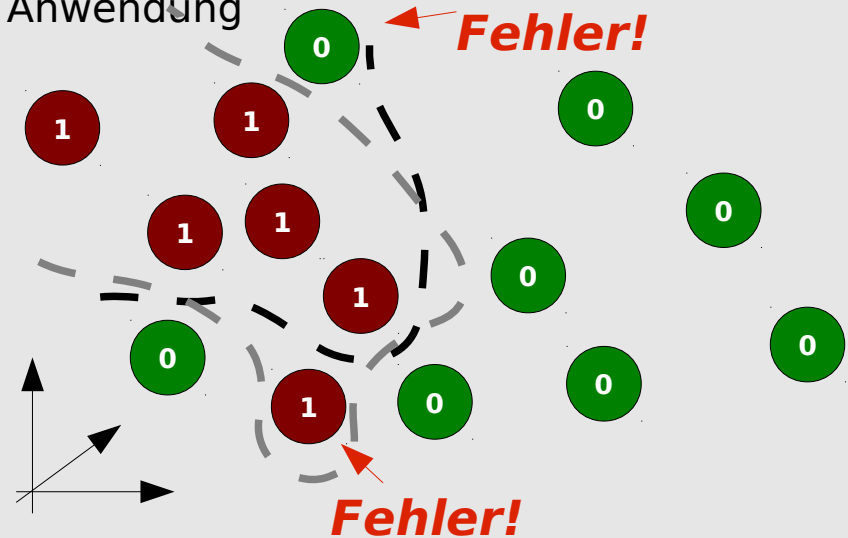
- **Ansatz:** Wir bestimmen für jede Klasse eine **Klassifikatorfunktion**  $\Phi_c^\theta: \mathbb{R}^n \rightarrow [0,1]$
- Gegeben ein Objekt  $\mathbf{x}$ , verwenden wir  $\Phi_c^\theta(\mathbf{x})$  als **Approximation** für  $P(\mathbf{c}|\mathbf{x})$ , um  $\mathbf{x}$  zu klassifizieren: Berechne  $\Phi_c^\theta(\mathbf{x})$  für **alle Klassen**  $\mathbf{c}$  und wähle die Klasse mit **maximalem Wert**.
- Die genaue Form von  $\Phi_c^\theta(\mathbf{x})$  wird von den **Parametern**  $\theta$  bestimmt. Die Parameter werden durch **Training** ermittelt
- Hierzu verwenden wir eine **Trainingsmenge**:
  - Beispiele (**Samples**)  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n \in \mathbb{R}^n$
  - Klassenzuordnungen (**Labels**)  $\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_n \in \{1, \dots, C\}$

Die Fähigkeit eines Klassifikators, von den Trainingsdaten auf (neue) Testdaten zu schließen, bezeichnen wir als **Generalisierung**.

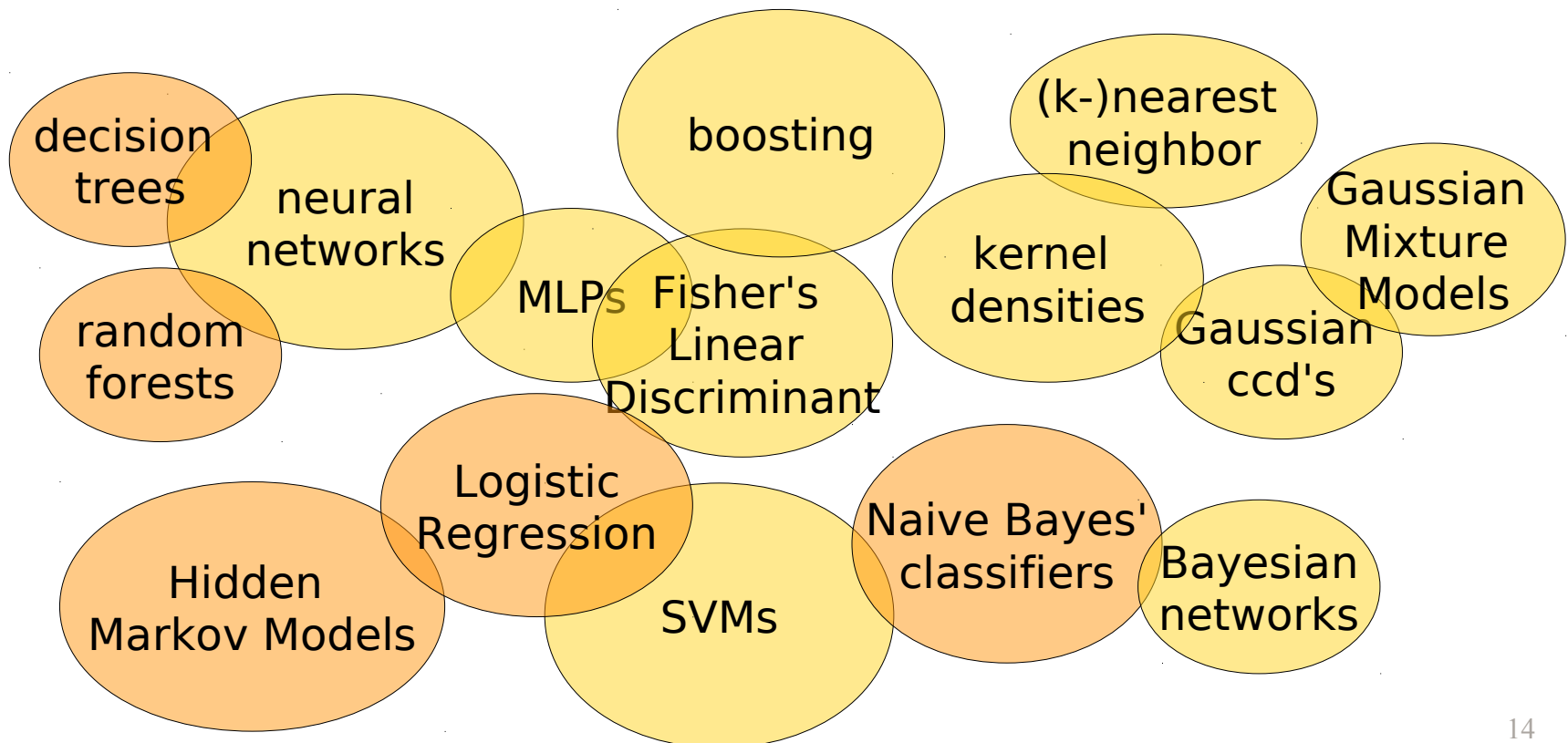
Training



Anwendung

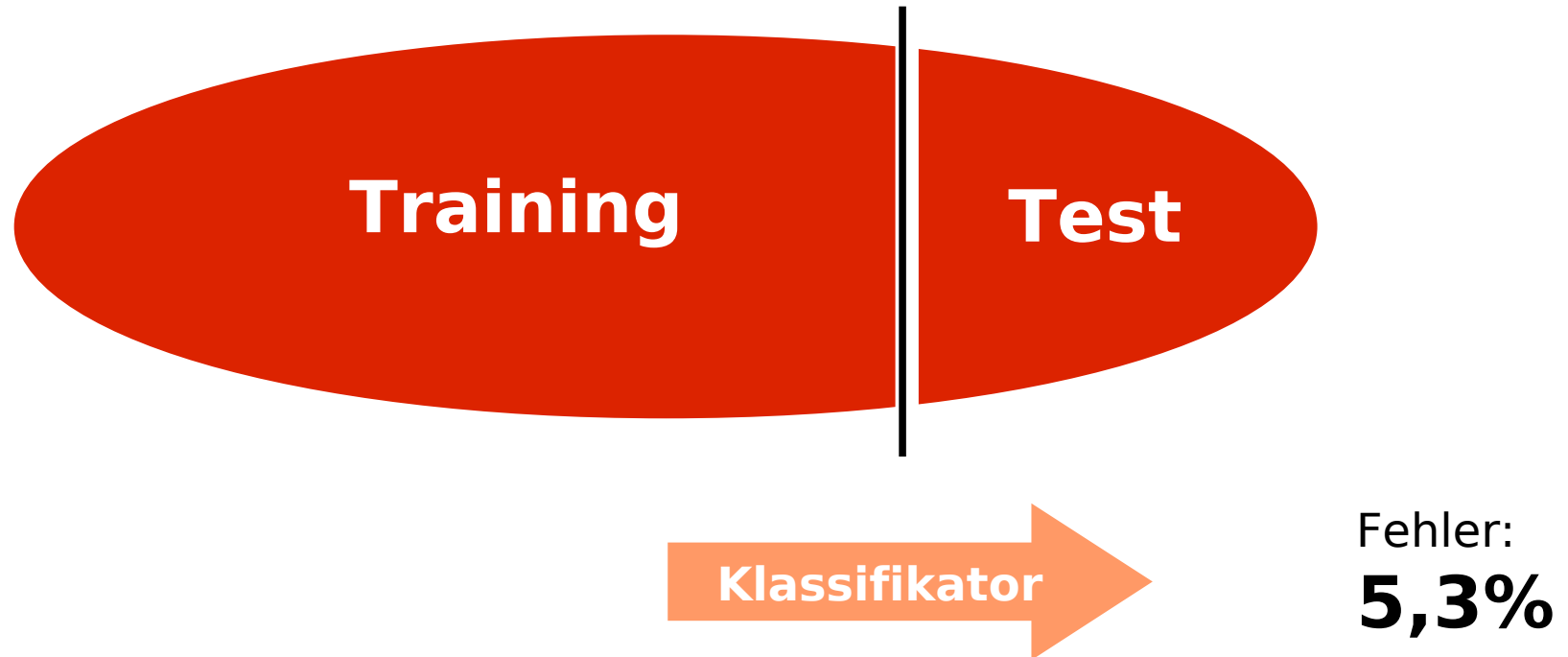


- **Zahlreiche Methoden** des maschinellen Lernens beschäftigen sich mit der Approximation von  $\mathbf{P}(\mathbf{y}|\mathbf{x})$  ...
- Den **universell besten Klassifikator** gibt es nicht  
(→ *No-free-lunch-Theorem*)





- **Design eines ML-Systems = *iterative Auswahl*** von Daten / Merkmalen / Modellen / Parametern
- **Schlüsseltreiber: *Benchmarking***



- Wir **unterteilen** die Datenmenge in Trainings- und Testdaten
- Wir trainieren auf den **Trainingsdaten** und benchmarken auf den **Testdaten**
- **Todsünde: „Testing on the Training Data“**<sup>16</sup>

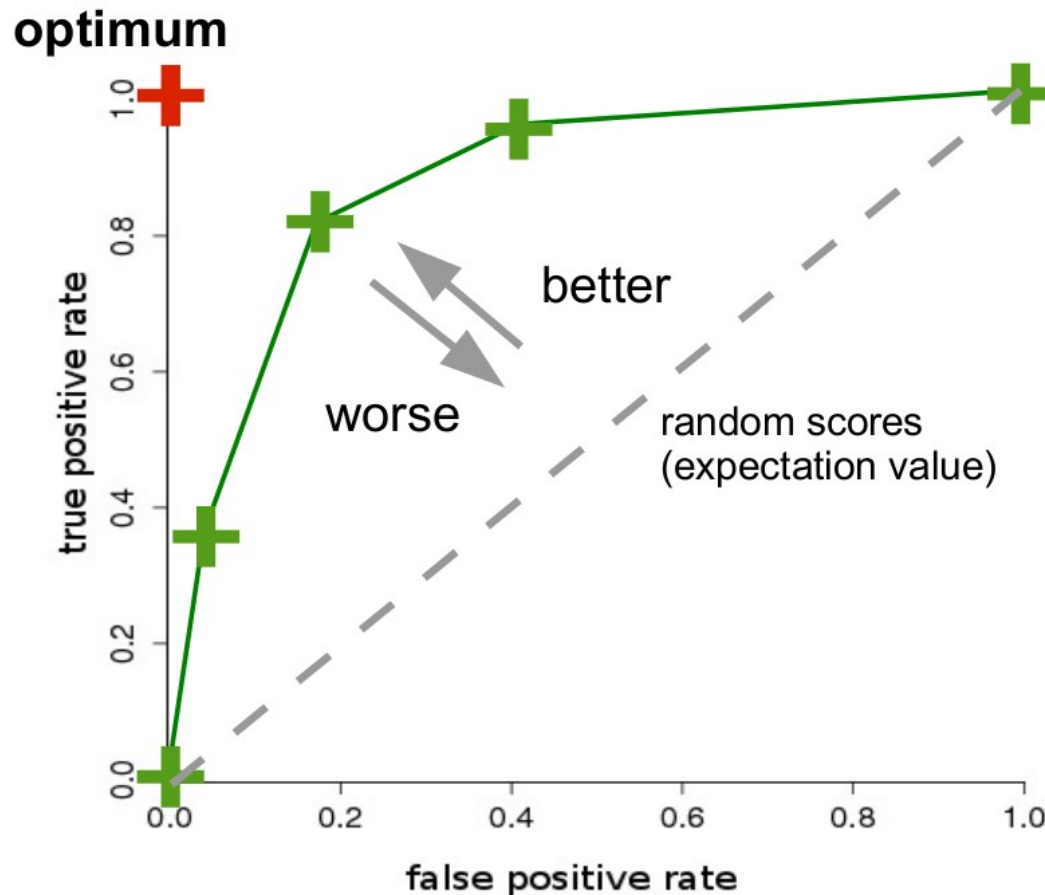


		Estimated Emotion							Emotion Recog. Rate
		Anger	Boredom	Disgust	Fear	Happiness	Sadness	Neutral	
True Emotion	Anger	19	0	2	0	3	0	0	79.2%
	Boredom	1	8	1	1	0	1	7	42.1%
	Disgust	0	1	6	0	1	0	3	54.5%
	Fear	1	3	2	7	2	0	1	43.8%
	Happiness	3	0	3	2	5	0	2	33.3%
	Sadness	0	0	0	0	0	14	0	100.0%
	Neutral	0	5	1	0	0	0	13	68.4%
HMM Recog. Rate		79.2%	47.1%	40.0%	70.0%	45.5%	93.3%	50.0%	

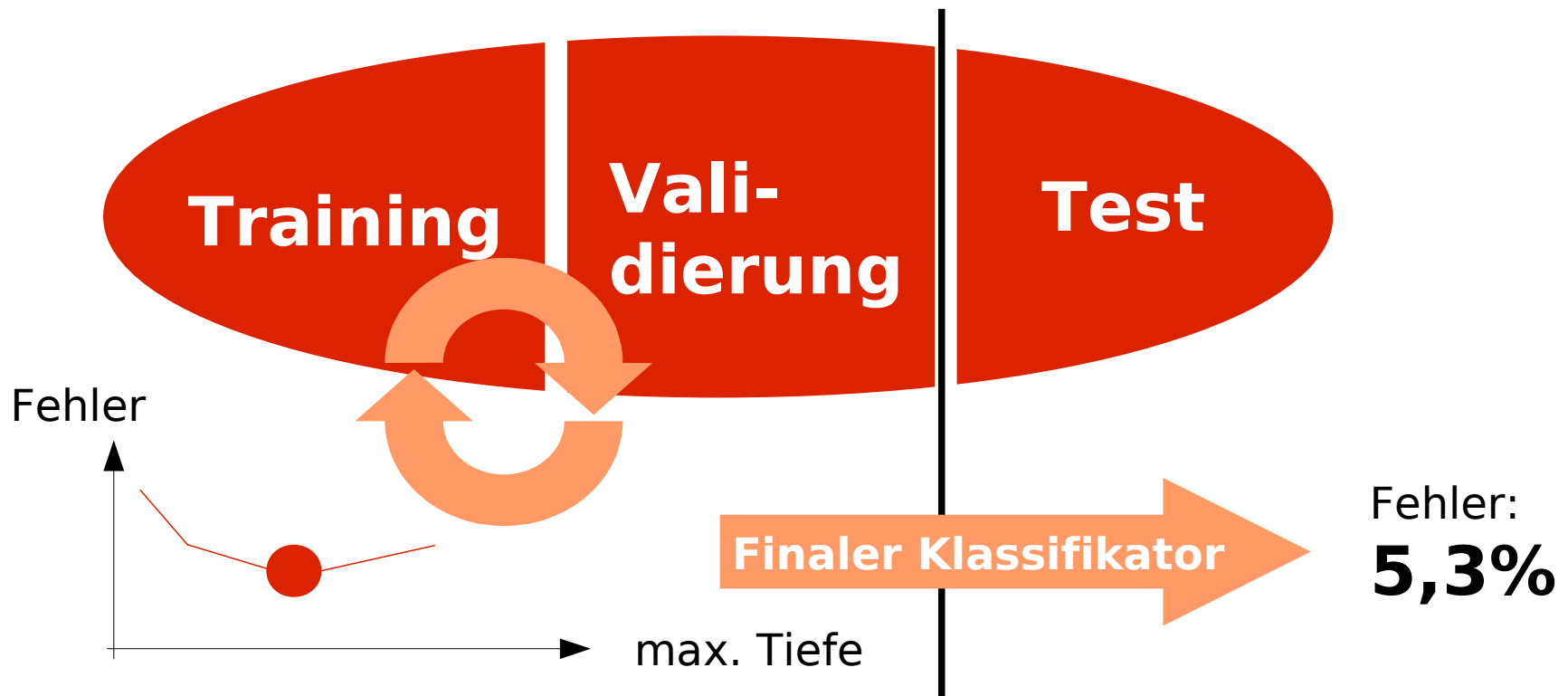
**Bsp.:** ***Confusion-Matrix*** (welche Klassen werden besonders häufig verwechselt?)

- Oft unterscheiden wir nur zwischen **zwei Klassen** (1 und 0)
- Wir unterscheiden vier Arten von Ergebnissen
  - *true/false positive*
  - *true/false negatives*
- Häufig ermitteln Klassifikatoren keine absolute Entscheidung, sondern einen **Score**  $s(x)$
- Wir wählen eine **Schwelle**  $T$  und entscheiden uns für Klasse 1, falls  $s(x) \geq T$
- $T$  bestimmt die **Sensitivität** des Klassifikators

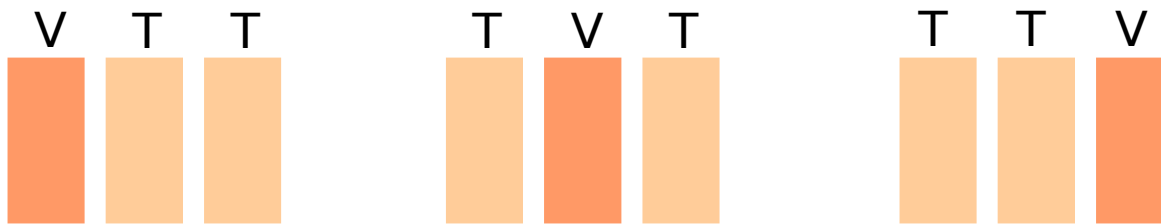
- Wir **variieren T** und erhalten eine Kurve, die sogenannte **ROC-Kurve**



- Einige **Parameter** des Klassifikators werden **gelernt**, andere (freie) Parameter werden **manuell** („trial-and-error“ / grid search) gesetzt
- **Beispiel Entscheidungsbaum**: maximale Tiefe
- **Notwendig**: Training → Validieren → Testen



- Sind die Trainingsdaten klein, teilt man sie in Teile (oder „**Folds**“) und trainiert/validiert **mehrfach**
- Die Ergebnisse werden über die Folds **gemittelt**

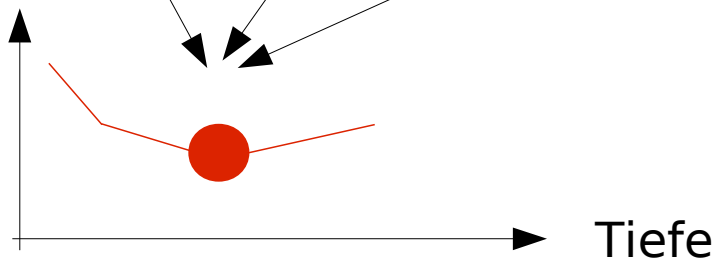


Fehler:  
**5,3%**

Fehler:  
**6,2%**

Fehler:  
**4,8%**

Fehler



Training +  
Validierung

Finaler Klassifikator