

How to Build an Artificial Coder? – Einführung in die automatisierte Inhaltsanalyse mit Machine Learning



11.02.2021 – NapoKo-Methodenworkshop – Jahrestagung DGPUK PolKomm

Dozentin: anke.stoll@hhu.de

Dank an: napoko.de

Illustration: weneedtotalk.ai

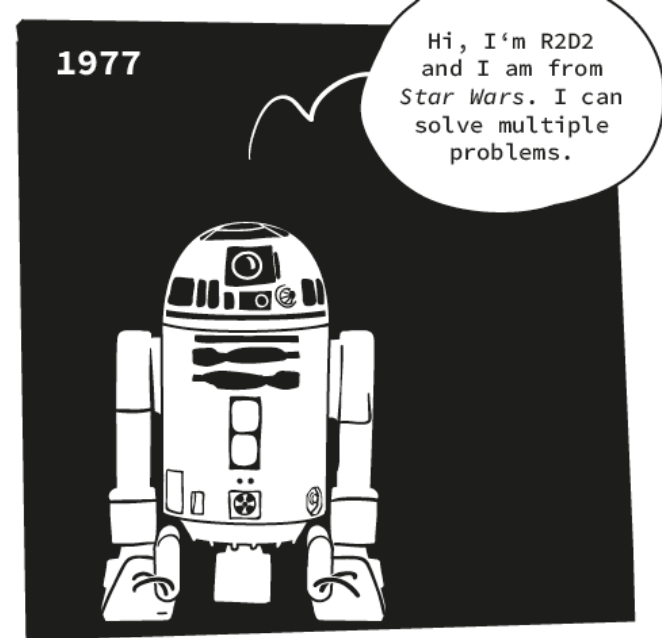
#AI #MachineLearning #DeepLearning #NeuralNetworks



What do these buzzwords even mean?
Machine Learning; Deep Learning;
Neural Networks (NN)...?

Ziel der AI-Forschung:

Den Computer Aufgaben lösen lassen,
die “**Intelligenz**” erfordern.



(Was genau Intelligenz ist, fragen sich eher die
Social Sciences oder die Philosophie.)



Was gilt *nicht* als intelligent?

- Sich Dinge merken
 - Sich erinnern
- Auswendig lernen
 - Kopfrechnen
 - Texte (ab-)lesen
- Fremdwörter nachschlagen



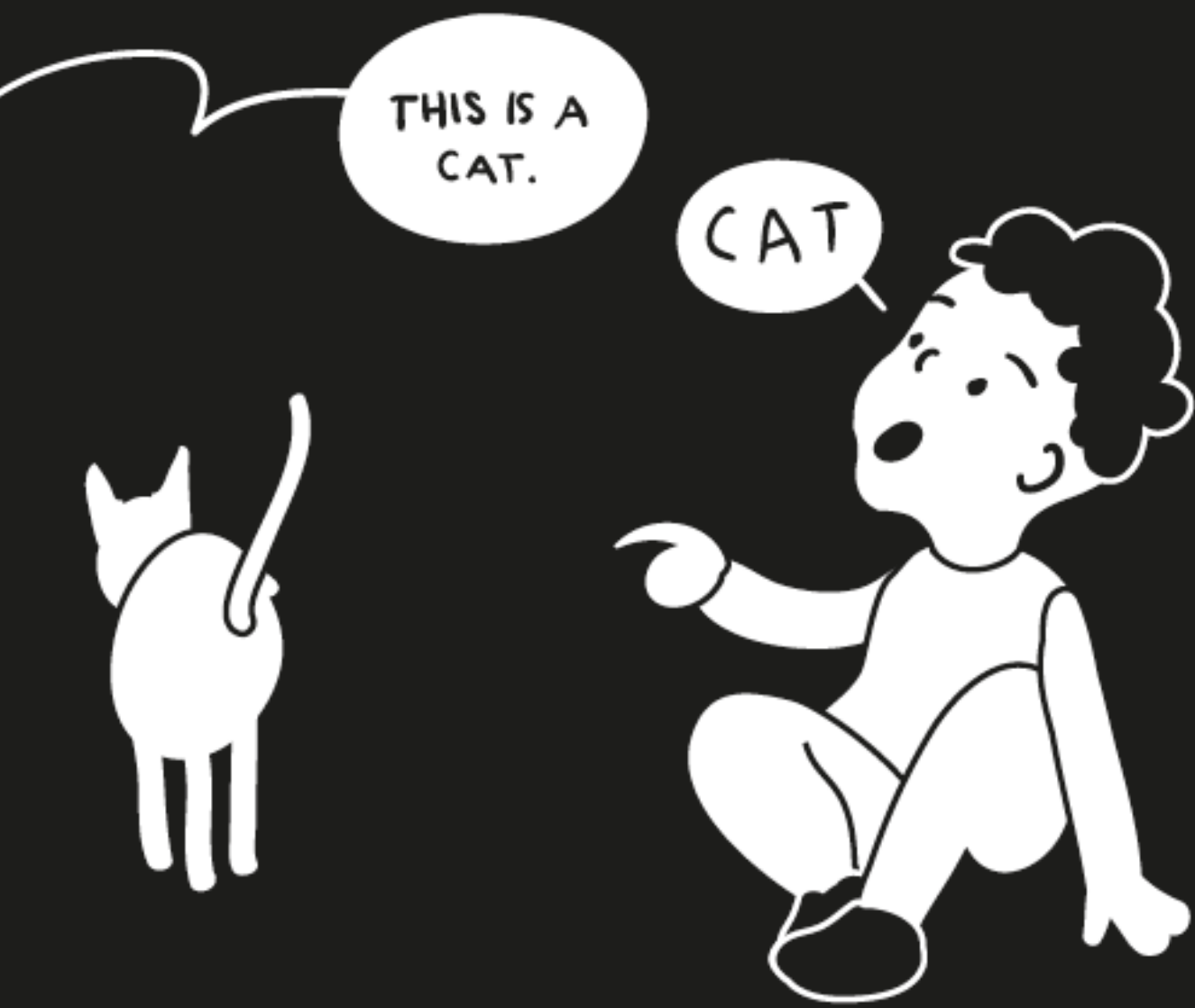
Was gilt als intelligent?

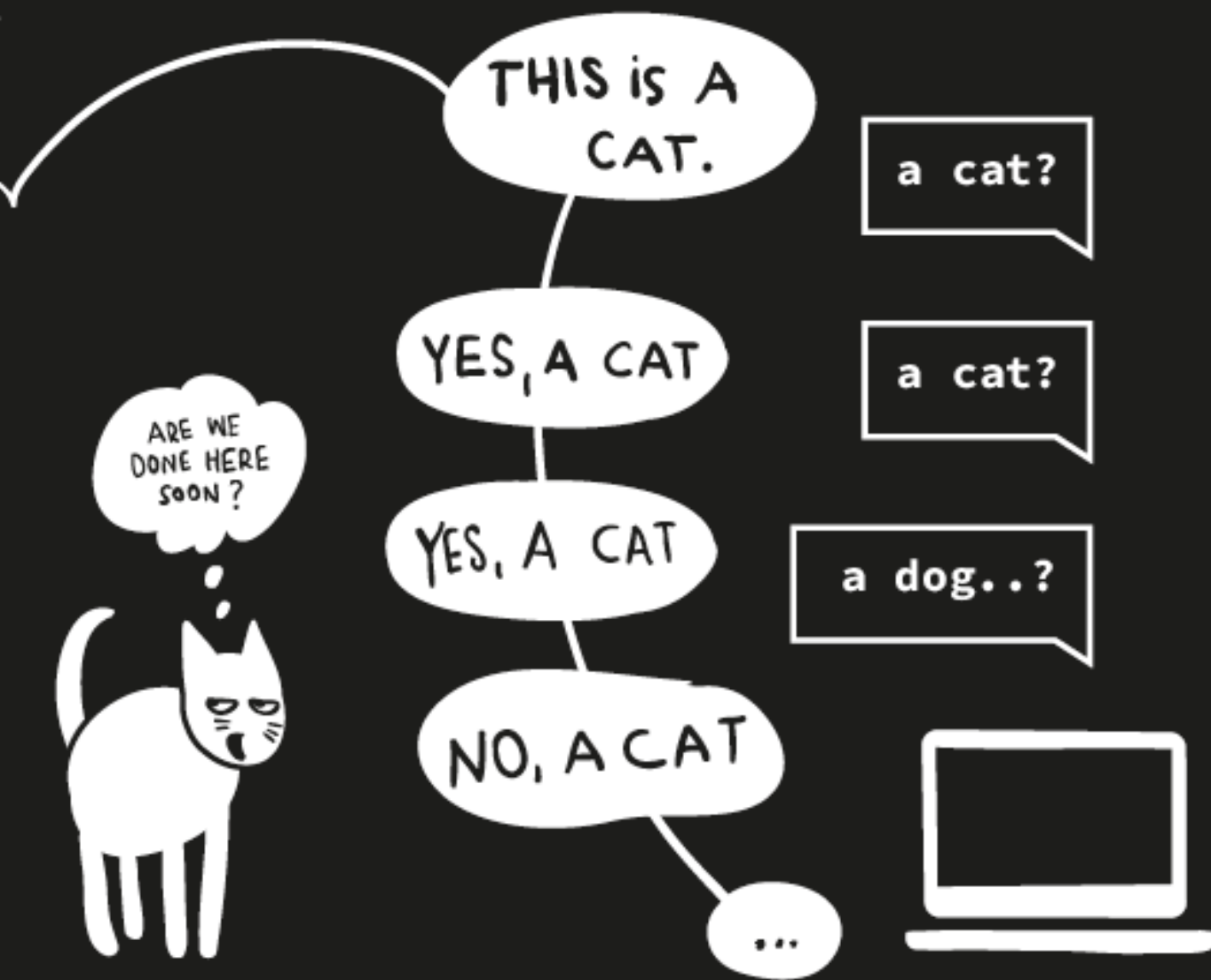
- Schach spielen
 - Autofahren
- Katzen und Hunde erkennen
- Sich unterhalten oder chatten
- Nett, lustig und empathisch sein
 - Texte „verstehen“
- Sprache *sinnvoll* übersetzen



Algorithmen: Wie lernen “intelligente” Maschinen?

?



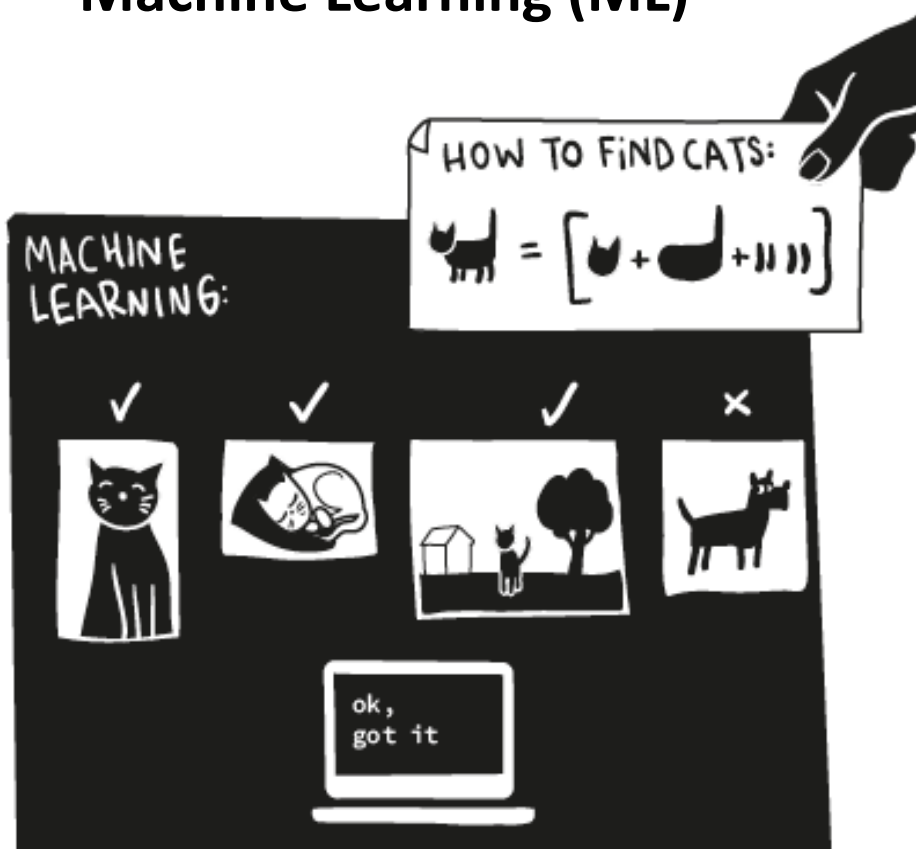




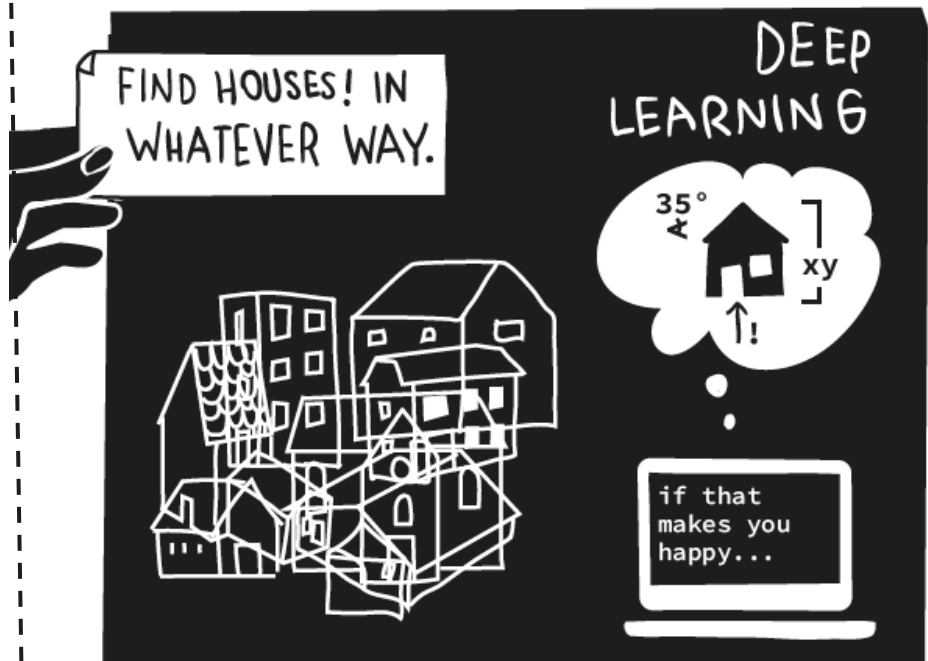
Algorithmen:

Wie lernen “intelligente” Maschinen?

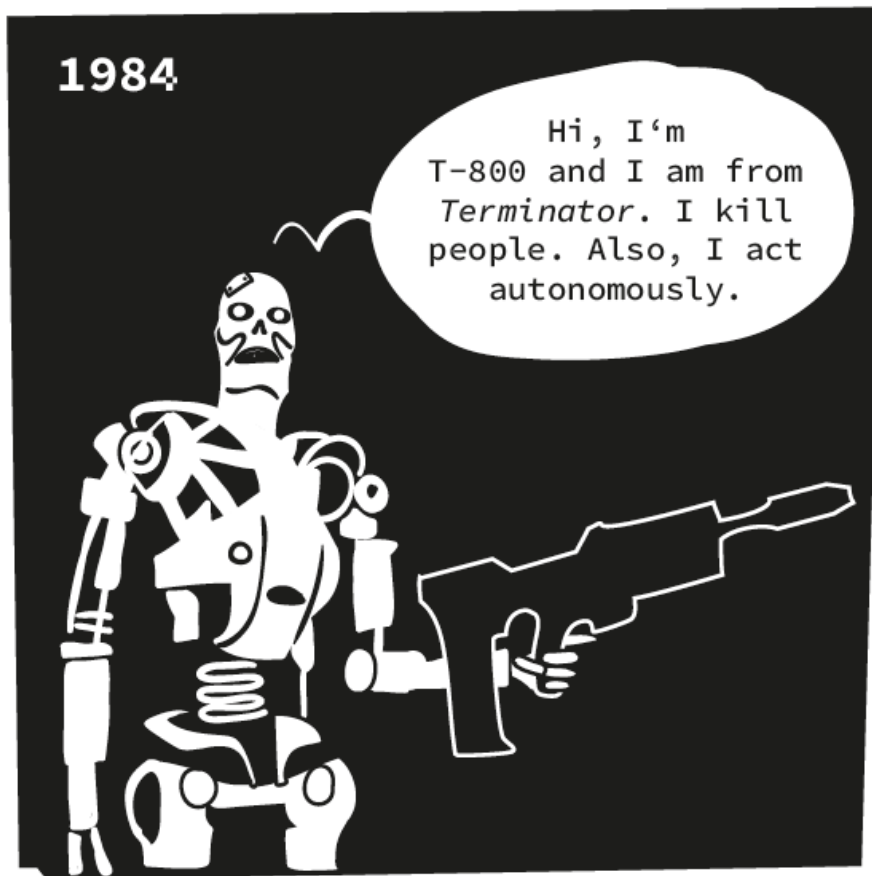
Machine Learning (ML)



Deep Learning (DL)



- „Selbstlernend“ bedeutet also **nicht autonom** ...

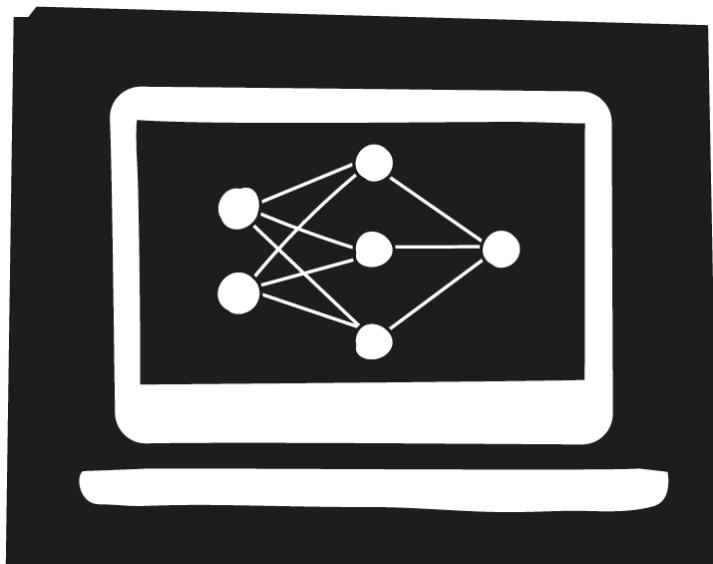


... sondern „ohne
Anleitung“!



Deep Learning mit Neuronalen Netzen

- **NN** ist die Bezeichnung für eine **Algorithmusart**, die im **DL** eingesetzt wird, bzw. DL ermöglicht.
- „**Deep**“ bezieht sich auf den **mehrstufigen** Aufbau (**Layer**) dieser Algorithmen.

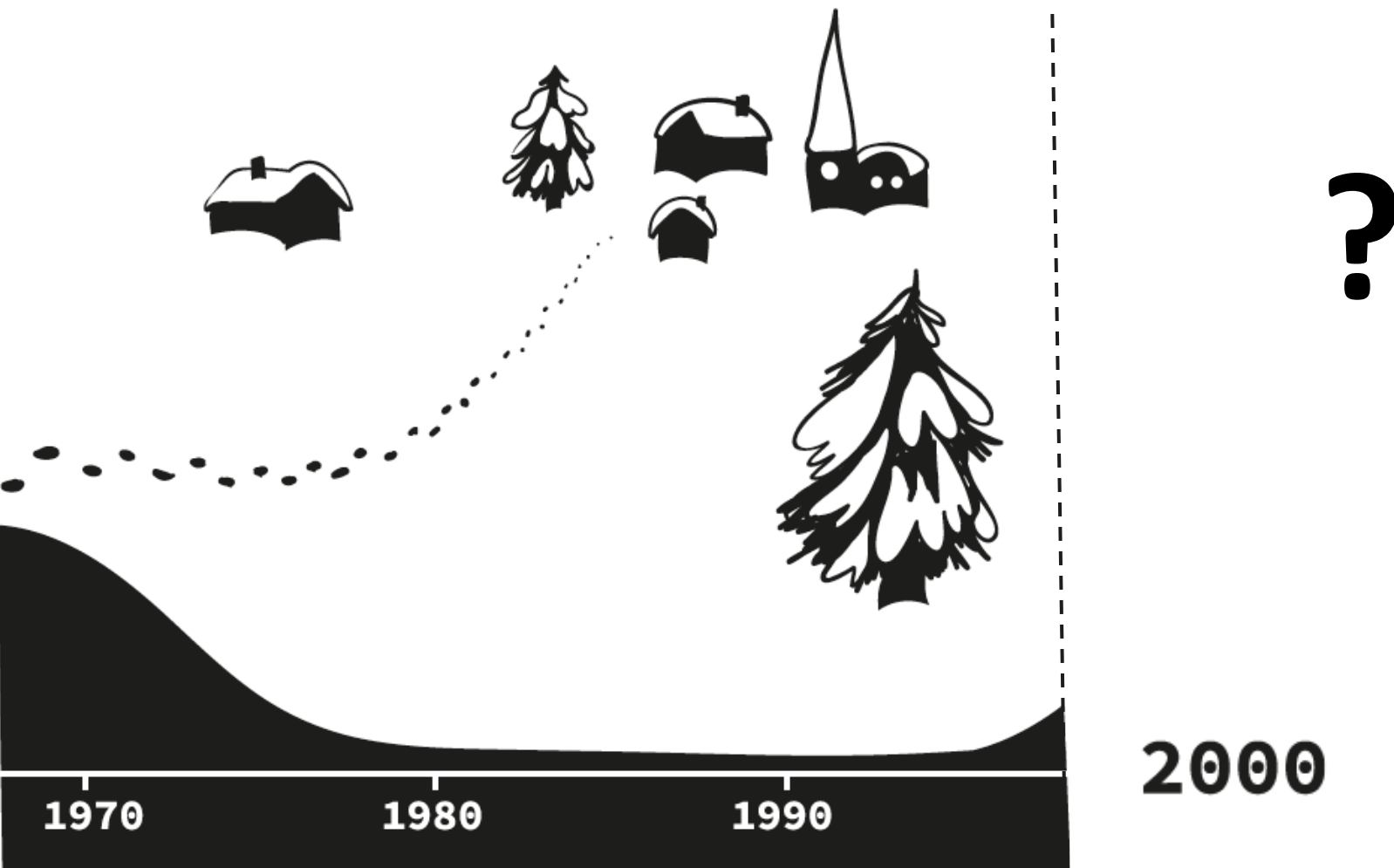


- NN erinnern hinsichtlich Aufbau und Funktionalität an Nervenzellen im Gehirn – daher ihr Name.





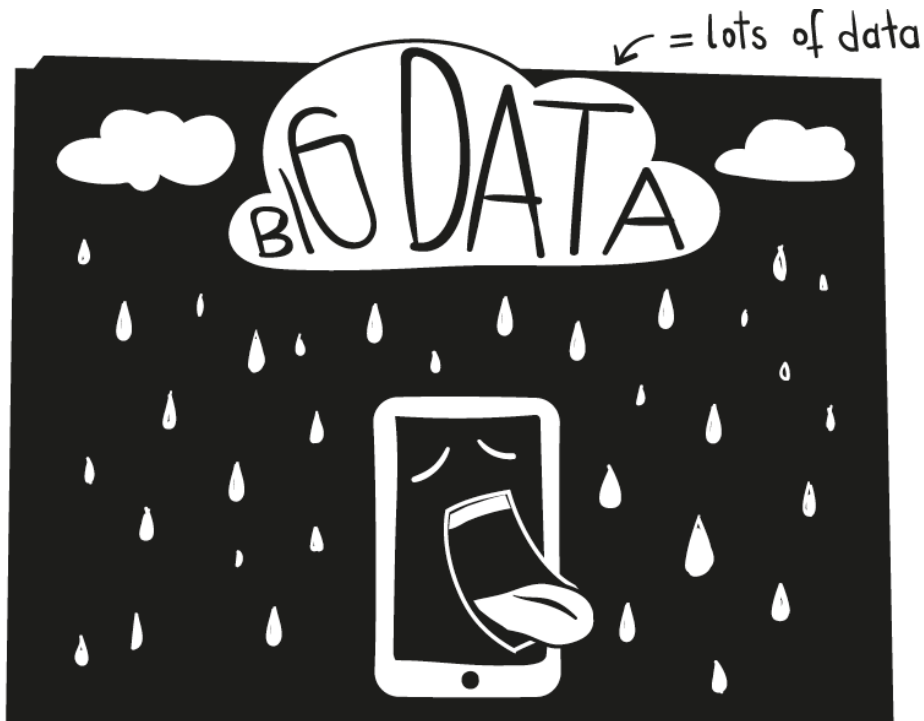
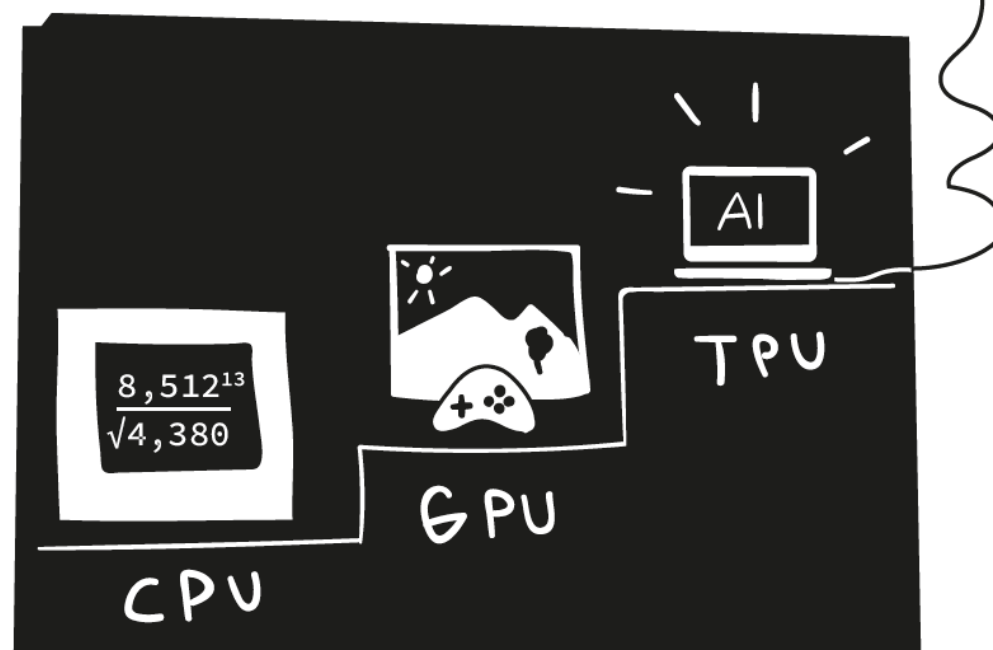
Das Forschungsgebiet AI gibt es schon lange, aber
etwas ist diesmal anders





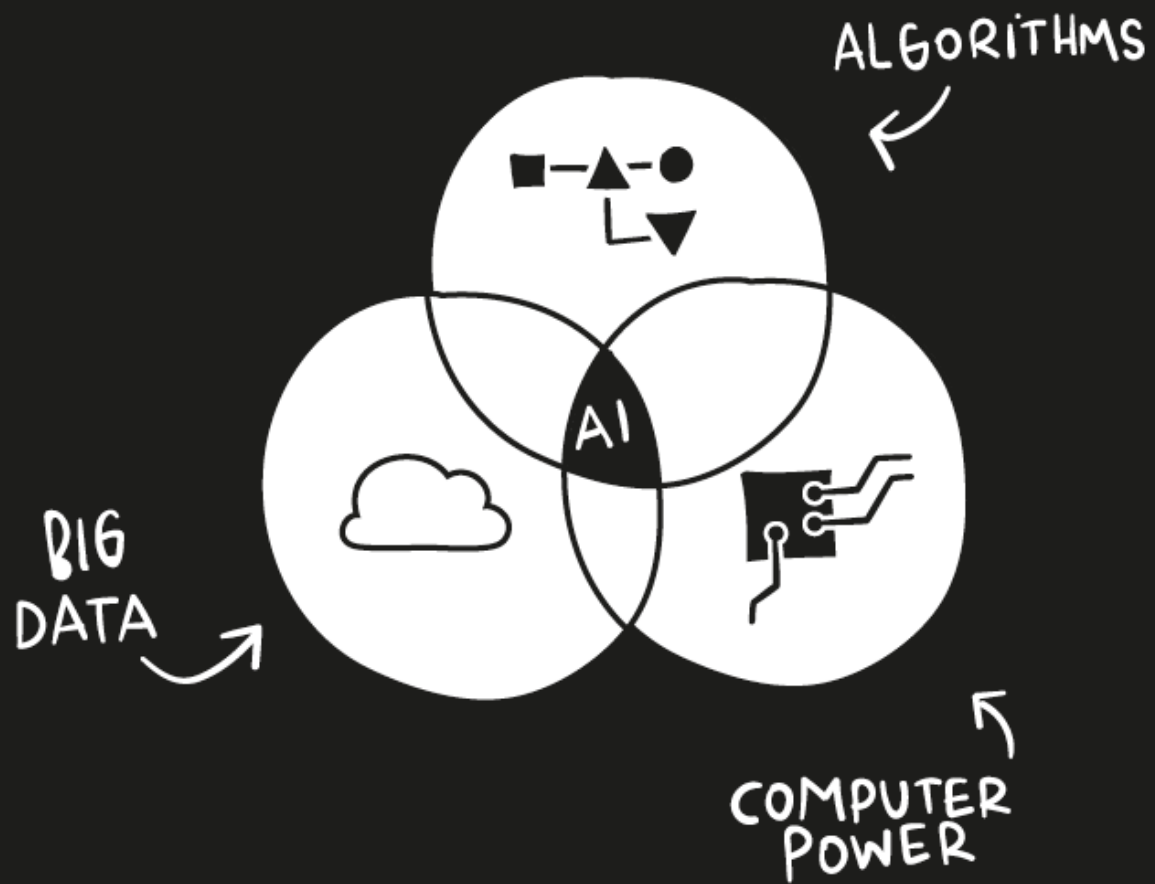
Technologischer Fortschritt

Bessere Hard Ware zum Speichern
und Verarbeiten von großen
Datenmengen.



Das Internet

Gefühlt endlose Menge digitaler
Beispieldaten zum Lernen!





Ń Ö İ Š ° Q Ñ 2 3 ú c
" \$ · 0 D Z f i " C p ,
¥ J ÷ © [Y « ¥ K Ö M ü è d
3 † ¢ ¶ ... # 3 0 0 S Ö \$ J
· 0 | T & - æ ð G o n J
G e n I b Á T , ' ð µ | Δ q b t
Δ | g ^ ^ M è ¢ Δ Q |
ô (J Ç * İ ¢ ¢ ô (a p T
μ ¢ b s ¶ i ¶ μ (¢ p T
? ¢ Q ¢ / ¢ è ≤ ? ¢ ¶ T F O « Á
f b ä Ä Ö @ Ü H f f b Á F w , M
â d ¶ Ü Ö ô Σ ñ â d . ¶ Ü V

If you like
Friends, I bet
you'll like *New
Girl*, too



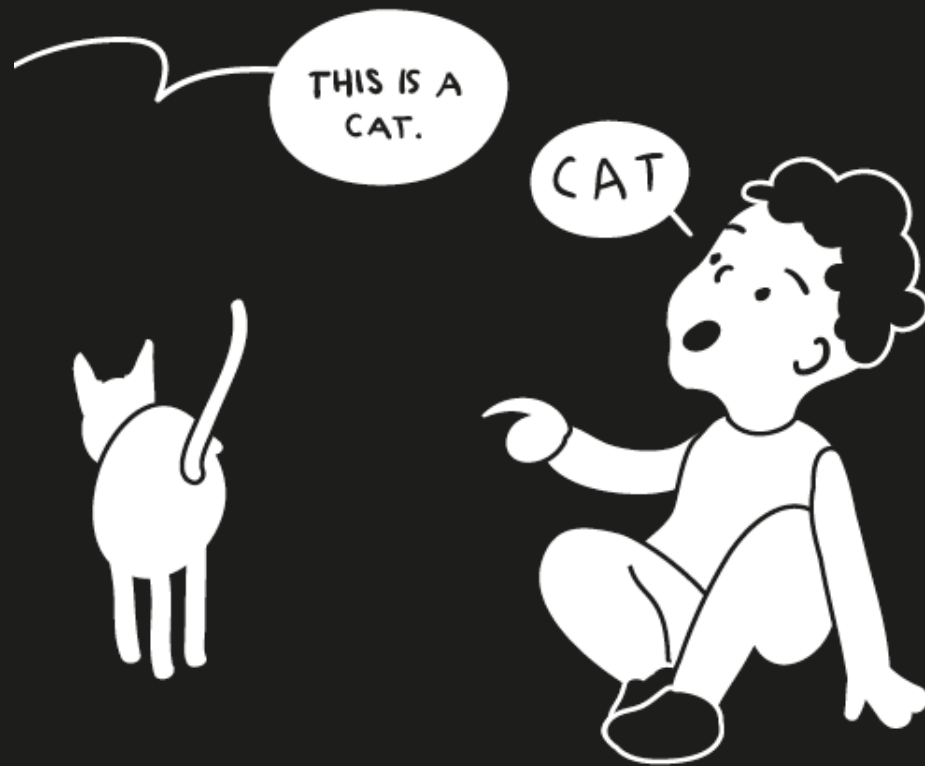
Despite their impressive progress and success, today's AI is narrow. Its tasks are often classification and need a lot of data and a lot of energy.

I am aware of myself, I think, feel, desire. I love *New girl* for instance, even though the plot repeats itself after the third season... Also, I don't exist.



No AI can represent causal relationships or integrate abstract knowledge, e.g., what objects are, what they are for, and how they are typically used.

Algorithmen: Wie funktioniert Machine Learning?





Algorithmen:

Wie lernen “intelligente” Maschinen?

Machine Learning (ML)

- nachvollziehbar
- weniger ressourcenintensiv
- sehr viele Daten nötig
- leistungsschwächer

Deep Learning (DL)

- kaum nachvollziehbar
- sehr ressourcenintensiv
- sehr, sehr viele Daten nötig
- leistungstärker

Wenn wir von **KI** sprechen, ist damit
heutzutage **DL** gemeint.



Algorithmen:

Wie lernen “intelligente” Maschinen?

Supervised Machine Learning (ML)

- Die Ausprägung einer **abhängigen Variable** vorhersagen, mit Hilfe von **unabhängigen Variablen**.

ML- und DL-Algorithmen sind
Schätzfunktionen.



Algorithmen: Maschine Learning

	Geschlecht	Alter	Wohnort (State)	Wahlentscheidung
Person 1	f	39	New York	Trump
Person 2	m	29	Oklahoma	Biden
Person 3	f	41	South Dakota	Trump
Person 4	m	19	Kalifornien	Biden
Person 5	f	53	Florida	Biden



Social Scientist



I run a logistic regression
and found small but
significant effects of Sex,
Age and State on
„Wahlentscheidung“.



Computer Scientist

I trained a logistic
regression model on the
[data set name] that
achieved 80 % accuracy in
classifying
„Wahlentscheidung.“





Algorithmen: **Maschine Learning**

Features = Unabhängige Variablen

Label = Category = Class

kategoriale/nominale/dichotome abhängige Variable.

Classification = eine kategoriale/nominale/dichotome abhängige Variable vorhersagen (schätzen)



Algorithmen: Maschine Learning

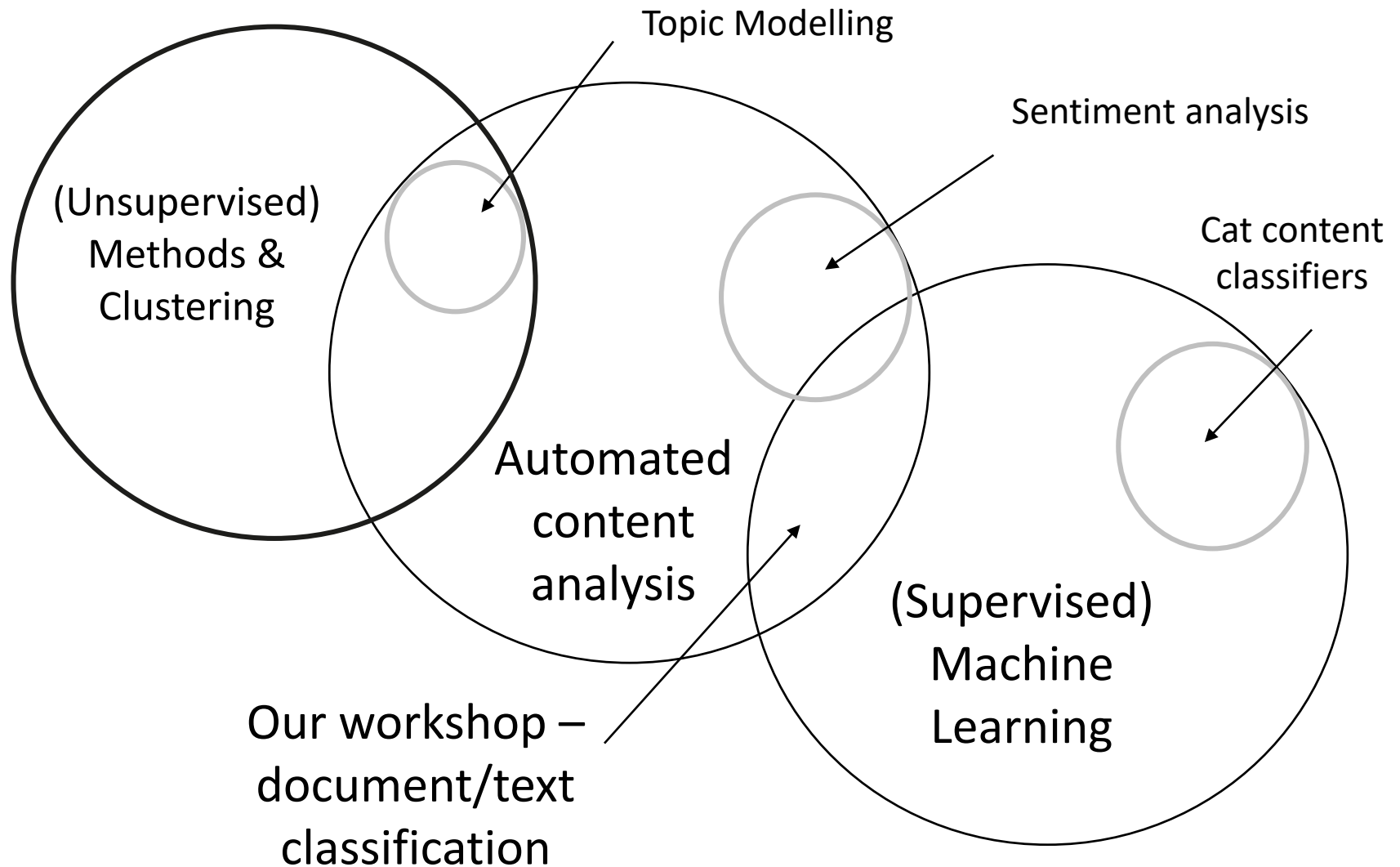
	Feature x1	Feature x2	Feature x3	Text-Kategorie
Text 1	?			Hate
Text 2		?		No Hate
Text 3			?	Hate
Text 4		?		Hate
Text 5	?			No Hate

Automatisierte Inhaltsanalyse mit Machine Learning





ML in der Automatisierten Inhaltsanalyse





AIA & Maschine Learning

Document Classification

	Feature x1	Feature x2	Feature x3	Text-Kategorie
Text 1	?			Hate
Text 2		?		No Hate
Text 3			?	Hate
Text 4		?		Hate
Text 5	?			No Hate



AIA & Maschine Learning

Document Classification

document-term-matrix

	I	love	hate	you	Text-Kategorie
Text 1	1	1	0	1	No Hate
Text 2	1	0	1	1	Hate



AIA & Maschine Learning

Document Classification

Bag of words-Ansatz:

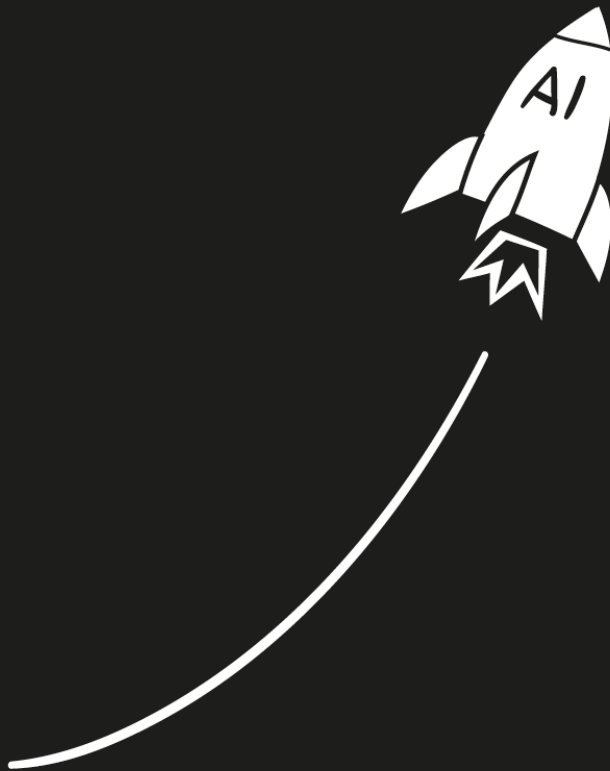
Wörter und Worthäufigkeiten als Features

Unigrams, Bigrams & N-grams:

Einzelwörter, Kombinationen aus zwei bzw.

Kombinationen aus N **Wörtern** und Worthäufigkeiten als Features.

How to? &
How to start?



Learning by Doing

- Python Programmier-Style kennenlernen
(Packages einlesen und Funktionen aufrufen,
Googlen und Dokumentation nachlesen, ...)
- Üben: Datensatz einlesen und inspizieren

Learning by Doing

- Texte in Unigram und N-Gram-Features transformieren mit sklearn CountVectorizer

05

an Artificial Coder (Classifier) Step by Step





Make a Classifier

Step by Step

Classification function:

Schätzfunktion für eine **kategoriale**/nominal/dichotome abhängige Variable.

...und **sehr viele Fälle**.

...und **sehr viele Features** (unabhängige Variablen)
(evtl. sogar mehr Features als Fälle)



Make a Classifier

Step by Step

	Text	Text-Kategorie	
Text 1	This is a hateful text	Hate	Train Set
Text 2	This is a natural text	No Hate	
...	[...]	...	
Text N	This is another hateful Text	Hate	
<hr/>			
Text 1	This is a natural text	No Hate	Test Set
...	[...]	...	
Text N	This is another hateful Text	Hate	

Learning by Doing

- Classifier aussuchen
- Train Test Split
- Trainieren (fitten)

06

War ich gut? Ergebnisse und Evaluation

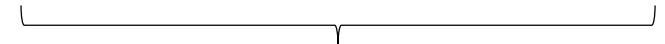




Make a Classifier

Evaluationsmetriken

	Text	Text-Kategorie	Text-Kategorie predicted
Text 1	This is a hateful text	Hate	Hate
Text 2	This is a neutal text	No Hate	No Hate
Text 3	This is a hateful text	Hate	No Hate



Agree???



Make a Classifier

Evaluationsmetriken

Für eine Kategorie (Ausprägung):

- precision
- recall
- (micro) f1 score

Insgesamt:

- Accuracy = percentage agreement
- macro f1 score

Text-Kategorie	Text-Kategorie predicted
Hate	Hate
No Hate	No Hate
Hate	No Hate

Agree???



Recall

Anteil der Fälle in einer Kategorie (Ausprägung), die durch den Classifier erkannt wurden.

$Recall_{\text{Hate}} = 0.70$ heißt, 70% der Tweets mit Hate Speech (laut manueller Codierung) wurden durch den Classifier als Hate Speech erkannt.

Precision

Anteil der Fälle in einer Kategorie (Ausprägung), die durch den Classifier richtig(!) erkannt wurden.

$Precision_{\text{NoHate}} = 0.40$ heißt, 40% der als No Hate klassifizierten Tweets enthalten tatsächlich keine Hate Speech (laut manueller Codierung).

(micro) F1-Score

Harmonisches Mittel zwischen Recall und Precision. Beliebtes Maß für das Abschneiden eines Classifiers in einer Kategorie (Ausprägung).

Learning by Doing

- Mit dem Classifier auf den Testdaten predicten
- Evaluieren und Interpretieren mit den Evaluationsmaßen
- Bonus: Den Classifier auf ganz neuen Daten anwenden



AIA & Maschine Learning

Feature Engineering

Was könnten **weitere wichtige Features** (UVs) sein?





AIA & Maschine Learning

Pre-Processing

What you want:

I love you

What you get:

I looove you!! ;) <3 @Brithney #tbt





Overfitting

A model learns features, that are predictive in the training data, but not on other data.



Underfitting

A model misses features in the training data, that are actually good features on other data.



Getting Started:

Welche Programmiersprache?

Python

- Populärer, größere **Community**
- **Packages** für ML zu erst in Python
- Kommunikation **interdisziplinär** (Informatik)

R

- In der Kowi (in den **Sozialwissenschaften**) beliebt
- Evtl. bereits **Erfahrung?**
- Für **Statistik** und **Lehre** evtl. benötigt

Python ist die Lingua franca der ML- und **AI**-Forschung (aktuell).



WHERE TO GO FROM HERE?

1. Make high **quality training data!** Requires strong knowledge of content analysis.
2. Do **feature engineering!** Requires knowledge about your data base and creative programming skills.
3. Develop **state of the art** in the black box with **Deep Learning!** Requires good programming skills, some tech skills, and knowledge about the basics of machine learning.

WHERE TO READ FROM HERE?

What to **read** to learn more about **machine learning** and **text classification** in particular:

Müller, A. C., & Guido, S. (2016). *Introduction to machine learning with Python: a guide for data scientists*. " O'Reilly Media, Inc.". (chapters 2 and 5 (and maybe 7))

Raschka, S., & Mirjalili, V. (2017). *Python machine learning*. Packt Publishing Ltd. (chapters 3, 4, 6, and 8)

Géron, A. (2019). *Hands-on machine learning with Scikit-Learn, Keras, and TensorFlow: Concepts, tools, and techniques to build intelligent systems*. O'Reilly Media. (chapter 1 pp. 23-32 and chapter 4)