

Machine Learning mit Python

Schulung der GfU

EXXETA
CONSULTING AND TECHNOLOGIES

Philipp Bongartz

10.-12. Mai 2021





Dr. Philipp Bongartz

Consultant

Dr. Philipp Bongartz programmiert seit 2013 unterschiedlichste Anwendungen in Python mit einem Fokus auf performante Algorithmen und Machine Learning. Er sammelte im Rahmen seiner Promotion Erfahrung im Deep Learning, im klassischen Algorithm Design und in der statistischen Analyse großer Datenmengen. Bei EXXETA ist er neben den Themen Python Programmierung und Machine Learning auch im Data Warehousing aktiv.

Auszug relevante Projekterfahrung

Lead Programmer, Applikation zum Auslesen von Personalausweisen, Produktentwicklung

- Entwicklung der Python Codebase für eine ML-Pipeline zum Auslesen von Personalausweisen mit Komponenten aus Bilderkennung und -bearbeitung, OCR, Optimierung, Dynamic Programming, Clustering, sowie Deep Learning.

Lead Programmer, Applikation zum automatischen Matchen von Consultants und Stellenausschreibungen

- Betreuung und Weiterentwicklung der Python Codebase einer serverbasierten KI-Anwendung für die automatische Auswertung von Stellenausschreibungen

DWH Entwickler, Betreuung und Weiterentwicklung einer Softwarelösung für die Datenversorgung von über 90 Applikationen aus verschiedenen SAP Systemen

- Anforderungsaufnahme und Analyse sowie Anpassen von ETL-Strecken, Datenbanken und Reports
- Sicherstellung von Datenkonsistenz & Performanceoptimierung
- Automatisierte ETL-Streckenimplementierung durch Scripting
- Onboarding Schulungen & Enabling für Projektübergabe zu Kunden
- Technologien: MS SQL Server, SSIS, Theobald XTract IS (für SAP BW, -CRM, -ECC, -SCM), SAP SolMan

Data Scientist, Predictive Maintenance auf Basis von Sensordaten, Schienen-Logistik Firma

- Entwicklung eines Self-Supervised Deep Learning Models zur Klassifizierung von Zeitreihen
- Entwicklung einer regelbasierten Zeitreihenklassifikation für Batteriespannungsverläufe im Rahmen von Monitoring und Predictive Maintenance

Biographie

- EXXETA, Mannheim
- PhD Bioinformatik, HITS, Heidelberg
- Mathematik-Diplom, Universität Bonn

Beratungskompetenz

- Softwareentwicklung Python, C
- Data Science & Deep Learning: Keras, Tensorflow, SciPy, Numpy, Matplotlib, Jupyter Notebook, Pandas, ScikitLearn
- DWH & ETL: T-SQL, SQL-Server, SSIS

Sprachen

- Deutsch
- Englisch



AGENDA

• Tag 1

- Machine Learning Einführung
- Panda Basics
- Data Preprocessing
- Lineare Regression
- Korrelation
- Decision Trees
- Clustering

• Tag 2

- Deep Learning Einführung
- Numpy Basics
- Gradient Descent & Generalisierung
- Convolutional Neural Networks
- Natural Language Processing

• Tag 3

- Visualisierung
- ML Ops



Inhalt

Daten ziehen und aufbereiten

- Dateien im Verzeichnis ansteuern
- Standardverfahren zum Lesen/Schreiben von Text- und CSV-Dateien
- SQL-Datenbanken ansteuern
- Arbeiten mit Datenmatrizen: Einführung in *Numpy* & *Pandas*
- Daten inspirieren und beschreiben
- Fehlende Werte behandeln

Machine Learning Grundlagen

- Grundlagen, Varianten und Techniken des Machine Learnings.
- Arbeiten mit der Machine Learning Bibliothek scikit-learn.
- Einfache Zusammenhänge zwischen stetigen Variablen modellieren: Lineare Regression
- Klassifizieren mit Logistic Regression, Softmax und Support Vector Machine.
- Modellen evaluieren: Accuracy, Precision, Recall & Confusion matrix
- Unterteilen der Daten in Trainings- und Testdaten

Feature-Extraction

- Kategoriale Daten vorbereiten: *One-Hot Codierung*
- Standardisierung von Daten
- Dimensionsreduktion mit PCA (Principle Component Analysis)
- Aufbereitung von Textdaten: Tokenizer und Bag-of-words.

Machine Learning Workflow

- Daten aufarbeiten und zusammenführen (DataMapper)
- Pipelines einrichten
- Speichern und laden trainierter Klassifizierer

Beschreibung

Das Seminar gibt eine Einführung in Konzepte und in die Praxis des Machine Learnings. Es zeigt Einsatzoptionen, erklärt statistische Grundlagen und schult anhand praktischer Beispiele die Anwendung der gängigsten Lernalgorithmen aus der *Scikit-Learn* Bibliothek. Neben der Klassifizierung und Prognostizierung von Daten, liegt der Schwerpunkt auf der Vorverarbeitung der Daten und der Extraktion relevanter Variablen für den Anlernprozess.

Schulungsziel

Sie lernen die Grundlagen und die praktische Anwendung der gängigsten Lernalgorithmen aus der *Scikit-Learn* Bibliothek. Sie bekommen ein Gespür für den Wert ihrer Daten und dafür, wie Sie diese Daten in Ihrem Geschäftsfeld einsetzen können, um Automatisierungsprozesse voranzutreiben. Sie verstehen die verschiedenen Lernalgorithmen in Theorie und Praxis (lineare und logistische Regression, Support Vector Machine, Decision Tree, Naive Bayes) und üben die Anwendung anhand praktischer Beispiele ein. Die Konzepte werden anhand von Folien anschaulich erklärt, an Beispielen verdeutlicht und gemeinsam in *Python* umgesetzt. Passgenaue Aufgabenstellungen ergänzen den Lernprozess und ermöglichen es Ihnen, die verschiedenen Lernszenarien kennenzulernen und einzuüben. Am Ende des Seminars sind Sie in der Lage, Daten zielsicher zu extrahieren, Algorithmen anzulernen und zur Klassifizierung oder Prognose einzusetzen.

https://github.com/PhilippBongartz/GfU_ML_mit_Python



Vorstellungsrunde

- Wie heißt Ihr?
- Was habt Ihr für einen Hintergrund?
- Was führt Euch hierher?
- Warum interessiert Euch Machine Learning?

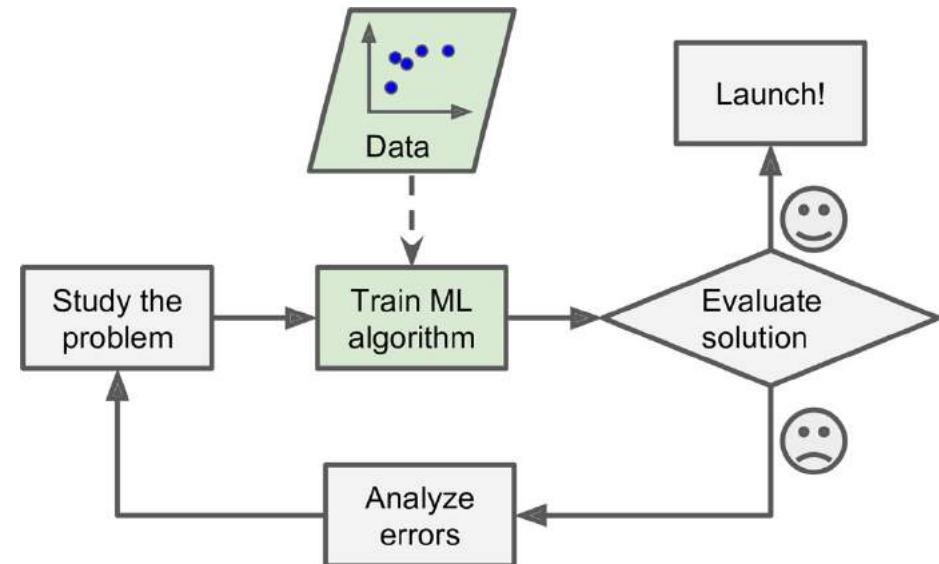
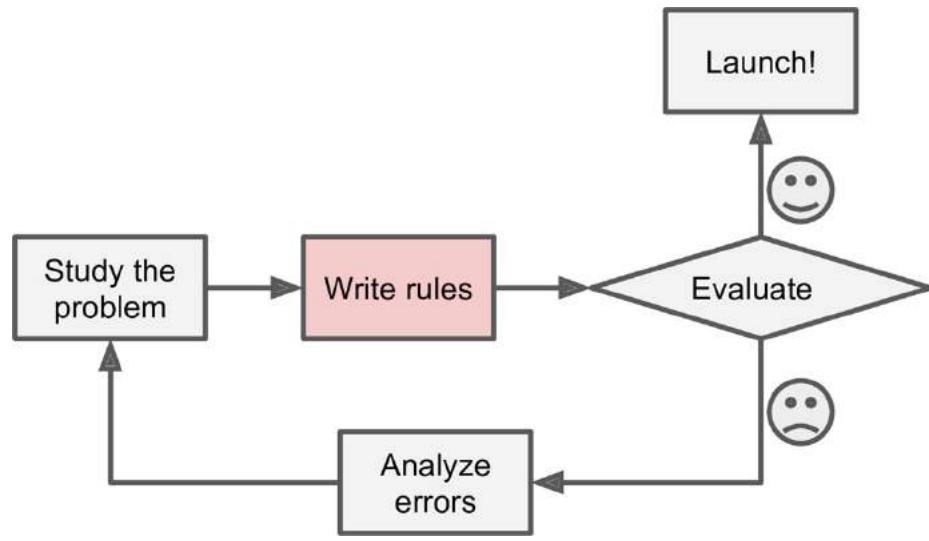


Machine Learning Überblick



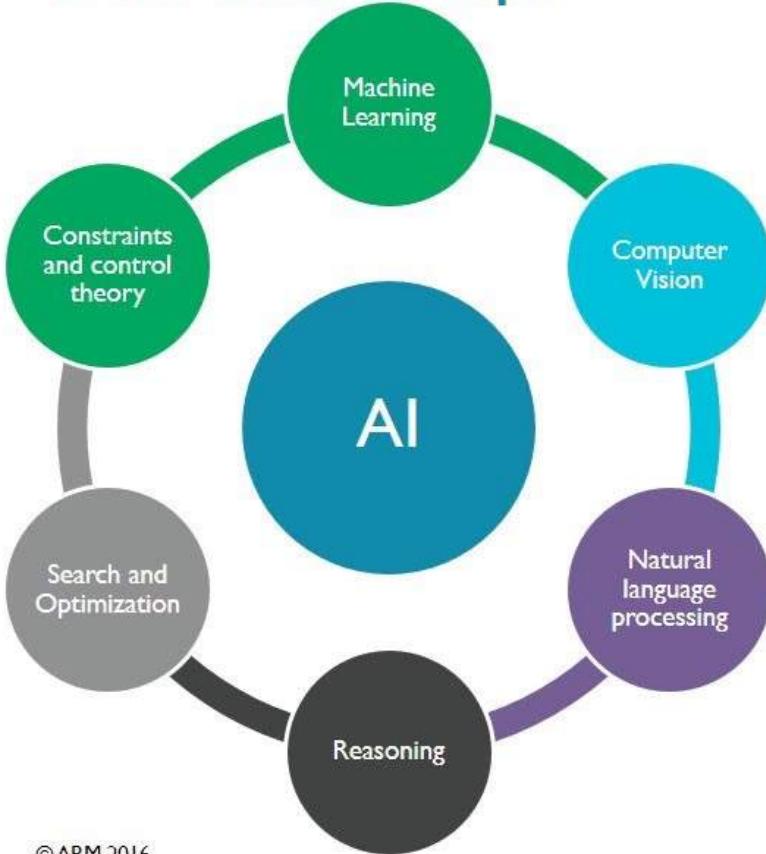
Was ist Machine Learning?

Algorithmen, die aus Daten lernen oder Struktur in Daten entdecken.



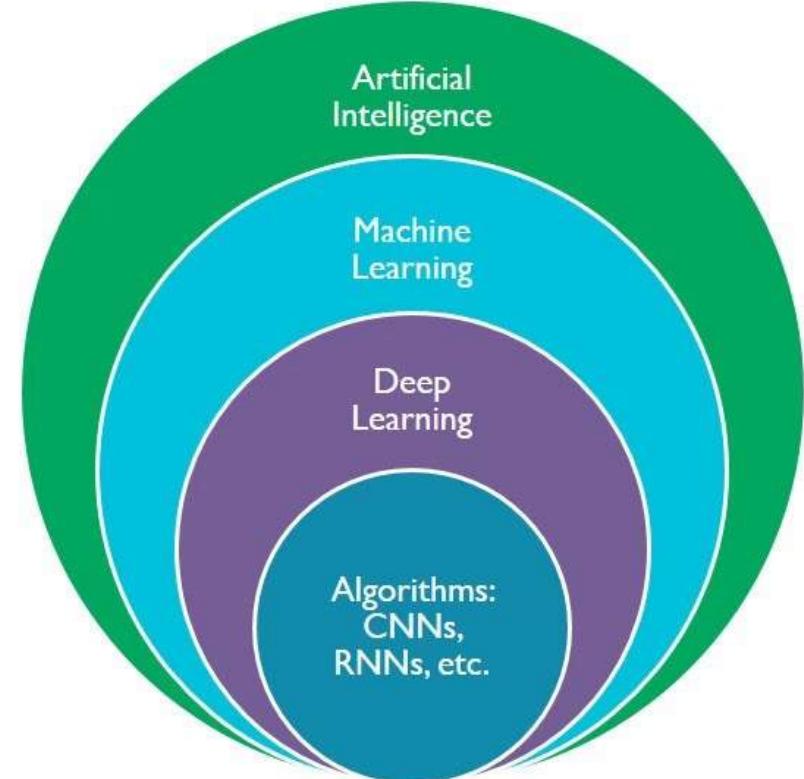
Artificial Intelligence and Machine Learning

The AI landscape



6

©ARM 2016



ARM



Startups and AI

The **AI 100** is CB Insights' annual ranking of the 100 most promising AI startups in the world.

2020

Healthcare



Finance & Insurance



Transportation



Construction



Retail & Warehousing



Govt. & City Planning



Legal



Mining



Food & Agriculture



CROSS-INDUSTRY TECH

AI Processors



AI Model Development



DevOps & Model Monitoring



NLP, NLG, & Computer Vision



Cybersecurity



BI & Ops Intel



Sales & CRM



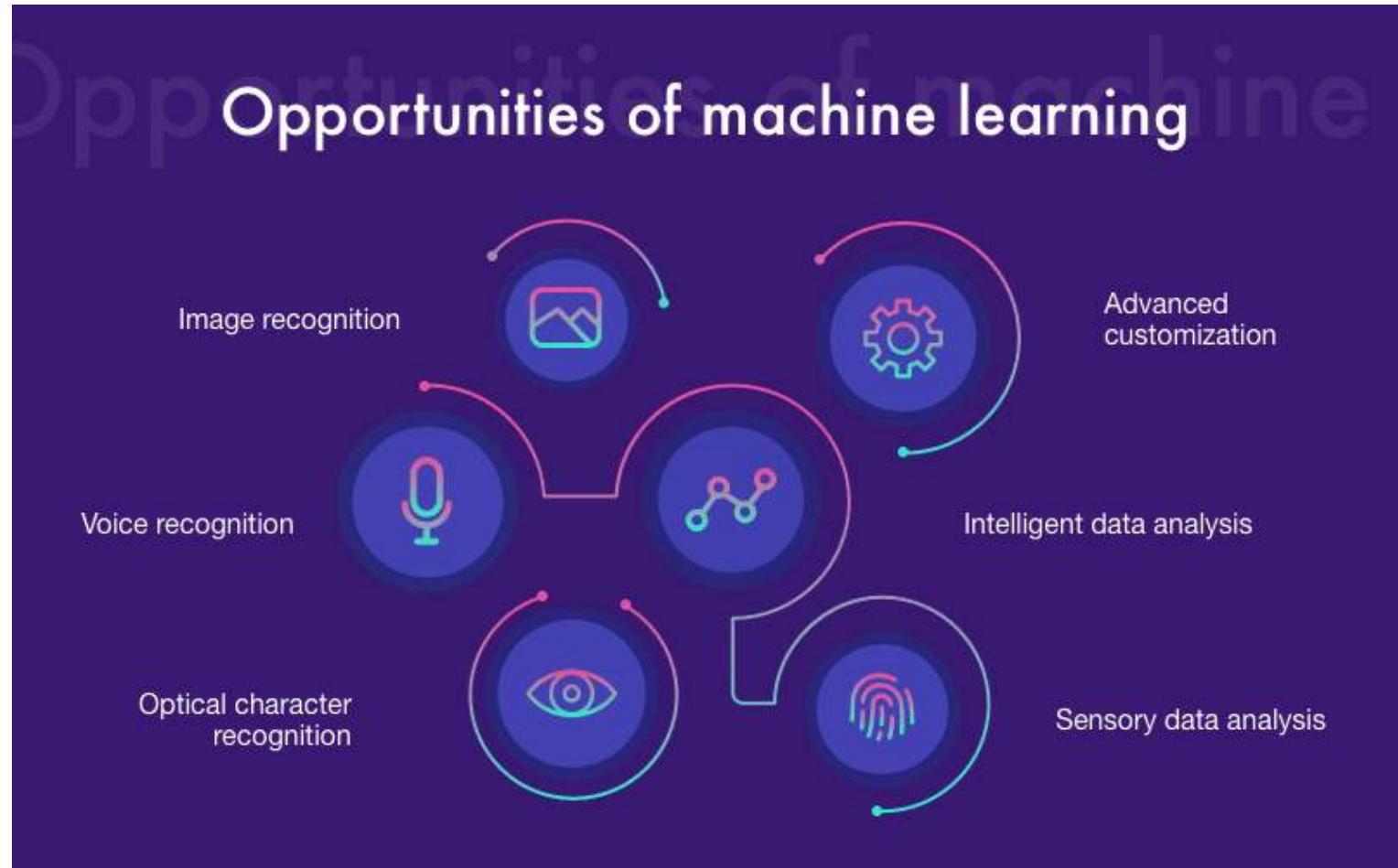
Other R&D



<https://www.cbinsights.com/research/artificial-intelligence-top-startups/>



Was kann man damit alles machen?



Arten des Machine Learning

- 1. Supervised Learning**
- 2. Unsupervised Learning**
- 3. Reinforcement Learning**

...



sind drei klar
unterschiedene
Hauptbereiche



Supervised Learning

Fragestellung Supervised Learning:

z.B.: Gibt es Kunden, welche eine hohe Wahrscheinlichkeit besitzen ihren Vertrag zu kündigen?

- Training erfolgt auf historischen Daten, welche alle gelabelt sind.
- Ziel ist es durch die vorliegenden Daten (Features), die Zielvariable so genau wie möglich vorhersagen zu können und ggf. eine Wahrscheinlichkeit der Klassenzugehörigkeit zu ermitteln.
- Vorhersage erfolgt mit Daten, bei denen das Target oder die Klasse fehlt.

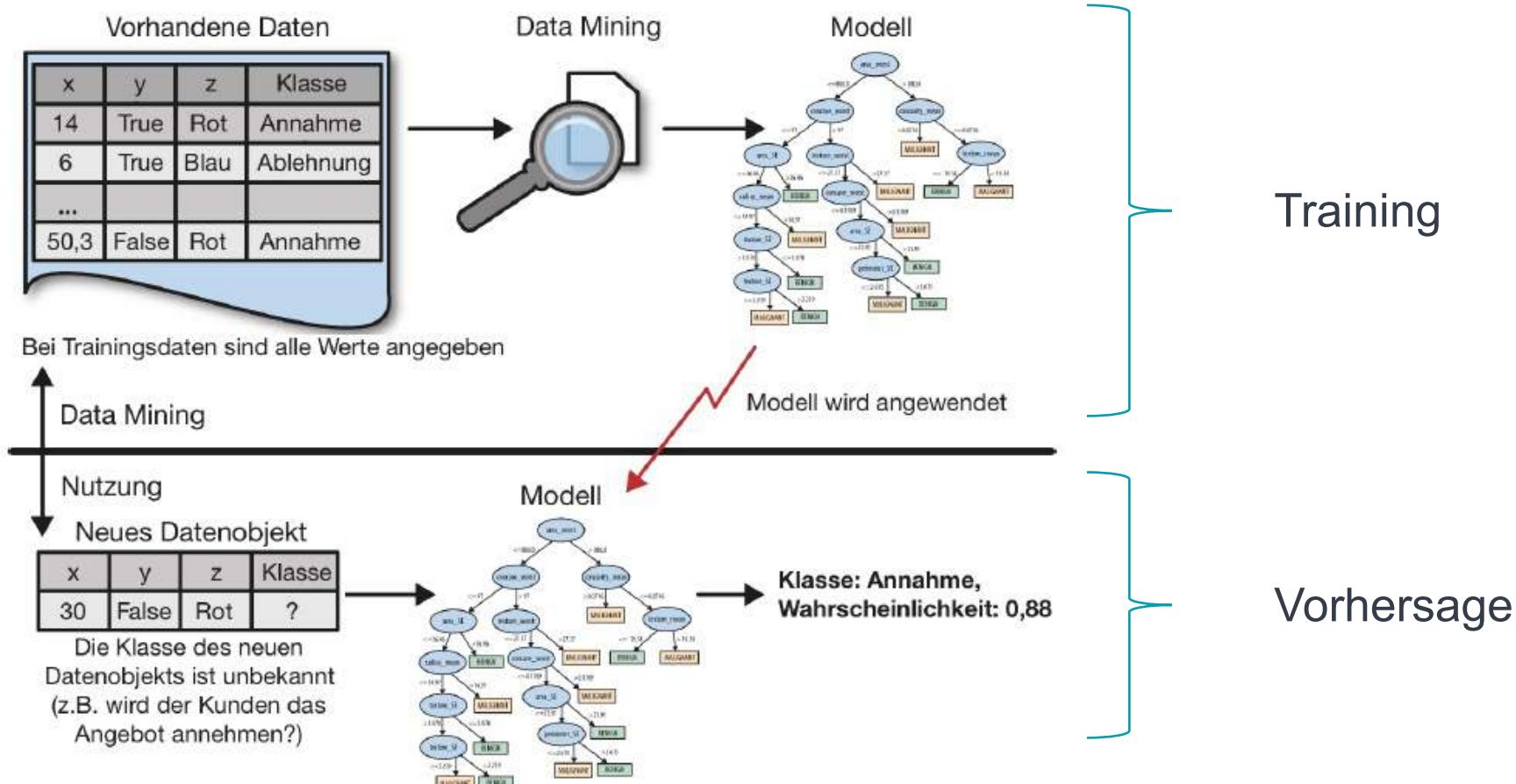


Live Demo

Teachable Machine



Supervised Learning



Supervised Learning – Schlüsselkonzept Train-, Test-, Validationsplit

Bevor man ein Machine Learning Modell trainiert, sollte man immer den Datensatz mindestens in Train- und Validationset unterteilen.

Performance auf dem Trainingsdatensatz garantiert nicht die Performance auf ungesehenen Daten!

Die Terminologie ist nicht immer sehr einheitlich. Manchmal spricht man auch von Training, Development und Testdatensätzen. Der dritte Datensatz ist dafür da, dass man nicht durch optimieren der Hyperparameter das Modell an den zweiten Datensatz überanpasst.

Der dritte Datensatz wird idealerweise nur ein **einziges** Mal ausgewertet!

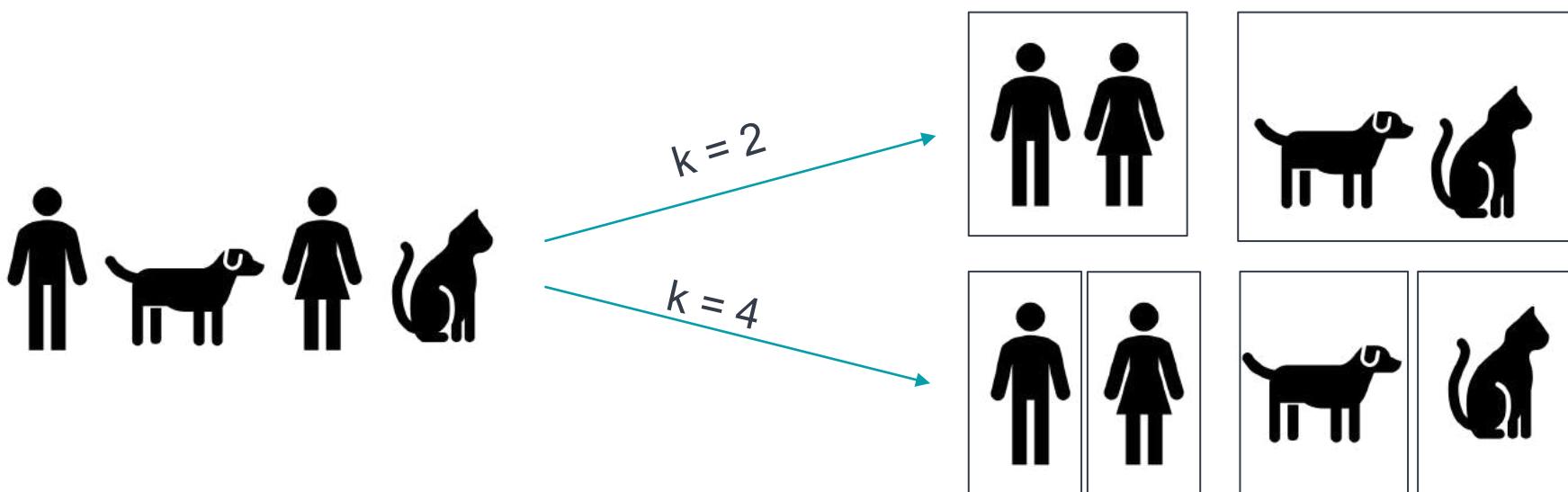


Unsupervised Learning

Fragestellung Unsupervised Learning:

z.B.: Können Accounts verschiedener Nutzer in bestimmte Gruppen aufgeteilt werden?

- Keine Angaben über ein bestimmtes Ziel in den Daten vorhanden
- Ziel wird allein anhand vorliegender Daten und deren Ähnlichkeit zueinander bestimmt



Unsupervised Learning

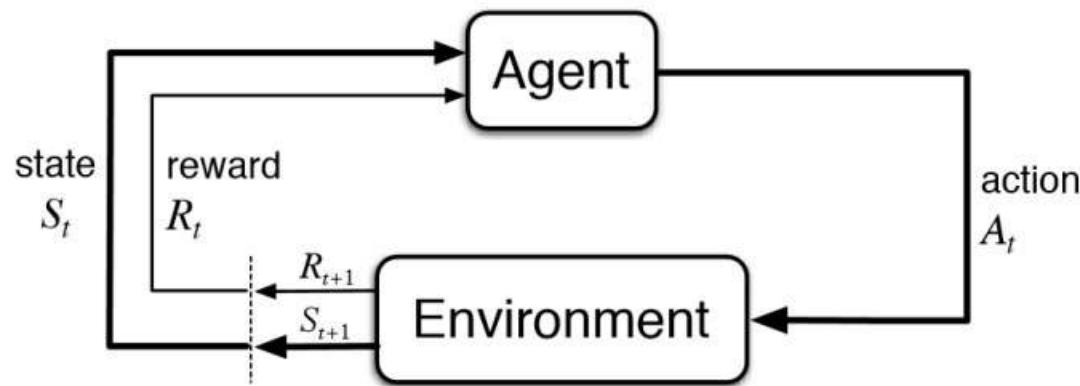
Arten des Unsupervised Learning:

- Clustering (z.B. kMeans oder DBSCAN)
- Dimensionsreduzierung (z.B. PCA)
- Visualisierung (z.b. t-SNE)
- Einbettungen oder Featurelearning (z.b. über Contrastive Loss)
- Lernen von Wahrscheinlichkeitsverteilungen über Daten (Self-Supervised)



Reinforcement Learning

- Sogenannter „Agent“ wird trainiert
- „Agent“ führt Aktionen aus und bekommt Rückmeldung von Umgebung → z.B.: Belohnung bei Bewegung in die richtige Richtung.
- Verwendung: Robotics, Gaming-Bots, AlphaGo,...



Source: https://www.google.com/search?q=reinforcement+learning&source=lnms&tbo=isch&sa=X&ved=2ahUKEwiwfG11IrwAhXUgf0HHdJxAwQQ_AUoAnoECAEQBA&biw=1920&bih=937#imgrc=ivmsD3gf5LO7WM

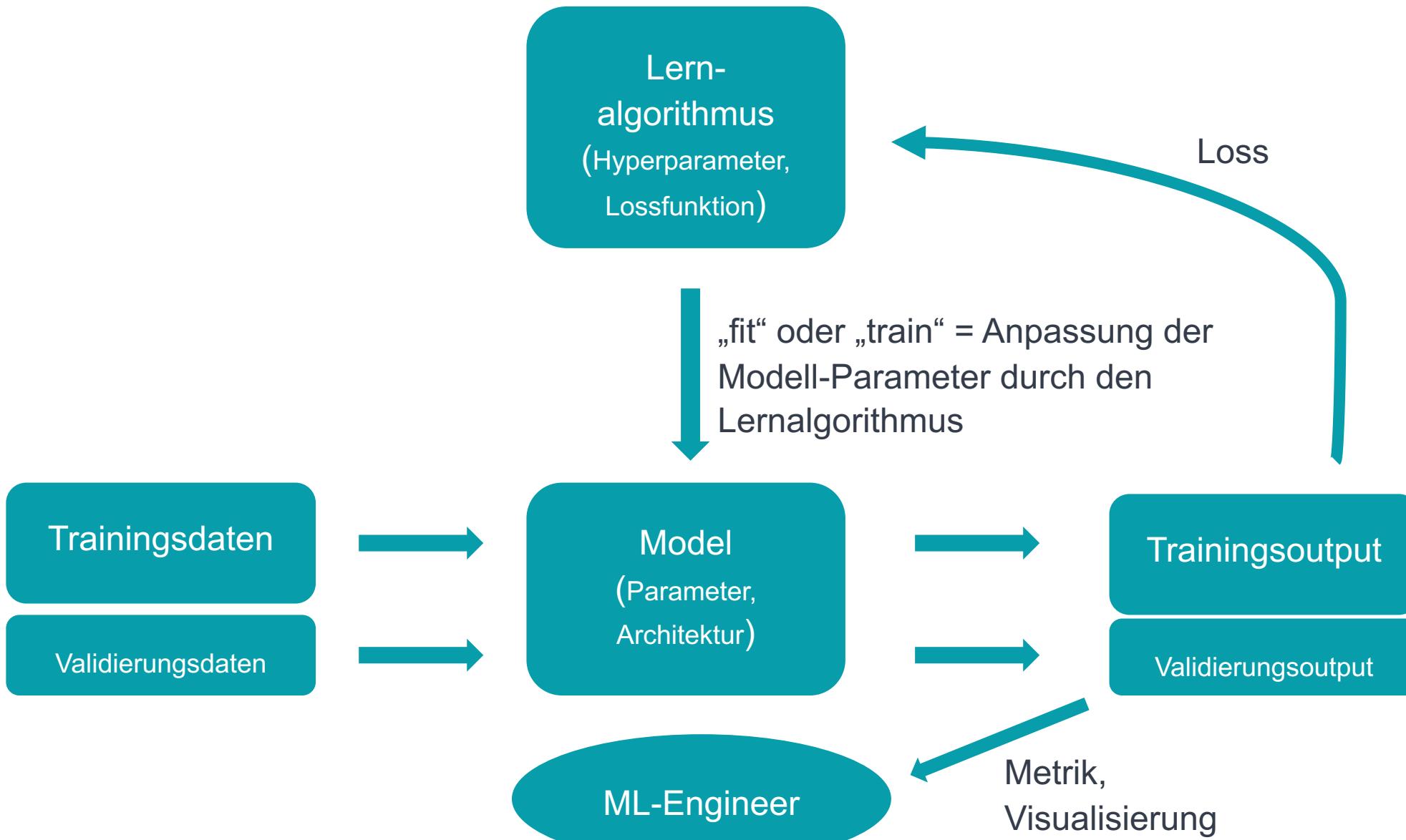


Herausforderungen Machine Learning

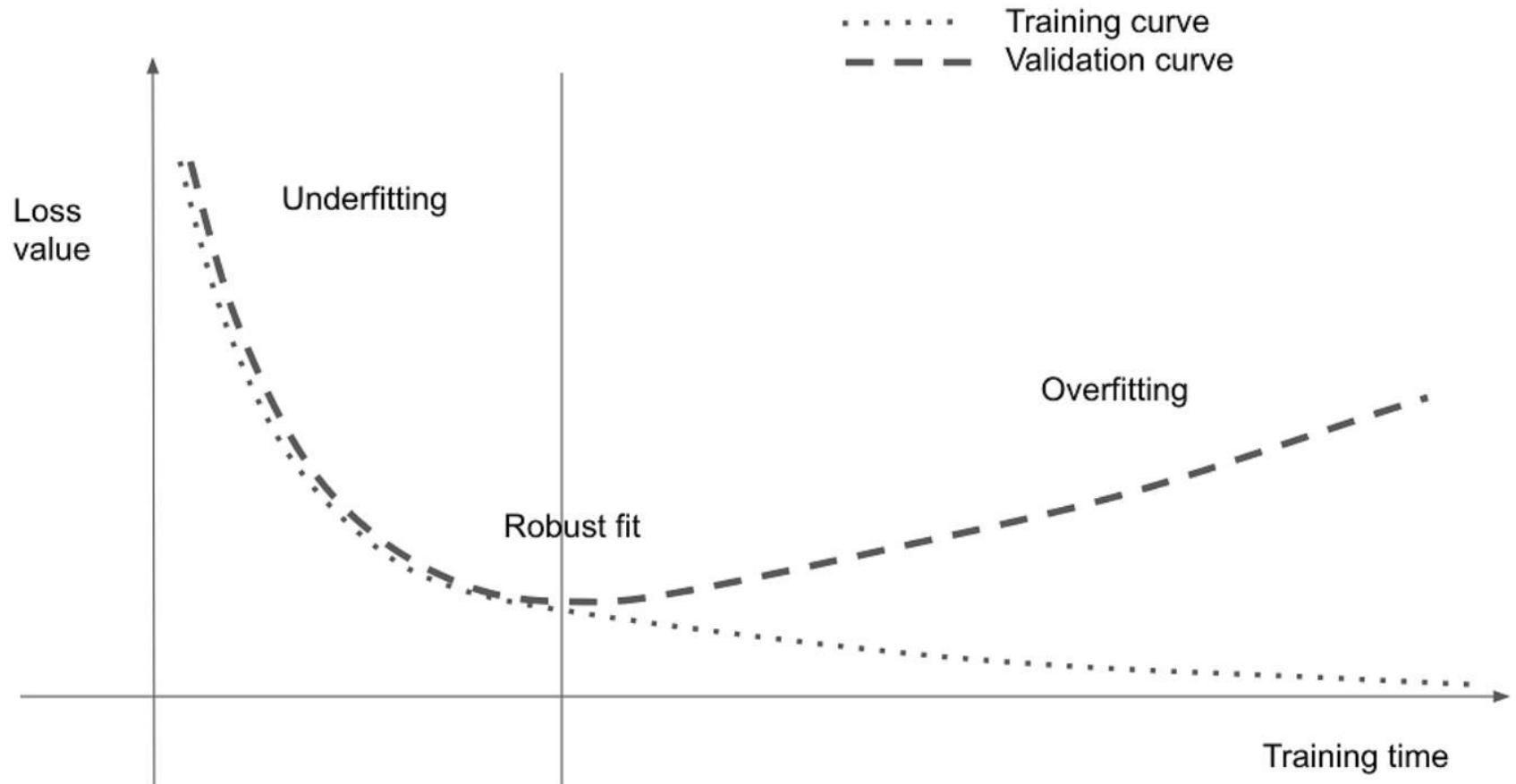
- Machine Learning und Predictive Applications bringen im Vergleich zu traditionellen Software Systemen **neue Herausforderungen** mit sich
 - Unzureichende Trainingsdaten
 - Trainingsdaten sind nicht repräsentativ (sampling bias)
 - Schlechte Datenqualität (Fehlwerte → oder Fehlwerte, die nicht als solche erkennbar sind Bsp.: als „unknown“ gekennzeichnet)
 - Irrelevante oder unzureichende Features
 - Overfitting → grundlegendes Problem ist Bias/Variance Trade-Off
 - Silent failure → Bugs zeigen sich oft nur indirekt durch etwas schlechtere Performance
- „Garbage in, garbage out“ ist zu beachten !
→ Besonders bei der Automatisierung von Geschäftsentscheidungen



Schlüsselbegriffe Machine Learning



Typische Overfitting Story

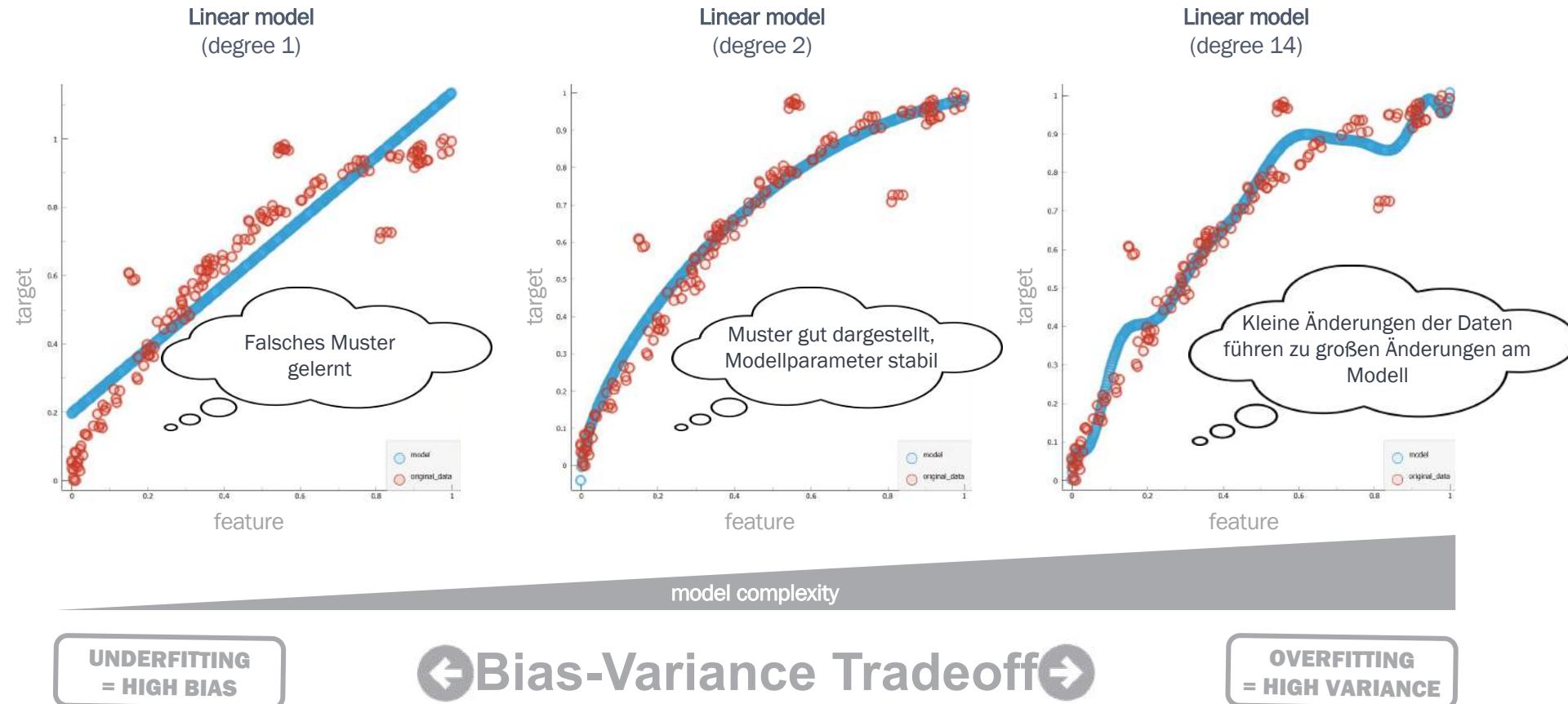


Live Demo

Orange – Bias Variance Tradeoff

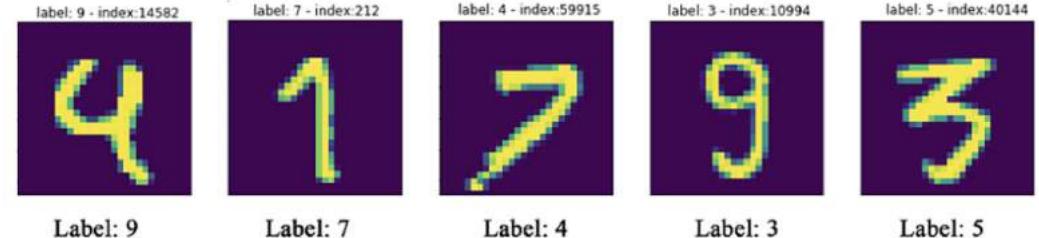
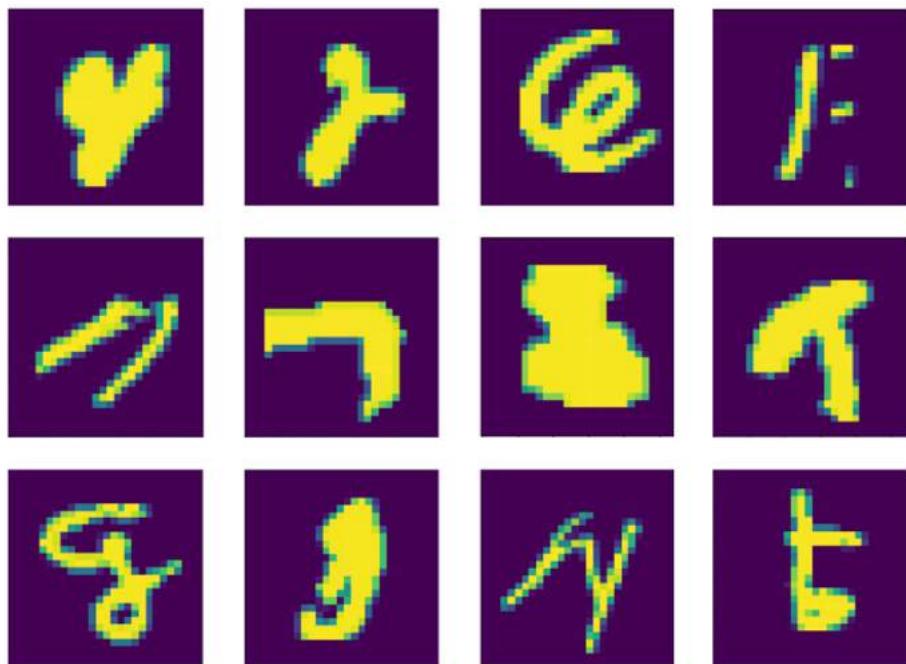


Bias Variance Tradeoff - Fazit

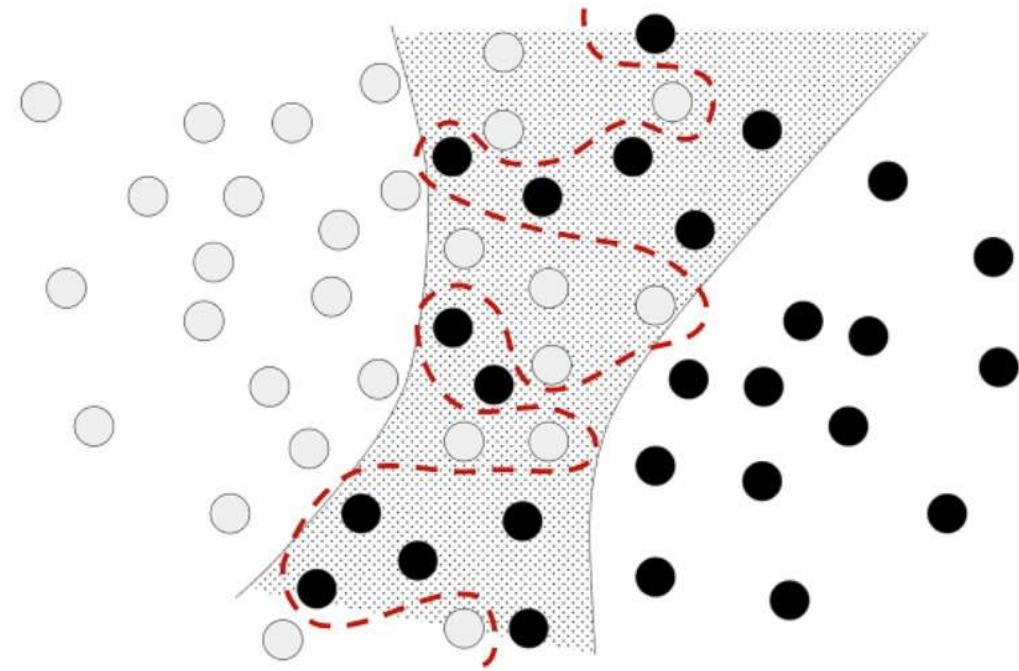
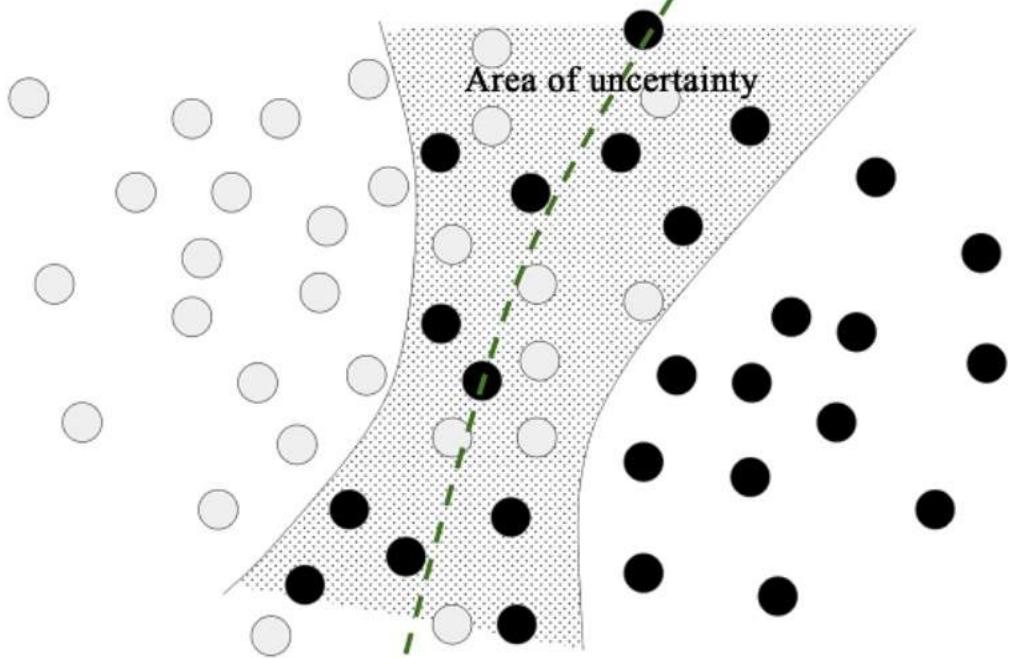


Noise auswendig lernen

- Falsche Labels
- Seltene Features
- Ambiguität
- Rauschen



Das Ergebnis: Schlechte Generalisierung



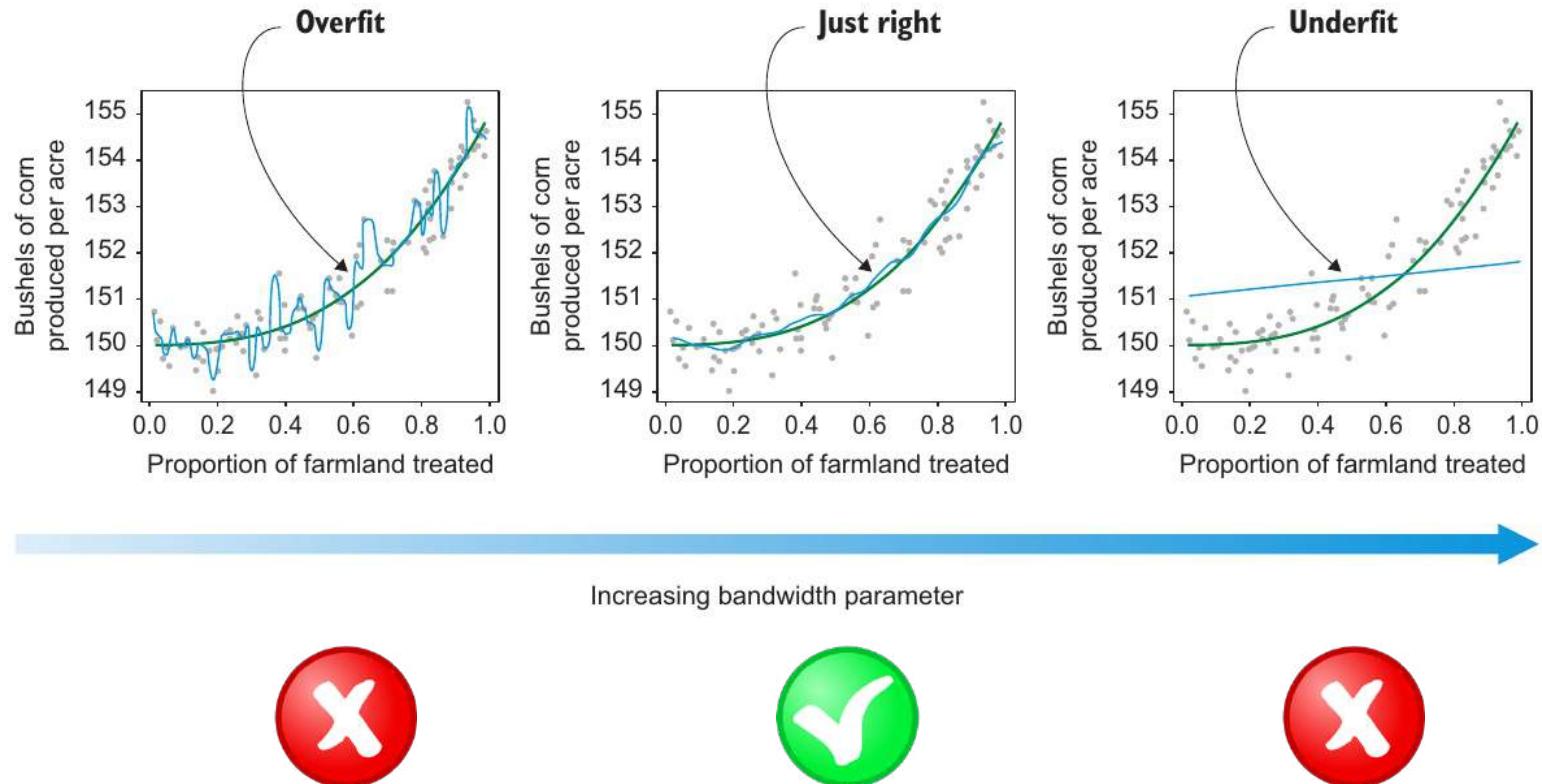
Einflüsse auf Overfitting

- Learning rate → Saturierung
- Datenmenge
- Datenkomplexität
- Modelgröße
- Regularisierung



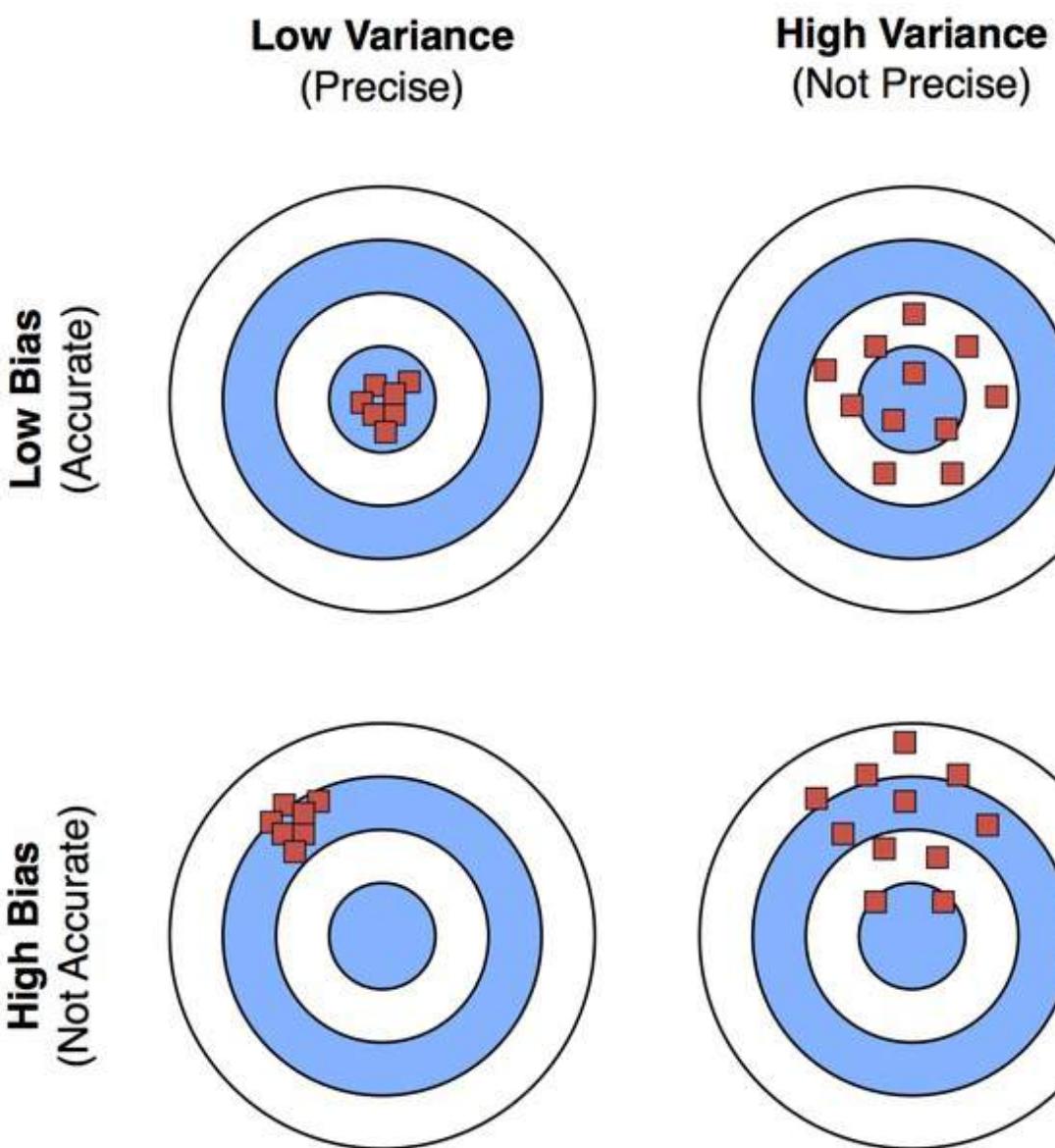
Overfitting und Underfitting

Lernt das Modell auswendig? Verallgemeinert es genau richtig? Oder lernt es zu wenig?



Bias-Variance-Tradeoff

Feste Struktur oder maximale Flexibilität?

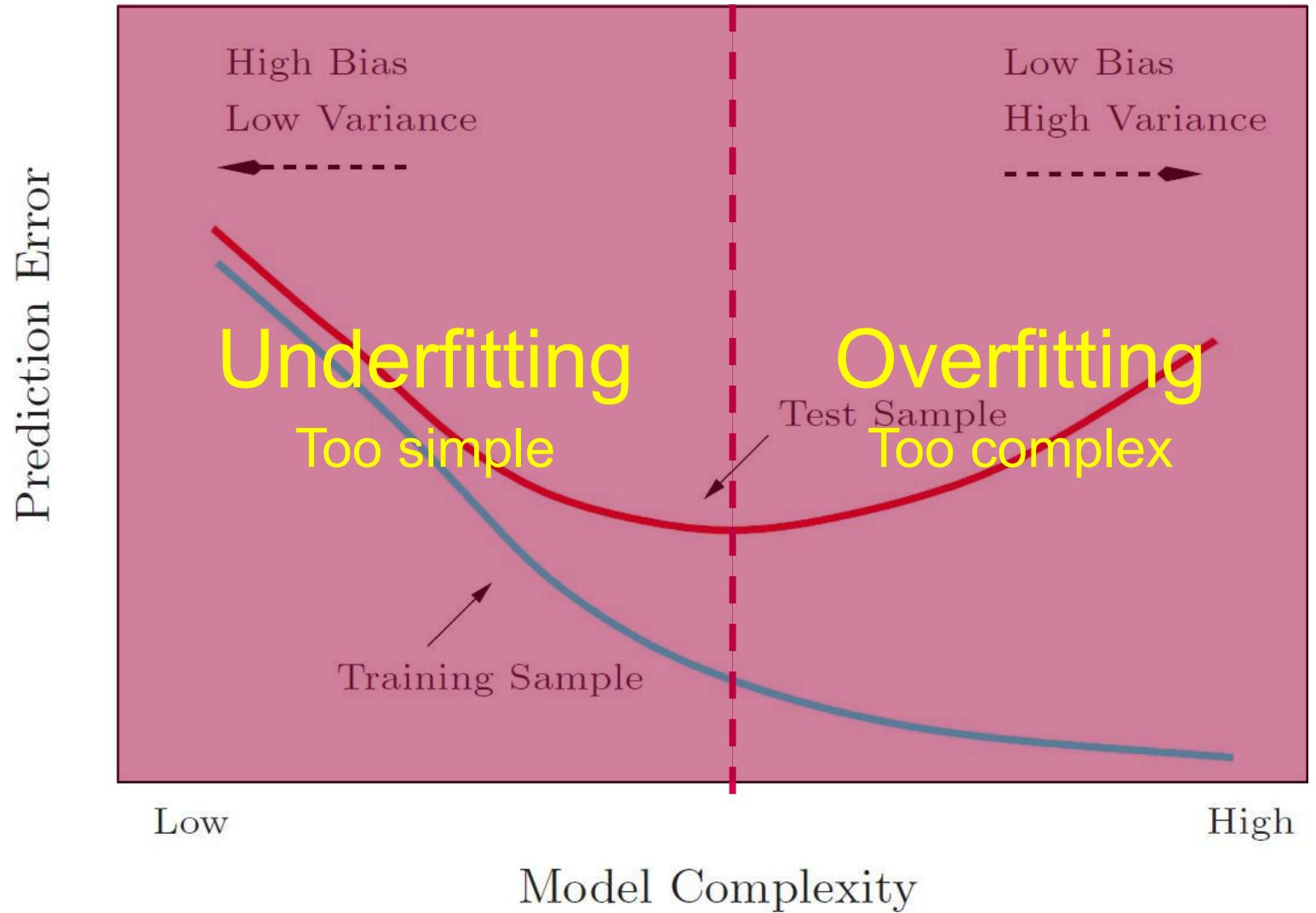


This work by Sebastian Raschka is licensed under a
Creative Commons Attribution 4.0 International License.



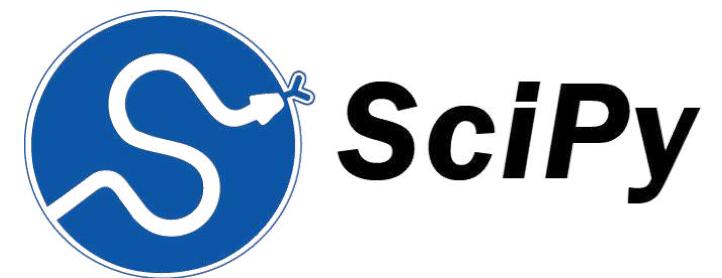
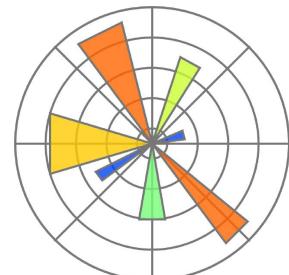
Bias-Variance-Tradeoff

Overfitting und andere Probleme



Der Machine Learning Tech Stack

Python und Python-Libraries die aufeinander aufbauen ...



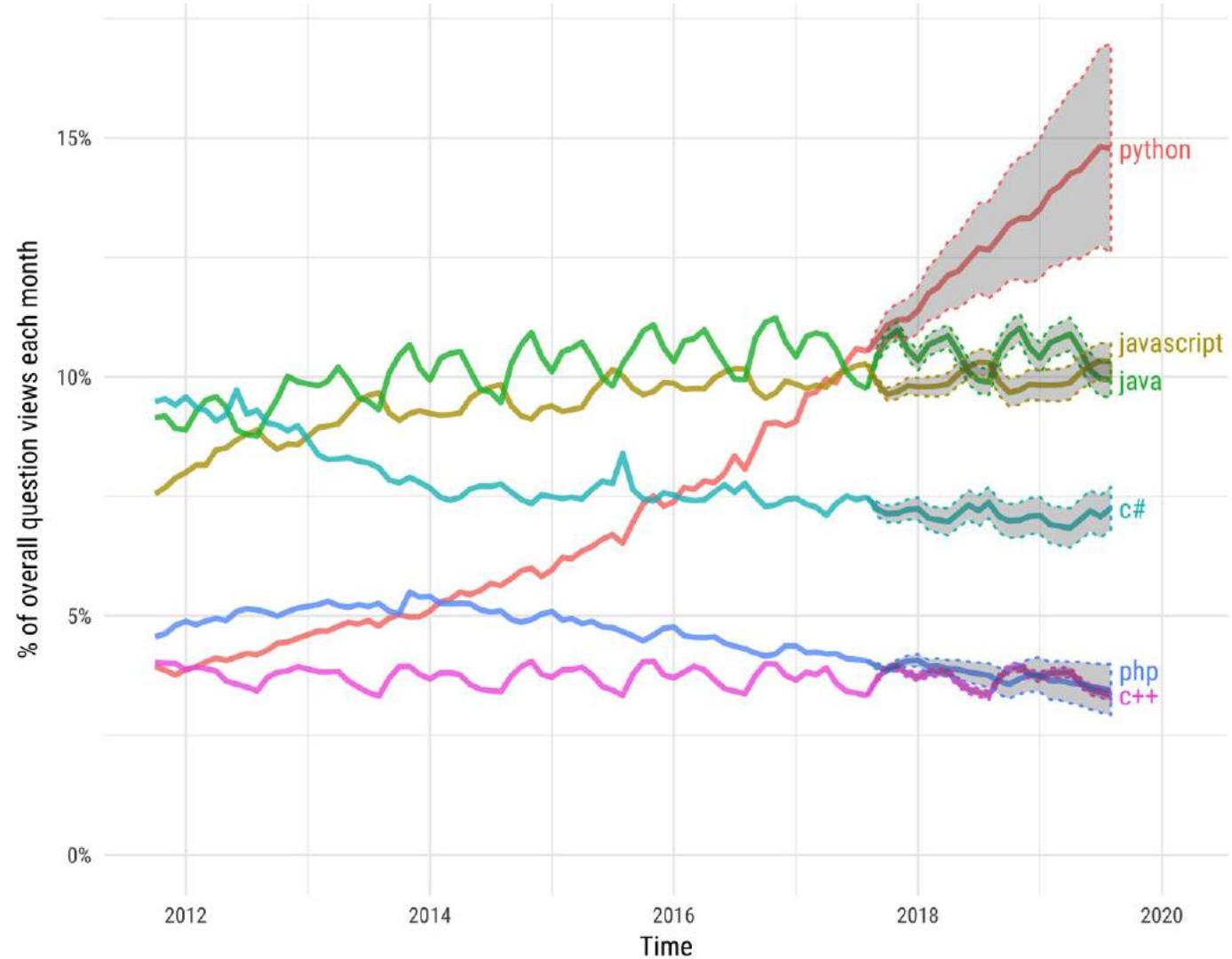
Quellen:
WikiCommons, Python.org,
fullstackpython.com,



Warum Python

Projections of future traffic for major programming languages

Future traffic is predicted with an STL model, along with an 80% prediction interval.



Source: Stack Overflow

<https://medium.com/@UdacityINDIA/why-use-python-for-machine-learning-e4b0b4457a77>



Was ist Data Science?

Ein relativ neues Berufsbild – Digitale Datenanalyse auf der Basis von Statistik und Machine Learning



Notebook für interaktive
Datenanalyse, Visualisierung
und Annotation.



Was ist Data Science?

Ein relativ neues Berufsbild – Digitale Datenanalyse auf der Basis von Statistik und Machine Learning



„Numerical Python“ – optimierte Algorithmen um mit Vektoren und Matrizen umzugehen.



Was ist Data Science?

Ein relativ neues Berufsbild – Digitale Datenanalyse auf der Basis von Statistik und Machine Learning



Library um tabellarische Daten zu
manipulieren. Excel in mächtig.



Was ist Data Science?

Ein relativ neues Berufsbild – Digitale Datenanalyse auf der Basis von Statistik und Machine Learning



Standard Visualisierungslibrary in Python



seaborn

Wenn es gut aussehen soll ...

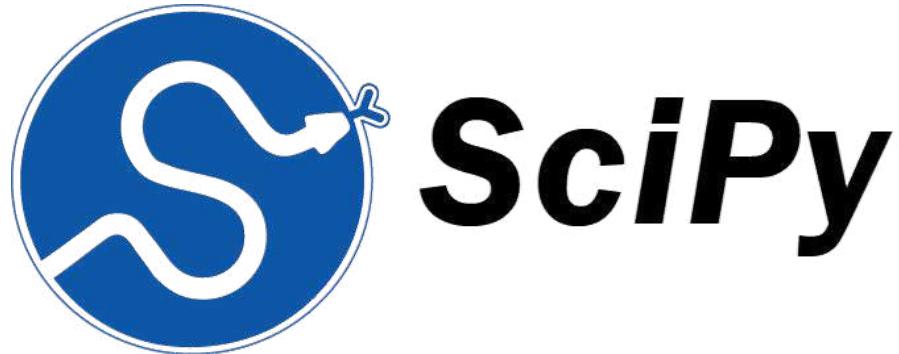


ScikitLearn und SciPy

Die beiden wichtigsten “reinen” Machine Learning Libraries



Machine Learning basiert,
d.h. um das Konzept eines
Modells herum aufgebaut.
Baut auf SciPy auf.



Scientific Computing basiert,
d.h. stärkerer Fokus auf Datenanalysen,
z.B. FFT, PCA, Integration, LinAlg,
Signalverarbeitung, Solver etc.



Scikit-learn Design

<https://arxiv.org/abs/1309.0238>

Scikit-Learn's API is remarkably well designed. These are the main design principles:

Consistency: All objects share a consistent and simple interface:

- **Estimators** Any object that can estimate some parameters based on a dataset is called an *estimator* (e.g., an imputer is an estimator). The estimation itself is performed by the `fit()` method, and it takes only a dataset as a parameter (or two for supervised learning algorithms; the second dataset contains the labels). Any other parameter needed to guide the estimation process is considered a hyperparameter (such as an imputer's strategy), and it must be set as an instance variable (generally via a constructor parameter).
- **Transformers** Some estimators (such as an imputer) can also transform a dataset; these are called *transformers*. Once again, the API is simple: the transformation is performed by the `transform()` method with the dataset to transform as a parameter. It returns the transformed dataset. This transformation generally relies on the learned parameters, as is the case for an imputer. All transformers also have a convenience method called `fit_transform()` that is equivalent to calling `fit()` and then `transform()` (but sometimes `fit_transform()` is optimized and runs much faster).
- **Predictors** Finally, some estimators, given a dataset, are capable of making predictions; they are called *predictors*. For example, the `LinearRegression` model in the previous chapter was a predictor: given a country's GDP per capita, it predicted life satisfaction. A predictor has a `predict()` method that takes a dataset of new instances and returns a dataset of corresponding predictions. It also has a `score()` method that measures the quality of the predictions, given a test set (and the corresponding labels, in the case of supervised learning algorithms).¹⁸

Inspection: All the estimator's hyperparameters are accessible directly via public instance variables (e.g., `imputer.strategy`), and all the estimator's learned parameters are accessible via public instance variables with an underscore suffix (e.g., `imputer.statistics_`).

Nonproliferation of classes Datasets are represented as NumPy arrays or SciPy sparse matrices, instead of homemade classes. Hyperparameters are just regular Python strings or numbers.

Composition Existing building blocks are reused as much as possible. For example, it is easy to create a Pipeline estimator from an arbitrary sequence of transformers followed by a final estimator, as we will see.

Sensible defaults Scikit-Learn provides reasonable default values for most parameters, making it easy to quickly create a baseline working system.



Entwicklungsumgebung

Python, Jupyter Notebook, Keras, Colab ...



Jupyter Notebook – Google Colab



- Notebook vereint Text und Code in einzeln ausführbaren Code-Zellen
- Interaktiv, gut debugbar und gut dokumentierbar

Colab ist ein kostenloses online Notebook von Google

<https://colab.research.google.com>

- Mit lokaler Laufzeit verbinden → Datei → Neues Notebook
- Laufzeit → Laufzeit ändern → GPU
- !pip install package_name , aber Keras und Tf sind vorinstalliert.



Der Plan

- Jeder öffnet ein Jupyter Notebook
 - Startmenü → Anaconda Prompt → >jupyter notebook
- Wir versuchen tensorflow zu importieren
- Außerdem numpy, keras, matplotlib ... sollte von der IT vorinstalliert sein
- Eventuell Zugriff auf Colab ausprobieren (Gmail Account vorhanden?)
- Alternativ: Direkt auf Eurem Computer arbeiten. Macht eventuell mehr Sinn.



Jupyter Basics und Markdown

- # Überschrift 1, ## Überschrift 1.1
- Code kann zur Illustration hinzugefügt werden, wird nicht ausgeführt
 - `''' ... '''` kann für Code verwendet werden.
- LaTeX-Gleichungen zwischen \$... \$ oder \$\$... \$\$
- *italic* **bold**
- HTML, Videoembeddings, etc.



Pandas Basics

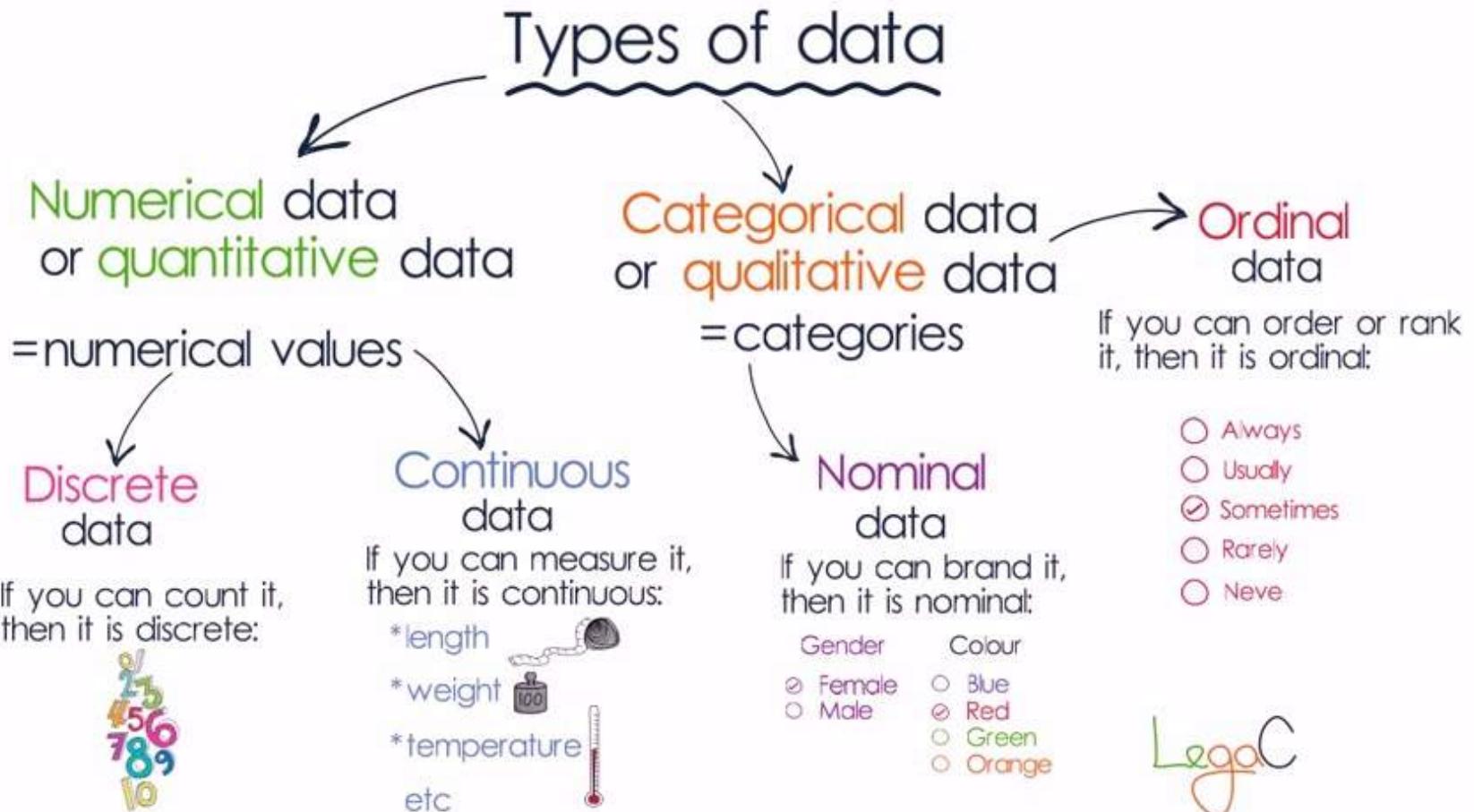
Library um tabellarische Daten zu manipulieren. Excel in mächtig.



Mit Daten umgehen



Wie sehen unsere Daten aus?



Pre-Processing - Überblick

- Daten sind meistens in vielfacher Hinsicht fehlerhaft, z.B. durch Messfehler, Menschliches Versagen oder Computerfehler, Übertragungsfehler, etc.
 - Unvollständig: Attribute fehlen
 - Noisy: verzerrte Daten, Error, Outliers (z.B.: Größe = -100cm)
 - Inkonsistenzen: Diskrepanzen zwischen verschiedenen Einträgen (Age vs. Birthday)
 - Intentional Errors: Vesteckte Fehlwerte
z.B.: jeder fehlende Eintrag für „birthday“ = „January 1st“

→ Ursprung für Fehlwerte kann vielerlei sein.



Warum Pre-Processing?

- Schlechte Datenqualität → geringe Qualität der darauf aufbauenden Data Mining und Machine Learning Ergebnisse
- Verbessert die Performance von Predictive Applications → z.B.: Accuracy
- Modellierung setzt gute Datenqualität voraus → Klassifikationsalgorithmen können grundsätzlich nicht mit Fehlwerten umgehen

→ **Datenaufbereitung, -säuberung und –transformation** beanspruchen den Hauptteil der Arbeit bei Predictive Applications und Data Science Projekten.



Hauptaufgaben des Pre-Processing

- **Datenbereinigung**
 - Füllen von Fehlerten, glätten von „noisy“ data, identifizieren und entfernen von Outliern und „noisy“ data, auflösen von Inkonsistenzen.
- **Datenintegration**
 - Integration von mehreren Datenbanken oder Files (Bsp.: Merging)
- **Datentransformation**
 - Normalisierung und Aggregation
 - Datendiskretisierung
- **Datenreduktion**
 - Reduzierung des Datenvolumens mit Beibehaltung der selben analytischen Ergebnisse



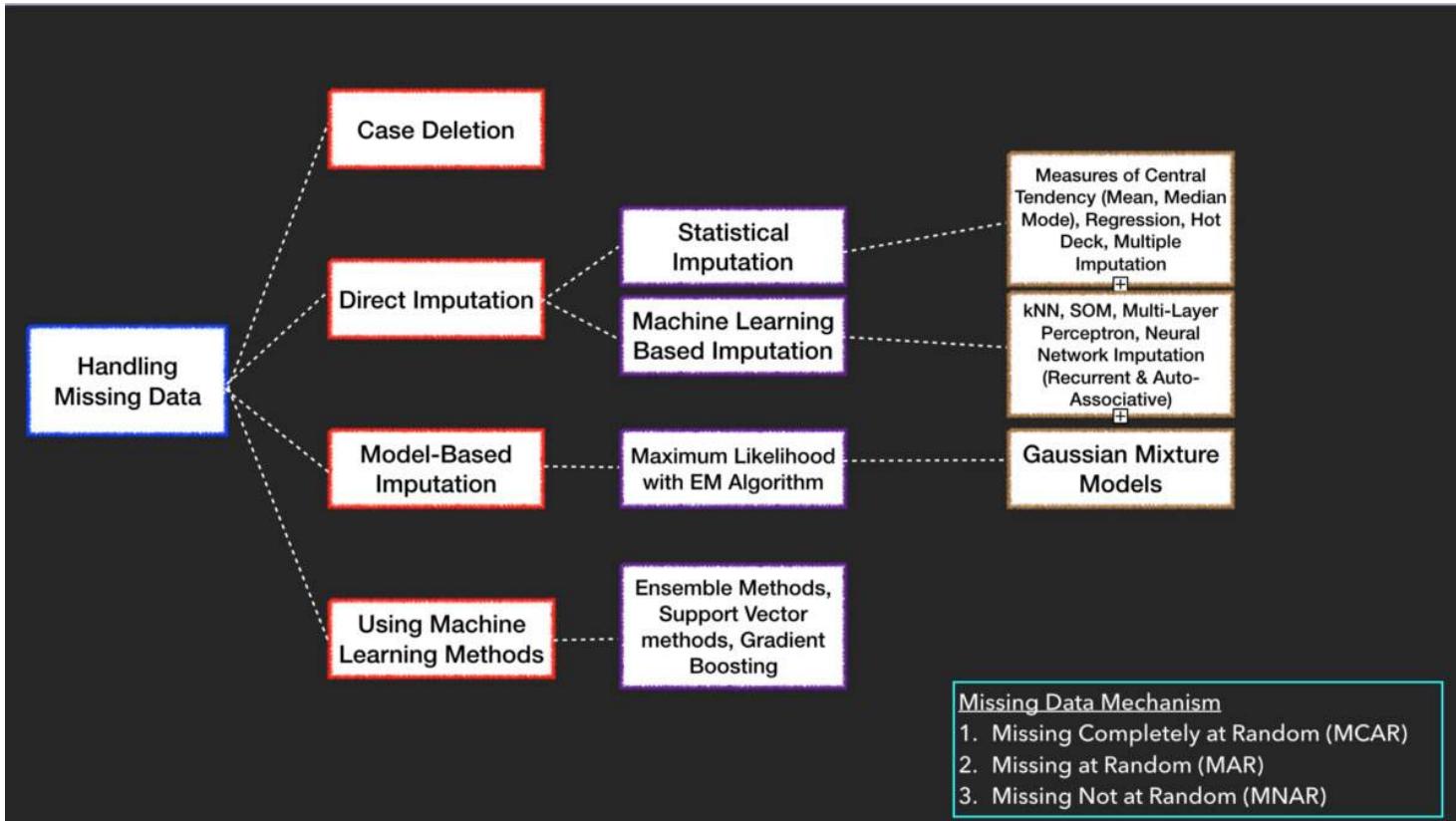
Datenbereinigung - Fehlwerte

- Aus der Statistik werden grundsätzlich 3 verschiedenen Klassen unterschieden
 - MCAR: Missing Completely at Random
 - Die Wahrscheinlichkeit, dass ein Wert fehlt ist nicht abhängig von den vorliegenden Datenwerten, noch von den fehlenden Datenwerten.
 - Beispiel: Ausfall eines Temperatursensors aufgrund technischer Probleme tritt zufällig auf
 - MAR: Missing at Random
 - Die Wahrscheinlichkeit, dass ein Wert fehlt ist teilweise von den anderen vorliegenden Daten abhängig aber nicht von irgendwelchen anderen Fehlwerten.
 - Beispiel: Ausfall eines Temperatursensors eher bei Nacht, trotzdem noch zufällig.
 - MNAR: Missing Not at Random
 - Die Wahrscheinlichkeit, dass eine Wert fehlt liegt an den Fehlwerten selbst.
 - Beispiel: Ausfall eines Temperatursensors wahrscheinlicher bei extremen Temperaturen.



Datenbereinigung – Umgang mit Fehlwerten

• Überblick

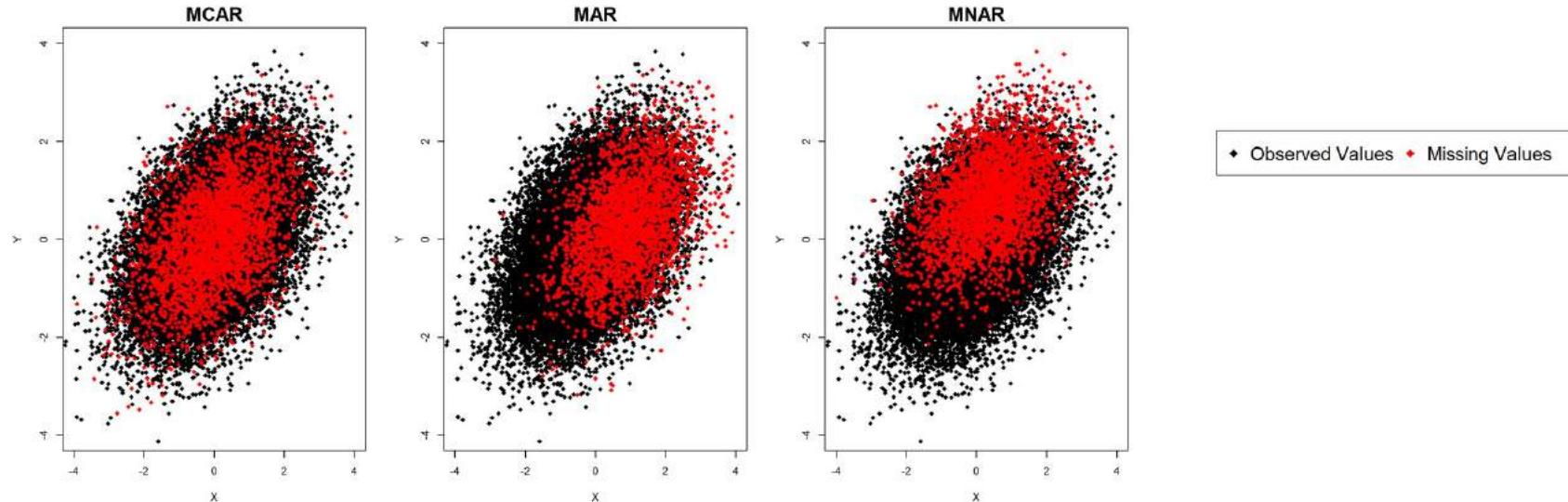


Source: <https://medium.com/ibm-data-science-experience/missing-data-conundrum-exploration-and-imputation-techniques-9f40abe0fd87>



Datenbereinigung – Kategorien von Fehlwerten

Beispiel: Studie/Umfrage



MCAR

- P hängt weder von X noch von Y ab.
- *Beispiel: Manche Antworten wurden ausversehen gelöscht.*

MAR

- P hängt von X ab, aber nicht von Y.
- *Beispiel: Befragte mit höherem Alter haben eine geringere Wahrscheinlichkeit ihre politische Meinung zu teilen*

MNAR

- P hängt von Y (und eventuell von X) ab.
- *Beispiel: Befragte mit höherem Einkommen geben ihr Einkommen seltener an.*



Datenbereinigung - MCAR

Arten von „Deletion“ (Lösung/Entfernung)

- **Listwise deletion**

User	Device	OS	Transactions
A	Mobile	NA	5
B	Mobile	Android	3
C	NA	iOS	2
D	Tablet	Android	1
E	Mobile	iOS	4

- **Pairwise deletion**

User	Device	OS	Transactions
A	Mobile	NA	5
B	Mobile	Android	3
C	NA	iOS	2
D	Tablet	Android	1
E	Mobile	iOS	4

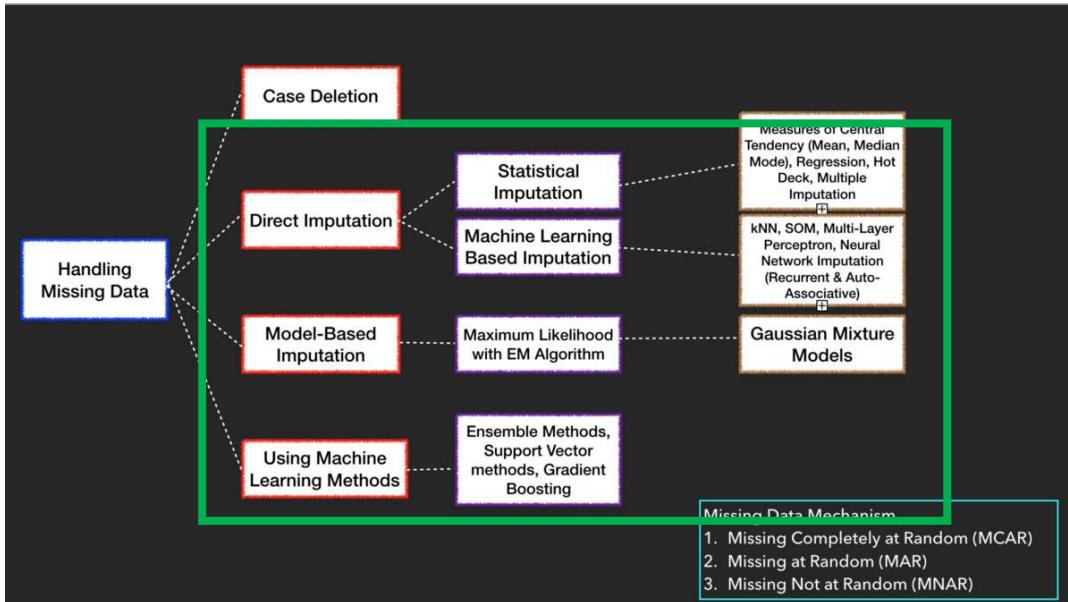
- **Löschen von Spalten (nicht empfohlen)**

- Nur falls >60% der Daten fehlen oder das Attribut nicht von Bedeutung ist



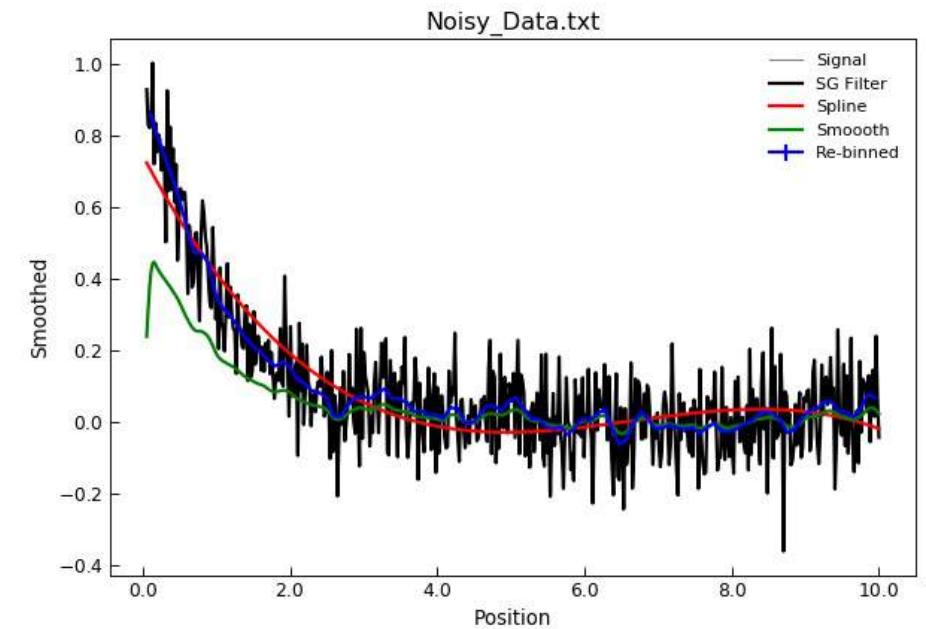
Datenaufbereitung - Imputation

- Imputation ist der Prozess des Ersetzens von Fehlwerten mit Ersatzwerten.
 - Beispiel: Ersetzen von Fehlwerten mit dem Mittelwert einer Spalte.



Noisy Data - Überblick

- Noise: Zufälliger Fehler oder Varianz in einer gemessenen Variable
- Fehlerhafte Daten resultieren beispielsweise aus:
 - Fehlerhafte Datensammlungsgeräte/-tools
 - Dateneingabeproblem (typo)
 - Datenübertragungsprobleme
 - Technologie-Limitierung



Sources: <https://stoner-pythoncode.readthedocs.io/en/latest/UserGuide/analysisfile.html>



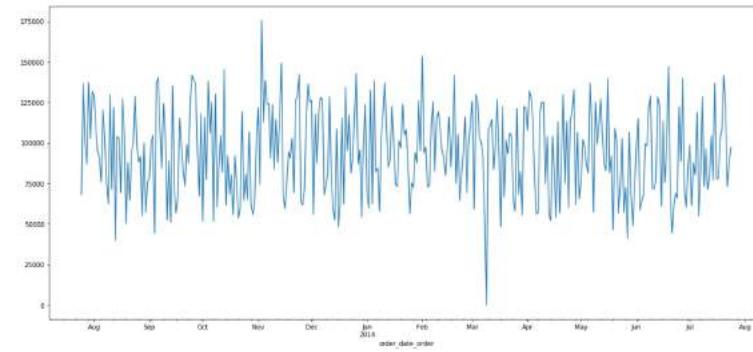
Noisy Data – Umgang mit Noisy Data

- Binning
 - Sortieren und Aufteilen der Daten in (gleichverteilte) bins
 - Dann kann nach dem Mittelwert/Median oder bin boundaries geglättet werden
- Gleitender Durchschnitt (Moving Average) & exponentielle Glättung (exponential smoothing)
- Regression
 - Glätten der Daten durch fitten von Regressionsfunktionen
- Untersuchen der Daten von Computer und Mensch
 - Entdecken von verdächtigen Werten und prüfen von Menschenhand (Bsp.: evtl. Outlier?)
 - EDA
- ...viele mehr

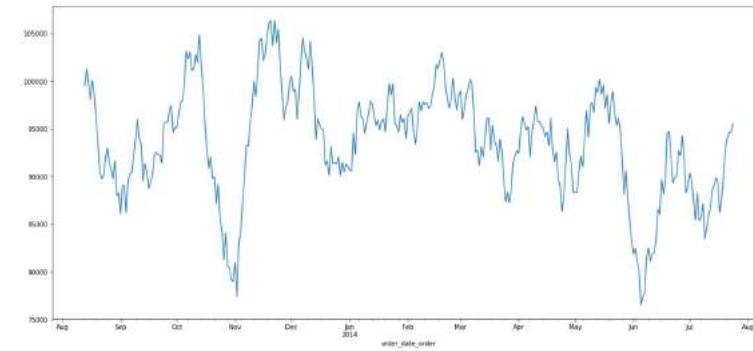


Noisy Data – Zeitreihen glätten

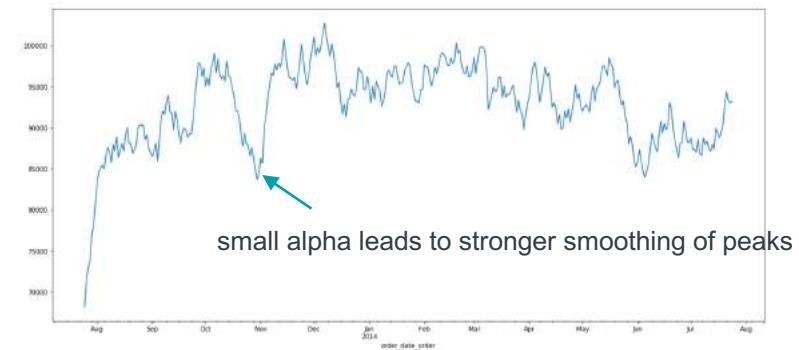
Original:



Rolling mean:



Simple exponential
smoothing (alpha=0.05)



Datentransformation - Normalisierung

- Normalisierung durch min-max Skalierung.
 - Normalisieren aller Werte zu einem Wertebereich von 0 bis 1.



Körpergröße (unskaliert)	Körpergröße (Min-Max skaliert)
190	1
175	0.57
155	0
176	0.6
164	0.257



Datentransformation - Normalisierung

- Normalisierung durch Standard Skalierung (z-Standardisierung)

STANDARDIZATION

Standardized feature value \tilde{x}_i = $\frac{x_i - \bar{x}}{\sigma}$

- Value of the i th observation
- Mean of the feature vector
- Standard deviation of the feature vector

Standardization is a common scaling method. \tilde{x}_i represents the number of standard deviations each value is from the mean value. It rescales a feature to have a mean of 0 and unit variance.

Chris Albon

<https://machinelearningflashcards.com/>

Standardization:

$$z = \frac{x - \mu}{\sigma}$$

with mean:

$$\mu = \frac{1}{N} \sum_{i=1}^N (x_i)$$

and standard deviation:

$$\sigma = \sqrt{\frac{1}{N} \sum_{i=1}^N (x_i - \mu)^2}$$

Unit variance bedeutet, dass die Standardabweichung der Variable gegen 1 geht wenn die sample size gegen unendlich geht.



Datentransformation - Encoding

Sample	Category
1	Human
2	Human
3	Penguin
4	Octopus
5	Alien
6	Octopus
7	Alien

LabelEncoding



Sample	Category	Numerical
1	Human	1
2	Human	1
3	Penguin	2
4	Octopus	3
5	Alien	4
6	Octopus	3
7	Alien	4

Achtung: Wir müssen vorsichtig sein mit Label Encoding! Machine Learning Algorithmen könnten das Attribut als Ordinal skaliert interpretieren.

OneHotEncoding



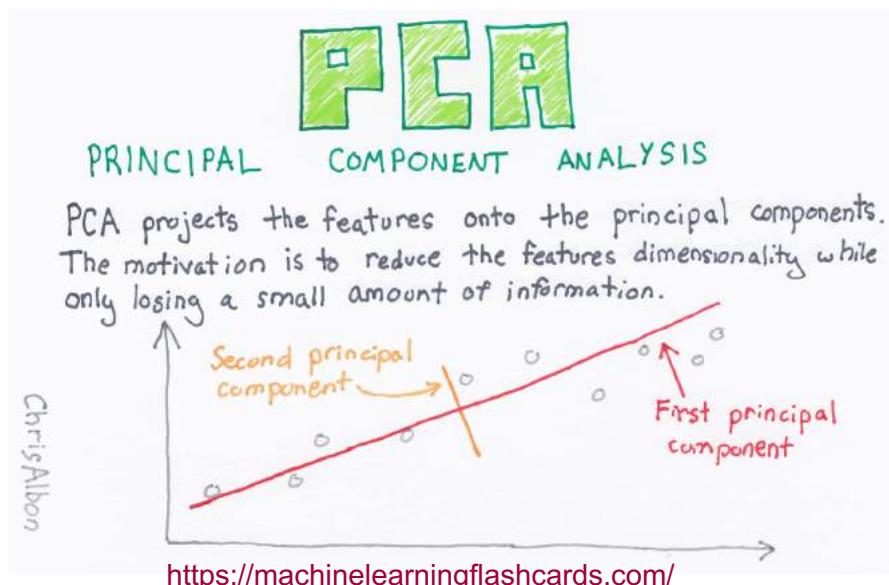
Sample	Human	Penguin	Octopus	Alien
1	1	0	0	0
2	1	0	0	0
3	0	1	0	0
4	0	0	1	0
5	0	0	0	1
6	0	0	1	0
7	0	0	0	1

One-Hot Encoding funktioniert besser mit normal nominalen/kategorischen Attributen für Machine Learning



Datenreduktion - PCA

- Ziel: Reduzieren der Anzahl von Features durch Identifikation von Korrelationen zwischen mehreren Features und reduzieren dieser zu einem (oder mehreren) Features, den sogenannten „Principle Components“
 - Reduzierter Rechenaufwand und Feature Dimensionen mit nur geringem Verlust an Informationen



Use Case Pre-Processing & Pre-Processing Pipelines



Machine Learning Basics

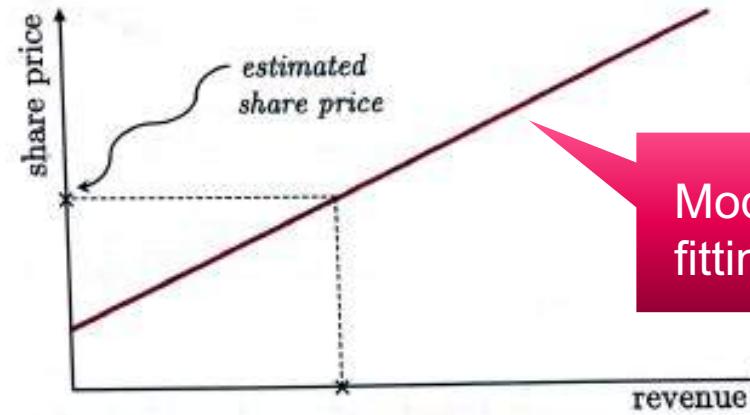
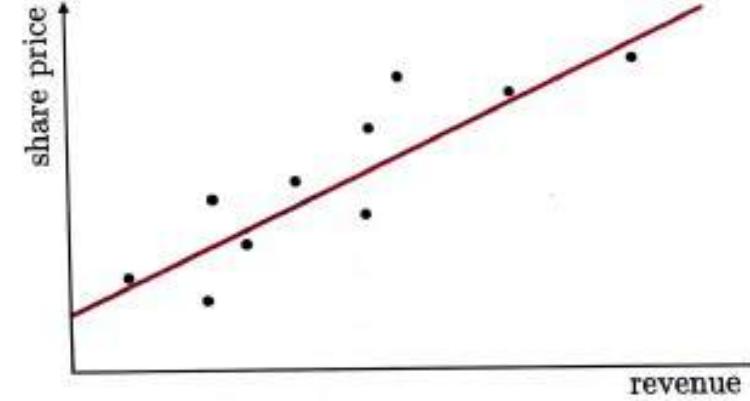
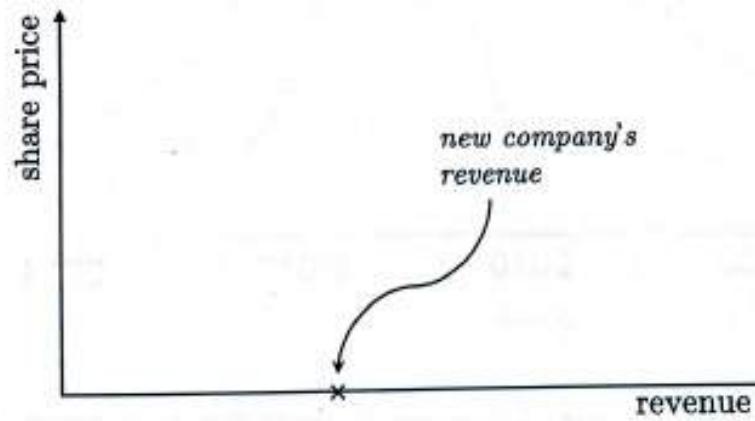
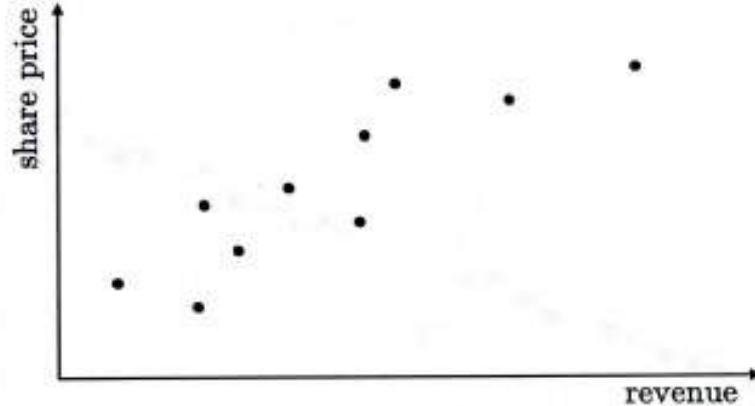


Regression vs Klassifizierung

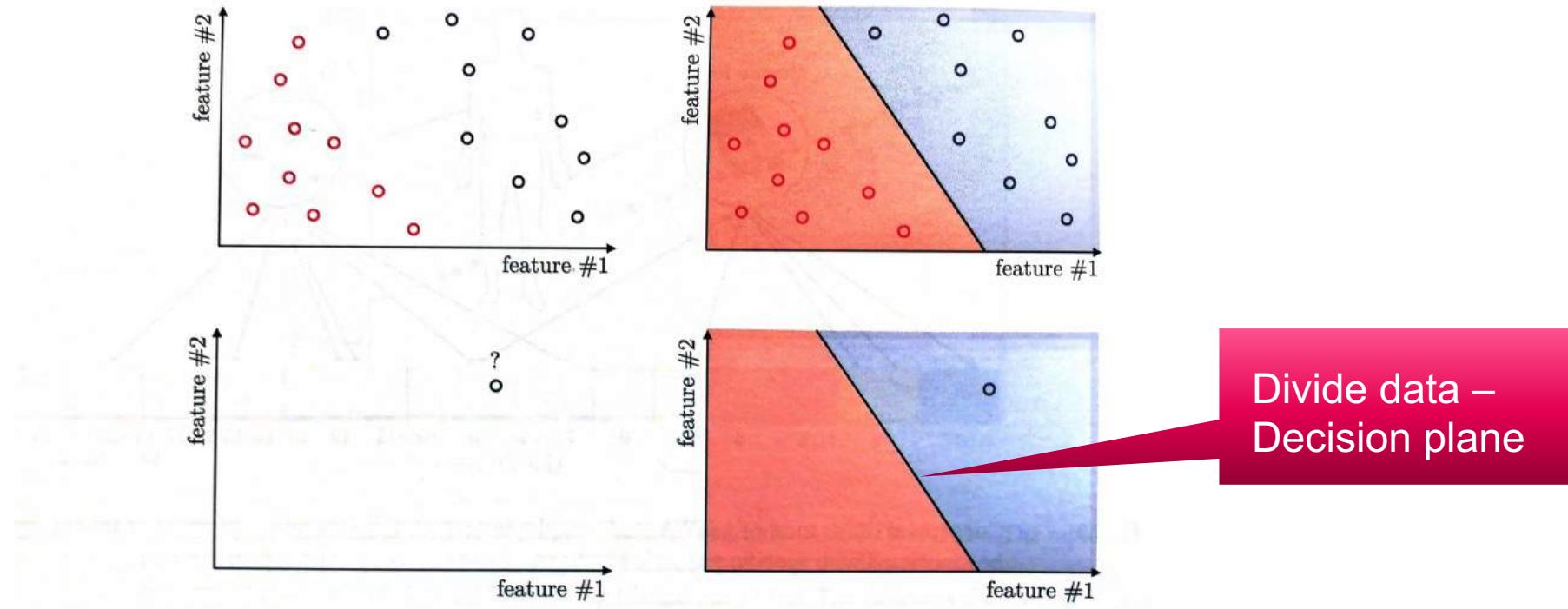
Klassen oder kontinuierliche Werte



Regression



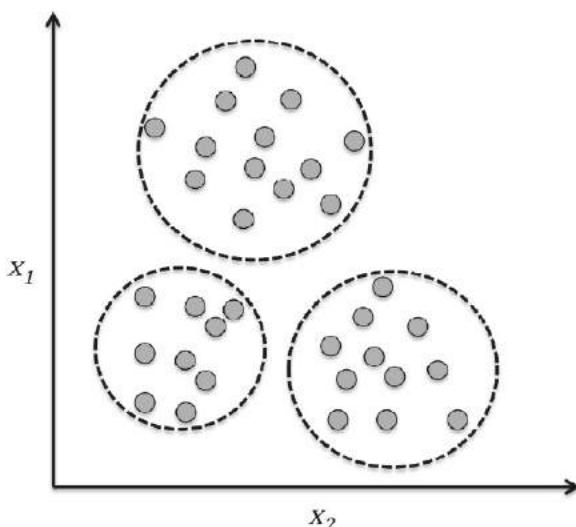
Klassifizierung



Unsupervised Learning - Clustering

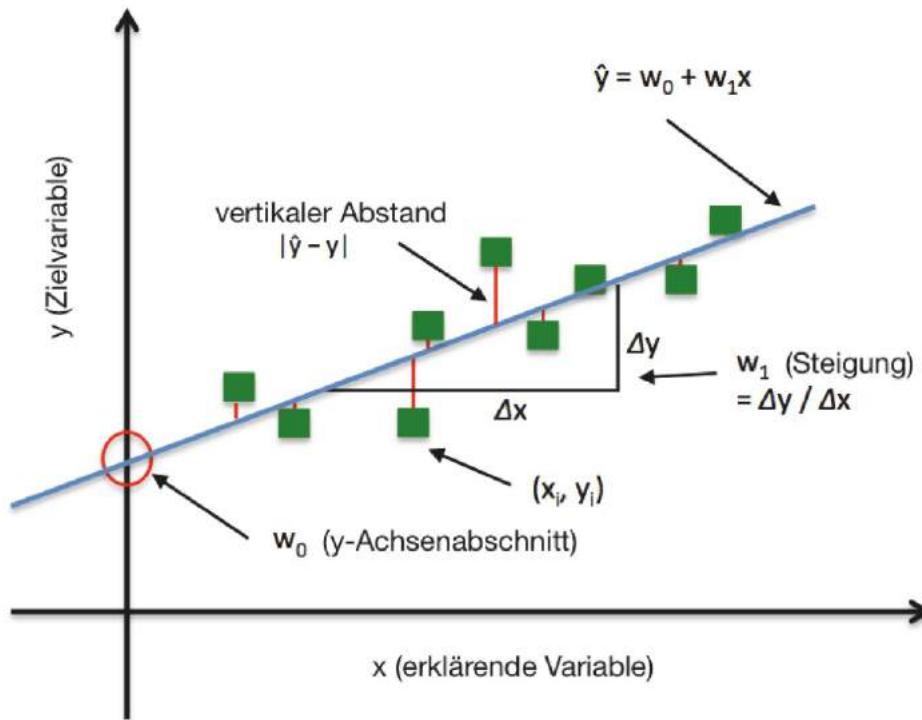
Clustering:

- Aufteilung von nicht gelabelten Daten in ähnliche Gruppen



Lineare Regression – Einfache Lineare Regression

- Dient zur Vorhersage eines Zielwerts y anhand eines Merkmals/Features x

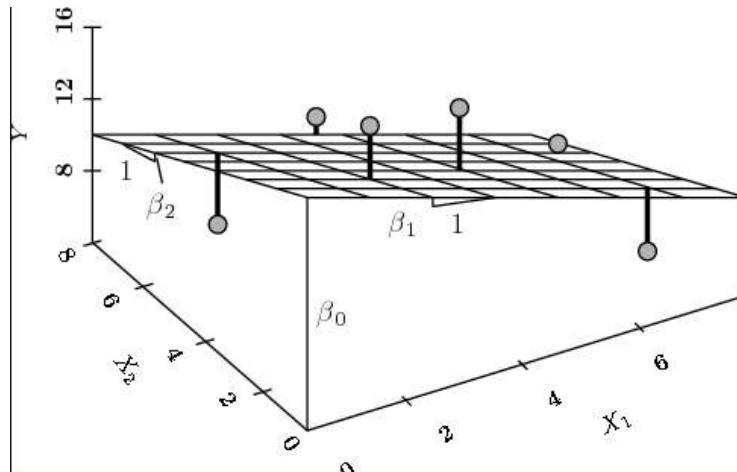


Lineare Regression – Multiple Lineare Regression

Es können mehrere Features zur Vorhersage der Zielvariable hinzugezogen werden

$$y = w_0x_0 + w_1x_1 + \dots + w_mx_m = \sum_{i=0}^m w_i x_i = w^T x$$

- Jedes Merkmal x hat seinen eigenen Koeffizienten und somit auch einen individuellen Einfluss auf die Zielvariable
- Aus Gerade wird eine Ebene:

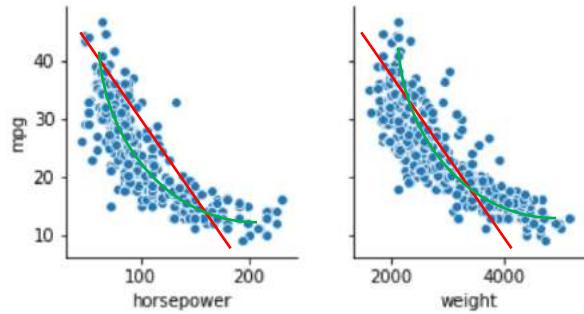


Source:https://www.google.com/search?q=Regressionsebene&source=lnms&tbo=isch&sa=X&ved=2ahUK EwiBn4GPv43wAhU4hf0HHVw6Av0Q_AUoA3oECAEQBQ&biw=1920&bih=937#imgrc=ICJ4cfQrU2z-FM



Polynomiale Regression

- Was wenn die Beziehung zwischen Target „y“ und unseren Features „komplexer“ ist?



- Polynomiale Regression erlaubt es die optimalen Koeffizienten für ein polynomial n-ten Grades zu ermitteln:

$$\hat{y} = w_0 + w_1x + w_2x^2 + \dots + w_nx^n$$



Polynomiale Regression

$$\hat{y} = w_0 + w_1x + w_2x^2 + \dots + w_nx^n$$

- Um die Koeffizienten w zu ermitteln, können wir die Werte für x^i berechnen und diese als normale numerische Features behandeln.
 - Es müssen keine weiteren Änderungen vorgenommen werden !

	horsepower	weight	horsepower^2	weight^2
0	130.0	3504	16900.0	12278016
1	165.0	3693	27225.0	13638249
2	150.0	3436	22500.0	11806096
3	150.0	3433	22500.0	11785489
4	140.0	3449	19600.0	11895601

Bei der Anwendung einer „normalen“ Linearen Regression werden diese Features dann verwendet und es werden bessere Resultate erzielt (falls eine nicht-lineare Beziehung besteht)



Use Case Lineare Regression



Korrelation – Pearson Korrelationskoeffizient

- Mithilfe der Korrelation lassen sich Beziehungen zwischen den vorliegenden Variablen finden
- Der Wert der Korrelation (Korrelationskoeffizient) zeigt wie stark eine Variable x, eine andere Variable y beeinflusst.
- Misst den **linearen Zusammenhang** zweier Variablen
- Korrelationskoeffizient variiert von -1 bis +1
 - +1 = Wenn x um 1 Einheit steigt, steigt y ebenfalls um 1 Einheit
 - 0 = keine Korrelation
 - -1 = Wenn x um 1 Einheit steigt, sinkt y um 1 Einheit

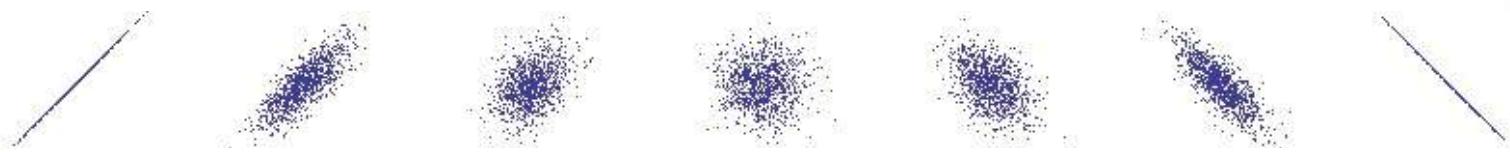
Berechnung beispielsweise durch Pearson Korrelationskoeffizienten:

$$\text{cor}_{x,y} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})} \sqrt{\sum_{i=1}^n (y_i - \bar{y})}} = \frac{\text{cov}_{x,y}}{\sigma_x \sigma_y}$$



Korrelation erkennen

a)



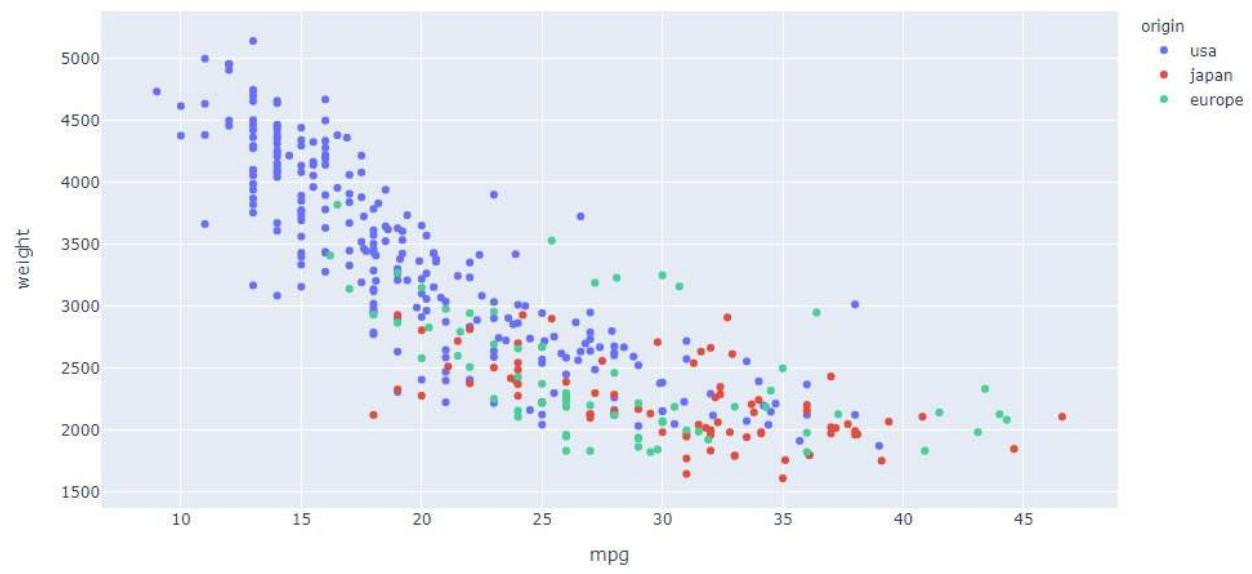
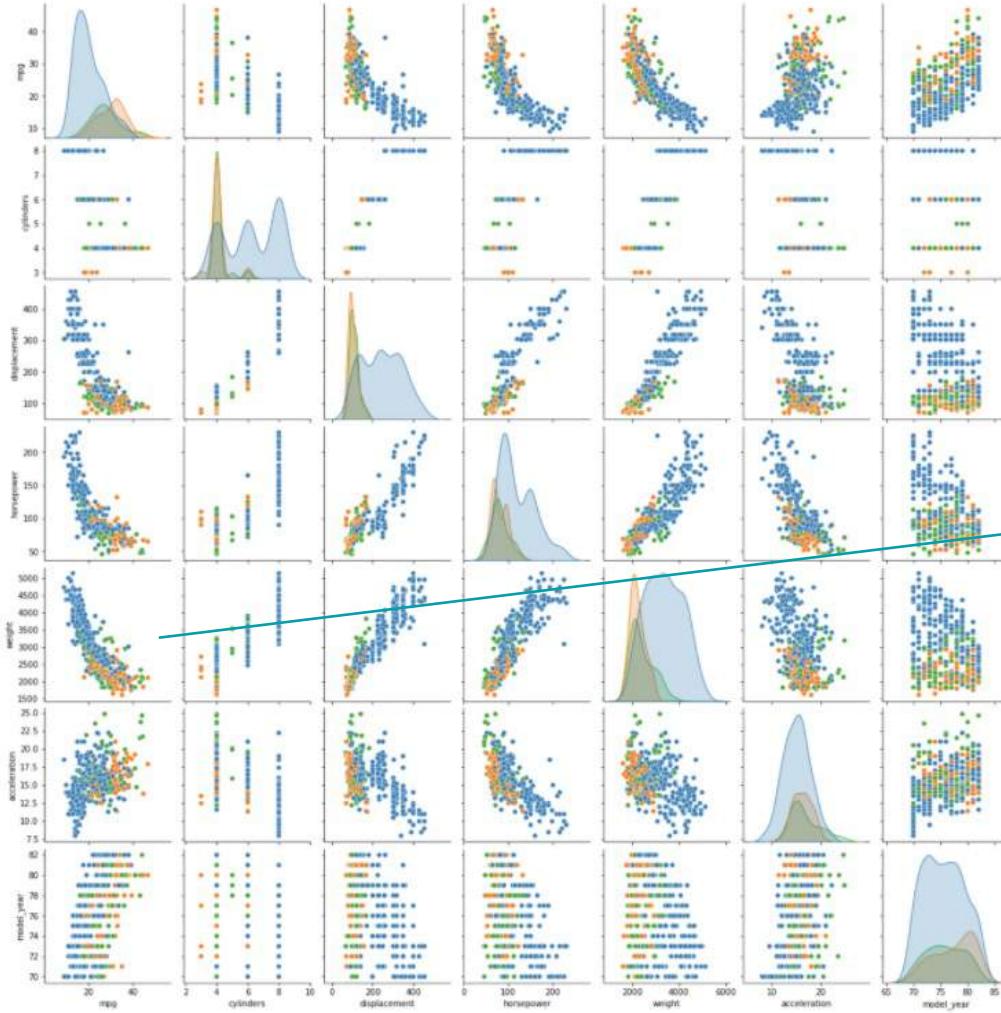
b)



c)

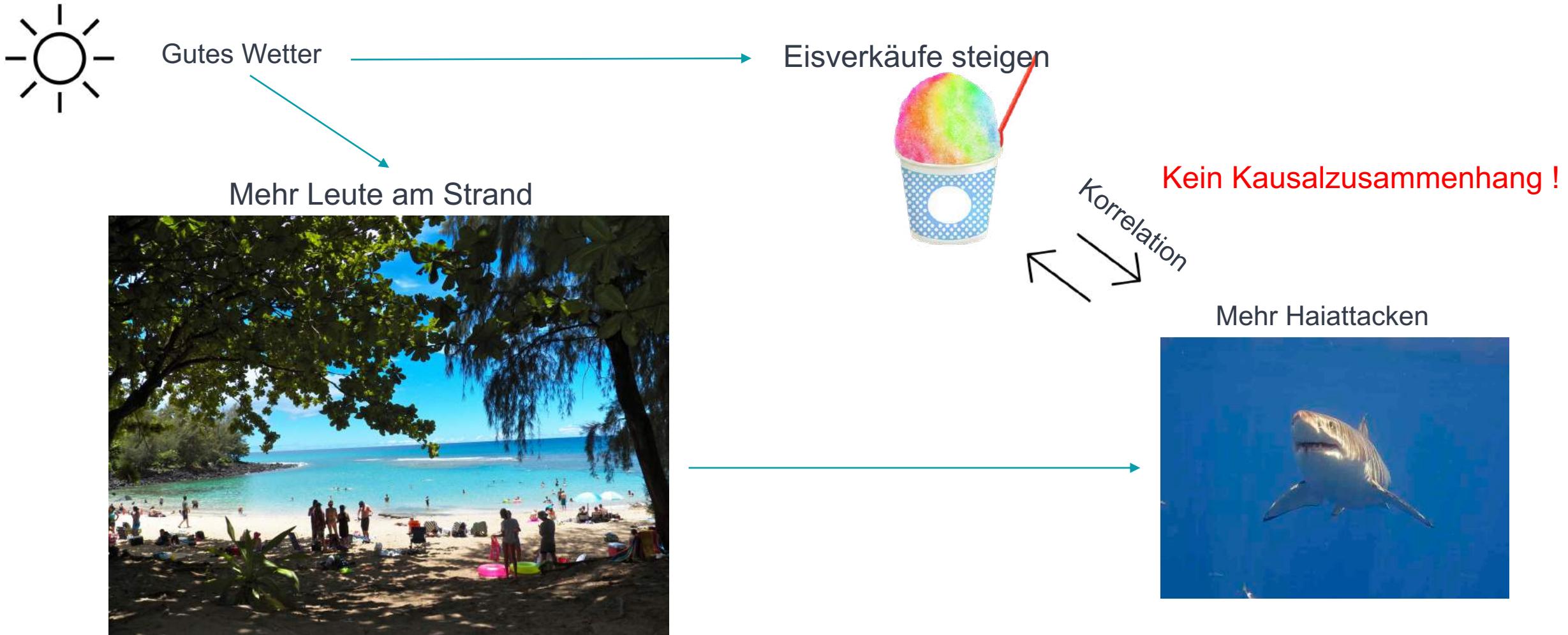


Korrelation – Pairplot zur Überprüfung von linearen Zusammenhängen



Korrelation ≠ Kausalität

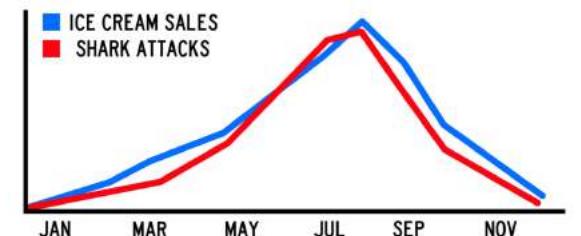
Beispiel Hawaii



Korrelation ≠ Kausalität

- Wenn ein Feature A und B korrelieren, kann man nicht zwangsläufig davon ausgehen, dass A der Auslöser/Treiber für B ist.
- Bsp.: Bei der Untersuchung einer demografischen Datenbank, wird eventuell ein Zusammenhang zwischen der „Anzahl von Krankenhäusern“ und der „Anzahl von Autodiebstählen in der Region“ gefunden, welche korrelieren.
 - Dies bedeutet nicht, dass das eine die Ursache für das andere ist.
 - Beide sind offensichtlich verbunden, jedoch durch ein drittes Attribut, nämlich „Bevölkerung“

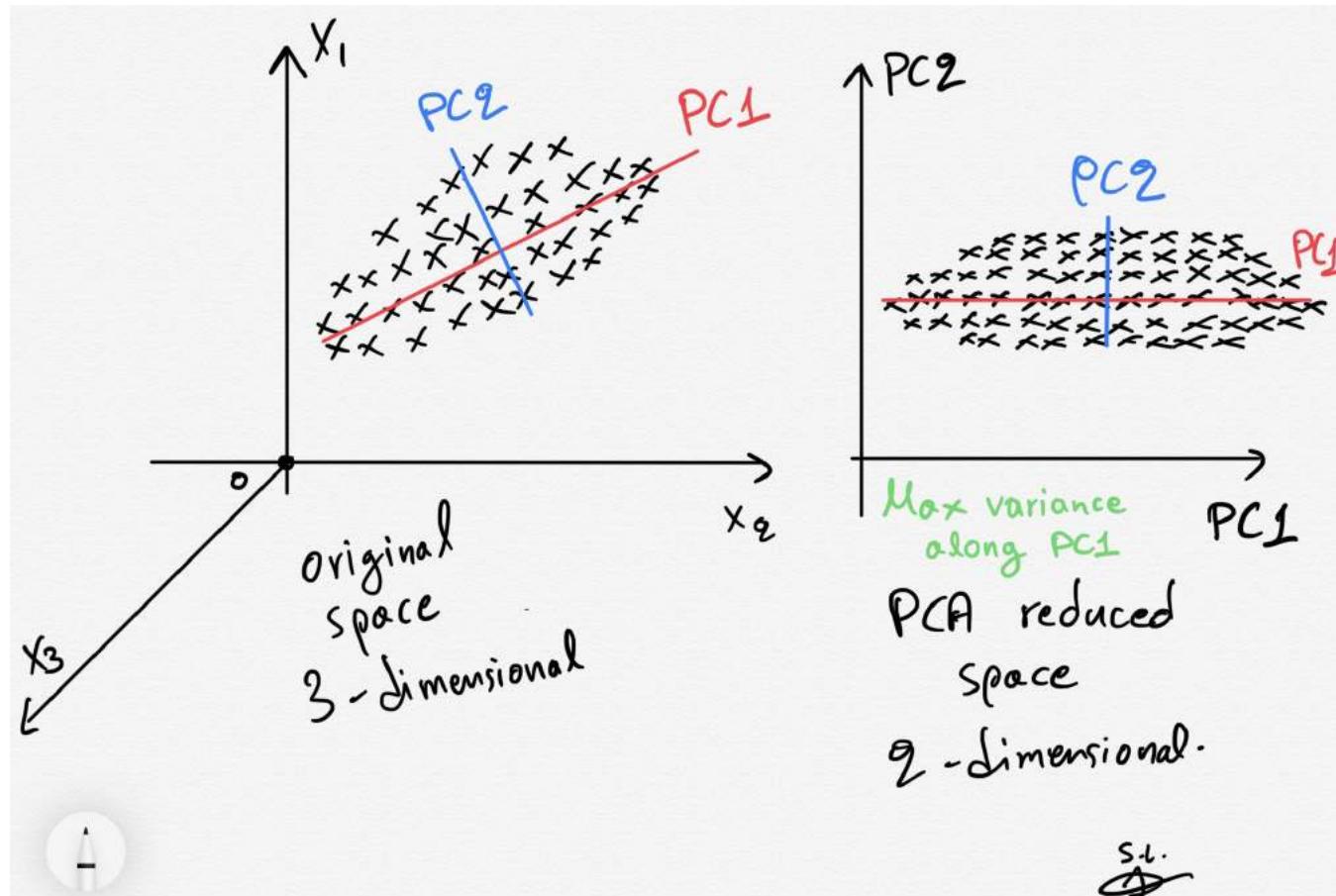
CORRELATION IS NOT CAUSATION!



Both ice cream sales and shark attacks increase when the weather is hot and sunny, but they are not caused by each other (they are caused by good weather, with lots of people at the beach, both eating ice cream and having a swim in the sea)



Lineare Regression, Pearson Korrelation und PCA sind drei mächtige Werkzeuge um linearen Zusammenhängen auf den Grund zu gehen



Source: <https://towardsdatascience.com/pca-clearly-explained-how-when-why-to-use-it-and-feature-importance-a-guide-in-python-7c274582c37e>

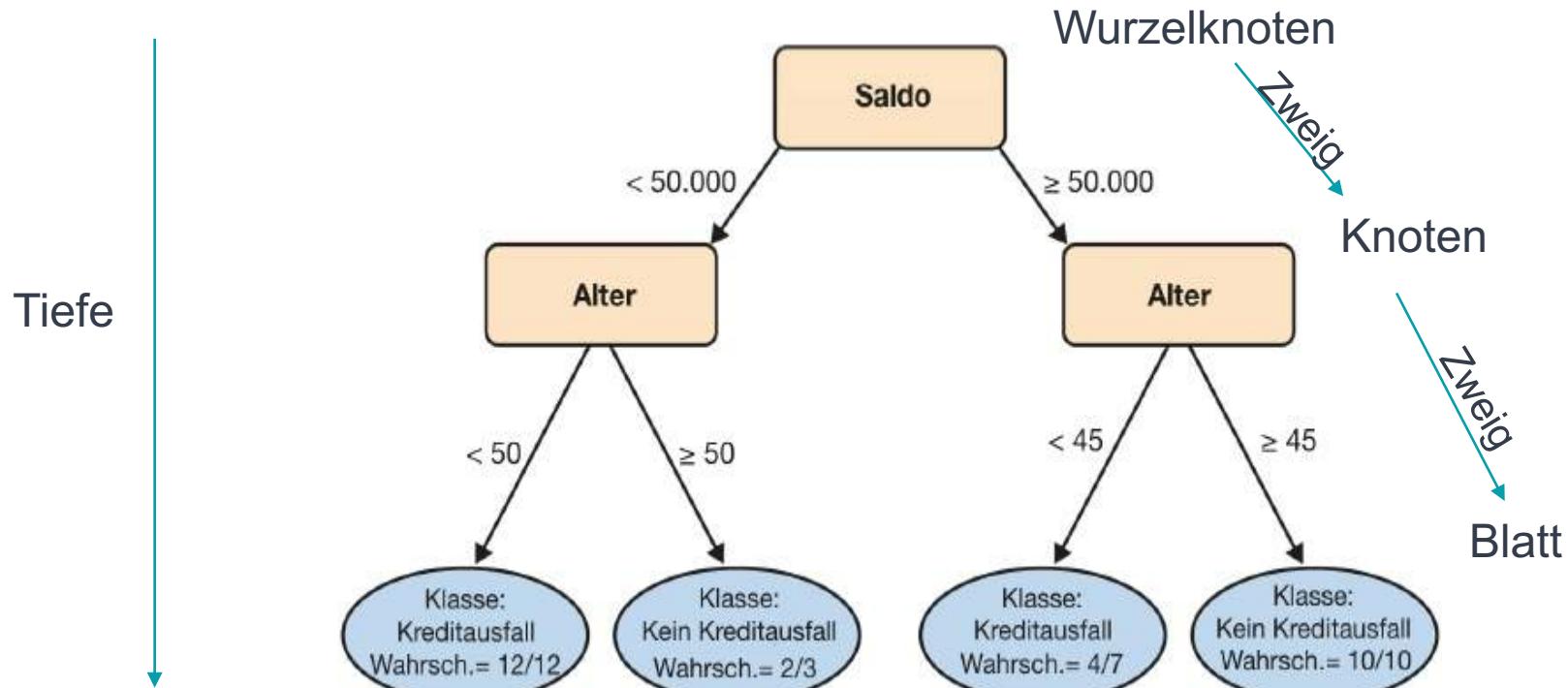


Use Case Korrelation (+PCA)



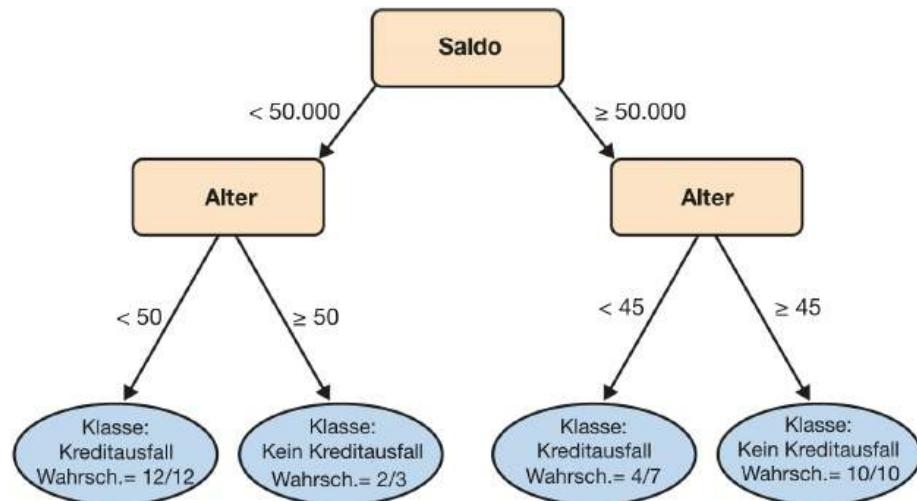
Klassifikation - Entscheidungsbaum

- Entscheidungsbaum zur Beurteilung der Kreditwürdigkeit von Kunden anhand von Merkmal Saldo und Alter

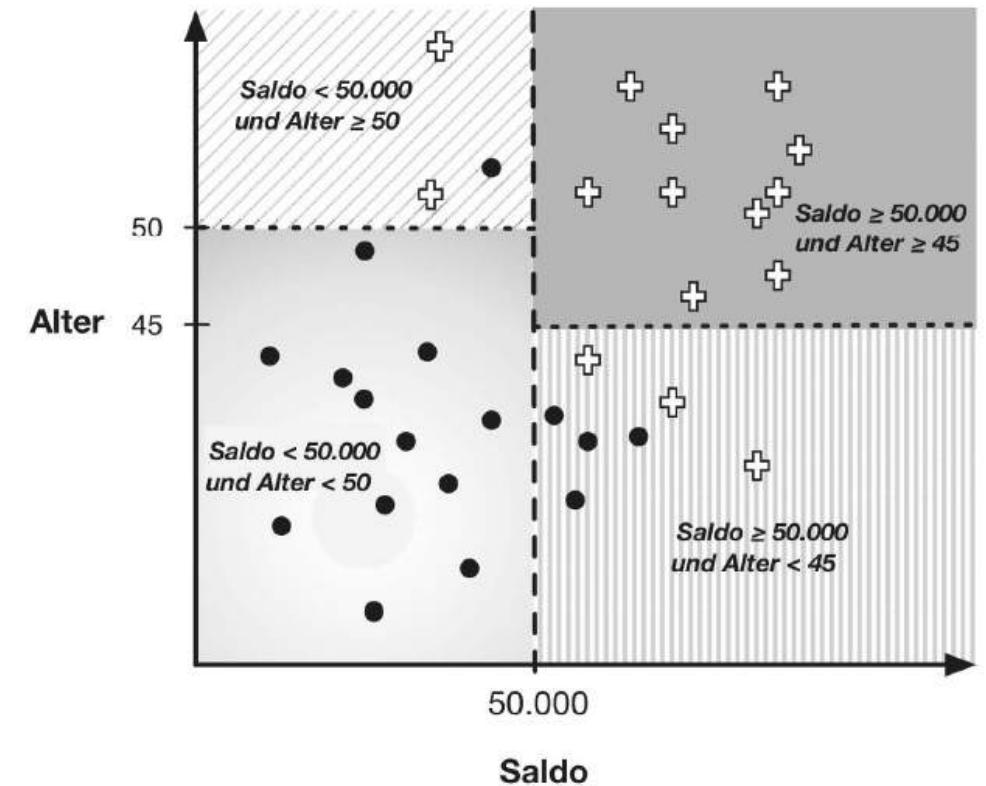


Klassifikation - Entscheidungsbaum

- Entscheidungsgrenzen sind keine Geraden oder Ebenen sondern Stufen/Bereiche



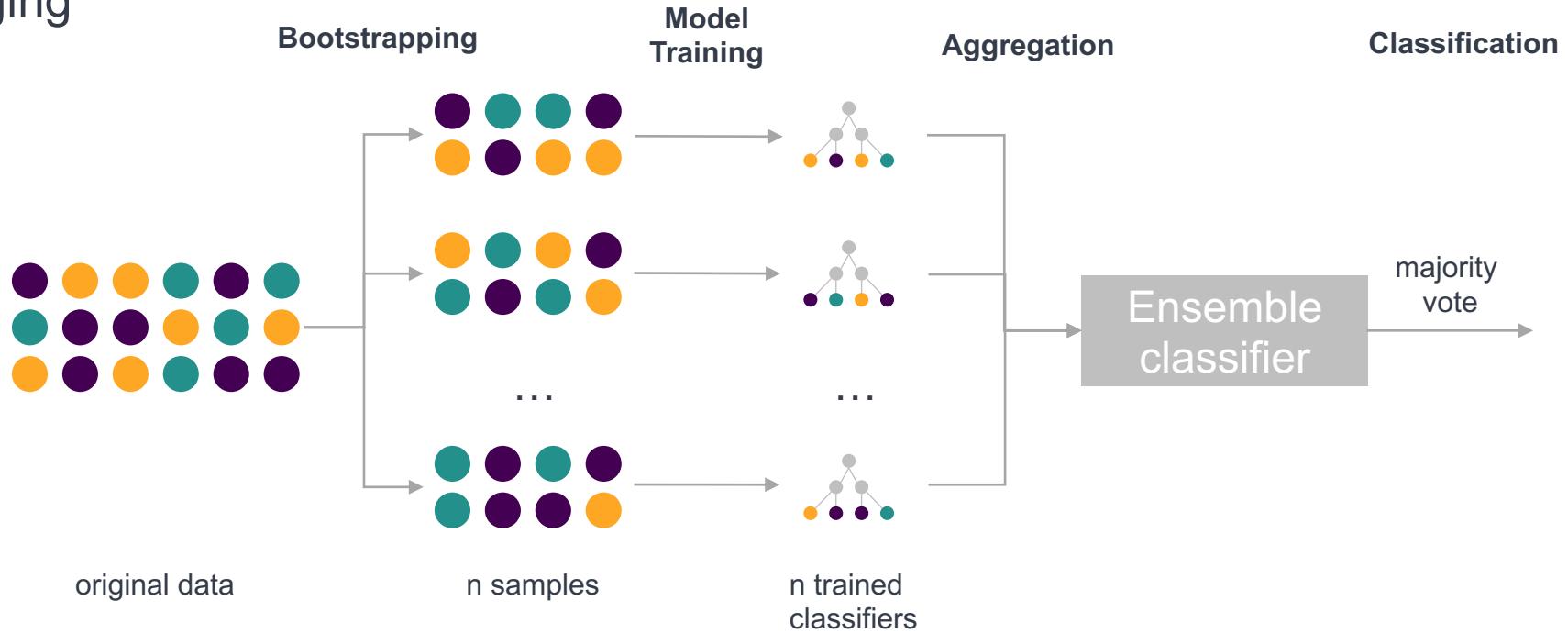
$$p(\text{Kreditausfall}) = \frac{\text{samples Kreditausfall}}{\text{samples Gesamt in Bereich}}$$



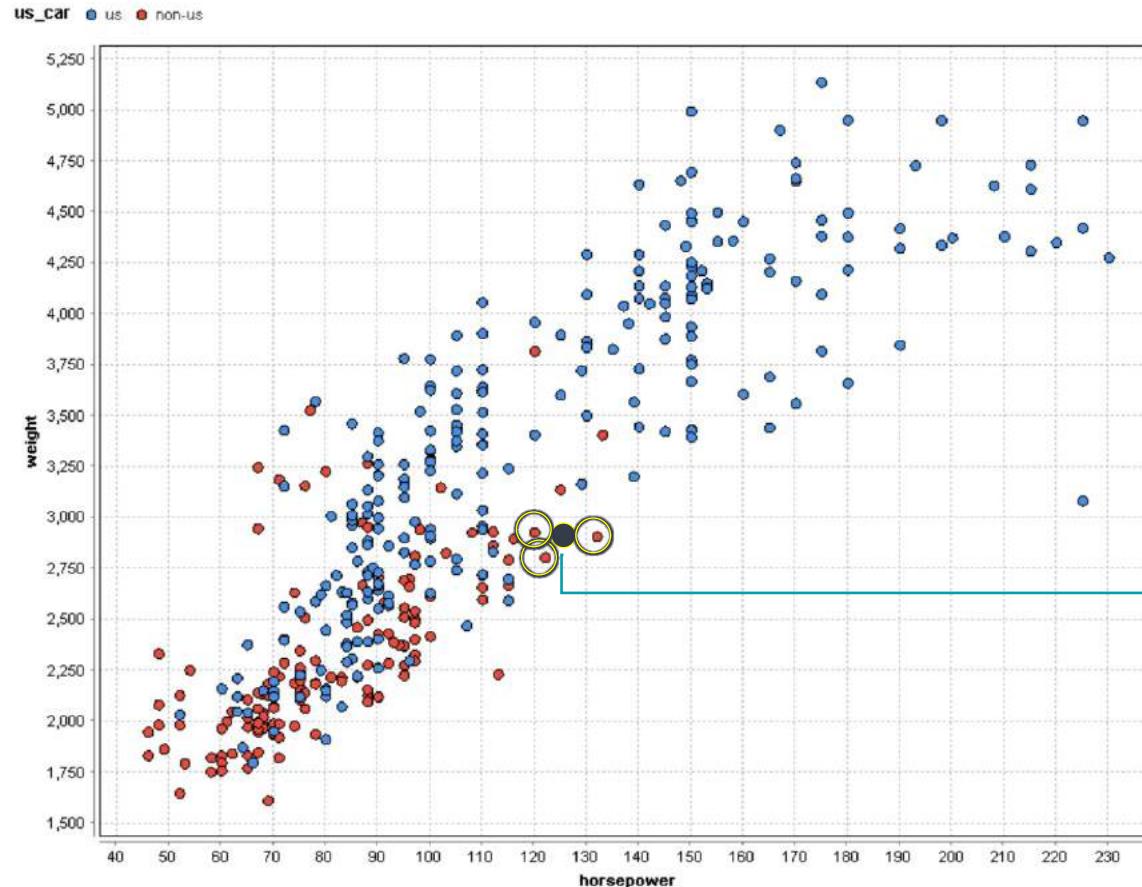
Ensemble Methoden - Bagging

- Um der Gefahr des Overfittings vorzubeugen können sogenannte Ensemble Methoden verwendet werden

Beispiel: Bagging



Klassifikation - K-Nearest-Neighbor



Für $k=3$ würde die
Vorhersage "non-us"
lauten



Klassifikation - K-Nearest-Neighbor

Distanzmaße:

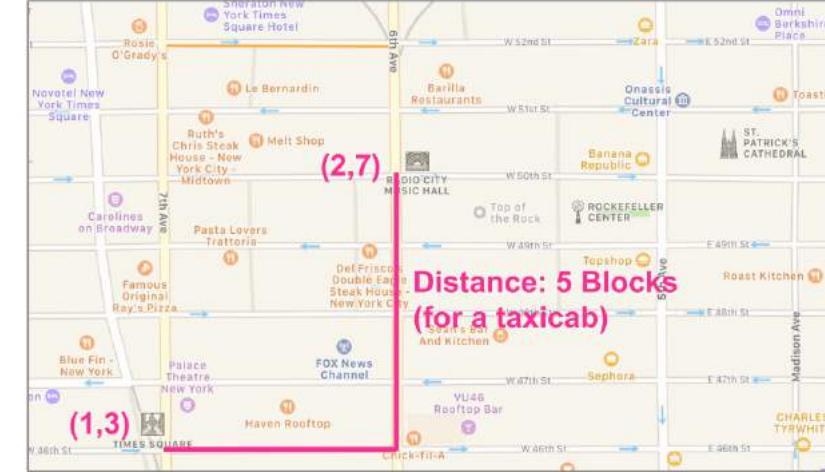
Euclidian distance

$$d(x, x') = \sqrt{\sum_{i=1}^m (x_i - x'_i)^2}$$



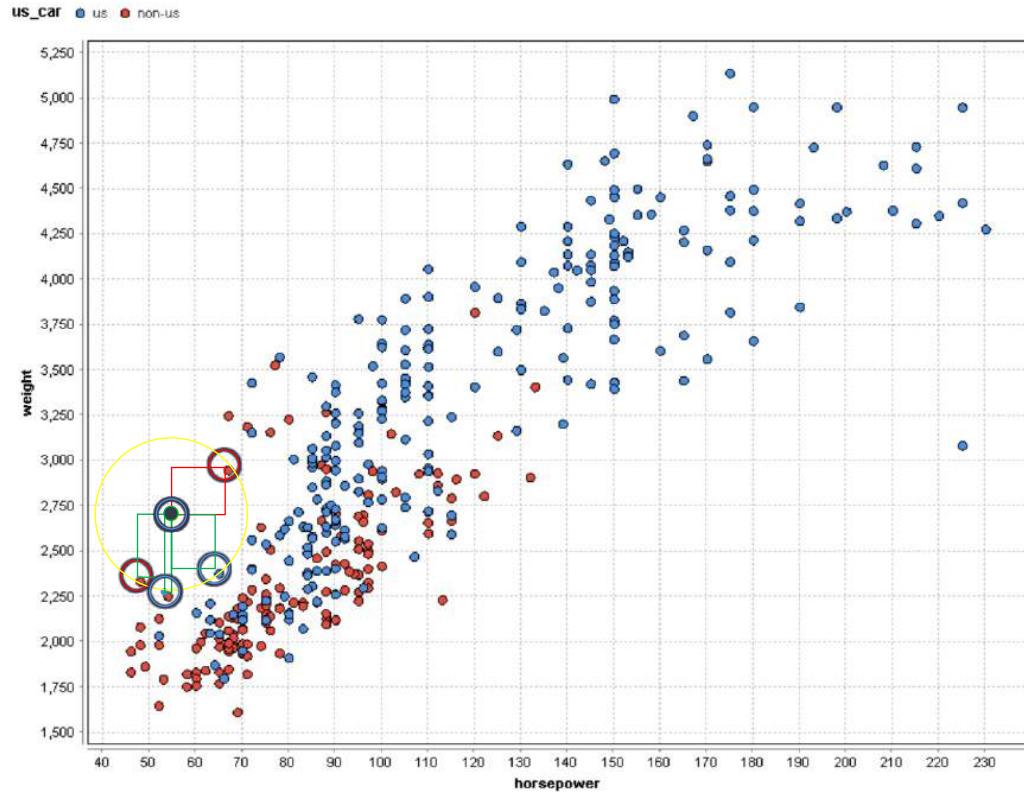
Manhattan distance

$$d(x, x') = \sum_{i=1}^m |x_i - x'_i|$$



Klassifikation - K-Nearest-Neighbor

- Beispiel, k=3
 - Euclidian distance
 - Manhattan distance
- Die Wahl des richtigen Distanzmaßes hängt von den Daten und den getroffenen Annahmen ab
 - Euclidian distance zieht eine „Homogenität“ des zu klassifizierenden Datenpunkts über alle Features vor (hier horsepower und weight)
 - Manhattan distance wird mit steigender Dimensionierung der Daten eher genutzt



Use Case Klassifikation



Clustering - KMeans

Algorithmus:

1. Auswahl aus Objekten k Zentroide als anfängliche Clusterzentren.
2. Alle Objekte dem nächsten Zentroiden zuweisen.
3. Neuberechnung des Zentroiden mit den aus Schritt 2 zugewiesenen Objekten.
4. Wiederholung von 2+3, bis sich die Zuordnung nicht mehr ändert (entweder Schwellenwert oder maximale Iterationen werden vorgegeben).

Aber was ist unser Ähnlichkeitsmaß ?



Clustering - KMeans

Euklidische Distanz:

- Datenobjekte die sich ähnlich sind nahe beieinander platziert

Formel: $d(x, y)^2 = \sum_{j=1}^m (x_j - y_j)^2 = \|x - y\|_2^2$

Durch euklidische Distanz kann k-Means Algorithmus als Optimierungsaufgabe formuliert werden

- Summe der quadrierten Abweichungen innerhalb eines Clusters soll minimiert werden

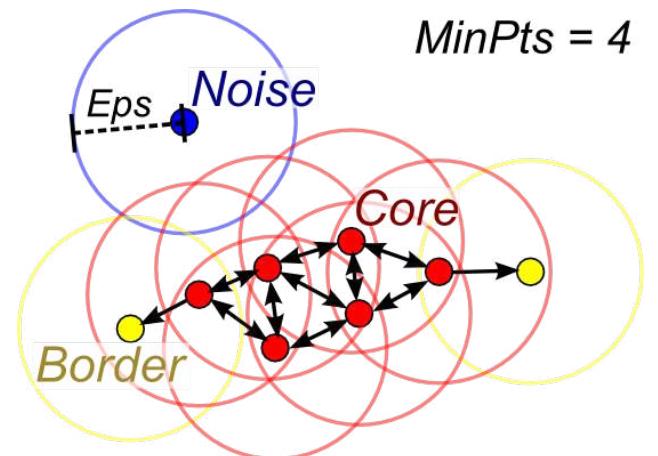
Formel: $SSE = \sum_{i=1}^n \sum_{j=1}^k w^{(i,j)} \|x^{(i)} - \mu^{(j)}\|_2^2$



Clustering - DBSCAN

(Density-Based Spatial Clustering of Applications with Noise)

1. Berechnen der Density für jeden Datenpunkt abhängig von epsilon → Density > MinPts = High-Density Area
2. Wenn Datenpunkt = High-Density Area → Core Point
Falls Datenpunkt ≠ High-Density Area, aber Core Point innerhalb Radius (epsilon) → Border Point
3. Falls Datenpunkt ≠ High-Density Area und kein Core Point innerhalb Radius (epsilon) → Noise



Live Demo

Kmeans & DBSCAN



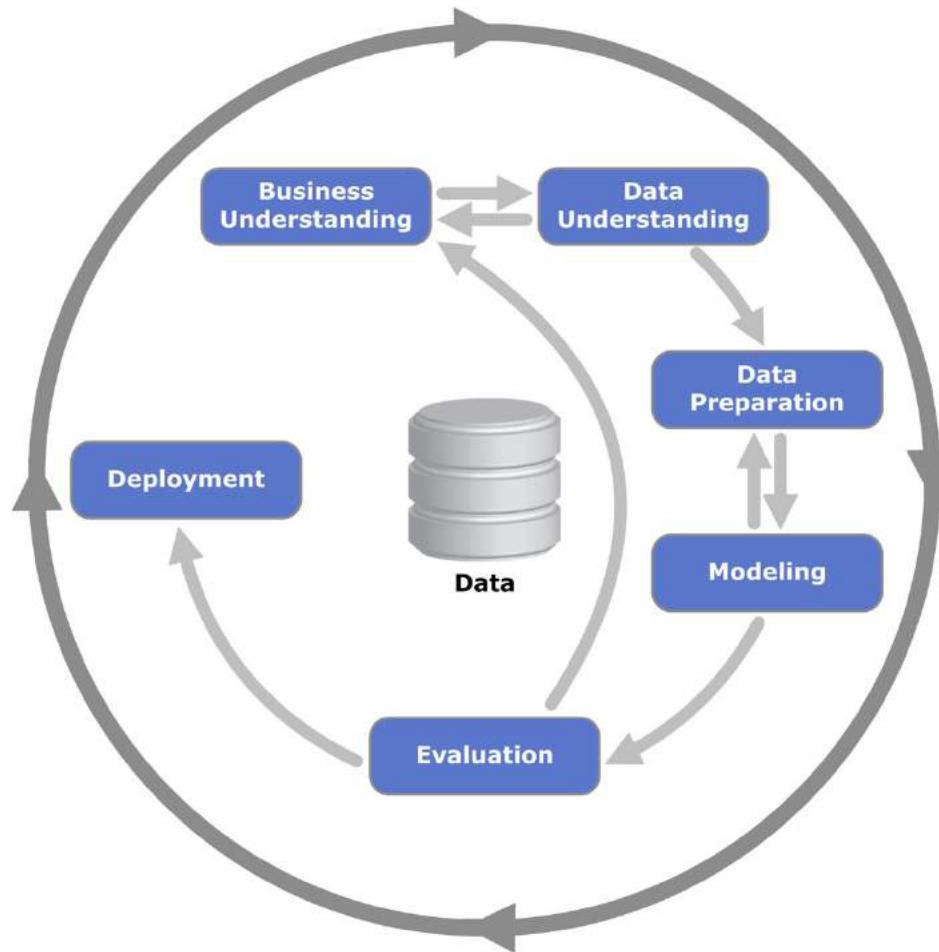
Use Case Clustering



Machine Learning Workflow



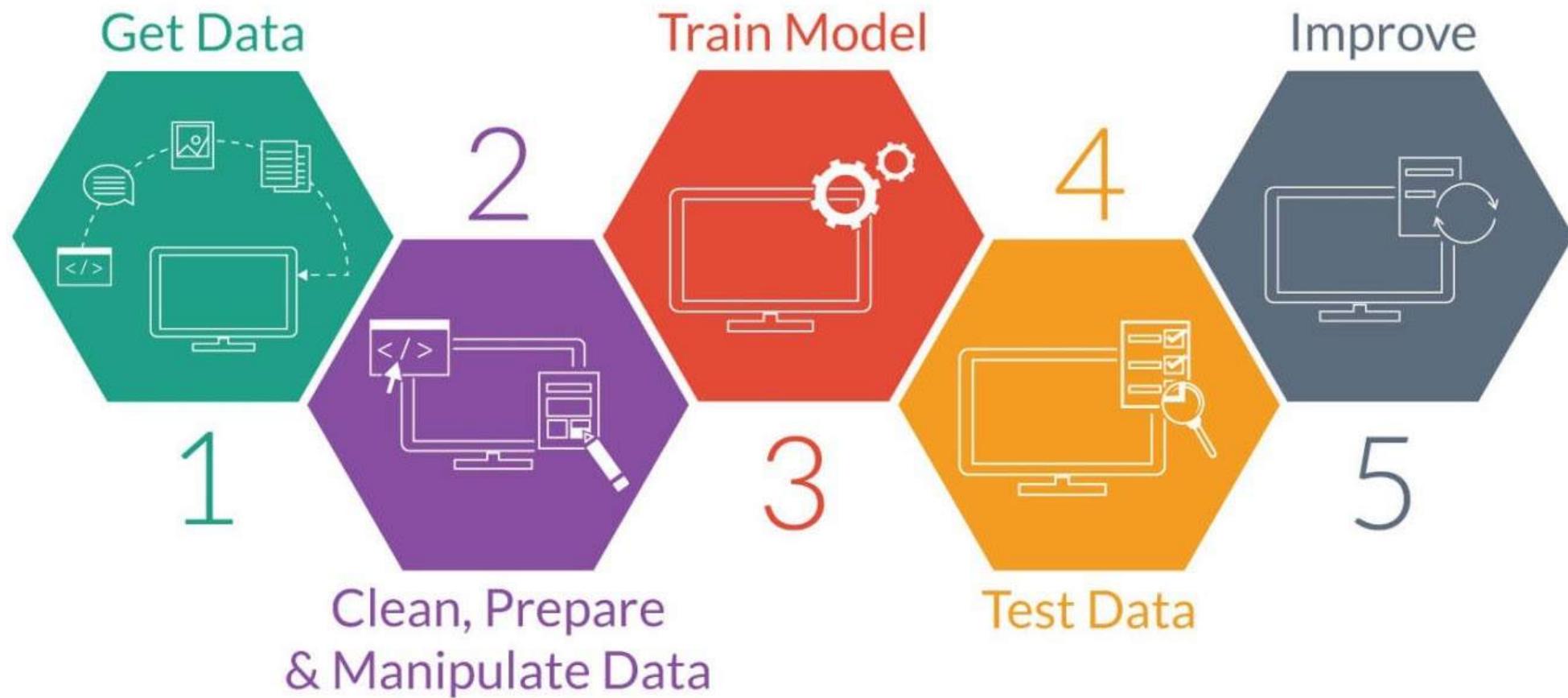
The bigger picture – CRISP-DM



https://en.wikipedia.org/wiki/Cross-industry_standard_process_for_data_mining

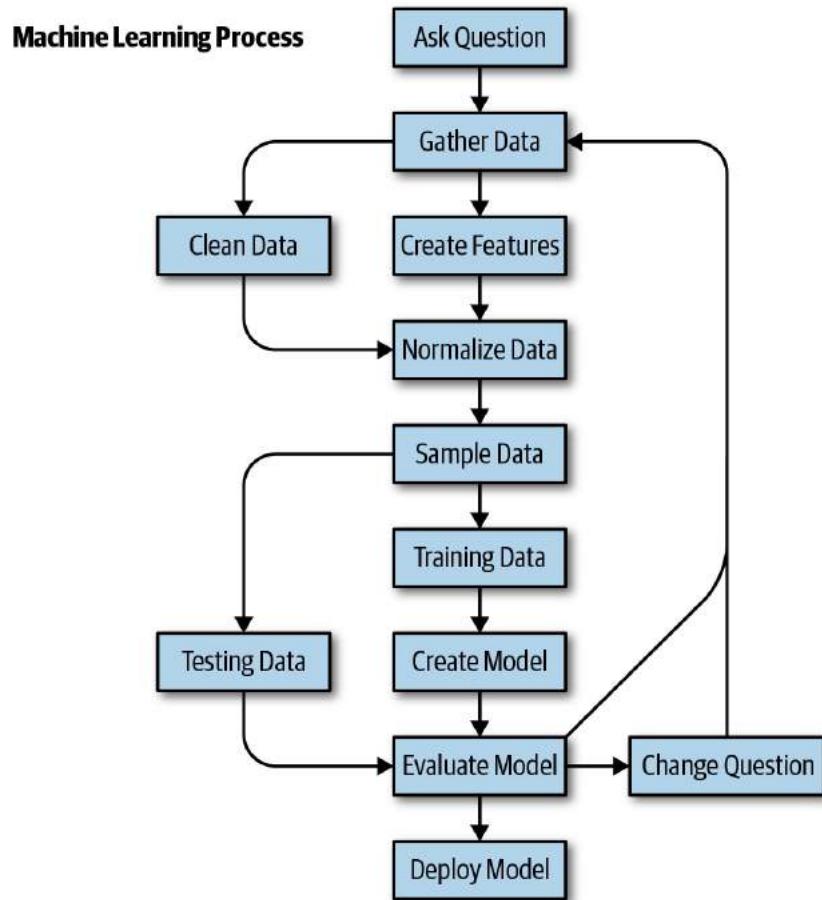


The Machine Learning workflow - simplified



A more detailed workflow

Machine Learning Pocket Reference by Matt Harrison



THE END





CONNECTING WORLDS

Kontakt



Data Scientist

philipp.bongartz@EXXETA.com



Cheat Sheets

Python For Data Science Cheat Sheet

Python Basics

Learn More Python for Data Science interactively at www.datacamp.com

Variables and Data Types

Variable Assignment

```
>>> x=5  
>>> x  
5
```

Calculations With Variables

>>> x+2	Sum of two variables
>>> x-2	Subtraction of two variables
>>> x*2	Multiplication of two variables
>>> x**2	Exponentiation of a variable
>>> x%2	Remainder of a variable
>>> x/float(2)	Division of a variable

Types and Type Conversion

str()	'5', '+3.45', 'True'	Variables to strings
int()	5, 3, 1	Variables to integers
float()	5.0, 1.0	Variables to floats
bool()	True, True, True	Variables to booleans

Asking For Help

```
>>> help(str)
```

Strings

```
>>> my_string = 'thisStringIsAwesome'  
>>> my_string  
'thisStringIsAwesome'
```

String Operations

```
>>> my_string * 2  
'thisStringIsAwesomethisStringIsAwesome'  
>>> my_string + 'Init'  
'thisStringIsAwesomeInit'  
>>> 'm' in my_string  
True
```

Lists

Also see NumPy Arrays

```
>>> a = 'is'  
>>> b = 'nice'  
>>> my_list = ['my', 'list', a, b]  
>>> my_list2 = [[4,5,6,7], [3,4,5,6]]
```

Selecting List Elements

Index starts at 0

Subset
>>> my_list[1]
>>> my_list[-3]
Slice
>>> my_list[1:3]
>>> my_list[1:
>>> my_list[:3]
>>> my_list[:]
Subset Lists of Lists
>>> my_list2[1][0]
>>> my_list2[1][:2]

List Operations

```
>>> my_list + my_list  
'my', 'list', 'is', 'now', 'my', 'list', 'is', 'nice'  
>>> my_list * 2  
'my', 'list', 'is', 'now', 'my', 'list', 'is', 'nice'  
>>> my_list2 > 4
```

List Methods

```
>>> my_list.index('a')  
>>> my_list.count('a')  
>>> my_list.append('!')  
>>> my_list.remove('!')  
>>> del(my_list[0:1])  
>>> my_list.reverse()  
>>> my_list.extend('!!')  
>>> my_list.pop(-1)  
>>> my_list.insert(0,'!')  
>>> my_list.sort()
```

String Operations

Index starts at 0

```
>>> my_string[3]  
>>> my_string[4:9]
```

String Methods

```
>>> my_string.upper()  
>>> my_string.lower()  
>>> my_string.count('w')  
>>> my_string.replace('o', '1')  
>>> my_string.strip()
```

Libraries

Import libraries
>>> import numpy
>>> import numpy as np
Selective import
>>> from math import pi

pandas Data analysis
NumPy
matplotlib

Install Python

ANACONDA
spyder
jupyter

Leading open data science platform
powered by Python
Free IDE that is included
with Anaconda
Create and share
documents with live code,
visualizations, text, ...

Numpy Arrays

Also see Lists
>>> my_list = [1, 2, 3, 4]
>>> my_array = np.array(my_list)
>>> my_2darray = np.array([[1,2,3], [4,5,6]])

Selecting Numpy Array Elements

Index starts at 0

Subset
>>> my_array[1]
2
Slice
>>> my_array[0:2]
array([1, 2])
Subset 2D Numpy arrays
>>> my_2darray[:,0]
array([1, 4])

Numpy Array Operations

```
>>> my_array > 3  
array([1, 0, 1, 0, 1, 0, 1, 0], dtype=bool)  
>>> my_array * 2  
array([2, 4, 6, 8])  
>>> my_array + np.array([5, 6, 7, 8])  
array([6, 8, 10, 12])
```

Numpy Array Functions

Get the dimensions of the array
>>> my_array.shape
>>> np.append(other_array)
Append items to an array
>>> np.insert(my_array, 1, 5)
Insert items in an array
>>> np.delete(my_array, [1])
Delete items in an array
>>> np.mean(my_array)
Mean of the array
>>> np.median(my_array)
Median of the array
>>> my_array.corrcoef()
Correlation coefficient
>>> np.std(my_array)
Standard deviation

DataCamp
Learn Python for Data Science interactively

Numpy Cheat Sheet



Pandas Cheat Sheet

Data Wrangling with pandas Cheat Sheet

<http://pandas.pydata.org>

Syntax – Creating DataFrames

```
df = pd.DataFrame([{"a": 4, "b": 5, "c": 6}, {"b": 7, "a": 8, "c": 9}, {"c": 10, "b": 11, "a": 12}], index=[1, 2, 3])
```

Specify values for each column.

```
df = pd.DataFrame([[4, 7, 10], [5, 8, 11], [6, 9, 12]], index=[1, 2, 3], columns=['a', 'b', 'c'])
```

Specify values for each row.

```
df = pd.DataFrame([{"n": 1, "v": 4, "b": 5, "t": 10}, {"n": 2, "v": 5, "b": 8, "t": 11}, {"n": 3, "v": 6, "b": 9, "t": 12}])
```

Create DataFrame with a MultiIndex

```
df = pd.DataFrame([(4, 5, 6), (7, 8, 9), (10, 11, 12)], index=pd.MultiIndex.from_tuples([('d', 1), ('d', 2), ('e', 2)], names=['n', 'v']))
```

Method Chaining

Most pandas methods return a DataFrame so that another pandas method can be applied to the result. This improves readability of code.

```
df = (pd.melt(df).rename(columns={'variable': 'var', 'value': 'val'})) .query('val >= 200')
```

Tidy Data – A foundation for wrangling in pandas

In a tidy data set:

&

Tidy data complements pandas's **vectorized operations**. pandas will automatically preserve observations as you manipulate variables. No other format works as intuitively with pandas.

M * A

Reshaping Data – Change the layout of a data set

pd.melt(df) Gather columns into rows.

pd.pivot(columns='var', values='val') Spread rows into columns.

pd.concat([df1, df2]) Append rows of DataFrames

pd.concat([df1, df2], axis=1) Append columns of DataFrames

Subset Observations (Rows)

df[df.Length > 7] Extract rows that meet logical criteria.

df.sample(frac=0.5) Randomly select fraction of rows.

df.sample(n=10) Randomly select n rows.

df.iloc[10:20] Select rows by position.

df.nlargest(n, 'value') Select and order top n entries.

df.nsmallest(n, 'value') Select and order bottom n entries.

Subset Variables (Columns)

df[['width', 'length', 'species']] Select multiple columns with specific names.

df['width'] or df.width Select single column with specific name.

df.filter(regex='regex') Select columns whose name matches regular expression regex.

regex (Regular Expressions) Examples

'.'	Matches strings containing a period ''.
'Length\$'	Matches strings ending with word 'Length'.
'Sepal'	Matches strings beginning with the word 'Sepal'.
'^x[1-5]\$'	Matches strings beginning with 'x' and ending with 1,2,3,4,5.
'^(?!(Species)).*'}	Matches strings except the string 'Species'.

df.loc[:, 'x2':'x4'] Select all columns between x2 and x4 (inclusive).

df.loc[:, [1, 2, 5]] Select columns in positions 1, 2 and 5 (first column is 0).

df.loc[df['a'] > 10, ['a', 'c']] Select rows meeting logical condition, and only the specific columns.

<http://pandas.pydata.org> This cheat sheet inspired by RStudio Data Wrangling Cheatsheet (<http://www.rstudio.com/ide/docs/guides/2015-07>Data-Wrangling-cheatsheet.pdf>) Written by Iván Uribe. <https://creativecommons.org/licenses/by-sa/4.0/>

Matplotlib Cheat Sheet

Python For Data Science Cheat Sheet

Matplotlib

Learn Python Interactively at www.DataCamp.com

Matplotlib

Matplotlib is a Python 2D plotting library which produces publication-quality figures in a variety of hardcopy formats and interactive environments across platforms.

1 Prepare The Data

Also see [Lists & NumPy](#)

1D Data

```
>>> import numpy as np  
>>> x = np.linspace(0, 10, 100)  
>>> y = np.sin(x)  
>>> V = 1 + x**2
```

from matplotlib.cbook import get_sample_data
img = np.loadtxt(get_sample_data('axes_grid/bivariate_normal.npy'))

2 Create Plot

```
>>> import matplotlib.pyplot as plt  
  
Figure  
>>> fig = plt.figure()  
>>> fig2 = plt.figure(figsize=plt.figaspect(2.0))  
  
Axes  
All plotting is done with respect to an Axes. In most cases, a subplot will fit your needs. A subplot is an axes on a grid system.  
>>> fig.add_axes()  
>>> ax1 = fig.add_subplot(221) # row-col-num  
>>> ax2 = fig.add_subplot(212)  
>>> fig3 = plt.subplots(rows=2, ncols=2)  
>>> fig4, axes2 = plt.subplots(ncols=3)
```

3 Plotting Routines

1D Data

```
>>> lines = ax.plot(x,y)  
>>> ax.scatter(x,y)  
>>> ax.patches[0].bar([1,2,3],[3,4,5])  
>>> axes[1,0].barh([0.5,1,2,5],[0.1,2])  
>>> axes[1,1].axhline(0.45)  
>>> axes[0,1].axvline(0.65)  
>>> ax.hill_between(y,y,color='blue')  
>>> ax.hill_between(y,y,color='yellow')
```

2D Data or Images

```
>>> fig, ax = plt.subplots()  
>>> im = ax.imshow(img, 'gist_earth',  
    interpolation='nearest',  
    vmin=-2,  
    vmax=2)
```

Colormapped or RGB arrays

Plot Anatomy & Workflow

Plot Anatomy

Workflow

The basic steps to creating plots with matplotlib are:

- 1 Prepare data
- 2 Create plot
- 3 Plot
- 4 Customize plot
- 5 Save plot
- 6 Show plot

```
>>> import matplotlib.pyplot as plt  
>>> x = [10,20,25,30] Step 1  
>>> fig = plt.figure() Step 2  
>>> ax = fig.add_subplot(111) Step 3  
>>> ax.plot(x, y, color='lightblue', linewidth=3) Step 4  
>>> ax.scatter([2,4,6],  
    [5,15,25],  
    color='darkgreen',  
    marker='*')  
>>> ax.set_xlim(1, 6.5)  
>>> plt.savefig('foo.png') Step 5  
>>> plt.show() Step 6
```

4 Customize Plot

Colors, Color Bars & Color Maps

```
>>> plt.plot(x, x, x**2, x, x**3)  
>>> ax.plot(x, y, alpha = 0.4)  
>>> fig.colorbar(im, orientation='horizontal')  
>>> im = ax.imshow(img,  
    cmap='seismic')
```

Markers

```
>>> fig, ax = plt.subplots()  
>>> ax.scatter(x,y,marker=".")  
>>> ax.plot(x,y,marker="o")
```

Linestyles

```
>>> plt.plot(x,y,linewidth=4.0)  
>>> plt.plot(x,y,ls="solid")  
>>> plt.plot(x,y,ls="--")  
>>> plt.plot(x,y,'-.',x**2,y**2,'-')  
>>> plt.setp(lines,color='r',linewidth=4.0)
```

Text & Annotations

```
>>> ax.text(1,  
    2.1,  
    'Example Graph',  
    style='italic',  
    transform='rotate(15deg)',  
    xytext=(8, 0),  
    xycoords='data',  
    textcoords='data',  
    arrowprops=dict(arrowstyle="->",  
    connectionstyle="arc3"))
```

Mathtext

```
>>> plt.title(r'$\sigma_i=15$', fontsize=20)
```

Limits, Legends & Layouts

Limits & Autoscaling
>>> ax.autoscale(x0=0.0, y0=0.1)
>>> ax.axis('equal')
>>> ax.set_xlim(-3,10.5), ylim=[-1.5,1.5])
>>> ax.set_xlim(0,10.5)

Legends
>>> ax.legend(loc='best')

Ticks
>>> ax.xaxis.set_ticks(range(1,5),
 ticklabels=[3,100,-12,"foo"])
>>> ax.tick_params(axis='y',
 direction='inout',
 length=10)

Subplot Spacing

```
>>> fig3.subplots_adjust(wspace=0.5,  
    hspace=0.3,  
    left=0.125,  
    right=0.9,  
    top=0.9,  
    bottom=0.1)
```

Axis Spines

```
>>> fig.tight_layout()  
>>> ax.spines['top'].set_visible(False)  
>>> ax.spines['bottom'].set_position((0,10))
```

5 Save Plot

Save figures
>>> plt.savefig('foo.png')
Save transparent figures
>>> plt.savefig('foo.png', transparent=True)

6 Show Plot

```
>>> plt.show()
```

Close & Clear

>>> plt.clf()
>>> plt.cla()
>>> plt.close()

DataCamp
Learn Python for Data Science Interactively!

Sscikit-Learn Cheat Sheet

Python For Data Science Cheat Sheet

Scikit-Learn

Learn Python for data science interactively at www.DataCamp.com

Scikit-learn

Scikit-learn is an open source Python library that implements a range of machine learning, preprocessing, cross-validation and visualization algorithms using a unified interface.

A Basic Example

```
>>> from sklearn import neighbors, datasets, preprocessing
>>> from sklearn.model_selection import train_test_split
>>> from sklearn.metrics import accuracy_score
>>> iris = datasets.load_iris()
>>> X, y = iris.data[:, :2], iris.target
>>> X_train, X_test, y_train, y_test = train_test_split(X, y, random_state=42)
>>> knn = neighbors.KNeighborsClassifier(n_neighbors=5)
>>> knn.fit(X_train, y_train)
>>> X_train = scaler.transform(X_train)
>>> X_test = scaler.transform(X_test)
>>> knn = neighbors.KNeighborsClassifier(n_neighbors=5)
>>> knn.fit(X_train, y_train)
>>> y_pred = knn.predict(X_test)
>>> accuracy_score(y_test, y_pred)
```

Loading The Data

Also see NumPy & Pandas

Your data needs to be numeric and stored as NumPy arrays or SciPy sparse matrices. Other types that are convertible to numeric arrays, such as Pandas DataFrames, are also acceptable.

```
>>> import numpy as np
>>> X = np.random.random((10, 5))
>>> X = np.array(['M', 'M', 'P', 'P', 'M', 'F', 'M', 'M', 'F', 'F'])
>>> X[X < 0.7] =
```

Training And Test Data

```
>>> from sklearn.model_selection import train_test_split
>>> X_train, X_test, y_train, y_test = train_test_split(X,
...                                                    y,
...                                                    random_state=0)
```

Preprocessing The Data

Standardization

```
>>> from sklearn.preprocessing import StandardScaler
>>> scaler = StandardScaler().fit(X_train)
>>> standardized_X_train = scaler.transform(X_train)
>>> standardized_X_test = scaler.transform(X_test)
```

Normalization

```
>>> from sklearn.preprocessing import Normalizer
>>> scaler = Normalizer().fit(X_train)
>>> normalized_X_train = scaler.transform(X_train)
>>> normalized_X_test = scaler.transform(X_test)
```

Binarization

```
>>> from sklearn.preprocessing import Binarizer
>>> binarizer = Binarizer(threshold=0.0).fit(X)
>>> binary_X = binarizer.transform(X)
```

Create Your Model

Supervised Learning Estimators

Linear Regression

```
>>> from sklearn.linear_model import LinearRegression
>>> lr = LinearRegression(normalize=True)
```

Support Vector Machines (SVM)

```
>>> from sklearn.svm import SVC
>>> svc = SVC(kernel='linear')
```

Naive Bayes

```
>>> from sklearn.naive_bayes import GaussianNB
>>> gnb = GaussianNB()
```

KNN

```
>>> from sklearn import neighbors
>>> knn = neighbors.KNeighborsClassifier(n_neighbors=5)
```

Unsupervised Learning Estimators

Principal Component Analysis (PCA)

```
>>> from sklearn.decomposition import PCA
>>> pca = PCA(n_components=0.95)
```

KMeans

```
>>> from sklearn.cluster import KMeans
>>> kmeans = KMeans(n_clusters=3, random_state=0)
```

Model Fitting

Supervised learning

```
>>> l = LinearRegression()
>>> l.fit(X_train, y_train)
```

Unsupervised Learning

```
>>> km = KMeans()
>>> km.fit(X_train)
```

Prediction

Supervised Estimators

```
>>> y_pred = svc.predict(np.random.random(2,5))
>>> y_pred = lr.predict(X_test)
>>> y_pred = knn.predict_proba(X_test)
```

Unsupervised Estimators

```
>>> y_pred = kmeans.predict(X_test)
```

Encoding Categorical Features

```
>>> from sklearn.preprocessing import LabelEncoder
>>> enc = LabelEncoder()
>>> y = enc.fit_transform(y)
```

Imputing Missing Values

```
>>> from sklearn.preprocessing import Imputer
>>> imp = Imputer(missing_values=0, strategy='mean', axis=0)
>>> imp.fit_transform(X_train)
```

Generating Polynomial Features

```
>>> from sklearn.preprocessing import PolynomialFeatures
>>> poly = PolynomialFeatures(5)
>>> poly.fit_transform(X)
```

Evaluate Your Model's Performance

Classification Metrics

Accuracy Score

```
>>> knn.score(X_test, y_test)
```

Classification Report

```
>>> from sklearn.metrics import classification_report
>>> print(classification_report(y_test, y_pred))
```

Confusion Matrix

```
>>> from sklearn.metrics import confusion_matrix
>>> print(confusion_matrix(y_test, y_pred))
```

Regression Metrics

Mean Absolute Error

```
>>> from sklearn.metrics import mean_absolute_error
>>> y_true = [3, -0.5, 2]
```

Mean Squared Error

```
>>> from sklearn.metrics import mean_squared_error
>>> y_pred
```

R² Score

```
>>> from sklearn.metrics import r2_score
>>> r2_score(y_true, y_pred)
```

Clustering Metrics

Adjusted Rand Index

```
>>> from sklearn.metrics import adjusted_rand_score
>>> adjusted_rand_score(y_true, y_pred)
```

Homogeneity

```
>>> from sklearn.metrics import homogeneity_score
>>> homogeneity_score(y_true, y_pred)
```

V-measure

```
>>> from sklearn.metrics import v_measure_score
>>> metrics.v_measure_score(y_true, y_pred)
```

Cross-Validation

```
>>> from sklearn.cross_validation import cross_val_score
>>> print(cross_val_score(knn, X_train, y_train, cv=4))
>>> print(cross_val_score(lr, X, y, cv=2))
```

Tune Your Model

Grid Search

```
>>> from sklearn.grid_search import GridSearchCV
>>> param = {'n_neighbors': np.arange(1, 3),
...           'weights': ['uniform', 'distance']}
>>> grid = GridSearchCV(estimator=knn,
...                      param_grid=param)
>>> grid.fit(X_train, y_train)
>>> print(grid.best_score_)
>>> print(grid.best_estimator_.n_neighbors)
```

Randomized Parameter Optimization

```
>>> from sklearn.grid_search import RandomizedSearchCV
>>> params = {'n_neighbors': range(1, 5),
...            'weights': ['uniform', 'distance']}
>>> research = RandomizedSearchCV(estimator=knn,
...                                 param_distributions=params,
...                                 cv=4,
...                                 n_iter=8,
...                                 random_state=5)
>>> research.fit(X_train, y_train)
>>> print(research.best_score_)
```

DataCamp
Learn Python for Data Science interactively



Scikit-learn Design

<https://arxiv.org/abs/1309.0238>

Scikit-Learn's API is remarkably well designed. These are the main design principles:

Consistency: All objects share a consistent and simple interface:

- **Estimators** Any object that can estimate some parameters based on a dataset is called an *estimator* (e.g., an imputer is an estimator). The estimation itself is performed by the `fit()` method, and it takes only a dataset as a parameter (or two for supervised learning algorithms; the second dataset contains the labels). Any other parameter needed to guide the estimation process is considered a hyperparameter (such as an imputer's strategy), and it must be set as an instance variable (generally via a constructor parameter).
- **Transformers** Some estimators (such as an imputer) can also transform a dataset; these are called *transformers*. Once again, the API is simple: the transformation is performed by the `transform()` method with the dataset to transform as a parameter. It returns the transformed dataset. This transformation generally relies on the learned parameters, as is the case for an imputer. All transformers also have a convenience method called `fit_transform()` that is equivalent to calling `fit()` and then `transform()` (but sometimes `fit_transform()` is optimized and runs much faster).
- **Predictors** Finally, some estimators, given a dataset, are capable of making predictions; they are called *predictors*. For example, the `LinearRegression` model in the previous chapter was a predictor: given a country's GDP per capita, it predicted life satisfaction. A predictor has a `predict()` method that takes a dataset of new instances and returns a dataset of corresponding predictions. It also has a `score()` method that measures the quality of the predictions, given a test set (and the corresponding labels, in the case of supervised learning algorithms).¹⁸

Inspection: All the estimator's hyperparameters are accessible directly via public instance variables (e.g., `imputer.strategy`), and all the estimator's learned parameters are accessible via public instance variables with an underscore suffix (e.g., `imputer.statistics_`).

Nonproliferation of classes Datasets are represented as NumPy arrays or SciPy sparse matrices, instead of homemade classes. Hyperparameters are just regular Python strings or numbers.

Composition Existing building blocks are reused as much as possible. For example, it is easy to create a Pipeline estimator from an arbitrary sequence of transformers followed by a final estimator, as we will see.

Sensible defaults Scikit-Learn provides reasonable default values for most parameters, making it easy to quickly create a baseline working system.



Where to go from here

1. Learn the language
2. Learn the tools
3. Work on projects!

Resources

- [Machine Learning Mastery](#)
- [Hands-On Machine Learning with Scikit-Learn and TensorFlow](#)
- [The Complete Python Course for Machine Learning Engineers](#)
- [Cheat Sheet of Machine Learning and Python \(and Math\) Cheat Sheets](#)



Webinar's Agenda

1. Machine Learning Basics

- Definition
- Relation to other fields
- Little terminology

2. The ML Workflow

- Get data
- Prepare data
- Model selection and training
- Model evaluation
- Model improvement

3. Resources



Machine Learning

What comes to your mind?



*“Every serious technology company now has an Artificial Intelligence team in place. These companies are **investing millions** into intelligent systems for situation assessment, prediction analysis, learning-based recognition systems, conversational interfaces, and recommendation engines.*

*Companies such as Google, Facebook, and Amazon aren’t just employing AI, but have made it a **central part of their core intellectual property.**” - Kristian J. Hammond*



Startups and AI

The **AI 100** is CB Insights' annual ranking of the 100 most promising AI startups in the world.

2020

Healthcare



Finance & Insurance



Transportation



Construction



Retail & Warehousing



Govt. & City Planning



Legal



Mining



Food & Agriculture



CROSS-INDUSTRY TECH

AI Processors



AI Model Development



DevOps & Model Monitoring



NLP, NLG, & Computer Vision



Cybersecurity



BI & Ops Intel



Sales & CRM



Other R&D

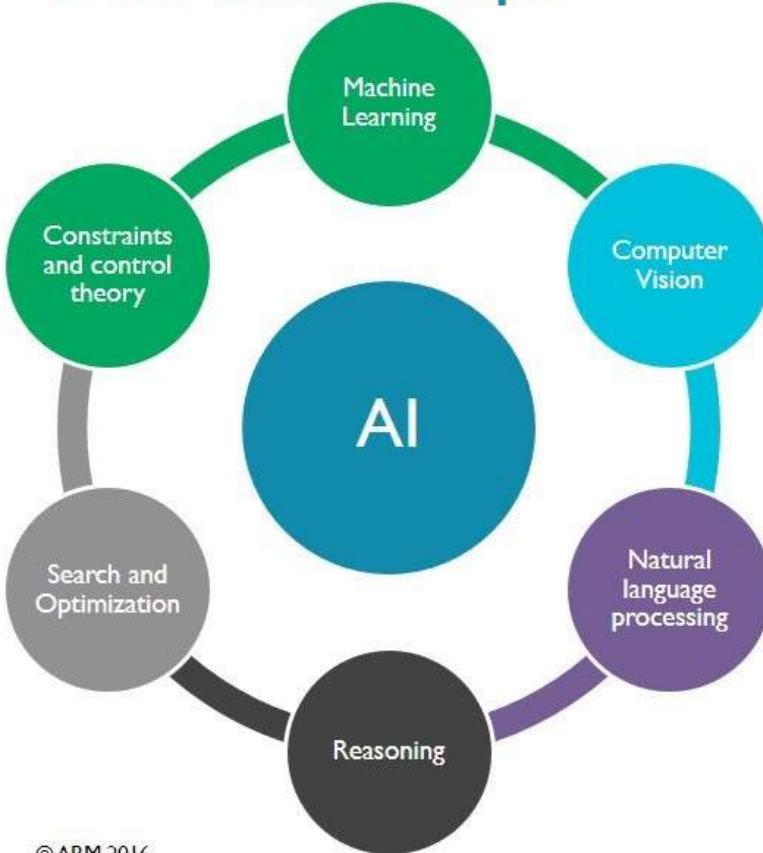


<https://www.cbinsights.com/research/artificial-intelligence-top-startups/>



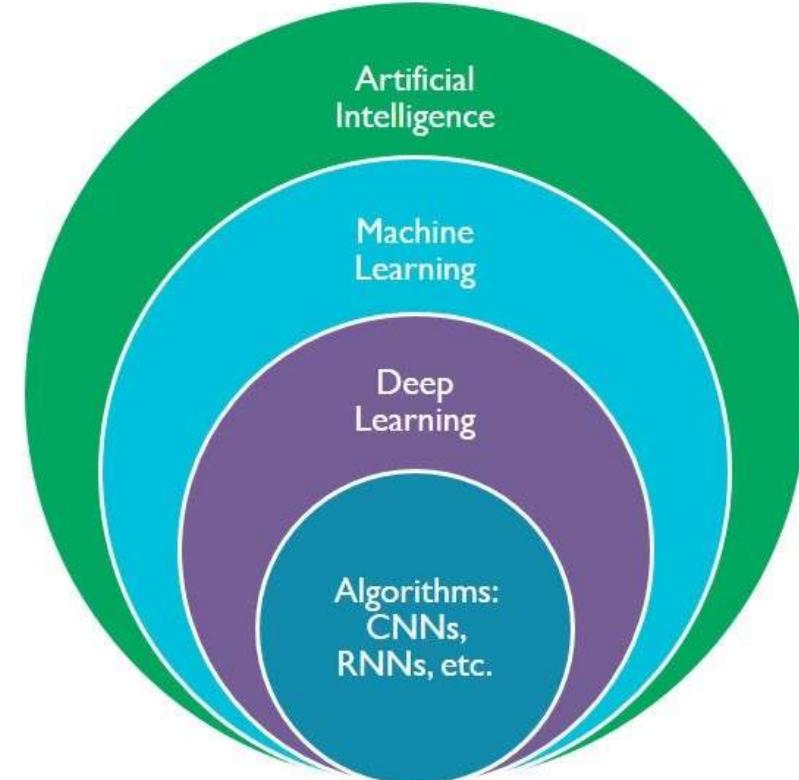
Artificial Intelligence and Machine Learning

The AI landscape



6

©ARM 2016



ARM



What does “learning” mean?



https://youtu.be/f_uwKZIAeM0



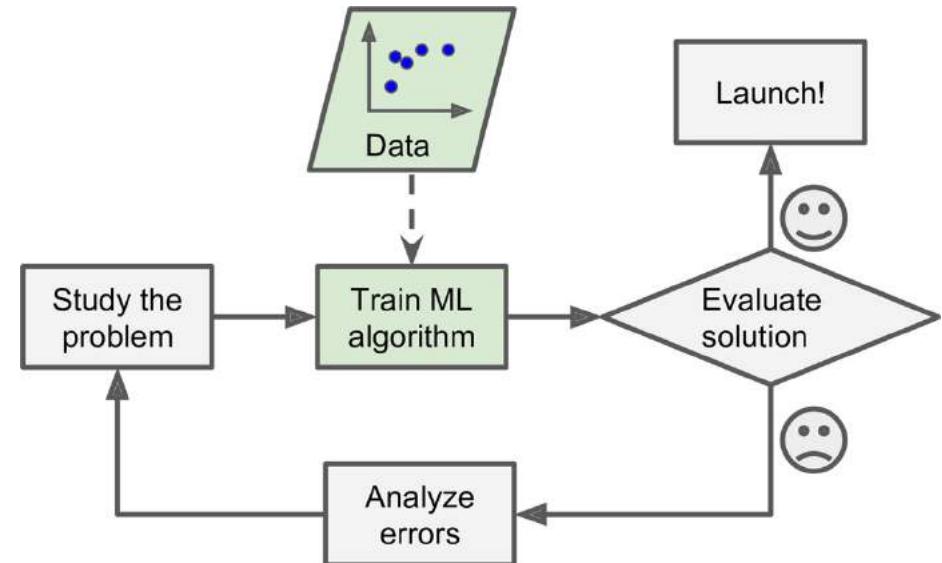
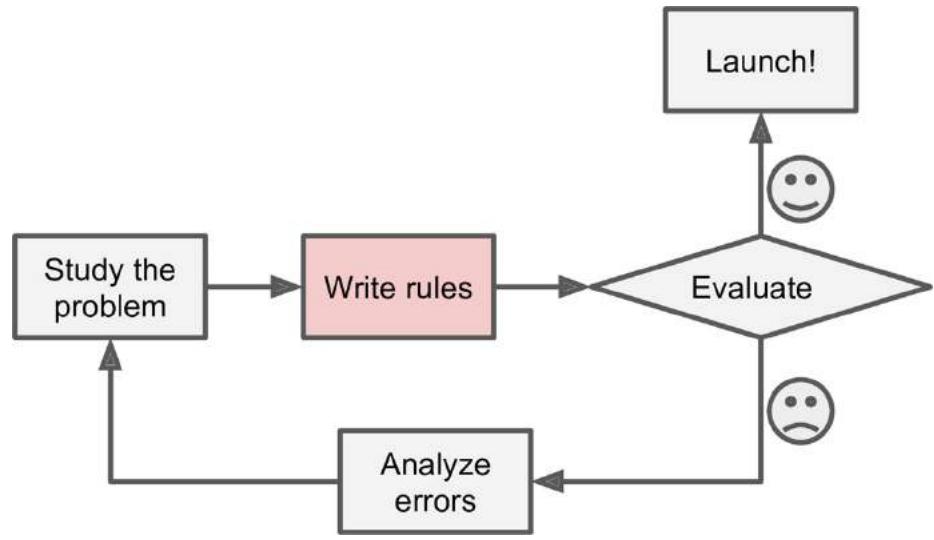
What does “learning” mean?

- 1.“Machine learning is the science of getting computers to act without being explicitly programmed.” – Andrew Ng, Stanford University (based on Arthur Samuel, 1959)
- 2.“Machine learning algorithms can figure out how to perform important tasks by generalizing from examples.” – University of Washington
- 3.“A computer program is said to learn from experience E with respect to some class of tasks T and performance measure P, if its performance at tasks in T, as measured by P, improves with experience E.” – Tom Mitchell



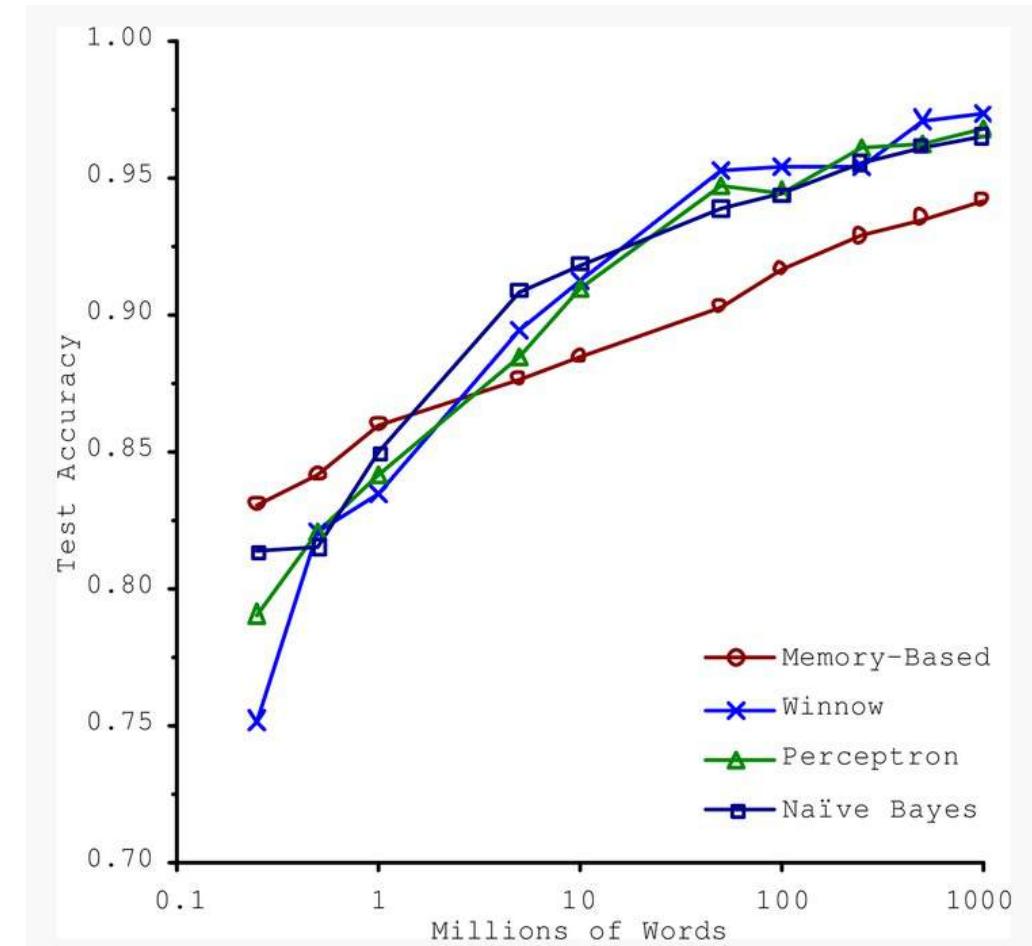
Why Use Machine Learning?

How would you implement a spam filter?



Main challenges of Machine Learning

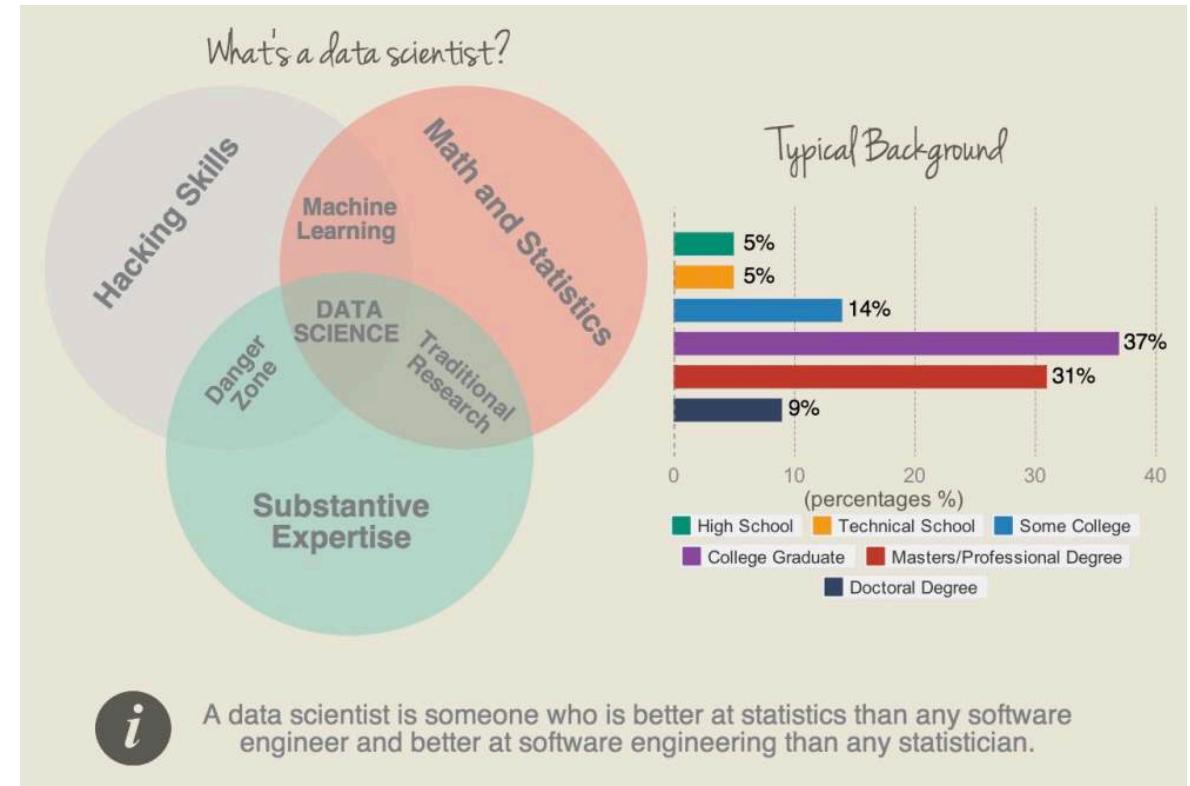
- Insufficient Quantity of Training Data
- Nonrepresentative Training Data
- Poor-Quality Data
- Irrelevant Features
- Overfitting the Training Data
- Underfitting the Training Data



Machine Learning Workflow



The bigger picture – skillset of a data scientist

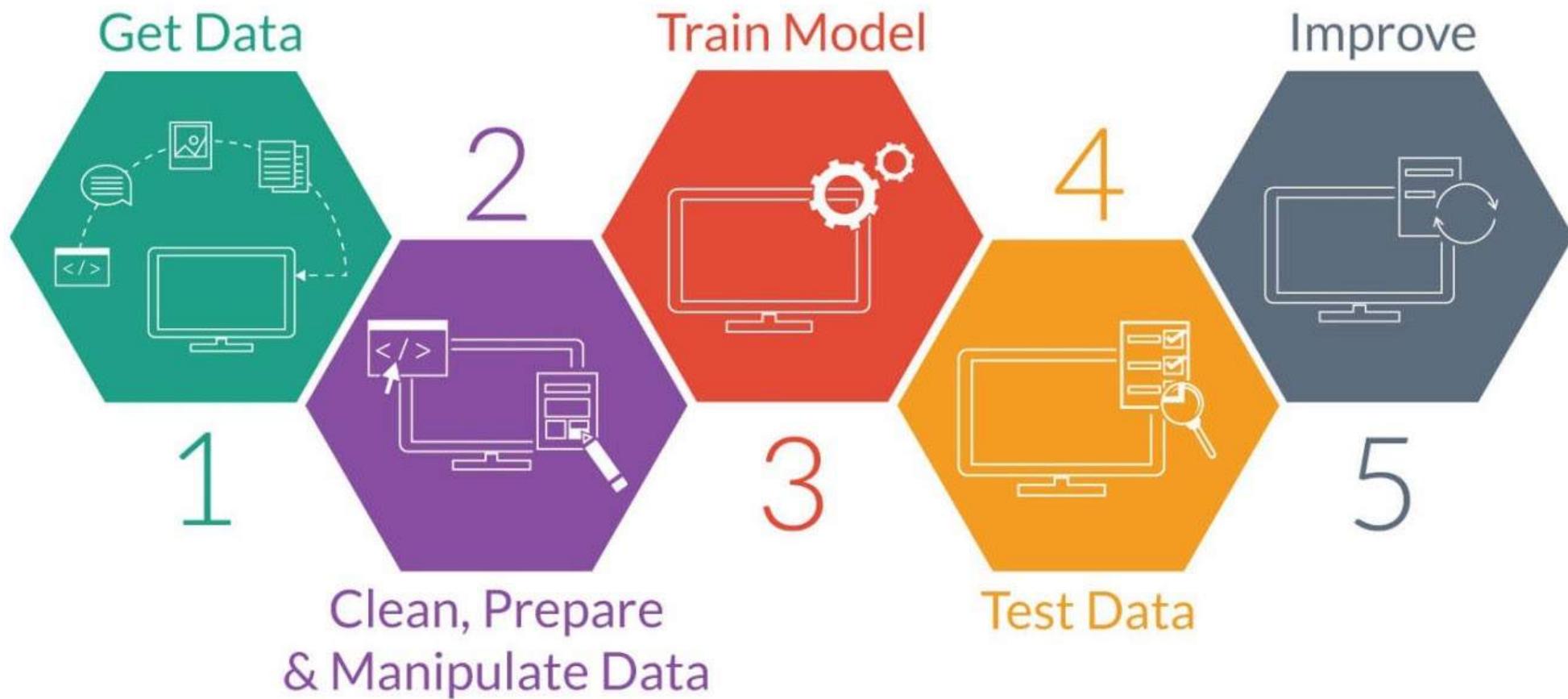


The bigger picture – AI vs ML vs DS

	 ARTIFICIAL INTELLIGENCE	 MACHINE LEARNING	 DATA SCIENCE
How does it work?	AI combines large amounts of data through iterative processing and intelligent algorithms to help computers learn automatically.	ML uses efficient programs that can use data to self learn without having to be instructed explicitly.	DS works by sourcing, cleaning, and processing data to extract meaning out of it for analytical purposes.
Popular Tools	<ul style="list-style-type: none">• TensorFlow• Scikit Learn• Keras	<ul style="list-style-type: none">• Amazon Lex• Scikit Learn• IBM Watson Studio• Microsoft Azure Machine Learning Studio	<ul style="list-style-type: none">• SAS• Tableau• Apache Spark• MATLAB• BigML
Top Applications	<ul style="list-style-type: none">• Chatbots• Cybersecurity• Voice assistant• Healthcare	<ul style="list-style-type: none">• Autonomous driving• Recommendation system• Facial recognition	<ul style="list-style-type: none">• Fraud & Risk detection• Healthcare analysis• Ecommerce



The Machine Learning workflow - simplified



The Machine Learning workflow - example

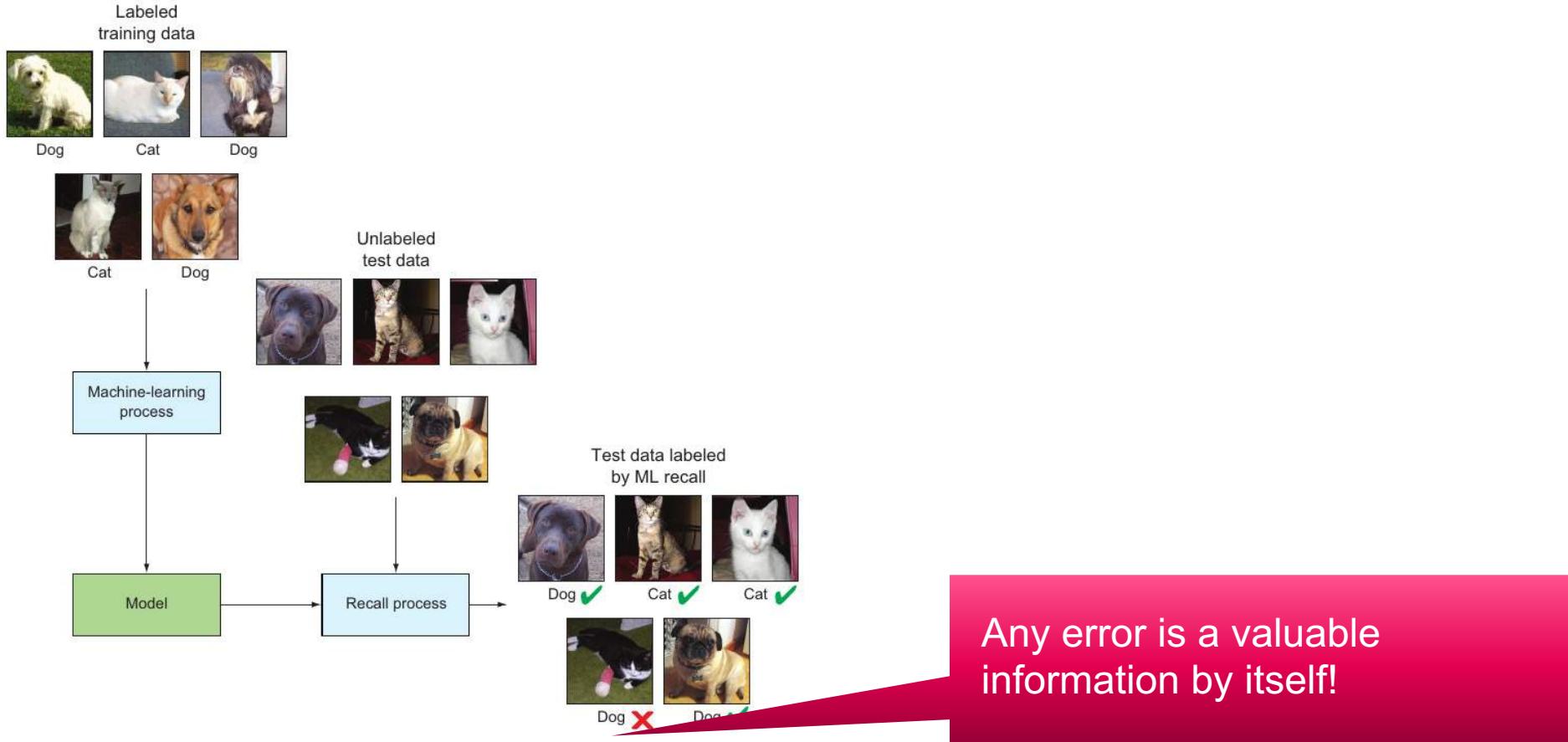


Figure 1.1 Machine-learning process for the cats and dogs competition



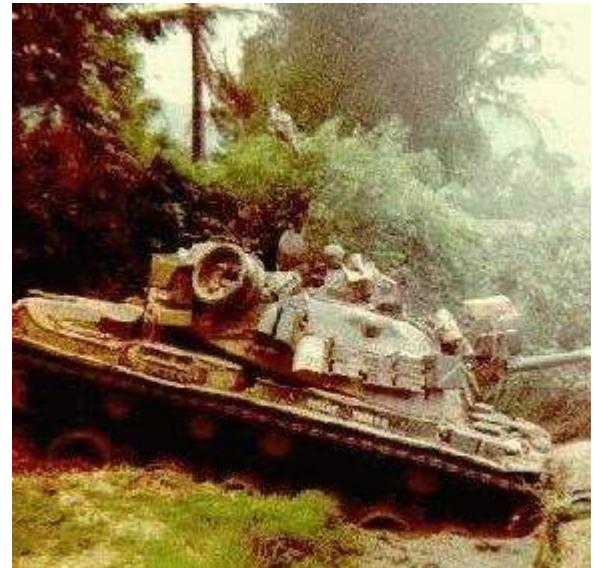
Why you should care about your data set

Detecting tanks

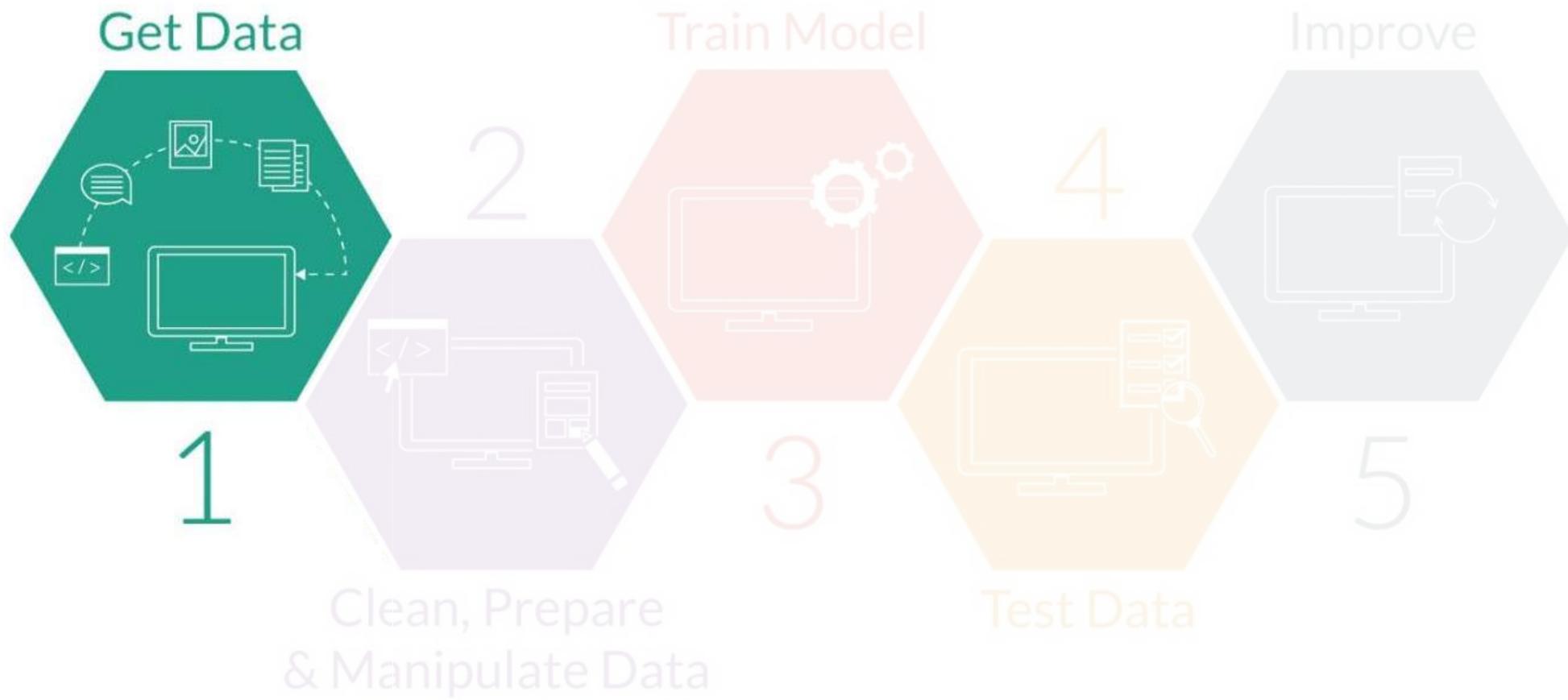
US Army wanted to use Neural Networks to detect tanks hidden behind trees

However...something went wrong

<https://neil.fraser.name/writing/tank/>



Get Data



Get Data

Machine Learning Terminology

Header	pclass	survived	name	sex	age	sibsp	parch	ticket	fare	cabin	embarked	boat	body	home.dest
Rows/ Examples/ Instances	0	1	Allen, Miss. Elisabeth Walton	female	29.0000	0	0	24160	211.3375	B5	S	2	NaN	St Louis, MO
	1	1	Allison, Master. Hudson Trevor	male	0.9167	1	2	113781	151.5500	C22 C26	S	11	NaN	Montreal, PQ / Chesterville, ON
	2	1	Allison, Miss. Helen Loraine	female	2.0000	1	2	113781	151.5500	C22 C26	S	NaN	NaN	Montreal, PQ / Chesterville, ON
	3	1	Allison, Mr. Hudson Joshua Creighton	male	30.0000	1	2	113781	151.5500	C22 C26	S	NaN	135.0	Montreal, PQ / Chesterville, ON
	4	1	Allison, Mrs. Hudson J C (Bessie Waldo Daniels)	female	25.0000	1	2	113781	151.5500	C22 C26	S	NaN	NaN	Montreal, PQ / Chesterville, ON
	5	1	Anderson, Mr. Harry	male	48.0000	0	0	19952	26.5500	E12	S	3	NaN	New York, NY
	6	1	Andrews, Miss. Kornelia Theodosia	female	63.0000	1	0	13502	77.9583	D7	S	10	NaN	Hudson, NY
	7	1	Andrews, Mr. Thomas Jr	male	39.0000	0	0	112050	0.0000	A36	S	NaN	NaN	Belfast, NI
	8	1	Appleton, Mrs. Edward Dale (Charlotte Lamson)	female	53.0000	2	0	11769	51.4792	C101	S	D	NaN	Bayside, Queens, NY
	9	1	Artagaveytia, Mr. Ramon	male	71.0000	0	0	PC 17609	49.5042	NaN	C	NaN	22.0	Montevideo, Uruguay
Index	Features/Attributes/Dimensions													



How to make sense of your data

You can visualize your data

- mosaic plot
- ...

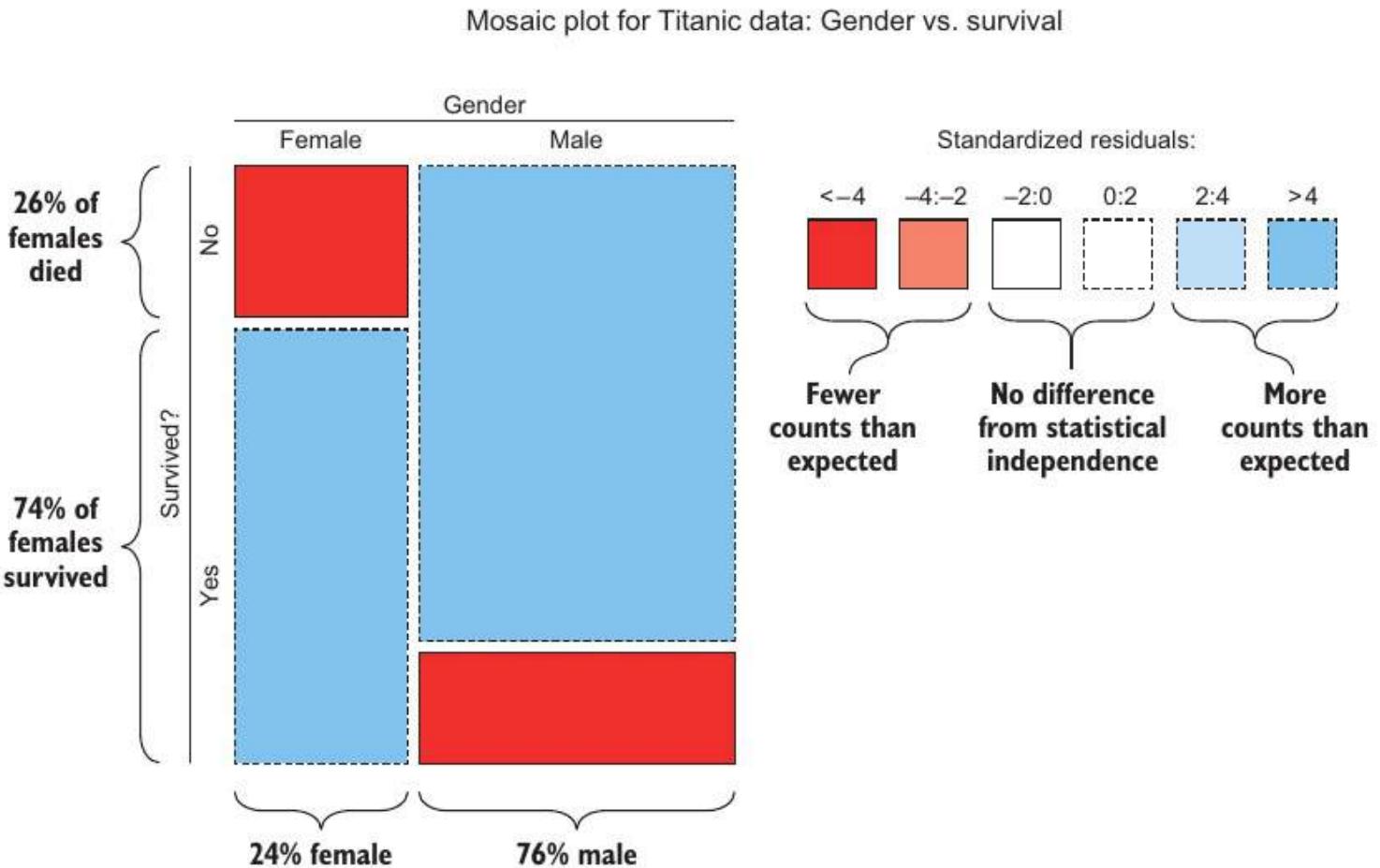


Figure 2.12 Mosaic plot showing the relationship between gender and survival on the Titanic. The visualization shows that a much higher proportion of females (and much smaller proportion of males) survived than would have been expected if survival were independent of gender. “Women and children first.”



How to make sense of your data

You can visualize your data

- Mosaic plot
- Box plot
- ...

Box plot for Titanic data: Passenger age vs. survival

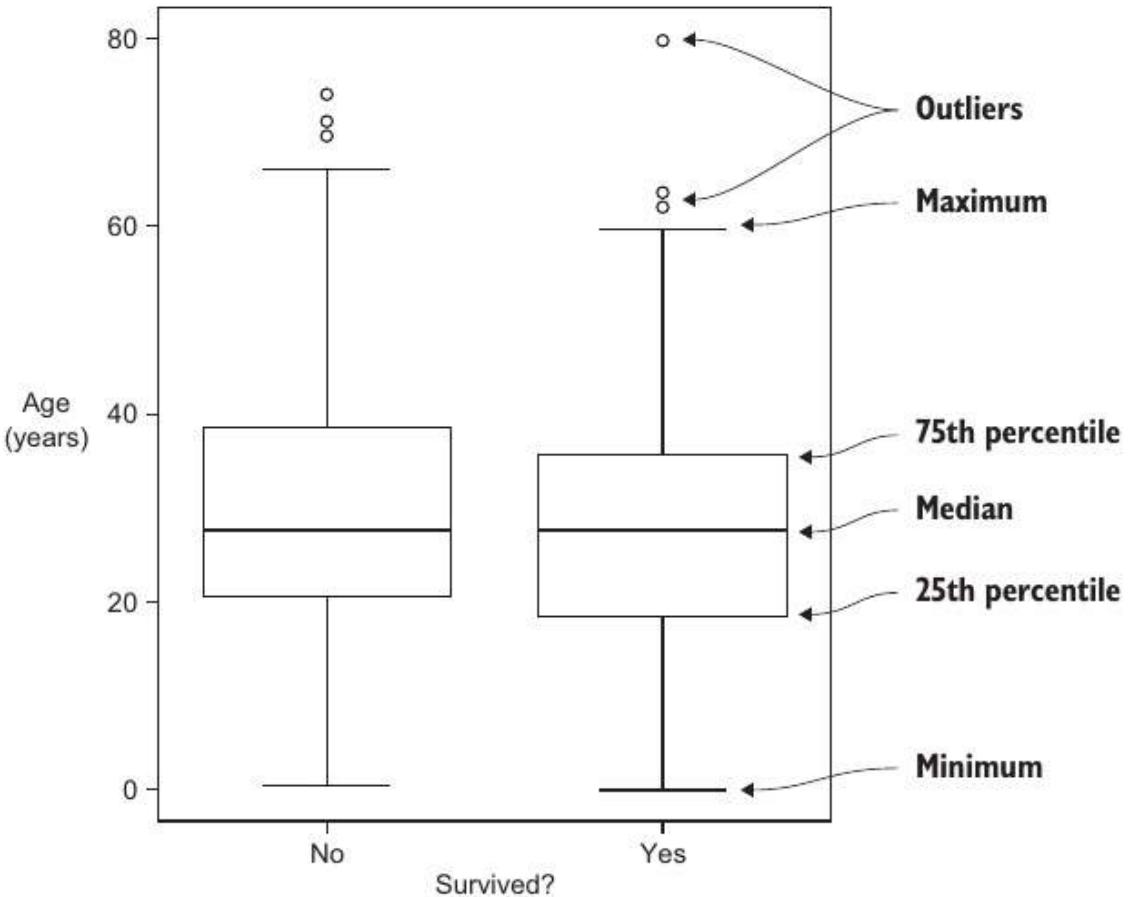


Figure 2.14 Box plot showing the relationship between passenger age and survival on the Titanic. No noticeable differences exist between the age distributions for survivors versus fatalities. (This alone *shouldn't* be a reason to exclude age from the ML model, as it may still be a predictive factor.)



How to make sense of your data

You can visualize your data

- Mosaic plot
- Box plot
- Density plot
- Scatter plot (matrix)
- ...

Scatterplots for MPG data

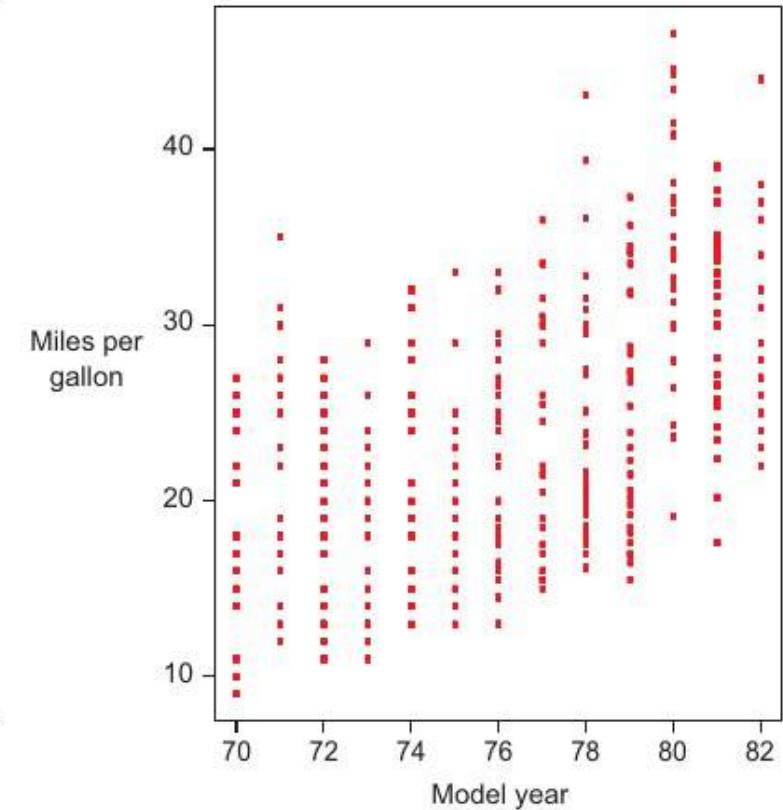
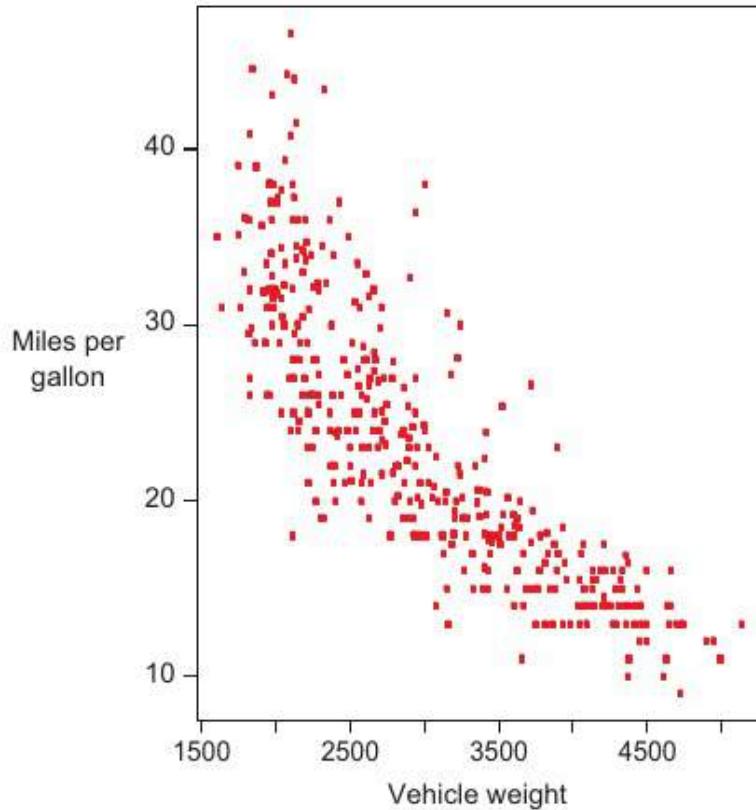


Figure 2.17 Scatter plots for the relationship of vehicle miles per gallon versus vehicle weight (left) and vehicle model year (right)



Clean and Prepare Data



Clean and Prepare Data

PassengerId	Survived	Pclass	Gender	Age	SibSp	Parch	Ticket	Fare	Cabin	Embarked
1	0	3	Male	22	1	0	A/5 21171	7.25	(Circled)	S
2	1	1	Female	38	1	0	PC 17599	71.2833	(Circled)	C
3	1	3	Female	26	0	0	STON/O2. 3101282	7.925	(Circled)	S
4	1	1	Female	35	1	0	113803	53.1	(Circled)	S
5	0	3	Male	35	0	0	373450	8.05	(Circled)	S
6	0	3	Male	(Circled)	0	0	330877	8.4583	(Circled)	Q

Missing values

Figure 2.7 The Titanic Passengers dataset has missing values in the Age and Cabin columns. The passenger information has been extracted from various historical sources, so in this case the missing values stem from information that couldn't be found in the sources.



Clean and Prepare Data

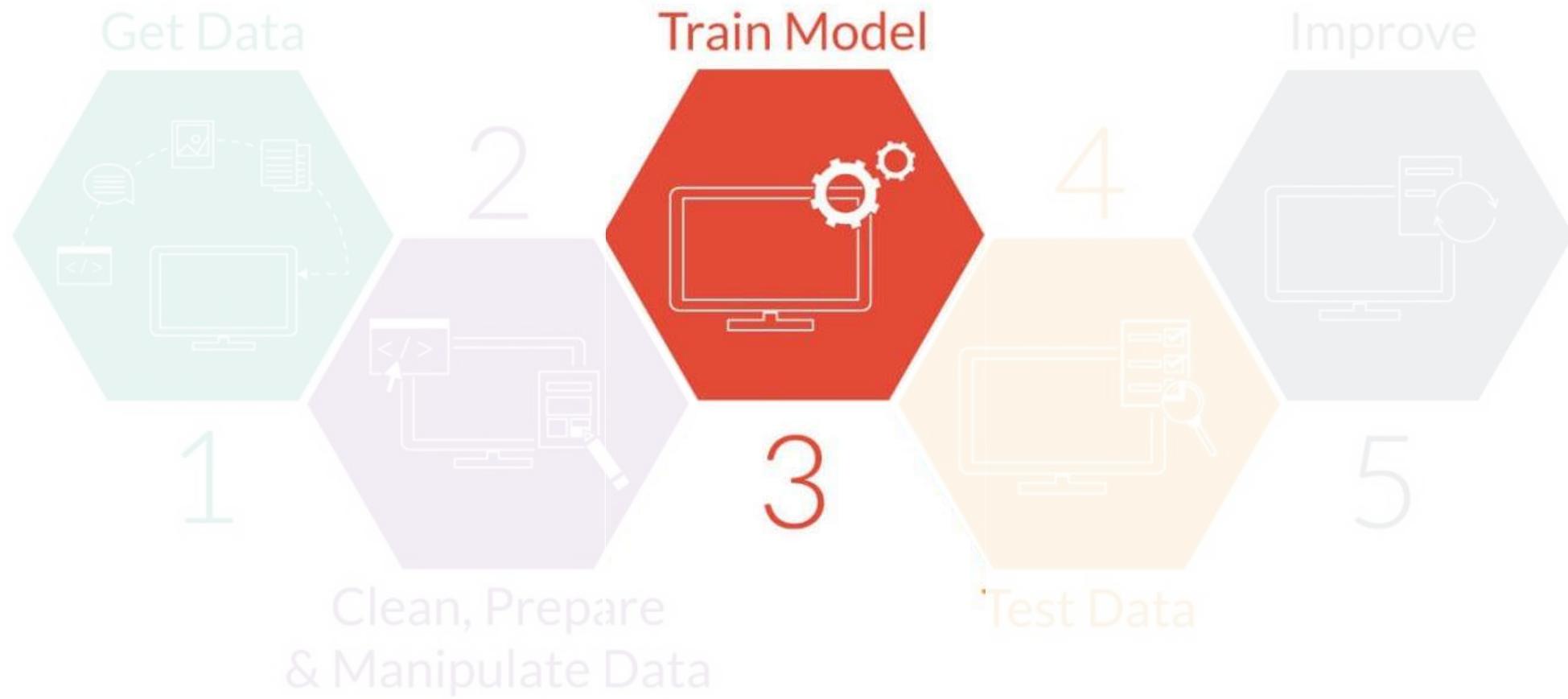
Person	Name	Age	Income	Marital status
1	Jane Doe	24	81,200	Single
2	John Smith	41	121,000	Married

PassengerId	Survived	Pclass	Gender	Age	SibSp	Parch	Ticket	Fare	Cabin	Embarked
1	0	3	Male	22	1	0	A/5 21171	7.25		S
2	1	1	Female	38	1	0	PC 17599	71.2833	C85	C
3	1	3	Female	26	0	0	STON/O2. 3101282	7.925		S
4	1	1	Female	35	1	0	113803	53.1	C123	S
5	0	3	Male	35	0	0	373450	8.05		S
6	0	3	Male		0	0	330877	8.4583		Q

Figure 2.4 Identifying categorical features. At the top is the simple Person dataset, which has a Marital Status categorical feature. At the bottom is a dataset with information about Titanic passengers. The features identified as categorical here are Survived (whether the passenger survived or not), Pclass (what class the passenger was traveling on), Gender (male or female), and Embarked (from which city the passenger embarked).



Clean and Prepare Data



Train Model

Features in columns				
Person	Name	Age	Income	Marital status
1	Jane Doe	24	81,200	Single
2	John Smith	41	121,000	Married

Figure 1.8 In a tabular dataset, rows are called *instances* and columns represent *features*.

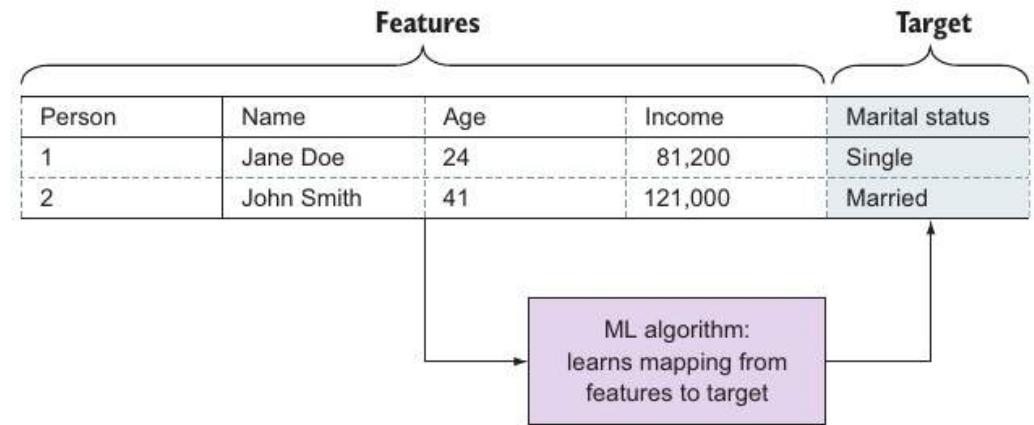
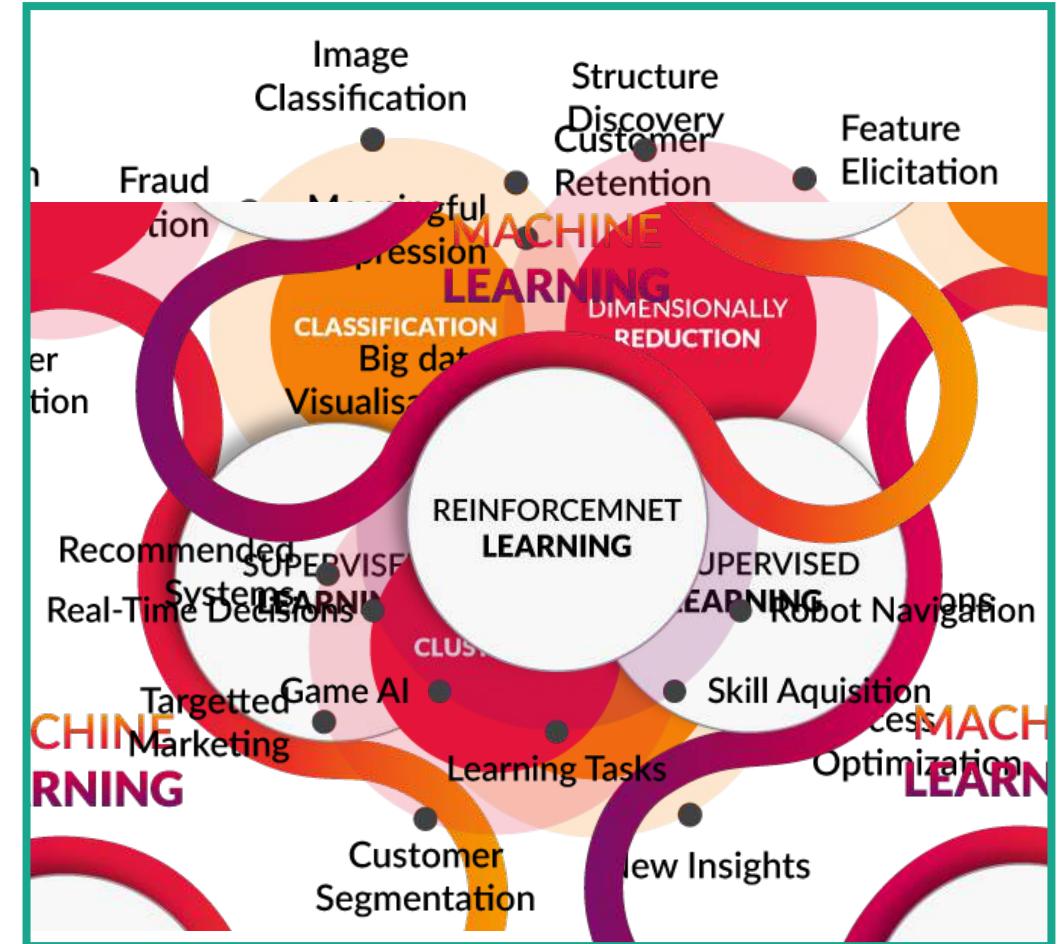
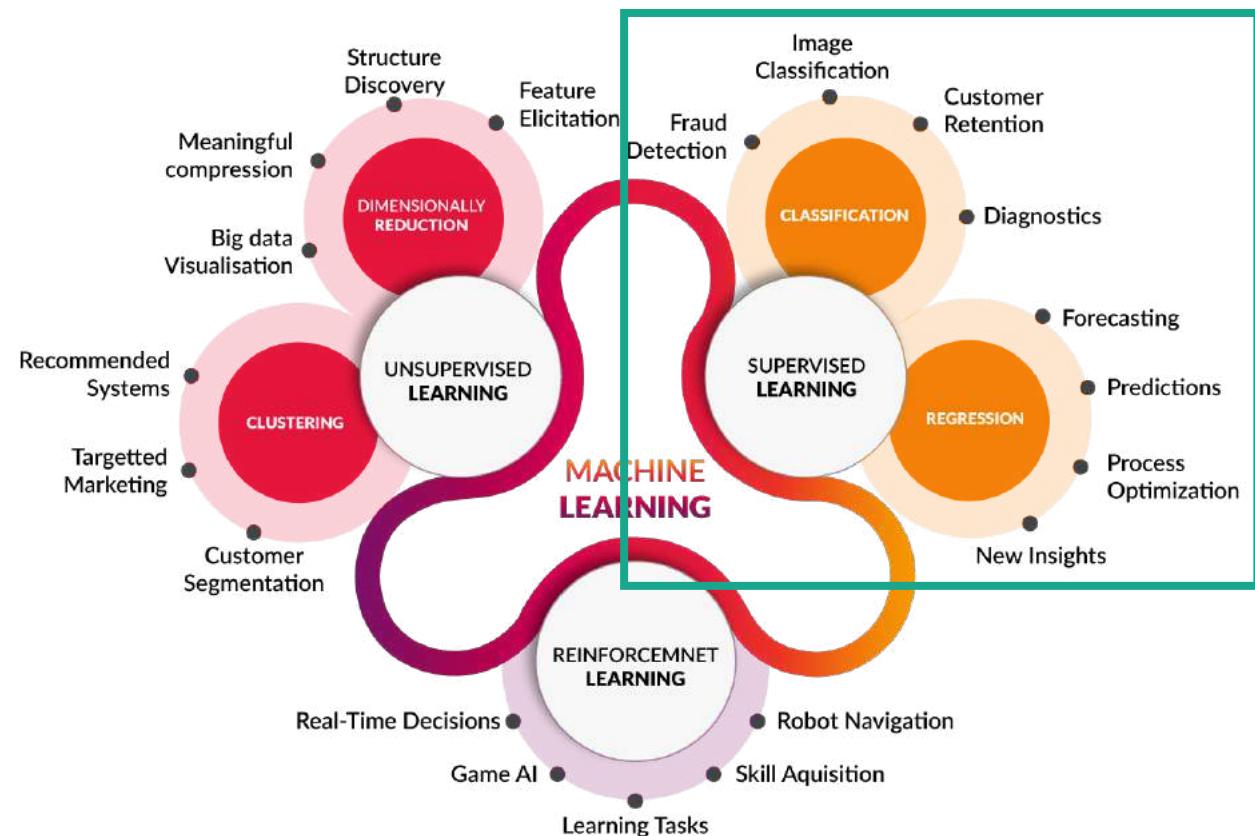


Figure 1.9 The machine-learning modeling process

But what algorithm?



The landscape of Machine Learning



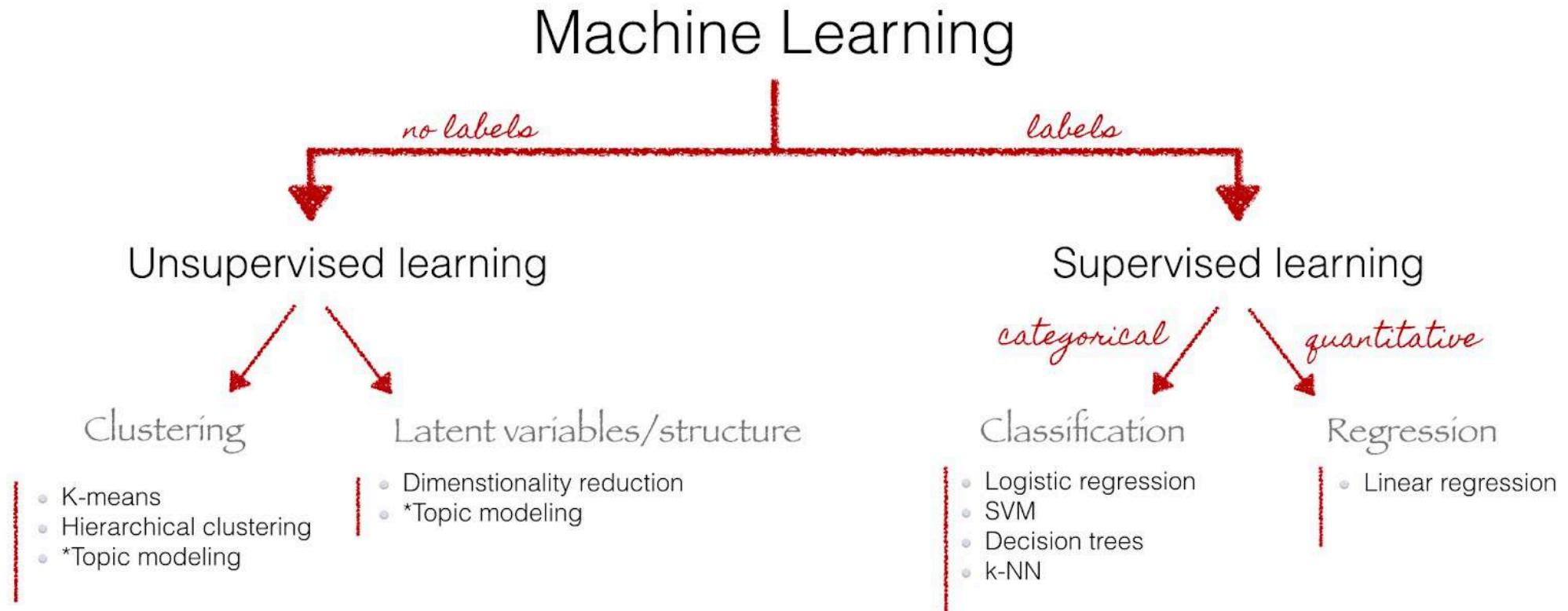
Source: Towards Data Science

Read further: <https://litslink.com/blog/an-introduction-to-machine-learning-algorithms>

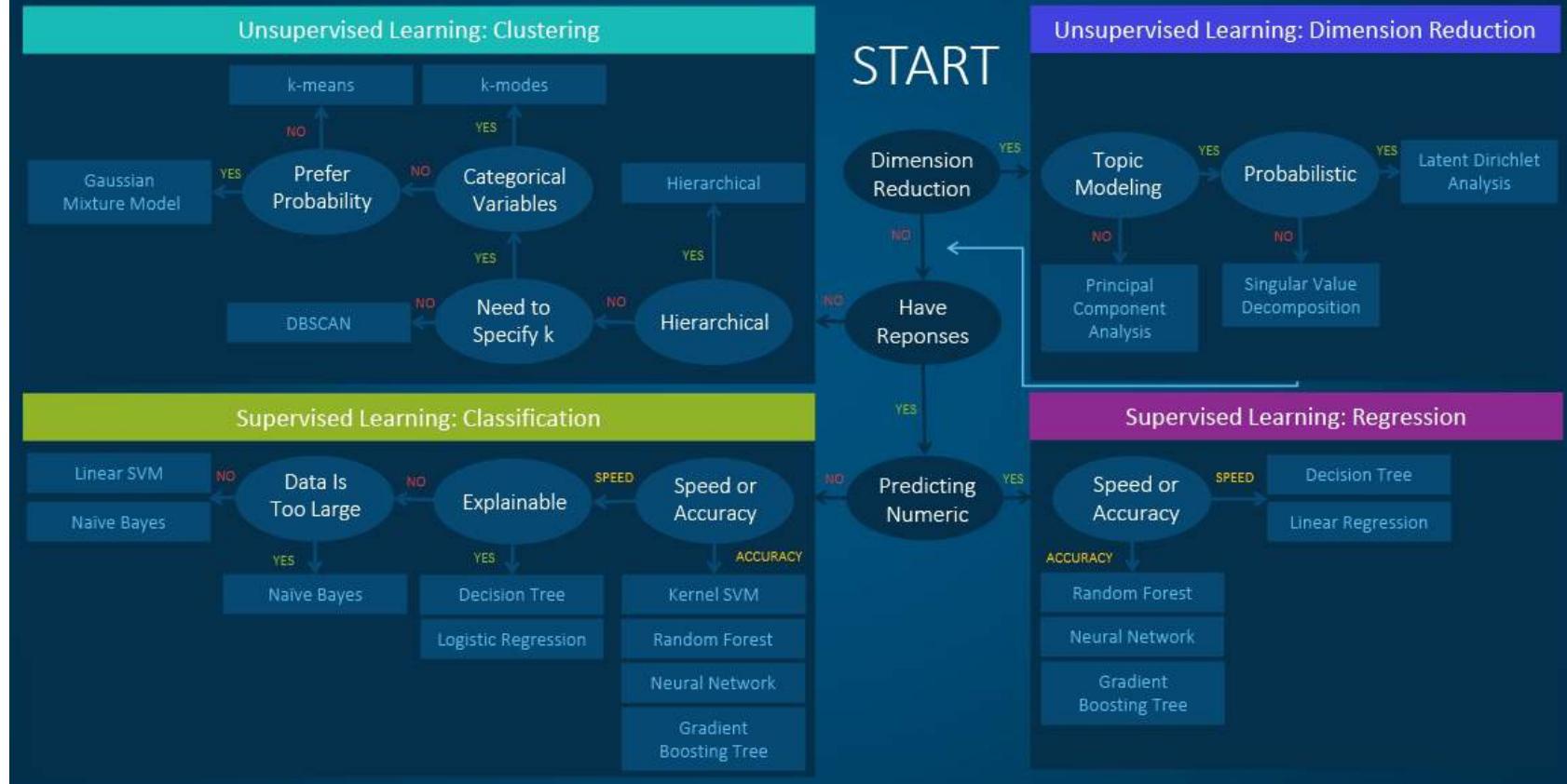


The landscape of Machine Learning

Common algorithms



Machine Learning Algorithms Cheat Sheet



Classification vs Regression

Regression

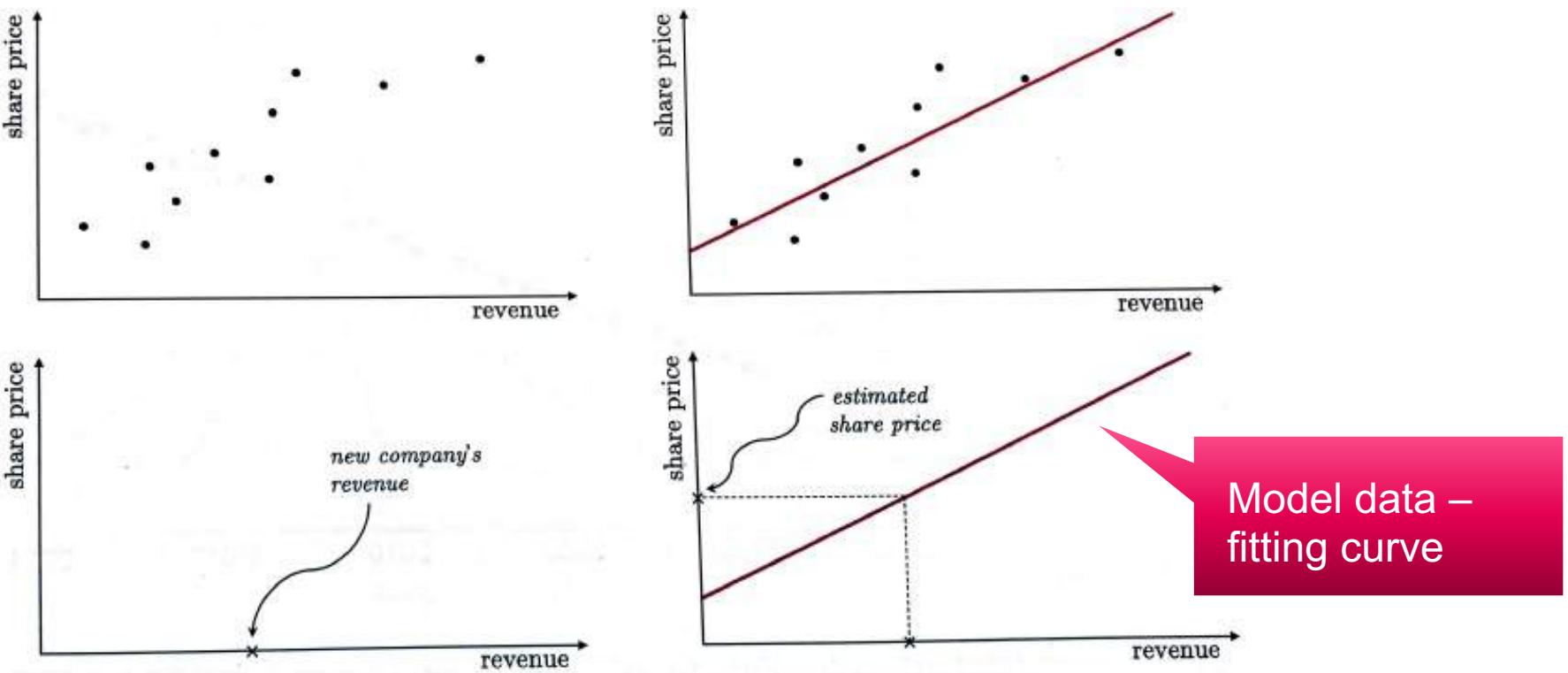


Fig. 1.7

(top left panel) A toy training dataset of ten corporations with their associated share price and revenue values. (top right panel) A linear model is fit to the data. This trend line models the overall trajectory of the points and can be used for prediction in the future as shown in the bottom left and bottom right panels.



Classification vs Regression

Classification

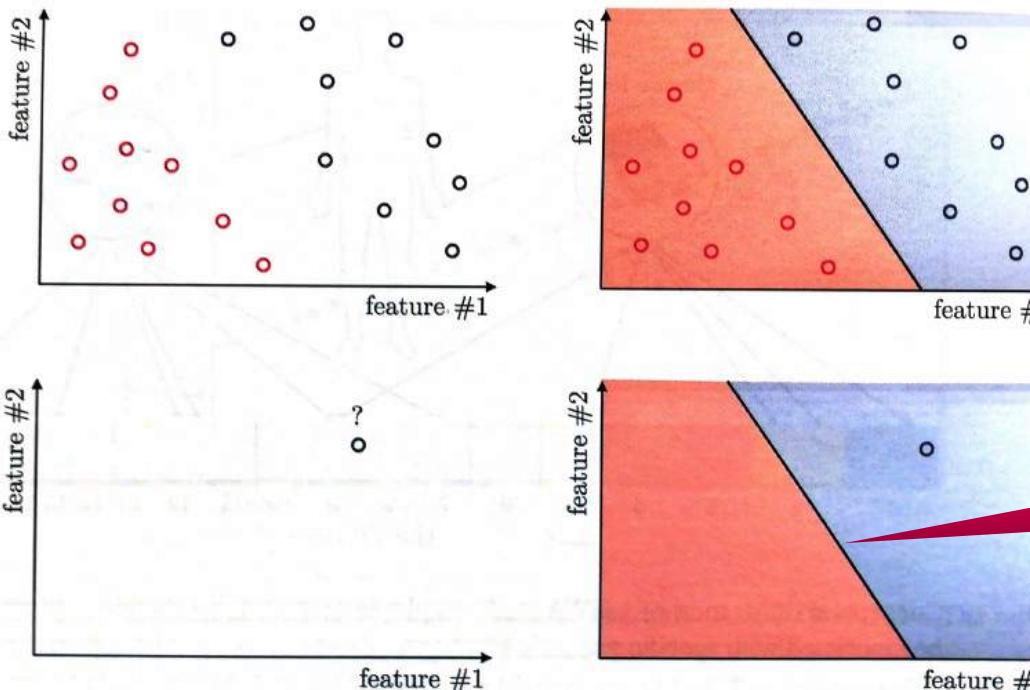


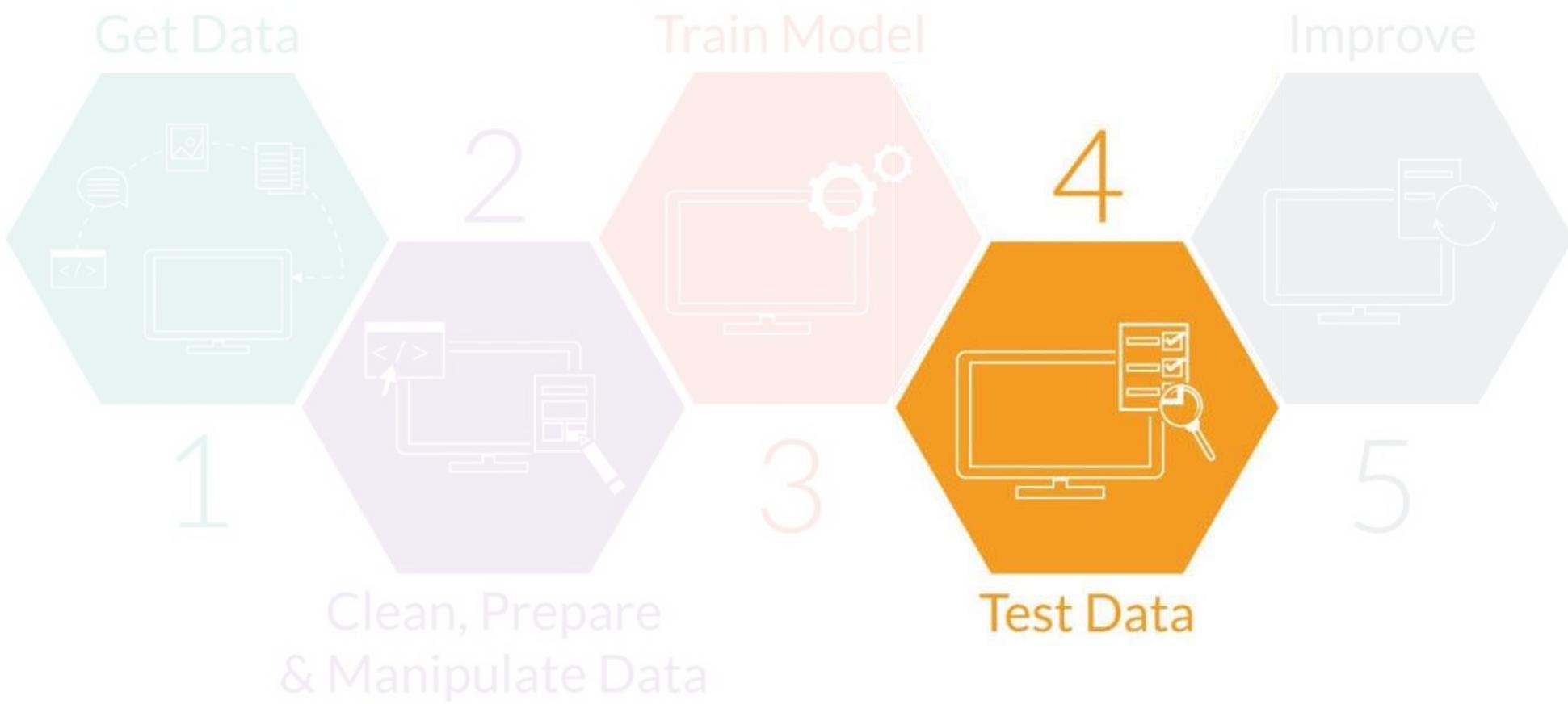
Fig. 1.10

(top left panel) A toy 2-dimensional training set consisting of two distinct classes, red and blue. (top right panel) A linear model is trained to separate the two classes. (bottom left panel) A test point whose class is unknown. (bottom right panel) The test point is classified as blue since it lies on the blue side of the trained linear classifier.

Divide data –
Decision plane



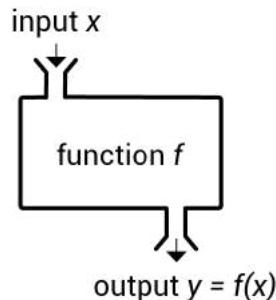
Test Data



Test Data

Features in columns					
Person	Name	Age	Income	Marital status	
1	Jane Doe	24	81,200	Single	
2	John Smith	41	121,000	Married	

Figure 1.8 In a tabular dataset, rows are called *instances* and columns represent *features*.



Testing data with target					
Person	Name	Age	Income	Marital status	
1	Trent Mosley	26	67,500	Single	
2	Lilly Peters	52	140,000	Married	

ML model: predicts the target variable on new data

Marital status
Single
Married

Figure 1.11 When using a testing set to evaluate model performance, you “pretend” that the target variable is unknown and compare the predictions with the true values.

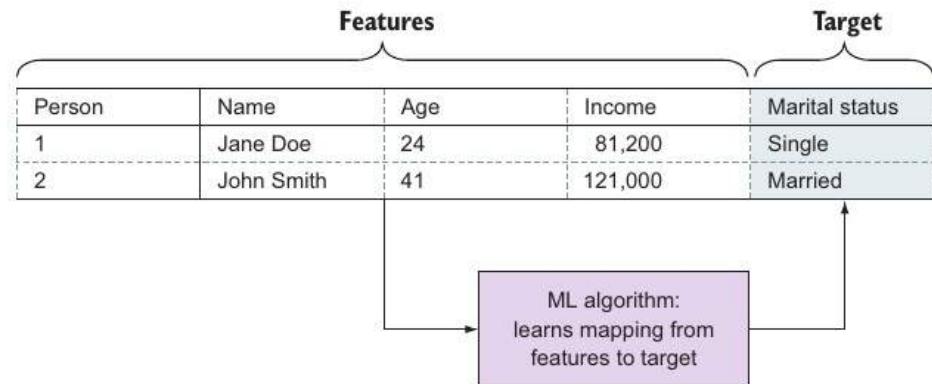


Figure 1.9 The machine-learning modeling process

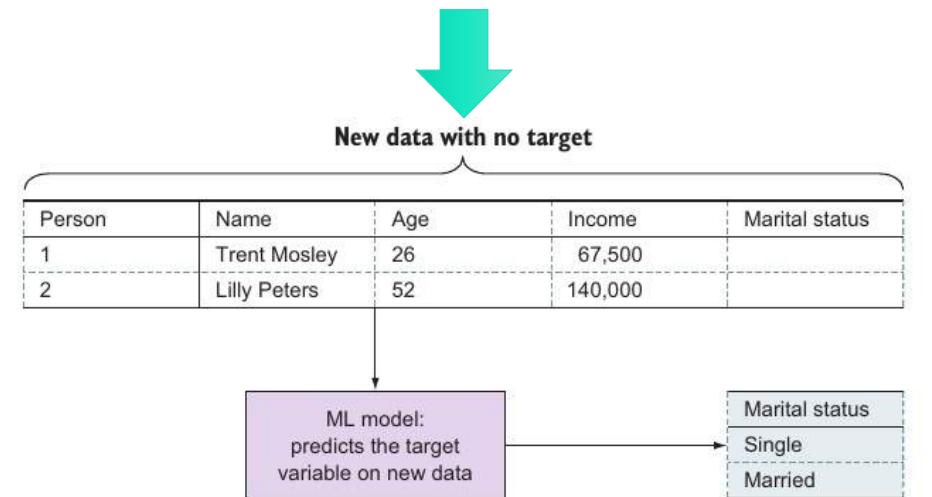


Figure 1.10 Using the model for prediction on new data



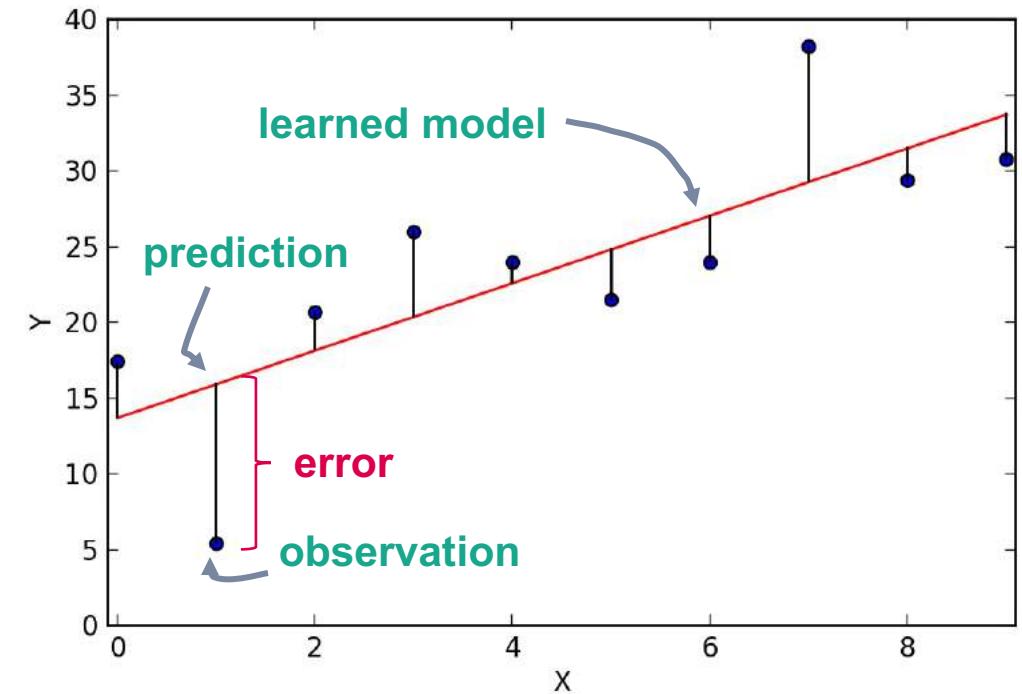
How do we measure our model's performance?

Regression

- (root) mean square error

$$RMSE = \sqrt{\sum_{i=1}^n \frac{(\hat{y}_i - y_i)^2}{n}}$$

Mean deviation



How do we measure our model's performance?

Classification

- (root) mean square error
- confusion matrix

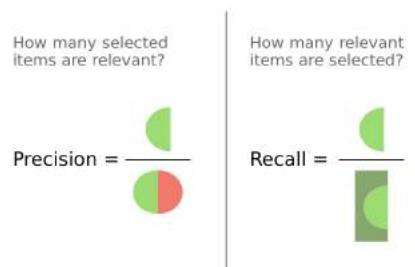
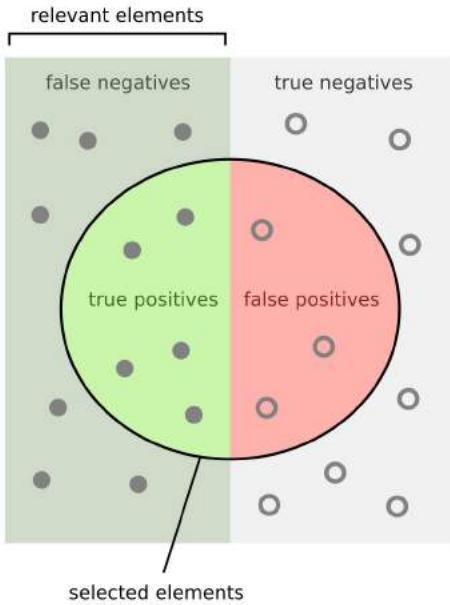
$$\text{Accuracy} = \frac{\text{True Positives} + \text{True Negatives}}{\text{All Samples}}$$

Percentage of correct predictions

		Predicted	
		Negative	Positive
Actual	Negative	True Negative	False Positive
	Positive	False Negative	True Positive



Beyond accuracy



		True condition		Prevalence = $\frac{\sum \text{Condition positive}}{\sum \text{Total population}}$	Accuracy (ACC) = $\frac{\sum \text{True positive} + \sum \text{True negative}}{\sum \text{Total population}}$
Total population	Predicted condition positive	Condition positive	Condition negative		
Predicted condition	Predicted condition positive	True positive	False positive, Type I error	Positive predictive value (PPV), Precision = $\frac{\sum \text{True positive}}{\sum \text{Predicted condition positive}}$	False discovery rate (FDR) = $\frac{\sum \text{False positive}}{\sum \text{Predicted condition positive}}$
	Predicted condition negative	False negative, Type II error	True negative	False omission rate (FOR) = $\frac{\sum \text{False negative}}{\sum \text{Predicted condition negative}}$	Negative predictive value (NPV) = $\frac{\sum \text{True negative}}{\sum \text{Predicted condition negative}}$
		True positive rate (TPR), Recall, Sensitivity, probability of detection, Power = $\frac{\sum \text{True positive}}{\sum \text{Condition positive}}$	False positive rate (FPR), Fall-out, probability of false alarm = $\frac{\sum \text{False positive}}{\sum \text{Condition negative}}$	Positive likelihood ratio (LR+) = $\frac{\text{TPR}}{\text{FPR}}$	Diagnostic odds ratio (DOR) = $\frac{\text{LR+}}{\text{LR-}}$
		False negative rate (FNR), Miss rate = $\frac{\sum \text{False negative}}{\sum \text{Condition positive}}$	Specificity (SPC), Selectivity, True negative rate (TNR) = $\frac{\sum \text{True negative}}{\sum \text{Condition negative}}$	Negative likelihood ratio (LR-) = $\frac{\text{FNR}}{\text{TNR}}$	
		F_1 score = $2 \cdot \frac{\text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}}$			

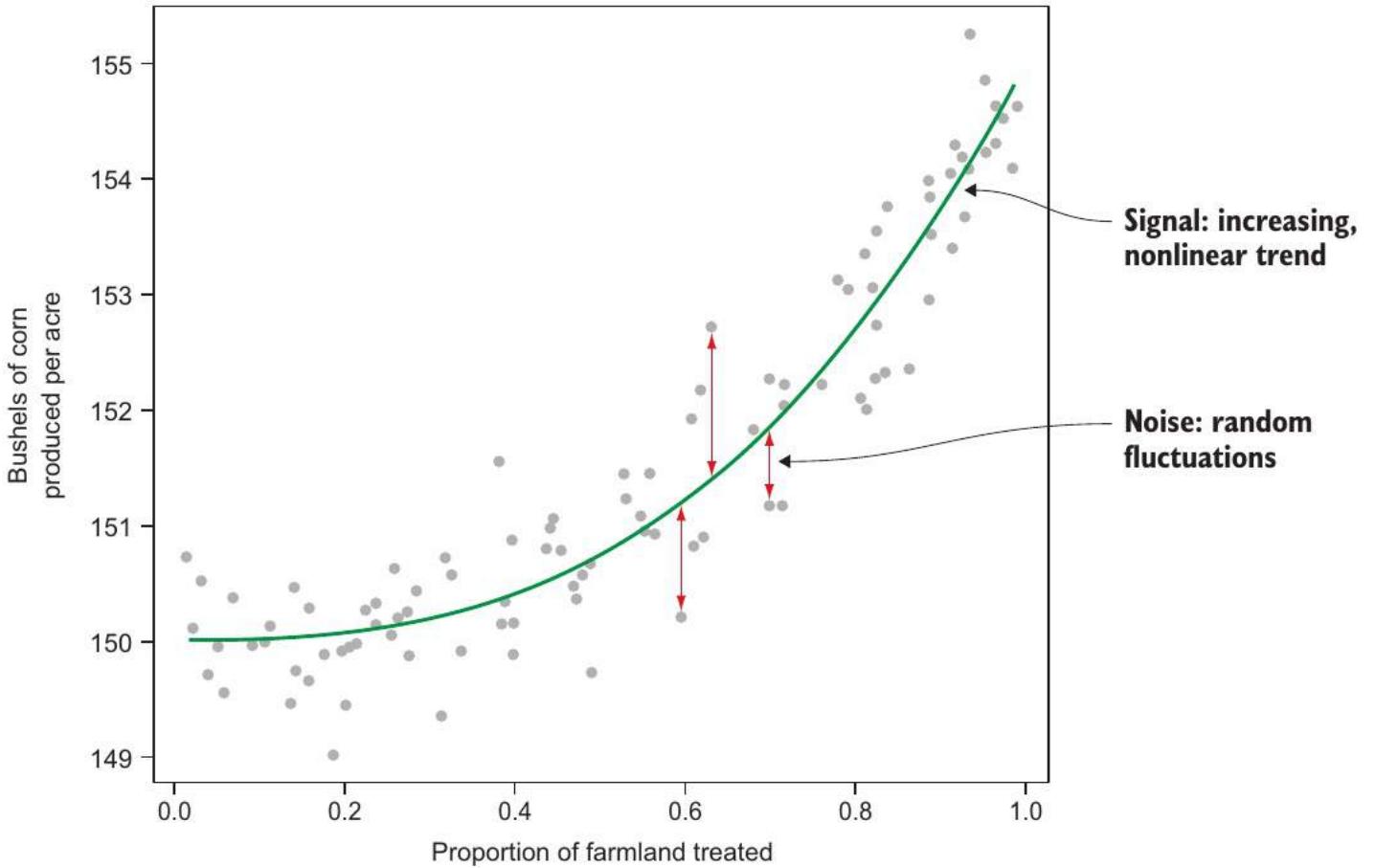


Improve Model



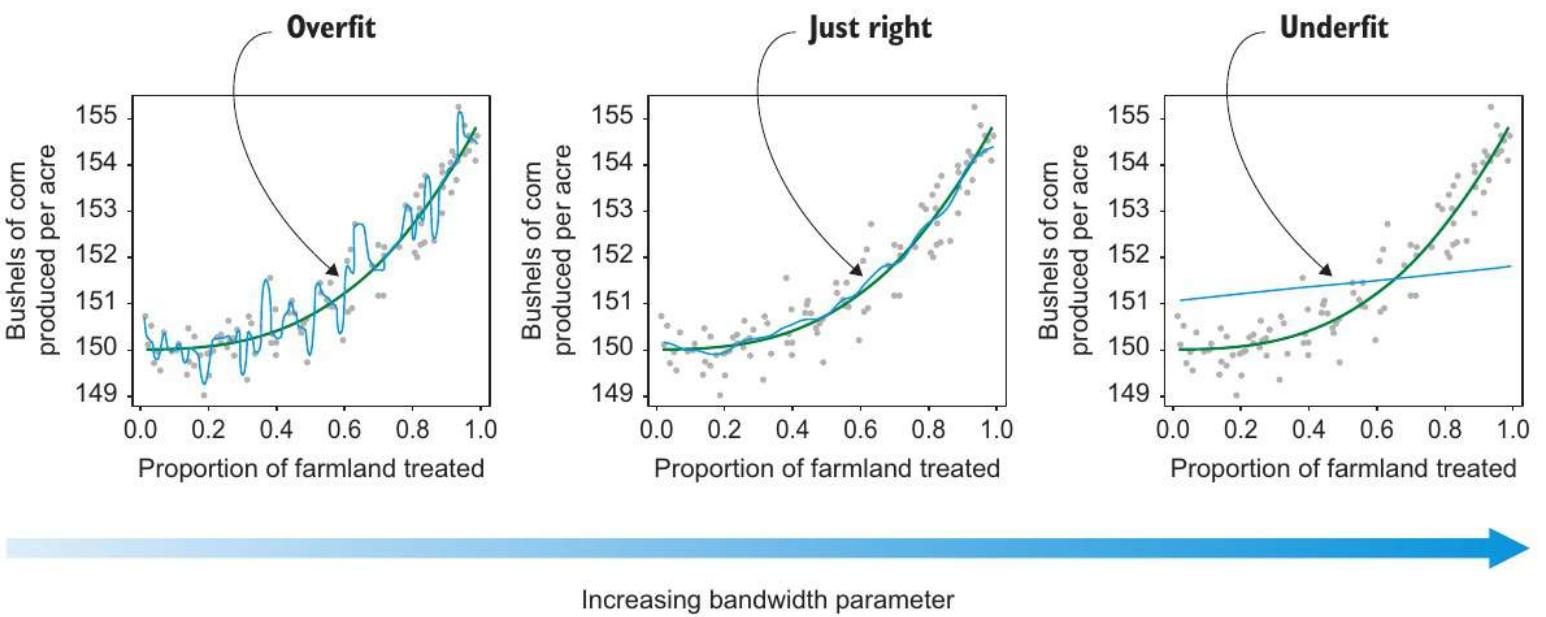
Real world data

signal-to-noise ratio



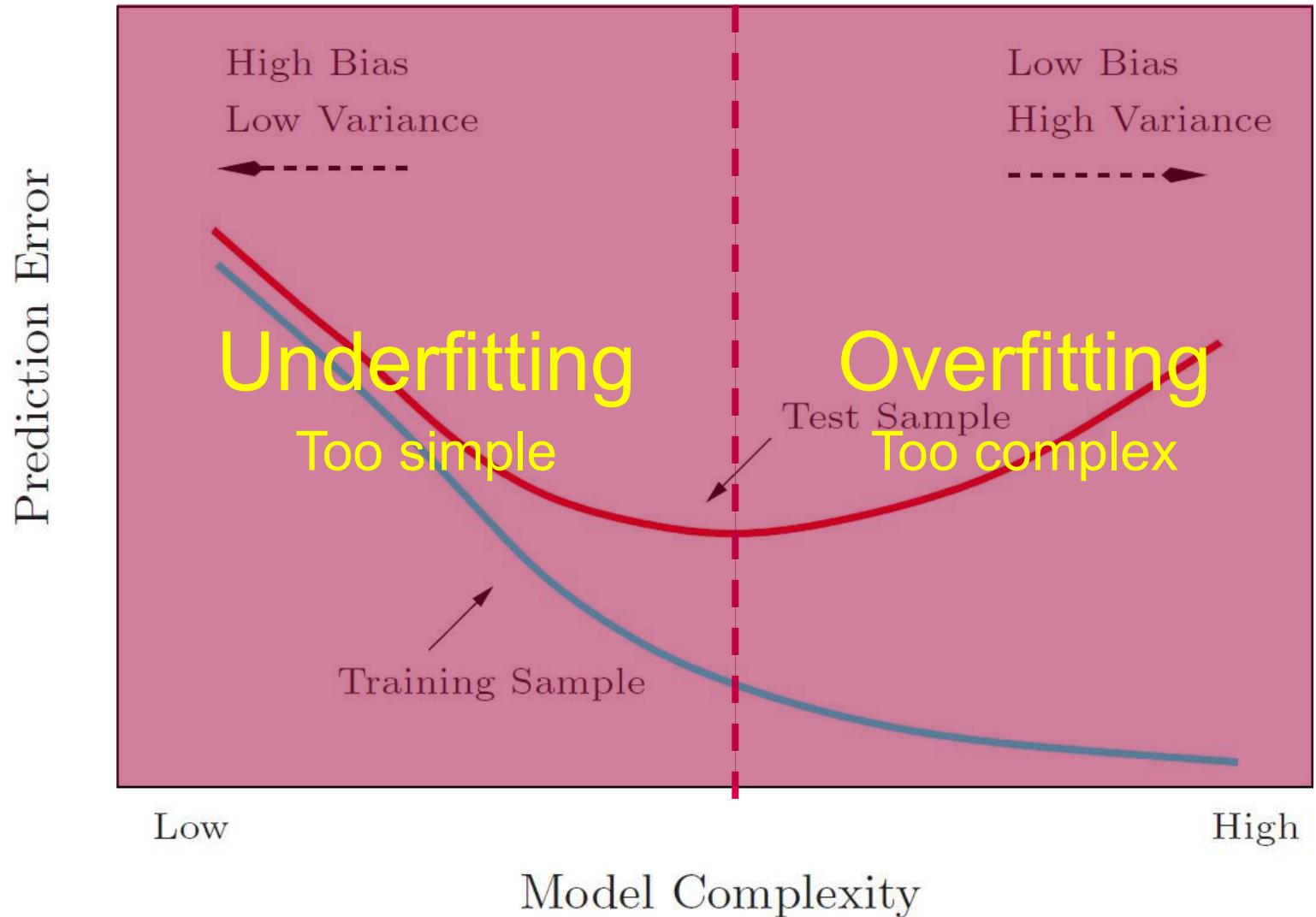
What makes a good model?

It may happen that the models we build are underfitted or overfitted



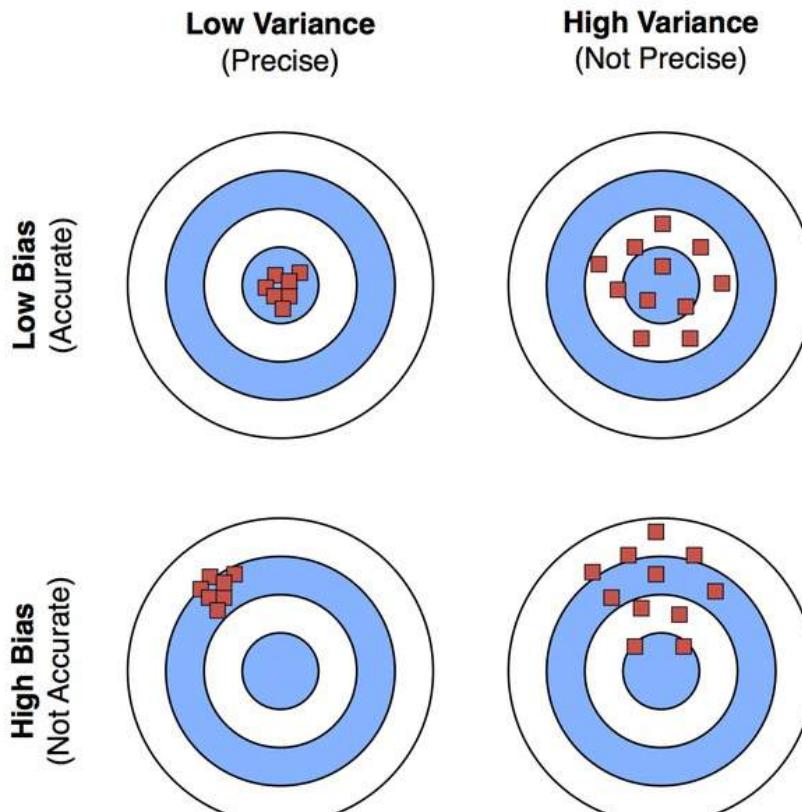
How do we measure our model's performance?

When we have a performance measure, we can compare different models and their **complexity**



How do we measure our model's performance?

Variance-Bias Trade-off

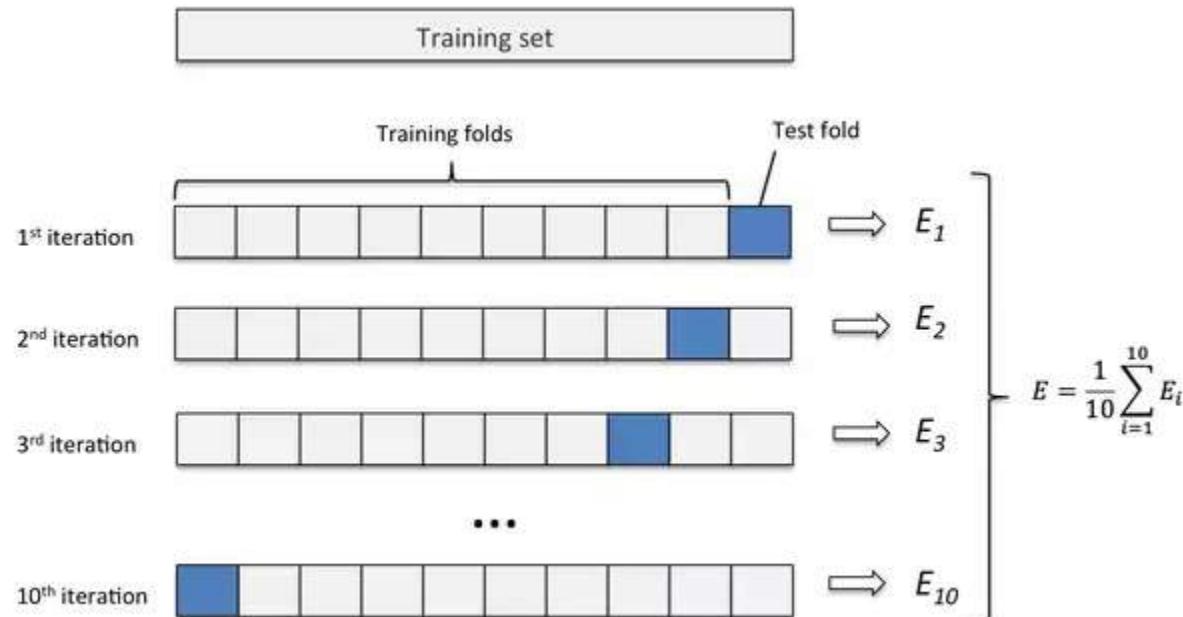


This work by Sebastian Raschka is licensed under a
Creative Commons Attribution 4.0 International License.



What about small data sets?

Cross-validation



What about large data sets?

Curse of dimensionality

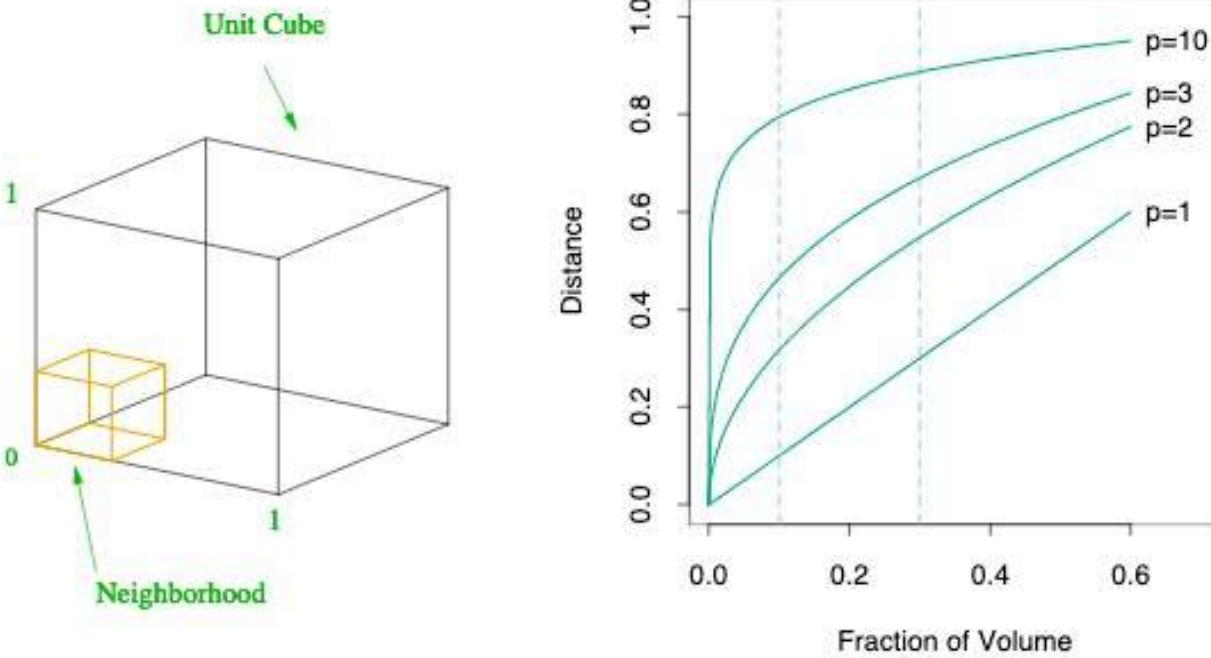


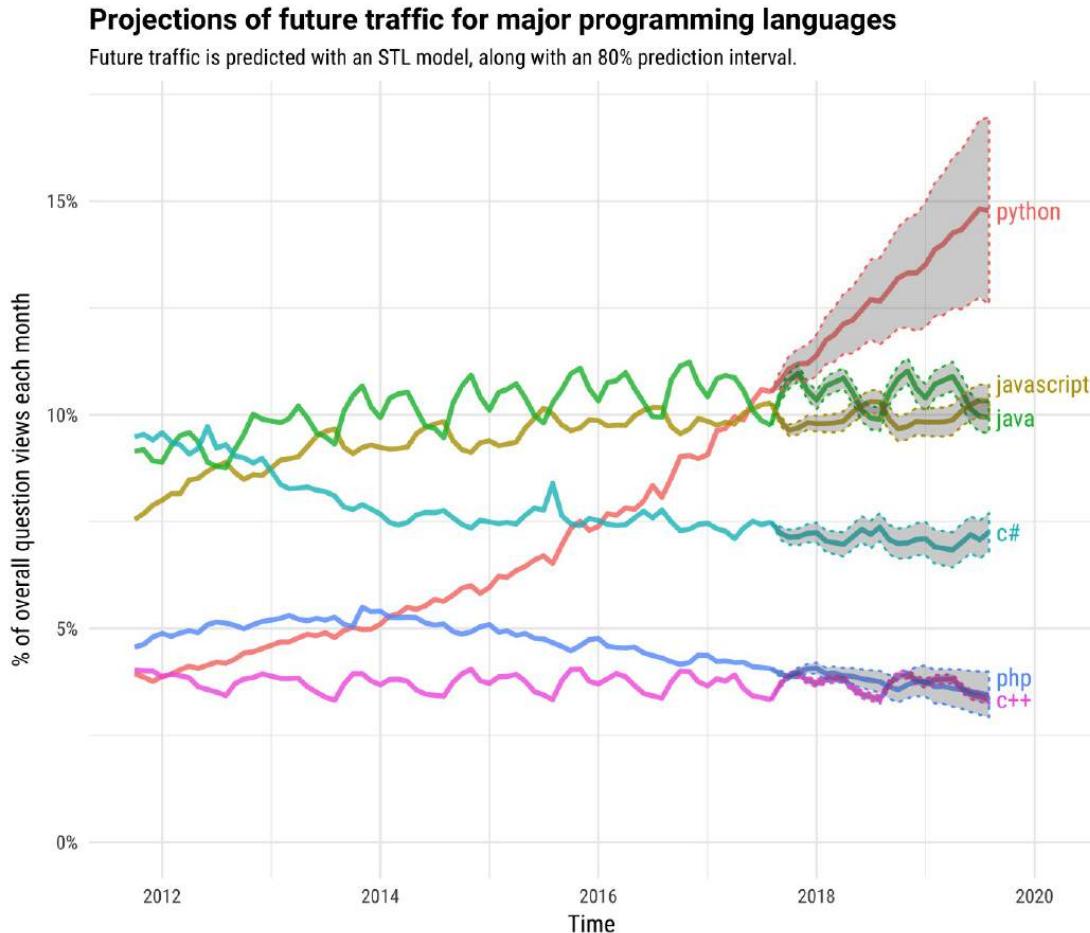
FIGURE 2.6. The curse of dimensionality is well illustrated by a subcubical neighborhood for uniform data in a unit cube. The figure on the right shows the side-length of the subcube needed to capture a fraction r of the volume of the data, for different dimensions p . In ten dimensions we need to cover 80% of the range of each coordinate to capture 10% of the data.



Practice



Why Python



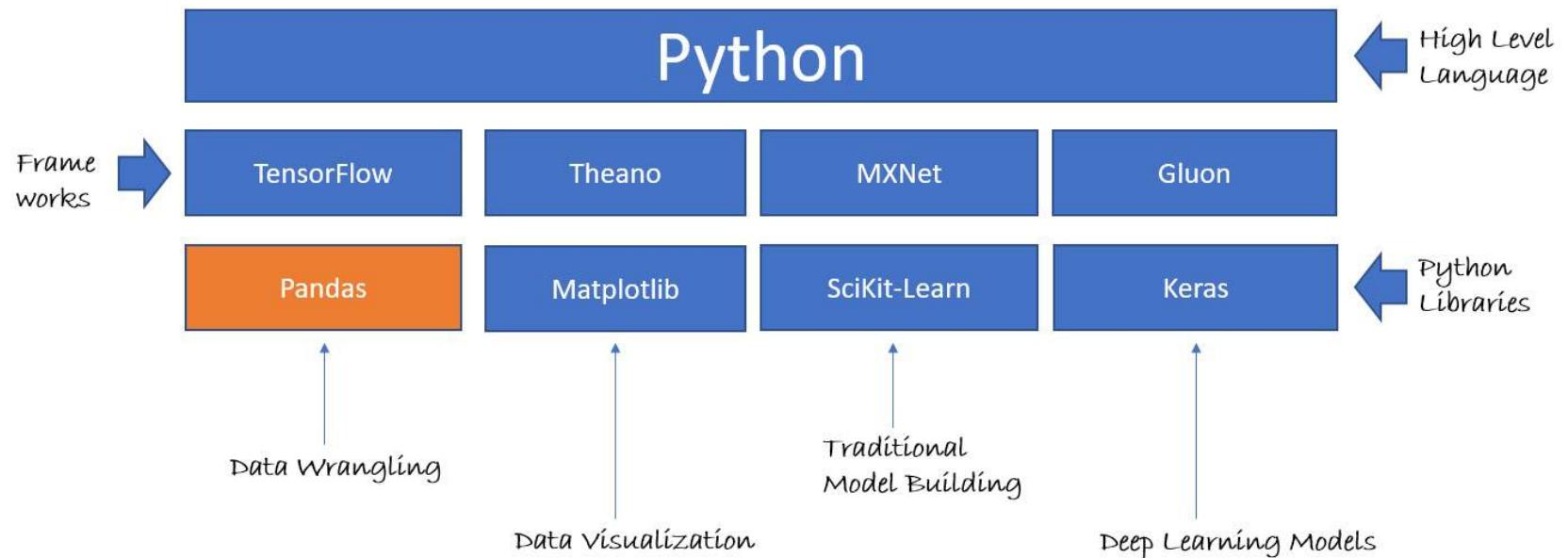
Source: Stack Overflow

<https://medium.com/@UdacityINDIA/why-use-python-for-machine-learning-e4b0b4457a77>



Python ecosystem

Machine learning is a **method of data analysis** that automates **analytical model building**. In the applied space most these models are built in a language called **Python**.



Python learning path



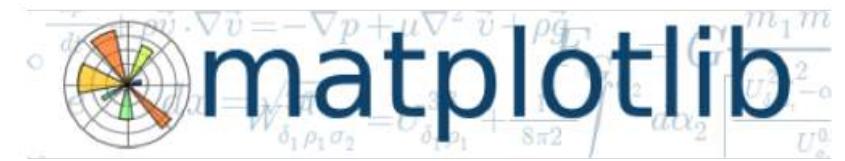
Packages, Packages everywhere !

- Want to work with images — numpy, opencv, scikit
- Want to work in text — nltk, numpy, scikit
- Want to work in audio — librosa
- Want to solve **machine learning problem** — pandas, scikit
- Want to see the data clearly — matplotlib, seaborn, scikit
- Want to use **deep learning** — tensorflow, pytorch
- Want to do scientific computing — scipy
- Want to integrate web applications — Django
- Want to take a shower Well

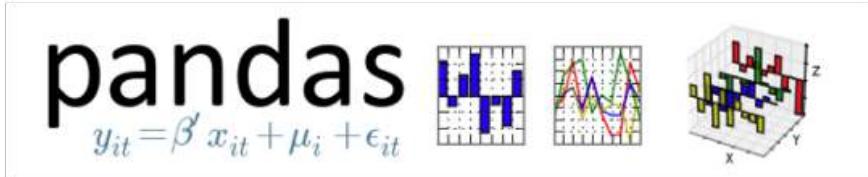


Machine Learning Landscape

Python ecosystem



IP[y]: IPython
Interactive Computing



Cheat Sheets

Python For Data Science Cheat Sheet

Python Basics

Learn More Python for Data Science interactively at www.datacamp.com

Variables and Data Types

Variable Assignment

```
>>> x=5  
>>> x  
5
```

Calculations With Variables

>>> x+2 7	Sum of two variables
>>> x-2 3	Subtraction of two variables
>>> x*2 10	Multiplication of two variables
>>> x**2 25	Exponentiation of a variable
>>> x%2 1	Remainder of a variable
>>> x/float(2) 2.5	Division of a variable

Types and Type Conversion

str()	'5', '+3.45', 'True'	Variables to strings
int()	5, 3, 1	Variables to integers
float()	5.0, 1.0	Variables to floats
bool()	True, False, True	Variables to booleans

Asking For Help

```
>>> help(str)
```

Strings

```
>>> my_string = 'thisStringIsAwesome'  
>>> my_string  
'thisStringIsAwesome'
```

String Operations

```
>>> my_string * 2  
'thisStringIsAwesomethisStringIsAwesome'  
>>> my_string + "Init"  
'thisStringIsAwesomeInit'  
>>> 'm' in my_string  
True
```

Lists

Also see NumPy Arrays

```
>>> a = 'is'  
>>> b = 'nice'  
>>> my_list = ['my', 'list', a, b]  
>>> my_list2 = [[4,5,6,7], [3,4,5,6]]
```

Selecting List Elements

Index starts at 0

Subset	Select item at index 1
>>> my_list[1]	Select 3rd last item
Slice	Select items at index 1 and 2
>>> my_list[1:]	Select items after index 0
>>> my_list[:3]	Select items before index 3
>>> my_list[:]	Copy my_list
Subset Lists of Lists	my_list[list][itemOfList]
>>> my_list2[1][0]	
>>> my_list2[1][:2]	

List Operations

```
>>> my_list + my_list  
'my', 'list', 'is', 'now', 'my', 'list', 'is', 'nice'  
>>> my_list * 2  
'my', 'list', 'is', 'now', 'my', 'list', 'is', 'nice'  
>>> my_list2 > 4
```

List Methods

>>> my_list.index('a')	Get the index of an item
>>> my_list.count('a')	Count an item
>>> my_list.append('!')	Append an item at a time
>>> my_list.remove('!')	Remove an item
>>> del(my_list[0:1])	Remove an item
>>> my_list.reverse()	Reverse the list
>>> my_list.extend('!')	Append an item
>>> my_list.pop(-1)	Remove an item
>>> my_list.insert(0,'!')	Insert an item
>>> my_list.sort()	Sort the list

Libraries

Import libraries
>>> import numpy
>>> import numpy as np
Selective import
>>> from math import pi

pandas Data analysis
Machine learning
NumPy Scientific computing
matplotlib 2D plotting

Install Python

anaconda Spyder Jupyter
Leading open data science platform Free IDE that is included Create and share powered by Python with Anaconda documents with live code, visualizations, text, ...

Numpy Arrays

Also see Lists

```
>>> my_list = [1, 2, 3, 4]  
>>> my_array = np.array(my_list)  
>>> my_2darray = np.array([(1,2,3), (4,5,6)])
```

Selecting Numpy Array Elements

Index starts at 0

Subset	Select item at index 1
>>> my_array[1]	
Slice	Select items at index 0 and 1
>>> my_array[0:2]	array([1, 2])
Subset 2D Numpy arrays	array([[1, 2], [3, 4]])
>>> my_2darray[:,0]	array([1, 4])
	my_2darray[rows, columns]

Numpy Array Operations

```
>>> my_array > 3  
array([False, False, False, True], dtype=bool)  
>>> my_array * 2  
array([2, 4, 6, 8])  
>>> my_array + np.array([5, 6, 7, 8])  
array([6, 8, 10, 12])
```

Numpy Array Functions

>>> my_array.shape	Get the dimensions of the array
>>> np.append(other_array)	Append items to an array
>>> np.insert(my_array, 1, 5)	Insert items in an array
>>> np.delete(my_array, [1])	Delete items in an array
>>> np.mean(my_array)	Mean of the array
>>> np.median(my_array)	Median of the array
>>> my_array.corrcoef()	Correlation coefficient
>>> np.std(my_array)	Standard deviation



Numpy Cheat Sheet



Pandas Cheat Sheet

Data Wrangling with pandas Cheat Sheet <http://pandas.pydata.org>

Syntax – Creating DataFrames

	a	b	c
1	4	7	10
2	5	8	11
3	6	9	12

```
df = pd.DataFrame(  
    {'a' : [4, 5, 6],  
     'b' : [7, 8, 9],  
     'c' : [10, 11, 12]},  
    index = [1, 2, 3])  
Specify values for each column.
```

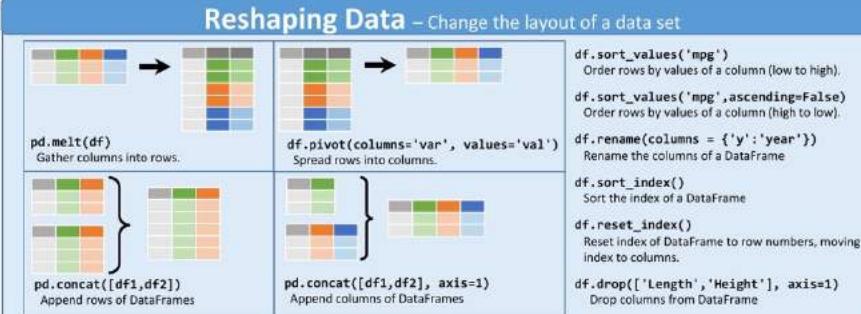
n	a	b	c
1	4	7	10
2	5	8	11
3	6	9	12

```
df = pd.DataFrame(  
    {'a' : [4, 5, 6],  
     'b' : [7, 8, 9],  
     'c' : [10, 11, 12]},  
    index = pd.MultiIndex.from_tuples(  
        [('d',1),('d',2),('e',2)],  
        names=['n', 'v']))  
Create DataFrame with a MultiIndex
```

Method Chaining

Most pandas methods return a DataFrame so that another pandas method can be applied to the result. This improves readability of code.

```
df = (pd.melt(df)  
      .rename(columns={  
          'variable' : 'var',  
          'value' : 'val'})  
      .query('val >= 200'))
```



Subset Observations (Rows)

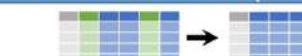


```
df.sample(frac=0.5) Randomly select fraction of rows.  
df.sample(n=10) Randomly select n rows.  
df.iloc[10:20] Select rows by position.  
df.head(n) Select first n rows.  
df.tail(n) Select last n rows.
```

Logic in Python (and pandas)		
<	Less than	!=
>	Greater than	df.column.isin(values)
==	Equals	Group membership
<=	Less than or equals	pd.isnull(obj)
>=	Greater than or equals	Is NaN
		pd.notnull(obj)
		Is not NaN
		pd.all()
		Logical and, or, not, xor, any, all

<http://pandas.pydata.org> | This cheat sheet inspired by RStudio Data Wrangling CheatSheet (<https://www.rstudio.com/wp-content/uploads/2015/07/data-wrangling-cheatsheet.pdf>) Written by Ivan Curic, [@ivancuric](#)

Subset Variables (Columns)



```
df.filter(regex='regex') Select columns whose name matches regular expression regex.
```

regex (Regular Expressions) Examples
'.' Matches strings containing a period ''
'Lengths' Matches strings ending with word 'Length'
'Sepal' Matches strings beginning with the word 'Sepal'
'^x[1-5]\$' Matches strings beginning with 'x' and ending with 1,2,3,4,5
'^(?i)Species\$' Matches strings except the string 'Species'

```
df.loc[:, 'x2':'x4'] Select all columns between x2 and x4 (inclusive).  
df.loc[:, [1,2,5]] Select columns in positions 1, 2 and 5 (first column is 0).  
df.loc[df['a'] > 10, ['a','c']] Select rows meeting logical condition, and only the specific columns .
```



Matplotlib Cheat Sheet

Python For Data Science Cheat Sheet

Matplotlib

Matplotlib is a Python 2D plotting library which produces publication-quality figures in a variety of hardcopy formats and interactive environments across platforms.



Prepare The Data

```
>>> import numpy as np
>>> x = np.linspace(0, 10, 100)
>>> y = np.cos(x)
>>> z = np.sin(x)
```

2D Data or Images

```
>>> data = 2 * np.random.random((10, 10))
>>> data2 = 3 * np.random.random((10, 10))
>>> Y, X = np.mgrid[-3:3:100j, -3:3:100j]
>>> V = np.sqrt(X**2 + Y**2)
>>> from matplotlib.cbook import get_sample_data
>>> img = np.loadtxt(get_sample_data('axes_grid/bivariate_normal.npy'))
```

Create Plot

```
>>> import matplotlib.pyplot as plt
```

Figure

```
>>> fig = plt.figure()
>>> fig2 = plt.figure(figsize=plt.figaspect(2.0))
```

Axes

All plotting is done with respect to an Axes. In most cases, a subplot will fit your needs. A subplot is an axes on a grid system.

```
>>> fig.add_axes()
>>> ax1 = fig.add_subplot(221) # row-col-num
>>> ax2 = fig.add_subplot(212)
>>> fig3, axes = plt.subplots(nrows=2, ncols=2)
>>> fig4, axes2 = plt.subplots(ncols=3)
```

Plotting Routines

1D Data

```
>>> lines = ax.plot(x, y)
>>> lines[0].set_color('red')
Draw points with lines or markers connecting them
Draw unconnected points, scaled or colored
```

```
>>> axes[0,0].bar([1,2,3],[3,4,5])
>>> axes[0,0].barh([0,1,2,3],[0,1,2,1])
>>> axes[1,0].axhline(0.45)
>>> axes[0,1].axvline(0.65)
>>> ax.hill(x,y,color='blue')
>>> ax.fill_between(x,y,color='yellow')
```

2D Data or Images

```
>>> fig, ax = plt.subplots()
>>> im = ax.imshow(data, interpolation='gist_stitch',
    interpolation='nearest',
    vmin=-2,
    vmax=2)
```

Colormapped or RGB arrays

Vector Fields

```
>>> axes[0,1].arrow(0,0,0.5,0.5)
>>> axes[1,1].quiver(y,x)
>>> axes[0,1].streamplot(X,Y,U,V)
```

Add an arrow to the axes
Plot a 2D field of arrows
Plot 2D vector fields

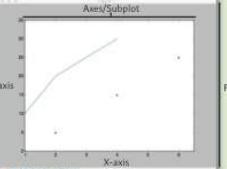
Data Distributions

```
>>> ax1.hist(y)
>>> ax3.boxplot(y)
>>> ax3.violinplot(z)
```

Plot a histogram
Make a box and whisker plot
Make a violin plot

Plot Anatomy & Workflow

Plot Anatomy



Figure

Workflow

The basic steps to creating plots with matplotlib are:

- 1 Prepare data
- 2 Create plot
- 3 Plot
- 4 Customize plot
- 5 Save plot
- 6 Show plot

```
>>> import matplotlib.pyplot as plt
>>> x = [1, 2, 3, 4, 5, 6, 7, 8, 9, 10]
>>> y = [1, 2, 3, 4, 5, 6, 7, 8, 9, 10]
>>> fig = plt.figure() Step 1
>>> ax = fig.add_subplot(111) Step 2
>>> ax.plot(x, y, color='lightblue', linewidth=3) Step 3
>>> ax.scatter([2, 4, 6], [5, 7, 25], color='darkgreen', marker='^') Step 4
>>> ax.set_xlim(1, 6.5) Step 5
>>> plt.savefig('foo.png') Step 6
>>> plt.show()
```

4) Customize Plot

Colors, Color Bars & Color Maps

```
>>> plt.plot(x, x, x**2, x, x**3)
>>> plt.plot(x, y, alpha= 0.4)
>>> ax.plot(x, y, c='k')
>>> fig.colorbar(im, orientation='horizontal')
>>> im = ax.imshow(img, cmap='seismic')
```

Markers

```
>>> fig, ax = plt.subplots()
>>> ax.scatter(x,y,marker="*")
>>> ax.plot(x,y,marker="o")
```

LineStyles

```
>>> plt.plot(x,y,linewidth=4.0)
>>> plt.plot(x,y,ls="solid")
>>> plt.plot(x,y,ls="dashed")
>>> plt.plot(x,y,ls="dashdot")
>>> plt.setp(lines,color='r',linewidth=4.0)
```

Text & Annotations

```
>>> ax.text(1, -2.1,
           'Example Graph',
           style='italic')
>>> ax.annotate("Sine", xy=(0, 0),
               xycoords='data',
               xytext=(1.5, 0),
               textcoords='data',
               arrowprops=dict(arrowsize=2,
                               connectionstyle="arc3",
```

MathText

```
>>> plt.title(r'$\sigma_i=15$', fontsize=20)
```

Limits, Legends & Layouts

Limits & Autoscaling

```
>>> ax.margins(0.0,0.0,y=0.1)
>>> ax.set_xlim(-1,1)
>>> ax.set_xlim(0,10.5), ylim=[-1.5,1.5])
>>> ax.set_xlim(0,10.5)
```

Legends

```
>>> ax.set_title("An Example Axes",
                 ylabel="Y-Axis",
                 xlabel="X-Axis")
>>> ax.legend(loc='best')
```

Ticks

```
>>> ax.xaxis.set_ticks(range(1,5),
                      ticklabels=[3,100,-12,"foo"])
>>> ax.tick_params(axis='y',
                  direction='inout',
                  length=10)
```

Subplot Spacing

```
>>> fig3.subplots_adjust(wspace=0.5,
                        hspace=0.3,
                        left=0.125,
                        right=0.9,
                        top=0.9,
                        bottom=0.1)
```

Axis Spines

```
>>> ax.spines['top'].set_visible(False)
>>> ax.spines['bottom'].set_position(('outward',10))
```

Add padding to a plot
Set the aspect ratio of the plot to 1
Set limits for x and y axes
Set limits for x-axis

Set a title and x and y-axis labels

No overlapping plot elements

Manually set x-ticks
Make y-ticks longer and go in and out

Adjust the spacing between subplots

Fit subplot(s) to the figure area

Make the top axis line for a plot invisible
Move the bottom axis line outward

5) Save Plot

Save figures

```
>>> plt.savefig('foo.png')
```

Save transparent figures

```
>>> plt.savefig('foo.png', transparent=True)
```

6) Show Plot

```
>>> plt.show()
```

Close & Clear

Close

```
>>> plt.close()
```

Clear

```
>>> plt.clf()
>>> plt.cla()
```

Clear an axis
Clear the entire figure
Close a window

DataCamp
Learn Python for Data Science Interactively!



Sscikit-Learn Cheat Sheet

Python For Data Science Cheat Sheet					
<h2>Scikit-Learn</h2> <p>Scikit-learn is an open source Python library that implements a range of machine learning, preprocessing, cross-validation and visualization algorithms using a unified interface.</p>	<h3>Create Your Model</h3> <h4>Supervised Learning Estimators</h4> <p>Linear Regression</p> <pre>>>> from sklearn.linear_model import LinearRegression >>> lr = LinearRegression(normalize=True)</pre> <p>Support Vector Machines (SVM)</p> <pre>>>> from sklearn.svm import SVC >>> svc = SVC(kernel='linear')</pre> <p>Naïve Bayes</p> <pre>>>> from sklearn.naive_bayes import GaussianNB >>> gnb = GaussianNB()</pre> <p>KNN</p> <pre>>>> from sklearn import neighbors >>> knn = neighbors.KNeighborsClassifier(n_neighbors=5)</pre> <h4>Unsupervised Learning Estimators</h4> <p>Principal Component Analysis (PCA)</p> <pre>>>> from sklearn.decomposition import PCA >>> pca = PCA(n_components=0.95)</pre> <p>K Means</p> <pre>>>> from sklearn.cluster import KMeans >>> kmeans = KMeans(n_clusters=3, random_state=0)</pre>				
<h3>A Basic Example</h3> <pre>>>> from sklearn import neighbors, datasets, preprocessing >>> from sklearn.model_selection import train_test_split >>> from sklearn.metrics import accuracy_score >>> iris = datasets.load_iris() >>> X = iris.data[:, 2:4] >>> y = iris.target >>> X_train, X_test, y_train, y_test = train_test_split(X, y, random_state=42) >>> scaler = preprocessing.StandardScaler().fit(X_train) >>> X_train = scaler.transform(X_train) >>> X_test = scaler.transform(X_test) >>> knn = neighbors.KNeighborsClassifier(n_neighbors=5) >>> knn.fit(X_train, y_train) >>> y_pred = knn.predict(X_test) >>> accuracy_score(y_test, y_pred)</pre>	<h3>Evaluate Your Model's Performance</h3> <h4>Classification Metrics</h4> <p>Accuracy Score</p> <pre>>>> knn.score(X_test, y_test) >>> from sklearn.metrics import accuracy_score >>> accuracy_score(y_test, y_pred)</pre> <p>Classification Report</p> <pre>>>> from sklearn.metrics import classification_report >>> print(classification_report(y_test, y_pred))</pre> <p>Confusion Matrix</p> <pre>>>> from sklearn.metrics import confusion_matrix >>> print(confusion_matrix(y_test, y_pred))</pre> <p>Estimator score method Metric scoring functions</p> <p>Precision, recall, f-score and support</p>				
<h3>Loading The Data</h3> <p>Also see NumPy & Pandas</p> <p>Your data needs to be numeric and stored as NumPy arrays or SciPy sparse matrices. Other types that are convertible to numeric arrays, such as Pandas DataFrame, are also acceptable.</p> <pre>>>> import numpy as np >>> X = np.array([[1, 2, 3], [4, 5, 6], [7, 8, 9]]) >>> y = np.array([0, 1, 0]) >>> X[X < 0.7] = 0</pre>	<h4>Model Fitting</h4> <table border="1"> <tr> <td> Supervised learning <pre>>>> lr.fit(X, y) >>> knn.fit(X_train, y_train) >>> svc.fit(X_train, y_train)</pre> </td><td> Fit the model to the data </td></tr> <tr> <td> Unsupervised Learning <pre>>>> kmeans.fit(X_train) >>> pca_model = pca.fit_transform(X_train)</pre> </td><td> Fit the model to the data Fit to data, then transform it </td></tr> </table>	Supervised learning <pre>>>> lr.fit(X, y) >>> knn.fit(X_train, y_train) >>> svc.fit(X_train, y_train)</pre>	Fit the model to the data	Unsupervised Learning <pre>>>> kmeans.fit(X_train) >>> pca_model = pca.fit_transform(X_train)</pre>	Fit the model to the data Fit to data, then transform it
Supervised learning <pre>>>> lr.fit(X, y) >>> knn.fit(X_train, y_train) >>> svc.fit(X_train, y_train)</pre>	Fit the model to the data				
Unsupervised Learning <pre>>>> kmeans.fit(X_train) >>> pca_model = pca.fit_transform(X_train)</pre>	Fit the model to the data Fit to data, then transform it				
<h3>Training And Test Data</h3> <pre>>>> from sklearn.model_selection import train_test_split >>> Xtrain, Xtest, y_train, y_test = train_test_split(X, >>> y, >>> random_state=0)</pre>	<h4>Prediction</h4> <table border="1"> <tr> <td> Supervised Estimators <pre>>>> y_pred = svc.predict(np.random.random((2,5))) >>> y_pred = lr.predict(X_test) >>> y_pred = knn.predict_proba(X_test)</pre> </td><td> Predict labels Predict labels Estimate probability of a label </td></tr> <tr> <td> Unsupervised Estimators <pre>>>> y_pred = kmeans.predict(X_test)</pre> </td><td> Predict labels in clustering algos </td></tr> </table>	Supervised Estimators <pre>>>> y_pred = svc.predict(np.random.random((2,5))) >>> y_pred = lr.predict(X_test) >>> y_pred = knn.predict_proba(X_test)</pre>	Predict labels Predict labels Estimate probability of a label	Unsupervised Estimators <pre>>>> y_pred = kmeans.predict(X_test)</pre>	Predict labels in clustering algos
Supervised Estimators <pre>>>> y_pred = svc.predict(np.random.random((2,5))) >>> y_pred = lr.predict(X_test) >>> y_pred = knn.predict_proba(X_test)</pre>	Predict labels Predict labels Estimate probability of a label				
Unsupervised Estimators <pre>>>> y_pred = kmeans.predict(X_test)</pre>	Predict labels in clustering algos				
<h3>Preprocessing The Data</h3> <h4>Standardization</h4> <pre>>>> from sklearn.preprocessing import StandardScaler >>> scaler = StandardScaler().fit(X_train) >>> standardized_X = scaler.transform(X_train) >>> standardized_X_test = scaler.transform(X_test)</pre>	<h4>Encoding Categorical Features</h4> <pre>>>> from sklearn.preprocessing import LabelEncoder >>> enc = LabelEncoder() >>> y = enc.fit_transform(y)</pre>				
<h4>Normalization</h4> <pre>>>> from sklearn.preprocessing import Normalizer >>> scaler = Normalizer().fit(X_train) >>> normalized_X = scaler.transform(X_train) >>> normalized_X_test = scaler.transform(X_test)</pre>	<h4>Imputing Missing Values</h4> <pre>>>> from sklearn.preprocessing import Imputer >>> imp = Imputer(missing_values=0, strategy='mean', axis=0) >>> imp.fit_transform(X_train)</pre>				
<h4>Binarization</h4> <pre>>>> from sklearn.preprocessing import Binarizer >>> binarizer = Binarizer(threshold=0.0).fit(X) >>> binary_X = binarizer.transform(X)</pre>	<h4>Generating Polynomial Features</h4> <pre>>>> from sklearn.preprocessing import PolynomialFeatures >>> poly = PolynomialFeatures(5) >>> poly.fit_transform(X)</pre>				
	<h4>Cross-Validation</h4> <pre>>>> from sklearn.cross_validation import cross_val_score >>> print(cross_val_score(knn, X_train, y_train, cv=4)) >>> print(cross_val_score(lr, X, y, cv=2))</pre>				
	<h3>Tune Your Model</h3> <h4>Grid Search</h4> <pre>>>> from sklearn.grid_search import GridSearchCV >>> params = {"n_neighbors": np.arange(1,3), >>> "weights": ["uniform", "distance"]} >>> grid = GridSearchCV(estimator=knn, >>> param_grid=params) >>> grid.fit(X_train, y_train) >>> print(grid.best_score_) >>> print(grid.best_estimator_.n_neighbors)</pre>				
	<h4>Randomized Parameter Optimization</h4> <pre>>>> from sklearn.grid_search import RandomizedSearchCV >>> params = {"n_neighbors": range(1,5), >>> "weights": ["uniform", "distance"]} >>> research = RandomizedSearchCV(estimator=knn, >>> param_distributions=params, >>> cv=4, >>> n_iter=8, >>> random_state=5} >>> research.fit(X_train, y_train) >>> print(research.best_score_)</pre>				
	<p>Learn Python for Data Science Interactively at www.DataCamp.com</p> 				



Scikit-learn Design

<https://arxiv.org/abs/1309.0238>

Scikit-Learn's API is remarkably well designed. These are the main design principles:

Consistency: All objects share a consistent and simple interface:

- **Estimators** Any object that can estimate some parameters based on a dataset is called an *estimator* (e.g., an imputer is an estimator). The estimation itself is performed by the `fit()` method, and it takes only a dataset as a parameter (or two for supervised learning algorithms; the second dataset contains the labels). Any other parameter needed to guide the estimation process is considered a hyperparameter (such as an imputer's strategy), and it must be set as an instance variable (generally via a constructor parameter).
- **Transformers** Some estimators (such as an imputer) can also transform a dataset; these are called *transformers*. Once again, the API is simple: the transformation is performed by the `transform()` method with the dataset to transform as a parameter. It returns the transformed dataset. This transformation generally relies on the learned parameters, as is the case for an imputer. All transformers also have a convenience method called `fit_transform()` that is equivalent to calling `fit()` and then `transform()` (but sometimes `fit_transform()` is optimized and runs much faster).
- **Predictors** Finally, some estimators, given a dataset, are capable of making predictions; they are called *predictors*. For example, the `LinearRegression` model in the previous chapter was a predictor: given a country's GDP per capita, it predicted life satisfaction. A predictor has a `predict()` method that takes a dataset of new instances and returns a dataset of corresponding predictions. It also has a `score()` method that measures the quality of the predictions, given a test set (and the corresponding labels, in the case of supervised learning algorithms).¹⁸

Inspection: All the estimator's hyperparameters are accessible directly via public instance variables (e.g., `imputer.strategy`), and all the estimator's learned parameters are accessible via public instance variables with an underscore suffix (e.g., `imputer.statistics_`).

Nonproliferation of classes Datasets are represented as NumPy arrays or SciPy sparse matrices, instead of homemade classes. Hyperparameters are just regular Python strings or numbers.

Composition Existing building blocks are reused as much as possible. For example, it is easy to create a Pipeline estimator from an arbitrary sequence of transformers followed by a final estimator, as we will see.

Sensible defaults Scikit-Learn provides reasonable default values for most parameters, making it easy to quickly create a baseline working system.



Where to go from here

1. Learn the language
2. Learn the tools
3. Work on projects!

Resources

- [Machine Learning Mastery](#)
- [Hands-On Machine Learning with Scikit-Learn and TensorFlow](#)
- [The Complete Python Course for Machine Learning Engineers](#)
- [Cheat Sheet of Machine Learning and Python \(and Math\) Cheat Sheets](#)



Questions

Hands-On Machine Learning with Scikit-Learn, Keras, and TensorFlow, 2nd Edition by Aurélien Géron, Publisher: O'Reilly Media, Inc. Release Date: September 201

1. How would you define Machine Learning?
2. Can you name four types of problems where it shines?
3. What is a labeled training set?
4. What are the two most common supervised tasks?
5. Can you name four common unsupervised tasks?
6. What type of Machine Learning algorithm would you use to allow a robot to walk in various unknown terrains?
7. What type of algorithm would you use to segment your customers into multiple groups?
8. Would you frame the problem of spam detection as a supervised learning problem or an unsupervised learning problem?
9. What is an online learning system?
10. What is out-of-core learning?
11. What type of learning algorithm relies on a similarity measure to make predictions?
12. What is the difference between a model parameter and a learning algorithm's hyperparameter?
13. What do model-based learning algorithms search for? What is the most common strategy they use to succeed? How do they make predictions?
14. Can you name four of the main challenges in Machine Learning?
15. If your model performs great on the training data but generalizes poorly to new instances, what is happening? Can you name three possible solutions?
16. What is a test set, and why would you want to use it?
17. What is the purpose of a validation set?
18. What is the train-dev set, when do you need it, and how do you use it?
19. What can go wrong if you tune hyperparameters using the test set?



Answers

1. Machine Learning is about building systems that can learn from data. Learning means getting better at some task, given some performance measure.
2. Machine Learning is great for complex problems for which we have no algorithmic solution, to replace long lists of hand-tuned rules, to build systems that adapt to fluctuating environments, and finally to help humans learn (e.g., data mining).
3. A labeled training set is a training set that contains the desired solution (a.k.a. a label) for each instance.
4. The two most common supervised tasks are regression and classification.
5. Common unsupervised tasks include clustering, visualization, dimensionality reduction, and association rule learning.
6. Reinforcement Learning is likely to perform best if we want a robot to learn to walk in various unknown terrains, since this is typically the type of problem that Reinforcement Learning tackles. It might be possible to express the problem as a supervised or semisupervised learning problem, but it would be less natural.
7. If you don't know how to define the groups, then you can use a clustering algorithm (unsupervised learning) to segment your customers into clusters of similar customers. However, if you know what groups you would like to have, then you can feed many examples of each group to a classification algorithm (supervised learning), and it will classify all your customers into these groups.
8. Spam detection is a typical supervised learning problem: the algorithm is fed many emails along with their labels (spam or not spam).
9. An online learning system can learn incrementally, as opposed to a batch learning system. This makes it capable of adapting rapidly to both changing data and autonomous systems, and of training on very large quantities of data.
10. Out-of-core algorithms can handle vast quantities of data that cannot fit in a computer's main memory. An out-of-core learning algorithm chops the data into mini-batches and uses online learning techniques to learn from these mini-batches.
11. An instance-based learning system learns the training data by heart; then, when given a new instance, it uses a similarity measure to find the most similar learned instances and uses them to make predictions.
12. A model has one or more model parameters that determine what it will predict given a new instance (e.g., the slope of a linear model). A learning algorithm tries to find optimal values for these parameters such that the model generalizes well to new instances. A hyperparameter is a parameter of the learning algorithm itself, not of the model (e.g., the amount of regularization to apply).
13. Model-based learning algorithms search for an optimal value for the model parameters such that the model will generalize well to new instances. We usually train such systems by minimizing a cost function that measures how bad the system is at making predictions on the training data, plus a penalty for model complexity if the model is regularized. To make predictions, we feed the new instance's features into the model's prediction function, using the parameter values found by the learning algorithm.
14. Some of the main challenges in Machine Learning are the lack of data, poor data quality, nonrepresentative data, uninformative features, excessively simple models that underfit the training data, and excessively complex models that overfit the data.
15. If a model performs great on the training data but generalizes poorly to new instances, the model is likely overfitting the training data (or we got extremely lucky on the training data). Possible solutions to overfitting are getting more data, simplifying the model (selecting a simpler algorithm, reducing the number of parameters or features used, or regularizing the model), or reducing the noise in the training data.
16. A test set is used to estimate the generalization error that a model will make on new instances, before the model is launched in production.
17. A validation set is used to compare models. It makes it possible to select the best model and tune the hyperparameters.
18. The train-dev set is used when there is a risk of mismatch between the training data and the data used in the validation and test datasets (which should always be as close as possible to the data used once the model is in production). The train-dev set is a part of the training set that's held out (the model is not trained on it). The model is trained on the rest of the training set, and evaluated on both the train-dev set and the validation set. If the model performs well on the training set but not on the train-dev set, then the model is likely overfitting the training set. If it performs well on both the training set and the train-dev set, but not on the validation set, then there is probably a significant data mismatch between the training data and the validation + test data, and you should try to improve the training data to make it look more like the validation + test data.
19. If you tune hyperparameters using the test set, you risk overfitting the test set, and the generalization error you measure will be optimistic (you may launch a model that performs worse than you expect).





CONNECTING WORLDS



Kontakt



Michael Aydinbas
Data Scientist

@ michael.aydinbas@EXXETA.com
m +49 174 99 50 913

