# Machine Learning

*Zusammenfassung* / Joel von Rotz / :octocat: [Quelldateien]

## Begriffe

**Artificial Intelligence** Software, die intelligentes Verhalten imitiert
↪ **Machine Learning** System, das von Erfahrung/Daten lernt
　　↪ **Deep Learning** basiert auf geschichteten neuralen Netzwerken

**Supervised Learning** Lernt anhand labeled Daten → predicts unseen data, *Dimensionality Reduction*
**Unsupervised Learning** unlabeled Daten → Daten-Struktur lernen und Kunden in Zielgruppen teilen (*Recommender Systems*)
**Semi-Supervised Learning** Mix Sup.+Unsup. → wenig *Labeled* & viele *Unlabeled*
**Reinforcement Learning** Belohnungsfunktion → sucht ideales Verhalten für maximale Belohnung

**Regression Problem** Label (vorherzusagendes Attribut) ist kontinuierlich
**Classification Problem** Label ist kategorial (binary -> Bsp. Ja, nein)
**Hyperparameters** vom Programmierer ausgewählte Parameter (Bsp. K bei K-NN)

## Data Quality Assessment

● Identify data sources and trustworthiness ● Check data ranges (e.g. negative salary) ● Interpret statistical key figures ● Validate plausibility of variable correlations ● Measure data redundancy ● Check for anomalies in sytax and semantics ● Explore Null values and duplicate data records

## Data Cleaning

Data quality can be improved (although document changes, inform,…):

● identify & remove duplicate data records ● replace null values (e.g. median or mean) or delete them (requires large data) ● investigate the origins of quality issues
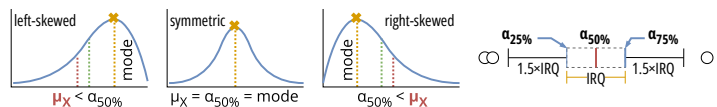
## Feat. Eng. - Vector Space Model

Data set consists **only** of numerical data (e.g. color → black, red, …). Categorical data is transformed into numerical data

## Statistics

**mean** $\mu$ easy calculate ; **median** robust against outliers ; $Var_{smpl}. \frac{1}{n-1}$ ; $Var_{pop}. \frac{1}{n}$

Var
$\sigma^2 = \frac{1}{n-1}\sum_{i=1}^{n}(x_i - \mu_X)^2$ ; $Cov(X,Y) = \frac{1}{n-1}\sum_{i=1}^{n}(x_i - \mu_X)(y_i - \mu_Y)$



**Pearson Correlation** $\rho(X,Y) = \frac{Cov(X,Y)}{\sigma_X \sigma_Y}$ normalizes ranges: -1=anti-corr., +1=corr

## Similariy Measures

**Euclidean Distance** $\sqrt{\sum_{i=1}^{n}(x_i - y_i)}$ (distance)

**Cosine Distance** $1 - CosineSimil. = 1 - \frac{\sum_{i=1}^{n} x_i \cdot y_i}{\sqrt{\sum_{i=1}^{n} x_i^2}\sqrt{\sum_{i=1}^{n} y_i^2}}$ (distance + angle)
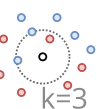
**Manhattan Distance** $\sum_i |x_i - y_i|$ (along a grid)

**Levenshtein/Edit Distance** for strings: $+1$ on delete, add ; $+2$ on change char

**Jacard Similarity**: $\frac{|X \cap^{and} Y|}{|X \cup^{or} Y|} \rightarrow [0,1]$; 1=same

## k-Nearest-Neighbor Classification

Data takes value of majority of $k$ (hyperparameter) neighbors + (optional) weight based on $1/d$ (distance $d$).
━ very slow with lots of data and most comput. is done in classification stage ➕ least data hungry algo, good baseline for other classifier


k=3

---

## Normalization

**Min-Max**: $\frac{x - \min X}{\max X - \min X} \rightarrow [0,1]$ ➕ interpretation in %, no negatives ━ no supervised

**Z-Score**: $\frac{x - \mu_X}{\sigma_X} \rightarrow \begin{matrix} \mu_{X'} = 0, \\ \sigma_{X'} = 1 \end{matrix}$ ➕ for (un)supervised ━ no direct interpr., negatives

## Regression Hypothesis & Evaluation Metrics

Goal is to fit a line or curve-function to a number of points. $y_i$ ground truth, $f_i$ prediction ; **MAE** (Mean Absolute Error), **MAPE** (Mean Absolute Percentage Error) human readable, **MSE** (Mean Squared Error) optimized for processing

$$\frac{1}{m}\left[\begin{matrix} MAE & MAPE & MSE \\ \sum_{i=1}^{m}|y_i - f_i| & \sum_{i=1}^{m}|\frac{y_i - f_i}{y_i}| & \sum_{i=1}^{m}(y_i - f_i)^2 \end{matrix}\right] \Big| \begin{matrix} R^2 = \rho^2 \\ 1 - \frac{\sum_{i=1}^{m}(y_i - f_i)^2}{\sum_{i=1}^{m}(y_i - \mu_Y)^2} \end{matrix}$$

$R^2 \le 1$ ; $R^2 < 0$ worse fit than horizontal hyperplane (mean) ; $R^2 > 1$ model performs worse than mean the fraction becomes

## Model Evaluation Workflow

**Workflow**: ① Mix data ② Split into training & test data ③ train ④ evaluate! NEVER reuse test set!

**Model Selection**: ① 60% Training, 20% Validation, 20% Test ② Train different models ③ Check via Validation and select model ④ Use selection on test data!

**K-Fold Cross Valid.**: not enough data → ① Split 80% data into $k$ parts ② Round 1 uses first part as **validation** ③ average all accuracy to get final accuracy

## Fitting Data

● Underfitting (too loose), Just Right (not perfect fit, but distribution is learnt), Overfitting (too perfect, no distribution)
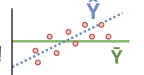
**Regression**: X regressor, $\varepsilon$ error $y_{real} - y_{lin}$ : $Y = X\theta + \varepsilon$ ($y_j = \theta_0 + \theta_1 \cdot x_j + \varepsilon_j$)

**Cost Function**: used to minimize MSE ; $J(\theta_0, \theta_1) = \frac{1}{N}\sum_{j=N}^{N} \varepsilon_j^2$

**Ordinary Least Squares** (OLS): $\frac{\partial J(\theta)}{\partial \theta} = 0 \rightarrow \theta_{opt} = (X^T X)^{-1} X^T y$ (**Gradient Descent** iteratively: ① random $\theta$ ② small steps down cost surface ● direction = negative gradient ; $\theta_{new} = \theta_{old} - \alpha \frac{\partial J(\theta)}{\partial \theta}$)

**Regression Performance**: $R^2 = 1 - \frac{\sum_j \varepsilon_j^2}{\sum_j \varepsilon_\delta^2} = 1 - \frac{\sum_j (y_j - \hat{Y})^2}{\sum_j (y_j - \bar{Y})^2}$

Plot data, as $R^2$ does not show the distribution (*Anscome's quartet*)!



To penalize large parameters and stop overfitting, **LASSO** and **Ridge** can be used. $J(\theta) = \underbrace{\frac{1}{N}\sum_{j=1}^{N}(y_j - \hat{y}_j)^2}_{OLS} + \underbrace{\lambda \sum_{k=1}^{M}\theta_k^2}_{Ridge} \Big/ \underbrace{\lambda \sum_{k=1}^{M}|\theta_k|}_{LASSO}$

**Polynomial Regression**: $h(\theta, X) = \theta_0 + \theta_1 X + \theta_2 X^2 + \cdots \Rightarrow \theta_0 + \theta_1 X_1 + \theta_2 X_2 + \cdots$

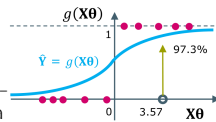## Classification

**Classification** predicts categories (*decision boundry*) *vs.* **Regression** predicts values (*data relationship*)
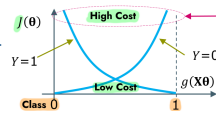
### Logistic Regression

$$\hat{Y} = P(Y = 1|X) = h_\theta(X) = g(X\theta) = \frac{1}{1 + e^{-X\theta}}$$

Instead of either exactly $1$ or $0$, a percentage $[0,1]$ is given. No OLS-form possible, good parameters calculated via *Gradient Descent*.



### Cross Entropy Cost Function

Classifying for one class **must** penalize the other!



---

## One Versus The Rest

Train model based on a specific class, against others. It is either in that class (1) or not (0) → **One vs. All** classification.

## Performance Analysis

$\text{Sensitivity} / \text{Recall} = \frac{TP}{TP + FN}$　$\text{Specifity} = \frac{TN}{TN + FP}$

$\text{Precision} = \frac{TP}{TP + FP}$　$\text{F1 Score} = \frac{2 \cdot prec \cdot reca}{prec + reca}$

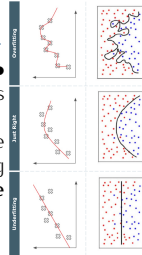| Conf. Matrix | | Predict | |
|---|---|---|---|
| | | No | Yes |
| **Actual** | No | TN | FP |
| | Yes | FN | TP |

Accuracy $= (TP + TN)/Total$ ; lower F1 indicates some kind of skewed data (unbalance).

## Dealing with Bias & Variance

**Dealing with bias is easier** (underfitting): ● larger set of feature ● different features ● feature engineering ● more complex algorithms (more data is not necessarily the answer)
**Dealing with variance** (overfitting): ● reduce features ● get more training data ● stop training early to avoid overfitting ● data cleaning
**Bias** means we are systematically off by a fixed amount. **Variance** means we are systematically off in a scattered way
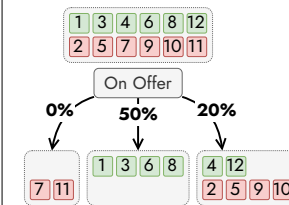


**Receiver Operating Characterstic (ROC)**:
　　● to compare different models
　　● largest area under curve is best



## Decision Trees

**Gini Impurity** measures the likelihood of an incorrect classification of a random data record. *Impure* features produce no significant information. Cont. Var. ① Sort Table ② Calculate avg. of adjacent values ($10 \rightarrow 15 \leftarrow 20$) ③ Averages are splitting criteria (histogramm binning reduces splits)



| Label | p(No) | p(Yes) | 1-p(No)²-p(Yes)² |
|---|---|---|---|
| 0% | 1 | 0 | 0 |
| 20% | 2/3 | 1/3 | 4/9 |
| 50% | 0 | 1 | 0 |

Gini Impurity of entire feature:
$2/12 \times 0 + 6/12 \times 4/9 + 4/12 \times 0 = 2/9$

**Regression Trees** ● Are grown like decision trees for classification but with different splitting criteria ● Instead of Gini impurity, **minimize variance of values in the same subset** ● Finally, predict average value of all instances in a leaf node

## Tree Construction Rules

① **Only** positive or negative nodes left → stop branch & assign corresponding decision as leaf node ② **Some** positive and negative remain → apply another feature ③ **No** instances left → look at parent node and decide according the more frequent label ④ Instances left but **no** features → training data is no good or contains hidden features ; decide on the more frequent label
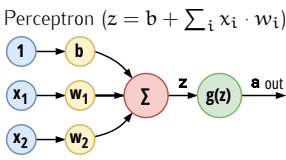
### Pros & Cons Decision & Regression Trees

➕ readable, can both candle numeric & categorial, (almost) no data preparation, perform well on larger datasets ━ not as accurate, danger of over-complex trees, very sensitive to data quality

**Random Forests** (collection of decision trees): Utilize **random** (row)sampling with replacement and **random** (column) feature selection (without replacement) to be less sensitive & no overfitting.

**Ensembles** complex ML model built from simpler models. *Building Bricks* → **Bragging** computes weighted average of votes from simple models, **Boosting** improves simple model by correcting its errors ($F_m(x) + h_m(x)$, e.g. Gradient Boosting).

# Neural Networks

**Input** $x_i$ from previous 'axons', contain some kind of information / **Weights** $b, w_i$ define the inputs' influence / **Sum** $\sum$ sum of the weighted inputs / **Activation** $g(z)$ can be **hard threshold** ($a = g(z) = 1$ if $\sum \geq 0$ else 0), **logistic regression** ($g(z) = \frac{1}{1+\exp(-z)}$)

Perceptron ($z = b + \sum_i x_i \cdot w_i$)

**Feed Forward** data moves one way (left → right)

**Hidden Layers**: allows for new features $a$ to be added

**Layer Count** $c_L = c_{hidden} + c_{output}$ (input is not counted)

| Feed Forward → | Back Propagation ← | Gradient Descent ↓ |
|---|---|---|
| • A neural net is the connection of many artificial neurons in parallel and series | • At the output of the ANN we get the estimate $\hat{y}$ of $y$ | • Back-propagation gives the gradient of $J$ with respect to weights **W** and biases **b** |
| • Each layer of a feed-forward ANN is defined by its weights **W** and biases **b** | • Because we start with random **W** and **b**, the estimate $\hat{y}$ is poor | • The gradient is a vector which points steeply up the error surface |
| | • Define the error $J$ as a measure of distance between $\hat{y}$ and $y$ | • So take a small step $\alpha$ in the opposite direction, reducing $J$ |
| • Given data with labels $y$ and features $x$, we feed $x$ forward through the ANN | • We seek the best **W** & **b** to minimize $J$. This occurs when the gradient of $J$ is small or zero | • "Take a step" means making small changes to **W** & **b** |
| | • So we need the gradient of $J$ | • Changes are defined exactly by $-\alpha$ times the gradient |
| | • Compute the gradient iteratively at each layer: more efficient! | • This is gradient descent by feed forward & back-prop using labelled data |
| | • We introduce deltas $\delta$ at each layer to simplify computation | |

## Activation Function

**ReLU** (rectified linear unit ;good for hidden) $\max(0, x)$ brings better perfomance and alleviates the vanishing gradients problem (flat line activation) → introduces **dying units** (bad for gradient training)
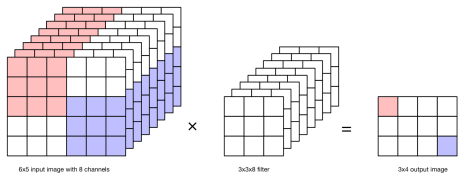
**Leaky ReLU** $\max(0.1x, x)$ introduces non-zero for negative values and fixes vanishing gradients + dying unit → results aren't always consistent

**Soft-Max** (good for output) $\text{softmax}_j = \frac{\exp z_j}{\sum_{k=1}^{K} \exp(z_k)}$ is a generalization of the logistic function to **multiple dimensions** and can be interpreted as probabilites (normalization).

## Convolution Neural Network

Features are extracted from images using **convolution** (Faltung). **Weights** are used as Masks/Filters on these images to extract certain features such as edges or colour. The final layer mostly uses **softmax**!

6x5 input image with 8 channels × 3x3x8 filter = 3x4 output image

**Pooling** applies operations such as max, min or avg to a pixel neighbour and the neighbourhood size is determined by the pool size parameter.

**Parameter Counting** Example on $20 \times 20$-RGB-Image ⓛ1 _2 Filters of $3 \times 3$_ (($3 \cdot 3 \cdot 3 + 1) \cdot \underline{2} = 56$ results $18 \times 18 \times 2$ Layers) ⓛ2 _3 Filters of $3 \times 3$ with 2 channels each_ (($3 \cdot 3 \cdot 2 + 1) \cdot \underline{3} = 57$ results $16 \times 16 \times 3$ Layers) ⓛ3 _Flatten result to a single column of neurons_ ($\underline{1} \cdot 16 \cdot 16 \cdot 3 = 1 \times 768$) ⓛ4 _Fully connected (aka. dense) layer and 2 output channels_ (#$\texttt{param} = 768 \cdot 2 + 2 = 1538$ (2 biases! & 2 neurons/output shape))

## Natural Language Processing

**Syntax** compare spelling or words / **Similarity** Levenshtein Distance

**Word Relatedness** Words are _related..._ • when they appear in the **same document**

(synonyms not related!!) • when they relate to the **same topic** (requires a lot of space) • when they can be replaced by each other

**Word Similarity** Words are _similar_ when they appear in the **same context**

**Term Frequency (Inverse Doucment Frequency)** $\text{TF-IDF}(x, y) = \text{TF}(x, y) \cdot \log\left(\frac{N}{\text{DF}(x)}\right)$ / $N$ number of documents in the corpus, TF counts the number of occurences of $x$ in $y$, DF counts the number of documents that contain $x$. _QUETZALCOATLUS_ has high score, due to being used a lot in a **single** document, compared to the low scoring word _AND_ (used everywhere).
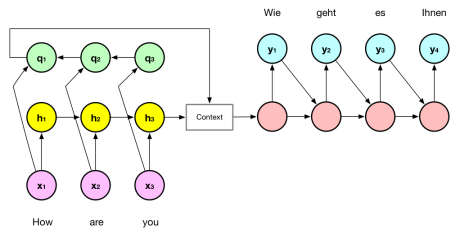
### Recurrent Neural Nets (RNN)

$h_t = (W_{xh} \cdot x_t + W_{hh} \cdot h_{t-1} + b_h) \quad y_t = \varphi(W_{hy} \cdot h_t + b_y)$

RNN are made for the supervised analysis of arbitrary time-series data (instead of fixed length, various sized strings can be analyzed) and **do not parallelize**! LSTM & GRU help with alleviating vanishing gradients to remember more of the context (Training an RNN by multiplying gradients backwards in time).

**GRU** (⋉) use hidden states as memory (faster) & **Long Short-Term Memory Model** have an additional memory state (better performance)

**Bidirectional RNN** reads left and right → produces together output

**Encoder-Decoder Architecture**: usually RNNs have one input and output, which makes it unusuable for i.e. translators. To allow in- and outputs of various lengths, en- & decoder architecture is used.

### Attention Mechanism

Mimics the behaviour of knowing which page the content can be found in. Words are assigned an _importance_ value which is used by the next generation for the output words.

In ⇨ Encoder ▷ Attention ▷ Decoder ⇨ **Out**

## k-Means Clustering

① Input $k$ Clusters and data points $x_n$ ② Randomly choose k clusters centers $\mu_k$ ③ Repeat until convergence (stability): (a) assign data points to nearest cluster center (b) update cluster center with mean of assigned datapoints

Two data points can meet in a cluster and separate again in the next step! → pattern always changes

### Clustering Distortion

**Total distortion** sum of squared distances of points and cluster center • **Average Distorion** average of total distortion (additional $1/n$) and allows for comparison with other data sets

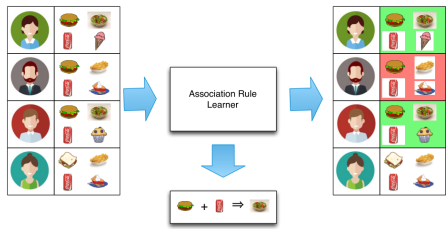**k-Means** only approximates optimal solutions and converges either to global **OR** local minimum

**Elbow Methods** increasing the amount of clusters leads to an 'elbow' in the distance-nCluster diagram → use cluster count at the 'elbow'

## Agglomerative Clustering

Outputs an hierarchy (_dendrogram_) that are informative and human interpretable. It does not scale with larger data, but cluster stays together forever! ① Initially each datapoint is a separate cluster ② While stop condit. (on no change) is not met: (a) calculate distance between pairs of cluster (b) merge pairs with shortest distance

## Association Rules

Implication ($X \rightarrow Y$): the higher confidence, the more likely $Y$ is also present in transaction with $X$

**Support** (Interestingness): how frequently an item set occurs in the data $\text{support}(\{i_1, \ldots, i_n\}) = \frac{\#\text{purchases of } \{i_1, \ldots, i_n\}}{\#\text{transactions}}$

**Support of an Association** $\text{support}(X \rightarrow Y) = \text{support}(Y \rightarrow X) = \text{support}(X \cup Y)$

**Confidence** (Trustworthiness): how frequently items in Y appear in transactions that contain X $\text{confidence}(X \rightarrow Y) = \frac{\text{support}(X \cup Y)}{\text{support}(X)}$

**BOTH Support** and **Confidence** need to be considered. ▬ Only Antecedent $X$ and co-occurence $X \cup Y$ is taken into account, not $Y$; sup. & conf. cannot filter coincidental occurrences out.

**Lift** (Association Strength): how many times more often X and Y show up together than expected $\text{lift}(X \rightarrow Y) = \frac{\text{support}(X \cup Y)}{\text{support}(X) \cdot \text{support}(Y)}$ ($= 1$ independent; $< 1$ X & Y appear less often together; $> 1$ X & Y appear often together)

### Apriori Algorithm

① Generate frequent item sets → ① frequent items leads subsets being frequent too ② infrequent leads supersets infrequent too. (shape **lattice of subsets**: first layer single node, second double connection,…)

② Extract rules from frequent item satisfying the confidence threshold → ① when an item violates the confidence bound all subsets can be eliminated. (shape **lattice of rules**)

### Recommender System - Non-Personalized

➕ simple, fast & cheap ; Associations are context aware ; Associations can be tailored to specific businesses (who viewed, who bought…) ; No new user cold start problem ▬ non-personalized

### Recommender System - Personalized

$P_{uj} = \frac{\sum_{i \in N_j} \text{similar}(j, i) \cdot r_{ui}}{\sum_{i \in N_j} \text{similar}(j, i)} \qquad P_{uj} = \bar{r}_u + \frac{\sum_{i \in N_u} \text{similar}(u, i) \cdot (r_{ij} - \bar{r}_i)}{\sum_{i \in N_u} \text{similar}(u, i)}$

N rated items (hyperparameter) ; $P_{uj}$ how much does user $u$ like item $j$ ; $N_j$ Neighbours of item $j$ ; $r_{ui}$ rating that user $u$ gave item $i$ ; $N_u$ Neighbors of user $u$ ; $r_{ij}$ from neighbor $i$ to item $j$

**Personalized Content-Based Recommendations** (left equ.): ➕ user independance, no new item cold start problem ▬ limited content analysis, overspecialization (bubble esque), severe new user cold start problem

**User-to-User Collaborative Filtering** (right equ.): ➕ cross-category recommendations ▬ many participants required, system must first learn user + item rating (cold start problem), gray sheep problem (unusual taste results in low accuracy)

### Hybridization

Combination of different recommender systems! • **Weighted Approach** item score is weighted sum score from different recommenders • **Mixed Approach** results of different recommender are presented • **Cascade Approach** one recommender refines the result of other method • **Switching Approach** use different methods ta different places (online shops)