

- X: observe clouds (0=no, 1=small, 2=big)
- Y: observe rain (0=no, 1=light, 2=moderate, 3=heavy)

Conditional (Marginal) Probability

- One variable is no longer random
- X is observed, its value is fixed
- Calculate the probabilities of Y given X: $P(Y|X)$

$$P(X, Y) = P(X|Y) * P(Y)$$

$$P(X, Y) = P(Y|X) * P(X)$$

$$P(Y|X) = \frac{P(X, Y)}{P(X)}$$

Bayes Rule

$$P(X|Y) * P(Y) = P(Y|X) * P(X)$$

Therefore:

$$P(Y|X) = \frac{P(X|Y) * P(Y)}{P(X)}$$

Derive

- Joint to Marginal: Ok
- Marginal to Joint: Not possible

4 Python

Chani alles

5 Data Visualization

- See trends, clusters and patterns in data
- Difficult to see in raw data
- Detect outliers and unusual groups
- Validate Hypothesis/Conjecture/Theory

Important in a Plot:

- X-Axis / Y-Axis
- Title
- Scale
- Dimensionality of the data 2D / 3D

5.1 Data Analysis Libraries

5.1.1 NumPy

- Package for scientific computing in Python
- Multidimensional array object
- Routines for fast array operations (sorting, selecting, FFT, linalg, etc)

5.1.2 pandas

- Built on top of NumPy
- Routines for accessing tabular data from files (.csv, xls, etc.)
- Supports 2-dimensional data (dataframe and series)
- Dataframes are something like database tables

5.1.3 Matplotlib

- Library for visualizing data
- Bargraphs, Histograms, Piecharts, Scatter plots, lines, boxplots, heatmaps, etc.

5.1.4 Seaborn

- Extension of Matplotlib, NumPy and pandas
- More user friendly
- Plots are aesthetically better

5.1.5 Chart types

Line Plots

- Bivariate, Continous
- Recognizes trend (pattern of change)

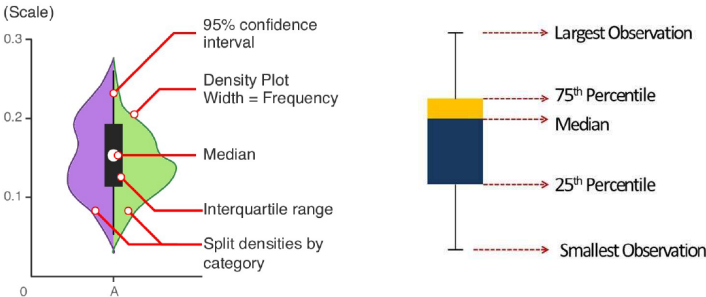
Bar Chart

- Used for categorical data
- Counting based on each category

Histogram

- Represents the empirical distribution of a variable
- Automatically creates bins (interval) along the range of values
- Shows vertical bars to indicate the number of observations / bin

Descriptive Statistics: Box Plots and Violin Plots



Scatter Plot

- Relationship between continous variables

- Helps to get an idea of the degree of correlation between variables

6 Regression

6.1 What is a model?

In ML, we use the term **model** for any mathematical function that explains the data:

$$y_i = f(x_i)$$

$$y_i = f(x_i) + \epsilon_i$$

where ϵ_i is unexplained noise. It is often assumed that ϵ_i follows a normal distribution.

Instead of approximating y_i , we calculate an **estimate** \hat{y}_i (y hat) of the usually unknown y_i :

$$\hat{y}_i = f(x)$$

6.1.1 Linear Regression

- Only considers a linear relationship between input and output
- In the simplest case, x and y are scalars and the linear model therefore has only two free parameters
- The goal is to identify a (slope) and b (intercept) for which the linear model best explains the data

$$\hat{y}_i = ax_i + b$$

6.1.2 Mean Squared Error (MSE)

- Loss we want to minimize
- Usually divided by 2

$$\hat{y}_i = ax_i + b$$

$$e_i = y_i - \hat{y}_i$$

The difference e_i , called residual

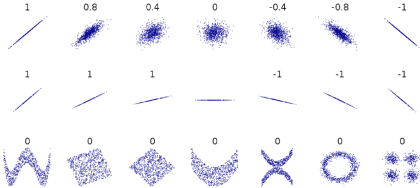
$$E = \frac{1}{2N} * \sum_{i=1}^N e_i^2$$

$$E = \frac{1}{2N} * \sum_{i=1}^N (\hat{y}_i - (a * x_i + b))^2$$

6.1.3 Correlation and Causality

- Correlation is not causality
- Correlation refers to the degree to which a pair of variables are linearly related
- Linear regression is a tool to detect correlations between two or more variables
- Correlation can be quantified using the Pearson correlation coefficient

Pearson Correlation Coefficient:



7 Optimization

- Training or learning in AI often suggests an algorithm performing some sort of optimization
- It is the problem of finding a set of inputs to an objective function that results in a maximum or minimum function evaluation
- In our examples the objective is to minimize the loss function

7.1 Gradient Descent

- Iterative Method
- Each iteration, the model parameters are updated such as that the Loss (MSE) is reduced

Parametervektor zum Zeitpunkt $t + 1$

$$\begin{bmatrix} a \\ b \end{bmatrix}_{t+1} = \begin{bmatrix} a \\ b \end{bmatrix}_t - \alpha \begin{bmatrix} \frac{\partial E}{\partial a} \\ \frac{\partial E}{\partial b} \end{bmatrix} \mid \begin{bmatrix} a \\ b \end{bmatrix}_t$$

$$\frac{\partial E}{\partial a} = \frac{1}{N} \sum_{i=1}^N (y_i - (ax_i + b)) * (-x)$$

$$\frac{\partial E}{\partial b} = \frac{1}{N} \sum_{i=1}^N (y_i - (ax_i + b)) * (-1)$$

7.2 Stochastic Gradient Descent (SGD)

- At each iteration, the gradient is calculated on a (randomly selected) subset of the data
- For a fixed learning rate, SGD does not converge
- Each iteration, the entire parameter vector gets optimized
- Only finds the local Minimum

7.2.1 Annealed SGD

- The learning rate alpha is reduced over time
- This is called (simulated) annealing
- There are different options (called schedules) how to reduce alpha over time

7.2.2 General remarks on SGD

- Gradient-based methods only work if we can express a Loss function as a differentiable function
- SGD is dealing with only a single datum at each iteration. This is very inefficient and rarely used.
- Batch- or mini-batch gradient-descent is usually used

8 Generalization & Regularization

8.1 Overfitting

- A model that perfectly fits the data does not have to be perfect
- In-Sample Error (Trainig error) was minimized (MSE = 0)
- Out-of-sample Error (Generalization Error, Test Error) is the MSE of new Data
- A good model has a low Generalization Error
- Overfitting happens if the MSE of Training Error is small thanks to a complex model but the Generalization Error is large

8.2 Underfitting

- Using a too simple model
- In-Sample Error is large
- Generalization Error is large

8.3 Training-Set, Test-Set, Model Evaluation

- The Generalization Error can't be calculated
- But Estimated
- Split the data into 2 sets
 - Training-Set (80% of data)
 - Test-Set (20% of data)

Training:

- Fit the model to the training set
- This minimizes the in-sample error

Evaluating

- Using the Test-Set
- Produces the Test-Error
- This is an estimate of the Generalization Error

8.4 Bias-Variance Trade-off

Variance: Difference of fits between data sets.

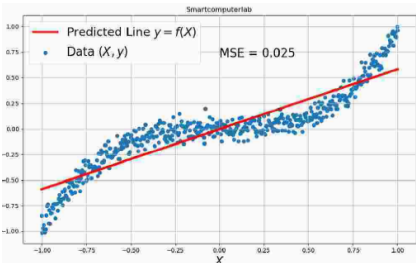
Bias: Results that are systematically prejudiced due to faulty assumptions.

High Bias

- A too simple model for the given data

Low Variance

- The model is relatively stable
- Very similar model if trained with new data

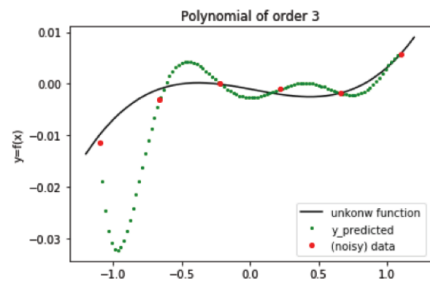


Low Bias

- A more complex model can better explain the data

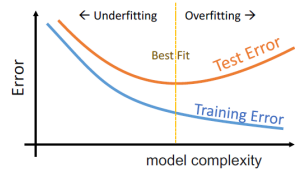
High Variance

- Given a new datapoint, the MSE can be very large
- For a different set with more datapoints, the model may be very different



8.4.1 Trade-off

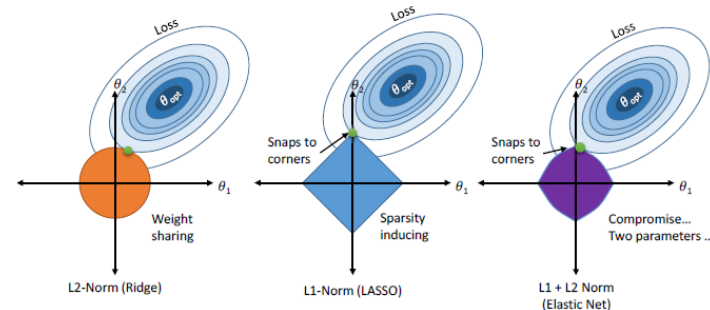
- Higher bias implies lower variance
- Lower bias implies higher variance
- In practice, all we want is low variance
- The model can only be as complex as the data permits
- You have to find an optimal balance between bias and variance



8.5 Regularization

- Technique to control the model complexity
 - Add a penalty term to the Loss
 - More complex models get a higher penalty
 - Add a constrain to the optimization process
 - *regularized loss = MSE + λ model-complexity*

$$\sum_{i=1}^n (y_i - \sum_{j=1}^p x_{ij} \beta_j)^2 + \lambda \sum_{j=1}^p \beta_j^2$$



9 Cross-Validation

Problem with 80/20 Data Separation

- Test Error depends on random set
- For different Set, the test error would be different

With Cross-Validation we can obtain a better estimate of the generalization error

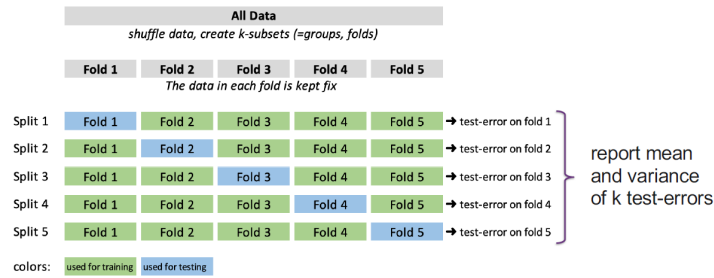
9.1 k-fold Cross-Validation

- Without cross-validation:



With k-Fold Cross-Validation

- The data is split once into k folds. Then train/test is repeated k-times. Each fold participates in k-1 training phases and is used once for testing:



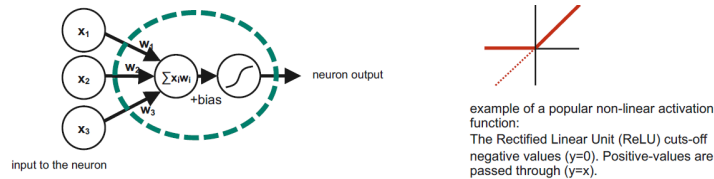
9.1.1 Some Comments

- Typical Values for k are 5,10 or N
- The data of a fold does not change during procedure
- Do not preprocess the whole dataset
- Apply the preprocessing pipe-line to each split

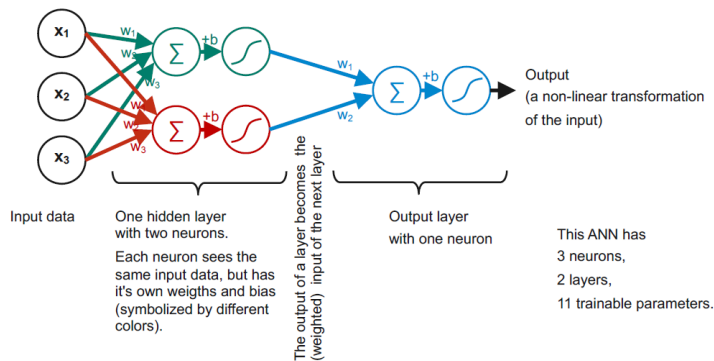
10 Artificial Neural Networks (ANN)

10.1 Artificial Neurons

- Receives an input vector $[x_1, x_2, \dots]$
- Each neuron has its own input weights $[w_1, w_2, \dots]$ and **bias** b
- Calculates the sum of the weighted input (dot product $\vec{x} \cdot \vec{w}$), adds a bias b , and passes it through a nonlinear activation function



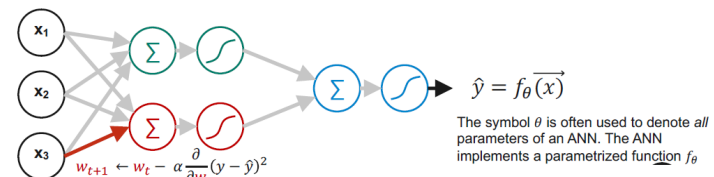
10.2 Simple ANN



10.3 Training an ANN

Supervised learning

- For each input \vec{x} we are given the output \vec{y}
- ANN is initialized with random weights
- An optimizer reduces a cost-function (e.g. MSE)
- At every iteration, and for every single weight w and bias b , the partial derivative needs to be calculated. (Backpropagation)



11 Classification & Logistic Regression

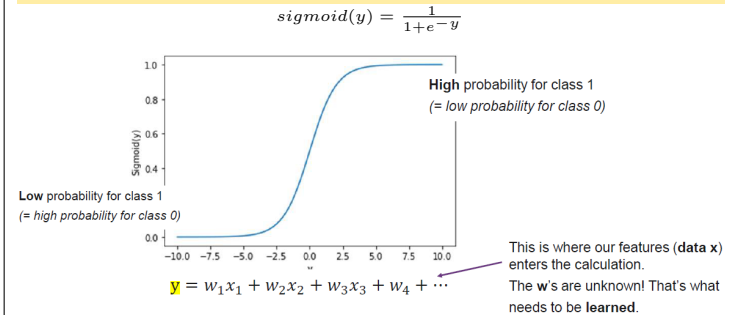
11.1 Binary Classification

- Decision with 2 possible outcomes
- Hail in Lausanne (yes/no)
- Master admission (admission / no admission)
- Based on different data / entity

11.1.1 Decision using Linear Regression

- Train the model with gradient descent
- **Bad Idea!**
- Models the response (y) and post process the response to compute the probability

11.1.2 The sigmoid function



Probabilities

- We can write the estimated probability
- For a prediction we can write

$$P(x) = \frac{1}{1+e^{-W^T x}}$$

11.1.3 Maximum Likelihood

- Given all the data points (X,Y) we want to maximize the probability that all the predictions are correct.
- For each of the training data, we want to maximize the likelihood of correct prediction
- We can use Gradient Descent to find W

Maximize Cost:

$$\text{MaximumCost}_2(W) = \sum_{y=1}^N \ln(p(x_i)) + \sum_{y=0}^N \ln(1 - p(x_i))$$

Loss Function (Minimize Cost):

$$\text{MinimizeCost}(W) = \frac{-1}{N} \sum_{i=1}^N (y_i * \ln(p_i)) + (1 - y_i) * \ln(1 - p_i)$$

12 Classifier Evaluation

12.1 Confusion Matrix

		Predicted condition		
		Positive (PP)	Negative (PN)	
Actual condition	Positive (P)	True positive (TP),	False negative (FN),	Prediction Correct Prediction Wrong
	Negative (N)	False positive (FP),	True negative (TN),	

• True Positive (t_p):
 • model predicted "yes/positive", and
 • the truth is also "yes/positive."

• True Negatives (t_n):
 • model predicted "no/negative", and
 • the truth is also "no/negative."

• False Positives (f_p):
 • model predicted "yes/positive", and
 • the truth is "no/negative".

• False Negatives (f_n):
 • model predicted "no/negative", and
 • the truth is "yes/positive".

Source: Wikipedia

Mean Accuracy:

- How often is the classifier correct?
- $A = (t_p + t_n) / n$

Mean Error:

- How often is the classifier wrong?
- $E = (f_p + f_n)/n$

Precision:

- When the prediction is 1, how often is it correct?
- $P = t_p / (t_p + f_p)$

Sensitivity, Recall, True Positive Rate (TPR):

- How often the prediction is 1 when it's actually 1
- $R = t_p / (t_p + f_n)$

Miss Rate, False Negative Rate (FNR)

- $MR = 1 - TPR$

False Positive Rate (FPR)

- $f_p / (f_p + t_n)$

True Negative Rate (TNR)

- $t_n / (t_n + f_p)$
- $1 - FPR$

12.2 Why Accuracy is not enough?

- If the prediction is constant the accuracy may still look decent
- E.g. always predict false
- 90% of the data is false
- Accuracy = 90% (decent)
- Precision = 0
- Recall = 0

12.3 Precision vs. Recall

- Increasing precision reduces Recall and vice versa
- Threshold is a business decision (depending on goals)

12.4 Receiver Operating Characteristics

- Multiple FPR / TPR needed to draw a curve
- Für eine Kurve müsste man den Classifier mit allen Schwellenwerten im Intervall auswerten
- Defined by FPR and TPR as x and y axes
- Visualizes tradeoff between TP (benefits) and FP (cost)



Area under the curve

- Area under the ROC curve
- Shows how well the TPR and FPR is looking in the aggregate
- The greater the area under the curve, the higher the quality of the model
- The greater the area, the higher the ratio of TP to FP

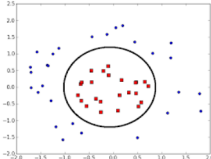
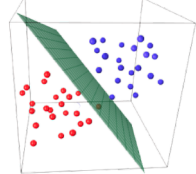
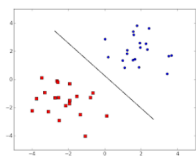
13 KNN

13.1 Linear Seperability

linearly separable in 2D

linearly separable in 3D

not linearly separable



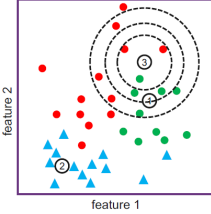
- Based on logistic regression model, you can draw a line
- This is the Linear decision boundary
- If a simple line perfectly separates the classes, then the classes are said to be linearly separable

13.2 Non-Linear decision boundary

- When classes are not linearly separable
- Resort to polynomial terms

13.3 k-Neares Neighbors (KNN)

- A datapoint is know by the company it keeps
- Computes k nearest neighbours
- Returns the most frequent class of the k neighbours



	k=3	k=5	k=10
sample 1	g	g	g
sample 2	b	b	b
sample 3			

- Parameter: how many neighbours? Choice of k!

13.3.1 Distance Metric

Cosine Distance (Skalar)

$$\frac{x_1 * x_2}{||x_1|| * ||x_2||}$$

Manhattan Distance (Quadrat um 90deg gedreht)

$$\sum_{i=1}^n |x_{1,n} - x_{2,n}|$$

Euclidean Distance (Kreis)

$$\sqrt{\sum_{i=1}^n (x_i - x_i)^2}$$

Minkowski Distance

13.3.2 Advantages

- Easy and simple ML model
- Few hyperparameters to tune

13.3.3 Disadvantages

- k should be wisely selected
- Large computation cost during runtime if sample size is large
- Not efficient for high dimensional datasets
- Proper scaling should be provided for fair treatment among features

13.3.4 Hyperparameters

- **K Value:** how many neighbours to participate in the KNN algo.
- **Distance Function:** Euclidean distance is most used

14 Clustering

14.1 Unsupervised Learning

- We are given Data (features, x) without labels (y)
- Can we still learn something from the data?
- Yes! Often the data has some structure
- **The goal** of unsupervised learning is to self-discover patterns from the data

14.2 Clusters

- Data points which have shared properties
- Fall into one cluster or one alike group
- Similar Data Points are close together

14.2.1 Applications

- Social Network Analysis
- Astronomical Data
- Marked segmentation
- Recommendation systems

14.3 Naive K-means

1. Let us assume we know the number of clusters k_c
2. Initialize the value of k cluster centres (aka, means, centroids) $(C_1, C_2, \dots, C_{k_c})$
3. **Assignment :**
 1. Find the **squared Euclidean distance** between the centres and **all the data points**.
 2. Assign each data point to the cluster of the **nearest centre**.
4. **Update:** Each cluster now potentially has a new centre (mean). Update the centre for each cluster
 1. New Centres $((C'_1, C'_2, \dots, C'_{k_c}) = \text{Average of all the data points in the cluster}(1, 2, \dots, k_c))$
5. **If some stopping criterion met, Done**
6. **Else, go to Assignment step 3**

14.3.1 Stopping Criterion

- When centres don't change (time consuming)
- The datapoints assigned to specific cluster remains the same (takes too much time)
- The distance of datapoints from their centres \geq threshold we have set
- Fixed number of iterations have reached (choose wisely)

14.3.2 Initialization

- Performance depends on the random initialization
- Some seeds can result in a poor convergence rate
- Some seeds can converge to suboptimal clustering
- If centres are very close, it takes a lot of iterations to converge
- Initialize randomly, run multiple times

14.3.3 Standardization of data

- Features with large values may dominate the distance value
- Features over small values will have no impact
- Normalize values!

14.3.4 Sklean k-means

Initialization

- Init = K-means++
- Only initialization of the centroids will change
- Chosen centroids should be far from each other

max_iter:

- Number of iterations before stopping

n_init:

- Number of time the k-means algorithm will be run with different centroid seeds

14.3.5 Evaluating Cluster Quality

- Make clusters so that for each cluster the distance of each cluster member from its center is minimizes

Inertia or within-cluster sum-of-squares (WCSS)

- Sum of squared distances to center (euclidian distance)
- As small as possible

Euclidian Distance:

- $(p_x - Center_x)^2 + (p_y - C1_y)^2$

Silhouette Score

- How far the datapoints in one cluster are from the datapoints in another cluster
- SS of a point: $\frac{b-a}{\max(a,b)}$
- a: average intra-cluster distance (distance between each point within)
- b: average inter-cluster distance (distance between a cluster and its nearest neighbour)
- Represents the Quality of the Clustering

15 Ensamble Methods

15.1 Wisdom of Crowd

- Suppose you have a difficult question
- Ask many people and aggregate the answer
- This might work very well instead of finding the best suited person

15.2 Ensemble

- Wisdom of Crowd can be applied to ML
- Instead of finding the best model, aggregate the results of weak models
- Aggregate predictions of regressors or classifiers
- Might get better accuracy than the best predictor
- Ensemble: group of predictors

15.3 Ensemble Method

- Suppose we have many different weak models (better than random)
- Get prediction from all of them and take a vote
- Class with most votes is the predicted class
- Commonly used towards the end of a project
- **Requirement:** enough models / diverse models

15.4 Bagging and Pasting

Bagging (Bootstrap Aggregating)

- Sampling with replacement
- Allows data points to be used several times
- Bootstrap: Reduces Variance
- Aggregating: reduces bias

Pasting

- Sampling without replacement

15.5 No free lunch theorem

No single machine learning algorithm is universally the best-performing algorithm for all problems

15.5.1 Out of Bag (oob) Evaluation

- Using Bagging
- Some Data Points may not be used at all
- Use them for evaluation