

Machine Learning FS2019

Alex Neher

June 4, 2019

Contents

1	Introduction	2
1.1	Disciplines	2
2	Data Quality	3
2.1	Data Quality Assessment	3
2.2	Approaches to Data Quality Assessment	4
2.3	Statistical Key Figures	5
2.3.1	Central Tendency	5
2.3.2	Skewdness	6
2.3.3	Quartile & Interquartile Range (IQR)	6
2.3.4	Five Number Summary	7
2.3.5	Boxplot	7
2.3.6	Variance	7
2.3.7	Covariance	8
2.3.8	Pearson Correlation	8
2.4	Normalization	8
3	Geometry of Data	9
3.1	Feature Engineering	9
3.2	Vector Space Model	9
3.3	Similarity of Data	10
3.3.1	Euclidean Distance	10
3.3.2	Cosine Similarity	11
3.3.3	Levenshtein / Edit Distance for Strings	11

1 Introduction

There are two popular definitions of Machine Learning:

“Field of study that gives computers the ability to learn without being explicitly programmed” (Arthur Samuel, IBM, 1959)

“A computer program is said to learn from experience E with respect to some task T and some performance measure P , if its performance on T , as measured by P , improves with experience E ” (Tom Mitchell, 1998)

So summarizing these two quotes, it can be said, that machine learning is defined as **the process in which machines learn something (mostly) on their own.**

1.1 Disciplines

There are different disciplines in machine learning:

Supervised Learning: The algorithm is given **labeled training data** and learns to **predict** the **labels** of yet unseen examples.

Unsupervised Learning: The algorithm is given **unlabeled data** and **creates labels by itself** based on the structure of the given data

Semi-Supervised Learning: A **mixture** of supervised and unsupervised learning. This approach is usually chosen if there is only **very little labeled test data**

Reinforcement Learning: No data is available, but the algorithm is **being rewarded**. The algorithm searches the ideal behaviour that maximizes its reward (Not subject of this lecture)

These classifications can be subdivided even more:

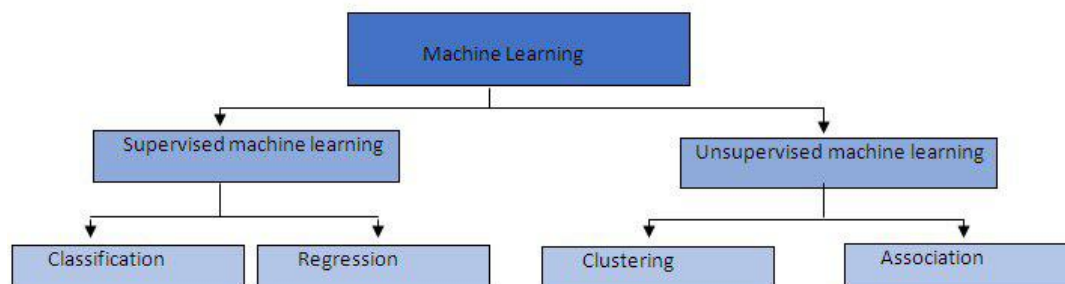


Figure 1.1: Distinction between supervised and unsupervised learning

The main difference between **classification** and **regression** is that when using classification, the result is **categorical**, whereas regression returns **numerical** results.

Clustering is similar to classification. However, while classification algorithms sort the given data into given groups, clustering algorithms determine these groups **by themselves**. This means, you can give a clustering algorithm a seemingly random dataset and the algorithm finds some kind of structure in it.

2 Data Quality

Data is categorized into **numerical** and **categorical** data.

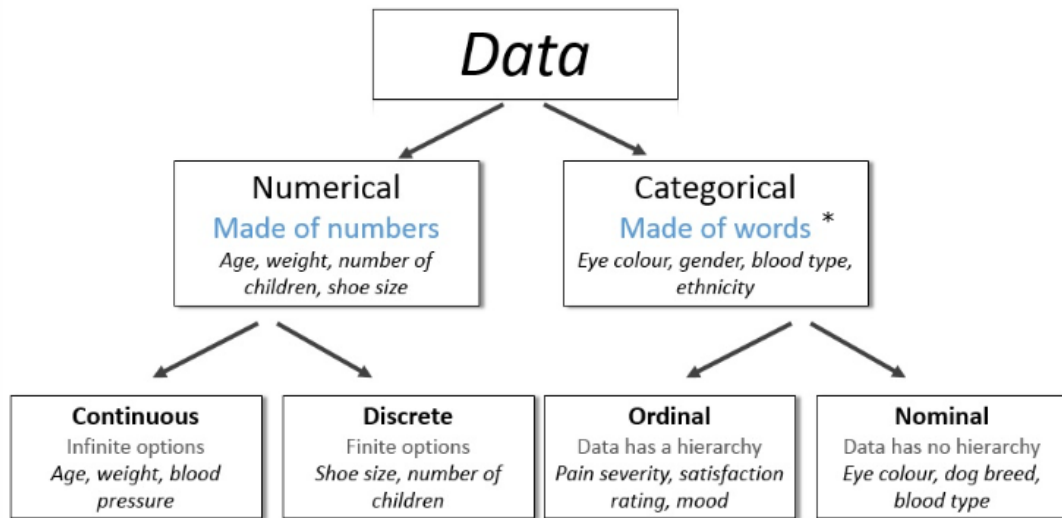


Figure 2.1: Classification of Data

Before any machine learning can take place, the quality of the given data has to be assessed and in some cases improved. Because every prediction made by machine learning algorithms is shit if the data quality is shit.

There are many reasons why the data quality could be poor:

- Ill-designed, inadequate or inconsistent data formats
- Programming errors or technical issues (e.g. sensor outage)
- Data decay (e.g. outdated e-mail addresses)
- Poorly designed data entry forms (e.g. data fields without verification)
- Human errors in data export or data pre-processing
- Deliberate errors and false information (e.g. due to privacy concerns everybody is called Hans Muster and lives at Musterstrasse 123)

2.1 Data Quality Assessment

Before even starting to assess the data-quality, it is seldomly a bad idea to **clean** the data first.

1. Identify and remove duplicates
2. Replace null-values (do not delete them because that might falsify the mean and median of the data)
3. Make data formats more machine-friendly (so-called *data-wrangling* e.g. store the gender as boolean)

If you change anything from the original data set, you should always

- Document all the changes
- Use a SVN (e.g. git)
- Let the data provider know that his data quality is shit (maybe they'll improve in the future)
- Investigate the origins of the poor data quality

2.2 Approaches to Data Quality Assessment

Identify data sources and their trustworthiness

Interpret statical key figures: See following sections

Visualize selected portions of the data: e.g. with Pair Plots (See Abb. 2.2)

Manually check data ranges Negative Salaries, People more than 200 years old...

Validate plausibility of attribute correlation: e.g. are mileage and number of seats in a core correlated? Can one of the columns be removed for redundancy?

Measure data redundancy: Can certain columns be removed due to not adding any real value to the data

Check for anomalies in syntax and semantics: Outliers can really distort a dataset and render the whole algorithm useless. Can be prevented by e.g. normalization of the data or removal of the outlier

Replace NULL Values and remove duplicate values

There are different ways to cope with NULL variables, but they have to be addressed, as most machine learning algorithms do not play well with them.

- Delete all rows with NULL values
Might be the easiest way if you have loads of data
- Fill in the missing values manually (e.g. from other sources)
Might be the hardest way if you have loads of data
- Fill in a global constant like N/A, UNKNOWN
- Use a measure for central tendency
e.g. take the mean if your data is symmetric or take the median if its skewed
- Use a measure for central tendency per class
e.g. take different values for healthy and sick people
- Use e.g Regression to 'guess' the missing values

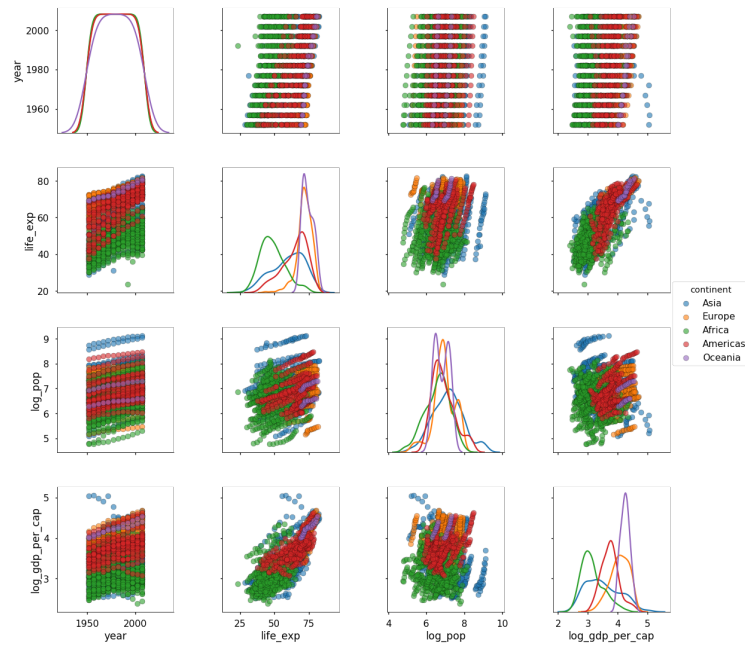


Figure 2.2: Visualisation of Data with Pair Plots

2.3 Statistical Key Figures

These figures can give you a rough overview about the whereabouts of your data-magnitude.

2.3.1 Central Tendency

Mean

This is the average in a set of numeric data. You add all data and divide it by the number of data points

$$\mu_x = \frac{1}{n} \sum_{i=1}^n x_i$$

Mode

This is the value that occurs the most in a given set of data

Median

This is the middlemost value of a sorted set of data. In contrast to the Mean, the Median can give information concerning the distribution of the data.

Given a dataset of 1, 2, 3, 4, 5, the median and mean are both 3. However, if we have 1, 2, 3, 1000, 10000, the mean is 2201.2 whereas the median is still 3

2.3.2 Skewness

All of these values can give information concerning the data's **skewness**

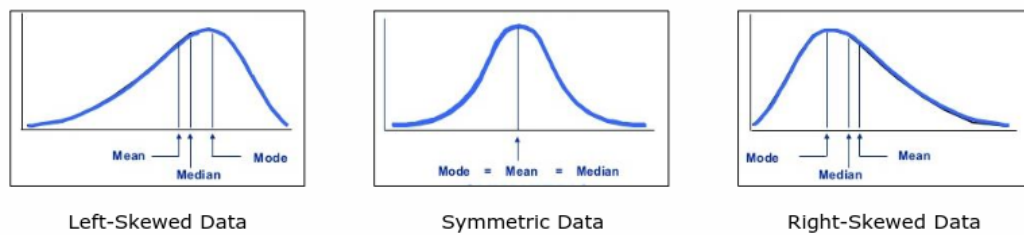


Figure 2.3: Skewness of data

$Mean - Mode > 0 \rightarrow$ Negative skewness / Left-skewed data

$Mean - Mode = 0 \rightarrow$ Symmetric Data

$Mean - Mode < 0 \rightarrow$ Positive skewness / Right-skewed data

2.3.3 Quartile & Interquartile Range (IQR)

The three quartiles divide your data into four equal-sized, consecutive subsets.

To calculate $Q1$, take the median of your data and then again the median of the left half of the data.

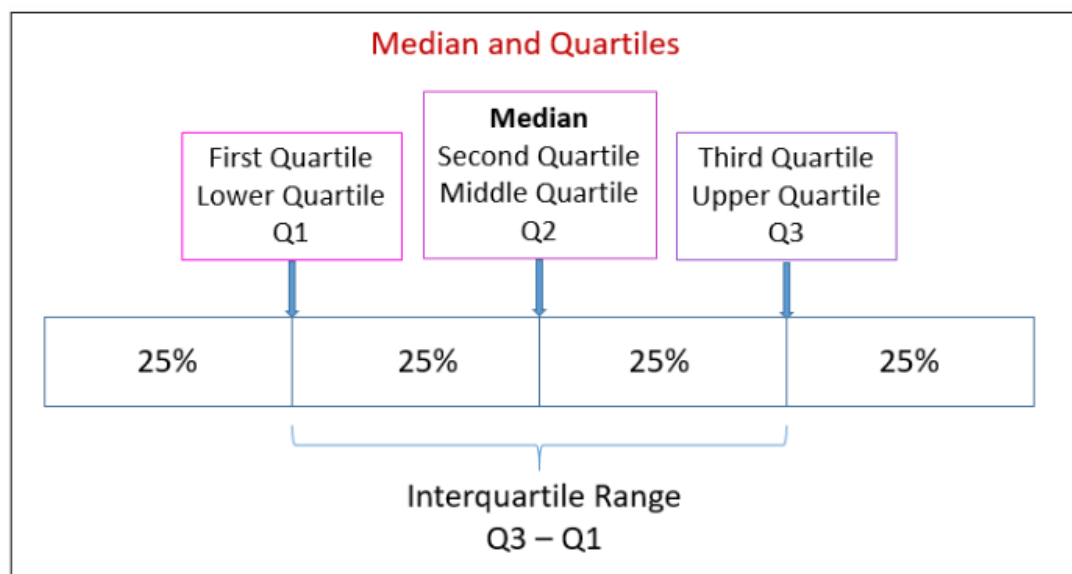


Figure 2.4: Quartiles of a dataset

2.3.4 Five Number Summary

With this method, you can get a pretty good overview of your data. The **Five Number Summary** of a dataset consists of:

- Median Q2
- Quartiles Q1 and Q2
- Smallest individual Value
- Largest individual Value

```
1 import numpy as np
2 import pandas as pd
3
4 s = pd.Series(np.random.rand(100))
5 s.describe()
```

Listing 2.1: Five Number Summary in Python

```
1 mean      0.524559
2 std       0.285565
3 min       0.003933
4 25%       0.298367
5 50%       0.530632
6 75%       0.765907
7 max       0.993293
8 dtype: float64
```

Listing 2.2: Output

2.3.5 Boxplot

This plot is a **visual representation of the five number summary** and can also give information on potential outliers.

Values $1.5 \cdot IQR$ above the 3rd or below the 1st Quartile can be considered outliers and are displayed with small circles.

2.3.6 Variance

The variance shows **how much the values are spread on average**. This is measured by squaring the sum of all deviations from the mean

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

The standard deviation is calculated as $\sqrt{\text{variance}}$

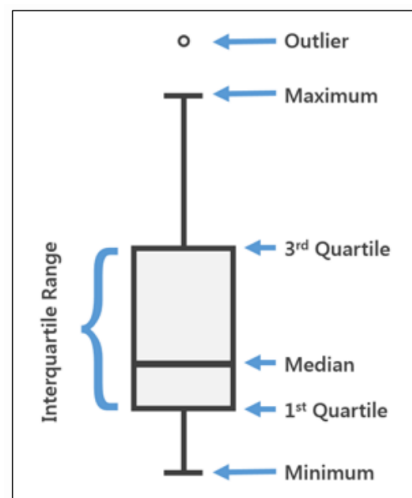


Figure 2.5: Boxplot

2.3.7 Covariance

The covariance is used to determine whether two variables are **connected** to each other.

If both variables are on the same side of the mean, the variance is **positive**, the variables are probably connected. Meaning if the value of one variable is rising, the other one will most likely rise as well.

If one is above and one is below the mean, the variance is **negative**, the variables are most likely **inversely connected** to each other. Meaning if the value of one variable is rising, the other is most likely falling.

If the variables are **independent** from each other, the covariance is zero, as they both cancel each other out.

$$Cov(x, y) = \frac{1}{n} \sum_{i=1}^n (x_i - \mu_x)(y_i - \mu_y)$$

The **covariance matrix** shows the covariance from all X with all Y . As $Cov(x, x) = Var(x)$, the covariance matrix has the variance of X in its diagonal

2.3.8 Pearson Correlation

Both the covariance and the variance are connected to the scale of the dataset, so the covariance of $X = [1, 2, 3, 4, 5]/Y = [6, 7, 8, 9, 10]$ is 2.5, whereas the covariance of $X = [1000, 2000, 3000, 4000, 5000]/Y = [6000, 7000, 8000, 9000, 10000]$ is 2'500'000'000. However, the Pearson Correlation is 1 in both examples.

$$\rho(X, Y) = \frac{Cov(X, Y)}{\sigma_x \sigma_y}$$

The Pearson Correlation is always between 1 and -1.

1 means the data is perfectly correlated, whereas -1 means that the data is perfectly uncorrelated

2.4 Normalization

It is immensely important that all data is normalized before we run a machine learning algorithm over it. Considering the data in figure 3.2, 'Mileage' and 'Price' are in a completely different scale. If the mileage of the first shown car goes up 500 miles, it's not really a big deal. However, a price increase by 500 would double the car's price.

Such differently scaled data can (and will) falsify the result of every machine learning algorithm you could find. Therefore, **normalization is really important**.

There are two popular normalization approaches: The Min-Max and the Z-Score normalization.

Min-Max normalization

All data is condensed to a value between 0 and 1. The smallest value becomes 0 and the largest one becomes 1.

$$x \rightarrow \frac{x - \min_x}{\max_x - \min_x}$$

Z-Score Normalization

The dataset is transformed in such a way, that the mean becomes 0 (so-called *mean-centering*) and the standard deviation is 1

$$x \rightarrow \frac{x - \mu_x}{\sigma_x}$$

3 Geometry of Data

3.1 Feature Engineering

Sometimes, data has to be modified to be better accessible/processable for machine learning algorithms. These algorithms can work the best with simple numbers, so that's the data we should be striving for:

Free	Date	Time	Free	Hour	Minute	Year	Month	Day
283	2015-09-27 00:00:00	06:26:46	283	6	26	2015	9	27
282	2015-09-11 00:00:00	05:18:55	282	5	18	2015	9	11
280	2015-09-20 00:00:00	21:14:49	280	21	14	2015	9	20
283	2015-09-25 00:00:00	01:22:47	283	1	22	2015	9	25
0	2015-10-15 00:00:00	08:12:35	0	8	12	2015	10	15
0	2015-10-27 00:00:00	10:02:28	0	10	2	2015	10	27
281	2015-09-13 00:00:00	12:20:54	281	12	20	2015	9	13
168	2015-10-14 00:00:00	08:07:35	168	8	7	2015	10	14
283	2015-09-25 00:00:00	05:42:47	283	5	42	2015	9	25
283	2015-09-18 00:00:00	22:57:50	283	22	57	2015	9	18
279	2015-09-10 00:00:00	20:26:55	279	20	26	2015	9	10
279	2015-10-04 00:00:00	18:37:40	279	18	37	2015	10	4
84	2015-09-17 00:00:00	17:17:51	84	17	17	2015	9	17
86	2015-09-11 00:00:00	08:28:55	86	8	28	2015	9	11
3	2015-10-26 00:00:00	13:51:28	3	13	51	2015	10	26
281	2015-09-30 00:00:00	00:44:44	281	0	44	2015	9	30
252	2015-10-15 00:00:00	07:19:35	252	7	19	2015	10	15
280	2015-09-15 00:00:00	00:41:52	280	0	41	2015	9	15
282	2015-09-09 00:00:00	06:05:56	282	6	5	2015	9	9
0	2015-10-29 00:00:00	12:16:27	0	12	16	2015	10	29

Figure 3.1: Turn 'complicated' data into easier data for better results

3.2 Vector Space Model

As described before, machine learning algorithms work best with **numeric** data. However, the real world isn't that easy and mostly throws categorical data at you. Therefore, you have to convert categorical data to numerical data.

Name	Price	Mileage	Color	Name	Price	Mileage	braun	gelb	grau	grün	rot	schwarz	silber	weiss
ALFA ROMEO 145 1.4 TS 16V L	500	187000	schwarz	ALFA ROMEO 145 1.4 TS 16V L	500	187000	0	0	0	0	0	1	0	0
ALFA ROMEO 145 1.8 TS 16V L	2600	182510	rot	ALFA ROMEO 145 1.8 TS 16V L	2600	182510	0	0	0	0	1	0	0	0
ALFA ROMEO 145 1.9 JTD	3500	116000	grau	ALFA ROMEO 145 1.9 JTD	3500	116000	0	0	1	0	0	0	0	0
ALFA ROMEO 145 2.0 TS 16V Quadrifoglio	4900	181000	rot	ALFA ROMEO 145 2.0 TS 16V Quadrifoglio	4900	181000	0	0	0	0	1	0	0	0
ALFA ROMEO 145 2.0 TS 16V Quadrifoglio	800	121000	rot	ALFA ROMEO 145 2.0 TS 16V Quadrifoglio	800	121000	0	0	0	0	1	0	0	0
ALFA ROMEO 145 2.0 TS 16V Quadrifoglio	3200	156000	schwarz	ALFA ROMEO 145 2.0 TS 16V Quadrifoglio	3200	156000	0	0	0	0	0	1	0	0
ALFA ROMEO 146 2.0 Ti 16V	770	158000	grau	ALFA ROMEO 146 2.0 Ti 16V	770	158000	0	0	1	0	0	0	0	0
ALFA ROMEO 146 2.0 Ti 16V	1200	119000	rot	ALFA ROMEO 146 2.0 Ti 16V	1200	119000	0	0	0	0	1	0	0	0
ALFA ROMEO 146 2.0 Ti 16V	4900	166000	schwarz	ALFA ROMEO 146 2.0 Ti 16V	4900	166000	0	0	0	0	0	1	0	0
ALFA ROMEO 146 2.0 Ti 16V	4900	102000	silber	ALFA ROMEO 146 2.0 Ti 16V	4900	102000	0	0	0	0	0	0	1	0
ALFA ROMEO 146 2.0 Ti 16V Kit Sport	5800	165000	schwarz	ALFA ROMEO 146 2.0 Ti 16V Kit Sport	5800	165000	0	0	0	0	0	1	0	0
ALFA ROMEO 147 1.6 16V Blackline	11500	46230	braun	ALFA ROMEO 147 1.6 16V Blackline	11500	46230	1	0	0	0	0	0	0	0

Figure 3.2: Turn categorical data into numerical data with the vector space model

This transformed data can also be visualized in a coordinate system and we can do math with it.

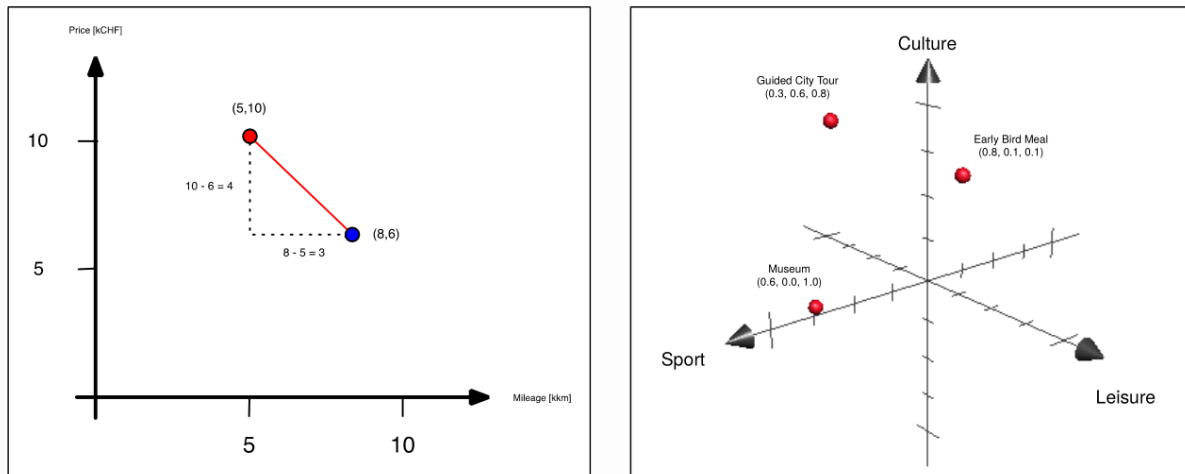


Figure 3.3: Transformed into vector space, data points can be interpreted as geometric points

3.3 Similarity of Data

The math we want to do is not even overly complicated: We just want to measure the distance between different points. Because **the smaller the distance between two points, the more similar they are.**

3.3.1 Euclidean Distance

The distance between two points is most easily calculated using the **euclidean distance**:

$$\text{dist}(X, Y) = \sqrt{\sum_{i=1}^n (x_i - y_i)^2}$$

So the distance between the points (5/10) and (8/6) can be calculated as

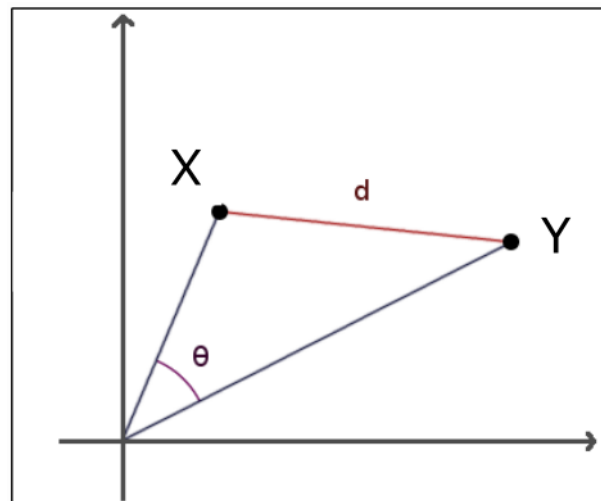
$$\begin{aligned} & \sqrt{(5 - 8)^2 + (10 - 6)^2} \\ & \quad \sqrt{-3^2 + 4^2} \\ & \quad \sqrt{9 + 16} \\ & \quad \sqrt{25} = 5 \end{aligned}$$

3.3.2 Cosine Similarity

If you want to compare two points that appear to be on a line (Pearson Correlation close to 1), but the euclidean distance is high, then the cosine similarity is probably pretty low.

The cosine similarity looks at the **angle** between point A and point B. However, it does also take the euclidean distance into consideration.

The cosine similarity is essentially just the scalar product of the two points.



$$sim(X, Y) = \frac{\sum_{i=1}^n x_i y_i}{\sqrt{\sum_{i=1}^n x_i^2} \sqrt{\sum_{i=1}^n y_i^2}}$$

Figure 3.4: Cosine Similarity

$$dist(X, Y) = 1 - sim(X, Y)$$

3.3.3 Levenshtein / Edit Distance for Strings

Count the minimal number of changes necessary to turn one string into another:

- count +1 when deleting a character [d]
- count +1 when adding a character [a]
- count +2 when changing a character [c]

1. Word	2. Word	Levenshtein Distance
Hello	Yellow	1 [c] + 1 [a] = 3
MacDonald	McDonalds	1 [d] + 1 [a] = 2
banana	ananas	? d+a=2

Figure 3.5: Examples for Levenshtein Distance