

Data Science

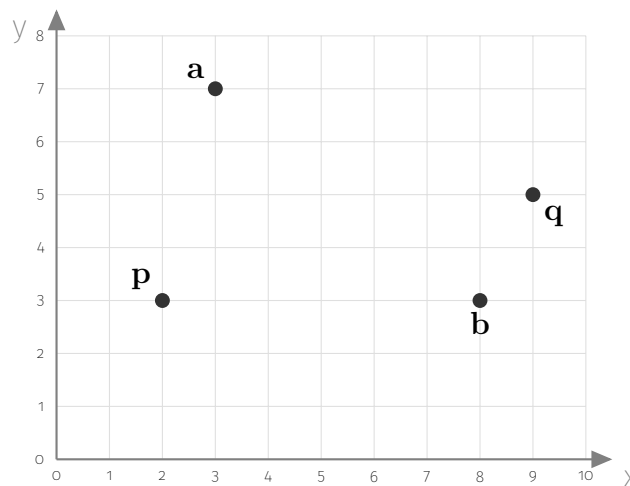
Sommersemester 2022

Übungsblatt 7

Dieses Übungsblatt befasst sich mit linearen Modellen und der damit verbundenen Thematik der Vektorräume. Vektoren lassen sich in Pandas sehr gut als **Series** Objekte darstellen. Damit kann man mit Python/Pandas sehr komfortabel Vektorraum Dinge umsetzen.

Aufgabe 1 (Vektoren mit Pandas)

Schauen Sie sich den folgenden Plot an.



1. Definieren Sie die Punkte **p**, **q**, **a** und **b** als **Series** Objekte.
2. Implementieren Sie eine Funktion **skalarprodukt(v, w)**, die das Skalarprodukt von zwei Vektoren **v** und **w** (gegeben jeweils als **Series** Objekt) berechnet.
3. Schreiben Sie eine Funktion **dist(v, w)**, die für zwei Vektoren (als **Series** Objekte) die Distanz zwischen den Vektoren berechnet. Welchen Abstand haben die Vektoren **p** und **q** aus der obigen Skizze?
4. Definieren Sie eine Python Funktion **f(x)**, die die Gerade beschreibt, die durch die Punkte **p** und **q** läuft.
5. Berechnen Sie den Vektor $\mathbf{v} = \mathbf{q} - \mathbf{p}$ als weiteres **Series** Objekt und berechnen Sie die Skalarprodukte $\langle \mathbf{b}, \mathbf{v} \rangle$ und $\langle \mathbf{a}, \mathbf{v} \rangle$.

Aufgabe 2 * (Support Vector Machine)

Schauen Sie sich das Notebook **Vorlesung/V7-Lineare-SVM.ipynb** auf dem Notebook-Server an.

1. Erzeugen Sie aus dem Iris-Datensatz mit den zwei Klassen *setosa* und *versicolor* einen Trainings- und einen Test-Datensatz. Nutzen Sie dafür 80% der Daten für die Trainingsmenge.
2. Trainieren Sie das SVM-Modell auf dem Trainingsdatensatz und berechnen Sie mit Hilfe der SkLearn Funktion **accuracy_score** den Trainingsfehler:

```
from sklearn.metrics import accuracy_score

y_pred = m.predict(X_train)
train_error = 1 - accuracy_score(y_train, y_pred)
```

3. Berechnen Sie auch den Testfehler.
4. Erzeugen Sie auch ein k -NN Modell mit $k = 5$ und trainieren Sie es auf den gleichen Trainingsdaten. Wie gut ist die Vorhersagegüte *accuracy* von k -NN im Vergleich zum SVM-Modell?
5. Schauen Sie sich für beide Modelle auch den **classification_report** an und vergleichen Sie die beiden Modelle auch anhand der anderen Metriken.