

Prof. Dr. Agnès Voisard
Nicolas Lehmann

Datenbanksysteme, SoSe 18

Übung 04

TutorIn: Toni Draßdo
Tutorium 014

Eduard Beiline, Mark Niehues, Antoen Oehler

21. Mai 2018

Task 1: ER-Modellierung

Task 2: Relationales Modell

Task 3: Reverse Engineering

Task 4: Data Mining

1 - K-Means

Der Log des K-Mean Aufrufs ist in Listing 1 angegeben. Der dazugehörige entwickelte Code in Listing 2.

Listing 1: Log File des K-Means Algorithmus

```
1 ::::Centers and their Clusters at Step 1
2 Center for this Cluster: [10 10 3]
3 Contains:
4 [10 10 3]
5 [7 9 1]
7 Center for this Cluster: [9 2 3]
8 Contains:
9 [9 2 3]
10 [9 4 1]
11 [6 6 8]
12 [4 3 5]
14 Center for this Cluster: [ 3 10 1]
15 Contains:
16 [ 3 10 1]
17 [4 5 1]
```

```

20 ::::Centers and their Clusters at Step 2
21 Center for this Cluster: [ 8.5  9.5  2. ]
22 Contains:
23 [10 10  3]
24 [7 9 1]

26 Center for this Cluster: [ 7.    3.75  4.25]
27 Contains:
28 [9 2 3]
29 [9 4 1]
30 [6 6 8]
31 [4 3 5]

33 Center for this Cluster: [ 3.5  7.5  1. ]
34 Contains:
35 [ 3 10  1]
36 [4 5 1]

39 ::::Final Centers and their Clusters:
40 Center for this Cluster: [ 8.5  9.5  2. ]
41 Contains:
42 [10 10  3]
43 [7 9 1]

45 Center for this Cluster: [ 7.    3.75  4.25]
46 Contains:
47 [9 2 3]
48 [9 4 1]
49 [6 6 8]
50 [4 3 5]

52 Center for this Cluster: [ 3.5  7.5  1. ]
53 Contains:
54 [ 3 10  1]
55 [4 5 1]

```

Listing 2: K-Means Implementierung

```

1  #!/usr/bin/env python3
2  # coding: utf-8
3  import numpy as np
4  import sys

6  DATA = np.array([
7      [3,9,9,10,6,7,4,4],
8      [10,2,4,10,6,9,5,3],
9      [1,3,1,3,8,1,1,5]
10 ])

13 def expectation(data, centers):
14     """
15     Assigns datapoints to centers
16     """
17     clusters = [[] for _ in range(centers.shape[0])]

19     for point in data:
20         # Calculate distances from centers
21         d = np.linalg.norm(centers - point, axis=1)

23         # Find the index according to the lowest distance
24         id_min = np.argmin(d)

26         # Assigning point to minimum distance
27         clusters[id_min].append(point)

29     for clt in clusters:

```

```

30         if len(clt) == 0:
31             print("Error: Empty cluster occurred, please re-run the program.")
32             sys.exit(1)
33         return clusters

36 def minimization(clusters):
37     """
38     Computes new cluster means
39     """
40     centers = [np.mean(cls, axis=0) for cls in clusters]

42     return np.vstack(centers)

45 def k_means(data, k, sigma):
46     # Initialize with k random points
47     centers = data[np.random.randint(data.shape[0], size=k)]

49     dist = sigma + 1
50     i = 0
51     while dist > sigma:
52         i+=1
53         # Assign data points to centers
54         clusters = expectation(data, centers)
55         print("::::Centers and their Clusters at Step " + str(i))
56         print_clusters(centers, clusters)
57         print()
58         # Calculate new centers
59         new_centers = minimization(clusters)

61         # Calc maximum center movement
62         dist = max(np.linalg.norm(centers - new_centers, axis=1))
63         centers = new_centers

65     return centers, clusters

68 def print_clusters(centers, clusters):
69     for center, clst in zip(centers, clusters):
70         print("Center for this Cluster: {}".format(center))
71         print("Contains:")
72         for _ in clst:
73             print(_)
74         print()

77 centers, clusters = k_means(DATA, k=3, sigma=3/4)
78 print("::::Final Centers and their Clusters:")
79 print_clusters(centers, clusters)

```

2 - Naive Bayes

1 - Wahrscheinlichkeit einer Grippe bei laufender Nase

Naive Bayes Formula:

$$P(C|x) = \frac{P(C) P(x|C)}{P(x)} \quad (1)$$

Aus Formel 1 folgt für die Wahrscheinlichkeit an einer Grippe zu leiden, bei laufender Nase:

$$P(\text{Grippe}|\text{Nase}) = \frac{P(\text{Grippe}) P(\text{Nase}|\text{Grippe})}{P(\text{Nase})} \quad (2)$$

wobei:

$$\begin{aligned}P(Nase) &= 4/8 = 1/2 \\P(Grippe) &= 1/2 \\P(Nase|Grippe) &= 3/4\end{aligned}$$

Durch einsetzen in Formel 2 erhält man:

$$P(Grippe|Nase) = 3/4$$

2 - Grippe, wenn X

Um die Frage zu beantworten, ob jemand eher Grippe oder keine Grippe besitzt wird, um die Rechnung zu Vereinfachen der Quotient aus $P(Grippe|x)$ und $P(\neg Grippe|x)$ gebildet, dadurch kürzt sich die aufwändig zu berechnende Evidenz $P(x)$ heraus:

$$Q = \frac{P(Grippe|x)}{P(\neg Grippe|x)} = \frac{P(x|Grippe) P(Grippe)}{P(x|\neg Grippe) P(\neg Grippe)} \quad (3)$$

wobei:

$$\begin{aligned}x &= \{Schttelfrost, schwacheKopfschmerzen, Fieber\} \\P(Schttelfrost|Grippe) &= 3/4 \\P(Schttelfrost|\neg Grippe) &= 1/2 \\P(schwacheKopfschmerzen|Grippe) &= 1/4 \\P(schwacheKopfschmerzen|\neg Grippe) &= 1/4 \\P(Fieber|Grippe) &= 1/2 \\P(Fieber|\neg Grippe) &= 1/2 \\P(Grippe) &= P(\neg Grippe) = 1/2\end{aligned}$$

Unter Annahme der (hinreichenden) Unabhängigkeit der Variablen, gilt $P(x_1, x_2|C) = P(x_1|C) P(x_2|C)$.

Daraus ergibt sich schließlich:

$$Q = \frac{\frac{3}{2} \frac{1}{4} \frac{1}{2}}{\frac{1}{2} \frac{1}{4} \frac{1}{2}} = 3/2$$

Aus der $Q > 1$ folgt, dass der Patient wahrscheinlicher Grippe hat als keine.

3 - Apriori

Die Supports werden anschließend wiederum kombiniert um Beziehungen der Form $\{A, B\} \rightarrow \{C\}$ zu bewerten. Die sogenannte *confidence* berechnet sich aus den Support als:

$$conf(X \rightarrow Y) = \frac{sup(X \cup Y)}{sup(X)} \quad (4)$$

Das Ergebnis der Berechnungen sei an dieser Stelle der Einfachheit halber als Textdatei angegeben:

Tabelle 1: Ausschnitt aus der Berechnung der Supports durch Kombinatorik und Abzählen

C_0	L_0	C_1	L_1	C_2
$\sup(\{A\}) = \frac{1}{2}$	$\{A\}$	$\sup(\{A,B\}) = \frac{1}{2}$	$\{A,B\}$...
$\sup(\{B\}) = \frac{5}{6}$	$\{B\}$	$\sup(\{A,C\}) = \frac{1}{3}$	$\{A,E\}$...
$\sup(\{C\}) = \frac{1}{2}$	$\{C\}$	$\sup(\{A,D\}) = \frac{1}{3}$	$\{B,C\}$...
$\sup(\{D\}) = \frac{5}{6}$	$\{D\}$	$\sup(\{A,E\}) = \frac{2}{3}$	$\{B,D\}$...
$\sup(\{E\}) = \frac{5}{6}$	$\{E\}$	$\sup(\{A,F\}) = \frac{1}{6}$	$\{B,E\}$...
$\sup(\{F\}) = \frac{1}{2}$	$\{F\}$	$\sup(\{B,C\}) = \frac{1}{2}$	$\{C,D\}$...
		$\sup(\{B,D\}) = \frac{2}{3}$	$\{D,E\}$...
		$\sup(\{B,E\}) = \frac{2}{3}$	$\{D,F\}$...
		$\sup(\{B,F\}) = \frac{1}{3}$	$\{E,F\}$...
		$\sup(\{C,D\}) = \frac{1}{2}$...
		$\sup(\{C,E\}) = \frac{1}{3}$...
		$\sup(\{C,F\}) = \frac{1}{6}$...
		$\sup(\{D,E\}) = \frac{2}{3}$...
		$\sup(\{D,F\}) = \frac{1}{2}$...
		$\sup(\{E,F\}) = \frac{1}{2}$...

Listing 3: Ergebnis des Apriori Algorithmus.

```

1 Supports:
2 -----
3 ('F') , 0.500
4 ('A') , 0.500
5 ('C') , 0.500
6 ('D') , 0.833
7 ('B') , 0.833
8 ('E') , 0.833
9 ('D', 'F') , 0.500
10 ('A', 'E') , 0.500
11 ('C', 'B') , 0.500
12 ('C', 'D') , 0.500
13 ('A', 'B') , 0.500
14 ('E', 'F') , 0.500
15 ('B', 'D') , 0.667
16 ('B', 'E') , 0.667
17 ('E', 'D') , 0.667
18 ('C', 'B', 'D') , 0.500
19 ('E', 'D', 'F') , 0.500
20 ('A', 'B', 'E') , 0.500
21 ('B', 'E', 'D') , 0.500

25 Regeln:
26 -----
27 ('B', 'D') ==> ('C') , 0.750
28 ('E', 'D') ==> ('F') , 0.750
29 ('B', 'E') ==> ('A') , 0.750
30 ('B', 'E') ==> ('D') , 0.750
31 ('B', 'D') ==> ('E') , 0.750
32 ('E', 'D') ==> ('B') , 0.750

```

```

33 ('B') ==> ('D') , 0.800
34 ('D') ==> ('B') , 0.800
35 ('B') ==> ('E') , 0.800
36 ('E') ==> ('B') , 0.800
37 ('E') ==> ('D') , 0.800
38 ('D') ==> ('E') , 0.800
39 ('F') ==> ('D') , 1.000
40 ('A') ==> ('E') , 1.000
41 ('C') ==> ('B') , 1.000
42 ('C') ==> ('D') , 1.000
43 ('A') ==> ('B') , 1.000
44 ('F') ==> ('E') , 1.000
45 ('C') ==> ('B', 'D') , 1.000
46 ('C', 'B') ==> ('D') , 1.000
47 ('C', 'D') ==> ('B') , 1.000
48 ('F') ==> ('E', 'D') , 1.000
49 ('E', 'F') ==> ('D') , 1.000
50 ('D', 'F') ==> ('E') , 1.000
51 ('A') ==> ('B', 'E') , 1.000
52 ('A', 'B') ==> ('E') , 1.000
53 ('A', 'E') ==> ('B') , 1.000

```

4 - Lineare Regression

Leider ist bei uns erst eine Korrektur eingetragen, daher nehmen wir für diese Aufgabe folgende Noten an:

Tabelle 2: Bisherige Notenverteilung

x	0	1	2
y	0,89	0,92	0,93

Aus Tabelle 2 folgt: $\bar{x} = 1$ und $\bar{y} = 0.91\bar{3}$.

Die lineare Regression beschreibt die Daten mit einer Funktion der Art:

$$f(x) = \beta_0 + \beta_1 x + \epsilon \quad (5)$$

mit

$$\beta_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} \quad (6)$$

$$\beta_0 = \bar{y} - \beta_1 \bar{x} \quad (7)$$

$$(8)$$

Aus Gleichung 6 folgt $\beta_1 = 0,02$ und damit $\beta_0 = 0.893$. Aus der linearen Gleichung 5 ergibt sich also für den nächsten Zettel $f(3) = 0,953$.