

Prof. Dr. Agn es Voisard, Nicolas Lehmann

Datenbanksysteme, SoSe 2017

 bungsblatt 5

Tutor: Nicolas Lehmann

Tutorium 10

Boyan Hristov, Julian Habib

8. Juni 2017

Link zum Git Repository: <https://github.com/BoyanH/Freie-Universitaet-Berlin/tree/master/Datenbanksysteme/Solutions/homework5>

1. Aufgabe

a)

$$\begin{aligned} \text{Minimal F} = \{ \\ & \{B\} \rightarrow \{C\} \\ & \{B\} \rightarrow \{D\} \\ & \{C\} \rightarrow \{E\} \\ & \{D\} \rightarrow \{E\} \\ & \{AE\} \rightarrow \{D\} \\ & \} \end{aligned}$$

$$\begin{aligned} & \Rightarrow \\ F+ = \{ \\ & \{B\} \rightarrow \{C\} \\ & \{B\} \rightarrow \{D\} \\ & \{B\} \rightarrow \{CD\} \\ & \{C\} \rightarrow \{E\} \\ & \{D\} \rightarrow \{E\} \\ & \{AE\} \rightarrow \{D\} \\ & \{B\} \rightarrow \{C\} \\ & \{B\} \rightarrow \{D\} \\ & \{BC\} \rightarrow \{CDE\} \\ & \{BD\} \rightarrow \{CDE\} \\ & \{CD\} \rightarrow \{E\} \\ & \{AB\} \rightarrow \{CDE\} \\ & \dots \\ & \} \end{aligned}$$

Damit ist AB ein Key / Schlüssel.

b) Pseudotransitivität wäre hier

$$\begin{aligned} & \{C\} \rightarrow \{E\} \wedge \{AE\} \rightarrow D \\ \Rightarrow & \{AC\} \rightarrow \{D\} \end{aligned}$$

c)

$$\begin{aligned} & \{BE\} \rightarrow \{B\} (\text{Reflexivity}) \\ \Rightarrow & \{A\} \rightarrow \{BE\} \wedge \{BE\} \rightarrow \{B\} \\ \Rightarrow & \{A\} \rightarrow \{B\} (\text{Transitivity}) \\ \Rightarrow & \{A\} \rightarrow \{B\} \wedge \{A\} \rightarrow \{A\} (\text{Reflexivity}) \\ \Rightarrow & \{A\} \rightarrow \{AB\} (\text{Union}) \wedge \{AB\} \rightarrow \{C\} \\ \Rightarrow & \{A\} \rightarrow \{C\} (\text{Transitivity}) \\ \Rightarrow & \{A\} \rightarrow \{C\} \wedge \{C\} \rightarrow \{D\} \\ \Rightarrow & \{A\} \rightarrow \{D\} (\text{Transitivity}) \end{aligned}$$

2. Aufgabe

a) Attributhülle von $\{A\}$

Es gibt keine funktionale Abhängigkeiten für den Attribut A, deswegen ist seine Attributhülle $\{A\}$ den Attribut selbst.

Attributhülle von $\{B\}$

$$\begin{aligned} & \{B\} \rightarrow \{AD\} \Rightarrow \{B\} \rightarrow \{D\} \text{ (Reflexivity)} \\ \Rightarrow & \{B\} \rightarrow \{D\} \wedge \{D\} \rightarrow \{BC\} \\ \Rightarrow & \{B\} \rightarrow \{C\} \text{ (Transitivity)} \end{aligned}$$

$$\begin{aligned} \Rightarrow & \{B\} \rightarrow \{C\} \wedge \{B\} \rightarrow \{D\} \\ \Rightarrow & \{B\} \rightarrow \{CD\} \text{ (Union)} \\ \Rightarrow & \{B\} \rightarrow \{CD\} \wedge \{CD\} \rightarrow \{AEF\} \\ \Rightarrow & \{B\} \rightarrow \{AEF\} \text{ (Transitivity)} \end{aligned}$$

$$\{B\} \rightarrow \{B\} \text{ (Reflexivity)}$$

$$\{B\} \rightarrow \{ABCDEF\}$$

Attributhülle von $\{B\} = B^+ = \{ABCDEF\}$. Deswegen ist B ein Schlüssel.

b) Wie wir bei Aufgabenteil a) gesehen haben ist B ein Schlüssel. Wegen $\{D\} \rightarrow \{BC\}$ ist D ein weiteres Schlüssel. Wegen $\{F\} \rightarrow \{ABD\}$ ist F auch ein Schlüssel. Von $\{E\} \rightarrow \{CD\}$ ist E auch Schlüssel. Weitere Kandidaten gibt es nicht, da per Definition ein Schlüssel das kleinste Superkey ist und alle andere Superkeys haben mindestens 2 Attributen.

Kandidaten für Schlüssel: B, D, E, F

c) Wir haben aus allen Kandidaten B für Schlüssel gewählt, da dieser am trivialsten zu sehen ist und auch alphabetisch der ersten Attribut ist. Da es keine feste Regeln gibt, welcher von den Kandidaten zu einem Schlüssel wird, war unsere Entscheidung nicht so schwer zu treffen.

3. Aufgabe

$$FD(R_2) = \{ \begin{array}{l} \{A\} \rightarrow \{D, E\} \\ \{B\} \rightarrow \{C, D\} \\ \{C\} \rightarrow \{A\} \\ \{D, E\} \rightarrow \{A, C\} \end{array} \}$$

$$\begin{array}{l} \{B\} \rightarrow \{CD\} \Rightarrow \{B\} \rightarrow \{C\} \text{(Reflexivity)} \\ \Rightarrow \{B\} \rightarrow \{C\} \wedge \{C\} \rightarrow \{A\} \\ \Rightarrow \{B\} \rightarrow \{A\} \text{(Transitivity)} \\ \Rightarrow \{B\} \rightarrow \{A\} \wedge \{A\} \rightarrow \{DE\} \\ \Rightarrow \{B\} \rightarrow \{DE\} \text{(Transitivity)} \\ \Rightarrow \{B\} \rightarrow \{D\} \text{(Reflexivity)} \\ \Rightarrow \text{Dependency } \{B\} \rightarrow \{CD\} \text{ kann durch } \{B\} \rightarrow \{C\} \text{ substituiert werden} \end{array}$$

$$\{C\} \rightarrow \{A\} \Rightarrow \text{Dependency } \{DE\} \rightarrow \{A\} \text{ ist wegen dependency } \{DE\} \rightarrow \{C\} \text{ nicht notwendig}$$

$$\Leftrightarrow \begin{array}{l} FD(R_2) = \{ \begin{array}{l} \{A\} \rightarrow \{D, E\} \\ \{B\} \rightarrow \{C\} \\ \{C\} \rightarrow \{A\} \\ \{D, E\} \rightarrow \{C\} \end{array} \} \end{array}$$

$$\Leftrightarrow \text{(Decomposition)}$$

$$FD(R_2) = \{ \begin{array}{l} \{A\} \rightarrow \{D\} \\ \{A\} \rightarrow \{E\} \\ \{B\} \rightarrow \{C\} \\ \{C\} \rightarrow \{A\} \\ \{D, E\} \rightarrow \{C\} \end{array} \}$$

Damit sind alle Anforderungen aus der Vorlesung erfüllt \rightarrow

1. Every right side of a dependency in F is a single attribute. (apply decomposition)

2. For no $X \rightarrow A$ in F is the set $F - \{X \rightarrow A\}$ equivalent to F .
3. For no $X \rightarrow A$ in F and subset Z of X is $F - \{X \rightarrow A\} \cup \{Z \rightarrow A\}$ equivalent to F .

Aufgabe 5

- a) Der Scraper liest mehrere Seiten von einem Web Application, und zwar auf mehrere Sprachen über die Jahren 2000 bis 2015 inklusive. Von den Seiten nimmt er immer die erste Tabelle, davon alle Zeilen mit rechtsbündigen Text. Für jede Zeile in der Tabelle speichert der Scraper eine Zeile in .csv Datei, wo Land, Jahr und die Zellen in der Tabelle nacheinander stehen.

Leider könnten wir den Scraper nicht ausführen, da diese Seiten nicht erreichbar waren. Laut Julian aber speichert der Scraper alle Informationen von Hunderassen der Jahre 2000-2016 in verschiedene Sprachen in einer CSV Datei.

Nach mehrere versuche und VPN (war es überhaupt notwendig?) und nachdem es spät genug war, damit nicht alle wild die Seite scrapen, haben wir gesehen, dass auf der Seite eine Tabelle von den besten Hunden in Wettbewerben steht, je mit Name, Geschlecht, abgeschlossene Wettbewerbe usw. Jede Zeile aus alle solchen Tabellen wird eine Zeile in unsere CSV Datei und dazu werden noch zwei Spalten für Land und Jahr eingefügt. Im großen und ganzen wird eine perfekte CSV Datei für das DBS Projekt im nächsten Semester :)

b)

```

1 from bs4 import BeautifulSoup
2 import requests
3 import re
4 from collections import Counter

6 PER_PAGE_NAVIGATION_CLASS = 'seitenweise_navigation'
7 PAGINATION_CLASS = 'pagination'
8 MAIN_PAGE_URL = 'https://www.heise.de/thema/https'
9 SITE_SUFFIX = '?seite='

11 # this function returns a soup page object
12 def getPage(url):
13     r = requests.get(url)
14     data = r.text
15     spobj = BeautifulSoup(data, "lxml")
16     return spobj

18 # pushes article soup objects into the passed articles array
19 def parsePage(soup, articles):
20     # Summary:
21     # In each page, get the <aside class="recommendations"> elements
22     # Inside, search for a <div class="recommendation">
23     # Every div of that kind is an article, add it to the articles array passed to the
        function

25     recommendationsAside = soup.find('aside', class_ = 'recommendations')
26     recommendations = recommendationsAside.find_all('div', class_ = 'recommendation')

28     for rec in recommendations:
29         articles.append(rec)

31 # returns array
32 def getHeadings(articles):
33     # Summary:
34     # get all <header> tags in an article,
35     # remove new lines and extra spaces and return the result

```

```

37     headings = []

39     for article in articles:
40         heading = article.find('header')
41         headingText = heading.text.replace('\n', '').strip()
42         headings.append(headingText)

44     return headings

46 # returns Dictionary<Word, Count>
47 def getPerWordCountInHeadings(headings):
48     # Summary:
49     # All symbol are removed from the article except for the one in the german
50     # language and dashes
51     # Words are then separated by dashes and spaces and counted

52     perWordCount = {}

54     for heading in headings:
55         #re.sub...
56         words = re.sub('[^a-zA-Z äöüÜß\-\`\' ]', '', heading).replace('-', ' ').lower()
57         .split(' ')

58         for word in words:
59             if len(word) < 3:
60                 continue
61             if not word in perWordCount:
62                 perWordCount[word] = 1
63             else:
64                 perWordCount[word] += 1

66     return perWordCount

68 # returns array
69 def getTopNWords(wordCountDict, n):
70     # Summary:
71     # Using the functionality from collections.Counter get the top 3 keys
72     # in the dictionary sorted by their value and push them to the returned array

74     counter = Counter(wordCountDict)
75     topNWords = []

77     for word in counter.most_common(n):
78         topNWords.append(word)

80     return topNWords

83 # scraper website: greyhound-data.com
84 def main():

86     articles = []

88     mainPage = getPage(MAIN_PAGE_URL)
89     navigation = mainPage.find_all('nav', class_ = PER_PAGE_NAVIGATION_CLASS)[0]
90     navItemsContainer = mainPage.find('span', class_ = PAGINATION_CLASS)
91     nonActivePages = navItemsContainer.find_all('a')
92     furtherPages = len(nonActivePages)

94     parsePage(mainPage, articles);

96     for page in range(1, furtherPages+1):
97         furtherPage = getPage(MAIN_PAGE_URL + SITE_SUFFIX + str(page))
98         parsePage(furtherPage, articles)

100     headings = getHeadings(articles)
101     perWordCount = getPerWordCountInHeadings(headings)

```

```

102     top3Words = getTopNWords(perWordCount, 3)

104     print("\nDie meist benutzten 3 Wörter mit mehr als 3 Buchstaben inklusive sind: \n")
105     for word in top3Words:
106         print("{0}: {1} mal".format(word[0], word[1]))
107     print("\n\nWichtig: Nur die Artikeln link die zu HTTPS releavnt sind wurden
        analysiert!\n")

110 # main program

112 if __name__ == '__main__':
113     main()

```