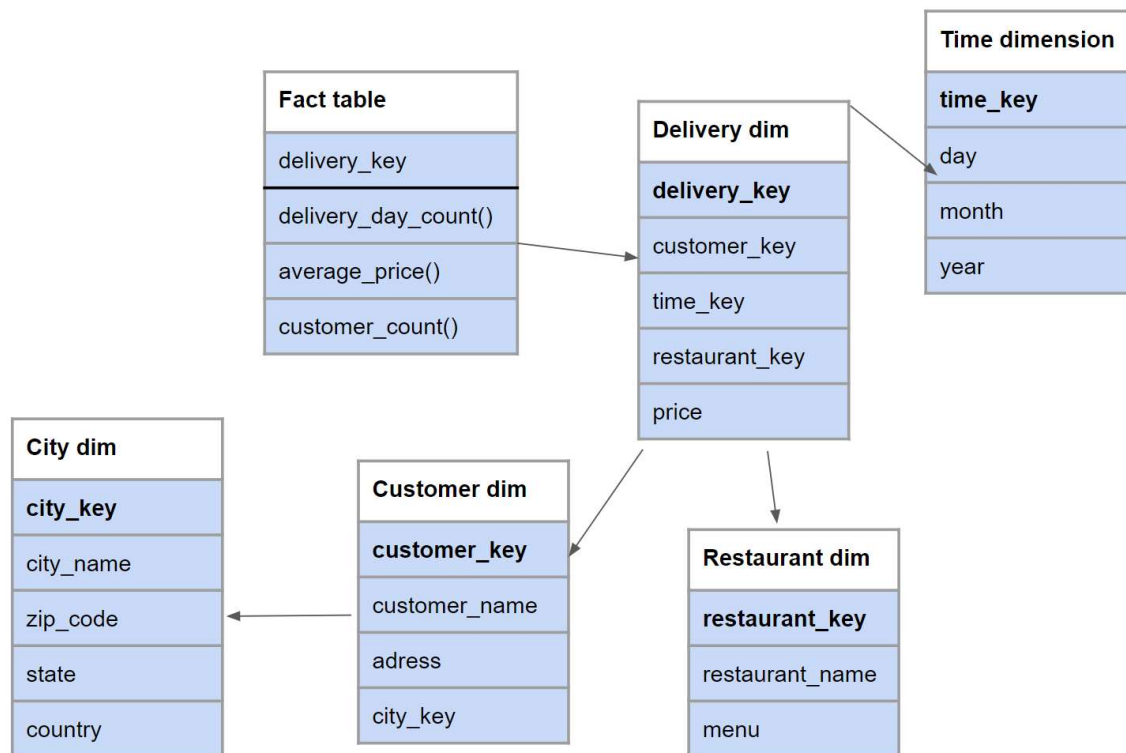


TDT4300 Exercise 5

Task 1: Datawarehousing

a)



- b) Only cuboid number 2 can be used to process the query. Since brand is less specific than item_name, the first cuboid cannot be used. Similarly, the third cuboid has country, which is less specific than city in the location dimension. Therefore you can't "roll up" to item_name and city in these two cuboids. In cuboid 2 you can "drill down" from street to city to get that information.

Task 2: Association Rules

Task 2: Apriori algorithm

A)

Minimum Support of 50% (i.e minimum support count is 4).
Using Fk-1 x Fk-1 method for candidate generation:

Initial Itemset		Initial C1 candidate set (K=1)		Frequent 1-element sets above min_sup		Frequent 2-element sets generated from the frequent 1-element set	
TID	Transaction	Itemsets	Sup_count	Itemsets	Sup_count	Itemsets	Sup_count
T1	BF	A	7	A	7	AB	5
T2	ABCDHF	B	6	B	6	AF	6
T3	ABF	C	1	F	7	AH	5
T4	ABFH	D	3	H	5	BF	4
T5	ADEF	E	2			BH	4
T6	ABFH	F	7			FH	4
T7	ABDEFH	G	1				
T8	AGH	H	5				

Frequent 3-element sets generated from the frequent 2-element sets. K = 3, so merging if their 3-2 = 1 first elements are identical.

		Itemsets	Sup_count
merge(AB,AF)	ABF	ABF	4
merge(AB,AH)	ABH	ABH	4
merge(AF,AH)	AFH	AFH	4
merge(BF, BH)	BFH	BFH	4

Frequent 4-element sets generated from the frequent 3-element sets. K = 4, so merging if their 4-2 = 2 first elements are identical.

		Itemset	Sup_count
merge(ABF,ABH)	ABFH	ABFH	4

B)

Rule generation for element set {ABH}

Rule	Confidence
{A} -> {BH}	4/7 = 57%
{B} -> {AH}	4/6 = 67%
{H} -> {AB}	4/5 = 80%
{AB} -> {H}	4/5 = 80%
{AH} -> {B}	4/5 = 80%
{BH} -> {A}	4/4 = 100%

Based on the confidence threshold $c = 75\%$, we can see 4 association rules (marked in green)

The numbers from the confidence calculation results from the formula

$$c(X \rightarrow Y) = \sigma(X \cup Y) / \sigma(X)$$

Task 3: Decision Trees

- a) Tree pruning is useful in decision tree induction because it can remove anomalies in the training dataset due to noise. This reduces the complexity of the decision tree and the amount of decisions that has to be made.
- One of the drawbacks of using a separate set of tuples to evaluate pruning is that there is a possibility that the tuples are not representative of the training tuples that were used in order to create the decision tree. If so, using

them to evaluate the accuracy of the pruned tree would not be a good indicator.

b) The stopping conditions in decision tree classification is as follows:

- i) Stop expanding a node when all the records belong to the same class.
- ii) Stop expanding a node when all the records have similar attribute values.
- iii) Early termination which is used to avoid too complex trees (i.e avoid overfitting). This can be done by:
 - 1) Stop if the number of instances is less than some user-specified threshold.
 - 2) Stop if class distribution of instances are independent of the available features, for example by using a χ^2 test.
 - 3) Stop if expanding the current node does not improve impurity measures. For example if the Gini-value or information gain is used as an impurity measure and it does not improve.

c)

Task 3: Decision Tree

D)

PC on credit:

Yes	12
No	8

$$\text{Entropy} = - \sum_{i=0}^{c-1} p_i(t) \cdot \log_2(t).$$

$E(\text{PC on credit}) = 0,971$

		PC on credit	
		Yes	No
Age	Young	4	3
	Middle	5	1
	Old	3	4

$E(\text{PC on credit, Age}) = P(\text{Young}) E(4,3) + P(\text{Middle}) E(5,1) + P(\text{Old}) E(3,4)$
 $= 7/20 * 0,985 + 6/20 * 0,650 + 7/20 * 0,985$
 $= 0,885$

		PC on credit	
		Yes	No
Income	Low	4	2
	Medium	6	3
	High	2	3

$E(\text{PC on credit, Income}) = 0,931$

		PC on credit	
		Yes	No
Married	Yes	9	3
	No	3	5

$E(\text{PC on Credit, Married}) = 0,868$

		PC on credit	
		Yes	No
Student	Yes	8	4
	No	4	4

$E(\text{PC on credit, Student}) = 0,951$

		PC on credit	
		Yes	No
Creditworth	Pass	6	4
	High	6	4

$E(\text{PC on credit, Creditworthiness}) = 0,971$

	Information gain
Age	$0,971 - 0,885 = 0,086$
Income	$0,971 - 0,931 = 0,040$
Married	$0,971 - 0,868 = 0,103$
Student	$0,971 - 0,951 = 0,020$
Creditworth	$0,971 - 0,971 = 0$

d)

- e) Since the attribute Married gives highest information gain, this attribute should be chosen as a splitting attribute.

Task 4: Data Types

(a) Time in terms of AM and PM.

- **Binary** (either AM or PM), **qualitative, ordinal** (since AM comes before PM on a specific day)

(b) Brightness as measured by a light meter.

- **Discrete** (not unlimited decimals or values), **quantitative, interval** (i.e 3.54 lumen defines interval from 3.535 to 3.5449)

(c) Brightness as measured by people's judgments.

- **Discrete** (not an unlimited amount of descriptions), **qualitative, nominal** ("very bright")

(d) Angles as measured in degrees between 0 and 360.

- **Discrete** (if you don't count decimals, there are 360 different values),
quantitative, ratio

(e) Bronze, Silver, and Gold medals as awarded at the Olympics.

- **Discrete, qualitative** (not measureable by nominal), **ordinal** (gold is better than silver)

(f) Height above sea level.

- Could argue for continuous and discrete, but in "day to day" speech or in writing it is often referred to in meters, so we'll say **discrete. Quantitative, ratio** (100m is twice the height of 50m etc).

(g) Number of patients in a hospital.

- **Discrete** (finite number). **Quantitative. Ratio.**

(h) ISBN numbers for books. (Look up the format on the Web.)

- **Discrete** (Limited digits), **qualitative, nominal** (code, numbers do not have a meaning)

(i) Ability to pass light in terms of the following values: opaque, translucent, transparent.

- **Discrete, qualitative, ordinal** (can rank by least to most transparent etc)

(j) Military rank.

- **Discrete** (finite number of ranks), **qualitative, ordinal** (major > leithenant)

(k) Distance from the center of campus.

- **Discrete** if measured in meters or kilometres. **Quantitative and ratio** (1km is twice as far as 500m).

(l) Density of a substance in grams per cubic centimeter.

- **Continuous** (infinite density measures), **quantitative, ratio** (water is twice as heavy as ...)

(m) Coat check number. (When you attend an event, you can often give your coat to someone who, in turn, gives you a number that you can use to claim your coat when you leave.)

- **Discrete** (numbers 1,2,3 etc.), **qualitative, nominal** (represents a code, not a meaningful order).

Task 5: Autocorrelation

Daily temperature likely has a higher temporal autocorrelation. It is unlikely that the temperature is doubling overnight, but that is quite common with rainfall on a day-to-day basis.

Task 6: Noise and Outliers

- a) Noise is defined as data that should be ignored, and does not contain any valuable information. It is desirable to reduce the noise as much as possible, if the model allows it. Outliers can be valid data points even though they do not fit within the scope of regular data points. Therefore outliers are most times not desirable, but could be interesting for further improvements, analyses etc to the model.
- b) It is possible that some instances of noise could appear as outliers, but it is not a common thing.
- c) No, sometimes noise objects are in the scope of the model (i.e appear as normal data) and are therefore not always outliers.
- d) No, sometimes legitimate data will not fit with the model and the rest of the data, and will therefore appear as an outlier, but it is still not noise.
- e) Yes, with enough noise, some typical values could seem unusual, or for instance outliers as typical data.

Task 7: Similarity Measures

Below you can see the code implementation of the different similarity measures and the result for the different tasks can be found in the last picture.

```
1  import math
2
3  # cosine similarity
4  def cos(x, y):
5      dot_product = 0
6      lenght_x = 0
7      lenght_y = 0
8      for val1, val2 in zip(x, y):
9          dot_product += val1*val2
10         lenght_x += val1**2
11         lenght_y += val2**2
12     return dot_product / (math.sqrt(lenght_x) * math.sqrt(lenght_y))
13
14 # Mean of a vector
15 def mean(vector):
16     n = len(vector)
17     summen = sum(vector)
18     return summen/n
19
20 # Standard deviation
21 def std_dev(vector):
22     s = 0
23     m = mean(vector)
24     for i in vector:
25         s += (i-m)**2
26     return math.sqrt(s/(len(vector)-1))
27
28
29 # Correlation
30 def correlation(x, y):
31     s_x = std_dev(x)
32     s_y = std_dev(y)
33     m_x = mean(x)
34     m_y = mean(y)
35     n = len(x) - 1
36     cov = 0
37
38     for x_i, y_i in zip(x, y):
39         c_x = (x_i-m_x)
40         c_y = (y_i - m_y)
41         cov += c_x * c_y
42     return cov/(n*s_x*s_y)
43
```

```
50
51 # Euclidean distance
52 v def euclidean(x, y):
53     d = 0
54 v     for i, j in zip(x, y):
55         d += (i-j)**2
56     return math.sqrt(d)
57
58 # Jaccard_coefficient
59 v def jaccard(x, y):
60     one_matches = 0
61     non_zero_matches = 0
62 v     for i, j in zip(x, y):
63 v         if i and j == 1:
64             one_matches += 1
65             non_zero_matches += 1
66 v         elif i and j == 0:
67             continue
68 v         else:
69             non_zero_matches += 1
70
71     return one_matches/non_zero_matches
72
```



```
76 # Defining vectors:
77 a_x, a_y = [1, 1, 1, 1], [2, 2, 2, 2]
78 b_x, b_y = [0, 1, 0, 1], [1, 0, 1, 0]
79 c_x, c_y = [0, -1, 0, 1], [1, 0, -1, 0]
80 d_x, d_y = [1, 1, 0, 1, 0, 1], [1, 1, 1, 0, 0, 1]
81 e_x, e_y = [2, -1, 0, 2, 0, -3], [-1, 1, -1, 0, 0, -1]
82
83 # Task a)
84 print("A ")
85 print("Cosine similarity = ", cos(a_x, a_y))
86 print("Correlation between X and Y is undefined")
87 print("Euclidean distance = ", euclidean(a_x, a_y))
88 print("")
89 print("-----")
90
91 # Task B)
92 print("B ")
93 print("Cosine similarity = ", cos(b_x, b_y))
94 print("Correlation = ", correlation(b_x, b_y))
95 print("Euclidean distance = ", euclidean(b_x, b_y))
96 print("Jaccard = ", jaccard(b_x, b_y))
97 print("")
98 print("-----")
99
100 # Task C)
101 print("Task C)")
102 print("Cosine similarity = ", cos(c_x, c_y))
103 print("Correlation = ", correlation(c_x, c_y))
104 print("Euclidean distance = ", euclidean(c_x, c_y))
105 print("")
106 print("-----")
107
108 # Task D)
109 print("Task D)")
110 print("Cosine similarity = ", cos(d_x, d_y))
111 print("Correlation = ", correlation(d_x, d_y))
112 print("Jaccard = ", jaccard(d_x, d_y))
113 print("")
114 print("-----")
115
116 # Task E)
117 print("Task E)")
118 print("Cosine similarity = ", cos(e_x, e_y))
119 print("Correlation = ", correlation(e_x, e_y))
120
```

```
A)
Cosine similarity = 1.0
Correlation between X and Y is undefined
Euclidean distance = 2.0
```

```
-----
B)
Cosine similarity = 0.0
Correlation = -1.0000000000000002
Euclidean distance = 2.0
Jaccard = 0.0
```

```
-----
Task C)
Cosine similarity = 0.0
Correlation = 0.0
Euclidean distance = 2.0
```

```
-----
Task D)
Cosine similarity = 0.75
Correlation = 0.24999999999999997
Jaccard = 0.6
```

```
-----
Task E
Cosine similarity = 0.0
Correlation = -5.73316704659901e-17
```