PROJECT JUPITER KONSEP DAN APLIKASI DATA MINING



OLEH: I'ROFUL BARIYAH 17.51.0004

STMIK PPKIA PRADNYA PARAMITA

MALANG

2020

1. Soal no1



```
In [14]: print(PLowYes)
          0.7407407407407407
In [15]: print(PMediumYes)
          0.18518518518518517
In [16]: print(PHigh)
          0.21568627450980393
In [17]: print(PLow)
          0.4117647058823529
In [18]: print(PMedium)
          0.37254901960784315
In [19]: #credit rating with student
pd.crosstab(df['Credit_rating'], df['Student'])
Out[19]:
               Student No Yes
          Credit_rating
           Excellent 8 12
                 Fair 16 15
In [20]: PExcellentNo = 8/24
PFairNo = 16/24
          PExcellentYes = 12/27
PFairYes = 15/27
         PExcellent = 20/51
PFair = 31/51
          print(PExcellentNo)
          0.3333333333333333
In [21]: print(PFairNo)
          0.66666666666666
In [22]: print(PExcellentYes)
          0.444444444444444
In [23]: print(PFairYes)
          0.55555555555556
In [24]: print(PExcellent)
          0.39215686274509803
In [25]: print(PFair)
          0.6078431372549019
In [26]: #income with class(buy_computer)
pd.crosstab(df['Income'], df['Class (buy_computer)'])
Out[26]:
          Class (buy_computer) No Yes
                      Income
           High 6 5
                       Low 11 10
          Medium 5 14
                                                                                                                            Activate Windows
```

```
In [27]: PHighNo = 6/22
PLowNo = 11/22
PMediumNo= 5/22
            PHighYes = 5/29
PLowYes = 10/29
PMediumYes = 24/29
            PHigh = 11/51
PLow = 21/51
PMedium = 19/51
             print(PHighNo)
             0.2727272727272727
 In [28]: print(PLowNo)
 In [29]: print(PMediumNo)
             0.22727272727272727
 In [30]: print(PHighYes)
             0.1724137931034483
In [31]: print(PLowYes)
           0.3448275862068966
In [32]: print(PMediumYes)
           0.8275862068965517
In [33]: #credit rating with class(buy_computer)
pd.crosstab(df['Credit_rating'], df['Class (buy_computer)'])
Out[33]:
            Class (buy_computer) No Yes
                   Credit_rating
            Excellent 8 12
                            Fair 14 17
In [34]: PExcellentNo = 8/22
PFairNo = 14/22
            PExcellentYes = 12/29
PFairYes = 17/29
            PExcellent = 20/51
PFair = 31/51
            print(PExcellentNo)
            0.36363636363636365
 In [35]: print(PFairNo)
            0.6363636363636364
 In [36]: print(PExcellentYes)
            0.41379310344827586
In [37]: print(PFairYes)
            0.5862068965517241
```

2. Soal no2

a. Apabila Cuaca buruk dengan nilai = 1, Weekday, dan Game = 0, maka berapa roti yang harus dibuat?

```
In [1]: import pandas as pd
  import numpy as np
  import matplotlib.pyplot as plt
       %matplotlib inline
In [2]: data=pd.read_excel('E:/KULIAH/SEMESTER 6/data mining/uas/data set2 ke 2/dataset_soal No. 2.xls')
In [3]: data
Out[3]:
          Category weatherv-1 holidayv-2 gamev-3 Qty
       0
              A 5 1
                                     0 250
              С
                                      0 75
       2
              D
                                      1 400
                                     0 150
                                      0 50
        5
                               0
       import math
       In [5]: data['dis'] = dis
Out[5]:
              A 5 1
                                       0 250 4.000000
               В
                        3
                                       1 200 2.236068
               D
                                       1 400 3.162278
               Е
                               0
                                       0 150 3.162278
                                       0 50 1.414214
In [6]: data.to_excel('E:/KULIAH/SEMESTER 6/data mining/uas/project uas/soal_no2.xls')
In [7]: data.to_excel('E:/KULIAH/SEMESTER 6/data mining/uas/project uas/soal_no2a.xls')
```

b. Apabila Cuaca baik dengan nilai 4, Weekend, dan Game =1, maka berapa roti yang harus dibuat?

```
In [1]: import pandas as pd
  import numpy as np
  import matplotlib.pyplot as plt
        %matplotlib inline
In [2]: data=pd.read_excel('E:/KULIAH/SEMESTER 6/data mining/uas/data set2 ke 2/dataset_soal No. 2.xls')
In [3]: data
Out[3]:
          Category weathery-1 holidayy-2 gamey-3 Qty
        0
               A 5 1
                                        0 250
                В
                                         1 200
        2 C
                                        0 75
               D
                                1
                                        1 400
        4 E 4 0 0 150
                                         0 50
```

```
In [5]: data['dis'] = dis
data
Out[5]:
        Category weatherv-1 holidayv-2 gamev-3 Qty
      0 A 5 1 0 250 1.414214
            В
                   3
                               1 200 1.000000
      1 B 3 1
2 C 1 1
                              0 75 3.162278
       3 D 4 1 1 400 0.000000
4 E 4 0 0 150 1.414214
      3
       5
                   2
                         0
                               0 50 2.449490
In [6]: data.to_excel('E:/KULIAH/SEMESTER 6/data mining/uas/project uas/soal_no2b.xls')
```

3. Jawaban No 3

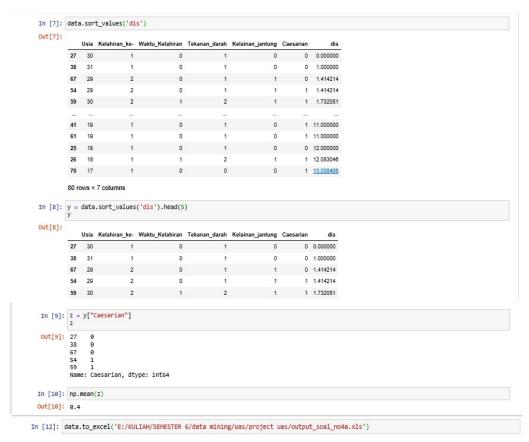
```
In [1]: import numpy as np import pandas as pd
           from apyori import apriori
In [2]: store_data = pd.read_excel('E:/KULIAH/SEMESTER 6/data mining/uas/dataset_soal No. 3.xls')
In [3]: store_data.head()
Out[3]:
                                         Item2
                                                     Item3
                                                                        Item4
                                                                                    Item5
                                                                                                  Item6
                                                                                                                  Item7
                                                                                                                                     Item8
                                                                                                                                                  Item9
                                                                                                                                                              Item10
           0
                    burgers
                                     meatballs
                                                     eggs
                                                                 low fat yogurt
                                                                                     NaN mineral water
                                                                                                                 salmon
                                                                                                                               low fat yogurt
                                                                                                                                                  NaN mineral water
                                                      NaN whole wheat pasta french fries mineral water
                                                                                                                  salmon whole wheat pasta french fries mineral water
                    chutney
                                  low fat yogurt
           3 mineral water
                                         soup light cream frozen vegetables
                                                                                 spaghetti
                                                                                               green tea
                                                                                                                   NaN frozen vegetables
                                                                                                                                             spaghetti
                                                                                                                                                             green tea
                                                                                                                                                            chocolate
            4 low fat yogurt frozen vegetables spaghetti
                                                                                             chocolate frozen smoothie
                                                                   french fries
                                                                                    eggs
                                                                                                                                french fries
                                                                                                                                                  eggs
In [4]: store_data.tail()
Out[4]:
                                                            Item3
                                                                                                                                  turkey
                                                                                                                                                               french fries
            2049
                             burgers
                                             eggs
                                                        french fries
                                                                     fresh tuna spaghetti
                                                                                              olive oil clothes accessories
            2050
                            burgers
                                             eggs frozen smoothie french wine
                                                                                     eaas french fries
                                                                                                              energy drink
                                                                                                                              french fries
                                                                                                                                                       NaN
                                                                                                                                                                 chocolate
                                                                                                                                                       milk herb & pepper
            2051 whole wheat pasta
                                                           melons champagne pancakes
                                                                                            light mayo
            2052
                        ground beef tomato sauce
                                                          spaghetti
                                                                       red wine
                                                                                   honey
                                                                                             hot doas
                                                                                                                   turkey herb & pepper whole wheat pasta
                                                                                                                                                             mineral water
            2053
                            burgers
                                            egas frozen smoothie
                                                                           milk
                                                                                                               french fries mineral water
                                                                                                                                                                   cookies
                                                                                   bacon
                                                                                                 eaas
                                                                                                                                                   avocado
In [ ]:
In [5]: store_data.shape
Out[5]: (2054, 10)
In [6]: records = []
           for i in range(0, 2054):
    records.append([str(store_data.values[i,j]) for j in range(0,10)])
In [7]: association_rules = apriori(records, min_support=0.2, min_confidence=0.2, min_lift=0.2, min_lenght=2) association_result = list(association_rules)
In [8]: print(len(association_result))
                                                                                                                                                         Go to Settings to activate
           61
In [9]: print(association result[0])
           Relation Record (items=frozenset (\{ 'avocado' \}), support=0.314508276533593, ordered\_statistics=[OrderedStatistic(items\_base=frozenset(), items\_add=frozenset(\{ 'avocado' \}), confidence=0.314508276533593, lift=1.0)])
In [13]: result =[]
for item in association_result:
                 pair = item[0]
items = [x for x in pair]
                 value0 = str(items[0])
                value0 = str(items[0])
value1 = str(item[1])
value2 = str(item[1])[:10]
value3 = str(item[2][0][2])[:10]
value4 = str(item[2][0][3])[:10]
                 rows = (value0,value1,value2,value3,value4)
                 result.append(rows)
                 label = ['title1', 'title2', 'support', 'confidence', 'lift']
                 store_suggestion = pd.DataFrame.from_records(result,columns=label)
                 print(store suggestion)
                            0.3145082/65335
                 title1
                              title2 support confidence lift
0.314508276533593 0.31450827 0.31450827 1.0
                avocado
                            0.24294060370009737
                                                      0.24294060 0.24294060
               burgers
                                title2 support confidence lift
0.314508276533593 0.31450827 0.31450827 1.0
                    title1
                   avocado
                                                                                                                                                          Activate Windows
                  burgers 0.24294060370009737
                                                         0.24294060
                                                                        0.24294060 1.0
                                             41382668 0.47565725 0.47565725 1.0
title2 support confidence lift
0.314508276533593 0.31450827 0.31450827 1.0
                               0.4756572541382668
title1
                                                                                                                                                          Go to Settings to activate
                               avocado
```

```
burgers 0.24294060370009737 0.24294060 0.24294060 1.0
title1 title2 support confidence lift
0 avocado 0.314508276533593 0.31450827 0.31450827 1.0
burgers 0.24294060370009737 0.24294060 0.24294060 1.0
            0 avocado
1 burgers
2 chocolate
                               25 1.0 confidence lift 0.31450827 1.0 0.24294060 1.0 0.47565725 1.0 confidence lift 0.31450827 1.0 0.24294060 1.0 0.47565725 1.0 0.35881207 1.0 0.35881207 1.0 confidence lift
            0
                                                          title2
                                                                          support
                                                                                          confidence lift
                                  title1
In [14]: store_suggestion.describe()
Out[14]:
                      title1
                                               title2
                                                         support confidence lift
                                                61
                                                                           61 61
              count 61
                                                             61
              unique
                         15
                                                 53
                                                              53
                                                                            53
             top nan 0.24294060370009737 0.24294060 0.24294060 1.0
                freq 19
                                                                            4 61
                                                                                                                                                                 Go to Settings to activate
In [15]: store_suggestion.to_excel('E:/KULIAH/SEMESTER 6/data mining/uas/project uas/output_soal3.xls')
```

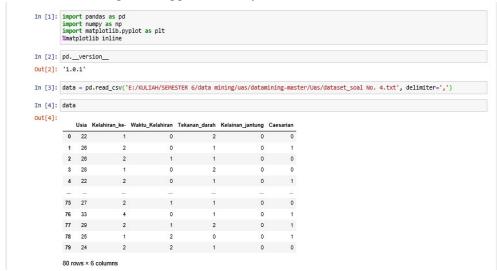
4. Jawaban No 4

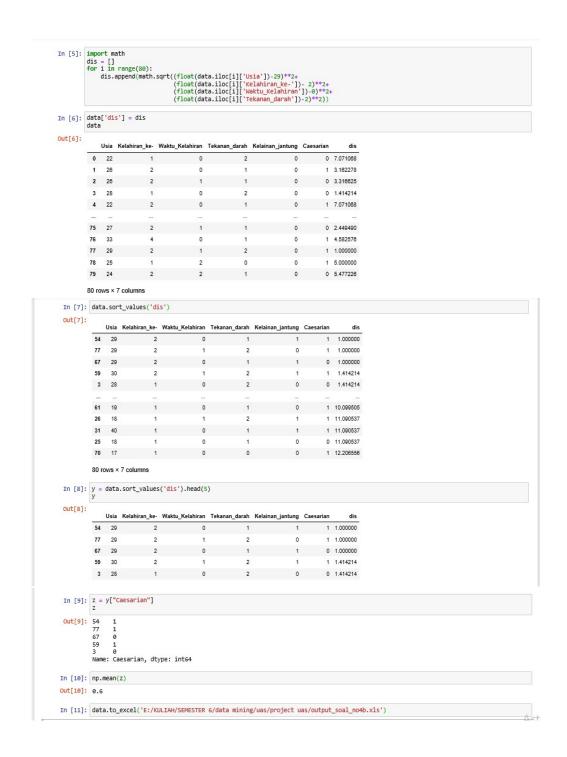
a. Berdasarkan data tersebut bagaimana perlakuan dengan kondisi Ibu hamil dengan Usia 30 Tahun, yang merupakan Kelahiran ke -1, dengan Waktu kelahiran sesuai dengan HPL, Memiliki tekanan darah Normal? Carilah KNN dengan menggunakan Key = 5





b. Bagaimana Apabila Ibu hamil dengan Usia 29 Tahun, yang merupakan Kelahiran ke 2, dengan Waktu kelahiran sesuai dengan HPL, Memiliki tekanan darah Tinggi?
 Carilah KNN dengan menggunakan Key =5





https://github.com/irafulbariyah/project uas