

MODUL PRAKTIKUM
MATA KULIAH: DATA MINING
MODUL VIII
[Algoritma Apriori dalam Analisis Asosiasi]



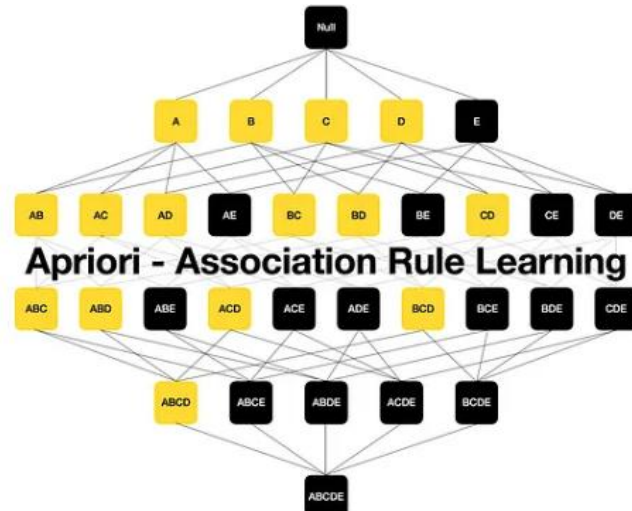
Program Studi Sains Data
Fakultas Sains
INSTITUT TEKNOLOGI SUMATERA
Tahun Ajaran Ganjil 2024/2025.

1. Tujuan Praktikum

- Dapat memahami algoritma apriori untuk *asosiasi rule learning*.
- Dapat mengimplementasikan metode apriori dalam algoritma python.

2. Landasan teori

Algoritma Apriori adalah metode data mining yang digunakan untuk menemukan aturan asosiasi dalam data transaksi, seperti pola belanja pelanggan. Algoritma Apriori digunakan untuk menemukan hubungan antar item dalam dataset besar.



Gambar 1 Skema apriori- Association Rule Learning.

Association Rule Learning

Association Rule Learning atau pembelajaran aturan asosiasi adalah metode pembelajaran mesin berbasis aturan untuk menemukan hubungan menarik antara variabel dalam basis data besar. Association rule mining adalah teknik yang digunakan untuk mengidentifikasi pola dalam kumpulan data besar. Teknik ini melibatkan pencarian hubungan antara variabel dalam data dan menggunakan hubungan tersebut untuk membuat prediksi atau keputusan. Tujuan penambangan aturan asosiasi adalah untuk mengungkap aturan yang menggambarkan hubungan antara berbagai item dalam dataset. Biasanya association rule learning terdiri dari algoritma untuk menganalisis data dan mengidentifikasi hubungan.

Algoritma ini dapat didasarkan pada metode statistik atau teknik pembelajaran mesin. Aturan yang dihasilkan sering kali dinyatakan dalam bentuk pernyataan "jika-maka", di mana bagian "jika" mewakili anteseden (kondisi yang diuji) dan bagian "maka" mewakili konsekuen (hasil yang terjadi jika kondisi terpenuhi). Hal ini dapat berguna untuk berbagai keperluan, seperti mengidentifikasi tren pasar, mendeteksi aktivitas penipuan, atau memahami perilaku pelanggan. Penambangan aturan asosiasi juga dapat digunakan sebagai batu loncatan untuk jenis analisis data lainnya, seperti memprediksi hasil atau mengidentifikasi pendorong utama fenomena tertentu.

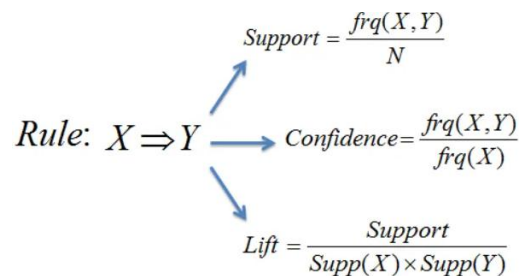
Use Cases of Association Rule Mining

Penambangan aturan asosiasi umumnya digunakan dalam berbagai aplikasi, beberapa yang umum adalah:

- Analisis Keranjang Belanja
- Customer Segmentation
- Fraud Detection
- Social network analysis
- Recommendation systems

Metrics for Evaluating Association Rules

Dalam penambangan aturan asosiasi, beberapa metrik umumnya digunakan untuk mengevaluasi kualitas dan pentingnya aturan asosiasi yang ditemukan. Dalam penambangan aturan asosiasi, beberapa metrik umumnya digunakan untuk mengevaluasi kualitas dan pentingnya aturan asosiasi yang ditemukan. Menafsirkan hasil metrik penambangan aturan asosiasi memerlukan pemahaman tentang makna dan implikasi setiap metrik, serta cara menggunakannya untuk mengevaluasi kualitas dan pentingnya aturan asosiasi yang ditemukan. Berikut adalah beberapa panduan untuk menafsirkan hasil metrik penambangan aturan asosiasi utama:



Gambar 2 matriks evaluasi untuk Association Rules

Dari gambar diatas dapat dijabarkan dibawah ini:

- Support

Support merupakan Ukuran seberapa sering item atau set item muncul dalam set data. Nilai ini dihitung sebagai jumlah transaksi yang memuat item dibagi dengan jumlah total transaksi dalam kumpulan data. Nilai dukungan yang tinggi menunjukkan bahwa suatu item atau kumpulan item umum dalam kumpulan data, sedangkan nilai dukungan yang rendah menunjukkan bahwa item atau kumpulan item tersebut jarang.

$$Support(\{X\} \rightarrow \{Y\}) = \frac{Transactions\ containing\ both\ X\ and\ Y}{Total\ number\ of\ transactions}$$

Contoh: Misalnya dari 1000 transaksi, 100 transaksi berisi Saus Tomat(*Ketchup*), maka Support untuk item Saus Tomat dapat dihitung sebagai berikut:

$$\begin{aligned} Support(Ketchup) &= (Transactions\ containing\ Ketchup) / (Total\ Transactions) \\ Support(Ketchup) &= 100/1000 \\ &= 10\% \end{aligned}$$

- Confidence(Keyakinan)

Confidence adalah Ukuran kekuatan hubungan antara dua item. Nilai ini dihitung sebagai jumlah transaksi yang memuat kedua item dibagi dengan jumlah transaksi yang memuat item pertama. Confidence yang tinggi menunjukkan bahwa keberadaan item pertama merupakan prediktor kuat keberadaan item kedua.

$$Confidence(\{X\} \rightarrow \{Y\}) = \frac{\text{Transactions containing both } X \text{ and } Y}{\text{Transactions containing } X}$$

Contoh: Kembali ke masalah kita, kita memiliki 50 transaksi di mana Burger dan Saus Tomat dibeli bersamaan. Sementara dalam 150 transaksi, burger dibeli. Maka kita dapat menemukan kemungkinan membeli saus tomat ketika burger dibeli dapat direpresentasikan sebagai confidence Burger -> Saus Tomat dan dapat ditulis secara matematis sebagai:

$$Confidence(Burger \rightarrow Ketchup) = \frac{\text{(Transactions containing both (Burger and Ketchup))}}{\text{(Transactions containing A)}}$$

$$\begin{aligned} Confidence(Burger \rightarrow Ketchup) &= 50/150 \\ &= 33.3\% \end{aligned}$$

- Lift

Lift adalah ukuran kekuatan hubungan antara dua item, dengan mempertimbangkan frekuensi kedua item dalam kumpulan data. Nilai ini dihitung sebagai Confidence hubungan dibagi dengan dukungan item kedua. Lift digunakan untuk membandingkan kekuatan hubungan antara dua item dengan kekuatan hubungan yang diharapkan jika item-item tersebut independen.

$$Lift(\{X\} \rightarrow \{Y\}) = \frac{\text{(Transactions containing both } X \text{ and } Y) / \text{(Transactions containing } X)}{\text{Fraction of transactions containing } Y}$$

Contoh: jika kita menggunakan data dari masalah diatas terkait dengan saus tomat dan burger maka kita akan dapat menghitung nilai lift nya yaitu:

$$\begin{aligned} Lift(Burger \rightarrow Ketchup) &= \text{(Confidence(Burger} \rightarrow \text{Ketchup))} / \text{(Support(Ketchup))} \\ Lift(Burger \rightarrow Ketchup) &= 33.3/10 \\ &= 3.33 \end{aligned}$$

3. Prosedur praktikum

Untuk dataset modul 8 dapat diakses, link: https://bit.ly/damin_datamodul8.

Install efficient-apriori.

```
1 pip install efficient-apriori
```

Setelah itu kita tuliskan library yang kita butuhkan untuk algoritma apriori ini

```

1 import pandas as pd # for data manipulation
2 import matplotlib.pyplot as plt # for plotting frequency distribution chart
3 from efficient_apriori import apriori # for association analysis
4 import sys
5 import os
6 # Assign main directory to a variable
7 main_dir=os.path.dirname(sys.path[0])

```

Masukan data [Market Basket Optimisation.csv](#) yang telah diunduh.

```

1 # Ingest the data
2 df = pd.read_csv(main_dir+'//content/Market_Basket_Optimisation.csv',
3                  encoding='utf-8', header=None)
4 # Show dataframe
5 df

```

Exploration

Sebelum kita menjalankan analisis aturan asosiasi, mari kita lihat terlebih dahulu distribusi frekuensi item.

```

1 # Put all transactions into a single list
2 txns=df.values.reshape(-1).tolist()
3 # Create a dataframe using this single list and add a column for count
4 df_list=pd.DataFrame(txns)
5 df_list['Count']=1
6 # Group by items and rename columns
7 df_list=df_list.groupby(by=[0], as_index=False).count().sort_values(by=['Count'], ascending=True) # count
8 df_list['Percentage'] = (df_list['Count'] / df_list['Count'].sum()) # percentage
9 df_list=df_list.rename(columns={0 : 'Item'})
10 # Show dataframe
11 df_list

```

Dari output yang diperoleh tahap exploitasi menunjukkan bahwa air mineral (1.788) sejauh ini merupakan barang paling populer yang dibeli di supermarket ini, diikuti oleh telur (1.348), spageti (1.306), kentang goreng (1.282), dan cokelat (1.230). Sementara itu, asparagus tua yang malang hanya dibeli satu kali. Penting juga untuk dicatat bahwa bahkan barang yang paling sering dibeli hanya muncul dalam lebih dari 6% transaksi. Kita dapat menggunakan informasi ini sebagai panduan saat menetapkan ambang batas dukungan minimum/ *minimum support threshold*.

Visualisasi distribusi frekuensi dalam diagram batang sebagai berikut:

```

1 # Draw a horizontal bar chart
2 plt.figure(figsize=(16,20), dpi=300)
3 plt.ylabel('Item Name')
4 plt.xlabel('Count')
5 plt.barh(df_list['Item'], width=df_list['Count'], color='black', height=0.8)
6 plt.margins(0.01)
7 plt.show()

```

Algoritma Apriori

Sebelum kita menjalankan algoritma, mari kita masukkan data kita dalam format yang diperlukan. Ini akan menjadi daftar daftar yang telah menghapus semua "NaN" yang Anda lihat di Pandas DataFrame.

```
1 # Create a list of lists from a dataframe
2 txns2=df.stack().groupby(level=0).apply(list).tolist()
3 # Show what it looks like
4 txns2
```

Terakhir, mari kita jalankan algoritma Apriori dan simpan itemset dan aturan asosiasi.

```
1 itemsets, rules = apriori(txns2, min_support=0.03, min_confidence=0.2, verbosity=1)
```

Algoritme tersebut telah menemukan 36 itemset berukuran 1 dan 18 itemset berukuran 2, yang memenuhi ambang batas dukungan minimum sebesar 0,03. Itemset tersebut ditampilkan di bawah ini.

Sekarang mari kita cetak aturan asosiasi yang ditemukan oleh algoritma. Perhatikan, beberapa set dikecualikan dalam tahap pembuatan aturan karena tidak memenuhi persyaratan keyakinan minimum yang kami tentukan (dalam contoh ini, min_confidence=0.2).

```
1 for item in sorted(rules, key=lambda item: (item.lift,item.conviction), reverse=True):
2     print(item)
```

Dari hasil algoritma diatas kita dapat melihat bahwa Seperti peningkatan tertinggi dalam data kehidupan nyata ini adalah untuk kombinasi daging sapi giling dan spageti. Namun, pembeli daging sapi giling lebih mungkin untuk juga membeli spageti(confidence: 0.399, conviction: 1.374) daripada sebaliknya(confidence: 0.225, conviction: 1.164).

Latihan

Dalam *Apriori Algorithm for Association Rule Learning* kita dapat membagi kedalam beberapa tahapan:

- Libraries

```
[53] 1 import numpy as np
      2 import pandas as pd
      3 import seaborn as sns
      4 import matplotlib.pyplot as plt
      5
      6 import itertools
      7 from mlxtend.preprocessing import TransactionEncoder
      8
      9 import warnings
     10 warnings.filterwarnings("ignore")
```

- Import DataSet

dalam Latihan modul 8 digunakan dataset GroceryStoreDataSet.csv yang dapat di unduh di https://bit.ly/damin_datamodul8.

```
1 data =pd.read_csv('GroceryStoreDataSet.csv',header=None)
2 data.columns=["Product"]
3 data

1 print("Number Of Rows is :",data.shape[0])
2 print("Number Of Columns is :",data.shape[1])

1 data.values
```

- data preparation

dalam modul ini kita akan menggunakan algoritma Apriori untuk melakukan analisis relevansi. Untuk mengonversi data, saya akan kembali menggunakan metode konversi data dari library mlxtend

```
1 data =list(data["Product"].apply(lambda x:x.split(",")))
2 data
```

- TransactionEncoder

Dengan menggunakan fungsi ini, kita menggunakan nilai Boolean untuk setiap transaksi dan memasukkan nama setiap produk dalam kolom.

```
1 te =TransactionEncoder()
2 te_data =te.fit_transform(data)
3 df =pd.DataFrame(te_data ,columns =te.columns_)
4 df
```

```
df_zero =df.replace(True ,1)
df_zero
```

- Visualisasi dengan menggunakan contoh Bar Plot()
Kita mevisualisasikan dengan bar plot untuk melihat pembelian terbanyak

```
1 ax = sns.barplot(df_zero,estimator="sum", errorbar=None)
2 ax.bar_label(ax.containers[0], fontsize=10)
3 plt.xticks(rotation =90)
4 plt.show()
```

- support: menunjukkan persentase setiap pembelian produk

```
1 first = pd.DataFrame(df.sum() / df.shape[0],
2                       columns = ["Support"]).sort_values("Support", ascending = False)
3 first
```

Untuk pembelian produk diatas 15%

```
[93] 1 first[first['Support'] >= 0.15]
```

- Second Iteration: Temukan nilai dukungan untuk kombinasi produk berpasangan.

```
1 second = list(itertools.combinations(first.index, 2))
2 second = [list(i) for i in second]
3 # Sample of combinations
4 second
```

Mencari nilai value dan mengeliminasi dari nilai support

```
1 value = []
2 for i in range(0, len(second)):
3     temp = df.T.loc[second[i]].sum()
4     temp = len(temp[temp == df.T.loc[second[i]].shape[0]]) / df.shape[0]
5     value.append(temp)
6 # Create a data frame
7 secondIteration = pd.DataFrame(value, columns = ["Support"])
8 secondIteration["index"] = [tuple(i) for i in second]
9 secondIteration['length'] = secondIteration['index'].apply(lambda x:len(x))
10 secondIteration = secondIteration.set_index("index").sort_values("Support", ascending = False)
11 # Elimination by Support Value
12 secondIteration
```

```
1 df_second=pd.DataFrame()
2 df_second.insert(0 , "Name Product" ,tuple(second))
3 df_second.insert(1 , "Support" ,[0.20,0.20,0.20,0.20,0.20,0.20,0.20,
4     0.15,0.15,0.15,0.15,0.15,
5     0.10,0.10,0.10,0.10,0.10,0.10,0.10,0.10,0.10,0.10,0.10,0.10,
6     0.05,0.05,0.05,0.05,0.05,0.05,0.05,0.05,0.05,0.05,0.05,0.05,
7     0.00,0.00,0.00,0.00,0.00,0.00,0.00,0.00,0.00,0.00,0.00,0.00,0.00,0.00,0.00,0.00,
8     0.00,0.00,0.00])
9 df_second
```



```

1 plt.figure(figsize =(15 ,5))
2 x =range(55)
3 y =df_second["Support"]
4 labels =[second]
5
6 plt.bar(x ,y)
7 plt.xticks(x ,labels[0] ,rotation =90)
8 plt.title("Product and Support")
9 plt.xlabel("Product")
10 plt.ylabel("Support Of Product")
11 plt.show()

```

Iterasi pertama

```

1 def ar_iterations(data, num_iter = 1, support_value = 0.1, iterationIndex = None):
2
3     # Next Iterations
4     def ar_calculation(iterationIndex = iterationIndex):
5         # Calculation of support value
6         value = []
7         for i in range(0, len(iterationIndex)):
8             result = data.T.loc[iterationIndex[i]].sum()
9             result = len(result[result == data.T.loc[iterationIndex[i]].shape[0]]) / data.shape[0]
10            value.append(result)
11
12            # Bind results
13            result = pd.DataFrame(value, columns = ["Support"])
14            result["index"] = [tuple(i) for i in iterationIndex]
15            result['length'] = result['index'].apply(lambda x:len(x))
16            result = result.set_index("index").sort_values("Support", ascending = False)
17            # Elimination by Support Value
18            result = result[result.Support > support_value]
19
20            return result

```

```

# First Iteration
first = pd.DataFrame(df.T.sum(axis = 1) / df.shape[0], columns = ["Support"]).sort_values("Support", ascending = False)
first = first[first.Support > support_value]
first["length"] = 1

if num_iter == 1:
    res = first.copy()

# Second Iteration
elif num_iter == 2:
    second = list(itertools.combinations(first.index, 2))
    second = [list(i) for i in second]
    res = ar_calculation(second)

# All Iterations > 2
else:
    nth = list(itertools.combinations(set(list(itertools.chain(*iterationIndex))), num_iter))
    nth = [list(i) for i in nth]
    res = ar_calculation(nth)
return res

```

Iterasi I

```

1 iteration1 = ar_iterations(df, num_iter=1, support_value=0.1)
2 iteration1

```

Iterasi II

```
1 iteration2 = ar_iterations(df, num_iter=2, support_value=0.1)
2 iteration2
```

Iterasi III

```
1 iteration3 = ar_iterations(df, num_iter=3, support_value=0.01,
2                       iterationIndex=iteration2.index)
3 iteration3
```

Iterasi IV

```
1 iteration4 = ar_iterations(df, num_iter=4, support_value=0.01,
2                       iterationIndex=iteration3.index)
3 iteration4
```