

# Microsoft

April 22, 2024

```
[1]: import os
import tempfile
import shutil
import urllib
import zipfile
import pandas as pd
```

1.

```
[2]: behaviors_path = os.path.join('./input', 'behaviors.tsv')
behaviors_df=pd.read_table(
    behaviors_path,
    header=None,
    names=['impression_id', 'user_id', 'time', 'history', 'impressions'])
behaviors_df
```

```
[2]:
```

|        | impression_id | user_id | time                   | \                               |
|--------|---------------|---------|------------------------|---------------------------------|
| 0      | 1             | U13740  | 11/11/2019 9:05:58 AM  |                                 |
| 1      | 2             | U91836  | 11/12/2019 6:11:30 PM  |                                 |
| 2      | 3             | U73700  | 11/14/2019 7:01:48 AM  |                                 |
| 3      | 4             | U34670  | 11/11/2019 5:28:05 AM  |                                 |
| 4      | 5             | U8125   | 11/12/2019 4:11:21 PM  |                                 |
| ...    | ...           | ...     | ...                    |                                 |
| 156960 | 156961        | U21593  | 11/14/2019 10:24:05 PM |                                 |
| 156961 | 156962        | U10123  | 11/13/2019 6:57:04 AM  |                                 |
| 156962 | 156963        | U75630  | 11/14/2019 10:58:13 AM |                                 |
| 156963 | 156964        | U44625  | 11/13/2019 2:57:02 PM  |                                 |
| 156964 | 156965        | U64800  | 11/14/2019 3:25:49 PM  |                                 |
|        |               |         |                        | history \                       |
| 0      | N55189        | N42782  | N34694                 | N45794 N18445 N63302 N104...    |
| 1      | N31739        | N6072   | N63045                 | N23979 N35656 N43353 N8129...   |
| 2      | N10732        | N25792  | N7563                  | N21087 N41087 N5445 N60384...   |
| 3      | N45729        | N2203   | N871                   | N53880 N41375 N43142 N33013 ... |
| 4      |               |         |                        | N10078 N56514 N14904 N33740     |
| ...    |               |         |                        | ...                             |
| 156960 | N7432         | N58559  | N1954                  | N43353 N14343 N13008 N28833...  |

```

156961  N9803 N104 N24462 N57318 N55743 N40526 N31726 ...
156962  N29898 N59704 N4408 N9803 N53644 N26103 N812 N...
156963  N4118 N47297 N3164 N43295 N6056 N38747 N42973 ...
156964                                     N22997 N48742

```

```

                                impressions
0                                N55689-1 N35729-0
1      N20678-0 N39317-0 N58114-0 N20495-0 N42977-0 N...
2      N50014-0 N23877-0 N35389-0 N49712-0 N16844-0 N...
3                                N35729-0 N33632-0 N49685-1 N27581-0
4      N39985-0 N36050-0 N16096-0 N8400-1 N22407-0 N6...
...
156960  N2235-0 N22975-0 N64037-0 N47652-0 N11378-0 N4...
156961  N3841-0 N61571-0 N58813-0 N28213-0 N4428-0 N25...
156962  N55913-0 N62318-0 N53515-0 N10960-0 N9135-0 N5...
156963  N6219-0 N3663-0 N31147-0 N58363-0 N4107-0 N457...
156964                                     N61233-0 N33828-1 N19661-0 N41934-0

```

[156965 rows x 5 columns]

```

[3]: news_path = os.path.join('./input', 'news.tsv')
news_df=pd.read_table(news_path,
                      header=None,
                      names=[
                          'id', 'category', 'subcategory', 'title', 'abstract', 'url',
                          'title_entities', 'abstract_entities'
                      ])
news_df

```

```

[3]:      id  category  subcategory \
0    N55528  lifestyle  lifestyleroyals
1    N19639   health    weightloss
2    N61837    news     newsworld
3    N53526   health    voices
4    N38324   health    medical
...
51277  N16909   weather  weathertopstories
51278  N47585  lifestyle  lifestylefamily
51279   N7482   sports    more_sports
51280  N34418   sports    soccer_epl
51281  N44276   autos     autossports

```

```

                                title \
0    The Brands Queen Elizabeth, Prince Charles, an...
1                                50 Worst Habits For Belly Fat
2    The Cost of Trump's Aid Freeze in the Trenches...
3    I Was An NBA Wife. Here's How It Affected My M...

```

4       How to Get Rid of Skin Tags, According to a De...  
 ...  
 51277   Adapting, Learning And Soul Searching: Reflect...  
 51278   Family says 13-year-old Broadway star died fro...  
 51279   St. Dominic soccer player tries to kick cancer...  
 51280                   How the Sounders won MLS Cup  
 51281                   Best Sports Car Deals for October

abstract \

0       Shop the notebooks, jackets, and more that the...  
 1       These seemingly harmless habits are holding yo...  
 2       Lt. Ivan Molchanets peeked over a parapet of s...  
 3       I felt like I was a fraud, and being an NBA wi...  
 4       They seem harmless, but there's a very good re...  
 ...  
 51277   Woolsey Fire Anniversary: A community is forev...  
 51278                   NaN  
 51279   Sometimes, what happens on the sidelines can b...  
 51280   Mark, Jeremiah and Casey were so excited they ...  
 51281                   NaN

url \

0       <https://assets.msn.com/labs/mind/AAGHOET.html>  
 1       <https://assets.msn.com/labs/mind/AAB19MK.html>  
 2       <https://assets.msn.com/labs/mind/AAJgNsz.html>  
 3       <https://assets.msn.com/labs/mind/AACk2N6.html>  
 4       <https://assets.msn.com/labs/mind/AAAKEkt.html>  
 ...  
 51277   <https://assets.msn.com/labs/mind/BBWzQJK.html>  
 51278   <https://assets.msn.com/labs/mind/BBWzQYV.html>  
 51279   <https://assets.msn.com/labs/mind/BBWzQnK.html>  
 51280   <https://assets.msn.com/labs/mind/BBWzQuK.html>  
 51281   <https://assets.msn.com/labs/mind/BBY5rVe.html>

title\_entities \

0       [{"Label": "Prince Philip, Duke of Edinburgh",...  
 1       [{"Label": "Adipose tissue", "Type": "C", "Wik...  
 2   []  
 3   []  
 4       [{"Label": "Skin tag", "Type": "C", "WikidataI...  
 ...  
 51277   [{"Label": "Woolsey Fire", "Type": "N", "Wikid...  
 51278   [{"Label": "Broadway theatre", "Type": "F", "W...  
 51279   []  
 51280   [{"Label": "MLS Cup", "Type": "U", "WikidataId...  
 51281   [{"Label": "Peugeot RCZ", "Type": "V", "Wikida...

```

                                abstract_entities
0                                []
1    [{"Label": "Adipose tissue", "Type": "C", "Wik...
2    [{"Label": "Ukraine", "Type": "G", "WikidataId...
3    [{"Label": "National Basketball Association", ...
4    [{"Label": "Skin tag", "Type": "C", "WikidataI...
...
51277 [{"Label": "Woolsey Fire", "Type": "N", "Wikid...
51278                                []
51279                                []
51280                                []
51281                                []

[51282 rows x 8 columns]

```

2. :

```

[20]: all_categories = news_df['subcategory']

#
def find_frequent_itemsets(categories, min_support=0.01):
    category_counts = categories.value_counts(normalize=True)
    frequent_categories = category_counts[category_counts >= min_support].index.
    ↪tolist()
    return frequent_categories

#
frequent_categories = find_frequent_itemsets(all_categories)

#
print("Frequent categories:")
for category in frequent_categories:
    print(category)

```

```

Frequent categories:
newsus
football_nfl
newspolitics
newscrime
weathertopstories
newsworld
football_ncaa
baseball_mlb
basketball_nba
newsscienceandtechnology
news
newstrends

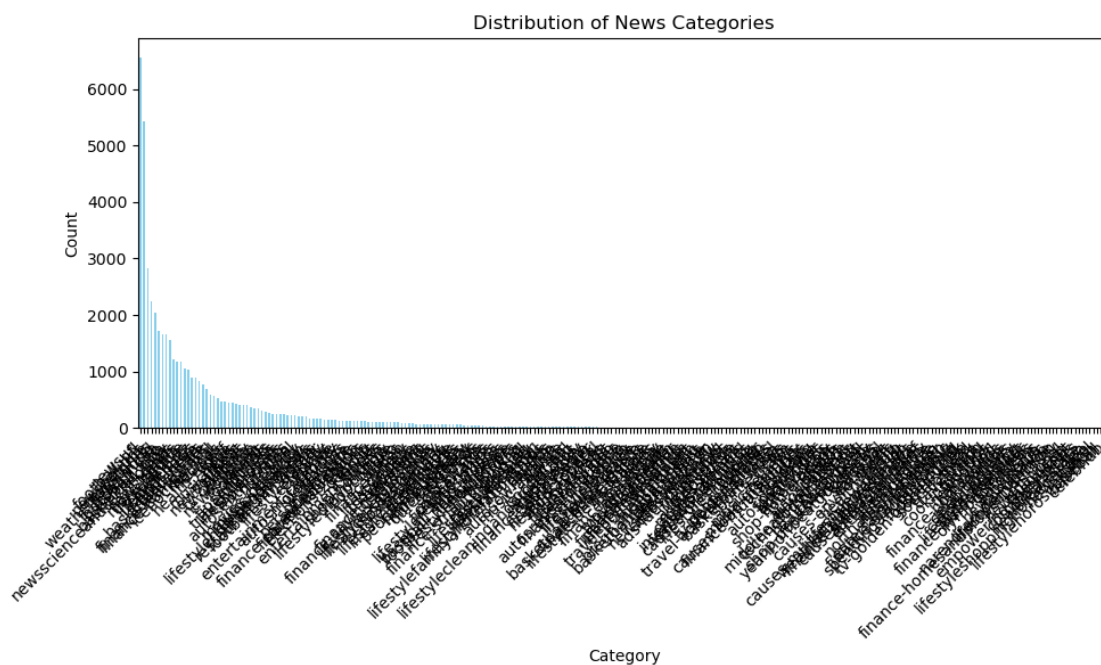
```

```
more_sports
travelarticle
travelnews
lifestylebuzz
autosnews
basketball_ncaa
financenews
finance-real-estate
finance-companies
icehockey_nhl
```

```
[21]: import matplotlib.pyplot as plt
category_counts = news_df['subcategory'].value_counts()

#
plt.figure(figsize=(10, 6))
category_counts.plot(kind='bar', color='skyblue')
plt.title('Distribution of News Categories')
plt.xlabel('Category')
plt.ylabel('Count')
plt.xticks(rotation=45, ha='right')
plt.tight_layout()

#
plt.show()
```



3. :

Apriori FP-growth

Apriori

```
[22]: from mlxtend.preprocessing import TransactionEncoder
from mlxtend.frequent_patterns import apriori, association_rules
# Apriori
user_data = []
last_user = '24186'
tmp = []

for index, row in news_df.iterrows():
    user_id = row['category']
    vroot_id = row['subcategory']

    if user_id == last_user:
        tmp.append(vroot_id)
    else:
        user_data.append(tmp)
        tmp = []
        tmp.append(vroot_id)
    last_user = user_id

user_data.append(tmp)
# user_data

te = TransactionEncoder()
data_encoded = te.fit_transform(user_data)
df = pd.DataFrame(data_encoded, columns=te.columns_)

df
```

```
[22]:
```

|       | ads-latingrammys | ads-lung-health | advice | animals | autosbuying | \ |
|-------|------------------|-----------------|--------|---------|-------------|---|
| 0     | False            | False           | False  | False   | False       |   |
| 1     | False            | False           | False  | False   | False       |   |
| 2     | False            | False           | False  | False   | False       |   |
| 3     | False            | False           | False  | False   | False       |   |
| 4     | False            | False           | False  | False   | False       |   |
| ...   | ...              | ...             | ...    | ...     | ...         |   |
| 41092 | False            | False           | False  | False   | False       |   |
| 41093 | False            | False           | False  | False   | False       |   |
| 41094 | False            | False           | False  | False   | False       |   |
| 41095 | False            | False           | False  | False   | False       |   |
| 41096 | False            | False           | False  | False   | False       |   |

|   | autoscartech | autosclassics | autoscompact | autosenthusiasts | \ |
|---|--------------|---------------|--------------|------------------|---|
| 0 | False        | False         | False        | False            |   |
| 1 | False        | False         | False        | False            |   |

|       |       |       |       |       |
|-------|-------|-------|-------|-------|
| 2     | False | False | False | False |
| 3     | False | False | False | False |
| 4     | False | False | False | False |
| ...   | ...   | ...   | ...   | ...   |
| 41092 | False | False | False | False |
| 41093 | False | False | False | False |
| 41094 | False | False | False | False |
| 41095 | False | False | False | False |
| 41096 | False | False | False | False |

|       |              |     |        |       |                       |   |
|-------|--------------|-----|--------|-------|-----------------------|---|
|       | autoshybrids | ... | voices | watch | weatherfullscreenmaps | \ |
| 0     | False        | ... | False  | False | False                 |   |
| 1     | False        | ... | False  | False | False                 |   |
| 2     | False        | ... | False  | False | False                 |   |
| 3     | False        | ... | False  | False | False                 |   |
| 4     | False        | ... | True   | False | False                 |   |
| ...   | ...          | ... | ...    | ...   | ...                   |   |
| 41092 | False        | ... | False  | False | False                 |   |
| 41093 | False        | ... | False  | False | False                 |   |
| 41094 | False        | ... | False  | False | False                 |   |
| 41095 | False        | ... | False  | False | False                 |   |
| 41096 | False        | ... | False  | False | False                 |   |

|       |                   |             |            |          |       |        |   |
|-------|-------------------|-------------|------------|----------|-------|--------|---|
|       | weathertopstories | weight-loss | weightloss | wellness | wines | wonder | \ |
| 0     | False             | False       | False      | False    | False | False  |   |
| 1     | False             | False       | False      | False    | False | False  |   |
| 2     | False             | False       | True       | False    | False | False  |   |
| 3     | False             | False       | False      | False    | False | False  |   |
| 4     | False             | False       | False      | False    | False | False  |   |
| ...   | ...               | ...         | ...        | ...      | ...   | ...    |   |
| 41092 | False             | False       | False      | False    | False | False  |   |
| 41093 | True              | False       | False      | False    | False | False  |   |
| 41094 | False             | False       | False      | False    | False | False  |   |
| 41095 | False             | False       | False      | False    | False | False  |   |
| 41096 | False             | False       | False      | False    | False | False  |   |

|       |                       |
|-------|-----------------------|
|       | yearinoffbeatgoodnews |
| 0     | False                 |
| 1     | False                 |
| 2     | False                 |
| 3     | False                 |
| 4     | False                 |
| ...   | ...                   |
| 41092 | False                 |
| 41093 | False                 |
| 41094 | False                 |
| 41095 | False                 |

41096 False

[41097 rows x 264 columns]

```
[25]: # Apriori
frequent_itemsets = apriori(df, min_support=0.01, use_colnames=True)

#
print("Frequent Itemsets:")
# print(frequent_itemsets)
frequent_itemsets
```

Frequent Itemsets:

```
[25]:      support      itemsets
0    0.006229      (animals)
1    0.002896      (autosclassics)
2    0.005548      (autosenthusiasts)
3    0.003090      (autosmotorcycles)
4    0.020099      (autosnews)
..      ...
132  0.001436      (newscime, newsscienceandtechnology, newsus)
133  0.001557      (newscime, newsworld, newsus)
134  0.001922      (newspolitics, newsscienceandtechnology, newsus)
135  0.002458      (newsworld, newspolitics, newsus)
136  0.001071      (newsworld, newsscienceandtechnology, newsus)
```

[137 rows x 2 columns]

```
[26]: rules = association_rules(frequent_itemsets, metric="confidence",
    ↪min_threshold=0.3)

#
print("\nAssociation Rules:")
# print(rules)
rules
```

Association Rules:

```
[26]:      antecedents      consequents \
0      (newsoffbeat)      (newsus)
1      (newsscienceandtechnology)      (newsus)
2      (baseball_mlb, football_ncaa)      (football_nfl)
3      (newspolitics, newscime)      (newsus)
4      (newsscienceandtechnology, newscime)      (newsus)
5      (newsworld, newscime)      (newsus)
6      (newspolitics, newsscienceandtechnology)      (newsus)
```



```

7          (newsworld, newspolitics)          (newsus)
8  (newsworld, newsscienceandtechnology)      (newsus)

    antecedent support    consequent support    support    confidence    lift \
0          0.009611          0.134803    0.003090    0.321519    2.385102
1          0.028469          0.134803    0.008881    0.311966    2.314234
2          0.003504          0.114047    0.001144    0.326389    2.861874
3          0.007519          0.134803    0.003285    0.436893    3.240975
4          0.003382          0.134803    0.001436    0.424460    3.148746
5          0.004648          0.134803    0.001557    0.335079    2.485690
6          0.004502          0.134803    0.001922    0.427027    3.167785
7          0.005694          0.134803    0.002458    0.431624    3.201886
8          0.002482          0.134803    0.001071    0.431373    3.200021

    leverage    conviction    zhangs_metric
0    0.001795    1.275197    0.586367
1    0.005044    1.257491    0.584533
2    0.000744    1.315229    0.652866
3    0.002271    1.536470    0.696689
4    0.000980    1.503280    0.684729
5    0.000931    1.301202    0.600488
6    0.001315    1.510014    0.687416
7    0.001690    1.522226    0.691622
8    0.000736    1.521553    0.689213

```

```

[27]: #
rules['lift'] = rules['lift'].apply(lambda x: round(x, 2))

#
observed = rules['support'] * len(df) #
expected = rules['antecedent support'] * rules['consequent support'] * len(df)
#
chi_squared = ((observed - expected) ** 2 / expected).sum() #

#
print("    :")
print(rules[['antecedents', 'consequents', 'lift']])
print("\n    :")
print("    :", chi_squared)

```

```

:
          antecedents    consequents    lift
0          (newsoffbeat)          (newsus)    2.39
1    (newsscienceandtechnology)          (newsus)    2.31
2    (baseball_mlb, football_ncaa)    (football_nfl)    2.86
3          (newspolitics, newscime)          (newsus)    3.24
4    (newsscienceandtechnology, newscime)          (newsus)    3.15
5          (newsworld, newscime)          (newsus)    2.49

```

```

6 (newspolitics, newsscienceandtechnology) (newsus) 3.17
7      (newsworld, newspolitics) (newsus) 3.20
8      (newsworld, newsscienceandtechnology) (newsus) 3.20

```

```

:
```

```

: 1120.7108790076857

```

4.

"support" "confidence" "lift"

```

[ ]: 0      (newsoffbeat) (newsus) 2.39
1      (newsscienceandtechnology) (newsus) 2.31
2      (baseball_mlb, football_ncaa) (football_nfl) 2.86
3      (newspolitics, newscime) (newsus) 3.24
4      (newsscienceandtechnology, newscime) (newsus) 3.15
5      (newsworld, newscime) (newsus) 2.49
6      (newspolitics, newsscienceandtechnology) (newsus) 3.17
7      (newsworld, newspolitics) (newsus) 3.20
8      (newsworld, newsscienceandtechnology) (newsus) 3.20

```