# Apriori Algorithm

Using Marketing Basket Data



## Intuition

The apriori algorithm is a type of association learning used to find groups of items that occur together frequently. This is typically done with market-based analysis.

## Algorithm

Count all items Filter for frequency

Count all pairs filter for frequent pairs

Count candidate triples Filter for frequent triples

## Implementation

```python
# Required Packages
import numpy as np
import matplotlib.pyplot as plt
import pandas as pd
import sklearn as sk
from collections import defaultdict
```

```python
1  # Load the data
2  data = pd.read_csv("../../../_resources/data/Market_Basket_Optimisation
3  data.head() # Display
```

Out[2]:

| | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | shrimp | almonds | avocado | vegetables mix | green grapes | whole weat flour | yams | cottage cheese | energy drink | tomato juice | low fat yogurt | gre |
| 1 | burgers | meatballs | eggs | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | N |
| 2 | chutney | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | N |
| 3 | turkey | avocado | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | N |
| 4 | mineral water | milk | energy bar | whole wheat rice | green tea | NaN | NaN | NaN | NaN | NaN | NaN | N |

```python
1  transactions = []
2  n = len(data)
3  for i in range(0, n):
4      transactions.append([str(data.values[i,j]) for j in range(0, 20)])
5
6  transactions
```

```
  'nan'],
 ['mineral water',
  'milk',
  'energy bar',
  'whole wheat rice',
  'green tea',
  'nan',
  'nan',
  'nan',
  'nan',

  'nan',
  'nan',
  'nan',
  'nan',
  'nan',
  'nan',
  'nan',
  'nan',
  'nan',
```

## Count all items and filter for frequency threshold

```python
In [4]:  1  # Occurences that makes an itemset 'frequent'
         2  threshold = 50
         3
         4  item_counts = defaultdict(int)
         5
         6  # Find candidate items
         7  for trans in transactions:
         8      for item in trans:
         9          item_counts[item] += 1
        10
        11  # filter for frequent items
        12  frequent_items = set()
        13  for key in item_counts:
        14      if item_counts[key] > threshold:
        15          frequent_items.add(key)
        16
        17  print("FREQUENT ITEMS IN BASKET : ")
        18  frequent_items
```

FREQUENT ITEMS IN BASKET :

Out[4]: {'almonds',
         'antioxydant juice',
         'avocado',
         'bacon',
         'barbecue sauce',
         'black tea',
         'blueberries',
         'body spray',
         'brownies',
         'bug spray',
         'burgers',
         'butter',
         'cake',
         'candy bars',
         'carrots',
         'cereals',
         'champagne',
         'chicken',
         'chocolate',
         'cider',
         'clothes accessories',
         'cookies',
         'cooking oil',
         'cottage cheese',
         'eggplant',
         'eggs',
         'energy bar',
         'energy drink',
         'escalope',
         'extra dark chocolate',
         'flax seed',
         'french fries',
         'french wine',
         'fresh bread',
         'fresh tuna',

         'fromage blanc'

'fromage blanc',
'frozen smoothie',
'frozen vegetables',
'gluten free bar',
'grated cheese',
'green beans',
'green grapes',
'green tea',
'ground beef',
'gums',
'ham',
'herb & pepper',
'honey',
'hot dogs',
'light cream',
'light mayo',
'low fat yogurt',
'magazines',
'meatballs',
'melons',
'milk',
'mineral water',
'mint',
'muffins',
'mushroom cream sauce',
'nan',
'nonfat milk',
'oil',
'olive oil',
'pancakes',
'parmesan cheese',
'pasta',
'pepper',
'protein bar',
'red wine',
'rice',
'salmon',
'salt',
'shallot',
'shrimp',
'soup',
'spaghetti',
'spinach',
'strawberries',
'strong cheese',
'tomato juice',
'tomato sauce',
'tomatoes',
'toothpaste',
'turkey',
'vegetables mix',
'white wine',
'whole weat flour',
'whole wheat pasta',
'whole wheat rice',
'yams',

'yogurt cake'

```
yoyuit caκe ,
'zucchini'}
```

## Count all item PAIRS and filter for frequency threshold

```
In [5]:   1  pair_counts = defaultdict(int)
          2  frequent_pairs = set()
```

```
In [9]:   1  def normalize_group(*args):
          2      return str(sorted(args))
          3
          4  # Get counts of candidate pairs
          5  for trans in transactions:
          6      for items in trans:
          7          for idx_1 in range(len(items) - 1):
          8              if items[idx_1] not in frequent_items:
          9                  continue
         10              for idx_2 in range(idx_1 + 1, len(items)):
         11                  if items[idx_2] not in frequent_items:
         12                      continue
         13                  pair = normalize_group(items[idx_1], items[idx_2])
         14                  pair_counts[pair] += 1
         15
         16
         17  # Get frequent pairs
         18  for key in pair_counts:
         19      if pair_counts[key] > threshold:
         20          frequent_pairs.add(key)
         21
         22  print("FREQUENT PAIRS IN BASKET : ")
         23  frequent_pairs
```

```
FREQUENT PAIRS IN BASKET :
```

```
Out[9]:   set()
```

## Count all item TRIPLES and filter for frequency threshold

```
In [12]:  1  triple_counts = defaultdict(int)
          2  frequent_triple = set()
```

```
In [13]:   1  def generate_pairs(*args):
           2      pairs = []
           3      for idx_1 in range(len(args) - 1):
           4          for idx_2 in range(idx_1+1, len(args)):
           5              pairs.append(normalize_group(args[idx_1], args[idx_1]))
           6      return pairs
           7
           8  # Get counts of candidate pairs
           9  for trans in transactions:
          10      for items in trans:
          11          for idx_1 in range(len(items) - 1):
          12              if items[idx_1] not in frequent_items:
          13                  continue
          14              for idx_2 in range(idx_1 + 1, len(items)):
          15                  if items[idx_2] not in frequent_items:
          16                      continue
          17                  first_pair = normalize_group(items[idx_1], items[idx_2]
          18                  if first_pair not in frequent_pairs:
          19                      continue
          20                  for idx_3 in range(idx_2 + 1, len(items)):
          21                      if items[idx_3] not in frequent_items:
          22                          continue
          23                      pairs = generate_pairs(items[idx_1], items[idx_2],
          24                      if any(pair not in frequent_pairs for pair in pairs
          25                          continue
          26                      triple = normalize_group(items[idx_1], items[idx_2]
          27                      triple_counts[triple] += 1
          28
          29
          30
          31  # Get frequent pairs
          32  for key in triple_counts:
          33      if triple_counts[key] > threshold:
          34          frequent_triple.add(key)
          35
          36  print("FREQUENT TRIPLE IN BASKET : ")
          37  frequent_triple
```

FREQUENT TRIPLE IN BASKET :

Out[13]:  set()


In [ ]:    1