

In [39]:

```

1 import pandas as pd
2 from mlxtend.frequent_patterns import apriori
3 from mlxtend.frequent_patterns import association_rules
4 import numpy as np
5 import matplotlib.pyplot as plt
6 from itertools import permutations
7 import seaborn as sns
8
9 from pandas.plotting import parallel_coordinates
10

```

In [2]:

```

1 # !pip install mlxtend
2 #https://goldinlocks.github.io/Market-Basket-Analysis-in-Python/

```

In [10]:

```
1 df = pd.read_csv("tel_samp_rec.csv",encoding="latin-1")
```

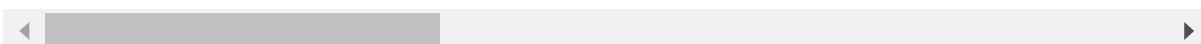
In [11]:

```
1 df.head()
```

Out[11]:

	Defence.date	Domains	Full.Text.Language	def.date	n.disc	these.id	disc1.lev1	
0	2010/09/23	Sciences du Vivant [q-bio] / Ecologie, Environ...	French	2010.0	1	tel-00662843v1	Sciences du Vivant [q-bio]	Env
1	2009/11/02	Sciences de l'Homme et Société	French	2009.0	1	tel-00491490v1	Sciences de l'Homme et Société	
2	1996/05/30	Sciences du Vivant [q-bio] / Alimentation et N...	French	1996.0	1	tel-01776364v1	Sciences du Vivant [q-bio]	Alin
3	2018/02/02	Informatique [cs] / Autre [cs.OH] \r\n\r\nInf...	French	2018.0	1	tel-02437294v1	Informatique [cs]	Al
4	2015/07/08	Informatique [cs] / Automatique \r\n\r\nInfor...	French	2015.0	1	tel-01245100v1	Informatique [cs]	A

5 rows × 25 columns



In [12]:

```
1 cols = ['disc1.rec.lev1', 'disc2.rec.lev1', 'disc3.rec.lev1']
2 #subset columns shown above and take columns where all the 3 columns are not null
3 df_sub = df[df[cols].notnull().all(axis=1)]
```

In [13]:

```
1 df_sub.head()
```

Out[13]:

	Defence.date	Domains	Full.Text.Language	def.date	n.disc	these.id	disc1.lev1	dis
53	1985/10/28	Planète et Univers [physics] / Sciences de la ...	French	1985.0	2	00711880v1 tel-	Planète et Univers [physics]	S de
104	2018/12/17	Sciences de l'ingénieur [physics] / Génie civi...	English	2018.0	2	02182014v1 tel-	Sciences de l'ingénieur [physics]	Gé
113	2003/06/17	Sciences de l'ingénieur [physics] / Traitement...	French	2003.0	3	00130932v1 tel-	Sciences de l'ingénieur [physics]	Tra d [e
193	1997/10/24	Planète et Univers [physics] / Sciences de la ...	French	1997.0	2	00675418v1 tel-	Planète et Univers [physics]	S de
212	2002/12/13	Sciences du Vivant [q-bio] / Autre [q-bio.OT] ...	French	2002.0	2	00008546v1 tel-	Sciences du Vivant [q-bio]	A

5 rows × 25 columns



In [14]:

```
1 df_sub = df_sub[cols]
```

In [15]:

```
1 df_sub.head()
```

Out[15]:

	disc1.rec.lev1	disc2.rec.lev1	disc3.rec.lev1
53	VIII	VIII	VIII
104	IX	IX	V
113	IX	VI	V
193	VIII	VIII	VIII
212	X	IX	IX

In [16]:

```
1 # Getting the list of transactions from the dataset
2 transactions = []
3 for i in range(0, len(df_sub)):
4     transactions.append([str(df_sub.values[i,j]) for j in range(0, len(df_sub.columns))])
```

In [17]:

```
1 #check transactions
2 transactions[:1]
3
```

Out[17]:

```
[['VIII', 'VIII', 'VIII']]
```

In [19]:

```
1
2 # Extract unique items.
3 flattened = [item for transaction in transactions for item in transaction]
4 items = list(set(flattened))
```

In [20]:

```
1 print('# of items:',len(items))
2 print(list(items))
```

of items: 13

```
['IV', 'pharmacie', 'VI', 'I', 'III', 'XII', 'II', 'VIII', 'X', 'V', 'I - Dr  
oit', 'VII', 'IX']
```

In [21]:

```
1 #remove nan if present in list
2 if 'nan' in items: items.remove('nan')
3 print(list(items))
```

```
['IV', 'pharmacie', 'VI', 'I', 'III', 'XII', 'II', 'VIII', 'X', 'V', 'I - Dr  
oit', 'VII', 'IX']
```

In [22]:

```

1 # Compute and print rules.
2 rules = list(permutations(items, 2))
3 print('# of rules:', len(rules))
4 print(rules[:5])

```

of rules: 156

```

[('IV', 'pharmacie'), ('IV', 'VI'), ('IV', 'I'), ('IV', 'III'), ('IV', 'XI
I')]

```

In [23]:

```

1 # Import the transaction encoder function from mlxtend
2 from mlxtend.preprocessing import TransactionEncoder
3
4 # Instantiate transaction encoder and identify unique items
5 encoder = TransactionEncoder().fit(transactions)
6
7 # One-hot encode transactions
8 onehot = encoder.transform(transactions)
9
10 # Convert one-hot encoded data to DataFrame
11 onehot = pd.DataFrame(onehot, columns = encoder.columns_)
12
13 # Print the one-hot encoded transaction dataset
14 onehot.head()

```

Out[23]:

	I	I- Droit	II	III	IV	IX	V	VI	VII	VIII	X	XII	pharmacie
0	False	False	False	False	False	False	False	False	False	True	False	False	False
1	False	False	False	False	False	True	True	False	False	False	False	False	False
2	False	False	False	False	False	True	True	True	False	False	False	False	False
3	False	False	False	False	False	False	False	False	False	True	False	False	False
4	False	False	False	False	False	True	False	False	False	False	True	False	False

In [24]:

```
1 def leverage(antecedent, consequent):
2     # Compute support for antecedent AND consequent
3     supportAB = np.logical_and(antecedent, consequent).mean()
4
5     # Compute support for antecedent
6     supportA = antecedent.mean()
7
8     # Compute support for consequent
9     supportB = consequent.mean()
10
11    # Return Leverage
12    return supportAB - supportB * supportA
13
14    # Define a function to compute Zhang's metric
15    def zhang(antecedent, consequent):
16        # Compute the support of each book
17        supportA = antecedent.mean()
18        supportC = consequent.mean()
19
20        # Compute the support of both books
21        supportAC = np.logical_and(antecedent, consequent).mean()
22
23        # Complete the expressions for the numerator and denominator
24        numerator = supportAC - supportA*supportC
25        denominator = max(supportAC*(1-supportA), supportA*(supportC-supportAC))
26
27        # Return Zhang's metric
28        return numerator / denominator
29
30    def conviction(antecedent, consequent):
31        # Compute support for antecedent AND consequent
32        supportAC = np.logical_and(antecedent, consequent).mean()
33
34        # Compute support for antecedent
35        supportA = antecedent.mean()
36
37        # Compute support for NOT consequent
38        supportnC = 1.0 - consequent.mean()
39
40        # Compute support for antecedent and NOT consequent
41        supportAnC = supportA - supportAC
42
43        # Return conviction
44        return supportA * supportnC / supportAnC
45
```

In [25]:

```

1 # Create rules DataFrame
2 rules_ = pd.DataFrame(rules, columns=['antecedents', 'consequents'])
3
4 # Define an empty list for metrics
5 zhangs, conv, lev, antec_supp, cons_supp, suppt, conf, lft = [], [], [], [], [], [], []
6
7 # Loop over lists in itemsets
8 for itemset in rules:
9     # Extract the antecedent and consequent columns
10    antecedent = onehot[itemset[0]]
11    consequent = onehot[itemset[1]]
12
13    antecedent_support = onehot[itemset[0]].mean()
14    consequent_support = onehot[itemset[1]].mean()
15    support = np.logical_and(onehot[itemset[0]], onehot[itemset[1]]).mean()
16    confidence = support / antecedent_support
17    lift = support / (antecedent_support * consequent_support)
18
19    # Complete metrics and append it to the list
20    antec_supp.append(antecedent_support)
21    cons_supp.append(consequent_support)
22    suppt.append(support)
23    conf.append(confidence)
24    lft.append(lift)
25    lev.append(leverage(antecedent, consequent))
26    conv.append(conviction(antecedent, consequent))
27    zhangs.append(zhang(antecedent, consequent))
28
29 # Store results
30 rules_['antecedent support'] = antec_supp
31 rules_['consequent support'] = cons_supp
32 rules_['support'] = suppt
33 rules_['confidence'] = conf
34 rules_['lift'] = lft
35 rules_['leverage'] = lev
36 rules_['conviction'] = conv
37 rules_['zhang'] = zhangs
38
39 # Print results
40 rules_.sort_values('zhang', ascending=False).head()

```

Out[25]:

	antecedents	consequents	antecedent support	consequent support	support	confidence	lift	leve
2	IV	I	0.143959	0.021994	0.020566	0.142857	6.495362	0.01
9	IV	I - Droit	0.143959	0.004856	0.003999	0.027778	5.720588	0.00
97	X	pharmacie	0.199372	0.019709	0.016281	0.081662	4.143453	0.01
45	I	I - Droit	0.021994	0.004856	0.001428	0.064935	13.372804	0.00
123	I - Droit	I	0.004856	0.021994	0.001428	0.294118	13.372804	0.00

In [26]:

```

1 # Function to convert rules to coordinates.
2 def rules_to_coordinates(rules):
3     rules['antecedent'] = rules['antecedents'].apply(lambda antecedent: list(antecedent))
4     rules['consequent'] = rules['consequents'].apply(lambda consequent: list(consequent))
5     rules['rule'] = rules.index
6     return rules[['antecedent', 'consequent', 'rule']]

```

In []:

1

In [32]:

```
1 rules_.head()
```

Out[32]:

	antecedents	consequents	antecedent support	consequent support	support	confidence	lift	leverage
0	IV	pharmacie	0.143959	0.019709	0.000857	0.005952	0.302019	-0.00198
1	IV	VI	0.143959	0.162239	0.000857	0.005952	0.036689	-0.02249
2	IV	I	0.143959	0.021994	0.020566	0.142857	6.495362	0.01739
3	IV	III	0.143959	0.032562	0.022279	0.154762	4.752820	0.01759
4	IV	XII	0.143959	0.049414	0.034276	0.238095	4.818332	0.02716

In [33]:

```

1 #remove rows where antecedent = consequent
2 rules_ = rules_[rules_['antecedents'] != rules_['consequents']]

```

In [35]:

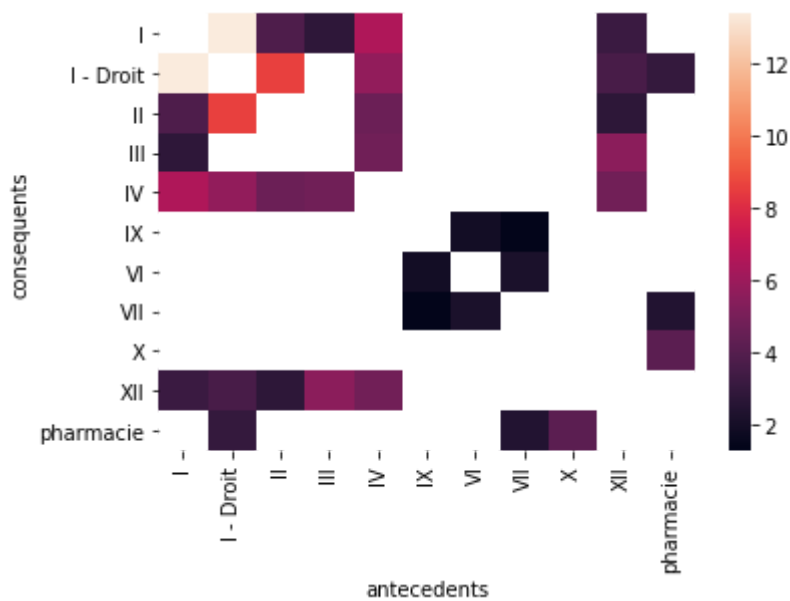
```

1 #filter rules with lift > 1
2 rules_ = rules_.query("lift>1")
3 #create support table based on Lift values > 1
4 support_table = rules_.pivot(index='consequents', columns='antecedents',
5 values='lift')
6
7
8
9
10 sns.heatmap(support_table)

```

Out[35]:

<AxesSubplot:xlabel='antecedents', ylabel='consequents'>

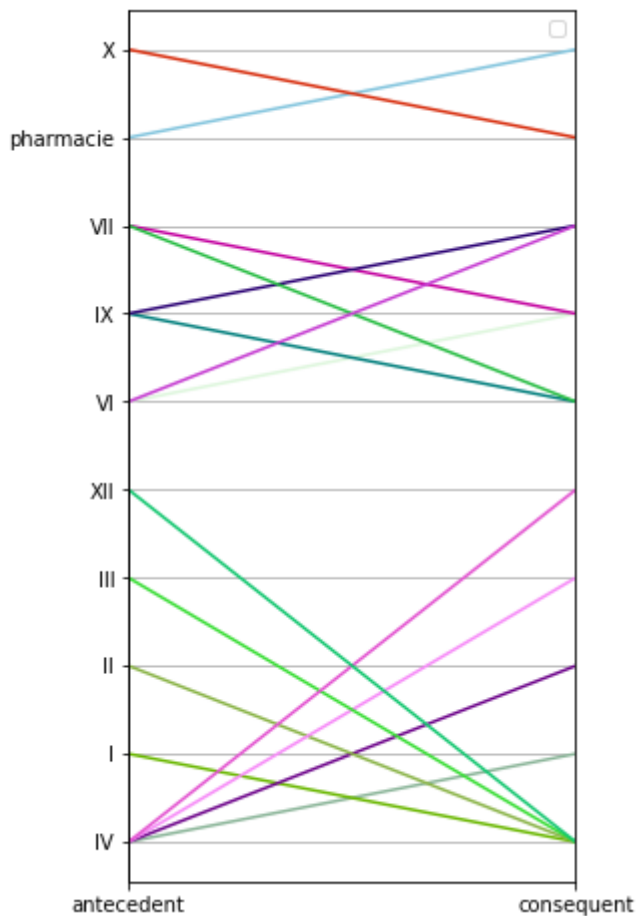


In [41]:

```

1  # Generate frequent itemsets
2  frequent_itemsets = apriori(onehot, min_support = 0.01, use_colnames = True, max_len =
3  # Generate association rules
4  rules = association_rules(frequent_itemsets, metric = 'lift', min_threshold = 1.00)
5  # Generate coordinates and print example
6  coords = rules_to_coordinates(rules)
7  # Generate parallel coordinates plot
8
9  plt.figure(figsize=(4,8))
10 parallel_coordinates(coords, 'rule')
11 plt.legend([])
12 plt.grid(True)
13 plt.show()

```



<https://www.galaxie.enseignementsup-recherche.gouv.fr/ensup/pdf/qualification/sections.pdf>
(<https://www.galaxie.enseignementsup-recherche.gouv.fr/ensup/pdf/qualification/sections.pdf>)

LINK : Reference what I , II etc means