

Q1: TF-IDF

1. Select texts.

Recently, there has been a rent crisis in Ireland, especially in Dublin. I am curious what the local people think about the rent crisis on Twitter. Therefore, I intercepted ten tweets and completed the normalization of them. Lowercase the text, remove punctuations and convert numbers to English words using `nltk.corpus.stopwords.words`. Here are the normalized tweets.

```
a = ['god', 'bless', 'politicians', 'oppose', 'homes', 'built', 'complain',
    ', 'housing', 'rent', 'crisis', 'even', 'think', 'could', 'troll', 'level',
    ', 'one', 'hundredk', 'year', 'make', 'lads']

b = ['protests', 'dublin', 'sharp', 'rise', 'costs', 'living', 'huge', 'ac
comodation', 'crisis', 'places', 'rent', 'people', 'afford', 'buy', 'anyth
ing']

c = ['entitled', 'opinion', 'young', 'person', 'dublin', 'anywhere', 'irel
and', 'living', 'crisis', 'impossible', 'save', 'monthly', 'rent', 'half',
    'wage', 'bills', 'extortionate', 'prices', 'basic', 'costs', 'rising']

d = ['dublin', 'student', 'accommodation', 'crisis', 'would', 'appear', 'b
ad', 'students', 'knocking', 'random', 'doors', 'asking', 'people', 'spare',
    ', 'room', 'rent', 'happened', 'us', 'hypothetical', 'effective', 'governm
ent', 'would', 'probably', 'something']

e = ['issue', 'high', 'cost', 'lack', 'supply', 'nowhere', 'even', 'rent',
    'dublin', 'talk', 'people', 'twentys', 'annoys', 'many', 'people', 'like',
    ', 'dan', 'touch', 'regards', 'rent', 'crisis', 'tragic', 'reality']

f = ['yes', 'many', 'people', 'putting', 'major', 'milestones', 'lives', '
tight', 'money', 'last', 'years', 'cost', 'living', 'crisis', 'going', 'co
st', 'rent', 'dublin']

g = ['rent', 'crisis', 'wild', 'like', 'seriously', 'feels', 'like', 'neve
r', 'find', 'flat']

h = ['doubt', 'someone', 'try', 'spin', 'vacant', 'houses', 'rural', 'part
s', 'west', 'reason', 'rent', 'crisis', 'dublin']

i = ['airbnbs', 'huge', 'factor', 'housing', 'crisis', 'dublin', 'rental',
    'prices', 'skyrocketed', 'amount', 'properties', 'actually', 'available',
    'rent', 'time', 'low', 'landlords', 'instead', 'airbnb', 'extortionate',
    'prices', 'need', 'address']

j = ['housing', 'crisis', 'reached', 'tipping', 'point', 'pretty', 'much',
    'rented', 'property', 'available', 'outside', 'dublin', 'people', 'buy',
    'due', 'cbi', 'rules', 'whose', 'landlords', 'selling', 'due', 'rent', 'co
ntrols', 'vulnerable', 'position']
```

Let `min_freq = 1`, here is the plot in `WordCloud`.



Figure 1: Repeated words

2. Compute TF and TF-IDF

The results obtained after the pre-processing are as follows. I have selected the TF scores for the first ten and last ten examples.

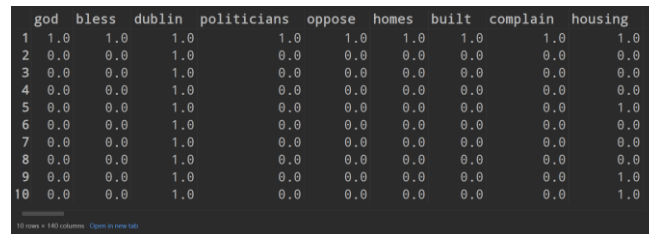


Figure 2: TF scores and matrix

The result of TF-IDF:

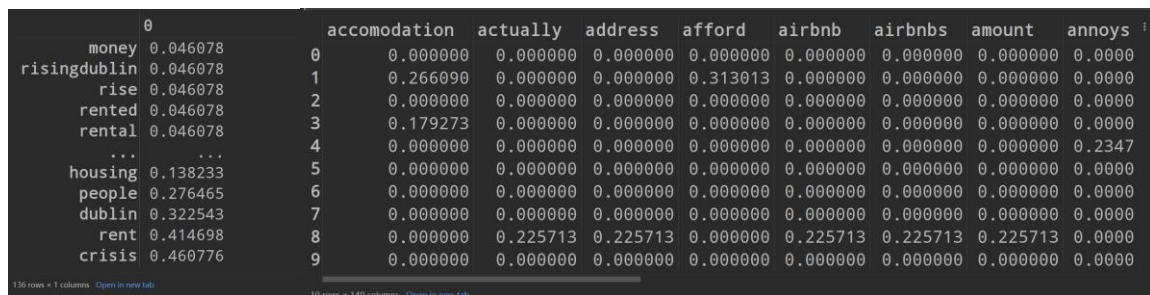


Figure 3: TF-IDF scores and matrix

There is little change in word ranking between TF and TF-IDF. The ranking of the last five is still in the last five. Furthermore, There are some internal changes in the ranking of the last three words, possibly due to changes in word combinations throughout the corpus.

Q2: PMI

Here are the top five PMI scores calculated by using `nlk.collocations`.

```
{('god', 'bless'): 7.584962500721156,  
 ('bless', 'politicians'): 7.584962500721156,  
 ('politicians', 'oppose'): 7.584962500721156,  
 ('oppose', 'homes'): 7.584962500721156,  
 ('homes', 'built'): 7.584962500721156}
```

The results do not make sense. Many words appear only one time but achieve a high PMI value, like 'homes' and 'built', which does not match the expected results. In this example, the expected value would be 'rent' and 'crisis' to have a maximum PMI value. Set this minimal cut-off frequency to 3 and calculate it again.

```
('rent', 'crisis') = 2.803602787196496
```

Q3: Entropy

I chose ten tweets from Tesla as the `spam_set`, all about its electric cars product. Apart from that, I also chose some other ten tweets as the `random_set` from the trend of Twitter which covers almost all aspects of life, including news, sports, entertainment, etc and has little correlation between them. Here are the sets.

```
spam_set = ['Your Tesla now shows energy consumed vs projected & gives range tips', 'The Tesla Model Y is officially the best selling vehicle in Germany, accounting for just over 4.5% of the car market.', 'Tesla Immersive, our multichannel audio upmixer, enables stereo content to be remixed in real time, optimizing the listening experience for our vehicles specifically', 'The brain of your Tesla: Neural networks with 1 billion parameters, completing 144 trillion operations per second', 'Enter your destination & your Tesla will automatically include Supercharging stops in your route', 'We are launching the Tesla Shareholder Platform – join the program to participate in Tesla events and hear more updates soon', 'Tesla will ask shareholders to vote at this year’s annual meeting to authorize additional shares in order to enable a stock split.', 'Non-Tesla vehicles can now charge at select Superchargers in Denmark, Finland, Germany, Luxembourg and Switzerland via the Tesla app.', 'Tesla navigation will now take predicted crosswind, headwind, humidity & temperature into account for calculating battery % on arrival', 'Tesla Vehicle Production & Deliveries and Date for Financial Results & Webcast for First Quarter 2022']
```

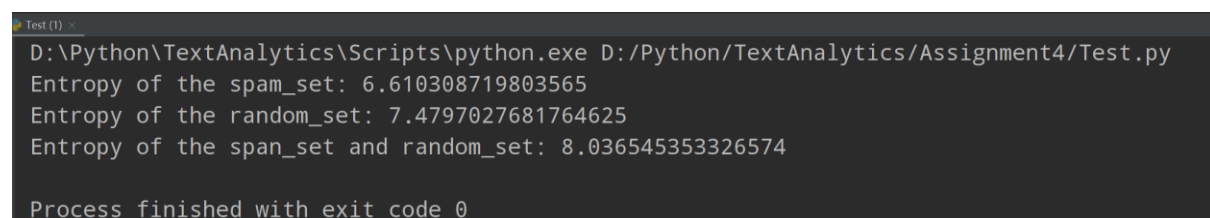
```
random_set = [ 'There are two groups who I will simply not accept any moral judgement or admonition from: Irish people feverishly supporting Ukraine but saying the girls singing Up The Ra is wrong+ anyone supporting a party who wanted to commemorate the Black and Tans. Hypocrisy is not for me.', 'On February 24, Russia launched an invasion of Ukraine by land, air and sea after months of tensions between Moscow and Kyiv. The attack triggered a chain of events over the next six months including unprecedented sancti
```

ons against Russia and the expansion of NATO.', 'Philly McMahon has hit out at Sky Sports News presenter Rob Wooton following his controversial questioning of Republic of Ireland star Chloe Mustaki.', 'Thoroughly enjoying Rob Wooton getting (deservedly) shat on from Irish Twitter. Typical English ignorance with the usual superiority complex. Lbu they'd do your head in.', 'Was told by my lecture today at 11.11am as we all walked out that the protest was a load of bullshit, he's not the one struggling to pay for accommodation, college fees, food, parking he's not on €10.50 an hour. The arrogance of some people.', 'Join dynamic alumni Nicai de Guzman (Wolfgang Digital), Emmet Daniel (Hubspot), and expert Dr. Linda Yang (Intercultural Development Programme, UCD) as they offer advice on thriving in multicultural work environments', 'Brendan Fraser hadn't played the lead in a major movie in 12 years, a gap on his résumé that's been attributed to personal issues, health problems and an assault allegation against the former head of the HFPA.', '58 years on the road in 2022, proud to have brought the story of Ireland around the world since 1964! Let The People Sing', 'Here's Xavi celebrating a group stage goal to avoid Europa league and here's Ancelotti reacting to the greatest come back of all time to make the UCL final. Levels to this game.', 'Venus Williams on sisterhood, securing equal pay for women in sport and building an empire off the court']

Here is the entropy method found in the previous PowerPoint. Before calculating the entropy value, I have done some normalization in the pre-processing section using the `pre_preprocessing()` method.

```
import nltk
import math
def entropy(labels):
    freqdist = nltk.FreqDist(labels)
    probs = [freqdist.freq(l) for l in freqdist]
    return -sum(p * math.log(p,2) for p in probs)

print("Entropy of the spam_set: {}".format(entropy(pre_preprocessing(spam_set))))
print("Entropy of the random_set: {}".format(entropy(pre_preprocessing(random_set))))
print("Entropy of the spam_set and random_set: {}".format(entropy(pre_preprocessing(spam_set + random_set))))
```



```
Test (1)
D:\Python\TextAnalytics\Scripts\python.exe D:/Python/TextAnalytics/Assignment4/Test.py
Entropy of the spam_set: 6.610308719803565
Entropy of the random_set: 7.4797027681764625
Entropy of the spam_set and random_set: 8.036545353326574

Process finished with exit code 0
```

Figure 4: Results of Entropy

Here it computes the results of three different entropy, corresponding to three cases. Entropy values symbolize a high or low state of disorder, even in text analytics, which can also reflect whether it is a disorder or not. The `spam_set` is primarily based on Tesla's product, which has lower entropy than the `random_set`. Moreover, the `random_set` is chosen entirely at random from the Twitter trends. It has a higher degree of disorder which is reflected in the result. At last, combining the two corpus sets has a higher disorder which gets the highest entropy value.