

Metin Madenciliğine Giriş
Metinden Tema Çıkarma

Mert Arslan

Doğuş Üniversitesi
İstanbul, Türkiye
Mayıs 2021

Dataset json formatında bir kitap.

```
def load_data(file):  
    with open (file, "r", encoding="utf-8") as f:  
        data = json.load(f)  
    return (data)  
  
def write_data(file, data):  
    with open (file, "w", encoding="utf-8") as f:  
        json.dump(data, f, indent=4)
```

Json tipinde datayı kolayca yüklemek ve yazmak için kullanılan kod bloğu.

```
1 stopwords = stopwords.words("english")  
  
1 print (stopwords)
```

Data set ingilizce olduğu için stopwordsü İngilizce seçildi.

```
1 def lemmatization(texts, allowed_postags=["NOUN", "ADJ", "VERB", "ADV"]):  
2     nlp = spacy.load("en_core_web_sm", disable=["parser", "ner"])  
3     texts_out = []  
4     for text in texts:  
5         doc = nlp(text)  
6         new_text = []  
7         for token in doc:  
8             if token.pos_ in allowed_postags:  
9                 new_text.append(token.lemma_)  
10        final = " ".join(new_text)  
11        texts_out.append(final)  
12    return (texts_out)  
13  
14  
15 lemmatized_texts = lemmatization(data)  
16 print (lemmatized_texts[0][0:90])
```

name be bear small town call be bear be very hard work child father mother have small mill

Bu kod bloğunda sadece İsim, sıfat, fiil, zarfların kalmasına izin verildi. Parser ve neri daha hızlı çalışması için kapatıldı. Sadece İsim, sıfat, fiil olan tokenlerin kökleri(lemma) bulundu.

```

1 def gen_words(texts):
2     final = []
3     for text in texts:
4         new = gensim.utils.simple_preprocess(text, deacc=True)
5         final.append(new)
6     return (final)
7
8 data_words = gen_words(lemmatized_texts)
9
10 print (data_words[0][0:20])

```

['name', 'be', 'bear', 'small', 'town', 'call', 'be', 'bear', 'be',

Stopwordler(my,if,or..) çıkarıldı ve kelimeler tek tek ayrıldı.

```

1 id2word = corpora.Dictionary(data_words)
2
3 corpus = []
4 for text in data_words:
5     new = id2word.doc2bow(text)
6     corpus.append(new)
7
8 print (corpus[0][0:20])
9
0 word = id2word[[0][:1][0]]
1 print (word)

```

Kelimeleri ve o kelimelerin frekanslarını gösteren bir sözlük oluşturuldu.

```

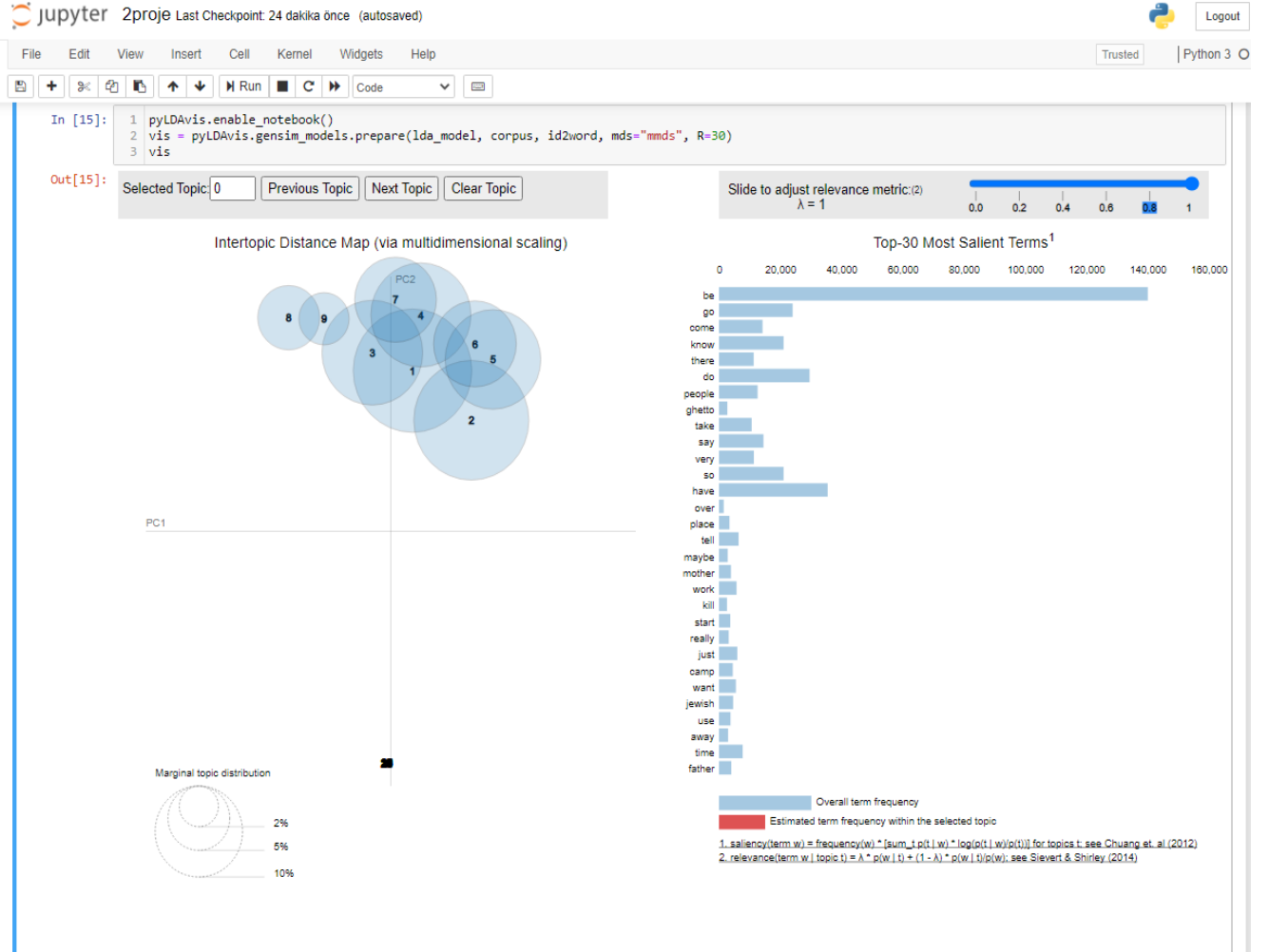
1 lda_model = gensim.models.ldamodel.LdaModel(corpus=corpus,
2                                             id2word=id2word,
3                                             num_topics=30,
4                                             random_state=100,
5                                             update_every=1,
6                                             chunksize=100,
7                                             passes=10,
8                                             alpha="auto")

```

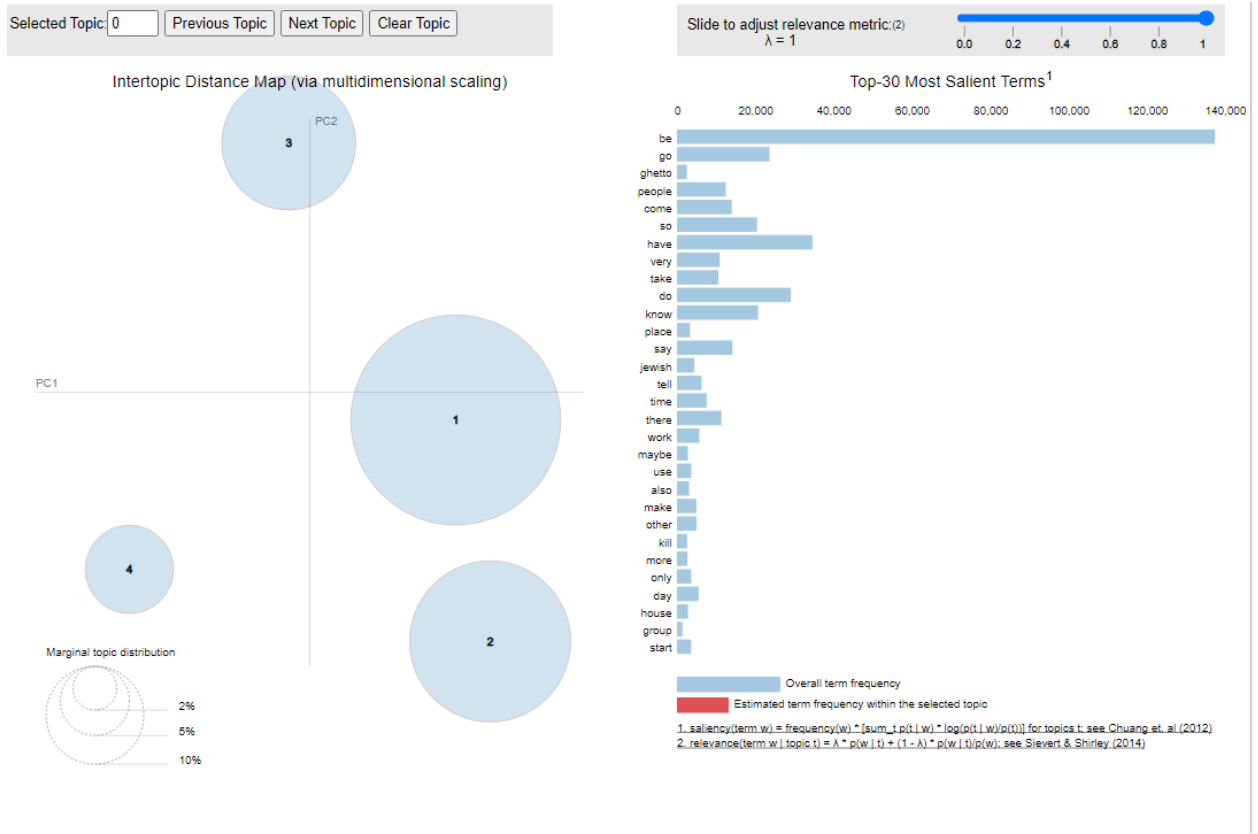
LDA topic modeli oluşturuldu.

```
1 pyLDavis.enable_notebook()
2 vis = pyLDavis.gensim_models.prepare(lda_model, corpus, id2word, mds="mmds", R=30)
3 vis
```

LDA topic modeli gösterildi.



LDA topic modelinde clusterlar üst üste bindi preproscing kısmında sorun olabilir. Tema olabilecek kelimeler know, go, take gibi kelimeler gösterdi bunlar filtrelenerek daha iyi bir sonuç alınabilir.



Number of topics kısmını 4 e düşürdüğümde clusterlar daha iyi bir şekilde ayrıldı. Tema olabilecek kelimelerde anlamsız olanlar filtrelenmeli.