

Exercise - 1

① Import dependencies

```
import gensim
from gensim.models.doc2vec import Doc2Vec, TaggedDocument
from nltk.tokenize import word_tokenize
from sklearn import utils
```

```
data = ["I love machine learning. It's awesome.",
        "I love coding in python",
        "I love building chatbots",
        "They chat amazingly well"]
```

```
import nltk
```

```
nltk.download('punkt')
```

```
tagged_data = [TaggedDocument(words_tokenize(d.lower()),
                               tags = [str(i)]), for i, d in enumerate(data)]
```

```
vec_size = 20
```

```
alpha = 0.025
```

```
model = Doc2Vec(vector_size=vec_size,
```

```
alpha=alpha
```

```
min_alpha=0.00025,
```

```
min_count=1
```

```
, dm=1)
```

```
model.build_vocab(tagged_data)
```

Natural Language Processing Lab

Lab4. Computing Document Similarity using Doc2Vec Model

EXERCISE-1

1. Import dependencies

```
from gensim.models.doc2vec import Doc2Vec, TaggedDocument
from nltk.tokenize import word_tokenize
from sklearn import utils
```

2. Create dataset

```
data = ["I love machine learning. Its awesome.",
        "I love coding in python",
        "I love building chatbots",
        "they chat amazingly well"]
```

3. Create TaggedDocument

```
tagged_data = [TaggedDocument(words=word_tokenize(d.lower()),
                               tags=[str(i)]) for i, d in enumerate(data)]
```

4. Train Model

```
# model parameters
vec_size = 20
alpha = 0.025

# create model
model = Doc2Vec(vector_size=vec_size,
                alpha=alpha,
                min_alpha=0.00025,
                min_count=1,
                dm=1)

# build vocabulary
model.build_vocab(tagged_data)

# shuffle data
tagged_data = utils.shuffle(tagged_data)

# train Doc2Vec model
model.train(tagged_data,
            total_examples=model.corpus_count,
            epochs=30)

model.save("d2v.model")
print("Model Saved")
```

5. Find Similar documents for the given document

```
from gensim.models.doc2vec import Doc2Vec

model = Doc2Vec.load("d2v.model")

# to find the vector of a document which is not in training data
test_data = word_tokenize("I love chatbots".lower())

v1 = model.infer_vector(test_data)
print("V1_infer", v1)

# to find most similar doc using tags
```

```
data_examples = model.ComplexCount,  
epochs = 30)
```

```
model.save("dev_model")  
print("Model saved")
```

```
from gensim.models.doc2vec import Doc2Vec  
model = Doc2Vec.load("dev_model")
```

```
test_data = word_tokenize("I love chocolate" lower())
```

```
v1 = model.infer_vector(test_data)
```

```
print("v1 - infer", v1)
```

```
similar_doc = model.docvecs.most_similar([1])
```

```
print(similar_doc)
```

```
print(model.docvecs[1])
```

```
docs = ["the house had a tiny little mouse",  
        "the cat saw the mouse",  
        "the mouse ran away from the house",  
        "the cat finally ate the mouse",  
        "the end of the mouse story"]
```

```
similar_doc = model.docvecs.most_similar('1')
print(similar_doc)
```

to find vector of doc in training data using tags or
in other words, printing the vector of document at index 1 in training data

```
print(model.docvecs['1'])
```

EXERCISE-2

Question1. Train the following documents using Doc2Vec model

```
docs=["the house had a tiny little mouse",
      "the cat saw the mouse",
      "the mouse ran away from the house",
      "the cat finally ate the mouse",
      "the end of the mouse story"
]
```

Question2. Find the most similar TWO documents for the query document "cat stayed in the house".

tagged_data = [TaggedDocument(words=word_tokenize
tags = [str(i)] for i, d in enumerate(docs))

vec_size = 20

alpha = 0.025

model = Doc2Vec(vector_size=vec_size, alpha=alpha,
min_alpha=0.00025,
min_count=1)

model.build_vocab(tagged_data)

tagged_data = utils.shuffle(tagged_data)

```
model.train (tagged, data_loader, examples=model, corpus  
              - count, epochs=20)
```

```
model.save ("doV.model")
```

```
print ("Model saved")
```

```
from gensim.models.docvec import Docvec.
```

```
model = Docvec.load ("doV.model")
```

```
- test_data = WordTokenizer ("cat stayed in the  
house".lower())
```

```
v1 = model.infer_vector (test_data)
```

```
print ("V1 - infer", v1)
```

```
similar_doc = model.docvecs.most_similar (v1)
```

```
print (similar_doc)
```