# Lab No. 13

## Natural Language Processing (NLP)

This laboratory session introduces students to Word2Vec, a popular technique in Natural Language Processing (NLP) used to represent words as meaningful dense vectors. Using textual data from the *Game of Thrones* series, students will learn how word embeddings capture semantic relationships between words based on their context. Through preprocessing, model training, and similarity analysis, this lab helps students understand how machines learn word meanings and relationships from large text corpora in an unsupervised manner.

**LAB Objectives:**

- Understand **distributional semantics**

- Learn how **Word2Vec** converts words into dense vectors

- Apply **Word2Vec (Skip-Gram / CBOW)** on real textual data (Game of Thrones)

- Explore **word similarity, analogy, and visualization**

**game-of-thrones-word2vec**

word2vec applied on game of thrones data

Dataset Link: https://www.kaggle.com/khulasasndh/game-of-thrones-books

Download the data set from Kaggle

▲ 29    <> Code    ⬇ Download

Data Card    Code (47)    Discussion (0)    Suggestions (0)

**001ssb.txt** (1.63 MB)    ⬇ ⛶ >

**About this file**    ✎ Suggest Edits

This file does not have a description yet.

⚠ This preview is truncated due to the large file size. Create a Notebook or download this file to see the full content.    **Download**    **Create Notebook**

```
A Game Of Thrones
Book One of A Song of Ice and Fire
By George R. R. Martin
PROLOGUE
"We should start back," Gared urged as the woods began to grow dark around them. "The wildlings are
dead."
```

**Data Explorer**

Version 1 (9.9 MB)

- 001ssb.txt
- 002ssb.txt
- 003ssb.txt
- 004ssb.txt
- 005ssb.txt

**Summary**
▸ 📁 5 files

Add the dataset in folder data where VS code directory present



**Code in jupter VS code:**

```
import numpy as np
import pandas as pd
```

```
!pip install gensim
```

```
import gensim
import os
```

```
!pip install nltk
```

```
data = "C:/Users/Syed Hamedoon/VScode Examples/data"
```

```
import nltk

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

```
import os
from nltk import sent_tokenize
from gensim.utils import simple_preprocess

DATA_PATH = r"C:\Users\Syed Hamedoon\VScode Examples\data"

story = []

for filename in os.listdir(DATA_PATH):
    if filename.endswith(".txt"):
        file_path = os.path.join(DATA_PATH, filename)

        try:
            with open(file_path, "r", encoding="utf-8") as f:
                corpus = f.read()
        except UnicodeDecodeError:
            with open(file_path, "r", encoding="cp1252") as f:
                corpus = f.read()

        for sent in sent_tokenize(corpus):
            story.append(simple_preprocess(sent))
```

```
print(len(story))
print(story[:2])
```

```
model = gensim.models.Word2Vec(
    window=10,
    min_count=2
)
```

```
model.build_vocab(story)
```

```
model.train(story, total_examples=model.corpus_count, epochs=model.epochs)
```

```
model.wv.most_similar('daenerys')
```

```
model.wv.doesnt_match(['jon','rikon','robb','arya','sansa','bran'])
```

```
model.wv.doesnt_match(['cersei', 'jaime', 'bronn', 'tyrion'])
```

```
model.wv['king']
```

```
model.wv.similarity('arya','sansa')
```

```python
model.wv.similarity('tywin','sansa')
```

```python
model.wv.get_normed_vectors()
```

```python
y = model.wv.index_to_key
```

```python
len(y)
```

```python
y
```

```python
from sklearn.decomposition import PCA
```

```python
pca = PCA(n_components=3)
```

```python
X = pca.fit_transform(model.wv.get_normed_vectors())
```

```python
X.shape
```

```python
!pip install --upgrade nbformat
```

```python
import pandas as pd
import plotly.express as px
import plotly.io as pio

pio.renderers.default = "browser"  # ← IMPORTANT

df = pd.DataFrame(X[200:300], columns=["x", "y", "z"])
df["label"] = y[200:300]
```
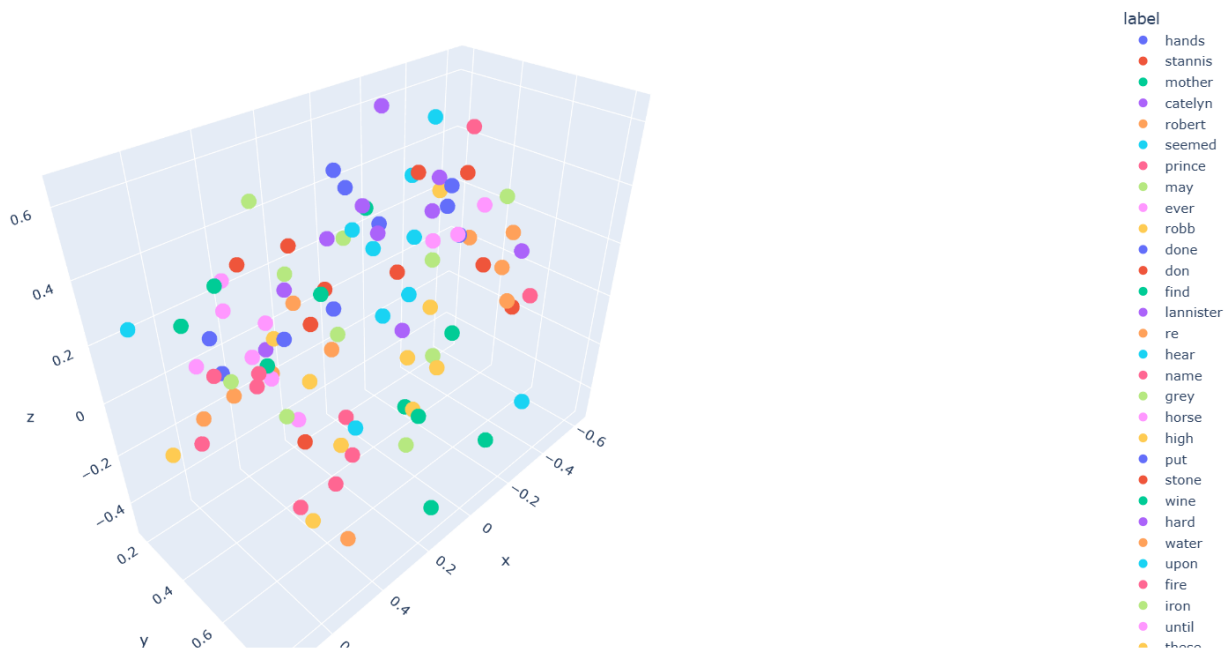
```
fig = px.scatter_3d(
    df,
    x="x",
    y="y",
    z="z",
    color="label"
)

fig.show()
```

Output:



## LAB Questions

1. What is the core idea behind Word2Vec?

2. Difference between CBOW and Skip-Gram?

3. Why is one-hot encoding inefficient?

4. Why do character names appear close in vector space?

5. How does window size affect semantic learning?

6. Why might rare characters have poor embeddings?

7. Which model performed better: CBOW or Skip-Gram? Why?

8. What happens if vector size is too small or too large?

9. Can Word2Vec understand word meaning without labels? Explain.

## LAB Assessment

| Student Name | | LAB Rubrics | CLO3 , P5, PLO5 |
|---|---|---|---|
| | | Total Marks | 10 |
| Registration No | | Obtained Marks | |
| | | Teacher Name | Dr. Syed M Hamedoon |
| Date | | Signature | |