

Precog Report

SANCHIT JALAN

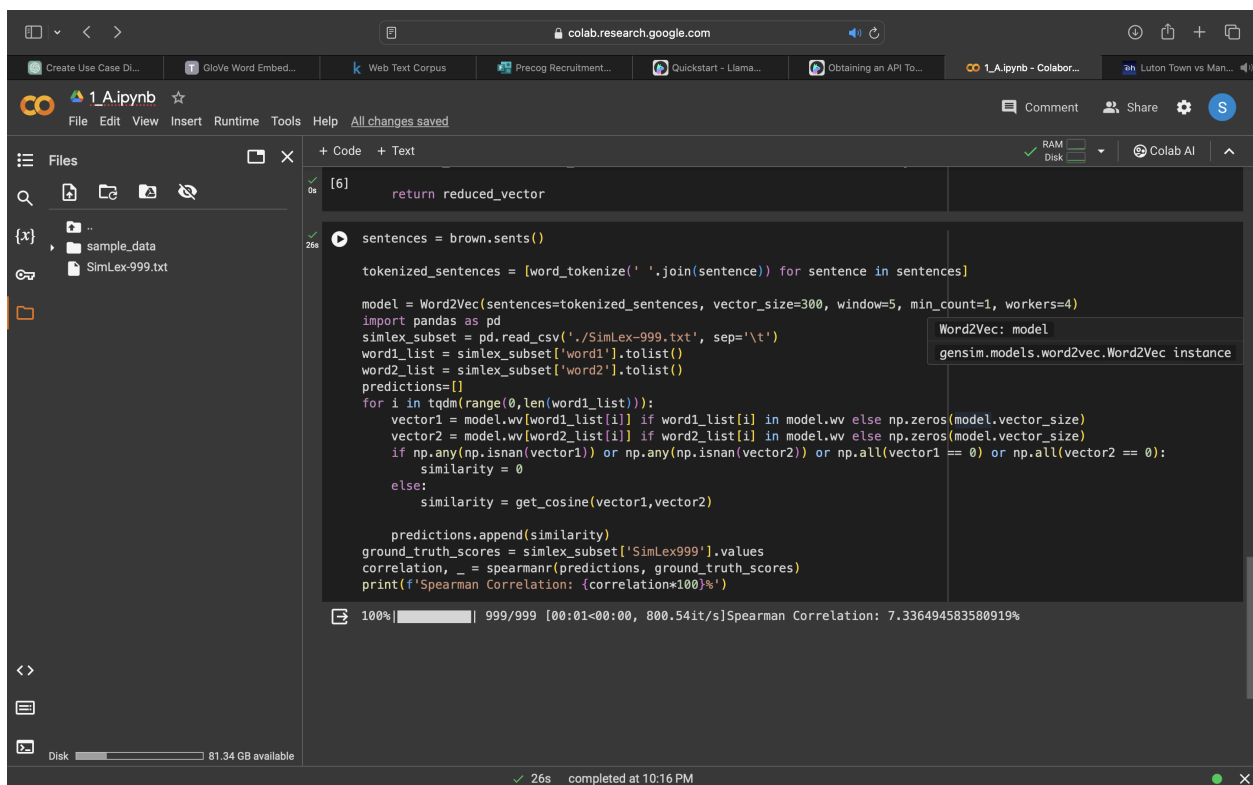
(2022101070)

Bonus part 3 was not done .

Task a

i) Constrained

In this task i imported a empty Word2Vec model . I trained my model on brown corpus of nltk library . Then i iterate over the dataset calculate the similarity of words using embeddings and cosine similarity and calculated the similarity using spearman coeffecient . It came out to be very less . This was not a very good way of training as my text corpus was not much .



The screenshot shows a Jupyter Notebook environment with a code cell containing the following Python code:

```
[6] return reduced_vector

sentences = brown.sents()

tokenized_sentences = [word_tokenize(' '.join(sentence)) for sentence in sentences]

model = Word2Vec(sentences=tokenized_sentences, vector_size=300, window=5, min_count=1, workers=4)

import pandas as pd
simlex_subset = pd.read_csv('./SimLex-999.txt', sep='\t')
word1_list = simlex_subset['word1'].tolist()
word2_list = simlex_subset['word2'].tolist()
predictions=[]
for i in tqdm(range(0,len(word1_list))):
    vector1 = model.wv[word1_list[i]] if word1_list[i] in model.wv else np.zeros(model.vector_size)
    vector2 = model.wv[word2_list[i]] if word2_list[i] in model.wv else np.zeros(model.vector_size)
    if np.any(np.isnan(vector1)) or np.any(np.isnan(vector2)) or np.all(vector1 == 0) or np.all(vector2 == 0):
        similarity = 0
    else:
        similarity = get_cosine(vector1,vector2)

    predictions.append(similarity)
ground_truth_scores = simlex_subset['SimLex999'].values
correlation, _ = spearmanr(predictions, ground_truth_scores)
print(f'Spearman Correlation: {correlation*100}%')
```

The output of the code cell shows a progress bar at 100% completion and the final result: Spearman Correlation: 7.336494583580919%.

A better approach would be implementing this model using CBOW using tensorflow and keras . Due to time constraint i was not able to implement this method .

ii) Unconstrained

In this method i imported Google Word2Vec pretrained model . I tested the dataset and found the spearman correlation as follows . It was better than my model and provided good results .

```
print("Mean Squared Error:", mse)
print("Mean Absolute Error:", mae)
print("Root Mean Squared Error:", rmse)
print("Pearson Correlation Coefficient:", pearson_corr)
print("Spearman Rank Correlation Coefficient:", spearman_corr)
```

✓ 0.0s

```
Mean Squared Error: 23.93555022540128
Mean Absolute Error: 4.1843771421902005
Root Mean Squared Error: 4.892397185981661
Pearson Correlation Coefficient: 0.45392821415513845
Spearman Rank Correlation Coefficient: 0.44196551091403796
```

Unconstrained works better because it was more fine tuned and more trained than the model i used .

Task b

- Similar approach was used for both phrases and sentences .

Imported the Word2Vec model .

Removed the stopwords from all phrases

Tried to use 2 different approaches for calculating embedding corresponding to a particular phrase .

Approach 1

Took average pooling of all the embeddings

Approach 2

Took a different approach of multiplying by a factor which was calculated using tfidf

In my conclusion Approach 1 was better , because we didn't have much data or corpus to make tfidf and same was concluded by testing it .

After making a embedding for a phrase , then cosine similarity was used .
A logistic regression model was trained on labels and cosine similarity .
Testing was done on the model and in the end accuracy was calculated .

Phrases

```
# Predict similarity using the trained logistic regression model
predicted_labels = classifier.predict(test_similarity.reshape(-1, 1))
# Evaluate performance
accuracy = accuracy_score(true_labels, predicted_labels)
print("Accuracy:", accuracy)
```

Test Dataset Loaded

Accuracy: 0.508

[/Library/Frameworks/Python.framework/Versions/3.10/lib/python3.10/site-packages/scipy](#)
 $\text{dist} = 1.0 - \frac{uv}{\sqrt{uu * vv}}$

Accuracy=50.8%

Sentences

```
# Convert lists to arrays
test_similarity = np.array(test_similarity)

# Predict similarity using the trained logistic regression model
predicted_labels = classifier.predict(test_similarity.reshape(-1, 1))
# Evaluate performance
accuracy = accuracy_score(true_labels, predicted_labels)
print("Accuracy:", accuracy)
```

[9] ✓ 7.9s

... Test Dataset Loaded

Accuracy: 0.558

Accuracy=55.8%

Bonus Task

i) Transformers

In this method sentence-transformer was used for calculating the embedding of a sentence . Xgb classifier was used as model for training and testing the data . This gave high accuracy because XGB and sentence transformer are more effecient .

```
# Evaluate performance
predicted_labels = xgb_classifier.predict(test_similarity.reshape(-1, 1))
accuracy = accuracy_score(true_labels, predicted_labels)
print("Accuracy:", accuracy)
```

```
100%|██████████| 8000/8000 [03:13<00:00, 41.29it/s]
Accuracy: 0.613125
```

Accuracy - 61.3%

ii) LLMs

Testing was only done on 100 phrases

1) Commercial LLMs

I used Gemini for my observation . It was equally effective for phrases as compared to fine tune model i made for phrases .

```
... 100%|██████████| 100/100 [07:19<00:00, 4.39s/it]
+ Code

test_accuracy = accuracy_score(ground_truth, results_list)
print("test accuracy: ", test_accuracy*100)

[18]

... test accuracy: 49.0
```

2) Open source LLMs

I used LLAMA for my observation . It was less effective for phrases as compared to fine tune model i made for phrases and also from Gemini API call .

```
except Exception as e:
    print(f"Error: {e}")
    results_list.append(0)

✓ 6m 58.1s

100%|██████████| 100/100 [06:58<00:00, 4.18s/it]

test_accuracy = accuracy_score(ground_truth, results_list)
print("test accuracy: ", test_accuracy*100)

✓ 0.0s

test accuracy: 44.0
```

Link for Paper Presentation

<https://docs.google.com/presentation/d/1GaQuf9ZYUANKMkO9oseWZOQzkwthAMYw3n7xQ0PCFuY/edit?usp=sharing>

