

CSE556: Natural Language Processing

Assignment-1

Himanshu Raj (2022216) | Ishita (2022224) | Ritika Thakur (2022408)

February 2, 2025

1 Task 1

1.1 Data Pre-processing

Our corpus `corpus.txt` had over 45k words which included all lower-case words with no number and punctuations. The corpus also included words like `http`, `src`, `img`, `href`, `www` which did not add anything to the contextual meaning of the words. The following pre-processing steps were applied:

- **Removing irrelevant words and trailing content:** Words such as `img`, `http`, `href`, and `src` along with everything after them were removed to eliminate unnecessary HTML-like tags and URLs. These words along with the trailing content were found to be irrelevant in the corpus and did not add to the context of the sentences.
- **Whitespace normalization:** Extra spaces were removed.

```
i feel like i am being punished for something that i didnt even do
the img file was completely broken and irrelevant something dumb
i really thought i was ok with how things are but here i am crying and feeling empty
nothing here makes sense but href them should be ignored
i left the hospital that night feeling helpless
check out this amazing site http nonsense words here
i still feel disappointed though
the href attribute was just pointing to garbage junk text
i was devastated and heartbroken never expecting this outcome
i found a link href to nowhere important gibberish words
i feel so worthless and ugly a href http afaerytaleinmakebelieve
why does every img have to be named something ridiculous
i was feeling extremely low this morning unable to focus on anything
an img tag with src pointing to absolutely nowhere something weird
i am excited to be introduced to a new kind of library environment but i feel stressed about it
this text contains src gibberish that means nothing
i was feeling pretty low about that but joan saw my disappointment and lifted my spirit
i clicked on an http link and it led to nonsense words
i do this if i allow myself to sit in this cycle today i will cause a nasty big blowup fight in public and feel humiliated
i mention this one doesnt feel fake
```

Figure 1: Example corpus before pre-processing

```
['i feel like i am being punished for something that i didnt even do', 'the', 'i really thought i was ok with how things are but here i am crying and feeling empty', 'nothing here makes sense but', 'i left the hospital that night feeling helpless', 'check out this amazing site', 'i still feel disappointed though', 'the', 'i was devastated and heartbroken never expecting this outcome', 'i found a link', 'i feel so worthless and ugly a', 'why does every', 'i was feeling extremely low i', 'is morning unable to focus on anything', 'an', 'i am excited to be introduced to a new kind of library environment but i feel stressed about it', 'this text contains', 'i was feeling pretty low about that but joan saw my disappointment and lifted my spirit', 'i clicked on an', 'i do this if i allow myself to sit in this cycle today i will cause a nasty big blowup fight in public and feel humiliated', 'i mention this one doesnt feel fake']
```

Figure 2: Example corpus after pre-processing

1.2 Vocabulary Generation

The vocabulary generation is based on the WordPiece subword tokenization approach (similar to BPE Algorithm).

1. **Splitting words into characters:** Each word is initially split into individual characters, with non-initial characters prefixed by `##`.

2. **Pair frequency calculation:** The frequency of adjacent character pairs is counted across the corpus.
3. **Pair scoring:** A score is calculated for each character pair based on its frequency relative to individual unit frequencies:
$$\text{freq}\{\text{pair}\} / (\text{freq}\{\text{word1}\} * \text{freq}\{\text{word2}\})$$
4. **Merging pairs iteratively:** The highest-scoring pairs are merged until the vocabulary size reaches the predefined limit.

[UNK] and [PAD] are also added to the vocabulary.

```
##vast
##vasta
##vastat
##w
##wup
##x
##xt
##y
##yc
##z
[PAD]
[UNK]
a
abou
about
allow
am
amaz
amazing
an
and
any
```

Figure 3: Snippet from vocabulary generated from example corpus with `size=500`

1.3 Tokenization

The tokenization process converts input text into subword units based on the generated vocabulary. After splitting words in the text into characters, as done during vocabulary generation, characters are merged based on the longest match in the vocabulary. If a word cannot be matched, it is replaced with [UNK].

```
[
  {
    "id": 0,
    "sentence": "i mention this one doesnt feel fake\n"
  },
  {
    "id": 1,
    "sentence": "not in the corpus zzz\n"
  }
]
```

Figure 4: Test json file

```
{
  "0": [
    "i",
    "mention",
    "this",
    "one",
    "doesnt",
    "feel",
    "fake"
  ],
  "1": [
    "not",
    "in",
    "the",
    "co",
    "##r",
    "##p",
    "##u",
    "##s",
    "[UNK]",
    "##Z",
    "##z"
  ]
}
```

Figure 5: Tokenized json file

2 Task 2

2.1 Introduction

We build a Word2Vec model using the Continuous Bag of Words (CBOW) approach. The goal is to preprocess the input text corpus, generate CBOW pairs, and use them to train the Word2Vec model. Also, we find similar and dissimilar word triplets using cosine similarity.

2.2 Dataset Preparation - Word2VecDataset Class

The **Word2VecDataset** class preprocesses the input text corpus for training the Word2Vec model using the **Continuous Bag of Words (CBOW)** approach.

1. **Initialization** : It initializes the window size for context words, vocabulary, and mappings between words and indices (word2idx and idx2word). It also initializes CBOW training pairs.
2. **Pre-processes Data** :
 - (a) Uses the WordPieceTokenizer from Task 1 to generate the vocabulary from the input text.

- (b) Tokenizes the input text corpus using the WordPieceTokenizer from Task 1.
- (c) Generates mapping from each word to a unique index and vice versa
- (d) Generates CBOW training pairs based on the tokenized corpus.

```
{
  "0": [
    "i",
    "stand",
    "h",
    "##e",
    "##r",
    "##e",
    "i",
    "c",
    "##e",
    "##e",
    "##l",
    "empty",
    "a",
    "class",
    "post",
    "count",
    "link"
  ],
  "1": [
    "i",
    "lit",
    "##e",
    "##rall",
    "##y",
    "just",

```

Figure 6: Tokenized Corpus

```
{
  "word2idx": {
    "##a": 0, "##aachan": 1, "##ab": 2, "##abb": 3, "##abbing": 4, "##abl": 5, "##ablsh": 6, "##ably": 7,
    "##abour": 8, "##abric": 9, "##abulous": 10, "##aby": 11, "##ac": 12, "##acation": 13, "##ach": 14, "##achan": 15,
    "##aching": 16, "##achm": 17, "##acious": 18, "##ack": 19, "##acking": 20, "##aconic": 21, "##acrific": 22, "##act": 23,
    "##actic": 24, "##actica": 25, "##actical": 26, "##acting": 27, "##action": 28, "##actions": 29, "##activ": 30, "##actu":
    31, "##acul": 32, "##ad": 33, "##ada": 34, "##adach": 35, "##aday": 36, "##adha": 37, "##adia": 38, "##adiation": 39,
    "##adiator": 40, "##adical": 41, "##ading": 42, "##adings": 43, "##adio": 44, "##adition": 45, "##adow": 46, "##adphon":
    "##adua": 48, "##adually": 49, "##aduation": 50, "##ady": 51, "##af": 52, "##aff": 53, "##affi": 54, "##affic": 55,

```

Figure 7: Word to Index Mapping

```
zoo": 8497, "zoom": 8498, "zooming": 8499},
  "idx2word": {
    "0": "##a", "1": "##aachan", "2": "##ab", "3": "##abb", "4":
    "##abbing", "5": "##abl", "6": "##ablsh", "7": "##ably", "8": "##abour", "9": "##abric", "10": "##abulous", "11": "##aby",
    "12": "##ac", "13": "##acation", "14": "##ach", "15": "##achan", "16": "##aching", "17": "##achm", "18": "##acious", "19":
    "##ack", "20": "##acking", "21": "##aconic", "22": "##acrific", "23": "##act", "24": "##actic", "25": "##actica", "26":
    "##actical", "27": "##acting", "28": "##action", "29": "##actions", "30": "##activ", "31": "##actu", "32": "##acul", "33":
    "##ad", "34": "##ada", "35": "##adach", "36": "##aday", "37": "##adha", "38": "##adia", "39": "##adiation", "40":
    "##adiator", "41": "##adical", "42": "##ading", "43": "##adings", "44": "##adio", "45": "##adition", "46": "##adow", "47":

```

Figure 8: Index to Word Mapping

2.3 Model Architecture - Word2VecModel Class

1. **Initialization** : Initializes the vocabulary size, embedding dimension for each word and weights using a uniform distribution. Bias terms are set to zero.
2. **Architecture Details**
 - **Embedding Layer**: Maps words to their corresponding vector representations.
 - **Linear Layer**: Predicts the target word given the context of word embeddings.

- **Log Softmax Activation:** Outputs probability distribution over the entire vocabulary.
3. **Forward Function :** The forward pass gets the embeddings for the context words, averages them, and sends the result to the output layer. The softmax function is then applied to give a probability distribution over all the words in the vocabulary.
 4. **Calculates Cosine Similarity :** Calculates cosine similarity between two vectors and also identifies the most and least similar vector to a given vector.

2.4 Training

1. **Initialization :** Initializes Adam optimizer, `NLLLoss()` for calculating loss and moves the model to the appropriate device (CPU/GPU). At the same time it also initializes history to keep track of loss and accuracy.
2. **Training :** The function trains the Word2Vec CBOW model. In every iteration, it decreases the learning rate and processes batches of context-target pairs. For each batch, the model performs a forward pass, computes the loss, and applies backpropagation to update the weights. Training progress, including batch loss and accuracy, is displayed using the progress bar.
3. **Model Validation :** After each epoch it computes validation loss, validation accuracy and average cosine similarity for the batch and updates the history.
4. **Updates Final Checkpoint and Saves the Final Model**

2.5 Finding Triplets

The `find_triplets` function creates triplets consisting of two similar words and one dissimilar word based on cosine similarity. It randomly selects a word, retrieves the two most similar words, and identifies the least similar word using the Word2Vec model.

Triplet 1:

- **High positive similarity:** "#run" has a similarity of **0.8747** with "haul" and **0.8696** with "#ashion" ighlighting their frequent occurrence in similar contexts.
- **Negative similarity:** "#run" and "uhur" have a similarity of **-0.9269**, suggesting they occur in opposite contexts or are rarely used together.

Triplet 2:

- **Very high positive similarity:** "#logic" has a similarity of **0.9312** with "##rful" and **0.9026** with "entry," highlighting their frequent occurrence in similar contexts.
- **Negative similarity:** "#logic" and "#aomi" have a similarity of **-0.8795**, suggesting they occur in opposite contexts or are rarely used together.

```

Triplet 1:
Token - ##run (index: 2296)
Similar tokens - haul (index: 5482), ##ashion (index: 285)
Dissimilar token - uhur (index: 8051)
Cosine similarity (##run, haul): 0.8747
Cosine similarity (##run, ##ashion): 0.8696
Cosine similarity (##run, uhur): -0.9269

Triplet 2:
Token - logic (index: 6054)
Similar tokens - ##rful (index: 2054), entry (index: 4856)
Dissimilar token - ##aomi (index: 211)
Cosine similarity (logic, ##rful): 0.9312
Cosine similarity (logic, entry): 0.9026
Cosine similarity (logic, ##aomi): -0.8795

```

Figure 9: Identifying Triplets

2.6 Experimental Results

After experimenting with the values of hyperparameters, the following hyperparameters gave the best results:

WINDOW_SIZE = 4

EMBEDDING_DIM = 10

BATCH_SIZE = 256

NUM_EPOCHS = 15

LEARNING_RATE = 0.02

TRAIN_SPLIT = 0.8

VOCAB_SIZE = 8500

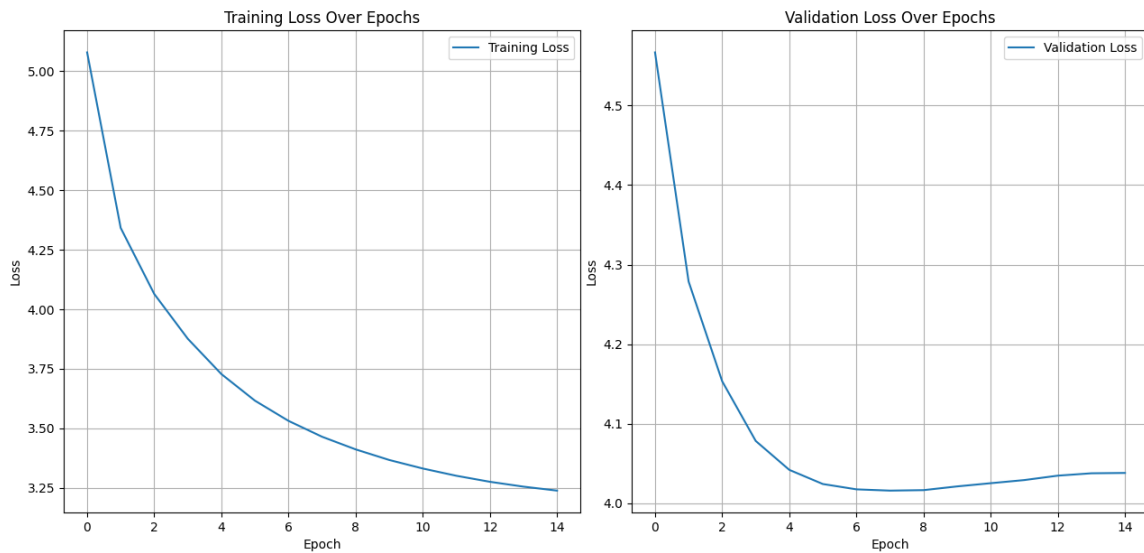


Figure 10: Training Loss and Validation Loss

```
Epoch 15/15:  
Average Loss: 3.2390  
Training Accuracy: 42.68%  
  
Validating model...  
Validation Loss: 4.0382  
Validation Accuracy: 39.04%  
Cosine Similarity: 0.4233977531685549  
Final model saved: task2-files/word2vec_checkpoints\word2vec_final_model.pt
```

Figure 11: Best Validation Accuracy using the given hyperparameters

3 Task 3

3.1 Introduction

For all architectures, input layer contains 40 neurons (contextwindow x embeddingdimension = 4x10) and output layer contains 8500 neurons (the vocab size used to train Word2Vec model). All models were trained for 5 epochs with a learning rate of 0.01 with a batch size of 1024 samples over a training set with 60,000+ samples. A common observation over these models was that validation loss increased significantly after 6th epoch, even with smaller learning rates. The training time and compute increased as we go from NeuralLM1 to NeuralLM2 to NeuralLM3.

3.2 Design Justifications and Performance comparisons

NeuralLM1: 1 hidden layer with 1024 neurons and tanh activation

Justification: This is a basic MLP based Neural Language Model architecture proposed by Bengio in 2003. This serves as a good baseline model to start with when training neural language models.

NeuralLM2: 1 hidden layer with 1024 neurons and ReLU activation and skip connection

Justification: The tanh activation outputs cap at a value of 1 which can lead to less/slower updates in learnable parameters. The ReLU activation output is unbounded and hence is used in an attempt to improve accuracy and perplexity. Skip connections help in tackling vanishing gradients, proposed by Bengio in 2003 but they were optional.

Improvement: This led to better training and validation accuracy than NeuralLM1. Training perplexity decreased, but validation perplexity increased signifying model is unable to generalize on unseen data.

NeuralLM3: 2 hidden layers with 512 and 2048 neurons and ReLU activation

Justification: In previous 2 architectures, validation loss increases after 6th epoch indicating model is unable to generalize on unseen data. Hence, another hidden layer is introduced in an attempt to capture the features in a better way, so that the model can generalise well on unseen data. We are not adding skip connection like NeuralLM2 because we saw an increase in validation perplexity because of which model is unable to generalize well on unseen data.

Improvement: This led to better training and validation accuracy than NeuralLM2. The validity perplexity is even less than NeuralLM1, signifying the best performance on unseen data out of all architectures. Training perplexity is a bit more than NeuralLM2, but still less than NeuralLM1.

3.3 Loss vs Epoch graphs

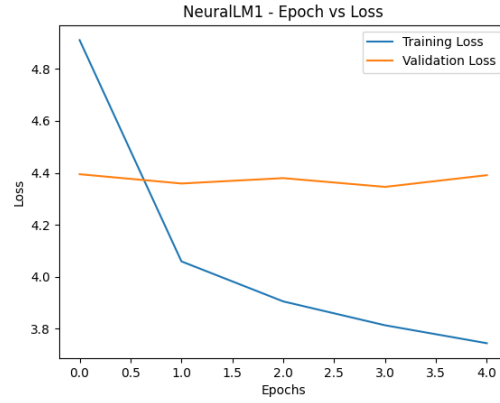


Figure 12: Loss vs Epoch curve for NeuralLM1

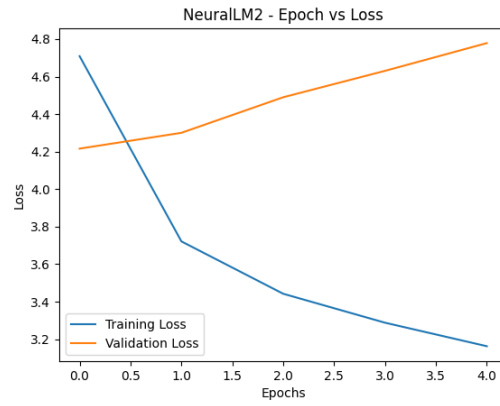


Figure 13: Loss vs Epoch curve for NeuralLM2

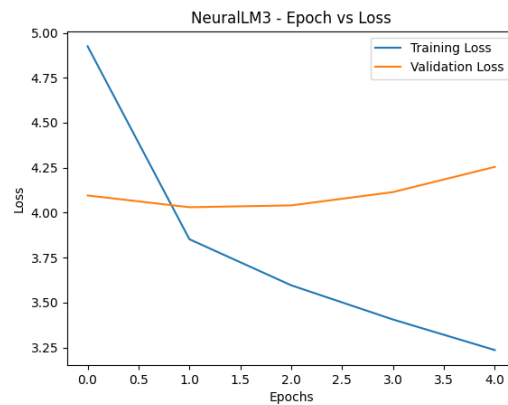


Figure 14: Loss vs Epoch curve for NeuralLM3

3.4 Accuracies and Perplexity scores

Model	Accuracy		Perplexity	
	Training	Validation	Training	Validation
NeuralLM1	36.01%	35.97%	42.27	80.71
NeuralLM2	39.73%	37.81%	23.63	118.94
NeuralLM3	39.17%	39.37%	25.43	70.43

```
(nlperv) PS D:\nlp\assignment1> python task3.py
Epoch 1/5: 100%|
Train Loss: 4.9109, Train Accuracy: 31.16%, Train Perplexity: 135.76
Val Loss: 4.3946, Val Accuracy: 34.24%, Val Perplexity: 81.01

Epoch 2/5: 100%|
Train Loss: 4.0591, Train Accuracy: 34.02%, Train Perplexity: 57.92
Val Loss: 4.3589, Val Accuracy: 34.83%, Val Perplexity: 78.17

Epoch 3/5: 100%|
Train Loss: 3.9050, Train Accuracy: 34.98%, Train Perplexity: 49.65
Val Loss: 4.3796, Val Accuracy: 35.28%, Val Perplexity: 79.81

Epoch 4/5: 100%|
Train Loss: 3.8131, Train Accuracy: 35.47%, Train Perplexity: 45.29
Val Loss: 4.3459, Val Accuracy: 36.34%, Val Perplexity: 77.16

Epoch 5/5: 100%|
Train Loss: 3.7440, Train Accuracy: 36.01%, Train Perplexity: 42.27
Val Loss: 4.3909, Val Accuracy: 35.97%, Val Perplexity: 80.71
```

Figure 15: Loss, Accuracy and Perplexity over epochs for NeuralLM1

```
Epoch 1/5: 100%|
Train Loss: 4.7089, Train Accuracy: 33.75%, Train Perplexity: 110.93
Val Loss: 4.2161, Val Accuracy: 37.03%, Val Perplexity: 67.77

Epoch 2/5: 100%|
Train Loss: 3.7215, Train Accuracy: 36.81%, Train Perplexity: 41.32
Val Loss: 4.3004, Val Accuracy: 37.81%, Val Perplexity: 73.73

Epoch 3/5: 100%|
Train Loss: 3.4423, Train Accuracy: 37.47%, Train Perplexity: 31.26
Val Loss: 4.4900, Val Accuracy: 37.80%, Val Perplexity: 89.12

Epoch 4/5: 100%|
Train Loss: 3.2880, Train Accuracy: 38.61%, Train Perplexity: 26.79
Val Loss: 4.6308, Val Accuracy: 38.09%, Val Perplexity: 102.60

Epoch 5/5: 100%|
Train Loss: 3.1627, Train Accuracy: 39.73%, Train Perplexity: 23.63
Val Loss: 4.7786, Val Accuracy: 37.81%, Val Perplexity: 118.94
```

Figure 16: Loss, Accuracy and Perplexity over epochs for NeuralLM2

Epoch 1/5: 100%	
Train Loss: 4.9253, Train Accuracy: 30.66%, Train Perplexity: 137.73	
Val Loss: 4.0956, Val Accuracy: 36.80%, Val Perplexity: 60.07	
Epoch 2/5: 100%	
Train Loss: 3.8522, Train Accuracy: 36.60%, Train Perplexity: 47.10	
Val Loss: 4.0302, Val Accuracy: 37.84%, Val Perplexity: 56.27	
Epoch 3/5: 100%	
Train Loss: 3.5962, Train Accuracy: 37.60%, Train Perplexity: 36.46	
Val Loss: 4.0404, Val Accuracy: 38.46%, Val Perplexity: 56.85	
Epoch 4/5: 100%	
Train Loss: 3.4058, Train Accuracy: 38.57%, Train Perplexity: 30.14	
Val Loss: 4.1145, Val Accuracy: 38.84%, Val Perplexity: 61.22	
Epoch 5/5: 100%	
Train Loss: 3.2359, Train Accuracy: 39.17%, Train Perplexity: 25.43	
Val Loss: 4.2547, Val Accuracy: 39.37%, Val Perplexity: 70.43	

Figure 17: Loss, Accuracy and Perplexity over epochs for NeuralLM3

4 Individual Contribution

- Himanshu Raj: Task 3 implementation and report
- Ishita: Task 2 implementation and report
- Ritika Thakur: Task 1 implementation and report

5 References

- WordPieceTokenizer Tutorial
- Word2Vec paper
- Custom dataset Tutorial
- PyTorch Implementations
- Skip-connections
- Lecture Slides - Language Modelling, Word Representation