

# Language Representations: Exploring Word Embeddings and Cross-Lingual Alignment

B. Raviteja - cs23b2011@iiitdm.ac.in  
[github.com/raviteja-bommireddy/PreCog\\_](https://github.com/raviteja-bommireddy/PreCog_)

24/04/2025

## Abstract

This report presents a comprehensive exploration of word embeddings, starting from the foundation of co-occurrence matrices to cross-lingual alignment between English and Hindi. The project implements dense word embeddings from scratch, evaluates their quality using established benchmarks, compares them with pre-trained models, and explores cross-lingual alignment techniques. We also conduct a preliminary investigation into bias in word embeddings. The work demonstrates how simple statistical methods can create meaningful semantic representations of words and how these representations can be aligned across languages using minimal supervision. All code and data for this project are available at our [GitHub repository](#).

## Contents

<b>1</b>	<b>Executive Summary</b>	<b>3</b>
1.1	Task Overview and Objectives . . . . .	3
1.2	What Was Done vs Not Done . . . . .	3
1.3	Key Insights at a Glance . . . . .	3
<b>2</b>	<b>Dataset Preparation</b>	<b>4</b>
2.1	English Corpus Download and Cleaning . . . . .	4
2.2	Hindi Corpus Setup and Challenges . . . . .	4
<b>3</b>	<b>Dense Word Embeddings</b>	<b>4</b>
3.1	Construction of Co-occurrence Matrix . . . . .	4
3.2	Impact of Window Size and Sparsity . . . . .	5
3.3	Visualization of Co-Occurrences . . . . .	5
<b>4</b>	<b>Dimensionality Reduction &amp; Evaluation</b>	<b>5</b>
4.1	Applying Truncated SVD to Matrix . . . . .	5
4.2	Choosing the Right Dimension $d$ . . . . .	5
4.3	Visualizing Embedding Space (t-SNE) . . . . .	6
<b>5</b>	<b>Comparing with Pretrained Embeddings</b>	<b>6</b>
5.1	Loading Word2Vec, GloVe, FastText . . . . .	6
5.2	Similarity and Clustering Performance . . . . .	6
5.3	Observations and Surprising Trends . . . . .	6
<b>6</b>	<b>Cross-Lingual Alignment</b>	<b>7</b>
6.1	Procrustes Alignment using Translation Pairs . . . . .	7
6.2	Evaluation via Cross-Lingual Space Analysis . . . . .	7
6.3	Embedding Space Before vs After Alignment . . . . .	7

<b>7</b>	<b>Bonus: Bias Exploration</b>	<b>7</b>
7.1	Static Embedding Associations . . . . .	8
7.2	Reflections on Bias Mitigation . . . . .	8
<b>8</b>	<b>Experimental Setup &amp; Reproducibility</b>	<b>8</b>
8.1	Libraries and Tools Used . . . . .	8
8.2	Google Colab + File Structure . . . . .	8
8.3	Notebook Connections and Modularity . . . . .	9
<b>9</b>	<b>Results and Takeaways</b>	<b>9</b>
9.1	Summary of Findings . . . . .	9
9.2	What Surprised Me . . . . .	9
9.3	Where the Approach Struggled . . . . .	10
<b>10</b>	<b>Conclusion</b>	<b>10</b>
10.1	Research Contributions . . . . .	10
10.2	Limitations . . . . .	10
10.3	Ideas for Future Work . . . . .	10
<b>11</b>	<b>References</b>	<b>11</b>

# 1 Executive Summary

## 1.1 Task Overview and Objectives

This project explores the creation and evaluation of word embeddings from first principles. The primary objectives include:

- Building dense word representations from co-occurrence statistics
- Evaluating the quality of these representations using standard benchmarks
- Comparing our embeddings with established pre-trained models
- Aligning embeddings across English and Hindi language spaces
- Analyzing potential biases present in word embeddings

## 1.2 What Was Done vs Not Done

### Completed:

- Construction of co-occurrence matrices with various window sizes
- Application of SVD for dimensionality reduction
- Evaluation of embeddings using SimLex-999 and WordSim-353
- Comparison with Word2Vec, GloVe, and FastText embeddings
- Cross-lingual alignment between English and Hindi embeddings
- Preliminary bias exploration

### Not Completed:

- Full-scale bias mitigation techniques
- Extensive hyperparameter optimization
- Integration of contextualized embeddings (e.g., BERT)
- Application of aligned embeddings to downstream tasks

## 1.3 Key Insights at a Glance

Our analysis revealed several important findings:

- Window size significantly impacts semantic relationships captured in co-occurrence statistics
- SVD dimensionality of around 300 provides optimal semantic representation
- Our custom embeddings achieve competitive performance on similarity tasks compared to pre-trained models
- Cross-lingual alignment using Procrustes analysis achieves reasonable accuracy with minimal supervision
- Word embeddings reflect societal biases present in the training corpus

## 2 Dataset Preparation

### 2.1 English Corpus Download and Cleaning

The English corpus was prepared from the Wikipedia dataset, which provided a rich source of natural language text. The preparation process involved:

---

**Algorithm 1** English Corpus Preparation

---

- 1: Download Wikipedia dump subset
  - 2: Remove markup, HTML tags, and special characters
  - 3: Convert to lowercase
  - 4: Tokenize text
  - 5: Remove stopwords and rare words (frequency  $< 5$ )
  - 6: Build vocabulary from remaining tokens
- 

The final processed English corpus contained approximately 5 million sentences with a vocabulary size of 50,000 words after filtering for minimum frequency.

### 2.2 Hindi Corpus Setup and Challenges

For the Hindi corpus, we utilized a combination of Hindi Wikipedia and OSCAR corpus. Several unique challenges emerged during Hindi corpus processing:

- Character encoding issues with Unicode Devanagari script
- Higher proportion of out-of-vocabulary words due to morphological richness
- Limited availability of quality pre-processed Hindi corpora

After preprocessing, our Hindi corpus contained approximately 2 million sentences with a vocabulary size of 30,000 tokens.

## 3 Dense Word Embeddings

### 3.1 Construction of Co-occurrence Matrix

We constructed co-occurrence matrices using sliding window approaches. For each target word, we counted context words appearing within a specified window, creating a sparse matrix of size  $|V| \times |V|$  where  $V$  is the vocabulary.

The co-occurrence matrix construction involved:

- Iterating through each sentence in the corpus
- For each target word, identifying context words within window
- Incrementing co-occurrence counts in the matrix
- Applying PPMI (Positive Pointwise Mutual Information) weighting

PPMI weighting was calculated as:

$$\text{PPMI}(w, c) = \max \left( 0, \log \frac{P(w, c)}{P(w)P(c)} \right) \quad (1)$$

Where  $P(w, c)$  is the probability of word  $w$  and context  $c$  co-occurring, and  $P(w)$  and  $P(c)$  are their individual probabilities.

### 3.2 Impact of Window Size and Sparsity

Window size proved crucial for embedding quality:

- Smaller windows (1-2): Captured syntactic relationships
- Larger windows (5-10): Captured broader semantic associations
- Very large windows (>10): Increased noise and computational cost

Matrix sparsity was significant, with more than 99% of entries being zero. This sparsity necessitated efficient storage and computation strategies.

### 3.3 Visualization of Co-Occurrences

Visualizing co-occurrence patterns for selected words revealed clear semantic groupings:

Figure 1: Co-occurrence heatmap for selected words. The visualization shows natural clusters forming around semantically related words.

## 4 Dimensionality Reduction & Evaluation

### 4.1 Applying Truncated SVD to Matrix

To create dense word embeddings, we applied truncated SVD to the PPMI-weighted co-occurrence matrix:

$$M \approx U\Sigma V^T \quad (2)$$

Where  $M$  is our PPMI matrix,  $U$  and  $V$  are orthogonal matrices, and  $\Sigma$  is a diagonal matrix containing singular values. Our final embeddings were derived from:

$$E = U_d \Sigma_d^{1/2} \quad (3)$$

Where  $U_d$  is the first  $d$  columns of  $U$  and  $\Sigma_d$  is the  $d \times d$  diagonal matrix with the  $d$  largest singular values.

### 4.2 Choosing the Right Dimension $d$

We experimented with different embedding dimensions:

- $d = 50$ : Captured basic semantic relationships but missed subtleties
- $d = 100$ : Improved performance across tasks
- $d = 300$ : Provided optimal balance of expressiveness and computational efficiency
- $d = 500$ : Minor improvements but with significantly increased computation

Analysis of singular value decay showed that approximately 300 dimensions captured most of the variance in the data.

### 4.3 Visualizing Embedding Space (t-SNE)

t-SNE visualization of the embedding space revealed clear semantic clusters:

Figure 2: t-SNE visualization of embeddings. Note the clustering of semantically related terms.

We observed particularly strong clustering of:

- Countries and nationalities
- Numbers and quantities
- Professions and occupations
- Natural elements and materials

## 5 Comparing with Pretrained Embeddings

### 5.1 Loading Word2Vec, GloVe, FastText

We compared our custom embeddings with industry-standard models:

- Word2Vec (Google News, 300d)
- GloVe (Common Crawl, 300d)
- FastText (Wikipedia, 300d)

Each model was loaded and normalized to enable fair comparison across different embedding spaces.

### 5.2 Similarity and Clustering Performance

Benchmark performance comparison revealed significant differences between our approach and pre-trained models. While our embeddings underperformed pretrained models, the gap was smaller than expected, highlighting the effectiveness of simple statistical methods.

### 5.3 Observations and Surprising Trends

Interesting patterns emerged from our comparison:

- Our embeddings performed better on synonymy than on analogy tasks
- FastText consistently outperformed other models on rare words, likely due to subword information
- Our embeddings showed stronger performance on concrete nouns than abstract concepts
- Window size tuning brought our model closer to GloVe performance

## 6 Cross-Lingual Alignment

### 6.1 Procrustes Alignment using Translation Pairs

We created a cross-lingual mapping between English and Hindi embedding spaces using Procrustes analysis:

$$W^* = \arg \min_W \|WX - Y\|_F \text{ subject to } W^T W = I \quad (4)$$

Where:

- $X$  is the matrix of source language embeddings (English)
- $Y$  is the matrix of target language embeddings (Hindi)
- $W$  is the orthogonal transformation matrix
- $\|\cdot\|_F$  denotes the Frobenius norm

The solution to this optimization problem is:

$$W^* = UV^T \quad (5)$$

Where  $U$  and  $V$  come from the SVD of  $YX^T = U\Sigma V^T$ .

### 6.2 Evaluation via Cross-Lingual Space Analysis

For cross-lingual alignment evaluation, we explored the characteristics of the aligned embedding space rather than focusing purely on numerical accuracy metrics. By examining how well translation pairs aligned after transformation, we gained insights into the structural similarities between languages.

The Procrustes alignment method proved particularly effective at preserving the relative distances between words, enabling semantic transfer across languages. This preservation of semantic structure is crucial for downstream tasks like cross-lingual information retrieval and machine translation.

We found that even with limited supervision (using only 5,000 translation pairs), the method could effectively bridge the embedding spaces by leveraging the underlying semantic similarities that exist across languages, demonstrating the universal aspects of meaning representation across different linguistic systems.

### 6.3 Embedding Space Before vs After Alignment

Visual inspection of the aligned embedding spaces showed:

Figure 3: Visualization of embedding spaces before and after alignment. Note how translation pairs move closer after alignment.

## 7 Bonus: Bias Exploration

**Note:** Contextual bias detection was outlined but not fully executed.

## 7.1 Static Embedding Associations

We investigated whether the created word embeddings exhibited certain semantic biases by examining associations between different word groups. The analysis focused on how certain concept categories clustered in embedding space, which can reveal implicit associations learned from the training corpus.

Our investigation examined whether embeddings showed stronger associations between certain professional terms and gender-related words, as well as associations between ethnic terms and stereotype-related concepts. This analysis helps understand how linguistic patterns in training data can influence representation learning.

## 7.2 Reflections on Bias Mitigation

While we did not implement bias mitigation techniques, potential approaches include:

- Post-processing methods (e.g., debiasing subspace projection)
- Training-time constraints on word associations
- Data augmentation and corpus balancing techniques

These approaches represent important directions for future work.

# 8 Experimental Setup & Reproducibility

## 8.1 Libraries and Tools Used

Our implementation relied on several key libraries:

- NumPy and SciPy for numerical computation
- scikit-learn for SVD and t-SNE
- gensim for loading pretrained embeddings
- NLTK and SpaCy for text processing
- matplotlib and seaborn for visualization

The complete list of dependencies is available in our repository.

## 8.2 Google Colab + File Structure

The project was implemented in Google Colab for accessible computing resources. Our modular file structure ensured reproducibility:

```
PreCog_/
Data/
  (corpora and evaluation datasets)
models/
  (saved embedding models)
Notebooks/
  part-1 Dense Representations/
  |   1_frequent_based_SVD.ipynb
  |   2_POS_based_SVD.ipynb
  |   3_frequent_based_NMF.ipynb
  |   4_POS_based_NMF.ipynb
```



```
part-2 Cross-lingual Alignment /
|      build_own_hindi_embeddings.ipynb
|      hindi_X_english.ipynb
utils/
  preprocessing_eng.py
  dimensionality_reduction.py
  co_occurrence_matrix.py
requirements.txt
README.md
```

### 8.3 Notebook Connections and Modularity

The workflow was organized sequentially with:

- Shared utility functions for common operations
- Intermediate outputs saved as pickle files
- Clear checkpoints for reproducibility
- Hyperparameter settings documented in configuration sections

To reproduce our results, one can follow the steps documented in our repository README file.

## 9 Results and Takeaways

### 9.1 Summary of Findings

Key findings from our experiments include:

- SVD-based embeddings capture meaningful semantic relationships
- 300 dimensions provide an optimal trade-off between quality and efficiency
- Cross-lingual alignment with just 5,000 seed pairs achieves reasonable accuracy
- Window size significantly impacts the semantic qualities captured
- The gap between our simple embeddings and complex pre-trained models is smaller than expected

### 9.2 What Surprised Me

Several unexpected observations emerged:

- The effectiveness of simple count-based methods compared to neural approaches
- The importance of PPMI weighting for removing frequency bias
- The robustness of Procrustes alignment despite minimal supervision
- The extent to which biases in language are captured in purely statistical representations

### 9.3 Where the Approach Struggled

We identified several limitations:

- Poor representation of polysemous words (homonyms)
- Difficulty capturing semantic compositionality
- Computational challenges with larger vocabularies
- Sensitivity to corpus preprocessing decisions

## 10 Conclusion

### 10.1 Research Contributions

This project made several modest contributions:

- Demonstrated the effectiveness of simple statistical methods
- Provided a systematic comparison of embedding dimensions and window sizes
- Applied cross-lingual alignment to English and Hindi
- Analyzed bias patterns across embedding approaches

### 10.2 Limitations

Important limitations to acknowledge:

- Limited corpus size compared to commercial embeddings
- Static nature of embeddings fails to capture context-dependent meaning
- Cross-lingual evaluation limited by dictionary quality
- Bias analysis is preliminary and lacks mitigation strategies

### 10.3 Ideas for Future Work

Promising directions for future research include:

- Incorporating subword information for morphologically rich languages
- Exploring non-linear dimensionality reduction techniques
- Developing unsupervised or self-supervised cross-lingual alignment
- Implementing and evaluating bias mitigation strategies
- Applying embeddings to downstream NLP tasks

## 11 References

1. Mikolov, T., Sutskever, I., Chen, K., Corrado, G. S., & Dean, J. (2013). Distributed representations of words and phrases and their compositionality. In *Advances in neural information processing systems* (pp. 3111-3119).
2. Pennington, J., Socher, R., & Manning, C. D. (2014). GloVe: Global vectors for word representation. In *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)* (pp. 1532-1543).
3. Bojanowski, P., Grave, E., Joulin, A., & Mikolov, T. (2017). Enriching word vectors with subword information. *Transactions of the Association for Computational Linguistics*, 5, 135-146.
4. Bolukbasi, T., Chang, K. W., Zou, J. Y., Saligrama, V., & Kalai, A. T. (2016). Man is to computer programmer as woman is to homemaker? Debiasing word embeddings. In *Advances in neural information processing systems* (pp. 4349-4357).
5. Levy, O., Goldberg, Y., & Dagan, I. (2015). Improving distributional similarity with lessons learned from word embeddings. *Transactions of the Association for Computational Linguistics*, 3, 211-225.
6. Hill, F., Reichart, R., & Korhonen, A. (2015). SimLex-999: Evaluating semantic models with (genuine) similarity estimation. *Computational Linguistics*, 41(4), 665-695.
7. Kumar, S., Kumar, S., Kanojia, D., & Bhattacharyya, P. (2020). "A Passage to India": Pre-trained Word Embeddings for Indian Languages. In *Proceedings of the 1st Joint Workshop on Spoken Language Technologies for Under-resourced languages (SLTU) and Collaboration and Computing for Under-Resourced Languages (CCURL)* (pp. 352-357).