

# Vector representation with a finance corpus

Student: Ritesh Bachhar, University of Rhode Island  
riteshbachhar@uri.edu

Mentor: Murat Aydogdu, Rhode Island College  
maydogdu@ric.edu

Date: August 10, 2022



# Vector representation with a finance corpus

- Timeframe
  - May 15th, 2022
  - August, 2022



# Vector representation with a finance corpus

- Goals
  - Vector representation & NLP
  - Wikipedia Articles<sup>1</sup>
  - 10-K documents of public traded companies<sup>2</sup>
  - GloVe<sup>3</sup> implementation of word vector
  - Setup workflow in HPC

<sup>1</sup> <https://dumps.wikimedia.org/enwiki/latest/enwiki-latest-pages-articles.xml.bz2>

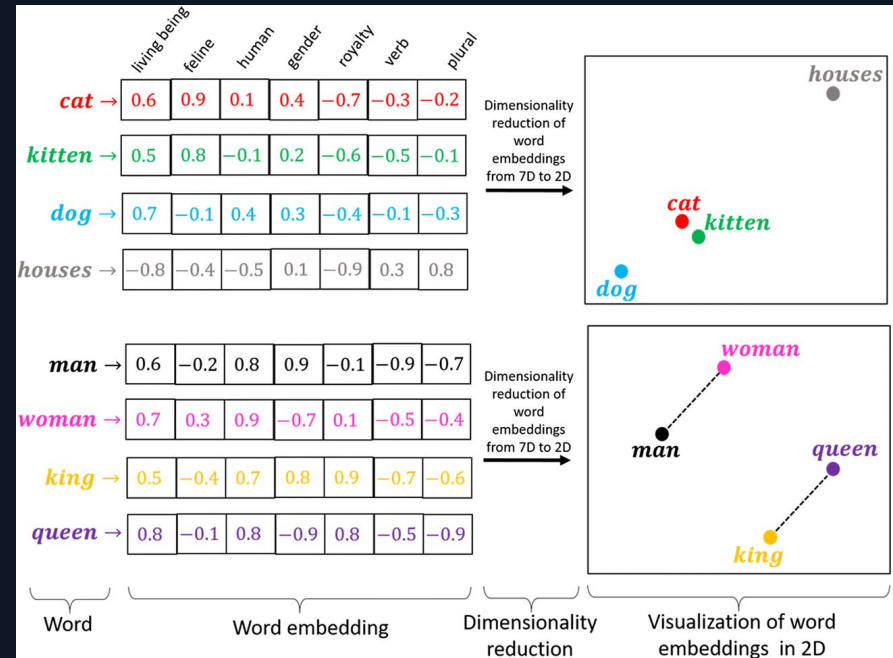
<sup>2</sup> <https://www.sec.gov/Archives/edgar/>

<sup>3</sup> <https://nlp.stanford.edu/projects/glove/>



# GloVe: Global Vector for word representation

- Word embedding
- Co-occurrence matrix
- Matrix Factorization Method (SVD ect.)
- Prediction based approach (Skip-Gram, CBOW)
- GloVe combines



<https://medium.com/@hari4om/word-embedding-d816f643140>

# GloVe: Global Vector for word representation

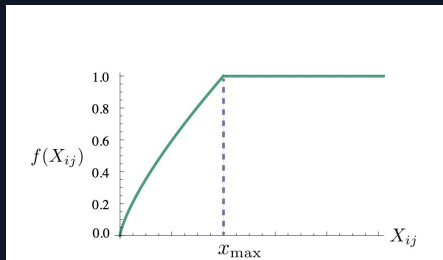
- $X_{ij}$ : word-word co-occurrence matrix
- Global statistics to predict the probability of word  $j$  appearing in the context of word  $i$

Model:

$$w_i^T \tilde{w}_k + b_i + \tilde{b}_k = \log(X_{ik}) .$$

Cost function:

$$J = \sum_{i,j=1}^V f(X_{ij}) (w_i^T \tilde{w}_j + b_i + \tilde{b}_j - \log X_{ij})^2 ,$$



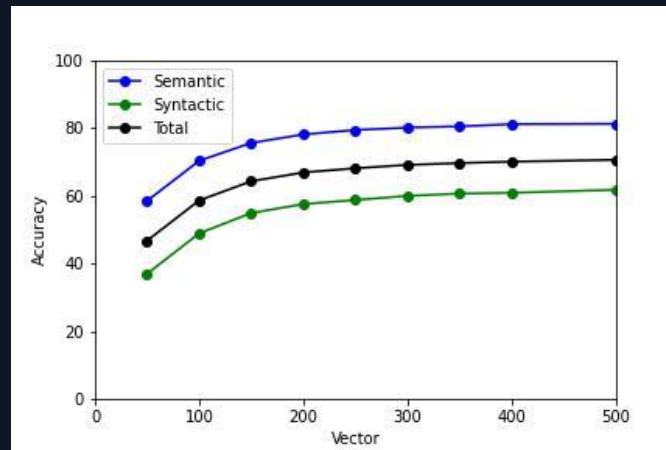
# GloVe on Wikipedia Articles

- ~ 5.8 mil article in english wikipedia
- Step 1: Download Wiki dump file (XML)
- Step 2: Gensim package to convert xml to json; much easier to parse
- Step 3: Batch of 10,000 article in a single json file
- Step 4: Process each batch using Spacy; json to txt file
- Step 5: Concatenate all article in single file

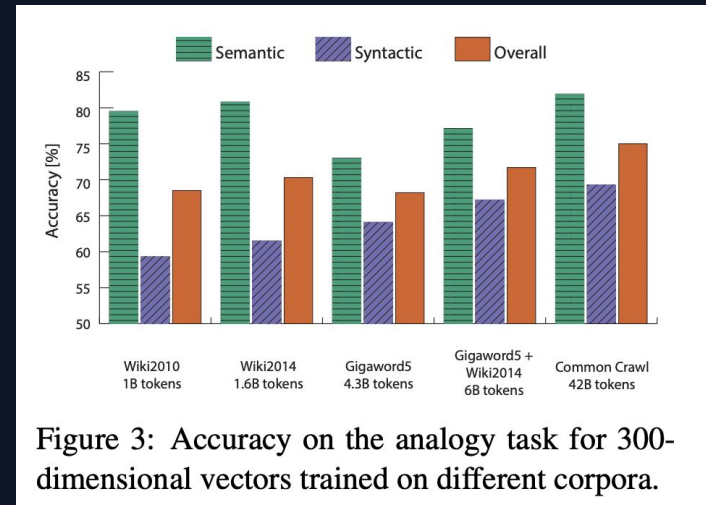
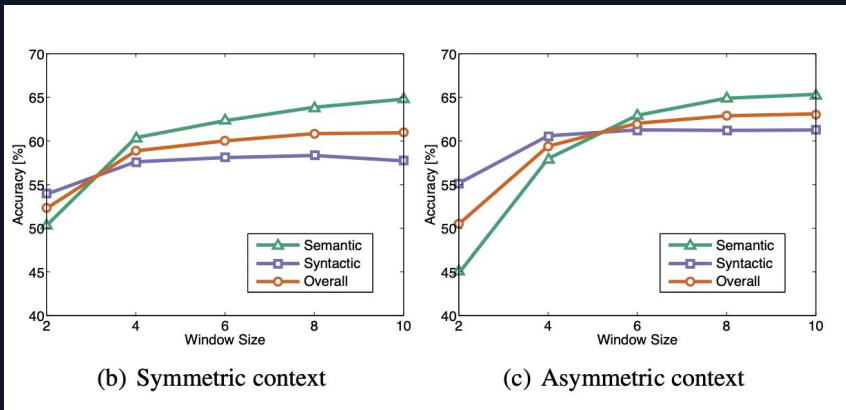


# GloVe on Wikipedia Articles

- Result: Word analogy tasks (Semantic and syntactic)
- “a” is to “b” as “c” is to \_\_\_?
- king - man + woman = queen
- Dancing - dance + play = playing



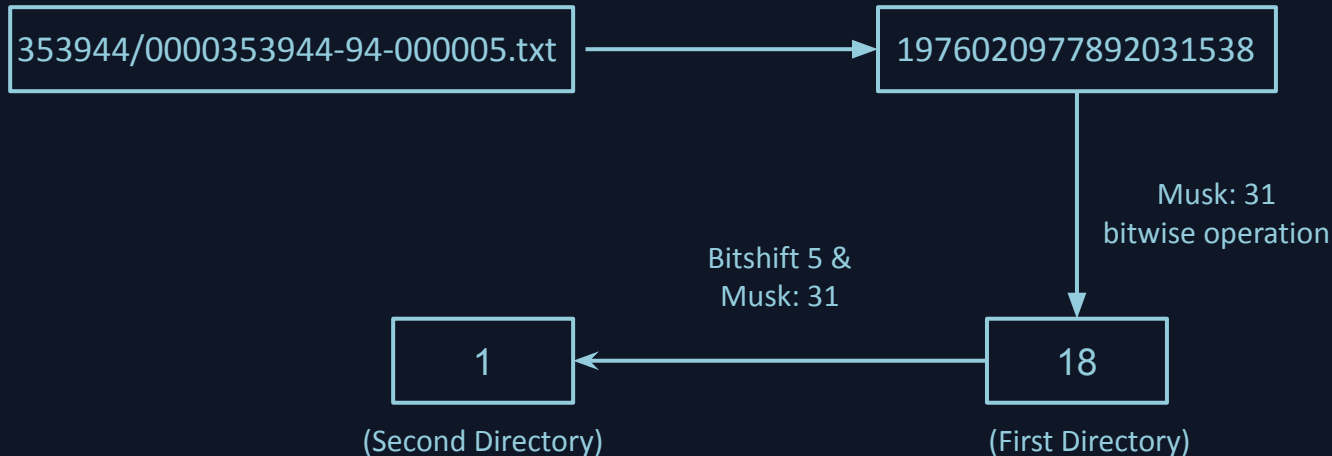
# GloVe on Wikipedia Articles





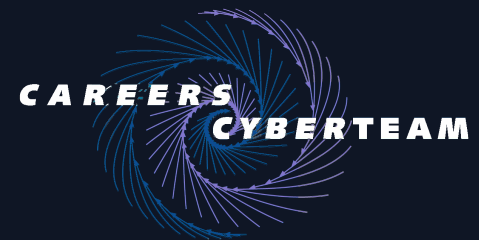
# 10K files

- Total files: ~200k
- Problem: Listing files, existence of file
- Hashed directory
- Two level directory structure



# Future work

- Processing of finance of document
- Identification of finance specific words
- Comparison between domain specific vectors



Thank You!

