

PGR 210 - Natural Language Processing Part

Kristiania University College

By Huamin Ren

Huamin.ren@kristiania.no



Outline of week 45

1. GloVe
2. Similary metric
3. K-means clustering

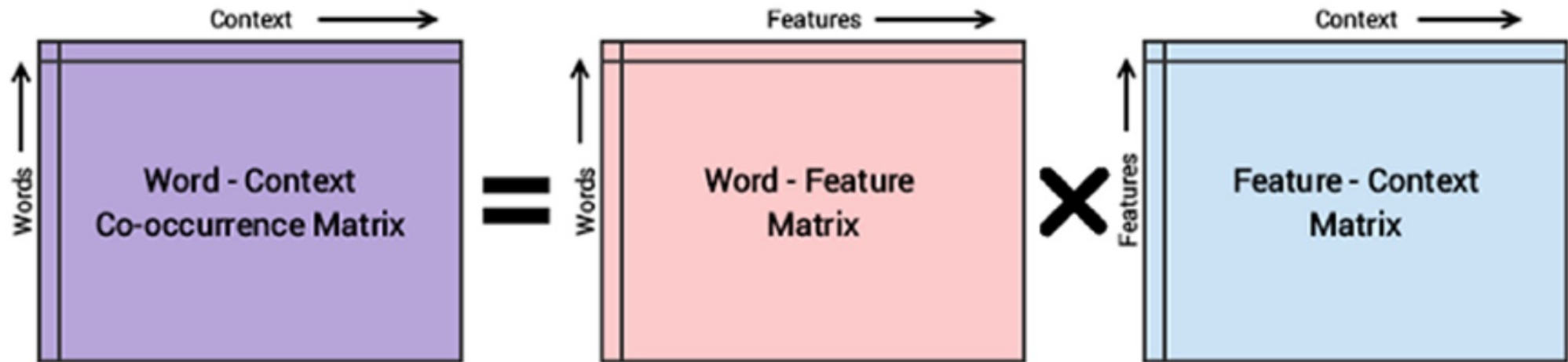
1. GloVe

1.1

- Word2vec was a breakthrough, but it relies on a neural network model that must be trained using backpropagation.
- GloVec is applying direct optimization of the global vectors of word co-occurrences (co-occurrences across the entire corpus)

<https://nlp.stanford.edu/pubs/glove.pdf?fileGuid=WyYwxqq8kWjKdWgd>

1.2 Conceptual model for the GloVe model's implementation



1.3

- Considering the Word-Context (WC) matrix, Word-Feature (WF) matrix, and Feature-Context (FC) matrix, we try to factorize

$$WC = WF FC$$

- Typically initialize WF and FC with some random weights and attempt to multiply them to get WC' (an approximation of WC) and measure how close it is to WC .
- Do this multiple times using *Stochastic Gradient Descent* (SGD) to minimize the error.

- WF (*Word-Feature*): word embeddings for each word, where F can be present to a specific number of dimensions.
- Similarity between Word2Vec and: build a vector space where the position of each word is influenced by its neighboring words based on their context and semantics.
- Difference in the two models: Word2Vec starts with local individual examples of word co-occurrence pairs; GloVe starts with global aggregated co-occurrence statistics across all words in the corpus.

2. Similarity metric

2.1 Similarity scores

- Similarity scores (and distances) tell on how similar or how far apart two documents are based on the similarity (or distance) of the vectors you used to represent them
- Some familiar distances:
 - Euclidean or Cartesian distance, or root mean square error (RMSE): 2-norm or L_2 , MSE,
 - Squared Euclidean distance, sum of squares distance (SSD): L_2
 - Cosine or angular or projected distance: normalized dot product
 - Minkowski distance: p-norm or L_p
 - Fractional distance, fractional norm: p-norm or L_p for $0 < p < 1$
 - City block, Manhattan, or taxicab distance; sum of absolute distance (SAD)

- V1, V2
- $\text{Sim}(V1, V2) = 0.999999999$
- $\text{Sim}(V1, V1) = 1$, 1 m documents, v1, v2, v1m
- $\text{Distance}(V1, V1) = 0$
- $\text{Sim}(V1, V2) = 0.49$
- $\text{Sim}(V1, V3) = 0.111$
- $\text{Sim}(V1, V4) = 0.47$
- $\text{Threshold} = 0.48$

Pairwise distances available in sklearn

- 'cityblock', 'cosine', 'euclidean', 'l1', 'l2', 'manhattan', 'braycurtis', 'canberra', 'chebyshev', 'correlation', 'dice', 'hamming', 'jaccard', 'kulsinski', 'mahalanobis', 'matching', 'minkowski', 'rogerstanimoto', 'russellrao', 'seuclidean', 'sokalmichener', 'sokalsneath', 'sqeuclidean', 'yule'

2.2 Distance

- Distance measures are often computed from similarity measures (scores) and vice versa such that distances are inversely proportional to similarity scores.
- For distances and similarity scores that range between 0 and 1, like probabilities, it's more common to use a formula like this:

similarity = 1. - distance

distance = 1. - similarity

$$hd(u,v) = \sum_{i=1}^n (u_i \neq v_i)$$

$$norm_hd(u,v) = \frac{\sum_{i=1}^n (u_i \neq v_i)}{n}$$

$U=[1\ 1\ 1\ 0\ 0]$ $v=[0\ 0\ 0\ 1\ 1]$

What is hd ? What is $norm_hd$?

2.3 Use similarity on?

- **Lexical similarity:** This involves observing the contents of the text documents with regards to its syntax, structure, and content and measuring their similarity based on these parameters.
- **Semantic similarity:** This involves determining the semantics, meaning, and context of the documents and then determining how close they are to each other.

- **Term similarity:** Similarity between individual tokens or words
- **Document similarity:** Similarity between entire text documents

3. K-means Clustering

- The k-means clustering algorithm is a centroid-based clustering model that tries to cluster data into groups or clusters of equal variance.
- Disadvantage of this algorithm is that the number of clusters (**k**) needs to be specified in advance.
- Advantage: perhaps the most popular clustering algorithm, due to its ease of use as well as it being scalable with large amounts of data.

Mathematical definition

- A dataset \mathbf{X} with \mathbf{N} data points or samples, the task is to group them into K clusters, where \mathbf{K} is a user-specified parameter.
- The k-means clustering algorithm will segregate the \mathbf{N} data points into \mathbf{K} disjoint separate clusters, C_k , and each of these clusters can be described by the means of the cluster samples.

$$X = \{x_1, x_2 \dots x_n\}$$

$$c_1 = \{x_1, x_2\} \quad c_2 = \{x_3\} \dots c_5 \dots$$

$$X_1 = [1 \ 0 \ 0 \ 1 \ 1 \ 1] \quad x_2 = [1 \ 0 \ 0 \ 1 \ 1 \ 0] \quad c_1 = [1 \ 0 \ 0 \ 1 \ 1 \ 0.5]$$

$$c_1 = \{\text{mean}(x_1, x_2)\}$$

- These means become the cluster centroids μ_k such that these centroids are not bound by the condition that they have to be actual data points from the \mathbf{N} samples in \mathbf{X} .

K-means procedure

Steps:

1. Choose initial k centroids μ_k by taking k random samples from the dataset \mathbf{X} .
2. Update clusters by assigning each data point or sample to its nearest centroid point.
3. Recalculate and update clusters based on the new cluster data points for each cluster obtained from Step 2. Mathematically this can be represented as follows:

1. $k=5$

$U_1(x_2) \dots U_5(x_1)$

2. For each data x_i :

Compute distances,

Find the closest cluster

Assign x_i to u_5

3. recalculate:

$u_1 := (x_2)$

$U_2 = (x_3)$

$$C_k = \{x_n : \|x_n - \mu_k\| \leq \text{all } \|x_n - \mu_l\|\}$$

$$\min \sum_{i=1}^K \sum_{x_n \in C_i} \|x_n - \mu_i\|^2$$