

Truncated SVD

Snowball Stemming
Porter stemming.

$$A = U * S * V^T$$

$$S = \begin{bmatrix} \alpha_1 & & & \\ & \alpha_2 & & \\ & & \alpha_3 & \\ & & & \ddots \\ & & & & \alpha_r \end{bmatrix}$$

first

$$\sigma_1 > \sigma_2 > \sigma_3 \dots > \sigma_r > 0$$

k-largest singular values \rightarrow keep them.
Others \rightarrow = zero.

Use only k columns of U, V.

k \rightarrow tradeoff b/w time and data.
(accuracy).

$$87773 \times \frac{11524}{\downarrow}$$

features

\rightarrow First step is to get words Co-occurrence matrix.
 \downarrow
How 2 or more words occur together in a given corpus.

Example

penny wise and pound foolish.
a penny saved is a penny earned

	a	and	earned	foolish	is	penny	pound	Saved	wise
a	0	0	0	0	0	2	0	0	0
penny	0	0	1	0	0	0	0	1	1

Bigram frequency.

Need $N \times N$ to represent bigram frequencies

\rightarrow this table is highly sparse.
(as most are 0s)

~~568454~~

568454 x 10

≈ 500k

⇒

After
EDA,

Data Preprocessing,
Featurization

87773 x $\left[\begin{array}{c} \text{Some} \\ \text{Columns} \end{array} + 11524 \text{ Word vector} \right]$

87773 x 11524

Need to reduce
dimensions.

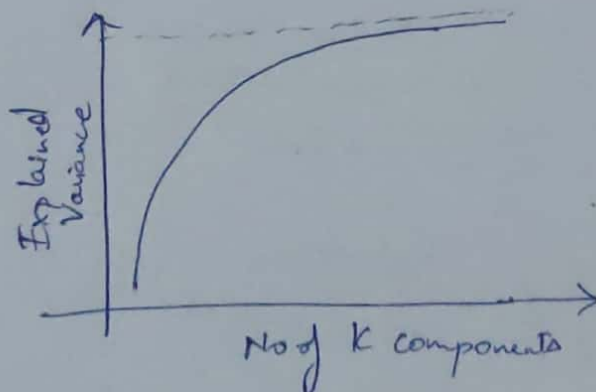
- 1) Bag of words
- 2) ~~TF~~ TFIDF
- 3) Word2Vec.

Co-occurrence matrix - $[11524 \times 11524]$

get the features from `TfidfVectorizer.get_feature_names()`

→ Do `svd(co-occurrence matrix)` for different n -components
(i.e. K -values)

→ get `svd.explained_variance_ratio_.sum()` for different K -values.
↳ ?



Select the optimal K .

→ Decompose 87773×11524 matrix.

$$A_{m \times n} = U_{m \times m} S_{m \times n} V_{n \times n}^T$$

$$B_{m \times n} = U_{m \times m} S_{m \times n} V_{n \times n}^T = \begin{bmatrix} \sigma_1 & \sigma_2 & \dots & \sigma_K & 0 & \dots & 0 \\ 0 & 0 & \dots & 0 & 0 & \dots & 0 \end{bmatrix}_{m \times n}$$

$$\begin{bmatrix} | & | & | & | & 0 & 0 & 0 & 0 \\ | & | & | & | & 0 & 0 & 0 & 0 \\ | & | & | & | & 0 & 0 & 0 & 0 \\ | & | & | & | & 0 & 0 & 0 & 0 \end{bmatrix}_{m \times n}$$