

Introduction to Natural Language Processing

Machine Learning Methods – Lecture 24



June 2021



A logo consisting of a red brushstroke border containing the text "fixel your pixels" in a cursive font. Below this, the URL "pixelalgorithms.gitlab.io" is written in a sans-serif font. The entire logo is set against a black background.

Natural Language Processing – Introduction

- Consider the following IMDb reviews:

“This movie is terrible but it has some good effects.”

positive

“I wouldn’t rent this one even on dollar rental night.”

negative

“Ming The Merciless does a little Bardwork and a movie most foul!”

“You’d better choose Paul Verhoeven’s even if you have watched it.”

“Adrian Pasdar is excellent in this film. He makes a fascinating woman.”

“Long, boring, blasphemous. Never have I been so glad to see ending credits roll.”

“I don’t know why I like this movie so well, but I never get tired of watching it.”

“no comment - stupid movie, acting average or worse... screenplay - no sense at all... SKIP IT!”

“A rating of ”1” does not begin to express how dull, depressing and relentlessly bad this movie is.”

“This is the definitive movie version of Hamlet.

Branagh cuts nothing, but there are no wasted moments.”

- A classic NLP task is to predict the sentiment of a given review.

Bag of Words – Introduction

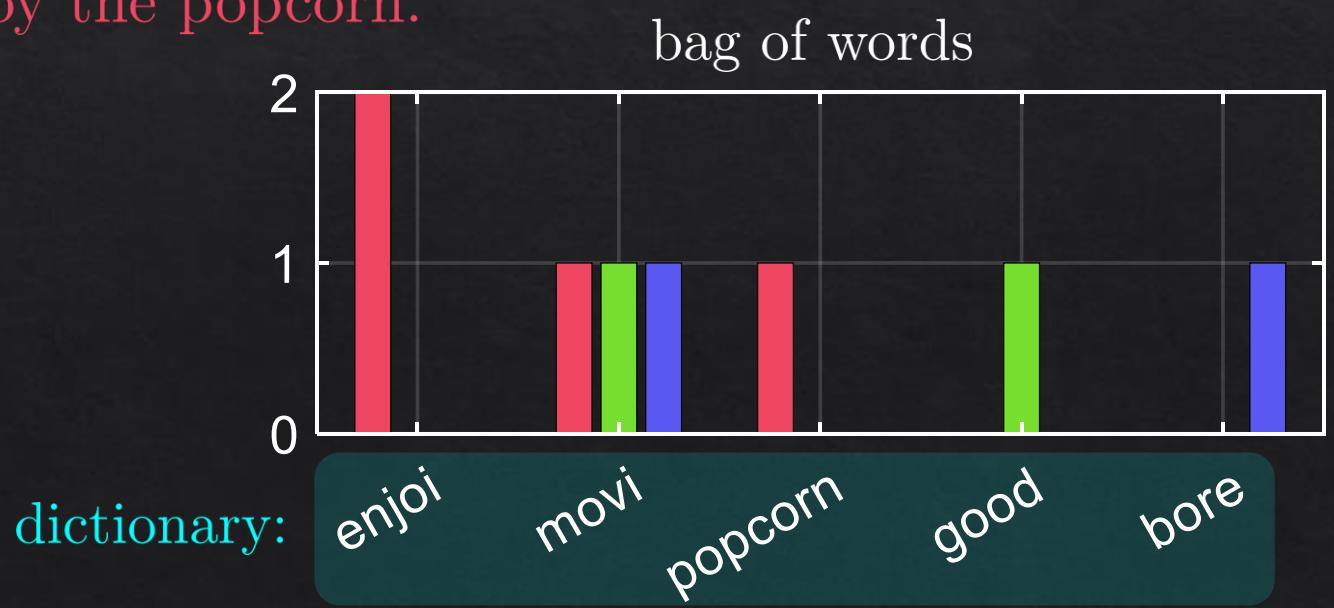
- Usually, the first task is to make the data numeric.
- In Bag of Words we:
 1. Apply some pre-processing.
 2. Build a **vocabulary** (a.k.a. **dictionary**).
 3. Replace each text with a corresponding histogram.

Example

“I enjoyed the movie, but I did not enjoy the popcorn.”

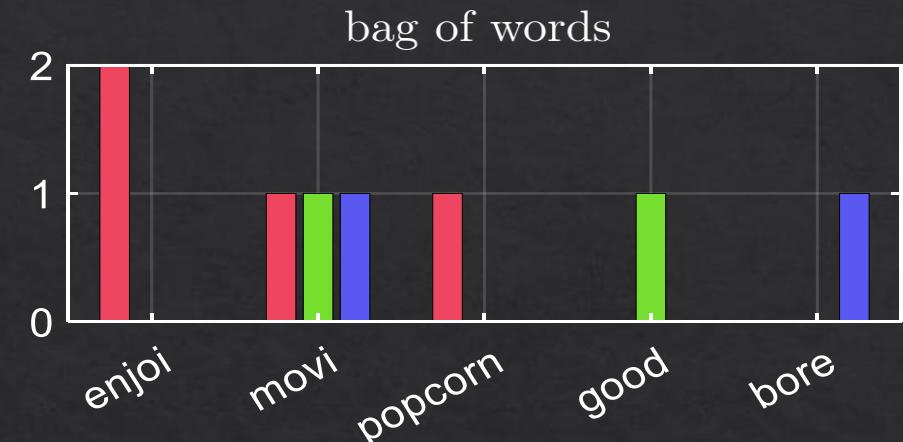
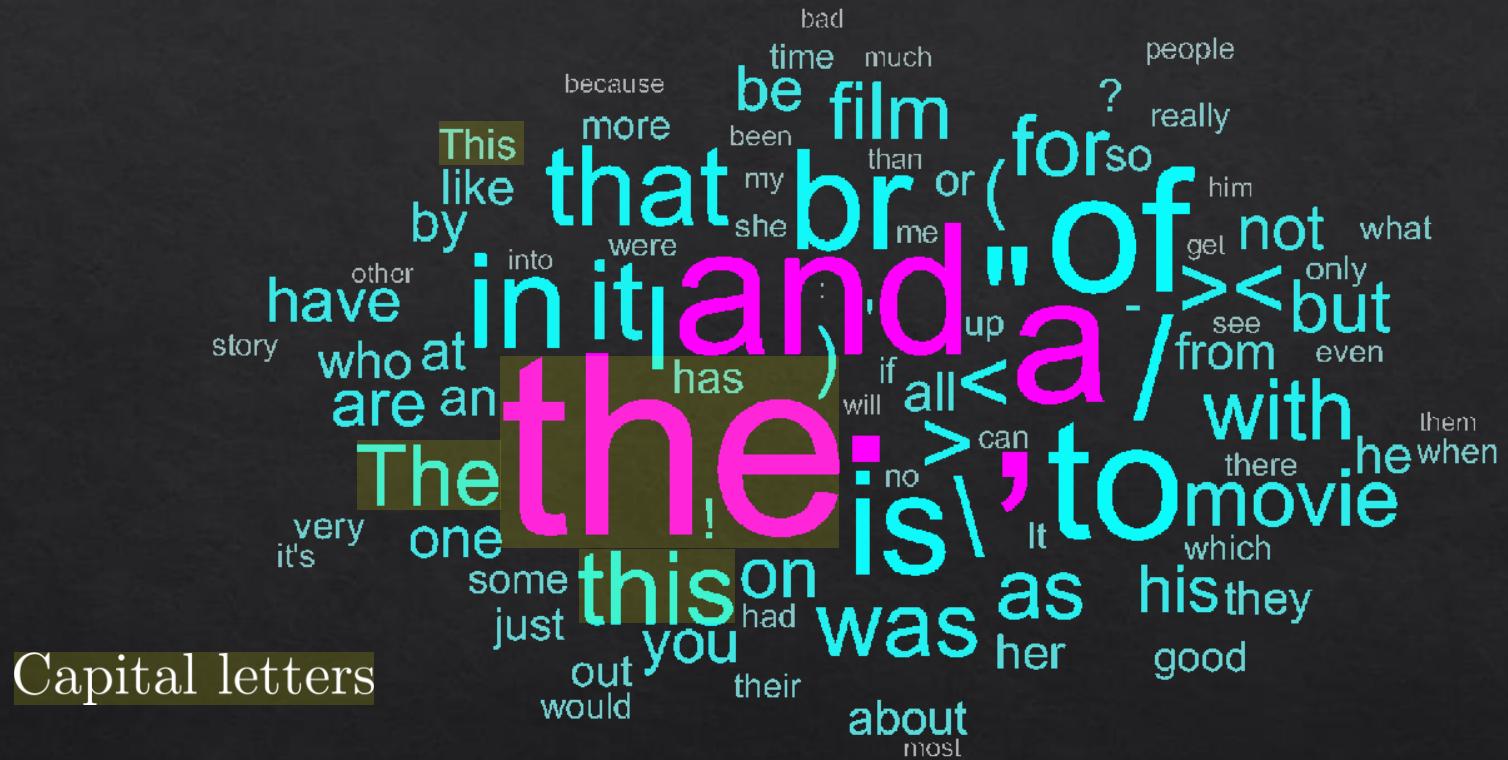
“This is a good movie.”

“The movie was boring.”



Bag of Words – Pre-processing

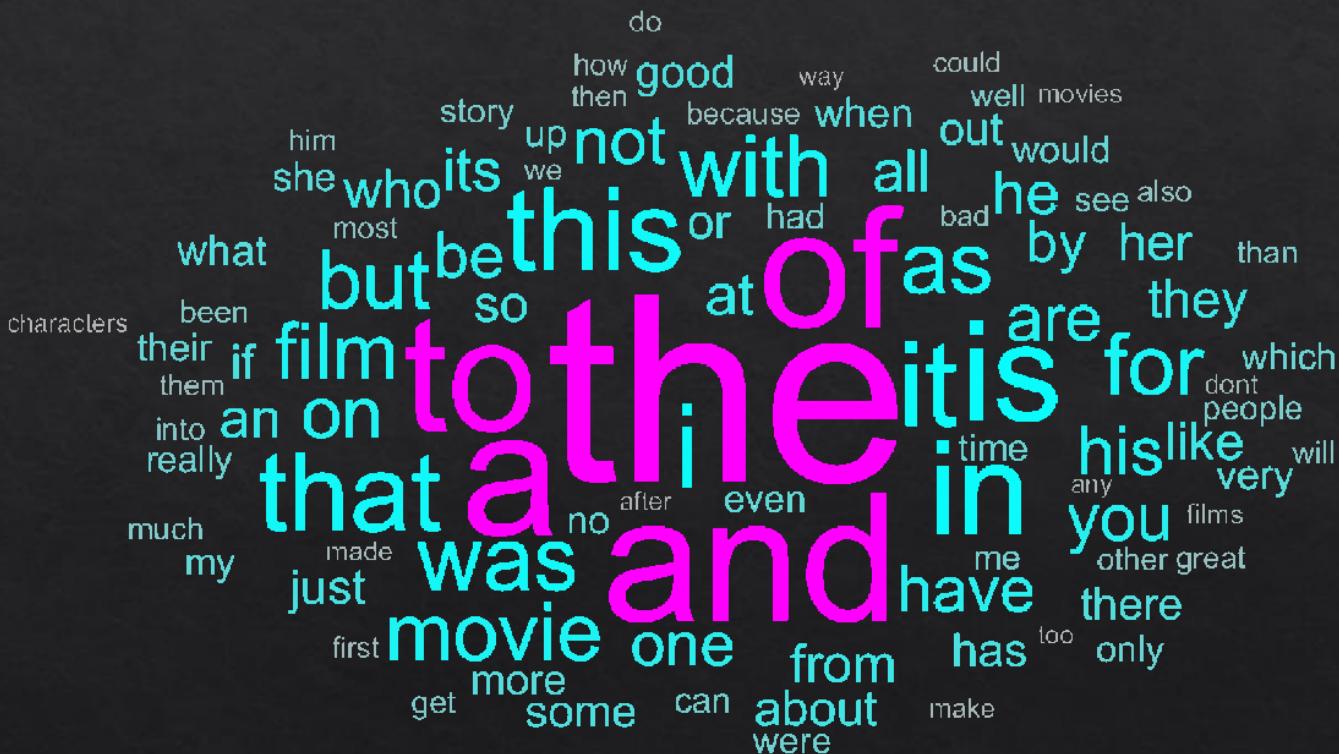
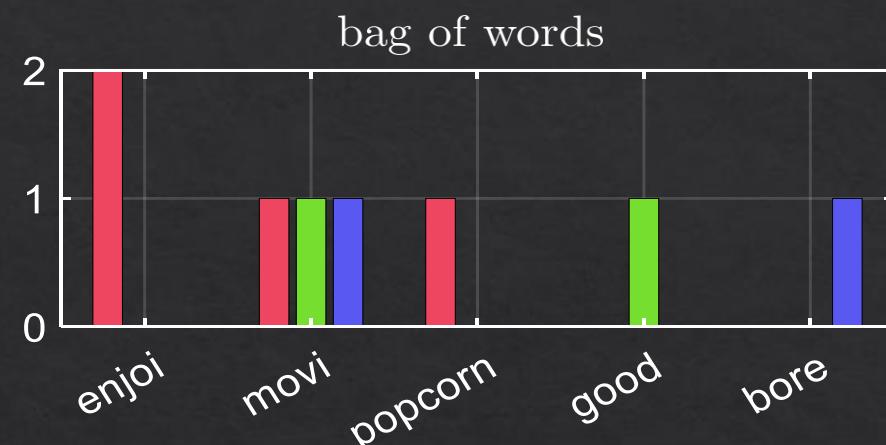
- The most common words among $N = 2,000$ IMDb reviews:



- First step is to remove spaces, punctuation, and change all characters to lower case.

Bag of Words – Pre-processing

- The most common words among $N = 2,000$ IMDb reviews:
 - First step is to remove spaces, punctuation, and change all characters to lower case.



- Second, we remove stop words: “the”, “of”, “to”, “and”, “a”, etc.

Bag of Words – Pre-processing

- The most common words among $N = 2,000$ IMDb reviews:

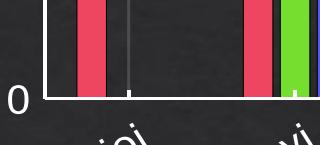
A bar chart with three bars. The y-axis ranges from 0 to 1. The x-axis labels are 'enjoy', 'movie', and 'good'. The bars are colored red, green, and blue respectively. The height of the red bar is approximately 0.8, the green bar is approximately 0.9, and the blue bar is approximately 1.0.
 - First step is to remove spaces, punctuation, and change all characters to lower case.
 - Second, we remove stop words: “the”, “of”, “to”, “and”, “a”, etc.



film movie like just bad off
watch time people even story
make show movies even first life
get many really up good films
plot quite give why makes pretty going saw know work though two way great
thought funny times never thought
look big new old real
lot find man still fact come
another plot
almost take years action
years action
character little part well got
own love back horror world
cast theres scenes
right enough down seen few
right

- Last step is stemming the words: movie(s) → movi, watch(ing) → watch.

Bag of Words – Dictionary

- The most common words among $N = 2,000$ IMDb reviews:

A bar chart with four bars representing word frequencies. The x-axis labels are 'enjoy', 'movi', 'watch', and 'popcorn'. The y-axis ranges from 0 to 1. The bars have heights approximately at 0.95, 0.95, 0.95, and 0.95 respectively. The bars are colored red, green, blue, and red.

Word	Frequency
enjoy	~0.95
movi	~0.95
watch	~0.95
popcorn	~0.95
 - First step is to remove spaces, punctuation, and change all characters to lower case.
 - Second, we remove stop words: “the”, “of”, “to”, “and”, “a”, etc.
 - Last step is stemming the words: movie(s) → movi, watch(ing) → watch.



after these steps,
we remain with 19,910 unique words

let $\mathcal{V} = \{w_j\}_{j=1}^{19,910}$ be the set (vocabulary) of all unique words.

Bag of Words – Data

- The numeric data set is $\mathbf{X} \in \mathbb{N}_0^{2,000 \times 19,910}$.
 - The i th row of \mathbf{X} is the histogram of the i th document.
 - Since most of the entries in \mathbf{X} are zero, it is efficient to use sparse representation.
 - For example:

$$\mathbf{X} = \begin{bmatrix} 2 & 1 & 1 & 0 & 0 \\ 0 & 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 & 1 \end{bmatrix}$$

$$X_{ij} = \sum_{w \in d_i} \mathbb{I}\{w_j = w\}$$

d_i is the i th review (document)

$w_j \in \mathcal{V}$ is the j th word (term)



$N = 2,000$ IMDb reviews:



after these steps,
we remain with 19,910 unique words

let $\mathcal{V} = \{w_j\}_{j=1}^{19,910}$ be the set (vocabulary) of all unique words.

Term Frequency Inverse Document Frequency (TF-IDF)

- TF-IDF is a weighted version of bag of words.

$$X_{ij} = \sum_{w \in d_i} \mathbb{I}\{w_j = w\}$$

- Term Frequency: counts the occurrences of the word $w \in \mathcal{V}$ in a document d :

$$\text{tf}(w, d) = \sum_{w_j \in d} \mathbb{I}\{w_j = w\}$$

tf is the bag of words: $X_{ij} = \text{tf}(w_j, d_i)$

- Inverse Document Frequency: measures whether the word $w \in \mathcal{V}$ is common or rare across all documents $\mathcal{D} = \{d_1, d_2, \dots, d_N\}$:
 \mathcal{D} is known as corpus

$$\text{idf}(w, \mathcal{D}) = \log \left(\frac{N}{\sum_{d_i \in \mathcal{D}} \mathbb{I}\{w \in d_i\}} \right) \quad N = |\mathcal{D}|$$

There are other possible definitions for both TF and IDF.

- TD-IDF is given by:

$$\text{tf-idf}(w, d, \mathcal{D}) = \text{tf}(w, d) \cdot \text{idf}(w, \mathcal{D})$$

$$\mathbf{X}_{\text{tf-idf}}[i, j] = \text{tf}(w_j, d_i) \cdot \text{idf}(w_j, \mathcal{D})$$

Bag of Words – IMDb Example

$$\text{tf}(w, d) = \sum_{w_j \in d} \mathbb{I}\{w_j = w\}$$

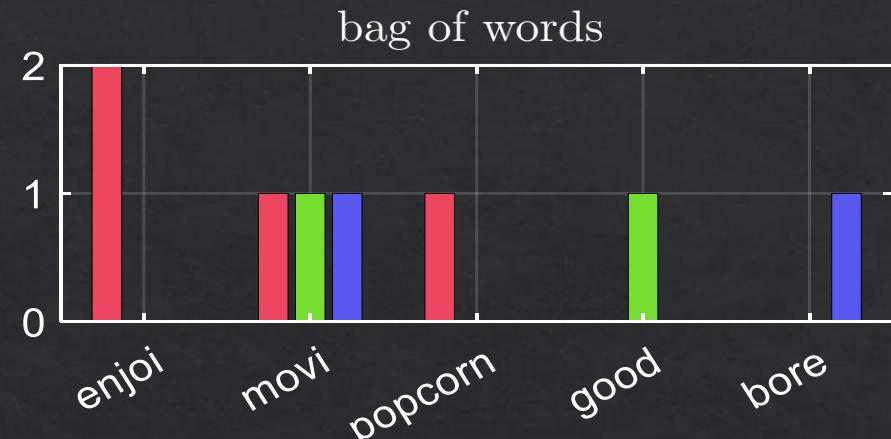
$$\text{idf}(w, \mathcal{D}) = \log \left(\frac{N}{\sum_{d_i \in \mathcal{D}} \mathbb{I}\{w \in d_i\}} \right)$$

$$\text{tf-idf}(w, d, \mathcal{D}) = \text{tf}(w, d) \cdot \text{idf}(w, \mathcal{D})$$

$$\mathbf{X}_{\text{tf-idf}}[i, j] = \text{tf}(w_j, d_i) \cdot \text{idf}(w_j, \mathcal{D})$$



Bag of Words – IMDb



let $\mathcal{V} = \{w_j\}_{j=1}^{19,910}$ be the set (vocabulary) of all unique words.

$$X_{ij} = \sum_{w \in d_i} \mathbb{I}\{w_j = w\}$$

Word2Vec – Introduction

- Word2Vec is a method to embed words into a vector (feature) space \mathbb{R}^d .
- Let $\mathbf{z}_i = \phi(\mathbf{w}_i)$ be the mapping of the word $\mathbf{w}_i \in \mathcal{V}$ to the vector $\mathbf{z}_i \in \mathbb{R}^d$, where \mathcal{V} is the vocabulary (set of all words).
- The feature space has a useful structure, for example:

$$\phi(\text{actor}) - \phi(\text{man}) + \phi(\text{woman}) \approx \phi(\text{actress})$$

- The word embedding is obtained by optimizing a linear transformation followed by a softmax activation:

$$\hat{\mathbf{y}}_i = \sigma_{\text{softmax}} \left(\begin{array}{c|c|c} \mathbf{V}_c \in \mathbb{R}^{|\mathcal{V}| \times d} & \mathbf{V}_w \in \mathbb{R}^{d \times |\mathcal{V}|} & \mathbf{x}_i \\ \hline & & \end{array} \right)$$

$$\text{movie} \xrightarrow{\phi} \begin{bmatrix} 0.7 \\ -2.1 \\ \vdots \\ 1.1 \end{bmatrix} \quad w_i \quad z_i$$

Word2Vec – Context Window

- Consider the following text:

“... an exciting summer blockbuster that was visually fantastic but also curiously ...”

context context word context context

- From the text we build pairs (word, context).

For example, given the word $w_i = \text{blockbuster}$ and a context radius of 2 we have:

$$(x_1, y_1) = (\text{blockbuster}, \text{exciting})$$

$$(x_{17}, y_{17}) = (\text{fantastic}, \text{was})$$

$$(x_2, y_2) = (\text{blockbuster}, \text{summer})$$

$$(x_{18}, y_{18}) = (\text{fantastic}, \text{visually})$$

$$(x_3, y_3) = (\text{blockbuster}, \text{that})$$

...

$$(x_{19}, y_{19}) = (\text{fantastic}, \text{but})$$

$$(x_4, y_4) = (\text{blockbuster}, \text{was})$$

$$(x_{20}, y_{20}) = (\text{fantastic}, \text{also})$$

...

- We iterate the entire text and construct all (word, context) pairs $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^N$.
- Word2Vec tries to predict the correct context from the given word.

Word2Vec – One-hot Encoding

- We encode each word using a one-hot encoding.

For example, if $x_i = \text{blockbuster}$ is the third word in the vocabulary \mathcal{V} , then:

$$\mathbf{x}_i = \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \in \mathbb{R}^{|\mathcal{V}|} \implies \mathcal{D} = \left\{ \underbrace{\left(\mathbf{x}_1, \mathbf{y}_1 \right)}_{\text{, }}, \underbrace{\left(\mathbf{x}_2, \mathbf{y}_2 \right)}_{\text{, }}, \dots, \underbrace{\left(\mathbf{x}_N, \mathbf{y}_N \right)}_{\text{, }} \right\}$$

$|\mathcal{V}|$ is the size of the vocabulary.

$$\begin{array}{c|c} \begin{array}{c} \text{---} \\ | \\ \text{---} \\ | \\ \text{---} \\ | \\ \text{---} \end{array} & \begin{array}{c} 0 \\ 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{array} \\ \hline \end{array} = \begin{array}{c} (x_1, y_1) = (\text{blockbuster}, \text{exciting}) \\ (x_2, y_2) = (\text{blockbuster}, \text{summer}) \\ (x_3, y_3) = (\text{blockbuster}, \text{that}) \\ (x_4, y_4) = (\text{blockbuster}, \text{was}) \\ \hline \mathcal{D} = \{(x_i, y_i)\}_{i=1}^N \end{array}$$

Architecture and Loss Function

- The prediction function of \mathbf{y}_i given \mathbf{x}_i :

$$\hat{\mathbf{y}}_i = \sigma_{\text{softmax}} (\mathbf{V}_c \mathbf{V}_w \mathbf{x})$$

$$\hat{\mathbf{y}}_i = \sigma_{\text{softmax}} \left(\begin{array}{c|c|c} \mathbf{V}_c \in \mathbb{R}^{|\mathcal{V}| \times d} & \mathbf{V}_w \in \mathbb{R}^{d \times |\mathcal{V}|} & \mathbf{x}_i \\ \hline & & \end{array} \right)$$

$$\mathcal{D} = \left\{ \left(\underbrace{(\mathbf{x}_1, \mathbf{y}_1)}_{\text{,}}, \underbrace{(\mathbf{x}_2, \mathbf{y}_2)}_{\text{,}}, \dots, \underbrace{(\mathbf{x}_N, \mathbf{y}_N)}_{\text{,}} \right) \right\}$$

$$\mathbf{p} = \sigma_{\text{softmax}} (\mathbf{v}) = \frac{\exp (\mathbf{v})}{\mathbf{1}^T \exp (\mathbf{v})}$$

\mathbf{p} is a probabilistic vector

- As in logistic regression,
the MLE of \mathbf{V}_w and \mathbf{V}_c are the minimizeres of the cross-entropy loss:

$$L = -\frac{1}{N} \sum_{i=1}^N \mathbf{y}_i^T \log (\hat{\mathbf{y}}_i) \implies \arg \min_{\mathbf{V}_c, \mathbf{V}_w} -\frac{1}{N} \sum_{i=1}^N \mathbf{y}_i^T \log (\hat{\mathbf{y}}_i)$$

the objective is not convex

- The final words embedding are given by the columns of \mathbf{V}_w^* .

Word2Vec – Interpretation

- Consider \mathbf{x}_i that corresponds to the k th word in \mathcal{V} :

$$\begin{bmatrix} | & | \\ z_1 & z_2 \\ | & | \end{bmatrix} \mathbf{V}_w \begin{bmatrix} | \\ z_{|\mathcal{V}|} \\ | \end{bmatrix} \mathbf{x}_i = \begin{bmatrix} | \\ z_2 \\ | \end{bmatrix}$$

$$\arg \min_{\mathbf{V}_c, \mathbf{V}_w} -\frac{1}{N} \sum_{i=1}^N \mathbf{y}_i^T \log (\hat{\mathbf{y}}_i)$$
$$\begin{bmatrix} | \\ \hat{\mathbf{y}}_i \\ | \end{bmatrix} = \sigma_{\text{softmax}} \left(\begin{bmatrix} | \\ \mathbf{V}_c \\ | \\ \mathbf{V}_w \\ | \\ \mathbf{x}_i \\ | \end{bmatrix} \right)$$

$\Rightarrow \mathbf{z}_k \in \mathbb{R}^d$ is a vector representation of $w_k \in \mathcal{V}$.

- In addition, the vector $\mathbf{c}_k = \mathbf{V}_c \mathbf{z}_k$ provide a context score between w_k and all other words (as context words).

$$\begin{bmatrix} | \\ \mathbf{V}_c \\ | \\ \mathbf{V}_c \\ | \\ \mathbf{V}_c \\ | \end{bmatrix} \mathbf{z}_k = \begin{bmatrix} | \\ \hat{\mathbf{c}}_k \\ | \end{bmatrix}$$

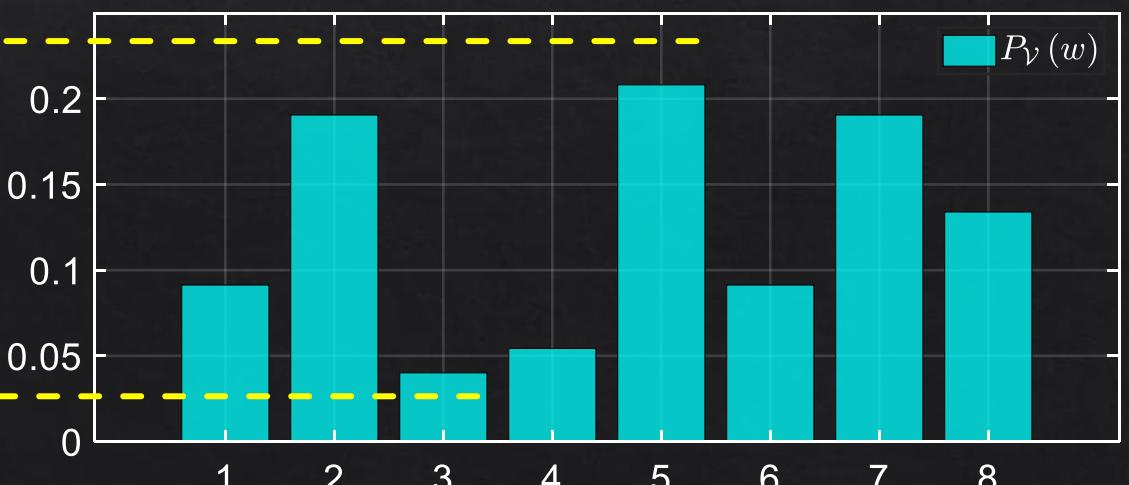
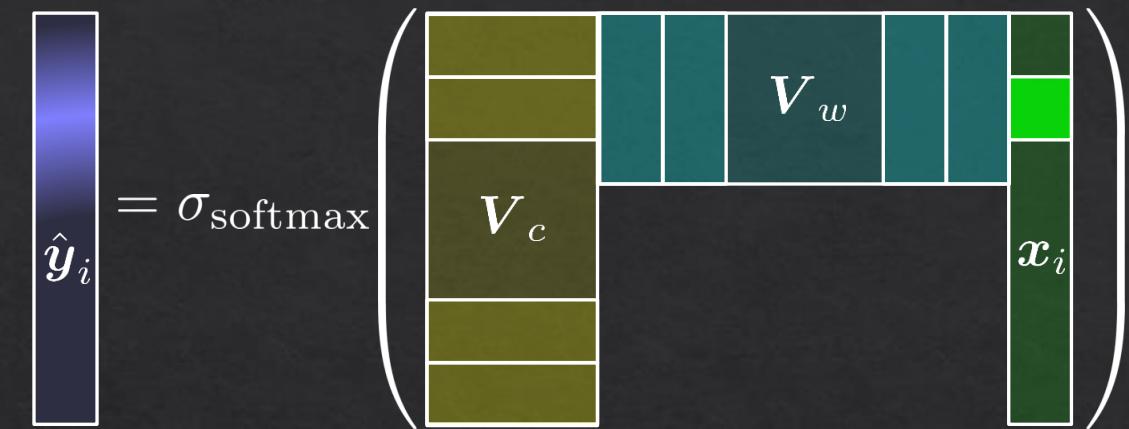
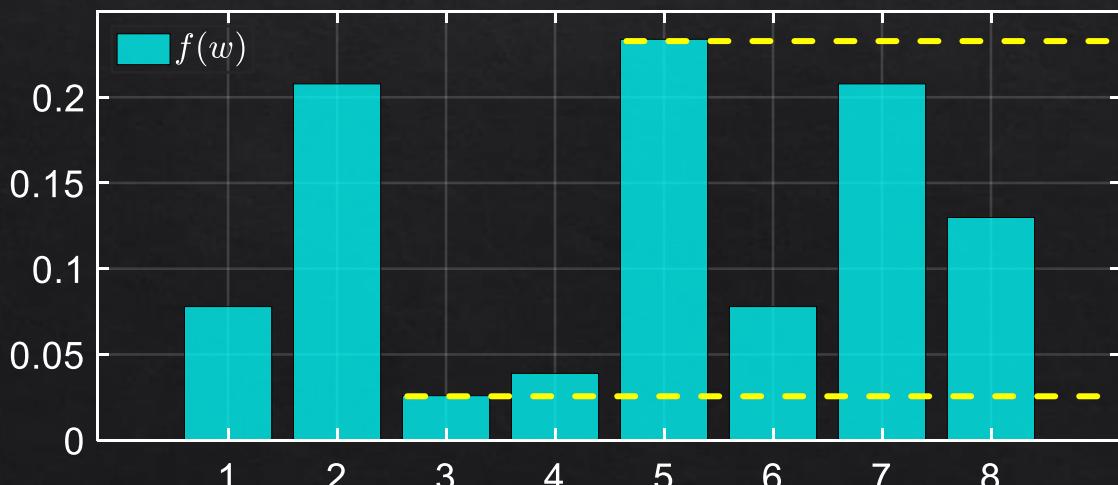
Word2Vec – Negative Sampling

- When the size of the vocabulary is large $|\mathcal{V}|$, the optimization is not trivial.
- Let $P_{\mathcal{V}}$ be a probability of the word w :

$$P_{\mathcal{V}}(w) = \frac{f^{\frac{3}{4}}(w)}{\sum_{w_i \in \mathcal{V}} f^{\frac{3}{4}}(w_i)}$$

where $f(w) \in [0, 1]$ is the frequency of w appears in the corpus.

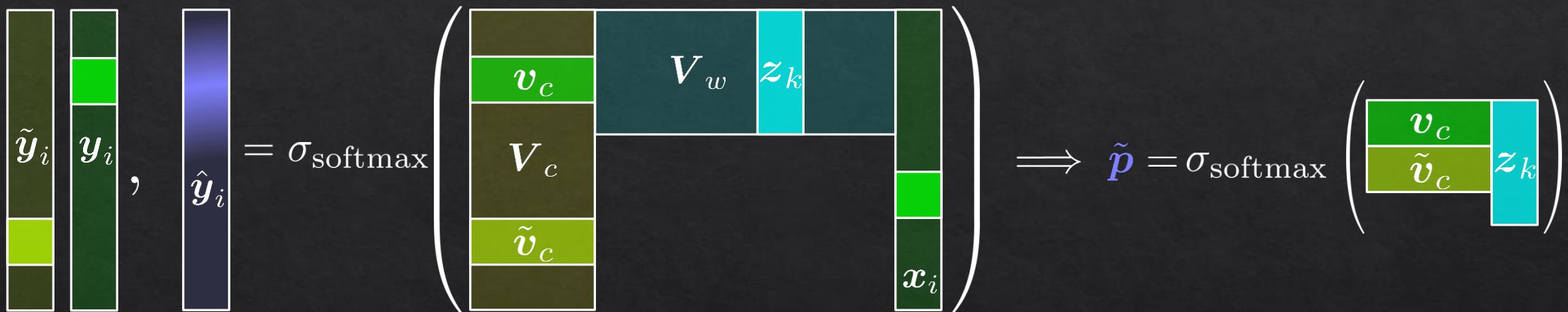
- The power $(\cdot)^{\frac{3}{4}}$ makes less frequent words to be sampled more often:



Word2Vec – Negative Sampling

$$P_{\mathcal{V}}(\mathbf{w}) = \frac{f^{\frac{3}{4}}(\mathbf{w})}{\sum_{w_i \in \mathcal{V}} f^{\frac{3}{4}}(w_i)}$$

- Given a true pair $(\mathbf{x}_i, \mathbf{y}_i)$, we create a fake (negative) pair $(\mathbf{x}_i, \tilde{\mathbf{y}}_i)$, where $\tilde{\mathbf{y}}_i$ is one-hot vector of a random word $\tilde{w} \sim P_{\mathcal{V}}$.
- Instead of consider the entire vocabulary \mathcal{V} , we reduce the problem to the words correspond to \mathbf{x}_i , \mathbf{y}_i , and $\tilde{\mathbf{y}}_i$:



- We optimize only w.r.t. $\mathbf{z}_k, \mathbf{v}_c$, and $\tilde{\mathbf{v}}_c$:

$$\arg \min_{\mathbf{z}_k, \mathbf{v}_c, \tilde{\mathbf{v}}_c} - \left\langle \begin{bmatrix} 1 \\ 0 \end{bmatrix}, \log(\tilde{\mathbf{p}}) \right\rangle \xleftarrow{\text{cross-entropy loss}}$$

$$\arg \min_{\mathbf{V}_c, \mathbf{V}_w} - \frac{1}{N} \sum_{i=1}^N \mathbf{y}_i^T \log(\hat{\mathbf{y}}_i)$$

- At each step (of the optimization), we draw a new negative sample.
- One can take more than one negative pair.

Word2Vec – Example



Word2Vec – IMDb

$$\hat{\mathbf{y}}_i = \sigma_{\text{softmax}} \left(\begin{array}{c} \mathbf{V}_c \\ \vdots \\ \mathbf{V}_c \\ \mathbf{V}_w \\ \vdots \\ \mathbf{V}_w \\ \mathbf{x}_i \end{array} \right)$$

$$\begin{bmatrix} z_1 & z_2 & \mathbf{V}_w & z_{|\mathcal{V}|} \end{bmatrix} \mathbf{x}_i = \begin{bmatrix} z_2 \end{bmatrix}$$

$$\begin{bmatrix} \mathbf{V}_c \\ \vdots \\ \mathbf{V}_c \\ z_k \end{bmatrix} = \begin{bmatrix} \hat{c}_k \end{bmatrix}$$

Questions



TF-IDF:

$$\text{tf}(\textcolor{teal}{w}, \textcolor{blue}{d}) = \sum_{w_j \in \textcolor{blue}{d}} \mathbb{I}\{w_j = w\}$$

$$\text{idf}(\textcolor{teal}{w}, \mathcal{D}) = \log \left(\frac{N}{\sum_{d_i \in \mathcal{D}} \mathbb{I}\{w \in d_i\}} \right)$$

$$\text{tf-idf}(w, d, \mathcal{D}) = \text{tf}(w, d) \cdot \text{idf}(w, \mathcal{D})$$

Word2Vec:

$$\hat{\mathbf{y}}_i = \sigma_{\text{softmax}} \left(\mathbf{V}_c \mathbf{V}_w \mathbf{x}_i \right)$$

