

Deep Learning

A course about theory & practice



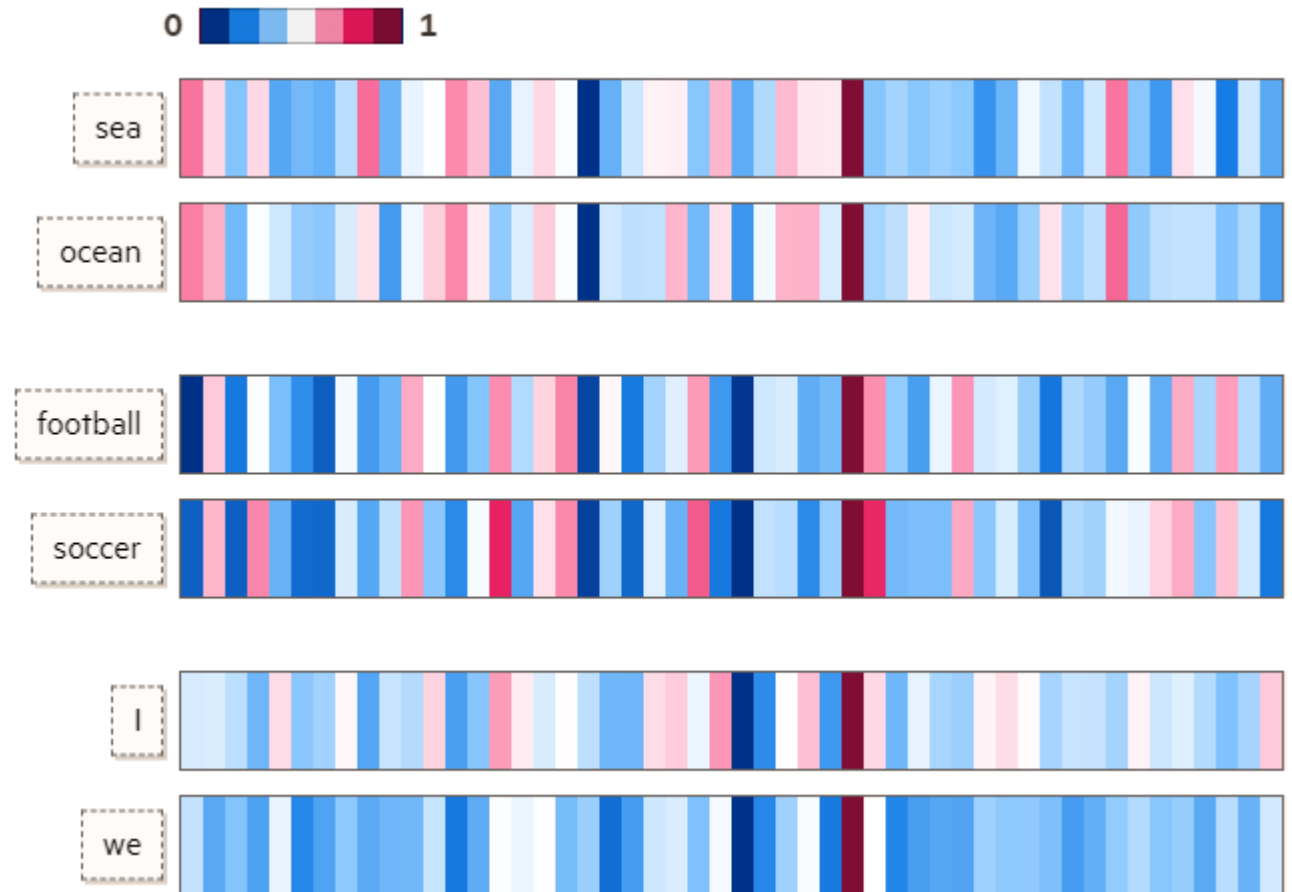
Word Embedding

Marco Piastra

Basic Idea

Embedding, in short

Words (=token) from natural language are each translated into a high-dimensional numerical vector

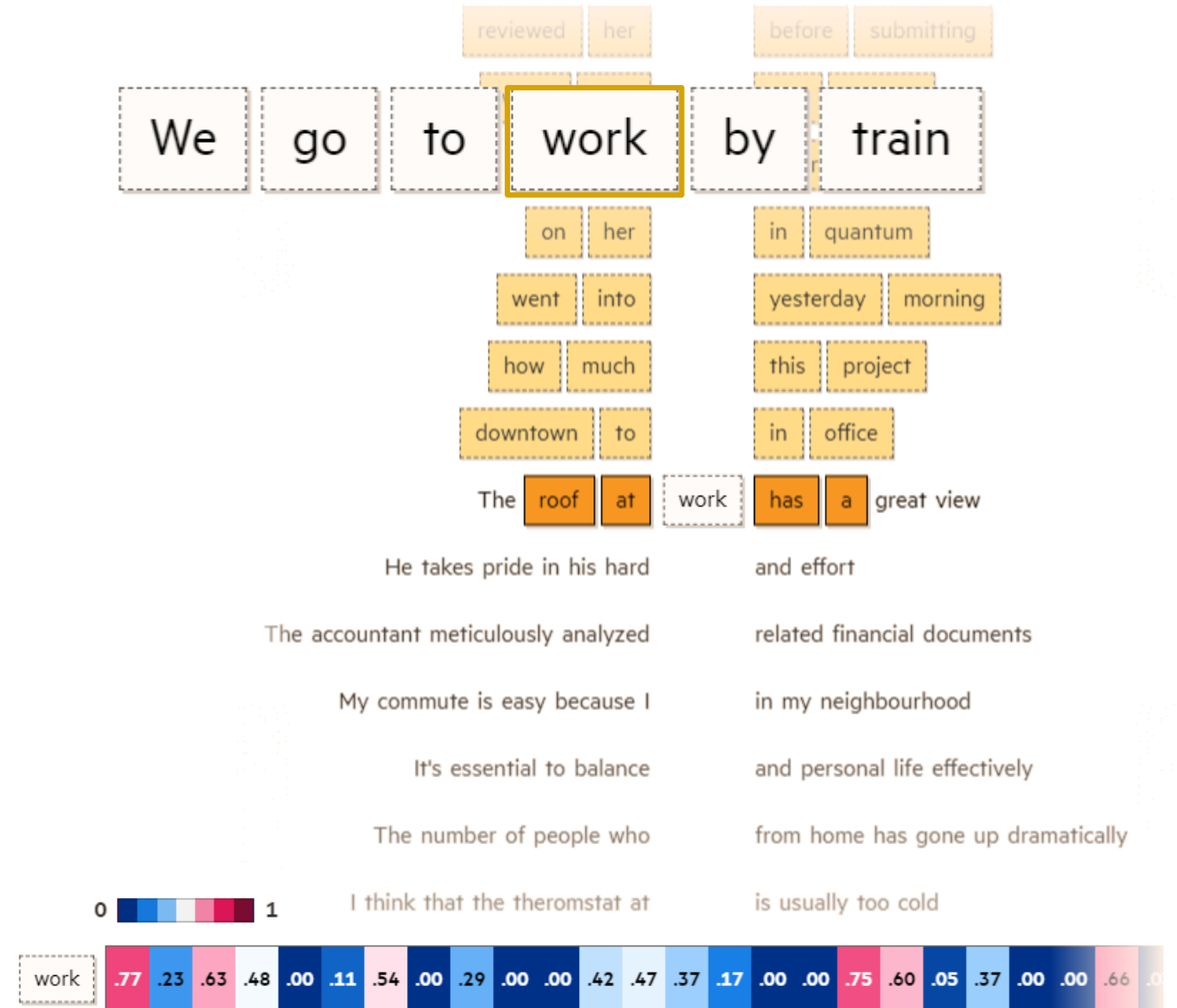


Images from <https://ig.ft.com/generative-ai/>

Embedding, in short

Words (=token) from natural language are each translated into a high-dimensional numerical vector

Such vector is computed by estimating the *probability of co-occurrence* in a context of other words in a (very) large text corpus



Images from <https://ig.ft.com/generative-ai/>

Embedding, in short

Words (=token) from natural language are each translated into a high-dimensional numerical vector

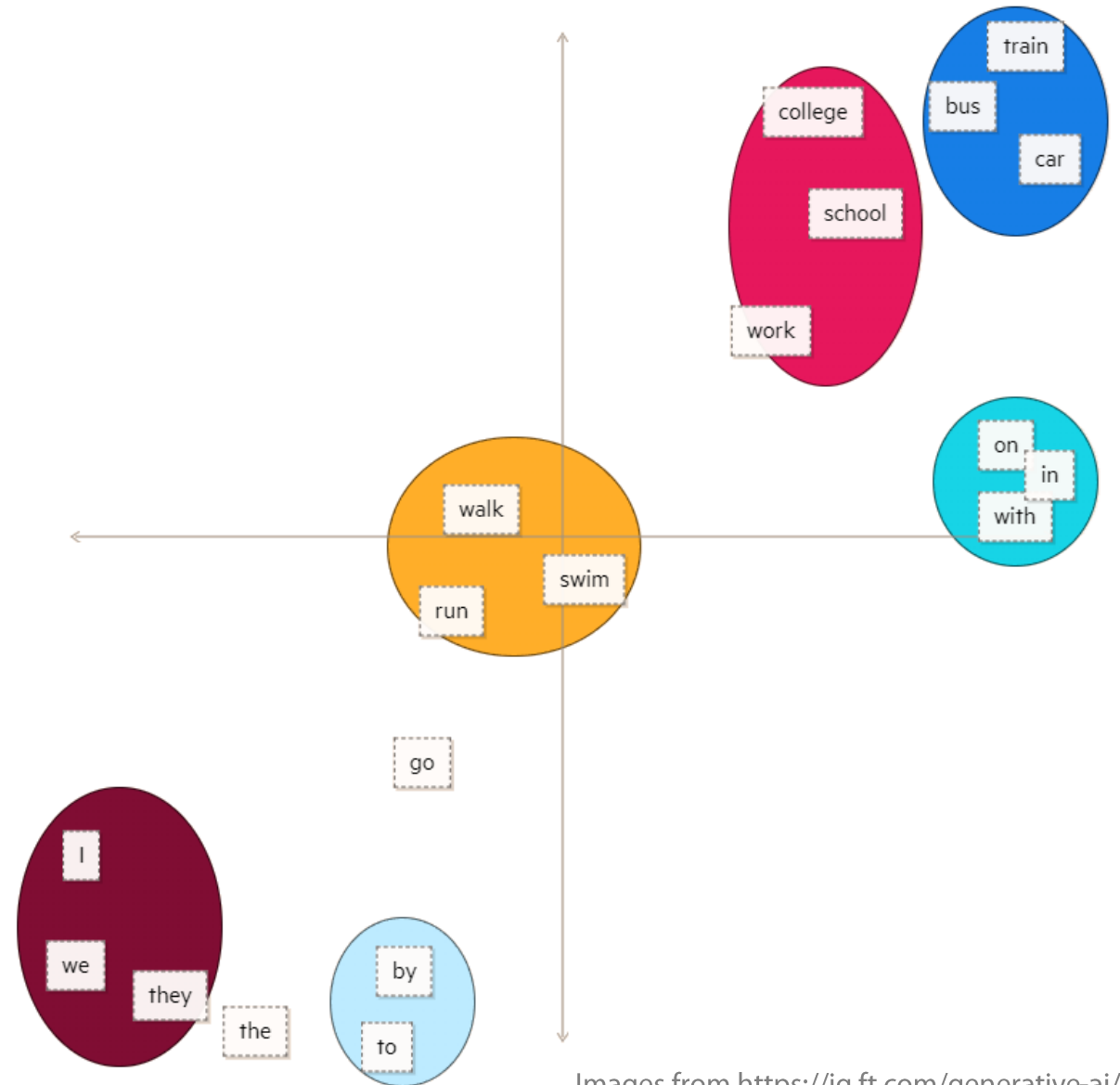
Such vector is computed by estimating the *probability of co-occurrence* in a context of other words in a (very) large text corpus

In this way, the *numerical similarity* among vectors is representative of words' affinity in terms of role or meaning (or both)

Distributional Semantics:

"You shall know a word by the company it keeps."

— J.R. Firth (1957)



Images from <https://ig.ft.com/generative-ai/>

Word Embedding

Representing sentences

▪ Natural Language

“The man loves his son”

Clearly, this is a *sequence*, of words

How can each word be represented by a numerical vector?

First idea: one hot encoding

Given a dictionary of W words, each word w could be assigned a unique vector

$$\mathbf{v}_w \in \{0, 1\}^W$$

- *Not particularly efficient: large vectors with almost entirely filled with zeros*
- *The ordering of components will be meaningless: similarities among words will not be represented at all*

Representing sentences

▪ Natural Language

“The man loves his son”

Clearly, this is a *sequence*, of words

How can each word be represented, effectively?

Nice-to-have: *similarity* among words

Cosine similarity between two vectors

$$\frac{\mathbf{v}_1 \cdot \mathbf{v}_2}{\|\mathbf{v}_1\| \|\mathbf{v}_2\|} \in [-1, 1]$$

- *Similar words (e.g., “son”, “daughter”) should have a high similarity value*

Representing words

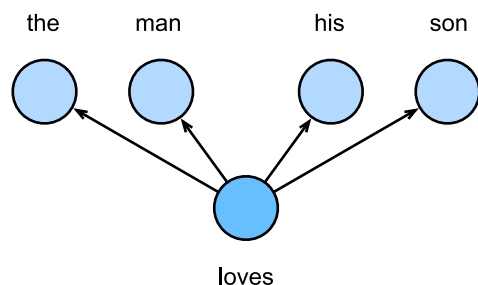
▪ The Skip-Gram Model (word 2vec)

“The man loves his son”

Basic idea: *representing words in relation to their context (in terms of conditional probability)*

$$P(\text{“the”}, \text{“man”}, \text{“his”}, \text{“son”} \mid \text{“loves”})$$

Assuming conditional independence (akin *Naïve Bayesian Classifier*):



the following factorization is correct:

$$= P(\text{“the”} \mid \text{“loves”})P(\text{“man”} \mid \text{“loves”})P(\text{“his”} \mid \text{“loves”})P(\text{“son”} \mid \text{“loves”})$$

Note that the ordering of context words is irrelevant

Representing words

▪ The Skip-Gram Model (word 2vec)

$$P(\text{“the”}, \text{“man”}, \text{“his”}, \text{“son”} \mid \text{“loves”}) = \\ P(\text{“the”} \mid \text{“loves”})P(\text{“man”} \mid \text{“loves”})P(\text{“his”} \mid \text{“loves”})P(\text{“son”} \mid \text{“loves”})$$

Conditional probability *factors* are defined via *softmax*

$$P(w_o \mid w_c) := \frac{\exp(\mathbf{u}_o \cdot \mathbf{v}_c)}{\sum_{i=1}^W \exp(\mathbf{u}_i \cdot \mathbf{v}_c)}$$

under these assumptions:

- each word i in the dictionary is associated to two vectors $\mathbf{u}_i, \mathbf{v}_i \in \mathbb{R}^d$
- \mathbf{u}_i is the vector for i as center word, whereas \mathbf{v}_i is the vector for i as context word
- the dimension d of vectors is an hyperparameter
- W is the vocabulary
- \mathbf{v}_i is the target of the training (embedding vector), whereas \mathbf{u}_i will be discarded after training

Representing words

■ The Skip-Gram Model

A **skip-gram** is a *context* of words in a sentence, corresponding to a '*center*' word

Each skip-gram is obtained from a fixed *window size*, that is, the number of words the context of the *center* word

Each skip-gram (a data item) is of the kind (**center_word, context word**)

Negative Sampling

A dataset for word embedding can be augmented using *negative sampling*: creating skip-grams for words that *do not* occur with the context of the center word in the sentence

Therefore, a skip-gram becomes

(**center_word, context word, label**)
where **label** is either 1 (positive) or 0 (negative)

Window Size	Text	Skip-grams
2	[The wide road shimmered] in the hot sun.	wide, the wide, road wide, shimmered
	The [wide road shimmered in the] hot sun.	shimmered, wide shimmered, road shimmered, in shimmered, the
	The wide road shimmered in [the hot sun].	sun, the sun, hot
3	[The wide road shimmered in] the hot sun.	wide, the wide, road wide, shimmered wide, in
	[The wide road shimmered in the hot] sun.	shimmered, the shimmered, wide shimmered, road shimmered, in shimmered, the shimmered, hot
	The wide road shimmered [in the hot sun].	sun, in sun, the sun, hot

Representing words

▪ Skip-gram: loss function

Given the independence conditions, the likelihood of a textual sentence of length T is:

$$\prod_{t=1}^T \prod_{j \in \text{Ctx}(t)} P(w^{(t)} | w^{(j)})$$

where $\text{Ctx}(t)$ is the context (of fixed length) of word t

Using log-probability:

$$\sum_{t=1}^T \sum_{j \in \text{Ctx}(t)} \log P(w^{(t)} | w^{(j)})$$

where:

$$\log P(w_o | w_c) = \mathbf{u}_o \cdot \mathbf{v}_c - \log \left(\sum_{i=1}^W \exp(\mathbf{u}_i \cdot \mathbf{v}_c) \right)$$

Representing words

▪ Skip gram: gradient

$$\begin{aligned}\frac{\partial \log P(w_o | w_c)}{\partial \mathbf{v}_c} &= \mathbf{u}_o - \frac{\sum_{j=1}^W \exp(\mathbf{u}_j \cdot \mathbf{v}_c) \mathbf{u}_j}{\sum_{i=1}^W \exp(\mathbf{u}_i \cdot \mathbf{v}_c)} \\ &= \mathbf{u}_o - \sum_{j=1}^W \left(\frac{\exp(\mathbf{u}_j \cdot \mathbf{v}_c)}{\sum_{i=1}^W \exp(\mathbf{u}_i \cdot \mathbf{v}_c)} \right) \mathbf{u}_j \\ &= \mathbf{u}_o - \sum_{j=1}^W P(w_j | w_c) \mathbf{u}_j\end{aligned}$$

Representing words

▪ Skip-gram: gradient

$$\begin{aligned}\frac{\partial \log P(w_o | w_c)}{\partial \mathbf{u}_o} &= \mathbf{v}_c - \frac{\exp(\mathbf{u}_o \cdot \mathbf{v}_c) \mathbf{v}_c}{\sum_{i=1}^W \exp(\mathbf{u}_i \cdot \mathbf{v}_c)} \\ &= \mathbf{v}_c - \frac{\exp(\mathbf{u}_j \cdot \mathbf{v}_c)}{\sum_{i=1}^W \exp(\mathbf{u}_i \cdot \mathbf{v}_c)} \mathbf{v}_c \\ &= \mathbf{v}_c - P(w_o | w_c) \mathbf{v}_c\end{aligned}$$

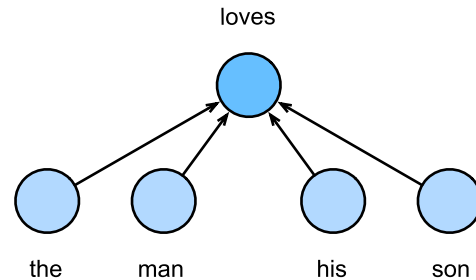
Representing words

■ The Continuous Bag of Words (CBOW) Model

“The man loves his son”

The basic idea is dual to skip-gram: predict center word starting from the context

$$P(\text{“loves”} \mid \text{“the”, “man”, “his”, “son”})$$



Mathematically, this is slightly more complex, since independence assumptions are in the *priors*

Once again, the ordering of context words is irrelevant

Representing words

▪ The Continuous Bag of Words (CBOW) Model

$$P(\text{“loves”} \mid \text{“the”, “man”, “his”, “son”})$$

Conditional probability *factors* are defined via a different *softmax*

$$P(w_c \mid w_{o_1}, \dots, w_{o_m}) = \frac{\exp\left(\frac{1}{m} \mathbf{u}_c \cdot (\mathbf{v}_{o_1} + \dots + \mathbf{v}_{o_m})\right)}{\sum_{i=1}^W \exp\left(\frac{1}{m} \mathbf{u}_i \cdot (\mathbf{v}_{o_1} + \dots + \mathbf{v}_{o_m})\right)}$$

From this point on, the derivation is similar.

Representing words

■ word2vec

- Word *vectors* are used to represent words, can also be considered as *feature vectors*
- The technique of mapping words to real vectors is called word *embedding*
- The **word2vec** tool contains both the skip-gram and continuous bag of words models
- The *skip-gram model* assumes that a word can be used to generate its surrounding words in a text sequence
- The *continuous bag of words* model assumes that a center word is generated based on its surrounding context words

Skip-gram or CBOW?

According to [Mikolov et al., 2013] Skip-Gram works well with small datasets and can better represent less frequent words

However, CBOW is considered to train faster than Skip-Gram and better in representing more frequent words

Say It with Tokens

Words vs Tokens

So far, we have assumed that binary vectors encode entire words, in natural language

Large Language Models (LLM) use *token* instead:

- sentences are pre-processed
- words and symbols are split apart, whit spaces removed
- individual words are further split apart

The man loves the sun!

Clear Show example

Tokens	Characters
6	23

The man loves the sun!

Text Token IDs

[791, 893, 16180, 279, 7160, 4999]

Text Token IDs

Punctuation marks and symbols translate into tokens

Note the different encoding of words 'The' and 'the'

Words vs Tokens

So far, we have assumed that binary vectors encode entire words, in natural language

Large Language Models (LLM) use *token* instead:

- sentences are pre-processed
- words and symbols are split apart, whit spaces removed
- individual words are further split apart

Occorre amare i propri figli!

Clear Show example

Tokens	Characters
9	30

Occorre amare i propri figli!

Text Token IDs

[22513, 93533, 1097, 548, 602, 21744, 4237, 747, 4999]

Text Token IDs

Typically, latin languages generate more sub-word tokens

Byte-Pair Encoding (BPE)

(see: <https://huggingface.co/learn/nlp-course/chapter6/5>)

Words: [("hug", 10), ("pug", 5), ("pun", 12), ("bun", 4), ("hugs", 5)]

Tokens: ["b", "g", "h", "n", "p", "s", "u"]

Corpus: ("h" "u" "g", 10), ("p" "u" "g", 5), ("p" "u" "n", 12), ("b" "u" "n", 4), ("h" "u" "g" "s", 5)

Most frequent pair: "u" + "g" 20

Tokens: ["b", "g", "h", "n", "p", "s", "u", "ug"]

Corpus: ("h" "ug", 10), ("p" "ug", 5), ("p" "u" "n", 12), ("b" "u" "n", 4), ("h" "ug" "s", 5)

Most frequent pair: "u" + "n" 16

Tokens: ["b", "g", "h", "n", "p", "s", "u", "ug", "un"]

Corpus: ("h" "ug", 10), ("p" "ug", 5), ("p" "un", 12), ("b" "un", 4), ("h" "ug" "s", 5)

Most frequent pair: "h" + "ug" 15

Tokens: ["b", "g", "h", "n", "p", "s", "u", "ug", "un", "hug"]

Corpus: ("hug", 10), ("p" "ug", 5), ("p" "un", 12), ("b" "un", 4), ("hug" "s", 5)

...

Tokens reaching zero frequency are removed from the vocabulary

frequency