

Natural Language Processing

Lecture 4 : Words and Meaning

Master Degree in Computer Engineering
University of Padua
Lecturer : Giorgio Satta

Lecture partially based on material originally developed by :
Marco Kuhlman, Linköping University
Mark-Jan Nederhof, University of St. Andrews
Elena Voita, University of Edinburgh

Words and meaning



What is the **meaning** of a word?

The linguistic study of word meaning is called **lexical semantics**.

A model of word meaning should allow us to relate different words and draw inferences to address meaning-related tasks.

We have already mentioned in lecture 2 the distinction between internal and external semantic structure of words.

A word-form can have multiple meanings; each meaning is called a **word sense**, or sometimes a **synset**.

Clipping of synonym and set.

Example : The word-form **mouse** can refer to the rodent or the cursor control device.

Word sense disambiguation (WSD) is the task of determining which sense of a word is being used in a particular context.

Lexical semantic relationship between words are important components of word meaning.

Two words are **synonyms** if they have a common word sense.

Example : **car** and **automobile**

Two words are **similar** if they have similar meanings.

Example : **car** and **bicycle**.

SimLex-999 is a dataset of word similarities for English.

Two words are **related** if they refer to related concepts.

Example : **car** and **gasoline**.

Words and meaning

Two words are **antonyms** if they define a binary opposition.

Example : 'hot' and 'cold'.

One word is an **hyponym** of another if the first has a more specific sense. Notions of **hypernym** or **hyperonym** are defined symmetrically.

Example : 'car' and 'vehicle' (subordinate).

Also called subordinate and superordinate relations.

Words can have **affective** meanings, implying positive or negative connotations/evaluation

Example : 'happy' and 'sad'; 'great' and 'terrible'.

Connection with word taxonomies and word ontologies.

Distributional semantics



It is very difficult to define the notion of word sense in a way that can be understood by computers.

And also in a way that can be understood by humans!

We take a radically different approach, already foreseen in the following works:

The meaning of a word is its use in the language

Ludwig Wittgenstein

Philosophical Investigations, 1953

You shall know a word by the company it keeps

John Rupert Firth

Selected papers, 1957

Suppose you don't know the meaning of the word 'tezgüino', but you have seen it in the following contexts

- a) A bottle of _____ is on the table.
- b) Everybody likes _____.
- c) Don't have _____ before you drive.
- d) We make _____ out of corn.

What other words fit into these contexts?

	a)	b)	c)	d)
tezgüino	1	1	1	1
motor oil	1	0	0	1
tortillas	0	1	0	1
choices	0	1	0	0
wine	1	1	1	0

Based on these vectors, we conclude that 'wine' is very similar to 'tezgüino'.

Distributional semantics develops methods for quantifying semantic similarities between words based on their distributional properties, meaning their **neighboring words**.

The basic idea lays in the so-called **distributional hypothesis**: linguistic items with similar distributions have similar meanings.

Thus the meaning of a word is defined by its distribution in language use.

The basic approach is to collect distributional information in **high-dimensional vectors**, and to define distributional/semantic similarity in terms of vector similarity.

Distributional semantics

Similar words are mapped into “close enough” vectors.

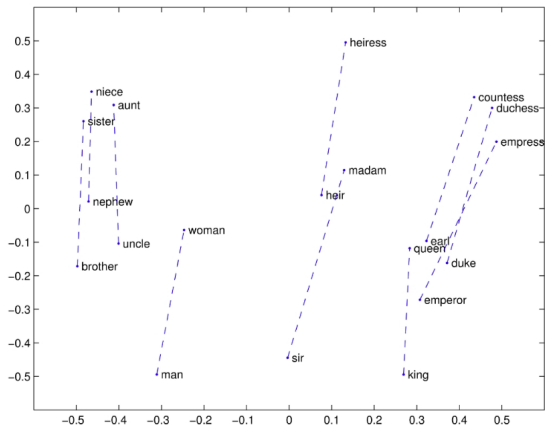
Example : Vectors from sentiment analysis application, projected into 2-dimensional space.



Distributional semantics

Lexical semantic relationships, represented as word vector differences, are generally preserved.

Example : Vectors projected into 2-dimensional space.



Distributional semantics

The obtained vectors are called **word embeddings**.

Each discrete word is embedded in a continuous vector space.

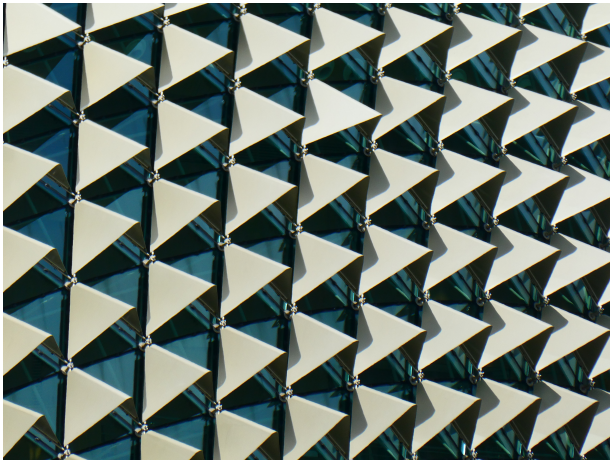
The approach is called **vector semantics** and is the standard way to represent word meaning in NLP. All the learning algorithms we report are **unsupervised**.

We discuss two families of word embeddings:

- **Sparse vectors**: Vector components are computed through some function of the counts of nearby words.
- **Dense vectors**: Vector components are computed through some optimisation or approximation process.

Alternative classifications: frequency-based vs. prediction-based embeddings, and static vs. dynamic or contextual embeddings, presented in a later lecture.

Review: vectors



Kevin Jiner from Unsplash

Review: vectors

We denote vectors as $\mathbf{v} = (v_1, \dots, v_N)$, where $v_1, \dots, v_N \in \mathbb{R}$.

Vectors can be added and subtracted “pointwise”:

$$\mathbf{v} + \mathbf{w} = (v_1 + w_1, \dots, v_N + w_N)$$

$$\mathbf{v} - \mathbf{w} = (v_1 - w_1, \dots, v_D - w_D)$$

Review: vectors

Vector **norm** is defined as

$$|\mathbf{v}| = \sqrt{\sum_{i=1}^N v_i^2}$$

Also known as Euclidean norm.

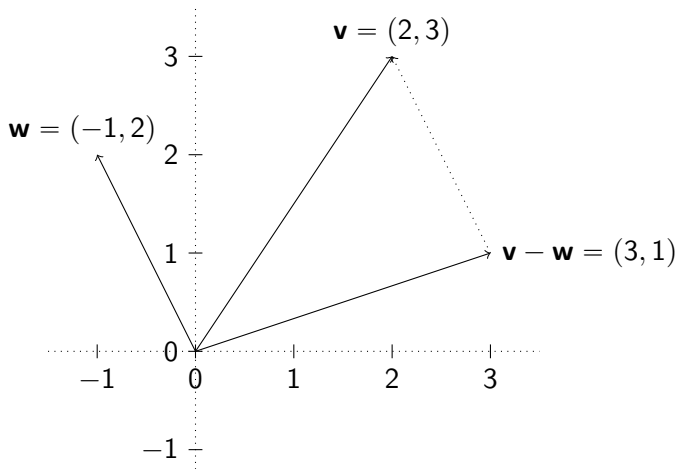
Vectors can be **normalised** to unit norm

$$\mathbf{w} = \mathbf{v}/|\mathbf{v}| = (v_1/|\mathbf{v}|, \dots, v_D/|\mathbf{v}|)$$

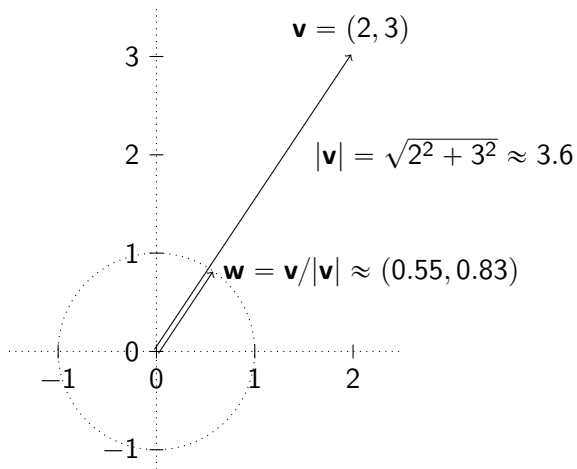
so that $|\mathbf{w}| = 1$.

Review: vectors

Special case of $N = 2$:



Review: vectors



Dot product (or internal product) of two vectors is defined as

$$\mathbf{v} \cdot \mathbf{w} = \sum_i v_i w_i$$

Also known as inner product.

We have:

$$\mathbf{v} \cdot \mathbf{w} = |\mathbf{v}| |\mathbf{w}| \cos \theta$$

where θ is the angle between \mathbf{v} and \mathbf{w} .

The **cosine similarity** metric between two vectors is computed as

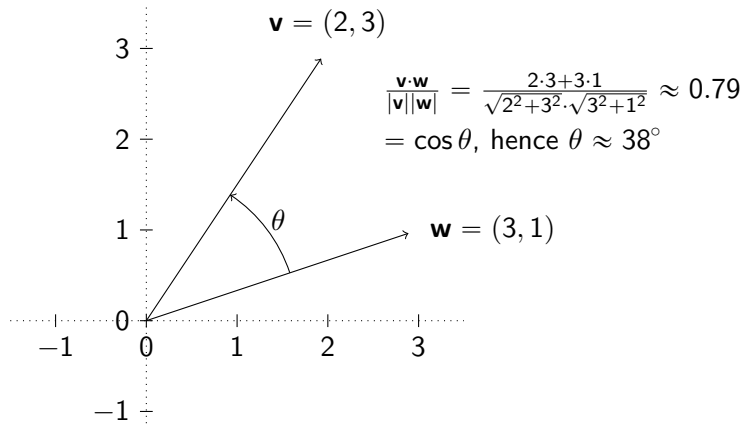
$$\text{cosine}(\mathbf{v}, \mathbf{w}) = \frac{\mathbf{v} \cdot \mathbf{w}}{|\mathbf{v}| |\mathbf{w}|} = \cos \theta$$

In general:

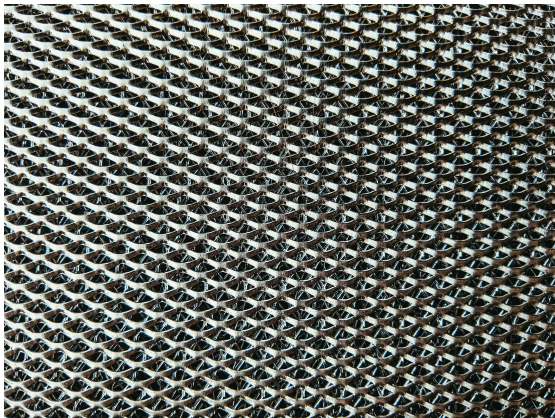
- $\text{cosine}(\mathbf{v}, \mathbf{v}) = 1$
- $\text{cosine}(\mathbf{v}, \mathbf{w}) = 0$ if \mathbf{v} and \mathbf{w} are orthogonal
- $\text{cosine}(\mathbf{v}, \mathbf{w}) = -1$ if \mathbf{v} and \mathbf{w} are in opposite directions

This will later be used to express the similarity between two words, normalization removes frequency information.

Review: vectors



Term-context matrix



Dimitry Zub on Unsplash

Term-context matrix

Let $\{w_1, \dots, w_W\}$ be a set of **term** (target) words and let $\{c_1, \dots, c_C\}$ be a set of **context** words.

For each term word w_i in the corpus, we consider context words appearing inside a **window** of size L , at the left or at the right of w_i .

A **term-context matrix** F is a matrix of size $W \times C$, where each element $f_{i,j}$ is the number of times context word c_j appears in the context of term word w_i in the corpus.

Example

Term words: 'cherry', 'digital', 'information'

Context words: 'pie', 'data', 'computer'

Term-context matrix F :

	pie	data	computer
cherry	442	8	2
digital	5	1683	1670
information	5	3982	3325

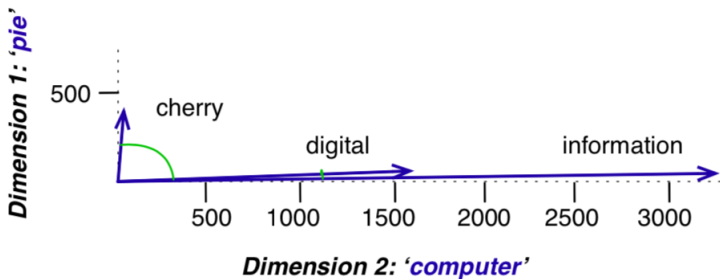
Cosine similarity between vectors associated with term words:

$$\cos(\text{cherry}, \text{information}) = \frac{442 * 5 + 8 * 3982 + 2 * 3325}{\sqrt{442^2 + 8^2 + 2^2} \sqrt{5^2 + 3982^2 + 3325^2}} = .018$$

$$\cos(\text{digital}, \text{information}) = \frac{5 * 5 + 1683 * 3982 + 1670 * 3325}{\sqrt{5^2 + 1683^2 + 1670^2} \sqrt{5^2 + 3982^2 + 3325^2}} = .996$$

Example

We restrict to dimensions 'pie' and 'computer' for convenience.



The model states that 'information' is way closer to 'digital' than it is to 'cherry'.

Observe that the difference in frequency between 'information' and 'digital' results in an undesired difference along the 'computer' dimension.



©tNotice_Group, Simone Favarin

Pointwise mutual information (PMI) is a measure of how often two events x and y occur, compared with what we would expect if they were independent

$$I(x, y) = \log_2 \frac{P(x, y)}{P(x)P(y)}$$

Properties of PMI:

- Symmetric: $I(x, y) = I(y, x)$
- $I(x, y)$ can take positive or negative values
- If the random variables are independent, $P(x, y) = P(x)P(y)$ and $I(x, y) = 0$

The mutual information of random variables X and Y is the expected value of the PMI w.r.t. distribution $P(X, Y)$.

Pointwise mutual information

The pointwise mutual information between w_i and c_j is

$$\text{PMI}(w_i, c_j) = \log_2 \frac{P(w_i, c_j)}{P(w_i)P(c_j)}$$

The numerator tells us how often we observed the two words together.

The denominator tells us how often we would expect the two words to co-occur assuming they each occurred independently.

The ratio gives us an estimate of how much more the two words co-occur than we expect by chance.

Pointwise mutual information

Using the definition of conditional probability

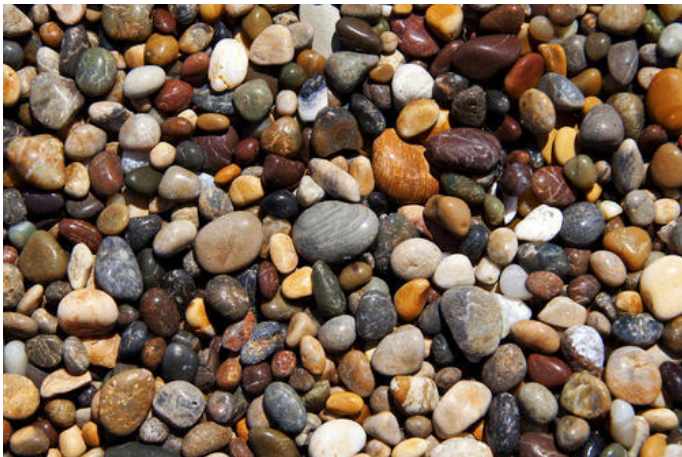
$$P(w_i | c_j) = \frac{P(w_i, c_j)}{P(c_j)}$$

we can rewrite PMI as

$$\text{PMI}(w_i, c_j) = \log_2 \frac{P(w_i | c_j)}{P(w_i)}$$

PMI may then be understood as quantifying how much our confidence in the outcome w_i increases after we observe c_j .

Probability estimation



©indiamart

How do we **estimate** probabilities $P(w_i, c_j)$, $P(w_i)$ and $P(c_j)$?

We estimate the probability $P(w_i, c_j)$ as

$$P(w_i, c_j) = \frac{f_{i,j}}{\sum_{i=1}^W \sum_{j=1}^C f_{i,j}}$$

Note that the denominator is the total number of possible target/context pairs (observations).

We estimate the probability $P(w_i)$ as

$$\begin{aligned} P(w_i) &= \sum_{j=1}^C P(w_i, c_j) \\ &= \frac{\sum_{j=1}^C f_{i,j}}{\sum_{i=1}^W \sum_{j=1}^C f_{i,j}} \end{aligned}$$

Similarly, we estimate the probability $P(c_j)$ as

$$\begin{aligned} P(c_j) &= \sum_{i=1}^W P(w_i, c_j) \\ &= \frac{\sum_{i=1}^W f_{i,j}}{\sum_{i=1}^W \sum_{j=1}^C f_{i,j}} \end{aligned}$$

It is difficult to estimate negative values of PMI:

- Assume two words with probabilities $\sim 10^{-6}$.
- We would need to be certain that the probability of the two words occurring together is **significantly smaller** than $\sim 10^{-12}$.

This kind of granularity would require an enormous corpus.

Use **positive PMI** (PPMI)

$$\text{PPMI}(w_t, w_c) = \max \left\{ \log_2 \frac{P(w_t, w_c)}{P(w_t)P(w_c)}, 0 \right\}$$

Example

Let us pretend that the ones reported below are the only term/context pairs that matter.

The textbook also implicitly assumes that each term token has exactly one context token (1-2-1 relation).

	computer	data	result	pie	sugar	count(w)
cherry	2	8	9	442	25	486
strawberry	0	0	1	60	19	80
digital	1670	1683	85	5	4	3447
information	3325	3982	378	5	13	7703
<u>count(context)</u>	4997	5673	473	512	61	11716

With the above assumptions, sum over rightmost column is the same as sum over bottom line, and equals the size of the corpus.

Example

We apply the relative frequency estimator:

$$P(w=\text{information},c=\text{data}) = \frac{3982}{11716} = .3399$$

$$P(w=\text{information}) = \frac{7703}{11716} = .6575$$

$$P(c=\text{data}) = \frac{5673}{11716} = .4842$$

$$\text{ppmi}(\text{information},\text{data}) = \log_2(.3399/ (.6575 * .4842)) = .0944$$

Example

The resulting word and joint probabilities:

	$p(w, context)$					$p(w)$
	computer	data	result	pie	sugar	$p(w)$
cherry	0.0002	0.0007	0.0008	0.0377	0.0021	0.0415
strawberry	0.0000	0.0000	0.0001	0.0051	0.0016	0.0068
digital	0.1425	0.1436	0.0073	0.0004	0.0003	0.2942
information	0.2838	0.3399	0.0323	0.0004	0.0011	0.6575
<u>$p(context)$</u>	0.4265	0.4842	0.0404	0.0437	0.0052	

The resulting PPMI term-context matrix:

	computer	data	result	pie	sugar
cherry	0	0	0	4.38	3.30
strawberry	0	0	0	4.10	5.51
digital	0.18	0.01	0	0	0
information	0.02	0.09	0.28	0	0

Practical issues



PMI has the problem of being **biased** toward infrequent events: very rare words tend to have very high PMI values.

One way to reduce this bias is to slightly change the computation for $P(w_c)$ in the PPMI

$$\text{PPMI}_\alpha(w_t, w_c) = \max \left\{ \log_2 \frac{P(w_t, w_c)}{P(w_t)P_\alpha(w_c)}, 0 \right\}$$
$$P_\alpha(w_c) = \frac{f(w_c)^\alpha}{\sum_{v \in V} (f(v))^\alpha}$$

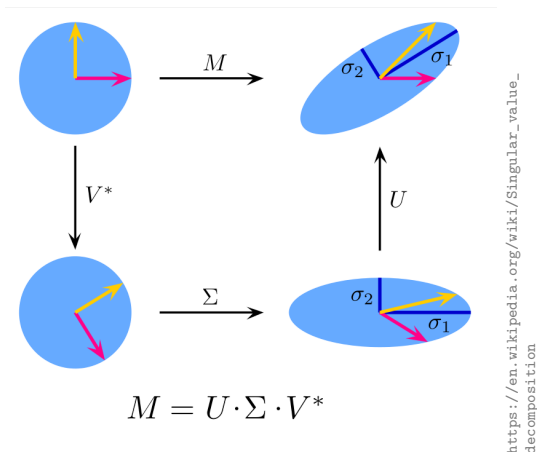
Levy *et al.* (2015) found that a setting of $\alpha = 0.75$ improved performance of embeddings on a wide range of tasks.

We can use the rows of the PPMI term-context matrix as word embeddings.

Notice that these vectors:

- have the size of the vocabulary, which can be quite **large**
- when viewed as arrays, they are very **sparse**

Dot product between these word embeddings needs to use specialised software libraries, implementing vectors as **dictionaries** rather than arrays.



We can obtain **dense** word embeddings from the PPMI term-context matrix.

This is done through a matrix approximation technique known as truncated singular value decomposition.

The approach is also used in information retrieval, under the name of latent semantic analysis (LSA).

For a matrix $\mathbf{A} \in \mathbb{R}^{n \times m}$, the **Frobenius norm** is

$$\|\mathbf{A}\|_F = \sum_{i=1}^n \sum_{j=1}^m a_{i,j}^2$$

where the $a_{i,j}$'s are the elements of the matrix.

Let $|V|$ be the vocabulary size, and let $K \ll |V|$ be the size of the desired embedding.

Let $\mathbf{P} \in \mathbb{R}^{|V| \times |V|}$ be the PPMI term-context matrix.

Let $\mathbf{U}, \mathbf{V} \in \mathbb{R}^{|V| \times K}$ be **learnable** parameters.

The **basic idea** is to approximate \mathbf{P} by means of a matrix $\tilde{\mathbf{P}}(\mathbf{U}, \mathbf{V}) \in \mathbb{R}^{|V| \times |V|}$ such that

$$\min_{\mathbf{U}, \mathbf{V}} \|\mathbf{P} - \tilde{\mathbf{P}}(\mathbf{U}, \mathbf{V})\|_F$$

We later ignore $\tilde{\mathbf{P}}(\mathbf{U}, \mathbf{V})$ and retain \mathbf{U} , which provides the desired word embeddings.

Truncated singular value decomposition (truncated SVD) solves the following constrained optimisation problem

$$\begin{aligned} & \min && \| \mathbf{P} - \mathbf{U}\mathbf{S}\mathbf{V}^{\top} \|_F \\ & \mathbf{U} \in \mathbb{R}^{|\mathcal{V}| \times K} \\ & \mathbf{S} \in \mathbb{R}^{K \times K} \\ & \mathbf{V} \in \mathbb{R}^{|\mathcal{V}| \times K} \end{aligned}$$

subject to conditions $\mathbf{U}^{\top}\mathbf{U} = \mathbb{I}$, $\mathbf{V}^{\top}\mathbf{V} = \mathbb{I}$, and \mathbf{S} is diagonal with its elements in decreasing order.

Orthonormality ensures that all pairs of columns in \mathbf{U} or \mathbf{V} are uncorrelated, so that each dimension conveys unique information.

\mathbf{U} are the target embeddings, $\mathbf{S}\mathbf{V}^{\top}$ are the context embeddings.

Let $\mathbf{U}_k, \mathbf{V}_k \in \mathbb{R}^{|\mathcal{V}| \times 1}$ be the k -th columns of matrices \mathbf{U} and \mathbf{V} , respectively. In truncated SVD, \mathbf{P} is approximated by the sum of K rank-1 matrices

$$\mathbf{P} \approx \sum_{k=1}^K s_k \mathbf{U}_k \mathbf{V}_k^T$$

$\mathbf{U}_k \mathbf{V}_k^T \in \mathbb{R}^{|\mathcal{V}| \times |\mathcal{V}|}$ is called the outer product.

Equivalently, each element $p_{i,j}$ of \mathbf{P} is approximated through the dot product between the embedding of the i -th target word and the embedding of the j -th context word, weighted by the singular values

$$p_{i,j} \approx \sum_{k=1}^K s_k u_{i,k} v_{j,k}$$

Truncated singular value decomposition

We can alternatively view the previous approximation approach as a factorisation problem for matrix \mathbf{P}

$$\mathbf{P} = \mathbf{U}\mathbf{S}\mathbf{V}^T$$

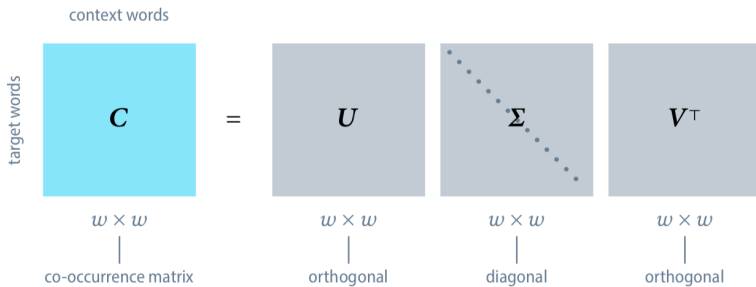
where $\mathbf{U}, \mathbf{S}, \mathbf{V} \in \mathbb{R}^{|V| \times |V|}$, $\mathbf{U}^T \mathbf{U} = \mathbb{I}$, $\mathbf{V}^T \mathbf{V} = \mathbb{I}$, and \mathbf{S} is diagonal with **singular values** $s_1 > s_2 > \dots > s_V$ in its diagonal.

The solution always exists if \mathbf{P} has full rank.

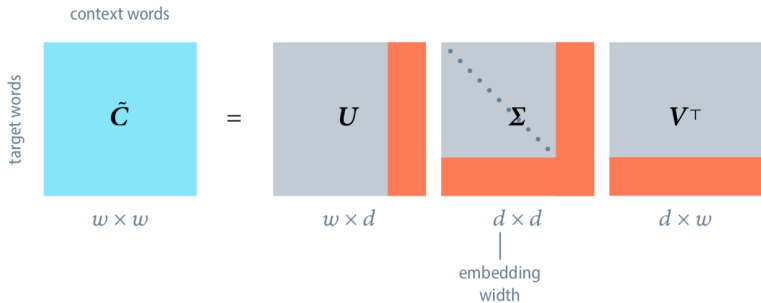
We then choose $K \ll |V|$ on the basis of **mass distribution** at the singular values, and truncate matrices $\mathbf{U}, \mathbf{S}, \mathbf{V}$ accordingly.

The equivalence of the two views is stated by the Eckart-Young-Mirsky theorem. Several techniques have been presented in the literature for the choice of K .

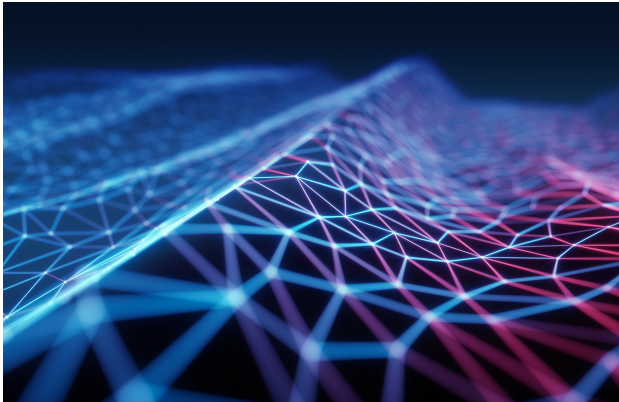
Truncated singular value decomposition



Truncated singular value decomposition



Neural word embeddings



©STAR7

Neural word embeddings

Word embeddings can also be computed using **neural networks**.

Generally speaking, neural word embeddings have the following properties:

- small number of dimensions, in the range [50..1000]
- vectors are dense
- the dimensions don't have a clear interpretation
- components can be negative

Neural word embeddings **work better** than previous embeddings in every NLP task.

Neural word embeddings

The neural word embeddings we introduce in this lecture are **static**, meaning that the presented methods learn one fixed embedding for each word in the vocabulary.

Later we will introduce dynamic (or contextual) embeddings such as BERT, which depend on specific sentences or documents.

Word2vec is a software package including two different algorithms for learning word embeddings:

- skip-gram with negative sampling (SGNS)
- continuous bag-of-words (CBOW)

We look into skip-gram in detail.

Idea : Instead of counting how often a context word appears near a target word, we train a classifier on the following binary

prediction task:

Is a given context word likely to appear near a given target word?

We don't really care about this prediction task: instead, we use the learned parameters as the word embeddings.

Skip-gram

More specifically, the basic steps of the **skip-gram** algorithm are as follows.

For each target word $w_t \in V$:

- treat w_t and any neighboring context word w_c as **positive examples**
- randomly sample other words $w_n \in V$, called **noise words**, to produce **negative examples** for w_t

Use **logistic regression** to train a classifier to distinguish positive and negative examples.

Use the learned weights as the embeddings.

For any pair of words $w_t, u \in V$, we need to define the probability that u is/is not a context word for target word w_t

$$P(+ \mid w_t, u)$$

$$P(- \mid w_t, u) = 1 - P(+ \mid w_t, u)$$

For each $w \in V$, we construct two complementary **embeddings**

- a target embedding $\mathbf{e}_t(w) \in \mathbb{R}^d$
- a context embedding $\mathbf{e}_c(w) \in \mathbb{R}^d$

Integer d is the size of the embedding. Operator $\mathbf{e}_t(w)$ is always assigned to target words, operator $\mathbf{e}_c(w)$ is always assigned to context or noise words.

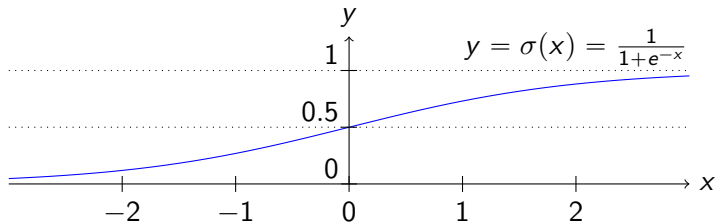
Logistic regression

We use the **dot product**

$$\mathbf{e}_t(w) \cdot \mathbf{e}_c(u) = |\mathbf{e}_t(w)| |\mathbf{e}_c(u)| \cos \theta$$

where θ is the angle between $\mathbf{e}_t(w)$ and $\mathbf{e}_c(u)$.

We also use the **logistic sigmoid function** σ , which maps real values to probabilities (values between 0 and 1)



We can now define

$$P(+ | w, u) = \sigma(\mathbf{e}_t(w) \cdot \mathbf{e}_c(u)) = \frac{1}{1 + \exp(-\mathbf{e}_t(w) \cdot \mathbf{e}_c(u))}$$

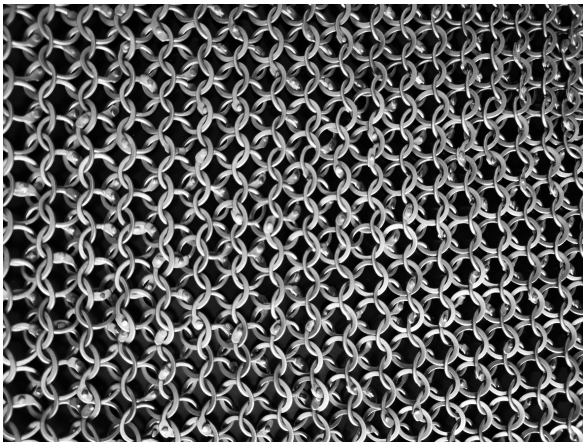
We then have (for fixed vector module)

- vectors $\mathbf{e}_t(w)$, $\mathbf{e}_c(u)$ with small angle: positive cosine similarity, $P(+ | w, u)$ close to one
- vectors $\mathbf{e}_t(w)$, $\mathbf{e}_c(u)$ with large angle: negative cosine similarity, $P(+ | w, u)$ close to zero

We also have

$$\begin{aligned} P(- | w, u) &= 1 - P(+ | w, u) = \sigma(-\mathbf{e}_t(w) \cdot \mathbf{e}_c(u)) \\ &= \frac{1}{1 + \exp(\mathbf{e}_t(w) \cdot \mathbf{e}_c(u))} \end{aligned}$$

Well-known property of the sigmoid function.



Yura Timoshenko on Unsplash

For simplicity, let us consider a dataset with only one target/context pair (w, u) , along with k noise words v_1, v_2, \dots, v_k (negative examples).

Negative examples are related to contrastive learning paradigm.

Skip-gram makes the simplifying assumption that all (positive and negative) context words are **independent**. Then the log-likelihood of the data is

$$\begin{aligned} LL_w &= \log \left(P(+ | w, u) \prod_{i=1}^k P(- | w, v_k) \right) \\ &= \log P(+ | w, u) + \sum_{i=1}^k \log P(- | w, v_k) \end{aligned}$$

$$\begin{aligned} LL_w &= \log P(+ | w, u) + \sum_{i=1}^k \log(1 - P(+ | w, v_k)) \\ &= \log \sigma(\mathbf{e}_t(w) \cdot \mathbf{e}_c(u)) + \sum_{i=1}^k \log(1 - \sigma(\mathbf{e}_t(w) \cdot \mathbf{e}_c(v_k))) \end{aligned}$$

We can alternatively minimise the inverse of LL_w , which in case of the logistic regression is the **cross-entropy** loss function

$$L_{CE} = -\log \sigma(\mathbf{e}_t(w) \cdot \mathbf{e}_c(u)) - \sum_{i=1}^k \log(1 - \sigma(\mathbf{e}_t(w) \cdot \mathbf{e}_c(v_k)))$$

We are omitting term $\frac{1}{(k+1)}$, regarding it as a constant.

The model parameters are the two $\mathbb{R}^{|V| \times d}$ matrices with the target and context embeddings $\mathbf{e}_t(w), \mathbf{e}_c(w), w \in V$, which are randomly initialised.

By making an **update** to minimise L_{CE} we

- force an increase in similarity (dot product) between $\mathbf{e}_t(w)$ and $\mathbf{e}_c(u)$
- force a decrease in similarity between $\mathbf{e}_t(w)$ and $\mathbf{e}_c(v_k)$, for all of the noise words v_k .

When several target words are involved, L_{CE} is not convex and there are several local minima.

We train the model with **stochastic gradient descent**, as usual for logistic regression.

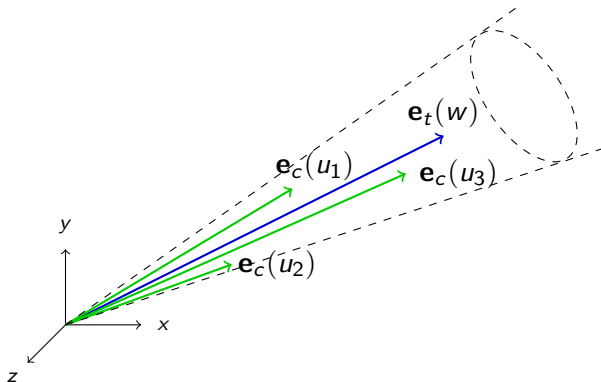
Textbook derives the gradient for the cross-entropy loss function.

In the end, we retain the target embeddings $\mathbf{e}_t(w)$ and ignore the context embeddings $\mathbf{e}_c(w)$.

In some applications, improvement in performance are reported by using $\mathbf{e}_t(w) + \mathbf{e}_c(w)$ as the final word embeddings.

Geometric interpretation

Target vector $\mathbf{e}_t(w)$ is surrounded by its context vectors $\mathbf{e}_c(u_j)$, which form a “thin” hypercone (dashed lines)



Geometric interpretation

Let w be a target word with context words u_j 's. Let also w' be a target word similar to w .

Then we have that

- $\mathbf{e}_c(u_j)$ form a thin cone surrounding vector $\mathbf{e}_t(w)$
- vector $\mathbf{e}_t(w')$ must be placed within that cone, because w' and w share their context words
- $\mathbf{e}_t(w)$ and $\mathbf{e}_t(w')$ are forced to be at a **small angle**

As a result, similar words cluster together.

With a similar reasoning, we can argue that different words repel each other.

Practical issues



Clean up text: convert any **punctuation** into token <PERIOD>.

Remove all words w such that $f(w) \leq D$ (usually $D = 5$):

- greatly reduces issues due to noise in the data
- improves the quality of the vector representation

Indexing: convert words in V to non-negative integers and back again, this makes your code simpler.

Usually, integers are assigned in descending frequency order (most frequent word is assigned integer 0).

Practical issues

Words that show up very often, such as 'the', 'of', and 'for', don't provide much context and are often regarded as noise.

The process of discarding occurrences of these words by means of some probability distribution is called **subsampling**.

We **discard** occurrence w_i with probability given by

$$P(w_i) = 1 - \sqrt{\frac{t}{z(w_i)}}$$

where $z(w_i)$ is the normalized frequency of w_i , and t is a threshold parameter depending on dataset size.

Example : With 10^7 word occurrences you may use $t = 10^{-5}$. Then if $z(w_i) = 0.1$ we get $P(w_i) = 0.99$, that is, we leave in the data 10^4 occurrences of w_i .

The ratio k between the positive and the negative examples is a **hyperparameter** and must be adjusted on the development set.

The noise words are chosen according to the relative frequency estimator and an α **corrector**

$$P_{\alpha}(v) = \frac{f(v)^{\alpha}}{\sum_{w \in V} (f(w))^{\alpha}}$$

Similarly to the PPMI case, $\alpha = .75$ gives better performance in most applications.

Choice of hyperparameter L , specifying context window size, will be discussed later.

Once hyperparameter L is fixed, context window size can be chosen randomly in the range $[1..L]$, with uniform distribution.

This results in **linear decay** for context words, depending on distance from target word:

- words at distance 1 from target will be available with $P = 1$
- words at distance L from target will be available with $P = \frac{1}{L}$

Should we estimate one embedding per word form or else estimate distinct embeddings for each **word sense**?

Intuitively, if word representations are to capture the meaning of individual words, then words with multiple meanings should have multiple embeddings.

However, in certain applications this is unnecessary, because the embedding of a word form is a linear combination of the embeddings of the underlying senses; see Arora *et al.* (2018).

A few other tricks are discussed in the word2vec original papers.



©allrecipes

FastText is an embedding that deals with unknown words and sparsity in languages with rich morphology, by using **subword models**.

FastText is an extension of word2vec.

Each word in fastText is represented by itself plus a bag of character N -grams, for bounded values of N .

A character N -gram is a sequence of N characters. We will introduce N -grams in some later lecture on language models.

Example : With $N = 3$ the word 'where' would be represented by the following strings (_ is a boundary marker)

where, _wh, whe, her, ere, re_

Then a skip-gram embedding is learned for each N -gram, and the word 'where' is represented as the sum of all of the resulting embeddings.

In **morphologically rich languages** (MRLs), words display productive internal structure. These regularity are then captured by fastText.

Global Vectors, or **GloVe** for short, is an embedding that accounts for **global corpus statistics**, which was somehow disregarded by word2vec.

GloVe combines the intuitions of count-based models like PPMI, while also capturing the linear structures used by methods like word2vec.

In the original paper for GloVe (Pennington et al., 2014), a mathematical comparison between this embedding and word2vec is provided.

Visualizing word embeddings

Word embeddings can be used to visualise the meaning of a word w .

Most commonly used methods:

- listing the words in V with highest cosine similarity with w .

Locality-sensitive hashing (LSH) can be used, that hashes similar input items into the same buckets with high probability.

- project the d dimensions of a word embedding down into 2 dimensions.

t -distributed stochastic neighbor embedding (t-SNE) is used, preserving metric properties.

Semantic properties

Both sparse and dense vectors are sensible to the **context window size** used to collect counts, controlled by hyperparameter L .

We generally work with 2 to 20 words in the context.

The choice of L depends on the goals of the representation:

- small values of L provide semantically similar words with the same parts of speech.

With a window ± 2 , 'Hogwarts' provides 'Sunnydale' and 'Evernight' (other fictional schools).

- larger values of L provide words that are topically related but not similar.

With a window ± 5 , 'Hogwarts' provides 'Dumbledore', 'Malfoy' and 'half-blood' (words topically related to the Harry Potter series).

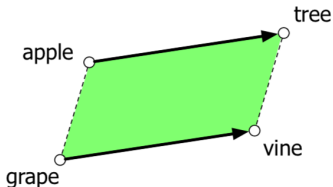
Semantic properties

Another semantic property of embeddings is their ability to capture **relational meanings**.

Word embeddings can be used to solve analogy problems, usually expressed in the form:

Apple is to tree as grape is to _____.

The **parallelogram model** can be used for solving such problems:



Unfortunately, word embeddings also reproduce the implicit biases and stereotypes that are latent in the text.

Example : On a standard dataset, embedding similarly suggests the analogy 'father' is to 'doctor' as 'mother' is to 'nurse'.

This could result in what is called **allocational harm**, when a system allocates resources (jobs or credit) unfairly to different groups.

Embeddings not only reflect the statistics of their input corpora, but have also been reported to amplify bias. This is a current research topic.



Extrinsic evaluation : Use the model to be evaluated in some end-to-end application (as for instance sentiment analysis, machine translation, etc.) and measure performance.

Difficult to do reliably, time consuming.

Intrinsic evaluation : Look at performance of model in isolation, with respect to a given evaluation measure.

For word embedding, several similarity datasets; see next slide.

The most common evaluation metric for embedding models is extrinsic evaluation on end-to-end tasks: machine translation, sentiment analysis, etc.

Three types of intrinsic evaluation

- word **similarity**: compute a numerical score for the semantic similarity between two words;
Example : car, automobile
- word **relatedness**: compute the degree of how much one word has to do with another word;
Example : car, highway
- word **analogy**: comparison of two things to show their similarities;
Example : tree : leaf :: flower : petal

Several datasets available for intrinsic evaluation:

- **WordSim-353** and **SimLex-999** are collection for measuring word relatedness and similarity for word pairs in isolation.
- **Stanford Contextual Word Similarity** (SCWS) and **Word-in-Context** (WiC) datasets provide pairs of words in their sentential context, offering richer evaluation scenarios.
- Several sets from **SemEval-2012 Task 2** provide analogy tasks covering morphology, lexicographic and encyclopedia relations.

Cross-lingual word embedding



©Cristinn – Adobe Stock

Cross-lingual word embeddings

- are very useful in **comparing** the meaning of words across languages
- are key to applications such as machine translation, multilingual lexicon induction, and cross-lingual information retrieval
- enable model **transfer** between resource-rich and low-resource languages, by providing a common representation space

Cross-lingual word embedding

The class of **objective functions** minimized by most cross-lingual word embedding methods can be formulated as

$$J = \mathcal{L}^1 + \dots + \mathcal{L}^\ell + \Omega$$

where \mathcal{L}^i is the monolingual loss of the i -th language and Ω a regularization term.

Joint optimization of multiple non-convex losses is **difficult**

- most approaches optimize one loss at a time, while keeping certain variables fixed
- this step-wise approach is approximate and does not guarantee to reach even a local optimum

The choice of multilingual **parallel data sources** is more important for the model performance than the actual underlying architecture.

Methods differ along the following two dimensions.

- Type of text **alignment**: at the level of words, sentences, or documents, which introduce stronger or weaker supervision.
- Type of text **comparability**: exact translations, weak translation, or similar meaning.

Cross-lingual word embedding

Extending the learning process from bilingual to multilingual settings **improves** results, as reported on many standard tasks.

Variation in ambiguity: ambiguity for word w in one language does not transfer to the translation of w in another language.

The general hypothesis is that variation in ambiguity acts as a form of naturally occurring **supervision**.

Research papers



Alfons Morales on Unsplash

Title: Neural Word Embedding as Implicit Matrix Factorization

Authors: Omer Levy, Yoav Goldberg

Conference: NIPS 2014

Content: This paper shows that skip-gram with negative-sampling implicitly factorises a word-context matrix with PMI shifted by a global constant.

[https://papers.nips.cc/paper/2014/file/
feab05aa91085b7a8012516bc3533958-Paper.pdf](https://papers.nips.cc/paper/2014/file/feab05aa91085b7a8012516bc3533958-Paper.pdf)

Title: Problems With Evaluation of Word Embeddings Using Word Similarity Tasks

Authors: Manaal Faruqui, Yulia Tsvetkov, Pushpendre Rastogi, Chris Dyer

Conference: Proceedings of the 1st Workshop on Evaluating Vector-Space Representations for NLP

Content: In this paper we present several problems associated with the evaluation of word vectors on word similarity datasets, and summarize existing solutions. Our study suggests that the use of word similarity tasks for evaluation of word vectors is not sustainable and calls for further research on evaluation methods.

<https://aclanthology.org/W16-2506/>

Title: A Survey of Cross-lingual Word Embedding Models

Authors: Sebastian Ruder, Ivan Vulić, Anders Sogaard

Conference: Proceedings of the 1st Workshop on Evaluating Vector-Space Representations for NLP

Content: This survey provides a comprehensive typology of cross-lingual word embedding models, and shows that many of the models presented in the literature optimize for the same objectives.

<https://www.jair.org/index.php/jair/article/view/11640>

Have fun!

Publicly available embedding projector.

