# IN5550: Neural Methods in Natural Language Processing
## Sub-lecture 4.1
### *One-hot and dense representations*

Andrey Kutuzov

University of Oslo
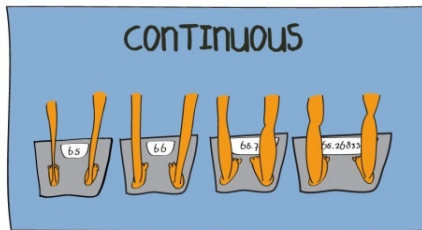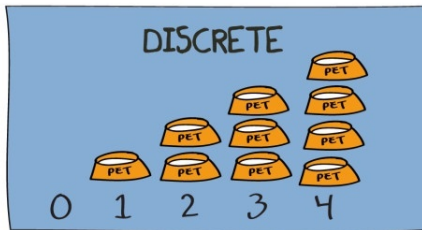
14 February 2023

# Contents

1. Dense Representations of Linguistic Features
   - One-hot representations: let's recall
   - Dense representations (embeddings)

How to make the world continuous?



Discrete and continuous variables

# Dense Representations of Linguistic Features

## Representations

▶ In the obligatory 1, you train neural document classifiers...

# Dense Representations of Linguistic Features

## Representations

▶ In the obligatory 1, you train neural document classifiers...

▶ ...using <span style="color:red">bags of words</span> as features.

# Dense Representations of Linguistic Features

## Representations

- ▶ In the obligatory 1, you train neural document classifiers...
- ▶ ...using bags of words as features.
- ▶ Documents are represented as sparse vocabulary vectors.

# Dense Representations of Linguistic Features

## Representations

▶ In the obligatory 1, you train neural document classifiers...

▶ ...using bags of words as features.

▶ Documents are represented as sparse vocabulary vectors.

▶ Core elements of this representation are words,

# Dense Representations of Linguistic Features

## Representations

- ▶ In the obligatory 1, you train neural document classifiers...
- ▶ ...using bags of words as features.
- ▶ Documents are represented as sparse vocabulary vectors.
- ▶ Core elements of this representation are words,
- ▶ and they are in turn represented with one-hot vectors.

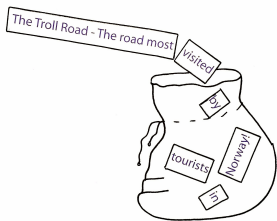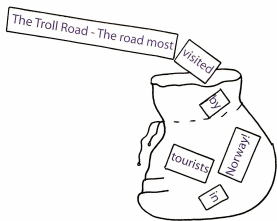▶ BOW feature vector of the document $i$ can be interpreted as a sum of one-hot vectors ($o$) for each token in it:

# One-hot representations: let's recall



▶ BOW feature vector of the document $i$ can be interpreted as a sum of one-hot vectors ($o$) for each token in it:

▶ Vocabulary $V$ from the picture above contains 10 words (lowercased):
['-', 'by', 'in', 'most', 'norway', 'road', 'the', 'tourists', 'troll', 'visited'].

# One-hot representations: let's recall



▶ BOW feature vector of the document $i$ can be interpreted as a sum of one-hot vectors ($o$) for each token in it:

  ▶ Vocabulary $V$ from the picture above contains 10 words (lowercased):
  ['-', 'by', 'in', 'most', 'norway', 'road', 'the', 'tourists', 'troll', 'visited'].
  ▶ $o^0 = [0, 0, 0, 0, 0, 0, 1, 0, 0, 0]$ ('The')
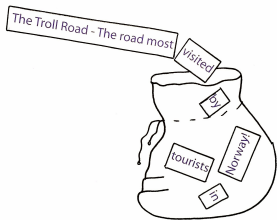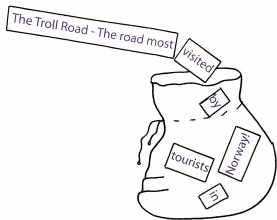
# One-hot representations: let's recall



- ▶ BOW feature vector of the document $i$ can be interpreted as a sum of one-hot vectors ($o$) for each token in it:
  - ▶ Vocabulary $V$ from the picture above contains 10 words (lowercased):
    ['-', 'by', 'in', 'most', 'norway', 'road', 'the', 'tourists', 'troll', 'visited'].
  - ▶ $o^0 = [0, 0, 0, 0, 0, 0, 1, 0, 0, 0]$ ('The')
  - ▶ $o^1 = [0, 0, 0, 0, 0, 0, 0, 0, 1, 0]$ ('Troll')
  - ▶ etc...

- BOW feature vector of the document $i$ can be interpreted as a sum of one-hot vectors ($o$) for each token in it:
    - Vocabulary $V$ from the picture above contains 10 words (lowercased):
      ['-', 'by', 'in', 'most', 'norway', 'road', 'the', 'tourists', 'troll', 'visited'].
    - $o^0 = [0, 0, 0, 0, 0, 0, 1, 0, 0, 0]$ ('The')
    - $o^1 = [0, 0, 0, 0, 0, 0, 0, 0, 1, 0]$ ('Troll')
    - etc...
    - $i = [1, 1, 1, 1, 1, 2, 2, 1, 1, 1]$ ('the' and 'road' occurred 2 times)

# One-hot representations: let's recall

▶ The network is trained on words represented with integer identifiers:
  ▶ '*the*' is the word number 6 in the vocabulary
  ▶ '*most*' is the word number 3 in the vocabulary
  ▶ '*visited*' is the word number 9 in the vocabulary
  ▶ etc.

# One-hot representations: let's recall

- The network is trained on words represented with integer identifiers:
    - '*the*' is the word number 6 in the vocabulary
    - '*most*' is the word number 3 in the vocabulary
    - '*visited*' is the word number 9 in the vocabulary
    - etc.
- Such features are discrete (categorical).
    - a.k.a. one-hot.

# One-hot representations: let's recall

▶ The network is trained on words represented with integer identifiers:
  ▶ '*the*' is the word number 6 in the vocabulary
  ▶ '*most*' is the word number 3 in the vocabulary
  ▶ '*visited*' is the word number 9 in the vocabulary
  ▶ etc.
▶ Such features are discrete (categorical).
  ▶ a.k.a. one-hot.
▶ Each word is a feature on its own, completely independent from other words.

- The network is trained on words represented with integer identifiers:
    - '*the*' is the word number 6 in the vocabulary
    - '*most*' is the word number 3 in the vocabulary
    - '*visited*' is the word number 9 in the vocabulary
    - etc.
- Such features are discrete (categorical).
    - a.k.a. one-hot.
- Each word is a feature on its own, completely independent from other words.
- Other NLP tasks: categorical features for PoS tags, dependency labels, etc.

### Why discrete features might be bad?

▶ Features for words are extremely sparse:

# One-hot representations: let's recall

## Why discrete features might be bad?

▶ Features for words are extremely sparse:
  ▶ the feature '*word form*' can take any of tens or hundreds of thousands categorical values...

### Why discrete features might be bad?

▶ Features for words are extremely sparse:
  ▶ the feature '*word form*' can take any of tens or hundreds of thousands categorical values...
  ▶ ...each absolutely unique and not related to each other.

## Why discrete features might be bad?

▶ Features for words are extremely sparse:
  ▶ the feature '*word form*' can take any of tens or hundreds of thousands categorical values...
  ▶ ...each absolutely unique and not related to each other.
▶ Classifiers have to learn weight matrices with $dim = |V|$.

### Why discrete features might be bad?

▶ Features for words are extremely sparse:
  ▶ the feature '*word form*' can take any of tens or hundreds of thousands categorical values...
  ▶ ...each absolutely unique and not related to each other.
▶ Classifiers have to learn weight matrices with $dim = |V|$.
▶ Not efficient:
  ▶ A 50-words text is $x \in \mathbb{R}^{100000}$, because 100K words in our vocabulary!

### Why discrete features might be bad?

▶ Features for words are extremely sparse:
  ▶ the feature '*word form*' can take any of tens or hundreds of thousands categorical values...
  ▶ ...each absolutely unique and not related to each other.
▶ Classifiers have to learn weight matrices with $dim = |V|$.
▶ Not efficient:
  ▶ A 50-words text is $x \in \mathbb{R}^{100000}$, because 100K words in our vocabulary!
▶ A bit easier for other linguistic entities (parts of speech, etc)...

## Why discrete features might be bad?

▶ Features for words are extremely sparse:
  ▶ the feature '*word form*' can take any of tens or hundreds of thousands categorical values...
  ▶ ...each absolutely unique and not related to each other.
▶ Classifiers have to learn weight matrices with $dim = |V|$.
▶ Not efficient:
  ▶ A 50-words text is $x \in \mathbb{R}^{100000}$, because 100K words in our vocabulary!
▶ A bit easier for other linguistic entities (parts of speech, etc)...
▶ ...but their feature combinations yield millions of resulting features.

## Why discrete features might be bad?

▶ Features for words are extremely sparse:
  ▶ the feature '*word form*' can take any of tens or hundreds of thousands categorical values...
  ▶ ...each absolutely unique and not related to each other.
▶ Classifiers have to learn weight matrices with $dim = |V|$.
▶ Not efficient:
  ▶ A 50-words text is $x \in \mathbb{R}^{100000}$, because 100K words in our vocabulary!
▶ A bit easier for other linguistic entities (parts of speech, etc)...
▶ ...but their feature combinations yield millions of resulting features.
▶ It's very difficult to learn good weights for them all.

## Why discrete features might be bad?

▶ Features for words are extremely sparse:
  ▶ the feature '*word form*' can take any of tens or hundreds of thousands categorical values...
  ▶ ...each absolutely unique and not related to each other.
▶ Classifiers have to learn weight matrices with $dim = |V|$.
▶ Not efficient:
  ▶ A 50-words text is $x \in \mathbb{R}^{100000}$, because 100K words in our vocabulary!
▶ A bit easier for other linguistic entities (parts of speech, etc)...
▶ ...but their feature combinations yield millions of resulting features.
▶ It's very difficult to learn good weights for them all.
▶ Feature extraction step haunted NLP practitioners for several decades.

Core Features +
Feature Combinations

As   McGwire   neared   ,   fans   went   wild

| [went] | [VBD] | [As] | [ADP] | [went] |
|---|---|---|---|---|
| [VERB] | [As] | [IN] | [went, VERB] | [As, ADP] |
| [went, As] | [VBD, ADP] | [went, VERB] | [As, IN] | [went, As] |
| [VERB, IN] | [VBD, As, ADP] | [went, As, ADP] | [went, VBD, ADP] | [went, VBD, As] |
| [ADJ, *, ADP] | [VBD, *, ADP] | [VBD, ADJ, *] | [VBD, ADJ, *] | [NNS, *, ADP] |
| [NNS, VBD, ADP] | [NNS, VBD, *] | [ADJ, ADP, NNP] | [VBD, ADP, NNP] | [VBD, ADJ, NNP] |
| [NNS, ADP, NNP] | [VERB, As, IN] | [went, left, 5] | [VBD, left, 5] | [As, left, 5] |
| [ADP, left, 5] | [VERB, *, IN] | [went, As, IN] | [went, VERB, IN] | [went, VERB, As] |
| [JJ, *, IN] | [NOUN, VERB, *] | [VERB, JJ, IN] | [VERB, JJ, *] | [NOUN, *, IN] |
| [NOUN, VERB, IN] | [NOUN, VERB, NOUN] | [JJ, IN, NOUN] | [VERB, IN, NOUN] | [VERB, JJ, NOUN] |
| [NOUN, IN, NOUN] | [went, VBD, As, ADP] | [went, left, 5] | [VERB, left, 5] | [As, left, 5] |
| [IN, left, 5] | [went, VBD, left, 5] | [VBD, ADJ, *, ADP] | [NNS, VBD, *, ADP] | [VBD, ADJ, ADP, NNP] |
| [NNS, VBD, ADP, NNP] | [VERB, JJ, *, IN] | [As, ADP, left, 5] | [went, As, left, 5] | [VBD, ADP, left, 5] |
| [went, VERB, As, IN] | [As, IN, left, 5] | [NOUN, VERB, *, IN] | [VERB, JJ, IN, NOUN] | [NOUN, VERB, IN, NOUN] |
| [went, VERB, left, 5] | [went, VBD, ADP, left, 5] | [went, As, left, 5] | [VERB, IN, left, 5] | [VBD, As, ADP, left, 5] |
| [went, As, ADP, left, 5] | [VBD, ADJ, *, left, 5] | [went, VBD, As, left, 5] | [ADJ, *, ADP, left, 5] | [VBD, *, ADP, left, 5] |
| [VBD, ADJ, ADP, left, 5] | [VBD, ADP, NNP, left, 5] | [NNS, *, ADP, left, 5] | [NNS, ADP, NNP, left, 5] | [NNS, *, left, 5] |
| [ADJ, ADP, NNP, left, 5] | [went, As, IN, left, 5] | [VBD, ADJ, NNP, left, 5] | [went, VERB, As, left, 5] | [NNS, VBD, NNP, left, 5] |
| [VERB, As, IN, left, 5] | [VERB, JJ, IN, left, 5] | [went, VERB, IN, left, 5] | [NOUN, *, IN, left, 5] | [JJ, *, IN, left, 5] |
| [VERB, *, IN, left, 5] | | [VERB, JJ, *, left, 5] | | [NOUN, VERB, IN, left, 5] |

Example from slides of Rush and Petrov (2012)

Feature model for parsing

'Is the 1st word to the right *wild*, and the 3rd word to the left a *verb*?'

### We can do better

▶ Is there a way to avoid using multitudes of discrete categorical features?

### We can do better
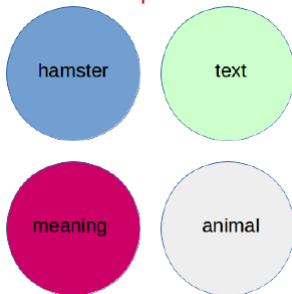
▶ Is there a way to avoid using multitudes of discrete categorical features?
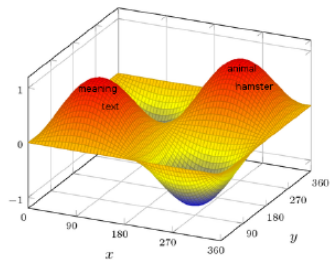
▶ Yes.

▶ Use dense continuous features.

# Dense representations (embeddings)

**Discrete representations**



**Continuous representations**

# Dense representations (embeddings)

### Discrete representations



### Continuous representations



- ► We would like linguistic entities to be represented with some meaningful 'coordinates'.
- ► It would allow our models to understand whether entities (for example, words) are more or less similar with respect to the current task at hand.

# Dense representations (embeddings)

## Vectors as coordinates

- A vector is a sequence or an array of $n$ real values:
  - $[0, 1, 2, 4]$ is a vector with 4 components/entries ($\in \mathbb{R}^4$);
  - $[200, 300, 1]$ is a vector with 3 components/entries ($\in \mathbb{R}^3$);

# Dense representations (embeddings)

## Vectors as coordinates

▶ A vector is a sequence or an array of $n$ real values:
  ▶ $[0, 1, 2, 4]$ is a vector with 4 components/entries ($\in \mathbb{R}^4$);
  ▶ $[200, 300, 1]$ is a vector with 3 components/entries ($\in \mathbb{R}^3$);

▶ Components can be viewed as coordinates in an $n$-dimensional space;

# Dense representations (embeddings)

## Vectors as coordinates

▶ A vector is a sequence or an array of $n$ real values:
  ▶ $[0, 1, 2, 4]$ is a vector with 4 components/entries ($\in \mathbb{R}^4$);
  ▶ $[200, 300, 1]$ is a vector with 3 components/entries ($\in \mathbb{R}^3$);

▶ Components can be viewed as coordinates in an $n$-dimensional space;

▶ then a vector is a point in this space.

# Dense representations (embeddings)
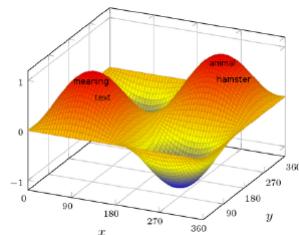
## Vectors as coordinates

▶ A vector is a sequence or an array of $n$ real values:
  ▶ $[0, 1, 2, 4]$ is a vector with 4 components/entries ($\in \mathbb{R}^4$);
  ▶ $[200, 300, 1]$ is a vector with 3 components/entries ($\in \mathbb{R}^3$);

▶ Components can be viewed as coordinates in an $n$-dimensional space;

▶ then a vector is a point in this space.

3-dimensional space:

# Dense representations (embeddings)

## Feature embeddings

▶ Say we have a vocabulary of size $|V|$;

# Dense representations (embeddings)

## Feature embeddings

▶ Say we have a vocabulary of size $|V|$;
▶ Let's embed these words into a $d$-dimensional space.

# Dense representations (embeddings)

### Feature embeddings

▶ Say we have a vocabulary of size $|V|$;

▶ Let's embed these words into a $d$-dimensional space.

▶ $d \ll |V|$

▶ e.g., $d = 100$ for a vocabulary of 100 000 words.

# Dense representations (embeddings)

## Feature embeddings

▶ Say we have a vocabulary of size $|V|$;

▶ Let's embed these words into a $d$-dimensional space.

▶ $d \ll |V|$

▶ e.g., $d = 100$ for a vocabulary of 100 000 words.

▶ Each word is associated with its own $d$-dimensional embedding vector;

# Dense representations (embeddings)

## Feature embeddings

▶ Say we have a vocabulary of size $|V|$;

▶ Let's embed these words into a $d$-dimensional space.

▶ $d \ll |V|$

▶ e.g., $d = 100$ for a vocabulary of 100 000 words.

▶ Each word is associated with its own $d$-dimensional embedding vector;

▶ These embeddings are part of $\theta$;

▶ can be trained together with the rest of the network.

# Dense representations (embeddings)

## Feature embeddings

- ► Say we have a vocabulary of size $|V|$;
- ► Let's embed these words into a $d$-dimensional space.
- ► $d \ll |V|$
- ► e.g., $d = 100$ for a vocabulary of 100 000 words.
- ► Each word is associated with its own $d$-dimensional embedding vector;
- ► These embeddings are part of $\theta$;
- ► can be trained together with the rest of the network.

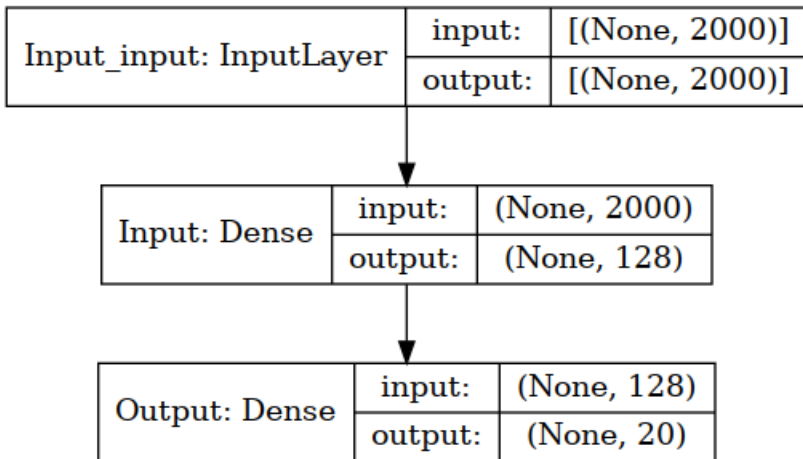Your models in the Obligatory 1 implicitly do this during the training.
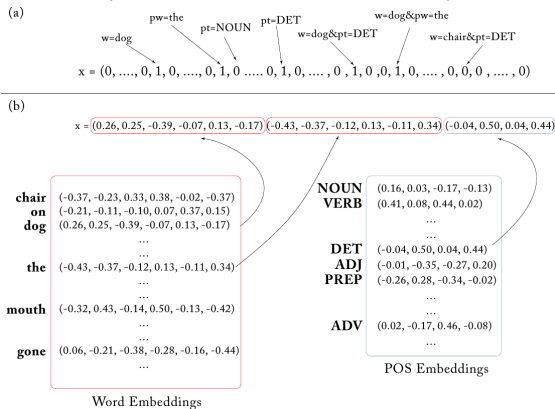This is representation learning.

# Representation learning

Sparse (a) and dense (b) feature representations for
'the_DET dog'
(part of speech tagging task)



(a)

pw=the
pt=NOUN    pt=DET
w=dog          w=dog&pw=the
w=dog&pt=DET
w=chair&pt=DET

x = (0, ...., 0, 1, 0, ...., 0, 1, 0 ..... 0, 1, 0, ... , 0 , 1, 0 ,0, 1, 0, ... , 0, 0, 0 , .... , 0)

(b)

x = (0.26, 0.25, -0.39, -0.07, 0.13, -0.17) (-0.43, -0.37, -0.12, 0.13, -0.11, 0.34) (-0.04, 0.50, 0.04, 0.44)

| | |
|---|---|
| **chair** | (-0.37, -0.23, 0.33, 0.38, -0.02, -0.37) |
| **on** | (-0.21, -0.11, -0.10, 0.07, 0.37, 0.15) |
| **dog** | (0.26, 0.25, -0.39, -0.07, 0.13, -0.17) |
| | ... |
| | ... |
| **the** | (-0.43, -0.37, -0.12, 0.13, -0.11, 0.34) |
| | ... |
| | ... |
| **mouth** | (-0.32, 0.43, -0.14, 0.50, -0.13, -0.42) |
| | ... |
| | ... |
| **gone** | (0.06, -0.21, -0.38, -0.28, -0.16, -0.44) |
| | ... |

Word Embeddings

| | |
|---|---|
| **NOUN** | (0.16, 0.03, -0.17, -0.13) |
| **VERB** | (0.41, 0.08, 0.44, 0.02) |
| | ... |
| | ... |
| **DET** | (-0.04, 0.50, 0.04, 0.44) |
| **ADJ** | (-0.01, -0.35, -0.27, 0.20) |
| **PREP** | (-0.26, 0.28, -0.34, -0.02) |
| | ... |
| | ... |
| **ADV** | (0.02, -0.17, 0.46, -0.08) |
| | ... |

POS Embeddings

Sparse (a) and dense (b) feature representations for
'the_DET dog'
(part of speech tagging task)

Q: what are the dimensionalities of word and PoS embeddings here?

# Dense representations (embeddings)

## Main benefits of continuous features

▶ Dimensionality of representations usually low (50-300).

# Dense representations (embeddings)

## Main benefits of continuous features

▶ Dimensionality of representations usually low (50-300).

▶ Feature vectors are dense, not sparse (more computationally efficient).

# Dense representations (embeddings)

## Main benefits of continuous features

▶ Dimensionality of representations usually low (50-300).

▶ Feature vectors are dense, not sparse (more computationally efficient).

▶ Generalization power: similar entities get similar embeddings (hopefully).

  ▶ '*town*' vector is closer to '*city*' vector than to '*banana*' vector (semantics);

# Dense representations (embeddings)

## Main benefits of continuous features

▶ Dimensionality of representations usually low (50-300).

▶ Feature vectors are dense, not sparse (more computationally efficient).

▶ Generalization power: similar entities get similar embeddings (hopefully).
  ▶ '*town*' vector is closer to '*city*' vector than to '*banana*' vector (semantics);
  ▶ *NOUN* vector is closer to *ADJ* vector than to *VERB* vector (grammar);

# Dense representations (embeddings)

## Main benefits of continuous features

► Dimensionality of representations usually low (50-300).

► Feature vectors are dense, not sparse (more computationally efficient).

► Generalization power: similar entities get similar embeddings (hopefully).
  - ► '*town*' vector is closer to '*city*' vector than to '*banana*' vector (semantics);
  - ► *NOUN* vector is closer to *ADJ* vector than to *VERB* vector (grammar);
  - ► *iobj* vector is closer to *obj* vector than to *punct* vector (syntax).

# Dense representations (embeddings)

## Main benefits of continuous features

▶ Dimensionality of representations usually low (50-300).

▶ Feature vectors are dense, not sparse (more computationally efficient).

▶ Generalization power: similar entities get similar embeddings (hopefully).
  - ▶ '*town*' vector is closer to '*city*' vector than to '*banana*' vector (semantics);
  - ▶ *NOUN* vector is closer to *ADJ* vector than to *VERB* vector (grammar);
  - ▶ *iobj* vector is closer to *obj* vector than to *punct* vector (syntax).

▶ Same features in different positions can share statistical strength:
  - ▶ A token 2 words to the right and a token 2 words to the left can be one and the same word. Would be good for the model to use this knowledge.
  - ▶ Not important for the Obligatory 1, but can be critical for other tasks.

## Demo web service

### Word vectors for English and Norwegian online

You can try the *WebVectors* service developed by our Language Technology group

# Demo web service

## Word vectors for English and Norwegian online

You can try the *WebVectors* service developed by our Language Technology group

`http://vectors.nlpl.eu/explore/embeddings/`



What words are related to **«computer»** in model «English Wikipedia»?

Word frequency

■ High  ■ Medium  □ Low

1. computr NOUN 0.7758
2. Computer NOUN 0.7665
3. microcomputer NOUN 0.6608
4. computing NOUN 0.6534
5. software NOUN 0.6387
6. mainframe NOUN 0.6324
7. compufl NOUN 0.6314
8. supercomputer NOUN 0.6026
9. workstation NOUN 0.5905
10. programmer NOUN 0.5837

- We show only the associates of the same part of speech as your query. All associates can be found at the *Similar Words* tab.

`0.6`  Similarity threshold  □ Show tags

This word in other models
- British National Corpus
- Google News
- English Gigaword
- Norsk Aviskorpus

Show the raw vector of «computer» in model