# Word Embeddings Tutorial

HILA GONEN

PHD STUDENT AT YOAV GOLDBERG'S LAB

BAR ILAN UNIVERSITY

5/3/18

# Outline

◦ NLP Intro

◦ Word representations and word embeddings

◦ Word2vec models

◦ Visualizing word embeddings

◦ Word2vec in Hebrew

◦ Similarity

◦ Analogies

◦ Evaluation

◦ A simple classification example

# NLP - Natural Language Processing

NLP is the field that includes **understanding** **processing** natural languages **analyzing** **generating**

We aim to create applicative models that perform as similar as possible to humans

# NLP

Applications in NLP:
- Translation
- Information Extraction
- Summarization
- Parsing
- Question Answering
- Sentiment Analysis
- Text Classification

And many more…

# NLP challenges

This field encounters numerous challenges:
◦ Polysemy
◦ Syntactic ambiguity
◦ Variability
◦ Co-reference resolution
◦ Lack of data / huge amounts of data

# NLP challenges - Polysemy

**Book**

Verb: **Book** a flight

Noun: He says it's a very good **book**

**Bank**

The edge of a river: He was strolling near the river **bank**

A financial institution: He works at the **bank**

**Solution**

An answer to a problem: Work out the **solution** in your head

From Chemistry: Heat the **solution** to 75° Celsius

# NLP challenges – Polysemy

## Kids make nutritious snacks

◦ Kids, when cooked well, can make nutritious snacks

## Kids make nutritious snacks

◦ Kids know how to prepare nutritious snacks

# NLP challenges – Syntactic Ambiguity

12 on their way to cruise among dead in plane crash

12 on their way to cruise among dead in plane crash

same words – different meanings

# NLP challenges – Syntactic Ambiguity

The cotton clothing is usually made of grows in Mississippi

The cotton clothing is usually made of grows in Mississippi

same words – different meanings

# NLP challenges – Syntactic Ambiguity

Fat people eat accumulates

Fat people eat accumulates

same words – different meanings

# NLP challenges – Variability

They allowed him to…

They let him …

He was allowed to…

He was permitted to…

Different words – same meaning

# NLP challenges – Co-Reference Resolution

Rachel had to wait for Dan because he said he wanted her advice.

This is a simple case...

There are more complex ones.

Dan called Bob to tell him about his surprising experience last week: "you won't believe it, I myself could not believe it".

# NLP challenges – Data-related issues

**A lot of data**

In some cases, we deal with huge amounts of data

Need to come up with models that can process a lot of data efficiently

**Lack of data**

Many problems in NLP suffer from lack of data:
◦ Non-standard platforms (code-switching)
◦ Expensive annotation (word-sense disambiguation, named-entity recognition)

Need to use methods to overcome this challenge (semi-supervised learning, multi-task learning…)

# Representation

We can represent objects in different hierarchy levels:
- Documents
- Sentences
- Phrases
- Words

We want the representation to be interpretable and easy-to-use

**Vector representation meets those requirements**

We will focus on word representation

# The Distributional Hypothesis

The Distributional Hypothesis:
- ◦ words that occur in the same contexts tend to have similar meanings
  (Harris, 1954)
- ◦ "You shall know a word by the company it keeps" (Firth, 1957)

Examples:
- ◦ Cucumber, sauce, pizza, ketchup
- ◦ Soundtrack, lyrics, sang, duet

tomato

song

# Vector Representation

We can define a word by a vector of counts over contexts, For Example:

|  | song | cucumber | meal | black |
|---|---|---|---|---|
| **tomato** | 0 | 6 | 5 | 0 |
| book | 2 | 0 | 2 | 3 |
| **pizza** | 0 | 2 | 4 | 1 |

◦ Each word is associated with a vector of dimension |V| (the size of the vocabulary)
◦ We expect similar words to have similar vectors
◦ Given the vectors of two words, we can determine their similarity (more about that later)

We can use different granularities of contexts: documents, sentences, phrases, n-grams

# Vector Representation

Raw counts are problematic:
◦ frequent words will characterize most words -> not informative

Except from raw counts, we can use other functions:
◦ TF-IDF (for term (t) – document (d)):

$$TF - IDF(t, d) = \frac{count(t, d)}{|d|} \cdot \log \frac{|D|}{|\{d \in D : t \in d\}|} \qquad D - \text{set of all documents}$$

◦ Pointwise Mutual Information:

$$PMI = \log \frac{p(x, y)}{p(x)p(y)}$$

# From Sparse to Dense

These vectors are:
- huge – each of dimension |V| (the size of the vocabulary ~ $100K +$)
- sparse – most entries will be 0

We want our vectors to be small and dense, two options:

1. Use a reduction algorithm such as SVD over a matrix of sparse vectors

2. Learn low-dimensional word vectors directly -

   usually referred as "word embeddings"

We will focus on the second option

# Word Embeddings

Each word in the vocabulary is represented by a low dimensional vector ($\sim 300d$)

All words are embedded into the same space

Similar words have similar vectors

(= their vectors are close to each other in the vector space)

Word embeddings are successfully used for various NLP applications

# Uses of word embeddings

Word embeddings are successfully used for various NLP applications (usually simply for initialization)

- ◦ Semantic similarity
- ◦ Word sense Disambiguation
- ◦ Semantic Role Labeling
- ◦ Named entity Recognition
- ◦ Summarization
- ◦ Question Answering
- ◦ Textual Entailment
- ◦ Coreference Resolution
- ◦ Sentiment analysis
- ◦ etc.

# Word2Vec

Models for efficiently creating word embeddings

Remember: our assumption is that similar words appear with similar context

Intuition: two words that share similar contexts are associated with vectors that are close to each other in the vector space

Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean, 2013. *Efficient estimation of word representations in vector space*. arXiv preprint arXiv:1301.3781.

Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean, 2013. *Distributed representations of words and phrases and their compositionality*. In Advances in neural information processing systems.
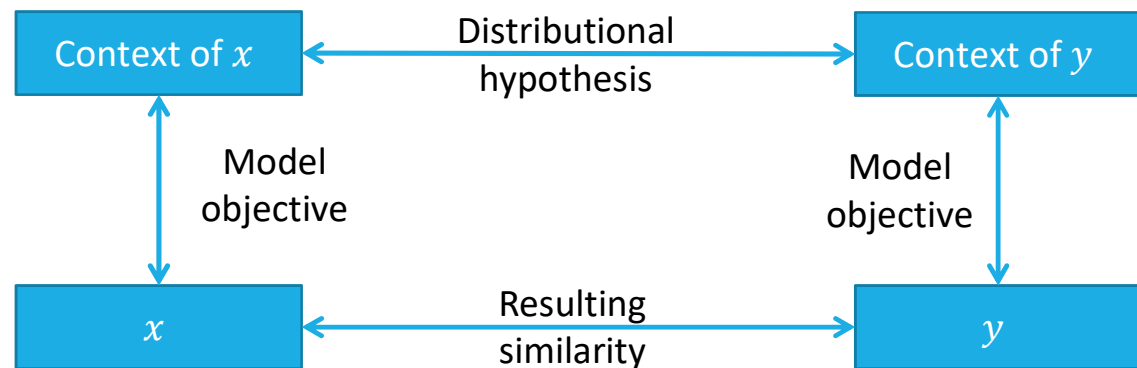
# Word2Vec

Models for efficiently creating word embeddings

Remember: our assumption is that similar words appear with similar context

Intuition: two words that share similar contexts are associated with vectors that are close to each other in the vector space

Let $x$ and $y$ be similar words

| Context of $x$ | ←  Distributional hypothesis  → | Context of $y$ |
| :---: | :---: | :---: |
| ↕ Model objective | | ↕ Model objective |
| $x$ | ←  Resulting similarity  → | $y$ |

# Word2Vec

The input: one-hot vectors
- bananas: $(1,0,0,0)$
- monkey: $(0,1,0,0)$
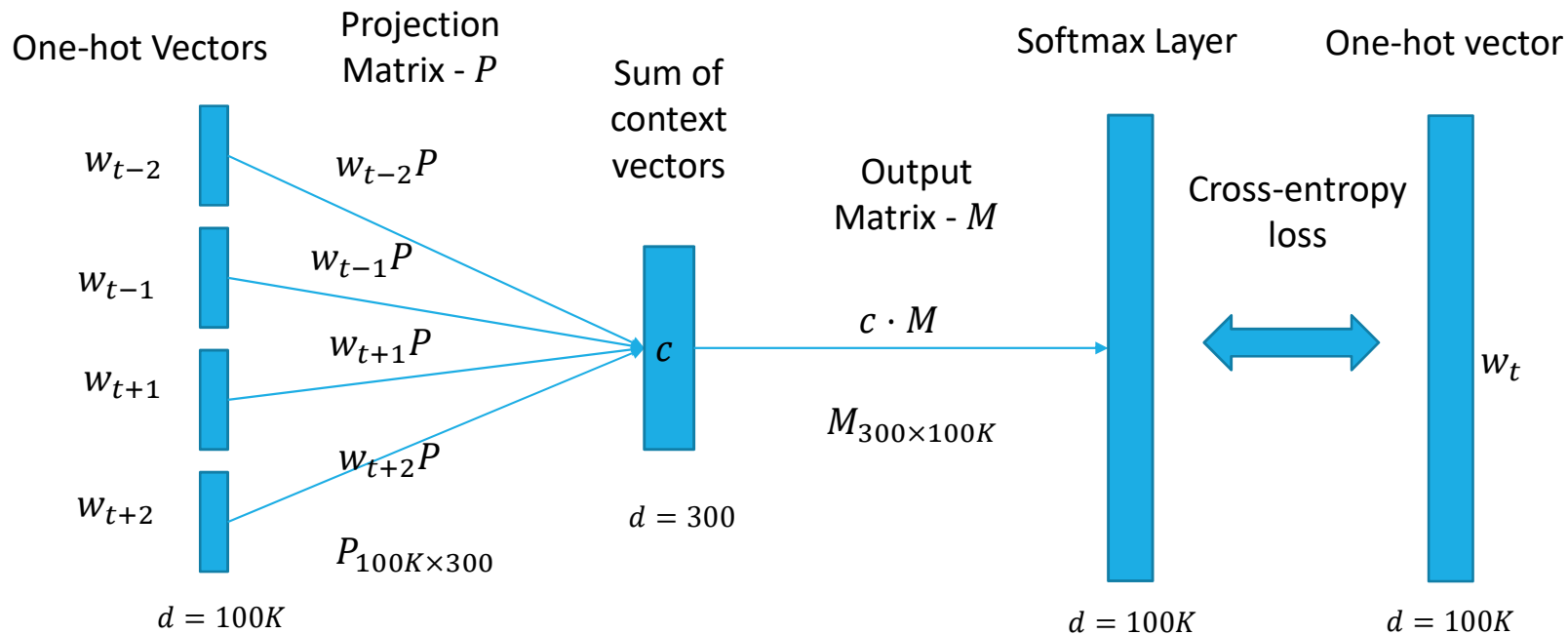- likes: $(0,0,1,0)$
- every: $(0,0,0,1)$

vocabulary size
$|V| = 4$

We are going to look at pairs of neighboring words:

Every monkey likes bananas $\longrightarrow$

$(every, monkey)$

$(likes, monkey)$

$(bananas, monkey)$

# CBOW – high level

Goal: Predict the middle word given the words of the context

One-hot Vectors

Projection Matrix - $P$

Sum of context vectors

Softmax Layer

One-hot vector

$w_{t-2}$

$w_{t-2}P$

$w_{t-1}$

$w_{t-1}P$

$w_{t+1}$

$w_{t+1}P$

Output Matrix - $M$

$c$

$c \cdot M$

Cross-entropy loss

$w_t$

$w_{t+2}$

$w_{t+2}P$

$M_{300 \times 100K}$

$P_{100K \times 300}$

$d = 300$

$d = 100K$

$d = 100K$

$d = 100K$

# Skip-gram – high level

Goal: Predict the context words given the middle word

The resulting projection matrix $P$ is the embedding matrix



One-hot Vector

Projection Matrix - $P$

Representation of $w_t$

Output Matrix - $M$

Softmax Layer

One-hot vectors

$w_t$

$w_t \cdot P$

$x$

$x \cdot M$

$x \cdot M$

$x \cdot M$

$x \cdot M$

$w_{t-2}$

$w_{t-1}$

$w_{t+1}$

$w_{t+2}$

$P_{100K \times 300}$

$d = 300$

$M_{300 \times 100K}$

$d = 100K$

$d = 100K$

Cross-entropy loss

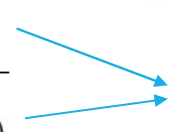$d = 100K$

# Skip-gram – details

Vector representations will be useful for predicting the surrounding words.

Formally:

Given a sequence of training words $w_1, w_2, \dots w_T$ , the objective of the Skip-gram model is to maximize the average log probability:

$$\frac{1}{T} \sum_{t=1}^{T} \sum_{-c \leq j \leq c, j \neq 0} \log p(w_{t+j}|w_t)$$

The basic Skip-gram formulation defines $p(w_{t+j}|w_t)$ using the softmax function:

$$p(w_{t+j}|w_t) = \frac{\exp(v'_{w_{t+j}} v_{w_t})}{\sum_{i=1}^{T} \exp(v'_{w_i} v_{w_t})}$$

$v$ - input vector representations
$v'$ - output vector representations

# Negative Sampling

Recall that for Skip-gram we want to maximize the average log probability:

$$\frac{1}{T}\sum_{t=1}^{T}\sum_{-c\leq j\leq c, j\neq 0}\log p(w_{t+j}|w_t)$$

Which is equivalent to minimizing the **cross-entropy loss**:

$$L = -\frac{1}{T}\sum_{t=1}^{T}\sum_{-c\leq j\leq c, j\neq 0}\log p(w_{t+j}|w_t) = -\frac{1}{T}\sum_{t=1}^{T}\sum_{-c\leq j\leq c, j\neq 0}\log \frac{\exp(v'_{w_{t+j}} v_{w_t})}{\sum_{i=1}^{T}\exp(v'_{w_i} v_{w_t})}$$

This is extremely computational-expensive, as we need to update all the parameters of the model for each training example…

# Negative Sampling

When looking at the loss obtained from a single training example, we get:

$$-\log p(w_{t+j}|w_t) = -\log \frac{\exp(v'_{w_{t+j}} v_{w_t})}{\sum_{i=1}^{T} \exp(v'_{w_i} v_{w_t})} = -v'_{w_{t+j}} v_{w_t} + \log \sum_{i=1}^{T} \exp(v'_{w_i} v_{w_t})$$

"positive" pair          "negative" pair

When using negative sampling, instead of going through all the words in the vocabulary for negative pairs, we sample a modest amount of $k$ words (around 5-20). The exact objective used:

$$\log \sigma(v'_{w_{t+j}} v_{w_t}) + \sum_{1=1}^{k} \log \sigma(-v'_{w_i} v_{w_t})$$

Replaces the term: $\log p(w_{t+j}|w_t)$ for each word in the training

$$\sigma(x) = \frac{1}{1 + \exp(-x)}$$

# Context Sampling

We want to give more weight to words closer to our target word

For a given window size C, we sample R in [1, C] and try to predict only R words before and after our target word

For each word in the training we need to perform 2*R word classifications

(R is not fixed)

# Subsampling of Frequent Words

In order to eliminate the negative effect of very frequent words such as "in", "the" etc. (that are usually not informative), a simple subsampling approach is used:

Each word $w_i$ in the training set is discarded with probability:   $P(w_i) = 1 - \sqrt{\dfrac{t}{f(w_i)}}$

This way frequent words are discarded more often

This method improves the training speed and makes the word representations significantly more accurate

# Phrases

These models are unable to represent phrases that are not compositions of the individual words

"New York"  != "New" + "York"

"Boston Globe" != "Boston" + "Globe"

The extension is simple:
- Find words that appear frequently together, and infrequently in other contexts
  - phrases are formed based on the unigram and bigram counts

$$\text{score}(w_i, w_j) = \frac{\text{count}(w_i w_j) - \delta}{\text{count}(w_i) \times \text{count}(w_j)}.$$

  - The bigrams with score above the chosen threshold are then used as phrases
  - "New York Times" will be replaced with a unique token, "this is" will remain unchanged
- Train word2vec as usual

# Use word2vec package

Using this package is extremely simple:
- ◦ Download the code from Mikolov's git repository:
- ◦ https://github.com/tmikolov/word2vec
- ◦ Compile the package
- ◦ Download the default corpus (wget http://mattmahoney.net/dc/text8.zip) or another corpus of your choice
- ◦ Train the model using the desired parameters

Jupyter: code for downloading, compiling, and training

# Importance of Parameters – window size

Word: **walk**

Window size = 3

| Word | Cosine distance |
| --- | --- |
| go | 0.488083 |
| snipe | 0.464912 |
| shoot | 0.456677 |
| fly | 0.449722 |
| sit | 0.449678 |
| pass | 0.442459 |
| climbs | 0.440931 |
| walked | 0.436502 |
| ride | 0.434034 |
| stumble | 0.426750 |
| bounce | 0.425577 |
| travelling | 0.419419 |
| walking | 0.412107 |
| walks | 0.410949 |
| trot | 0.410418 |
| leaping | 0.406744 |
| sneak | 0.401918 |
| climb | 0.399793 |
| move | 0.396715 |
| wait | 0.394463 |
| going | 0.391639 |
| shouted | 0.388382 |
| roam | 0.388073 |
| thrown | 0.384087 |
| get | 0.383894 |

Window size = 30

| Word | Cosine distance |
| --- | --- |
| walking | 0.486317 |
| walked | 0.430764 |
| walks | 0.406772 |
| stairs | 0.401518 |
| go | 0.399274 |
| sidewalk | 0.385786 |
| stand | 0.380480 |
| cortege | 0.371033 |
| wheelchair | 0.362877 |
| strapped | 0.360179 |
| hollywood | 0.356544 |
| carousel | 0.356187 |
| grabs | 0.356007 |
| swim | 0.355027 |
| breathe | 0.354314 |
| tripped | 0.352899 |
| cheer | 0.352477 |
| moving | 0.350943 |
| inductees | 0.347791 |
| walkway | 0.347164 |
| shout | 0.346229 |
| pounding | 0.340554 |
| blvd | 0.339121 |
| crowd | 0.338731 |
| levada | 0.334899 |

# Importance of Parameters – iterations

Word: **walk**

No. of iterations = 1

| Word | Cosine distance |
| --- | --- |
| walking | 0.851438 |
| walks | 0.846485 |
| bat | 0.843796 |
| ride | 0.830734 |
| crowd | 0.821692 |
| quiet | 0.812538 |
| spot | 0.802777 |
| steal | 0.787917 |
| door | 0.787571 |
| doors | 0.786485 |
| bed | 0.773686 |
| dinner | 0.772160 |
| shadow | 0.769573 |
| luck | 0.768221 |
| baby | 0.767862 |
| shoot | 0.765968 |
| walked | 0.765739 |
| sitting | 0.765394 |
| shirt | 0.759116 |
| rides | 0.759047 |
| watching | 0.755140 |
| watch | 0.750808 |
| gehrig | 0.741494 |
| shoots | 0.740971 |
| looking | 0.740904 |

No. of iterations = 100

| Word | Cosine distance |
| --- | --- |
| walked | 0.483473 |
| ride | 0.470925 |
| walks | 0.470889 |
| stand | 0.449993 |
| walking | 0.449071 |
| go | 0.430172 |
| shoot | 0.421110 |
| get | 0.404258 |
| move | 0.403757 |
| live | 0.403347 |
| fly | 0.400929 |
| climbs | 0.396346 |
| throw | 0.391768 |
| climb | 0.384038 |
| wiggle | 0.380892 |
| thrown | 0.380426 |
| pull | 0.375478 |
| goes | 0.375406 |
| moving | 0.374447 |
| pass | 0.372463 |
| conversing | 0.364413 |
| sit | 0.362765 |
| crowd | 0.361651 |
| kiss | 0.359883 |
| stay | 0.357015 |

# Importance of Parameters – dimensions

Word: **walk**

No. of dimensions = 5

No. of dimensions = 1000

| Word | Cosine distance |
|---|---|
| catcher | 0.998074 |
| shirt | 0.996589 |
| lechuck | 0.995313 |
| bullseye | 0.994644 |
| bowler | 0.994381 |
| punter | 0.993154 |
| lovell | 0.992815 |
| heels | 0.992255 |
| whip | 0.992085 |
| outfit | 0.992047 |
| tore | 0.991924 |
| steals | 0.991524 |
| guybrush | 0.991166 |
| gigs | 0.990291 |
| hanging | 0.990201 |
| burns | 0.990043 |
| backing | 0.989966 |
| orser | 0.989960 |
| torch | 0.989747 |
| beat | 0.989435 |
| showdown | 0.989381 |
| feat | 0.989242 |
| cheers | 0.988951 |
| clad | 0.988646 |
| lunch | 0.988326 |

| Word | Cosine distance |
|---|---|
| walks | 0.304954 |
| walked | 0.303322 |
| snipe | 0.287221 |
| walking | 0.272690 |
| ride | 0.266770 |
| canter | 0.251025 |
| bandleaders | 0.246454 |
| climbs | 0.233725 |
| catapulted | 0.230075 |
| climb | 0.229263 |
| trot | 0.228362 |
| shouted | 0.227306 |
| stand | 0.223288 |
| seagulls | 0.221745 |
| fly | 0.216602 |
| fences | 0.216366 |
| lifts | 0.215402 |
| pray | 0.214977 |
| paws | 0.214865 |
| bounces | 0.214449 |
| shoot | 0.213457 |
| grabs | 0.212018 |
| walkway | 0.211136 |
| swim | 0.209120 |
| tumble | 0.207765 |

# How does the file usually look like

Word embedding files in readable format usually have a row for each word in the vocabulary

In each row, the specific word is followed by the values of the respected vector

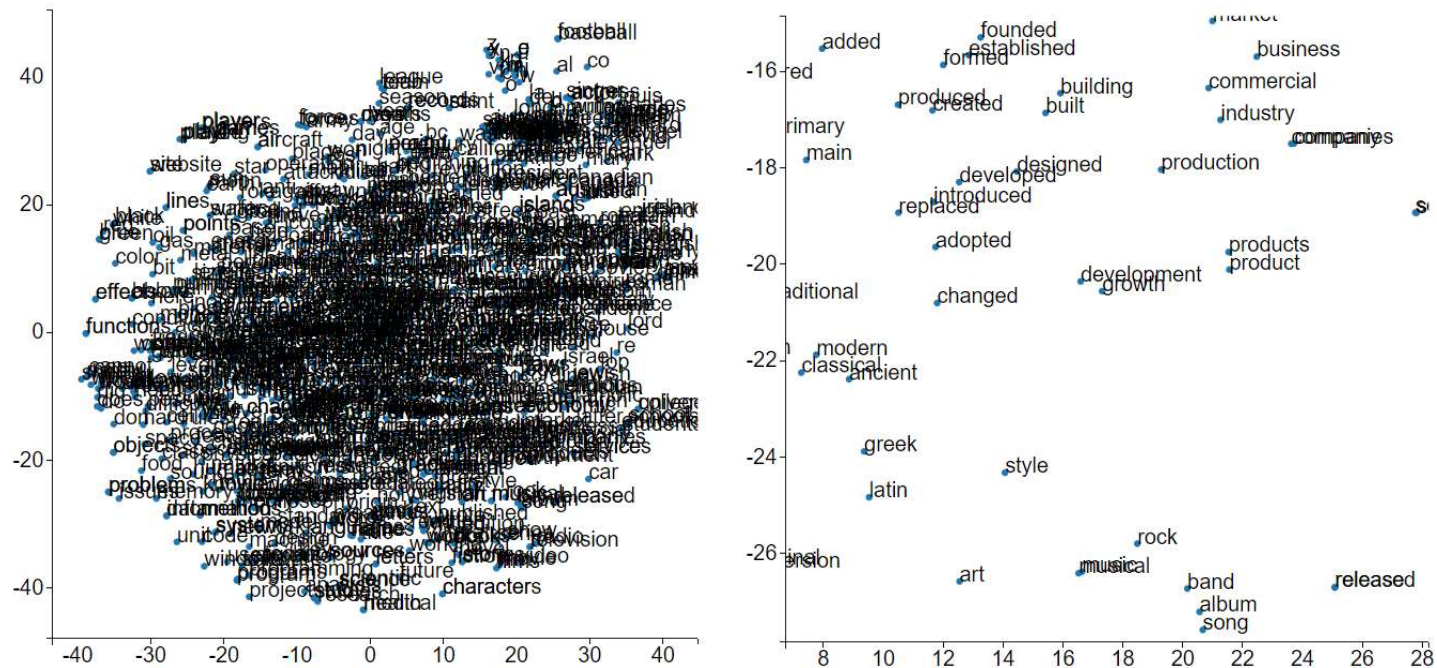Possibly some additional information in the first rows
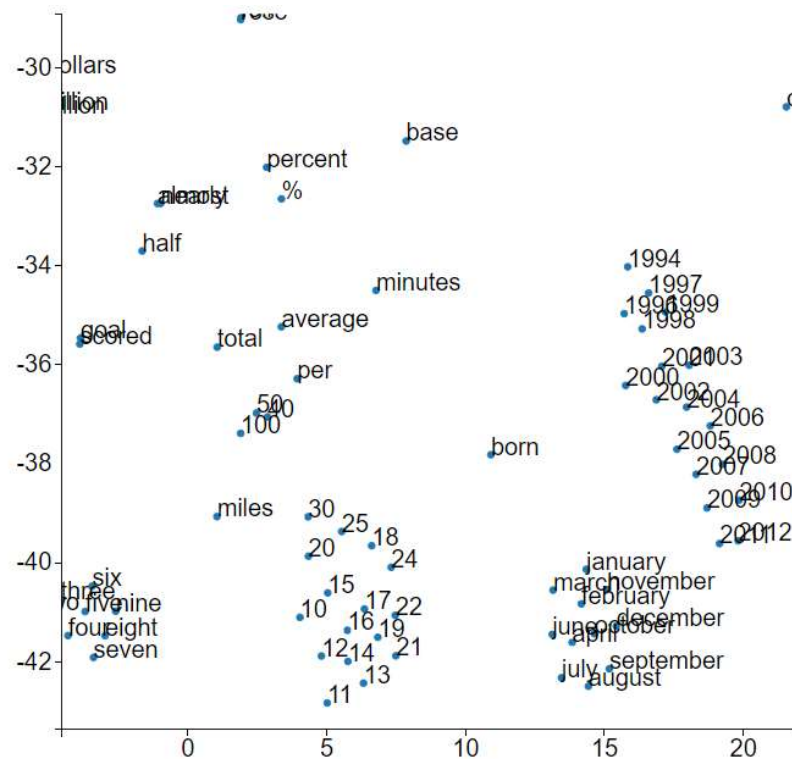
# Visualizing word embeddings

Using tSNE

◦ *Visualizing Data using t-SNE*, Maaten and Hinton, 2008

Jupyter: Loading and Visualizing word vectors

# Visualizing word embeddings– Word2Vec

# Visualizing word embeddings - Glove

# Wevi: word embedding visual inspector

A tool that visualizes the basic working mechanism of word2vec

https://ronxin.github.io/wevi/

# We...ctor

A too...d2vec

https:...

# Word2Vec - Hebrew

Link: http://u.cs.biu.ac.il/~yogo/tw2v/similar/ (By Ron Shemesh and Yoav Goldberg)

Based on tweets in Hebrew

# Word Embeddings

Link: http .../tw2... nesh

Based on

| שכיחות | דימיון מילה | | שכיחות | דימיון מילה | | שכיחות | דימיון מילה |
|---|---|---|---|---|---|---|---|
| 16065 | 😀 100.0 | | 327100 | הולך 100.0 | | 20864 | שפה 100.0 |
| 14738 | 🙂 94.7 | | 639232 | בא 75.9 | | 125 | הגיה 78.1 |
| 20020 | ☺ 93.5 | | 680624 | יכול 75.3 | | 255 | הגייה 74.3 |
| 16020 | 😄 92.7 | | 155222 | מתחיל 76.3 | | 288 | אסכולה 72.9 |
| 62507 | 😊 90.2 | | 166958 | הולכת 75.3 | | 292 | אינטונציה 71.7 |
| 68331 | 😌 87.1 | | 87292 | אמר 75.3 | | 44 | דיסציפלינה 71.2 |
| 5438 | 😋 86.6 | | 944253 | צריך 71.0 | | 2287 | לוגיקה 70.2 |
| 3556 | 🙂 85.8 | | 199438 | חייב 73.7 | | 261 | משפה 70.6 |
| 4577 | 😇 82.7 | | 1120959 | רוצה 68.9 | | 3938 | תיאוריה 68.7 |
| 44223 | 😛 81.5 | | 105244 | מוכן 73.1 | | 496 | תפישה 69.3 |
| 30838 | 😁 80.8 | | 29930 | מפסיק 72.8 | | 391 | מיתולוגיה 69.2 |
| 19796 | ☺ 79.9 | | 14106 | מחליט 72.7 | | 1099 | גנטיקה 68.9 |
| 24873 | 😎 79.6 | | 31218 | אלך 72.2 | | 187 | ורסיה 69.1 |
| 2327 | 😊😊 78.8 | | 13029 | יכל 72.2 | | 450 | כשפה 69.0 |
| 13376 | 😇 76.5 | | 5077 | ממהר 72.0 | | 179 | מתודולוגיה 68.9 |
| 1236 | 😻 76.7 | | 28144 | שהולך 71.9 | | 3684 | פלטפורמה 67.6 |
| 2197 | 😲 76.6 | | 60953 | ממשיך 70.6 | | 589 | התחכמות 68.2 |
| 3758 | 😊 76.0 | | 83324 | עומד 69.4 | | 9034 | הגדרה 66.2 |
| 749 | 😀😀 75.4 | | 37413 | מסוגל 69.2 | | 11916 | שיטה 65.3 |
| 2258 | ^_^ 75.3 | | 78116 | מצליח 68.9 | | 53 | סכמה 68.1 |
| 13669 | 💜 74.5 | | 11052 | יצטרך 69.1 | | 1579 | פרקטיקה 67.7 |
| | | | 17427 | יתחיל 69.0 | | 582 | טרמינולוגיה 67.9 |
| | | | 177308 | מנסה 67.9 | | 4492 | תפיסה 66.9 |
| | | | 47496 | רצה 68.1 | | 50 | עגה 67.9 |

# Word2Vec - Hebrew

Ok... Nice.

But:

What about יעד vs. מטרה which are synonyms?

Noun genders dramatically affect results –

We do not want that, or at least not for arbitrary gender

# Word2Vec

Ok… Nice.

But:

What about יעד vs. ... onym

Noun genders dran ... ; –

We do not want th ... arbi

| | חפש | יעד |
|---|---|---|

| שכיחות | דימיון מילה | |
|---|---|---|
| 8136 | יעד | 100.0 |
| 535 | כיעד | 70.6 |
| 1476 | קריטריון | 65.2 |
| 6366 | נכס | 63.6 |
| 10784 | מודל | 61.6 |
| 158 | ויעד | 61.6 |
| 3626 | תרחיש | 60.5 |
| 11644 | היעד | 60.0 |
| 17600 | הישג | 59.2 |
| 46 | כסיכון | 59.4 |
| 2528 | סטנדרט | 58.7 |
| 38259 | פתרון | 57.1 |
| 1470 | תמריץ | 56.8 |
| 278883 | מקום | 51.3 |
| 37 | קרטריון | 56.2 |
| 95 | ובינלאומי | 56.0 |
| 13 | ואקסקלוסיבי | 56.0 |
| 237 | ייתרון | 55.8 |
| 23446 | שטח | 55.3 |
| 5116 | רף | 55.1 |
| 12953 | יתרון | 54.6 |
| 24049 | מהלך | 54.1 |
| 11696 | פוטנציאל | 54.2 |
| 31019 | כיוון | 53.9 |

| | חפש | מטרה |
|---|---|---|

| שכיחות | דימיון מילה | |
|---|---|---|
| 19358 | מטרה | 100.0 |
| 2658 | שאיפה | 74.5 |
| 9828 | משימה | 71.9 |
| 61495 | הזדמנות | 69.1 |
| 105895 | סיבה | 66.2 |
| 8805 | אלטרנטיבה | 70.5 |
| 1613 | נוסחה | 70.7 |
| 11916 | שיטה | 70.0 |
| 237 | נוסחא | 70.5 |
| 4714 | אסטרטגיה | 70.0 |
| 7000 | אג'נדה | 67.9 |
| 681 | יישות | 67.4 |
| 496 | תפישה | 67.4 |
| 3684 | פלטפורמה | 67.0 |
| 39317 | נקודה | 65.1 |
| 19288 | גישה | 65.6 |
| 196 | דוקטרינה | 66.4 |
| 44 | אסטרטגייה | 66.0 |
| 21334 | בחירה | 64.9 |
| 4492 | תפיסה | 65.6 |
| 3168 | חלופה | 65.6 |
| 9034 | הגדרה | 65.3 |
| 39166 | אפשרות | 62.9 |
| 1900 | ישות | 64.6 |

# Word2Vec - Hebrew

Another example:

בצרה

בצל

Prefixes and suffixes are not always handled correctly

Also, not always clear what the wanted behavior is

# Word2V...ew

Another examp...

בצרה

בצל

Prefixes and su...ys handled...

Also, not alway...anted beh...

חפש | בצל

| דימיון מילה | שכיחות |
|---|---|
| בצל 100.0 | 23556 |
| ובצל 79.8 | 1574 |
| וסלרי 63.3 | 103 |
| מרק 56.8 | 27207 |
| ובזיליקום 61.7 | 314 |
| ופטרוזיליה 61.5 | 286 |
| מגורד 61.4 | 183 |
| ברוטב 59.8 | 5301 |
| וזוקיני 60.7 | 45 |
| בחרדל 60.2 | 127 |
| וכוסברה 59.9 | 335 |
| וזעתר 60.0 | 185 |
| סלק 59.1 | 3523 |
| רוטב 58.2 | 7020 |
| ולימון 59.1 | 1723 |
| קונפי 58.9 | 298 |
| וברוקולי 58.9 | 203 |
| וסלמון 58.6 | 196 |
| ואורגנו 58.6 | 77 |
| ואספרגוס 58.6 | 88 |
| ופרמזן 58.3 | 147 |
| כתוש 58.1 | 270 |
| וסלק 58.1 | 235 |
| פורטבלו 58.0 | 135 |

חפש | בצרה

| דימיון מילה | שכיחות |
|---|---|
| בצרה 100.0 | 1624 |
| במצוקה 60.6 | 4622 |
| בבעיה 50.2 | 8617 |
| בסביבה 44.5 | 13472 |
| בצרות 47.2 | 3026 |
| בהכרה 46.2 | 3730 |
| בסכנה 42.7 | 9296 |
| ומרעננה 46.1 | 85 |
| בקרבתם 46.0 | 181 |
| באחווה 45.8 | 184 |
| בקירבה 45.5 | 97 |
| ונטועה 45.5 | 16 |
| משמאלם 45.5 | 12 |
| בחובות 44.8 | 1742 |
| בטובתם 45.1 | 96 |
| במיעוט 44.1 | 1852 |
| בגלות 44.2 | 1610 |
| מאחוריך 43.0 | 4082 |
| בבועה 43.1 | 2503 |
| בבוטניקה 43.7 | 27 |
| בפרנויה 43.5 | 150 |
| מאופקים 43.2 | 775 |
| ובשביה 43.5 | 53 |
| ובריאותו 43.3 | 61 |

# Word2Vec – Hebrew

Word embeddings include inherent biases, as a result of biases in the corpus (not only in Hebrew…)

For example, רופא vs. רופאה

# Word2Ve...w

Word embeddings ... biases, a ... s in the corpus
(not only in Hebre...
For example, רופא ...



| | | רופאה | חפש |
|---|---|---|---|
| שכיחות | דימיון מילה | | |
| 4307 | רופאה | 100.0 | |
| 3064 | הרופאה | 74.4 | |
| 2407 | מטופלת | 73.2 | |
| 64 | פציינטית | 72.3 | |
| 457 | וטרינרית | 72.1 | |
| 1242 | קוסמטיקאית | 71.9 | |
| 80 | עו"סית | 71.6 | |
| 453 | כאחות | 71.3 | |
| 2470 | פקידה | 70.9 | |
| 3860 | תלמידה | 69.5 | |
| 440 | מיילדת | 69.5 | |
| 1748 | מדריכה | 69.1 | |
| 304 | רוקחת | 69.1 | |
| 530 | מדענית | 68.6 | |
| 43 | קרדיולוגית | 68.3 | |
| 2529 | יולדת | 67.9 | |
| 1237 | שיננית | 67.5 | |
| 5575 | סטודנטית | 67.4 | |
| 4402 | גננת | 67.2 | |
| 749 | פיליפינית | 67.0 | |
| 329 | הוטרינרית | 66.7 | |
| 25801 | אחות | 66.0 | |
| 168 | פרופסורית | 66.5 | |
| 182017 | אישה | 61.3 | |

| | | רופא | חפש |
|---|---|---|---|
| שכיחות | דימיון מילה | | |
| 31223 | רופא | 100.0 | |
| 368 | אורתופד | 85.1 | |
| 511 | כירורג | 82.0 | |
| 3270 | פסיכיאטר | 79.4 | |
| 207 | נוירולוג | 78.0 | |
| 17732 | הרופא | 74.6 | |
| 1772 | וטרינר | 76.3 | |
| 218 | אורולוג | 76.4 | |
| 5093 | פסיכולוג | 73.1 | |
| 298 | קרדיולוג | 72.6 | |
| 454 | אורטופד | 69.8 | |
| 12820 | לרופא | 66.1 | |
| 6271 | מטופל | 66.7 | |
| 175 | כירופרקט | 66.7 | |
| 50 | סקסולוג | 66.6 | |
| 47 | פנימאי | 66.3 | |
| 106 | סטאז'ר | 66.0 | |
| 1140 | גניקולוג | 65.3 | |
| 309 | פתולוג | 64.9 | |
| 52 | ג.א.א. | 64.4 | |
| 980 | הוטרינר | 64.3 | |
| 4307 | רופאה | 63.7 | |
| 450 | ורופא | 64.0 | |
| 1630 | פסיכיאטרי | 63.4 | |

# Other pre-trained word embeddings

Glove (Pennington et al.):
◦ Based on ratios of co-occurrence probabilities
◦ https://nlp.stanford.edu/projects/glove/


Fast-text (Bojanowski et al.):
◦ Each word is represented as a bag of character n-grams. A vector representation is associated to each character n-gram, and words are represented as the sum of these representations
◦ https://fasttext.cc/

# Similarity

In order to evaluate word embeddings on similarity tasks, we first need to define "similarity"

There are many different ways to define "similarity" and "correlation" between words...

◦ walk – walking, walk – run, walk – stroll

◦ Germany – Berlin, Germany – England

◦ dog – cat, dog – Labrador, dog – leash

This is still an open issue...

# Similarity measure

The distance between two vectors is not a good measure

◦ We do not want to take the length of the vector into account

Most popular similarity measure is Cosine similarity:

◦ The similarity between two vectors $v$ and $w$ is: $\dfrac{v \cdot w}{||v|| \, ||w||}$

# Similarity

Jupiter: find topK

# Similarity

You can't always get what you want…

walk - Top5 similar words:

walked, walks, walking, climbs, ride

book – top5 similar words:

books, chapter, novel, abridged, autobiography

# Analogies

These models are capable of learning linguistic regularities

For example,

vector("king") - vector("man") + vector("woman") $\cong$ vector("queen")

vector("mice") - vector("mouse") + vector("door") $\cong$ vector("doors")

Jupyter: Analogies

# Analogies

How does it work?

Given the analogy $a : a^*, b : b^*$ , where word $b^*$ is to be found, we try to maximize the following objective:

$$\arg\max_{b^* \in V}(sim(b^*, b - a + a^*)) \longrightarrow \arg\max_{b^* \in V}(cos(b^*, b - a + a^*))$$

When vectors are normalized, this is equivalent to:

$$\arg\max_{b^* \in V}(cos(b^*, b) - cos(b^*, a) + cos(b^*, a^*))$$

We actually search for a word that is similar to *b*, and *a\**, but different from *a*

*Linguistic Regularities in Sparse and Explicit Word Representations, Levy and Goldberg, 2014*

# Analogies

This does not always work that well…

Naturally, depends on the corpus and the hyper-parameters

Additionally, in some cases, specific aspect of relations might dominate others:

**London : England , Baghdad : ?** ⟶   Mosul (instead of Iraq)

Here, even though Iraq is more similar to England than Mosul, the similarity of Mosul to Baghdad dominates, making Mosul the best candidate

A possible solution – use multiplication instead of summation (equivalent to using *log* values):

$$\arg\max_{b^*\in V} \frac{cos(b^*,b)cos(b^*,a)}{cos(b^*,a^*)+\epsilon}$$

*Linguistic Regularities in Sparse and Explicit Word Representations, Levy and Goldberg, 2014*

# Evaluation

Intrinsic Evaluation:

1. Syntactic and semantic analogies:
   ◦ Athens : Greece ; Oslo : ?
   ◦ think : thinking ; read : ?
   ◦ mouse : mice ; door : ?

2. Word correlation benchmarks with human scores (wordsim353, simLex999):

| word1 | word2 | Human score |
|-------|-------|-------------|
| train | car | 6.31 |
| drink | ear | 1.31 |
| gem | jewel | 8.96 |

Extrinsic evaluation:
   ◦ Show Improvement on downstream tasks when using word embeddings

# Classification with word embedding

An optional use of word embedding is a simple classification:

Say we have a short list of professions, and we want to elaborate it

We can run a simple classification model with sklearn
- Use the short list as positive examples
- Add random negative example
- Learn a classification model
- Predict True/False for new words from the vocabulary

Jupyter: Classification example

# Biases in word embeddings

A very nice tutorial about word embedding biases:

How to make a racist AI without really trying

https://gist.github.com/rspeer/ef750e7e407e04894cb3b78a82d66aed

# Conclusion

Word embedding:
◦ A powerful word representation
◦ Easy to incorporate into different models

Can capture word similarities and linguistic regularities

Existing models have their limitations

Need to custom training parameters according to the desired properties and similarities

# Questions?

| Word | Cosine distance |
|---|---|
| question | 0.744367 |
| answers | 0.577580 |
| doubts | 0.535749 |
| answer | 0.521965 |
| concerns | 0.492096 |
| issues | 0.487333 |
| unanswered | 0.457177 |
| discussions | 0.454697 |
| answering | 0.448383 |
| matters | 0.446861 |
| debates | 0.439423 |
| statements | 0.433208 |
| objections | 0.429943 |
| issue | 0.424728 |
| debate | 0.421695 |
| conclusions | 0.421594 |
| inquiries | 0.419190 |
| problem | 0.419025 |
| answered | 0.414844 |
| responses | 0.412600 |
| arguments | 0.410255 |
| problems | 0.406869 |

# Thank you!