

Deep Learning NLP Introduction

M. Vazirgiannis

<https://bit.ly/2rwmvQU>

LIX, Ecole Polytechnique

March 2022

OUTLINE

- **Representation Learning for Text**
 - Latent Semantic Indexing (LSI)
 - Word2Vec
- CNN for text classification

Language model

- Goal: determine $P(s = w_1 \dots w_k)$ in some domain of interest

$$P(s) = \prod_{i=1}^k P(w_i | w_1 \dots w_{i-1})$$

e.g., $P(w_1 w_2 w_3) = P(w_1) P(w_2 | w_1) P(w_3 | w_1 w_2)$

- Traditional n-gram language model assumption:
“probability of a word depends only on **context** of $n - 1$ previous words”

$$\Rightarrow \hat{P}(s) = \prod_{i=1}^k P(w_i | w_{i-n+1} \dots w_{i-1})$$

- i.e. “Paris is the capital of France located in Ile de”
- Typical ML-smoothing learning process (e.g., Katz 1987):
 1. compute $\hat{P}(w_i | w_{i-n+1} \dots w_{i-1}) = \frac{\#w_{i-n+1} \dots w_{i-1} w_i}{\#w_{i-n+1} \dots w_{i-1}}$ on training corpus
 2. smooth to avoid zero probabilities

Representing Words

➤ One-hot vector

- high dimensionality
- sparse vectors
- dimensions= $|V|$ ($10^6 < |V|$)
- unable to capture semantic similarity between words



<i>eat</i>							■						
<i>food</i>											■		
<i>news</i>		■											

➤ Distributional vector

- words that occur in similar contexts, tend to have similar meanings
- each word vector contains the frequencies of all its neighbors
- dimensions= $|V|$
- computational complexity for ML algorithms

<i>eat</i>				■			■			■			
<i>food</i>					■			■			■		■
<i>news</i>		■						■			■		

Representing Words

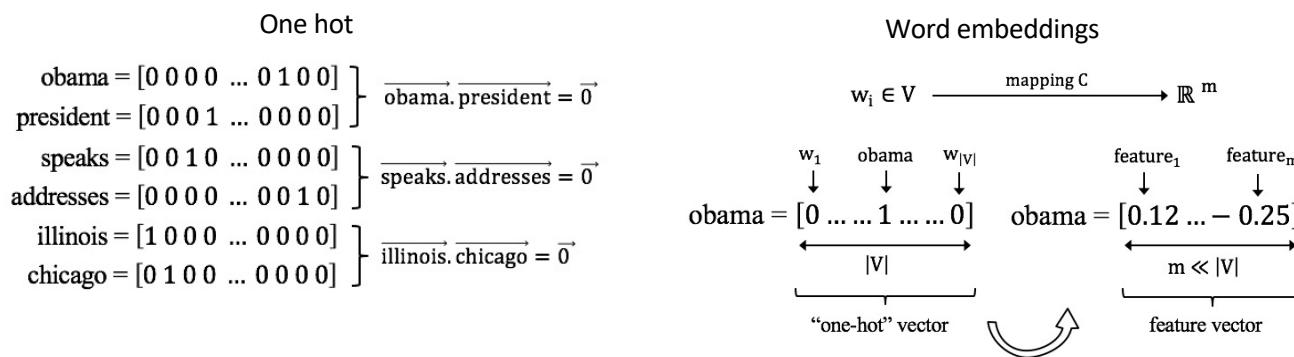
➤ Word embeddings

- store the same contextual information in a low-dimensional vector
- **densification** (sparse to dense)
- **compression**
 - dimensionality reduction
 - dimensions=m
 $100 < m < 500$
- able to capture semantic similarity between words
- learned vectors (unsupervised)
- Learning methods
 - **SVD**
 - **word2vec**
 - **GloVe**

<i>eat</i>									
<i>food</i>									
<i>news</i>									

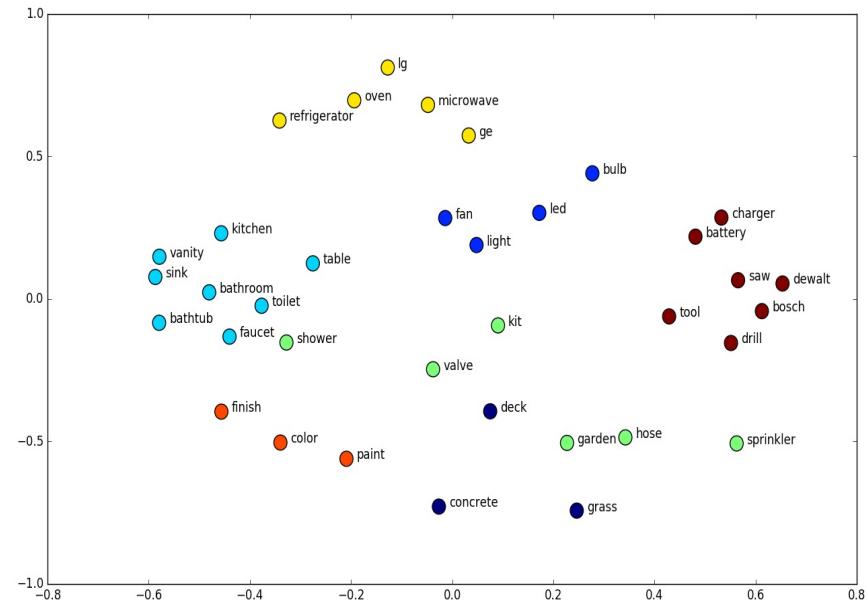
Text Similarity

- We should assign similar probabilities (discover similarity) to *Obama speaks to the media in Illinois* and the *President addresses the press in Chicago*
- This does not happen because of the “one-hot” vector space representation



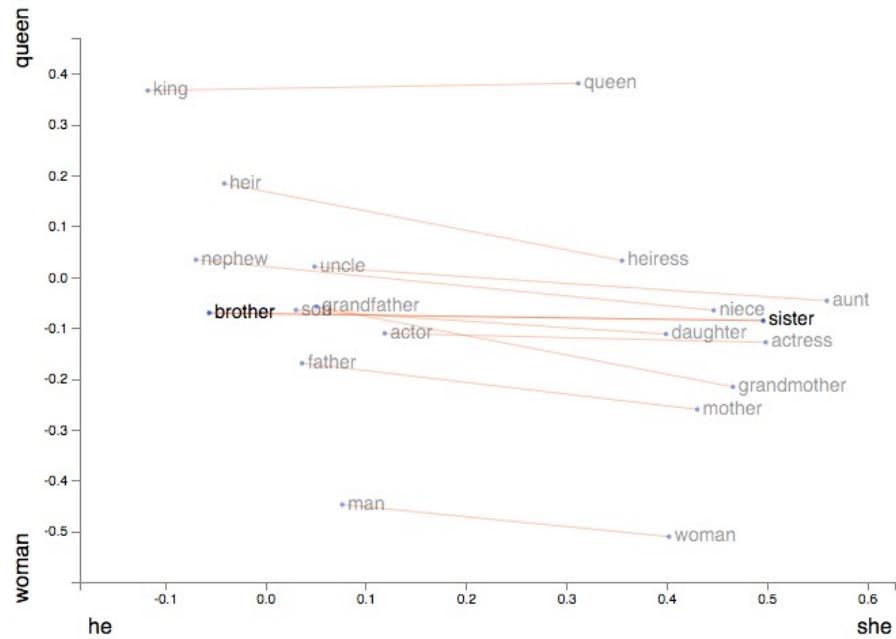
Representation Learning for Text

- “**a word is defined by “the company it keeps” (Firth, 1957)**
- Word embeddings are a class of algorithms where each word is represented as real-valued vector.
- The learning process of these vectors is either joint with a neural network model on some task or is an unsupervised process.
- Similar words in meaning have similar representation.



Representation Learning for Text

- Words with similar meaning end up close to each other
- Words sharing similar contexts may be analogous
 - Synonyms
 - Antonyms
 - Names
 - Colors
 - Places
 - Interchangeable words
- Vector arithmetics to work with analogies
- i.e. **king - man + woman = queen**



<https://lamiowce.github.io/word2viz/>

OUTLINE

- Representation Learning for Text
 - **Latent Semantic Indexing (LSI)**
 - Word2Vec
- CNN for text classification

SVD word embeddings

- Dimensionality reduction on co-occurrence matrix
- Create a $|V| \times |V|$ word co-occurrence matrix X
- Apply SVD $X = USV^T$
- Take first k columns of U
- Use the k -dimensional vectors as representations for each word
- Able to capture semantic and syntactic similarity

LSI – an example

LSI application on a term – document matrix

- C1: Human machine Interface for Lab ABC computer application
 - C2: A survey of user opinion of computer system response time
 - C3: The EPS user interface management system
 - C4: System and human system engineering testing of EPS
 - C5: Relation of user-perceived response time to error measurements
 - M1: The generation of random, binary unordered trees
 - M2: The intersection graph of path in trees
 - M3: Graph minors IV: Widths of trees and well-quasi-ordering
 - M4: Graph minors: A survey
- The dataset consists of 2 classes, 1st: “human – computer interaction” (c1-c5) 2nd: related to graph (m1-m4). After feature extraction the titles are represented as follows.

LSI – an example

$$A = U L V^T$$

A =

LSI – an example

$U =$

$$A = U \Lambda V^T$$

0.22	-0.11	0.29	-0.41	-0.11	-0.34	0.52	-0.06	-0.41
0.20	-0.07	0.14	-0.55	0.28	0.50	-0.07	-0.01	-0.11
0.24	0.04	-0.16	-0.59	-0.11	-0.25	-0.30	0.06	0.49
0.40	0.06	-0.34	0.10	0.33	0.38	0.00	0.00	0.01
0.64	-0.17	0.36	0.33	-0.16	-0.21	-0.17	0.03	0.27
0.27	0.11	-0.43	0.07	0.08	-0.17	0.28	-0.02	-0.05
0.27	0.11	-0.43	0.07	0.08	-0.17	0.28	-0.02	-0.05
0.30	-0.14	0.33	0.19	0.11	0.27	0.03	-0.02	-0.17
0.21	0.27	-0.18	-0.03	-0.54	0.08	-0.47	-0.04	-0.58
0.01	0.49	0.23	0.03	0.59	-0.39	-0.29	0.25	-0.23
0.04	0.62	0.22	0.00	-0.07	0.11	0.16	-0.68	0.23
0.03	0.45	0.14	-0.01	-0.30	0.28	0.34	0.68	0.18

LSI – an example

$$A = U L V^T$$

L =

LSI – an example

$$A = U L V^T$$

$V =$

0.20	-0.06	0.11	-0.95	0.05	-0.08	0.18	-0.01	-0.06
0.61	0.17	-0.50	-0.03	-0.21	-0.26	-0.43	0.05	0.24
0.46	-0.13	0.21	0.04	0.38	0.72	-0.24	0.01	0.02
0.54	-0.23	0.57	0.27	-0.21	-0.37	0.26	-0.02	-0.08
0.28	0.11	-0.51	0.15	0.33	0.03	0.67	-0.06	-0.26
0.00	0.19	0.10	0.02	0.39	-0.30	-0.34	0.45	-0.62
0.01	0.44	0.19	0.02	0.35	-0.21	-0.15	-0.76	0.02
0.02	0.62	0.25	0.01	0.15	0.00	0.25	0.45	0.52
0.08	0.53	0.08	-0.03	-0.60	0.36	0.04	-0.07	-0.45

LSI – an example

Choosing the 2 largest singular values we have

0.22	-0.11
0.20	-0.07
0.24	0.04
0.40	0.06
0.64	-0.17
0.27	0.11
0.27	0.11
0.30	-0.14
0.21	0.27
0.01	0.49
0.04	0.62
0.03	0.45

$L_k =$

3.34	0
0	2.54

$V_k^T =$

0.20	0.61	0.46	0.54	0.28	0.00	0.02	0.02	0.08
-0.06	0.17	-0.13	-0.23	0.11	0.19	0.44	0.62	0.53

LSI reconstruction (2 singular values)

$A_k =$

	C1	C2	C3	C4	C5	M1	M2	M3	M4
human	0.16	0.40	0.38	0.47	0.18	-0.05	-0.12	-0.16	-0.09
Interface	0.14	0.37	0.33	0.40	0.16	-0.03	-0.07	-0.10	-0.04
Computer	0.15	0.51	0.36	0.41	0.24	0.02	0.06	0.09	0.12
User	0.26	0.84	0.61	0.70	0.39	0.03	0.08	0.12	0.19
System	0.45	1.23	1.05	1.27	0.56	-0.07	-0.15	-0.21	-0.05
Response	0.16	0.58	0.38	0.42	0.28	0.06	0.13	0.19	0.22
Time	0.16	0.58	0.38	0.42	0.28	0.06	0.13	0.19	0.22
EPS	0.22	0.55	0.51	0.63	0.24	-0.07	-0.14	-0.20	-0.11
Survey	0.10	0.53	0.23	0.21	0.27	0.14	0.31	0.44	0.42
Trees	-0.06	0.23	-0.14	-0.27	0.14	0.24	0.55	0.77	0.66
Graph	-0.06	0.34	-0.15	-0.30	0.20	0.31	0.69	0.98	0.85
Minors	-0.04	0.25	-0.10	-0.21	0.15	0.22	0.50	0.71	0.62

LSI Example

- Query: “human computer interaction” retrieves documents: c_1, c_2, c_4 but *not* c_3 and c_5 .
- If we submit the same query (based on the transformation shown before) to the transformed matrix we retrieve (using cosine similarity) all c_1-c_5 even if c_3 and c_5 have no common keyword to the query.
- According to the transformation for the queries we have:

Query transformation

	query
human	1
Interface	0
computer	1
User	0
System	0
Response	0
Time	0
EPS	0
Survey	0
Trees	0
Graph	0
Minors	0

q =

Query transformation

$$q^T = \begin{bmatrix} 1 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

$$U_k = \begin{bmatrix} 0.22 & -0.11 \\ 0.20 & -0.07 \\ 0.24 & 0.04 \\ 0.40 & 0.06 \\ 0.64 & -0.17 \\ 0.27 & 0.11 \\ 0.27 & 0.11 \\ 0.30 & -0.14 \\ 0.21 & 0.27 \\ 0.01 & 0.49 \\ 0.04 & 0.62 \\ 0.03 & 0.45 \end{bmatrix}$$

$$L_k^{-1} = \begin{bmatrix} 0.3 & 0 \\ 0 & 0.39 \end{bmatrix}$$

$$q_n = q^T U_k L_k^{-1} = \begin{bmatrix} 0.138 & -0.0273 \end{bmatrix}$$

Query transformation

Map docs to
the 2 dim
space $V_k L_k =$

0.20	-0.06
0.61	0.17
0.46	-0.13
0.54	-0.23
0.28	0.11
0.00	0.19
0.01	0.44
0.02	0.62
0.08	0.53

$$\begin{array}{|c|c|} \hline 3.34 & 0 \\ \hline 0 & 2.54 \\ \hline \end{array}$$

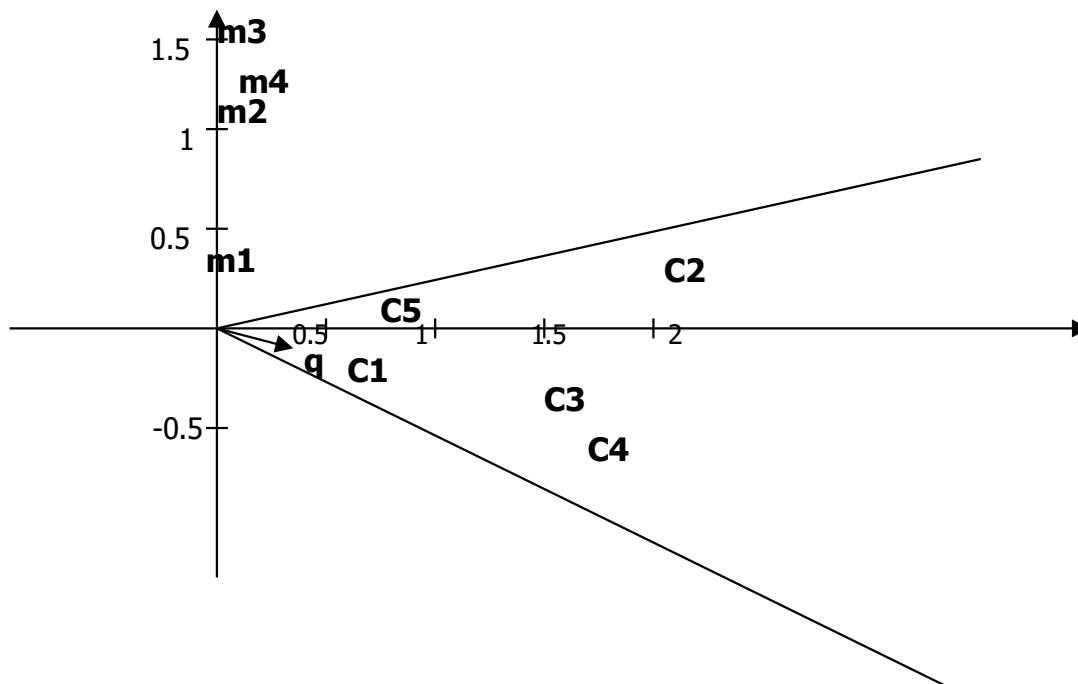
0.67	-0.15
2.04	0.43
1.54	-0.33
1.80	-0.58
0.94	0.28
0.00	0.48
0.03	1.12
0.07	1.57
0.27	1.35

$$q_n L_k = \begin{array}{|c|c|} \hline 0.138 & -0.0273 \\ \hline \end{array} \begin{array}{|c|c|} \hline 3.34 & 0 \\ \hline 0 & 2.54 \\ \hline \end{array} = \begin{array}{|c|c|} \hline 0.46 & -0.069 \\ \hline \end{array}$$

Query transformation

- The cosine similarity matrix of query vector to the documents is:

	query
C1	0.99
C2	0.94
C3	0.99
C4	0.99
C5	0.90
M1	-0.14
M2	-0.13
M3	-0.11
M4	0.05



SVD problems

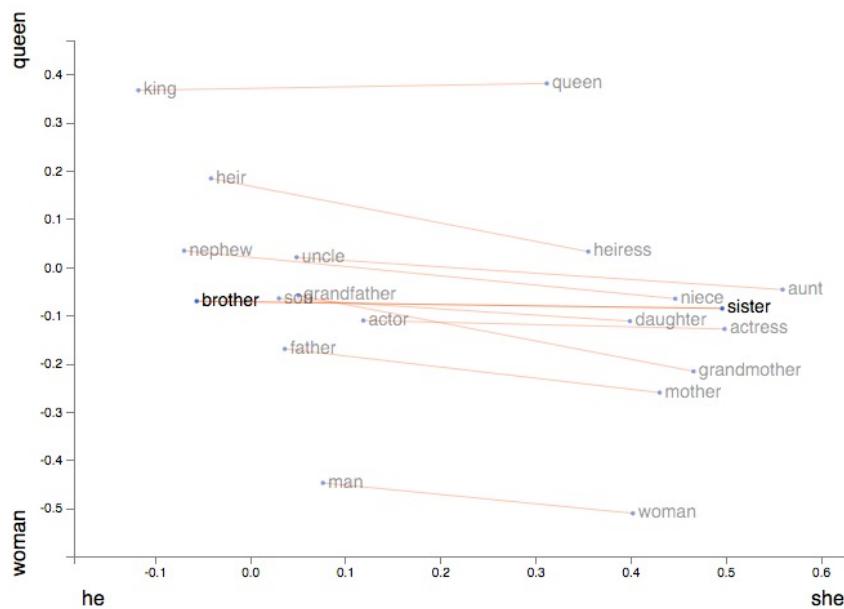
- The dimensions of the matrix change when dictionary changes
- The whole decomposition must be re-calculated when we add a word
- Sensitive to the imbalance of word frequency
- Very high dimensional matrix
- Not suitable for large corpora or dictionaries
- Quadratic cost to perform SVD
- Solution: Directly calculate a low-dimensional representation

OUTLINE

- Representation Learning for Text
 - Latent Semantic Indexing (LSI)
 - **Word2Vec**
- CNN for text classification

Word analogy

- Words with similar meaning end up close to each other
- Words sharing similar contexts may be analogous
 - Synonyms
 - Antonyms
 - Names
 - Colors
 - Places
 - Interchangeable words
- Vector arithmetics to work with analogies
- i.e. **king - man + woman = queen**



<https://lamyiowce.github.io/word2viz/>

But why?

- what's an analogy?

$$\frac{p(w'|man)}{p(w'|woman)} \approx \frac{p(w'|king)}{p(w'|queen)}$$

Assume PMI is approximated by a low rank approximation of the co-occurrence matrix.

1. $PMI(w', w) \approx v_w v_{w'}^T$ *inner product*
2. Isotropic: $E_{w'}[(v_{w'} v_u)^T]^2 = \|v_u\|^2$

Then

3. $\operatorname{argmin}_w E_{w'} [\ln \frac{p(w'|w)}{p(w'|queen)} - \ln \frac{p(w'|man)}{p(w'|woman)}]^2$
4. $\operatorname{argmin}_w E_{w'} [(PMI(w'|w) - PMI(w'|queen)) - (PMI(w'|man) - PMI(w'|woman))]^2$
5. $\operatorname{argmin}_w \|(v_w - v_{queen}) - (v_{man} - v_{woman})\|^2$
6. $v_w \approx v_{queen} - v_{woman} + v_{man}$ which is an analogy!

- Arora et al (ACL 2016) shows that if (2) holds then (1) holds as well
- So we need to construct vectors from co-occurrence that satisfy (2)
- $d < |V|$ in order to have isotropic vectors

Learning Word Vectors

- Corpus containing a sequence of T training words

- Objective: $f(w_t, \dots, w_{t-n+1}) = \hat{P}(w_t | w_{t-n+1} \dots w_{t-1})$

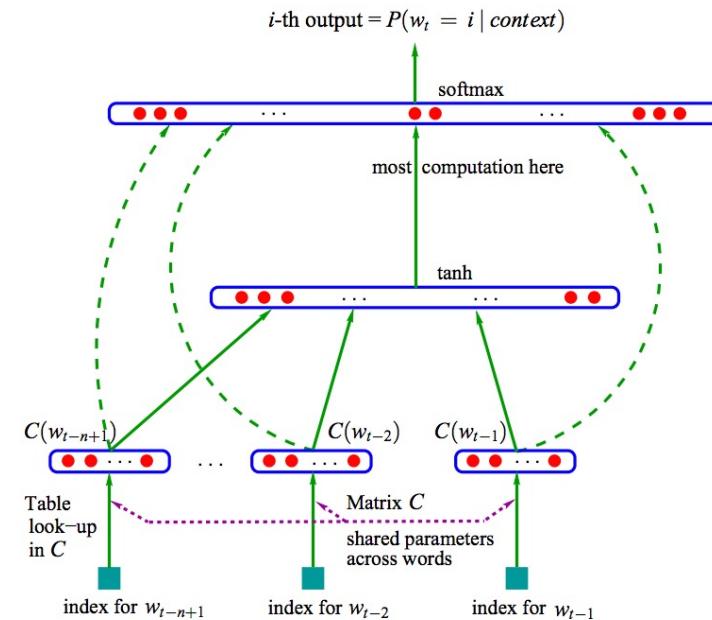
- Decomposed in two parts:

$$w_i \in V \xrightarrow{\text{mapping } C} \mathbb{R}^m$$

- Mapping **C** (1-hotv => lower dimensions)
- Mapping any **g** s.t. (estimate prob $t+1 | t$ previous)

$$f(w_{t-1}, \dots, w_{t-n+1}) = g(C(w_{t-1}), \dots, C(w_{t-n+1}))$$

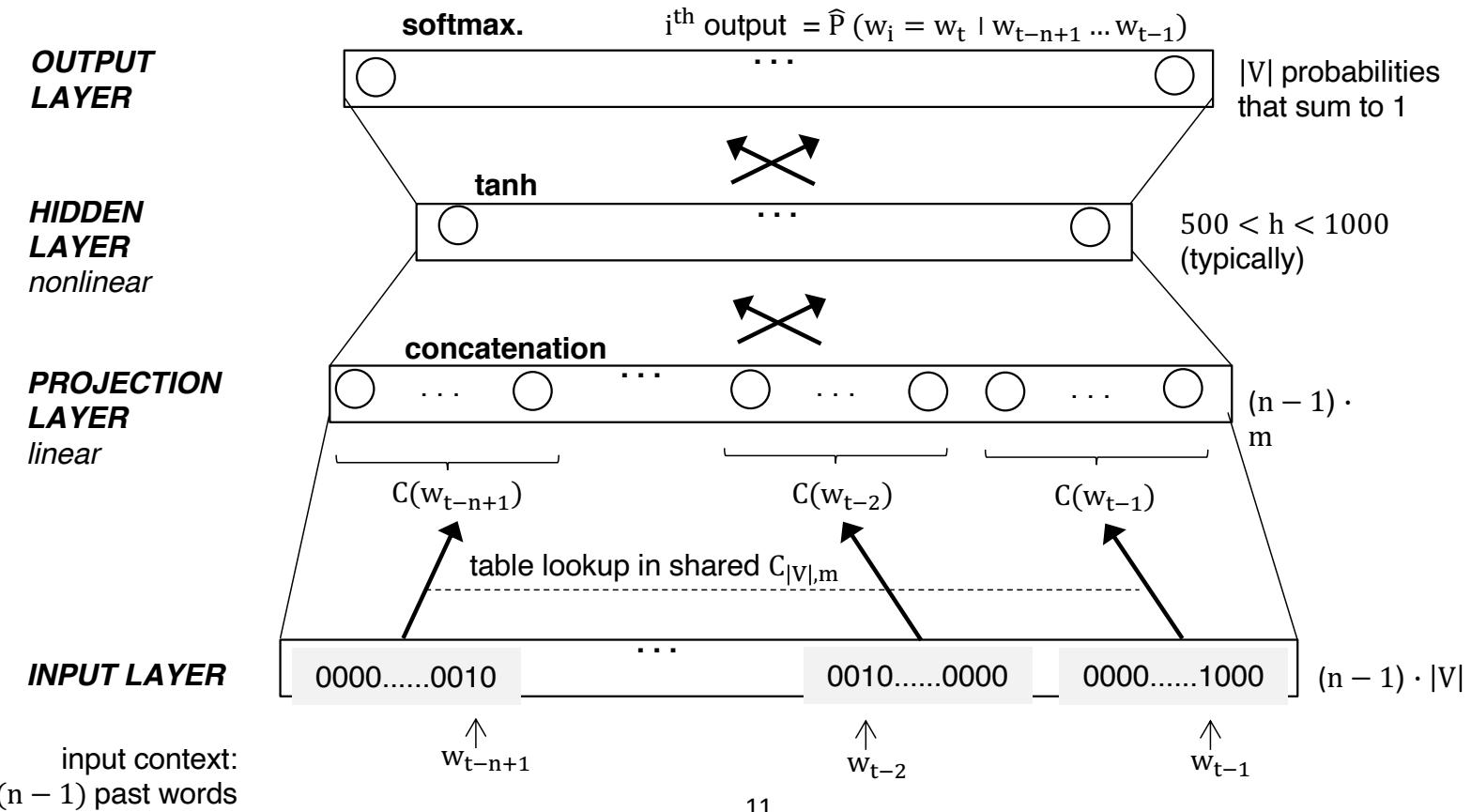
- $C(i)$ is the i-th word feature vector (Word embedding)
- Objective function: $J = \frac{1}{T} \sum f(w_t, \dots, w_{t-n+1})$



[Bengio, Yoshua, et al. "A neural probabilistic language model."](#)
[The Journal of Machine Learning Research 3 \(2003\): 1137-1155.](#)

Neural Net Language Model

For each training sequence:
 input = (context, target) pair: $(w_{t-n+1} \dots w_{t-1}, w_t)$
 objective: minimize $E = -\log \hat{P}(w_t | w_{t-n+1} \dots w_{t-1})$



Objective function

- $E = -\log \hat{P}(w_t | w_{t-n+1} \dots w_{t-1})$
- a probability between 0 and 1.
- On this support, the log is negative $\Rightarrow -\log$ term positive.
- makes sense to try to minimize it.
 - Probability of word given the context be as high as possible (1 for a perfect prediction).
 - case the error is equal to 0 (global minimum).

p	log(p)	-log(p)
0,7	-0,15490196	0,15490196
0,2	-0,698970004	0,698970004

NNLM facts

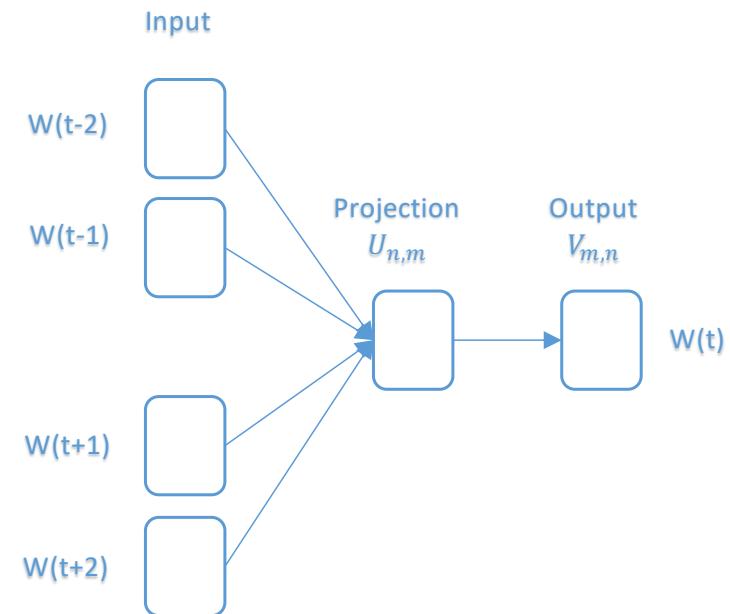
- tested on Brown (1.2M words, $|V| \approx 16K$) and AP News (14M words, $|V| \approx 150K$ reduced to 18K) corpuses
- Brown: $h = 100$, $n = 5$, $m = 30$
- AP News: $h = 60$, $n = 6$, $m = 100$, **3 week** training using **40 cores**
- 24% and 8% relative improvement (resp.) over traditional smoothed n-gram LMs
 - in terms of test set *perplexity*: geometric average
$$1/\widehat{P}(w_t \mid w_{t-n+1} \dots w_{t-1})$$
- Due to **complexity**, NNLM can't be applied to large data sets → poor performance on rare words
- Bengio et al. (2003) initially thought their main contribution was a more accurate LM. They let the interpretation and use of the word vectors as **future work**
- On the opposite, Mikolov et al. (2013) focus on the **word vectors**

Word2Vec

- Mikolov et al. in 2013
- Word2vec key idea: achieve better performance not by using a more complex model (i.e., with more layers), but by allowing a **simpler (shallower) model** to be trained on **much larger amounts of data**
- no hidden layer (leads to 1000X speedup)
- projection layer is shared (not just the weight matrix) - C
- context: words from both history & future:
- Two algorithms for learning words vectors:
 - **CBOW**: from context predict target
 - **Skip-gram**: from target predict context

W2V: Continuous Bag Of Words – CBOW

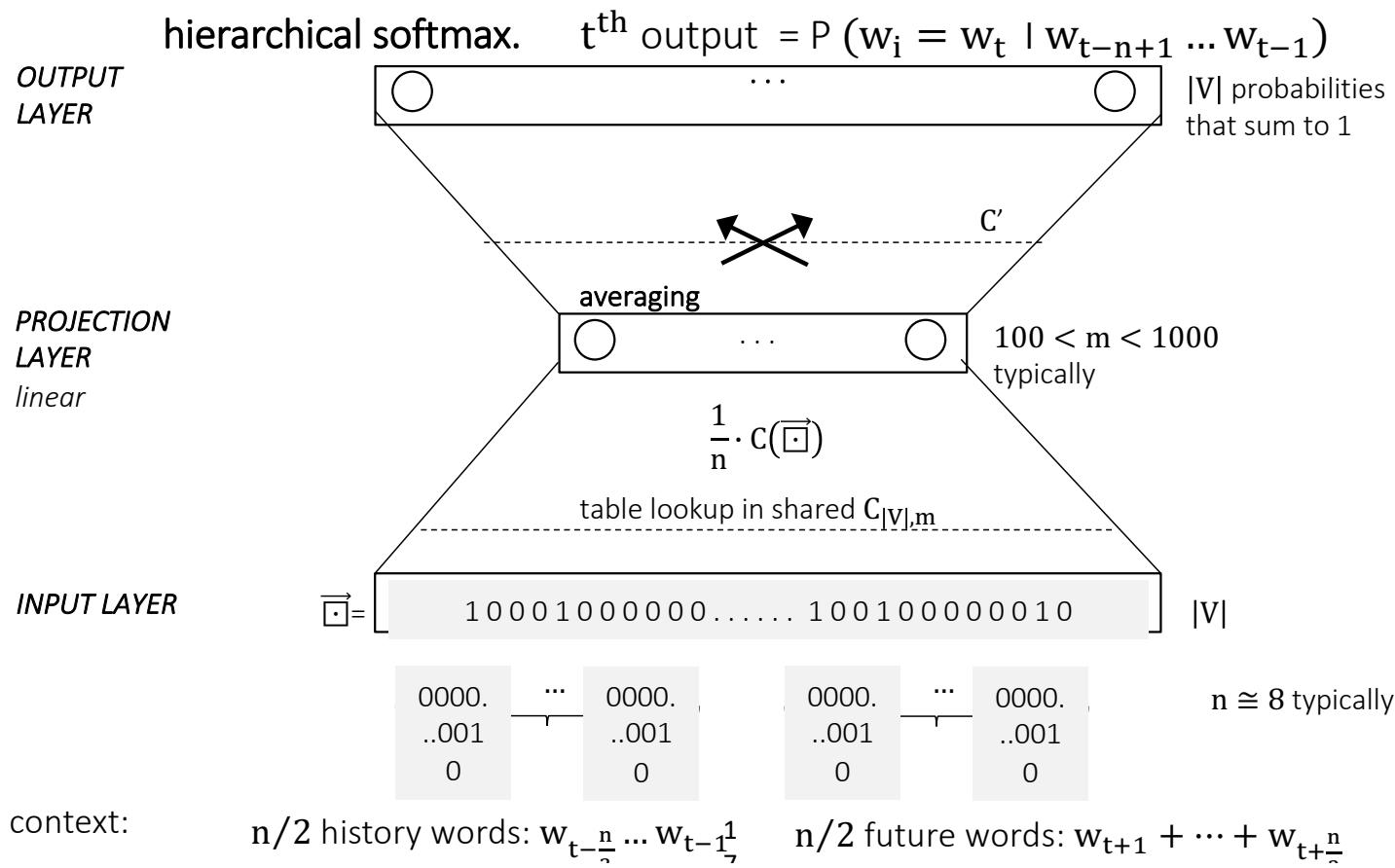
- An unsupervised technique to learn word embeddings.
- CBOW learns the embeddings by predicting the target word (the center one) based on the context words (surrounding words).
- i.e. “Paris is the capital of France located in Ile de France”



Continuous Bag-of-Words (CBOW)

For each training sequence input = (context, target) pair: $(w_{t-\frac{n}{2}} \dots w_{t-1} w_{t+1} \dots w_{t+\frac{n}{2}}, w_t)$

objective: minimize $-\log \hat{P}(w_t | w_{t-n+1} \dots w_{t-1})$



W2V: Continuous Bag Of Words – CBOW : Forward

- Each word **W(t)** is represented by one-hot vector of size **n** (vocabulary size).
- The **w** context words are averaged and forwarded to the projection layer to produce an embedded vector **z** of size **m**:

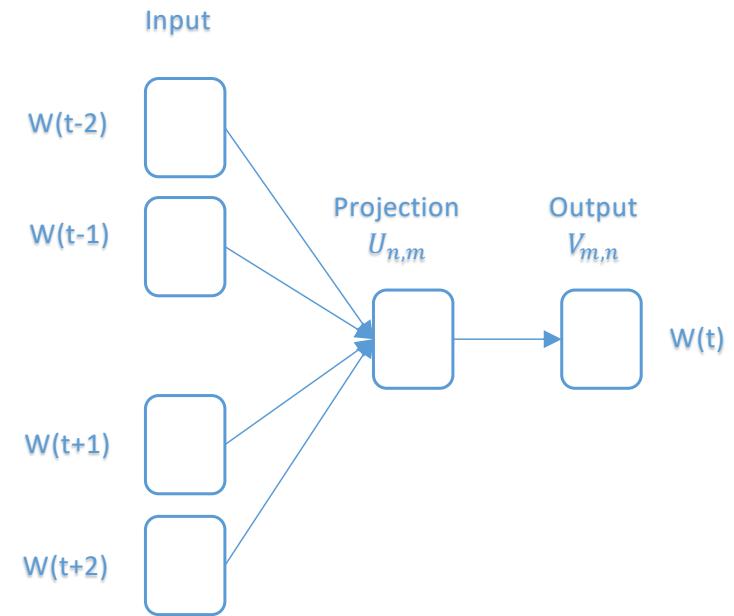
$$z_{1,m} = \frac{1}{w} \sum_{W \in \text{context}} W_{1,n} U_{n,m}$$

- The vector **z** is forwarded to an output projection layer that produce the out vector **y** of size **n**.

$$y_{1,n} = z_{1,m} V_{m,n}$$

- Finally, a soft-max activation function is applied to the output to find a vector of probabilities for each word:

$$\sigma(y_i) = \frac{e^{y_i}}{\sum_j e^{y_j}}$$

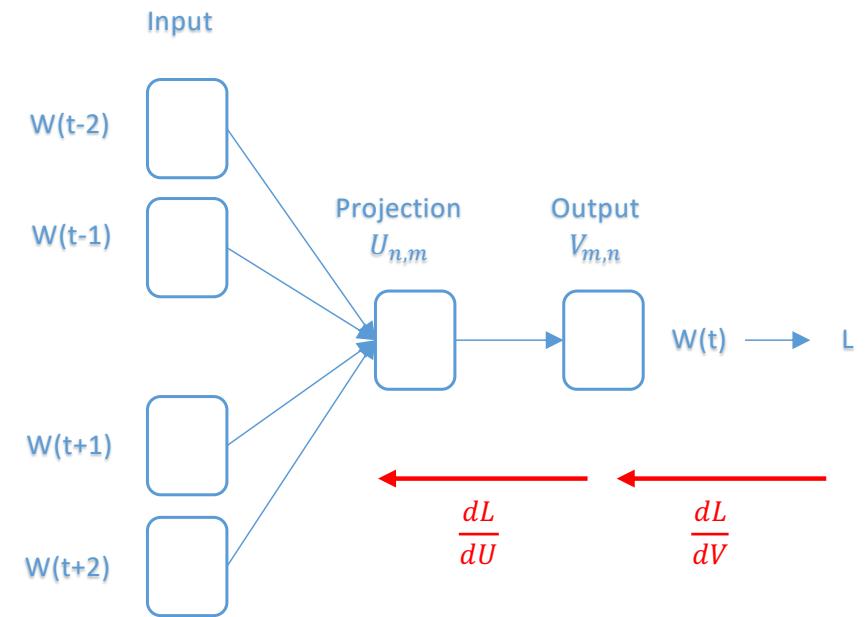


W2V: Continuous Bag Of Words – CBOW : Backward

- To update the weights, first we compute the log loss function (cross-entropy):

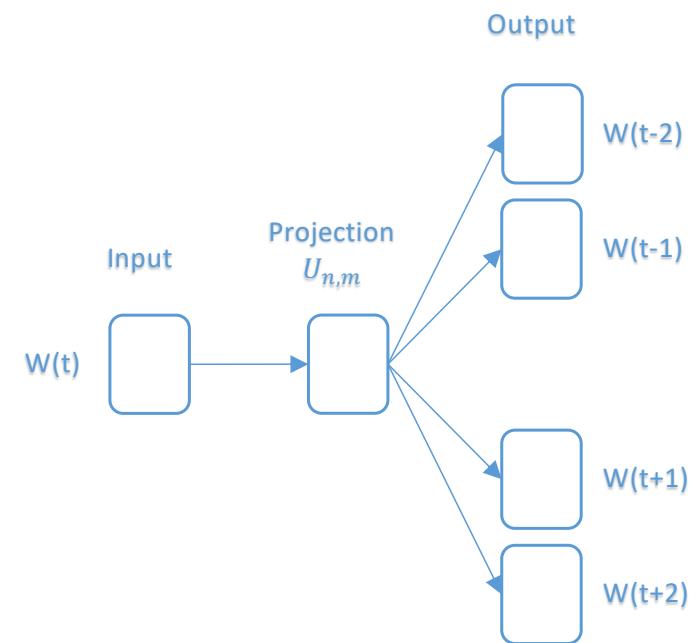
$$L = -\frac{1}{n} \sum_{i=1}^n W(t)_i \log(\sigma(y_i))$$

- The weights (matrices U and V) are now updated using gradient descent with learning rate α .
- $V = V - \alpha \frac{dL}{dV}$
- $U = U - \alpha \frac{dL}{dU}$
- Finally, after multiple passes through the corpus, **U** is the final **Word Embeddings Matrix** where each row represent a vector of size **m** for a specific word.



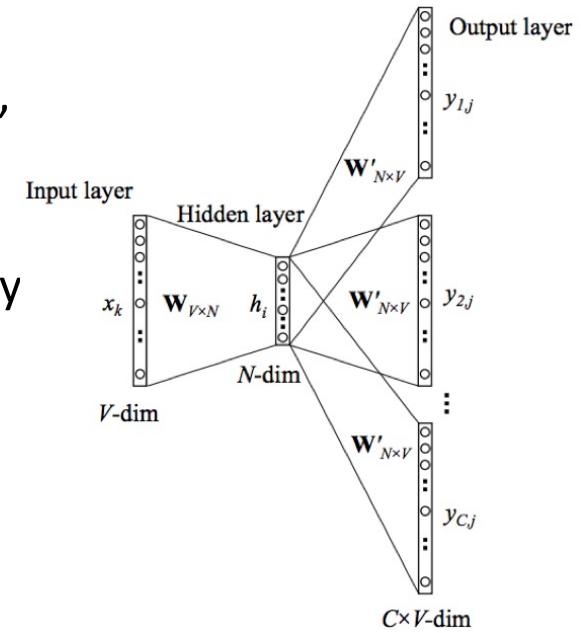
W2V: Skip-Gram

- An unsupervised technique to learn word embeddings.
- Skip-gram learn the embeddings by predicting the context of the word.
- The used loss function is cross-entropy as in CBOW.



Skip-gram

- skip-gram uses the context's center word to predict the surrounding words
- i.e. “Paris is the capital of France located in Ile de France”
- instead of computing the probability of the target word w_t given its previous words, we calculate the probability of the surrounding word w_{t+j} given w_t
- $p(w_{t+j}|w_t) = \frac{\exp(v_{w_t}^T v'_{w_{t+j}})}{\sum_{w \in V} \exp(v_{w_t}^T v'_{w_{t+j}})}$
- $v_{w_t}^T$ is a column of W_{VxN} and $v'_{w_{t+j}}$ is a column of W'_{NxV}
- Objective function $J = \frac{1}{T} \sum_{t=1}^T \sum_{-n \leq j \leq n} \log p(w_{t+j}|w_t)$



Word2vec facts

- Complexity is $n * m + m * \log|V|$ (Mikolov et al. 2013a)
- n : size of context window (~ 10) $n \times m$: dimensions of the projection layer, $|V|$ size of the vocabulary
- On Google news 6B words training corpus, with $|V| \sim 10^6$:
 - CBOW with $m = 1000$ took **2 days** to train on **140 cores**
 - Skip-gram with $m = 1000$ took **2.5 days** on **125 cores**
 - NNLM (Bengio et al. 2003) took **14 days** on **180 cores**, for $m = 100$ only!
- word2vec training speed $\cong 100K\text{-}5M$ words/s
- Quality of the word vectors:
 - \nearrow significantly with **amount of training data** and **dimension of the word vectors** (m)
 - measured in terms of accuracy on 20K semantic and syntactic association tasks.
 - e.g., words in **bold** have to be returned:

Capital-Country	Past tense	Superlative	Male-Female	Opposite
Athens: Greece	walking: walked	easy: easiest	brother: sister	ethical: unethical

- Best NNLM: 12.3% overall accuracy. Word2vec (with Skip-gram): 53.3%
- References: <http://www.scribd.com/doc/285890694/NIPS-DeepLearningWorkshop-NNforText#scribd> ---- <https://code.google.com/p/word2vec/>

GloVe

Probability and Ratio	$k = solid$	$k = gas$	$k = water$	$k = fashion$
$P(k ice)$	1.9×10^{-4}	6.6×10^{-5}	3.0×10^{-3}	1.7×10^{-5}
$P(k steam)$	2.2×10^{-5}	7.8×10^{-4}	2.2×10^{-3}	1.8×10^{-5}
$P(k ice)/P(k steam)$	8.9	8.5×10^{-2}	1.36	0.96

- Ratio of co-occurrence probabilities best distinguishes relevant words

$$F(w_i, w_j, \tilde{w}_k) = \frac{P_{ik}}{P_{jk}}$$



$$w_i^T \tilde{w}_k + b_i + \tilde{b}_k = \log(X_{ik})$$

- Cast this into a least square problem:

- X co-occurrence matrix
- f weighting function,
- b bias terms
- w_i = word vector
- \tilde{w}_j = context vector

$$J = \sum_{i,j=1}^V f(X_{ij}) (w_i^T \tilde{w}_j + b_i + \tilde{b}_j - \log X_{ij})^2$$

$$f(x) = \begin{cases} (x/x_{\max})^\alpha & \text{if } x < x_{\max} \\ 1 & \text{otherwise} \end{cases} .$$

model that utilizes

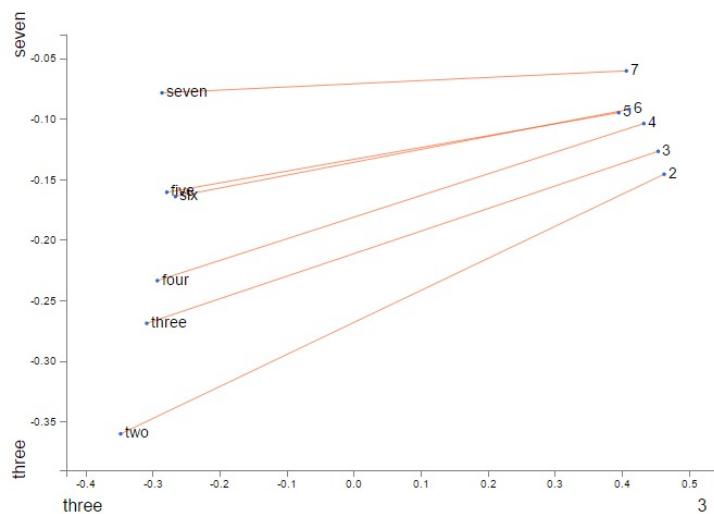
- count data
- bilinear prediction-based methods like word2vec

Which is better?

- Open question
- SVD vs word2vec vs GloVe
- All based on co-occurrence
- *Levy, O., Goldberg, Y., & Dagan, I. (2015)*
 - SVD performs best on similarity tasks
 - Word2vec performs best on analogy tasks
 - *No single algorithm consistently outperforms the other methods*
 - *Hyperparameter tuning is important*
 - 3 out of 6 cases, tuning hyperparameters is more beneficial than increasing corpus size
 - word2vec outperforms GloVe on all tasks
 - *CBOW is worse than skip-gram on all tasks*

Applications

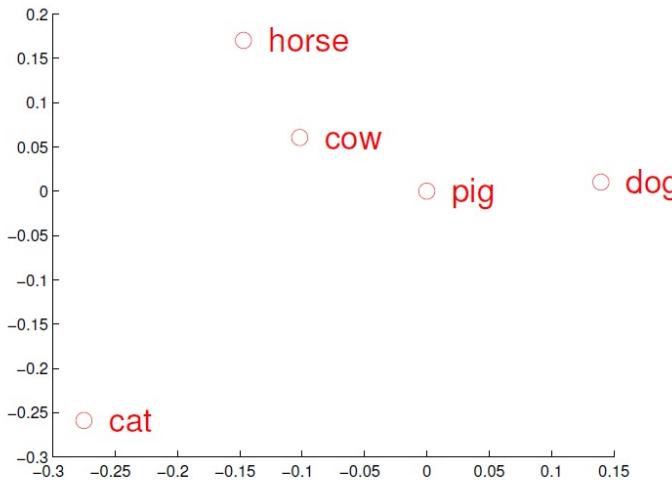
- Word analogies
- Find similar words
 - Semantic similarity
 - Syntactic similarity
- POS tagging
- Similar analogies for different languages
- Document classification



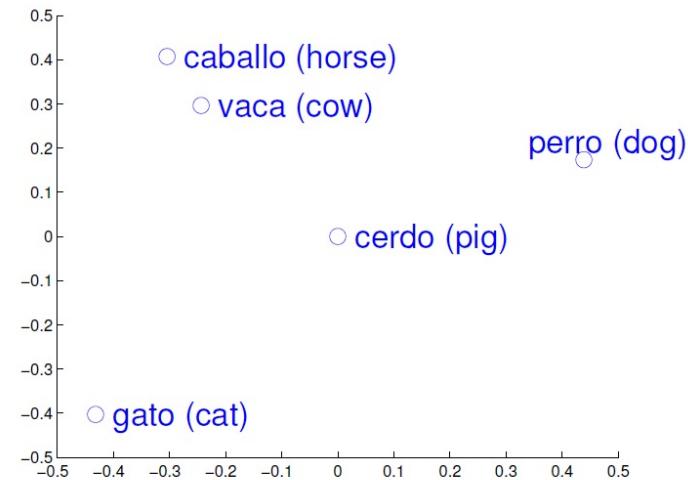
<https://lamyiowce.github.io/word2viz/>

Applications

- High quality word vectors boost performance of all NLP tasks, including document classification, machine translation, information retrieval...
- Example for English to Spanish machine translation:

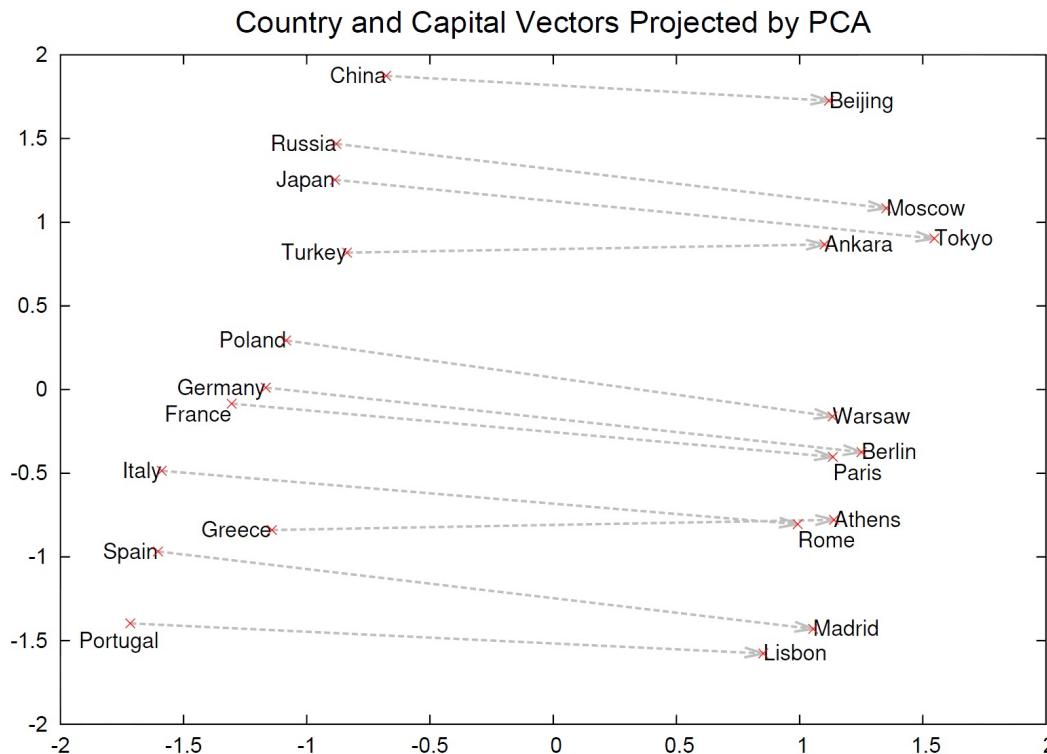


About 90% reported accuracy (Mikolov et al. 2013c)



[Mikolov, T., Le, Q. V., & Sutskever, I. \(2013\). Exploiting similarities among languages for machine translation. arXiv preprint arXiv:1309.4168.](https://arxiv.org/abs/1309.4168)

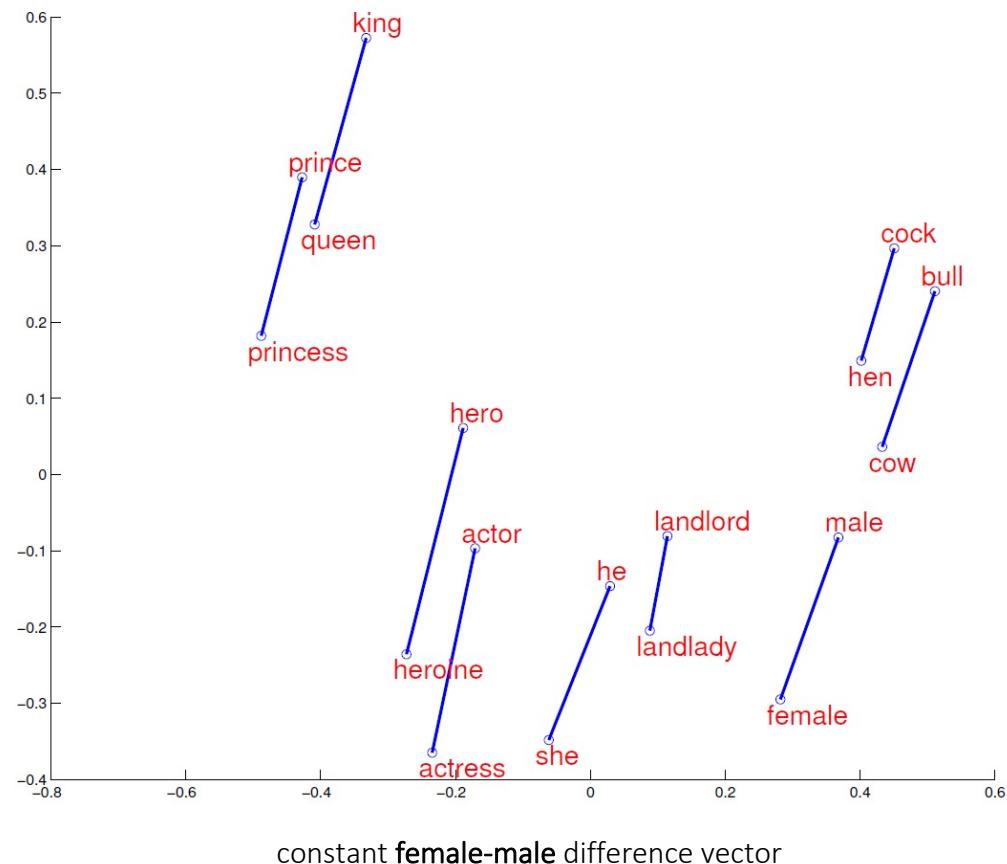
Remarkable properties of word vectors



regularities between words are encoded in the difference vectors
e.g., there is a constant **country-capital** difference vector

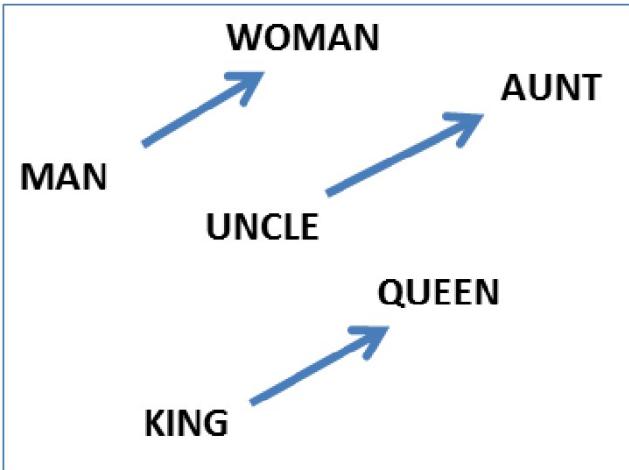
Mikolov et al. (2013b)
Distributed representations of
words and phrases and their
compositionality

Remarkable properties of word vectors

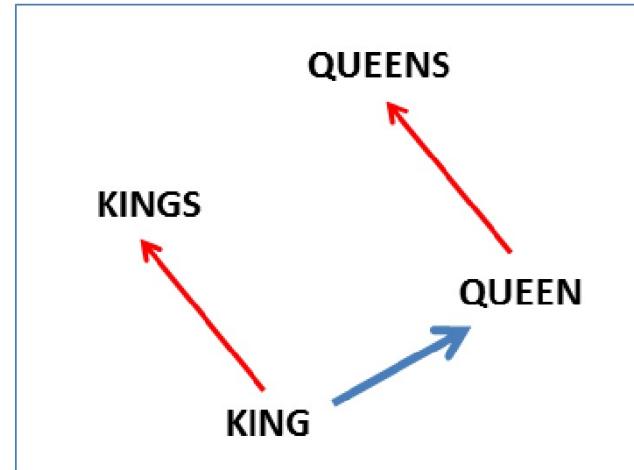


<http://www.scribd.com/doc/285890694/NIPS-DeepLearningWorkshop-NNforText#scribd>

Remarkable properties of word vectors



constant **male-female** difference vector



constant **singular-plural** difference vector

- Vector operations are supported and make intuitive sense:

$$w_{king} - w_{man} + w_{woman} \cong w_{queen}$$

$$w_{einstein} - w_{scientist} + w_{painter} \cong w_{picasso}$$

$$w_{paris} - w_{france} + w_{italy} \cong w_{rome}$$

$$w_{his} - w_{he} + w_{she} \cong w_{her}$$

$$w_{windows} - w_{microsoft} + w_{google} \cong w_{android}$$

$$w_{cu} - w_{copper} + w_{gold} \cong w_{au}$$

- Online [demo](#) (scroll down to end of tutorial)

<http://rare-technologies.com/word2vec-tutorial/>

OUTLINE

- Representation Learning for Text
 - Latent Semantic Indexing (LSI)
 - Word2Vec
- **CNN for text classification**

CNN for Text Classification

- Use the word embeddings of the document terms as input for Convolutional Neural Network
- Input must be fixed size
- Applies multiple filters to concatenated word vectors
- Produces new features for every filter
- picks the max as a feature for the CNN

CNN architecture for document classification

- Use the high quality embeddings as input for Convolutional Neural Network
- Applies multiple filters to concatenated word vector

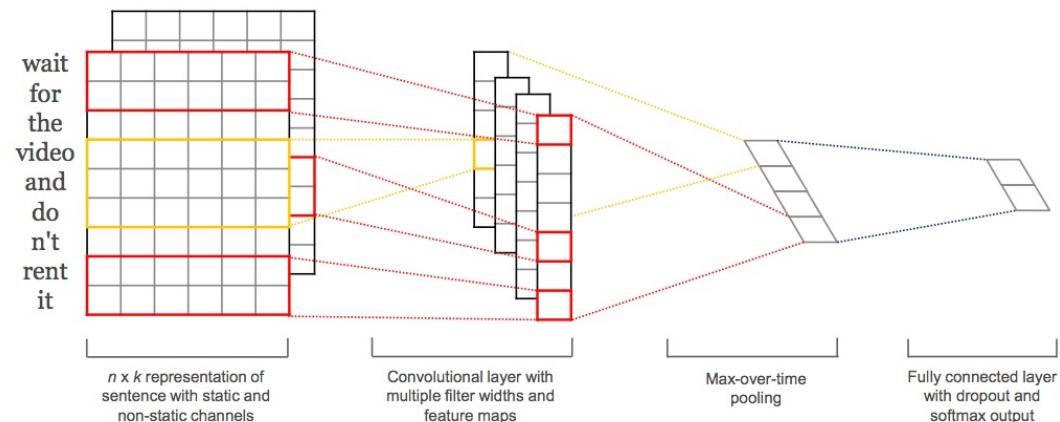
$$\mathbf{x}_{1:n} = \mathbf{x}_1 \oplus \mathbf{x}_2 \oplus \dots \oplus \mathbf{x}_n$$

- Produces new features for every filter

$$c_i = f(\mathbf{w} \cdot \mathbf{x}_{i:i+h-1} + b)$$

- And picks the max as a feature for the CNN

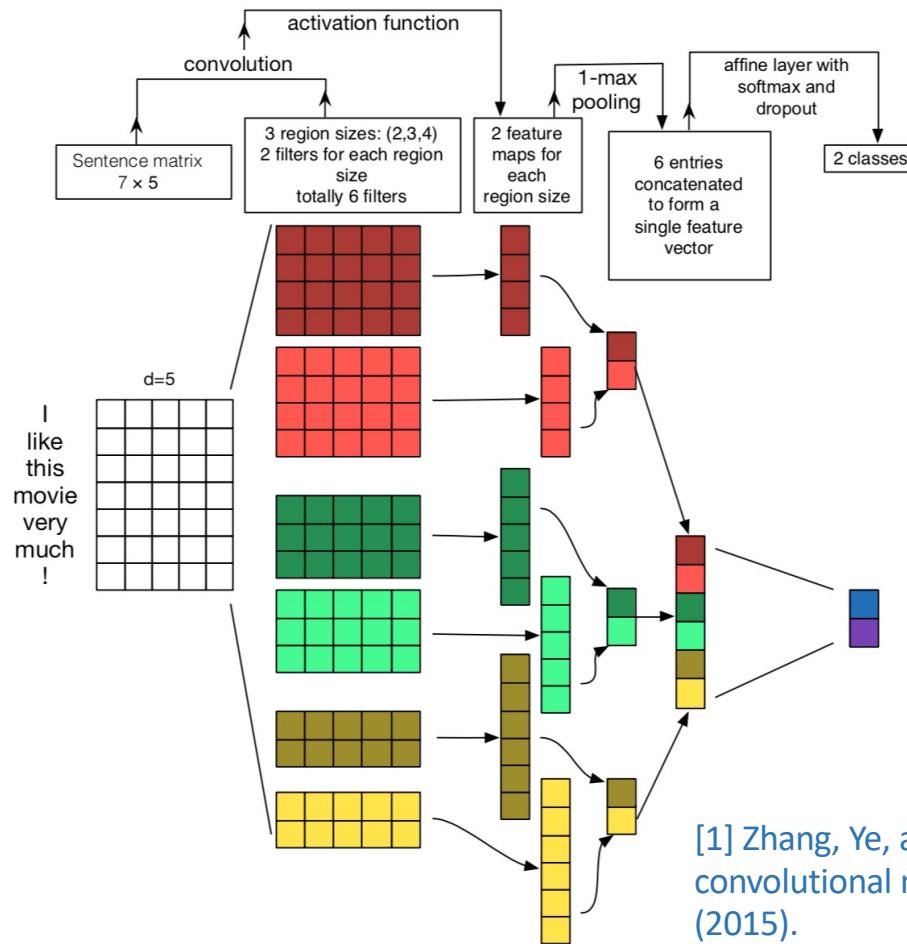
$$\mathbf{c} = [c_1, c_2, \dots, c_{n-h+1}] \quad \hat{c} = \max\{\mathbf{c}\}$$



[Yoon Kim - Convolutional Neural Networks for Sentence Classification](#)

CNN architecture for document classification

[1]



- Data (text) only 1st column of input
- Rest of each row: embedding (in images 2D+RGB dimension)
- Filters of different sizes (4x5, 3x5 etc.)
 - Each size captures different features (need $\sim 10^2$ filters/size)
- Feature maps:
 - As many as the times filter fits on data matrix
 - Max pooling maintains the “best features”
 - Global feature map => classification via softmax

[1] Zhang, Ye, and Byron Wallace. "A sensitivity analysis of (and practitioners' guide to) convolutional neural networks for sentence classification." arXiv preprint arXiv:1510.03820 (2015).

CNN for text classification

Many variations of the model [1]

- use existing vectors as input (CNN-static)
- learn vectors for the specific classification task through backpropagation (CNN-rand)
- Modify existing vectors for the specific task through backpropagation(CNN-non-static)

[1] Y. Kim, Convolutional Neural Networks for Sentence Classification, EMNLP 2014

CNN for text classification

- Combine multiple word embeddings
- Each set of vectors is treated as a ‘channel’
- Filters applied to all channels
- Gradients are back-propagated only through one of the channels
- Fine-tunes one set of vectors while keeping the other static

CNN for text classification

Model	MR	SST-1	SST-2	Subj	TREC	CR	MPQA
CNN-rand	76.1	45.0	82.7	89.6	91.2	79.8	83.4
CNN-static	81.0	45.5	86.8	93.0	92.8	84.7	89.6
CNN-non-static	81.5	48.0	87.2	93.4	93.6	84.3	89.5
CNN-multichannel	81.1	47.4	88.1	93.2	92.2	85.0	89.4
RAE (Socher et al., 2011)	77.7	43.2	82.4	—	—	—	86.4
MV-RNN (Socher et al., 2012)	79.0	44.4	82.9	—	—	—	—
RNTN (Socher et al., 2013)	—	45.7	85.4	—	—	—	—
DCNN (Kalchbrenner et al., 2014)	—	48.5	86.8	—	93.0	—	—
Paragraph-Vec (Le and Mikolov, 2014)	—	48.7	87.8	—	—	—	—
CCAE (Hermann and Blunsom, 2013)	77.8	—	—	—	—	—	87.2
Sent-Parser (Dong et al., 2014)	79.5	—	—	—	—	—	86.3
NBSVM (Wang and Manning, 2012)	79.4	—	—	93.2	—	81.8	86.3
MNB (Wang and Manning, 2012)	79.0	—	—	93.6	—	80.0	86.3
G-Dropout (Wang and Manning, 2013)	79.0	—	—	93.4	—	82.1	86.1
F-Dropout (Wang and Manning, 2013)	79.1	—	—	93.6	—	81.9	86.3
Tree-CRF (Nakagawa et al., 2010)	77.3	—	—	—	—	81.4	86.1
CRF-PR (Yang and Cardie, 2014)	—	—	—	—	—	82.7	—
SVM _S (Silva et al., 2011)	—	—	—	—	95.0	—	—

Accuracy scores (Kim et al vs others)

CNN architecture for (short) document classification – T-SNE visualization (see Lab notes)

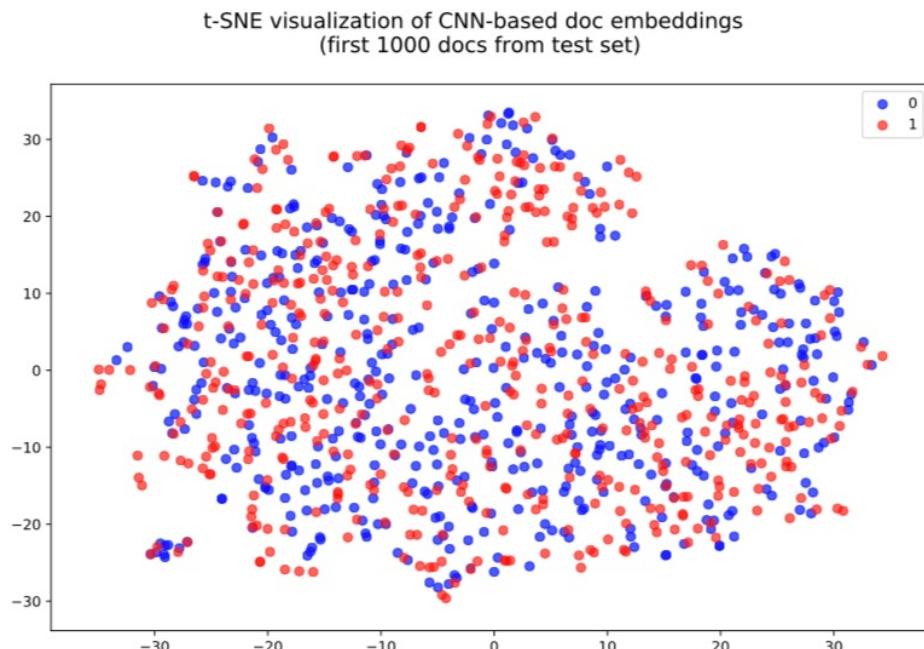


Figure 2: Doc embeddings before training.

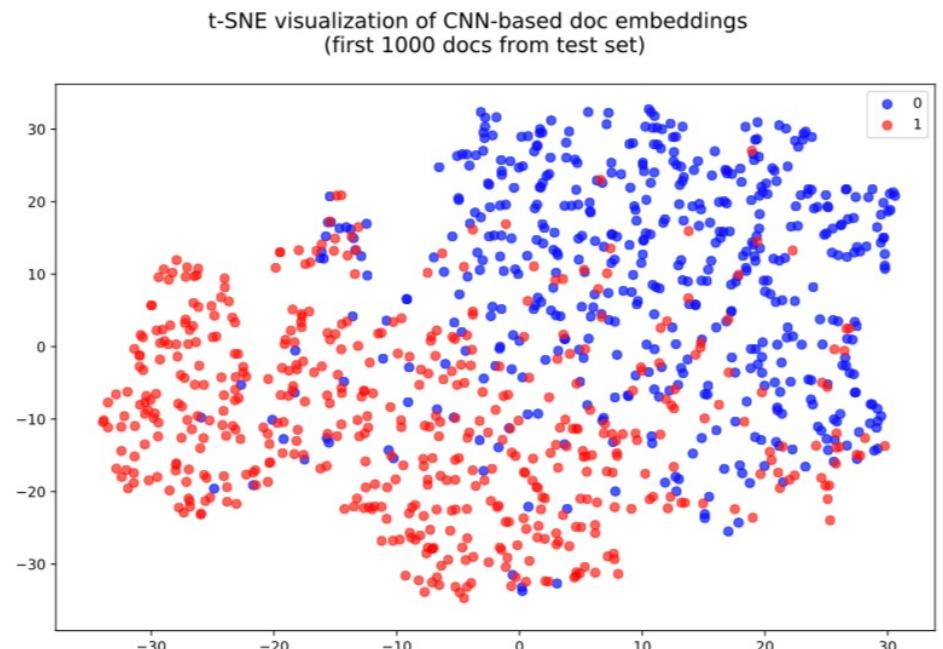


Figure 3: Doc embeddings after 2 epochs.

CNN architecture for document classification - Saliency maps (see Lab notes)

- words are most related to changing the doc classification
- $A \in R^{sxd}$, s :# sentence words, d :size of embeddings saliency(a) = $\left| \frac{\partial(\text{CNN})}{\partial a} \right|_a$

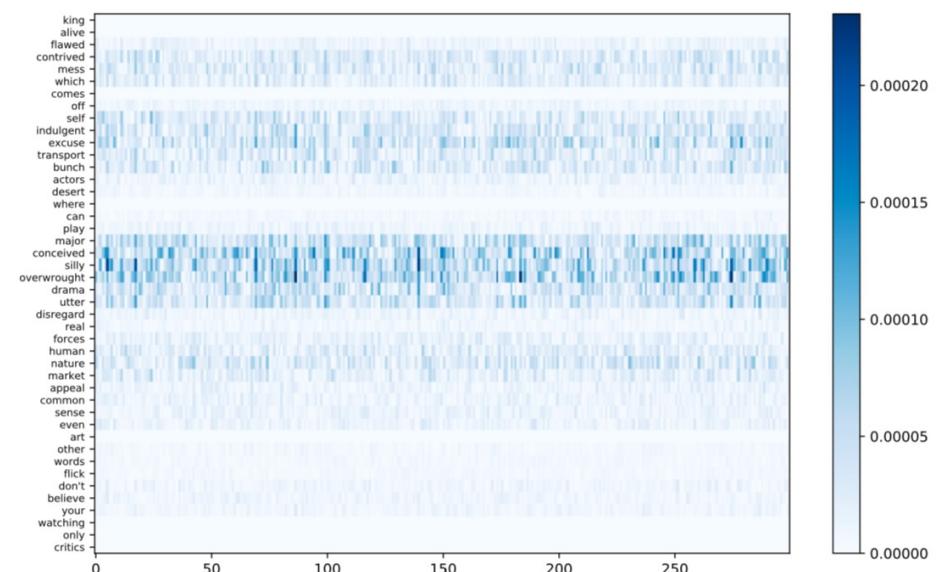
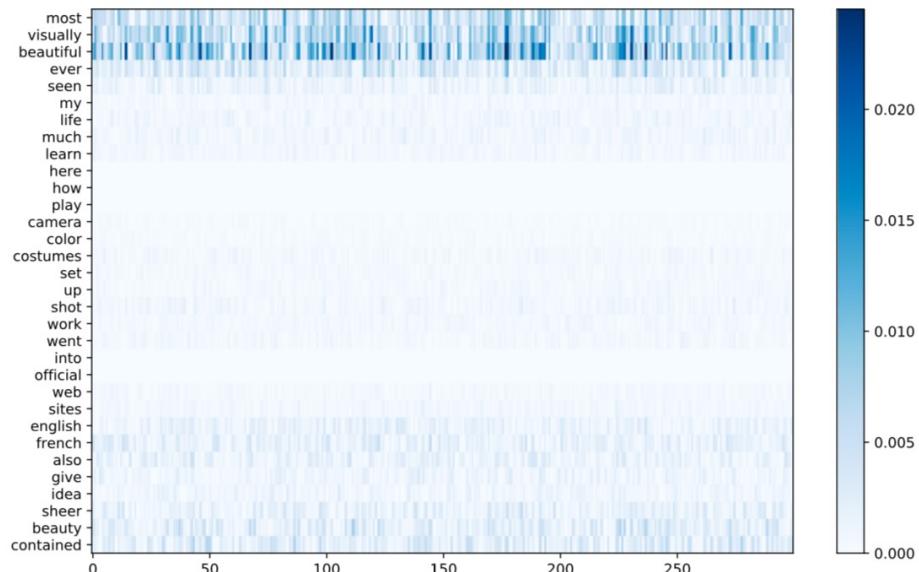
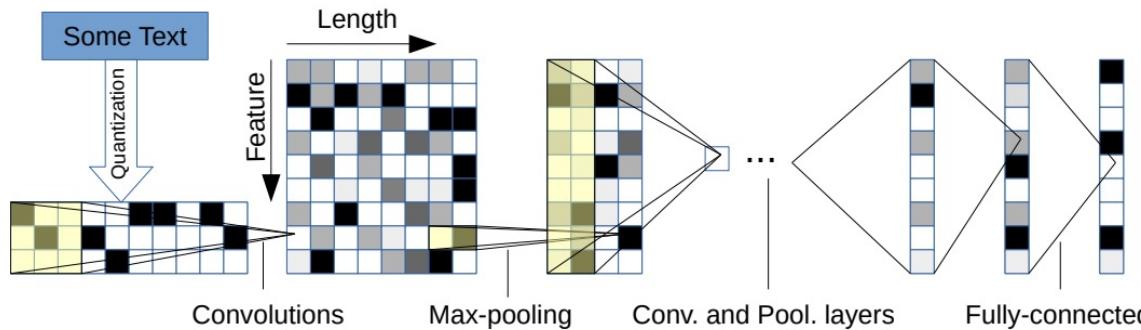


Figure 4: Saliency map for document 1 of the IMDB test set (true label: positive) Figure 5: Saliency map for document 15 of the IMDB test set (true label: negative)

Character-level CNN for Text Classification

- Input: sequence of encoded characters
- quantize each character using “one-hot” encoding
- input feature length is 1014 characters
- 1014 characters able capture most of the texts of interest
- Also perform Data Augmentation using Thesaurus as preprocessing step

Model Architecture



- 9 layers deep
- 6 convolutional layers
- 3 fully-connected layers
- 2 dropout modules in between the fully-connected layers for regularization

Model Comparison

Model	AG	Sogou	DBP.	Yelp P.	Yelp F.	Yah. A.	Amz. F.	Amz. P.
BoW	11.19	7.15	3.39	7.76	42.01	31.11	45.36	9.60
BoW TFIDF	10.36	6.55	2.63	6.34	40.14	28.96	44.74	9.00
ngrams	7.96	2.92	1.37	4.36	43.74	31.53	45.73	7.98
ngrams TFIDF	7.64	2.81	1.31	4.56	45.20	31.49	47.56	8.46
Bag-of-means	16.91	10.79	9.55	12.67	47.46	39.45	55.87	18.39
LSTM	13.94	4.82	1.45	5.26	41.83	29.16	40.57	6.10
Lg. w2v Conv.	9.92	4.39	1.42	4.60	40.16	31.97	44.40	5.88
Sm. w2v Conv.	11.35	4.54	1.71	5.56	42.13	31.50	42.59	6.00
Lg. w2v Conv. Th.	9.91	-	1.37	4.63	39.58	31.23	43.75	5.80
Sm. w2v Conv. Th.	10.88	-	1.53	5.36	41.09	29.86	42.50	5.63
Lg. Lk. Conv.	8.55	4.95	1.72	4.89	40.52	29.06	45.95	5.84
Sm. Lk. Conv.	10.87	4.93	1.85	5.54	41.41	30.02	43.66	5.85
Lg. Lk. Conv. Th.	8.93	-	1.58	5.03	40.52	28.84	42.39	5.52
Sm. Lk. Conv. Th.	9.12	-	1.77	5.37	41.17	28.92	43.19	5.51
Lg. Full Conv.	9.85	8.80	1.66	5.25	38.40	29.90	40.89	5.78
Sm. Full Conv.	11.59	8.95	1.89	5.67	38.82	30.01	40.88	5.78
Lg. Full Conv. Th.	9.51	-	1.55	4.88	38.04	29.58	40.54	5.51
Sm. Full Conv. Th.	10.89	-	1.69	5.42	37.95	29.90	40.53	5.66
Lg. Conv.	12.82	4.88	1.73	5.89	39.62	29.55	41.31	5.51
Sm. Conv.	15.65	8.65	1.98	6.53	40.84	29.84	40.53	5.50
Lg. Conv. Th.	13.39	-	1.60	5.82	39.30	28.80	40.45	4.93
Sm. Conv. Th.	14.80	-	1.85	6.49	40.16	29.84	40.43	5.67

Testing errors for all models
 Blue->best, Red->worst

links

- <http://yann.lecun.com/exdb/publis/pdf/lecun-01a.pdf>
- <https://arxiv.org/pdf/1509.01626.pdf>
- <http://www.aclweb.org/anthology/D14-1181>
- <http://cs231n.github.io/convolutional-networks/>
- <http://ufldl.stanford.edu/tutorial/supervised/Pooling/>

THANK YOU