

Word2Vec

Hung-Hsuan Chen

1/11/21

1

1

Attribute of a text document

- One-hot encoding the words

	Document 1	Document 2	Document 3	Document 4	Document 5	Document 6	Document 7	Document 8
Term(s) 1	10	0	1	0	0	0	0	2
Term(s) 2	0	2	0	0	0	18	0	2
Term(s) 3	0	0	0	0	0	0	0	2
Term(s) 4	6	0	0	4	6	0	0	0
Term(s) 5	0	0	0	0	0	0	0	2
Term(s) 6	0	0	1	0	0	1	0	0
Term(s) 7	0	1	8	0	0	0	0	0
Term(s) 8	0	0	0	0	0	3	0	0

← Word Vector (Passage Vector)

Document Vector

1/11/21

2

2

Term-document count matrices

- Consider the number of occurrences of a term in a document:
 - Each document is a count vector in \mathbb{N}^v : a column below

	Antony and Cleopatra	Julius Caesar	The Tempest	Hamlet	Othello	Macbeth
Antony	157	73	0	0	0	0
Brutus	4	157	0	1	0	0
Caesar	232	227	0	2	1	1
Calpurnia	0	10	0	0	0	0
Cleopatra	57	0	0	0	0	0
mercy	2	0	3	5	5	1
worser	2	0	1	1	1	0

1/11/21

3

3

Bag of words model

- Vector representation doesn't consider the ordering of words in a document
- John is quicker than Mary and Mary is quicker than John have the same vectors*
- This is called the bag of words model.

1/11/21

4

4

Term frequency tf

- The term frequency $tf_{t,d}$ of term t in document d is defined as the number of times that t occurs in d .
- We want to use tf when computing query-document match scores. But how?
- Raw term frequency is not what we want:
 - A document with 10 occurrences of the term is more relevant than a document with 1 occurrence of the term.
 - But not 10 times more relevant.
- Relevance does not increase proportionally with term frequency.

1/11/21

5

5

Log-frequency weighting

- The log frequency weight of term t in d is

$$w_{t,d} = \begin{cases} 1 + \log_{10} tf_{t,d}, & \text{if } tf_{t,d} > 0 \\ 0, & \text{otherwise} \end{cases}$$
- $0 \rightarrow 0, 1 \rightarrow 1, 2 \rightarrow 1.3, 10 \rightarrow 2, 1000 \rightarrow 4$, etc.
- Score for a document-query pair: sum over terms t in both q and d :
- $\text{score} = \sum_{t \in q \cap d} (1 + \log tf_{t,d})$
- The score is 0 if none of the query terms is present in the document.

1/11/21

6

6

Document frequency

- Rare terms are more informative than frequent terms
 - Recall stop words
- Consider a term in the query that is rare in the collection (e.g., *arachnocentric*)
- A document containing this term is very likely to be relevant to the query *arachnocentric*
 - We want a high weight for rare terms like *arachnocentric*.
- We will use document frequency (df) to capture this

1/11/21

7

7

idf weight

- df_t is the document frequency of t : the number of documents that contain t
 - df_t is an inverse measure of the informativeness of t
 - $df_t \leq N$
- We define the idf (inverse document frequency) of t by

$$idf_t = \log_{10} (N/df_t)$$
 - We use $\log (N/df_t)$ instead of N/df_t to “dampen” the effect of idf.

1/11/21

8

8

idf example, suppose $N = 1$ million

term	df_t	idf_t
calpurnia	1	
animal	100	
sunday	1,000	
fly	10,000	
under	100,000	
the	1,000,000	

$$idf_t = \log_{10} (N/df_t)$$

There is one idf value for each term t in a collection.

1/11/21

9

9

tf-idf weighting

- The tf-idf weight of a term is the product of its tf weight and its idf weight.

$$w_{t,d} = \log(1 + tf_{t,d}) \times \log_{10}(N / df_t)$$

- Best known weighting scheme in information retrieval**
 - Note: the “-” in tf-idf is a hyphen, not a minus sign!
 - Alternative names: tf.idf, tf x idf**
- Increases with the number of occurrences within a document
- Increases with the rarity of the term in the collection**
- Many variants, e.g., BM25

1/11/21

10

10

These are all bag-of-words model

- The feature sets are the same
 - Each word is a feature
- The feature values are generated in different ways
 - Based on term counts
 - Based on log of term frequency
 - Based on tf-idf
- Compare document similarity by vectors
- Compare word similarity by vectors

1/11/21

11

11

N-grams

Word-level unigrams

Text	Token Sequence	Token Value
One Two Three Four	1	One
One Two Three Four	2	Two
One Two Three Four	3	Three
One Two Three Four	4	Four

Word-level bigrams

Text	Token Sequence	Token Value
One Two Three Four	1	One Two
One Two Three Four	2	Two Three
One Two Three Four	3	Three Four

Word-level trigrams

Text	Token Sequence	Token Value
One Two Three Four	1	One Two Three
One Two Three Four	2	Two Three Four

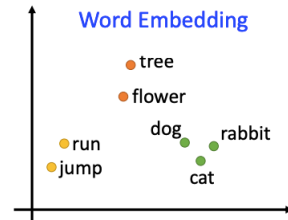
1/11/21

12

12

Word2Vec

- Learning the word vectors from a large collection of documents
- Resulting word representation
 - Similar words are closer
 - The word vectors capture the semantic of the words
- The output vectors are also known as word embedding



1/11/21

13

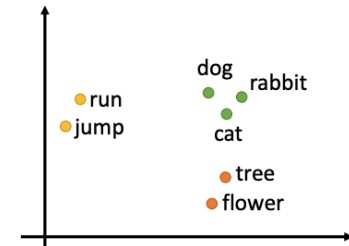
13

One-hot encoding vs word embedding

1-of-N Encoding

apple = [1 0 0 0 0]
 bag = [0 1 0 0 0]
 cat = [0 0 1 0 0]
 dog = [0 0 0 1 0]
 elephant = [0 0 0 0 1]

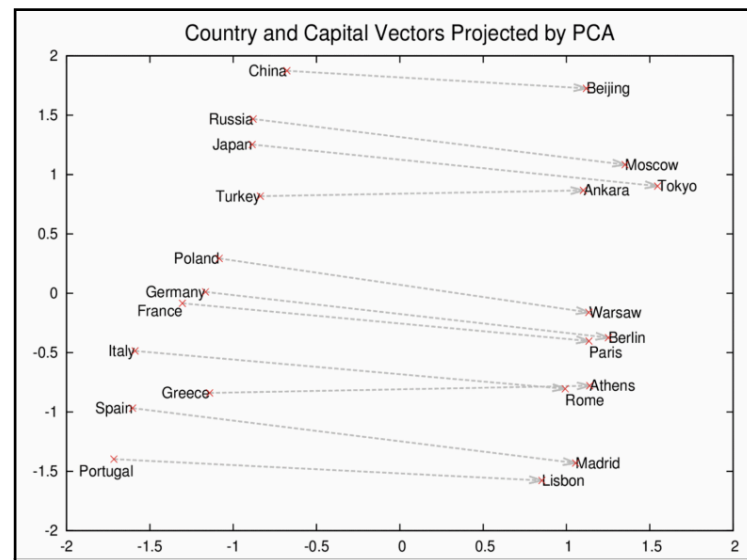
Word Embedding



1/11/21

14

14



15

Word2Vec (training data: Chinese Wikipedia)

- 台灣: 台北 = 法國: ?
⇒ 巴黎 (0.9350)
- 國民黨: 馬英九 = 民進黨: ?
⇒ 陳水扁 (0.9844)
- 海賊王: 魯夫 = 火影忍者: ?
⇒ 自來也 (0.7876)
- 爵士樂: 紐奧良 = 鄉村音樂: ?
⇒ 納許維爾 (0.9718)
- 中研院: 李遠哲 = 工研院: ?
⇒ 林同炎 (0.9381)
- 台灣: 台灣大學 = 美國: ?
⇒ 約翰霍普金斯大學 (0.9976)

1/11/21

16

16

Singular Value Decomposition

$$\underbrace{\begin{bmatrix} * & * & * \\ * & * & * \\ * & * & * \\ * & * & * \end{bmatrix}}_A = \underbrace{\begin{bmatrix} * & * & * & * \\ * & * & * & * \\ * & * & * & * \\ * & * & * & * \end{bmatrix}}_U \underbrace{\begin{bmatrix} \bullet & & \\ & \bullet & \\ & & \bullet \end{bmatrix}}_{\Sigma} \underbrace{\begin{bmatrix} * & * & * \\ * & * & * \\ * & * & * \end{bmatrix}}_{V^T}$$

- U and V are unitary matrices
 - i.e., $U^T U = I$; $V^T V = I$
- Σ is a diagonal matrix
 - The diagonal entries σ_i of Σ are known as the singular values of A , commonly listed in descending order

1/11/21

17

17

Reduced SVD

- If we retain **only k largest** singular values, and set the rest to 0
 - Then Σ is $k \times k$, U is $M \times k$, V^T is $k \times N$, and A_k is $M \times N$
- This is referred to as the **reduced SVD**
- It is the convenient (space-saving) and usual form for computational applications
- Among matrices with rank k , A_k is the best approximation of A (in terms of Frobenius norm)

1/11/21

18

18

A

$$\begin{matrix} M & D_1 & D_2 & D_3 & D_4 & D_5 & D_6 & \dots & D_n \\ \begin{matrix} T_1 \\ T_2 \\ T_3 \\ T_4 \\ T_5 \\ T_6 \\ \vdots \\ T_m \end{matrix} & \begin{pmatrix} 0.00060 & 0.00012 & 0.00003 & 0.00003 & 0.00333 & 0.00048 & \dots & a_{1n} \\ 0 & 0 & 0 & 0 & 0 & 0 & \dots & a_{2n} \\ 0 & 2.98862 & 0 & 0 & 0 & 1.49431 & \dots & a_{3n} \\ 0 & 0 & 0 & 13.32555 & 0 & 0 & \dots & a_{4n} \\ 0 & 0 & 0 & 0 & 0 & 0 & \dots & a_{5n} \\ 1.03442 & 1.03442 & 0 & 0 & 0 & 3.10326 & \dots & a_{6n} \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ a_{m1} & a_{m2} & a_{m3} & a_{m4} & a_{m5} & a_{m6} & \dots & a_{mn} \end{pmatrix} \end{matrix}$$

B

$$U = \begin{bmatrix} C_1 & C_2 & C_3 & \dots & C_m \\ T_1 & a_{11} & a_{12} & a_{13} & \dots & a_{1m} \\ T_2 & a_{21} & a_{22} & a_{23} & \dots & a_{2m} \\ T_3 & a_{31} & a_{32} & a_{33} & \dots & a_{3m} \\ T_4 & a_{41} & a_{42} & a_{43} & \dots & a_{4m} \\ T_5 & a_{51} & a_{52} & a_{53} & \dots & a_{5m} \\ T_6 & a_{61} & a_{62} & a_{63} & \dots & a_{6m} \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ T_m & a_{m1} & a_{m2} & a_{m3} & \dots & a_{mm} \end{bmatrix}$$

$$\Sigma = \begin{bmatrix} D_1 & D_2 & D_3 & \dots & D_n \\ T_1 & a_{11} & 0 & 0 & \dots & 0 \\ T_2 & 0 & a_{22} & 0 & \dots & 0 \\ T_3 & 0 & 0 & a_{33} & \dots & 0 \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ T_m & 0 & 0 & 0 & \dots & a_{mm} \end{bmatrix}$$

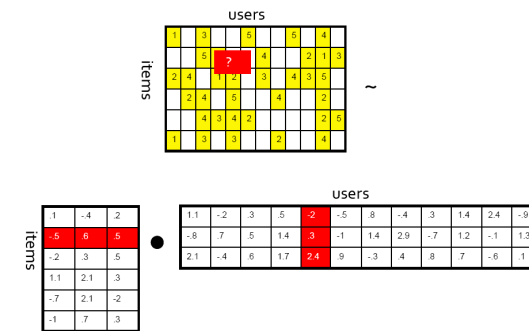
$$V^T = \begin{bmatrix} D_1 & D_2 & D_3 & \dots & D_n \\ C_1 & a_{11} & a_{12} & a_{13} & \dots & a_{1n} \\ C_2 & a_{21} & a_{22} & a_{23} & \dots & a_{2n} \\ C_3 & a_{31} & a_{32} & a_{33} & \dots & a_{3n} \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ C_n & a_{n1} & a_{n2} & a_{n3} & \dots & a_{nn} \end{bmatrix}$$

1/11/21:

19

19

Relationship with matrix factorization (in recommender systems)



1/11/21

20

20

Counting-based vs prediction-based model

- Counting-based
 - Generate the word representation based on the co-occurrence counts
 - Perform matrix factorization
- Prediction-based
 - Predict the words based on the neighboring words
 - The word representation is a by-product of the prediction task

1/11/21

21

21

Two popular types of prediction-based model

- Continuous bag-of-words (CBOW)
 - Predict current word based on the context
- Skip-gram
 - Predict the context words based on the current word
- A word's context

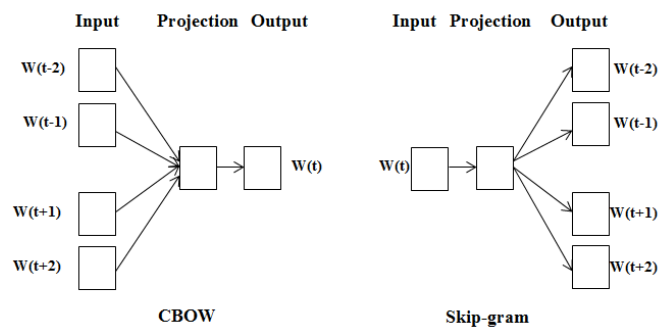


1/11/21

22

22

CBOW vs Skip-gram

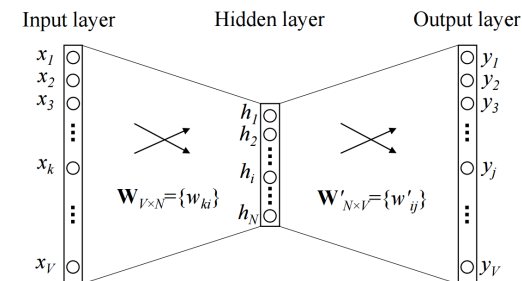


1/11/21

23

23

Simple CBOW



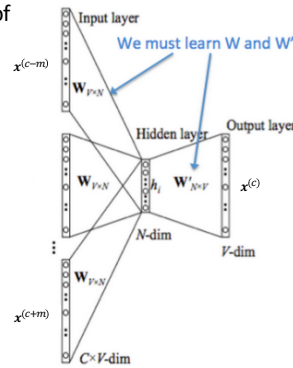
Given (\mathbf{x}, \mathbf{y}) as two neighboring words:

$$\hat{\mathbf{y}} = \text{softmax}(\mathbf{W}'^T (\mathbf{W}^T \mathbf{x}))$$

24

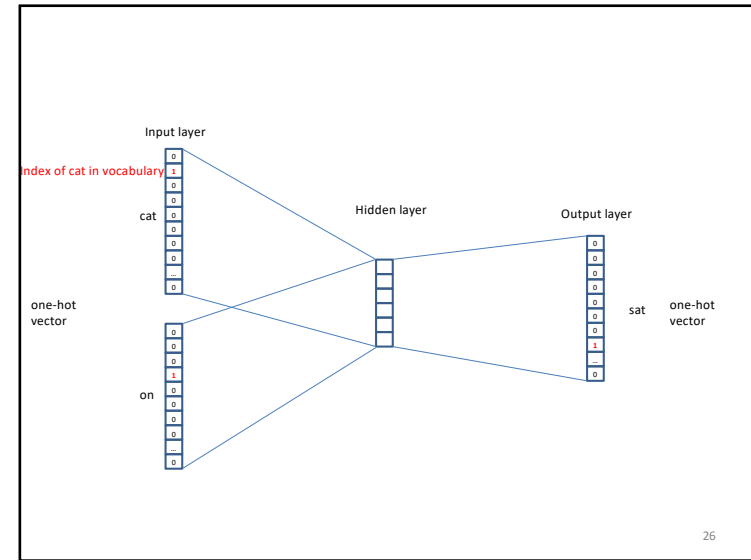
Continuous bag-of-words (CBOW)

- Called bag-of-words, because the order of words are lost
- Forward
 - generate one-hot word vectors $(x^{(c-m)}, \dots, x^{(c-1)}, x^{(c+1)}, \dots, x^{(c+m)})$ for the input context of size m
 - generate word vectors of the context $(v^{(c-m)}, \dots, v^{(c-1)}, v^{(c+1)}, \dots, v^{(c+m)})$, where $v^{(i)} = W^T x^{(i)}$
 - $\tilde{v} = \frac{v^{(c-m)} + \dots + v^{(c-1)} + v^{(c+1)} + \dots + v^{(c+m)}}{2m}$
 - $z = W'^T \tilde{v}$
 - $\hat{x}^{(c)} = \text{softmax}(z)$
- Backprop
 - Compute error and perform back-prop to update W and W'

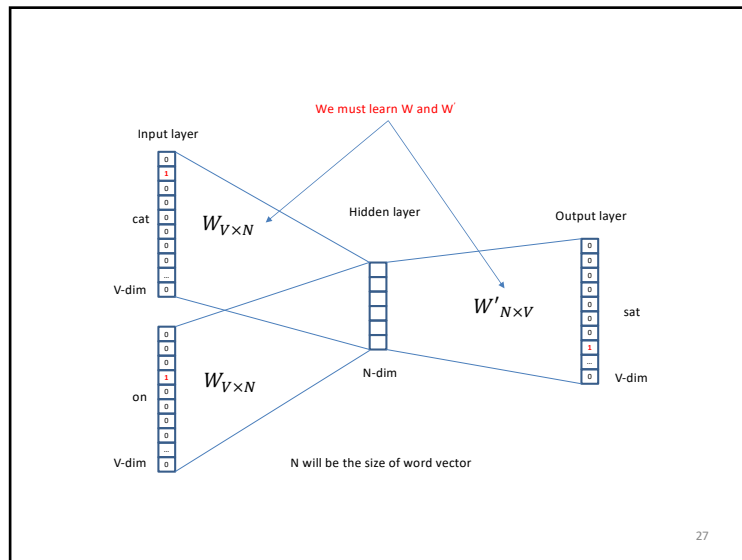


1/11/21

25

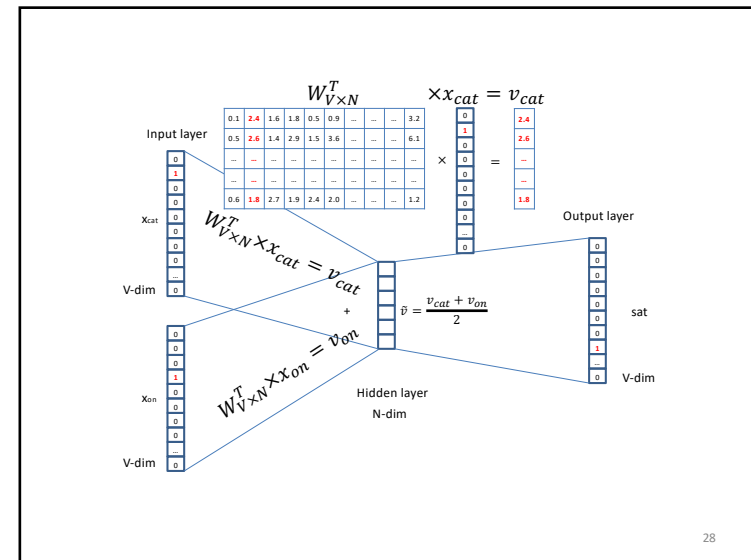


26



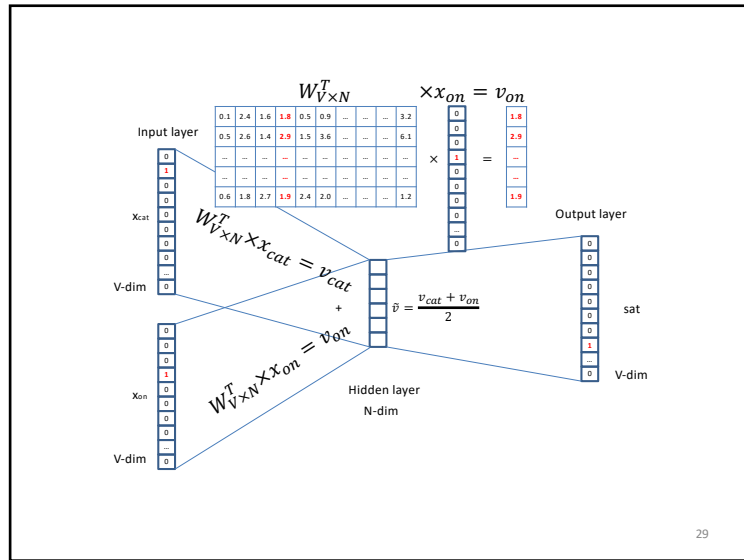
27

27

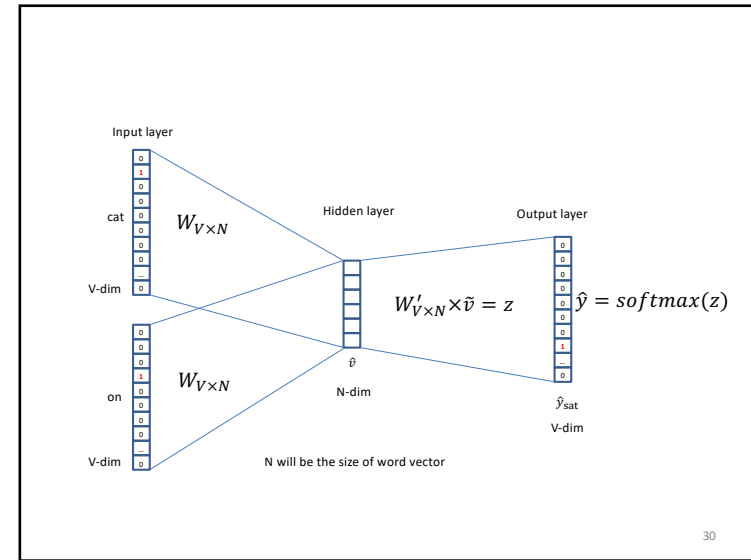


28

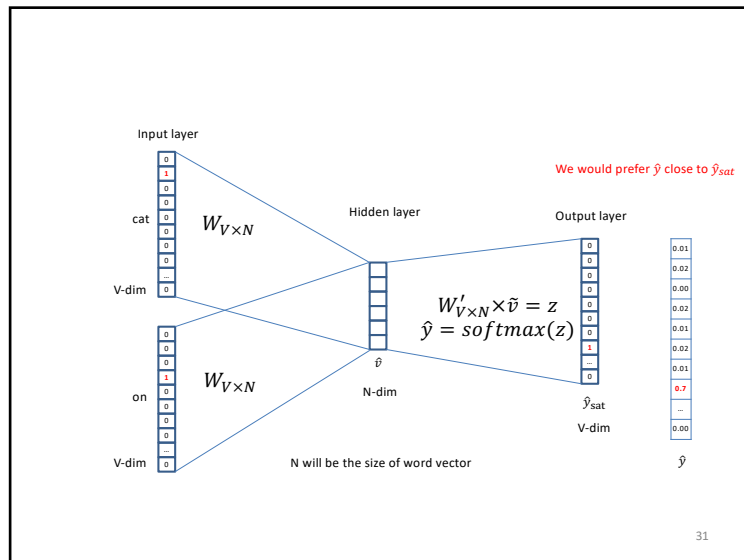
28



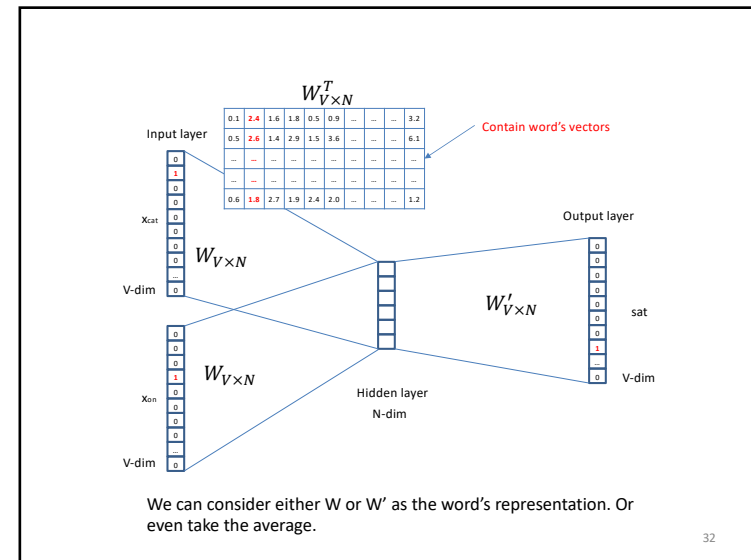
29



30



31



32

Skip-gram

- Predict context words by the current word
- Or more precisely, if there are k context words, we predict one context word based on the current word for k times
- E.g., “this is a wonderful day”, window size=1
 - When current word is “is”, the context words are “this” and “a”
 - Use “is” to predict “this”
 - Use “is” to predict “a”
 - When the current word is “wonderful”, the context words are “a” and “day”
 - Use “wonderful” to predict “a”
 - Use “wonderful” to predict “day”

1/11/21

33

33

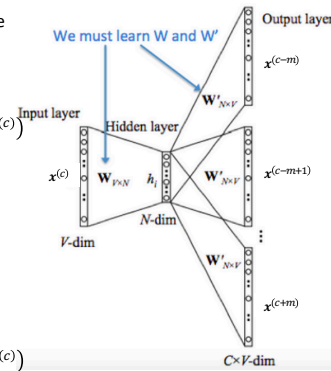
Skip-gram (cont')

- For each (target, context) pair $x^{(c)}$ and $x^{(d)}$, we want:

$$x^{(d)} \approx \text{softmax}(W'^T W x^{(c)})$$

- Objective:

$$\begin{aligned} J &= -\log P(x^{(c-m)}, \dots, x^{(c-1)}, x^{(c+1)}, \dots, x^{(c+m)} | x^{(c)}) \\ &= -\log \prod_{j=0, j \neq m}^{2m} P(x^{(c-m+j)} | x^{(c)}) \\ &= -\log \prod_{j=0, j \neq m}^{2m} P(u^{(c-m+j)} | v^{(c)}) \\ &= -\log \prod_{j=0, j \neq m}^{2m} \frac{\exp(u^{(c-m+j)T} v^{(c)})}{\sum_{k=1}^{|V|} \exp(u_k^T v^{(c)})} \\ &= - \sum_{j=0, j \neq m}^{2m} u^{(c-m+j)T} v^{(c)} + 2m \log \sum_{k=1}^{|V|} \exp(u_k^T v^{(c)}) \end{aligned}$$



1/11/21

34

34

Negative sampling

- In the objective function, summing over $|V|$ is a huge cost
 - Every update and every objective function evaluation take $O(|V|)$
- Can we approximate the term? → Negative sampling!
 - “Sample” negative examples from $P_n(w)$, whose probability matches the ordering of the frequency of the vocabulary

1/11/21

35

35

Turning prediction function of skip gram into logistic regression

- Positive instances: (w, c) pairs where w is a word and c is a context word
- Negative instances: (w, c') where c' is from negative sampling
- Training: find v_w , v_c , and $v_{c'}$ such that

$$P(w, c) = \frac{1}{1 + \exp(-v_c^T v_w)} \approx 1, \text{ and}$$

$$P(w, c') = \frac{1}{1 + \exp(-v_{c'}^T v_w)} \approx 0$$

1/11/21

36

36

We may apply similar concept to other domains

- As long as you can identify the “relevant object pairs”, it is possible to apply similar concept (X2vec) to your domain

1/11/21

37

37

Example: Prod2Vec

- Build a recommendation algorithm for a e-commerce website based on clickstream
- Generate the embedding for the products and make recommendations based on the distance between product embeddings

1/11/21

38

38

Example: Node2Vec

- Generate node embeddings for the nodes in a graph

1/11/21

39

39

Example: Behavior2Vec

- On a e-commerce website, “viewing” a product X and “purchasing” a product X are very different behaviors
- Instead of generating embedding for products, we generate two embeddings for each product
 - Embedding of “buying” a product
 - Embedding of “viewing” a product
- Use different embeddings in different scenarios

Behavior2Vec: Generating Distributed Representations of Users' Behaviors on Products for Recommender Systems
Hung-Hsuan Chen TKDD ACM Transactions on Knowledge Discovery from Data 12(4), 2018

1/11/21

40

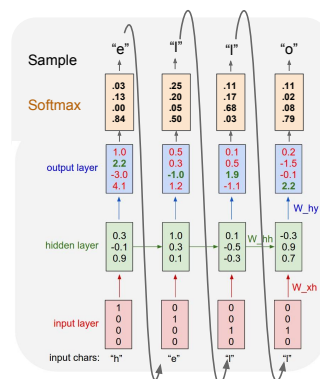
40

Previously, we use “one-hot-encoded” words as input for a NN

**Example:
Character-level
Language Model
Sampling**

Vocabulary:
[h,e,l,o]

At test-time sample
characters one at a time,
feed back to model



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 40 May 4, 2017

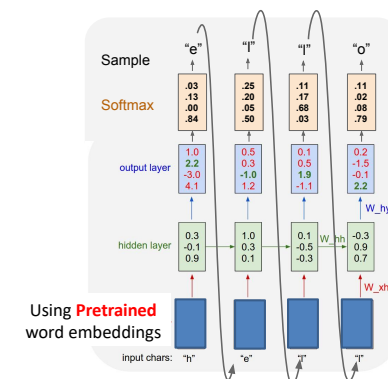
41

But we can use word embedding as the input for a NN

**Example:
Character-level
Language Model
Sampling**

Vocabulary:
[h,e,l,o]

At test-time sample
characters one at a time,
feed back to model



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 10 - 40 May 4, 2017

42

Summary

- Word2Vec generates dense vector (instead of one-hot encoded vector) for each word
- Relationship between words may be captured by Word2Vec
- Word embeddings are sometimes pre-trained and regarded as the input features to other tasks
- We may apply similar concept to many fields

1/11/21

43

43