

CSE 4392 Special Topics Natural Language Processing

1

Word Embedding

2025 Spring

HOW TO REPRESENT WORDS? • N-gram language models: $P(w \mid \text{it is 76 F and})$ *It is 76 F and* _____. [0.0001, 0.1, 0, 0, 0.002, ..., 0.3, ..., 0] red sunny • Text classification: $P(y = 1 | \mathbf{x}) = \sigma(\mathbf{w} \cdot \mathbf{x} + b)$ $\stackrel{\checkmark}{=} x^{(1)} [0, 1, 0, 0, 0, ..., 1, ..., 1]$ I like this movie. I don't like this movie. $\P_{x^{(2)}}$ [0, 1, 0, 1, 0, ..., 1, ..., 1] don't *Feature: don't*

REPRESENTING WORDS AS DISCRETE SYMBOLS

• In traditional NLP, we regard words as discrete symbols: hotel, conference, motel — a localist representation

one 1, the rest O's

 • Words can be represented by one-hot vectors: hotel = [00000000000000] motel = [00010000000000]
 • Vector dimension = num of words in vocabulary (e.g., 500,000)

• No way to encode similarity between words!

QUIZ: ONE-HOT VECTORS

• Using one-hot vectors to represent words, why is there no way to encode similarity between the words?

Represent Words by their Context

• **Distributional hypothesis**: words that occur in similar contexts tend to have similar meanings



J.R. Firth 1957

"You shall know a word by the company it keeps" One of the most successful ideas of modern statistical NLP!

...government debt problems turning into **banking** crises as happened in 2009...

...saying that Europe needs unified **banking** regulation to replace the hodgepodge...

...India has just given its banking system a shot in the arm...

These context words will represent banking.

DISTRIBUTIONAL HYPOTHESIS

"tejuino"



• C1: A bottle of _____ is on the table.

- C2: Everybody likes ____.
- C3: Don't have <u>before you drive</u>.
- C4: We make ____ out of corn.

DISTRIBUTIONAL HYPOTHESIS

C1: A bottle of _____ is on the table. C2: Everybody likes ____.

C3: Don't have ____ before you drive.

C4: We make ____ out of corn.

	C1	C2	C3	C4
tejuino	1	1	1	1
loud	0	0	0	0
motor-oil	1	0	0	0
tortillas	0	1	0	1
choices	0	1	0	0
wine	1	1	1	0

"words that occur in similar contexts tend to have similar meanings"

WORDS AS VECTORS

- We'll build a new model of meaning focusing on similarity
 - Each word is a vector
 - Similar words are "nearby in space"
- A first solution: we can just use context vectors to represent the meaning of words!

• Word-word co-occurrence matrix:

	aardvark	computer	data	pinch	result	sugar	
apricot	0	0	0	1	0	1	
pineapple	0	0	0	1	0	1	
digital	0	2	1	0	1	0	
information	0	1	6	0	4	0	





The range of $\cos(.)$ is [0, 1]

Quiz: Compute the cosine similarity between "digital" [1, 1] and "information" [6, 4].

WORDS AS VECTORS

- Problem: using raw frequency counts is not always very good..
 - Solution: let's weight the counts!
 - PPMI = Positive Pointwise Mutual Information

$PPMI(w,c) = \max(\log_2 \frac{1}{P})$	$\frac{P(w,c)}{P(w)P(c)},0)$
--	------------------------------

w and c are two words in the same text window

	computer	data	result	pie	sugar
cherry	2	8	9	442	25
strawberry	0	0	1	60	19
digital	1670	1683	85	5	4
information	3325	3982	378	5	13

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

https://gyan-mittal.com/nlp-ai-ml/nlp-word-embedding-svd-based-methods/

Sparse vs. Dense Vectors

• Still, the vectors we get from word-word occurrence matrix are sparse (most are 0's) & long (vocabulary size)

- Alternative: we want to represent words as **short** (50-300 dimensional) & **dense** (real-valued) vectors
 - The focus of this lecture
 - The basis of all the modern NLP systems

DENSE VECTORS



WHY DENSE VECTORS?

- Short vectors are easier to use as features in ML systems
- Dense vectors may generalize better than storing explicit counts
- They do better at capturing synonymy
 - *w*₁ co-occurs with "car", *w*₂ co-occurs with "automobile"
 → *w*₁ and *w*₂ are synonyms

DIFFERENT METHODS FOR GETTING DENSE VECTORS

• Singular value decomposition (SVD)

$$\begin{bmatrix} X \\ X \end{bmatrix} = \begin{bmatrix} W \\ W \end{bmatrix} \begin{bmatrix} \sigma_1 & 0 & 0 & \dots & 0 \\ 0 & \sigma_2 & 0 & \dots & 0 \\ 0 & 0 & \sigma_3 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & \sigma_k \end{bmatrix} \begin{bmatrix} C \\ k \times |V| \end{bmatrix}$$

- Limitations:
 - Vocab usually big, so matrix is big and difficult to train
 - Imbalance due to high frequency words
 - As new words are added to the vocab, matrix size changes, need to re-train!
- <u>https://gyan-mittal.com/nlp-ai-ml/nlp-word-embedding-svd-based-methods/</u>
- word2vec and friends: "learn" the vectors!
 - Much more elegant than SVD

WORD2VEC AND FRIENDS



• (Mikolov et al, 2013): Distributed Representations of Words and Phrases and their Compositionality



WORD2VEC

- Input: a large text corpus, V, d
 - V: a pre-defined vocabulary
 - d: dimension of word vectors (e.g. 300)
 - Text corpora:
 - Wikipedia + Gigaword 5: 6B
 - o Twitter: 27B
 - o Common Crawl: 840B

• Output:

$$f: V \to \mathbb{R}^d$$

$$v_{\text{cat}} = \begin{pmatrix} -0.224\\ 0.130\\ -0.290\\ 0.276 \end{pmatrix} \quad v_{\text{dog}} = \begin{pmatrix} -0.124\\ 0.430\\ -0.200\\ 0.329 \end{pmatrix}$$
$$v_{\text{dog}} = \begin{pmatrix} 0.430\\ 0.200\\ 0.329 \end{pmatrix} \quad v_{\text{language}} = \begin{pmatrix} 0.290\\ -0.441\\ 0.762\\ 0.982 \end{pmatrix}$$

(0191)

WORD2VEC

	Word	Cosine distance
	norway	0.760124
• Word = "sweden"	finland	0.620022
	belgium	0.585835
	netherlands iceland	0.574631 0.562368
	estonia slovenia	0.547621 0.531408

WORD2VEC

INPUT PROJECTION OUTPUT INPUT PROJECTION OUTPUT



SKIP-GRAM

- Idea: Use the *center* word to predict its *context* words.
- Context: a fixed window of 2m



SKIP-GRAM

• Next time step *t*:



SKIP-GRAM: OBJECTIVE FUNCTION

• For each position t = 1, 2, ..., T, predict context words within context size m, given center word w_i :

all the parameters to be optimized
$$\mathcal{L}(heta) = \prod_{t=1}^T \prod_{-m \leq j \leq m, j \neq 0} P(w_{t+j} \mid w_t; heta)$$

• The objective function $J(\theta)$ is the (average) negative log likelihood:

$$J(\theta) = -\frac{1}{T} \log \mathcal{L}(\theta) = -\frac{1}{T} \sum_{t=1}^{T} \sum_{-m \le j \le m, j \ne 0}^{T} \log P(w_{t+j} \mid w_t; \theta)$$

HOW TO DEFINE: $P(W_{T+J} | W_T; \Theta)$?

• There are two word embeddings (vectors) for each word in the vocabulary:

 $\mathbf{u}_i \in \mathbb{R}^d$: embedding for target word i $\mathbf{v}_{i'} \in \mathbb{R}^d$: embedding for context word i'

• Use inner product $u_i \cdot v_{i'}$ to measure how likely a word *i* appears with context word *i*', the larger the better.

"softmax" we learned last time!

$$P(w_{t+j} \mid w_t) = \frac{\exp(\mathbf{u}_{w_t} \cdot \mathbf{v}_{w_{t+j}})}{\sum_{k \in V} \exp(\mathbf{u}_{w_t} \cdot \mathbf{v}_k)}$$

 $\theta = \{\{\mathbf{u}_k\}, \{\mathbf{v}_k\}\}$ are all the parameters in this model!

HOW TO TRAIN THE MODEL

Λ

• Calculate all the gradients together!

$$\theta = \{\{\mathbf{u}_k\}, \{\mathbf{v}_k\}\}$$
$$J(\theta) = -\frac{1}{T} \sum_{t=1}^T \sum_{-m \le j \le m, j \ne 0} \log P(w_{t+j} \mid w_t; \theta) \quad \nabla_{\theta} J(\theta) = ?$$

Quiz: Suppose the vocab is *V* in the corpus, how many parameters in θ to train in total?

• We can apply stochastic gradient descent (SGD)!

(() ())

$$\theta^{(t+1)} = \theta^{(t)} - \eta \nabla_{\theta} J(\theta)$$

• Let's go through the math.

WARM-UP

$$f(x) = \exp(x)$$

$$f(x) = \log(x)$$

 $f(x) = f_1(f_2(x))$

 $f(\mathbf{x}) = \mathbf{x} \cdot \mathbf{a}$

$$\frac{df}{dx} = \exp(x)$$
$$\frac{df}{dx} = \frac{1}{x}$$

chain rule:

$$\frac{df}{dx} = \frac{df_1(z)}{dz} \frac{df_2(x)}{dx} \quad z = f_2(x)$$

$$\frac{\partial f}{\partial \mathbf{x}} = \mathbf{a}$$

$$\frac{\partial f}{\partial \mathbf{x}} = [\frac{\partial f}{\partial x_1}, \frac{\partial f}{\partial x_2}, \dots, \frac{\partial f}{\partial x_n}]$$

COMPUTING THE GRADIENTS

 $\partial \mathbf{v}_k$

• Consider one pair of target/context words (t, c):

$$\begin{split} y &= -\log\left(\frac{\exp(\mathbf{u}_t \cdot \mathbf{v}_c)}{\sum_{k \in V} \exp(\mathbf{u}_t \cdot \mathbf{v}_k)}\right) \\ \frac{\partial y}{\partial \mathbf{u}_t} &= \frac{\partial \left(-\mathbf{u}_t \cdot \mathbf{v}_c + \log\left(\sum_{k \in V} \exp\left(\mathbf{u}_t \cdot \mathbf{v}_k\right)\right)\right)}{\partial \mathbf{u}_t} \\ &= -\mathbf{v}_c + \frac{\sum_{k \in V} \frac{\partial \exp(\mathbf{u}_t \cdot \mathbf{v}_k)}{\partial \mathbf{u}_t}}{\sum_{k \in V} \exp\left(\mathbf{u}_t \cdot \mathbf{v}_k\right)} \\ &= -\mathbf{v}_c + \frac{\sum_{k \in V} \exp\left(\mathbf{u}_t \cdot \mathbf{v}_k\right)\mathbf{v}_k}{\sum_{k \in V} \exp\left(\mathbf{u}_t \cdot \mathbf{v}_k\right)} \\ &= -\mathbf{v}_c + \sum_{k \in V} P(k \mid t)\mathbf{v}_k \\ \frac{\partial y}{\partial \mathbf{v}_t} &= -1(k = c)\mathbf{u}_t + P(k \mid t)\mathbf{u}_t \quad \text{(Try to derive this!)} \end{split}$$

PUTTING IT TOGETHER

- o Input: text corpus, context size m, embedding size d, vocab V
- Initialize $\boldsymbol{u}_i, \boldsymbol{v}_i$ randomly
- Work through the training corpus and collection training data (*t*, *c*):
 - Update:

$$\begin{split} \mathbf{u}_{t} &\leftarrow \mathbf{u}_{t} - \eta \frac{\partial y}{\partial \mathbf{u}_{t}} \\ \mathbf{v}_{k} &\leftarrow \mathbf{v}_{k} - \eta \frac{\partial y}{\partial \mathbf{v}_{k}}, \forall k \in V \end{split}$$

• Any problem here?

SKIP-GRAM WITH NEGATIVE SAMPLING (SGNS)

• Problem: for each training data pair (*t*, *c*), you need to update the embedding of every word in the vocab. That's too expensive!

$$\frac{\partial y}{\partial \mathbf{u}_t} = -\mathbf{v}_c + \sum_{k \in V} P(k \mid t) \mathbf{v}_k$$
$$\frac{\partial y}{\partial \mathbf{v}_k} = -1(k = c) \mathbf{u}_t + P(k \mid t) \mathbf{u}_t$$

• *Negative Sampling*: instead of considering all the words in V, let's randomly sample K (5-20) negative examples.

softmax:
$$y = -\log\left(\frac{\exp(\mathbf{u}_t \cdot \mathbf{v}_c)}{\sum_{k \in V} \exp(\mathbf{u}_t \cdot \mathbf{v}_k)}\right)$$

NS:
$$y = -\log(\sigma(\mathbf{u}_t \cdot \mathbf{v}_c)) - \sum_{i=1}^{K} \mathbb{E}_{j \sim P(w)} \log(\sigma(-\mathbf{u}_t \cdot \mathbf{v}_j))$$

SKIP-GRAM WITH NEGATIVE SAMPLING (SGNS)

$$y = -\log(\sigma(\mathbf{u}_t \cdot \mathbf{v}_c)) - \sum_{i=1}^{K} \mathbb{E}_{j \sim P(w)} \log(\sigma(-\mathbf{u}_t \cdot \mathbf{v}_j))$$

positive examples +

с

apricot tablespoon

apricot of

apricot a

apricot jam

negative examples -

tctcapricotaardvarkapricotsevenapricotmyapricotforeverapricotwhereapricotdearapricotcoaxialapricotif

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



• Same as training a **logistic regression** for binary classification!

$$P(D=1 \mid t, c) = \sigma(\mathbf{u}_t \cdot \mathbf{v}_c)$$

• Compute the gradient in assignment!

CONTINUUS BAG OF WORDS (CBOW)

INPUT PROJECTION OUTPUT



GLOVE: GLOBAL VECTORS

• Skip-gram and CBoW uses local context

- Slow to train when the corpus is very large
- Let's take the global co-occurrence statistics X_{ij} : $J = \sum_{i,j=1}^{V} f(X_{ij}) \left(w_i^T \tilde{w}_j + b_i + \tilde{b}_j - \log X_{ij} \right)^2$
 - X_{ij} tabulate the number of times word *j* occurs in the context of word *i*.



GLOVE: GLOBAL VECTORS

- Nearest word to frog:
 - frogs
 - toad
 - litoria
 - leptodactylidae
 - rana
 - lizard
 - eleutherodactylus



litoria



leptodactylidae



rana



eleutherodactylus



(Pennington et al, 2014): GloVe: Global Vectors for Word Representation

FASTTEXT: SUB-WORD EMBEDDINGS

Similar as Skip-gram, but break words into n-grams with n = 3 to 6

where: 3-grams: <wh, whe, her, ere, re>
 4-grams: <whe, wher, here, ere>
 5-grams: <wher, where, here>
 6-grams: <where, where>

```
• Replace \boldsymbol{u}_i \cdot \boldsymbol{v}_j by
```



• More to come: contextualized word embeddings

(Bojanowski et al, 2017): Enriching Word Vectors with Subword Information



PRE-TRAINED WORD EMBEDDINGS AVAILABLE

- word2vec: <u>https://code.google.com/archive/p/word2vec/</u>
- GloVe: <u>https://nlp.stanford.edu/projects/glove/</u>
- FastText: <u>https://fasttext.cc/</u>

NLPL word embeddings repository

brought to you by Language Technology Group at the University of Oslo

We feature models trained with clearly stated hyperparametes, on clearly described and linguistically pre-processed corpora. More information and hints at the NLPL wiki page. You can also download the JSON file containing metadata for all the models in the repository.

Filter your search by:

Language

Select one or more languages:

English [eng] (models: 43) Estonian [est] (models: 2) Basque [eus] (models: 2) Persian [fas] (models: 2)

Algorithms:

Global Vectors BERT Embeddings from Language Models (ELMo) CfastText Skipgram Word2Vec Continuous Skipgram Gensim Continuous Skipgram

Lemmatization:

✓True ✓False

• Differ in algorithms, text corpora, dimensions, cased/uncased...

EVALUATING WORD EMBEDDINGS

EXTRINSIC VS INTRINSIC EVALUATION

• Extrinsic evaluation

- Plug these word embeddings into a real NLP system and see if this improves the performance
- Could take a long time but still the most important evaluation metric
- Intrinsic evaluation
 - Evaluate on a specific/intermediate subtask
 - Fast to compute
 - Not clear if it really helps the downstream task



INTRINSIC EVALUATION

• Word similarity task

- Example dataset: WordSim-353 353 pairs of words with human judgement
- <u>http://www.cs.technion.ac.il/~gabr/resources/data/wordsim353/</u>

Word 1	Word 2	Human (mean)
tiger	cat	7.35
tiger	tiger	10
book	paper	7.46
computer	internet	7.58
plane	car	5.77
professor	doctor	6.62
stock	phone	1.62
stock	CD	1.31
stock	jaguar	0.92

Cosine similarity:

$$\cos(oldsymbol{u}_i,oldsymbol{u}_j) = rac{oldsymbol{u}_i\cdotoldsymbol{u}_j}{||oldsymbol{u}_i||_2 imes||oldsymbol{u}_j||_2}$$

Metric: Spearman Rank Correlation

INTRINSIC EVALUATION

• Word Similarity Results (Spearman correlations):

Model	Size	WS353	MC	RG	SCWS	RW
SVD	6B	35.3	35.1	42.5	38.3	25.6
SVD-S	6B	56.5	71.5	71.0	53.6	34.7
SVD-L	6B	65.7	<u>72.7</u>	75.1	56.5	37.0
$CBOW^{\dagger}$	6B	57.2	65.6	68.2	57.0	32.5
SG^\dagger	6B	62.8	65.2	69.7	<u>58.1</u>	37.2
GloVe	6B	<u>65.8</u>	<u>72.7</u>	<u>77.8</u>	53.9	<u>38.1</u>
SVD-L	42B	74.0	76.4	74.1	58.3	39.9
GloVe	42B	<u>75.9</u>	<u>83.6</u>	<u>82.9</u>	<u>59.6</u>	<u>47.8</u>
CBOW*	100B	68.4	79.6	75.4	59.4	45.5

INTRINSIC EVALUATION

• Word analogy

```
man:woman \approx king: ?
```

$$rg\max_i \left(\cos(\mathbf{u}_i, \mathbf{u}_b - \mathbf{u}_a + \mathbf{u}_c) \right)$$

Semantic

Syntactic

Austin:Texas ≈ ? :California

bad:worst \approx hot: ?

More examples at:

http://download.tensorflow.org/data/questions-words.txt