



**CSE 4392 SPECIAL TOPICS  
NATURAL LANGUAGE PROCESSING**

# Word Embedding

1

2024 Spring

# HOW TO REPRESENT WORDS?

- N-gram language models:

*It is 76 F and \_\_\_\_\_.*


$$P(w \mid \text{it is 76 F and})$$


$$[0.0001, 0.1, 0, 0, 0.002, \dots, 0.3, \dots, 0]$$

*red**sunny*

- Text classification:

$$P(y = 1 \mid \mathbf{x}) = \sigma(\mathbf{w} \cdot \mathbf{x} + b)$$

*I like this movie.*   $\mathbf{x}^{(1)}$  [0, 1, 0, 0, 0, ..., 1, ..., 1]

*I don't like this movie.*   $\mathbf{x}^{(2)}$  [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 0's



- Words can be represented by **one-hot** vectors:

**hotel** = [00000000000010000]

**motel** = [00010000000000000]

- 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 \_\_\_ before you drive.
- 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”

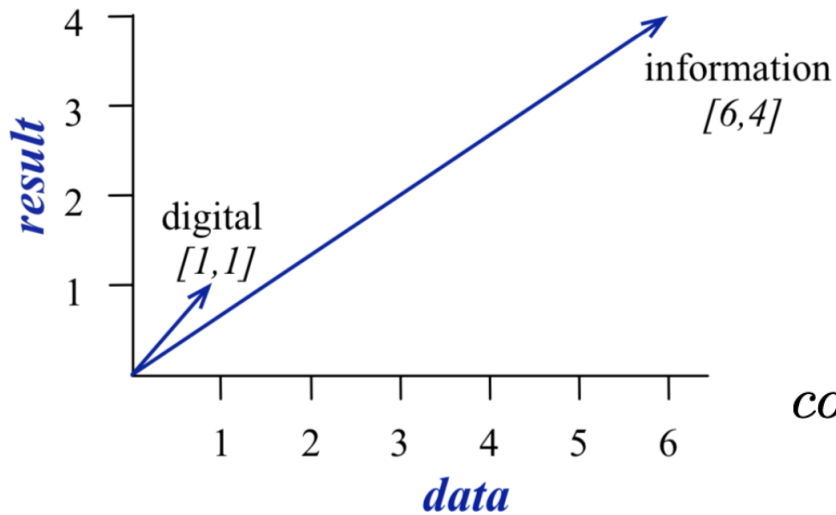
# WORDS AS VECTORS

- 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	



# WORDS AS VECTORS



$$\cos(\mathbf{u}, \mathbf{v}) = \frac{\mathbf{u} \cdot \mathbf{v}}{\|\mathbf{u}\| \|\mathbf{v}\|}$$

$$\cos(\mathbf{u}, \mathbf{v}) = \frac{\sum_{i=1}^V u_i v_i}{\sqrt{\sum_{i=1}^V u_i^2} \sqrt{\sum_{i=1}^V v_i^2}}$$

The range of  $\cos(\cdot)$  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

$$\text{PPMI}(w, c) = \max\left(\log_2 \frac{P(w, c)}{P(w)P(c)}, 0\right)$$

*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

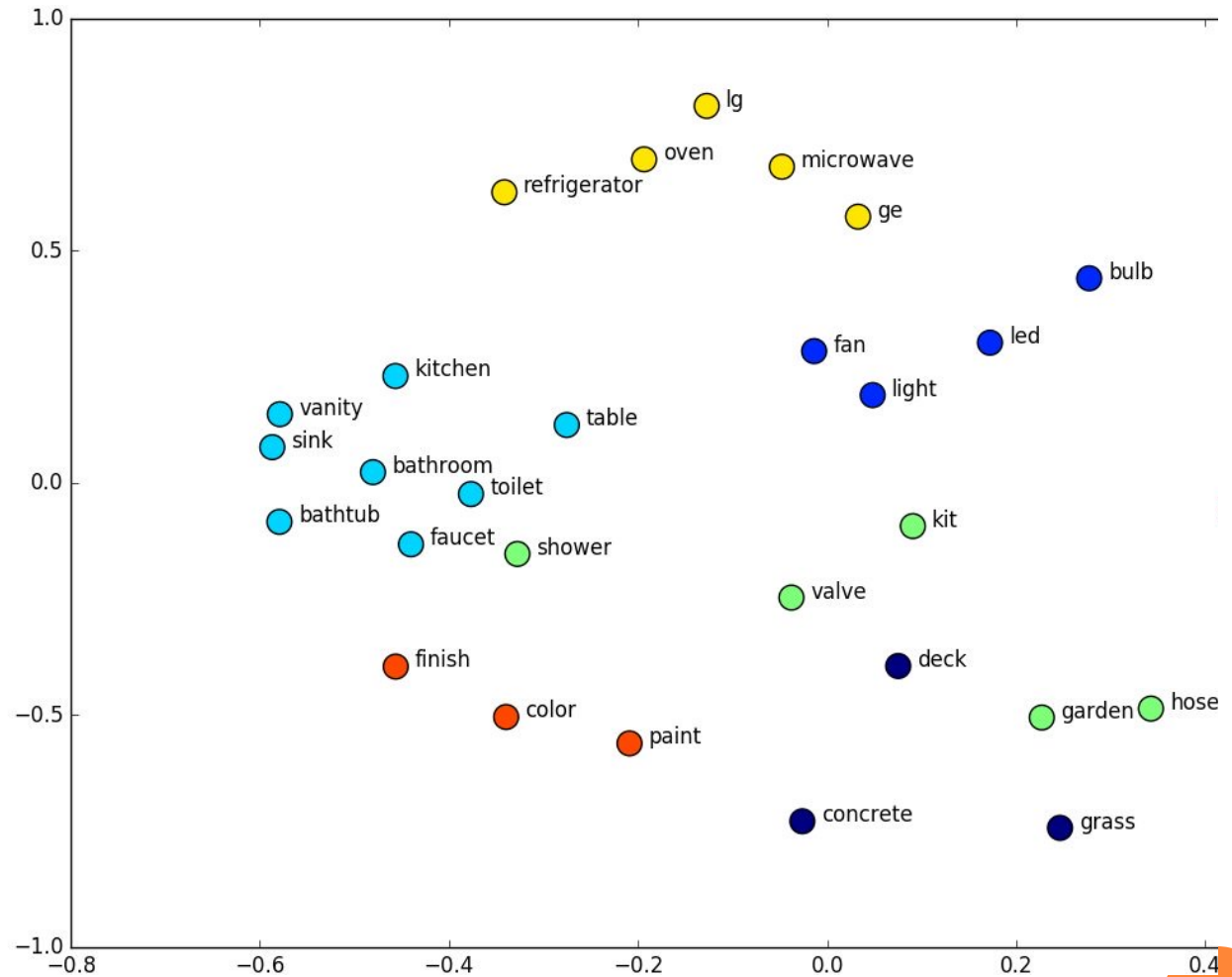
	computer	data	result	pie	sugar
cherry	0	0	0	4.38	3.30
strawberry	0	0	0	4.10	5.51
digital	0.18	0.01	0	0	0
information	0.02	0.09	0.28	0	0

# 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

$$\text{faucet} = \begin{pmatrix} 0.246 \\ 0.793 \\ -0.177 \\ -2.212 \\ 12.322 \\ 3.292 \\ -0.768 \\ 3.211 \end{pmatrix}$$



# 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

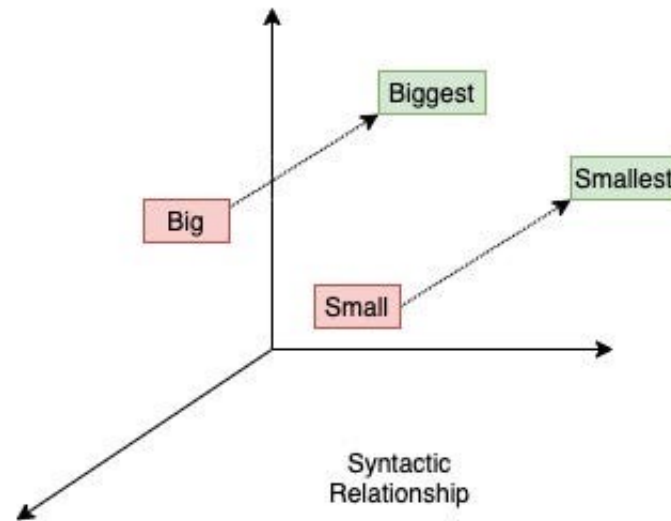
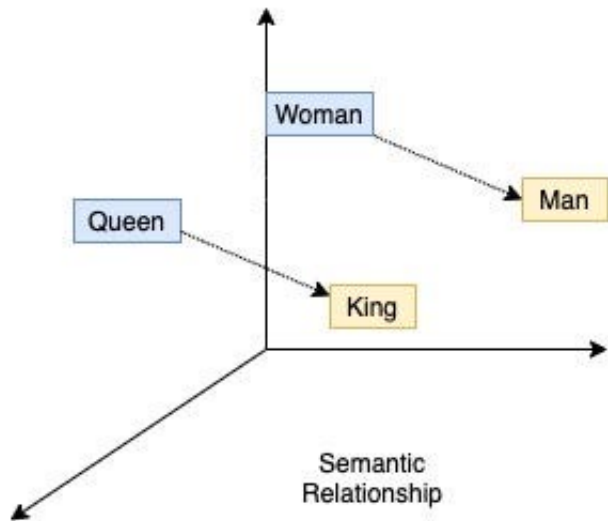
# DIFFERENT METHODS FOR GETTING DENSE VECTORS

- Singular value decomposition (SVD)

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

- 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!
- <https://gyan-mittal.com/nlp-ai-ml/nlp-word-embedding-svd-based-methods/>
- 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
    - Twitter: 27B
    - Common Crawl: 840B

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

$$v_{\text{the}} = \begin{pmatrix} 0.234 \\ 0.266 \\ 0.239 \\ -0.199 \end{pmatrix} \quad v_{\text{language}} = \begin{pmatrix} 0.290 \\ -0.441 \\ 0.762 \\ 0.982 \end{pmatrix}$$

- Output:

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



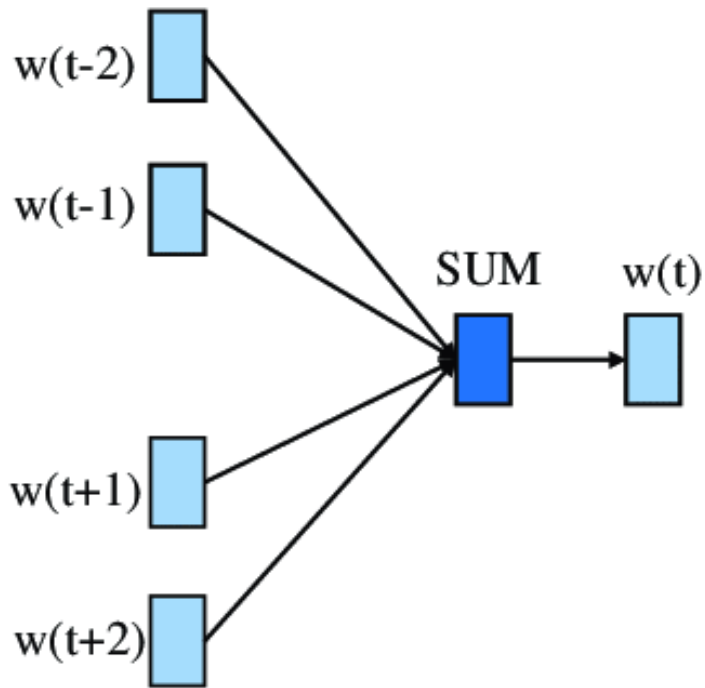
# WORD2VEC

- Word = “sweden”

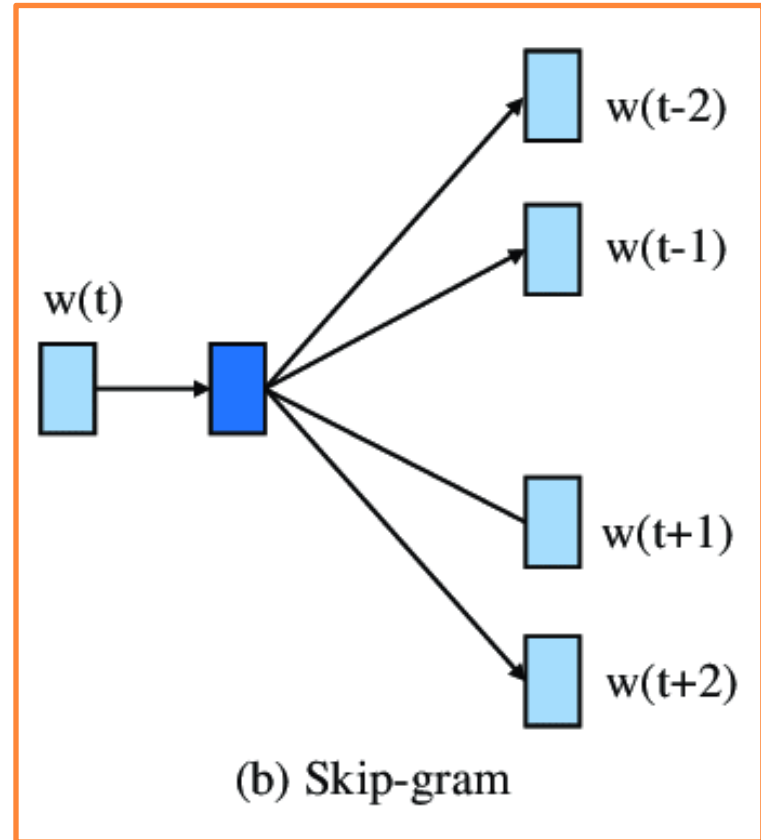
Word	Cosine distance
norway	0.760124
denmark	0.715460
finland	0.620022
switzerland	0.588132
belgium	0.585835
netherlands	0.574631
iceland	0.562368
estonia	0.547621
slovenia	0.531408

# WORD2VEC

INPUT PROJECTION OUTPUT INPUT PROJECTION OUTPUT



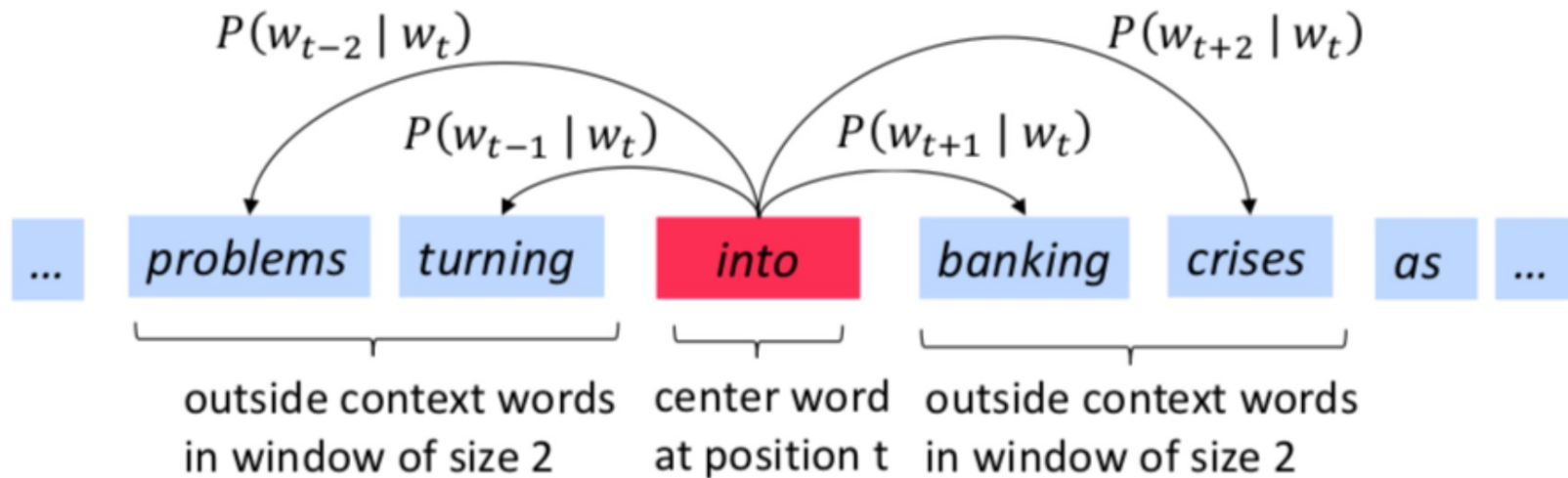
(a) CBOW



(b) Skip-gram

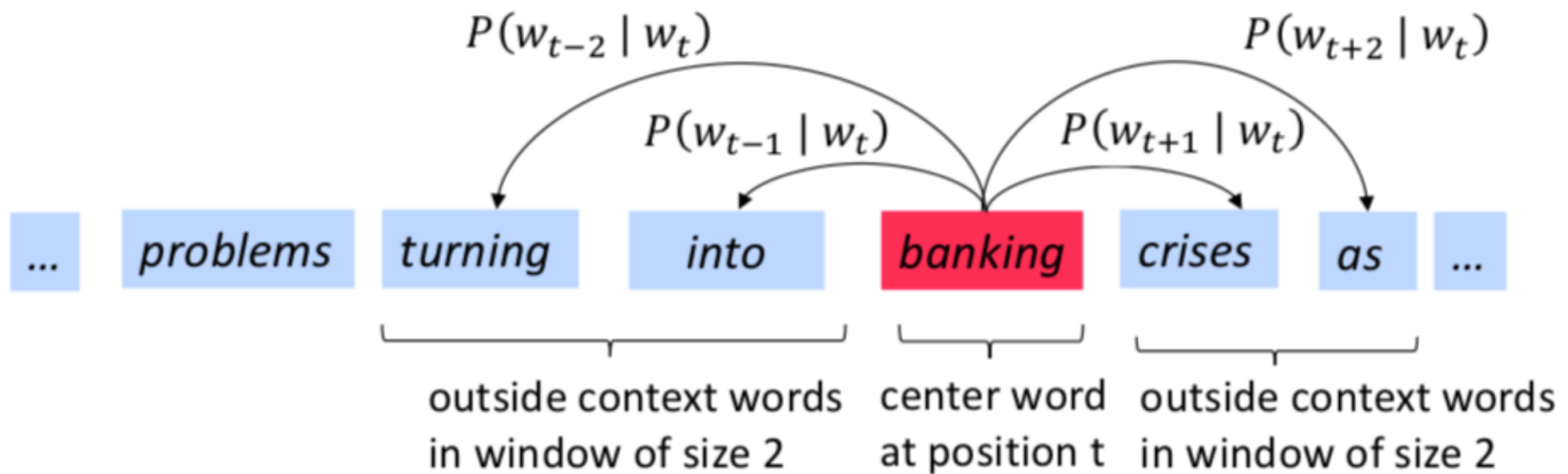
# SKIP-GRAM

- 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  $t = 1, 2, \dots, T$ , predict context words within context size  $m$ , given center word  $w_j$ :

$$\mathcal{L}(\theta) = \prod_{t=1}^T \prod_{-m \leq j \leq m, j \neq 0} P(w_{t+j} \mid w_t; \theta)$$

all the parameters to be optimized 

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

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

## HOW TO DEFINE: $P(w_{T+J} \mid w_T; \Theta)$ ?


- 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  $\mathbf{u}_i \cdot \mathbf{v}_{i'}$  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

- Calculate all the gradients together!

$$\theta = \{\{\mathbf{u}_k\}, \{\mathbf{v}_k\}\}$$

$$J(\theta) = -\frac{1}{T} \sum_{t=1}^T \sum_{-m \leq j \leq m, j \neq 0} \log P(w_{t+j} | w_t; \theta) \quad \nabla_{\theta} J(\theta) = ?$$

**Quiz:** Suppose the vocab is  $V$  in the corpus, how many parameters in  $\theta$  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

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

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

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

$$\frac{df}{dx} = \frac{1}{x}$$

**chain rule:**

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

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

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

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

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



# COMPUTING THE GRADIENTS

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

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

$$\frac{\partial y}{\partial \mathbf{u}_t} = \frac{\partial (-\mathbf{u}_t \cdot \mathbf{v}_c + \log(\sum_{k \in V} \exp(\mathbf{u}_t \cdot \mathbf{v}_k)))}{\partial \mathbf{u}_t}$$

$$= -\mathbf{v}_c + \frac{\sum_{k \in V} \frac{\partial \exp(\mathbf{u}_t \cdot \mathbf{v}_k)}{\partial \mathbf{u}_t}}{\sum_{k \in V} \exp(\mathbf{u}_t \cdot \mathbf{v}_k)}$$

$$= -\mathbf{v}_c + \frac{\sum_{k \in V} \exp(\mathbf{u}_t \cdot \mathbf{v}_k) \mathbf{v}_k}{\sum_{k \in V} \exp(\mathbf{u}_t \cdot \mathbf{v}_k)}$$

$$= -\mathbf{v}_c + \sum_{k \in V} P(k | t) \mathbf{v}_k$$

$$\frac{\partial y}{\partial \mathbf{v}_k} = -1(k = c) \mathbf{u}_t + P(k | t) \mathbf{u}_t$$

(Try to derive this!)

# PUTTING IT TOGETHER

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

- Update:

$$\mathbf{u}_t \leftarrow \mathbf{u}_t - \eta \frac{\partial y}{\partial \mathbf{u}_t}$$
$$\mathbf{v}_k \leftarrow \mathbf{v}_k - \eta \frac{\partial y}{\partial \mathbf{v}_k}, \forall k \in V$$

- 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!

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

$$\frac{\partial y}{\partial \mathbf{v}_k} = -1(k = c) \mathbf{u}_t + P(k | t) \mathbf{u}_t$$

- 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)

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

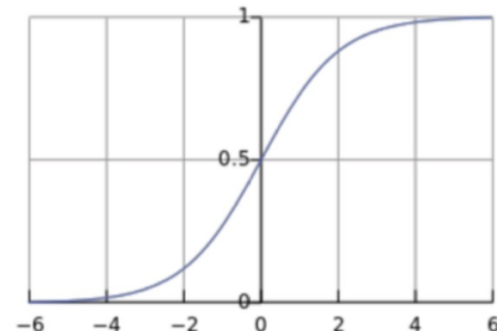
**positive examples +**

t	c
apricot	tablespoon
apricot	of
apricot	jam
apricot	a

**negative examples -**

t	c	t	c
apricot	aardvark	apricot	seven
apricot	my	apricot	forever
apricot	where	apricot	dear
apricot	coaxial	apricot	if

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



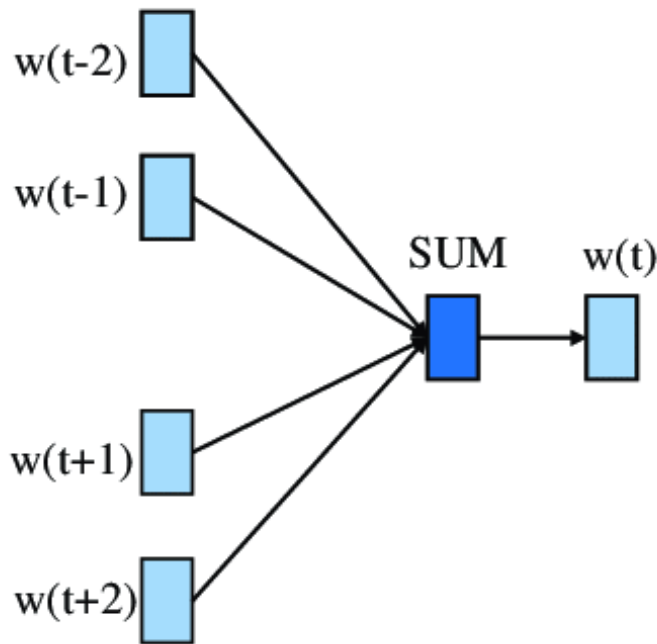
- Same as training a **logistic regression** for binary classification!

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

- Compute the gradient in assignment!

# CONTINUOUS BAG OF WORDS (CBOW)

INPUT PROJECTION OUTPUT



$$L(\theta) = \prod_{t=1}^T P(w_t | \{w_{t+j}\}, -m \leq j \leq m, j \neq 0)$$

$$\bar{\mathbf{v}}_t = \frac{1}{2m} \sum_{-m \leq j \leq m, j \neq 0} \mathbf{v}_{t+j}$$

(average context vector for  $w_t$ )

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

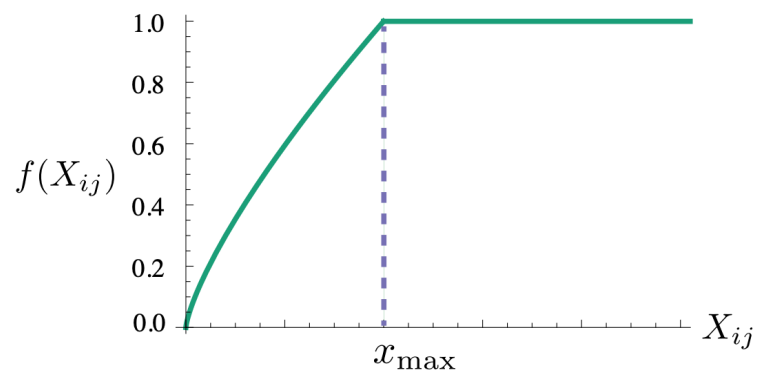
# GLOVE: GLOBAL VECTORS

- Skip-gram and CBoW uses local context
  - Slow to train when the corpus is very large
- Let's take the global co-occurrence statistics  $X_{ij}$ :

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

- $X_{ij}$  tabulate the number of times word  $j$  occurs in the context of word  $i$ .

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



# 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  $\mathbf{u}_i \cdot \mathbf{v}_j$  by

$$\sum_{g \in \text{ngrams}(w_i)} \mathbf{u}_g \cdot \mathbf{v}_j$$

- More to come: contextualized word embeddings



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



# PRE-TRAINED WORD EMBEDDINGS AVAILABLE

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

## NLPL word embeddings repository

brought to you by [Language Technology Group](#) at the University of Oslo

We feature models trained with clearly stated hyperparameters, 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)  fastText Skipgram  Word2Vec Continuous Skipgram  Gensim Continuous Skipgram  
 Gensim Continuous Bag-of-Words  fastText Continuous Bag-of-Words

#### Lemmatization:

True  False

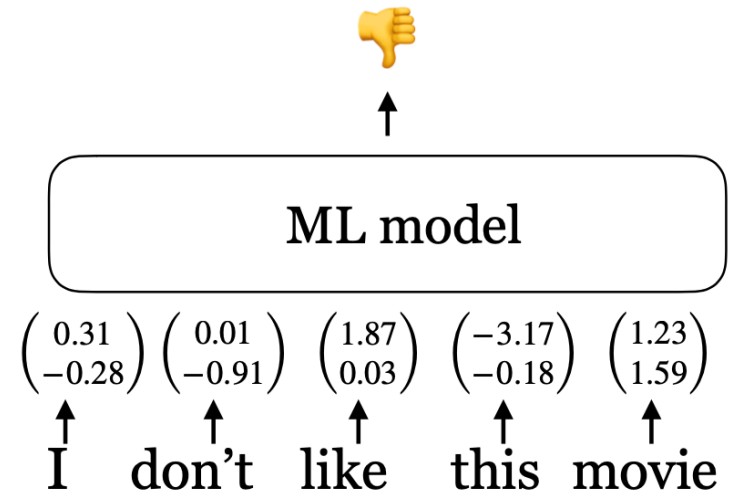
- Differ in algorithms, text corpora, dimensions, cased/uncased...

# 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
- <http://www.cs.technion.ac.il/~gabr/resources/data/wordsim353/>

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:

$$\cos(\mathbf{u}_i, \mathbf{u}_j) = \frac{\mathbf{u}_i \cdot \mathbf{u}_j}{\|\mathbf{u}_i\|_2 \times \|\mathbf{u}_j\|_2}$$

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 <sup>†</sup>	6B	57.2	65.6	68.2	57.0	32.5
SG <sup>†</sup>	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	<b><u>75.9</u></b>	<b><u>83.6</u></b>	<b><u>82.9</u></b>	<b><u>59.6</u></b>	<b><u>47.8</u></b>
CBOW*	100B	68.4	79.6	75.4	59.4	45.5

# INTRINSIC EVALUATION

- Word analogy

man:woman  $\approx$  king: ?

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

## Semantic

Austin:Texas  $\approx$  ? :California

## Syntactic

bad:worst  $\approx$  hot: ?

More examples at:

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