

# lecture 8/13 (LAST CLASS, no class FRI 8/15)

- word embedding
- Glove + HW6 (optional) demo

logistics: HWS done (closed)

[HW6] due Fri incl done  
PROJ: submitted, look for discussion?

# Introduction to Deep Learning

## 22. Embeddings, Word2vec, fastText, GloVe

STAT 157, Spring 2019, UC Berkeley

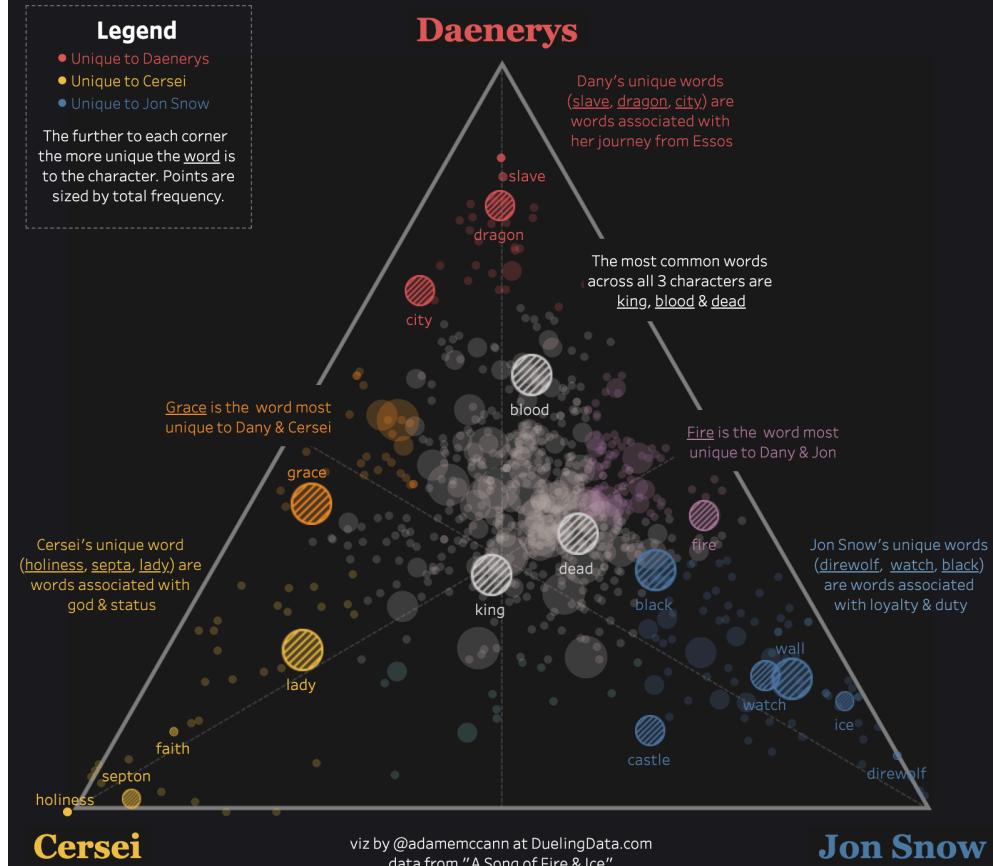
Alex Smola and Mu Li

[courses.d2l.ai/berkeley-stat-157](https://courses.d2l.ai/berkeley-stat-157)

# GAME OF THRONES™ IN WORDS

This viz shows the most unique words by character for each chapter in the 5 Game of Thrones books

# Word2Vec

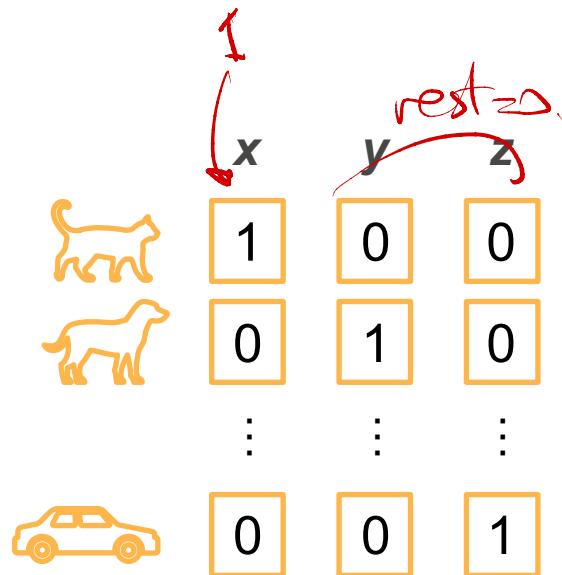


# Motivation

$\text{Enc} \Rightarrow \text{word ID in vector}^{\text{dim}} = 300 \text{ words} \Rightarrow \text{dim} = \# \text{vocab}$

- One-hot vectors map objects/words into fixed-length vectors
- These vectors only contain the identity information, not semantic meaning, e.g. no easy sim()

$$\langle \mathbf{x}, \mathbf{y} \rangle = \langle \mathbf{z}, \mathbf{y} \rangle = 0$$



# Word2vec

300 dim real numbers

- Learn an embedding vector for each word
- Use  $\langle x, y \rangle$  to measure the similarity

dot prod for sim

$$\text{sim}(x, y) = \langle x, y \rangle > \langle z, y \rangle$$

- Build a probability model
- Maximize the likelihood function to learn the model

loss  $\equiv$  meaningful  
min  $\equiv$  word vectors

- Tricks for fast computation

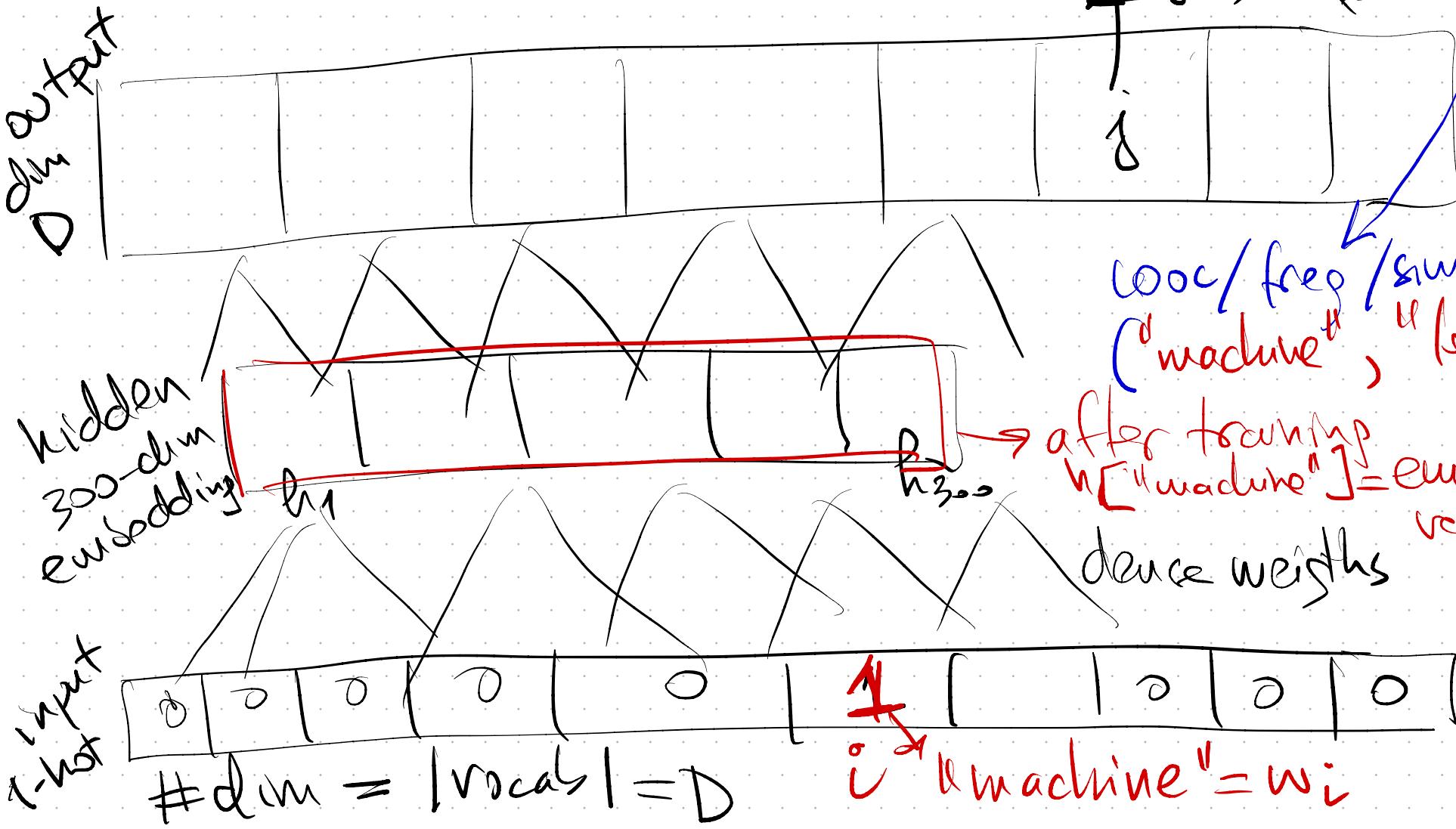
x	y	z
	1	0
	0	1
:	:	:
	0	0
		1

## Simple/basic idea

- one-hot representation

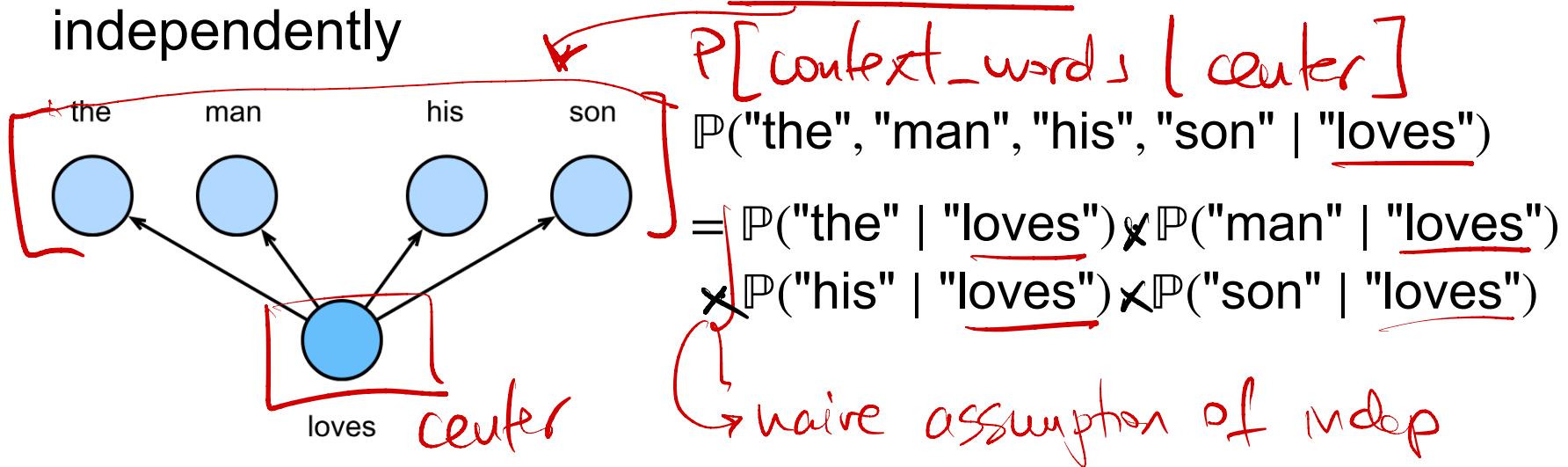
- NN-aufencoder/similr

$W_j = \text{learning}^j$

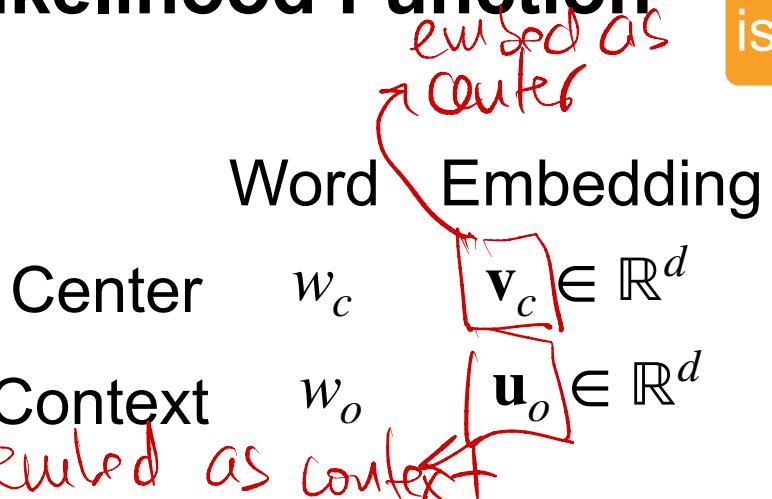


# The Skip-Gram Model

- A word can be used to generate the words surround it
- Given the center word, the context words are generated independently



# Likelihood Function



Summing over all words  
is too expensive

$$P(\text{context } w_o | \text{center}) = \frac{\exp(u_o^\top v_c)}{\sum_{i \in \mathcal{V}} \exp(u_i^\top v_c)}$$

$\mathcal{V}$  : all context words

Normalization  
over all words

- Given length  $T$  sequence, context window  $m$ , the likelihood function:

each word

naive assumption

$$\prod_{t=1}^T \prod_{m \leq j \leq m, j \neq 0} \mathbb{P}(w^{(t+j)} | w^{(t)})$$

## TRICKS for computation

- Simplify/approx SOFTMAX; do not compute normalized denominator  $\text{Sigmoid}(U_c^T \cdot V_o)$  instead.
  - much easier to compute prob(in window)
- $P(D=1 | w^t, w^{t+j})$  = prob that  $w^t, w^{t+j}$  are in same context window
- Negative Sampling: include words  $w_n$  NOT in context window.
  - select some (not all)
  - use probab  $P(\quad) = 1 - \text{Sigmoid}(U_c^T \cdot V_n)$
- not include all counts  $\text{prob}(\text{not in window})$
- Subsampling frequent words  $(P_{\text{pos}}^{\text{neg}}) P_{w_i} = \max\left(1 - \frac{\sum_{\text{fw.}}}{n}\right)$

# Negative Sampling

• Simplify/approx SOFTMAX

- Treat a center word and a context word appear in the same context window as an event

$$\mathbb{P}(D = 1 | w_c, w_o) = \sigma(\mathbf{u}_c^T \mathbf{v}_o)$$
$$\sigma(x) = \frac{1}{1 + \exp(-x)}$$

- Change the likelihood function from  $\prod_{t=1}^T \prod_{-m \leq j \leq m, j \neq 0} \mathbb{P}(w^{(t+j)} | w^{(t)})$  to

$$\prod_{t=1}^T \prod_{-m \leq j \leq m, j \neq 0} \mathbb{P}(D = 1 | w^{(t)}, w^{(t+j)})$$

Naive solution: infinity

# Negative Sampling

- Sample noise word  $w_n$  that doesn't appear in the window

$$\mathbb{P}(D = 0 | w_c, w_n) = 1 - \sigma(\mathbf{u}_n^T \mathbf{v}_c)$$

- Add into the likelihood function as well
- Maximizing the likelihood equals to solve a binary classification problem with a binary logistic regression loss

The skip-gram model parameters are the center word vector and context word vector for each word in the vocabulary. In training, we learn the model parameters by maximizing the likelihood function (i.e., maximum likelihood estimation). This is equivalent to minimizing the following loss function:

$$-\sum_{t=1}^T \sum_{-m \leq j \leq m, j \neq 0} \log P(w^{(t+j)} | w^{(t)}). \quad (15.1.6)$$

When using stochastic gradient descent to minimize the loss, in each iteration we can randomly sample a shorter subsequence to calculate the (stochastic) gradient for this subsequence to update the model parameters. To calculate this (stochastic) gradient, we need to obtain the gradients of the log conditional probability with respect to the center word vector and the context word vector. In general, according to (15.1.4) the log conditional probability involving any pair of the center word  $w_c$  and the context word  $w_o$  is

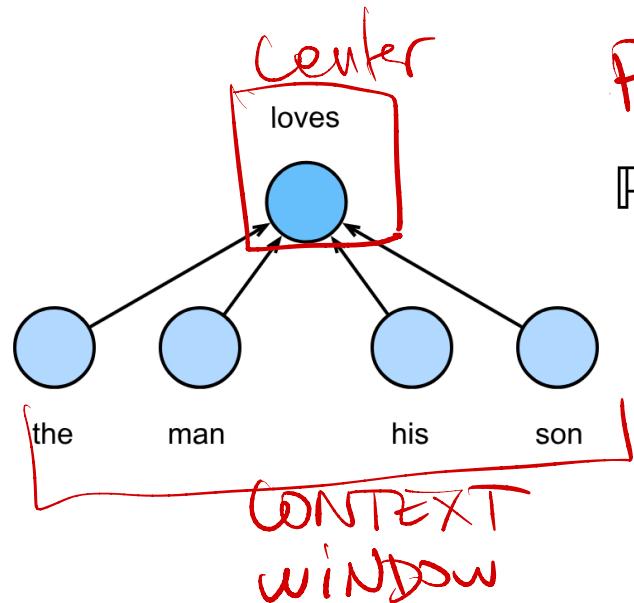
$$\log P(w_o | w_c) = \mathbf{u}_o^\top \mathbf{v}_c - \log \left( \sum_{i \in \mathcal{V}} \exp(\mathbf{u}_i^\top \mathbf{v}_c) \right). \quad (15.1.7)$$

Through differentiation, we can obtain its gradient with respect to the center word vector  $\mathbf{v}_c$  as

$$\begin{aligned} \frac{\partial \log P(w_o | w_c)}{\partial \mathbf{v}_c} &= \mathbf{u}_o - \frac{\sum_{j \in \mathcal{V}} \exp(\mathbf{u}_j^\top \mathbf{v}_c) \mathbf{u}_j}{\sum_{i \in \mathcal{V}} \exp(\mathbf{u}_i^\top \mathbf{v}_c)} \\ &= \mathbf{u}_o - \sum_{j \in \mathcal{V}} \left( \frac{\exp(\mathbf{u}_j^\top \mathbf{v}_c)}{\sum_{i \in \mathcal{V}} \exp(\mathbf{u}_i^\top \mathbf{v}_c)} \right) \mathbf{u}_j \\ &= \mathbf{u}_o - \sum_{j \in \mathcal{V}} P(w_j | w_c) \mathbf{u}_j. \end{aligned} \quad (15.1.8)$$

# Continuous Bag Of Words (CBOW) REVERSE COND PROB

- The center word is generated based on the context words



$P(\text{center} \mid \text{context\_words})$

$P(\text{"loves"} \mid \text{"the"}, \text{"man"}, \text{"his"}, \text{"son"})$

one naive product-proba / indep

instead evaluate score  
based on all context

# Likelihood Function

- Compute the probability

$$\mathbb{P}(w_c \mid w_{o_1}, \dots, w_{o_{2m}}) =$$

Dot product on whole window.

$$\frac{\exp\left(\frac{1}{2m}\mathbf{u}_c^\top(\mathbf{v}_{o_1} + \dots + \mathbf{v}_{o_{2m}})\right)}{\sum_{i \in \mathcal{V}} \exp\left(\frac{1}{2m}\mathbf{u}_i^\top(\mathbf{v}_{o_1} + \dots + \mathbf{v}_{o_{2m}})\right)}$$

normal.

SOFTMAX

- Likelihood

each word

$$\prod_{t=1}^T \mathbb{P}(w^{(t)} \mid w^{(t-m)}, \dots, w^{(t-1)}, w^{(t+1)}, \dots, w^{(t+m)})$$

min ↑

## Training

Training continuous bag of words models is almost the same as training skip-gram models. The maximum likelihood estimation of the continuous bag of words model is equivalent to minimizing the following loss function:

$$-\sum_{t=1}^T \log P(w^{(t)} | w^{(t-m)}, \dots, w^{(t-1)}, w^{(t+1)}, \dots, w^{(t+m)}). \quad (15.1.13)$$

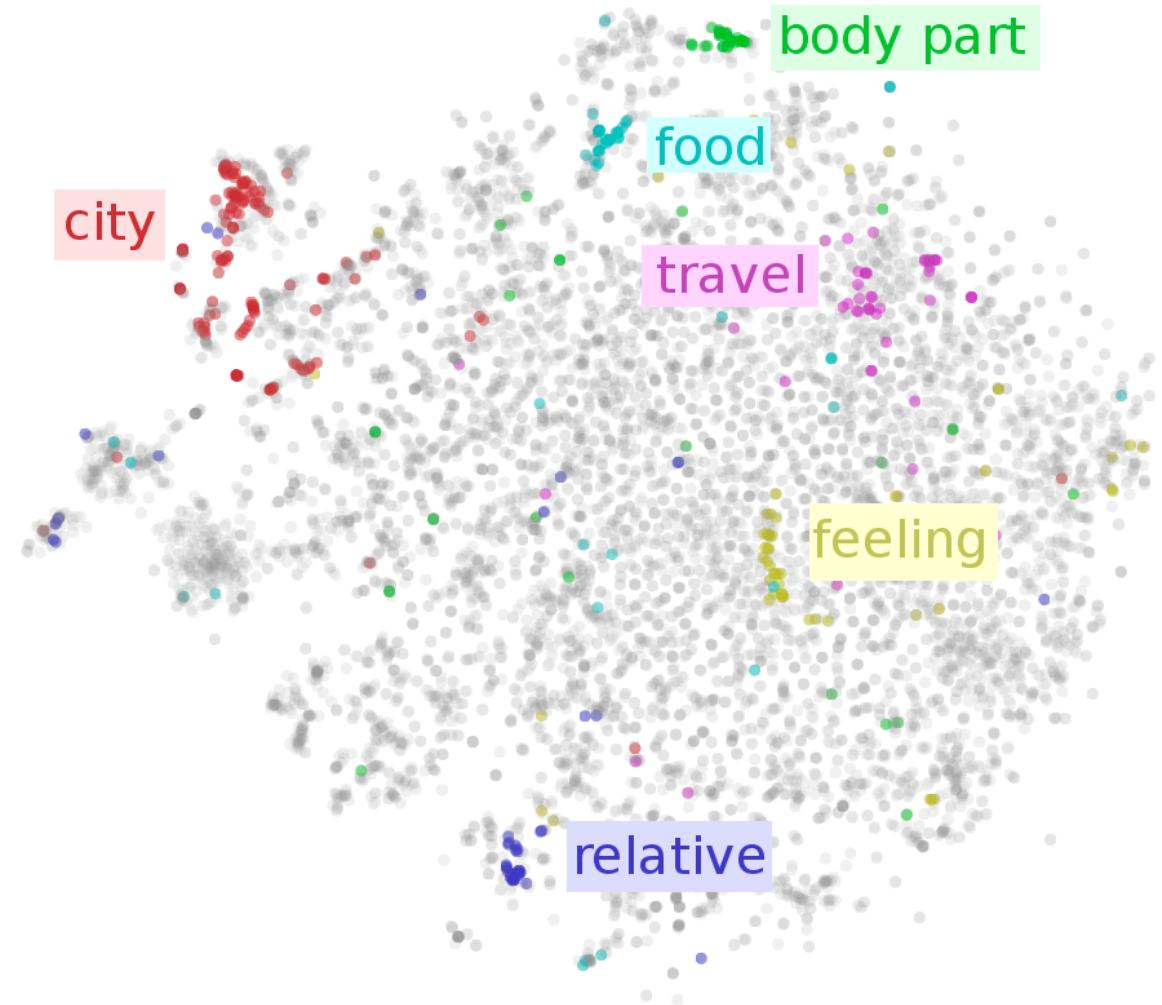
Notice that

$$\log P(w_c | \mathcal{W}_o) = \mathbf{u}_c^\top \bar{\mathbf{v}}_o - \log \left( \sum_{i \in \mathcal{V}} \exp(\mathbf{u}_i^\top \bar{\mathbf{v}}_o) \right). \quad (15.1.14)$$

Through differentiation, we can obtain its gradient with respect to any context word vector  $\mathbf{v}_{o_i}$  ( $i = 1, \dots, 2m$ ) as

$$\frac{\partial \log P(w_c | \mathcal{W}_o)}{\partial \mathbf{v}_{o_i}} = \frac{1}{2m} \left( \mathbf{u}_c - \sum_{j \in \mathcal{V}} \frac{\exp(\mathbf{u}_j^\top \bar{\mathbf{v}}_o) \mathbf{u}_j}{\sum_{i \in \mathcal{V}} \exp(\mathbf{u}_i^\top \bar{\mathbf{v}}_o)} \right) = \frac{1}{2m} \left( \mathbf{u}_c - \sum_{j \in \mathcal{V}} P(w_j | \mathcal{W}_o) \mathbf{u}_j \right).$$

# More Embedding Models



FastText := stemming / lexicographics  
"friend" = "friends" = "friendship" = "friendness". . .

- English words usually have internal structures and formation methods
  - dog, dogs, dogcatcher
- Each center word is represented as a set of subwords
  - “where” -> “<where>” ->  $n$ -gram
  - $n=3$ : “<wh”, “whe”, “her”, “ere”, “re>”
- Useful for long but infrequent words
  - e.g. pneumonoultramicroscopicsilicovolcanoconiosis



$$\text{embedding} = \sum \text{embed}(grows)$$

# FastText

- For word  $w$ ,  $\mathcal{G}_w$  is the union of subwords with length from 3 to 6
- The center vector is then

$$\mathbf{u}_w = \sum_{g \in \mathcal{G}_w} \mathbf{u}_g$$

↑ sum (grams)

- The rest model is same as skip-gram

# Word Embedding with Global Vectors (GloVe)

general  
wordvec

pretrained

(large corpus)

- Denote by

$$q_{ij} = \frac{\exp(\mathbf{u}_j^\top \mathbf{v}_i)}{\sum_{k \in \mathcal{V}} \exp(\mathbf{u}_k^\top \mathbf{v}_i)}$$

- Rewrite the negative log-likelihood function of skip-gram

$$-\log \left[ \prod_{t=1}^T \prod_{-m \leq j \leq m, j \neq 0} \mathbb{P}(w^{(t+j)} | w^{(t)}) \right]$$

tune for your  
specific corpus.

- as  $-\sum_{i \in \mathcal{V}} \sum_{j \in \mathcal{V}} x_{ij} \log q_{ij}$  with proper counts  $x_{ij}$   
*MULTICOUNT  
Count(i,j) in context window*

# Glove

download  
various sizes

not approx, same  
as original.

- Further rewrite

$$-\sum_{i \in \mathcal{V}} \sum_{j \in \mathcal{V}} x_{ij} \log q_{ij} = -\sum_{i \in \mathcal{V}} x_i \sum_{j \in \mathcal{V}} p_{ij} \log q_{ij}$$

with  $x_i = \sum_j x_{ij}$ ,  $p_{ij} = x_i / x_{ij}$

Cross entropy

$x_i$  / ratio - normalized

# Glove

- Replace the cross entropy with a log square loss

$$\sum_{j \in \mathcal{V}} p_{ij} \log q_{ij} \rightarrow \sum_{j \in \mathcal{V}} (\log p_{ij} - \log q'_{ij})^2$$

approx  
easier to  
compute

with an easy to compute  $q'_{ij} = \exp(\mathbf{u}_j^\top \mathbf{v}_i)$

- Add bias term for center and context words
- Replace the weights  $x_i$  with a monotone increasing function in  $[0, 1]$

$$\sum_{i \in \mathcal{V}} \sum_{j \in \mathcal{V}} h(x_{ij}) \left( \mathbf{u}_j^\top \mathbf{v}_i + b_i + c_j - \log x_{ij} \right)^2$$

monotonic func(x) control scale  
bias and sensitivity.