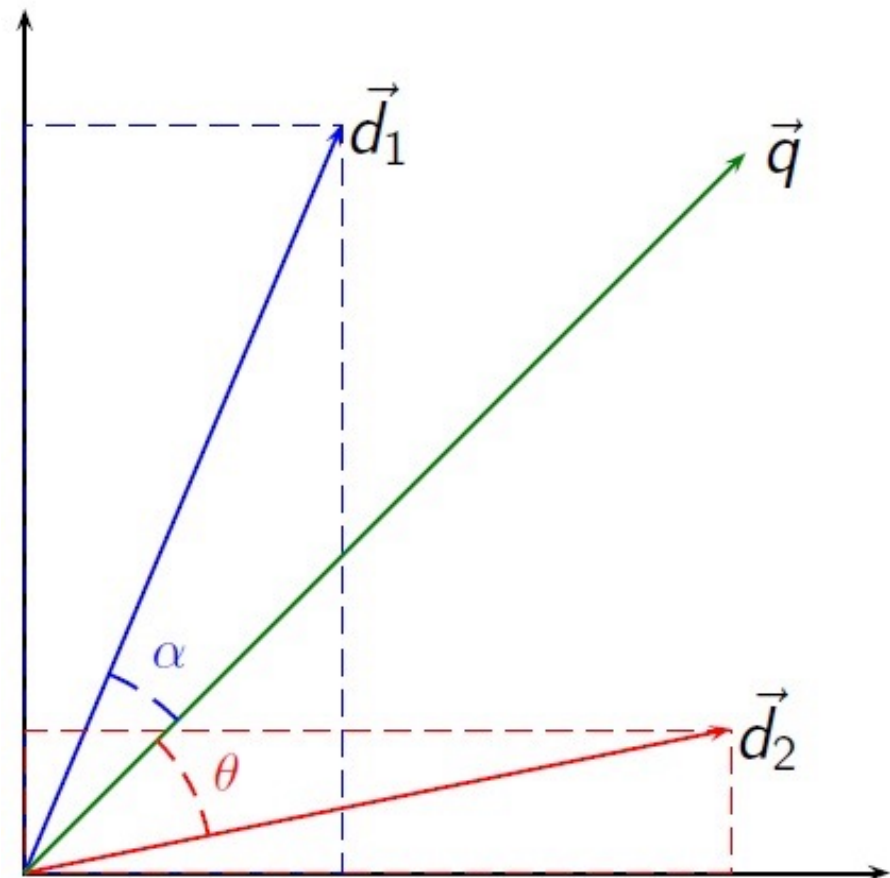


# Lecture 9: Vector Semantics

Mark Hasegawa-Johnson

CC0: Public domain. Re-  
use, Remix, Redistribute at  
will.



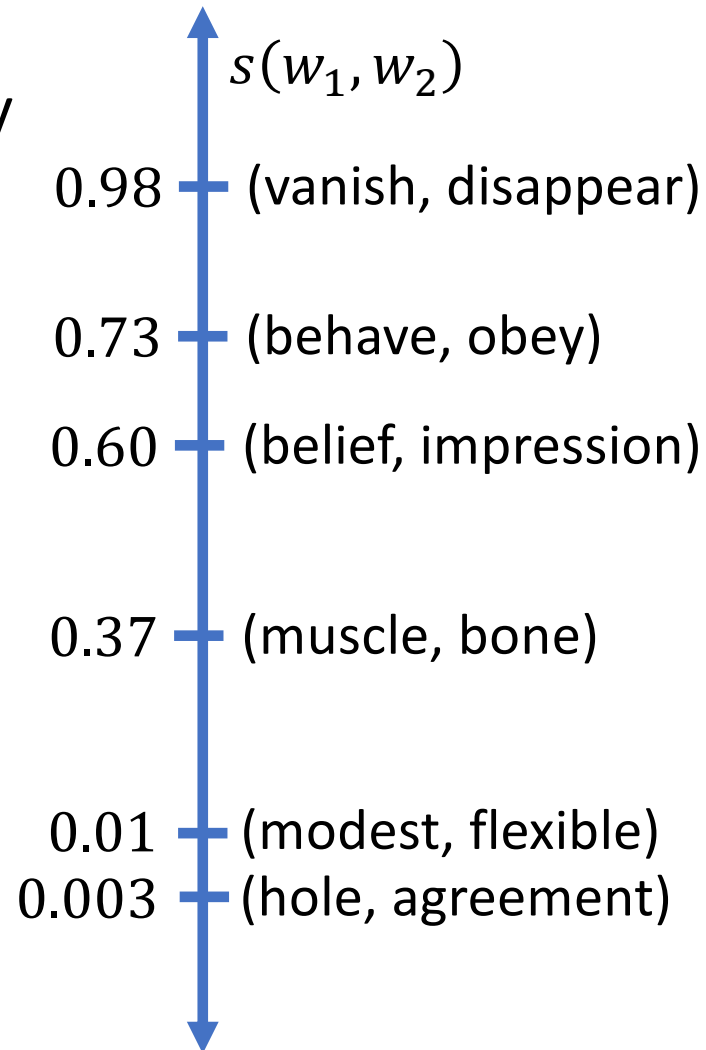
By Riclas - Own work, CC BY 3.0,  
<https://commons.wikimedia.org/w/index.php?curid=9076846>

# Outline

- Vector semantics: CBOW and skip-gram
- Calculating probabilities using the logistic sigmoid
- Dot-product similarity
- Noise contrastive estimation (NCE)
- Visualizations

# Synonymy and similarity

- Words are “synonyms” if they have exactly the same meaning.
- No words ever have **exactly** the same meaning, so no two words are ever exactly synonyms.
- We prefer to talk about word similarity,  $0 \leq s(w_1, w_2) \leq 1$ 
  - $s(w_1, w_2) = 1$ :  $w_1$  and  $w_2$  are perfect synonyms. Never happens in practice, but sometimes close.
  - $s(w_1, w_2) = 0$ :  $w_1$  and  $w_2$  are completely different.



# Skip-gram: “You shall know a word by the company it keeps”

- The key idea of skip-gram is that “You shall know a word by the company it keeps” (J.R. Firth, 1957)
- The words “vanish” and “disappear” should be considered similar if they can occur in the same contexts:

“The ship will vanish into the mists.”

“The ship will disappear into the mists.”

# Review: Naïve Bayes: the “Bag-of-words” model

We can estimate the likelihood of an e-mail by pretending that the e-mail is just a bag of words (order doesn't matter).

With only a few thousand spam e-mails, we can get a pretty good estimate of these things:

- $P(W = \text{"hi"}|Y = \text{spam}), P(W = \text{"hi"}|Y = \text{ham})$
- $P(W = \text{"vitality"}|Y = \text{spam}), P(W = \text{"vitality"}|Y = \text{ham})$
- $P(W = \text{"production"}|Y = \text{spam}), P(W = \text{"production"}|Y = \text{ham})$

Then we can approximate  $P(X|Y)$  by assuming that the words,  $W$ , are **conditionally independent of one another given the category label**:

$$P(X = x|Y = y) \approx \prod_{i=1}^n P(W = w_i|Y = y)$$



# Similarity: The Internet is the database

Similarity = words can be used interchangeably in most contexts

How do we measure that in practice?

Answer: extract examples of word  $w_1$ , +/- C words (C=2 or 3):

...hot, although iced coffee is a popular...

...indicate that moderate coffee consumption is benign...

...and of  $w_2$ :

...consumed as iced tea. Sweet tea is...

...national average of tea consumption in Ireland...

The words “iced” and “consumption” appear in both contexts, so we can conclude that  $s(\text{coffee}, \text{tea}) > 0$ . No other words are shared, so we can conclude  $s(\text{coffee}, \text{tea}) < 1$ .

# Continuous Bag of Words (CBOW)

“Context bag of words” (CBOW) approximates each word’s probability, in context, by the product of single-word context probabilities:

$$p(w_t | w_{t-3}, \dots, w_{t+3}) \approx \prod_{\substack{c \neq 0 \\ c=-3 \\ c=3}}^3 p(w_t | w_{t+c})$$

Using this model, two words are considered similar if they have similar CBOW probabilities:

$$\begin{aligned} p(\text{vanish} | \text{the ship will} - \text{into the mists}) &\approx \\ p(\text{disappear} | \text{the ship will} - \text{into the mists}) \end{aligned}$$

# Skip-gram

“Skip-gram” approximates the probability of the context given the center word:

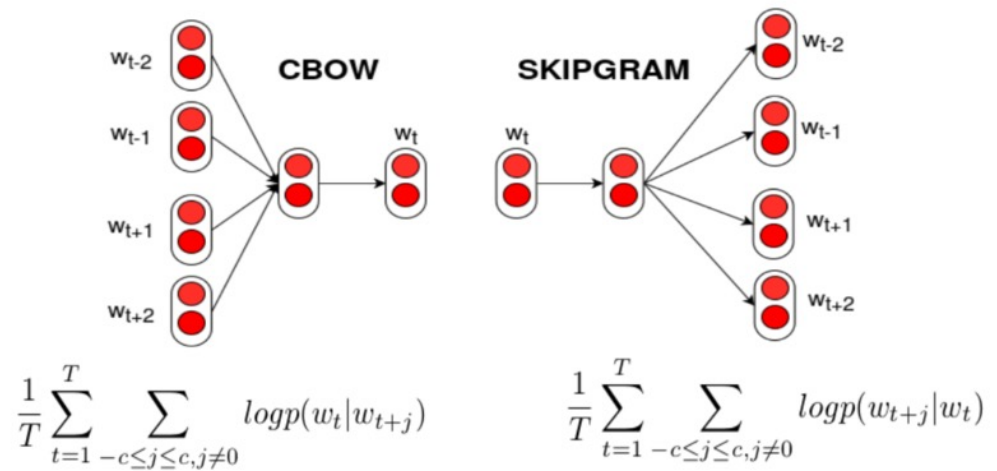
$$p(w_{t-3}, \dots, w_{t+3} | w_t) \approx \prod_{\substack{j \neq 0 \\ j=-3}}^3 p(w_{t+j} | w_t)$$

Using this model, two words are considered similar if they have similar skip-gram probabilities:

$$\begin{aligned} p(\text{the ship will} - \text{into the mists} | \text{vanish}) &\approx \\ p(\text{the ship will} - \text{into the mists} | \text{disappear}) \end{aligned}$$



Neural net loss =  
negative log probability



Log probability is easier to differentiate than probability, so:

CBOW:

$$\mathcal{L} = -\frac{1}{T} \sum_{t=0}^{T-1} \sum_{j=-c, j \neq 0}^c \ln P(w_t | w_{t+j})$$

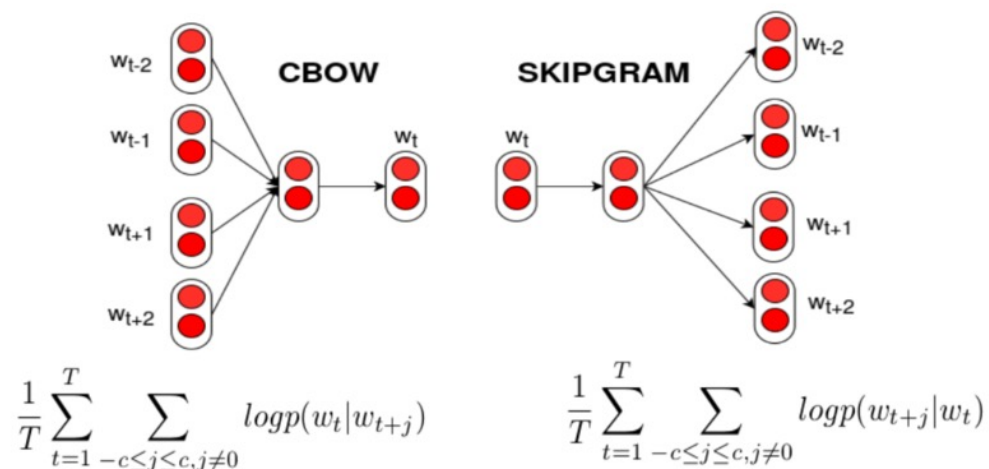
Skip-gram:

$$\mathcal{L} = -\frac{1}{T} \sum_{t=0}^{T-1} \sum_{j=-c, j \neq 0}^c \ln P(w_{t+j} | w_t)$$

# Outline

- Vector semantics: CBOW and skip-gram
- Calculating probabilities using the logistic sigmoid
- Dot-product similarity
- Noise contrastive estimation (NCE)
- Visualizations

# Neural net probability estimate?



Let's focus on skip-gram:

$$\mathcal{L} = -\frac{1}{T} \sum_{t=0}^{T-1} \sum_{j=-c, j \neq 0}^c \ln P(w_{t+j} | w_t)$$

Suppose we want to use a neural net to estimate  $P(w_{t+j} | w_t)$ . How can we do that?

# The logistic sigmoid function

- Remember the logistic sigmoid function:

$$\sigma(z) = \frac{1}{1 + e^{-z}}$$

- Since  $0 < \sigma(z) < 1$ , we can interpret  $\sigma(z)$  as a probability.

$$P(w_{t+j}|w_t) = \sigma(z)$$

- Apparently, we want the logit,  $z$ , to be a function that measures the similarity or predictability of  $w_{t+j}$  and  $w_t$ :
  - If  $w_{t+j}$  never occurs in the context of  $w_t$ , we want  $P(w_{t+j}|w_t) \approx 0$ , so  $z(w_{t+j}, w_t) \approx -\infty$
  - If  $w_{t+j}$  always occurs in the context of  $w_t$ , we want  $P(w_{t+j}|w_t) \approx 1$ , so  $z(w_{t+j}, w_t) \approx +\infty$

# Training the neural net

- Now we can write:

$$\mathcal{L} = -\frac{1}{T} \sum_{t=0}^{T-1} \sum_{j=-c, j \neq 0}^c \ln \sigma(z(w_{t+j}, w_t))$$

- For any pair of words  $w_1$  and  $w_2$ , the gradient is:

$$\frac{\partial \mathcal{L}}{\partial z(w_1, w_2)} = -\frac{1}{T} \sum_{w_t=w_2} \sum_{w_{t+j}=w_1} (1 - \sigma(z(w_1, w_2)))$$

- The more often  $w_1$  and  $w_2$  occur in context, the more that gradient descent will push  $z(w_1, w_2) \rightarrow \infty$

# Outline

- Vector semantics: CBOW and skip-gram
- Calculating probabilities using the logistic sigmoid
- Dot-product similarity
- Noise contrastive estimation (NCE)
- Visualizations

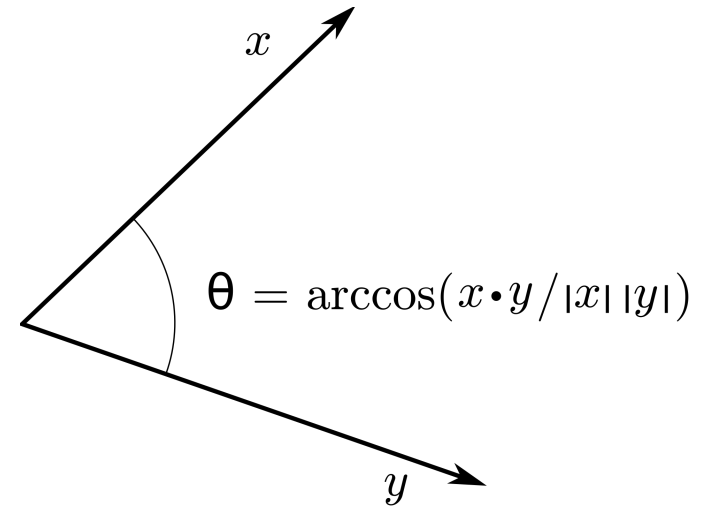
# Dot product between vectors

The dot product between two vectors is:

$$\mathbf{v}_1^T \mathbf{v}_2 = |\mathbf{v}_1| |\mathbf{v}_2| \cos \theta$$

The dot product is a pretty good measure of the similarity between the vectors.

- If  $\theta < 90$  degrees, then  $\mathbf{v}_1^T \mathbf{v}_2 > 0$
- If  $\theta = 90$  degrees, then  $\mathbf{v}_1^T \mathbf{v}_2 = 0$
- If  $\theta > 90$  degrees, then  $\mathbf{v}_1^T \mathbf{v}_2 < 0$



By BenFrantzDale at the English Wikipedia, CC BY-SA 3.0,  
<https://commons.wikimedia.org/w/index.php?curid=49972362>

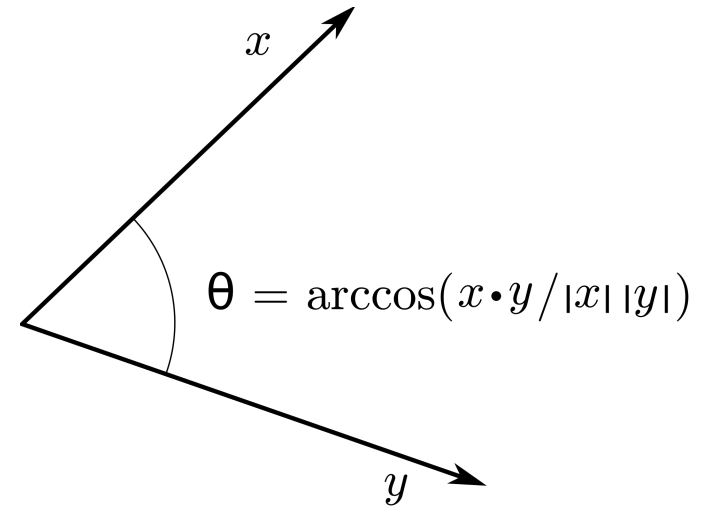
# Words as vectors

- Suppose that, for every word in the dictionary  $w_1$ , we train a vector  $\mathbf{v}_1$
- The similarity between two words should be measured by the dot product between their vectors:

$$z(w_1, w_2) = \mathbf{v}_1^T \mathbf{v}_2$$

- The probability of one word, given the other, is the sigmoid of their dot product:

$$P(w_1 | w_2) = \frac{1}{1 + e^{-\mathbf{v}_1^T \mathbf{v}_2}}$$

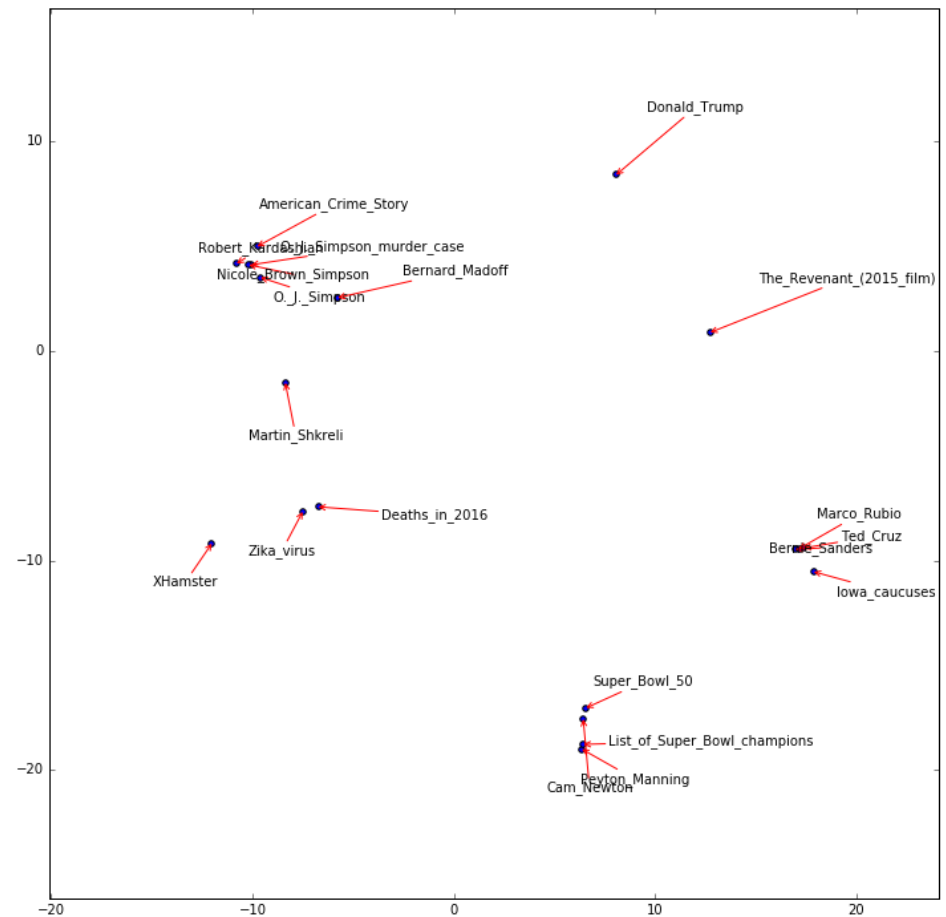


By BenFrantzDale at the English Wikipedia, CC BY-SA 3.0,  
<https://commons.wikimedia.org/w/index.php?curid=49972362>



# Words as vectors

t-SNE projection of Word2vec  
embedding, ewulczyn 2016,  
[https://commons.wikimedia.org/wiki/File:2016\\_02\\_mini\\_embedding.png](https://commons.wikimedia.org/wiki/File:2016_02_mini_embedding.png)



# Training the neural net

- Now we can write:

$$\mathcal{L} = -\frac{1}{T} \sum_{t=0}^{T-1} \sum_{j=-c, j \neq 0}^c \ln \sigma(\mathbf{v}_1^T \mathbf{v}_2)$$

- For every word  $w$  with vector  $\mathbf{v}_t$ , the loss gradient is:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{v}} = -\frac{1}{T} \sum_{w_t=w} \sum_{j=-c, j \neq 0}^c (1 - \sigma(\mathbf{v}_{t+j}^T \mathbf{v})) \mathbf{v}_{t+j}$$

# Training the neural net

- Gradient descent updates  $\mathbf{v}$  as:

$$\begin{aligned}\mathbf{v} &\leftarrow \mathbf{v} - \eta \frac{\partial \mathcal{L}}{\partial \mathbf{v}} \\ &= \mathbf{v} + \frac{\eta}{T} \sum_{w_t=w} \sum_{j=-c, j \neq 0}^c (1 - \sigma(\mathbf{v}_{t+j}^T \mathbf{v})) \mathbf{v}_{t+j}\end{aligned}$$

- In other words,  $\mathbf{v}$  becomes a weighted average of the vectors  $\mathbf{v}_{t+j}$  for words  $w_{t+j}$  that occur in the context of  $w$ .
- The vectors start out random, then vectors that frequently co-occur get pushed together, to become more similar

# Outline

- Vector semantics: CBOW and skip-gram
- Calculating probabilities using the logistic sigmoid
- Dot-product similarity
- Noise contrastive estimation (NCE)
- Visualizations

# Training the neural net

- This formula has a problem:

$$\mathbf{v} \leftarrow \mathbf{v} + \frac{\eta}{T} \sum_{w_t=w} \sum_{j=-c, j \neq 0}^c (1 - \sigma(\mathbf{v}_{t+j}^T \mathbf{v})) \mathbf{v}_{t+j}$$

- The problem is that  $(1 - \sigma(\mathbf{v}_{t+j}^T \mathbf{v}))$  is always positive
- Therefore,  $\mathbf{v}$  gets gradually bigger and bigger – it never stops growing!!

## Fixing the problem

To keep all the vectors from growing forever, we need to subtract something from the right-hand side:

$$\mathbf{v} \leftarrow \mathbf{v} + \frac{\eta}{T} \sum_{w_t=w} \sum_{j=-c, j \neq 0}^c (1 - \sigma(\mathbf{v}_{t+j}^T \mathbf{v})) \mathbf{v}_{t+j} - \text{SOMETHING}$$

...such that, on average, the update will be zero (keeping all vectors small).

- If we still want to say that

$$\mathbf{v} \leftarrow \mathbf{v} - \eta \frac{\partial \mathcal{L}}{\partial \mathbf{v}}$$

- ...then we need to add something to  $\mathcal{L}$  whose derivative will be “SOMETHING”

# Noise contrastive estimation

So far, we've said that the probability of  $\mathbf{v}_{t+j}$  occurring in the context of  $\mathbf{v}_t$  is:

$$P(\mathbf{v}_{t+j}|\mathbf{v}_t) = \sigma(\mathbf{v}_{t+j}^T \mathbf{v}_t)$$

For some randomly chosen “noise” sample,  $\mathbf{v}_i$ , the probability it DOESN'T occur in the context of  $\mathbf{v}_t$  is:

$$1 - P(\mathbf{v}_i|\mathbf{v}_t) = 1 - \sigma(\mathbf{v}_i^T \mathbf{v}_t)$$

We can create a well-normalized loss function by adding these two log probabilities:

$$\mathcal{L} = -\ln \sigma(\mathbf{v}_{t+j}^T \mathbf{v}_t) - \ln(1 - \sigma(\mathbf{v}_i^T \mathbf{v}_t))$$

# Skipgram NCE

Skipgram NCE uses the following update step:

1. Choose a pair  $\mathbf{v}_t$  and  $\mathbf{v}_{t+j}$  from the training data
2. Choose  $k$  different “noise” examples,  $\mathbf{v}_i$ , uniformly at random from the vocabulary
3. Update  $\mathbf{v}_t$  according to

$$\mathcal{L} = -\ln \sigma(\mathbf{v}_{t+j}^T \mathbf{v}_t) - \frac{1}{k} \sum_{i=1}^k \ln(1 - \sigma(\mathbf{v}_i^T \mathbf{v}_t))$$

$$\mathbf{v}_t \leftarrow \mathbf{v}_t - \eta \frac{\partial \mathcal{L}}{\partial \mathbf{v}}$$



# Quiz

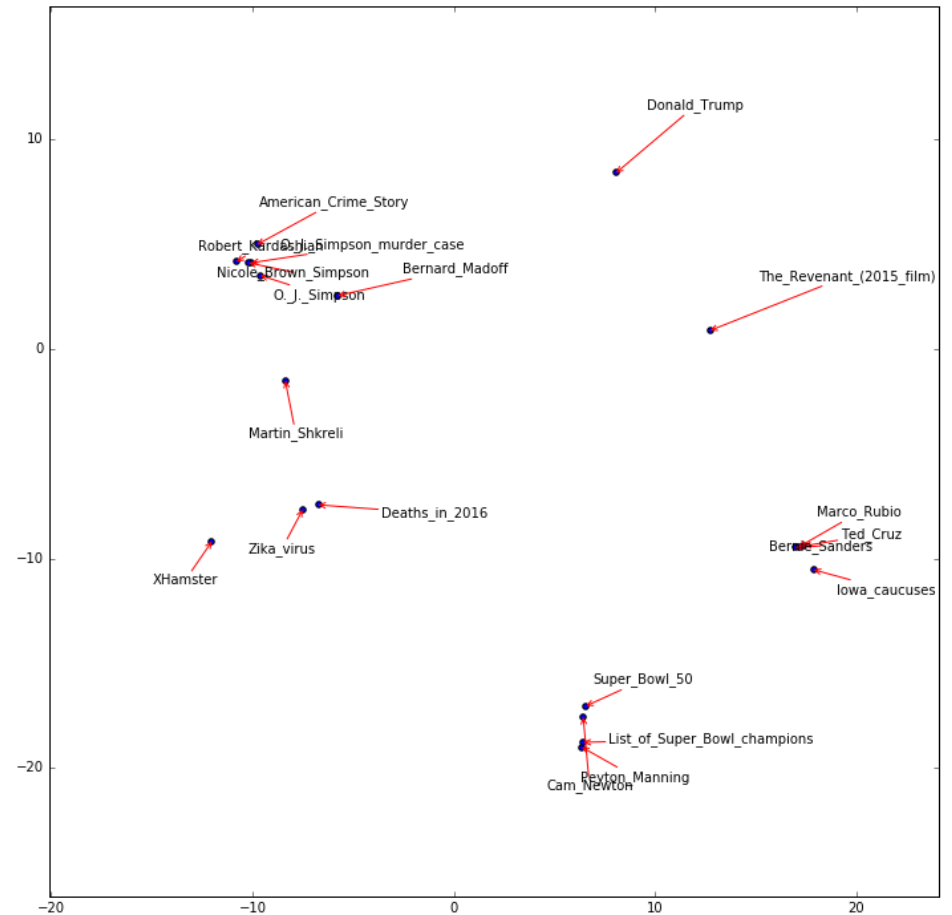
- Try the quiz!

# Outline

- Vector semantics: CBOW and skip-gram
- Calculating probabilities using the logistic sigmoid
- Dot-product similarity
- Noise contrastive estimation (NCE)
- **Visualizations**

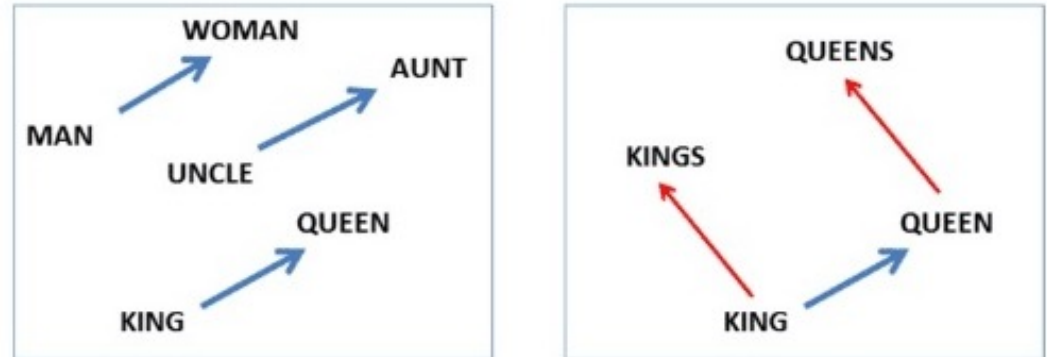
# Visualizations: Vectors

One way to visualize the results is with a scatter plot of the trained vectors:



# Visualizations: Relatedness

$$\text{vec}(\text{"woman"}) - \text{vec}(\text{"man"}) + \text{vec}(\text{"king"}) = \text{vec}(\text{"queen"})$$



Christian S. Perone, "Voynich Manuscript: word vectors and t-SNE visualization of some patterns," in *Terra Incognita*, 16/01/2016, <http://blog.christianperone.com/2016/01/voynich-manuscript-word-vectors-and-t-sne-visualization-of-some-patterns/>.

Mikolov (2013) showed that word2vec captures similarity relationships among words. For example, the difference between the vectors for “woman” and “man” is roughly the same as the difference between the vectors for “queen” and “king.” Perone (2016) showed that this effect works differently depending on the training corpus: in his blog post, he looks at word relatedness in the 15<sup>th</sup> century Voynich manuscript.

# Visualizations: Similarity

Another way to visualize the results is to see which words are nearest neighbors in the vector space:

Model (training time)	Redmond	Havel	ninjutsu	graffiti	capitulate
Collobert (50d) (2 months)	conyers lubbock keene	plauen dzerzhinsky osterreich	reiki kohona karate	cheesecake gossip dioramas	abdicate accede rearm
Turian (200d) (few weeks)	McCarthy Alston Cousins	Jewell Arzu Ovitz	- - -	gunfire emotion impunity	- - -
Mnih (100d) (7 days)	Podhurst Harlang Agarwal	Pontiff Pinochet Rodionov	- - -	anaesthetics monkeys Jews	Mavericks planning hesitated
Skip-Phrase (1000d, 1 day)	Redmond Wash. Redmond Washington Microsoft	Vaclav Havel president Vaclav Havel Velvet Revolution	ninja martial arts swordsmanship	spray paint grafitti taggers	capitulation capitulated capitulating

Table 6: Examples of the closest tokens given various well known models and the Skip-gram model trained on phrases using over 30 billion training words. An empty cell means that the word was not in the vocabulary.

# Learning biased analogies from data

- It's useful that algorithms like word2vec learn appropriate analogies, like "kings → king as queens → queen."
- Unfortunately, it also learns biases. For example, it learns that some occupations are "male," and some are "female!"
- This paper demonstrated the problem, and proposed a really interesting solution:  
    "Man is to Computer Programmer as Woman is to Homemaker?  
    Debiasing Word Embeddings," Bolukbasi et al., 2016

# The Male-Female vs. Neutral-Specific Space

Here's the resulting 2D space, from Bolukbasi et al., 2016:



# Summary

- CBOW:

$$\mathcal{L} = -\frac{1}{T} \sum_{t=0}^{T-1} \sum_{j=-c, j \neq 0}^c \ln P(w_t | w_{t+j})$$

- Skip-gram:

$$\mathcal{L} = -\frac{1}{T} \sum_{t=0}^{T-1} \sum_{j=-c, j \neq 0}^c \ln P(w_{t+j} | w_t)$$

- Dot-product similarity:

$$\mathcal{L} = -\frac{1}{T} \sum_{t=0}^{T-1} \sum_{j=-c, j \neq 0}^c \ln \sigma(\mathbf{v}_1^T \mathbf{v}_2)$$

- Skip-gram noise contrastive estimation:

$$\mathcal{L} = -\ln \sigma(\mathbf{v}_{t+j}^T \mathbf{v}_t) - \frac{1}{k} \sum_{i=1}^k \ln(1 - \sigma(\mathbf{v}_i^T \mathbf{v}_t))$$