

Word2Vec and Sentiment Analysis

An implementation from scratch, using only plain python

S.Shankar
2 May 2017

Overview

- This project demonstrates a sentiment analysis application.
- Here is the catch: word2vec is implemented from scratch, using only plain python primitives.

Stanford Sentiment Treebank (SST)

- SST offers a fine grained evaluation of sentiment analysis algorithms.
- Plausible phrases are parsed by the Stanford parser.
- Sample output of parse tree:

an artistl48897

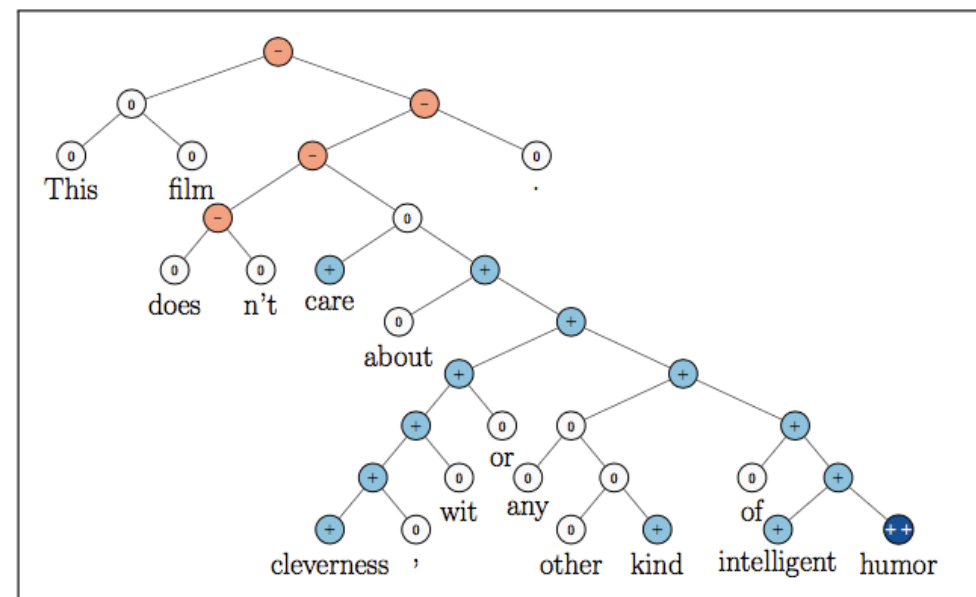
an artist who is simply tired --l192486

an artist who is simply tired -- of fighting the same fightsl192487

an artist who is simply tired -- of fighting the same fights , of putting the weight of the world on his shouldersl192489

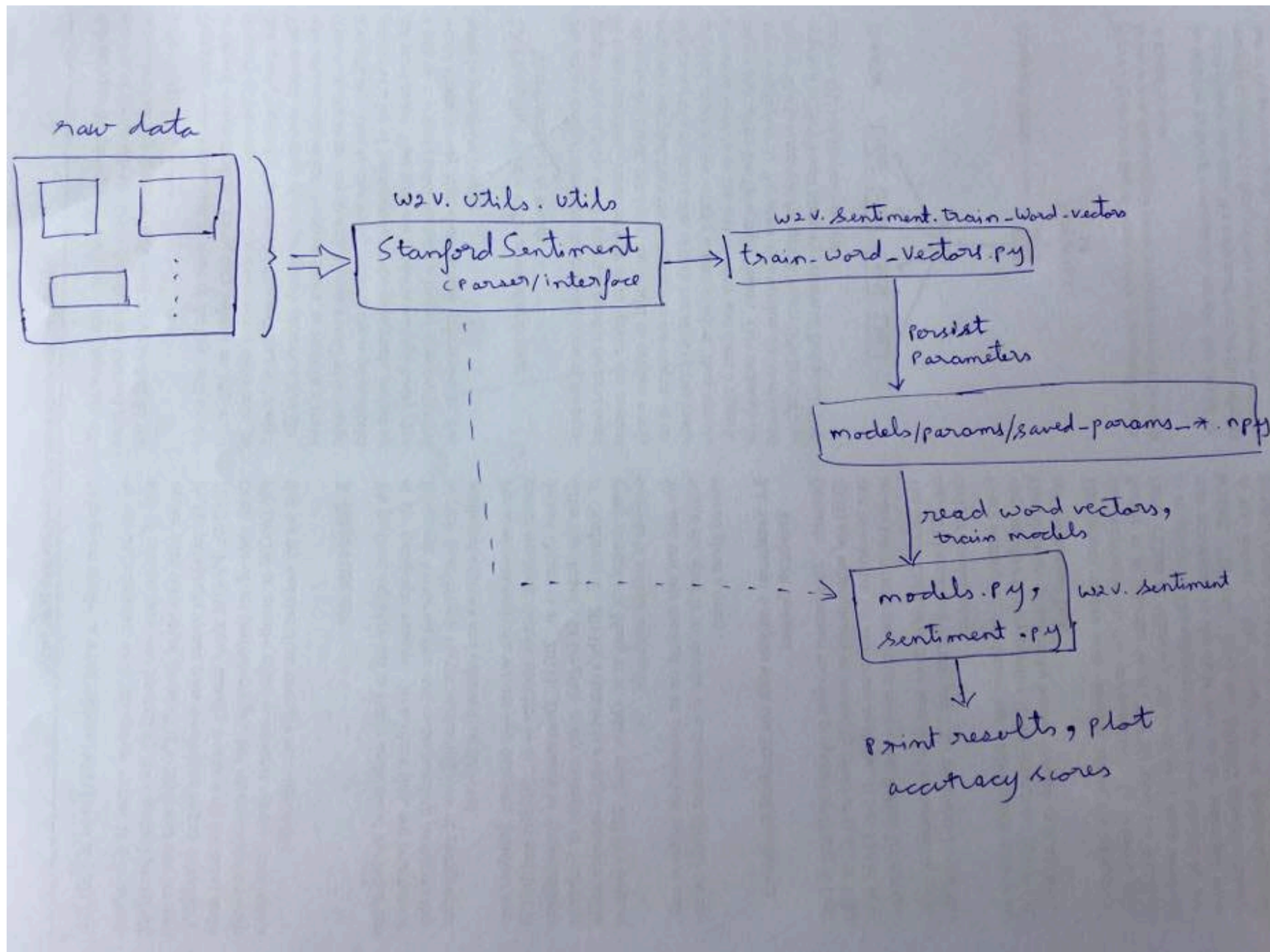
an artist who is simply tired -- of fighting the same fights , of putting the weight of the world on his shoulders , of playing with narrative forml192491

an artistic collaborationl48900



- Reference: "Recursive Deep Models for Semantic Compositionality Over a Sentiment Treebank"
See references for paper
See w2v/utils/utils.py for parsing code

Sentiment analysis pipeline



Sentiment analysis:

- Train word vectors using the word2vec package
- Learn a softmax-regression model: map from word vectors to sentiment scores
- Use L2 regularisation!
- Softmax regression cost function: $\frac{1}{N} * \sum(\text{cross_entropy}(x_i, y_i)) + \frac{1}{2} * |w|^2$

Sentiment analysis pipeline

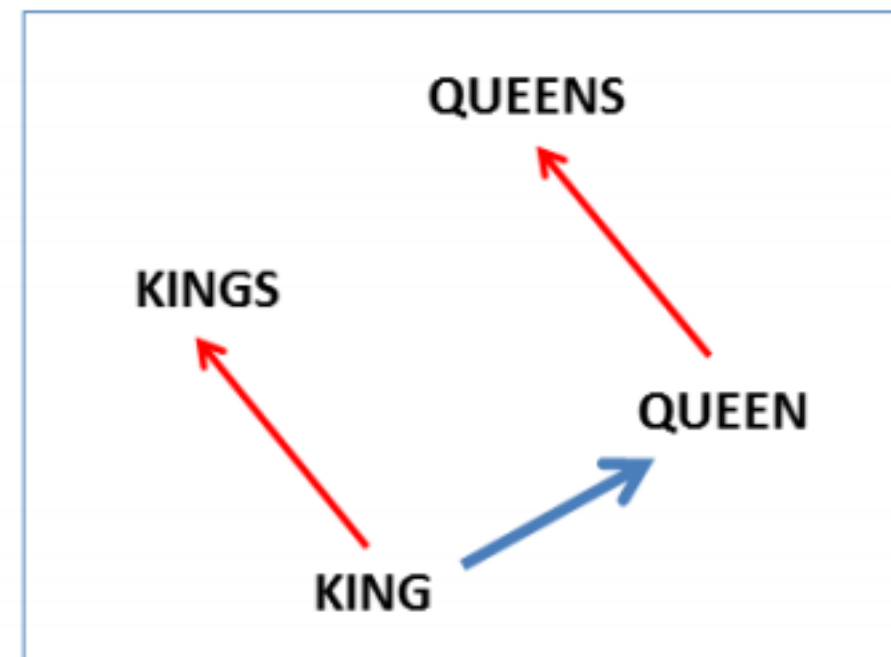
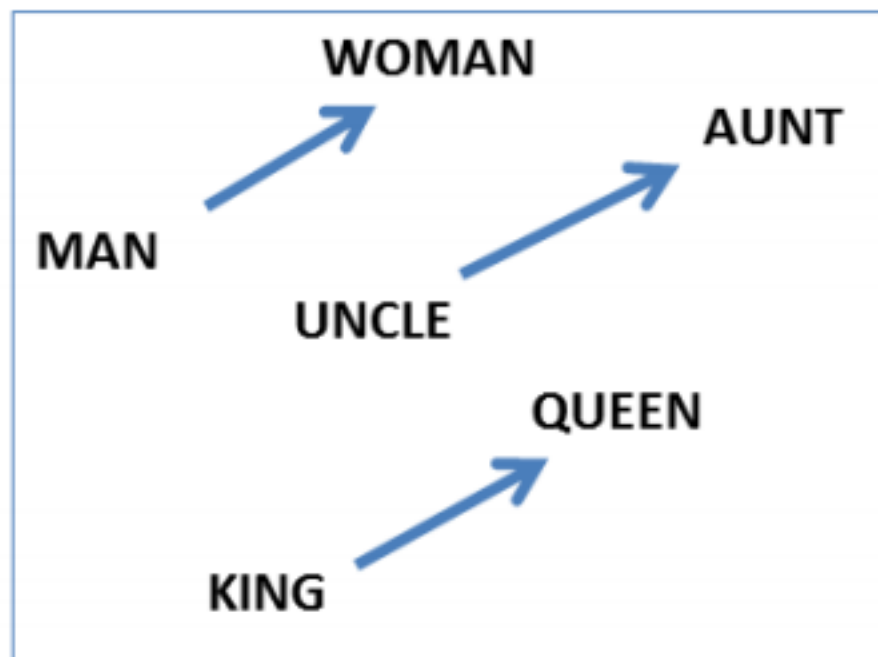
- 1-button run project:
 - make sentiment

Word2Vec

- Word2Vec: an algorithm to represent words with vectors:
- Fundamental step in NLP: first step in many, many downstream application such as search, classification, sentiment, translation, question-answering, etc.

Word2Vec

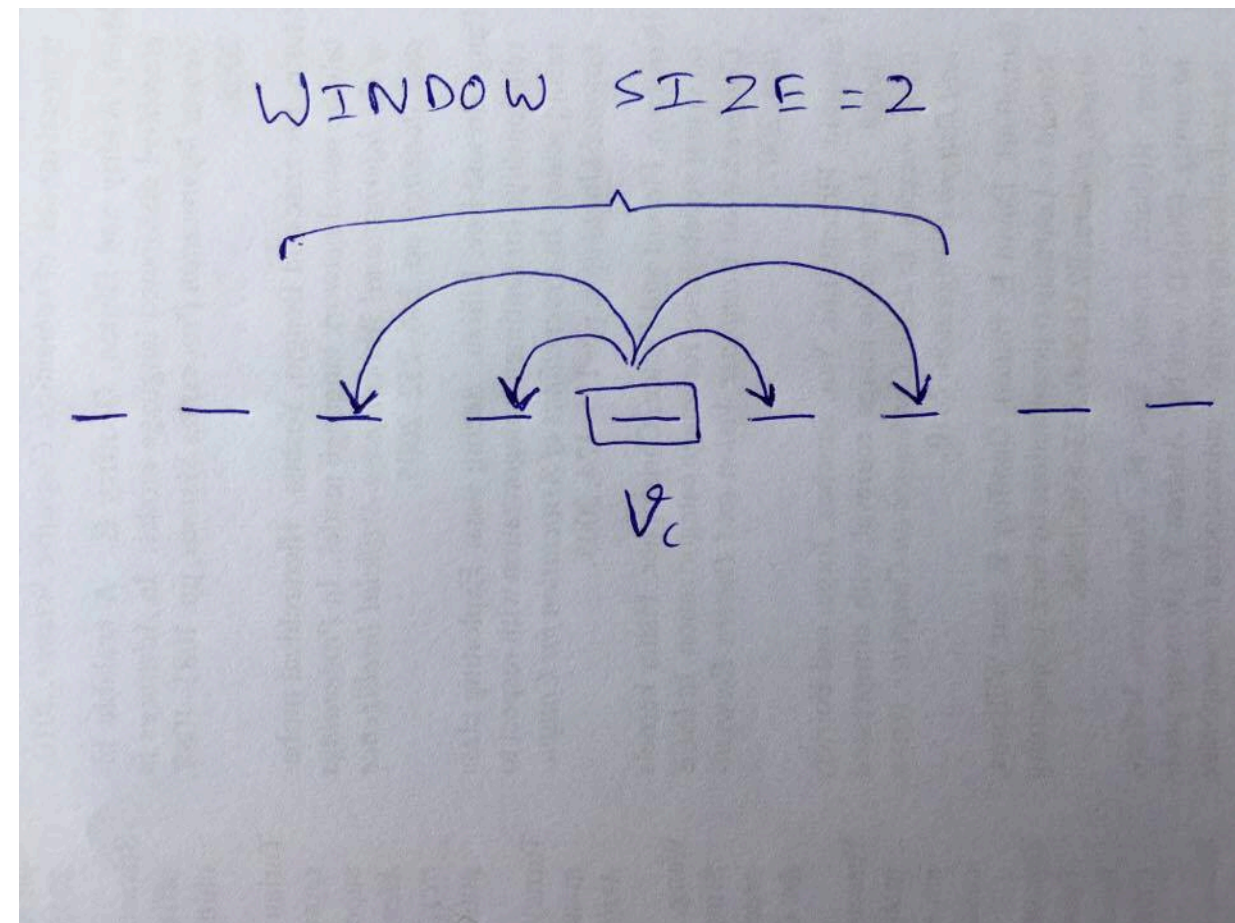
- man:woman :: king: ? \rightarrow queen



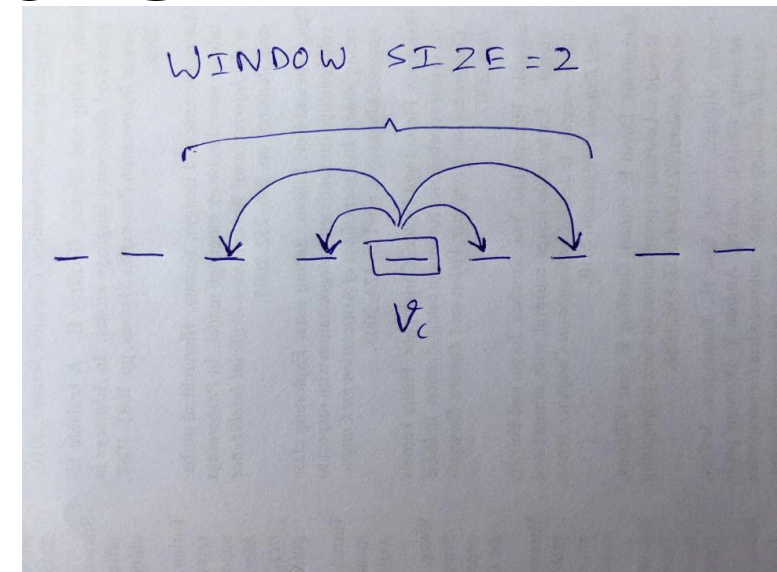
(Mikolov et al., NAACL HLT, 2013)

Word2Vec: Intuition

- Window based approach
- “You shall know a word by the company it keeps”



Word2Vec: Intuition



- Increase the probability of correct word co-occurring, decrease the probability of other words co-occurring.
- Probability of words co-occurring:

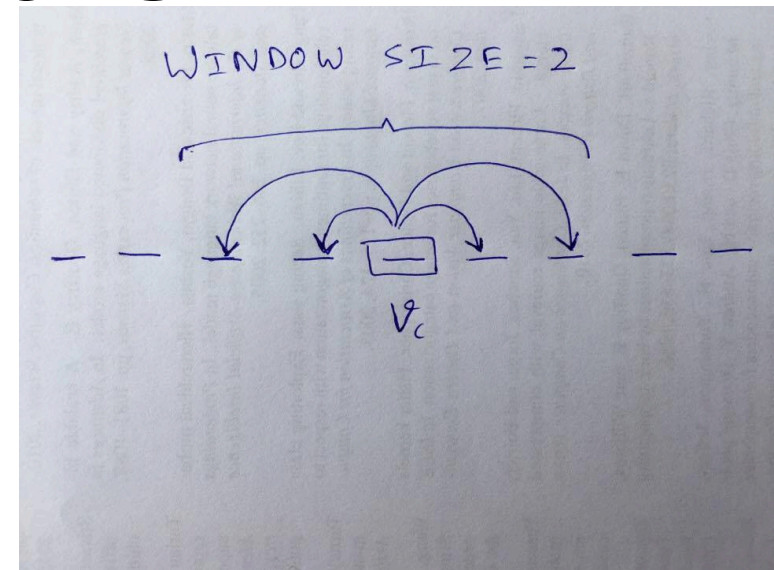
$$p(o|c) = \frac{\exp(u_o^T v_c)}{\sum_{w=1}^W \exp(u_w^T v_c)}$$

Word2Vec: Intuition

- Express our intuition in terms of a cost function :

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

- Minimising the above cost function will result in word vectors that “best” express our intuition.



Word2Vec: Details

- Actually, we have 2 variations of the cost function:
- Softmax-CE:

$$\hat{y}_o = p(\mathbf{o} \mid \mathbf{c}) = \frac{\exp(\mathbf{u}_o^\top \mathbf{v}_c)}{\sum_{w=1}^W \exp(\mathbf{u}_w^\top \mathbf{v}_c)}$$

$$CE(\mathbf{y}, \hat{\mathbf{y}}) = - \sum_i y_i \log(\hat{y}_i)$$

$$J_{softmax-CE}(\mathbf{o}, \mathbf{v}_c, \mathbf{U}) = CE(\mathbf{y}, \hat{\mathbf{y}})$$

Word2Vec: Details

- Actually, we have 2 variations of the cost function:
- Negative sampling:

$$J_{neg-sample}(\mathbf{o}, \mathbf{v}_c, \mathbf{U}) = -\log(\sigma(\mathbf{u}_o^\top \mathbf{v}_c)) - \sum_{k=1}^K \log(\sigma(-\mathbf{u}_k^\top \mathbf{v}_c))$$

Word2Vec: Details

- We also have 2 models for word2vec:
- Skipgram:

$$J_{\text{skip-gram}}(\text{word}_{c-m\dots c+m}) = \sum_{-m \leq j \leq m, j \neq 0} F(\mathbf{w}_{c+j}, \mathbf{v}_c)$$

Word2Vec: Details

- We also have 2 models for word2vec:
- CBOW:

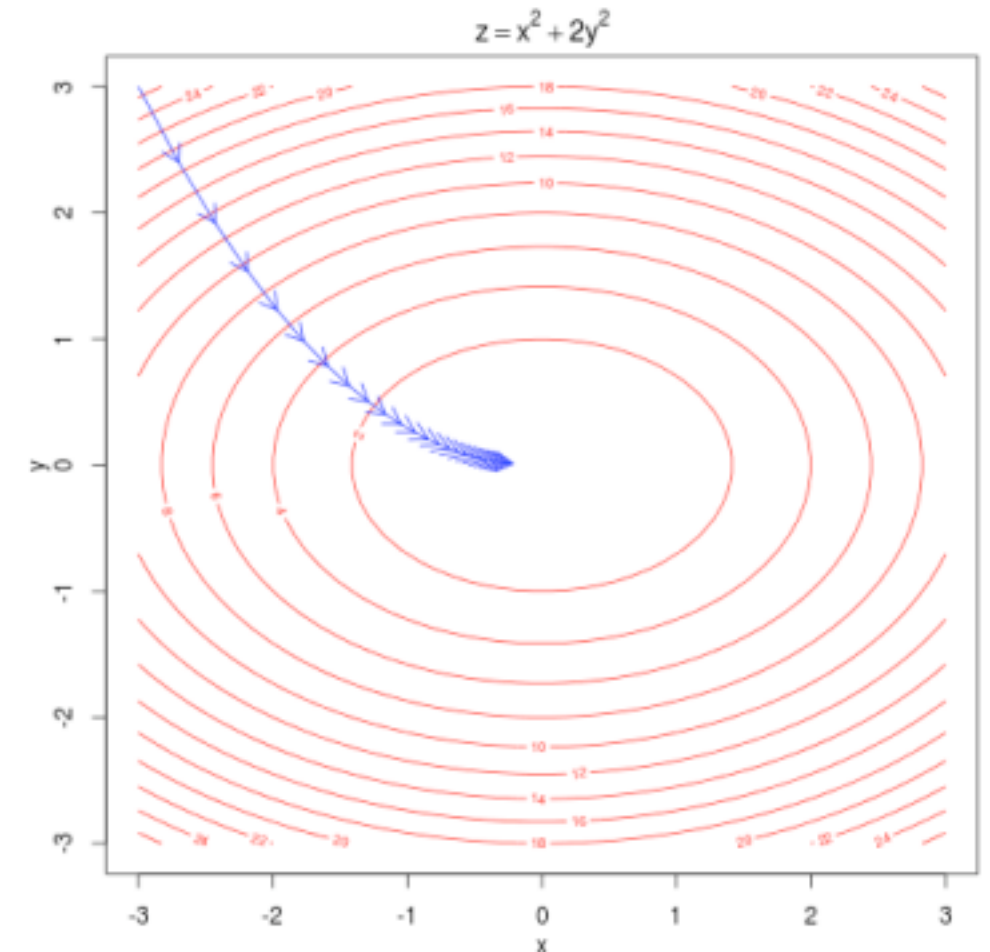
$$\hat{\mathbf{v}} = \sum_{-m \leq j \leq m, j \neq 0} \mathbf{v}_{c+j}$$

$$J_{\text{CBOW}}(\text{word}_{c-m \dots c+m}) = F(\mathbf{w}_c, \hat{\mathbf{v}})$$

Word2Vec: SGD

- SGD: numerically minimise a differentiable function

$$\theta^{new} = \theta^{old} - \alpha \nabla_{\theta} J(\theta)$$



```
while True:
    theta_grad = evaluate_gradient(J, corpus, theta)
    theta = theta - alpha * theta_grad
```


Differentiation

- We compute gradients from scratch
- Empirical gradient checking

Word2Vec: Differentiation

- Gradients for softmax-CE:

$$\frac{\partial J}{\partial \mathbf{v}_c} = U^T (\hat{\mathbf{y}} - \mathbf{y}).$$

$$\frac{\partial J}{\partial U} = \mathbf{v}_c (\hat{\mathbf{y}} - \mathbf{y})^\top$$

Word2Vec: Differentiation

- Gradients for negative sampling:

$$\frac{\partial J}{\partial \mathbf{v}_c} = (\sigma(\mathbf{u}_o^\top \mathbf{v}_c) - 1)\mathbf{u}_o - \sum_{k=1}^K (\sigma(-\mathbf{u}_k^\top \mathbf{v}_c) - 1)\mathbf{u}_k$$

$$\frac{\partial J}{\partial \mathbf{u}_o} = (\sigma(\mathbf{u}_o^\top \mathbf{v}_c) - 1)\mathbf{v}_c$$

$$\frac{\partial J}{\partial \mathbf{u}_k} = -(\sigma(-\mathbf{u}_k^\top \mathbf{v}_c) - 1)\mathbf{v}_c, \quad \text{for all } k = 1, 2, \dots, K$$

Word2Vec: Differentiation

- Gradients for skipgram:

$$\frac{\partial J_{\text{skip-gram}}(\text{word}_{c-m \dots c+m})}{\partial \mathbf{U}} = \sum_{-m \leq j \leq m, j \neq 0} \frac{\partial F(\mathbf{w}_{c+j}, \mathbf{v}_c)}{\partial \mathbf{U}},$$

$$\frac{\partial J_{\text{skip-gram}}(\text{word}_{c-m \dots c+m})}{\partial \mathbf{v}_c} = \sum_{-m \leq j \leq m, j \neq 0} \frac{\partial F(\mathbf{w}_{c+j}, \mathbf{v}_c)}{\partial \mathbf{v}_c},$$

$$\frac{\partial J_{\text{skip-gram}}(\text{word}_{c-m \dots c+m})}{\partial \mathbf{v}_j} = \mathbf{0}, \text{ for all } j \neq c.$$

Word2Vec: Differentiation

- Gradients for CBOW:

$$\frac{\partial J_{\text{CBOW}}(\text{word}_{c-m \dots c+m})}{\partial \mathbf{U}} = \frac{\partial F(\mathbf{w}_c, \hat{\mathbf{v}})}{\partial \mathbf{U}},$$

$$\frac{\partial J_{\text{CBOW}}(\text{word}_{c-m \dots c+m})}{\partial \mathbf{v}_j} = \frac{\partial F(\mathbf{w}_c, \hat{\mathbf{v}})}{\partial \hat{\mathbf{v}}},$$

$$\frac{\partial J_{\text{CBOW}}(\text{word}_{c-m \dots c+m})}{\partial \mathbf{v}_j} = \mathbf{0}, \quad \text{for all } j.$$