

# Word Embeddings, Sentiment Classification

Paulina Grnarova, Andrew An Bian

April 6-7, 2017

# Overview

## Word Embeddings

What are word embeddings?

GloVe

Evaluation

Tips and Extensions

## Pen and Paper

## Sentiment classification

What is sentiment classification?

Baseline

Classifiers

Possible improvement of the baseline

# What is a word embedding?

Suppose you are given a dictionary of words.

The  $i$ -th word  $w_i$  in the dictionary is represented by an embedding  $\mathbf{x}_{w_i} \in \mathbb{R}^d$ , i.e. a  $d$ -dimensional latent vector, which is learned.

- ▶ Captures the meaning of the word
- ▶ Similar words should have similar embeddings (share latent features)
- ▶ Angles and distances between embeddings should relate to meaning

The embedding is a way of representing (word) meaning.

# Discrete Representation

Represent a vector by its index in the vocabulary

$[0\ 0\ 0\ 1\ 0\ 0\ 0\ \dots\ 0\ 0\ 0]$  - "one-hot" vector representation

Problems:

- ▶ Dimensionality

Wikipedia + Gigaword (400K vocab), English Language (1M vocab), Twitter-2B tweets (1.2M vocab)

- ▶ Do not capture similarity of words

*good* =  $[0\ 1\ 0\ 0\ 0\ 0\ 0\ 0\ 0\ 0\ 0\ 0\ 0\ 0\ 0]$

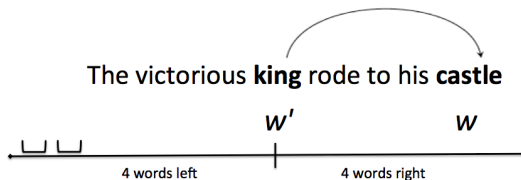
*great* =  $[0\ 0\ 0\ 0\ 1\ 0\ 0\ 0\ 0\ 0\ 0\ 0\ 0\ 0\ 0]$

*milk* =  $[0\ 0\ 1\ 0\ 0\ 0\ 0\ 0\ 0\ 0\ 0\ 0\ 0\ 0\ 0]$

*good* AND *great* = *good* AND *milk* = 0

# Distributional Similarity Based Representations

Represent a word by its neighbors - key idea in modern NLP



- ▶ Words with similar meanings have similar neighbors.
- ▶ How to represent a vector by its context?
  - ▶ Use a co-occurrence matrix

# Latent Vector Model: Basic Model

- ▶ Latent vector representation of words = embedding

$$w \mapsto (\mathbf{x}_w, b_w) \in \mathbb{R}^{D+1}, \quad (\text{vector} + \text{bias})$$

- ▶ Define **log-bilinear** model

$$\log p_{\theta}(w \mid w') = \langle \mathbf{x}_w, \mathbf{x}_{w'} \rangle + b_w + \text{const.}$$

- ▶ symmetric bilinear form fitted to log-probabilities
- ▶ normalization constant (see below)

# Latent Vector Model: Basic Model (cont'd)

- ▶ Exponentiating  
⇒ **soft-max**

$$p_{\theta}(w \mid w') = \frac{\exp [\langle \mathbf{x}_w, \mathbf{x}_{w'} \rangle + b_w]}{Z_{\theta}(w')}$$

- ▶ partition function (normalization constant):

$$Z_{\theta}(w') := \sum_{v \in \mathcal{V}} \exp [\langle \mathbf{x}_v, \mathbf{x}_{w'} \rangle + b_v]$$

- ▶ model parameters:

$$\theta = ((\mathbf{x}_w, b_w)_{w \in \mathcal{V}}) \in \mathbb{R}^{(D+1) \cdot |\mathcal{V}|}$$

# Latent Vector Model: Challenges

- Log-likelihood of basic model

$$\mathcal{L}(\theta; \mathbf{w}) = \sum_{t=1}^T \sum_{\Delta \in \mathcal{I}} \left[ \begin{aligned} & b_{w(t+\Delta)} \quad \text{ok} \\ & + \langle \mathbf{x}_{w(t+\Delta)}, \mathbf{x}_{w(t)} \rangle \quad \text{bi-linear} \leftarrow \#1 \\ & - \log \sum_{v \in \mathcal{V}} \exp [\langle \mathbf{x}_v, \mathbf{x}_{w(t)} \rangle + b_v] \quad \text{large cardinality} \leftarrow \#2 \end{aligned} \right]$$



# GloVe: Co-occurrence Matrix

Summarize data in **co-occurrence matrix**

$$\mathbf{N} = (n_{ij}) \in \mathbb{R}^{|\mathcal{V}| \cdot |\mathcal{C}|},$$

$n_{ij} = \#$  occurrences of  $w_i \in \mathcal{V}$  in context of  $w_j \in \mathcal{C}$

Example corpus:

- ▶ The king rode to his castle.
- ▶ The king lives in the castle.

counts	the	king	rode	lives	to	in	his	castle
the	0	2	0	0	0	1	0	1
king	2	0	1	1	0	0	0	0
rode	0	1	0	0	1	0	0	0
lives	0	1	0	0	0	1	0	0
to	0	0	1	0	0	0	1	0
in	1	0	0	1	0	0	0	0
his	0	0	0	0	1	0	0	1
castle	1	0	0	0	0	0	1	0

SVD?

# GloVe Objective

- ▶ Weighted least squares fit of log-counts

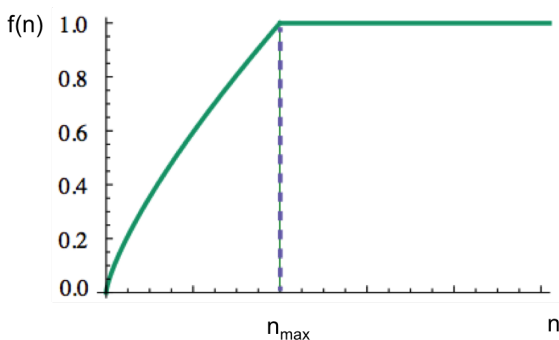
$$\mathcal{H}(\theta; \mathbf{N}) = \sum_{i,j} f(n_{ij}) \left( \underbrace{\log n_{ij}}_{\text{target}} - \underbrace{\log \tilde{p}_{\theta}(w_i|w_j)}_{\text{model}} \right)^2,$$

with **unnormalized** distribution

$$\tilde{p}_{\theta}(w_i|w_j) = \exp [\langle \mathbf{x}_i, \mathbf{y}_j \rangle + b_i + d_j]$$

and **weighting function**  $f$

# GloVe Weighting



- ▶ Scalable to large corpora
- ▶ Fast training
- ▶ Limit influence of large counts (very frequent words)

# How to optimize the objective of GloVe?

- ▶ Non-convex problem: hard to find global minimum
- ▶ Goal: Minimize the objective/cost function
- ▶ Use gradient descent method

Trivial example: Find a local minimum of the function

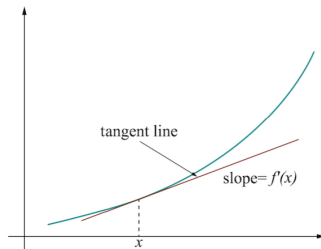
$$f(x) = x^2 - 6x, \text{ where } f'(x) = 2x - 6.$$

```
x_old = 0
x_new = 6
eps = 0.01
precision = 0.00001
```

```
def f_derivative(x):
    return 2*x - 6
```

```
while abs(x_new - x_old) > precision:
    x_old = x_new
    x_new = x_old - eps * f_derivative(x_old)
```

```
print("Local minimum occurs at ", x_new)
```



# GloVe Optimization (no!)

- ▶ Non-convex problem: hard to find global minimum
- ▶ Gradient descent (aka **steepest descent**)

$$\theta^{\text{new}} \leftarrow \theta^{\text{old}} - \eta \nabla_{\theta} \mathcal{H}(\theta; \mathbf{N}), \quad \eta > 0 \text{ (step size)}$$

- ▶  $\theta = ((\mathbf{x}_w, b_w)_{w \in \mathcal{V}}, (\mathbf{y}_w, d_w)_{w \in \mathcal{C}})$ , embeddings = parameters

To minimize over the full batch (the entire training data) would require us to compute gradients for all entries in the co-occurrence matrix w.r.t. all parameters

# GloVe Optimization (yes!)

There might be billions of entries in the co-occurrence matrix!  
Long waiting time before a single update

- ▶ Non-convex problem: hard to find global minimum
- ▶ Use stochastic optimization to find local minimum
- ▶ Stochastic gradient descent (SGD):
  - ▶ sample  $(i, j)$  such that  $n_{ij} > 0$  uniformly at random
  - ▶ perform "cheap" update (single entry and sparse)

$$\mathbf{x}_i^{\text{new}} \leftarrow \mathbf{x}_i + 2\eta f(n_{ij}) (\log n_{ij} - \langle \mathbf{x}_i, \mathbf{y}_j \rangle) \mathbf{y}_j$$

$$\mathbf{y}_j^{\text{new}} \leftarrow \mathbf{y}_j + 2\eta f(n_{ij}) (\log n_{ij} - \langle \mathbf{x}_i, \mathbf{y}_j \rangle) \mathbf{x}_i$$

# How to evaluate word vectors?

## Intrinsic versus extrinsic evaluation

- ▶ Intrinsic
  - ▶ Evaluate on a specific or intermediate task
  - ▶ Fast to compute
  - ▶ Leads to better understanding of the system
  - ▶ Usefulness depends on correlation to the realistic task
- ▶ Extrinsic
  - ▶ Evaluation on a real task
  - ▶ Can take a long time to compute accuracy
  - ▶ Hard to debug

# Practical Tips

- ▶ What to do with the two sets of vectors?
  - ▶ Word vectors  $\mathbf{x}_i$
  - ▶ Context vectors  $\mathbf{y}_i$
- ▶ Both capture similar co-occurrence
- ▶ In order to get the final embeddings, a simple and efficient way is to sum them up

$$X_{final} = X + Y$$

- ▶ Different variations explored in Global Vectors for Word Representation (Pennington et al. (2014))





# Pen and Paper Assignment

## Problem 1:

1) The objective of GloVe is

$$\mathcal{H}(\theta; \mathbf{N}) = \sum_{i,j} f(n_{ij}) \left( \underbrace{\log n_{ij}}_{\text{target}} - \underbrace{\log \tilde{p}_{\theta}(w_i|w_j)}_{\text{model}} \right)^2.$$

Suppose  $f(.) = 1$  for all arguments, and  $m_{ij} = \log n_{ij}$ .

- ▶ a) Derive the gradient of the objective function of GloVe w.r.t.  $\mathbf{x}_i$  and  $\mathbf{y}_j$
- ▶ b) Derive the stochastic gradient of the objective function of GloVe w.r.t.  $\mathbf{x}_i$  and  $\mathbf{y}_j$

# Pen and Paper Assignment

## Problem 1:

2) Show that GloVe with  $f(n_{ij}) := \begin{cases} 1 & \text{if } n_{ij} > 0, \\ 0 & \text{otherwise.} \end{cases}$  solves a **matrix completion** problem

$$\min_{\mathbf{X}, \mathbf{Y}} \sum_{ij: n_{ij} > 0} \left( m_{ij} - (\mathbf{X}^\top \mathbf{Y})_{ij} \right)^2 .$$

3) Derive the gradient and stochastic gradient of  $H$  for any weighting function  $f$ .

# What is sentiment classification?

- ▶ Given a tweet predict whether it has a positive or negative opinion  
e.g.: if a tweet message contains a :) or :(
- ▶ Examples (the smiley has been intentionally removed):
  - ▶ “i know android sucks :(”
  - ▶ “twitter is dead right now :(”
  - ▶ “my sis made apple crisp with extra crisp ! it's awesome :)”
  - ▶ “i hope your wednesday was awesome :)”

## Project 2: Dataset provided

- ▶ Twitter data:
  - ▶ `train_pos_full.txt` ~ 1M tweets that contained :)
  - ▶ `train_neg_full.txt` ~ 1M tweets that contained :(
  - ▶ `train_pos.txt` - 10% from the positive tweets for training
  - ▶ `train_neg.txt` - 10% from the negative tweets for training
  - ▶ `test_data.txt` - 10K unlabeled tweets
- ▶ Each tweet contains at most 140 characters
- ▶ All tweets are tokenized - words are separated by a single whitespace
- ▶ All labels (smileys) are removed
- ▶ User mentions replaced with `<user>`
- ▶ Links replaced with `<url>`

# Simple baseline using word embeddings

- ▶ Average the word embeddings to get an embedding for the whole tweet message
- ▶ Feed the resulting word embedding to a classifier
- ▶ If the tweet has positive sentiment, its embedding is close to the words that represent positive meaning (close to good, great, amazing and far from bad, horrible etc.)

# Which classifiers to use?

- ▶ Any classifier you would like
- ▶ Logistic regression, Support Vector Machine (SVM), Gaussian process classifier, neural network, etc
- ▶ Even using ensemble methods.
- ▶ No need to go deep into classifiers so far, can use library :-)

# Logistic regression from Scikit

- ▶ Logistic regression classification:
  - ▶ a linear model for classification rather than regression.
  - ▶ a.k.a. logit regression, maximum-entropy classification (MaxEnt) or the log-linear classifier
  - ▶ The probabilities describing the possible outcomes of a single trial are modeled using a logistic function.
- ▶ Example using it from Scikit-learn
  - ▶ Scikit-learn: Machine learning in Python
  - ▶ Exemplar code for logistic regression of the Iris dataset



# Python source code: plot\_iris\_logistic.py

## Import data

```
print(__doc__)  
# Code source: Gael Varoquaux  
  
import numpy as np  
import matplotlib.pyplot as plt  
from sklearn import linear_model, datasets  
  
# import some data to play with  
iris = datasets.load_iris()  
X = iris.data[:, :2] # we only take the first two features.  
Y = iris.target  
h = .02 # step size in the mesh
```

## Train the model

```
logreg = linear_model.LogisticRegression(C=1e5)  
logreg.fit(X, Y)
```

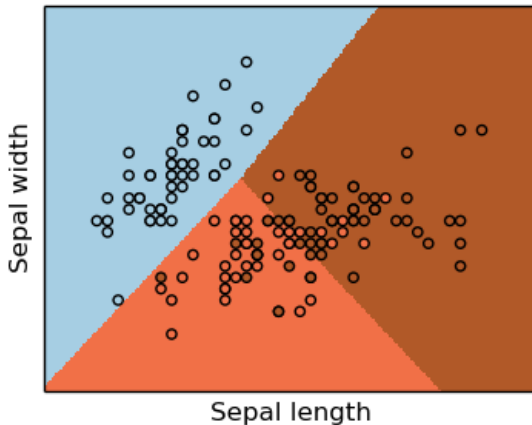
## Python source code: show results

```
# Plot the decision boundary.
x_min, x_max = X[:, 0].min() - .5, X[:, 0].max() + .5
y_min, y_max = X[:, 1].min() - .5, X[:, 1].max() + .5
xx, yy = np.meshgrid(np.arange(x_min, x_max, h), np.arange(y_min, y_max, h))
Z = logreg.predict(np.c_[xx.ravel(), yy.ravel()])

# Put the result into a color plot
Z = Z.reshape(xx.shape)
plt.figure(1, figsize=(4, 3))
plt.pcolormesh(xx, yy, Z, cmap=plt.cm.Paired)
# Plot also the training points
plt.scatter(X[:, 0], X[:, 1], c=Y, edgecolors='k', cmap=plt.cm.Paired)
plt.xlabel('Sepal length')
plt.ylabel('Sepal width')
plt.xlim(xx.min(), xx.max())
plt.ylim(yy.min(), yy.max())
plt.xticks(())
plt.yticks(())

plt.show()
```

# Python source code: The decision boundary



# Why is the baseline overly simplistic?

- ▶ If the tweet consists of only positive (negative) words it may correctly predict the sentiment.
- ▶ However, it will fail with double negation or mixed words. e.g.  
“It was not horrible and definitely not boring.”,  
”It started pretty bad, but it turned out to be amazing”.
- ▶ The task is to improve the baseline taking this into consideration.

# Possible improvement–Retraining the word embeddings

- ▶ Init word embeddings using GloVe
- ▶ Treat them as parameters
- ▶ Retrain the word embeddings during the classification

