

Deep Learning for Natural Language Processing

Cornelia Caragea

Computer Science

University of Illinois at Chicago

Credits for slides: Manning, Socher

Word Vectors - Glove

Today

Lecture Structure:

- The GloVe model of word vectors
- Evaluating Glove word vectors

Recommended reading:

- GloVe: Global Vectors for Word Representation

Count based vs. direct prediction

- LSA, HAL (Lund & Burgess),
- COALS, Hellinger-PCA (Rohde et al, Lebrecht & Collobert)

- Fast training (For small matrices)
- Efficient usage of statistics
- Primarily used to capture word similarity
- Disproportionate importance given to large counts

- Skip-gram/CBOW (Mikolov et al)
- NNLM, HLBL, RNN (Bengio et al; Collobert & Weston; Huang et al; Mnih & Hinton)

- Scales with corpus size
- Inefficient usage of statistics
- Generate improved performance on other tasks
- Can capture complex patterns beyond word similarity

LSA - term-doc matrices; HAL (Hyperspace Analogue to Language) - term-term matrices; HPCA (Hellinger PCA) - a square root type transformation; COALS - co-occurrence matrix is transformed, e.g., by a correlation-based normalization.

The GloVe model

Notation:

- X the matrix of word-word co-occurrence counts;
 X_{ij} number of times word j occurs in the context of word i ;
 $X_i = \sum_k X_{ik}$ the number of times any word appears in the context of word i .
- $P_{ij} = P(j|i) = \frac{X_{ij}}{X_i}$ - the probability that word j appears in the context of word i .

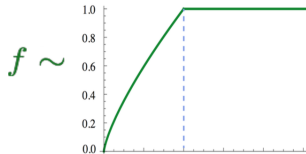
Combining the best of both worlds GloVe

[Pennington, Socher, and Manning, EMNLP 2014]

$$J = \sum_{i,j=1}^V f(P_{ij})(w_i \cdot \tilde{w}_j - \log P_{ij})^2$$

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

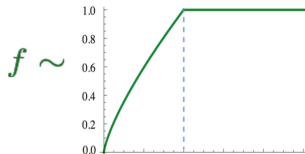
- Fast training
- Scalable to huge corpora
- Good performance even with small corpus and small vectors



Combining the best of both worlds GloVe

[Pennington, Socher, and Manning, EMNLP 2014]

1. $f(0) = 0$. If f is viewed as a continuous function, it should vanish as $x \rightarrow 0$ fast enough that the $\lim_{x \rightarrow 0} f(x) \log^2 x$ is finite.
2. $f(x)$ should be non-decreasing so that rare co-occurrences are not overweighted.
3. $f(x)$ should be relatively small for large values of x , so that frequent co-occurrences are not overweighted.



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

GloVe results

Nearest words to frog:

- ① frogs
- ② toad
- ③ litoria
- ④ leptodactylidae
- ⑤ rana
- ⑥ lizard
- ⑦ eleutherodactylus



litoria



leptodactylidae



rana



eleutherodactylus

How to evaluate word vectors?

- Related to general evaluation in NLP: Intrinsic vs extrinsic
- Intrinsic:
 - Evaluation on a specific/intermediate subtask
 - Fast to compute
 - Helps to understand that system
 - Not clear if really helpful unless correlation to real task is established
- Extrinsic:
 - Evaluation on a real task
 - Can take a long time to compute accuracy
 - Unclear if the subsystem is the problem or its interaction with other subsystems
 - If replacing exactly one subsystem with another improves accuracy → Winning!

Intrinsic word vector evaluation

- Word Vector Analogies

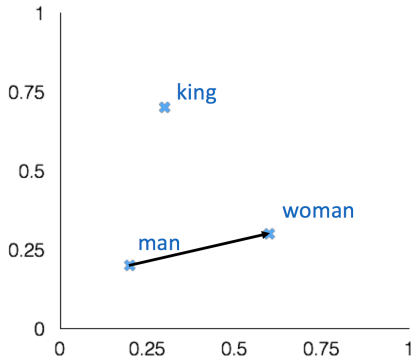
a:b :: c:?



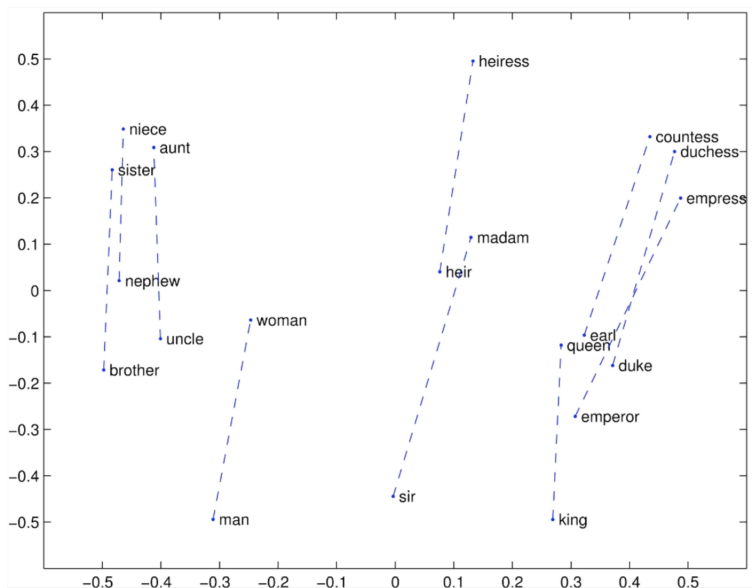
$$d = \arg \max_i \frac{(x_b - x_a + x_c)^T x_i}{\|x_b - x_a + x_c\|}$$

man:woman :: king:?

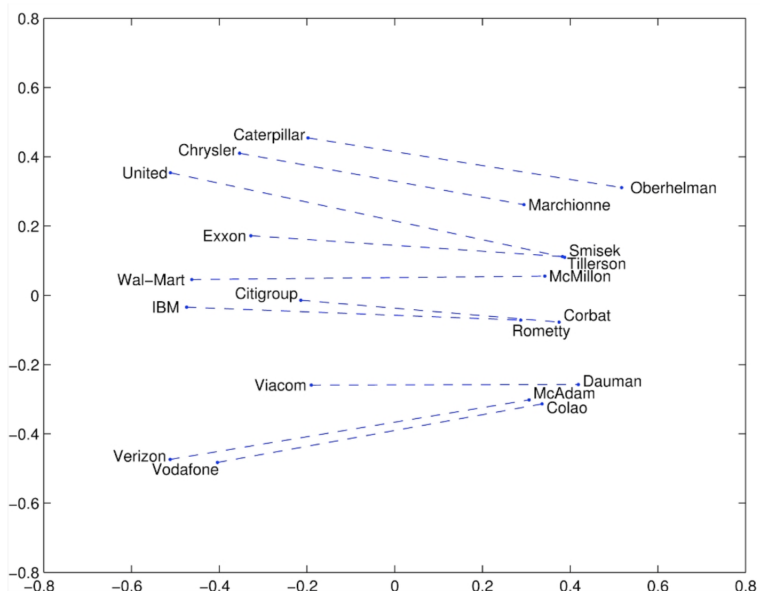
- Evaluate word vectors by how well their cosine distance after addition captures intuitive semantic and syntactic analogy questions
- Discarding the input words from the search!



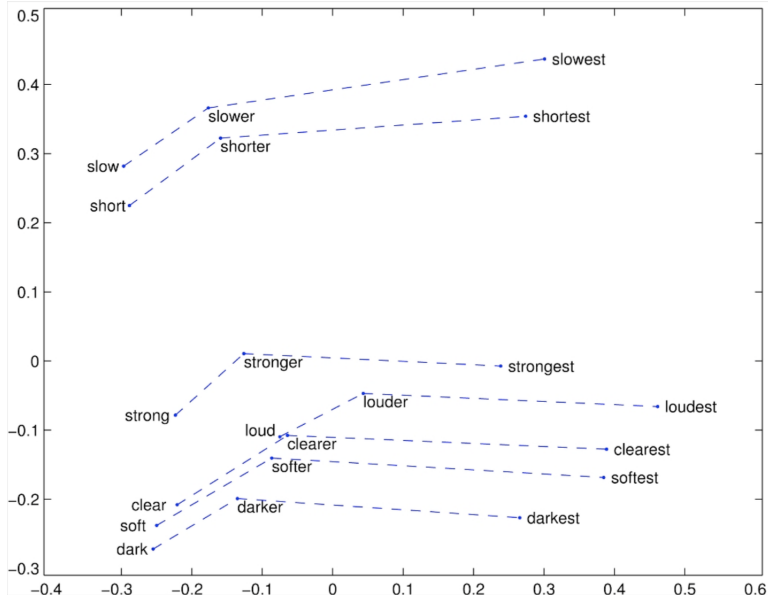
Glove Visualizations - Similar Vector Offsets



Glove Visualizations: Company - CEO



Glove Visualizations: Superlatives



Details of intrinsic word vector evaluation

- Word Vector Analogies: Syntactic and **Semantic** examples from <http://code.google.com/p/word2vec/source/browse/trunk/questions-words.txt>

: city-in-state

Chicago Illinois Houston Texas

Chicago Illinois Philadelphia Pennsylvania

Chicago Illinois Phoenix Arizona

Chicago Illinois Dallas Texas

Chicago Illinois Jacksonville Florida

Chicago Illinois Indianapolis Indiana

Chicago Illinois Austin Texas

Chicago Illinois Detroit Michigan

Chicago Illinois Memphis Tennessee

Chicago Illinois Boston Massachusetts

Problem: different cities may have same name

Details of intrinsic word vector evaluation

- Word Vector Analogies: **Syntactic** and Semantic examples from <http://code.google.com/p/word2vec/source/browse/trunk/questions-words.txt>

: gram4-superlative

bad worst big biggest

bad worst bright brightest

bad worst cold coldest

bad worst cool coolest

bad worst dark darkest

bad worst easy easiest

bad worst fast fastest

bad worst good best

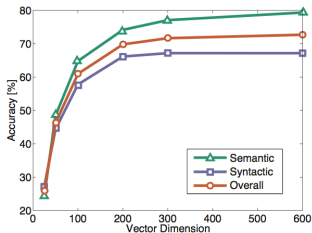
bad worst great greatest

Analogy evaluation and hyperparameters

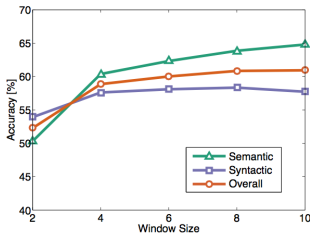
- Glove word vectors evaluation

Model	Dim.	Size	Sem.	Syn.	Tot.
ivLBL	100	1.5B	55.9	50.1	53.2
HPCA	100	1.6B	4.2	16.4	10.8
GloVe	100	1.6B	<u>67.5</u>	<u>54.3</u>	<u>60.3</u>
SG	300	1B	61	61	61
CBOW	300	1.6B	16.1	52.6	36.1
vLBL	300	1.5B	54.2	<u>64.8</u>	60.0
ivLBL	300	1.5B	65.2	63.0	64.0
GloVe	300	1.6B	<u>80.8</u>	61.5	<u>70.3</u>
SVD	300	6B	6.3	8.1	7.3
SVD-S	300	6B	36.7	46.6	42.1
SVD-L	300	6B	56.6	63.0	60.1
CBOW [†]	300	6B	63.6	<u>67.4</u>	65.7
SG [†]	300	6B	73.0	66.0	69.1
GloVe	300	6B	<u>77.4</u>	67.0	<u>71.7</u>
CBOW	1000	6B	57.3	68.9	63.7
SG	1000	6B	66.1	65.1	65.6
SVD-L	300	42B	38.4	58.2	49.2
GloVe	300	42B	<u>81.9</u>	<u>69.3</u>	<u>75.0</u>

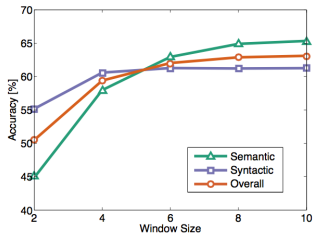
Analogy evaluation and hyperparameters



(a) Symmetric context



(b) Symmetric context



(c) Asymmetric context

Dimensionality

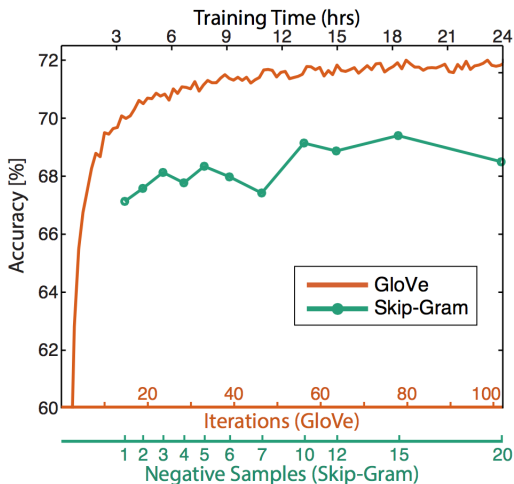
Window size

Window size

- Good dimension is ~ 300
- Asymmetric context (only words to the left) are not as good
- But this might be different for downstream tasks!
- Window size of 8 around each center word is good for Glove vectors

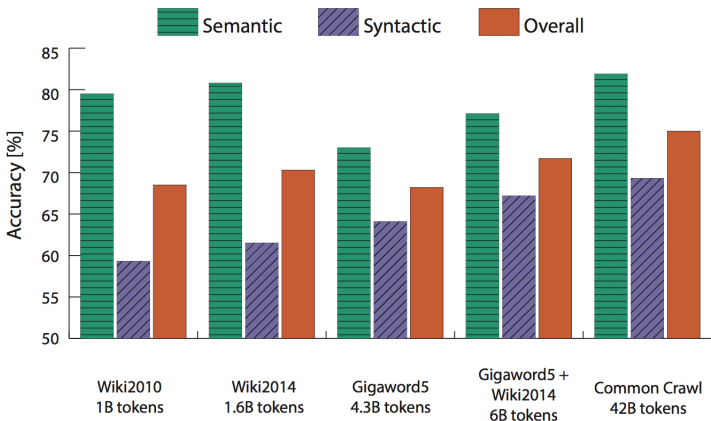
Analogy evaluation and hyperparameters

- More training time helps



Analogy evaluation and hyperparameters

- More data helps, Wikipedia is better than news text!



Neural Net Fundamentals

- We concentrate on understanding (deep, multi-layer) neural networks and how they can be trained (learned from data) using backpropagation (the judicious application of calculus)
- We will look at an NLP classifier that adds context by taking in windows around a word and classifies the center word (not just representing it across all windows)!