

ENGR-E 533 “Deep Learning Systems”

Lecture 15: Word Embedding

Minje Kim

Department of Intelligent Systems Engineering

Email: minje@indiana.edu

Website: <http://minjekim.com>

Research Group: <http://saige.sice.indiana.edu>

Meeting Request: <http://doodle.com/minje>



INDIANA UNIVERSITY

**SCHOOL OF INFORMATICS,
COMPUTING, AND ENGINEERING**

How to Represent Words?

- Once again we seek latent representations

- Say that we're interested in these words:
 - Illinois, Chicago, Champaign, Indiana, Indianapolis, Bloomington
- How do we represent them for computer algorithms to process?
 - One-hot vectors have been common
- Any problems with this?
 - Too high dimensional
 - There are about 700k words in English
 - Not really a feature vector
 - Very discriminant but that's it
 - Difficult to compare them

Illinois	1	0	0	0	0	0

Chicago	0	1	0	0	0	0

Champaign	0	0	1	0	0	0
Indiana	0	0	0	1	0	0

Indianapolis	0	0	0	0	1	0

Bloomington	0	0	0	0	0	1



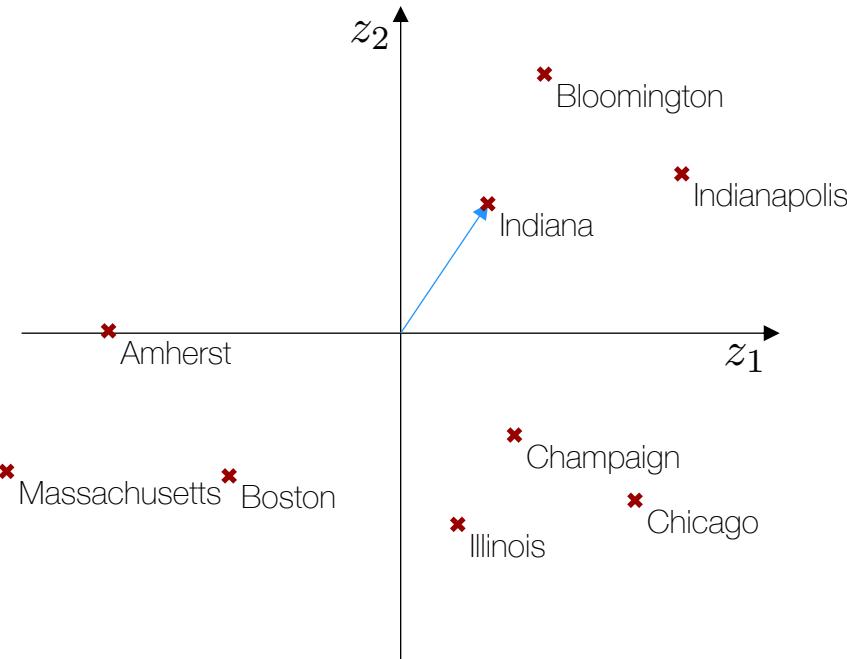
INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

How to Represent Words?

- Once again we seek latent representations

- A better latent representation



INDIANA UNIVERSITY

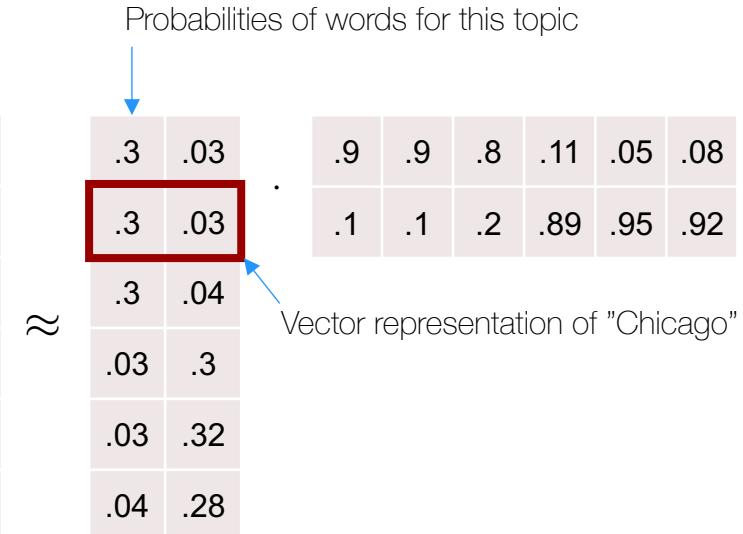
SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

Vector Representation of Words

- A matrix factorization approach

- Similarity matrix
 - e.g. co-occurrence
- Matrix factorization
 - Eigendecomposition
 - SVD
 - NMF
 - Topic modeling
 - You name it
- Implication
 - Each word has a K-d vector representation
 - K is the low rank of your choice
 - Similar words often appear in the same topic together
 - Each word has a unique activation pattern over the latent variables

	Illinois	Chicago	Champaign	Indiana	Indianapolis	Bloomington
Illinois	1	.9	.8	.2	.05	.04
Chicago	.9	1	.7	.02	.03	.07
Champaign	.8	.7	1	.02	.1	.3
Indiana	.2	.02	.1	1	.95	.88
Indianapolis	.05	.03	.1	.95	1	.82
Bloomington	.04	.07	.3	.88	.82	1



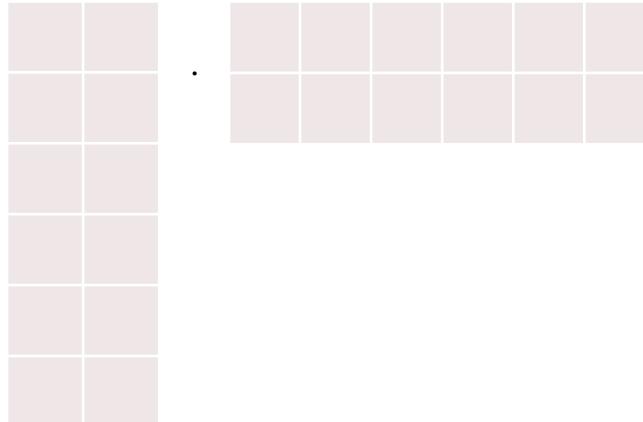
Vector Representation of Words

- A matrix factorization approach

- Same story for the Term-Frequency (TF) matrix

	d1	d2	d3	d4	...	dN
Illinois	.03
Chicago	.02
Champaign	.04
Indiana	.001
Indianapolis	.000
Bloomington	.000

≈



- An equation for this procedure $X \approx WH$
- Where am I going?
 - Autoencoder



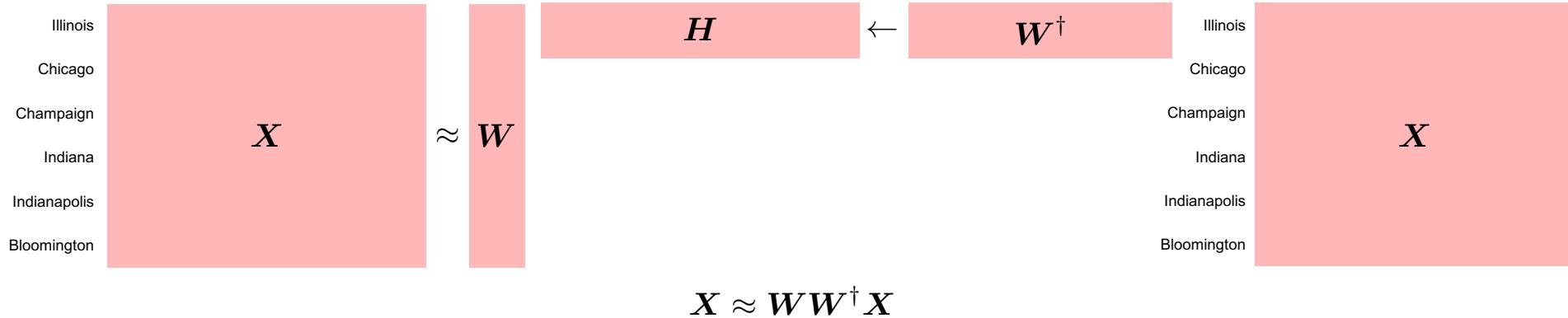
INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

Autoencoder for Word Embedding

- The basic structure

- The autoencoding process



- You'll see this is just a shallow neural network
- In order to make it practical, we need to define
 - The error function
 - Input data
 - Targets
 - Nonlinearity



INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

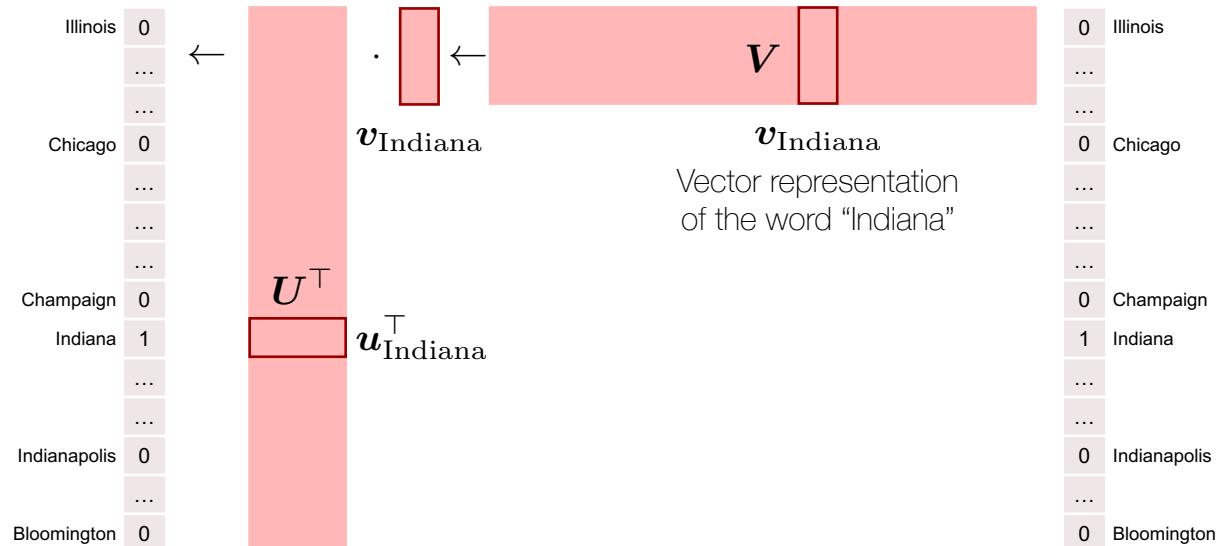
Word2vec

- Skip-gram

- Basic autoencoding network
- Activation function for the output
 - Softmax

$$p(w_t | w_t; \mathbf{V}, \mathbf{U}) = \frac{\exp(\mathbf{u}_t^\top \mathbf{v}_t)}{\sum_{t'} \exp(\mathbf{u}_{t'}^\top \mathbf{v}_t)}$$

- If this autoencoding job is successful
 - $\mathbf{v}_{\text{Indiana}}$ is totally different from any other vectors in \mathbf{U} , but $\mathbf{u}_{\text{Indiana}}$
- Is this good enough?
What's the problem with this?
 - Is "Indiana" closer to "Indianapolis" than to "Illinois"?



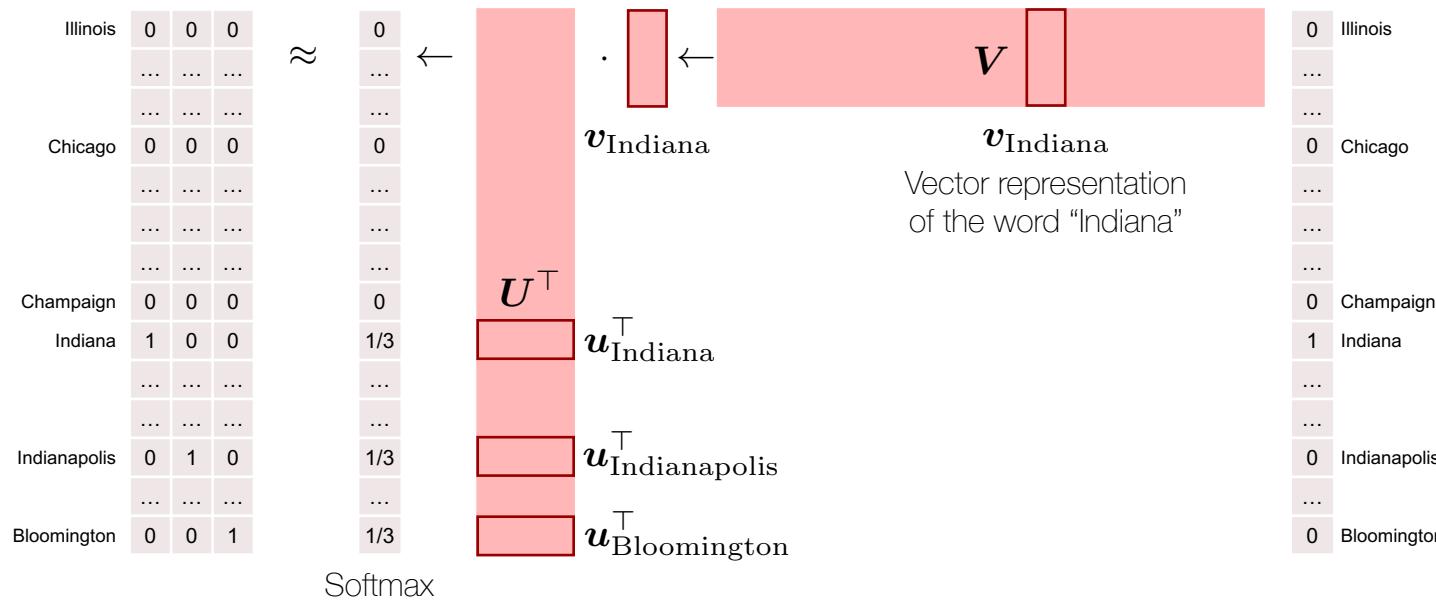
INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

Word2vec

- Skip-gram

- Another prediction mechanism that encodes co-occurrence



Word2vec

- Negative sampling

- Average negative log-likelihood

$$\mathcal{J}(\mathbf{U}, \mathbf{V}) = -\frac{1}{T} \sum_{t=1}^T \sum_{|\tau| \leq m, \tau \neq 0} \log p(w_{t+\tau} | w_t) = -\frac{1}{T} \sum_{t=1}^T \sum_{|\tau| \leq m, \tau \neq 0} \log \frac{\exp(\mathbf{u}_{t+\tau}^\top \mathbf{v}_t)}{\sum_{t'} \exp(\mathbf{u}_{t'}^\top \mathbf{v}_t)}$$

context window

- For a given input word word2vec predicts its surrounding words
- Note that word2vec doesn't predict the input word

- Any problem?

- Denominator involves all words

- Negative sampling

- Word-wise binary classification [Prob. of predicting the word of interest] versus [A randomly sampled word outside of the context]

- Suppose softmax with two classes

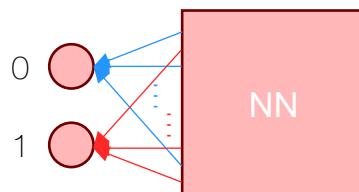
- You can maximize this instead $\log \sigma(\mathbf{u}_{t+\tau}^\top \mathbf{v}_t) + (1 - \log \sigma(\mathbf{u}_j^\top \mathbf{v}_t))$

Prob. of predicting one of the context word

Prob. of predicting a negative word

- Or, we could check on a few more negative samples

$$\log p(w_{t+\tau} | w_t) \approx \log \sigma(\mathbf{u}_{t+\tau}^\top \mathbf{v}_t) + \sum_{j \in P(w)} (\log \sigma(-\mathbf{u}_j^\top \mathbf{v}_t))$$



INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

The Other Things & Reading

- GloVe
 - Pennington, Jeffrey, Richard Socher, and Christopher Manning. "Glove: Global vectors for word representation." *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*. 2014.
- FastText <https://fasttext.cc>
 - Takes the sum of the context words as the input
 - Supported by Facebook AI Research
- Word2vec
 - Mikolov, Tomas, et al. "Efficient estimation of word representations in vector space." *arXiv preprint arXiv:1301.3781* (2013).
 - Mikolov, Tomas, et al. "Distributed representations of words and phrases and their compositionality." *Advances in neural information processing systems*. 2013.
- <https://www.tensorflow.org/tutorials/word2vec>



INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING



Thank You!



INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING