

Word Embeddings

Neural Embeddings for Political Science Research

Fridolin Linder

November 8, 2016

Methods Comprehensive Exam Presentation

Table of contents

1. Introduction
2. Representing Meaning in Text Analysis
3. Word Embeddings
4. Neural Networks
5. Continuous Bag of Words Model

Introduction

- A lot of data in political science is text

- A lot of data in political science is text
- Every corpus is analyzed in isolation

Representing Meaning in Text Analysis

Meaning of Documents

- Three documents:

Meaning of Documents

- Three documents:
 1. *The cat jumps over the dog*

Meaning of Documents

- Three documents:
 1. *The cat jumps over the dog*
 2. *The kitten hops ovqer the hound*

Meaning of Documents

- Three documents:
 1. *The cat jumps over the dog*
 2. *The kitten hops ovqer the hound*
 3. *Die Katze springt über den Hund*

Meaning of Documents

- Three documents:
 1. *The cat jumps over the dog*
 2. *The kitten hops ovqer the hound*
 3. *Die Katze springt über den Hund*
- Take the vocabulary $\mathcal{V} = \{\text{The, cat, ...}\}$ and represent them as vectors of dimension $|\mathcal{V}|$ in a Document-Term-Matrix $\mathbf{X}_{|\mathcal{D}| \times |\mathcal{V}|}$:

	the	cat	jump	over	dog	kitten	ovqer	hound	die	katze	springt	über	den	hund
D_1	2	1	1	1	1	0	0	0	0	0	0	0	0	0
D_2	2	0	0	0	0	1	1	1	0	0	0	0	0	0
D_3	0	0	0	0	0	0	0	0	1	1	1	1	1	1

Meaning of Documents

- Three documents:
 1. *The cat jumps over the dog*
 2. *The kitten hops ovqer the hound*
 3. *Die Katze springt über den Hund*
- Take the vocabulary $\mathcal{V} = \{\text{The, cat, ...}\}$ and represent them as vectors of dimension $|\mathcal{V}|$ in a Document-Term-Matrix $\mathbf{X}_{|\mathcal{D}| \times |\mathcal{V}|}$:

	the	cat	jump	over	dog	kitten	ovqer	hound	die	katze	springt	über	den	hund
D_1	2	1	1	1	1	0	0	0	0	0	0	0	0	0
D_2	2	0	0	0	0	1	1	1	0	0	0	0	0	0
D_3	0	0	0	0	0	0	0	0	1	1	1	1	1	1

- Each word is a dimension of meaning

Meaning of Documents

- Three documents:
 1. *The cat jumps over the dog*
 2. *The kitten hops ovqer the hound*
 3. *Die Katze springt über den Hund*
- Take the vocabulary $\mathcal{V} = \{\text{The, cat, ...}\}$ and represent them as vectors of dimension $|\mathcal{V}|$ in a Document-Term-Matrix $\mathbf{X}_{|\mathcal{D}| \times |\mathcal{V}|}$:

	the	cat	jump	over	dog	kitten	ovqer	hound	die	katze	springt	über	den	hund
D_1	2	1	1	1	1	0	0	0	0	0	0	0	0	0
D_2	2	0	0	0	0	1	1	1	0	0	0	0	0	0
D_3	0	0	0	0	0	0	0	0	1	1	1	1	1	1

- Each word is a dimension of meaning
- Similarity: Angle between these vectors ($\cos(D_i, D_j)$)

You shall know a word by the company it keeps (Firth, 1957)

- Define a context (e.g. a document)

Meaning of Words

You shall know a word by the company it keeps (Firth, 1957)

- Define a context (e.g. a document)
- How often do two words appear in the same context?

	the	cat	jump	over	dog	kitten	...
the	2	1	1	1	1	1	...
cat	1	1	1	1	1	0	...
jump	1	1	1	1	1	0	...
over	1	1	1	1	1	0	...
dog	1	1	1	1	1	0	...
kitten	1	0	0	0	0	1	...
...

Meaning of Words

You shall know a word by the company it keeps (Firth, 1957)

- Define a context (e.g. a document)
- How often do two words appear in the same context?

	the	cat	jump	over	dog	kitten	...
the	2	1	1	1	1	1	...
cat	1	1	1	1	1	0	...
jump	1	1	1	1	1	0	...
over	1	1	1	1	1	0	...
dog	1	1	1	1	1	0	...
kitten	1	0	0	0	0	1	...
...

- What we want: $\cos(\text{cat}, \text{jump}) > \cos(\text{cat}, \text{kitten})$

Problem

- Many redundant dimensions

	the	cat	jump	over	dog	kitten	ovqer	hound	die	katze	springt	über	den	hund
Doc. 1	2	1	1	1	1	0	0	0	0	0	0	0	0	0
Doc. 2	2	0	0	0	0	1	1	1	0	0	0	0	0	0
Doc. 3	0	0	0	0	0	0	0	0	1	1	1	1	1	1

Problem

- Many redundant dimensions
- Short documents have many 0's (open ended surveys, social media data)

	the	cat	jump	over	dog	kitten	ovqer	hound	die	katze	springt	über	den	hund
Doc. 1	2	1	1	1	1	0	0	0	0	0	0	0	0	0
Doc. 2	2	0	0	0	0	1	1	1	0	0	0	0	0	0
Doc. 3	0	0	0	0	0	0	0	0	1	1	1	1	1	1

Problem

- Many redundant dimensions
- Short documents have many 0's (open ended surveys, social media data)
- Things that should be similar are not

	the	cat	jump	over	dog	kitten	ovqer	hound	die	katze	springt	über	den	hund
Doc. 1	2	1	1	1	1	0	0	0	0	0	0	0	0	0
Doc. 2	2	0	0	0	0	1	1	1	0	0	0	0	0	0
Doc. 3	0	0	0	0	0	0	0	0	1	1	1	1	1	1

Problem

- Many redundant dimensions
- Short documents have many 0's (open ended surveys, social media data)
- Things that should be similar are not
- Computational cost: A usual corpus has 10^5 unique words sometimes more

	the	cat	jump	over	dog	kitten	ovqer	hound	die	katze	springt	über	den	hund
Doc. 1	2	1	1	1	1	0	0	0	0	0	0	0	0	0
Doc. 2	2	0	0	0	0	1	1	1	0	0	0	0	0	0
Doc. 3	0	0	0	0	0	0	0	0	1	1	1	1	1	1

Word Embeddings

- Documents / Words are 'embedded' in lower dimensional space

- Documents / Words are 'embedded' in lower dimensional space
- Our example:

- Documents / Words are 'embedded' in lower dimensional space
- Our example:
 - 'cat' \sim {cat, kitten, Katze}

- Documents / Words are ‘embedded’ in lower dimensional space
- Our example:
 - ‘cat’ \sim {cat, kitten, Katze}
 - ‘dog’ \sim {dog, hound, Hund}

- Documents / Words are 'embedded' in lower dimensional space
- Our example:
 - 'cat' \sim {cat, kitten, Katze}
 - 'dog' \sim {dog, hound, Hund}
 - 'article' \sim {the, The, die, den}

Basic Idea

- Documents / Words are ‘embedded’ in lower dimensional space
- Our example:
 - ‘cat’ \sim {cat, kitten, Katze}
 - ‘dog’ \sim {dog, hound, Hund}
 - ‘article’ \sim {the, The, die, den}
- New Document Matrix:

	‘cat’	‘dog’	‘jump’	‘article’
D_1	1	1	1	2
D_2	1	1	1	2
D_3	1	1	1	2

- Methods to find a good low dimensional representation of a matrix

- Methods to find a good low dimensional representation of a matrix
 - SVD, PCA, Factor Analysis

- Methods to find a good low dimensional representation of a matrix
 - SVD, PCA, Factor Analysis
 - IRT Models

- Methods to find a good low dimensional representation of a matrix
 - SVD, PCA, Factor Analysis
 - IRT Models
 - Topic Models

- Methods to find a good low dimensional representation of a matrix
 - SVD, PCA, Factor Analysis
 - IRT Models
 - Topic Models
 - **Neural Networks**

- Continuous Bag of Words (CBOW) Model (Mikolov et al., 2013; Le and Mikolov, 2014)

Word Embeddings with Neural Nets

- Continuous Bag of Words (CBOW) Model (Mikolov et al., 2013; Le and Mikolov, 2014)
- Shallow neural net

Word Embeddings with Neural Nets

- Continuous Bag of Words (CBOW) Model (Mikolov et al., 2013; Le and Mikolov, 2014)
- Shallow neural net
- Can use very large training corpora

Word Embeddings with Neural Nets

- Continuous Bag of Words (CBOW) Model (Mikolov et al., 2013; Le and Mikolov, 2014)
- Shallow neural net
- Can use very large training corpora
- Produces good embeddings:

$$\mathbf{v}_{queen} - \mathbf{v}_{woman} \approx \mathbf{v}_{king}$$

$$\mathbf{v}_{paris} - \mathbf{v}_{france} \approx \mathbf{v}_{rome}$$

Meaning and Prediction

context word context
The cat jumps over the dog . The kitten hops over the hound.

- Intuition, let's model:

$$P(\text{word}|\text{context})$$

For all words in \mathcal{V}

Meaning and Prediction

context word context
The cat jumps over the dog . The kitten hops over the hound.

- Intuition, let's model:

$$P(\text{word}|\text{context})$$

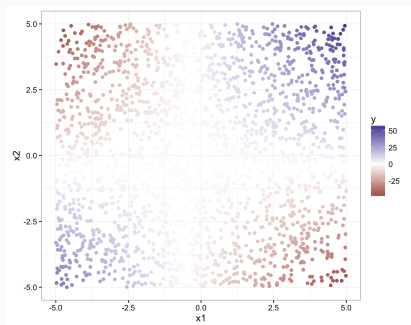
For all words in \mathcal{V}

- The parameters of this model will contain meaning

Neural Networks

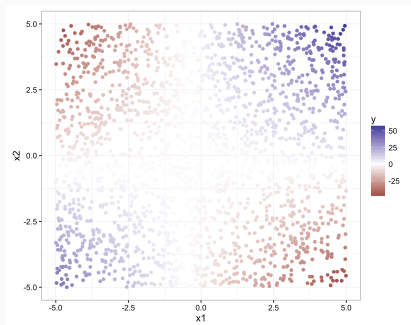
A Simple Neural Net

- Learn the function $y = f(\mathbf{x})$



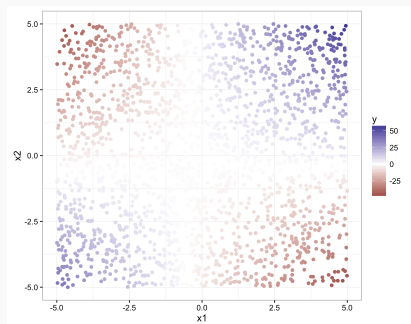
A Simple Neural Net

- Learn the function $y = f(\mathbf{x})$
- $y \in \mathbb{R}$, $\mathbf{x} \in \mathbb{R}^2$

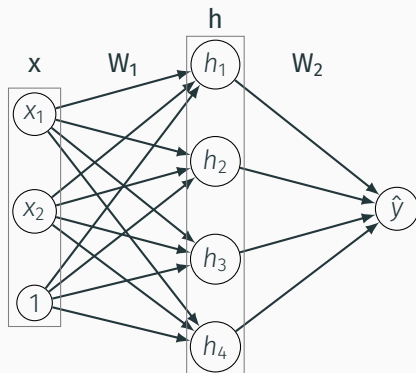


A Simple Neural Net

- Learn the function $y = f(\mathbf{x})$
- $y \in \mathbb{R}, \mathbf{x} \in \mathbb{R}^2$
- $y = x_1 + x_2 + 2x_1x_2 + \epsilon$

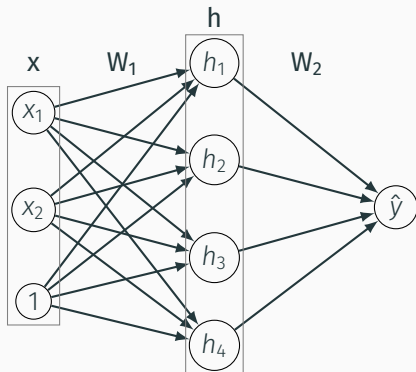


Architecture



- $p = 3$ inputs
- One hidden layer
- $k = 4$ hidden nodes
- Weight Matrices W_1, W_2
 $p \times k \quad 1 \times k$
- Activation function:
$$\sigma(x) = \frac{1}{1+e^{-x}}$$

Feed Forward Prediction



Hidden Layer

$$z = x \times W_1$$

$1 \times k \quad 1 \times p \quad p \times k$

$$h = \sigma(z)$$

$1 \times k \quad 1 \times k$

Output Layer

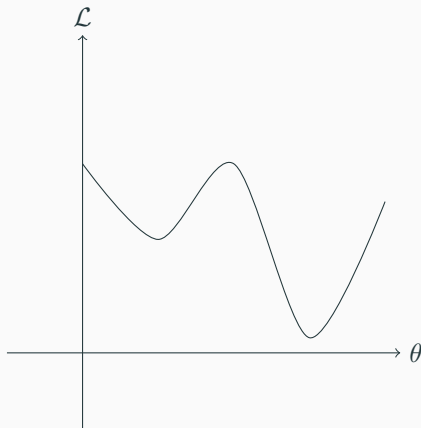
$$\hat{y} = h \times W_2$$

$1 \times 1 \quad 1 \times k \quad k \times 1$

Loss

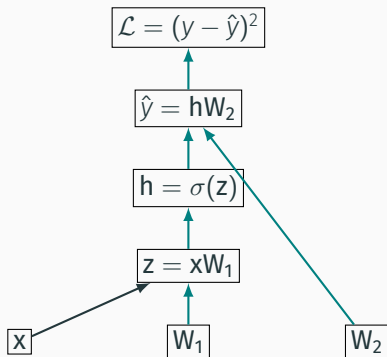
$$\mathcal{L} = (y - \hat{y})^2$$

Estimation - Error Back Propagation



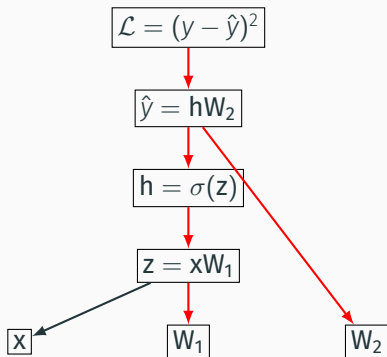
- Minimize the loss through the parameters $\theta = \{\mathbf{W}_1, \mathbf{W}_2\}$
- Need to find $\frac{\delta \mathcal{L}}{\delta \theta}$
- Then \mathcal{L} can be minimized with gradient descent

Estimation - Backpropagation I



- Make a prediction given inputs and parameters

Estimation - Backpropagation II



- Propagate the errors back through the network

- Feed Forward – > Back Propagate – > Adjust θ with $\frac{\delta \mathcal{L}}{\delta \theta}$ – > Repeat

Continuous Bag of Words Model

Overview

- Context size C
- Vocabulary $\mathcal{V} = \{\mathbf{x}_1, \mathbf{x}_2, \dots\}$
- $\mathbf{x}_j = [0, 0, 0, \dots, 1, 0, \dots, 0]$
 $1 \times |\mathcal{V}|$
- Focus word $\mathbf{y} \in \mathcal{V}$
 $1 \times |\mathcal{V}|$
- Context words $\mathbf{X} \subset \mathcal{V}$
 $C \times |\mathcal{V}|$
- Model $P(\mathbf{y}|\mathbf{X})$

“

I'm telling you, I used to use the word incompetent. Now I just call them stupid. I went to an Ivy League school. I'm very highly educated. I know words, I have the best words...but there is no better word than stupid. Right?

”

Architecture

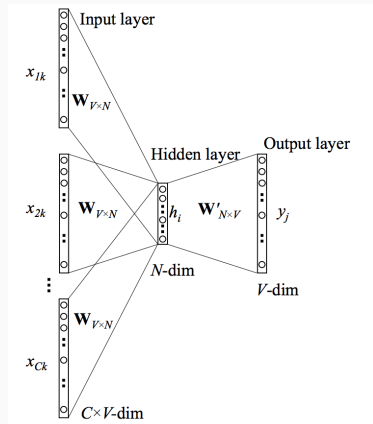


Figure 1: CBOW architecture. Figure from Rong (2014)

Hidden Layer (N units)

$$\mathbf{z}_{1 \times N} = (\mathbf{x}_1 + \mathbf{x}_1 + \dots + \mathbf{x}_C)_{1 \times |\mathcal{V}|} \times \mathbf{W}_1_{|\mathcal{V}| \times N}$$

$$\mathbf{h}_{1 \times N} = \frac{\mathbf{z}}{C}$$

Output Layer

$$\mathbf{u}_{1 \times |\mathcal{V}|} = \mathbf{h}_{1 \times N} \times \mathbf{W}_2_{N \times |\mathcal{V}|}$$

$$\hat{\mathbf{y}}_{1 \times |\mathcal{V}|} = \frac{e^{\mathbf{u}}}{\sum_{i=1}^{|\mathcal{V}|} e^{u_i}}$$

Loss

$$\mathcal{L} = -\log P(\mathbf{y}|\mathbf{X}) = -\log(\hat{\mathbf{y}}_{1 \times |\mathcal{V}|} \times \mathbf{y}_{|\mathcal{V}| \times 1}) \quad 24$$

Meaning of Words?

- W_1 and W_2 encode knowledge about $P(y|X)$
 $|\mathcal{V}| \times N$ $N \times |\mathcal{V}|$

Meaning of Words?

- \mathbf{W}_1 and \mathbf{W}_2 encode knowledge about $P(\mathbf{y}|\mathbf{X})$
 $|\mathcal{V}| \times N$ $N \times |\mathcal{V}|$
- Words are represented by vectors

Meaning of Words?

- \mathbf{W}_1 and \mathbf{W}_2 encode knowledge about $P(\mathbf{y}|\mathbf{X})$
 $|\mathcal{V}| \times N$ $N \times |\mathcal{V}|$
- Words are represented by vectors
- Dot product is a function of the angle (i.e. the similarity) of two vectors

Meaning of Words?

- \mathbf{W}_1 and \mathbf{W}_2 encode knowledge about $P(\mathbf{y}|\mathbf{X})$
 $|\mathcal{V}| \times N$ $N \times |\mathcal{V}|$
- Words are represented by vectors
- Dot product is a function of the angle (i.e. the similarity) of two vectors
- Training objective:

Meaning of Words?

- \mathbf{W}_1 and \mathbf{W}_2 encode knowledge about $P(\mathbf{y}|\mathbf{X})$
 $|\mathcal{V}| \times N$ $N \times |\mathcal{V}|$
- Words are represented by vectors
- Dot product is a function of the angle (i.e. the similarity) of two vectors
- Training objective:
 - Frequent context - word pairs: Large dot products ($\mathbf{u} = \mathbf{h}\mathbf{W}_2$)

Meaning of Words?

- \mathbf{W}_1 and \mathbf{W}_2 encode knowledge about $P(\mathbf{y}|\mathbf{X})$
 $|\mathcal{V}| \times N$ $N \times |\mathcal{V}|$
- Words are represented by vectors
- Dot product is a function of the angle (i.e. the similarity) of two vectors
- Training objective:
 - Frequent context - word pairs: Large dot products ($\mathbf{u} = \mathbf{h}\mathbf{W}_2$)
 - Infrequent context - word pairs: Small dot products ($\mathbf{u} = \mathbf{h}\mathbf{W}_2$)

- In actual training not all ϕ are calculated (Negative Sampling, Hierarchical Softmax)
- There is much debate about why these models produce useful embeddings (Arora et al. (2015) give intuition)
- It is not clear why a lower dimensional representation does better in NLP tasks
- → Much theoretical work to be done

Why is this valuable for political science research?

- Short documents
- Comparative Research - Multiple languages
- Use general knowledge about data in every task

Questions?

References

- Arora, S., Y. Li, Y. Liang, T. Ma, and A. Risteski (2015). Random walks on context spaces: Towards an explanation of the mysteries of semantic word embeddings. *CoRR abs/1502.03520*.
- Firth, J. R. (1957). *Papers in linguistics, 1934-1951*. Oxford University Press.
- Le, Q. V. and T. Mikolov (2014). Distributed representations of sentences and documents. *arXiv preprint arXiv:1405.4053*.
- Mikolov, T., I. Sutskever, K. Chen, G. S. Corrado, and J. Dean (2013). Distributed representations of words and phrases and their compositionality. In *Advances in neural information processing systems*, pp. 3111–3119.

Rong, X. (2014). word2vec parameter learning explained.
CoRR abs/1411.2738.

Appendix

Estimation - Error Back Propagation

$$\mathcal{L} = (y - \hat{y})^2$$

$$\hat{y} = \mathbf{h} \mathbf{W}_2$$

$$\mathbf{h} = \sigma(\mathbf{z})$$

$$\mathbf{z} = \mathbf{x} \mathbf{W}_1$$

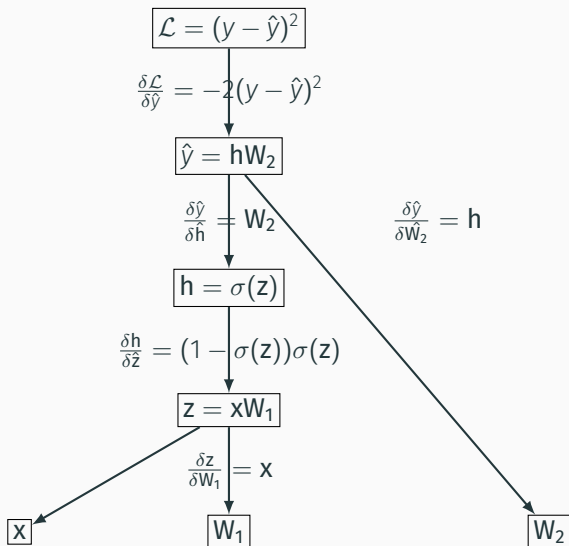
- How do the parameters affect \mathcal{L}
- Decompose \mathcal{L} into a computational graph

\mathbf{x}

\mathbf{W}_1

\mathbf{W}_2

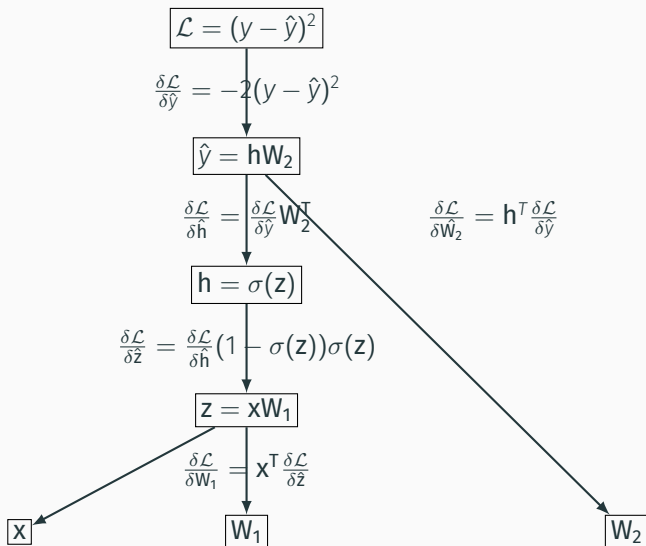
Simple Network - Backpropagation



- Calculate node derivatives

- How do the parameters affect \mathcal{L} ?

Simple Network - Backpropagation



- Calculate node derivatives
- Use the chain rule to propagate the errors through the network

CBOW Backprop Equations

Node Derivatives

$$\begin{aligned}\frac{\delta \mathcal{L}}{\delta u} &= \hat{y} - y = e \\ \frac{\delta u}{\delta W_2} &= h \\ \frac{\delta u}{\delta h} &= W_2 \\ \frac{\delta h}{\delta W_1} &= \left(\frac{x_1 + \dots + x_C}{C} \right)^T\end{aligned}$$

Chained Derivatives w.r.t. \mathcal{L}

$$\begin{aligned}\frac{\delta \mathcal{L}}{\delta u} &= \hat{y} - y = e \\ \frac{\delta \mathcal{L}}{\delta W_2} &= h^T e \\ \frac{\delta \mathcal{L}}{\delta h} &= \frac{\delta \mathcal{L}}{\delta u} W_2^T \\ \frac{\delta \mathcal{L}}{\delta W_1} &= \left(\frac{x_1 + \dots + x_C}{C} \right)^T \frac{\delta \mathcal{L}}{\delta h}\end{aligned}$$

Updating the Word Vectors

Intuition

- If the model is wrong for word i (i.e. $|e_i|$ is large) the corresponding column in \mathbf{W}_2 is adjusted
- If $\hat{y}_i > y_i$ (overestimated) \mathbf{w}_{2i} is pushed away from \mathbf{h} (the context)
- If $\hat{y}_i < y_i$ (underestimated) \mathbf{w}_{2i} is drawn towards \mathbf{h} (the context)

Updating Equations

$$\begin{aligned}\frac{\delta \mathcal{L}}{\delta \mathbf{u}} &= \hat{y} - y = \mathbf{e} \\ \frac{\delta \mathcal{L}}{\delta \mathbf{W}_2} &= \mathbf{h}^\top \mathbf{e}\end{aligned}$$

Example - Setup

- $V = \{\text{know, words, have}\}$
- $|V| = 3, N = 2, C = 1$
- Focus word: 'words'
- Context word: 'know'
- Initialize parameters randomly:

$$W_1 = \left[\begin{array}{c|cc} \text{know} & 0.3 & 0.02 \\ \text{words} & 0.01 & -0.03 \\ \text{have} & -0.4 & 0.001 \end{array} \right]$$

$$W_2 = \left[\begin{array}{ccc} \text{know} & \text{words} & \text{have} \\ \hline 0.02 & -0.002 & 0.1 \\ 0.04 & 0.008 & -0.03 \end{array} \right]$$

Example - Feed Forward

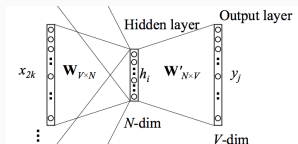


Figure 2: CBOW architecture.
Figure from Rong (2014)

- $\mathbf{h} = \mathbf{x}_1 \mathbf{W}_1 = [0.3, 0.02]$
- $\mathbf{u} = \mathbf{h} \mathbf{W}_2 = [0.0068, -0.00044, 0.0294]$
- $\hat{\mathbf{y}} = [0.331, 0.329, 0.339]$
- $\mathbf{e} = \hat{\mathbf{y}} - \mathbf{y} = [0.331, 0.329, 0.339] - [0, 1, 0] = [0.331, -0.671, 0.339]$

Example - Update

$$\begin{aligned}W_2^{\text{new}} &= W_2 + \eta \mathbf{h}^T \mathbf{e} \\&= W_2 + \eta \begin{bmatrix} \textit{know} & \textit{words} & \textit{have} \\ 0.010 & -0.201 & 0.101 \\ 0.006 & -0.013 & 0.006 \end{bmatrix}\end{aligned}$$