

# Distributed representation of text

He He



**NEW YORK UNIVERSITY**

January 29, 2025

# Table of Contents

Review

Introduction

Count-based word embeddings

Prediction-based word embeddings

Neural networks

## Last week

### **Generative vs discriminative models** for text classification

- (Multinomial) naive Bayes

What's the key assumption?

## Last week

### **Generative vs discriminative models** for text classification

- (Multinomial) naive Bayes What's the key assumption?
  - Assumes conditional independence
  - Very efficient in practice (closed-form solution)

## Last week

### **Generative vs discriminative models** for text classification

- (Multinomial) naive Bayes What's the key assumption?
  - Assumes conditional independence
  - Very efficient in practice (closed-form solution)
- Logistic regression What's the main advantage?

## Last week

### **Generative vs discriminative models** for text classification

- (Multinomial) naive Bayes What's the key assumption?
  - Assumes conditional independence
  - Very efficient in practice (closed-form solution)
- Logistic regression What's the main advantage?
  - Works with all kinds of features
  - Wins with more data [Ng and Jordan, 2001]

## Last week

### **Generative vs discriminative models** for text classification

- (Multinomial) naive Bayes What's the key assumption?
  - Assumes conditional independence
  - Very efficient in practice (closed-form solution)
- Logistic regression What's the main advantage?
  - Works with all kinds of features
  - Wins with more data [Ng and Jordan, 2001]

### **Feature vector** of text input

- BoW representation
- N-gram features (usually  $n \leq 3$ )

# Table of Contents

Review

Introduction

Count-based word embeddings

Prediction-based word embeddings

Neural networks



# Objective

**Goal:** come up with a **good representation** of text

- What is a representation?

# Objective

**Goal:** come up with a **good representation** of text

- What is a representation?
  - Feature map:  $\phi: \text{text} \rightarrow \mathbb{R}^d$ , e.g., BoW, handcrafted features

# Objective

**Goal:** come up with a **good representation** of text

- What is a representation?
  - Feature map:  $\phi: \text{text} \rightarrow \mathbb{R}^d$ , e.g., BoW, handcrafted features
  - “Representation” often refers to **learned** features of the input

# Objective

**Goal:** come up with a **good representation** of text

- What is a representation?
  - Feature map:  $\phi: \text{text} \rightarrow \mathbb{R}^d$ , e.g., BoW, handcrafted features
  - “Representation” often refers to **learned** features of the input
- What is a good representation?

# Objective

**Goal:** come up with a **good representation** of text

- What is a representation?
  - Feature map:  $\phi: \text{text} \rightarrow \mathbb{R}^d$ , e.g., BoW, handcrafted features
  - “Representation” often refers to **learned** features of the input
- What is a good representation?
  - Enables **a notion of distance** over text:  $d(\phi(a), \phi(b))$  is small for semantically similar texts  $a$  and  $b$ 
    - Abstraction over individual words

# Objective

**Goal:** come up with a **good representation** of text

- What is a representation?
  - Feature map:  $\phi: \text{text} \rightarrow \mathbb{R}^d$ , e.g., BoW, handcrafted features
  - “Representation” often refers to **learned** features of the input
- What is a good representation?
  - Enables **a notion of distance** over text:  $d(\phi(a), \phi(b))$  is small for semantically similar texts  $a$  and  $b$ 
    - Abstraction over individual words
  - Leads to **good task performance** (with less training data)
    - Contains useful features

# Objective

**Goal:** come up with a **good representation** of text

- What is a representation?
  - Feature map:  $\phi: \text{text} \rightarrow \mathbb{R}^d$ , e.g., BoW, handcrafted features
  - “Representation” often refers to **learned** features of the input
- What is a good representation?
  - Enables **a notion of distance** over text:  $d(\phi(a), \phi(b))$  is small for semantically similar texts  $a$  and  $b$ 
    - Abstraction over individual words
  - Leads to **good task performance** (with less training data)
    - Contains useful features

# Objective

**Goal:** come up with a **good representation** of text

- What is a representation?
  - Feature map:  $\phi: \text{text} \rightarrow \mathbb{R}^d$ , e.g., BoW, handcrafted features
  - “Representation” often refers to **learned** features of the input
- What is a good representation?
  - Enables **a notion of distance** over text:  $d(\phi(a), \phi(b))$  is small for semantically similar texts  $a$  and  $b$ 
    - Abstraction over individual words
  - Leads to **good task performance** (with less training data)
    - Contains useful features



What text representation have we learned?



# Table of Contents

Review

Introduction

Count-based word embeddings

Prediction-based word embeddings

Neural networks

# The distributional hypothesis

How do we come up with a representation without handcrafting features?

## The distributional hypothesis

How do we come up with a representation without handcrafting features?

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

# The distributional hypothesis

How do we come up with a representation without handcrafting features?

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

Word guessing! (example from Eisenstein's book)

Everybody likes [tezgüino](#).

# The distributional hypothesis

How do we come up with a representation without handcrafting features?

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

Word guessing! (example from Eisenstein's book)

Everybody likes **tezgüino**.

We make **tezgüino** out of corn.

# The distributional hypothesis

How do we come up with a representation without handcrafting features?

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

Word guessing! (example from Eisenstein's book)

Everybody likes **tezgüino**.

We make **tezgüino** out of corn.

A bottle of **tezgüino** is on the table.

# The distributional hypothesis

How do we come up with a representation without handcrafting features?

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

Word guessing! (example from Eisenstein's book)

Everybody likes **tezgüino**.

We make **tezgüino** out of corn.

A bottle of **tezgüino** is on the table.

Don't have **tezgüino** before you drive.

# The distributional hypothesis

How do we come up with a representation without handcrafting features?

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

Word guessing! (example from Eisenstein's book)

Everybody likes **tezgüino**.

We make **tezgüino** out of corn.

A bottle of **tezgüino** is on the table.

Don't have **tezgüino** before you drive.

**Idea:** Represent a word by its neighbors.



## Step 1: Choose the context

What are the neighbors?

Example:

- word  $\times$  document

	As You Like It	Twelfth Night	Julius Caesar	Henry V
battle	1	0	7	13
good	114	80	62	89
fool	36	58	1	4
wit	20	15	2	3

**Figure 6.2** The term-document matrix for four words in four Shakespeare plays. Each cell contains the number of times the (row) word occurs in the (column) document.

Figure: Jurafsky and Martin.

## Step 1: Choose the context

What are the neighbors?

Example:

- word  $\times$  document
- word  $\times$  word

	As You Like It	Twelfth Night	Julius Caesar	Henry V
battle	1	0	7	13
good	114	80	62	89
fool	36	58	1	4
wit	20	15	2	3

**Figure 6.2** The term-document matrix for four words in four Shakespeare plays. Each cell contains the number of times the (row) word occurs in the (column) document.

Figure: Jurafsky and Martin.

## Step 1: Choose the context

What are the neighbors?

Example:

- word  $\times$  document
- word  $\times$  word
- person  $\times$  movie

	As You Like It	Twelfth Night	Julius Caesar	Henry V
battle	1	0	7	13
good	114	80	62	89
fool	36	58	1	4
wit	20	15	2	3

**Figure 6.2** The term-document matrix for four words in four Shakespeare plays. Each cell contains the number of times the (row) word occurs in the (column) document.

Figure: Jurafsky and Martin.

## Step 1: Choose the context

What are the neighbors?

Example:

- word  $\times$  document
- word  $\times$  word
- person  $\times$  movie
- note  $\times$  song

	As You Like It	Twelfth Night	Julius Caesar	Henry V
battle	1	0	7	13
good	114	80	62	89
fool	36	58	1	4
wit	20	15	2	3

**Figure 6.2** The term-document matrix for four words in four Shakespeare plays. Each cell contains the number of times the (row) word occurs in the (column) document.

Figure: Jurafsky and Martin.

## Step 1: Choose the context

What are the neighbors?

Example:

- word  $\times$  document
- word  $\times$  word
- person  $\times$  movie
- note  $\times$  song

	As You Like It	Twelfth Night	Julius Caesar	Henry V
battle	1	0	7	13
good	114	80	62	89
fool	36	58	1	4
wit	20	15	2	3

**Figure 6.2** The term-document matrix for four words in four Shakespeare plays. Each cell contains the number of times the (row) word occurs in the (column) document.

Figure: Jurafsky and Martin.

Construct a matrix where

- Row and columns represent two sets of objects
- Each entry is the co-occurrence counts of the two objects

## Distance functions

Now we already have a representation for each document / word!

	As You Like It	Twelfth Night	Julius Caesar	Henry V
<b>battle</b>	1	0	7	13
<b>good</b>	114	80	62	89
<b>fool</b>	36	58	1	4
<b>wit</b>	20	15	2	3

**Figure 6.2** The term-document matrix for four words in four Shakespeare plays. Each cell contains the number of times the (row) word occurs in the (column) document.

How do we compute similarities between documents?

## Distance functions

Now we already have a representation for each document / word!

	As You Like It	Twelfth Night	Julius Caesar	Henry V
battle	1	0	7	13
good	114	80	62	89
fool	36	58	1	4
wit	20	15	2	3

**Figure 6.2** The term-document matrix for four words in four Shakespeare plays. Each cell contains the number of times the (row) word occurs in the (column) document.

How do we compute similarities between documents?

### Euclidean distance

For  $a, b \in \mathbb{R}^d$ ,

$$d(a, b) = \sqrt{\sum_{i=1}^d (a_i - b_i)^2}.$$

$$b_i = 2a_i$$

Assume  $a$  and  $b$  are BoW vectors. What if  $b$  repeats each sentence in  $a$  twice?

## Distance functions

Now we already have a representation for each document / word!

	As You Like It	Twelfth Night	Julius Caesar	Henry V
battle	1	0	7	13
good	114	80	62	89
fool	36	58	1	4
wit	20	15	2	3

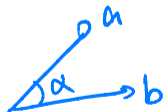
**Figure 6.2** The term-document matrix for four words in four Shakespeare plays. Each cell contains the number of times the (row) word occurs in the (column) document.

How do we compute similarities between documents?

### Cosine similarity

For  $a, b \in \mathbb{R}^d$ ,

$$\text{sim}(a, b) = \frac{a \cdot b}{\|a\| \|b\|}$$



Defines angle between two vectors (=  $\cos \alpha$  in 2D)



## Example application: information retrieval

Given a set of documents and a query, use the BoW representation and cosine similarity to find the most relevant document.

## Example application: information retrieval

Given a set of documents and a query, use the ~~BoW~~ representation and cosine similarity to find the most relevant document.

What are potential problems?

Q: What are **good** books on great **battles**?

	<b>As You Like It</b>	<b>Twelfth Night</b>	<b>Julius Caesar</b>	<b>Henry V</b>
<b>battle</b>	1	0	7	13
<b>good</b>	114	80	62	89
<b>fool</b>	36	58	1	4
<b>wit</b>	20	15	2	3

**Figure 6.2** The term-document matrix for four words in four Shakespeare plays. Each cell contains the number of times the (row) word occurs in the (column) document.

## Example application: information retrieval

Given a set of documents and a query, use the BoW representation and cosine similarity to find the most relevant document.

What are potential problems?

Q: What are good books on great battles?

	<b>As You Like It</b>	<b>Twelfth Night</b>	<b>Julius Caesar</b>	<b>Henry V</b>
<b>battle</b>	1	0	7	13
<b>good</b>	114	80	62	89
<b>fool</b>	36	58	1	4
<b>wit</b>	20	15	2	3

**Figure 6.2** The term-document matrix for four words in four Shakespeare plays. Each cell contains the number of times the (row) word occurs in the (column) document.

Similarity may be dominated by **common words**!

## Step 2: Reweight counts

**Key idea:** upweight words that carry more information about the document

Construct a feature map  $\phi$ : document  $\rightarrow \mathbb{R}^{|\mathcal{V}|}$

TFIDF weight for token  $w$ :

## Step 2: Reweight counts

**Key idea:** upweight words that carry more information about the document

Construct a feature map  $\phi: \text{document} \rightarrow \mathbb{R}^{|\mathcal{V}|}$

TFIDF weight for token  $w$ :

$$\phi_w(d) = \underbrace{\text{count}(w, d)}_{\text{tf}(w, d)} \times$$

- **Term frequency (TF):** count of the word in the document (same as BoW)

## Step 2: Reweight counts

**Key idea:** upweight words that carry more information about the document

Construct a feature map  $\phi: \text{document} \rightarrow \mathbb{R}^{|\mathcal{V}|}$

TFIDF weight for token  $w$ :

$$\phi_w(d) = \underbrace{\text{count}(w, d)}_{\text{tf}(w, d)} \times \log \underbrace{\frac{\# \text{ documents}}{\# \text{ documents containing } w}}_{\text{idf}(w)} .$$

- **Term frequency (TF):** count of the word in the document (same as BoW)
- Reweight by **inverse document frequency (IDF):** how **specific** is the word to any particular document

## Step 2: Reweight counts

**Key idea:** upweight words that carry more information about the document

Construct a feature map  $\phi: \text{document} \rightarrow \mathbb{R}^{|\mathcal{V}|}$

TFIDF weight for token  $w$ :

$$\phi_w(d) = \underbrace{\text{count}(w, d)}_{\text{tf}(w, d)} \times \log \underbrace{\frac{\# \text{ documents}}{\# \text{ documents containing } w}}_{\text{idf}(w)} .$$

- **Term frequency (TF):** count of the word in the document (same as BoW)
- Reweight by **inverse document frequency (IDF):** how **specific** is the word to any particular document
- Higher weight on **frequent** words that only **occur in a few documents**

## TFIDF example

	As You Like It	Twelfth Night	Julius Caesar	Henry V
battle	0.074	0	0.22	0.28
good	0	0	0	0
fool	0.019	0.021	0.0036	0.0083
wit	0.049	0.044	0.018	0.022

**Figure 6.9** A tf-idf weighted term-document matrix for four words in four Shakespeare plays, using the counts in Fig. 6.2. For example the 0.049 value for *wit* in *As You Like It* is the product of  $tf = \log_{10}(20 + 1) = 1.322$  and  $idf = .037$ . Note that the idf weighting has eliminated the importance of the ubiquitous word *good* and vastly reduced the impact of the almost-ubiquitous word *fool*.

Figure: From Jurafsky and Martin.

Why do some words have zero weights?



## An alternative way to reweighting using pointwise mutual information

$$\text{PMI}(x; y) \stackrel{\text{def}}{=} \log \frac{p(x, y)}{p(x)p(y)} = \log \frac{p(x | y)}{p(x)} = \log \frac{p(y | x)}{p(y)}$$

## An alternative way to reweighting using pointwise mutual information

$$\text{PMI}(x; y) \stackrel{\text{def}}{=} \log \frac{p(x, y)}{p(x)p(y)} = \log \frac{p(x | y)}{p(x)} = \log \frac{p(y | x)}{p(y)}$$

- Symmetric:  $\text{PMI}(x; y) = \text{PMI}(y; x)$

## An alternative way to reweighting using pointwise mutual information

$$\text{PMI}(x; y) \stackrel{\text{def}}{=} \log \frac{p(x, y)}{p(x)p(y)} = \log \frac{p(x | y)}{p(x)} = \log \frac{p(y | x)}{p(y)}$$

- Symmetric:  $\text{PMI}(x; y) = \text{PMI}(y; x)$
- Estimates:

$$\hat{p}(x | y) = \frac{\text{count}(x, y)}{\sum_{x' \in \mathcal{X}} \text{count}(x', y)} \quad \text{how often does } x \text{ occur in the neighborhood of } y$$

$$\hat{p}(x) = \frac{\text{count}(x)}{\sum_{x' \in \mathcal{X}} \text{count}(x')} \quad \text{how often does } x \text{ occur in the corpus}$$

## An alternative way to reweighting using pointwise mutual information

$$\text{PMI}(x; y) \stackrel{\text{def}}{=} \log \frac{p(x, y)}{p(x)p(y)} = \log \frac{p(x | y)}{p(x)} = \log \frac{p(y | x)}{p(y)}$$

- Symmetric:  $\text{PMI}(x; y) = \text{PMI}(y; x)$
- Estimates:

$$\hat{p}(x | y) = \frac{\text{count}(x, y)}{\sum_{x' \in \mathcal{X}} \text{count}(x', y)} \quad \text{how often does } x \text{ occur in the neighborhood of } y$$

$$\hat{p}(x) = \frac{\text{count}(x)}{\sum_{x' \in \mathcal{X}} \text{count}(x')} \quad \text{how often does } x \text{ occur in the corpus}$$

- **Positive PMI:**  $\text{PPMI}(x; y) \stackrel{\text{def}}{=} \max(0, \text{PMI}(x; y))$

## An alternative way to reweighting using pointwise mutual information

$$\text{PMI}(x; y) \stackrel{\text{def}}{=} \log \frac{p(x, y)}{p(x)p(y)} = \log \frac{p(x | y)}{p(x)} = \log \frac{p(y | x)}{p(y)}$$

- Symmetric:  $\text{PMI}(x; y) = \text{PMI}(y; x)$
- Estimates:

$$\hat{p}(x | y) = \frac{\text{count}(x, y)}{\sum_{x' \in \mathcal{X}} \text{count}(x', y)} \quad \text{how often does } x \text{ occur in the neighborhood of } y$$

$$\hat{p}(x) = \frac{\text{count}(x)}{\sum_{x' \in \mathcal{X}} \text{count}(x')} \quad \text{how often does } x \text{ occur in the corpus}$$

- **Positive PMI:**  $\text{PPMI}(x; y) \stackrel{\text{def}}{=} \max(0, \text{PMI}(x; y))$
- Application in NLP: measure association between words

## Step 3: Dimensionality reduction

What's the size of this matrix?

	<b>As You Like It</b>	<b>Twelfth Night</b>	<b>Julius Caesar</b>	<b>Henry V</b>
<b>battle</b>	0.074	0	0.22	0.28
<b>good</b>	0	0	0	0
<b>fool</b>	0.019	0.021	0.0036	0.0083
<b>wit</b>	0.049	0.044	0.018	0.022

**Figure 6.9** A tf-idf weighted term-document matrix for four words in four Shakespeare

Figure: Jurafsky and Martin.

## Step 3: Dimensionality reduction

What's the size of this matrix?

	As You Like It	Twelfth Night	Julius Caesar	Henry V
battle	0.074	0	0.22	0.28
good	0	0	0	0
fool	0.019	0.021	0.0036	0.0083
wit	0.049	0.044	0.018	0.022

**Figure 6.9** A tf-idf weighted term-document matrix for four words in four Shakespeare

Figure: Jurafsky and Martin.

**Motivation:** want a lower-dimensional, dense representation for efficiency

### Step 3: Dimensionality reduction

Recall **SVD**: a  $m \times n$  matrix  $A_{m \times n}$  (e.g., a word-document matrix), can be decomposed to

$$U_{m \times m} \Sigma_{m \times n} V_{n \times n}^T,$$

where  $U$  and  $V$  are orthogonal matrices, and  $\Sigma$  is a diagonal matrix.



## Step 3: Dimensionality reduction

Recall **SVD**: a  $m \times n$  matrix  $A_{m \times n}$  (e.g., a word-document matrix), can be decomposed to

$$U_{m \times m} \Sigma_{m \times n} V_{n \times n}^T,$$

where  $U$  and  $V$  are orthogonal matrices, and  $\Sigma$  is a diagonal matrix.

**Interpretation:**

$$AA^T = (U\Sigma V^T)(V\Sigma U^T) = U\Sigma^2 U^T.$$

- $\sigma_i^2$  are eigenvalues of  $AA^T$
- Connection to PCA: If columns of  $A$  have zero mean (i.e.  $AA^T$  is the covariance matrix), then columns of  $U$  are principle components of the column space of  $A$ .

## SVD for the word-document matrix

$$A = \begin{matrix} & d_1 & d_2 & \cdots & d_n & \\ \begin{pmatrix} 12 & 5 & \cdots & 8 \\ 7 & 10 & \cdots & 3 \\ 4 & 8 & \cdots & 6 \\ 9 & 3 & \cdots & 7 \end{pmatrix} & \text{US} \\ & \text{gov} \\ & \text{gene} \\ & \text{lab} \end{matrix}$$

# SVD for the word-document matrix

$$A = \begin{pmatrix} d_1 & d_2 & \dots & d_n \\ 12 & 5 & \dots & 8 \\ 7 & 10 & \dots & 3 \\ 4 & 8 & \dots & 6 \\ 9 & 3 & \dots & 7 \end{pmatrix} \begin{matrix} \text{US} \\ \text{gov} \\ \text{gene} \\ \text{lab} \end{matrix}$$

$$= \underbrace{\begin{pmatrix} 0.50 & 0.02 & \dots \\ 0.60 & 0.03 & \dots \\ 0.01 & 0.72 & \dots \\ 0.02 & 0.84 & \dots \end{pmatrix}}_U \underbrace{\begin{pmatrix} 15 & 0 & 0 & 0 & \dots & 0 \\ 0 & 10 & 0 & 0 & \dots & 0 \\ 0 & 0 & 5 & 0 & \dots & 0 \\ 0 & 0 & 0 & 2 & \dots & 0 \end{pmatrix}}_\Sigma \underbrace{\begin{pmatrix} 0.50 & 0.60 & \dots \\ 0.64 & 0.48 & \dots \\ 0.12 & 0.24 & \dots \\ 0.36 & 0.12 & \dots \end{pmatrix}}_V^T$$

$4 \times 4$        $4 \times n$        $n \times n$

## SVD for the word-document matrix

$$A = \begin{pmatrix} d_1 & d_2 & \cdots & d_n \\ 12 & 5 & \cdots & 8 \\ 7 & 10 & \cdots & 3 \\ 4 & 8 & \cdots & 6 \\ 9 & 3 & \cdots & 7 \end{pmatrix} \begin{matrix} \text{US} \\ \text{gov} \\ \text{gene} \\ \text{lab} \end{matrix}$$

$$= \begin{pmatrix} 0.50 & 0.02 & \cdots \\ 0.60 & 0.03 & \cdots \\ 0.01 & 0.72 & \cdots \\ 0.02 & 0.84 & \cdots \end{pmatrix}_{4 \times 4} \begin{pmatrix} 15 & 0 & 0 & 0 & \cdots & 0 \\ 0 & 10 & 0 & 0 & \cdots & 0 \\ 0 & 0 & 5 & 0 & \cdots & 0 \\ 0 & 0 & 0 & 2 & \cdots & 0 \end{pmatrix}_{4 \times n} \begin{pmatrix} 0.50 & 0.60 & \cdots \\ 0.64 & 0.48 & \cdots \\ 0.12 & 0.24 & \cdots \\ 0.36 & 0.12 & \cdots \end{pmatrix}_{n \times n}^T$$

- $u_i$  are document clusters and  $v_i$  are word clusters
- Take top- $k$  components to obtain word vectors:  $W = U_k \Sigma_k$  (or just  $U_k$ )

# SVD for the word-document matrix

Computing the dense word vectors:

- Run **truncated SVD** of the word-document matrix  $A_{m \times n}$
- Each row of  $U_{m \times k} \Sigma_k$  corresponds to a word vector of dimension  $k$
- Each coordinate of the word vector corresponds to a cluster of documents (i.e., topics such as politics, music, etc.)

# Summary

## Count-based word embeddings

1. Design the matrix, e.g. word  $\times$  document, people  $\times$  movie.
2. Reweight the raw counts, e.g. TFIDF, PPMI.
3. Reduce dimensionality by truncated SVD.
4. Use word/person/etc. vectors in downstream tasks.

## Key idea:

- Intuition: Represent an object by its connection to other objects.
- Lexical semantics: the word meaning can be represented by the context it occurs in.
- Linear algebra: Infer clusters (e.g., concepts, topics) using co-occurrence statistics

# Table of Contents

Review

Introduction

Count-based word embeddings

Prediction-based word embeddings

Neural networks

# Learning word embeddings

**Goal:** map each word to a vector in  $\mathbb{R}^d$  such that *similar* words also have *similar* word vectors.



# Learning word embeddings

**Goal:** map each word to a vector in  $\mathbb{R}^d$  such that *similar* words also have *similar* word vectors.

Can we directly optimize the goal by formalizing a **prediction problem**?

- Needs to be self-supervised since our data is unlabeled.

# Learning word embeddings

**Goal:** map each word to a vector in  $\mathbb{R}^d$  such that *similar* words also have *similar* word vectors.

Can we directly optimize the goal by formalizing a **prediction problem**?

- Needs to be self-supervised since our data is unlabeled.

**Distributional hypothesis:** Similar words occur in similar contexts

- Predict the context given a word  $f : \text{word} \rightarrow \text{context}$
- Words that tend to occur in same contexts will have similar representation

# The skip-gram model

**Task:** given a **word**, predict its **neighboring words** within a window

The **quick brown fox jumps over** the lazy dog

# The skip-gram model

**Task:** given a word, predict its neighboring words within a window

The quick brown fox jumps over the lazy dog

Assume **conditional independence** of the context words:

$$p(w_{i-k}, \dots, w_{i-1}, w_{i+1}, \dots, w_{i+k} \mid w_i) = \prod_{j=i-k, j \neq i}^{i+k} p(w_j \mid w_i)$$

# The skip-gram model

**Task:** given a **word**, predict its **neighboring words** within a window

The **quick brown fox jumps over** the lazy dog

Assume **conditional independence** of the context words:

$$p(w_{i-k}, \dots, w_{i-1}, w_{i+1}, \dots, w_{i+k} \mid w_i) = \prod_{j=i-k, j \neq i}^{i+k} p(w_j \mid w_i)$$

How to model  $p(w_j \mid w_i)$ ?

# The skip-gram model

**Task:** given a **word**, predict its **neighboring words** within a window

The **quick brown fox jumps over** the lazy dog

Assume **conditional independence** of the context words:

$$p(w_{i-k}, \dots, w_{i-1}, w_{i+1}, \dots, w_{i+k} \mid w_i) = \prod_{j=i-k, j \neq i}^{i+k} p(w_j \mid w_i)$$

How to model  $p(w_j \mid w_i)$ ?

Multiclass classification

## The skip-gram model

Use the softmax function to predict **context words** from the **center word**

$$p(w_j | w_i) = \frac{\exp[\phi_{\text{ctx}}(w_j) \cdot \phi_{\text{wrd}}(w_i)]}{\sum_{w \in \mathcal{V}} \exp[\phi_{\text{ctx}}(w) \cdot \phi_{\text{wrd}}(w_i)]}$$

## The skip-gram model

Use the softmax function to predict **context words** from the **center word**

$$p(w_j | w_i) = \frac{\exp[\phi_{\text{ctx}}(w_j) \cdot \phi_{\text{wrd}}(w_i)]}{\sum_{w \in \mathcal{V}} \exp[\phi_{\text{ctx}}(w_j) \cdot \phi_{\text{wrd}}(w_i)]} \propto \exp(\theta_j \cdot \phi(x))$$



What's the difference from multinomial logistic regression?



# The skip-gram model

Use the softmax function to predict **context words** from the **center word**

$$p(w_j | w_i) = \frac{\exp[\phi_{\text{ctx}}(w_j) \cdot \phi_{\text{wrđ}}(w_i)]}{\sum_{w \in \mathcal{V}} \exp[\phi_{\text{ctx}}(w_j) \cdot \phi_{\text{wrđ}}(w_i)]}$$



What's the difference from multinomial logistic regression?

Implementation:

- $\phi: \mathcal{V} \rightarrow \mathbb{R}^d$ .  $\phi$  can be implemented as a dictionary from words (ids) to vectors.

# The skip-gram model

Use the softmax function to predict **context words** from the **center word**

$$p(w_j | w_i) = \frac{\exp[\phi_{\text{ctx}}(w_j) \cdot \phi_{\text{wrđ}}(w_i)]}{\sum_{w \in \mathcal{V}} \exp[\phi_{\text{ctx}}(w_j) \cdot \phi_{\text{wrđ}}(w_i)]}$$



What's the difference from multinomial logistic regression?

Implementation:

- $\phi: \mathcal{V} \rightarrow \mathbb{R}^d$ .  $\phi$  can be implemented as a dictionary from words (ids) to vectors.
- For each word  $w$ , learn two vectors.

# The skip-gram model

Use the softmax function to predict **context words** from the **center word**

$$p(w_j | w_i) = \frac{\exp[\phi_{\text{ctx}}(w_j) \cdot \phi_{\text{wrđ}}(w_i)]}{\sum_{w \in \mathcal{V}} \exp[\phi_{\text{ctx}}(w_j) \cdot \phi_{\text{wrđ}}(w_i)]}$$



What's the difference from multinomial logistic regression?

Implementation:

- $\phi: \mathcal{V} \rightarrow \mathbb{R}^d$ .  $\phi$  can be implemented as a dictionary from words (ids) to vectors.
- For each word  $w$ , learn two vectors.
- Learn parameters by MLE and SGD (Is the objective convex?)

# The skip-gram model

Use the softmax function to predict **context words** from the **center word**

$$p(w_j | w_i) = \frac{\exp[\phi_{\text{ctx}}(w_j) \cdot \phi_{\text{wrd}}(w_i)]}{\sum_{w \in \mathcal{V}} \exp[\phi_{\text{ctx}}(w) \cdot \phi_{\text{wrd}}(w_i)]}$$

$-\log p(w_j | w_i) = -\phi_{\text{ctx}}(w_j) \cdot \phi_{\text{wrd}}(w_i) + \log \sum_{w \in \mathcal{V}} \exp(\phi_{\text{ctx}}(w) \cdot \phi_{\text{wrd}}(w_i))$



What's the difference from multinomial logistic regression?

Implementation:

- $\phi: \mathcal{V} \rightarrow \mathbb{R}^d$ .  $\phi$  can be implemented as a dictionary from words (ids) to vectors.
- For each word  $w$ , learn two vectors.
- Learn parameters by MLE and SGD (Is the objective convex?)
- $\phi_{\text{wrd}}$  is taken as the word embedding

## Negative sampling

Challenge in MLE: computing the normalizer is expensive (try calculate the gradient)!

## Negative sampling

Challenge in MLE: computing the normalizer is expensive (try calculate the gradient)!

Key idea: solve a binary classification problem instead

Is the (word, context) pair real or fake?

### positive examples +

$w$	$c_{\text{pos}}$
apricot	tablespoon
apricot	of
apricot	jam
apricot	a

### negative examples -

$w$	$c_{\text{neg}}$	$w$	$c_{\text{neg}}$
apricot	aardvark	apricot	seven
apricot	my	apricot	forever
apricot	where	apricot	dear
apricot	coaxial	apricot	if

## Negative sampling

Challenge in MLE: computing the normalizer is expensive (try calculate the gradient)!

Key idea: solve a binary classification problem instead

Is the (word, context) pair real or fake?

**positive examples +**

$w$	$c_{\text{pos}}$
apricot	tablespoon
apricot	of
apricot	jam
apricot	a

**negative examples -**

$w$	$c_{\text{neg}}$	$w$	$c_{\text{neg}}$
apricot	aardvark	apricot	seven
apricot	my	apricot	forever
apricot	where	apricot	dear
apricot	coaxial	apricot	if

$$- \log p_{\theta}(\text{real} \mid w, c) = \frac{1}{1 + e^{-\phi_{\text{ctx}}(c) \cdot \phi_{\text{word}}(w)}}$$

Large dot product between  $w$  and  $c$  if they co-occur.

## Negative sampling: loss function

Let  $s(w, c) = \phi_{\text{ctx}}(c) \cdot \phi_{\text{word}}(w)$  be the score of the context and target word.

NLL loss:

$$\log(1 + e^{-s(w, c)}) + \sum_{c_n \in \mathcal{N}_{w, c}} \log(1 + e^{s(w, c_n)})$$

Let  $\ell(x) = \log(1 + e^{-x})$ , we have the loss for negative sampling

$$\ell(s(w, c)) + \sum_{c_n \in \mathcal{N}_{w, c}} \ell(-s(w, c_n))$$



# Fasttext

Improvements over skipgram:

- Negative sampling
- Character n-gram, e.g., <ap, app, ppl, ple, le>

$$s(w, c) = \sum_{v \in \text{ngram}(w)} s(v, c)$$

- Faster to train and maintains similar quality to skipgram
- Well-maintained open source project: <https://fasttext.cc>

# The continuous bag-of-words model

**Task:** given the context, predict the word in the middle

The quick brown fox jumps over the lazy dog

Similarly, we can use multiclass classification for the prediction

$$p(w_i \mid w_{i-k}, \dots, w_{i-1}, w_{i+1}, \dots, w_{i+k})$$

How to represent the context (input)?

# The continuous bag-of-words model

The context is a sequence of words.

$$c = w_{i-k}, \dots, w_{i-1}, w_{i+1}, \dots, w_{i+k}$$

$$p(w_i | c) = \frac{\exp[\phi_{\text{word}}(w_i) \cdot \phi_{\text{BoW}}(c)]}{\sum_{w \in \mathcal{V}} \exp[\phi_{\text{word}}(w) \cdot \phi_{\text{BoW}}(c)]}$$

*Handwritten note:*  $= \sum_{w \in \mathcal{V}} \phi_{\text{aux}}(w)$

# The continuous bag-of-words model

The context is a sequence of words.

$$c = w_{i-k}, \dots, w_{i-1}, w_{i+1}, \dots, w_{i+k}$$

$$\begin{aligned} p(w_i | c) &= \frac{\exp[\phi_{\text{word}}(w_i) \cdot \phi_{\text{BoW}}(c)]}{\sum_{w \in \mathcal{V}} \exp[\phi_{\text{word}}(w) \cdot \phi_{\text{BoW}}(c)]} \\ &= \frac{\exp[\phi_{\text{word}}(w_i) \cdot \sum_{w' \in c} \phi_{\text{ctx}}(w')]}{\sum_{w \in \mathcal{V}} \exp[\phi_{\text{word}}(w) \cdot \sum_{w' \in c} \phi_{\text{ctx}}(w')]} \end{aligned}$$

- $\phi_{\text{BoW}}(c)$  sums over representations of each word in  $c$
- Implementation is similar to the skip-gram model.

## Surprising observation of word embeddings

Find similar words: top- $k$  nearest neighbors using cosine similarity

- Size of window influences the type of similarity
- Shorter window produces **syntactically similar** words, e.g., Hogwarts and Sunnydale (fictional schools)
- Longer window produces **topically related** words, e.g., Hogwarts and Dumbledore (Harry Potter entities)

## Semantic ~~properties~~ of word embeddings

Solve word analogy problems: a is to b as a' is to what?

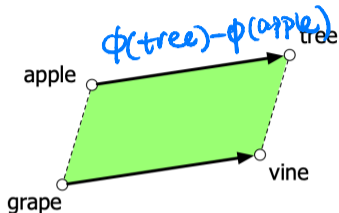


Figure: Parallelogram model (from J&H).

- man : woman :: king : queen  
 $\phi_{\text{word}}(\text{man}) - \phi_{\text{word}}(\text{king}) \approx \phi_{\text{word}}(\text{woman}) - \phi_{\text{word}}(\text{queen})$  ?
- Caveat: must exclude the three input words
- Does **not** work for general relations

## Comparison

### Count-based

matrix factorization  
fast to compute  
interpretable components

### Prediction-based

prediction problem  
slow (with large corpus)  
hard to interpret but has intriguing properties

- Both uses the **distributional hypothesis**.
- Both generalize beyond text: using co-occurrence between any types of objects
  - Learn product embeddings from customer orders
  - Learn region embeddings from images

# Evaluate word vectors

## **Intrinsic evaluation**

- Evaluate on the proxy task (related to the learning objective)
- Word similarity/analogy datasets (e.g., WordSim-353, SimLex-999)

## **Extrinsic evaluation**

- Evaluate on the real/downstream task we care about
- Use word vectors as features in applications, e.g., text classification.



# Summary

**Key idea:** formalize word representation learning as a self-supervised prediction problem

Prediction problems:

- Skip-gram: Predict context from words
- CBOW: Predict word from context
- Other possibilities:
  - Predict  $\log \hat{p}(\text{word} \mid \text{context})$ , e.g. GloVe
  - Contextual word embeddings (later)

# Table of Contents

Review

Introduction

Count-based word embeddings

Prediction-based word embeddings

Neural networks

# Feature learning

Linear predictor with **handcrafted features**:  $h(x) = w \cdot \phi(x)$ .

Can we **learn features** from data?

# Feature learning

Linear predictor with **handcrafted features**:  $h(x) = w \cdot \phi(x)$ .

Can we **learn features** from data?

Example:

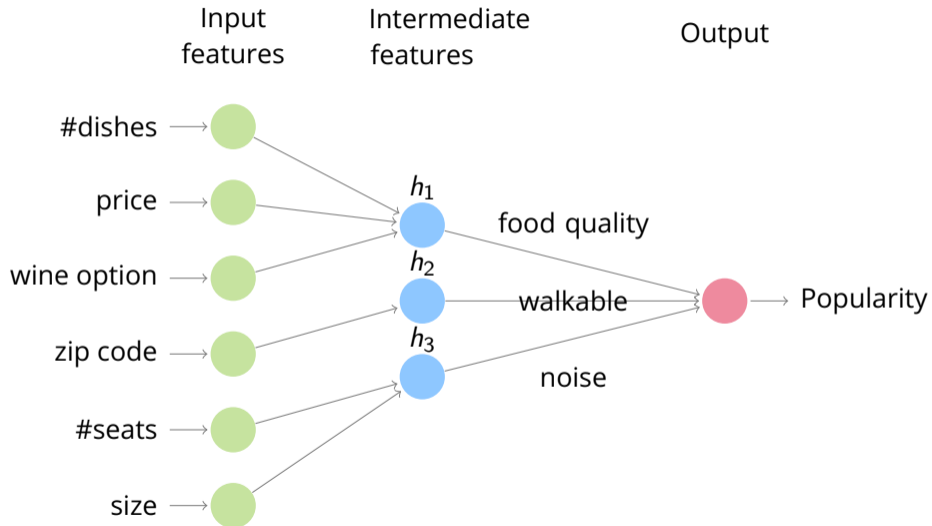
- Predict popularity of restaurants.
- Raw input: #dishes, price, wine option, zip code, #seats, size
- Decompose into subproblems:

$$h_1([\text{\#dishes, price, wine option}]) = \text{food quality}$$

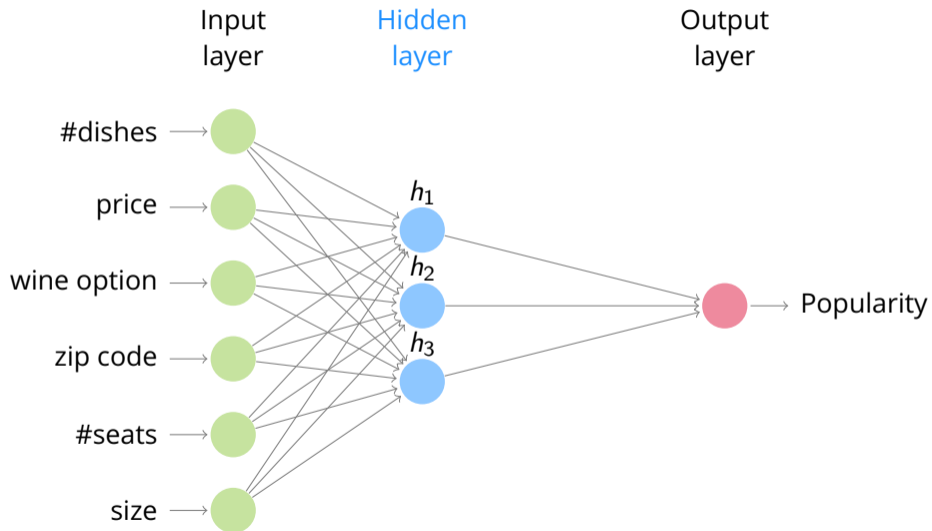
$$h_2([\text{zip code}]) = \text{walkable}$$

$$h_3([\text{\#seats, size}]) = \text{nosie}$$

# Predefined subproblems



# Learning intermediate features



# Neural networks

**Key idea:** automatically learn the intermediate features.

**Feature engineering:** Manually specify  $\phi(x)$  based on domain knowledge and learn the weights:

$$f(x) = w^T \phi(x).$$

**Feature learning:** Automatically learn both the features ( $K$  hidden units) and the weights:

$$h(x) = [h_1(x), \dots, h_K(x)], \quad f(x) = w^T h(x)$$

## Activation function

- How should we parametrize  $h_i$ 's?



## Activation function

- How should we parametrize  $h_i$ 's?

$$h_i(x) = \sigma(v_i^T x). \quad (1)$$

- $\sigma$  is the **activation function**.

# Activation function

- How should we parametrize  $h_i$ 's?

$$h_i(x) = \sigma(v_i^T x). \quad (1)$$

- $\sigma$  is the **activation function**.
- What might be some activation functions we want to use?

# Activation function

- How should we parametrize  $h_i$ 's?

$$h_i(x) = \sigma(v_i^T x). \quad (1)$$

- $\sigma$  is the **activation function**.
- What might be some activation functions we want to use?
  - sign function? **Non-differentiable**.

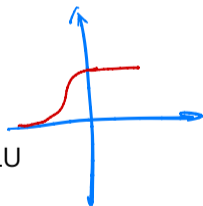


# Activation function

- How should we parametrize  $h_i$ 's?

$$h_i(x) = \sigma(v_i^T x). \quad (1)$$

- $\sigma$  is the **activation function**.
- What might be some activation functions we want to use?
  - sign function? **Non-differentiable**.
  - *Differentiable* approximations: sigmoid functions.
    - E.g., logistic function, hyperbolic tangent function, ReLU



# Activation function

- How should we parametrize  $h_i$ 's?

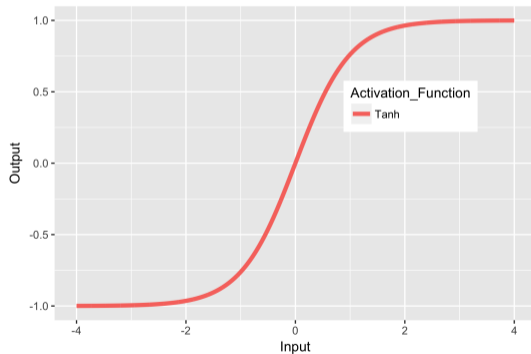
$$h_i(x) = \sigma(v_i^T x). \quad (1)$$

- $\sigma$  is the **activation function**.
- What might be some activation functions we want to use?
  - sign function? **Non-differentiable**.
  - *Differentiable* approximations: sigmoid functions.
    - E.g., logistic function, hyperbolic tangent function, ReLU
  - Non-linearity

# Activation Functions

- The **hyperbolic tangent** is a common activation function:

$$\sigma(x) = \tanh(x).$$

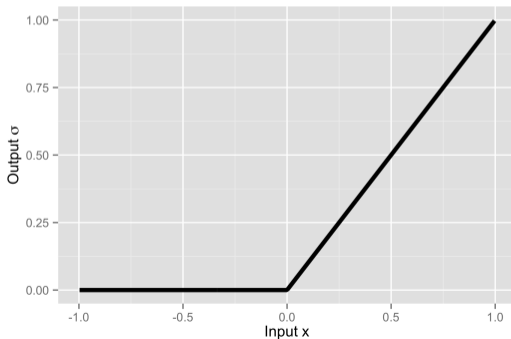


# Activation Functions

- The **rectified linear (ReLU)** function:

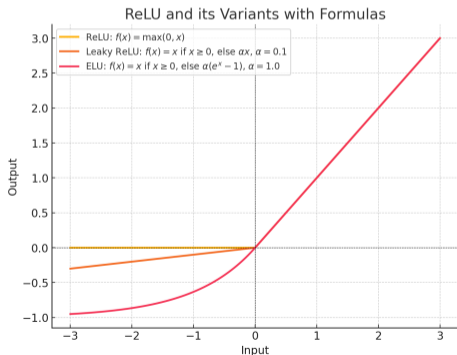
$$\sigma(x) = \max(0, x).$$

- Efficient backpropagation, sparsity, avoiding vanishing gradient
- Work well empirically.



# Activation Functions

- The dying ReLU problem: neurons become inactive
- Solution: allow small gradients when the neuron is not active
  - Still need to tune the hyperparameter  $\alpha$



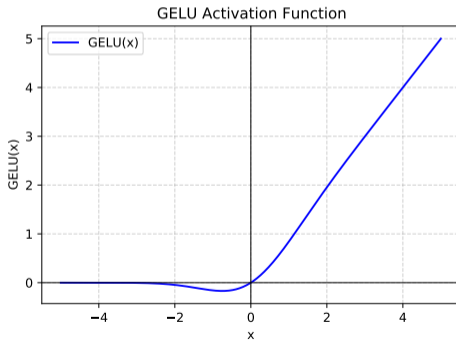


# Activation Functions

- GELU: smooth transition around the origin

$$\text{GELU}(x) = x \cdot \Phi(x)$$

$\Phi(x)$  is the CDF of the normal distribution

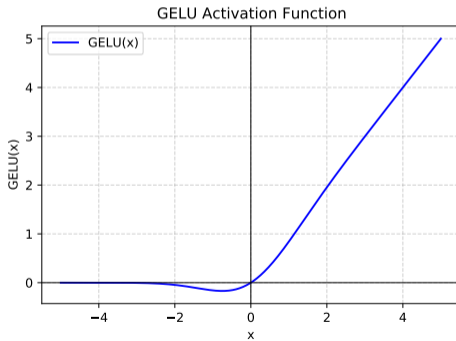


# Activation Functions

- GELU: smooth transition around the origin

$$\text{GELU}(x) = x \cdot \Phi(x)$$

$\Phi(x)$  is the CDF of the normal distribution



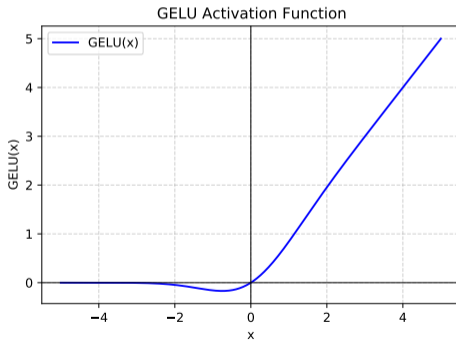
- Downside?

# Activation Functions

- GELU: smooth transition around the origin

$$\text{GELU}(x) = x \cdot \Phi(x)$$

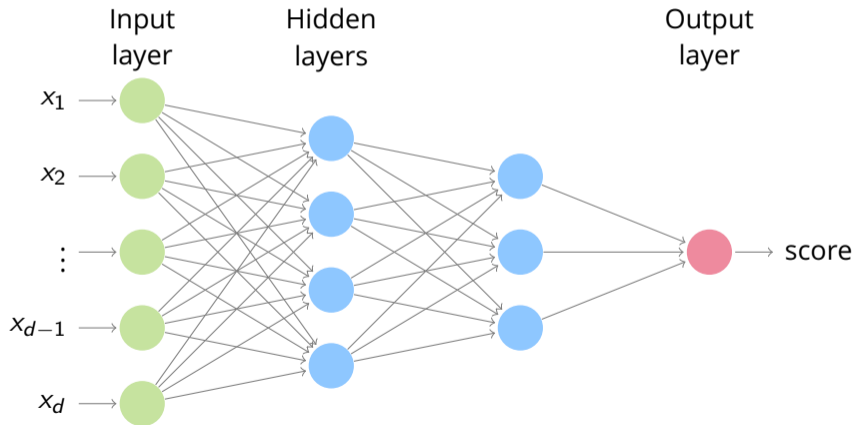
$\Phi(x)$  is the CDF of the normal distribution



- Downside? More compute (but negligible)

# Multilayer perceptron / Feed-forward neural networks

- Wider: more hidden units.
- Deeper: more hidden layers.



## Multilayer Perceptron: Standard Recipe

- Each subsequent hidden layer takes the **output  $o \in \mathbb{R}^m$  of previous layer** and produces

$$h^{(j)}(o^{(j-1)}) = \sigma \left( W^{(j)} o^{(j-1)} + b^{(j)} \right), \text{ for } j = 2, \dots, L$$

where  $W^{(j)} \in \mathbb{R}^{m \times m}$ ,  $b^{(j)} \in \mathbb{R}^m$ .

## Multilayer Perceptron: Standard Recipe

- Each subsequent hidden layer takes the **output  $o \in \mathbb{R}^m$  of previous layer** and produces

$$h^{(j)}(o^{(j-1)}) = \sigma \left( W^{(j)} o^{(j-1)} + b^{(j)} \right), \text{ for } j = 2, \dots, L$$

where  $W^{(j)} \in \mathbb{R}^{m \times m}$ ,  $b^{(j)} \in \mathbb{R}^m$ .

- Last layer is an *affine* mapping (no activation function):

$$a(o^{(L)}) = W^{(L+1)} o^{(L)} + b^{(L+1)},$$

where  $W^{(L+1)} \in \mathbb{R}^{k \times m}$  and  $b^{(L+1)} \in \mathbb{R}^k$ .

## Multilayer Perceptron: Standard Recipe

- Each subsequent hidden layer takes the **output**  $o \in \mathbb{R}^m$  of previous layer and produces

$$h^{(j)}(o^{(j-1)}) = \sigma \left( W^{(j)} o^{(j-1)} + b^{(j)} \right), \text{ for } j = 2, \dots, L$$

where  $W^{(j)} \in \mathbb{R}^{m \times m}$ ,  $b^{(j)} \in \mathbb{R}^m$ .

- Last layer is an *affine* mapping (no activation function):

$$a(o^{(L)}) = W^{(L+1)} o^{(L)} + b^{(L+1)},$$

where  $W^{(L+1)} \in \mathbb{R}^{k \times m}$  and  $b^{(L+1)} \in \mathbb{R}^k$ .

- The full neural network function is given by the *composition* of layers:

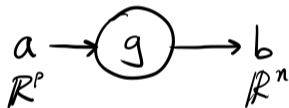
$$f(x) = \left( a \circ h^{(L)} \circ \dots \circ h^{(1)} \right) (x) \tag{2}$$

# Computation graphs

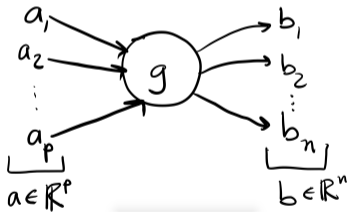
(adapted from David Rosenberg's slides)

Function as a *node* that takes in *inputs* and produces *outputs*.

- Typical computation graph:



- Broken out into components:

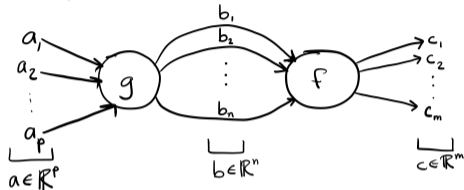




# Compose multiple functions

(adpated from David Rosenberg's slides)

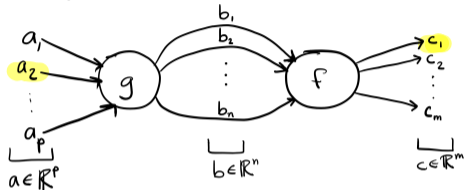
Compose two functions  $g : \mathbb{R}^p \rightarrow \mathbb{R}^n$  and  $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ :  $c = f(g(a))$



# Compose multiple functions

(adpated from David Rosenberg's slides)

Compose two functions  $g : \mathbb{R}^p \rightarrow \mathbb{R}^n$  and  $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ :  $c = f(g(a))$

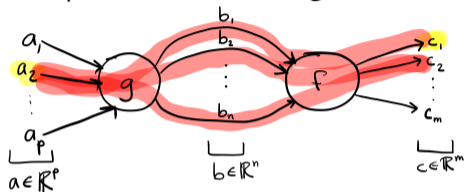


- Derivative: How does change in  $a_j$  affect  $c_i$ ?

# Compose multiple functions

(adpated from David Rosenberg's slides)

Compose two functions  $g : \mathbb{R}^p \rightarrow \mathbb{R}^n$  and  $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ :  $c = f(g(a))$



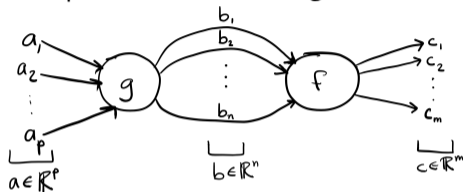
- Derivative: How does change in  $a_j$  affect  $c_i$ ?

$$\frac{\partial c_i}{\partial a_j} = \sum_{k=1}^n \frac{\partial c_i}{\partial b_k} \frac{\partial b_k}{\partial a_j}$$

# Compose multiple functions

(adpated from David Rosenberg's slides)

Compose two functions  $g : \mathbb{R}^p \rightarrow \mathbb{R}^n$  and  $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ :  $c = f(g(a))$



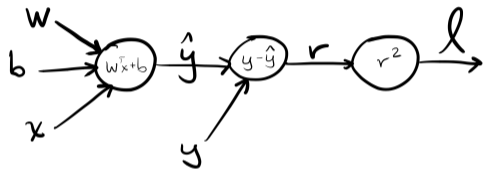
- Derivative: How does change in  $a_j$  affect  $c_i$ ?

$$\frac{\partial c_i}{\partial a_j} = \sum_{k=1}^n \frac{\partial c_i}{\partial b_k} \frac{\partial b_k}{\partial a_j}.$$

- Visualize the multivariable **chain rule**:
  - **Sum** changes induced on all paths from  $a_j$  to  $c_i$ .
  - Changes on one path is the **product** of changes across each node.

# Computation graph example

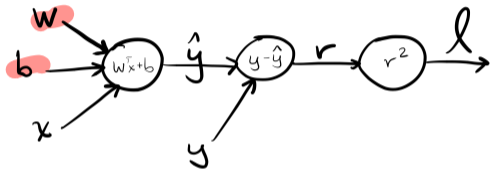
(adapted from David Rosenberg's slides)



(What is this graph computing?)

# Computation graph example

(adapted from David Rosenberg's slides)

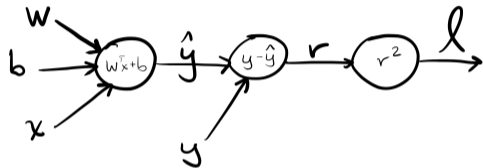


(What is this graph computing?)

$$\frac{\partial l}{\partial b} = \frac{\partial l}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial b} = (-2r)(1) = -2r$$
$$\frac{\partial l}{\partial w_j} = \frac{\partial l}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial w_j} = (-2r)x_j = -2rx_j$$

# Computation graph example

(adapted from David Rosenberg's slides)



(What is this graph computing?)

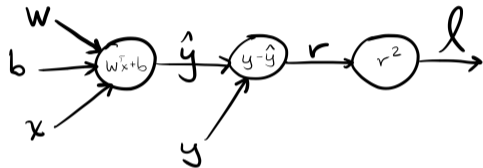
$$\frac{\partial l}{\partial b} = \frac{\partial l}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial b} = (-2r)(1) = -2r$$

$$\frac{\partial l}{\partial w_j} = \frac{\partial l}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial w_j} = (-2r)x_j = -2rx_j$$

Computing the derivatives in certain order allows us to save compute!

# Computation graph example

(adapted from David Rosenberg's slides)



(What is this graph computing?)

$$\begin{aligned}\frac{\partial l}{\partial r} &= 2r \\ \frac{\partial l}{\partial \hat{y}} &= \frac{\partial l}{\partial r} \frac{\partial r}{\partial \hat{y}} = (2r)(-1) = -2r \\ \frac{\partial l}{\partial b} &= \frac{\partial l}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial b} = (-2r)(1) = -2r \\ \frac{\partial l}{\partial w_j} &= \frac{\partial l}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial w_j} = (-2r)x_j = -2rx_j\end{aligned}$$

Computing the derivatives in certain order allows us to save compute!



# Backpropogation

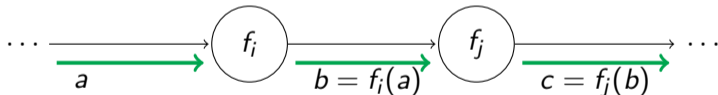
Backpropogation = chain rule + dynamic programming on a computation graph

# Backpropogation

Backpropogation = chain rule + dynamic programming on a computation graph

Forward pass

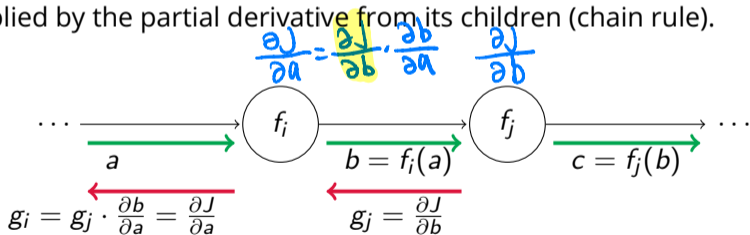
- **Topological order:** every node appears before its children
- For each node, compute the output given the input (from its parents).



# Backpropagation

## Backward pass

- **Reverse topological order:** every node appear after its children
- For each node, compute the partial derivative of its output w.r.t. its input, multiplied by the partial derivative from its children (chain rule).



# Summary

Key idea in neural nets: feature/representation learning

Building blocks:

- Input layer: raw features (no learnable parameters)
- Hidden layer: perceptron + nonlinear activation function
- Output layer: linear (+ transformation, e.g. softmax)

Optimization:

- Optimize by SGD (implemented by back-propagation)
- Objective is non-convex, may not reach a global minimum