

CS769 Advanced NLP

Word Embeddings and Text Classification

Junjie Hu



Slides adapted from Noah, Yulia
<https://junjiehu.github.io/cs769-fall25/>

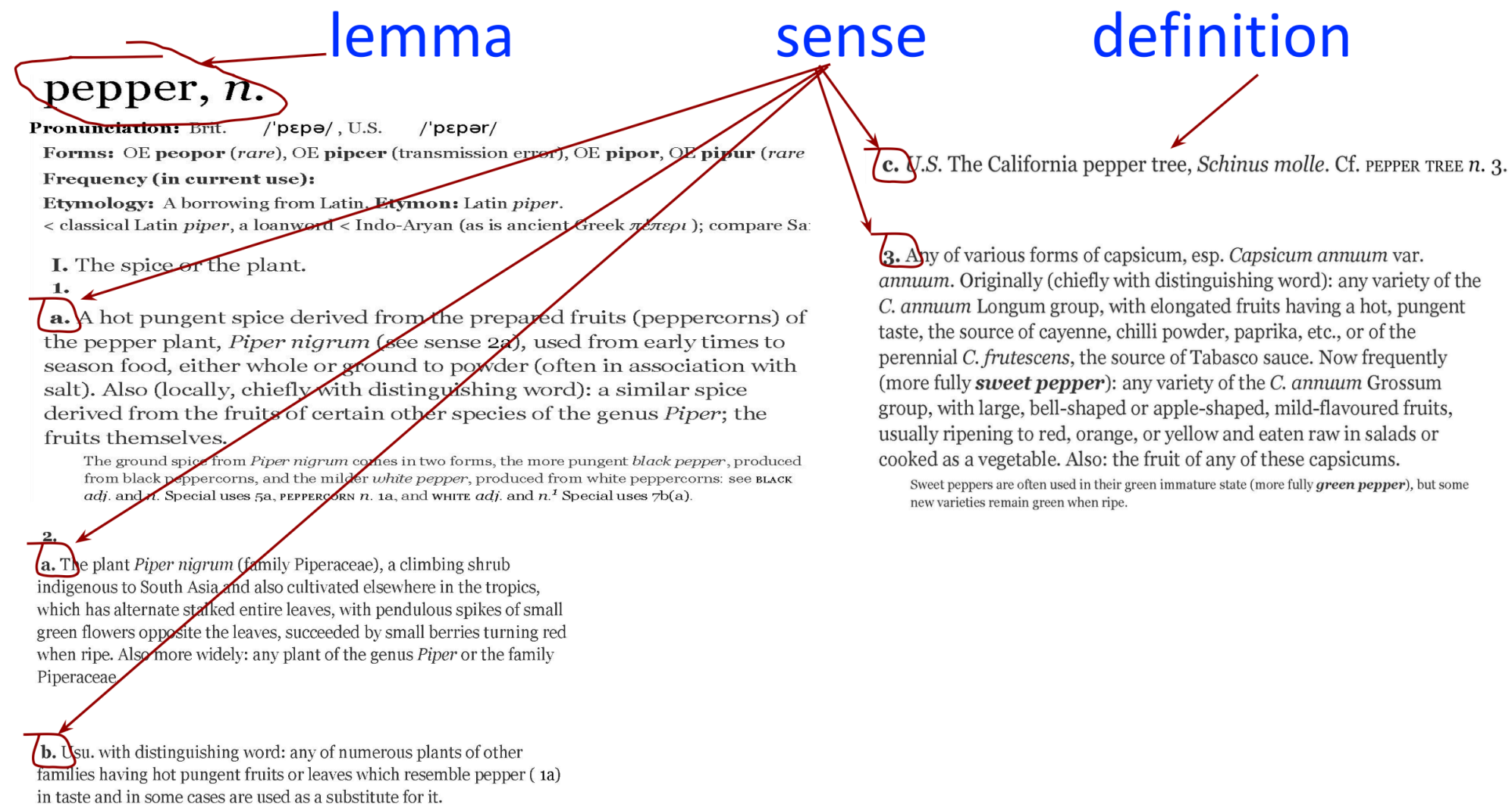
Goals for Today

- Lexical Semantics and Distributional Semantics
- **Word Embeddings** (e.g., Skip-gram, CBOW)
- Evaluation (intrinsic and extrinsic)
- Text Classification

How should we represent the
meaning of the word?

Lexical Semantics

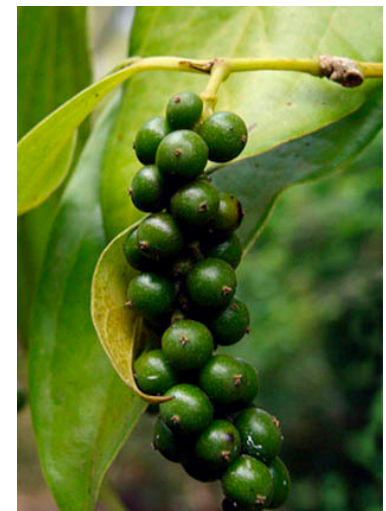
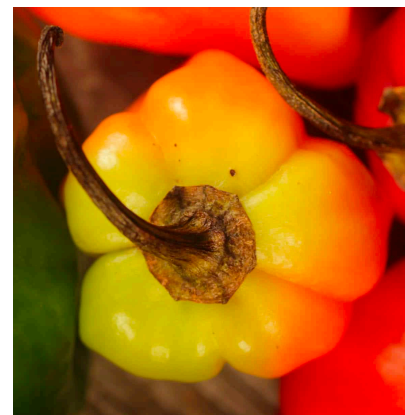
- How should we represent the meaning of the word?
- Words, lemmas, senses, definition



Oxford English Dictionary: <https://www.oed.com/>

Lemma pepper

- Sense 1: spice from pepper plant
- Sense 2: the pepper plant itself
- Sense 3: another similar plant (Jamaican pepper)
- Sense 4: plant with peppercorns (California pepper)
- Sense 5: capsicum (i.e., chili, paprika, bell pepper, etc)



Lexical Semantics

- How should we represent the meaning of the word?
 - Words, lemmas, senses, definition
 - Relationships between words or senses
 1. Synonymy: same meaning, e.g., couch/sofa
 2. Antonymy: opposite senses, e.g., hot/cold
 3. Similarity: similar meanings, e.g., car/bicycle
 4. Relatedness: association, e.g., car/gasoline
 5. Superordinate/Subordinate: e.g., car/vehicle, mango/fruit

Lexical Semantics

- How should we represent the meaning of the word?
 - Words, lemmas, senses, definition
 - Relationships between words or senses
 - Taxonomy: abstract -> concrete

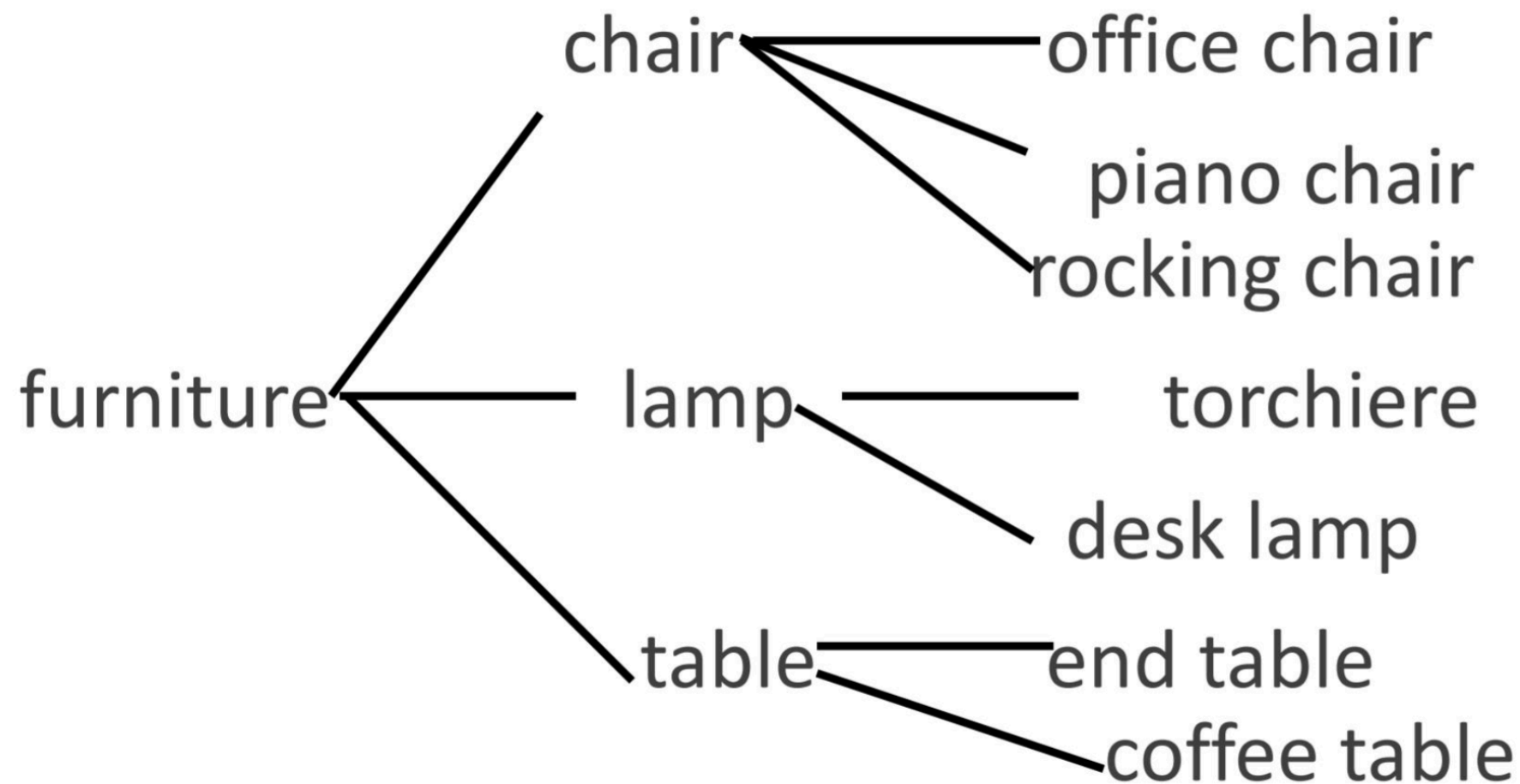
Taxonomy

- abstract -> concrete

Superordinate

Basic

Subordinate



Lexical Semantics

- How should we represent the meaning of the word?
 - Words, lemmas, senses, definition
 - Relationships between words or senses
 - Taxonomy: abstract -> concrete
 - Semantic frames and roles

Semantic Frame

- A set of words that denote perspectives or participants in an event

- Tom brought a book from Bill.

Tom brought a book from Bill
buyer event from the perspective of the buyer

- Bill sold a book to Tom.

Bill sold a book to Tom
seller event from the perspective of the seller

Mismatch

- Theories of language tend to view the data (words, sentences, documents) and abstractions over it as *symbolic* or categorical.
 - Uses *symbols* to represent linguistic information
- Machine learning algorithms built on optimization rely more on *continuous* data.
 - Uses *floating-point numbers (vectors)*

Problems with Discrete Representations

- Too coarse: *expert* \leftrightarrow *skillful*
- Sparse
- Subjective
- Expensive
- Hard to compute word relationships

One-hot vector:

<i>expert</i>	[0	0	0	1	0	0	0	0	0	0	0	0	0	0]
<i>skillful</i>	[0	0	0	0	0	0	0	0	0	0	1	0	0	0]

Distributional Hypothesis

“The meaning of a word is its use in the language”

[Wittgenstein 1943]

“You shall know a word by the company it keeps”

[Firth 1957]

“If A and B have almost identical environments we say that they are synonyms.”

[Harris 1954]

Example

- What does “Ong Choy” mean?
 - Suppose you see these sentences:
 - Ong Choy is delicious **sautéed with garlic**
 - Ong Choy is superb **over rice**
 - Ong Choy **leaves** with salty sauces
 - And you’ve also seen these:
 - ... water spinach **sautéed with garlic over rice**
 - Chard stems and **leaves** are delicious
 - Collard greens and other **salty leafy** greens

Ong Choy \approx “Water Spinach”?

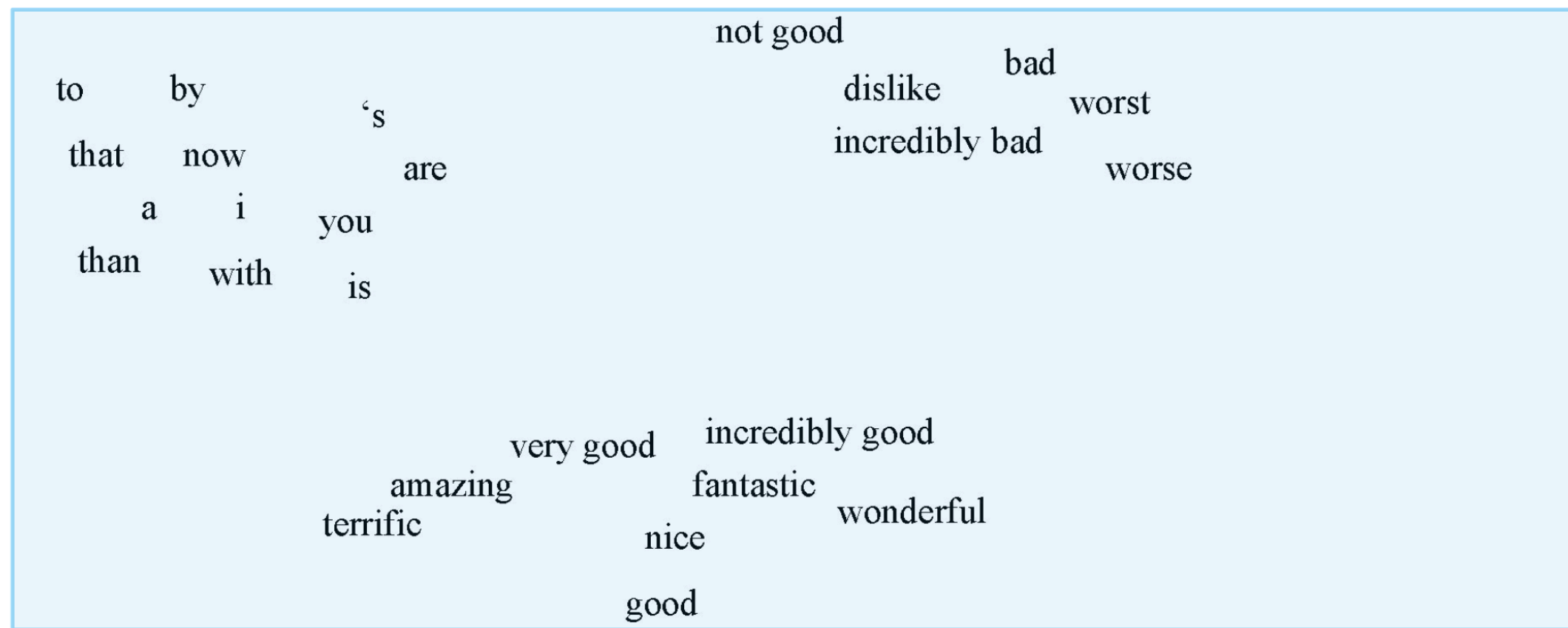
- Ong Choy is a leafy green like spinach, chard, or collard greens



Ong Choy: pronunciation of “薺菜” in Cantonese

Model of Meaning Focusing on Similarity

- Each word = a vector
 - Similar words are “nearby in space”
 - the standard way to represent meaning in NLP



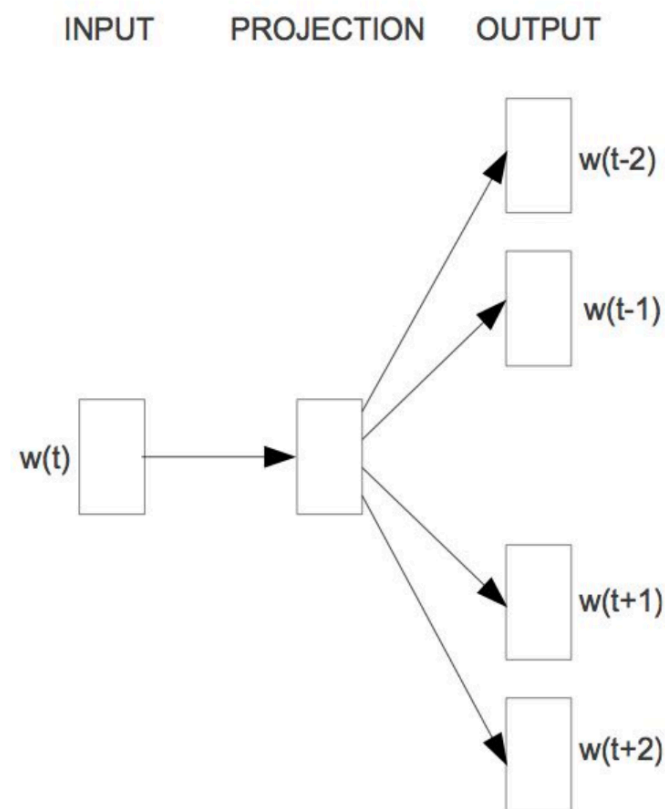
Distributed Word Embeddings

Word Vector Models

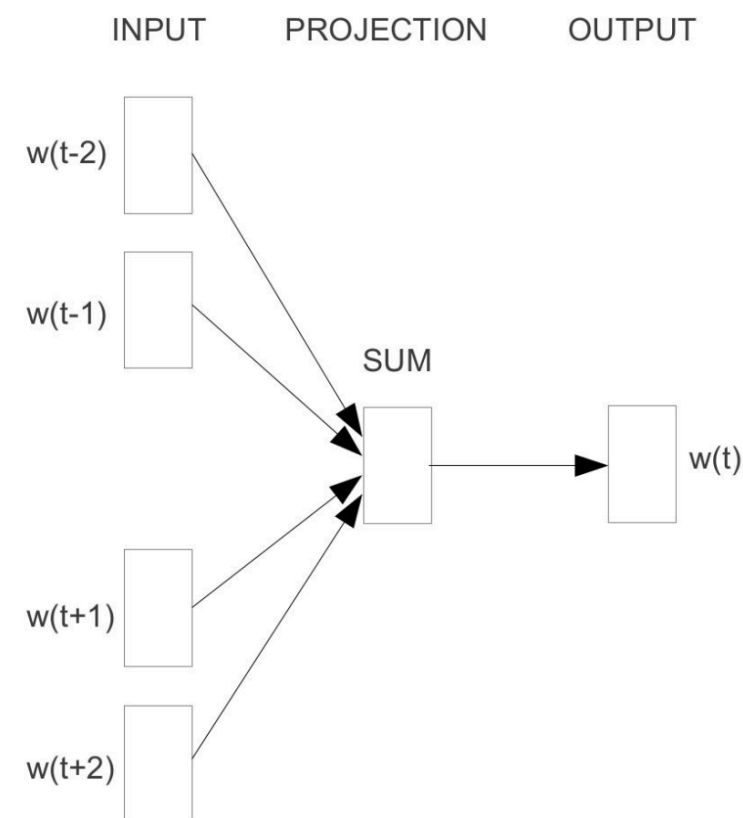
- These models are designed to “guess” a word at position i given a word at a position in $\{i - w, \dots, i - 1\} \cup \{i + 1, \dots, i + w\}$
- “Pre-train” word vectors are used in other larger models (e.g., neural LM)

Word2vec

- Continuous bag of words (CBOW): $p(v \mid c)$
 - Similar to feedforward neural LM w/o the feedforward layers in Lecture 3.
- Skip-gram: $p(c \mid v)$



Skip-gram



CBOW

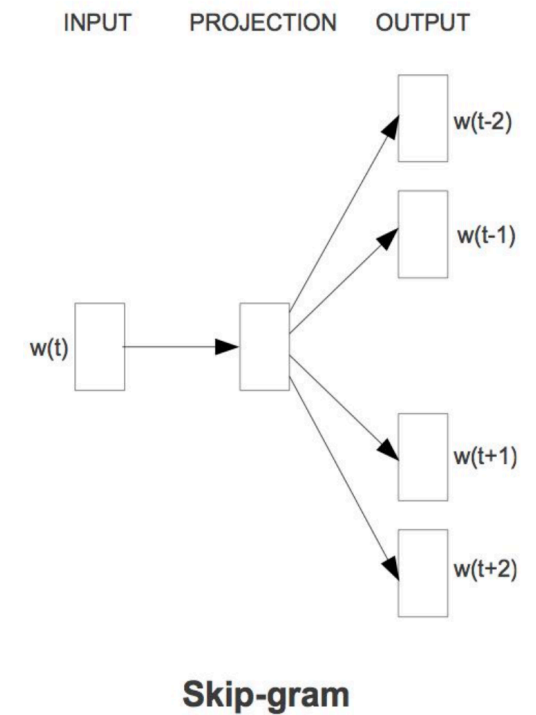
Skip-gram Prediction

- Predict vs Count

the cat sat on the mat



context size = 2



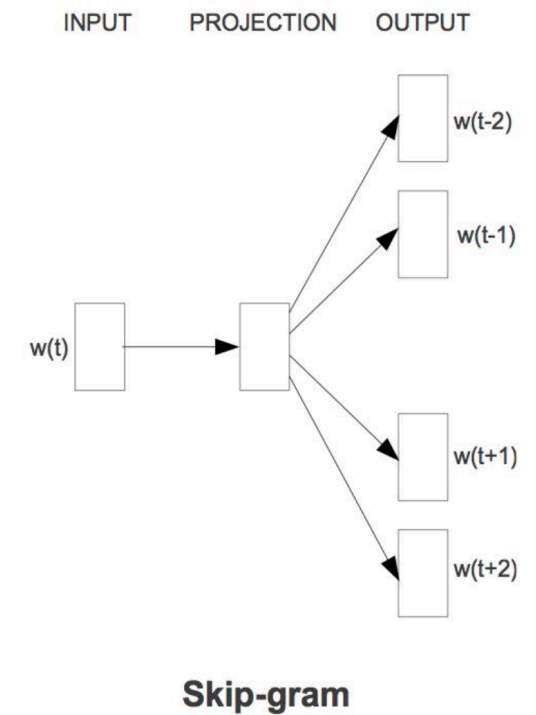
Skip-gram Prediction

- Predict vs Count

the cat sat on the mat



context size = 2



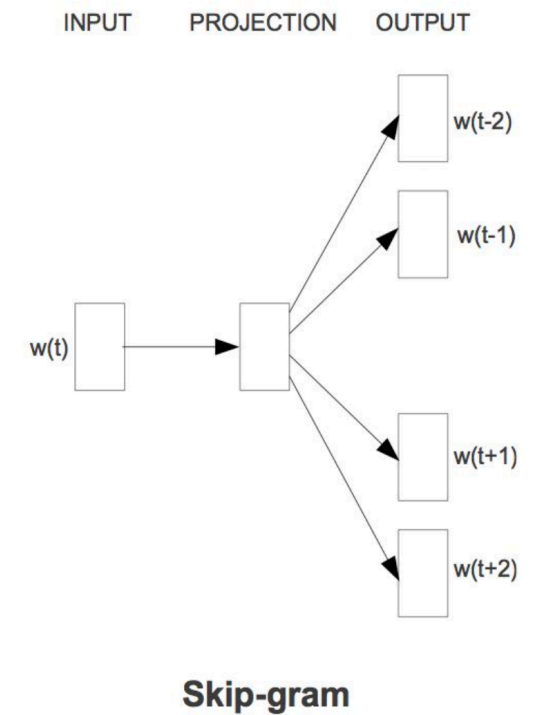
Skip-gram Prediction

- Predict vs Count

the cat sat on the mat



context size = 2



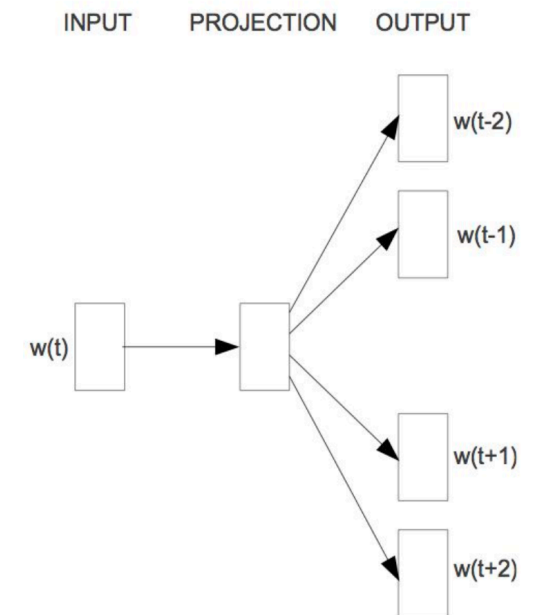
Skip-gram Prediction

- Predict vs Count

the cat sat on the mat



context size = 2

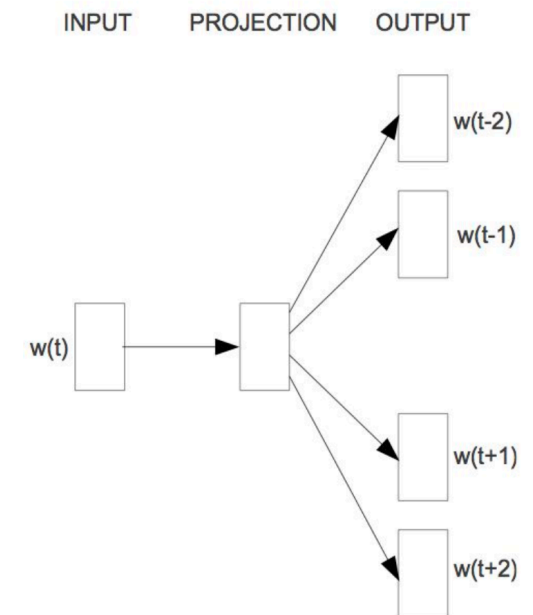


Skip-gram

Skip-gram Prediction

- Predict vs Count

the cat sat on the mat



Skip-gram



context size = 2

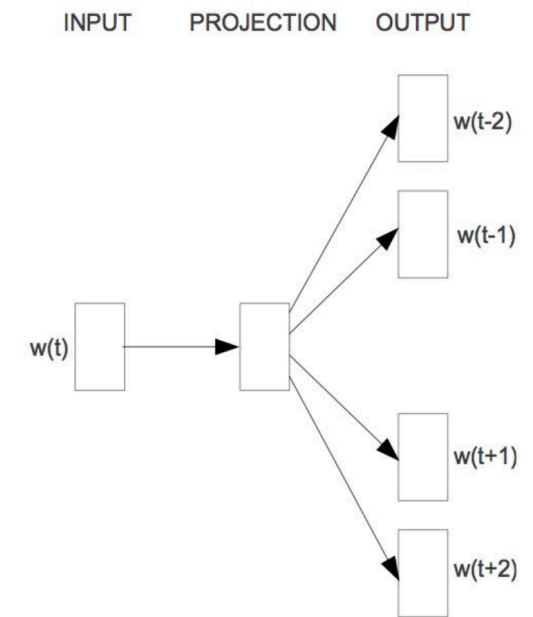
Skip-gram Prediction

- Predict vs Count

the cat sat on the mat



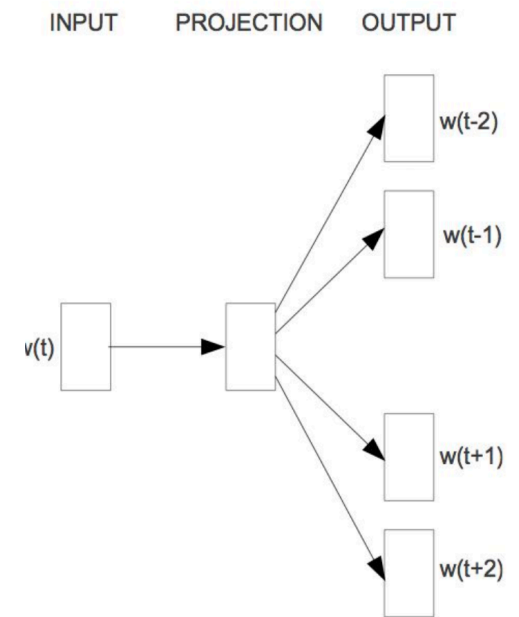
context size = 2



Skip-gram

Skip-gram Prediction

- The same word can appear in different context.



Skip-gram

context size = 2

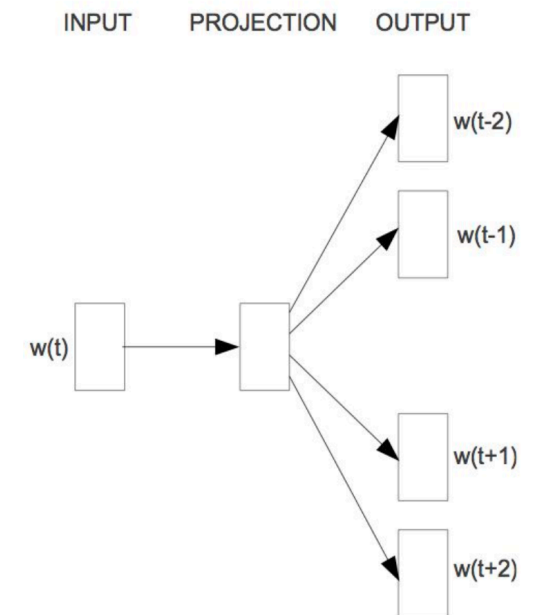
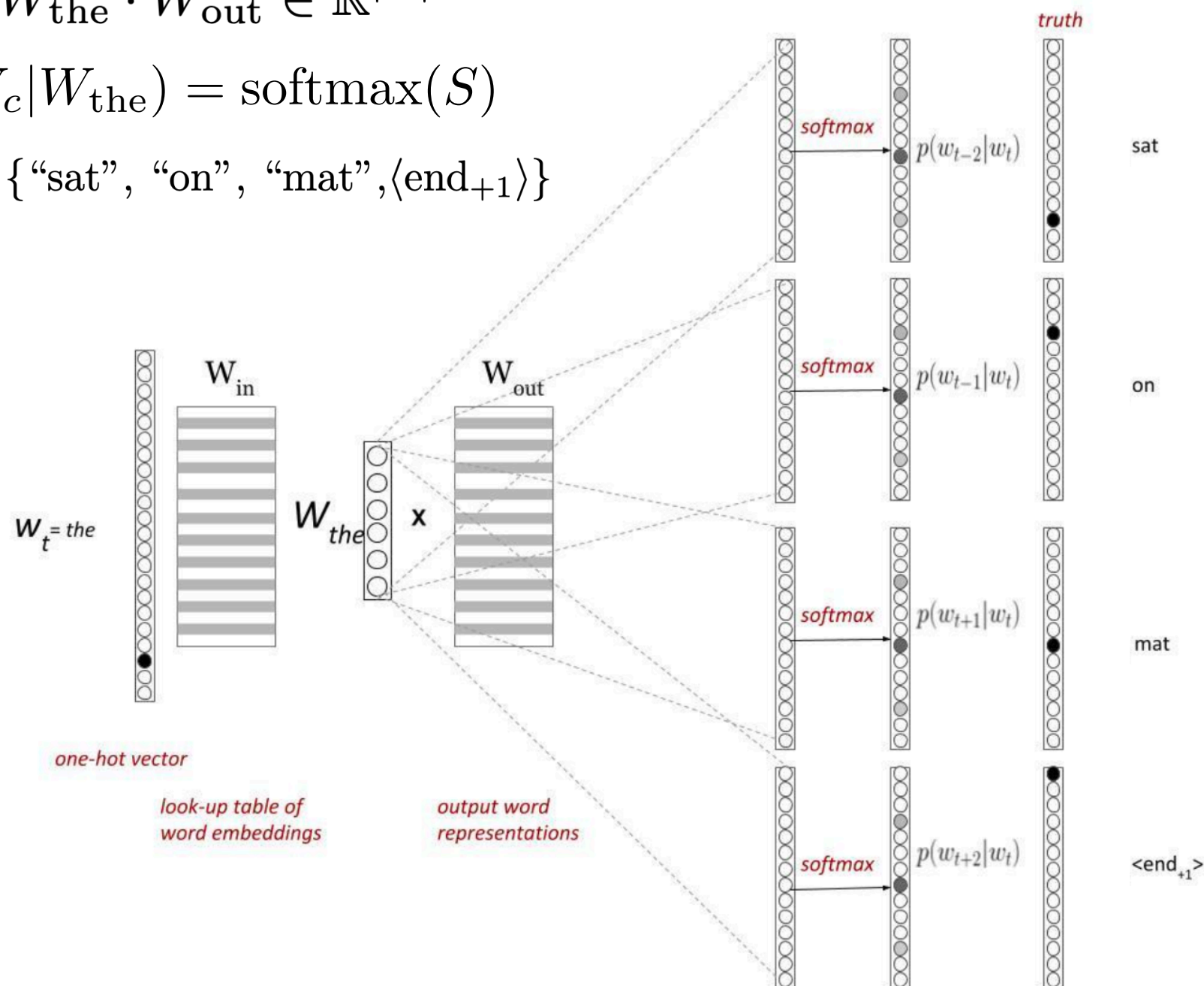
Skip-gram Prediction

$$W_{\text{the}} = \text{LookUp}(W_{\text{in}}, \text{"the"}) \in \mathbb{R}^d, W_{\text{in}} \in \mathbb{R}^{|V| \times d}$$

$$S = W_{\text{the}} \cdot W_{\text{out}} \in \mathbb{R}^{|V|}$$

$$P(W_c | W_{\text{the}}) = \text{softmax}(S)$$

$$W_c = \{\text{"sat"}, \text{"on"}, \text{"mat"}, \langle \text{end}_{+1} \rangle\}$$



Skip-gram

Skip-gram Objective

- For each word in the corpus

$$J(\Theta) = \prod_{t=1}^T \prod_{-m \leq j \leq m, j \neq 0} p(w_{t+j} | w_t; \Theta)$$

$$J(\Theta) = -\frac{1}{T} \sum_{t=1}^T \sum_{-m \leq j \leq m, j \neq 0} \log p(w_{t+j} | w_t; \Theta)$$

Maximize the probability of any context window given the current center word

Skip-gram Objective

- For each word in the corpus

$$J(\Theta) = -\frac{1}{T} \sum_{t=1}^T \sum_{-m \leq j \leq m, j \neq 0} \log p(w_{t+j} | w_t; \Theta)$$

$$p(w_{t+j} | w_t) = p(o|c) = \frac{\exp(u_o^\top v_c)}{\sum_{i=1}^V \exp(u_i^\top v_c)}$$

dot product
(similarity)
between outside
and center word
vectors

Notation simplification:

o = index of outside (context) word

c = index of center word (w_t)

V = vocab size, V can be large 50K - 30M

Skip-gram w/ negative sampling

- $V=50K-30M$, too large!

$$p(w_{t+j}|w_t) = p(o|c) = \frac{\exp(u_o^\top v_c)}{\sum_{i=1}^V \exp(u_i^\top v_c)}$$

- Negative sampling:
 - Treat the center word and a neighboring context word as positive examples.
 - Randomly sample other words in the lexicon to get negative samples.

(banking, regulation)

(banking, aardvark)

Skip-gram w/ negative sampling

- Convert the task to binary classification rather than multiclass:

$$P(o \mid c) = \frac{\exp(u_o^T v_c)}{\sum_{i=1}^V \exp(u_i^T v_c)} \longrightarrow P(o \mid c) = \frac{1}{1 + \exp(-u_o^T v_c)} = \sigma(u_o^T v_c)$$

- New objective (single context word, k negative samples):

$$\log P(o_+ \mid c) + \sum_{i=1}^k \log(1 - P(o_i \mid c))$$

Choosing negative samples

- Pick negative samples according to unigram frequency $P(w)$
- More common to choose according to:

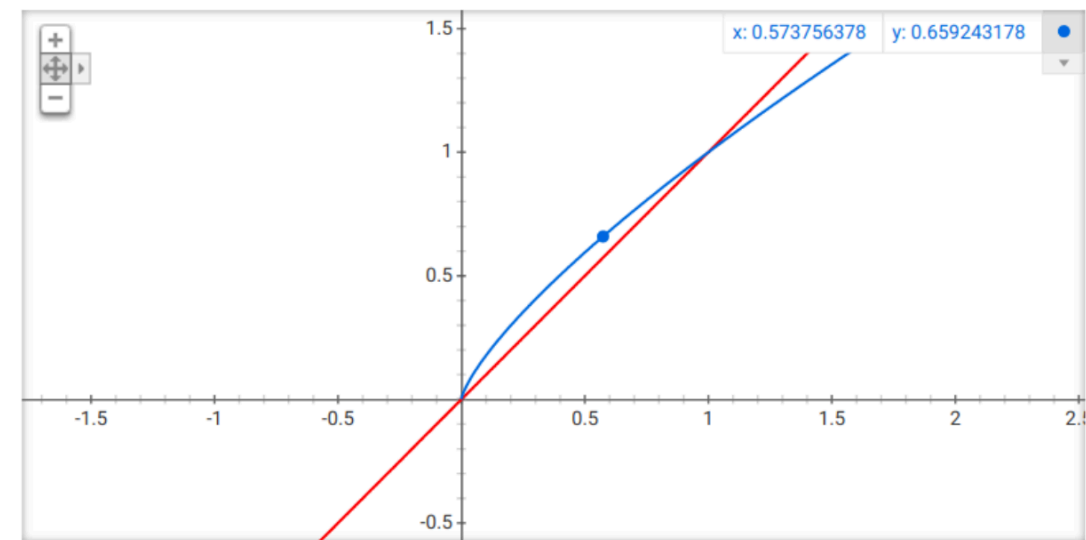
$$P_{\alpha}(w) = \frac{\text{count}(w)^{\alpha}}{\sum_w \text{count}(w)^{\alpha}}$$

- $\alpha = 0.75$ works well empirically
- Gives rare words slightly higher probability
- e.g., $P(a) = 0.99$, $P(b) = 0.01$

$$P_{\alpha}(a) = \frac{0.99^{0.75}}{0.99^{0.75} + 0.01^{0.75}} = 0.97$$

$$P_{\alpha}(b) = \frac{0.01^{0.75}}{0.99^{0.75} + 0.01^{0.75}} = 0.03$$

Graph for $x^{3/4}$, x



Available dense embeddings

- Word2vec (Mikolov et al. 2013)
 - <https://code.google.com/archive/p/word2vec/>
- GloVe (Pennington et al. 2014)
 - <http://nlp.stanford.edu/projects/glove/>
- Fasttext (Bojanowski et al. 2017)
 - <http://www.fasttext.cc/>

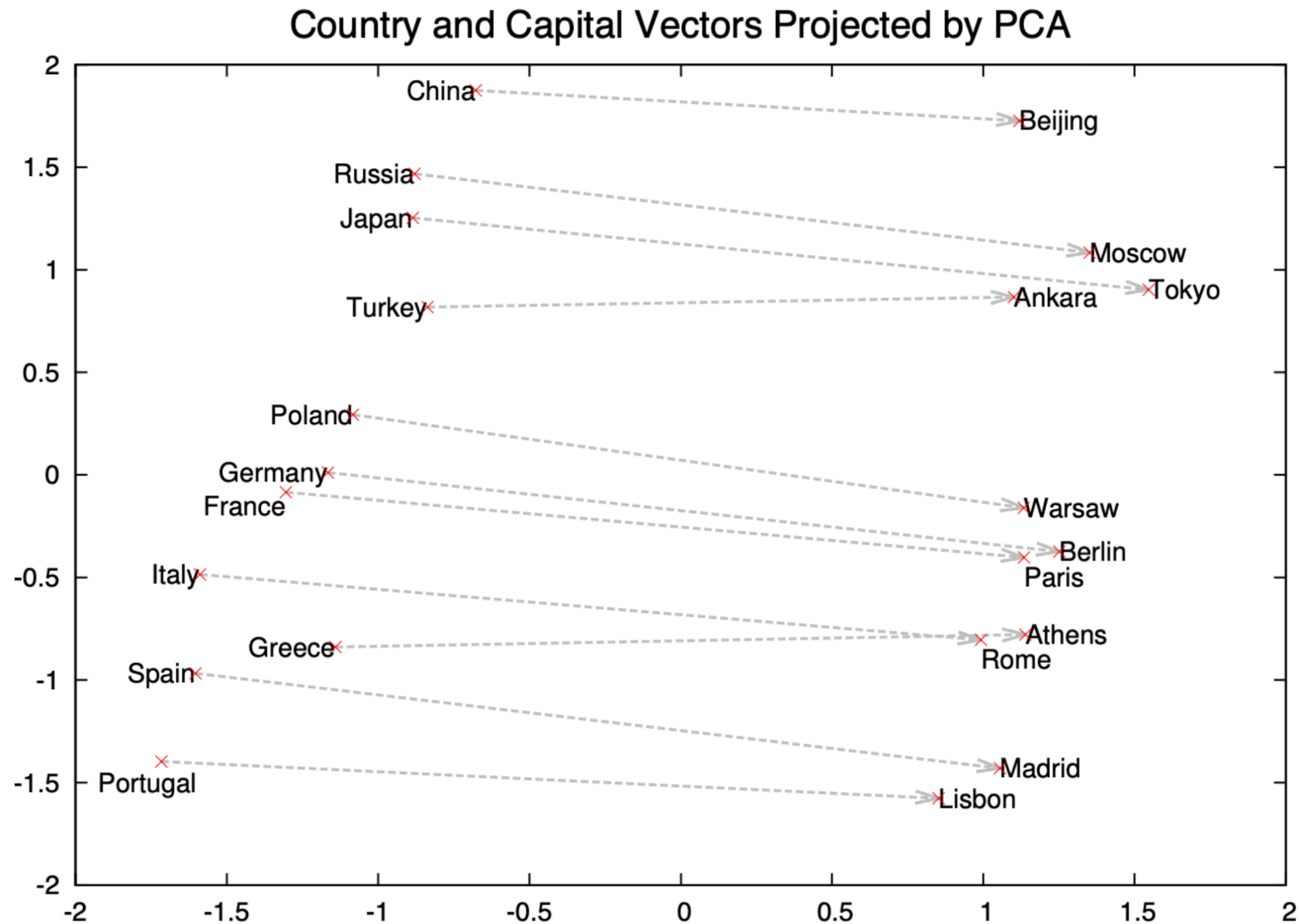
Evaluation

— how well do word vectors capture embedding similarity?

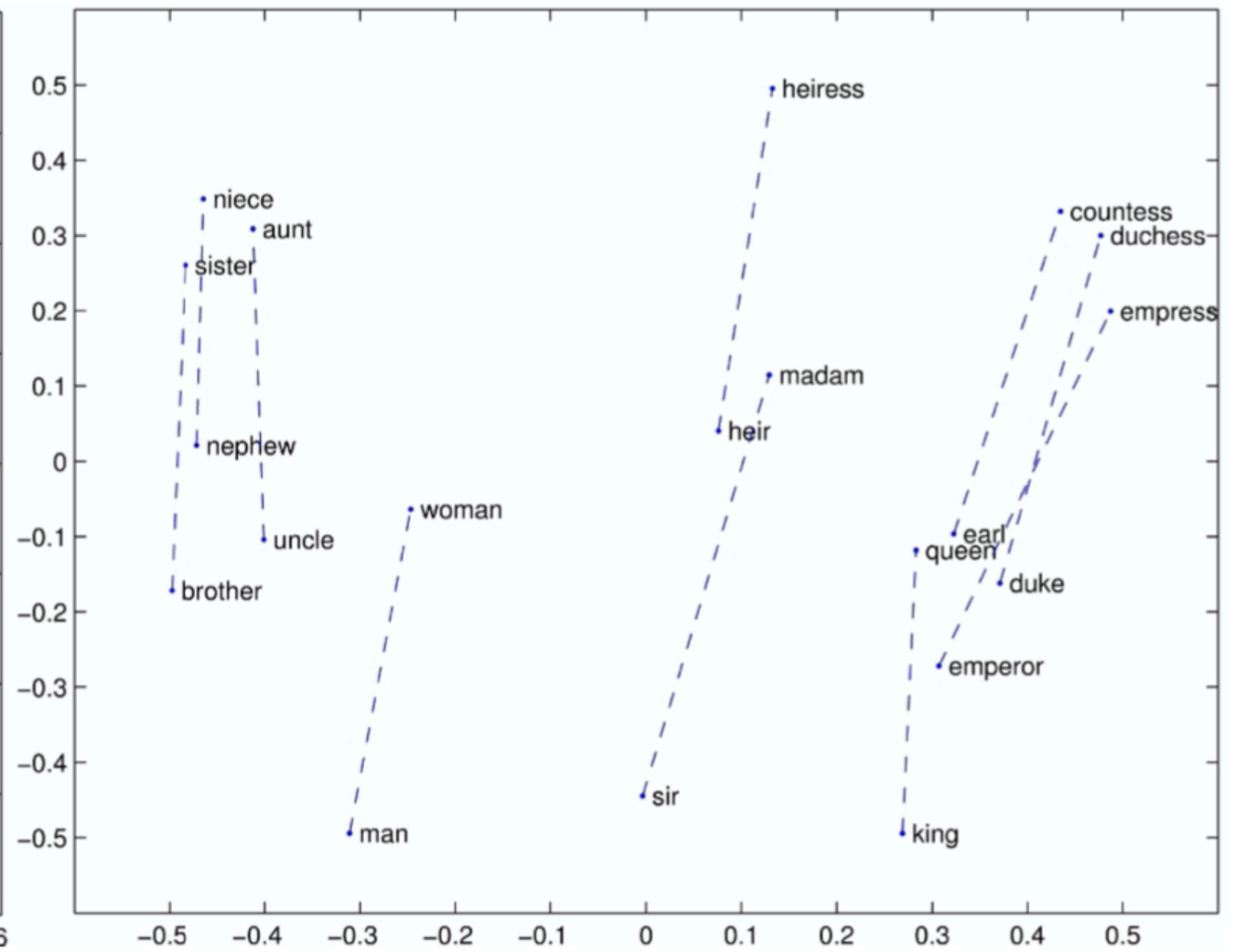
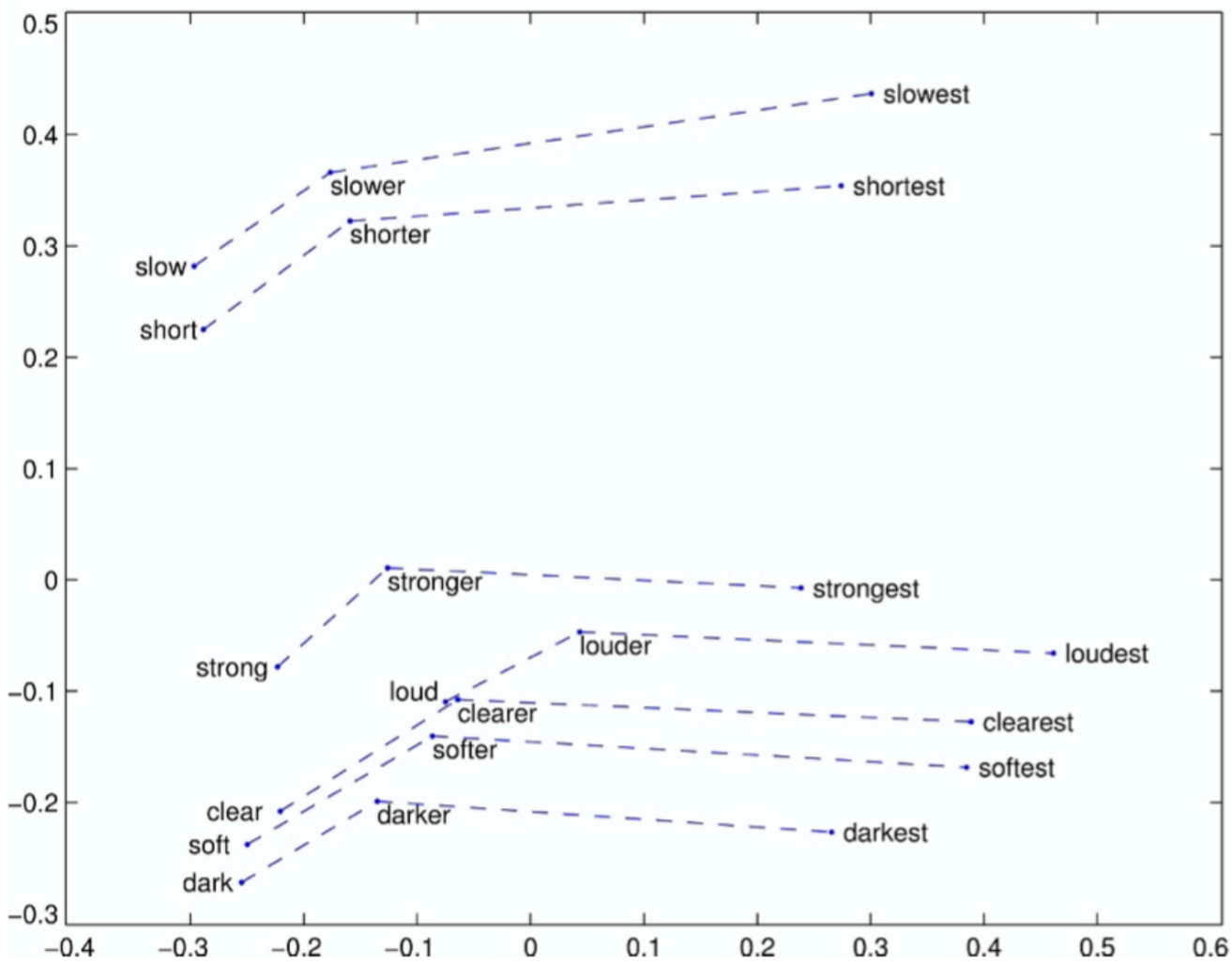
Evaluating word vectors

- **Intrinsic evaluation:** test whether the representations align with our intuitions about word meaning.
 - How well does cosine similarity of word embeddings correlate with human judgements?
 - Completing analogies: $a:b \leftrightarrow c: ?$
- **Extrinsic evaluation:** test whether the representations are useful for downstream tasks, such as tagging, parsing, QA, ...
 - Provide embeddings as input to the same classifier, how well does a model w/ pre-trained embeddings perform?

A:B <-> C:?



A:B <-> C:?




Text Classification

Text Classification

- Classify sentences according to various traits
- Topic, sentiment, subjectivity/objectivity, etc.

I hate this movie



The diagram shows a black arrow pointing from the end of the sentence "I hate this movie" to a vertical list of three sentiment labels: "positive" (green), "neutral" (black), and "negative" (red). The "negative" label is positioned at the tip of the arrow, indicating the classification result.

positive
neutral
negative

Example: Movie Review

- Predict sentiment from IMDB movie review:
{positive, neutral, negative}

★ 10/10

So glad to have you back Guiliana!
[Love_Life_Laughter](#) 2 September 2018

I watched the ugly breakup of the E! channel circle around Joan Rivers, after her tragic and unexpected demise, with such sadness. I am so glad to have Guiliana back. She has such a unique persona, equally comfortable talking about fashion and cancer recovery, that somehow provides the perfect balance between appropriate gravitas and celebratory girliness. Bravo E!

positive

★ 2/10

Not funny or entertaining at all
[gtamaniak-16300](#) 5 November 2019

Some times nitpicking films and other media might be interesting to watch if the are some clever jokes made by the editor. Here, while some plot holes of certain films are revealed which I wouldn't even think of noticing and that's a good thing, there are some observations that are born from the poor imagination of the narrator and don't have anything to do with the certain media reviewed. Absolutely unfunny moments that don't even make me chuckle. At least the narrator's voice is OK and not like one of those annoying British accents from whatculture.

negative

IMDb Menu All Search IMDb

Sort by IMDb Rating - Highest Rated Movies and TV Shows tagged with keyword "movie-review"

Refine Movie Review

1 to 50 of 659 titles | Next » Sort by: IMDb Rating View:

1. **Freaking Pause pra TV** (2012-)
Adventure, Comedy
★ 9.5 Rate this
Add a Plot
Star: [Wilson Rafael de Azevedo](#)
Votes: 40

2. **Best of the Worst** (2013-)
Comedy, Talk-Show
★ 9.4 Rate this
Mike, Jay, Rich and the rest of the Red Letter Media crew brave some of the worst movies ever created by man.
Stars: [Jay Bauman](#), [Rich Evans](#), [Mike Stoklasa](#), [Jack Packard](#)
Votes: 3,155

Example: Customer Rating

- Predict Amazon customer rating: {1, 2, 3, 4, 5}



Patrick Montoya **Top Contributor: Baby**

★★★★★ **Great batteries**

Reviewed in the United States on March 23, 2019

Size: 100 Count | **Verified Purchase**

The batteries last forever. It's nice to have a huge box like this.

43 people found this helpful



N. Caruso **VINE VOICE**

★☆☆☆☆ **Poor quality!**

Reviewed in the United States on May 24, 2018

Size: 48 Count | **Verified Purchase**

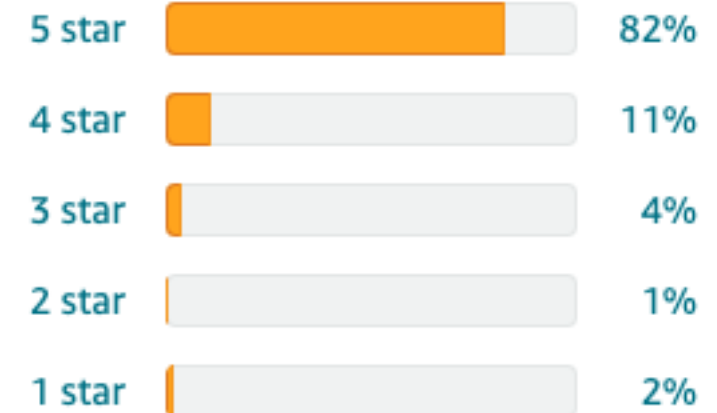
I loved these batteries when I first started buying them. They were cheap and seemed to last. Then I bought 48! I figured - cheaper by the (4) dozen - Big mistake! Not only don't they hold a charge in a device, they are actually dead coming out of the shrink wrap. I actually had a couple literally fall apart in a camera. Luckily there was no damage. I expect high quality from Amazon basic products - they usually are great - but these batteries are terrible!.

854 people found this helpful

Customer reviews

★★★★☆ 4.7 out of 5

412,923 global ratings



Generative and Discriminative Models

Generative vs. Discriminative Models

- **Generative model:** a model that calculates the probability of the input data itself

$$P(X)$$

stand-alone

$$P(X, y)$$

joint

- **Discriminative model:** a model that calculates the probability of the output given the input data

$$P(y | X)$$

conditional

Application to Text Classification

- **Generative text classification:** Learn a model of the joint $P(\textcolor{blue}{X}, \textcolor{red}{y})$, and find

$$\hat{y} = \arg \max_y P(X, y)$$

- **Discriminative text classification:** Learn a model of the conditional $P(\textcolor{red}{y} \mid \textcolor{blue}{X})$, and find

$$\hat{y} = \arg \max_y P(y|X)$$

Discriminative Text Classification

Why Discriminative Classifiers?

- Generative models are somewhat roundabout
→ spend lots of capacity modeling the input
- Discriminative models directly model the probability of the output → what we care about
- However, discriminative models **don't have an easy count-based decomposition!**

BOW Generative:

$$P(X, y) = P(y) \prod_{i=1}^{|X|} P(x_i|y) = \frac{c(y)}{\sum_{\tilde{y}} c(\tilde{y})} \prod_{i=1}^{|X|} \frac{c(x_i, y)}{\sum_{\tilde{x}} c(\tilde{x}, y)}$$

BOW Discriminative:

$$P(y|X) = \frac{P(y, X)}{P(X)} \quad ?? \quad \text{Sentence space is infinite!}$$

Discriminative Model Training

- Instead, define model that calculates probability directly based on parameters θ

$$P(y|X; \theta)$$

- Define a **loss function** that is lower if the model is better, such as **negative log likelihood** over training data

$$\mathcal{L}_{\text{train}}(\theta) = - \sum_{\langle X, y \rangle \in \mathcal{D}_{\text{train}}} \log P(y|X; \theta)$$

- And **optimize the parameters directly** to minimize loss

$$\hat{\theta} = \operatorname{argmin}_{\tilde{\theta}} \mathcal{L}_{\text{train}}(\tilde{\theta})$$

Logistic Regression

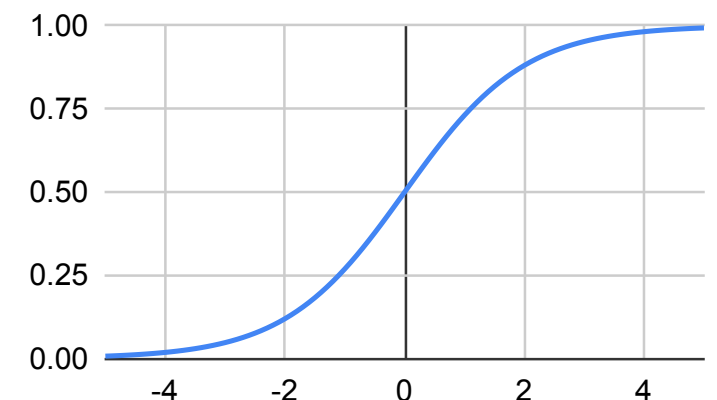
- For **binary classification** of positive/negative, first calculate score

$$s_{y|X} = \theta_{y|X} \cdot \underline{f(X)}$$

Learn a feature extractor

- Convert into a **probability**, e.g. using *sigmoid* function

$$P(y|X; \theta) = \text{sigmoid}(s_{y|X}) = \frac{1}{1 + e^{-s_{y|X}}}$$



- Learning**: maximize log likelihood of training data

$$\mathcal{L}_{\text{train}}(\theta) = \sum_{\langle X, y \rangle \sim \mathcal{D}_{\text{train}}} \log P(y|X; \theta), \quad \theta^* = \arg \max_{\theta} \mathcal{L}(\theta)$$

Multi-class Classification: Softmax

- Sigmoid can be used for binary decisions
- For multi-class decisions, calculate score for each class and use **softmax**

$$P(y|X; \theta) = \frac{e^{s_{y|X}}}{\sum_{\tilde{y}} e^{s_{\tilde{y}|X}}}$$

$$s = \begin{pmatrix} -3.2 \\ -2.9 \\ 1.0 \\ 2.2 \\ 0.6 \\ \dots \end{pmatrix} \longrightarrow p = \begin{pmatrix} 0.002 \\ 0.003 \\ 0.329 \\ 0.444 \\ 0.090 \\ \dots \end{pmatrix}$$

Gradient Descent

- Calculate the **gradient of the loss function** with respect to the parameters

$$\frac{\partial \mathcal{L}_{\text{train}}(\theta)}{\partial \theta}$$

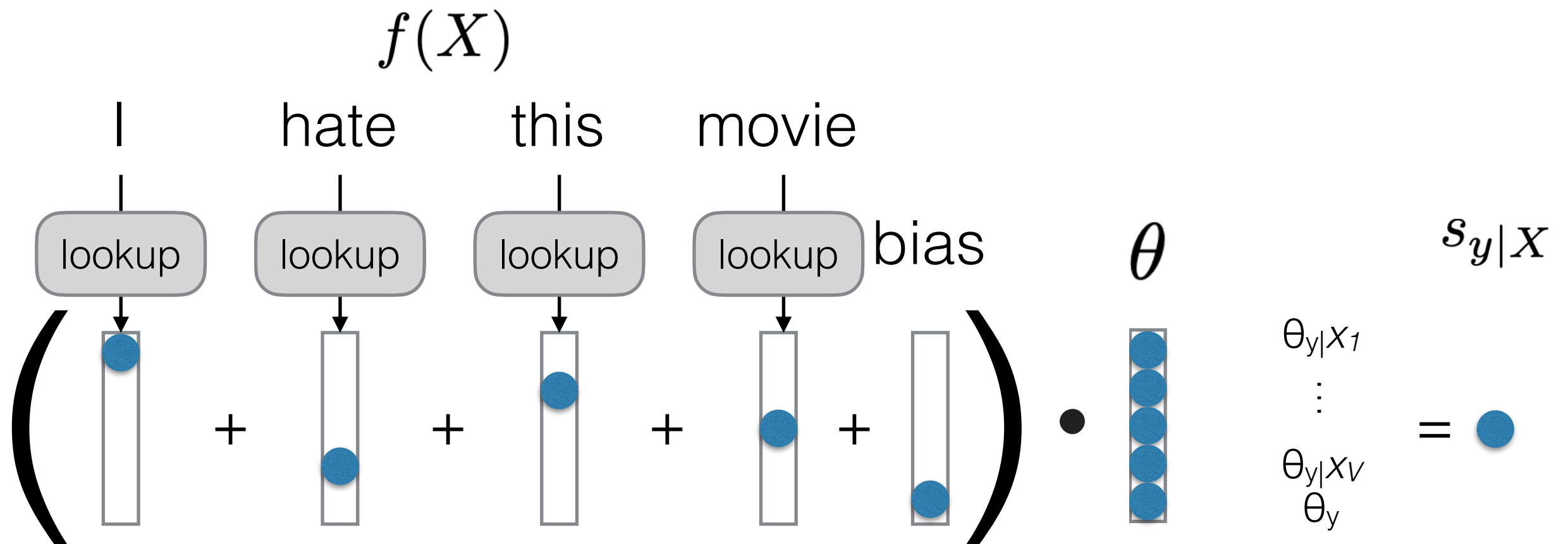
- How? Use the chain rule - more in later lectures.
- **Update** to move in a direction that decreases the loss

$$\theta \leftarrow \theta - \alpha \frac{\partial \mathcal{L}_{\text{train}}(\theta)}{\partial \theta}$$

- α is a **learning rate** dictating speed of movement
- This is the *first-order* gradient descent
- Others, e.g. Newton's method, consider *second-order* (curvature) information and converge more quickly

BOW Discriminative Model

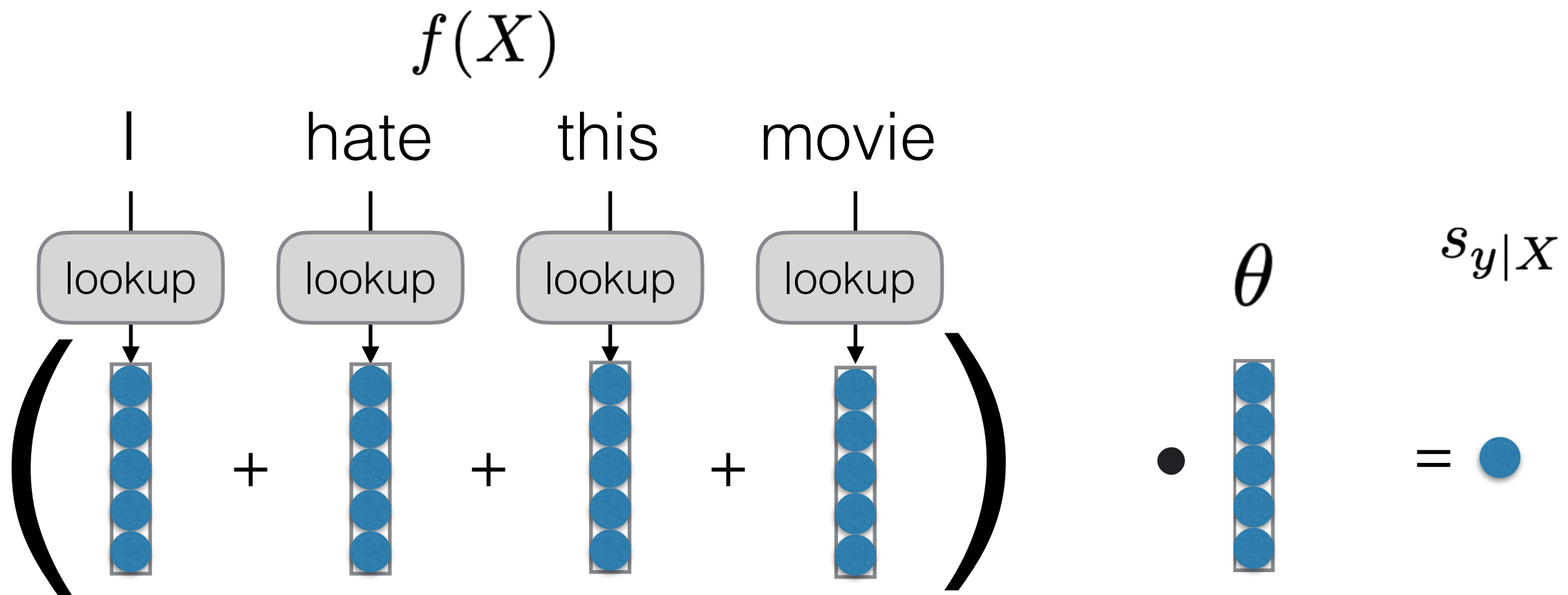
- Use BOW representations for $f(X)$



$$s_{y|X} = \theta_y + \sum_{i=1}^{|X|} \theta_{y|x_i}$$

CBOW Discriminative Model

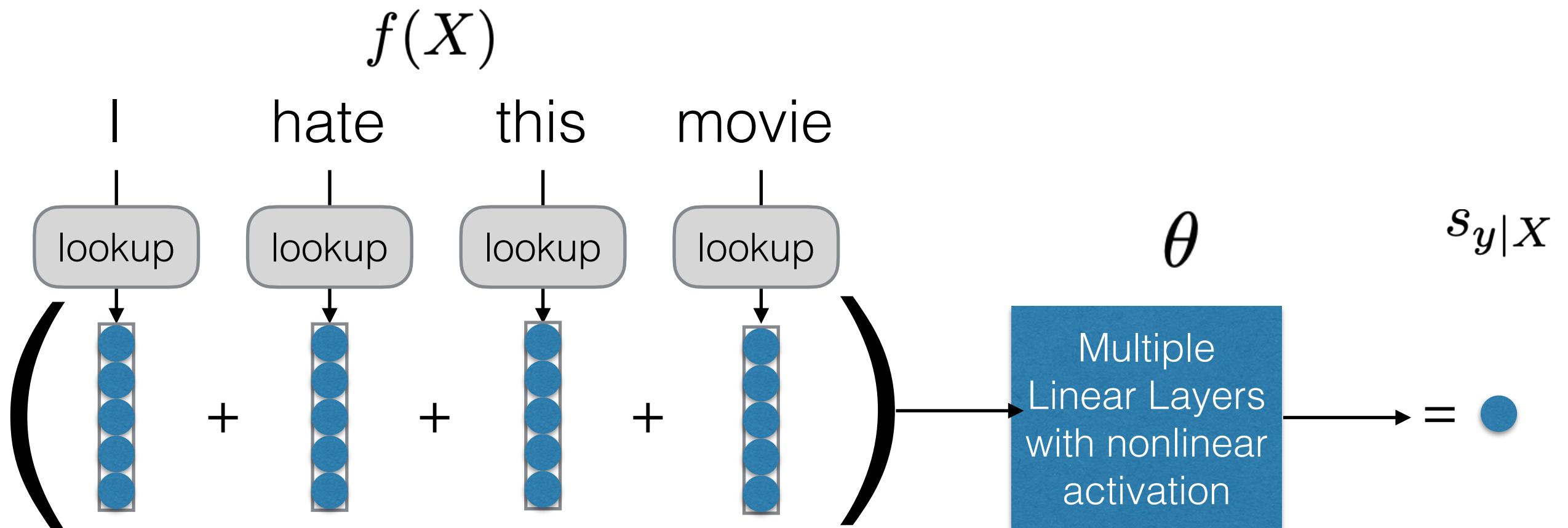
- Use CBOW representations for encoding a sentence $f(X)$



$$s_{y|X} = \theta_y + \sum_{i=1}^{|X|} \theta_{y|x_i}$$

CBOW Discriminative Model

- Use CBOW representations for encoding a sentence $f(X)$



$$s_{y|X} = \theta_y + \sum_{i=1}^{|X|} \theta_{y|x_i}$$

Covert Scores to Probabilities

- Using Softmax function

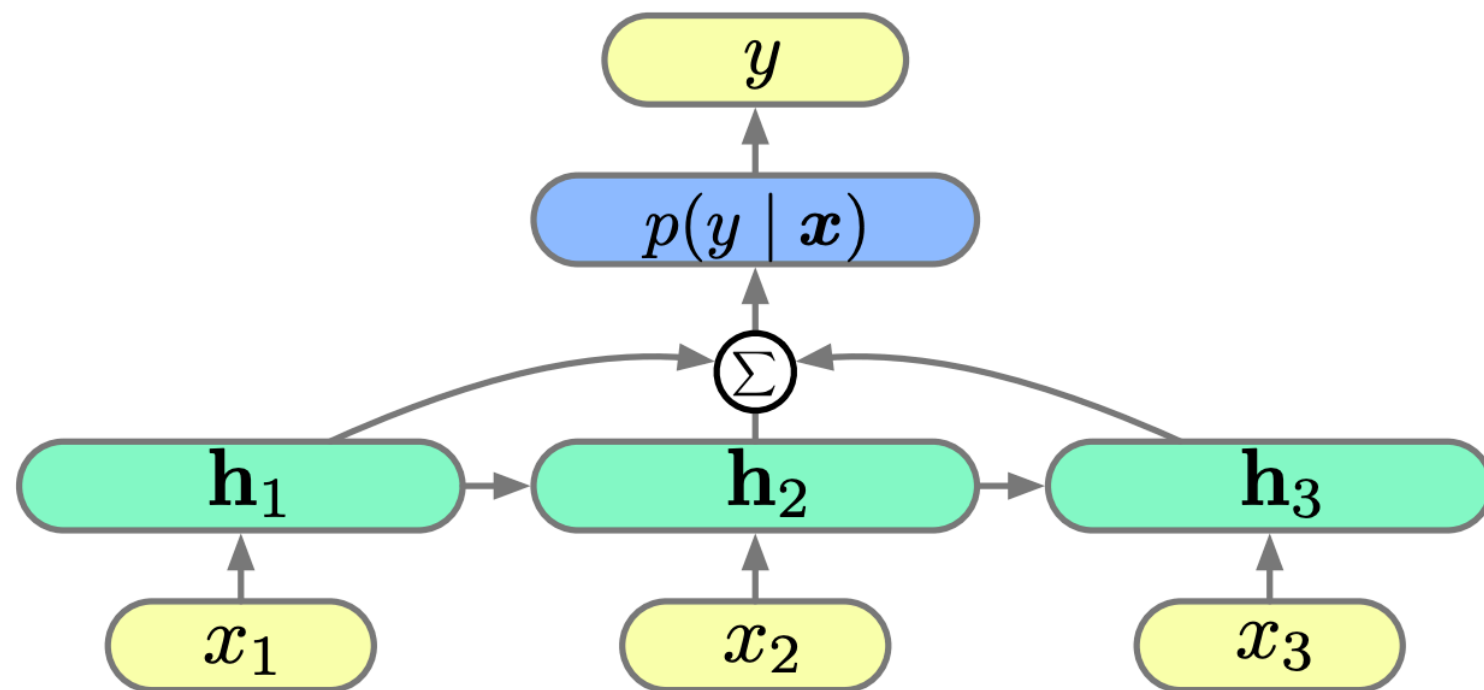
$$P(y|X; \theta) = \text{SoftMax}(s_{y|X}) = \frac{e^{s_{y|X}}}{\sum_{y' \in Y} e^{s_{y'|X}}}$$

- Optimize by gradient descent with an regularization

$$\mathcal{L}_{\text{train}}(\theta) = \sum_{\langle X, y \rangle \sim \mathcal{D}_{\text{train}}} \log P(y|X; \theta) + \lambda \|\theta\|^2$$

Neural Network Discriminative Model

- Use neural network (e.g., LSTM) to learn features $f(X)$



$$s_{y|X} = \text{NN}(X)$$

Any Questions?