

Neural NLP Tutorial

Language Modeling: Models of $P(\text{text})$

Are These Sentences

OK?

- Jane went to the store.
- store to Jane went the.
- Jane went store.
- Jane goed to the store.
- The store went to Jane.

Engineering Solutions

- Jane went to the store.
- store to Jane went the.
- Jane went store.
- Jane goed to the store.
- The store went to Jane.

} Create a grammar
of the language

} morphology and
Consider exceptions
Semantic
} preference
categories,
s

What Can we Do w/

LMs?

- Score sentences:

Jane went to the store . →

high store to Jane went the .

→ low

(same as calculating loss for training)

- Generate sentences:


while didn't choose end-of-sentence symbol:

calculate probability

sample a new word from the probability

distribution

Calculating the Probability of a Sentence

$$P(X) = \prod_{i=1}^I P(x_i \mid x_1, \dots, x_{i-1})$$


The diagram shows the formula $P(X) = \prod_{i=1}^I P(x_i \mid x_1, \dots, x_{i-1})$. Below the formula, there are two horizontal lines. The first line is red and is positioned under the x_i term. A red arrow points from the text "Next Word" below it to this red line. The second line is blue and is positioned under the x_1, \dots, x_{i-1} term. A blue arrow points from the text "Context" below it to this blue line.

This is a classification problem over the next word!

$$P(x_i \mid x_1, \dots, x_{i-1})$$

How do we do
this?!?!

Count-based Language Models

- Count up the frequency and divide:

$$P_{ML}(x_i \mid x_{i-n+1}, \dots, x_{i-1}) := \frac{c(x_{i-n+1}, \dots, x_i)}{c(x_{i-n+1}, \dots, x_{i-1})}$$

<https://books.google.com/ngrams>

Any problem?

Problems w/ Count-based Models

- Cannot share strength among **similar**

words	bought	a	car	she	bought	a	bicycle
she	purchased	a		she	purchased	a	
car				bicycle			

- Cannot handle **long-distance**

for tennis class	he wanted to buy his own	racquet
for programming class	he wanted to buy his own	computer

A Slight Simplification:
Sentence
Classification Models
of $P(\text{label} \mid \text{text})$

An Example Prediction Problem: Sentence Classification

I hate this movie

I love this movie

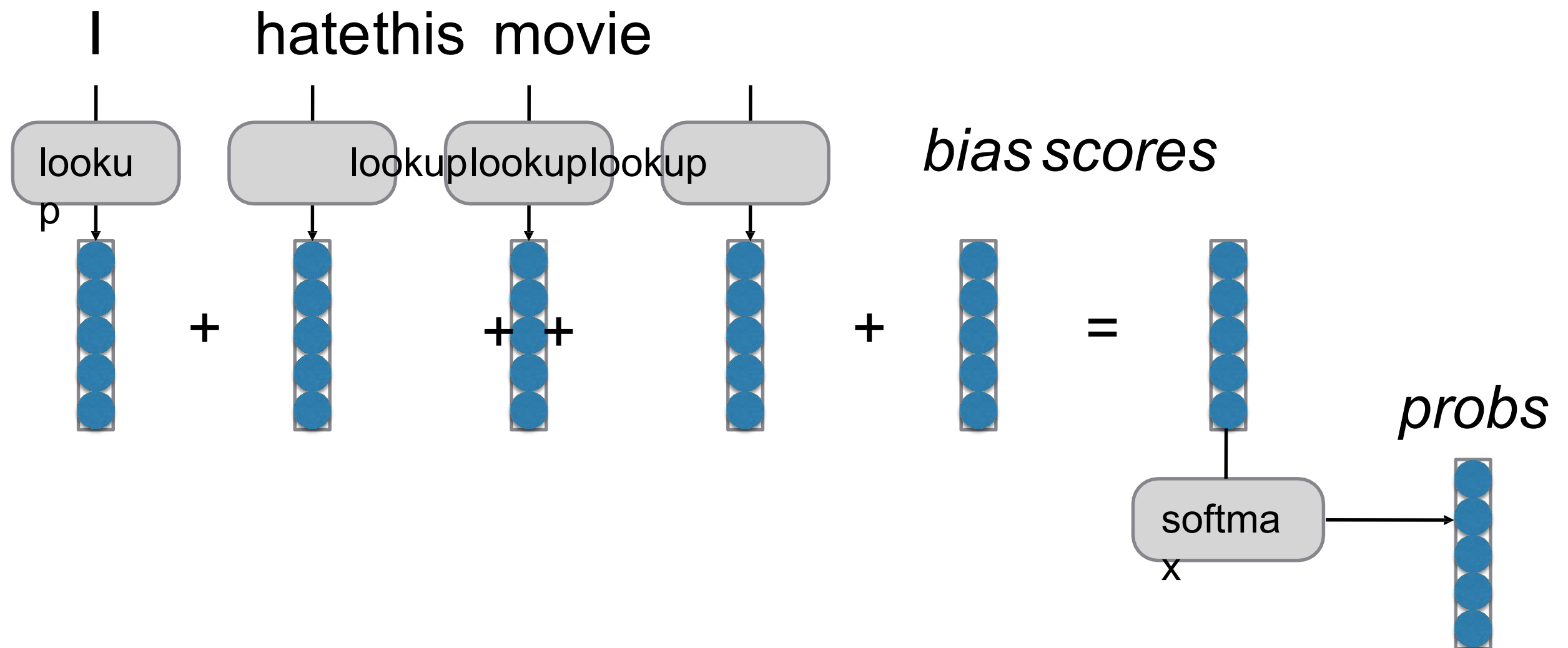
very
good
neutral
bad
very
bad
very
good
good
neutral
bad
very bad

Word Representation (task-specific)

- Each word has its own 5 elements corresponding to [very good, good, neutral, bad, very bad]
- “hate” will have a high value for “very bad”, etc.



A First Try: Bag of Words (BOW)



Adversarial Examples

I don't love this
movie

very
good
good
neutral
bad
very bad

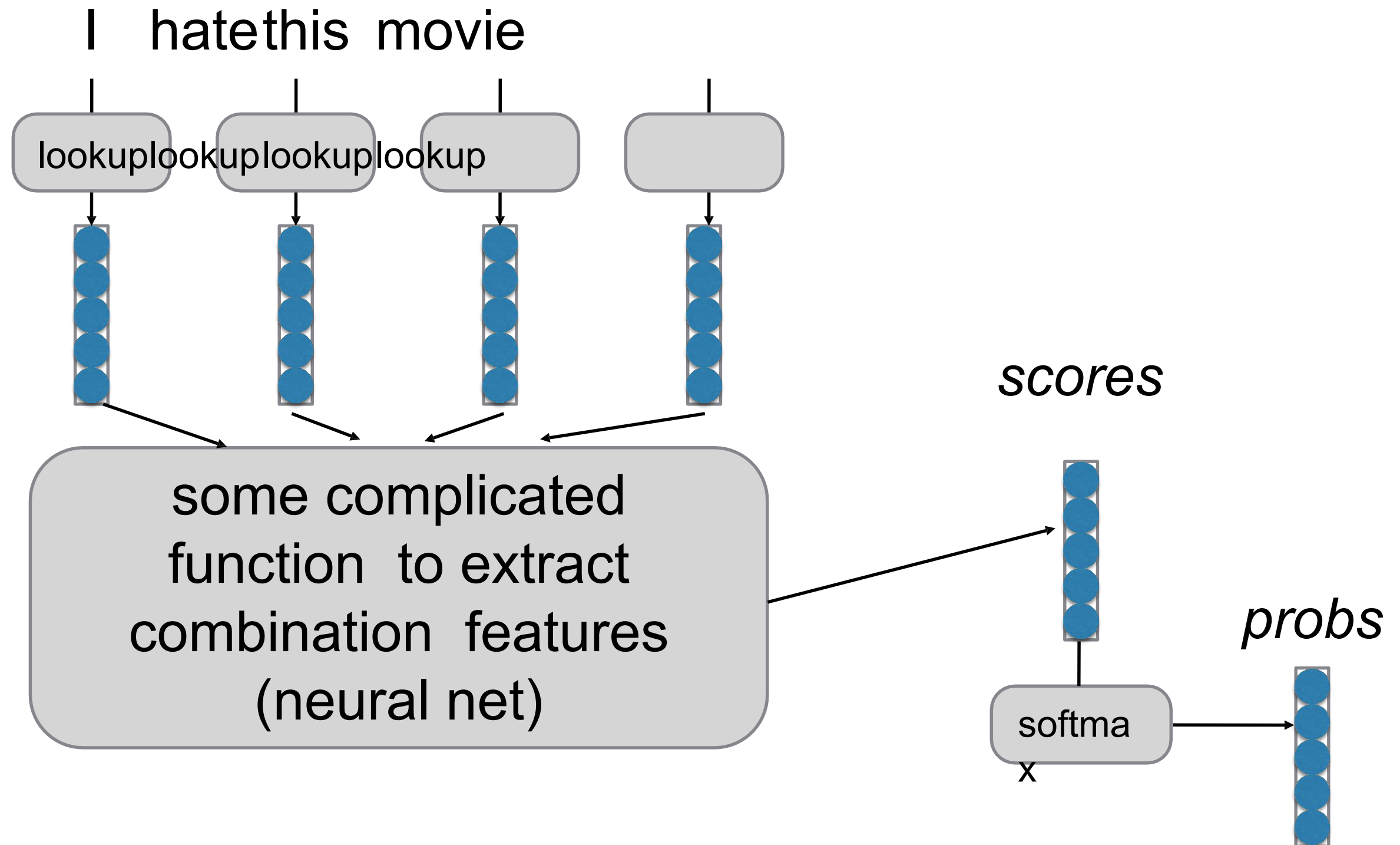
There's nothing I
don't
love about this movie

very
good
good
neutral
bad
very bad

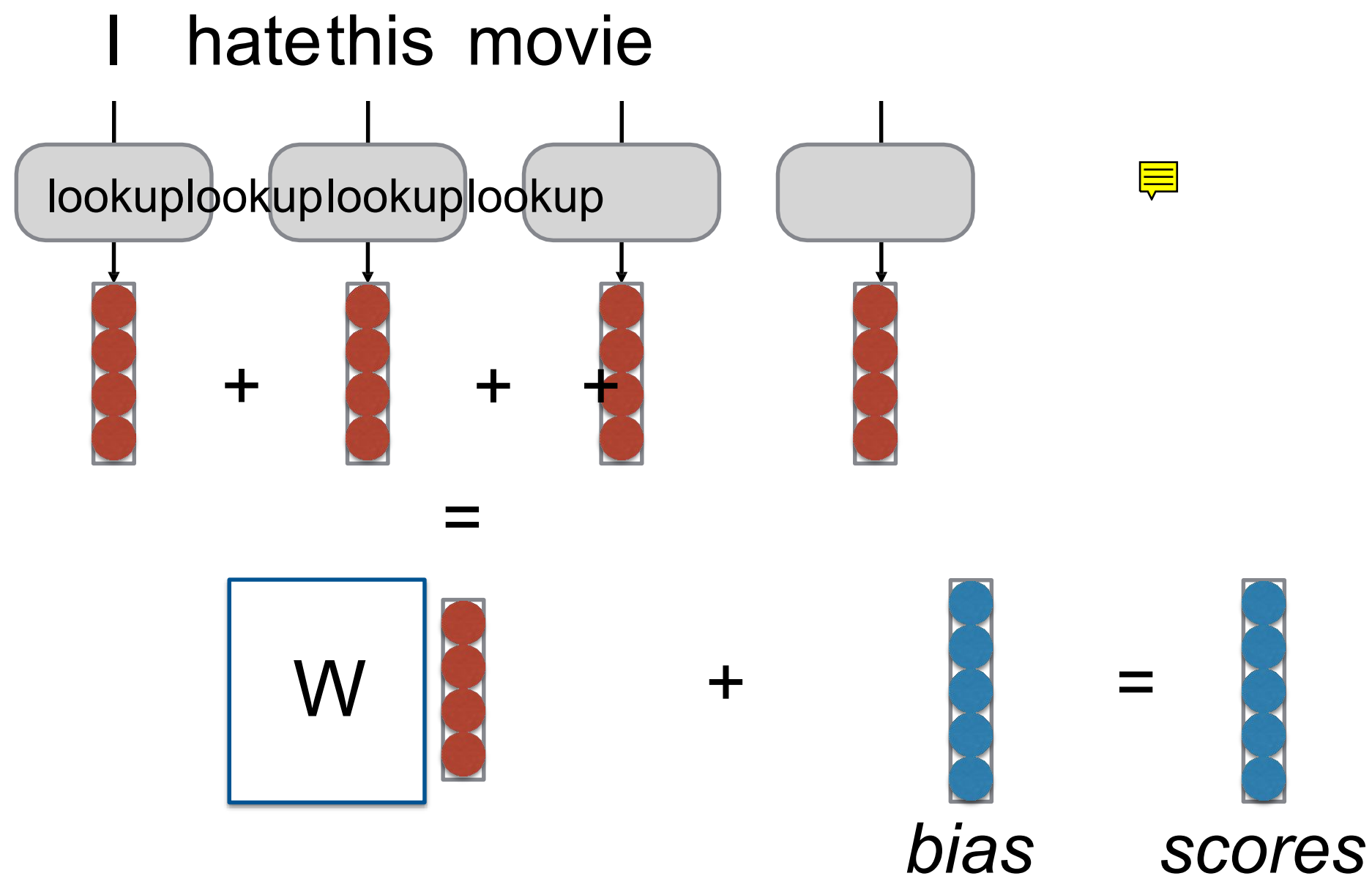
Combination Features

- Does it contain “don’t” and “love”?
- Does it contain “don’t”, “i”, “love”, and “nothing”?

Idea 1: Neural Networks



Idea 2: Continuous Bag of Words (CBOW)



What do Our Vectors Represent?

- Each vector has “features” (e.g. is this an animate object? is this a positive word, etc.)
- We sum these features, then use these to make predictions
- Still no combination features: only the expressive power of a linear model, but dimension reduced

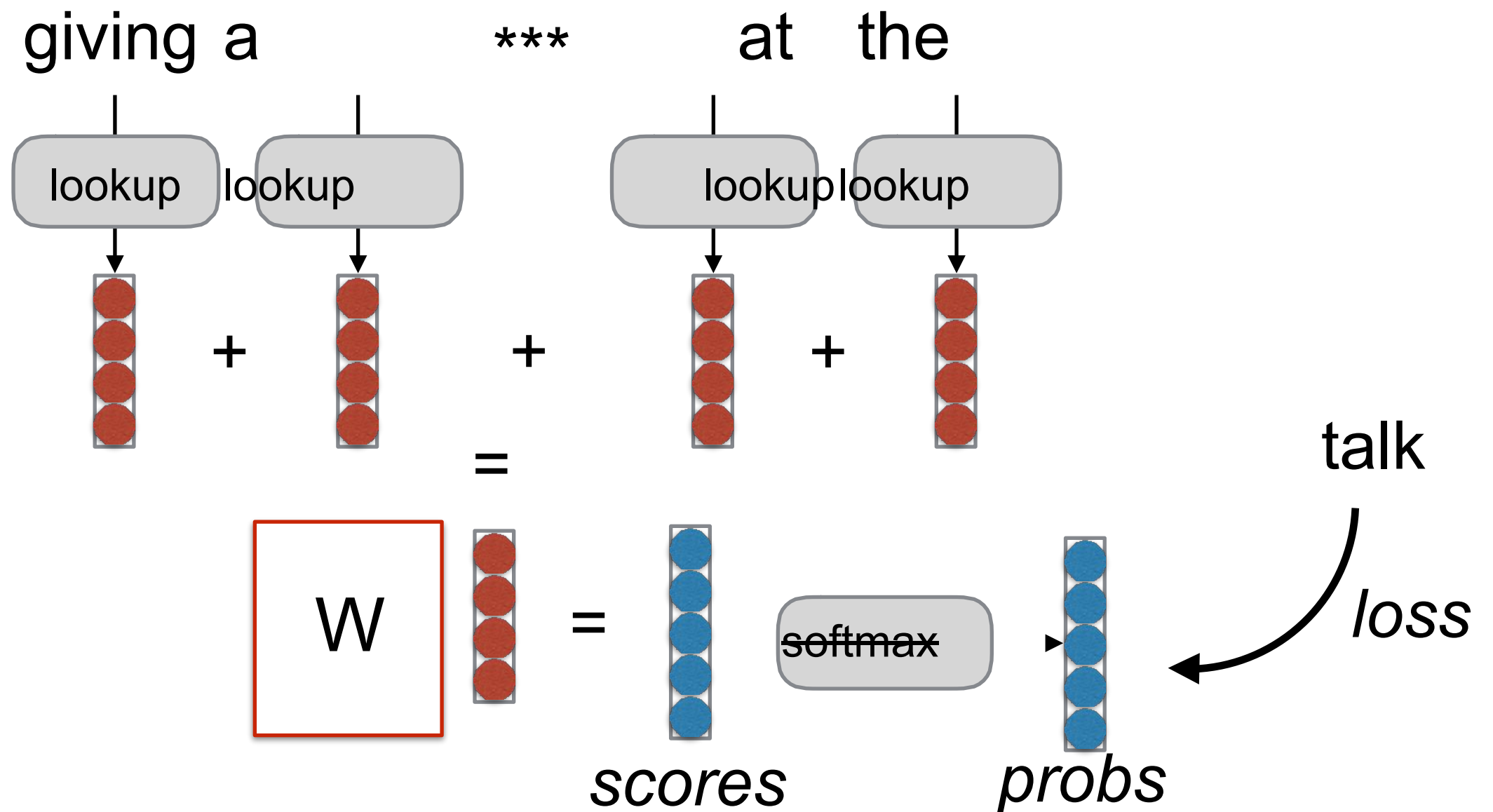
Revisiting LM

LM as a classification problem into words

CBO

(Mikolov et al. 2013)

- Predict word based on sum of surrounding embeddings



Softma

- Convert scores ~~into~~ probabilities by taking the exponent and normalizing (softmax)

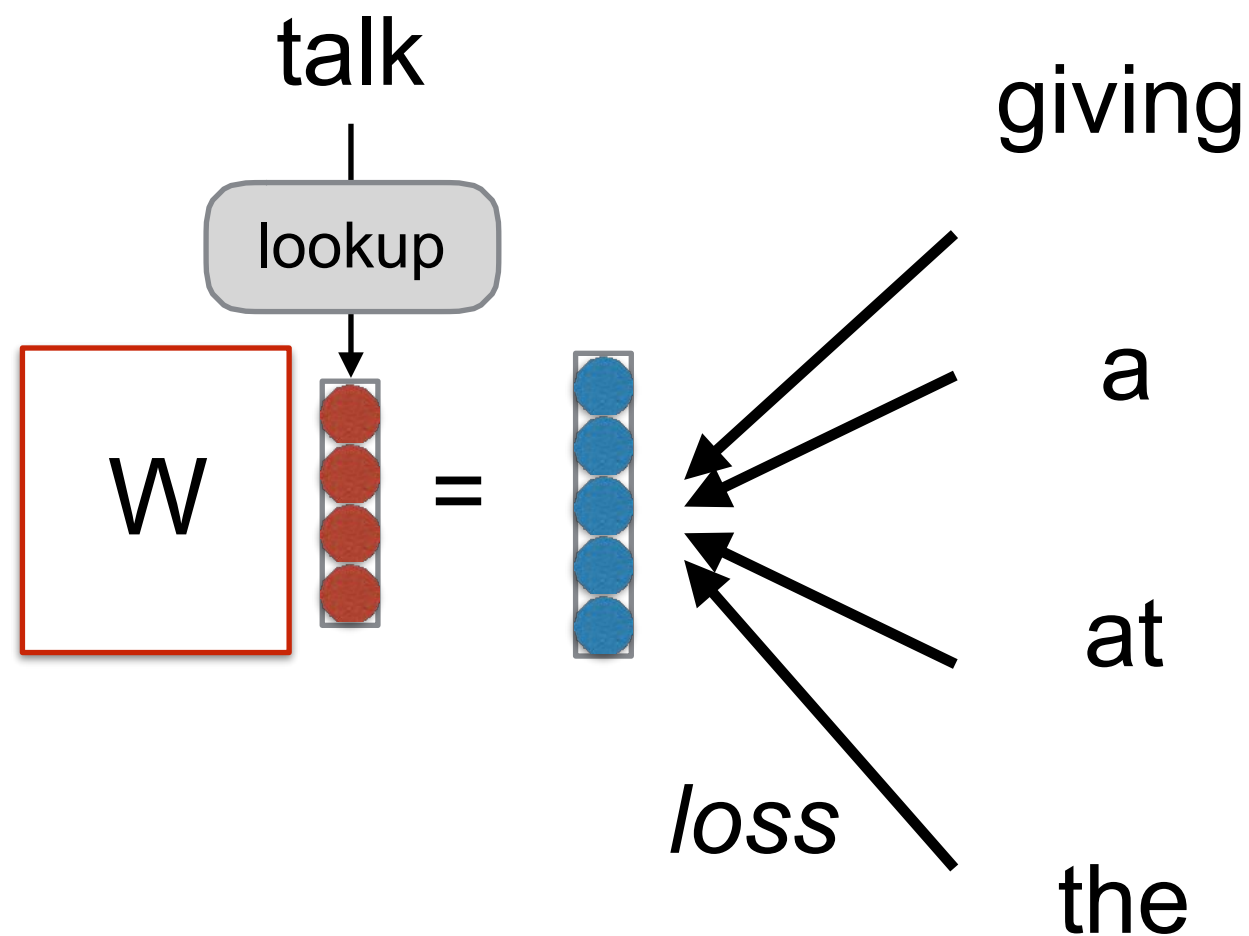
$$P(x_i \mid x_{i-n+1}^{i-1}) = \frac{e^{s(x_i \mid x_{i-n+1}^{i-1})}}{\sum_{\tilde{x}_i} e^{s(\tilde{x}_i \mid x_{i-n+1}^{i-1})}}$$

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

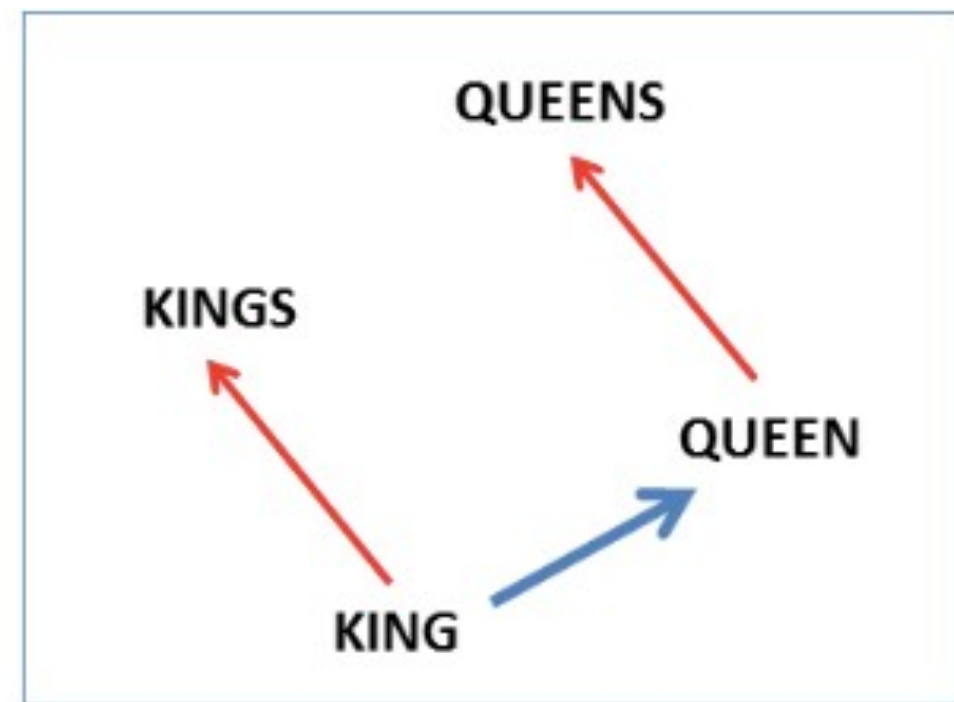
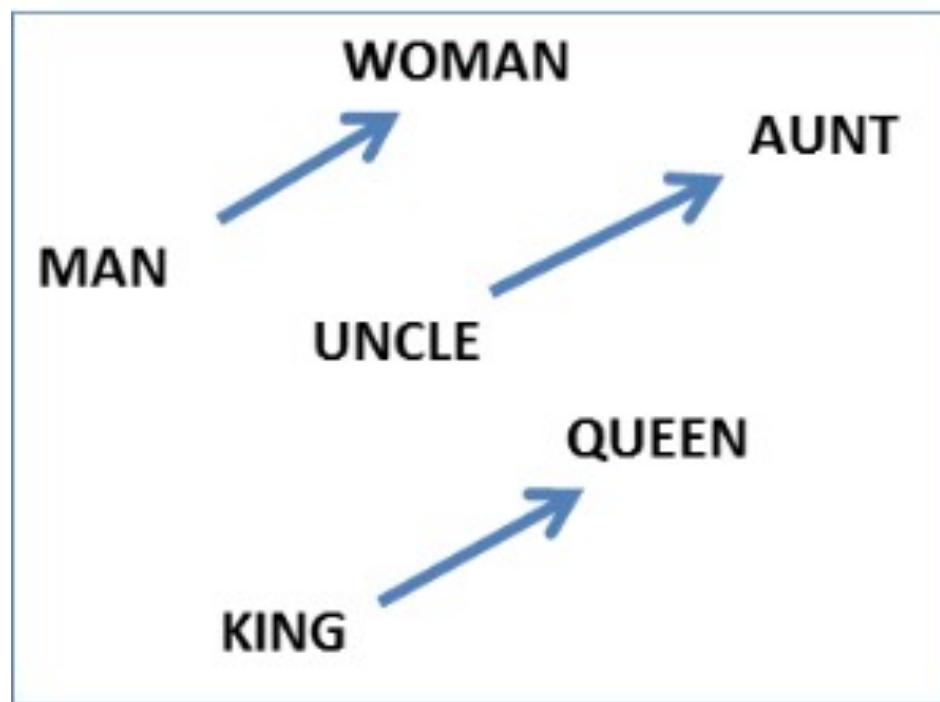
Skip-gram

(Mikolov et al. 2013)

- Predict each word in the context given the word



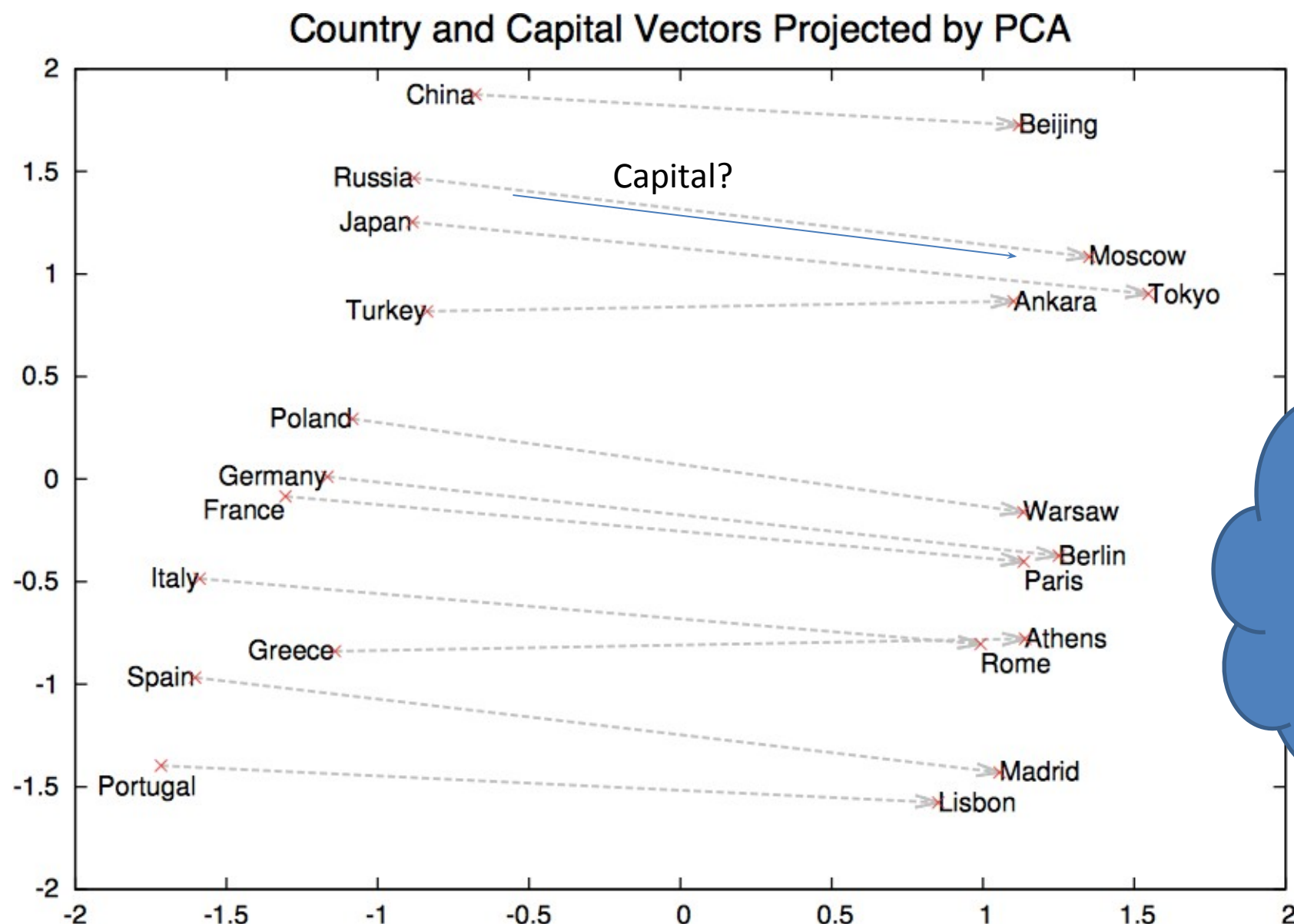
Word Embeddings



- “What is the female equivalent of king?” is not easily accessible in many traditional resources

Visualization of Embeddings

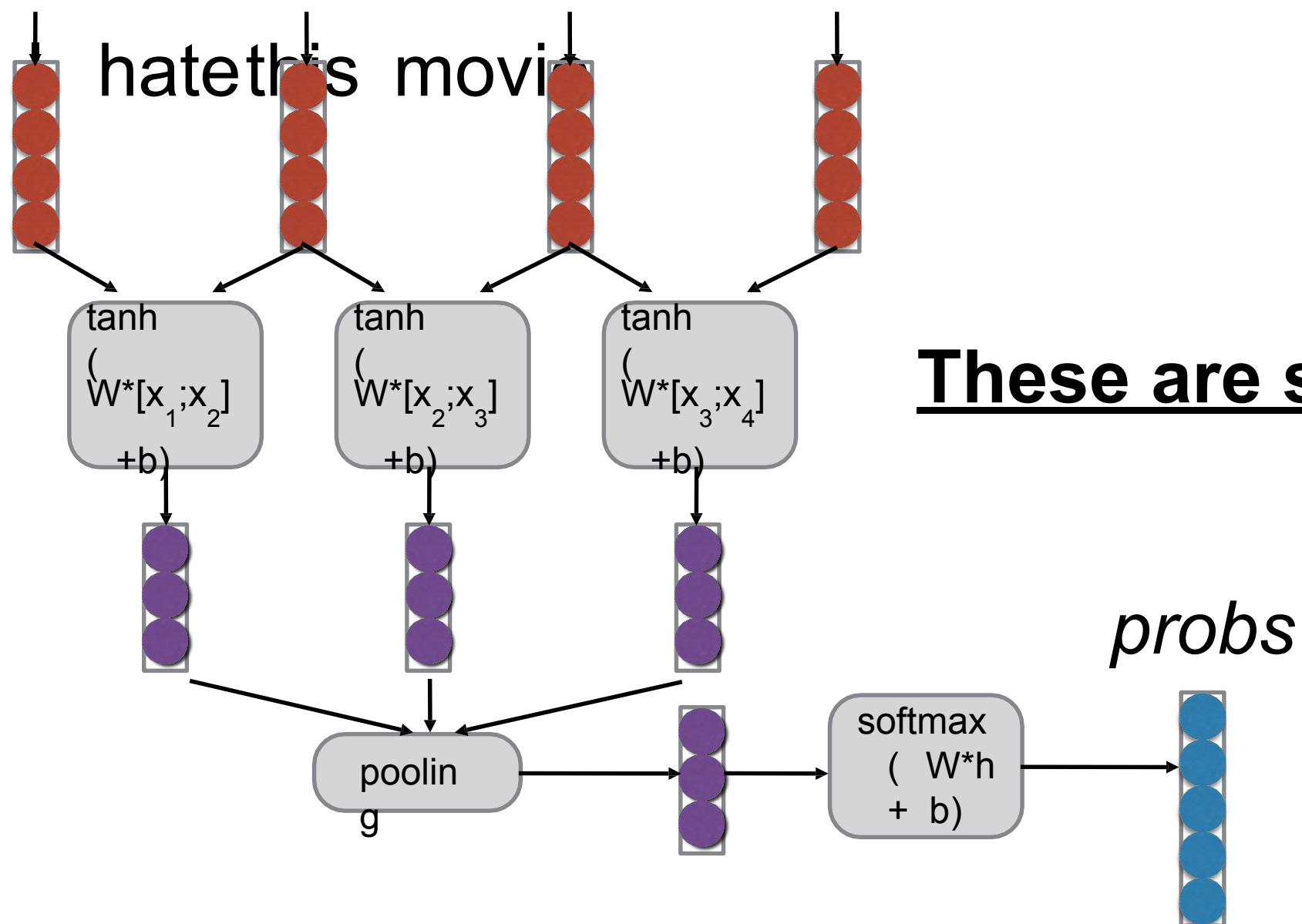
- Reduce high-dimensional embeddings into 2/3D for visualization (e.g. Mikolov et al. 2013)



(Korea, ?, Seoul) =
True
If $h + ? = t$

Time Delay Neural Networks

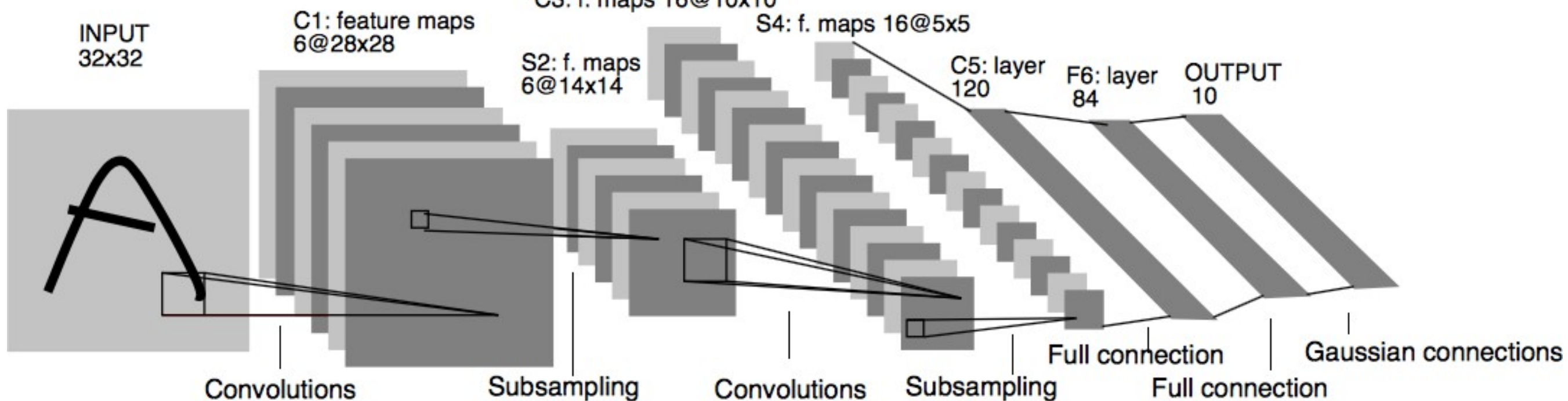
(Waibel et al.
1989)



These are soft 2-grams!

Convolutional Networks

(LeCun et al 1997)



Parameter extraction performs a 2D sweep, not 1D

Poolin

- Calculate some reduction function feature-wise
- **Max pooling:** “Did you see this feature anywhere in the range?” (most common)
- **Average pooling:** “How prevalent is this feature over the entire range”
- **k-Max, Dynamic, ...**

Weaknesses of CNNs

- CNNs are great for short-distance feature extractors
- But don't have holistic view of the sentence to capture long-distance dependencies

Long-distance Dependencies in Language

- Agreement in number, gender, etc.

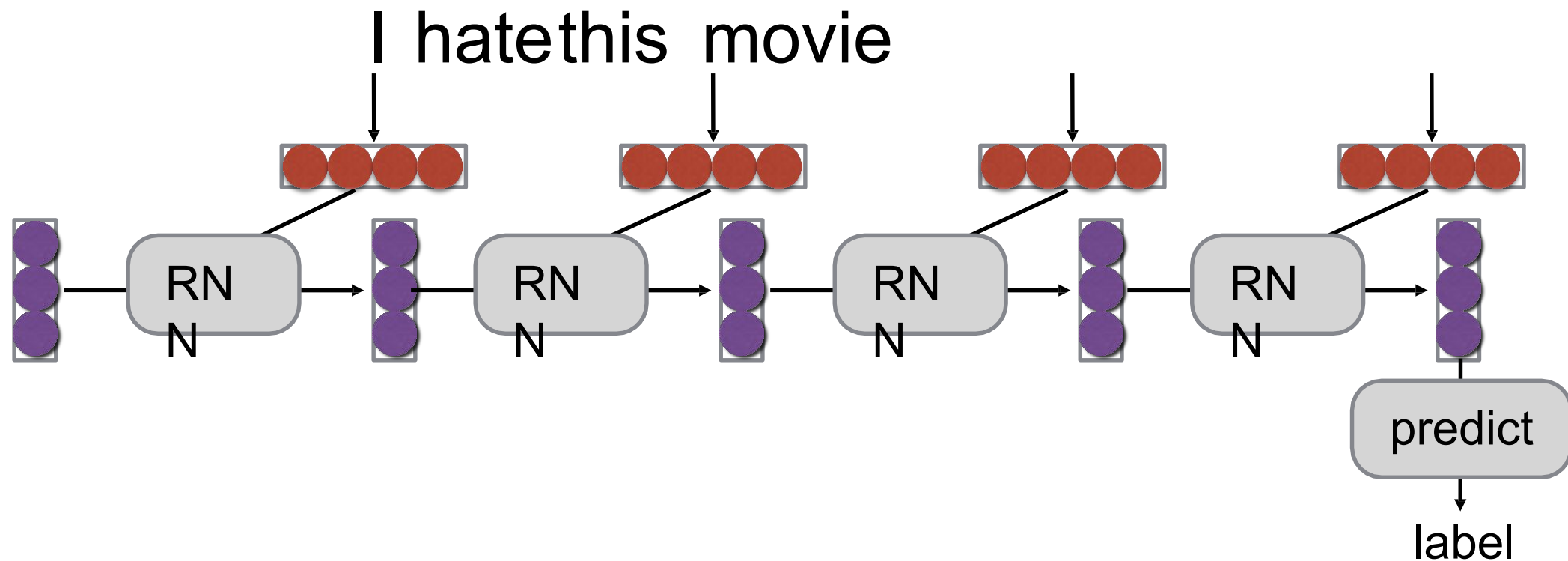
He does not have very much confidence in **himself**. **She** does not have very much confidence in **herself**.

- Selectional preference

The **reign** has lasted as long as the life of the **queen**. The **rain** has lasted as long as the life of the **clouds**.

Remember History w/ RNNs

- What does processing a sequence look like?



Weaknesses of RNNs

- Indirect passing of information, credit assignment more difficult
 - Made better by LSTMs/GRUs/etc. but not perfect
- Can be slow, due to incremental processing