

DS 207: Introduction to Natural Language Processing



Text Classification II

Danish Pruthi

Slides courtesy: Graham Neubig

Example: topic classification

- Sports: "Kohli scores another remarkable hundred"
- Politics: "Minister announces new metro plans ahead of elections"
- Entertainment: "A sleeper hit, 12th fail, is praised by many Bollywood actors"
- Finance: "Stocks for Indian delivery startups plummet"

Generative Naive Bayes

ith example

$$P(X^{(i)}, y^{(i)}) = P(y^{(i)}) P(X^{(i)} | y^{(i)})$$

$$P(X | y) = P(w_1, w_2, w_3, \dots w_t | y)$$

$$P(X | y) = \prod_i^t P(w_i | y)$$

Assumed independent

$$P(w_i = "Kohli" | y = "sports")$$

Estimating parameters

$$P(w_i = \text{"Kohli"} \mid y = \text{"sports"})$$

$$= \frac{\text{count}(w_i = \text{Kohli} \in y = \text{sports})}{\sum_{w \in |V|} \text{count}(w \in y = \text{sports})}$$

Add- α smoothing (or Laplace smoothing)

$$P(w_i = \text{"Kohli"} \mid y = \text{"finance"})$$

$$= \frac{\text{count}(w_i = \text{Kohli} \in y = \text{finance}) + \alpha}{\sum_{w \in |V|} (\text{count}(w \in y = \text{finance}) + \alpha)}$$

Evaluation → 2 classes

TP FN

- Accuracy $\frac{TP + TN}{P + N}$

- Precision = $TP / (\text{Predicted}) P$

- Recall = $TP / (\text{Actual}) P$

- F1 score

$\sqrt{P \cdot R}$

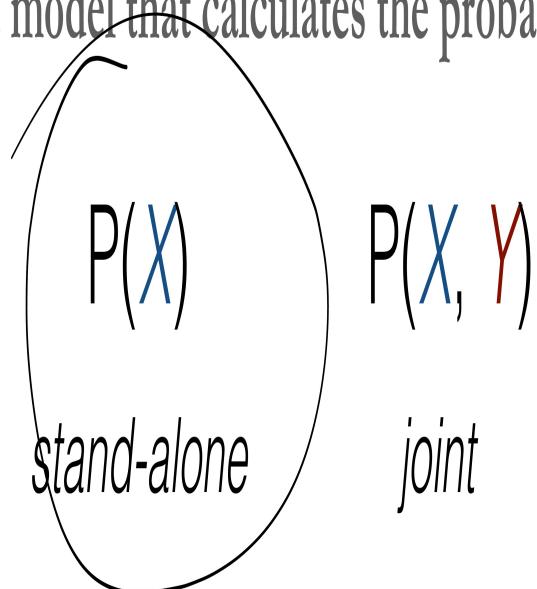
		Predicted condition	
		Total population $= P + N$	Predicted positive (PP)
		Positive (P) [a]	True positive (TP), hit [b]
Actual condition		Negative (N) [d]	False positive (FP), false alarm, overestimation
			True negative (TN), correct rejection [e]

What metrics are best suited

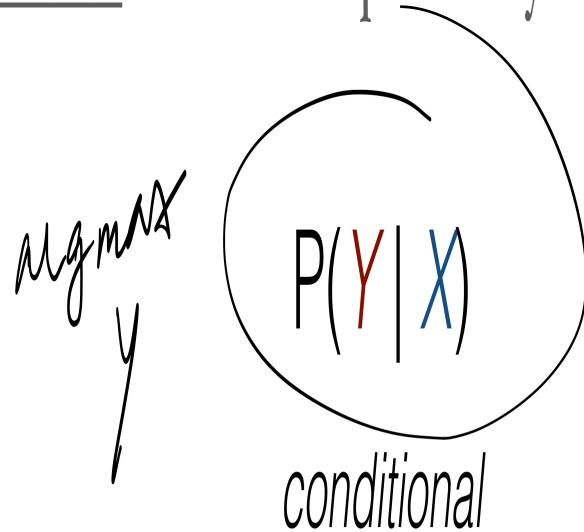
- Diagnosing rare type of cancer → recall
- Criminal punishment ✗ ↗ ↓
- Detecting spam
- Recruitment/Filtering based on text in the CVs
- Recommending products/songs/movies ↗ ↘~~~~~

Generative vs Discriminative models

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



- Discriminative model calculates the probability of a class (or trait) given the data



Rule based systems: is the headline sports or entertainment?

Kohli scores a ton. India Win.

- **Feature extraction:** Extract the salient features for making the predictions from text

1. Does the headline contain the word "Kohli"
2. Does the headline contain the word "win"
3. Does the headline contain the word "cricket"
4. Does the headline contain the word "actor"
5. Does the headline contain the word "show"
6. Does the headline contain the word "blockbuster"

Rule based systems: is the headline sports or entertainment?

- **Feature extraction:** Extract the salient features for making the predictions from text

{	1. Does the headline contain the word "Kohli"	<input checked="" type="checkbox"/>	1	Kohli scores a ton, India win.
	2. Does the headline contain the word "win"	<input type="checkbox"/>	1	
	3. Does the headline contain the word "cricket" \times	$f(x) =$	0	
	4. Does the headline contain the word "actor" \times		0	
	5. Does the headline contain the word "show" \checkmark		0	
	6. Does the headline contain the word "blockbuster" \times		0	

Rule based systems: is the headline sports or entertainment?

$$f(w) \cdot w$$

- **Feature extraction:** Extract the salient features for making the predictions from text

1. Does the headline contain the word "Kohli"

1	1
1	+1
0	+1
0	-1
0	-1
0	-1

2. Does the headline contain the word "win"

3. Does the headline contain the word "cricket"

$$f(x) =$$

4. Does the headline contain the word "actor"

5. Does the headline contain the word "show"

6. Does the headline contain the word "blockbuster"

A three step process for making predictions

- Feature extraction: Extract the salient features for making the decision from text

$$h = f(x) \quad \text{test} \quad R^d \quad 4 \times d$$

- Score calculation: Calculate a score for one or more possibilities

$$s = w \cdot h \quad \text{binary} \quad s = w \cdot h \quad \text{multi-class} \quad R \quad R^4 \quad R$$

- Decision function: Choose one of the several possibilities

$$\hat{y} = \text{decide}(s) \quad HAT \quad S ? t \quad y^* \quad t \quad -$$

Another discriminative classifier

Another discriminative classifier

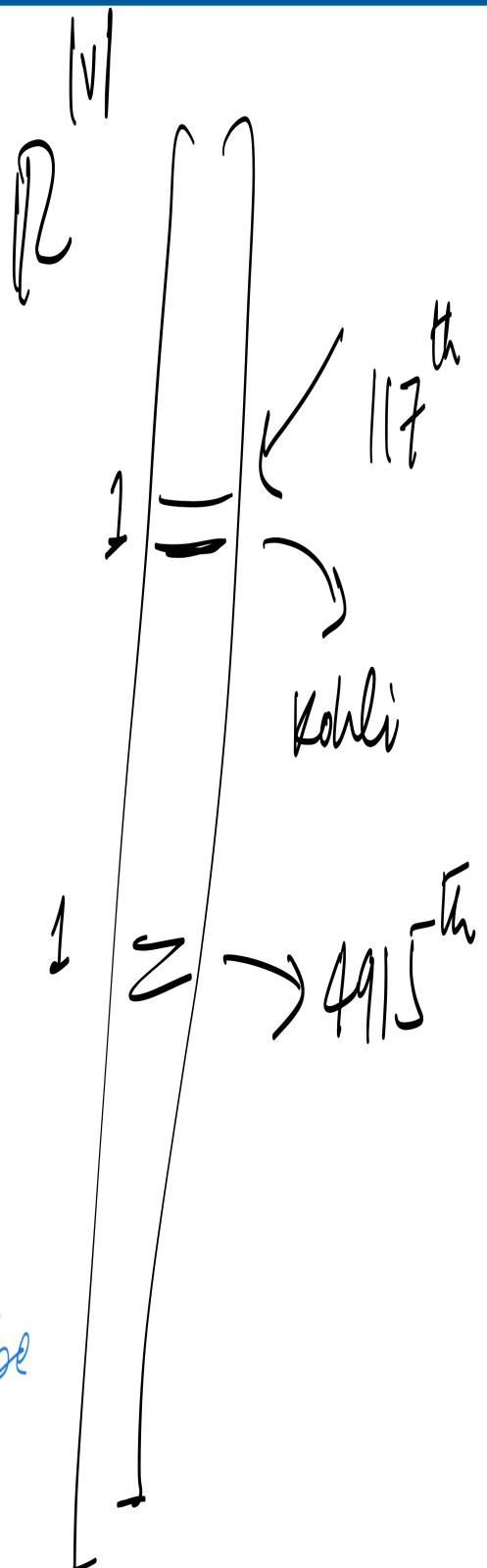
- Feature Extraction: $h = f(x)$

One hot vectors ("bag of words")

- Count every vector 201
- Order ignored.

- "count of words" instead

"bag of words" could be



Another discriminative classifier

- Feature Extraction: $h = f(x) \rightarrow \mathbb{R}^{|V|}$

One hot vectors ("bag of words")

- Score Calculation: binary or multi-class

$$s = \mathbf{w} \cdot \mathbf{h}$$

$s = \mathbf{W} \mathbf{h}$ $\mathbf{X} \in \mathbb{R}^{N \times |V|}$

Another discriminative classifier

- **Feature Extraction:** $h = f(x)$

One hot vectors ("bag of words")

- **Score Calculation:** binary or multi-class

$$s = \mathbf{w} \cdot \mathbf{h} \quad s = \mathbf{W}\mathbf{h}$$

- **Decision:** Convert to a probability:

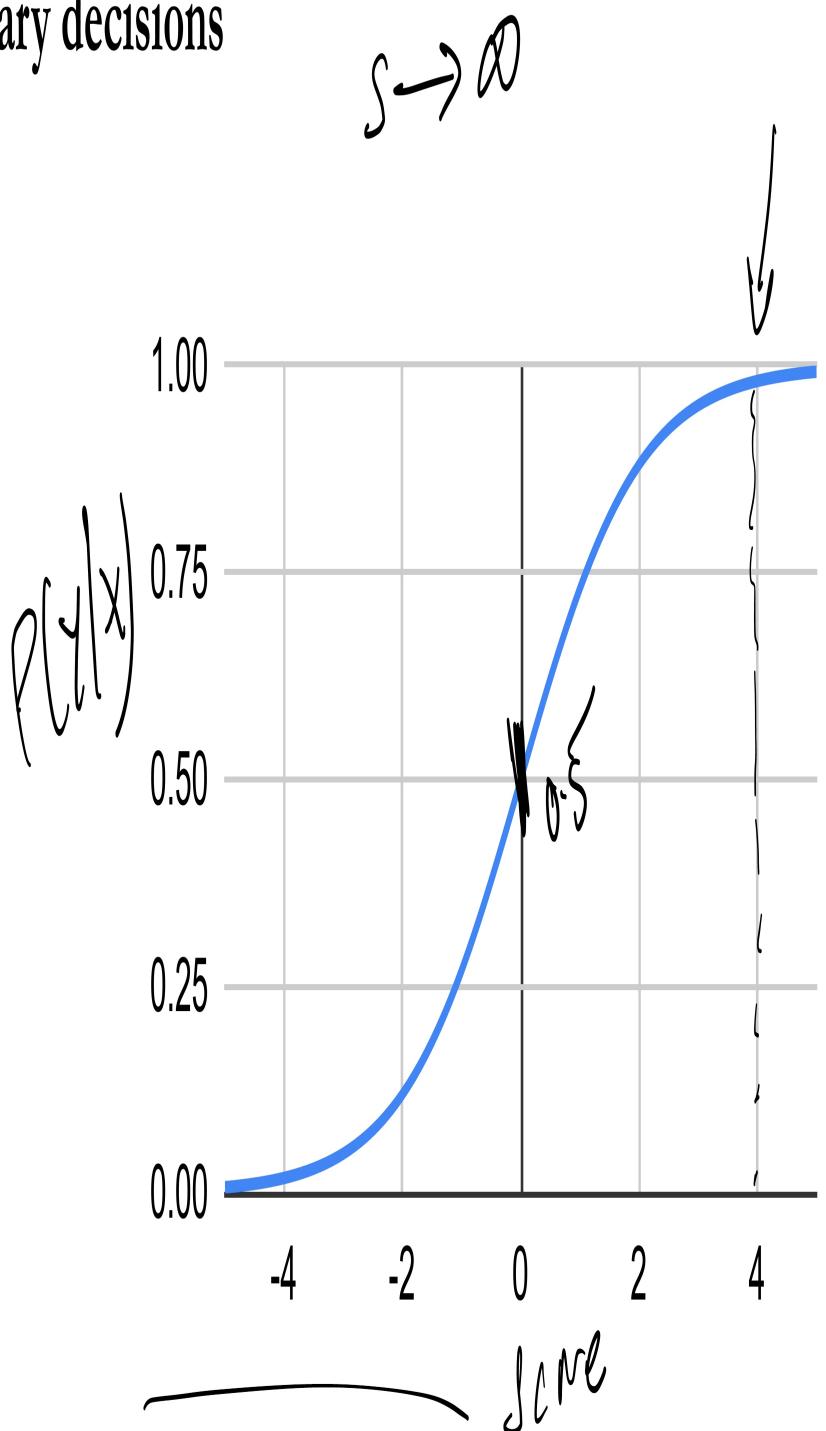
$$P(y|x) = \sigma(s) \text{ or } \text{softmax}(\mathbf{s})$$



Sigmoid function

- Sigmoid can be used for binary decisions

$$\sigma(s) = \frac{1}{1 + e^{-s}}$$



Softmax

- Softmax is used for multi-label classification

$$\text{softmax}(s) = \frac{e^{s_i}}{\sum_i^d e^{s_i}}$$

↑ to malce v
positive.

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

why not $\sigma(w_2 + b)$?

Softmax

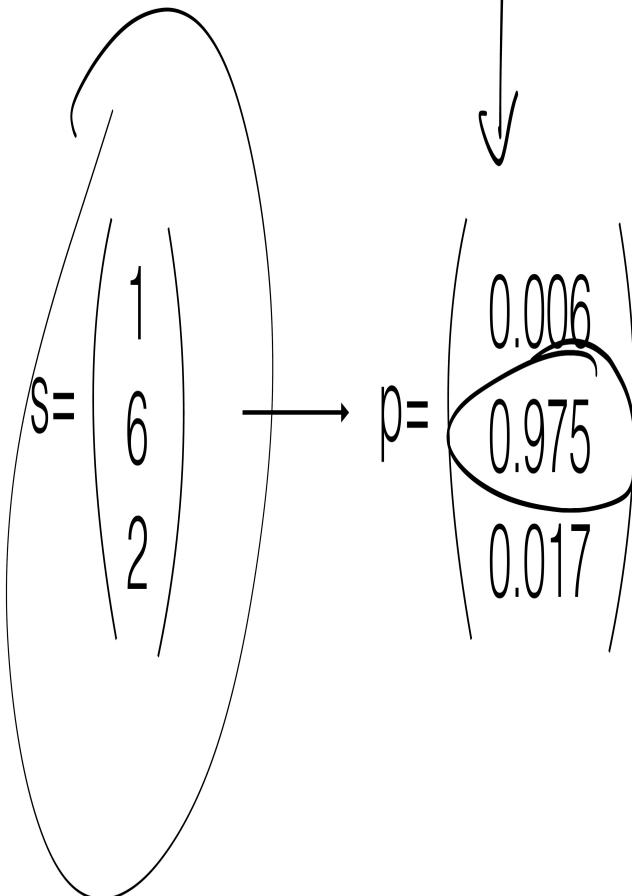
$$p(y|x)$$

- Softmax is used for multi-label classification

↓
3 class

$$\text{softmax}(s) = \frac{e^{s_i}}{\sum_i^d e^{s_i}}$$

↑
points to s_{\max}



Recap: Discriminative models

- Define a model that calculates probability directly based on parameters W

$$\begin{aligned} P(Y|X; W) \\ = \sigma(W f(X)) \\ = \text{Softmax}(W \cdot f(X)) \end{aligned}$$

Computing loss (or error/cost) function

Computing loss (or error/cost) function

Define a loss function (a function which is lower for better models):

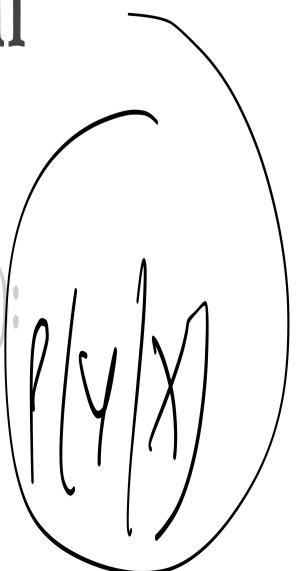
Computing loss (or error/cost) function

Define a loss function (a function which is lower for better models):

Computing loss (or error/cost) function

Define a loss function (a function which is lower for better models):

$$f(y, x) = \underbrace{\dots}_{z}$$



$$L(W) = - \sum_{x, y \in D} \log P(y|x; W) \quad [\text{negative log likelihood}]$$

- training

Train to find W what values of W gives the minimum loss?

Computing loss (or error/cost) function

Define a loss function (a function which is lower for better models):

$$L(W) = - \sum_{x, y \in D} \log P(y|x; W) \quad [\text{negative log likelihood}]$$

Computing loss (or error/cost) function

Define a loss function (a function which is lower for better models):

$$(1 - y_i) \log \left(1 - P(y=1|x_i; w) \right)$$

$$L(W) = - \sum_{x, y \in D} \log P(y|x; W) \quad [\text{negative log likelihood}]$$

algm

$$W\left[L(W)\right] = - \sum_i^n \left(y_i \log P(y=1|x_i; W) + (1 - y_i) \log(P(y=0|x_i; W)) \right)$$

[binary cross entropy]

$y_i = 1/0$

$y_i = 0$

$y_i = 1$ class

-ve class

Learn parameters through gradient descent

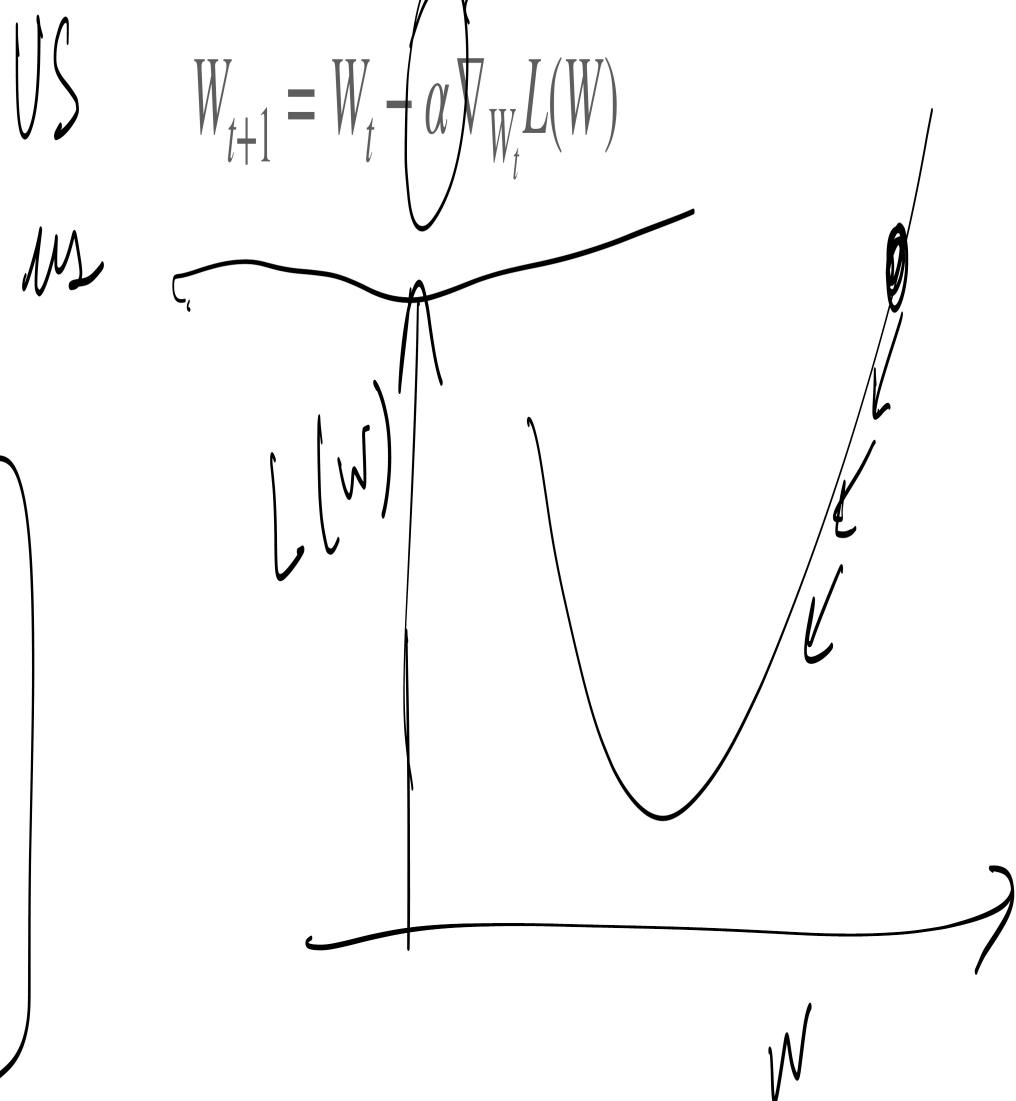
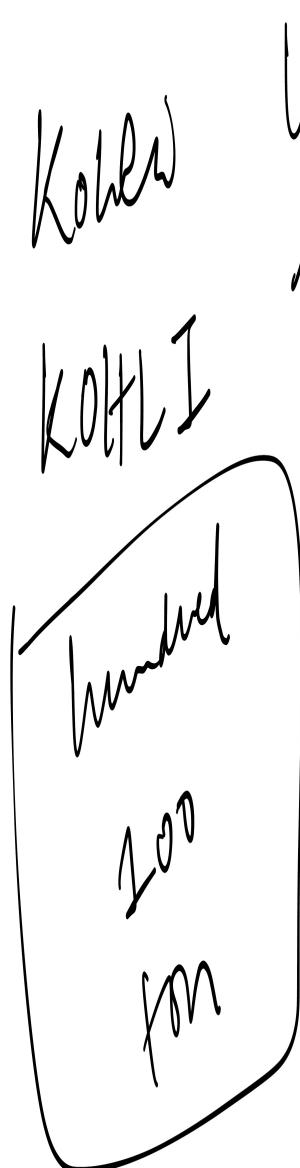
A big win for India!

GREAT! great

- Compute the gradient of the loss function with respect to the parameters

learning rate

- Keep updating the parameters to move in a direction that decreases the loss

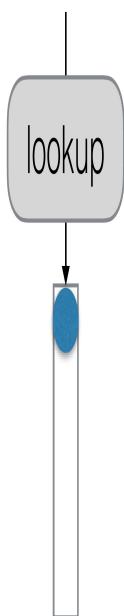


Computational Graph View

Kohli scored a hundred

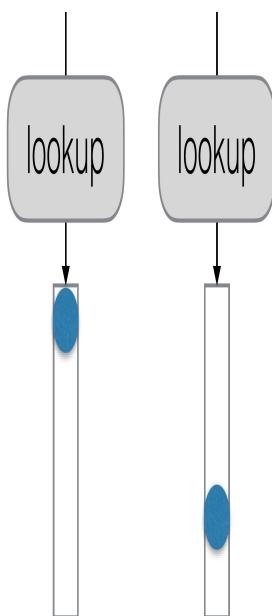
Computational Graph View

Kohli scored a hundred



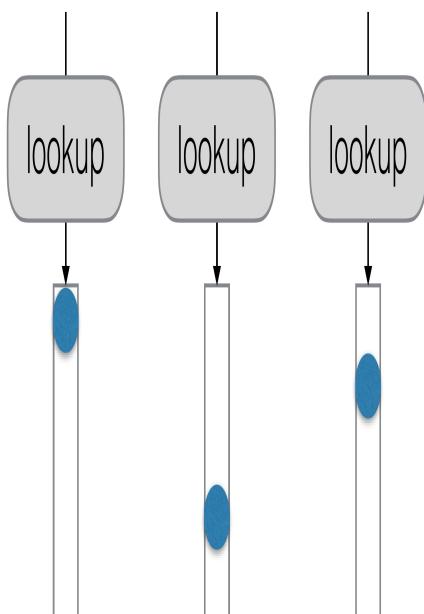
Computational Graph View

Kohli scored a hundred



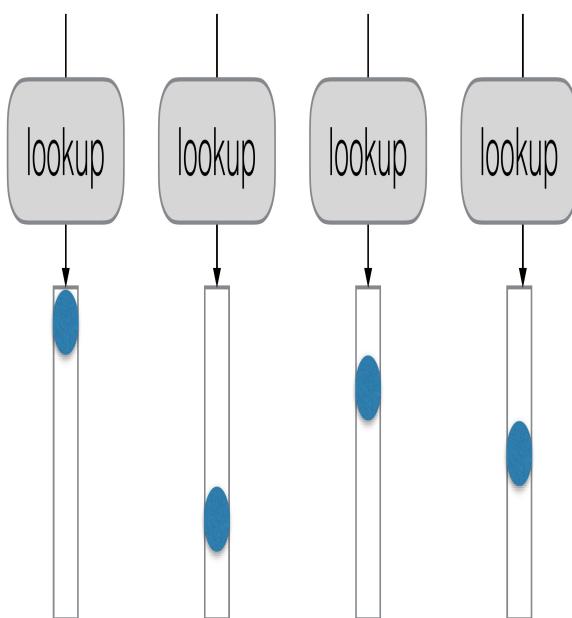
Computational Graph View

Kohli scored a hundred



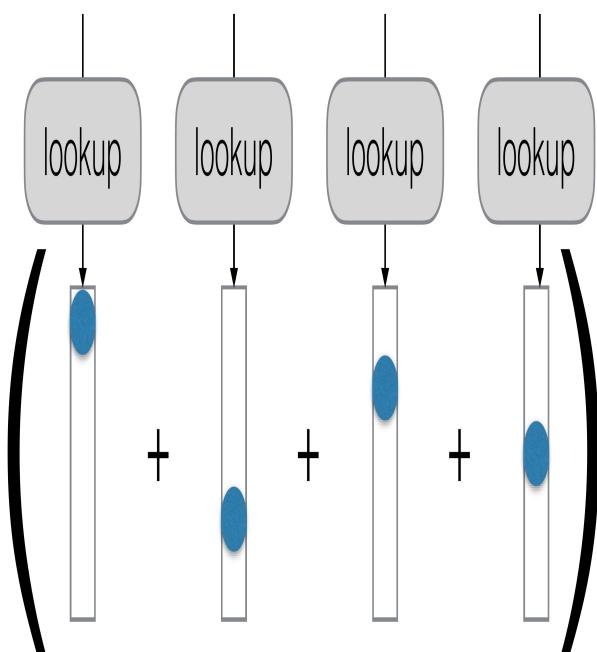
Computational Graph View

Kohli scored a hundred



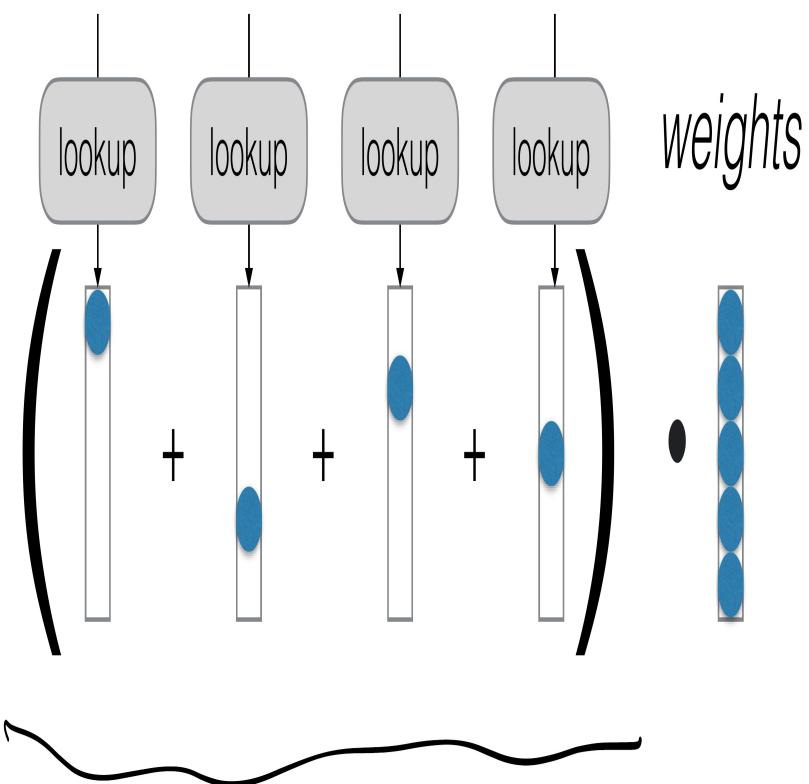
Computational Graph View

Kohli scored a hundred



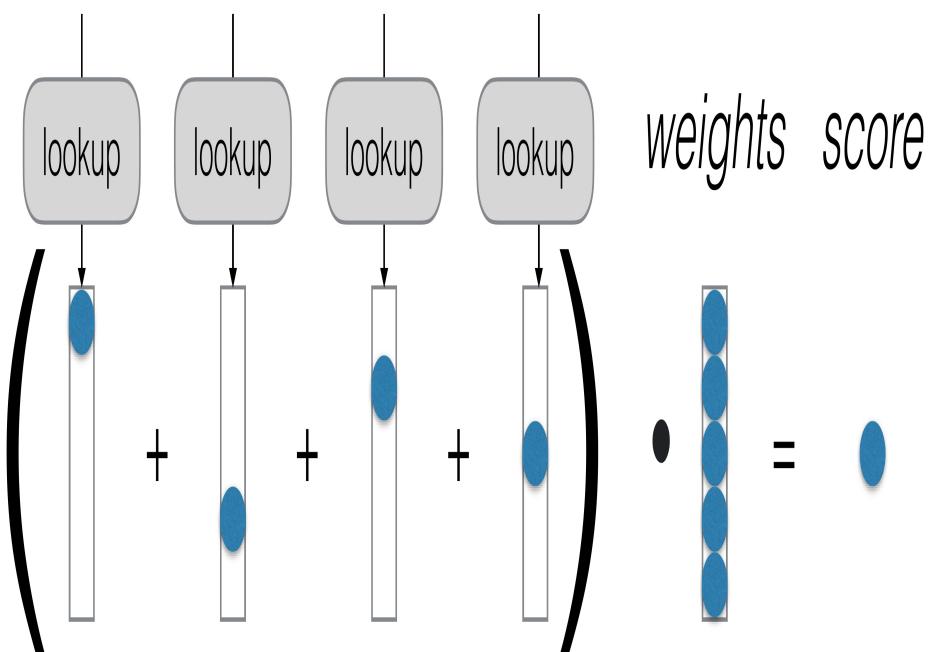
Computational Graph View

Kohli scored a hundred



Computational Graph View

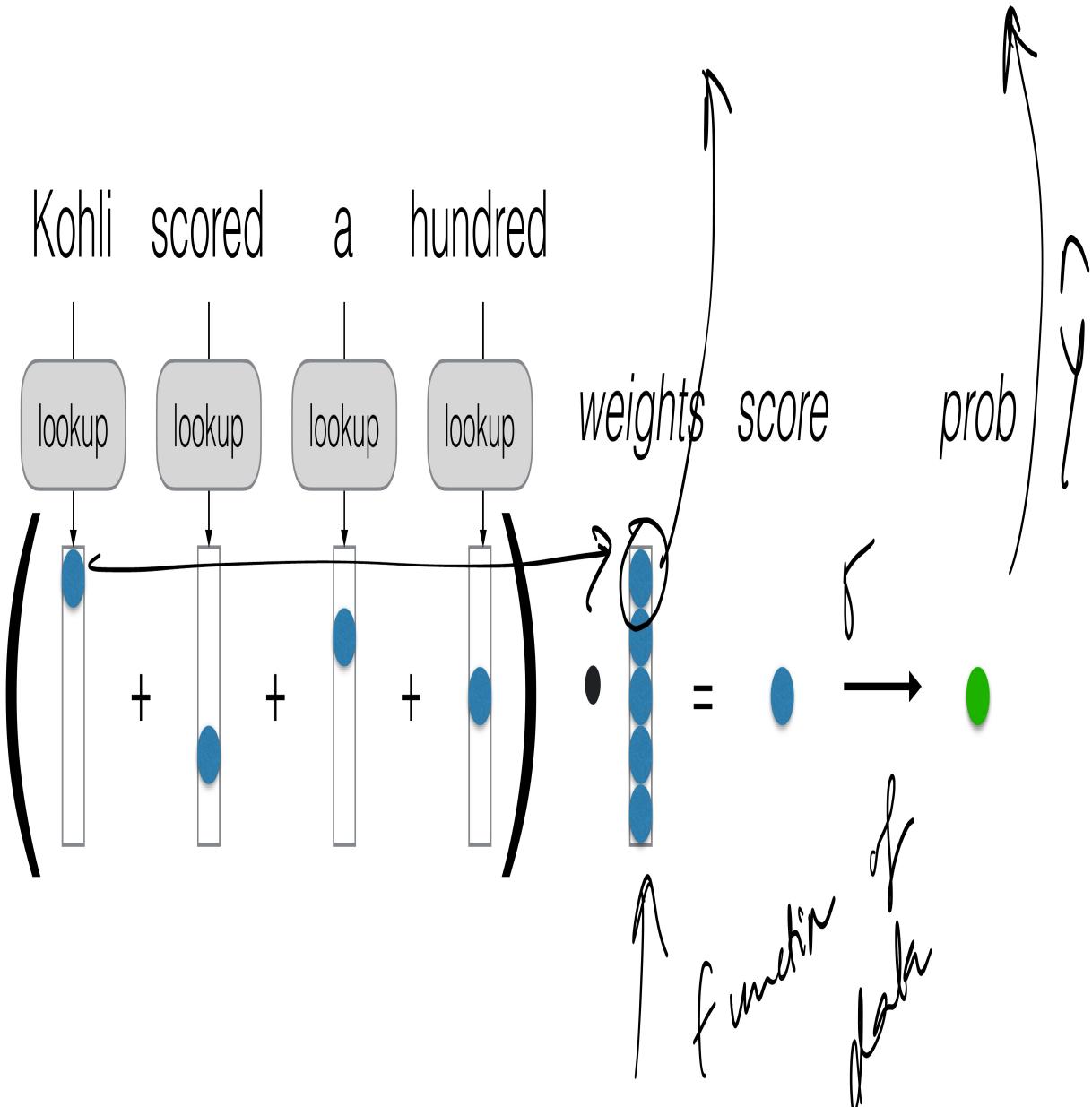
Kohli scored a hundred



Computational Graph View

$$U(y^*, y)$$

Kohli scored a hundred



"WordNet" "JPN"

Questions?

Next class: Word2vec