# DS 207: Introduction to Natural Language Processing

## Text Classification

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

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

$$P(X \mid y) = P(w_1, w_2, w_3, \ldots w_t \mid y)$$

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

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

# 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-$\alpha$ smoothing (or Laplace smoothing)

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

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

# Evaluation

Accuracy $= \dfrac{TP + TN}{P + N}$

- **Accuracy**

- **Precision = TP / (Predicted) P**

- **Recall     = TP / (Actual) P**

- **F1 score**

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

Actual condition

# What metrics are best suited

- Diagnosing rare type of cancer

- 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

$$P(X) \qquad P(X,\ Y)$$

*stand-alone*      *joint*

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

$$P(Y \mid X)$$

*conditional*

# 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"

  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"
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"

$$f(x) = \begin{bmatrix} 1 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

Kohli scores a ton, India win.

# 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"

  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"

$$f(x) = \begin{bmatrix} 1 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix} \quad W = \begin{bmatrix} +1 \\ +1 \\ +1 \\ -1 \\ -1 \\ -1 \end{bmatrix}$$

# A three step process for making predictions

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

$$\mathbf{h} = f(\mathbf{x})$$

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

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

$$\textit{binary} \qquad \textit{multi-class}$$

- Decision function: Choose one of the several possibilities

$$\hat{y} = \text{decide}(\mathbf{s})$$

# Another discriminative classifier

# Another discriminative classifier

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

  One hot vectors ("bag of words")

  - $h$ is very sparse.
  - Order ignored.
  - "count of words" instead of "bag of words" could be used.

# Another discriminative classifier

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

  One hot vectors ("bag of words")

- **Score Calculation:** binary    or   multi-class
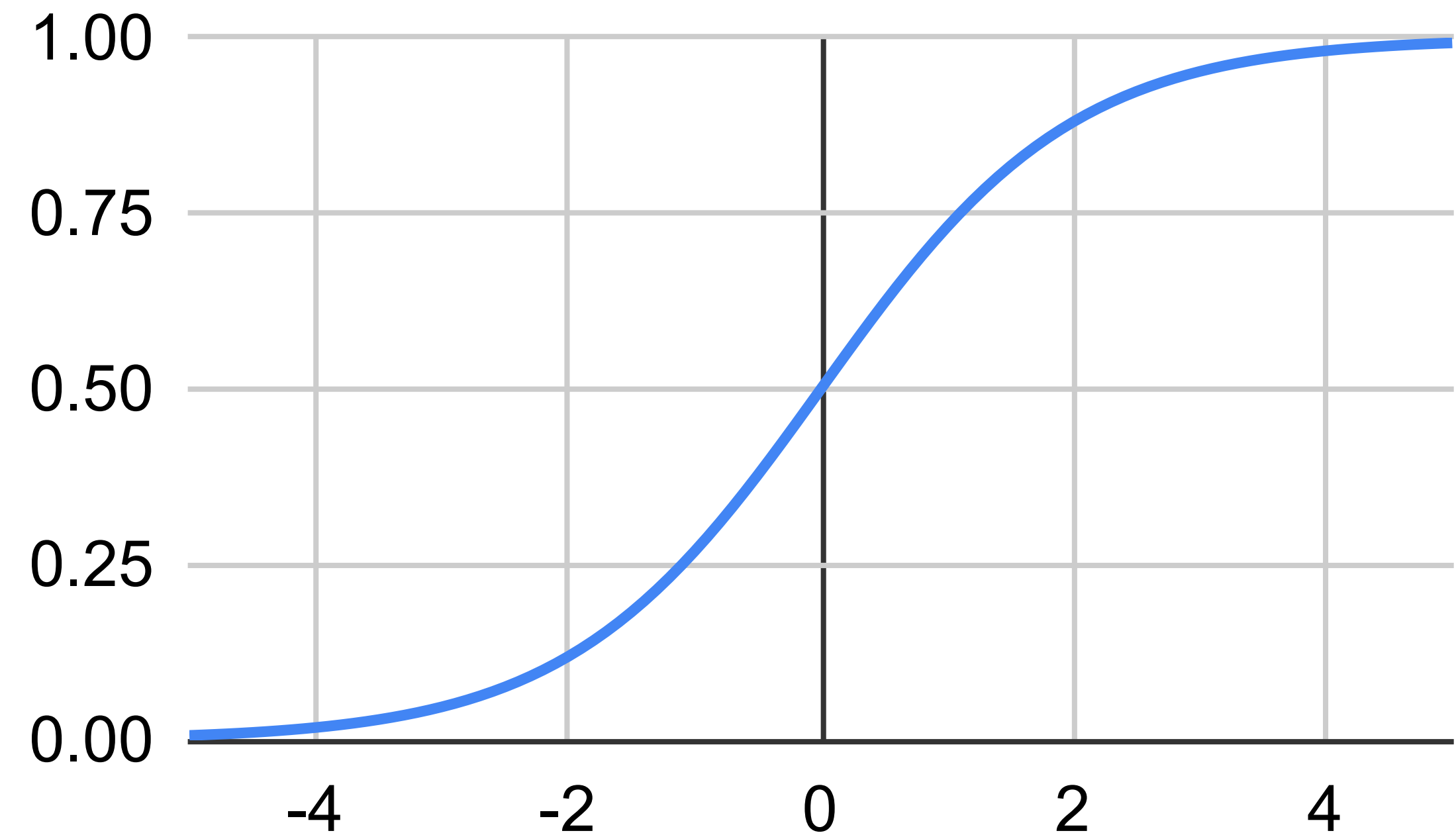
$$s = \mathbf{w} \cdot \mathbf{h} \qquad \mathbf{s} = W\mathbf{h}$$

11

# Another discriminative classifier

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

  One hot vectors ("bag of words")

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

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

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

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

11

# 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 make values 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 $\sigma(wx+b)$ ?

# Softmax

- **Softmax is used for multi-label classification**

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

↑
points to "man"

$$s = \begin{pmatrix} 1 \\ 6 \\ 2 \end{pmatrix} \longrightarrow p = \begin{pmatrix} 0.006 \\ 0.975 \\ 0.017 \end{pmatrix}$$

14

# Recap: Discriminative models

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

$$P(Y \mid X; W)$$

# 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):

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

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):

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

16

# Computing loss (or error/cost) function

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

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

$$L(W) = - \sum_{i}^{n} \left( y_i \ \log P(y = 1 \,|\, x_i; W) \quad + \quad (1 - y_i) \ \log(P(y = 0 \,|\, x_i; W) \right)$$

[binary cross entropy]

16

# Learn parameters through gradient descent

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


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

$$W_{t+1} = W_t - \alpha \nabla_{W_t} L(W)$$

# Computational Graph View

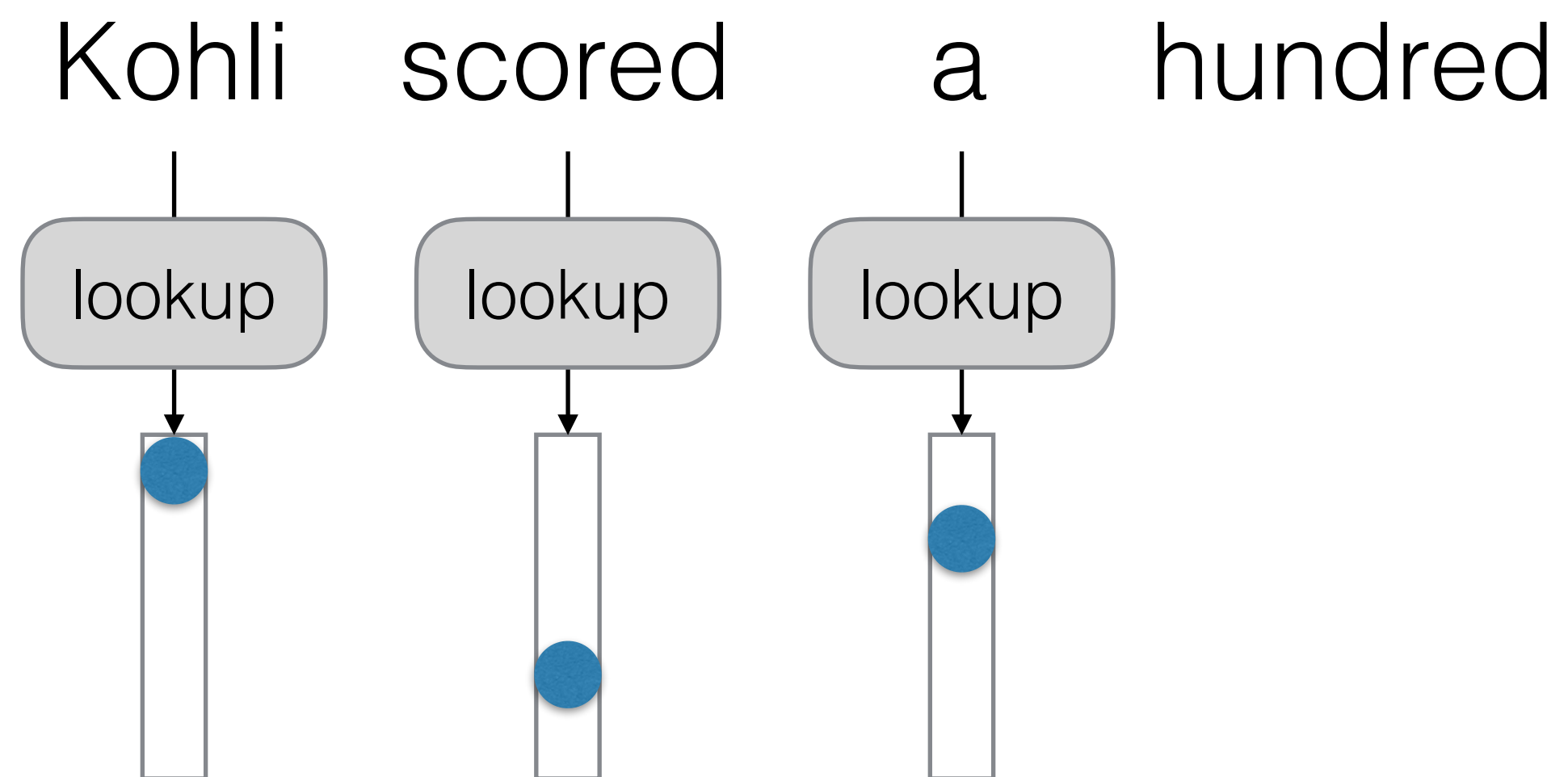Kohli   scored      a     hundred

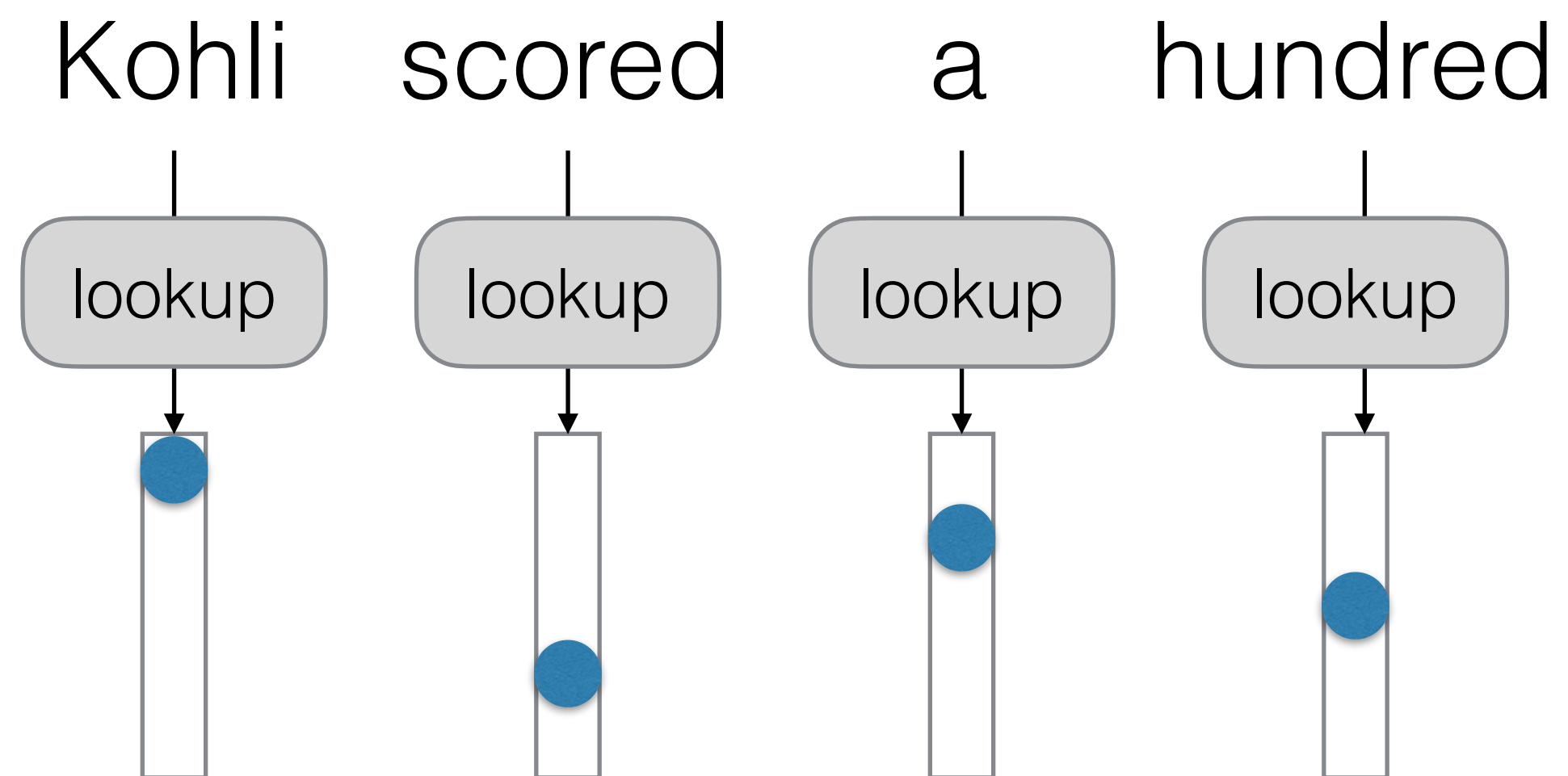# Computational Graph View

Kohli    scored    a    hundred

lookup

# Computational Graph View

Kohli    scored      a     hundred

# Computational Graph View

Kohli    scored    a    hundred

lookup    lookup    lookup

# Computational Graph View

Kohli    scored       a      hundred

# Computational Graph View

Kohli    scored    a    hundred

# Computational Graph View

Kohli    scored    a    hundred



weights

18

# Computational Graph View

# Computational Graph View

Kohli    scored    a    hundred

weights    score    prob

# Questions?

Next class: Word2vec