

# Classification Methods

December 13<sup>th</sup>, 2023



# Logistics - Reminder

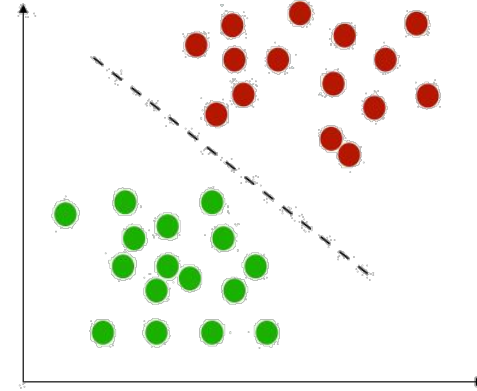
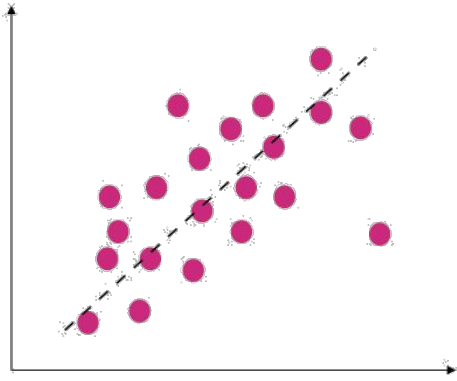
- Everything on course website / Moodle
- 4 credit points
- 2 hrs. lecture, 2 hr. tutorial (Tutorial covers new material)
- 4 homework assignments + Final Project
- All communication through Moodle

# Supervised Learning

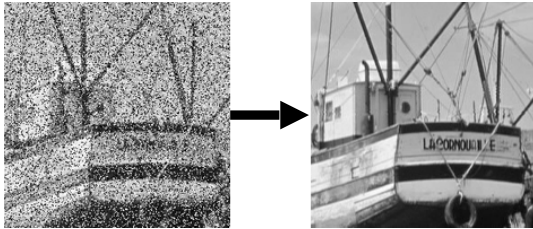
Regression

$$\{x_i, y_i\}_{i=1}^m$$

Classification



E.g. Image denoising



E.g. Image classification



E.g. Object localization



# Binary Classification

New email

בעקבות המצב:  
עד ש50,000 באופן מיידית לכל מטרה  
לשכירים/עצמאיים.  
באישור משרד האוצר  
למחזיקי כרטיס אשראי  
לבדיקת זכאות ללא עלות לחצו:  
<http://bit.ly/2Zrtplp>

**Spam**

**Not Spam**

# Binary Classification

- Samples:  $x_i \in \mathcal{R}^d$ ,  $y_i \in \{0, 1\}$ 
  - $x_i$  - text email features,  $y_i = \begin{cases} 1, & \textit{spam} \\ 0, & \textit{not spam} \end{cases}$
- Training Set:  $\{x_i, y_i\}_{i=1}^m$
- Hypothesis:  $h: \mathcal{R}^d \rightarrow \{0, 1\}$ , parameterized by  $\theta$

# Linear Classifier

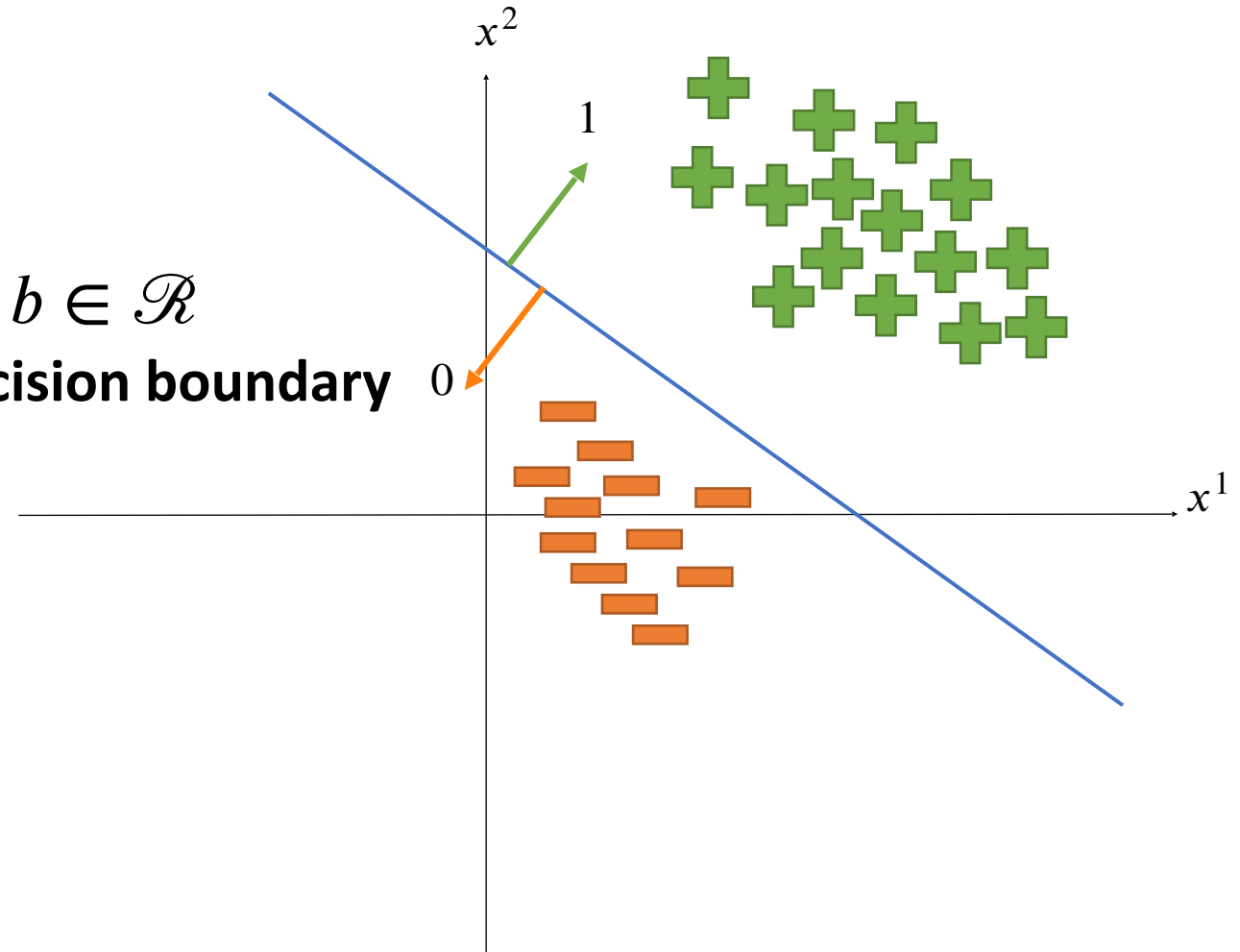
Indicator function

- $\theta = [W, b], h_{W,b}(x_i) = 1(Wx_i + b > 0)$

- Assume:

$$x_i = \begin{pmatrix} x_i^1 \\ x_i^2 \end{pmatrix} \in \mathcal{R}^2, W \in \mathcal{R}^{1 \times 2}, b \in \mathcal{R}$$

Decision boundary



# Linear Classifier

$$h_{\theta}(x_i) = \mathbf{1}(Wx_i + b > 0) \quad \text{+ } x_1$$

- How to find good  $W, b$ ?
- How to interpret the results?

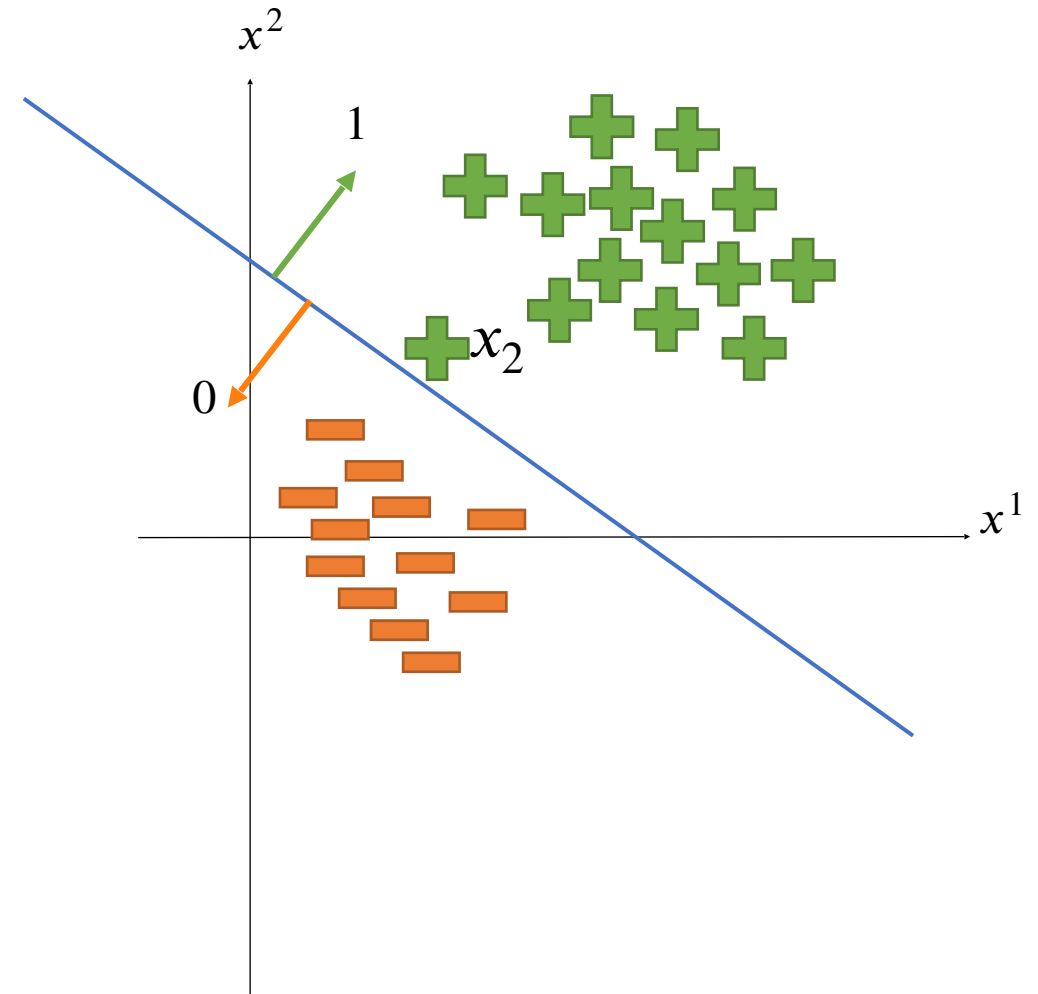
$$Wx_1 + b = \mathbf{1000}$$
$$\Rightarrow h_{W,b}(x_1) = 1$$

**For sure!**

$$Wx_2 + b = \mathbf{0.6}$$
$$\Rightarrow h_{W,b}(x_2) = 1$$

**Possibly...**

- Only intuitive, not measurable



# Logistic Regression (Classification)

- Interpretable results
- Linear decision boundary
- Easy to find  $W, b$





# Sigmoid (Logistic Function)

$Wx + b \rightarrow$  probability

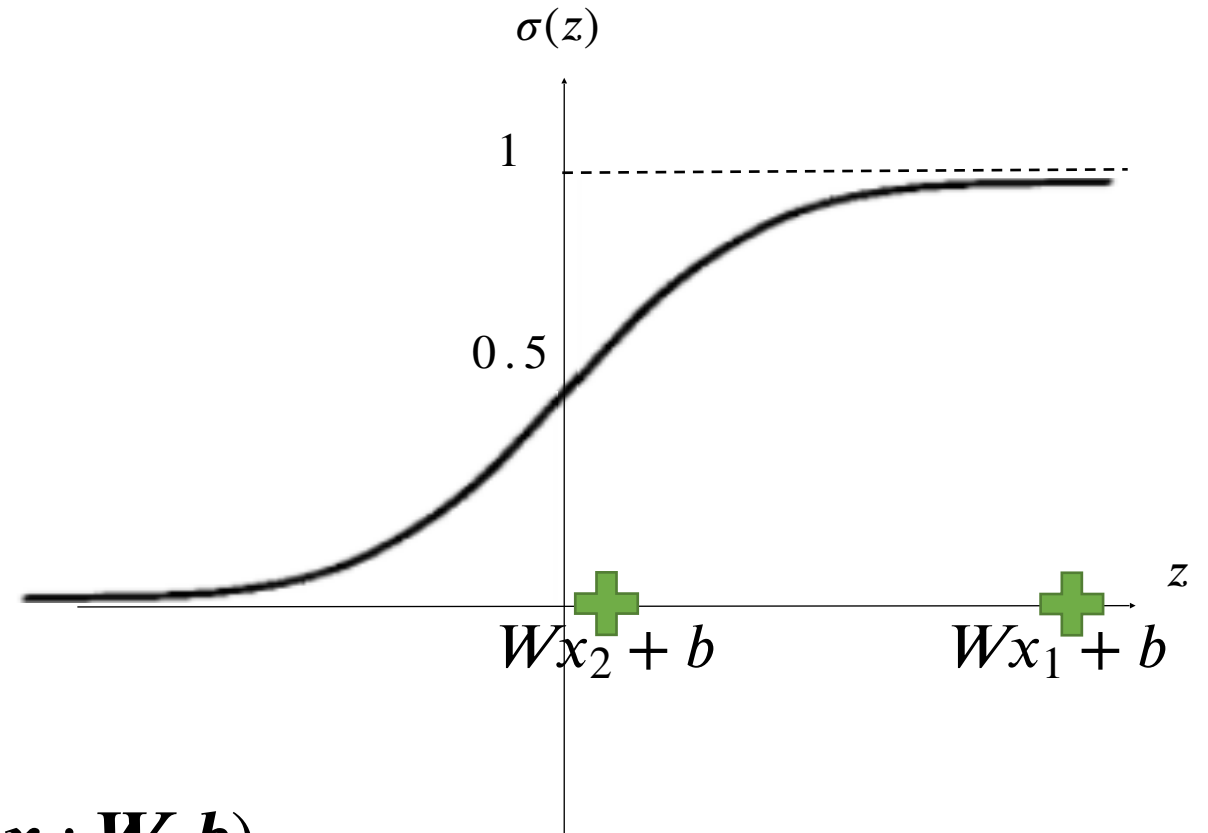
$$z = Wx + b \Rightarrow \sigma(z) = \frac{1}{1 + e^{-z}}$$

Now, our final prediction:

$$h_{W,b}(x_i) = 1(\sigma(z) > 0.5)$$

Interpretation:  $\sigma(z) = \Pr(y_i = \mathbf{1} \mid x_i; \mathbf{W}, \mathbf{b})$

Note:  $1 - \sigma(z) = \Pr(y_i = \mathbf{0} \mid x_i; \mathbf{W}, \mathbf{b})$



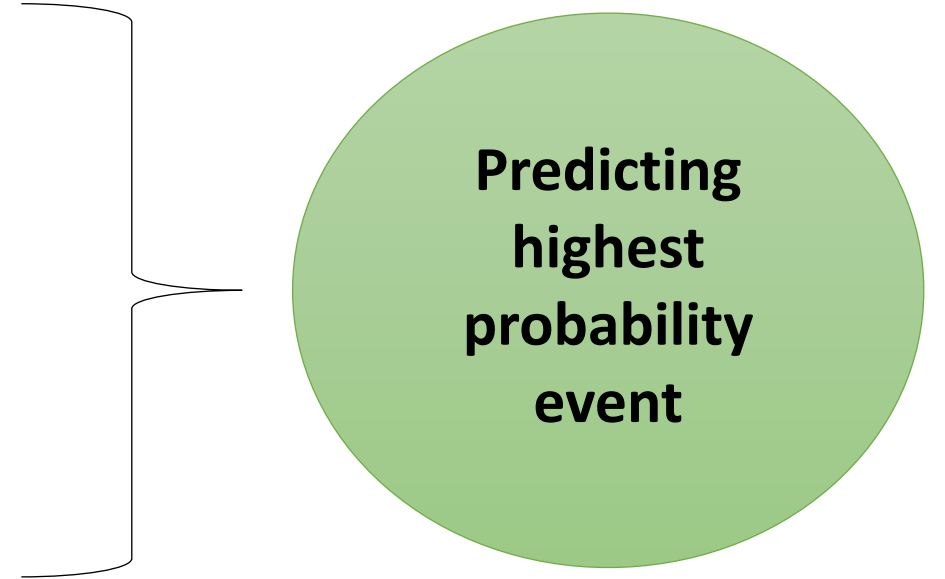
# Sigmoid (Logistic Function)

Final prediction:

$$h_{W,b}(x_i) = 1(\sigma(z) > 0.5)$$

$$\sigma(z) = \Pr(y_i = \mathbf{1} \mid x_i; \mathbf{W}, b)$$

$$1 - \sigma(z) = \Pr(y_i = \mathbf{0} \mid x_i; \mathbf{W}, b)$$



# Logistic Regression

- Interpretable results
- Linear decision boundary
- Easy to find  $W, b$

# Sigmoid (Logistic Function)

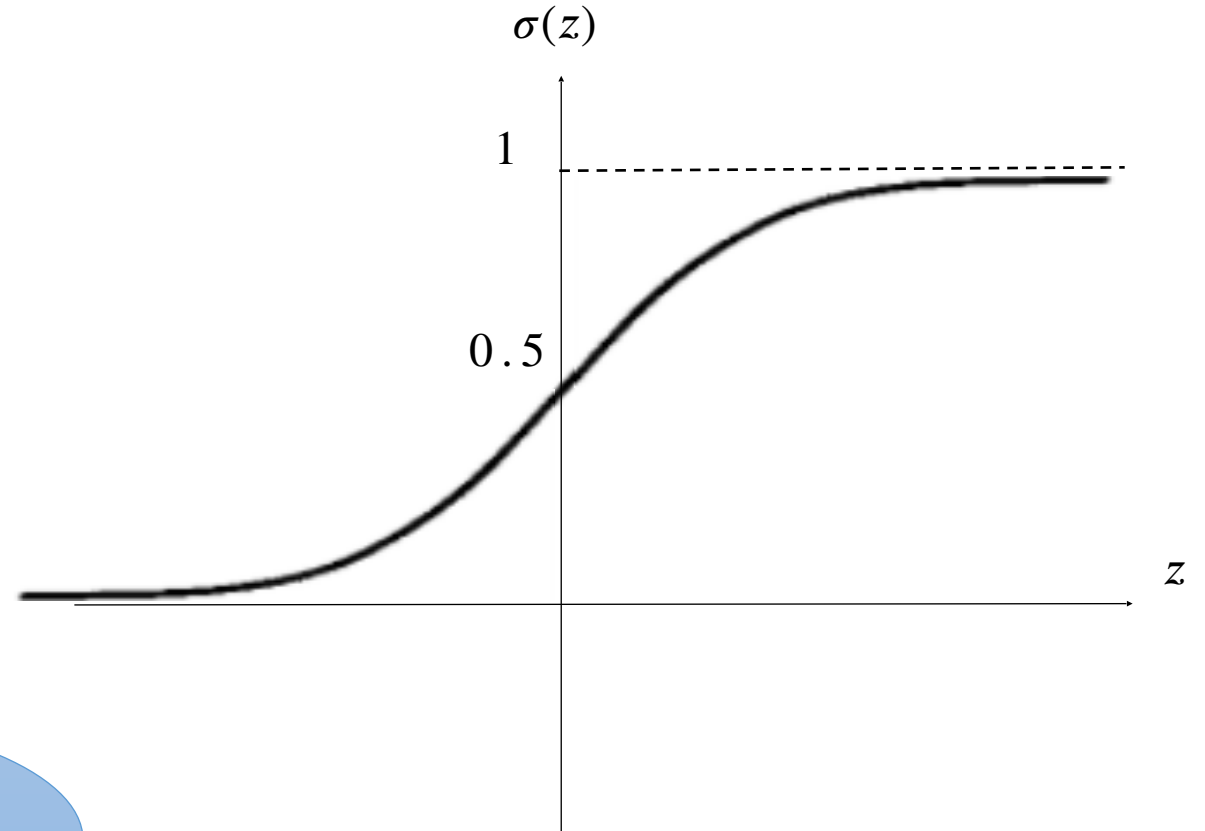
$$h_{W,b}(x_i) = 1(\sigma(z) > 0.5)$$

$$\sigma(z) > 0.5$$

$$\iff z > 0$$

$\Rightarrow$  Linear decision boundary


$$z = Wx + b$$



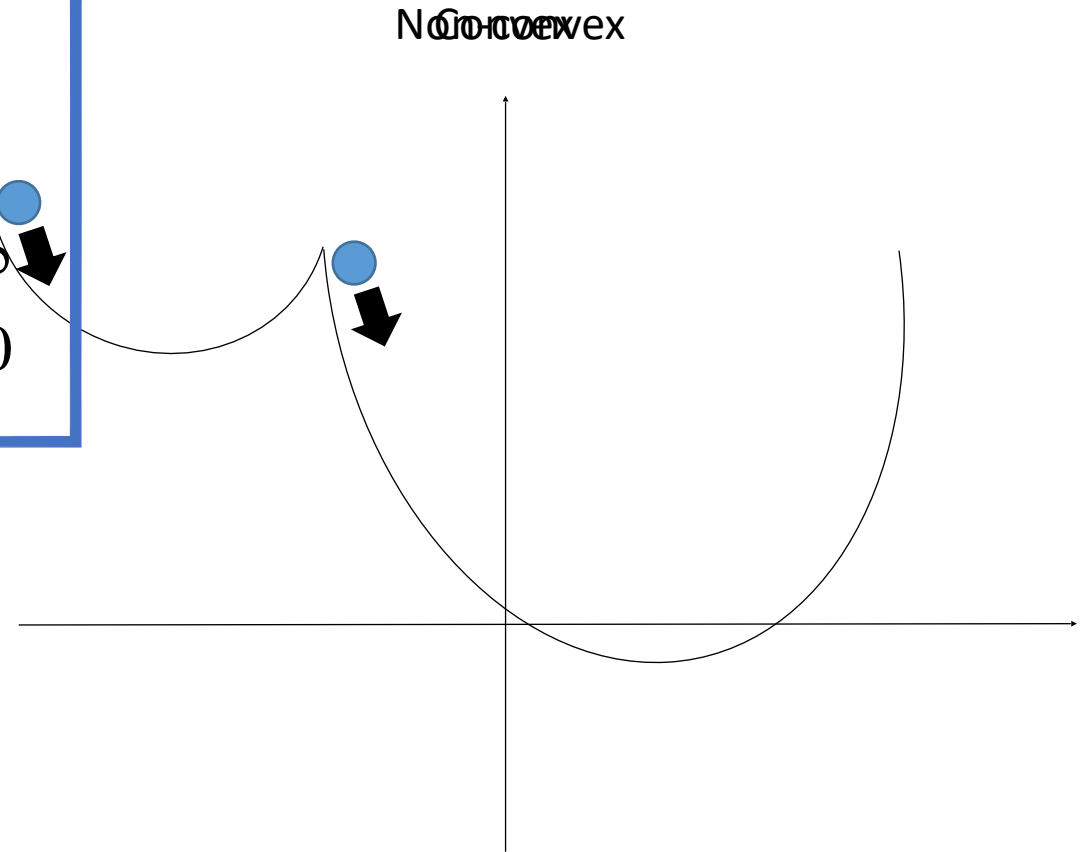
# Logistic Regression

- Interpretable results
- Linear decision boundary
- Easy to find  $W, b$

# How to find good $W, b$ ?

- Define a **convex** cost function  $\mathcal{L}$
- $\mathcal{L}$  penalizes errors
  - Wrong prediction, large penalty:  $\mathcal{L} \rightarrow \infty$
  - Correct prediction, small penalty:  $\mathcal{L} \rightarrow 0$

- Use gradient descent to update  $W, b$



# Cross-Entropy Loss

$$\Pr[y_i = 1 | x_i]$$

- **Prediction:**  $\hat{y}_i = \sigma(z) = \sigma(Wx + b)$
- **Label (GT):**  $y_i$

$$\mathcal{L}_i = y_i \cdot (-\log \hat{y}_i) + (1 - y_i) \cdot (-\log(1 - \hat{y}_i))$$

- Note that  $y_i \in \{0, 1\}$ ,  $\hat{y}_i \in (0, 1)$
- Is this cost function good?

# Cross-Entropy Loss

$$\mathcal{L}_i = y_i \cdot (-\log \hat{y}_i) + (1 - y_i) \cdot (-\log(1 - \hat{y}_i))$$

	$y_i = 0$	$y_i = 1$
$\hat{y}_i \rightarrow 0$	$\mathcal{L} \rightarrow 0$	$\mathcal{L} \rightarrow \infty$
$\hat{y}_i \rightarrow 1$	$\mathcal{L} \rightarrow \infty$	$\mathcal{L} \rightarrow 0$



# Cross-Entropy Loss

- **Prediction:**  $\hat{y}_i = \sigma(z) = \sigma(Wx + b)$
- **Label (GT):**  $y_i$

$$\mathcal{L}_i = y_i \cdot (-\log \hat{y}_i) + (1 - y_i) \cdot (-\log(1 - \hat{y}_i))$$

- Note that  $y_i \in \{0, 1\}$ ,  $\hat{y}_i \in (0, 1)$
- Is this cost function good?
  - Convex w.r.t  $W, b$  ✓
  - Penalizes wrong predictions ✓

# Cross-Entropy Loss - Intuition

- Entropy (information theory) – measures “uncertainty”:
  - Uncertain – all events have the same probability
  - Certain – only one event possible

- Formally:  $H(X) = - \sum_{x \in X} p(x) \log p(x)$

- Expected number of bits to represent an event

- Cross-entropy: have 2 distributions

- How many bits are needed

- Formally:  $H(p, q) =$

- Minimal when  $q = p$



# Cross-Entropy Loss - Intuition

$$H(p, q) = - \sum_{x \in X} p(x) \log q(x)$$

- How is it related to us?
- Want to learn  $W, b$ , such that:  
 $\sigma(Wx_i + b) = \hat{y}_i = q(y_i = 1 | x_i; W, b)$
- Note:  $q(y_i = 0 | x_i; W, b) = 1 - q(y_i = 1 | x_i; W, b)$



# Cross-Entropy Loss - Intuition

$$H(p, q) = - \sum_{x \in X} p(x) \log q(x)$$

- During training, we have a target distribution  $p$  and a predicted distribution  $q$ .

$$p(y_i)$$

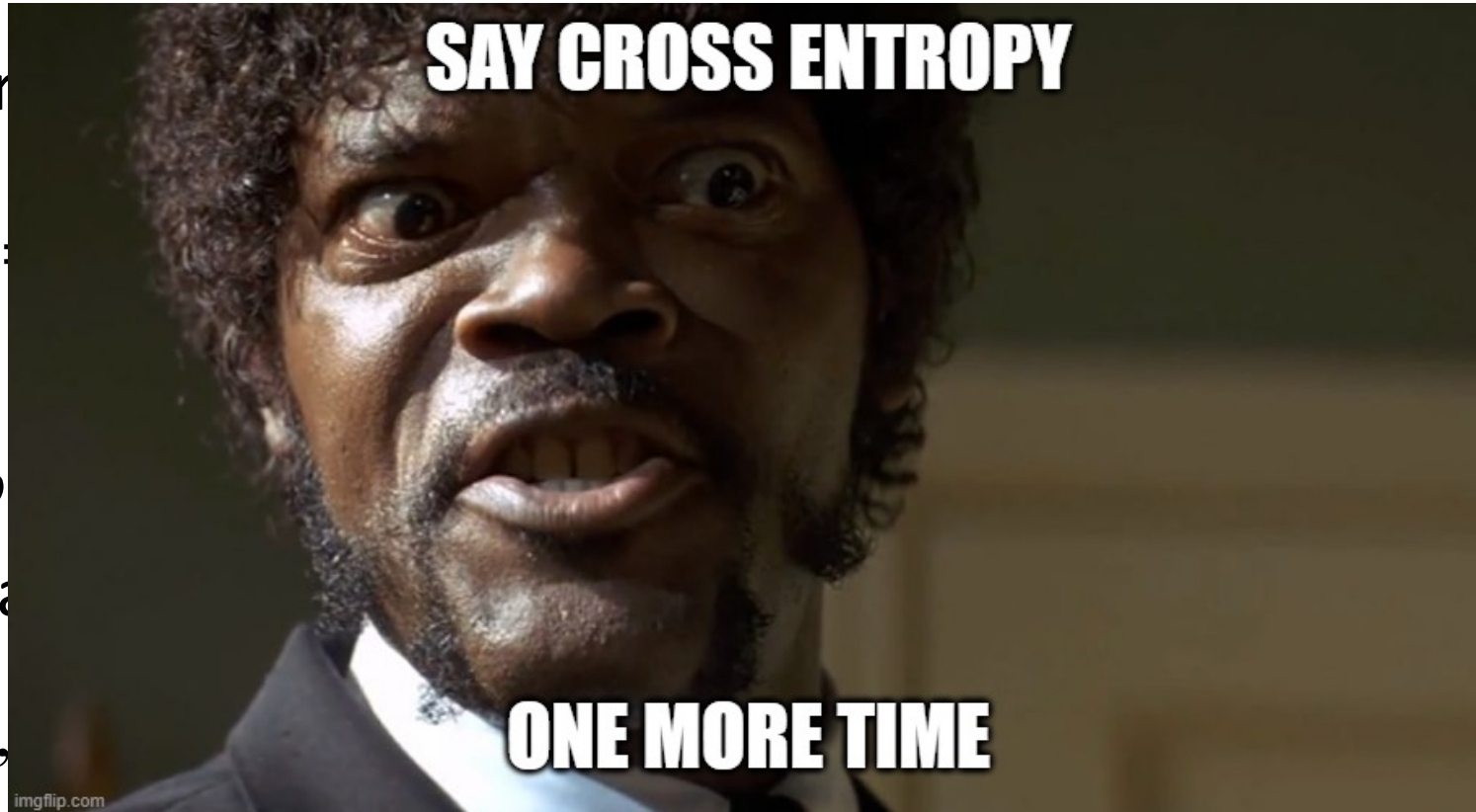
$$\Rightarrow p$$

- Finally, the cross-entropy loss is defined as:

$$H(p, q)$$

$$= y_i \cdot (-\log \hat{y}_i) + (1 - y_i) \cdot (-\log(1 - \hat{y}_i))$$

$$= \mathcal{L}_i$$



$$(1 - \hat{y}_i)$$

# How to find good $W, b$ ?

- Define a **convex** cost function  $\mathcal{L}$
- $\mathcal{L}$  penalizes errors
  - Wrong prediction, large penalty:  $\mathcal{L} \rightarrow \infty$
  - Correct prediction, small penalty:  $\mathcal{L} \rightarrow 0$

- Use gradient descent to update  $W, b$

# Gradient Descent – Single Sample

$$x_i \in \mathcal{R}^{d \times 1}, y_i \in \{0, 1\}$$

Find prediction

$$x_i \in \mathcal{R}^{1 \times d}, b \in \mathcal{R} \Rightarrow \hat{y}_i \in (0, 1)$$

Find prediction

Calculate Loss

$$+ (1 - y_i) \cdot \left( -\log(1 - \hat{y}_i) \right)$$

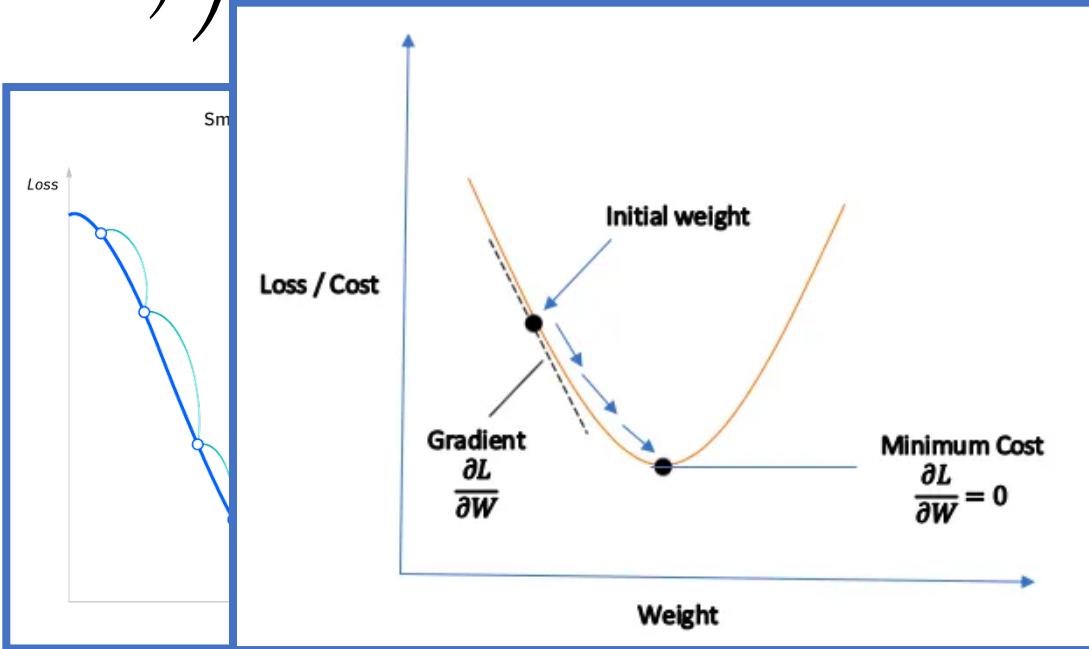
Calculate Loss

Update weights

$$W \leftarrow W - \alpha \frac{\partial \mathcal{L}_i}{\partial W} = W - \alpha \cdot x_i^T (\hat{y}_i - y_i)$$

Update bias

$$b \leftarrow b - \alpha \frac{\partial \mathcal{L}_i}{\partial b} = b - \alpha \cdot (\hat{y}_i - y_i)$$



# Stochastic Mini-Batch Gradient Descent

Repeat until convergence:

Sample a mini batch  $\{x_i, y_i\}_{i=1}^m$ :

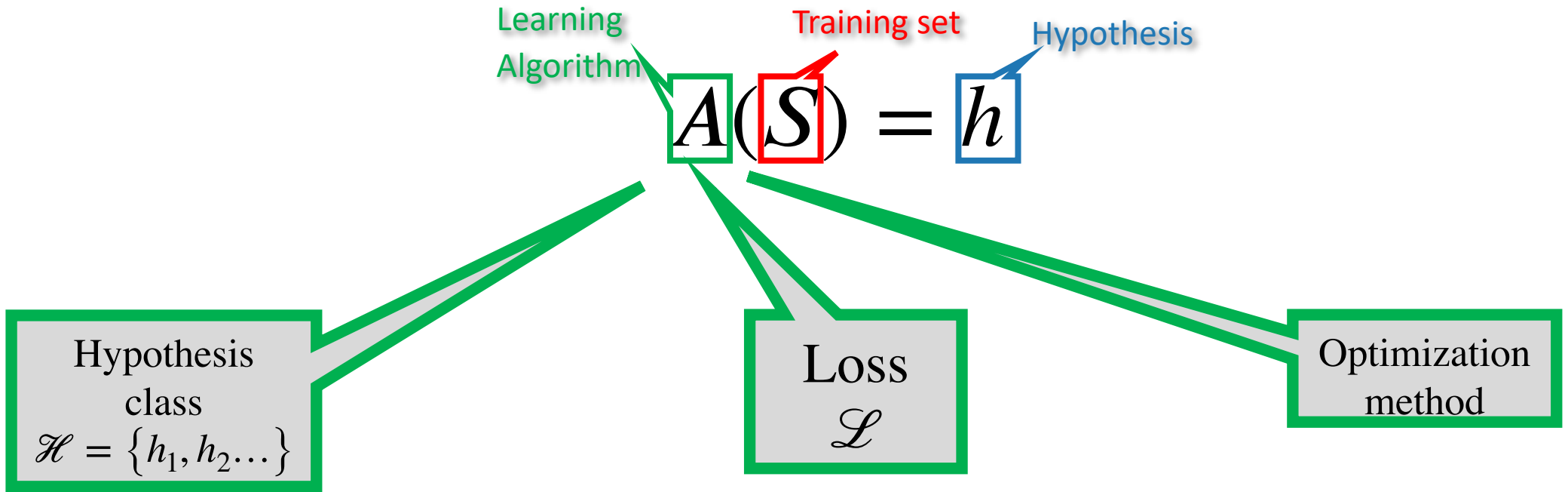
Calculate  $\mathcal{L}_i$  for each pair

$$\mathcal{L} = \frac{1}{m} \sum_i \mathcal{L}_i$$

$$W \leftarrow W - \alpha \cdot \frac{\partial \mathcal{L}}{\partial W}$$

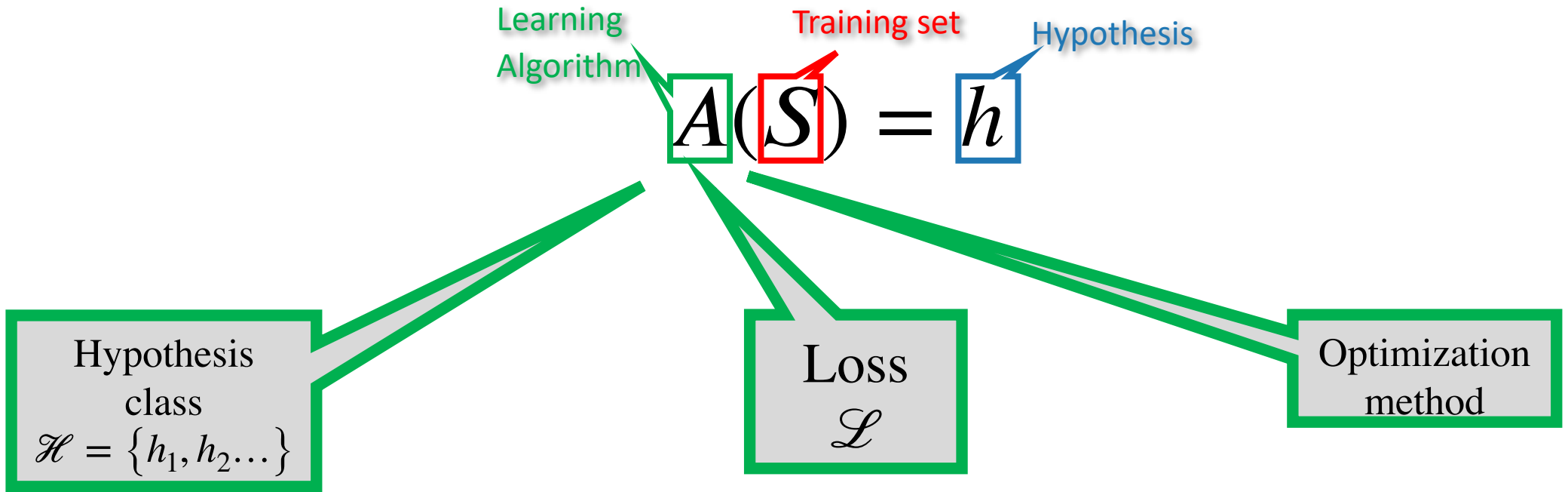
$$b \leftarrow b - \alpha \cdot \frac{\partial \mathcal{L}}{\partial b}$$

# Supervised Learning





# Supervised Learning – Logistic Regression



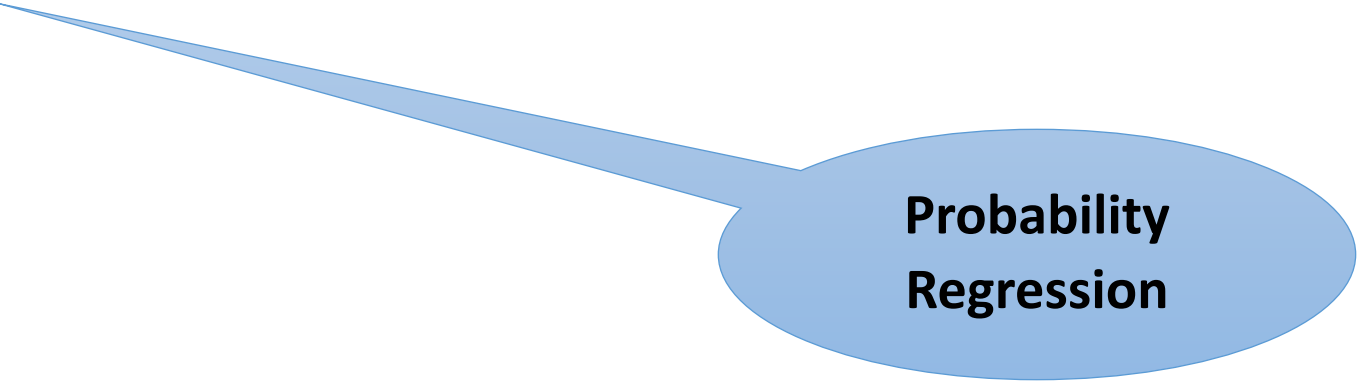
$$\mathcal{H} = \{h_{W,b}(x_i) = 1(\sigma(Wx_i + b) > 0.5)\}$$

Cross-entropy loss

SGD

# Logistic Regression - Summary

- **Binary** classification
- **Linear** ( $Wx + b$ )
- **Sigmoid** for interpretable (probability) output
- **Predict** highest probability event
- **Learning:**
  - **Cross-entropy** loss (convex, penalizes mistakes)
  - **Gradient descent** to update weights and bias



**Probability  
Regression**

# Multi Class Classification

New email

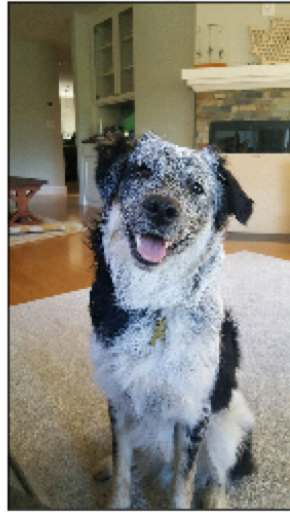
בעקבות המצב:  
עד ש50,000 באופן מיידית לכל מטרה  
לשכירים/עצמאיים.  
באישור משרד האוצר  
למחזיקי כרטיס אשראי  
לבדיקת זכאות ללא עלות לחצו:  
<http://bit.ly/2Zrtplp>

**Spam**

**Maybe**

**Not Spam**

# Multi Class Classification



# Multi Class Classification

- $C$  classes
- Samples:  $x_i \in \mathcal{R}^d$ ,  $y_i \in [C]$ 
  - $x_i$  - text email features,  $y_i = \begin{cases} 2, & \text{maybe} \\ 1, & \text{spam} \\ 0, & \text{not spam} \end{cases}$
- Training Set:  $\{x_i, y_i\}_{i=1}^m$
- Hypothesis:  $h: \mathcal{R}^d \rightarrow [C]$ , parameterized by  $\theta$

# Softmax Classifier

- Generalizing Logistic Regression, **binary** → **multi-class**
- Begin with **Linear Transformation** ( $z = Wx + b$ )

Similar Learning:  
1. Cross-entropy loss  
2. SGD

Logistic Regression	Softmax classifier
$W \in \mathbb{R}^{1 \times d}, b \in \mathbb{R}$	$W \in \mathbb{R}^{C \times d}, b \in \mathbb{R}^C$

- Apply **non-linearity** to transform  $z$  to  $\hat{y}$  (probability)

Sigmoid (logistic function)	Softmax
-----------------------------	---------

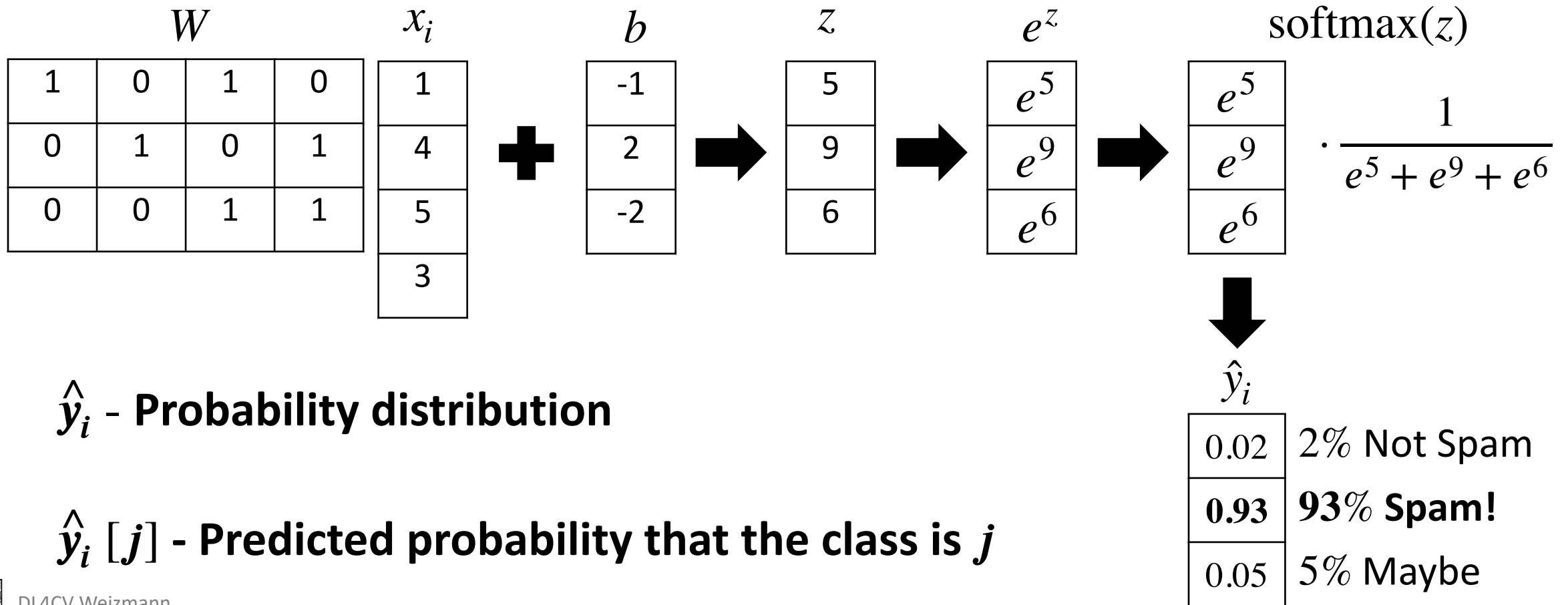
$C$  dimensional  $z$   
 $z^j = j$ 'th class score

- Final prediction – pick the event with **highest probability**

$1(\hat{y} > 0.5)$	$\operatorname{argmax}(\hat{y})$
--------------------	----------------------------------

# Softmax Function - Example

Given a text  $x_i$ : **Not Spam (0)**, **Spam(1)** or **Maybe (2)**



# Softmax Function - Formally

Converting  $z$  to probability distribution over the classes  $C$ :

$$z = Wx + b, \quad z \in \mathcal{R}^C$$

$$\Rightarrow f(z)^i = \frac{e^{z^i}}{\sum_{j \in [C]} e^{z^j}}$$

$$\Rightarrow \hat{y} = \left[ f(z)^0, \dots, f(z)^{C-1} \right], \quad \sum_{j \in C} \hat{y}^j = 1$$



# Softmax Classifier

- Generalizing Logistic Regression, **binary** → **multi-class**
- Begin with **Linear Transformation** ( $z = Wx + b$ )

Similar Learning:  
1. Cross-entropy loss  
2. SGD

Logistic Regression	Softmax classifier
$W \in \mathbb{R}^{1 \times d}, b \in \mathbb{R}$	$W \in \mathbb{R}^{C \times d}, b \in \mathbb{R}^C$

- Apply **non-linearity** to transform  $z$  to  $\hat{y}$  (probability)

Sigmoid (logistic function)	Softmax
-----------------------------	---------

- Final prediction – pick the event with **highest probability**

$1(\hat{y} > 0.5)$	$\operatorname{argmax}(\hat{y})$
--------------------	----------------------------------

# Cross-Entropy Loss – Softmax Classifier

$$H(p, q) = - \sum_{x \in X} p(x) \log q(x)$$

$p(y_i = c | x)$  – Ground Truth (known),  $\hat{y}_i^c = q(y_i = c | x)$  – Our Prediction

0
<b>1</b>
0

Spam  
**Not spam!**  
Maybe

0.02
<b>0.93</b>
0.05

2% Spam  
**93% Not spam!**  
5% Maybe

$$H(p, q) = - \sum_{c \in \{0, 1, 2\}} p(y_i = c | x_i) \cdot \log(\hat{y}_i^c) = - \log(\hat{y}_i^1)$$

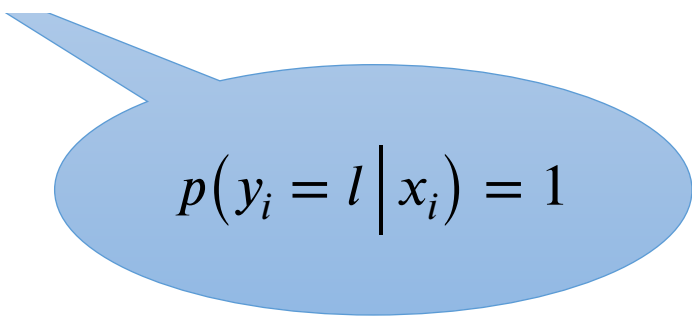
**Goal: Prediction should be closer to GT.**

# Cross-Entropy Loss – Softmax Classifier

- Generally, with  $C$  classes, training example  $\{x_i, l\}$

$$\mathcal{L}_i = H(p, q) = - \sum_{c \in C} p(y_i = c | x_i) \cdot \log(\hat{y}_i^c)$$

$$= - \log \frac{e^{z^l}}{\sum_{j \in [C]} e^{z^j}}$$


$$p(y_i = l | x_i) = 1$$

- Wish to make the loss **smaller** (log argument larger):
  - **Increase numerator** → higher confidence of class  $l$
  - **Decrease denominator** → higher confidence it's not other classes!

# Softmax Classifier

- Generalizing Logistic Regression, **binary** → **multi-class**
- Begin with **Linear Transformation** ( $z = Wx + b$ )

Similar Learning:  
1. Cross-entropy loss  
2. SGD

Logistic Regression	Softmax classifier
$W \in \mathbb{R}^{1 \times d}, b \in \mathbb{R}$	$W \in \mathbb{R}^{C \times d}, b \in \mathbb{R}^C$

- Apply **non-linearity** to transform  $z$  to  $\hat{y}$  (probability)

Sigmoid (logistic function)	Softmax
-----------------------------	---------

- Final prediction – pick the event with **highest probability**

$1(\hat{y} > 0.5)$	$\operatorname{argmax}(\hat{y})$
--------------------	----------------------------------

# Gradient Descent – Logistic Regression

$$x_i \in \mathcal{R}^{d \times 1}, y_i \in \{0, 1\}$$

$$\hat{y}_i = \sigma(Wx_i + b) \quad (W \in \mathcal{R}^{1 \times d}, b \in \mathcal{R} \Rightarrow \hat{y}_i \in [0, 1])$$

$$\mathcal{L}_i = y_i \cdot \left( -\log \hat{y}_i \right) + (1 - y_i) \cdot \left( -\log (1 - \hat{y}_i) \right)$$

$$W \leftarrow W - \alpha \cdot \frac{\partial \mathcal{L}_i}{\partial W} = W - \alpha \cdot x_i^T (\hat{y}_i - y_i)$$

$$b \leftarrow b - \alpha \cdot \frac{\partial \mathcal{L}_i}{\partial b} = b - \alpha \cdot (\hat{y}_i - y_i)$$

Find prediction

Calculate Loss

Update weight

Update bias

# Gradient Descent – Softmax Classifier

$$x_i \in \mathcal{R}^{d \times 1}, y_i \in \{\mathbf{0}, \mathbf{1}, \dots, \mathbf{C} - \mathbf{1}\}$$

$$\hat{y}_i = \sigma(Wx_i + b) \quad (W \in \mathcal{R}^{1 \times d}, b \in \mathcal{R} \Rightarrow \hat{y}_i \in [0, 1])$$

$$\mathcal{L}_i = y_i \cdot \left(-\log \hat{y}_i\right) + (1 - y_i) \cdot \left(-\log(1 - \hat{y}_i)\right)$$

$$W \leftarrow W - \alpha \cdot \frac{\partial \mathcal{L}_i}{\partial W} = W - \alpha \cdot x_i^T (\hat{y}_i - y_i)$$

$$b \leftarrow b - \alpha \cdot \frac{\partial \mathcal{L}_i}{\partial b} = b - \alpha \cdot (\hat{y}_i - y_i)$$

Find prediction

Calculate Loss

Update weight

Update bias

# Gradient Descent – Softmax Classifier

$$x_i \in \mathcal{R}^{d \times 1}, y_i \in \{0, 1, \dots, C - 1\}$$

$$\hat{y}_i = f(Wx_i + b) \quad (W \in \mathcal{R}^{C \times d}, b \in \mathcal{R}^C, \sum_{j \in C} \hat{y}_i[j] = 1)$$

$$\mathcal{L}_i = y_i \cdot \left(-\log \hat{y}_i\right) + (1 - y_i) \cdot \left(-\log(1 - \hat{y}_i)\right)$$

$$W \leftarrow W - \alpha \cdot \frac{\partial \mathcal{L}_i}{\partial W} = W - \alpha \cdot x_i^T (\hat{y}_i - y_i)$$

$$b \leftarrow b - \alpha \cdot \frac{\partial \mathcal{L}_i}{\partial b} = b - \alpha \cdot (\hat{y}_i - y_i)$$

Find prediction

Calculate Loss

Update weight

Update bias

# Gradient Descent – Softmax Classifier

$$x_i \in \mathcal{R}^{d \times 1}, y_i \in \{0, 1, \dots, C - 1\}$$

$$\hat{y}_i = f(Wx_i + b) \quad (W \in \mathcal{R}^{C \times d} \quad b \in \mathcal{R}^C, \quad \sum_{j \in C} \hat{y}_i[j] = 1)$$

$$\mathcal{L}_i = - \sum_{c \in C} p(y_i = c | x_i) \cdot \log(f(z)[c])$$

$$W \leftarrow W - \alpha \cdot \frac{\partial \mathcal{L}_i}{\partial W} = W - \alpha \cdot x_i^T (\hat{y}_i - y_i)$$

$$b \leftarrow b - \alpha \cdot \frac{\partial \mathcal{L}_i}{\partial b} = b - \alpha \cdot (\hat{y}_i - y_i)$$

Find prediction

Calculate Loss

Update weight

Update bias



# Gradient Descent – Softmax Classifier

$$x_i \in \mathcal{R}^{d \times 1}, y_i \in \{0, 1, \dots, C - 1\}$$

$$\hat{y}_i = f(Wx_i + b) \quad (W \in \mathcal{R}^{C \times d} \quad b \in \mathcal{R}^C \quad \sum_{c=0}^{C-1} \Delta_c = 1)$$

$$\mathcal{L}_i = - \sum_{c \in C} p(y_i = c | x_i) \cdot \log(f(z) [c])$$

$$\delta_{y_i}[c] = \begin{cases} 1, & c = y_i \\ 0, & o.w. \end{cases}$$

Find prediction

Calculate Loss

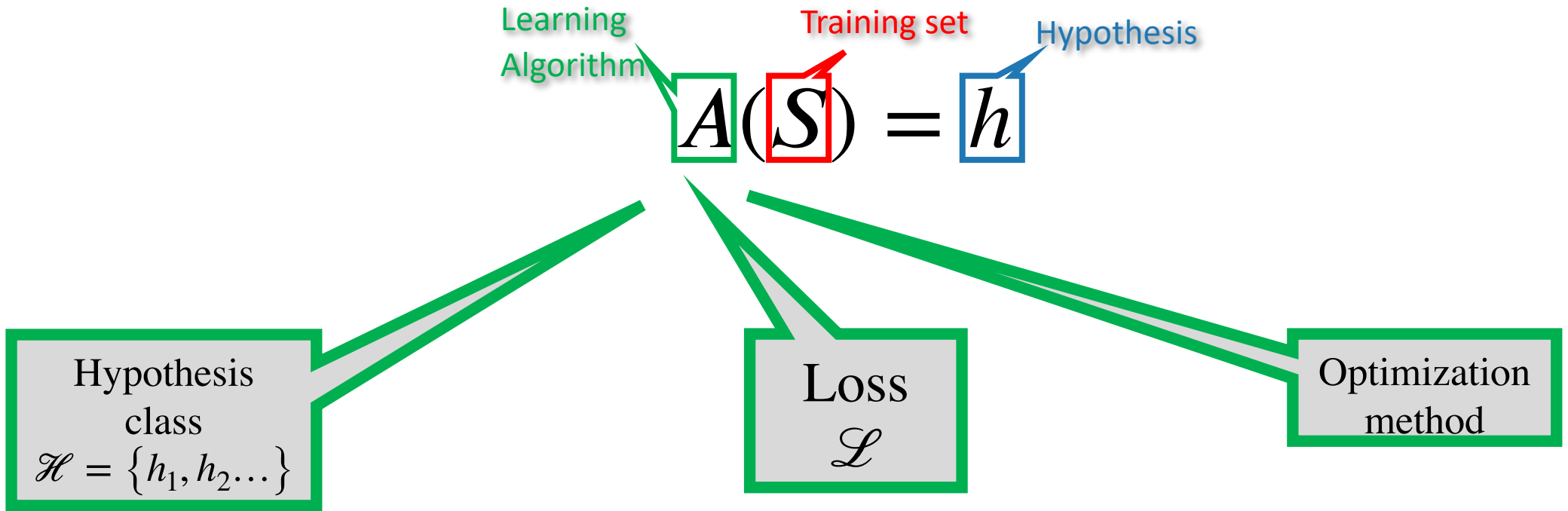
Update weight

Update bias

$$W \leftarrow W - \alpha \cdot \frac{\partial \mathcal{L}_i}{\partial W} = W - \alpha \cdot (\hat{y}_i - \delta_{y_i}) x_i^T$$

$$b \leftarrow b - \alpha \cdot \frac{\partial \mathcal{L}_i}{\partial b} = b - \alpha \cdot (\hat{y}_i - \delta_{y_i})$$

# Supervised Learning – Softmax Classifier



$$\mathcal{H} = \{h_{W,b}(x_i) = \operatorname{argmax}(f(Wx_i + b))\}$$

Cross-entropy loss

SGD

# Conclusion

- Begin with **Linear Transformation** ( $z = Wx + b$ ):

Logistic Regression	Softmax classifier
$W \in \mathbb{R}^{1 \times d}, b \in \mathbb{R}$	$W \in \mathbb{R}^{C \times d}, b \in \mathbb{R}^C$

- Apply **non-linearity** to transform  $z$  to  $\hat{y}$  (probability)

Sigmoid (logistic function)	Softmax
-----------------------------	---------

- Final prediction – the event with **highest probability**

$1(\hat{y} > 0.5)$	$\operatorname{argmax}(\hat{y})$
--------------------	----------------------------------

Similar Learning:

1. Cross-entropy loss
2. SGD

# Practical Considerations

- **Mini Batches**
- **Numerical Stability**

# Stochastic Mini-Batch Gradient Descent

Repeat until convergence:

Sample a mini batch  $\{x_i, y_i\}_{i=1}^m$ :

Calculate  $\mathcal{L}_i$  for each pair

$$\mathcal{L} = \frac{1}{m} \sum_i \mathcal{L}_i$$

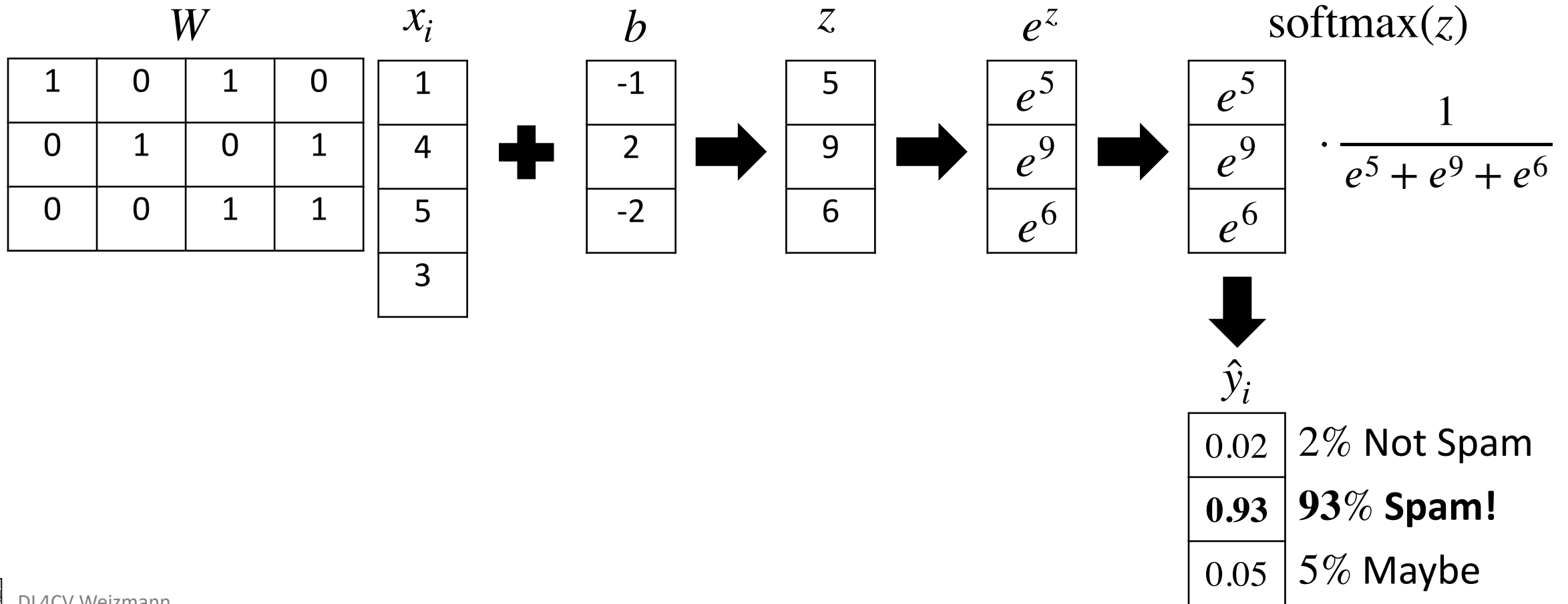
$$W \leftarrow W - \alpha \cdot \frac{\partial \mathcal{L}}{\partial W}$$

$$b \leftarrow b - \alpha \cdot \frac{\partial \mathcal{L}}{\partial b}$$

**Not efficient  
sequentially**

# Softmax Function - Reminder

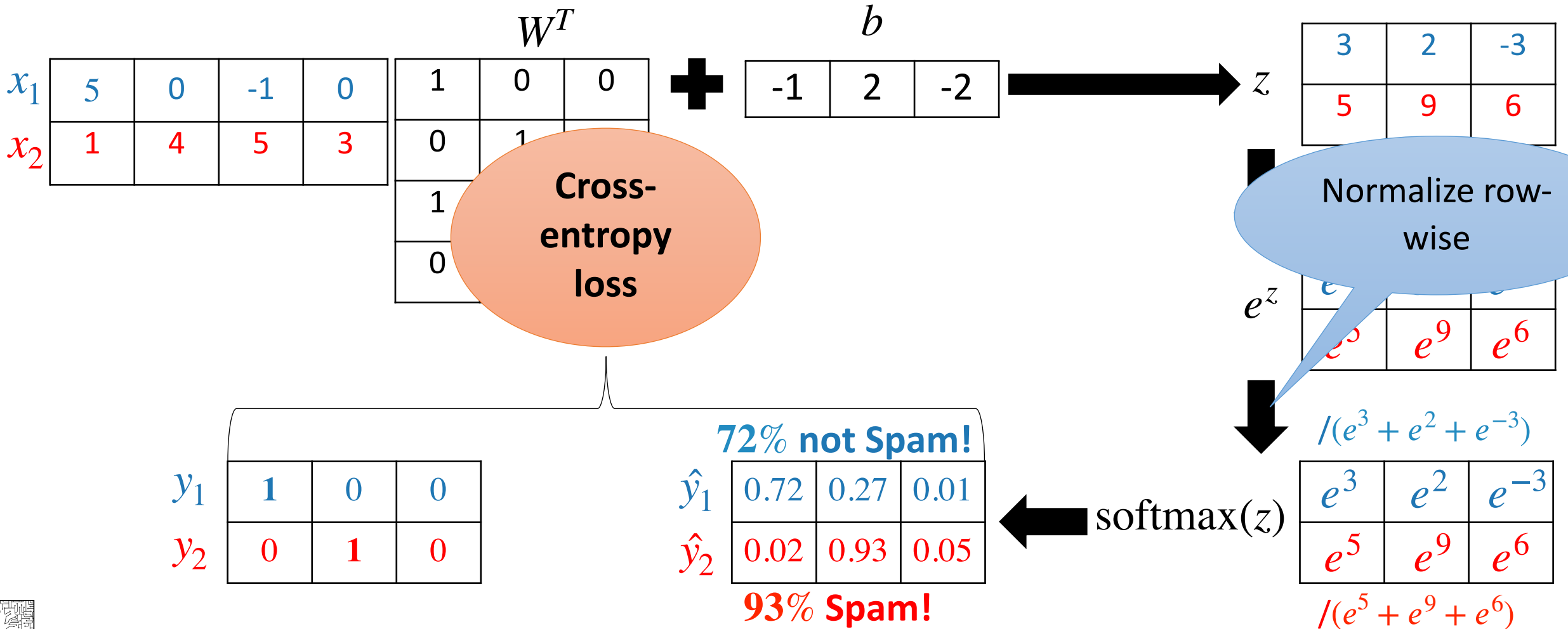
Given a text  $x_i$ : **Not Spam (0)**, **Spam(1)** or **Maybe (2)**



# Softmax Classifier - Batched Example

Given two pairs:  $\{x_1, 0\}$ ,  $\{x_2, 1\}$

spam



# Mini Batches - Formally

Define  $X = [x_1, \dots, x_m]^T \in \mathcal{R}^{m \times d}$  - samples matrix

$Y \in \mathcal{R}^{m \times C}$  - GT labels matrix (Each row in  $Y$  is one-hot vector)

Reminder:  $W \in \mathcal{R}^{C \times d}$ ,  $b \in \mathcal{R}^C$

$$Z = \underbrace{XW^T}_{\in \mathcal{R}^{m \times c}} + b$$

$\in \mathcal{R}^{m \times c}$

( $b$  is added to each row)



# Mini Batches - Formally

$$Z = \underbrace{XW^T}_{\in \mathcal{R}^{m \times c}} + b, \quad Z \in \mathcal{R}^{m \times C}$$

Softmax & Cross-entropy applied row-wise

$$\hat{Y} = f(Z)$$

$$\mathcal{L} = \frac{1}{m} \sum_i \mathcal{L}_i$$

Reminder: Weights update (single-sample):

$$W \leftarrow W - \alpha \cdot \frac{\partial \mathcal{L}_i}{\partial W} = W - \alpha \cdot \left( \hat{y}_i - \delta_{y_i} \right) x_i^T$$

Weights update (mini-batch):

$$W \leftarrow W - \alpha \cdot \frac{\partial \mathcal{L}}{\partial W} = W - \frac{\alpha}{m} \cdot \left( \hat{Y} - Y \right)^T X$$

Similarly derive  
cross-entropy  
rule

# Numerical Stability

$$f(z)^i = \frac{e^{z^i}}{\sum_{j \in [C]} e^{z^j}}$$

- Softmax:  $e^z$  can be extremely large

- Cross-entropy:  $\hat{y}$  possibly close to zero  $\Rightarrow \log(\hat{y}_i) \rightarrow -\infty$

- Many possible solutions

$$H(p, q) = - \sum_{c \in [C]} p(y_i = c | x_i) \cdot \log(\hat{y}_i^c)$$

**QUESTIONS?**



imgflip.com

