# Basics of Neural Networks

## Deep Learning



Biological Neuron

Artificial Neuron / Perceptron

Recap:

# Logistic Regression!

Learned Model Parameters

## Linear Predictive Model



$x_{i1}$ $x_{i2}$      features of data      $x_{iM}$

## Linear Predictive Model



$x_{i1}$ $x_{i2}$      features of data      $x_{iM}$

## Linear Predictive Model



$x_{i1}$  $x_{i2}$          features of data          $x_{iM}$

## Linear Predictive Model

$$(b_1 \times x_{i1})$$



$x_{i1}$ $x_{i2}$       features of data       $x_{iM}$

## Linear Predictive Model

$$(b_1 \times x_{i1}) + (b_2 \times x_{i2}) + \cdots + (b_M \times x_{iM})$$

$x_{i1}$  $x_{i2}$          features of data          $x_{iM}$

## Linear Predictive Model

$$(b_1 \times x_{i1}) + (b_2 \times x_{i2}) + \cdots + (b_M \times x_{iM}) + b_0$$
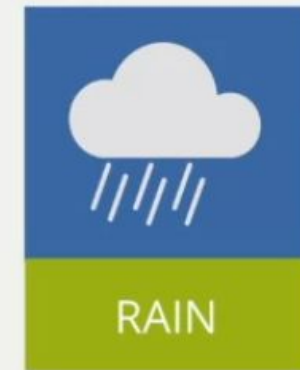
bias



$x_{i1}$  $x_{i2}$    features of data    $x_{iM}$

## Will it Rain?

$x_i$ = features for day $i$

features

outcome

$y_i$ = 1, yes
$y_i$ = 0, no

**WEATHER**

| Cloud Cover | Humidity | Temperature | Air Pressure | Did it Rain |
|-------------|----------|-------------|--------------|-------------|
| 0.5 | 80% | 75 | 1.2 | 1 |
| 0.2 | 95% | 83 | 1.3 | 0 |

**RAIN**

$z_1 = (b_1 \times 0.5) + (b_2 \times 0.8) + (b_3 \times 75) + (b_4 \times 1.2) + b_0$    $y_1 = 1$

$z_2 = (b_1 \times 0.2) + (b_2 \times 0.95) + (b_3 \times 83) + (b_4 \times 1.3) + b_0$    $y_2 = 0$
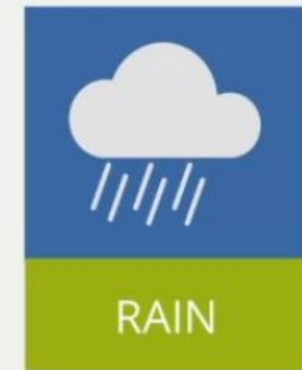
## Will it Rain?

$x_i$ = features for day $i$

features

outcome

$y_i = 1$, yes
$y_i = 0$, no

| Cloud Cover | Humidity | Temperature | Air Pressure | Did it Rain |
|---|---|---|---|---|
| 0.5 | 80% | 75 | 1.2 | 1 |
| 0.2 | 95% | 83 | 1.3 | 0 |

WEATHER

RAIN

$z_1 = (b_1 \times 0.5) + (b_2 \times 0.8) + (b_3 \times 75) + (b_4 \times 1.2) + b_0 \quad y_1 = 1$

$z_2 = (b_1 \times 0.2) + (b_2 \times 0.95) + (b_3 \times 83) + (b_4 \times 1.3) + b_0 \quad y_2 = 0$

sigma

$p(y_i = 1 | x_i) = \sigma(z_i)$
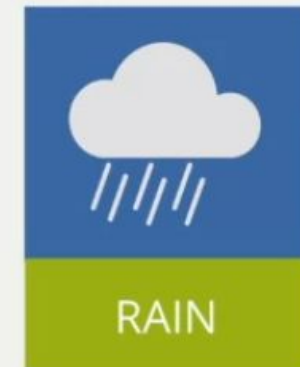
## Will it Rain?

$x_i$ = features for day $i$

features

| Cloud Cover | Humidity | Temperature | Air Pressure |
|---|---|---|---|
| 0.5 | 80% | 75 | 1.2 |
| 0.2 | 95% | 83 | 1.3 |

outcome

| Did it Rain |
|---|
| 1 |
| 0 |

$y_i$ = 1, yes
$y_i$ = 0, no

WEATHER

RAIN

$z_1 = (b_1 \times 0.5) + (b_2 \times 0.8) + (b_3 \times 75) + (b_4 \times 1.2) + b_0 \quad y_1 = 1$

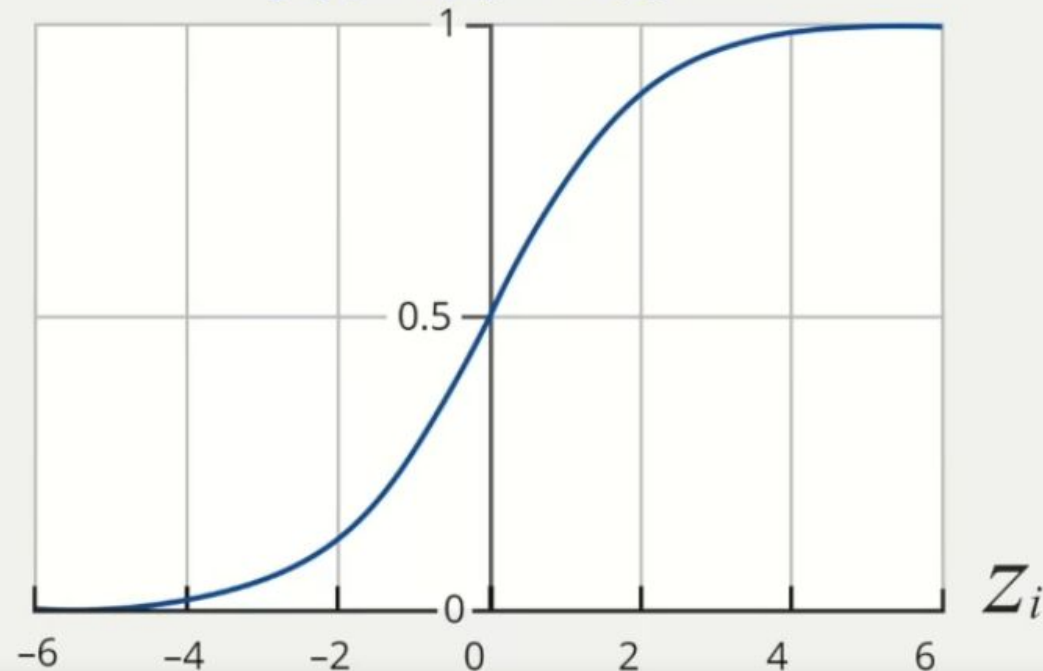$z_2 = (b_1 \times 0.2) + (b_2 \times 0.95) + (b_3 \times 83) + (b_4 \times 1.3) + b_0 \quad y_2 = 0$
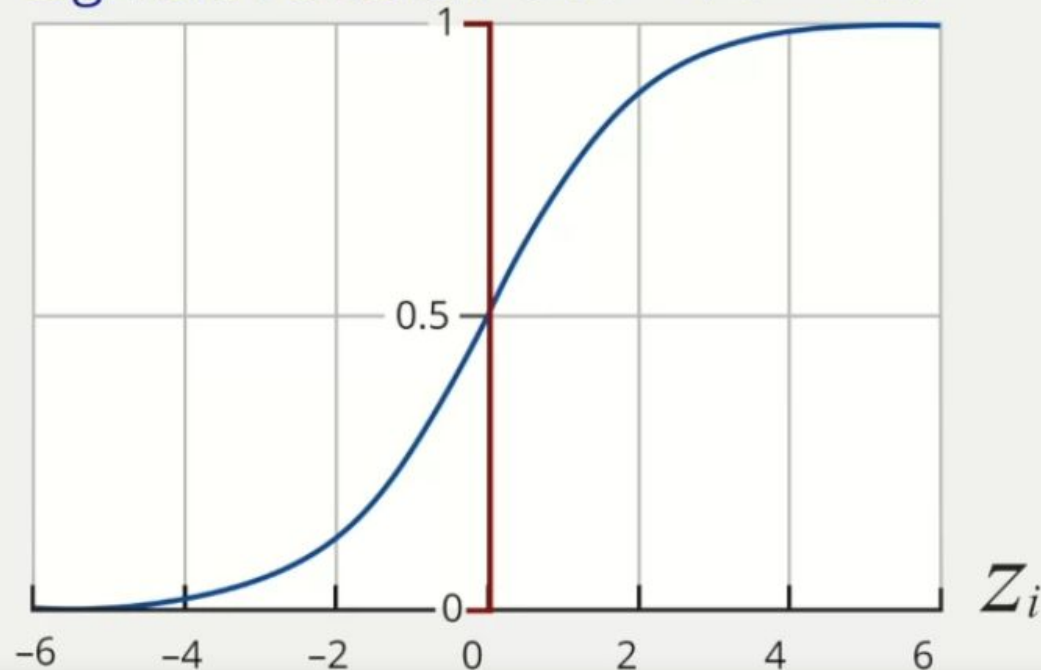
sigma

$p(y_i = 1 | x_i) = \sigma(z_i)$

## Convert to a Probability

$$z_i = (b_1 \times x_{i1}) + (b_2 \times x_{i2}) + \cdots + (b_M \times x_{iM}) + b_0$$
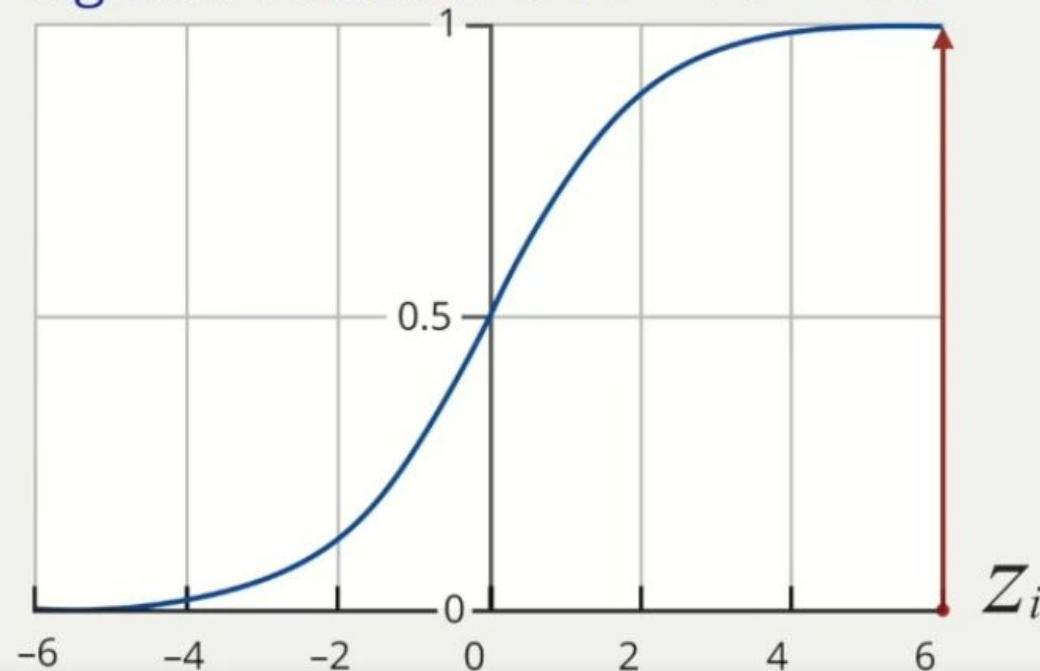
$$p(y_i = 1 | x_i) = \sigma(z_i)$$

## Convert to a Probability

$$z_i = (b_1 \times x_{i1}) + (b_2 \times x_{i2}) + \cdots + (b_M \times x_{iM}) + b_0$$

Sigmoid Function $\quad p(y_i = 1 | x_i) = \sigma(z_i)$

## Convert to a Probability

$$z_i = (b_1 \times x_{i1}) + (b_2 \times x_{i2}) + \cdots + (b_M \times x_{iM}) + b_0$$

Sigmoid Function $\quad p(y_i = 1 | x_i) = \sigma(z_i)$

## Convert to a Probability

$$z_i = (b_1 \times x_{i1}) + (b_2 \times x_{i2}) + \cdots + (b_M \times x_{iM}) + b_0$$
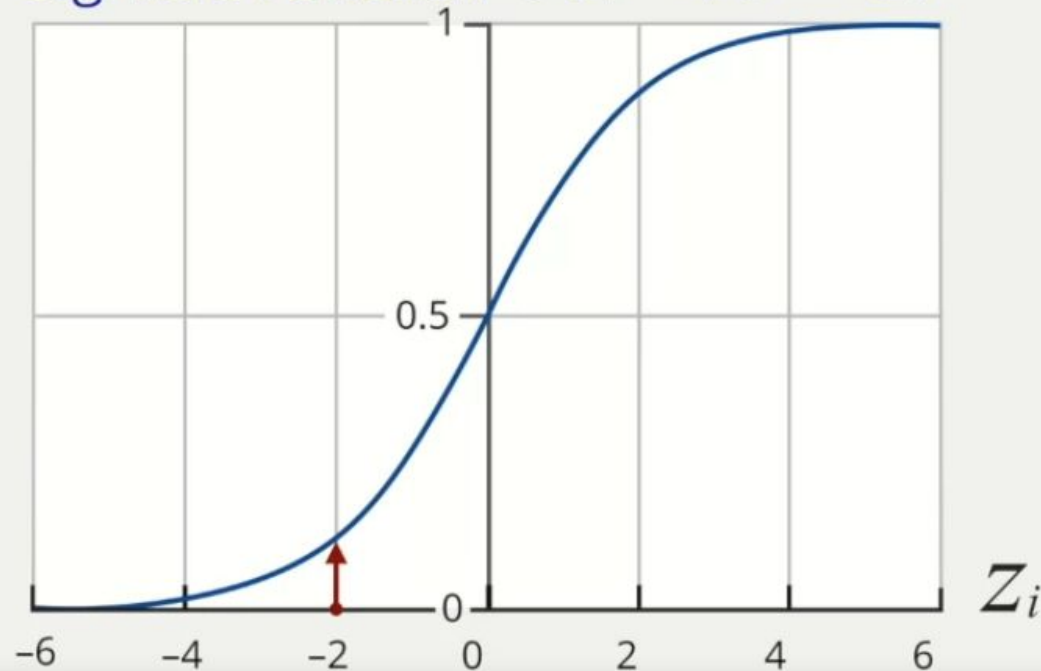
Sigmoid Function $\quad p(y_i = 1 | x_i) = \sigma(z_i)$
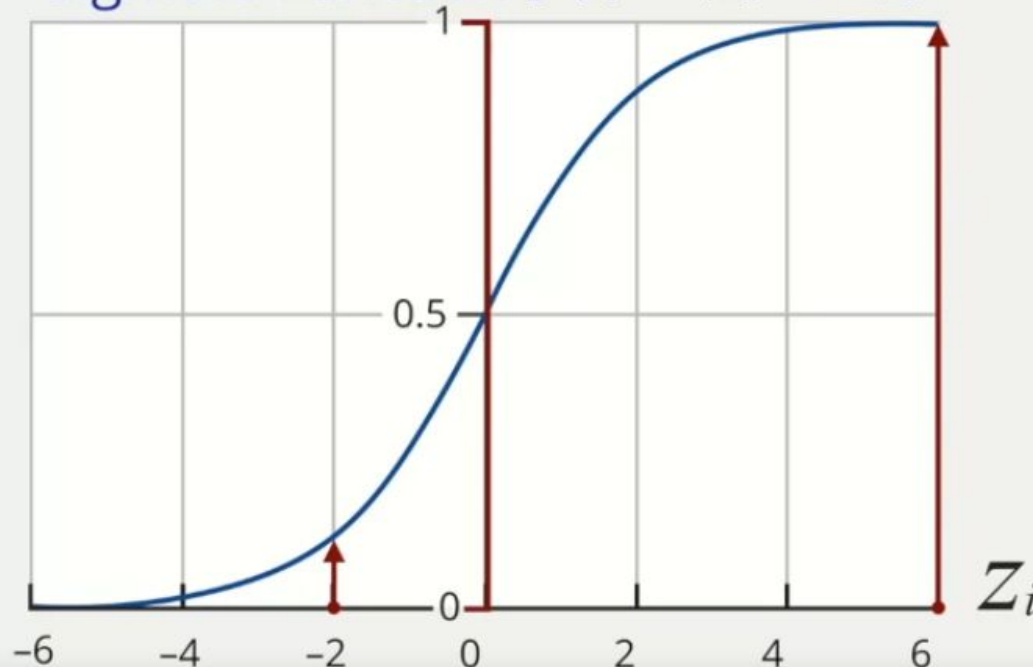
# Basics of Neural Networks

Sigmoid Function is a way to convert predictions to a probabilistic perspective

$$f(x) = \frac{1}{1 + e^{-(x)}}$$

## Convert to a Probability

$$z_i = (b_1 \times x_{i1}) + (b_2 \times x_{i2}) + \cdots + (b_M \times x_{iM}) + b_0$$
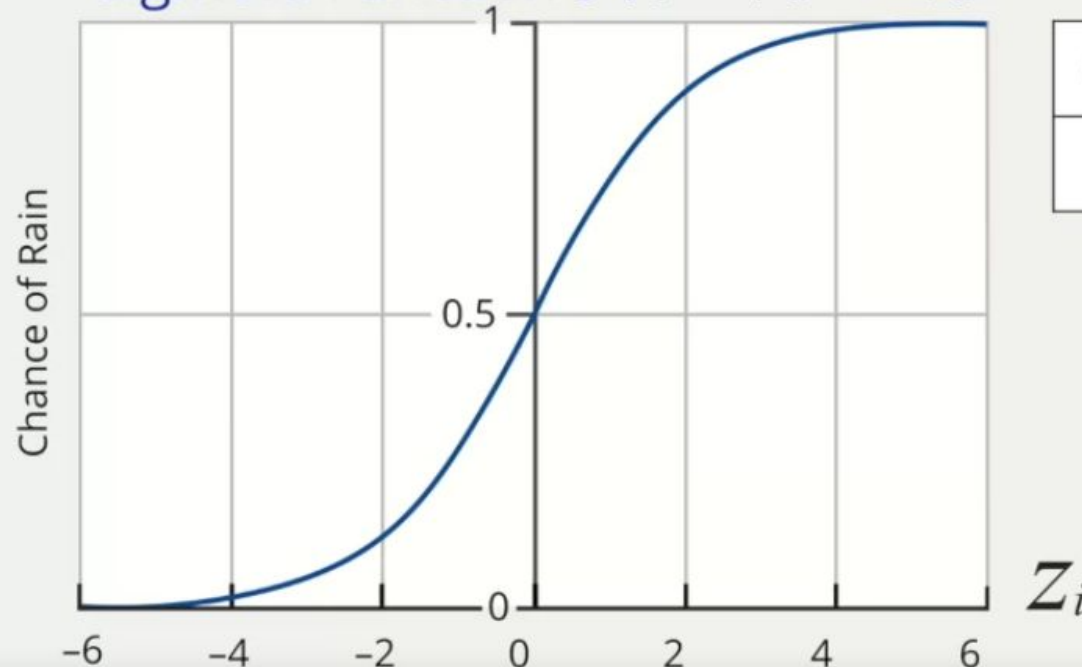
Sigmoid Function $p(y_i = 1 | x_i) = \sigma(z_i)$



Outcome of $Z$

- $Z_i$ = Large and positive indicates $y_i$ = 1 is likely

- $Z_i$ = Large and negative indicates $y_i$ = 0 is likely

## Convert to a Probability

$$z_i = (b_1 \times 0.5) + (b_2 \times 0.8) + (b_3 \times 75) + (b_4 \times 1.2) + b_0$$
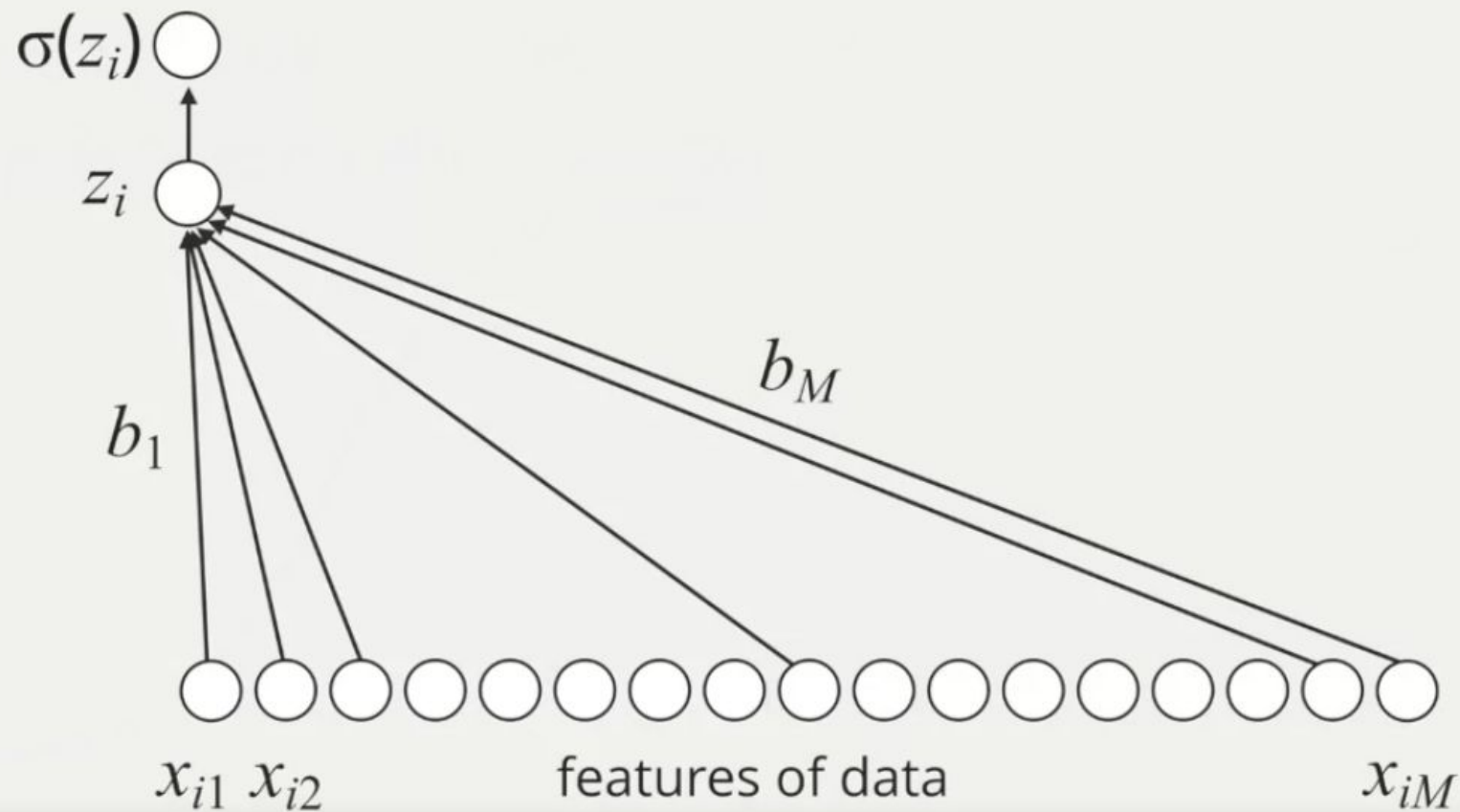
Sigmoid Function $p(y_i = 1 | x_i) = \sigma(z_i)$

features
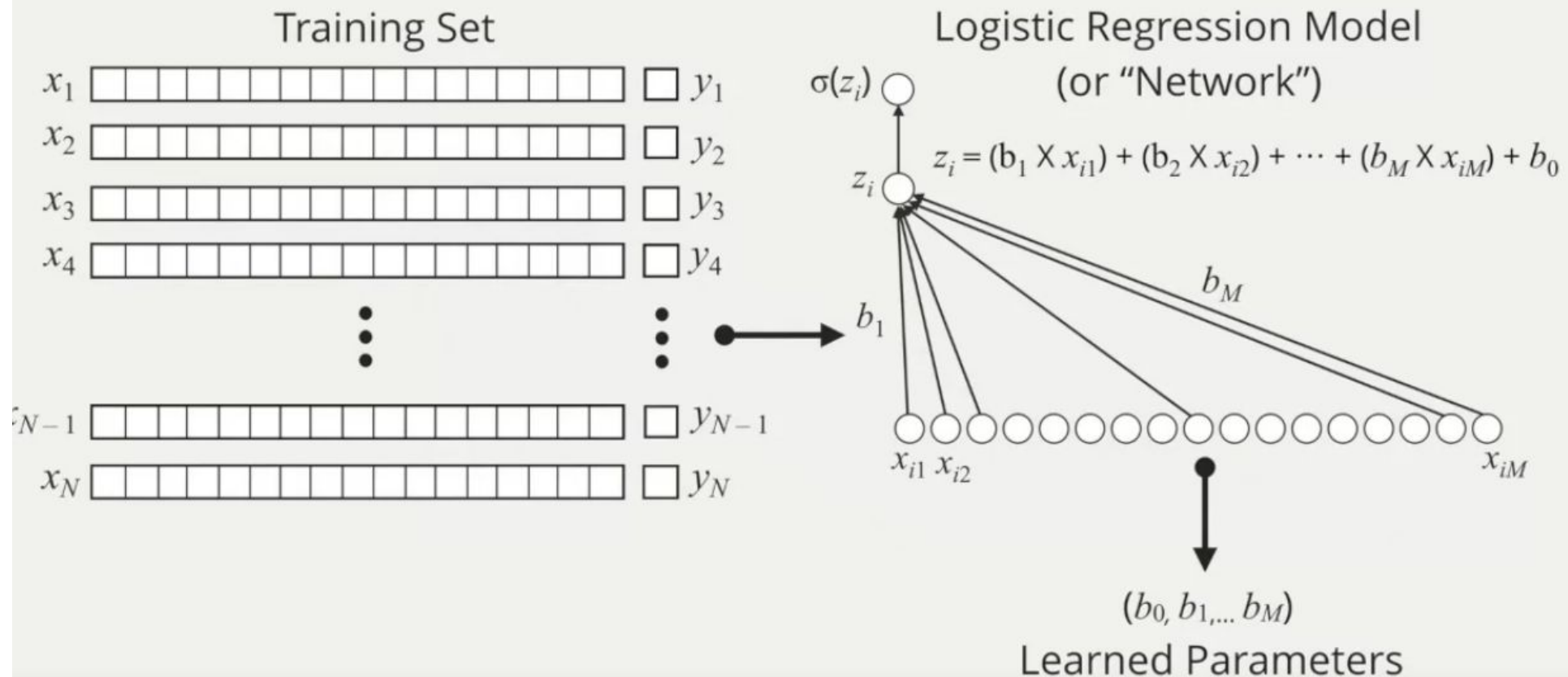
| Cloud Cover | Humidity | Temperature | Air Pressure |
|---|---|---|---|
| 0.5 | 80% | 75 | 1.2 |



Chance of Rain vs $Z_i$

b parameters tell us how important
data variables are to the prediction

Logistic Regression

$\sigma(z_i)$

$z_i$

$b_M$

$b_1$

$x_{i1}$ $x_{i2}$   features of data   $x_{iM}$

## Learned Model Parameters

Training Set

Logistic Regression Model (or "Network")

$\sigma(z_i)$

$z_i = (b_1 \times x_{i1}) + (b_2 \times x_{i2}) + \cdots + (b_M \times x_{iM}) + b_0$
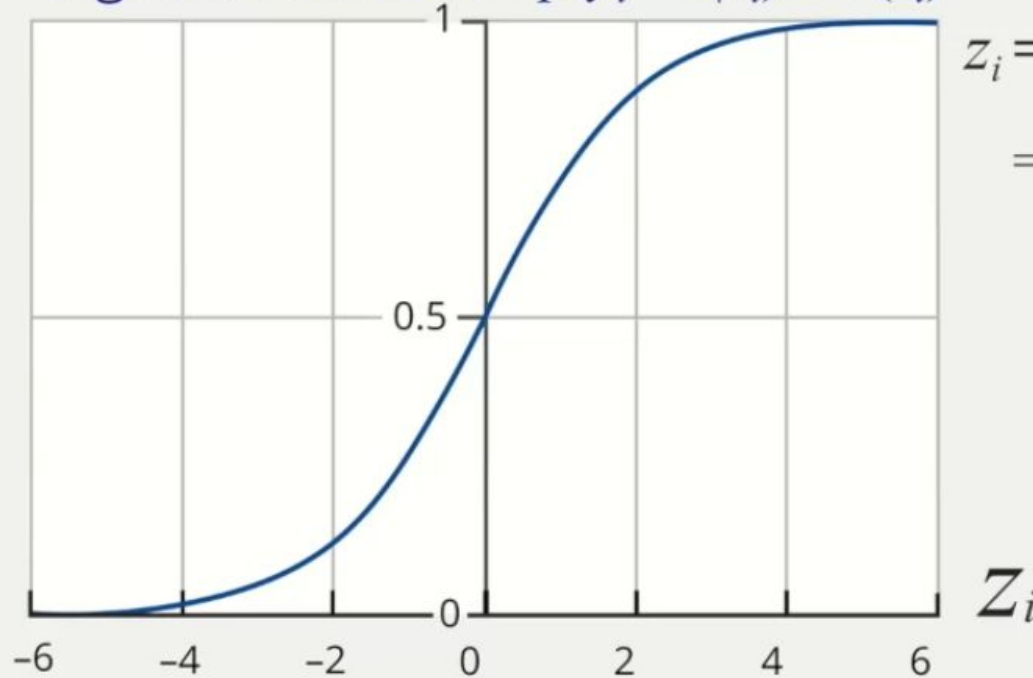
$(b_0, b_1, \ldots b_M)$

Learned Parameters

# Logistic Regression

Sigmoid Function $p(y_i = 1|x_i) = \sigma(z_i)$

$$z_i = (b_1 \times x_{i1}) + (b_2 \times x_{i2}) + \cdots + (b_M \times x_{iM}) + b_0$$
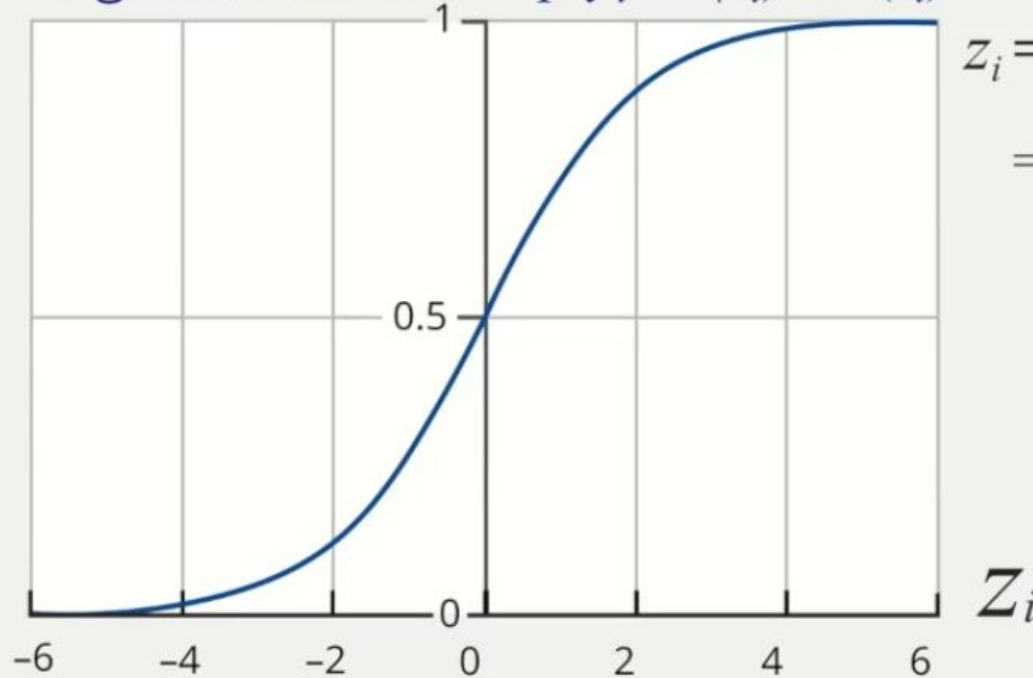
$$= b_0 + x_i \odot b$$

bias    inner product

$Z_i$

## Logistic Regression

Sigmoid Function $p(y_i = 1 | x_i) = \sigma(z_i)$



$$z_i = (b_1 \times x_{i1}) + (b_2 \times x_{i2}) + \cdots + (b_M \times x_{iM}) + b_0$$
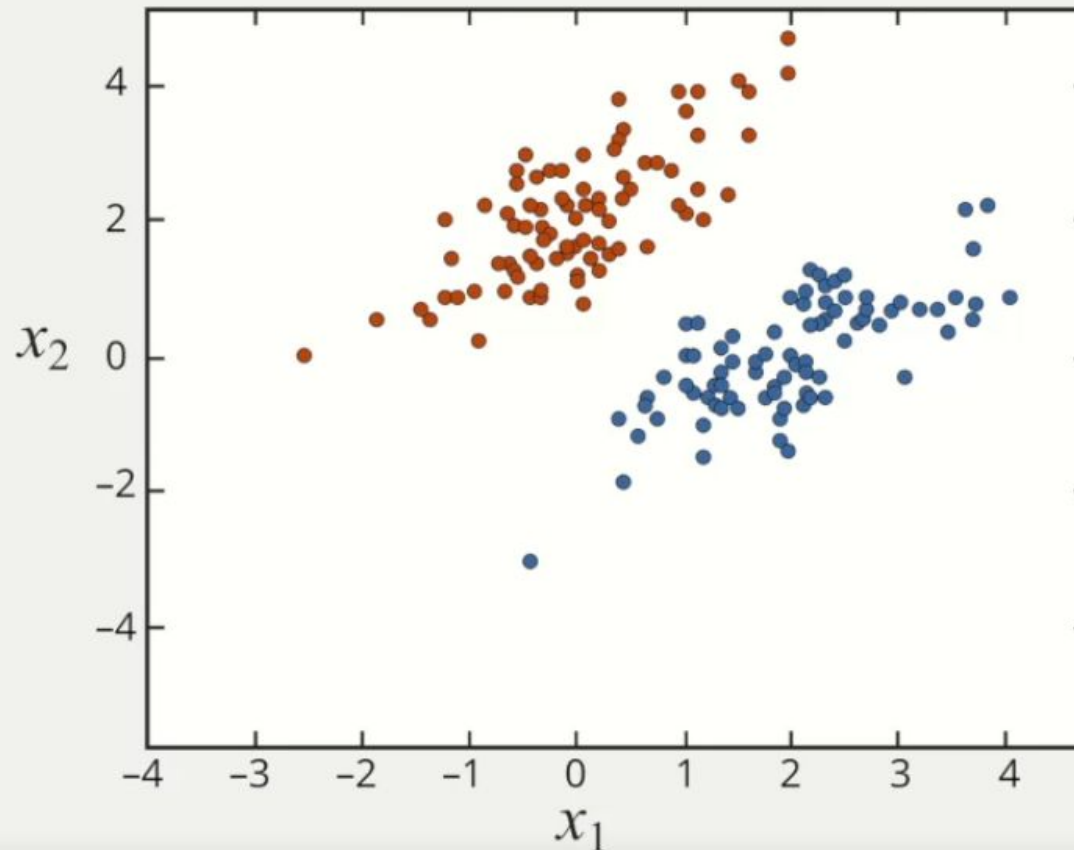
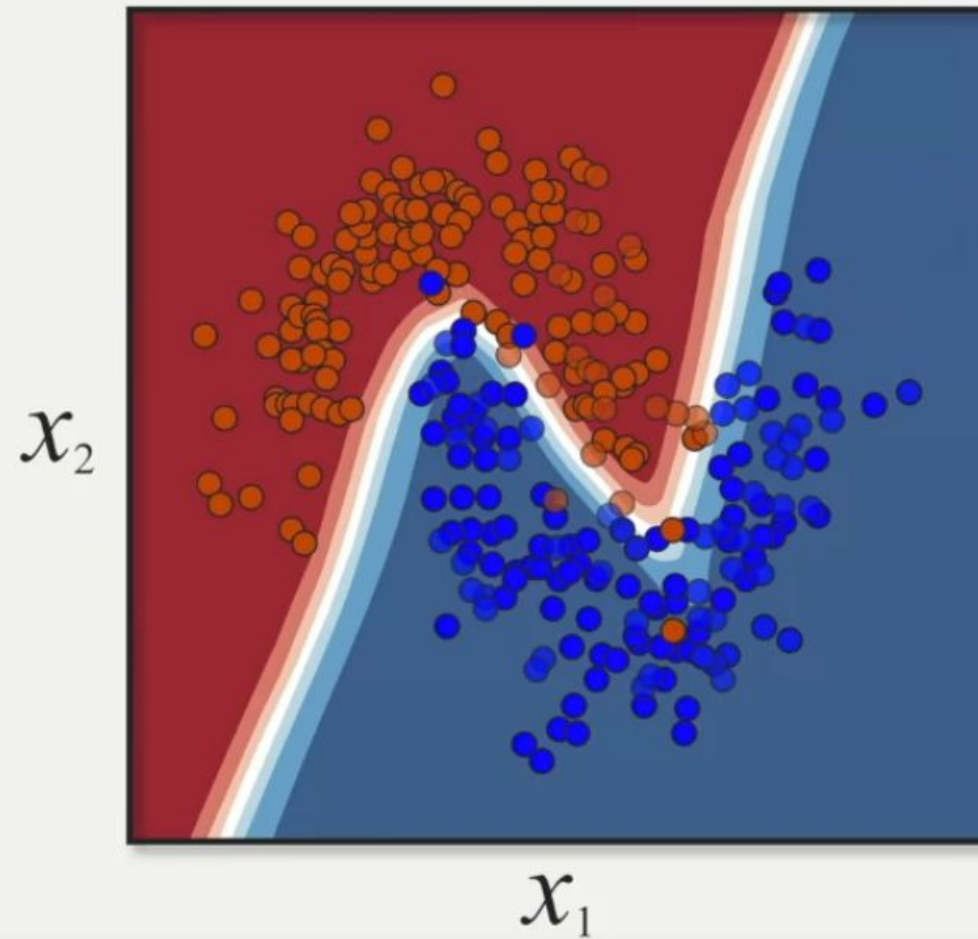$$= b_0 + x_i \odot b$$

bias    inner product
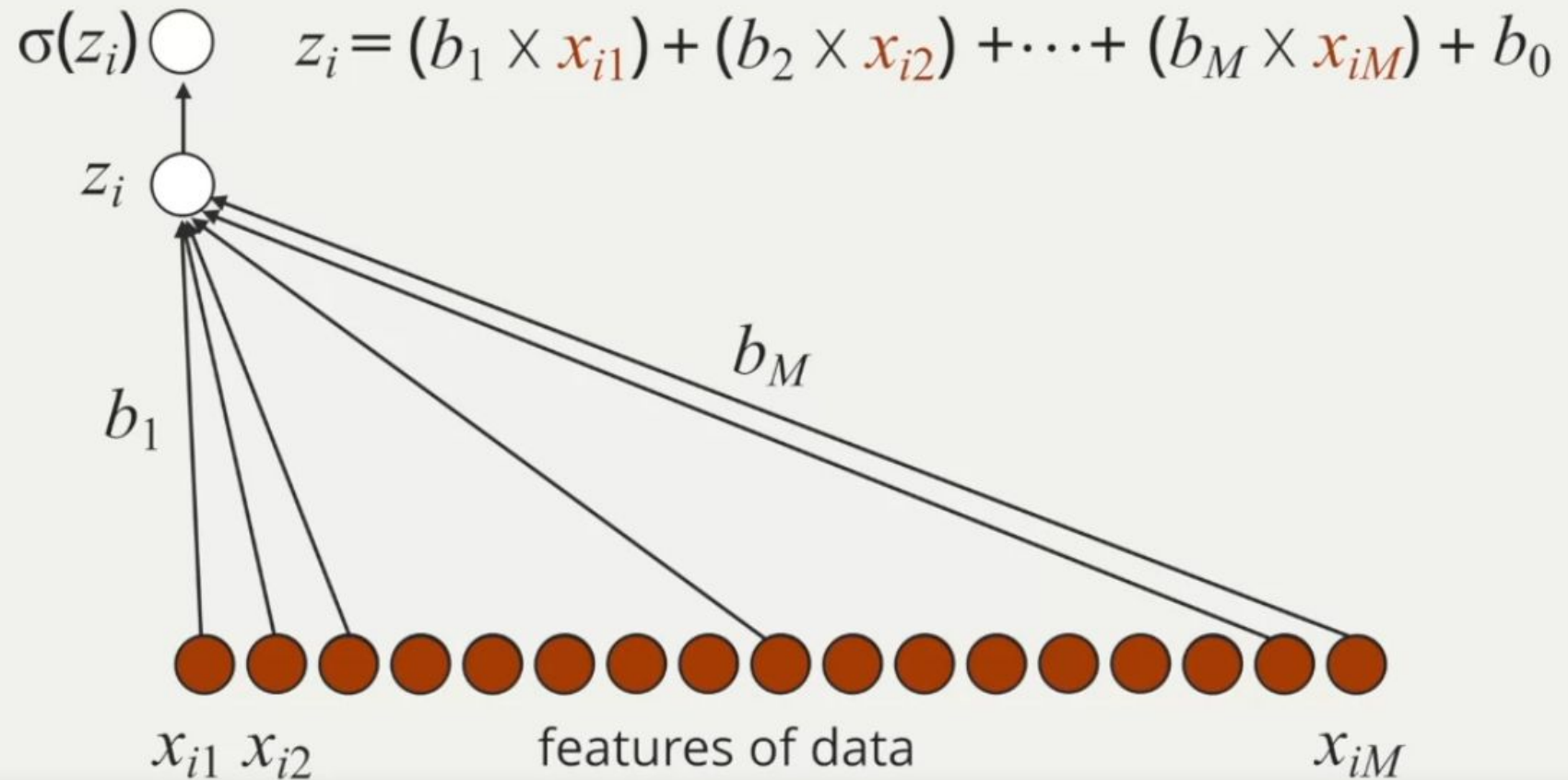
**Limitations of Logistic Regression**

Linear

- Linear classifiers can only represent limited relationships
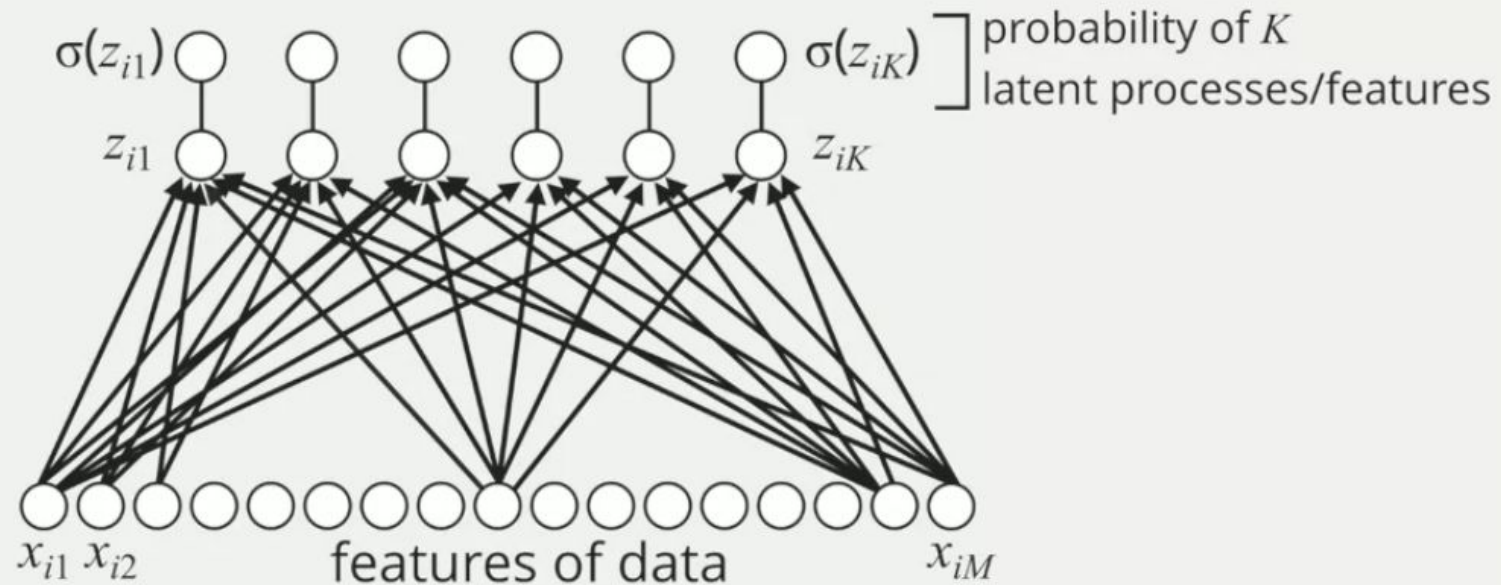- Often want to use a classifier that can handle non-linearities

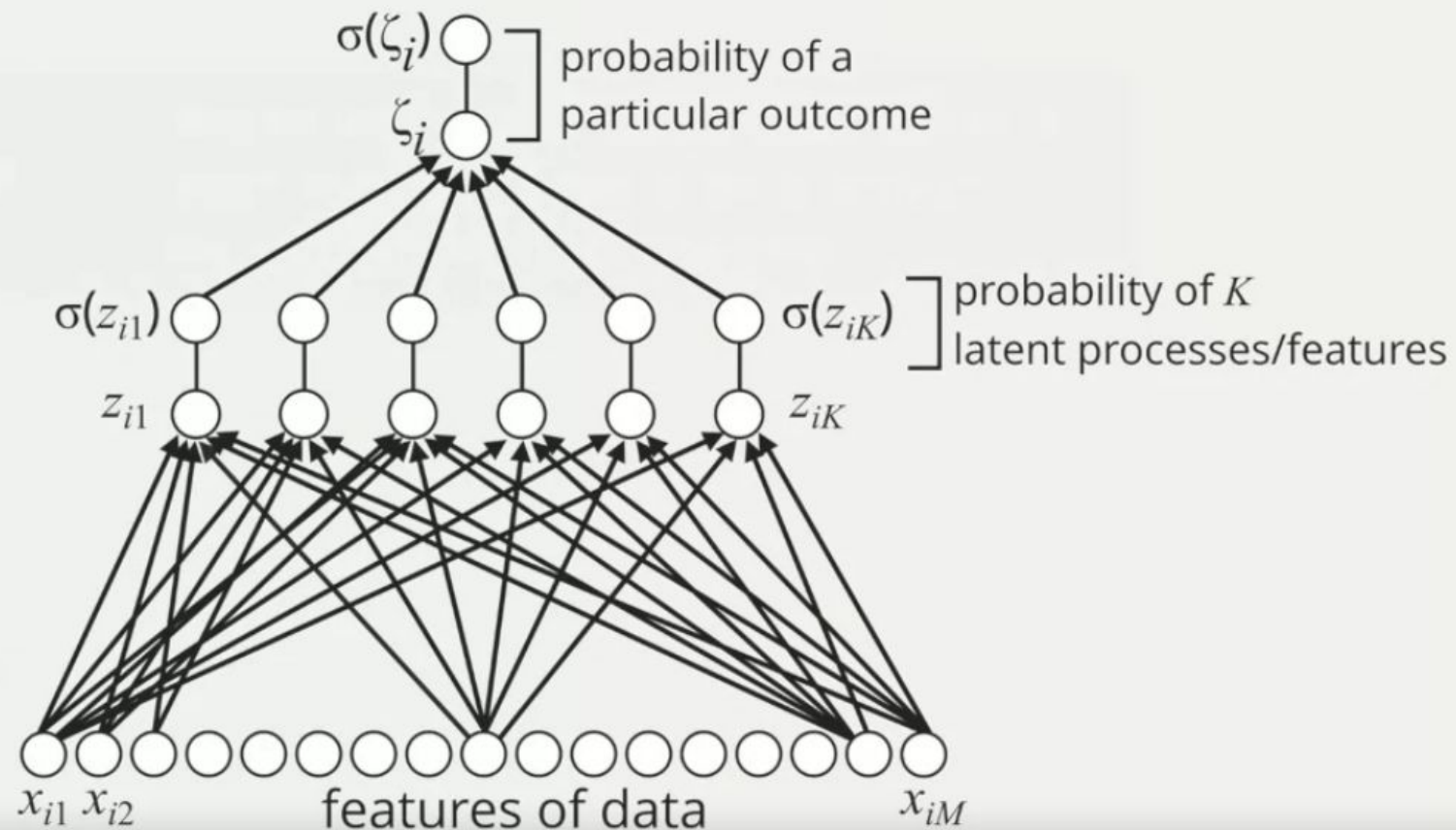**Generalization of Logistic Regression: Learned Features**

**Logistic Regression**
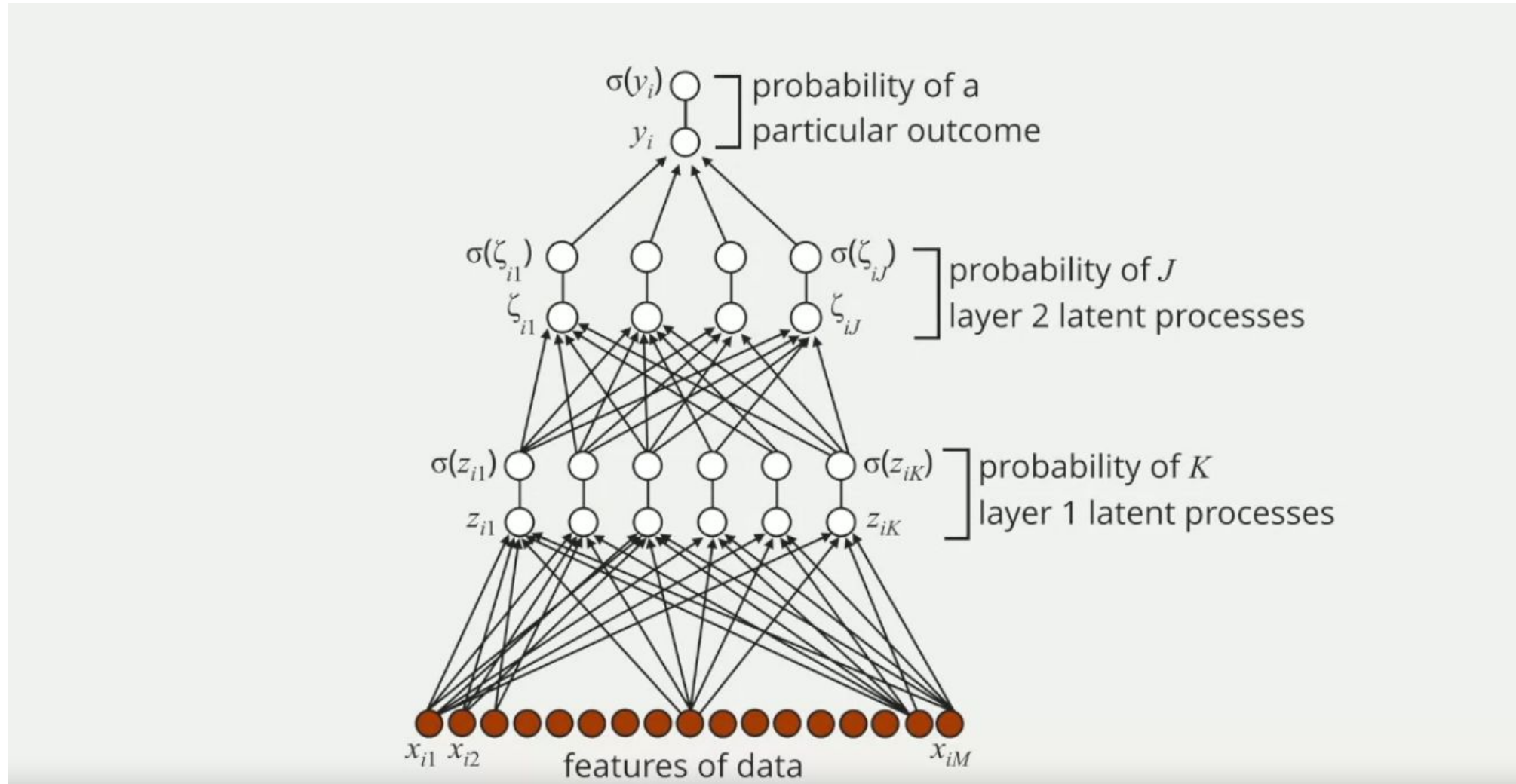
$$\sigma(z_i) \qquad z_i = (b_1 \times x_{i1}) + (b_2 \times x_{i2}) + \cdots + (b_M \times x_{iM}) + b_0$$

$z_i$

$b_1$

$b_M$

$x_{i1} \quad x_{i2}$      features of data      $x_{iM}$

Generalization of Logistic Regression: Learned Features

Extended Logistic Regression

$\sigma(\zeta_i)$ — probability of a particular outcome

$\zeta_i$

$\sigma(z_{i1})$ ... $\sigma(z_{iK})$ — probability of $K$ latent processes/features

$z_{i1}$ ... $z_{iK}$

$x_{i1}$ $x_{i2}$ features of data $x_{iM}$

**Analysis of Documents**

$x_i$ = features for document $i$

features

outcome

$y_i = 1$, like
$y_i = 0$, dislike

| Word 1 | Word 2 | Word 3 | • | • | • | • | Word V |
|--------|--------|--------|---|---|---|---|--------|
| 11 | 20 | 10 | • | • | • | • | 32 |

| Liked/ Disliked |
|-----------------|
| 1 |

number of times each word appears in document

WORDS

LIKED DOCUMENT

Training Set

$x$ = data    $y$ = outcome

$x_1$

$x_2$

$x_3$

$x_4$

$x_{N-1}$

$x_N$

$y_1$

$y_2$

$y_3$

$y_4$

$y_{N-1}$

$y_N$

LIKED DOCUMENT

number of times each word appears in document
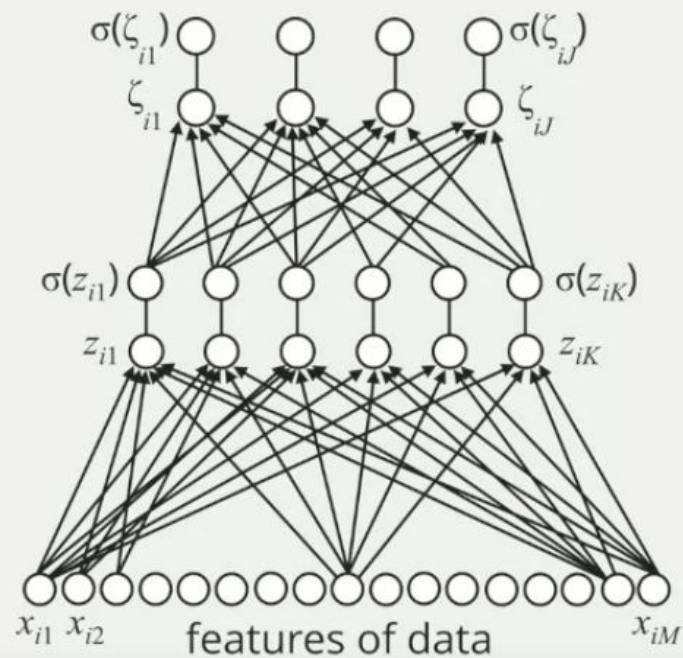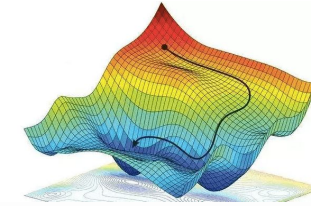
$$z_{i1} = b_{01} + x_i \odot b_1$$

$$z_{i2} = b_{02} + x_i \odot b_2$$

$$\vdots$$

$$z_{iK} = b_{0K} + x_i \odot b_K$$

38

## Gradient Descent

Network parameters $\theta = \{w_1, w_2, \cdots, b_1, b_2, \cdots\}$

Starting Parameters $\theta^0 \longrightarrow \theta^1 \longrightarrow \theta^2 \longrightarrow \cdots\cdots$

$$\nabla L(\theta) = \begin{bmatrix} \partial L(\theta)/\partial w_1 \\ \partial L(\theta)/\partial w_2 \\ \vdots \\ \partial L(\theta)/\partial b_1 \\ \partial L(\theta)/\partial b_2 \\ \vdots \end{bmatrix}$$

Compute $\nabla L(\theta^0)$     $\theta^1 = \theta^0 - \eta \nabla L(\theta^0)$

Compute $\nabla L(\theta^1)$     $\theta^2 = \theta^1 - \eta \nabla L(\theta^1)$

Millions of parameters ......

To compute the gradients efficiently, we use **backpropagation**.

# Gradient Descent

Network parameters $\theta = \{w_1, w_2, \cdots, b_1, b_2, \cdots\}$

Starting Parameters $\theta^0 \longrightarrow \theta^1 \longrightarrow \theta^2 \longrightarrow \cdots\cdots$

$$\nabla L(\theta) = \begin{bmatrix} \partial L(\theta)/\partial w_1 \\ \partial L(\theta)/\partial w_2 \\ \vdots \\ \partial L(\theta)/\partial b_1 \\ \partial L(\theta)/\partial b_2 \\ \vdots \end{bmatrix}$$

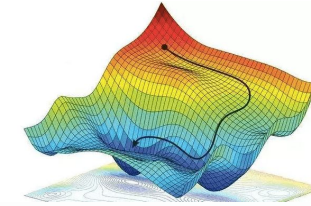Compute $\nabla L(\theta^0)$     $\theta^1 = \theta^0 - \eta \nabla L(\theta^0)$

Compute $\nabla L(\theta^1)$     $\theta^2 = \theta^1 - \eta \nabla L(\theta^1)$

Millions of parameters ……

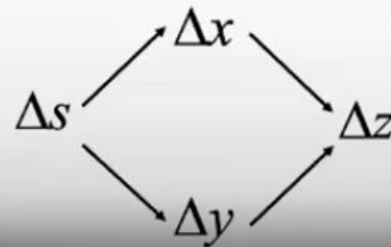To compute the gradients efficiently, we use **_backpropagation_**.

## Chain Rule

**Case 1**     $y = g(x)$     $z = h(y)$

$$\Delta x \rightarrow \Delta y \rightarrow \Delta z \qquad \frac{dz}{dx} = \frac{dz}{dy}\frac{dy}{dx}$$

**Case 2**

$$x = g(s) \qquad y = h(s) \qquad z = k(x, y)$$

$$\frac{dz}{ds} = \frac{\partial z}{\partial x}\frac{dx}{ds} + \frac{\partial z}{\partial y}\frac{dy}{ds}$$

(diagram: $\Delta s \nearrow \Delta x \searrow \Delta z$ ; $\Delta s \searrow \Delta y \nearrow \Delta z$)

$$\frac{\partial E_{total}}{\partial w_1} = \frac{\partial E_{total}}{\partial out_{h1}} * \frac{\partial out_{h1}}{\partial net_{h1}} * \frac{\partial net_{h1}}{\partial w_1}$$

$$\frac{\partial E_{total}}{\partial out_{h1}} = \frac{\partial E_{o1}}{\partial out_{h1}} + \frac{\partial E_{o2}}{\partial out_{h1}}$$



$$E_{total} = E_{o1} + E_{o2}$$