



# LOGISTIC REGRESSION

CS6140

Predrag Radivojac

KHOURY COLLEGE OF COMPUTER SCIENCES

NORTHEASTERN UNIVERSITY

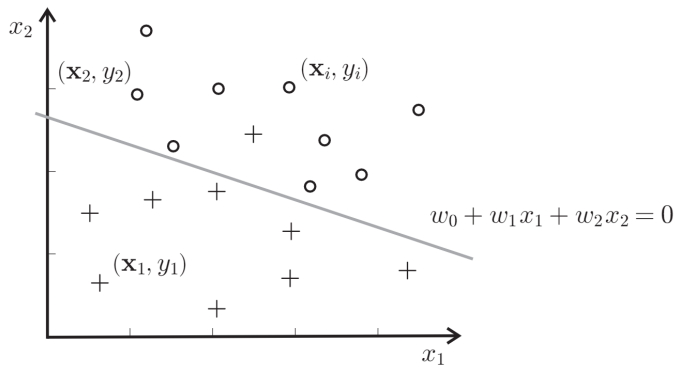
Spring 2021

# LINEAR CLASSIFICATION

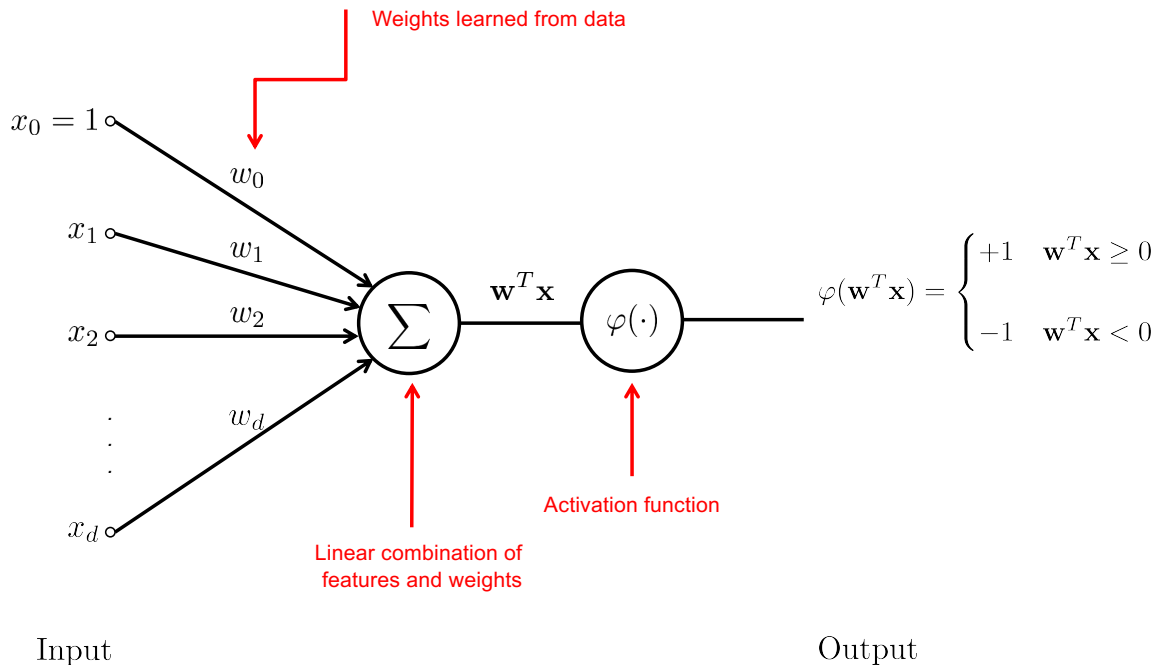
**Given:** a set of observations  $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$ ,  $(\mathbf{x}_i, y_i) \in \mathbb{R}^d \times \{0, 1\}$

**Objective:** find best linear separator  $f(\mathbf{x}) = 0$ , where  $f(\mathbf{x}) = w_0 + \sum_{j=1}^d w_j x_j$

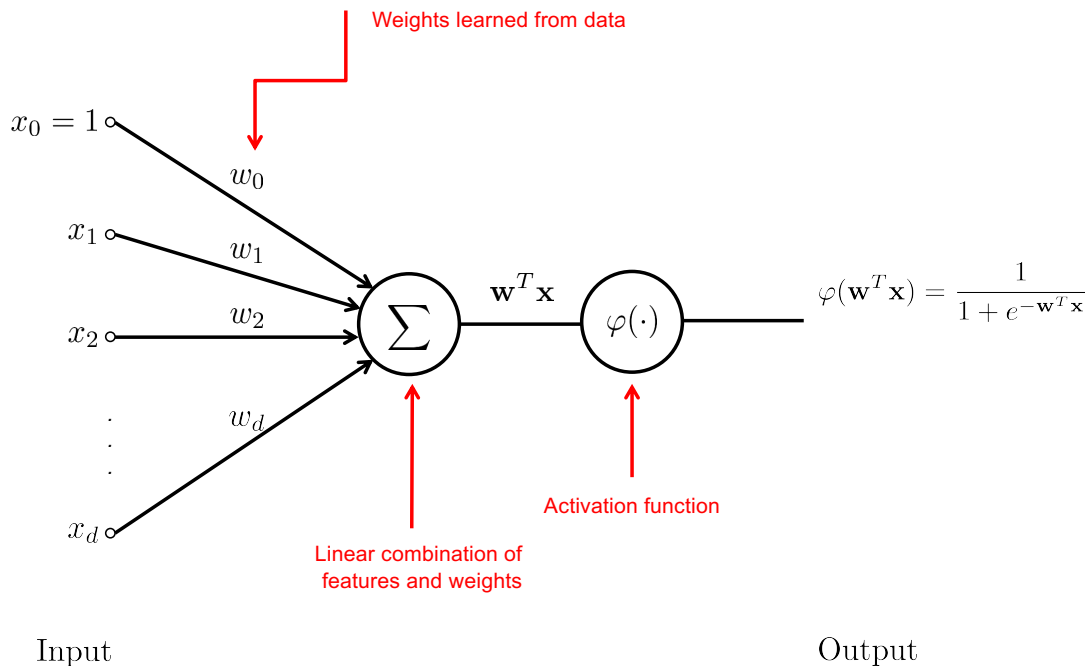
$$\mathcal{X} = \mathbb{R} \times \mathbb{R}, \mathcal{Y} = \{0, 1\}$$



# PERCEPTRON

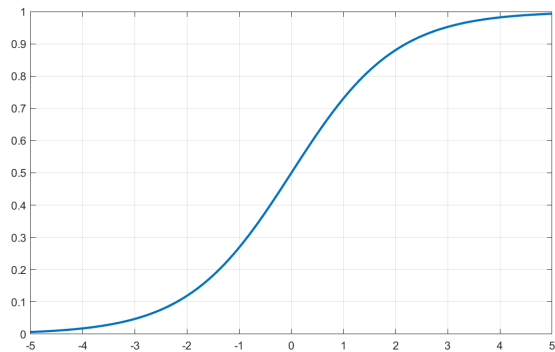


# LOGISTIC REGRESSION



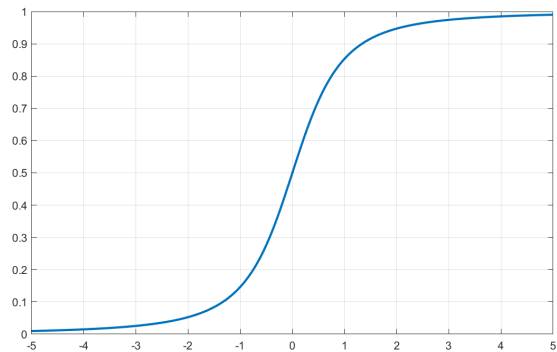
# ACTIVATION FUNCTIONS

$$\varphi(t) = \frac{1}{1 + e^{-t}} \quad \leftarrow \text{sigmoid function}$$



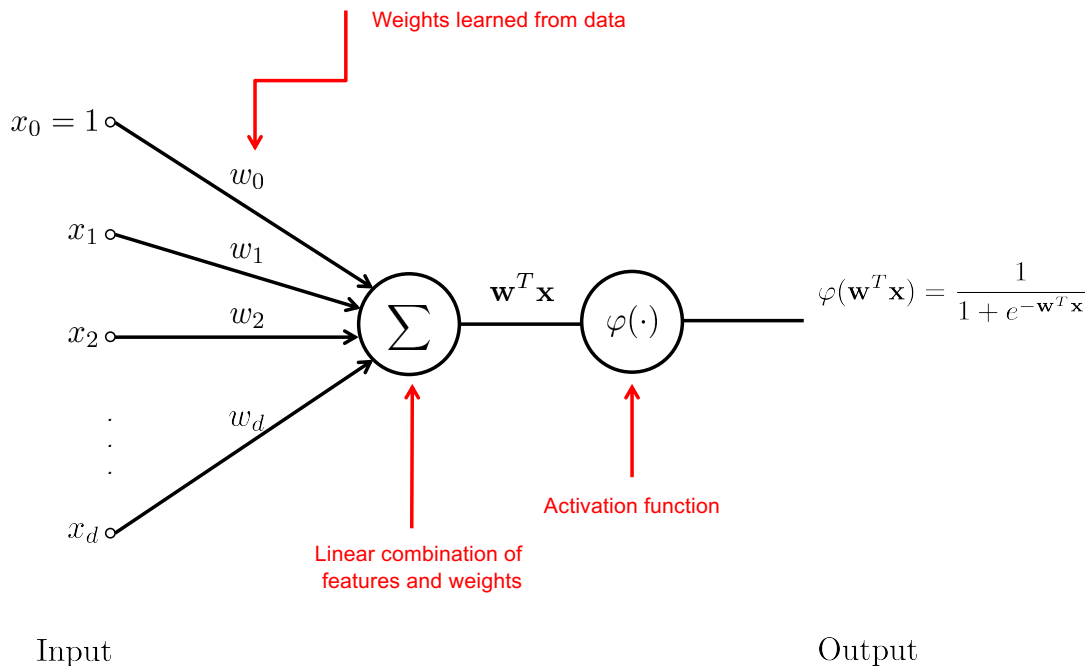
$t$

$$\varphi(t) = \frac{1}{2} \cdot \left( 1 + \frac{t}{\sqrt{1 + t^2}} \right)$$



$t$

# LOGISTIC REGRESSION



# THINKING ABOUT LOGISTIC REGRESSION

Can we model the  $(X, Y)$  dependence using a linear combination?

$$\mathcal{X} = \{1\} \times \mathbb{R}^d, \mathcal{Y} = \{0, 1\}$$

$$P(Y = 1|\mathbf{x}) = \mathbf{w}^T \mathbf{x}$$

How about the odds function?

$$\underbrace{\frac{P(Y = 1|\mathbf{x})}{1 - P(Y = 1|\mathbf{x})}}_{\text{odds function}} = \mathbf{w}^T \mathbf{x}$$

How about the log odds function?

$$\log \frac{P(Y = 1|\mathbf{x})}{1 - P(Y = 1|\mathbf{x})} = \mathbf{w}^T \mathbf{x}$$

# MAXIMIZING LIKELIHOOD

$$P(Y = 1|\mathbf{x}, \mathbf{w}) = \frac{1}{1+e^{-\mathbf{w}^T \mathbf{x}}}$$

$$\mathcal{X} = \{1\} \times \mathbb{R}^d, \mathcal{Y} = \{0, 1\}$$

Let's express the probability mass function as:

$$p(y|\mathbf{x}, \mathbf{w}) = \begin{cases} \left( \frac{1}{1+e^{-\mathbf{w}^T \mathbf{x}}} \right)^y & \text{for } y = 1 \\ \left( 1 - \frac{1}{1+e^{-\mathbf{w}^T \mathbf{x}}} \right)^{1-y} & \text{for } y = 0 \end{cases}$$

And make it more compact:

$$p(y|\mathbf{x}, \mathbf{w}) = \left( \frac{1}{1 + e^{-\mathbf{w}^T \mathbf{x}}} \right)^y \left( 1 - \frac{1}{1 + e^{-\mathbf{w}^T \mathbf{x}}} \right)^{1-y}$$



# MAXIMIZING LIKELIHOOD

**Given:** a set of observations  $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$ ,  $(\mathbf{x}_i, y_i) \in \mathcal{X} \times \mathcal{Y}$   $\mathcal{X} = \{1\} \times \mathbb{R}^d$ ,  $\mathcal{Y} = \{0, 1\}$

Let's express the conditional likelihood as:

$$l(\mathbf{w}) = \prod_{i=1}^n p(y_i | \mathbf{x}_i, \mathbf{w}) \quad \longrightarrow \quad \mathbf{w}^* = \arg \max_{\mathbf{w}} \left\{ \prod_{i=1}^n p(y_i | \mathbf{x}_i, \mathbf{w}) \right\}$$

# UPDATE RULES TO MAXIMIZE LIKELIHOOD

$\mathbf{w}^{(0)}$  = something

$$\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} + \Delta \mathbf{w}^{(t)}$$

**Maximum Likelihood:**

$$\mathbf{w}^{(0)} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

$$\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} + \eta \left( \mathbf{X}^T \mathbf{P}^{(t)} \left( \mathbf{I} - \mathbf{P}^{(t)} \right) \mathbf{X} \right)^{-1} \mathbf{X}^T \left( \mathbf{y} - \mathbf{p}^{(t)} \right)$$

$$\mathbf{p} = (p_1, p_2, \dots, p_n)$$

$$\mathbf{P} = \text{diag} \{ \mathbf{p} \}$$

$$p_i = P(Y_i = 1 | \mathbf{x}_i, \mathbf{w}) = \frac{1}{1 + e^{-\mathbf{w}^T \mathbf{x}_i}}$$

$$\mathbf{y} = (y_1, y_2, \dots, y_n)$$

$\mathbf{I}$  = identity matrix

# MINIMIZING SUM OF SQUARED ERRORS

**Given:** a set of observations  $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$ ,  $(\mathbf{x}_i, y_i) \in \mathcal{X} \times \mathcal{Y}$   $\mathcal{X} = \{1\} \times \mathbb{R}^d$ ,  $\mathcal{Y} = \{0, 1\}$

Let's express the sum of squared errors as:

$$E(\mathbf{w}) = \sum_{i=1}^n (y_i - p_i)^2 = \sum_{i=1}^n e_i^2$$

$$p_i = P(Y_i = 1 | \mathbf{x}_i, \mathbf{w}) = \frac{1}{1 + e^{-\mathbf{w}^T \mathbf{x}_i}}$$

# UPDATE RULES TO MINIMIZE SUM OF SQUARED ERRORS

$$\mathbf{w}^{(0)} = \text{something}$$

$$\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} + \Delta \mathbf{w}^{(t)}$$

Minimum Sum of Squared Errors:

$$\mathbf{w}^{(0)} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

$$\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} + \eta \left( \mathbf{J}^{(t)T} \mathbf{J}^{(t)} + \mathbf{J}^{(t)T} \mathbf{E}^{(t)} (2\mathbf{P}^{(t)} - \mathbf{I}) \mathbf{X} \right)^{-1} \mathbf{X}^T \mathbf{P}^{(t)} (\mathbf{I} - \mathbf{P}^{(t)}) (\mathbf{y} - \mathbf{p}^{(t)})$$

$$\mathbf{p} = (p_1, p_2, \dots, p_n)$$

$$\mathbf{P} = \text{diag} \{ \mathbf{p} \}$$

$$\mathbf{e} = (y_1 - p_1, y_2 - p_2, \dots, y_n - p_n)$$

$$\mathbf{y} = (y_1, y_2, \dots, y_n)$$

$$\mathbf{J} = \mathbf{P}(\mathbf{I} - \mathbf{P})\mathbf{X}$$

$$\mathbf{E} = \text{diag} \{ \mathbf{e} \}$$

$\mathbf{I}$  = identity matrix

## UPDATE RULES, BATCH MODE

$$\mathbf{w}^{(0)} = \text{something}$$

$$\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} + \Delta \mathbf{w}^{(t)}$$

**Maximum Likelihood:**

$$\mathbf{w}^{(0)} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

$$\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} + \eta \left( \mathbf{X}^T \mathbf{P}^{(t)} \left( \mathbf{I} - \mathbf{P}^{(t)} \right) \mathbf{X} \right)^{-1} \mathbf{X}^T \left( \mathbf{y} - \mathbf{p}^{(t)} \right)$$

**Minimum Sum of Squared Errors:**

$$\mathbf{w}^{(0)} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

$$\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} + \eta \left( \mathbf{J}^{(t)T} \mathbf{J}^{(t)} + \mathbf{J}^{(t)T} \mathbf{E}^{(t)} (2\mathbf{P}^{(t)} - \mathbf{I}) \mathbf{X} \right)^{-1} \mathbf{X}^T \mathbf{P}^{(t)} \left( \mathbf{I} - \mathbf{P}^{(t)} \right) \left( \mathbf{y} - \mathbf{p}^{(t)} \right)$$

$$\mathbf{p} = (p_1, p_2, \dots, p_n)$$

$$\mathbf{P} = \text{diag} \{ \mathbf{p} \}$$

$$\mathbf{e} = (y_1 - p_1, y_2 - p_2, \dots, y_n - p_n)$$

$$\mathbf{y} = (y_1, y_2, \dots, y_n)$$

$$\mathbf{J} = \mathbf{P}(\mathbf{I} - \mathbf{P})\mathbf{X}$$

$$\mathbf{E} = \text{diag} \{ \mathbf{e} \}$$

$\mathbf{I}$  = identity matrix

## UPDATE RULES, BATCH MODE

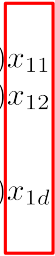
$\mathbf{w} \leftarrow \text{something}$

$\mathbf{w} \leftarrow \mathbf{w} + \Delta \mathbf{w}$

**Gradient Descent:**

$$\Delta \mathbf{w} = \eta \mathbf{X}^T \mathbf{P} (\mathbf{I} - \mathbf{P}) (\mathbf{y} - \mathbf{p})$$

$$\Delta \mathbf{w} = \eta \cdot \begin{bmatrix} p_1(1-p_1)x_{11} & p_2(1-p_2)x_{21} & \cdots & p_n(1-p_n)x_{n1} \\ p_1(1-p_1)x_{12} & p_2(1-p_2)x_{22} & & \\ \vdots & & \ddots & \\ p_1(1-p_1)x_{1d} & & & p_n(1-p_n)x_{nd} \end{bmatrix} \cdot \begin{bmatrix} e_1 \\ e_2 \\ \vdots \\ e_n \end{bmatrix}$$

  $\mathbf{x}_1$

$$\Delta \mathbf{w} = \eta \sum_{i=1}^n e_i p_i (1 - p_i) \mathbf{x}_i,$$

## UPDATE RULES, STOCHASTIC MODE

$\mathbf{w} \leftarrow \text{something}$

$\mathbf{w} \leftarrow \mathbf{w} + \Delta \mathbf{w}$

**Gradient Descent:**

$$\Delta \mathbf{w} = \eta \sum_{i=1}^n e_i p_i (1 - p_i) \mathbf{x}_i$$

**Stochastic Gradient Descent:**

$$\Delta \mathbf{w} = \eta e_i p_i (1 - p_i) \mathbf{x}_i$$

## PREDICTION

**Given:** a predictor defined by  $\mathbf{w}$  and a new data point  $\mathbf{x}$

$$P(Y = 1|\mathbf{x}, \mathbf{w}) = \frac{1}{1+e^{-\mathbf{w}^T \mathbf{x}}}$$

$$r = \frac{w_0 + \sum_{j=1}^d w_j x_j}{\sqrt{\sum_{j=1}^d w_j^2}} \quad \longleftarrow \quad \text{distance from } \mathbf{x} \text{ to hyperplane}$$

$$P(Y = 1|\mathbf{x}) = \frac{1}{1+e^{-r \cdot \sqrt{\sum_{j=1}^d w_j^2}}}$$