

# Machine learning for signal processing [5LSL0]

Ruud van Sloun, Rik Vullings

# Activation functions:

From optimal linear filtering to nonlinear classification

Weight vector of  $n$  samples

$$\mathbf{w} = [w_1, w_2, \dots, w_n]^T.$$

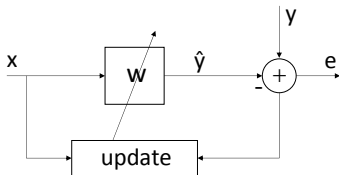
$m$  observations consisting of data vectors  $\mathbf{x}$  and outputs  $\mathbf{y}$   
collected in an  $m \times n$  input data matrix:

$$\mathbf{X} = [\mathbf{x}^{(0)}, \mathbf{x}^{(1)}, \dots, \mathbf{x}^{(i)}, \dots, \mathbf{x}^{(m)}]^T.$$

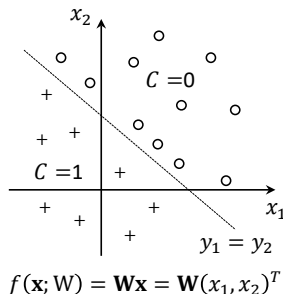
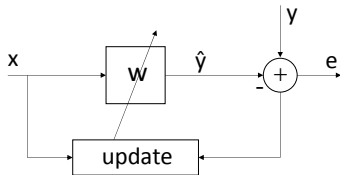
$$\mathbf{x}^{(i)} = [x_1^{(i)}, x_2^{(i)}, \dots, x_n^{(i)}]^T, \text{ with associated output } \mathbf{y}^{(i)} \text{ or } y^{(i)}.$$

Set of input data vectors:  $\mathbb{X}$

So far: optimal linear operations given some cost criterion.

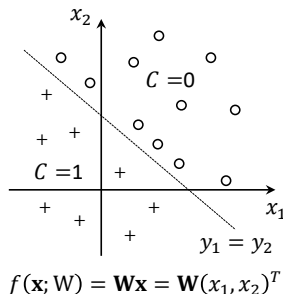
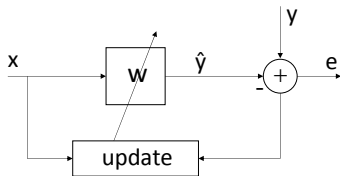


So far: optimal linear operations given some cost criterion.



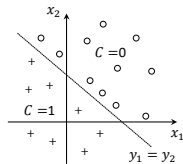
Same framework also enables classification:

So far: optimal linear operations given some cost criterion.



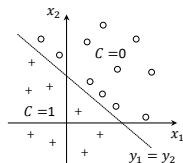
Same framework also enables classification:  
[ $y_1, y_2$ ] =  $\mathbf{W}\mathbf{x}$ . Given  $\mathbf{y}$ , classification through:

$$C(\mathbf{x}) = \begin{cases} 1 & y_2 > y_1 \\ 0 & \text{else} \end{cases} \quad (1)$$



$$f(\mathbf{x}; \mathbf{W}) = \mathbf{W}\mathbf{x} = \mathbf{W}(x_1, x_2)^T$$

Regression problem through MSE cost criterion:



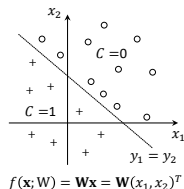
$$f(\mathbf{x}; \mathbf{W}) = \mathbf{W}\mathbf{x} = \mathbf{W}(x_1, x_2)^T$$

Regression problem through MSE cost criterion:

$$J(\theta) = \frac{1}{N} \sum_{\mathbf{x} \in \mathbb{X}} (f^*(\mathbf{x}) - f(\mathbf{x}; \theta))^2,$$

where  $f^*(\mathbf{x})$  is some known/target output on  $\mathbf{x} \in \mathbb{X}$ , and  
 $f(\mathbf{x}; \theta) = f(\mathbf{x}; \mathbf{W}) = \mathbf{W}^T \mathbf{x}$



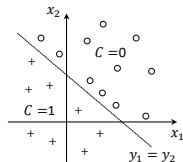


Regression problem through MSE cost criterion:

$$J(\theta) = \frac{1}{N} \sum_{\mathbf{x} \in \mathbb{X}} (f^*(\mathbf{x}) - f(\mathbf{x}; \theta))^2,$$

where  $f^*(\mathbf{x})$  is some known/target output on  $\mathbf{x} \in \mathbb{X}$ , and  
 $f(\mathbf{x}; \theta) = f(\mathbf{x}; \mathbf{W}) = \mathbf{W}^T \mathbf{x}$

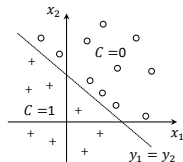
As before, we can minimize  $J(\theta)$  w.r.t. the coefficients in matrix  $\mathbf{W}$ .



$$f(\mathbf{x}; \mathbf{W}) = \mathbf{W}\mathbf{x} = \mathbf{W}(x_1, x_2)^T$$

Regression problem through MSE cost criterion:

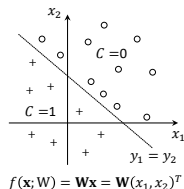
$$J(\theta) = \frac{1}{N} \sum_{\mathbf{x} \in \mathbb{X}} (f^*(\mathbf{x}) - f(\mathbf{x}; \theta))^2$$



$$f(\mathbf{x}; \mathbf{W}) = \mathbf{W}\mathbf{x} = \mathbf{W}(x_1, x_2)^T$$

Minimum MSE cost criterion for vector  $\mathbf{w}$ :

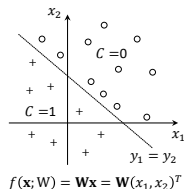
$$J(\theta) = \sum_{\mathbf{x} \in \mathbb{X}} (f^*(\mathbf{x}) - f(\mathbf{x}; \theta))^2 \quad // \text{ omitting scalar}$$



Minimum MSE cost criterion for vector  $\mathbf{w}$ :

$$J(\theta) = \sum_{\mathbf{x} \in \mathbb{X}} (f^*(\mathbf{x}) - f(\mathbf{x}; \theta))^2 \quad // \text{ omitting scalar}$$

$$J(\mathbf{w}) = (\mathbf{w}^T \mathbf{X} - \mathbf{y})(\mathbf{w}^T \mathbf{X} - \mathbf{y})^T \quad // \text{ rewriting in matrix form}$$
$$\mathbf{X} = [\mathbf{x}^{(0)}, \mathbf{x}^{(1)}, \dots, \mathbf{x}^{(i)}, \dots, \mathbf{x}^{(m)}]$$
$$\mathbf{y} = [y^{(0)}, y^{(1)}, \dots, y^{(i)}, \dots, y^{(m)}]$$



Minimum MSE cost criterion for vector  $\mathbf{w}$ :

$$J(\theta) = \sum_{\mathbf{x} \in \mathbb{X}} (f^*(\mathbf{x}) - f(\mathbf{x}; \theta))^2 \quad // \text{ omitting scalar}$$

$$J(\mathbf{w}) = (\mathbf{w}^T \mathbf{X} - \mathbf{y})(\mathbf{w}^T \mathbf{X} - \mathbf{y})^T \quad // \text{ rewriting in matrix form}$$
$$\mathbf{X} = [\mathbf{x}^{(0)}, \mathbf{x}^{(1)}, \dots, \mathbf{x}^{(i)}, \dots, \mathbf{x}^{(m)}]$$
$$\mathbf{y} = [y^{(0)}, y^{(1)}, \dots, y^{(i)}, \dots, y^{(m)}]$$

$$\Rightarrow \partial_{\mathbf{w}} J(\mathbf{w}) = 2\mathbf{X}\mathbf{X}^T \mathbf{w} - 2\mathbf{X}\mathbf{y}^T = 0$$

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

Example:

OR function:  $\mathbb{X} = \{[0, 0], [0, 1], [1, 0], [1, 1]\}$ ; target function  $f^*(\mathbf{x})$  gives  $\{0, 1, 1, 1\}$  on  $\mathbb{X}$ .

Example:

OR function:  $\mathbb{X} = \{[0, 0], [0, 1], [1, 0], [1, 1]\}$ ; target function  $f^*(\mathbf{x})$  gives  $\{0, 1, 1, 1\}$  on  $\mathbb{X}$ .

GOAL: Find an optimal (minimum MSE) function  $f(\mathbf{x}; \theta)$  that can fit these data points.

Example:

OR function:  $\mathbb{X} = \{[0, 0], [0, 1], [1, 0], [1, 1]\}$ ; target function  $f^*(\mathbf{x})$  gives  $\{0, 1, 1, 1\}$  on  $\mathbb{X}$ .

GOAL: Find an optimal (minimum MSE) function  $f(\mathbf{x}; \theta)$  that can fit these data points.

Define a linear function  $f(\mathbf{x}; \mathbf{w}) = \mathbf{w}^T \mathbf{x}$

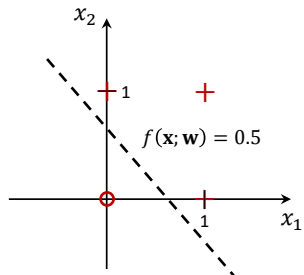


Example:

OR function:  $\mathbb{X} = \{[0, 0], [0, 1], [1, 0], [1, 1]\}$ ; target function  $f^*(\mathbf{x})$  gives  $\{0, 1, 1, 1\}$  on  $\mathbb{X}$ .

GOAL: Find an optimal (minimum MSE) function  $f(\mathbf{x}; \theta)$  that can fit these data points.

Define a linear function  $f(\mathbf{x}; \mathbf{w}) = \mathbf{w}^T \mathbf{x}$



Solving the normal equations leads to:  $\mathbf{w} = [0.6667, 0.6667]$

## Example 2:

XOR (exclusive or):  $\mathbb{X} = \{[0, 0], [0, 1], [1, 0], [1, 1]\}$ ; target function  $f^*(\mathbf{x})$  gives  $\{0, 1, 1, 0\}$  on  $\mathbb{X}$ .

## Example 2:

XOR (exclusive or):  $\mathbb{X} = \{[0, 0], [0, 1], [1, 0], [1, 1]\}$ ; target function  $f^*(\mathbf{x})$  gives  $\{0, 1, 1, 0\}$  on  $\mathbb{X}$ .

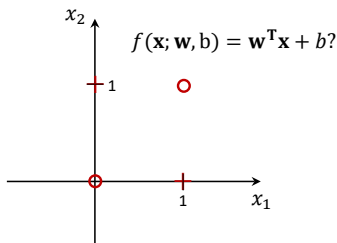
GOAL: Find an optimal (minimum MSE) function  $f(\mathbf{x}; \theta)$  that can fit these data points.

## Example 2:

XOR (exclusive or):  $\mathbb{X} = \{[0, 0], [0, 1], [1, 0], [1, 1]\}$ ; target function  $f^*(\mathbf{x})$  gives  $\{0, 1, 1, 0\}$  on  $\mathbb{X}$ .

GOAL: Find an optimal (minimum MSE) function  $f(\mathbf{x}; \theta)$  that can fit these data points.

Define a linear function  $f(\mathbf{x}; \mathbf{w}, b) = \mathbf{w}^T \mathbf{x} + b$

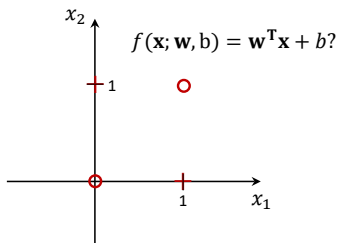


## Example 2:

XOR (exclusive or):  $\mathbb{X} = \{[0, 0], [0, 1], [1, 0], [1, 1]\}$ ; target function  $f^*(\mathbf{x})$  gives  $\{0, 1, 1, 0\}$  on  $\mathbb{X}$ .

GOAL: Find an optimal (minimum MSE) function  $f(\mathbf{x}; \theta)$  that can fit these data points.

Define a linear function  $f(\mathbf{x}; \mathbf{w}, b) = \mathbf{w}^T \mathbf{x} + b$



Solving the normal equations leads to:  $\mathbf{w} = \mathbf{0}$  and  $b = \frac{1}{2}$  (work this out yourselves).

Problem: A linear model  $f(\mathbf{x}; \mathbf{w}, b) = \mathbf{w}^T \mathbf{x} + b$  is not able to represent all desirable functions (e.g XOR, and in practice, many detection/estimation signal processing functions).

Problem: A linear model  $f(\mathbf{x}; \mathbf{w}, b) = \mathbf{w}^T \mathbf{x} + b$  is not able to represent all desirable functions (e.g XOR, and in practice, many detection/estimation signal processing functions).

Solution: Define a mapping function that (nonlinearly) transforms  $\mathbf{x}$  into a space  $\mathbf{h}$  in which the function can be represented with a linear model:

Problem: A linear model  $f(\mathbf{x}; \mathbf{w}, b) = \mathbf{w}^T \mathbf{x} + b$  is not able to represent all desirable functions (e.g XOR, and in practice, many detection/estimation signal processing functions).

Solution: Define a mapping function that (nonlinearly) transforms  $\mathbf{x}$  into a space  $\mathbf{h}$  in which the function can be represented with a linear model:

$$f(\mathbf{x}; \theta) = f^{(2)}(\mathbf{h}; \mathbf{w}^{(2)}, b^{(2)}), \quad (2)$$



Problem: A linear model  $f(\mathbf{x}; \mathbf{w}, b) = \mathbf{w}^T \mathbf{x} + b$  is not able to represent all desirable functions (e.g XOR, and in practice, many detection/estimation signal processing functions).

Solution: Define a mapping function that (nonlinearly) transforms  $\mathbf{x}$  into a space  $\mathbf{h}$  in which the function can be represented with a linear model:

$$f(\mathbf{x}; \theta) = f^{(2)}(\mathbf{h}; \mathbf{w}^{(2)}, b^{(2)}), \quad (2)$$

where

$$\mathbf{h} = f^{(1)}(\mathbf{x}; \mathbf{W}^{(1)}, \mathbf{b}^{(1)}). \quad (3)$$

Solution: Define a mapping function that (nonlinearly) transforms  $\mathbf{x}$  into a space  $\mathbf{h}$  in which the function can be represented with a linear model:

$$f(\mathbf{x}; \theta) = f^{(2)}(\mathbf{h}; \mathbf{w}^{(2)}, b^{(2)}), \quad (4)$$

where

$$\mathbf{h} = f^{(1)}(\mathbf{x}; \mathbf{W}^{(1)}, b^{(1)}). \quad (5)$$

Solution: Define a mapping function that (nonlinearly) transforms  $\mathbf{x}$  into a space  $\mathbf{h}$  in which the function can be represented with a linear model:

$$f(\mathbf{x}; \theta) = f^{(2)}(\mathbf{h}; \mathbf{w}^{(2)}, b^{(2)}), \quad (4)$$

where

$$\mathbf{h} = g(\mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)}). \quad (6)$$

Solution: Define a mapping function that (nonlinearly) transforms  $\mathbf{x}$  into a space  $\mathbf{h}$  in which the function can be represented with a linear model:

$$f(\mathbf{x}; \theta) = f^{(2)}(\mathbf{h}; \mathbf{w}^{(2)}, b^{(2)}), \quad (4)$$

where

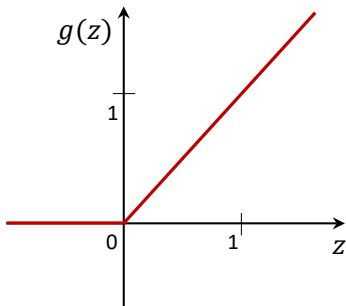
$$\mathbf{h} = g(\mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)}). \quad (6)$$

Let's consider a rectifying nonlinearity (in the machine learning community known as a rectified linear unit, or ReLU), such that

$$\mathbf{h} = \max(0, \mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)}). \quad (7)$$

$$f(\mathbf{x}; \theta) = f^{(2)}(\mathbf{h}; \mathbf{w}^{(2)}, b^{(2)}), \quad (8)$$

$$\mathbf{h} = \max(0, \mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)}). \quad (9)$$



$$f(\mathbf{x}; \theta) = f^{(2)}(\mathbf{h}; \mathbf{w}^{(2)}, b^{(2)}), \quad (10)$$

$$\mathbf{h} = \max(0, \mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)}). \quad (11)$$

The complete network is then:

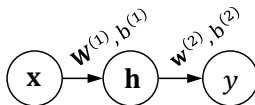
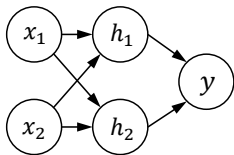
$$f(\mathbf{x}; \mathbf{W}^{(1)}, \mathbf{b}^{(1)}, \mathbf{w}^{(2)}, b^{(2)}) = \left(\mathbf{w}^{(2)}\right)^T \max(0, \mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)}) + b^{(2)}. \quad (12)$$

$$f(\mathbf{x}; \theta) = f^{(2)}(\mathbf{h}; \mathbf{w}^{(2)}, b^{(2)}), \quad (10)$$

$$\mathbf{h} = \max(0, \mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)}). \quad (11)$$

The complete network is then:

$$f(\mathbf{x}; \mathbf{W}^{(1)}, \mathbf{b}^{(1)}, \mathbf{w}^{(2)}, b^{(2)}) = \left(\mathbf{w}^{(2)}\right)^T \max(0, \mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)}) + b^{(2)}. \quad (12)$$



Given:

$$f(\mathbf{x}; \mathbf{W}^{(1)}, \mathbf{b}^{(1)}, \mathbf{w}^{(2)}, b^{(2)}) = \left(\mathbf{w}^{(2)}\right)^T \max(0, \mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)}) + b^{(2)}, \quad (13)$$

what parameter values should  $\mathbf{W}^{(1)}$ ,  $\mathbf{b}^{(1)}$ ,  $\mathbf{w}^{(2)}$  and  $b^{(2)}$  have?



Given:

$$f(\mathbf{x}; \mathbf{W}^{(1)}, \mathbf{b}^{(1)}, \mathbf{w}^{(2)}, b^{(2)}) = \left(\mathbf{w}^{(2)}\right)^T \max(0, \mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)}) + b^{(2)}, \quad (13)$$

what parameter values should  $\mathbf{W}^{(1)}$ ,  $\mathbf{b}^{(1)}$ ,  $\mathbf{w}^{(2)}$  and  $b^{(2)}$  have? Consider the following values:

$$\mathbf{W}^{(1)} = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}, \mathbf{b}^{(1)} = \begin{bmatrix} 0 \\ -1 \end{bmatrix}, \mathbf{w}^{(2)} = \begin{bmatrix} 1 \\ -2 \end{bmatrix} \text{ and } b^{(2)} = 0.$$

Given:

$$f(\mathbf{x}; \mathbf{W}^{(1)}, \mathbf{b}^{(1)}, \mathbf{w}^{(2)}, b^{(2)}) = \left(\mathbf{w}^{(2)}\right)^T \max(0, \mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)}) + b^{(2)}, \quad (13)$$

what parameter values should  $\mathbf{W}^{(1)}$ ,  $\mathbf{b}^{(1)}$ ,  $\mathbf{w}^{(2)}$  and  $b^{(2)}$  have? Consider the following values:

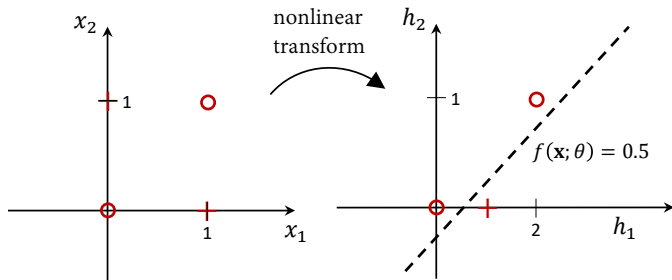
$$\mathbf{W}^{(1)} = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}, \mathbf{b}^{(1)} = \begin{bmatrix} 0 \\ -1 \end{bmatrix}, \mathbf{w}^{(2)} = \begin{bmatrix} 1 \\ -2 \end{bmatrix} \text{ and } b^{(2)} = 0.$$

$$\text{Then } \mathbf{X} = \begin{bmatrix} 0 & 0 & 1 & 1 \\ 0 & 1 & 0 & 1 \end{bmatrix}, \text{ maps to outputs: } \mathbf{y} = [0 \quad 1 \quad 1 \quad 0].$$

$$f(\mathbf{x}; \mathbf{W}^{(1)}, \mathbf{b}^{(1)}, \mathbf{w}^{(2)}, b^{(2)}) = \left(\mathbf{w}^{(2)}\right)^T \max(0, \mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)}) + b^{(2)}, \quad (14)$$

where

$$\mathbf{W}^{(1)} = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}, \mathbf{b}^{(1)} = \begin{bmatrix} 0 \\ -1 \end{bmatrix}, \mathbf{w}^{(2)} = \begin{bmatrix} 1 \\ -2 \end{bmatrix} \text{ and } b^{(2)} = 0.$$



Optimal parameter values?

Optimal parameter values?

⇒ Maximize the likelihood of the observations ( $\mathbf{y}$ ) given the model ( $\theta$ ) and the input data ( $\mathbf{x}$ ).

Negative log-likelihood cost function:

$$J(\theta) = -\mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim \hat{p}_{\text{data}}} \log p_{\text{model}}(\mathbf{y} | \mathbf{x}, \theta), \quad (15)$$

Optimal parameter values?

⇒ Maximize the likelihood of the observations ( $\mathbf{y}$ ) given the model ( $\theta$ ) and the input data ( $\mathbf{x}$ ).

Negative log-likelihood cost function:

$$J(\theta) = -\mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim \hat{p}_{\text{data}}} \log p_{\text{model}}(\mathbf{y}|\mathbf{x}, \theta), \quad (15)$$

$\hat{p}_{\text{data}}$ : Empirical data distribution defined by the training set  $\mathbb{X}$   
 $\hat{p}_{\text{model}}(\mathbf{y}|\mathbf{x}; \theta)$ : Error model for the predictions.

Optimal parameter values?

⇒ Maximize the likelihood of the observations ( $\mathbf{y}$ ) given the model ( $\theta$ ) and the input data ( $\mathbf{x}$ ).

Negative log-likelihood cost function:

$$J(\theta) = -\mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim \hat{p}_{\text{data}}} \log p_{\text{model}}(\mathbf{y}|\mathbf{x}, \theta), \quad (15)$$

$\hat{p}_{\text{data}}$ : Empirical data distribution defined by the training set  $\mathbb{X}$   
 $\hat{p}_{\text{model}}(\mathbf{y}|\mathbf{x}; \theta)$ : Error model for the predictions.

$$J(\theta) = -\log \prod_{i=0}^{m-1} \hat{p}_{\text{model}}(\mathbf{y}^{(i)}|\mathbf{x}^{(i)}, \theta), \quad (16)$$

Optimal parameter values?

⇒ Maximize the likelihood of the observations ( $\mathbf{y}$ ) given the model ( $\theta$ ) and the input data ( $\mathbf{x}$ ).

Negative log-likelihood cost function:

$$J(\theta) = -\mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim \hat{p}_{\text{data}}} \log p_{\text{model}}(\mathbf{y}|\mathbf{x}, \theta), \quad (15)$$

$\hat{p}_{\text{data}}$ : Empirical data distribution defined by the training set  $\mathbb{X}$   
 $\hat{p}_{\text{model}}(\mathbf{y}|\mathbf{x}; \theta)$ : Error model for the predictions.

$$J(\theta) = -\log \prod_{i=0}^{m-1} \hat{p}_{\text{model}}(\mathbf{y}^{(i)}|\mathbf{x}^{(i)}, \theta), \quad (16)$$

$$J(\theta) = -\sum_{i=0}^{m-1} \log \hat{p}_{\text{model}}(\mathbf{y}^{(i)}|\mathbf{x}^{(i)}, \theta), \quad (17)$$



Optimal parameter values?

⇒ Maximize the likelihood of the observations ( $\mathbf{y}$ ) given the model ( $\theta$ ) and the input data ( $\mathbf{x}$ ).

Negative log-likelihood cost function:

$$J(\theta) = - \sum_{i=0}^{m-1} \log \hat{p}_{\text{model}}(\mathbf{y}^{(i)} | \mathbf{x}^{(i)}, \theta), \quad (18)$$

Optimal parameter values?

⇒ Maximize the likelihood of the observations ( $\mathbf{y}$ ) given the model ( $\theta$ ) and the input data ( $\mathbf{x}$ ).

Negative log-likelihood cost function:

$$J(\theta) = - \sum_{i=0}^{m-1} \log \hat{p}_{\text{model}}(\mathbf{y}^{(i)} | \mathbf{x}^{(i)}, \theta), \quad (18)$$

$$\hat{\theta} = \underset{\theta}{\operatorname{argmin}} J(\theta) = \underset{\theta}{\operatorname{argmin}} \left( - \sum_{i=0}^{m-1} \log \hat{p}_{\text{model}}(\mathbf{y}^{(i)} | \mathbf{x}^{(i)}, \theta) \right), \quad (19)$$

Optimal parameter values:

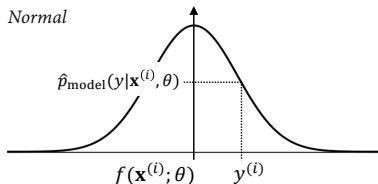
$$\hat{\theta} = \underset{\theta}{\operatorname{argmin}} \left( - \sum_{i=0}^{m-1} \log \hat{p}_{\text{model}}(\mathbf{y}^{(i)} | \mathbf{x}^{(i)}, \theta) \right), \quad (20)$$

Choice of  $\hat{p}_{\text{model}}$  depends on the error distribution, e.g.:

Optimal parameter values:

$$\hat{\theta} = \underset{\theta}{\operatorname{argmin}} \left( - \sum_{i=0}^{m-1} \log \hat{p}_{\text{model}}(\mathbf{y}^{(i)} | \mathbf{x}^{(i)}, \theta) \right), \quad (20)$$

Choice of  $\hat{p}_{\text{model}}$  depends on the error distribution, e.g.:

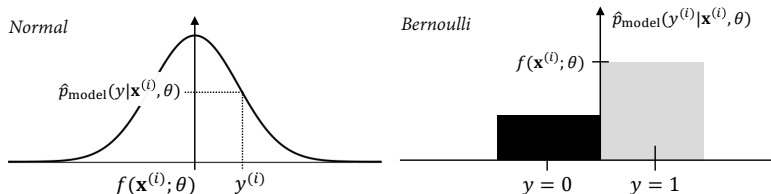


- ▶ Normal (Gaussian) distribution
  - Regression of continuous variables

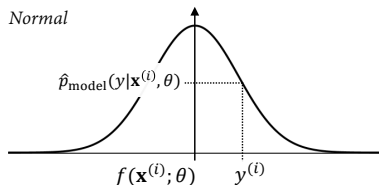
Optimal parameter values:

$$\hat{\theta} = \underset{\theta}{\operatorname{argmin}} \left( - \sum_{i=0}^{m-1} \log \hat{p}_{\text{model}}(\mathbf{y}^{(i)} | \mathbf{x}^{(i)}, \theta) \right), \quad (20)$$

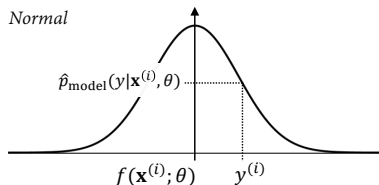
Choice of  $\hat{p}_{\text{model}}$  depends on the error distribution, e.g.:



- ▶ Normal (Gaussian) distribution
  - Regression of continuous variables
- ▶ Bernoulli/Categorical distribution
  - Binary/multi-class classification

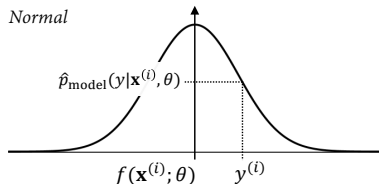


Optimal parameter values for Gaussian error distribution,  
 $\hat{p}_{\text{model}} \sim \mathcal{N}\{\mathbf{y}; f(\mathbf{x}; \theta), \sigma \mathbf{I}\}$ :



Optimal parameter values for Gaussian error distribution,  
 $\hat{p}_{\text{model}} \sim \mathcal{N}\{\mathbf{y}; f(\mathbf{x}; \theta), \sigma^2 \mathbf{I}\}$ :

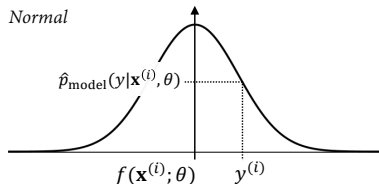
“The probability of observing  $y^{(i)}$  follows a normal distribution with the model prediction  $f(\mathbf{x}^{(i)}; \theta)$  as its mean.”



Optimal parameter values for Gaussian error distribution,  
 $\hat{p}_{\text{model}} \sim \mathcal{N}\{\mathbf{y}; f(\mathbf{x}; \theta), \sigma \mathbf{I}\}$ :

$$\hat{\theta} = \underset{\theta}{\operatorname{argmin}} \left( -\frac{1}{m} \sum_{i=0}^{m-1} \log \hat{p}_{\text{model}}(\mathbf{y}^{(i)} | \mathbf{x}^{(i)}, \theta) \right). \quad (21)$$

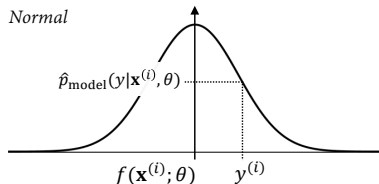




Optimal parameter values for Gaussian error distribution,  
 $\hat{p}_{\text{model}} \sim \mathcal{N}\{\mathbf{y}; f(\mathbf{x}; \theta), \sigma \mathbf{I}\}$ :

$$\hat{\theta} = \underset{\theta}{\operatorname{argmin}} \left( -\frac{1}{m} \sum_{i=0}^{m-1} \log \hat{p}_{\text{model}}(\mathbf{y}^{(i)} | \mathbf{x}^{(i)}, \theta) \right), \quad (22)$$

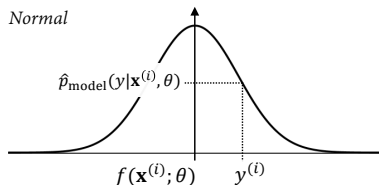
$$\hat{\theta} = \underset{\theta}{\operatorname{argmin}} \left( -\frac{1}{m} \sum_{i=0}^{m-1} \log e^{-\frac{1}{2\sigma^2}} [\mathbf{y}^{(i)} - f(\mathbf{x}^{(i)}; \theta)] [\mathbf{y}^{(i)} - f(\mathbf{x}^{(i)}; \theta)]^T \right), \quad (23)$$



Optimal parameter values for Gaussian error distribution,  
 $\hat{p}_{\text{model}} \sim \mathcal{N}\{\mathbf{y}; f(\mathbf{x}; \theta), \sigma^2 \mathbf{I}\}$ :

$$\hat{\theta} = \underset{\theta}{\operatorname{argmin}} \left( -\frac{1}{m} \sum_{i=0}^{m-1} \log e^{-\frac{1}{2\sigma^2}} [\mathbf{y}^{(i)} - f(\mathbf{x}^{(i)}; \theta)] [\mathbf{y}^{(i)} - f(\mathbf{x}^{(i)}; \theta)]^T \right), \quad (24)$$

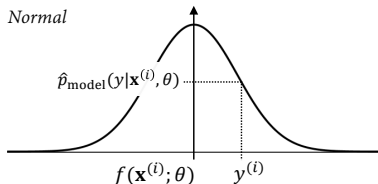
$$\hat{\theta} = \underset{\theta}{\operatorname{argmin}} \left( \frac{1}{m} \sum_{i=0}^{m-1} [\mathbf{y}^{(i)} - f(\mathbf{x}^{(i)}; \theta)] [\mathbf{y}^{(i)} - f(\mathbf{x}^{(i)}; \theta)]^T \right), \quad (25)$$



Optimal parameter values for Gaussian error distribution,  
 $\hat{p}_{\text{model}} \sim \mathcal{N}\{\mathbf{y}; f(\mathbf{x}; \theta), \sigma \mathbf{I}\}$ :

$$\hat{\theta} = \underset{\theta}{\operatorname{argmin}} \left( \frac{1}{m} \sum_{i=0}^{m-1} [\mathbf{y}^{(i)} - f(\mathbf{x}^{(i)}; \theta)] [\mathbf{y}^{(i)} - f(\mathbf{x}^{(i)}; \theta)]^T \right), \quad (26)$$

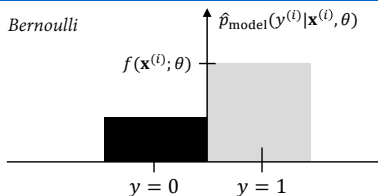
$$\hat{\theta} = \underset{\theta}{\operatorname{argmin}} \frac{1}{m} \sum_{i=0}^{m-1} \left\| \mathbf{y}^{(i)} - f(\mathbf{x}^{(i)}; \theta) \right\|_2^2, \quad (27)$$



Optimal parameter values for Gaussian error distribution,  
 $\hat{p}_{\text{model}} \sim \mathcal{N}\{\mathbf{y}; f(\mathbf{x}; \theta), \sigma \mathbf{I}\}$ :

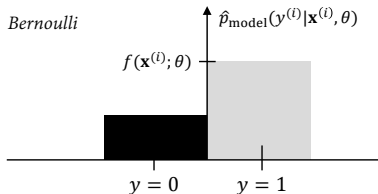
$$\hat{\theta} = \underset{\theta}{\operatorname{argmin}} \frac{1}{m} \sum_{i=0}^{m-1} \left\| \mathbf{y}^{(i)} - f(\mathbf{x}^{(i)}; \theta) \right\|_2^2, \quad (28)$$

$$y^{(i)} \in \mathbb{R}^1 \Rightarrow \hat{\theta} = \underset{\theta}{\operatorname{argmin}} \underbrace{\frac{1}{m} \sum_{i=0}^{m-1} \left( y^{(i)} - f(\mathbf{x}^{(i)}; \theta) \right)^2}_{J(\theta) = \text{mean squared error}}, \quad (29)$$



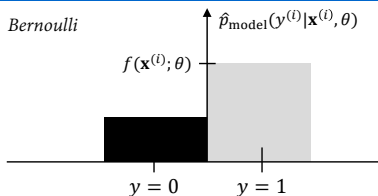
Optimal parameter values for Bernoulli error distribution,

$$\hat{p}_{\text{model}}(y^{(i)}|\mathbf{x}^{(i)}, \theta) = \begin{cases} p & \text{for } y^{(i)} = 1 \\ 1 - p & \text{for } y^{(i)} = 0 \end{cases}$$



Optimal parameter values for Bernoulli error distribution,

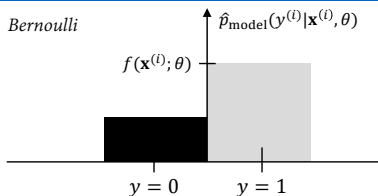
$$\hat{p}_{\text{model}}(y^{(i)}|\mathbf{x}^{(i)}, \theta) = p^{y^{(i)}}(1-p)^{(1-y^{(i)})} \text{ for } y^{(i)} \in \{0, 1\}$$



Optimal parameter values for Bernoulli error distribution,

$$\hat{p}_{\text{model}}(y^{(i)}|\mathbf{x}^{(i)}, \theta) = p^{y^{(i)}}(1-p)^{(1-y^{(i)})} \text{ for } y^{(i)} \in \{0, 1\}$$

“The probability of observing  $y^{(i)} = 1$  follows a Bernoulli distribution with the model prediction  $f(\mathbf{x}^{(i)}; \theta)$  determining its probability  $p$ .”

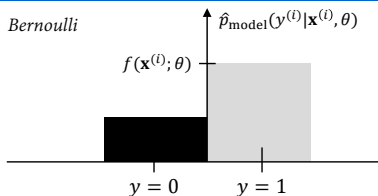


Optimal parameter values for Bernoulli error distribution,

$$\hat{p}_{\text{model}}(y^{(i)}|\mathbf{x}^{(i)}, \theta) = p^{y^{(i)}}(1-p)^{(1-y^{(i)})} \text{ for } y^{(i)} \in \{0, 1\}$$

$$\hat{\theta} = \underset{\theta}{\operatorname{argmin}} \left( -\frac{1}{m} \sum_{i=0}^{m-1} \log \hat{p}_{\text{model}}(y^{(i)}|\mathbf{x}^{(i)}, \theta) \right). \quad (30)$$



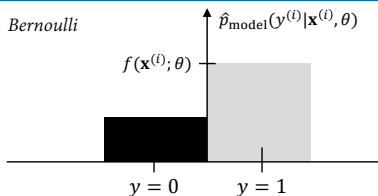


Optimal parameter values for Bernoulli error distribution,

$$\hat{p}_{\text{model}}(y^{(i)}|\mathbf{x}^{(i)}, \theta) = p^{y^{(i)}}(1-p)^{(1-y^{(i)})} \text{ for } y^{(i)} \in \{0, 1\}$$

$$\hat{\theta} = \underset{\theta}{\operatorname{argmin}} \left( - \sum_{i=0}^{m-1} \log \hat{p}_{\text{model}}(y^{(i)}|\mathbf{x}^{(i)}, \theta) \right). \quad (31)$$

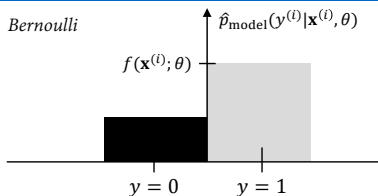
$$\hat{\theta} = \underset{\theta}{\operatorname{argmin}} \left( - \sum_{i=0}^{m-1} \log \left( \left( p^{(i)} \right)^{y^{(i)}} \left( 1 - p^{(i)} \right)^{(1-y^{(i)})} \right) \right), \quad (32)$$



Optimal parameter values for Bernoulli error distribution,

$$\hat{p}_{\text{model}}(y^{(i)}|\mathbf{x}^{(i)}, \theta) = p^{y^{(i)}}(1-p)^{(1-y^{(i)})} \text{ for } y^{(i)} \in \{0, 1\}$$

$$\Rightarrow \hat{\theta} = \underset{\theta}{\operatorname{argmin}} \left( - \sum_{i=0}^{m-1} y^{(i)} \log(p^{(i)}) + (1 - y^{(i)}) \log(1 - p^{(i)}) \right), \quad (33)$$



Optimal parameter values for Bernoulli error distribution,

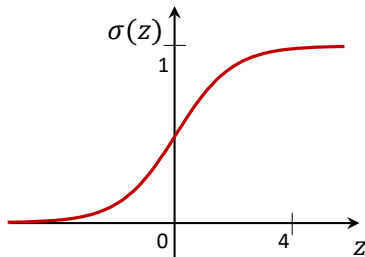
$$\hat{p}_{\text{model}}(y^{(i)}|\mathbf{x}^{(i)}, \theta) = p^{y^{(i)}}(1-p)^{(1-y^{(i)})} \text{ for } y^{(i)} \in \{0, 1\}$$

$$\Rightarrow \hat{\theta} = \underset{\theta}{\operatorname{argmin}} \underbrace{\left( - \sum_{i=0}^{m-1} y^{(i)} \log(p^{(i)}) + (1 - y^{(i)}) \log(1 - p^{(i)}) \right)}_{J(\theta)=\text{Binary cross entropy}}, \quad (34)$$

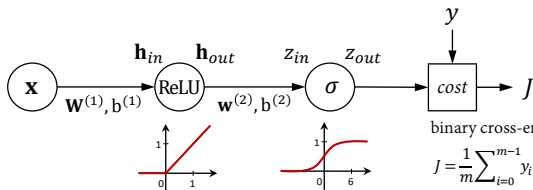
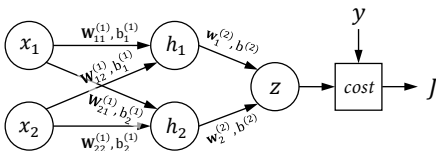
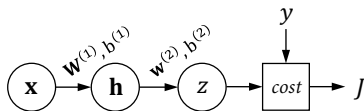
We require model that maps inputs to probabilities  $p$ .

Consider a nonlinearity that squeezes all model output values between 0 and 1:

$$p = \sigma(f(z)) = \frac{1}{1 + e^{-z}}. \quad (35)$$



The full binary classification model becomes:



binary cross-entropy:

$$J = \frac{1}{m} \sum_{i=0}^{m-1} y_i \log z_{out} + (1 - y_i) \log(1 - z_{out})$$

## Finding optimal parameter values given the cost

⇒ Gradient-based learning (e.g. Gradient-Descent algorithm):

- ▶  $\mathbf{W}_{n+1}^{(1)} = \mathbf{W}_n^{(1)} - \mu \partial_{\mathbf{W}^{(1)}} J(\theta_n)$
- ▶  $\mathbf{w}_{n+1}^{(2)} = \mathbf{w}_n^{(2)} - \mu \partial_{\mathbf{w}^{(2)}} J(\theta_n)$
- ▶  $\mathbf{b}_{n+1}^{(1)} = \mathbf{b}_n^{(1)} - \mu \partial_{\mathbf{b}^{(1)}} J(\theta_n)$
- ▶  $b_{n+1}^{(2)} = b_n^{(2)} - \mu \partial_{b^{(2)}} J(\theta)$

With learning rate  $\mu$ .

## Finding optimal parameter values given the cost

⇒ Gradient-based learning (e.g. Gradient-Descent algorithm):

- ▶ Linear models; Convex → guaranteed global convergence.
- ▶ Nonlinear models; Often non-convex → no global convergence guarantees.

## Finding optimal parameter values given the cost

⇒ Gradient-based learning (e.g. Gradient-Descent algorithm):

- ▶ Linear models; Convex → guaranteed global convergence.
- ▶ Nonlinear models; Often non-convex → no global convergence guarantees.

