

# Logistic Regression

Data Mining

Prof. Sujee Lee

Department of Systems Management Engineering

Sungkyunkwan University

# Logistic Regression

# Logistic Regression

---

## ■ (Recap.) Linear Regression

- In linear regression,  $y$  is predicted by

$$\hat{y} = f(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b = w_1 x_1 + \cdots + w_d x_d + b$$

- here,  $y, \hat{y} \in \mathbb{R}$

## ■ Logistic Regression

- Logistic Regression extends the ideas of linear regression for classification problem
  - i.e., the labels are binary  $y = 0$  or  $1$
- we will use linear model  $\mathbf{w}^T \mathbf{x} + b$ , but  $f(\mathbf{x})$  to be a probability

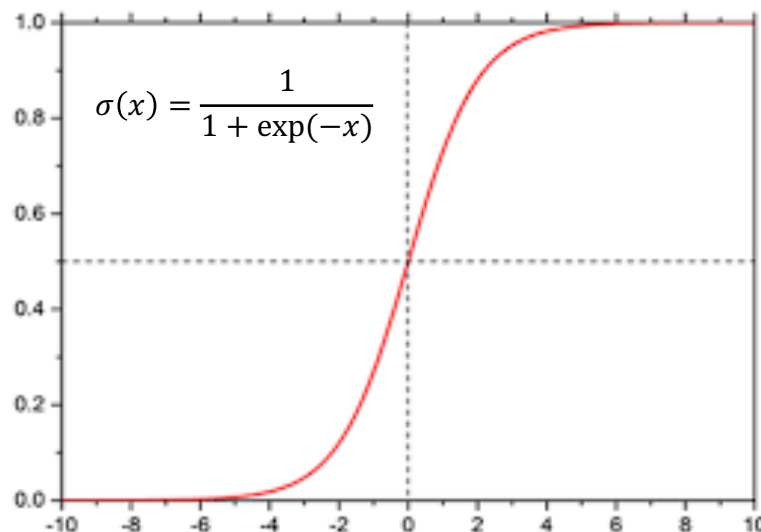
$$\hat{y} = f(\mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x} + b) = \frac{1}{1 + \exp(-\mathbf{w}^T \mathbf{x} - b)}$$

$$\mathbf{x} = (x_1, \dots, x_d) \in \mathbb{R}^d, \quad y \in \mathbb{B}, \quad 0 \leq \hat{y} \leq 1$$

# Logistic Regression

$$\hat{y} = f(\mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x} + b) = \frac{1}{1 + \exp(-\mathbf{w}^T \mathbf{x} - b)}$$

- here  $\sigma(\cdot)$  is called as a **sigmoid function** or a **logistic function**, and results in between 0 and 1

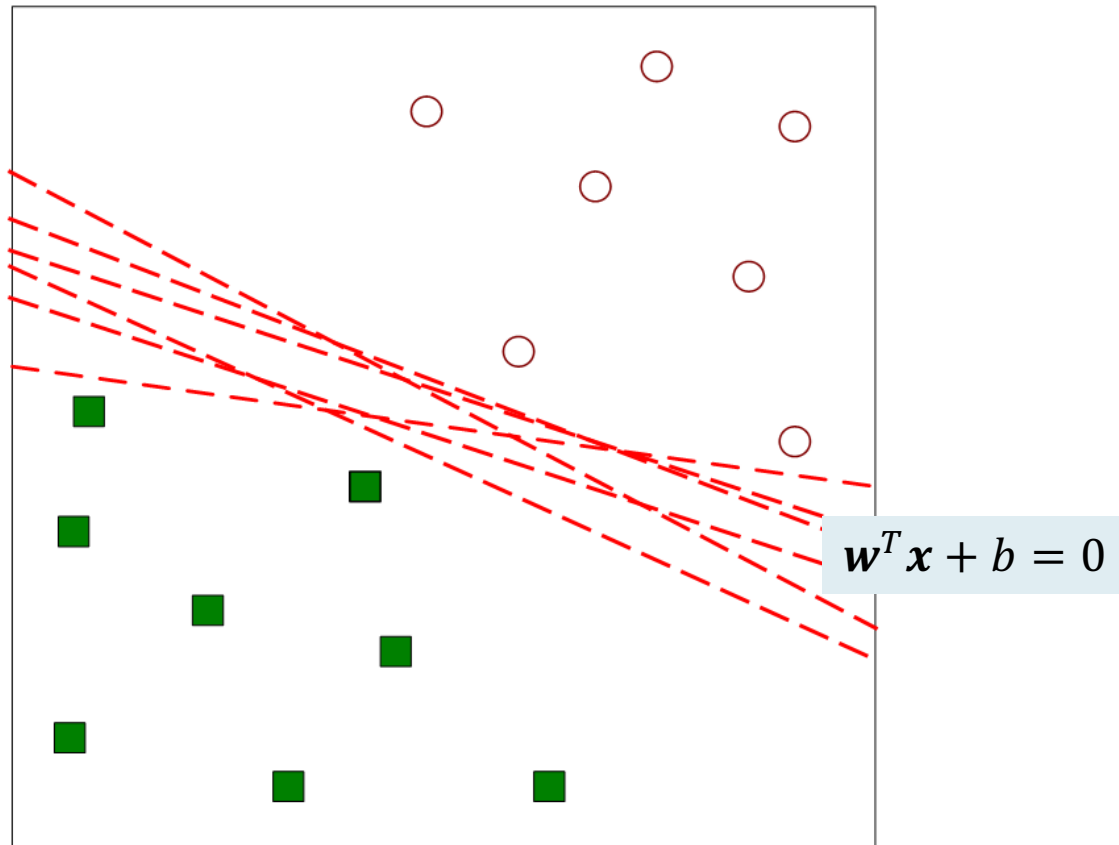


- then, we can consider the class probabilities as

$$P(y = 1) = \frac{1}{1 + \exp(-\mathbf{w}^T \mathbf{x} - b)}, \quad P(y = 0) = \frac{\exp(-\mathbf{w}^T \mathbf{x} - b)}{1 + \exp(-\mathbf{w}^T \mathbf{x} - b)}$$

# Logistic Regression

- This algorithm is a linear classifier
  - For linear models for binary classification, the **decision boundary** (hyperplane) that separates two classes is a **linear function** of input features.



# Learning a Logistic Regression Model

## ■ Logistic Regression

- Goal : Given data pairs  $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$ , find the optimal parameter  $\mathbf{w}^*$  that minimizes the training error. (as we did in linear regression)
- here, we use “*binary cross-entropy*” loss to define classification error

$$\min_{\mathbf{w}} \sum_{i=1}^n \{-y_i \log \hat{y}_i - (1 - y_i) \log(1 - \hat{y}_i)\}$$

$$-y_i \log \hat{y}_i - (1 - y_i) \log(1 - \hat{y}_i): \quad \text{If } y_i = 1, -\log \hat{y}_i \quad \text{If } y_i = 0, -\log(1 - \hat{y}_i)$$

$$\hat{y}_i = 1$$

$$\hat{y}_i = 0$$

- There is no “closed-form”. How can we solve the above problem? => “*gradient descent*”

# Gradient Descent

# Optimization Problem

---

- Optimization Problem

- (Recap) Linear Regression

- Ordinary Least Squares :  $\min_w \|\mathbf{y} - \mathbf{X}\mathbf{w}\|^2$

- (Further)

- Logistic Regression :  $\min_w \sum_{i=1}^n \{-y_i \log \hat{y}_i - (1 - y_i) \log(1 - \hat{y}_i)\}$
    - Neural Network : ...

- All are about **minimizing loss function** (cost function or error function)

$$\min_w L(w; \mathbf{X}, \mathbf{y})$$



# Optimization Problem

---

- **Function Minimization / Maximization Problem**

$$\min_x f(x)$$

- When can we find an (global) **optimal solution**? => **when  $f$  is convex**

- Note: How to check the convexity :  $f''(x)$  or  $\nabla_x^2 f(x)$  (Hessian matrix)
  - $f$  is convex if and only if  $f''(x) \geq 0$  for all  $x$
  - $f$  is convex if and only if  $\nabla_x^2 f(x) \geq 0$  (positive semi-definite) for all  $x$

# Optimization Problem

---

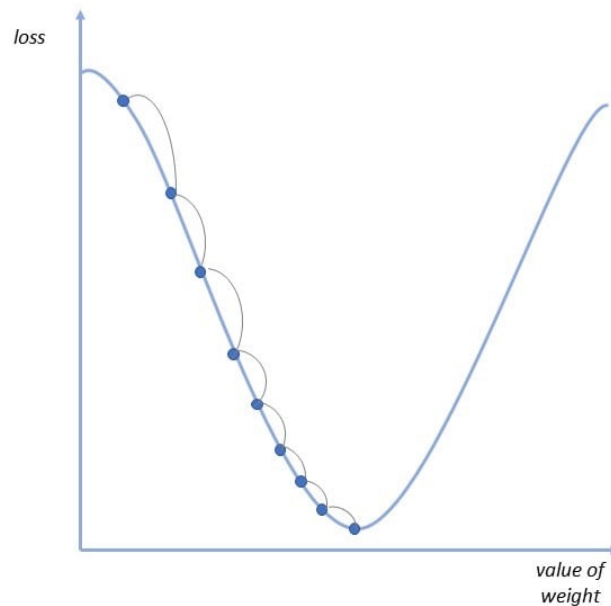
- When  $f$  is differentiable and convex, a necessary and sufficient condition for a point  $x^*$  to be optimal is  $\nabla f(x^*) = 0$ .
- If you can find an *analytical solution* (closed-form solution) for  $\nabla f(x^*) = 0$ , please do!
- If you can't,  $\nabla f(x^*) = 0$  usually can be solved by an *iterative algorithm*
  - it means, you will update your solutions as the following until convergence
$$x^{(k+1)} = x^{(k)} + \Delta x^{(k)}$$

# Gradient Descent

## ■ Gradient Descent

- Gradient descent (GD) is an iterative first-order optimization algorithm used to find a local minimum/maximum of a given function.
- To find  $\mathbf{x}^* = \operatorname{argmin}_{\mathbf{x}} f(\mathbf{x})$ , iteratively update

$$\mathbf{x} \leftarrow \mathbf{x} - \alpha \nabla_{\mathbf{x}} f(\mathbf{x}), \text{ i.e., } x_j \leftarrow x_j - \alpha \frac{\partial}{\partial x_j} f(\mathbf{x})$$



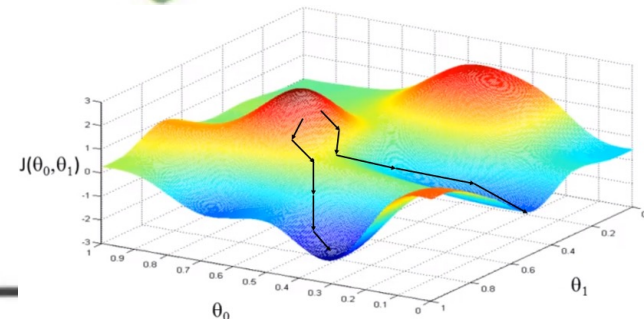
# Gradient Descent

- The gradient descent algorithm guides the search for values that minimize the function at a local/global minimum

$$\mathbf{x} \leftarrow \mathbf{x} - \alpha \nabla_{\mathbf{x}} f(\mathbf{x})$$

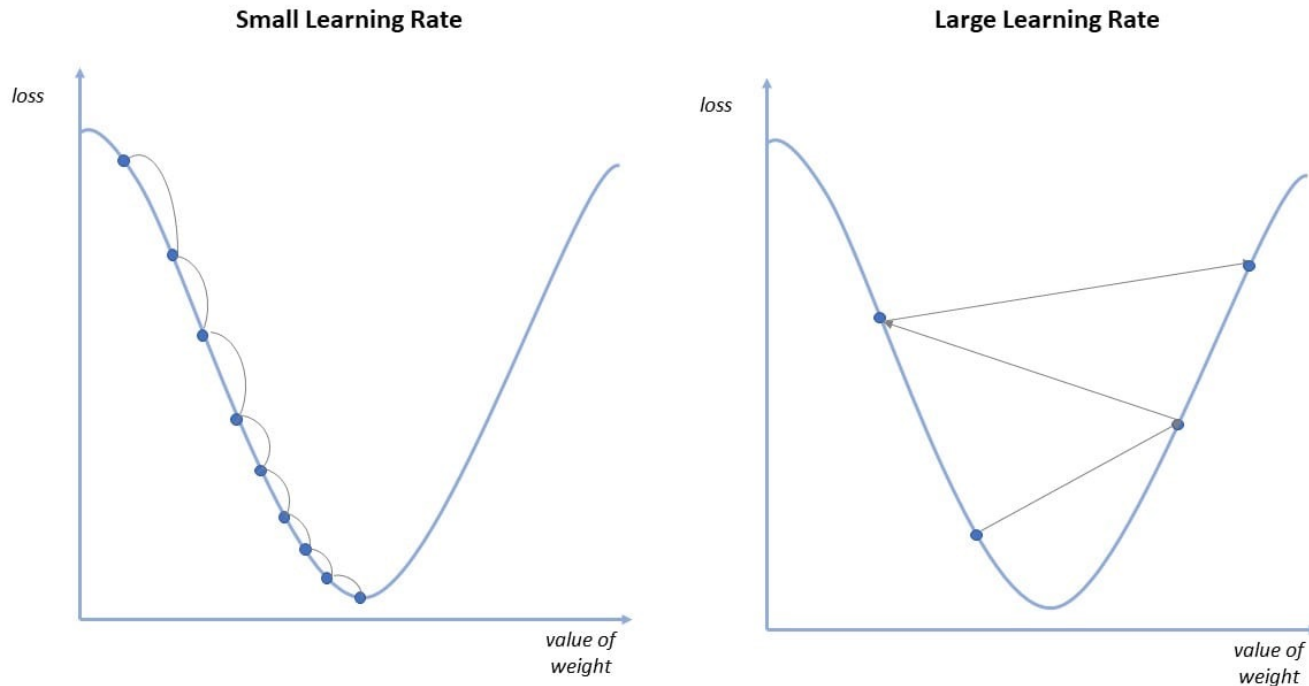


- **Local Minimum:** The minimum parameter values within a specified range or sector of the cost function.
- **Global Minimum:** This is the smallest parameter value within the entire cost function domain.



# Gradient Descent

- Role of the learning rate  $\alpha$ 
  - The learning rate determines the size of the steps taken towards the minimum of the loss function. A higher learning rate means larger steps, while a lower learning rate means smaller steps.



\* warning : gradient descent might “overshoot” if step size is chosen incorrectly ( $\alpha$  has to be small enough relative to curvature of the function)

# **Logistic Regression (again)**

# Learning a Logistic Regression Model

---

## ■ Logistic Regression

- Goal : Given data pairs  $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$ , find the optimal parameter  $\mathbf{w}^*$  that minimizes the training error. *(as we did in linear regression)*
- here, we use “*binary cross-entropy*” loss to define classification error

$$\min_{\mathbf{w}} \sum_{i=1}^n \{-y_i \log \hat{y}_i - (1 - y_i) \log(1 - \hat{y}_i)\}$$

- There is no “closed-form”. How can we solve the above problem? => “*gradient descent*”

# Learning a Logistic Regression Model

---

- Check the convexity

- Let  $L(\mathbf{w})$  be

$$L(\mathbf{w}) = \sum_{i=1}^n \{-y_i \log \hat{y}_i - (1 - y_i) \log(1 - \hat{y}_i)\}$$

- You want to obtain a global optimal solution  $\mathbf{w}^*$  for  $\min_{\mathbf{w}} L(\mathbf{w})$  using Gradient Descent

- To make sure that you can find it, you have to check that  $L(\mathbf{w})$  is convex

- Note: To do so,

- Calculate  $\nabla_{\mathbf{w}}^2 L(\mathbf{w})$
- Show that  $\nabla_{\mathbf{w}}^2 L(\mathbf{w})$  is positive semi-definite.
- If  $\nabla_{\mathbf{w}}^2 L(\mathbf{w})$  is positive semi-definite  $\Rightarrow$  then,  $L(\mathbf{w})$  is convex



# Learning a Logistic Regression Model

- Find the optimal  $\mathbf{w}^*$

- By the gradient decent on  $L(\mathbf{w})$ ,

$$\mathbf{w} \leftarrow \mathbf{w} - \alpha \nabla_{\mathbf{w}} L(\mathbf{w}) \quad (\alpha : \text{learning rate})$$

- Repeat the above until convergence
- $L(\mathbf{w}) = \sum_{i=1}^n \{-y_i \log \hat{y}_i - (1 - y_i) \log(1 - \hat{y}_i)\}$   
 $= \sum_{i=1}^n \{-y_i \log \sigma(\mathbf{w}^T \mathbf{x}) - (1 - y_i) \log(1 - \sigma(\mathbf{w}^T \mathbf{x}))\}$
- $\nabla_{\mathbf{w}} L(\mathbf{w}) =$

- it will converge to the optimal  $\mathbf{w}^*$

# Another View

# Probabilistic Approach for Logistic Regression

## ■ Probabilistic Approach for Logistic Regression

- Start from a generative model:
  - Assume that  $y_i$ 's are generated probabilistically from  $x_i$ 's and a parameter  $\mathbf{w}$  via

$$P(y = 1) = \frac{1}{1 + \exp(-\mathbf{w}^T \mathbf{x} - b)}, \quad P(y = 0) = \frac{\exp(-\mathbf{w}^T \mathbf{x} - b)}{1 + \exp(-\mathbf{w}^T \mathbf{x} - b)}$$

- Let's perform MLE to find a formula for  $\mathbf{w}$ :

$$L_{\mathbf{w}}(\mathbf{X}, \mathbf{y}) = \prod_{i=1}^n \left( \frac{1}{1 + \exp(-\mathbf{w}^T \mathbf{x}_i - b)} \right)^{y_i} \left( \frac{\exp(-\mathbf{w}^T \mathbf{x}_i - b)}{1 + \exp(-\mathbf{w}^T \mathbf{x}_i - b)} \right)^{(1-y_i)}$$

- Maximizing  $L_{\mathbf{w}}(\mathbf{X}, \mathbf{y})$  with respect to  $\mathbf{w}$  is equivalent to

$$\min_{\mathbf{w}} \sum_{i=1}^n \left\{ -y_i \log \left( \frac{1}{1 + \exp(-\mathbf{w}^T \mathbf{x}_i - b)} \right) - (1 - y_i) \log \left( \frac{\exp(-\mathbf{w}^T \mathbf{x}_i - b)}{1 + \exp(-\mathbf{w}^T \mathbf{x}_i - b)} \right) \right\}$$