



MLE, Logistic Regression

Iran University of Science and Technology

M. S. Tahaei, PhD.

Fall 2024

Courtesy: slides are adopted partly from Dr. Soleymani, Sharif University

Outline

- Maximum-Likelihood (ML) estimation
- Logistic Regression

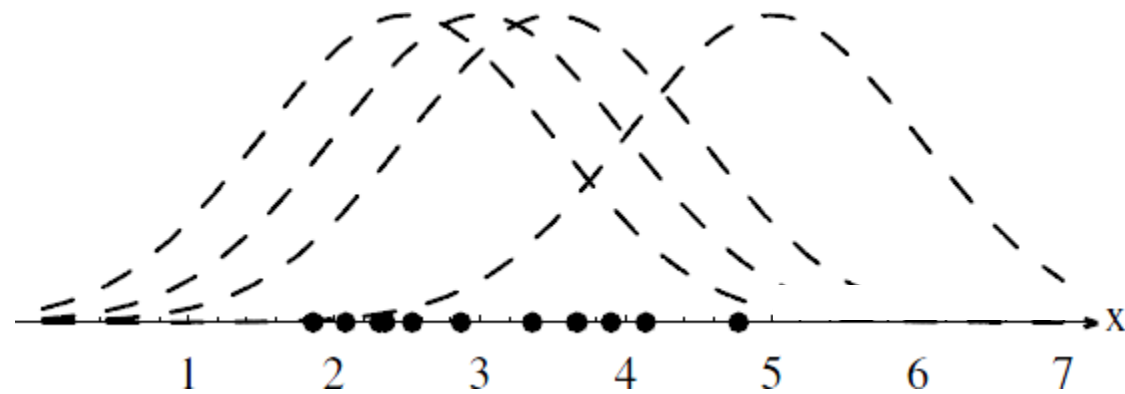
Parametric density estimation

- ▶ Estimating the probability density function $p(\boldsymbol{x})$, given a set of data points $\{\boldsymbol{x}^{(i)}\}_{i=1}^N$ drawn from it.
- ▶ Assume that $p(\boldsymbol{x})$ in terms of a specific functional form which has a number of adjustable parameters.
- ▶ Methods for parameter estimation
 - ▶ Maximum likelihood estimation
 - ▶ Maximum A Posteriori (MAP) estimation

Parametric density estimation

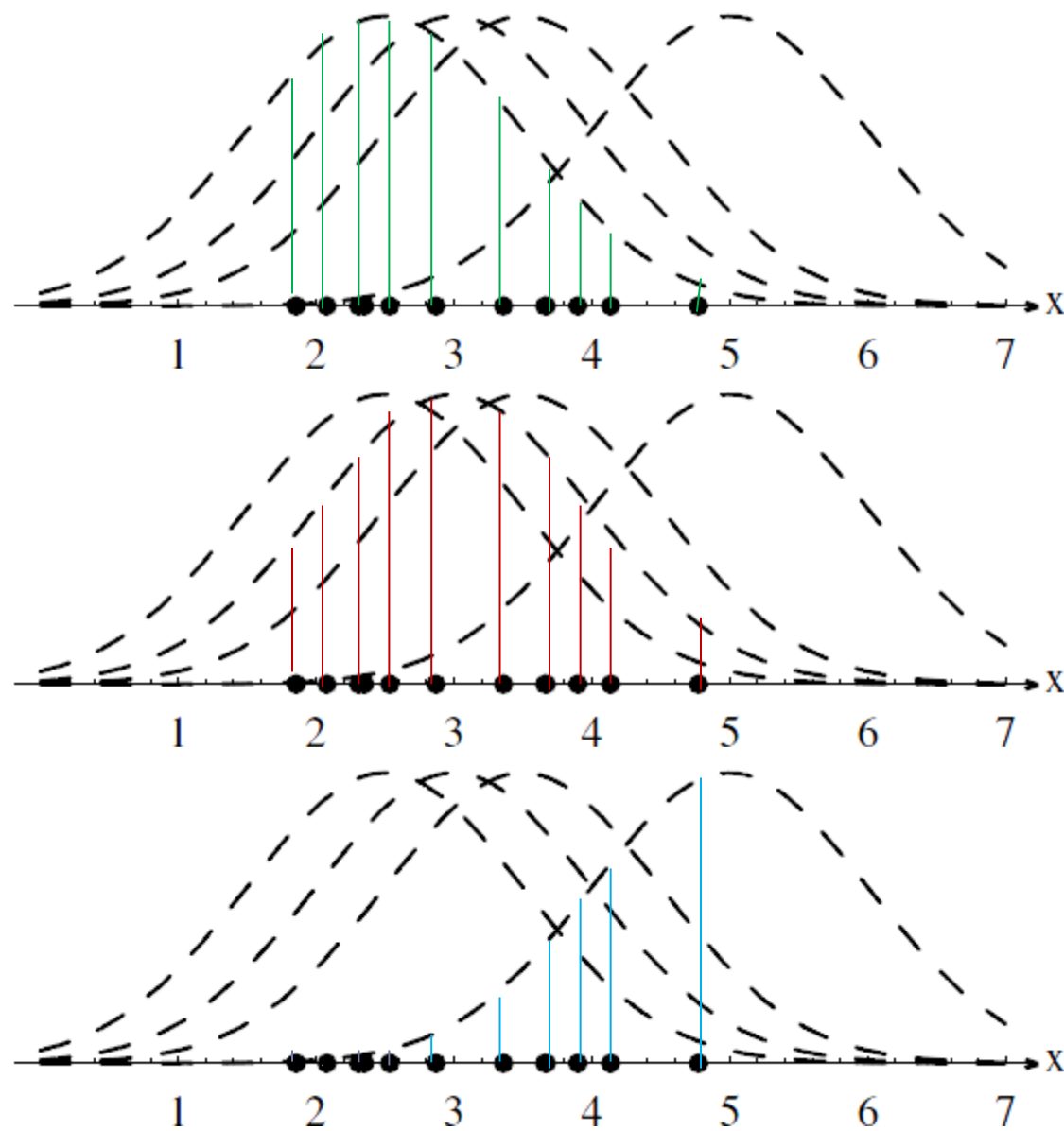
- ▶ Goal: estimate parameters of a distribution from a dataset $\mathcal{D} = \{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}\}$
 - ▶ \mathcal{D} contains N independent, identically distributed (i.i.d.) training samples.
- ▶ We need to determine $\boldsymbol{\theta}$ given $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}\}$
 - ▶ How to represent $\boldsymbol{\theta}$?
 - ▶ $\boldsymbol{\theta}^*$ or $p(\boldsymbol{\theta})$?

Example



$$P(x|\mu) = N(x|\mu, 1)$$


Example



Maximum Likelihood Estimation (MLE)

- ▶ Maximum-likelihood estimation (MLE) is a method of estimating the parameters of a statistical model given data.
- ▶ Likelihood is the conditional probability of observations $\mathcal{D} = \{\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(N)}\}$ given the value of parameters $\boldsymbol{\theta}$
 - ▶ Assuming i.i.d. observations:

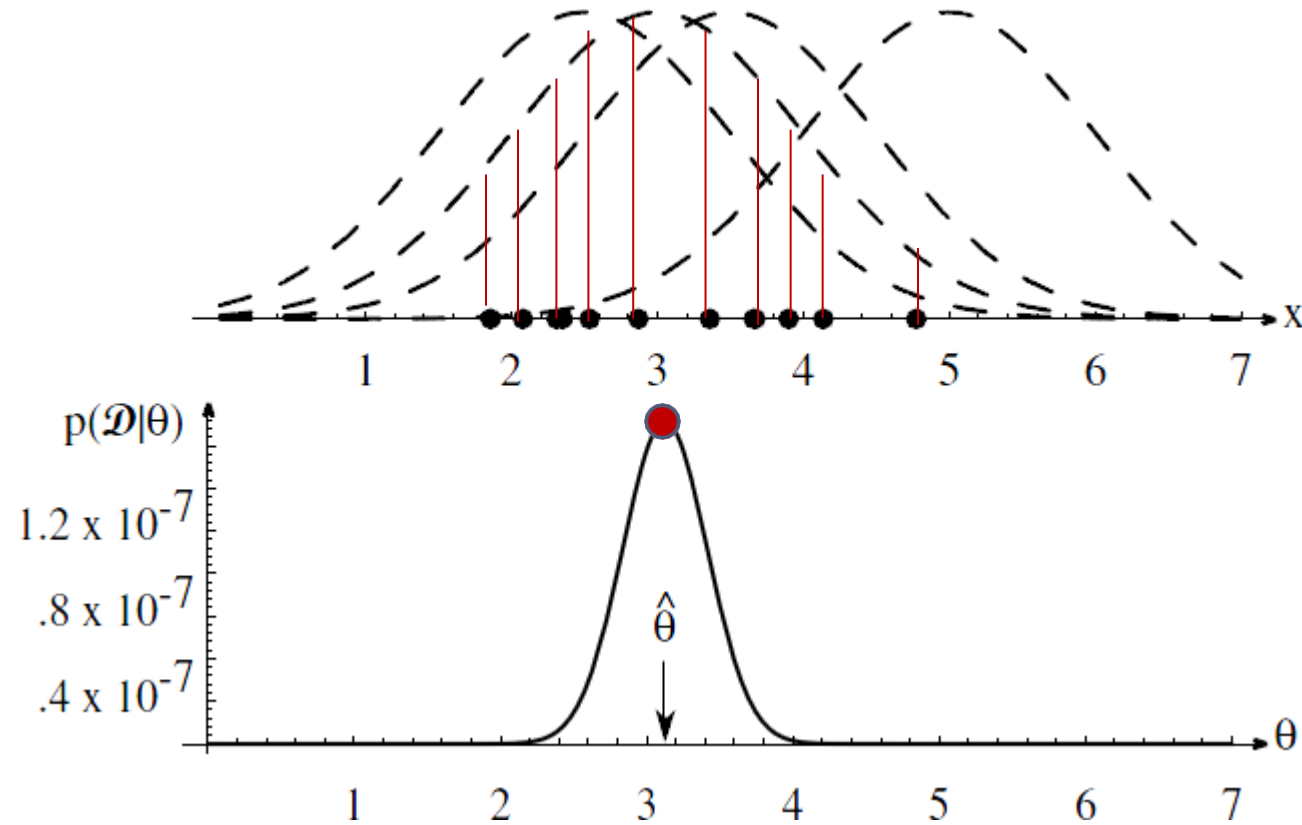
$$p(\mathcal{D}|\boldsymbol{\theta}) = \sum_{i=1}^N p(\mathbf{x}^{(i)}|\boldsymbol{\theta})$$


likelihood of $\boldsymbol{\theta}$ w.r.t. the samples

- ▶ Maximum Likelihood estimation

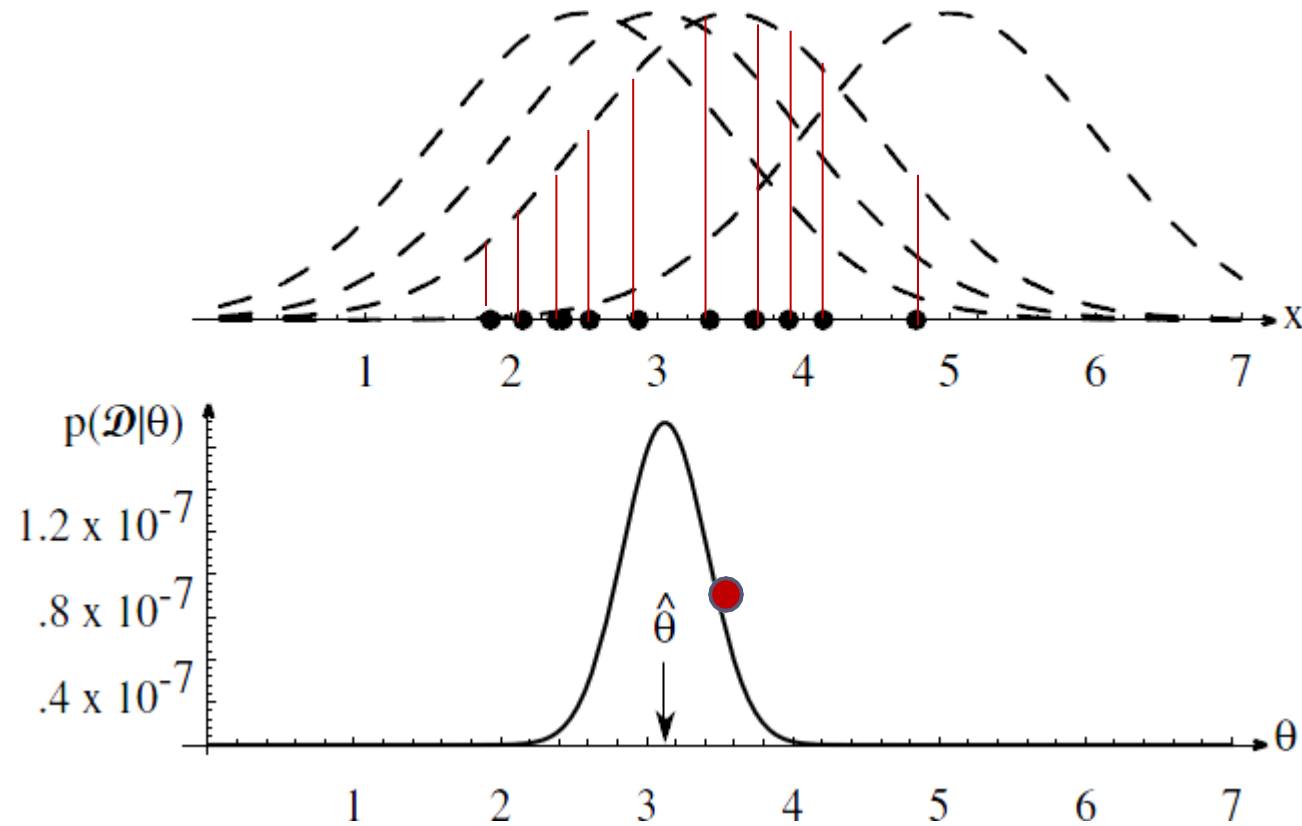
$$\boldsymbol{\theta}_{ML} = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} p(\mathcal{D}|\boldsymbol{\theta})$$

Maximum Likelihood Estimation (MLE)



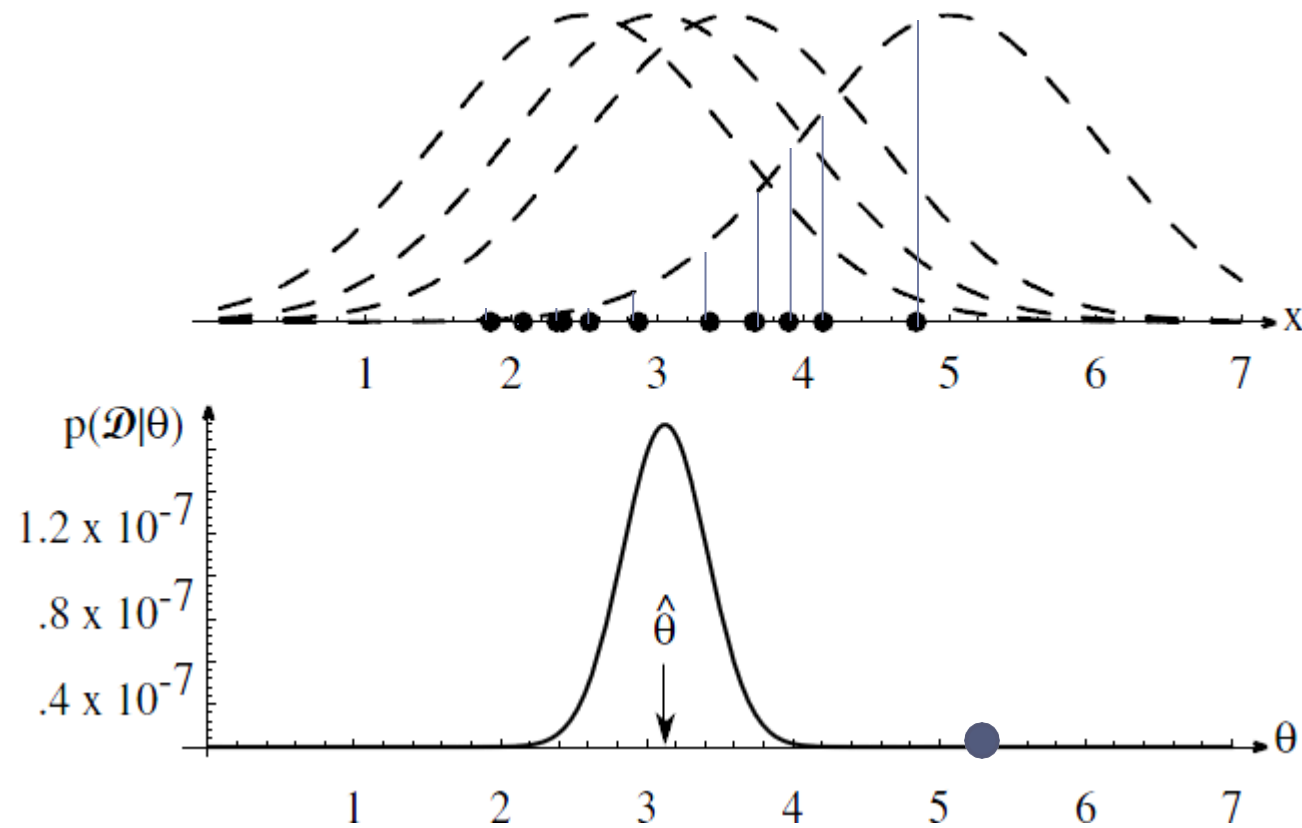
θ best agrees with the observed samples

Maximum Likelihood Estimation (MLE)



θ best agrees with the observed samples

Maximum Likelihood Estimation (MLE)



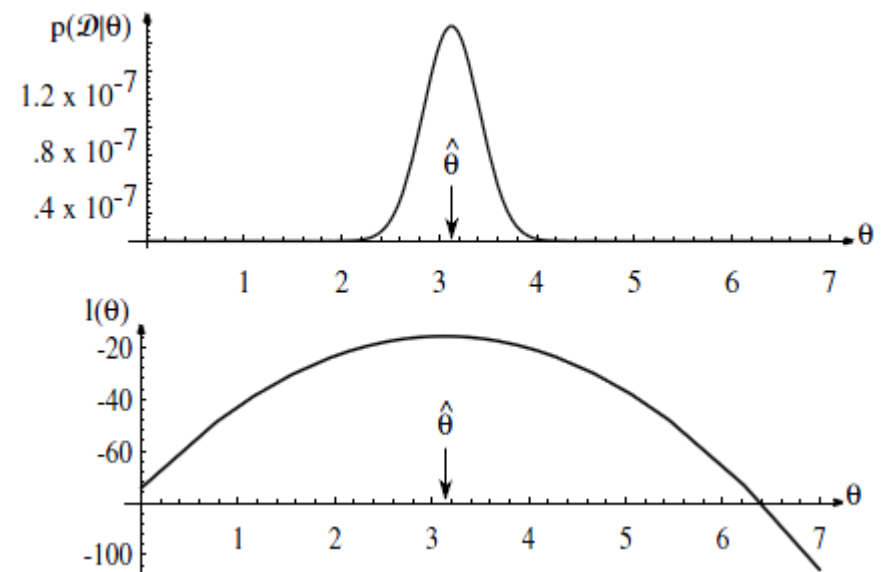
θ best agrees with the observed samples

Maximum Likelihood Estimation (MLE)

$$\mathcal{L}(\boldsymbol{\theta}) = \ln p(\mathcal{D}|\boldsymbol{\theta}) = \ln \prod_{i=1}^N p(\mathbf{x}^{(i)}|\boldsymbol{\theta}) = \sum_{i=1}^N \ln p(\mathbf{x}^{(i)}|\boldsymbol{\theta})$$

$$\boldsymbol{\theta}_{ML} = \operatorname{argmax}_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}) = \operatorname{argmax}_{\boldsymbol{\theta}} \sum_{i=1}^N \ln p(\mathbf{x}^{(i)}|\boldsymbol{\theta})$$

- ▶ Thus, we solve $\nabla_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}) = \mathbf{0}$ to find global optimum



MLE. Bernoulli

- ▶ Given: $\mathcal{D} = \{x^{(1)}, x^{(2)}, \dots, x^{(N)}\}$, m heads (1), $N - m$ tails (0)

$$p(x|\theta) = \theta^x (1 - \theta)^{1-x}$$

$$p(\mathcal{D}|\theta) = \prod_{i=1}^N p(x^{(i)}|\theta) = \prod_{i=1}^N \theta^{x^{(i)}} (1 - \theta)^{1-x^{(i)}}$$

$$\ln p(\mathcal{D}|\theta) = \sum_{i=1}^N \ln p(x^{(i)}|\theta) = \sum_{i=1}^N \{x^{(i)} \ln \theta + (1 - x^{(i)}) \ln(1 - \theta)\}$$

$$\frac{\partial \ln p(\mathcal{D}|\theta)}{\partial \theta} = 0 \Rightarrow \theta_{ML} = \frac{\sum_{i=1}^N x^{(i)}}{N} = \frac{m}{N}$$

MLE. Bernoulli: example

- ▶ Example: $\mathcal{D} = \{1,1,1\}$, $\theta_{ML} = \frac{3}{3} = 1$
 - ▶ Prediction: all future tosses will land heads up
- ▶ Overfitting to \mathcal{D}

Logistic regression

- ▶ More general than discriminant functions:
 - ▶ $f(\mathbf{x}; \mathbf{w})$ predicts posterior probabilities $P(y = 1|\mathbf{x})$

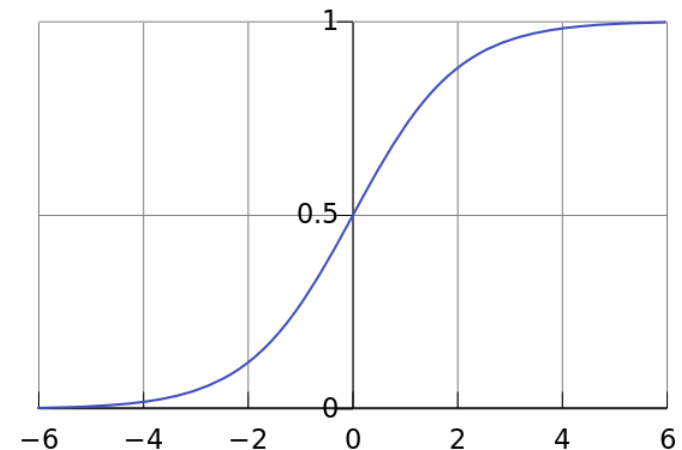
$$f(\mathbf{x}; \mathbf{w}) = \sigma(\mathbf{w}^T \mathbf{x})$$

$\mathbf{x} = [1, x_1, \dots, x_d]$
 $\mathbf{w} = [w_0, w_1, \dots, w_d]$

$\sigma(\cdot)$ is an activation function

- ▶ Sigmoid (logistic) function
 - ▶ Activation function

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



Logistic regression

- ▶ $f(\mathbf{x}; \mathbf{w})$: probability that $y = 1$ given \mathbf{x} (parameterized by \mathbf{w})

$$P(y = 1 | \mathbf{x}, \mathbf{w}) = f(\mathbf{x}; \mathbf{w})$$

$K = 2$
 $y \in \{0, 1\}$

$$P(y = 0 | \mathbf{x}, \mathbf{w}) = 1 - f(\mathbf{x}; \mathbf{w})$$

$$f(\mathbf{x}; \mathbf{w}) = \sigma(\mathbf{w}^T \mathbf{x})$$

$$0 \leq f(\mathbf{x}; \mathbf{w}) \leq 1$$

estimated probability of $y = 1$ on input \mathbf{x}

- ▶ Example: Cancer (Malignant, Benign)

- ▶ $f(\mathbf{x}; \mathbf{w}) = 0.7$

- ▶ 70% chance of tumor being malignant

Logistic regression: Decision surface

- ▶ Decision surface $f(\mathbf{x}; \mathbf{w}) = \text{constant}$
 - ▶ $f(\mathbf{x}; \mathbf{w}) = \sigma(\mathbf{w}^T \mathbf{x}) = \frac{1}{1+e^{-(\mathbf{w}^T \mathbf{x})}} = 0.5$
- ▶ Decision surfaces are linear functions of \mathbf{x}

if $f(\mathbf{x}; \mathbf{w}) \geq 0.5$ then $y = 1$
else $y = 0$

Equivalent to

if $\mathbf{w}^T \mathbf{x} + w_0 \geq 0$ then $y = 1$
else $y = 0$

Logistic regression: ML estimation

- ▶ Maximum (conditional) log likelihood:

$$\mathbf{w} = \underset{\mathbf{w}}{\operatorname{argmax}} \quad \operatorname{Log} \sum_{i=1}^n p(y^{(i)} | \mathbf{w}, \mathbf{x}^{(i)})$$

$$p(y^{(i)} | \mathbf{w}, \mathbf{x}^{(i)}) = f(\mathbf{x}^{(i)}; \mathbf{w})^{y^{(i)}} \left(1 - f(\mathbf{x}^{(i)}; \mathbf{w})\right)^{(1-y^{(i)})}$$

$$\begin{aligned} & \log p(\mathbf{y} | \mathbf{X}, \mathbf{w}) \\ &= \sum_{i=1}^n \left[y^{(i)} \log \left(f(\mathbf{x}^{(i)}; \mathbf{w}) \right) + (1 - y^{(i)}) \log \left(1 - f(\mathbf{x}^{(i)}; \mathbf{w}) \right) \right] \end{aligned}$$

Logistic regression: cost function

$$\mathbf{w} = \underset{\mathbf{w}}{\operatorname{argmin}} J(\mathbf{w})$$

$$\begin{aligned} J(\mathbf{w}) &= - \sum_{i=1}^n \log p(y^{(i)} | \mathbf{w}, \mathbf{x}^{(i)}) \\ &= \sum_{i=1}^n -y^{(i)} \log(f(\mathbf{x}^{(i)}; \mathbf{w})) - (1 - y^{(i)}) \log(1 - f(\mathbf{x}^{(i)}; \mathbf{w})) \end{aligned}$$

- ▶ No closed form solution for

$$\nabla_{\mathbf{w}} J(\mathbf{w}) = 0$$

- ▶ However $J(\mathbf{w})$ is convex.

Logistic regression: Gradient descent

$$\mathbf{w}^{t+1} = \mathbf{w}^t - \eta \nabla_{\mathbf{w}} J(\mathbf{w}^t)$$

$$\nabla_{\mathbf{w}} J(\mathbf{w}) = \sum_{i=1}^n (f(\mathbf{x}^{(i)}; \mathbf{w}) - y^{(i)}) \mathbf{x}^{(i)}$$

- Is it similar to gradient of SSE for linear regression?

$$\nabla_{\mathbf{w}} J(\mathbf{w}) = \sum_{i=1}^n (\mathbf{w}^T \mathbf{x}^{(i)} - y^{(i)}) \mathbf{x}^{(i)}$$

Logistic regression: loss function

$$\text{Loss}(y, f(\mathbf{x}; \mathbf{w})) = -y \times \log(f(\mathbf{x}; \mathbf{w})) - (1 - y) \times \log(1 - f(\mathbf{x}; \mathbf{w}))$$

$$\text{Since } y = 1 \text{ or } y = 0 \Rightarrow \text{Loss}(y, f(\mathbf{x}; \mathbf{w})) = \begin{cases} -\log(f(\mathbf{x}; \mathbf{w})) & \text{if } y = 1 \\ -\log(1 - f(\mathbf{x}; \mathbf{w})) & \text{if } y = 0 \end{cases}$$

How is it related to zero-one loss?

$$\text{Loss}(y, y) = \begin{cases} 1 & y \neq y \\ 0 & y = y \end{cases}$$

$$f(\mathbf{x}; \mathbf{w}) = \frac{1}{1 + \exp(-\mathbf{w}^T \mathbf{x})}$$

Logistic regression: cost function

- ▶ Logistic Regression (LR) has a more proper cost function for classification than SSE and Perceptron
- ▶ Why is the cost function of LR also more suitable than?

$$J(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n \left(y^{(i)} - f(\mathbf{x}^{(i)}; \mathbf{w}) \right)^2$$

where $f(\mathbf{x}; \mathbf{w}) = \sigma(\mathbf{w}^T \mathbf{x})$

The conditional distribution $p(y|\mathbf{x}, \mathbf{w})$ in the classification problem is not Gaussian (it is Bernoulli)

The cost function of LR is also convex

Logistic Regression (LR): summary

- ▶ LR is a linear classifier
- ▶ LR optimization problem is obtained by maximum likelihood
 - ▶ when assuming Bernoulli distribution for conditional probabilities whose mean is $\frac{1}{1+e^{-(\mathbf{w}^T \mathbf{x})}}$
- ▶ No closed-form solution for its optimization problem
 - ▶ But convex cost function and global optimum can be found by gradient ascent

Resource

- Yaser S. Abu-Mostafa, MalikMaghdon-Ismael, and Hsuan Tien Lin, “**Learning from Data**”, 2012.
- C. Bishop, “Pattern Recognition and Machine Learning”.