

Lecture 5: Linear Models Review

Rui Xia

School of Computer Science and Engineering
Nanjing University of Science and Technology

Linear Models (We Learnt So Far)

- Linear Regression
- Logistic Regression
- Perceptron Algorithm

3 Key Concepts in Machine Learning

- Hypothesis
 - Math models with (unknown) **parameters** (or structures)
- Learning (**to estimate the parameters**)
 - Maximum Likelihood Estimation (MLE), MAP, Bayesian Estimation
 - Cost Function Optimization
- Decision
 - Bayes decision rule
 - Direct prediction function

Model Hypothesis

- Linear Regression

$$h_{\theta}(x) = \theta^T x$$

- Perceptron Algorithm

$$h_{\theta}(x) = \begin{cases} 1 & \text{if } \theta^T x \geq 0 \\ 0 & \text{if } \theta^T x < 0 \end{cases}$$

- Logistic Regression

$$h_{\theta}(x) = \delta(\theta^T x) = \frac{1}{1 + e^{-\theta^T x}}$$

$$P(y = 1|x; \theta) = h_{\theta}(x)$$

$$P(y = 0|x; \theta) = 1 - h_{\theta}(x)$$

Learning Criteria (Cost Functions)

- Linear Regression

$$J_l(\theta) = \frac{1}{2} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

**Maximum Likelihood \Leftrightarrow
Least Mean Square**

Learning Criteria (Cost Functions)

- Perceptron Algorithm

$$\begin{aligned} J_p(\theta) &= \sum_{\substack{x^{(i)} \in M_0 \\ m}} \theta^T x^{(i)} - \sum_{x^{(j)} \in M_1} \theta^T x^{(j)} \\ &= \sum_{\substack{i=1 \\ m}} \left((1 - y^{(i)}) h_{\theta}(x^{(i)}) - y^{(i)} (1 - h_{\theta}(x^{(i)})) \right) \theta^T x^{(i)} \\ &= \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) \theta^T x^{(i)} \end{aligned}$$

Perceptron Criterion

Learning Criteria (Cost Functions)

- Logistic Regression

$$J_c(\theta) = \sum_{i=1}^m y^{(i)} \log h_{\theta}(x^{(i)}) + (1 - y^{(i)}) \log (1 - h_{\theta}(x^{(i)}))$$

**Maximum Likelihood \Leftrightarrow
Minimum Cross Entropy Error**

Gradient Descent Optimization

- Linear Regression

$$\begin{aligned}\frac{\partial}{\partial \theta} J_l(\theta) &= \frac{1}{2} \frac{\partial}{\partial \theta} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 \\ &= \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x^{(i)}\end{aligned}$$



$$\theta := \theta - \alpha \frac{\partial}{\partial \theta} J_l(\theta) = \theta - \alpha \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x^{(i)}$$

(Stochastic) Gradient Descent Optimization

- Perceptron Algorithm

$$\frac{\partial}{\partial \theta} J_p(\theta) = \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})x^{(i)}$$



$$\omega := \omega + \alpha(y - h_{\theta}(x))x$$

$$= \begin{cases} \omega + \alpha x & \text{if } y = 1 \text{ and } h_{\theta}(x) = 0 \\ \omega - \alpha x & \text{if } y = 0 \text{ and } h_{\theta}(x) = 1 \\ \omega & \text{others} \end{cases}$$

Gradient Descent Optimization

- Logistic Regression

$$\begin{aligned}\frac{\partial J_c(\theta)}{\partial \theta} &= \sum_{i=1}^m \left(y^{(i)} \frac{1}{h_\theta(x^{(i)})} - (1 - y^{(i)}) \frac{1}{1 - h_\theta(x^{(i)})} \right) \frac{\partial}{\partial \theta} h_\theta(x^{(i)}) \\ &= \sum_{i=1}^m \left(y^{(i)} \frac{1}{h_\theta(x^{(i)})} - (1 - y^{(i)}) \frac{1}{1 - h_\theta(x^{(i)})} \right) h_\theta(x^{(i)}) (1 - h_\theta(x^{(i)})) \frac{\partial}{\partial \theta} \theta^T x^{(i)} \\ &= \sum_{i=1}^m \left(y^{(i)} (1 - h_\theta(x^{(i)})) - (1 - y^{(i)}) h_\theta(x^{(i)}) \right) x^{(i)} \\ &= \sum_{i=1}^m \left(y - h_\theta(x^{(i)}) \right) x^{(i)}\end{aligned}$$



$$\theta := \theta + \alpha \sum_{i=1}^m \left(y^{(i)} - h_\theta(x^{(i)}) \right) x^{(i)}$$



Any Questions?