

Lecture 6

Some Concepts Revisit

Rui Xia

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

<http://www.nustm.cn/~rxia>

Outline

- Review of Linear Models
 - Linear Regression
 - Logistic Regression
 - Perceptron
- Generative vs. Discriminative
 - Hypothesis
 - Decision
 - Learning
- Over-fitting
 - ML - MAP
 - Regularization

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

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

- Perceptron Algorithm

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

- Logistic Regression

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

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

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

Learning Criteria (Cost Functions)

- Linear Regression

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

Maximum Likelihood \Leftrightarrow Least Mean Square

Learning Criteria (Cost Functions)

- Perceptron Algorithm

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

Perceptron Criterion

Learning Criteria (Cost Functions)

- Logistic Regression

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

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

Gradient Descent Optimization

- Linear Regression

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



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

(Stochastic) Gradient Descent Optimization

- Perceptron Algorithm

$$\frac{\partial}{\partial \mathbf{w}} J_p(\mathbf{w}) = \sum_{i=1}^N (h_{\mathbf{w}}(\mathbf{x}^{(i)}) - y^{(i)}) \mathbf{x}^{(i)}$$



$$\begin{aligned} \mathbf{w} &:= \mathbf{w} + \alpha(y - h_{\mathbf{w}}(\mathbf{x}))\mathbf{x} \\ &:= \begin{cases} \mathbf{w} + \alpha\mathbf{x}, & \text{if } y = 1 \text{ and } h_{\mathbf{w}}(\mathbf{x}) = 0 \\ \mathbf{w} + \alpha\mathbf{x}, & \text{if } y = 0 \text{ and } h_{\mathbf{w}}(\mathbf{x}) = 1 \\ \mathbf{w}, & \text{otherwise} \end{cases} \end{aligned}$$

Gradient Descent Optimization

- Logistic Regression

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



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

Outline

- Review of Linear Models
 - Linear Regression
 - Logistic Regression
 - Perceptron
- Generative vs. Discriminative
 - Hypothesis
 - Decision
 - Learning
- Over-fitting
 - ML - MAP
 - Regularization

Hypothesis - Learning - Decision

- Discriminative Model

- Directly Modeling Predictive Function

$$y = f(\mathbf{x})$$

Example:
Perceptron, SVMs

- Modeling Conditional Distribution

$$p(y|\mathbf{x})$$

Example:
Logistic/Softmax Regression

- Generative Model (Modeling Joint Distribution)

$$p(\mathbf{x}, y) = p(y)p(\mathbf{x}|y)$$

Examples:
Naïve Bayes, GMM

Hypothesis - Learning - Decision

- Discriminative Model
 - Modeling Predictive Function

$$\theta^* = \arg \max_{\theta} J(\theta)$$

Optimizing some loss functions, such as least mean square (LMS), cross entropy (CE), Maximum Margin, etc.

- Modeling Posterior Distribution

$$\theta^* = \arg \max_{\theta} \sum_i \log p(y^{(i)} | \mathbf{x}^{(i)})$$

ML, MAP (for posterior distribution)
 \Leftrightarrow Equivalent to some criteria in some cases

- Generative Model (Modeling Joint/Marginal Distribution)

$$\theta^* = \arg \max_{\theta} \sum_i \log p(\mathbf{x}^{(i)}, y^{(i)})$$

ML, MAP, Bayesian Inference (for joint or marginal distribution)

Hypothesis - Learning - Decision

- Discriminative Model
 - Conditional Distribution
 - Predictive Function

$$\arg \max_y p(y|\mathbf{x})$$

$$y = f(\mathbf{x})$$

- Generative Model
 - Bayes Formula

$$p(y|\mathbf{x}) = \frac{p(\mathbf{x}, y)}{p(\mathbf{x})}$$



$$\arg \max_y p(y|\mathbf{x}) = \arg \max_y p(\mathbf{x}, y) = \arg \max_y p(\mathbf{x}|y)p(y)$$

Generative Models for Classification

- Modeling Joint Distribution

$$p(\mathbf{x}, y = c_j) = p(c_j)p(\mathbf{x}|c_j)$$

Class-conditional
probability

Class prior probability

- Different Class-Conditional Distribution

- Multinomial Distribution

$$\begin{aligned} p(\mathbf{x}, c_j | \boldsymbol{\theta}) &= p(c_j | \boldsymbol{\theta}) p(\mathbf{x} | c_j; \boldsymbol{\theta}) \\ &= p_j \prod_{t=1}^M \theta_{t,j}^{N(w_t, \mathbf{x})} \end{aligned}$$

- Gaussian Distribution

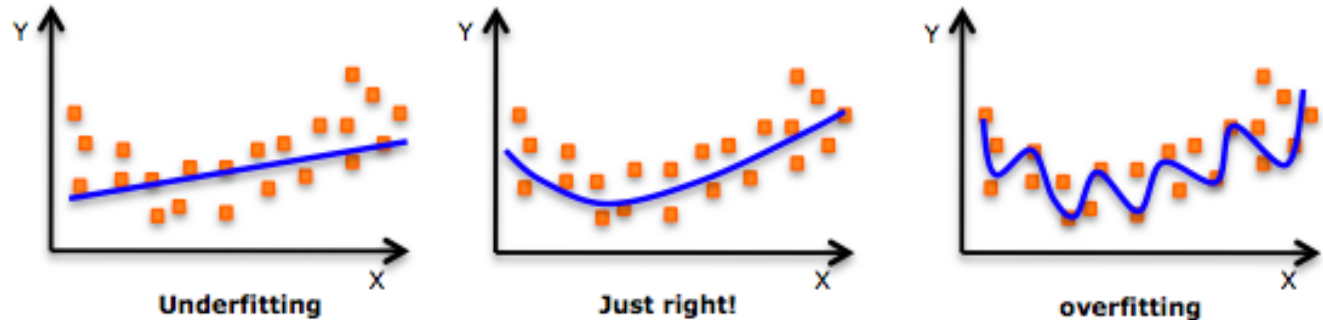
$$\begin{aligned} p(\mathbf{x}, c_j | \boldsymbol{\theta}) &= p(c_j | \boldsymbol{\theta}) p(\mathbf{x} | c_j; \boldsymbol{\theta}) \\ &= p_j N(\mathbf{x} | \mu_j, \Sigma_j) \end{aligned}$$

Outline

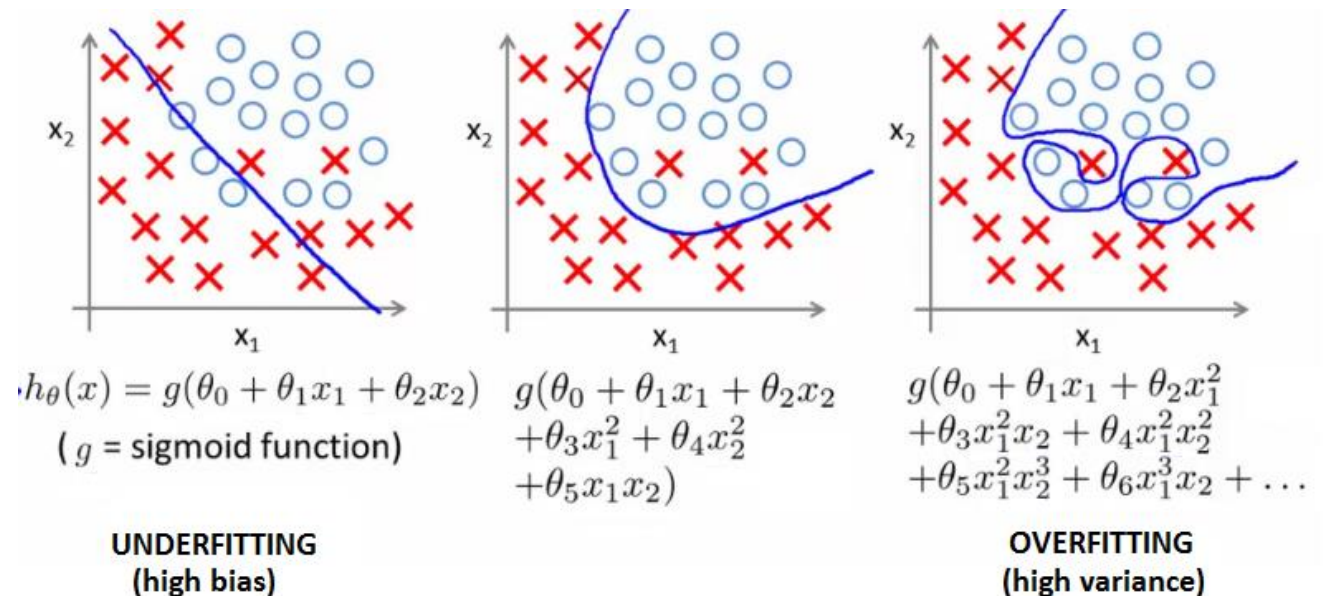
- Review of Linear Models
 - Linear Regression
 - Logistic Regression
 - Perceptron
- Generative vs. Discriminative
 - Hypothesis
 - Decision
 - Learning
- Over-fitting
 - ML - MAP
 - Regularization

Over-fitting

- Regression



- Classification



ML - MAP

- Maximum Likelihood (ML)

$$\begin{aligned}\theta_{ML}^* &= \arg \max_{\theta} L(\theta) = \arg \max_{\theta} p(X|\theta) \\ &= \arg \max_{\theta} \sum_{x \in X} \log p(x|\theta)\end{aligned}$$

likelihood

- Maximum A Posteriori (MAP)

$$\begin{aligned}\theta_{MAP}^* &= \arg \max_{\theta} p(\theta|X) = \arg \max_{\theta} \frac{p(X|\theta)p(\theta)}{p(X)} \\ &= \arg \max_{\theta} p(X|\theta) p(\theta) \\ &= \arg \max_{\theta} \sum_{x \in X} \log p(x|\theta) + \log p(\theta)\end{aligned}$$

likelihood • prior

Regularization

- ML - MAP

$$\theta_{ML}^* = \arg \max_{\theta} \sum_{x \in X} \log p(x|\theta)$$



$$\theta_{MAP}^* = \arg \max_{\theta} \sum_{x \in X} \log p(x|\theta) + \log p(\theta)$$

regularization term

- Loss function plus regularization

$$\theta^* = \arg \max_{\theta} J(\theta)$$



$$\theta^* = \arg \max_{\theta} \hat{J}(\theta) = \arg \max_{\theta} J(\theta) + \lambda R(\theta)$$

regularization term

Example: Polynomial Curve Fitting

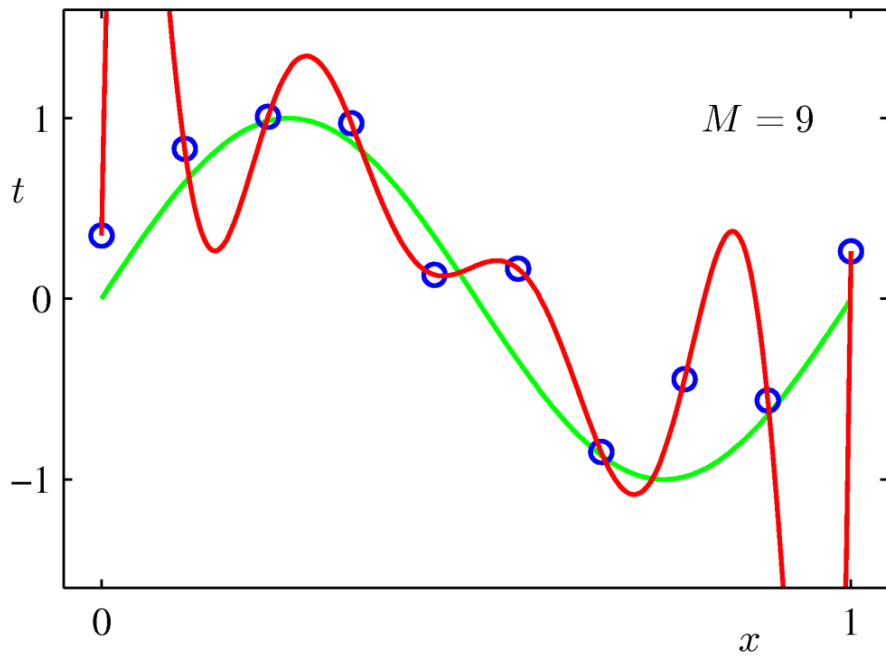
- ML (PRML Equation 1.62)

$$\ln p(t|\mathbf{x}, \mathbf{w}, \beta) = -\frac{\beta}{2} \sum_{n=1}^N \{y(\mathbf{x}_n, \mathbf{w}) - t_n\}^2 + \frac{N}{2} \ln \beta - \frac{N}{2} \ln(2\pi)$$

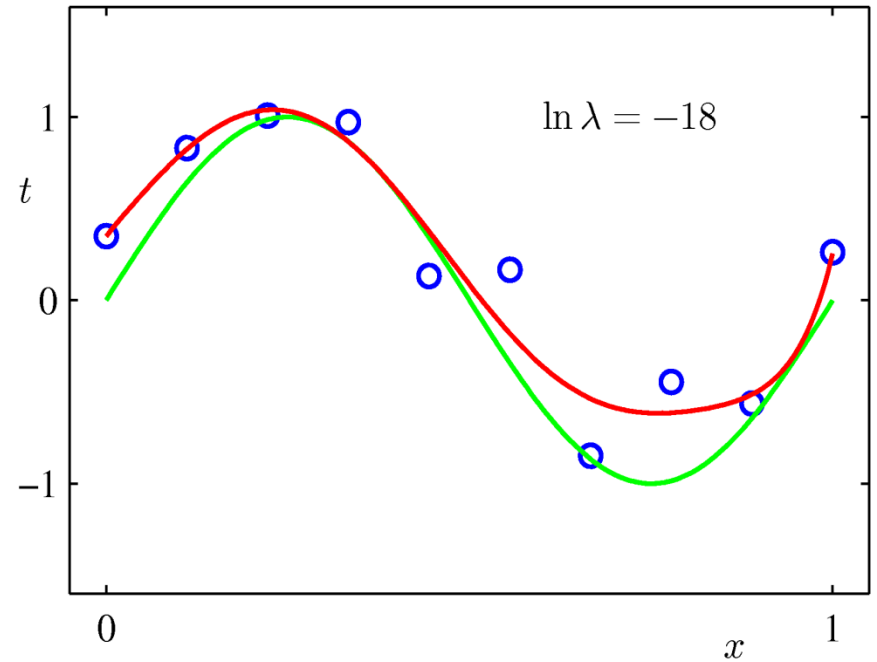
- MAP (PRML Equation 1.67)

$$\frac{\beta}{2} \sum_{n=1}^N \{y(\mathbf{x}_n, \mathbf{w}) - t_n\}^2 + \frac{\alpha}{2} \mathbf{w}^T \mathbf{w}$$

Example: Polynomial Curve Fitting



ML



MAP

Example: Logistic Regression

- ML

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

$$\text{where } h_{\boldsymbol{\theta}}(\mathbf{x}) = \frac{1}{1 + \exp^{-\boldsymbol{\theta}^T \mathbf{x}}}$$

- MAP

$$\hat{J}_c(\boldsymbol{\theta}) = \sum_{i=1}^N y^{(i)} \log h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) + (1 - y^{(i)}) \log (1 - h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)})) + \frac{1}{2} \|\boldsymbol{\theta}\|^2$$

Example: Bernoulli Experiments

- Bernoulli Distribution

$$\text{Bern}(x|\mu) = \mu^x (1 - \mu)^{1-x}, \quad x \in \{0,1\}$$

- Log-likelihood

$$\begin{aligned} \log L(\mu|X) &= \log \prod_{i=1}^N p(x_i|\mu) = \sum_{i=1}^N \log p(x_i|\mu) \\ &= m_1 \log p(1|\mu) + m_0 \log p(0|\mu) \\ &= m_1 \log \mu + m_0 \log(1 - \mu) \end{aligned}$$

- ML solution

$$\frac{\partial \log L}{\partial \mu} = \frac{m_1}{\mu} - \frac{m_0}{1 - \mu} = 0 \Leftrightarrow \hat{\mu}_{\text{ML}} = \frac{m - 1}{N}$$

Example: Bernoulli Experiments

- Prior Distribution

$$p(\mu|\alpha, \beta) = \frac{1}{B(\alpha, \beta)} \mu^{\alpha-1} (1 - \mu)^{\beta-1} \triangleq \text{Beta}(\mu|\alpha, \beta)$$

$$\text{where } B(\alpha, \beta) = \frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha + \beta)} \quad \Gamma(x) = \int_0^\infty u^{x-1} e^{-u} du$$

- MAP solution

$$\begin{aligned} \frac{\partial}{\partial \mu} \log L(\mu|X) + \log p(\mu) &= \frac{m_1}{\mu} - \frac{m_0}{1 - \mu} + \frac{\alpha - 1}{\mu} - \frac{\beta - 1}{1 - \mu} = 0 \\ \Leftrightarrow \hat{\mu}_{\text{MAP}} &= \frac{m_1 + \alpha - 1}{N + \alpha + \beta - 2} \end{aligned}$$

Outline

- Review of Linear Models
 - Linear Regression
 - Logistic Regression
 - Perceptron
- Generative vs. Discriminative
 - Hypothesis
 - Decision
 - Learning
- Over-fitting
 - ML - MAP
 - Regularization



Any Questions?