Regression Models for this week and next week:

Regression models
- Linear Regression → One variable ⎫
- Polynomial Regression → Multiple variable ⎬ → $L_2$ Regularization ⇒ Ridge Regression
- Logistic Regression ⎭ → $L_1$ Regularization ⇒ LASSO Regression

Four Components for Regression Models:

① Parameters : $\theta$

② Hypothesis : $h_\theta(X)$

③ Cost Function : $J(\theta)$

④ Goal : $\underset{\theta}{\text{minize}} \{J(\theta)\}$

Hierarchical structure

Something important about Gradient Descent:

① Updat simultaneously

② Learning Rate $\alpha$
- Too small: slow
- Too large: Diverge
- No need to decrease over time $\Rightarrow$ Auto take small steps
- Batch Gradient Descent:

- Debugging $\Rightarrow$ Make sure that $J(\theta)$ decrease on every iteration.

How to choose : 0.003
0.03
0.3

③ Normalize Features
- Feature Scaling ($-1 \leq X \leq 1$)
- Mean Normalization ($-0.5 \leq X \leq 0.5$)

④ In specific cases: Normal equation.

$$\theta = (X^T \cdot X)^{-1} \cdot X^T \cdot y$$

In Octave/Matlab: $\theta = pinv(X' * X) * X' * y$

When $(X^T \cdot X)$ is not invertible:

How to solve:  ① Delete some features $\Rightarrow$ avoid $\underbrace{(\text{linearly}}_{\text{A.}}$

$\underbrace{\text{dependent}}$ ) or $(n > \overset{B}{m}$ case)

$\Downarrow$ known as feature redundancy

② Use Regularization $\Rightarrow$

$$\theta = \left[ X^T X + \lambda \cdot \underbrace{\begin{bmatrix} 0 & & & \\ & 1 & & \\ & & \ddots & \\ & & & 1 \end{bmatrix}}_{(n+1) \times (n+1)} \right]^{-1} \cdot X^T y \text{ , where } \lambda > 0$$

Comparison with Gradient Descent:

Advantage:

① No need to choose $\alpha$

② No need to iterate

Disadvantage:  ① Slow when $n$ is large $\Rightarrow$

$O(n^3)$ while GD is $O(n^2)$

② $X^T X$ may non-invertible.

# Linear Regressions:

Cost function:

$$\frac{1}{2m} \cdot \left[ \sum_{i=1}^{m} \cdot \left( h_\theta(X^{(i)}) - y^{(i)} \right)^2 + \lambda \cdot \sum_{i=1}^{n} \cdot \theta_i^2 \right]$$

$$\underbrace{\qquad}_{} \qquad \underbrace{\qquad}_{\lambda \cdot \text{从} \theta \text{开始}}$$

$$\theta_0 + \theta_1 \cdot X_1^{(i)} + \theta_2 \cdot X_2^{(i)} \text{---} \theta_n \cdot X_n^{(i)}$$

Error function:

$$\frac{1}{2m_{test}} \cdot \sum_{i=1}^{m_{test}} \left( h_\theta(X^{(i)}) - y^{(i)} \right)^2$$

# Logistic Regressions:

Cost function: $\Rightarrow$ comes from cross entropy

$$\underline{\bigcirc} - \frac{1}{m} \cdot \sum_{i=1}^{m} \cdot \left[ y^{(i)} \cdot \log h_\theta(X^{(i)}) + \left(1 - y^{(i)}\right) \log\left( 1 - h_\theta(X^{(i)}) \right) \right]$$

Error function:

$$\frac{1}{m_{test}} \cdot \sum_{i=1}^{m_{test}} \underbrace{error\left( h_\theta(X^i), \; y^{(i)} \right)}_{}$$

$$h_\theta(X^{(i)}) > 0.5 \;\&\& \; y^{(i)} = 0 \qquad\qquad 0 \;, \text{otherwise}$$

$$\| \; h_\theta(X^{(i)}) < 0.5 \;\&\& \; y^{(i)} = 1$$

$$\Rightarrow \quad 1$$

Cost :

$$L(y, x, P_\theta) = -[y \cdot \log(\underbrace{h_\theta(x)}) + (1-y) \cdot \log(1-h_\theta(x))]$$

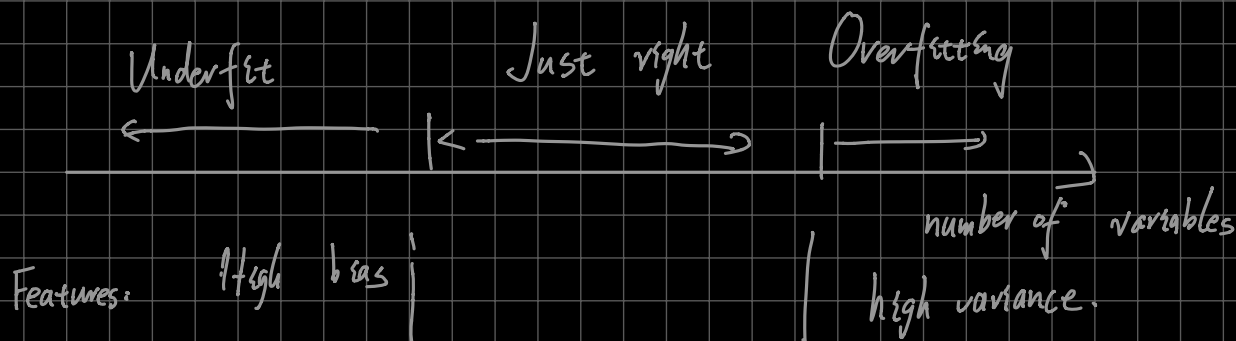$$= -[y \cdot \log(\sigma(\theta^T x)) + (1-y) \cdot \log(1-\sigma(\theta^T x))]$$

Gradient :

$$\frac{\partial L(y, x, P_\theta)}{\partial \theta} = -y \cdot \frac{1}{\sigma(\theta^T x)} \cdot \sigma(\theta^T x)(1-\sigma(\theta^T x)) \cdot x$$

$$- (1-y) \cdot \frac{1}{1-\sigma(\theta^T x)} \cdot -1 \cdot \sigma(\theta^T x)(1-\sigma(\theta^T x)) x$$

$$= (\sigma(\theta^T x) - y) \cdot x$$

Overfitting:

Causes: When the samples is not enough while the features are too much, it will cause overfitting. The curve try too hard to wiggly through all the training examples.

Underfit                Just right           Overfitting

$\leftarrow$ ———————— |$\leftarrow$ ————————— $\rightarrow$| |———————— $\rightarrow$|————————— $\rightarrow$

number of variables

Features:        High bias |                | high variance.

## Addressing Overfitting:

Method:

A. Reduce features ⟨ Manually pick features

Model selection Algorithms.

B. Regularization

⟹ keep all the features, but reduce (magnitude) of parameters $\theta_i$. This method is more suitable for handling multi-feature problems.

each feature is slightly useful.

Weight regularization
— L1 norm
— L2 norm

Drop out regularization
— Normal Dropout
— Data Augmentation
— Early stopping