# Coursera-Stanford-ML-Notes

Quentin Truong

20 June 2017 - ? July 2017

# Contents

# 1 Week 1: Introduction

## 1.1 Overview

– Machine Learning: "A computer program is said to learn from experience E with respect to some class of tasks T and performance measure P, if its performance at tasks in T, as measured by P, improves with experience E."

– Supervised Learning: know what our correct output looks like
  • Regression: want continuous output
  • Classification: want discrete output

– Unsupervised Learning: little or no idea what our results should look like
  • Clustering: find groups according to similarity in various variables
  • Nonclustering: find structure in chaos

# 2 Week 2: Linear Regression with Multiple Variables

## 2.1 Overview

– Use linear regression for continuous output

– Choose gradient descent if many features (million+) because the inverse matrix required for the normal equation can become expensive to compute

– Normal equation will directly compute theta

– Normalize features if using gradient descent

## 2.2 Symbols

$$m = number\ of\ samples$$
$$n = number\ of\ feature$$
$$x = (n \times 1)$$
$$X = (m \times n)$$
$$X_j = (m \times 1)$$
$$\theta = (n \times 1)$$
$$\theta_j = (1 \times 1)$$

## 2.3 Gradient Descent

| | |
|---|---|
| Hypothesis Function | $h_\theta(x) = \theta^\mathsf{T} \times x$ |
| Vectorized Hypothesis Function | $h_\theta(X) = X \cdot \theta$ |
| Linear Regression Cost Function | $J(\theta) = \dfrac{1}{2m} \sum (h_\theta(X) - y)^2$ |
| Derivative of Linear Regression CF wrt $\theta_j$ | $\dfrac{\partial}{\partial \theta_j} J(\theta) = \dfrac{1}{m} \sum (h_\theta(X) - y) \mathbin{.*} X_j$ |
| Change in $\theta_j$ | $\theta_j = \theta_j - \alpha \dfrac{\partial}{\partial \theta_j}$ |
| | $= \theta_j - \alpha \dfrac{1}{m} \sum (h_\theta(X) - y) \mathbin{.*} X_j$ |
| Vectorized Change in $\theta$ | $\theta = \theta - \alpha \dfrac{1}{m} X^\mathsf{T}(X \cdot \theta - y)$ |

## 2.4 Normal Equation

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

# 3 Week 3: Logistic Regression

## 3.1 Overview

– Use logistic regression for discrete output (classification)
  - $h_\theta(x) = (y = 1|x; \theta)$; gives probability that the output is 1
  - For multi-class classification, use one-vs-all
  - Sigmoid/Logistic function maps any real number to (0, 1)
  - Pick class i that maximizes $h_\theta^i(x)$

– Overfitting is when learned hypothesis fits training data well but fails to generalize
  - Underfitting is when doesn't fit training data

– Address overfitting by reducing number of features, model selection, and regularization
  - Regularization results in simpler hypothesis and less overfitting
  - Extremely large $\lambda$ will result in underfitting and gradient descent will fail to converge

– Use other prewritten optimization algorithims (conjugate gradient, BFGS, L-BFGS) because they are faster

## 3.2 Logistic Regression Hypothesis Function

$$\text{Sigmoid/Logistic Function} \qquad g(z) = \frac{1}{1 + e^{-z}}$$

$$\text{Hypothesis Function} \qquad h_\theta(x) = g(\theta^\mathsf{T} x)$$

$$= \frac{1}{1 + e^{-\theta^\mathsf{T} x}}$$

## 3.3 Logistic Regression Cost Function

$$Cost(h_\theta(x), y) = \begin{cases} -\log(h_\theta(x)) \text{ if } y = 1 \\ -\log(1 - h_\theta(x)) \text{ if } y = 0 \end{cases}$$

$$= -y\log(h_\theta(x)) - (1 - y)\log(1 - h_\theta(x))$$

$$J(\theta) = \frac{1}{m} \sum_{i=1}^{m} \text{Cost}(h_\theta(x^i), y^i)$$

$$J(\theta) = \frac{-1}{m} \sum_{i=1}^{m} \left[ y^i \log(h_\theta(x^i)) + (1 - y^i)\log(1 - h_\theta(x^i)) \right]$$

## 3.4 Proof of Logistic Regression Cost Function Derivative

$$J(\theta) = \frac{-1}{m} \sum_{i=1}^{m} [y^i \log(h_\theta(x^i)) + (1 - y^i) \log(1 - h_\theta(x^i))]$$

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

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

$$J(\theta) = -\frac{1}{m} \sum_{i=1}^{m} \left[ -y^i (\log(1 + e^{-\theta x^i})) + (1 - y^i)(-\theta x^i - \log(1 + e^{-\theta x^i})) \right]$$

$$= -\frac{1}{m} \sum_{i=1}^{m} \left[ y^i \theta x^i - \theta x^i - \log(1 + e^{-\theta x^i}) \right]$$

$$= -\frac{1}{m} \sum_{i=1}^{m} \left[ y^i \theta x^i - \log(e^{\theta x^i}) - \log(1 + e^{-\theta x^i}) \right]$$

$$= -\frac{1}{m} \sum_{i=1}^{m} \left[ y^i \theta x^i - \log(1 + e^{\theta x^i}) \right]$$

$$\frac{\partial}{\partial \theta_j} y^i \theta x^i = y^i x_j^i$$

$$\frac{\partial}{\partial \theta_j} \log(1 + e^{\theta x^i}) = \frac{x_j^i e^{\theta x^i}}{1 + e^{\theta x^i}}$$

$$= \frac{x_j^i}{1 + e^{-\theta x^i}}$$

$$= x_j^i h_\theta(x^i)$$

$$\frac{\partial}{\partial \theta_j} J(\theta) = -\frac{1}{m} \sum_{i=1}^{m} \left[ y^i x_j^i - x_j^i h_\theta(x^i) \right]$$

$$\frac{\partial}{\partial \theta_j} J(\theta) = \frac{1}{m} \sum_{i=1}^{m} \left[ h_\theta(x^i) - y^i \right] x_j^i$$

## 3.5 Regularization

$$\text{Regularizing Term} \quad \lambda \sum_{j=1}^{n} \theta_j^2$$

$$\text{Regularized Linear Regression CF} \quad J(\theta) = \frac{1}{2m} \sum_{i=1}^{m} (h_\theta(x^i) - y^i)^2 + \lambda \sum_{j=1}^{n} \theta_j^2$$

$$\text{Regularized Logistic Regression CF} \quad J(\theta) = \frac{-1}{m} \sum_{i=1}^{m} \left[ y^i \log(h_\theta(x^i)) + (1 - y^i) \log(1 - h_\theta(x^i)) \right] + \frac{\lambda}{2m} \sum_{j=1}^{n} \theta_j^2$$

$$\text{Regularized GD (Lin/Log Regression)} \quad \begin{cases} \theta_0 = \theta_j - \alpha \left[ \frac{1}{m} \sum_{i=1}^{m} (h_\theta(x^i) - y^i) x_0^i \right] \\ \theta_j = \theta_j - \alpha \left[ \frac{1}{m} \sum_{i=1}^{m} (h_\theta(x^i) - y^i) x_j^i + \frac{\lambda}{m} \theta_j \right] \text{ (j=1,2,...,n)} \end{cases}$$

$$\text{Regularized Normal Equation} \quad \theta = (X^\intercal X + \lambda \begin{bmatrix} 0 & 0 & \cdots & 0 \\ 0 & 1 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & 1 \end{bmatrix}_{n+1,n+1})^{-1} X^\intercal y$$