

# Machine Learning

## Equations and Definitions

Daniel Hamacher

August 6, 2018

## Contents

<b>1</b>	<b>Linear Regression with one Variable</b>	<b>2</b>
1.1	Notation . . . . .	2
1.2	Function Definitions . . . . .	2
1.3	Updating Parameters . . . . .	2
<b>2</b>	<b>Linear Regression with Multiple Variables</b>	<b>3</b>
2.1	Notation . . . . .	3
2.2	Function Definitions . . . . .	3
2.3	Updating Parameters . . . . .	4
2.4	Using the Normal Equation instead . . . . .	4
<b>3</b>	<b>Logistic Regression</b>	<b>5</b>
3.1	Notation . . . . .	5
3.2	Function Definitions . . . . .	6
3.3	Updating Parameters . . . . .	6
3.4	Further Readings . . . . .	6

# 1 Linear Regression with one Variable

## 1.1 Notation

$$m = \text{The number of training samples} \quad (1)$$

$$x = \text{Input variables / features} \quad (2)$$

$$y = \text{Output variable / target variable} \quad (3)$$

$$\alpha = \text{The learning rate} \quad (4)$$

$$h_{\theta}(x) = \text{The hypothesis function} \quad (5)$$

$$J(\theta_0, \theta_1) = \text{The cost function} \quad (6)$$

$$\theta_0, \theta_1 = \text{The parameters (gradients)} \quad (7)$$

## 1.2 Function Definitions

$$h_{\theta}(x) = \theta_0 + \theta_1 x \quad (8)$$

$$J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 \quad (9)$$

## 1.3 Updating Parameters

$$\theta_0 = \theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1) = \theta_0 - \alpha \left( \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) \right) \quad (10)$$

$$\theta_1 = \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1) = \theta_1 - \alpha \left( \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x^{(i)} \right) \quad (11)$$

**Note** that the parameters must be updated simultaneously.

## 2 Linear Regression with Multiple Variables

### 2.1 Notation

$m$  = The number of training samples

$n$  = The number of features/variables

$x^{(i)}/\theta^{(i)}$  = Input variables/parameters of  $i^{th}$  training sample

$x_j^{(i)}/\theta_j^{(i)}$  = Value of feature/parameter  $j$  in  $i^{th}$  training sample

$x_0^{(i)}$  = Set to 1 for convenience

$$x = \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix} \in \mathbb{R}^{n+1} = (n+1, 1) \text{ input vector}$$

$$X = \begin{bmatrix} x_0^T \\ x_1^T \\ \vdots \\ x_m^T \end{bmatrix} \in \mathbb{R}^{n+1} = (m, n+1) \text{ design matrix}$$

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix} \in \mathbb{R}^{n+1} = (n+1, 1) \text{ parameter vector}$$

$y$  = Output variable / target variable

$\alpha$  = The learning rate

$h_\theta(x)$  = The hypothesis function

$J(\theta)$  = The cost function

### 2.2 Function Definitions

$$h_\theta(x) = \theta_0 x_0 + \theta_1 x_1 + \cdots + \theta_n x_n = \theta^T x \quad (12)$$

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2 \quad (13)$$

## 2.3 Updating Parameters

$$\theta = \begin{cases} \theta_0 =: & \theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\theta) = \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_0^{(i)} \\ \theta_1 =: & \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta) = \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_1^{(i)} \\ \theta_2 =: & \theta_2 - \alpha \frac{\partial}{\partial \theta_2} J(\theta) = \theta_2 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_2^{(i)} \\ \dots & \\ \theta_{n+1} =: & \theta_{n+1} - \alpha \frac{\partial}{\partial \theta_{n+1}} J(\theta) = \theta_{n+1} - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_{n+1}^{(i)} \end{cases}$$

**Note** to utilize vectorization to increase performance and effectiveness.

## 2.4 Using the Normal Equation instead

The normal equation can be used to find the minimum value for  $J(\theta)$ . However, the normal equation will not work for design matrices that are not invertible i.e. singular. The normal equation is suitable for machine learning solutions that only have a small set of features. The normal equation uses the design matrix in the 2.1 and looks like this:

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

With the normal equation there is no need to choose  $\alpha$  and finding the minimum value requires a matrix operation rather than iteration. However, this approach is slow when the model has more than  $10^6$  features. It is recommended to use gradient descent for those models. It is important to remember that the expression  $X^T X$  could end up to be a non-invertible matrix. In that case, investigate the design matrix on linear dependent features (redundant) and see if there are some features that can be deleted or use regularization.

## 3 Logistic Regression

### 3.1 Notation

$m$  = The number of training samples

$n$  = The number of features/variables

$x^{(i)}/\theta^{(i)}$  = Input variables/parameters of  $i^{th}$  training sample

$x_j^{(i)}/\theta_j^{(i)}$  = Value of feature/parameter  $j$  in  $i^{th}$  training sample

$x_0^{(i)}$  = Set to 1 for convenience

$x = \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix} \in \mathbb{R}^{n+1} = (n+1, 1)$  input vector

$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix} \in \mathbb{R}^{n+1} = (n+1), 1)$  parameter vector

$y$  = Output variable / target variable

$\alpha$  = The learning rate

$h_\theta(x)$  = The hypothesis function

$J(\theta)$  = The cost function

if  $h_\theta(x) \geq 0.5$ , predict " $y = 1$ "

if  $h_\theta(x) < 0.5$ , predict " $y = 0$ "

### 3.2 Function Definitions

The hypothesis is based on the sigmoid function

$$\begin{aligned}g(z) &= \frac{1}{1 + e^{-z}} \\h_{\theta}(x) &= g(\theta^T x) \\J(h_{\theta}(x), y) &= \begin{cases} -\log(h_{\theta}(x)) & \text{if } y = 1 \\ -\log(1 - h_{\theta}(x)) & \text{if } y = 0 \end{cases} \\J(\theta) &= -\frac{1}{m} \left[ \sum_{i=1}^m y^{(i)} \log h_{\theta}(x^{(i)}) + (1 - y^{(i)}) \log(1 - h_{\theta}(x^{(i)})) \right] \\J(\theta) &= -\frac{1}{m} \left[ \sum_{i=1}^m y^{(i)} \log h_{\theta}(x^{(i)}) + (1 - y^{(i)}) \log(1 - h_{\theta}(x^{(i)})) \right] + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2 (\text{Regularization Term})\end{aligned}$$

**NOTE:** There is no need to regularize  $j = 0$  for  $\theta_j^2$ , since it is used for the bias term.

### 3.3 Updating Parameters

$$\begin{aligned}\theta_j &= \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta), \text{ where} \\ \frac{\partial}{\partial \theta} J(\theta) &= \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)} \\ \frac{\partial}{\partial \theta} J(\theta) &= \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)} - \frac{\lambda}{m} \theta_j (\text{regularized})\end{aligned}$$

**NOTE:** There is no need to regularize  $j = 0$  for  $\theta_j^2$ , since it is used for the bias term.

### 3.4 Further Readings

This section contains links to further read about machine learning