

Machine Learning

**Nagubandi Krishna Sai
MS20BTECH11014**

June 2021

Contents

1	Fundamentals of Machine Learning	2
1.1	Supervised Learning	2
1.1.1	Linear Regression	2
1.1.2	Hypothesis	3
1.1.3	Cost function	4
1.1.4	Gradient descent	4
1.1.5	Linear Regression for multivariables	5
1.1.6	Polynomial regression	5
1.1.7	Normal Equation	6
1.1.8	Classification	7
1.1.9	Logistic Regression Model	8

Chapter 1

Fundamentals of Machine Learning

1.1 Supervised Learning

Supervised Learning gives "correct answers", the output values are same as real life values.

In Supervised Learning, we are given a set of data and we know what our correct output should look like, having an idea that there is relationship between the input and the output.

Supervised Learning problems has two types of problems,

1. Regression.
2. Classification.

1.1.1 Linear Regression

In regression type of problems, we are trying to predict results within a continuous output, means that we are trying to map input variables to some continuous function.

Linear regression has real-valued output, but the output will be same or near valued to the actual output.

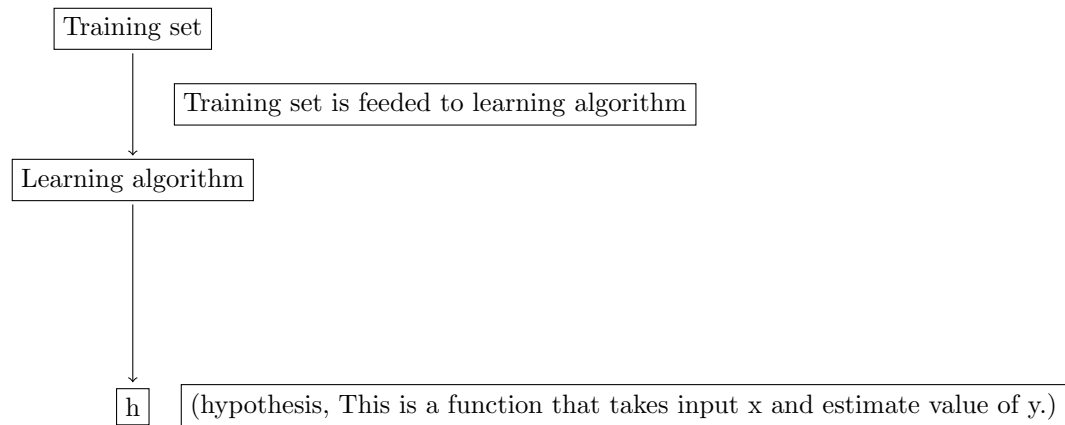
1. m = Number of training examples.
2. x 's = "input" variable (or) feature.
3. y 's = "output (or) target" variable.
4. (x,y) = one training example.
5. $(x^{(i)}, y^{(i)}) = i^{th}$ training example.

Example :

Size of feet ² (x)	Price(\$) in 1000's (y)
2104	460
1416	232
1534	315
852	178
\vdots	\vdots

Table 1.1: Training set of housing prices.

1. $m = 47$.
2. $x^{(1)} = 2104$ and $y^{(1)} = 460$.
3. $x^{(2)} = 1416$ and $y^{(2)} = 232$.



1.1.2 Hypothesis

$$h_{\theta}(x) = \theta_0 + \theta_1 \cdot x \quad (1.1)$$

θ_i 's = parameters.

This type of hypothesis mode is called "Linear regression with one variable" (or) "Univariate linear regression."

1.1.3 Cost function

Cost function helps us know that how well to fit the best possible straight line over the given data.

Q) How to choose θ_i 's ?

1. Choose θ_i 's so that $h_\theta(x)$ is close to y for our training example (x,y) .
2. minimise θ_0, θ_1 so, that $[h_\theta(x) - y]$ is small.

$$J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m [h_\theta(x^{(i)}) - y^{(i)}]^2 \quad (1.2)$$

$J(\theta_0, \theta_1)$ = Cost function (or) Squared error function.

1.1.4 Gradient descent

Gradient descent is used to minimise cost function(J) in linear regression.

Gradient descent is used in many areas to minimise many functions in ML/AI.

Gradient descent algorithm,

$$\text{Repeat until convergence (minimum)} \left\{ \theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1), \text{ for } j = 0, j = 1. \right. \quad (1.3)$$

1. $:=$ is Assignment operator.
2. α is learning rate.
3. $\frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1)$ is derivative.

Gradient descent is nothing but the derivative of the Cost function.

$$\text{Slope of cost function curve} = \frac{\partial J(\theta_1)}{\partial \theta_1}, \text{ when } \theta_0 = 0. \quad (1.4)$$

Learning rate,

1. If α is too small, gradient descent can be slow. After many such operations(can be infinite times), the ' θ_1 ' could reach "global minimum".
2. If α is too large, gradient descent can overshoot the minimum. It may "fail to converge (or) even diverge".
3. If θ_1 is at the local optima itself when we started or taken θ_1 , then there is no use of " α (or) gradient descent".

1.1.5 Linear Regression for multivariables

Hypothesis,

$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n. \quad (1.5)$$

In the total context of Supervised learning, hypothesis is just predicting the output.

For convenience of notation, declare $x_0 = 1$ ($x_0^{(i)} = 1$).

$$x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \in \mathbb{R}^{n+1} \quad \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_n \end{bmatrix} \in \mathbb{R}^{n+1}$$

The above matrix is 0 - indexed.

$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n. h_{\theta}(x) = \theta^{\top} x. \quad (1.6)$$

Cost function,

$$J(\theta) = J(\theta_0, \theta_1, \theta_2, \dots, \theta_n) = \frac{1}{2m} \sum_{i=1}^m [h_{\theta}(x^{(i)}) - y^{(i)}]^2 \quad (1.7)$$

Gradient descent,

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1, \theta_2, \dots, \theta_n) \quad (1.8)$$

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta) \quad (1.9)$$

$$\theta_j := \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m [h_{\theta}(x^{(i)}) - y^{(i)}] x_j^{(i)} \quad (1.10)$$

Feature scaling,

Get every feature into approximately $-1 \leq x_i \leq 1$ range.

Mean normalization,

Replace x_i with $x_i - \mu_i$ to make features have approximately zero mean.

1.1.6 Polynomial regression

Hypothesis,

$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3. \quad (1.11)$$

Housing price prediction
$x_1 = size$
$x_2 = (size)^2$
$x_3 = (size)^3$

Table 1.2: Features be-like in Polynomial regression.

1.1.7 Normal Equation

Intuition,

$$\frac{d}{d\theta} J(\theta) = 0.$$

Cost function,

$$\theta \in \mathbb{R}^{n+1}, J(\theta_0, \theta_1, \theta_2, \dots, \theta_n) = \frac{1}{2m} \sum_{i=1}^m [h_{\theta}(x^{(i)}) - y(i)]^2 \quad (1.12)$$

$$\frac{\partial}{\partial \theta_j} J(\theta) = 0, \text{ (solve for } \theta_0, \theta_1, \dots, \theta_n) \quad (1.13)$$

Example,

	Size (feet ²)	No.of Bed rooms	No.of floors	Age of home (years)	Price(\$1000)
x_0	x_1	x_2	x_3	x_4	y
1	2104	5	1	45	460
1	1416	3	2	40	232
1	1534	3	2	30	315
1	1852	2	1	36	178

Table 1.3: Sample Training set for Multi-Variate Linear regression.

1. n = number of Features.
2. $x_j^{(i)}$ = value of j in the i^{th} training example.
3. $x_{(i)}$ = the input(features) of the i^{th} training example.

$$X = \begin{bmatrix} 1 & 2104 & 5 & 1 & 45 \\ 1 & 1416 & 3 & 2 & 40 \\ 1 & 1534 & 3 & 2 & 30 \\ 1 & 1852 & 2 & 1 & 36 \end{bmatrix} \quad Y = \begin{bmatrix} 460 \\ 232 \\ 315 \\ 178 \end{bmatrix}$$

X is $m \times (n+1)$ -dimensional matrix and Y is a m -dimensional vector.

$$\theta = (X^T X)^{-1} X^T Y \quad (1.14)$$

The above θ value is optimal θ value.

For Normal equation method, then no need to use **feature scaling**.

We use 'Gradient descent' and 'Normal equation' methods to minimise cost function.

Gradient descent	Normal equation
1) Need to choose ' α '.	1) No need to choose ' α '.
2) Need many iterations.	2) Don't need to iterate.
3) Works well even, when 'n' is large. ($n \leq 10000$)	3) Need to compute $n \times n$ matrix inverse $(X^T X)^{-1}$
	4) Works Now if n is very large.

Table 1.4: Why should we use the particular method? Advantages and Disadvantages of two methods.

$$\theta = (X^T X)^{-1} X^T Y \quad (1.15)$$

Q) What is $X^T X$ is non-invertible(singular/degenerate) ?
Reasons,

1. Redundant features (linearly dependent)
 - $x_1 = \text{Size in feet}^2$
 - $x_2 = \text{Size in } m^2$
 - $x_2 = (3.28)^2 x_1$, $1m = 3.28\text{feet}$.
2. Too many features. ($m \leq n$)
 - $m = 10$
 - $n = 100$, $\theta \in \mathbb{R}^{101}$, Delete some features (or) Regularization.

1.1.8 Classification

The output value 'y' is **discrete value**.

The algorithm used is **logistic regression**.

1.1.9 Logistic Regression Model

We want $0 \leq h_\theta(x) \leq 1$.

$$h_\theta(x) = g(\theta^\top x) \quad (1.16)$$

$$g(z) = \frac{1}{1 + e^{-z}} \quad (1.17)$$

$$h_\theta(x) = g(\theta^\top x) \quad (1.18)$$

$$= \frac{1}{1 + e^{-\theta^\top x}} \quad (1.19)$$

The above $g(z)$ is called sigmoid function (or) logistic function.

Graph,

