# Contents

Machine Learning
(Andrew Ng)

# CHAPTER 1

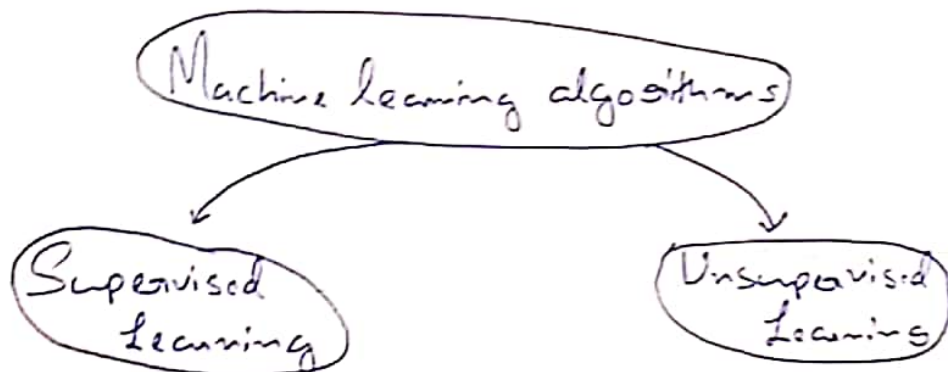# Introduction

# Introduction

**1.1) Welcome**

→ Machine learning grew out of work in AI.

**1.2) What is machine Learning**

→ Field of study that gives computers the ability to learn without being explicitly programmed.

<div align="right">Arthur Samuel (1959)</div>

Machine Learning algorithms

Supervised Learning ← → Unsupervised Learning

**1.3) Supervised learning**

→ ( We gave the algorithm the data set with right answers )

↓

( and the task of the algorithm was to predict the right answer for the next question )

⇒ **Regression** : Predict continuous value output.

⇒ **Classification** :: Discrete valued output (0 or 1)

(Support Vector machine)

→ Can have more than two values.

## 1.4) Unsupervised - Learning

↳

{ Algorithm that learns from test data that has not been labeled, classified or categorized }

↳ Clustering

{ Give the algorithm a ton of data ask to find structure in it. }

————✗————✗————

# CHAPTER 2

## Linear Regression With One Variable

# Linear regression with one variable

## 2.1) Model representation

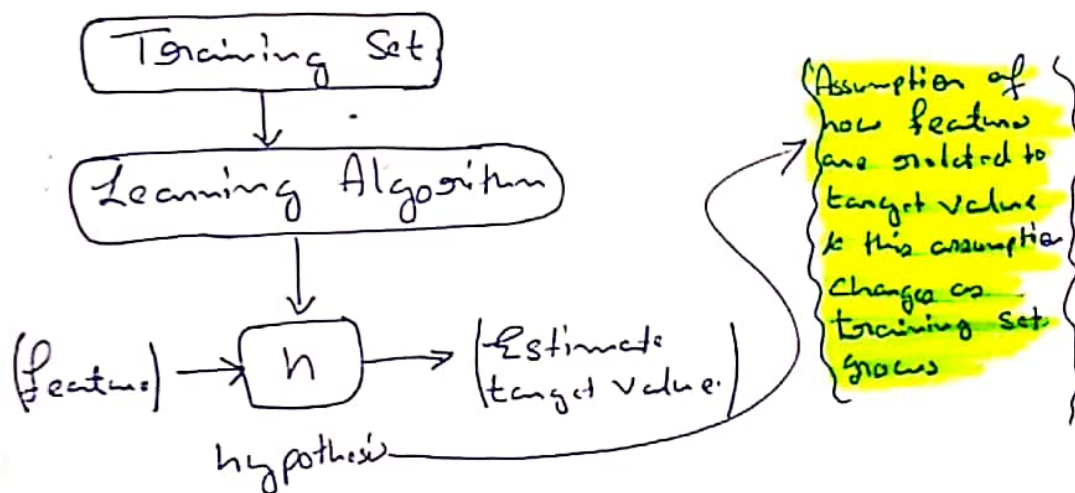⟹ Data set is called **training set.**

**Notation:**

$m$ ⟹ Number of training example

$x's$ ⟹ "input" variables / features

$y's$ ⟹ "output" variables / "target" variable
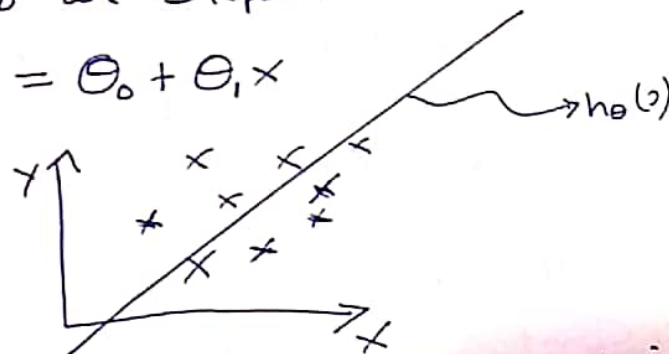
$(x, y)$ ⟹ To denote single training example.

$(x^{(i)}, y^{(i)})$ ⟹ To denote $i^{th}$ training example.

Training Set
↓
Learning Algorithm
↓
(feature) → $h$ → Estimate target value.
hypothesis

(Assumption of how features are related to target value & this assumption changes as training set grows)

⟹ $h$ is a function which maps $x's$ to $y's$

⟹ How do we represent $h$?

$$h_\theta(x) = \theta_0 + \theta_1 x$$

→ $h_\theta(x)$

⟹ Linear regression with one variable.
(i.e. Univariate linear regression)

## 2.2> Cost function

Hypothesis: $h_\theta(x) = \theta_0 + \theta_1 x$

$$\{ \theta_i's \Rightarrow Parameters \}$$

Objective: Choose $\theta_0, \theta_1$ so that $h_\theta(x)$ is close to $y$ for our training examples $(x, y)$.

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

↑

Cost function

$\Big\downarrow$

$\{$ Squared error cost function $\}$

Goal: minimize $J(\theta_0, \theta_1)$
$\theta_0\ \theta_1$

$\Rightarrow$ $\Big\{$ Minimize $J(\theta_0, \theta_1)$ to obtain value of $\theta_0$ & $\theta_1$ $\Big\}$

Contour line $\Rightarrow$ A contour line of a function of two variables is a curve along which the function has a constant value, so that the curve joins points of equal value.

$\hookrightarrow$ different colors are used to represent contour lines at different height.

# Gradient descent

→ {An algorithm for minimizing the cost function J}

↓

$$\min_{\theta_0 \cdots \theta_n} J(\theta_0, \cdots \theta_n)$$

⇒ First Start from a random value,

**# Gradient descent algorithm**

repeat until convergence {

$$\theta_j := \theta_j - \alpha \frac{\delta}{\delta \theta_j} J(\theta_0, \theta_1) \quad \left(\forall \; j=0 \; \text{and} \; j=1\right)$$

}

⇓ {Correct simultaneous update}

$$temp0 := \theta_0 - \alpha \frac{\delta}{\delta \theta_0} J(\theta_0, \theta_1)$$

$$temp1 := \theta_1 - \alpha \frac{\delta}{\delta \theta_1} J(\theta_0, \theta_1)$$

$$\theta_0 := temp0$$

$$\theta_1 := temp1$$

→ {Assignment operator}

$\alpha$ ⇒ learning rate

(i.e. how big the step is)

$= \rightarrow$ {Truth assertion}

**Issue:** It can be susceptible to local optima.

⟹ The algorithm is also called (**Batch**) Gradient Descent.

{ Each step of gradient descent uses all the training examples }

(Two extension)

(Exact method)    (Larger number of features)

{ $x_1$ $x_2$ ---- $x_n$ can be different feature }

# CHAPTER 3

# Linear Regression With Multiple Variable

# Linear regression with multiple Variables

## 4.1) Multiple features

### Notation:

$M$ = Number of features

$x^{(i)}$ = input (features) of $i^{th}$ training example

$\left( x_j^{(i)} \right.$ = Value of feature $j$ in $i^{th}$ training example

$\longrightarrow$ (will be a number)

{ Will be $n$ dimentional Vector }

$x^{(i)} \in \mathbb{R}^n$

$\Rightarrow$ New hypothesis function : $h_\theta(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \cdots \theta_n x_n$

$\Rightarrow$ For convenience of notation define $x_0^{(i)} = 1$

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

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

$$h_\theta(x) = \theta^T x$$ $\left\{ x = \begin{bmatrix} x_0 \\ x^{(i)} \end{bmatrix} \right\}$

$\Rightarrow$ Multivariate linear regression

## 4.2> Gradient descent for multiple variables

### Cost function

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

Vector        Vector

### Gradient descent

Repeat {

$$\theta_j := \theta_j - \alpha \frac{\delta}{\delta \theta_j} J(\theta) \quad \left\{ \forall \; j = 0, 1 \cdots m \right\}$$

}

$\Downarrow$

Repeat {

$$\theta_j := \theta_j - \alpha \frac{1}{m} \sum_{i=1}^{m} \left( h_\theta(x^{(i)}) - y^{(i)} \right) x_j^{(i)}$$

$$\left( \begin{array}{l} \text{Simultaneously update } \theta_j \\ \text{for } j = 0 \cdots m \end{array} \right)$$

}

# 15) Feature Scaling

If different features takes different range of values then it becomes very difficult for algorithm to find minimum or it takes way more time for algorithm to find minimum.

     ↳ To avoid this we scale every feature so that they are in similar range.

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

## Mean Normalization

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

General rule $\quad x_i \longleftrightarrow \dfrac{x_i - \mu_i}{S_i}$ → avg value of $x_i$ in training set

               ↳ range of $x_i$ (i.e. $\max x_i - \min x_i$)

⟹ To make sure Gradient descent is working correctly.

     ↳ Plot $J(\theta)$ vs (Number of iteration)
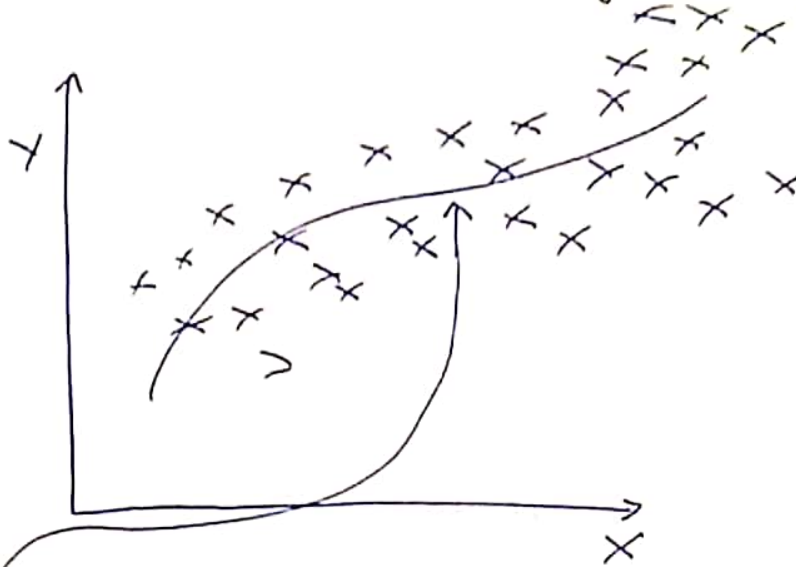
         ↓

{ If it is continuously decreasing than its going well. }

         ↓

{ If it is almost flat than you can assume that it has converged }

$\Rightarrow$ To choose $\alpha$, try

.... 0.001 , 0.01 , 0.1 , 1 ...

## 4.5> Polynomial regression



$h_\theta(x) = \theta_0 + \theta_1 x$ $\{$ doesn't give good fit $\}$

So, $h_\theta(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3$

$\{$ Very good fit $\}$

So $x_1 = x$
$x_2 = x^2$ $\{$ define three new features $\}$
$x_3 = x^3$

## 4.6) Normal Equations

Method to Solve for $\theta$ analytically.

Let $X = \begin{bmatrix} \vdots & x'_1 & x'_2 \cdots x'_n \\ \vdots & x^2_1 & x^2_2 - x^2_n \\ \vdots & \vdots & \vdots \\ \vdots & x^m_1 & x^m_2 \cdots x^m_n \end{bmatrix}_{m \times (n+1)}$

→ number of feature

{ number of training Set }

{ Matrix that Contain all the training example }

Similarly $Y = \begin{bmatrix} y' \\ y^2 \\ \vdots \\ y^m \end{bmatrix}_{m \times 1}$

{ Collection of all the target Variable in training example }

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

→ This gives you the value of $\theta$ that minimizes the Cost function

{ 
## Note

$$J(\theta) = \frac{1}{2m} (X\theta - Y)^T (X\theta - Y)$$
}

| Gradient Descent | Normal Equation |
|---|---|
| ⇒ Need to choose $\alpha$. | ⇒ No need to chose $\alpha$. |
| ⇒ Need many iterations. | ⇒ Don't need to iterate. |
| ⇒ Works well even when $n$ is large | ⇒ Slow if $n$ is very large. |
| | $O(n^3)$ |
| $n > 10000$ | $n = 100$    OK |
| | $n = 1000$    OK |
| | $n = 10000$ ← Not greater than this |

# CHAPTER 4

# Logistic Regression

# Logistic Regression

**6.1> Classification**

→ Logistic Regression is a classification algorithm.

$$0 \le h_\theta(x) < 1$$

**6.2> hypothesis - representation**

$$h_\theta(x) = g(\theta^T x)$$

where $g(z) = \dfrac{1}{1+e^{-z}}$

→ Sigmoid function
or
logistic function

$$\Rightarrow \quad h_\theta(x) = \dfrac{1}{1+e^{\theta^T x}}$$

$\left. \begin{array}{l} \text{estimated probability that} \\ y=1 \text{ on input } x \end{array} \right\}$

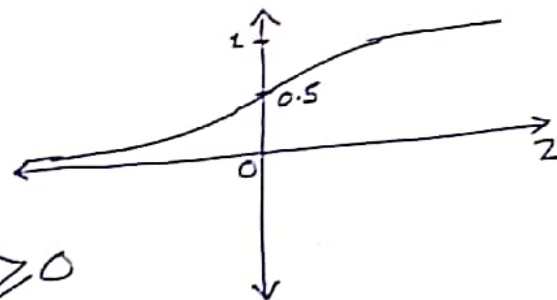**6.3> Decision boundary**

Predict "$y=1$" if $h_\theta(x) \ge 0.5$

& Predict "$y=0$" if $h_\theta(x) < 0.5$
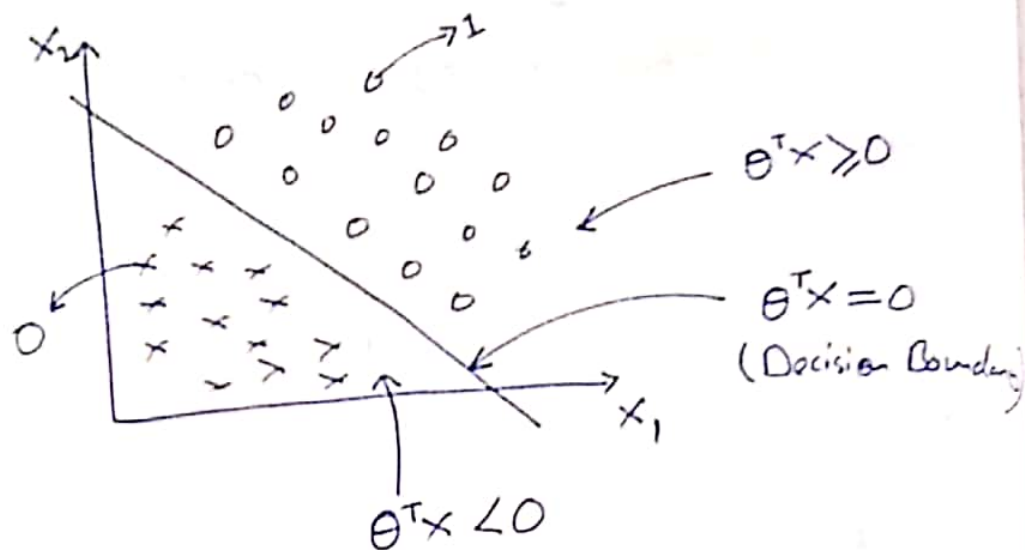
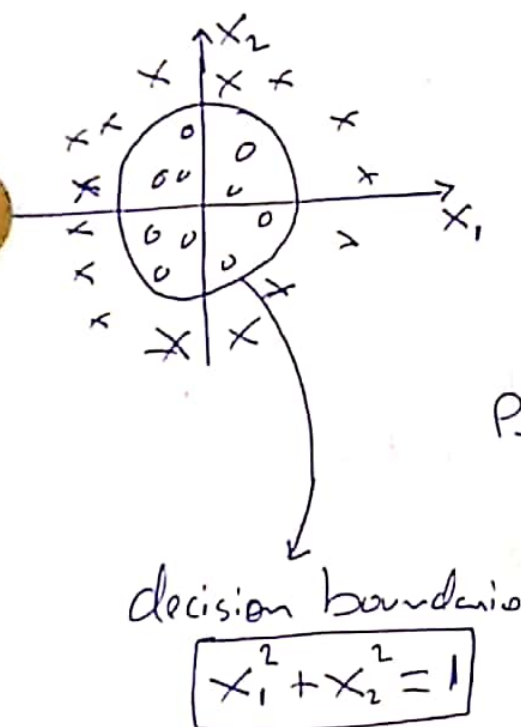$$h_\theta(z) = \dfrac{1}{1-e^{-z}}$$

when, $z = \theta^T x$

$h_\theta(z) \ge 0.5 \Rightarrow z \ge 0$

$h_\theta(z) < 0.5 \Rightarrow z < 0$

$\theta^T x \geqslant 0$

$\theta^T x = 0$
(Decision Boundary)

$\theta^T x < 0$

# Non linear decision boundaries



$h_\theta(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_2^2)$

$\theta = \begin{bmatrix} -1 \\ 0 \\ 0 \\ 1 \\ 1 \end{bmatrix}$

Predict "$y=1$" if $-1 + x_1^2 + x_2^2 \geqslant 0$

$\Rightarrow x_1^2 + x_2^2 \geqslant 0$

decision boundary

$$\boxed{x_1^2 + x_2^2 = 1}$$

6.4) Cost function

Training set : $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots (x^{(m)}, y^{(m)})\}$

$\{m \text{ examples}\}$

$x = \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix}$  $x_0 = 1$  $y \in \{0, 1\}$

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

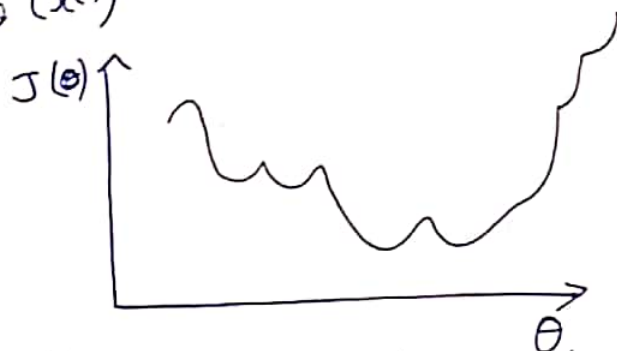* ## Cost function

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

6)

⇒ If $Cost(h_\theta(x^{(i)}), y) = \frac{1}{2}(h_\theta(x^{(i)}) - y)^2$

{ Same. as Linear regression }

⇒ then, $J(\theta)$ is not a Convex function because of highly non linear term $h_\theta(x^{(i)})$

$J(\theta)$ ↑



θ

⇒ So using gradient descent will not guarantee global minima.

⇒ So to avoide this problem we will use the following as cost function.

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

( → It gives very high penalty if prediction is wrong and zero penalty if prediction is write.

→ It guarantee a convex function.
( Proving this is out of scope of this course) )

6.5) **Simplified Cost function and gradient descent**

⇒ Simplified Cost function:

$$\text{Cost}(h_\theta(x), y) = -y(\log(h_\theta(x)))$$
$$- (1-y)\log(1 - h_\theta(x))$$

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

\* **Gradient Descent**

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

6

Want $\min_\theta J(\theta)$:

Repeat {

$$\theta_j : \theta_j - \alpha \frac{\delta}{\delta \theta_j} J(\theta)$$

}

{ Simultaneously update all $\theta_j$ }

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

## 6.6) Advanced Optimization

→ Other optimization algorithms: { other than gradient descent }

- Conjugate gradient
- BFGS
- L-BFGS

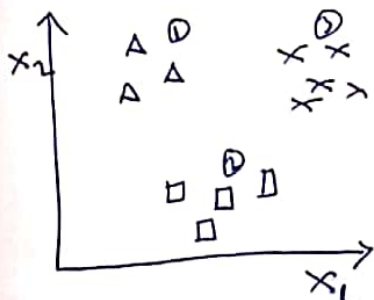**Advantages**

→ No need to manually pick $\alpha$.

→ Often faster than gradient descent

**Disadvantages**

→ More complex

## 6.7) Multiclass - classification (one-vs-all)

⇒ Classification problem with more than one classes.



⇒ Turn this into three seperate binary classification problem:

$h_\theta^{(1)}(x) \rightarrow \triangle$ Vs rest

$h_\theta^{(2)}(x) \rightarrow \square$ Vs rest

$h_\theta^{(3)}(x) \rightarrow \times$ Vs rest

① Train a logistic regression classifier $h_\theta^{(i)}(x)$ for each class $i$ to predict the probability that $y = i$.

② On a new input $x$, to make a prediction, pick the class $i$ that maximizes
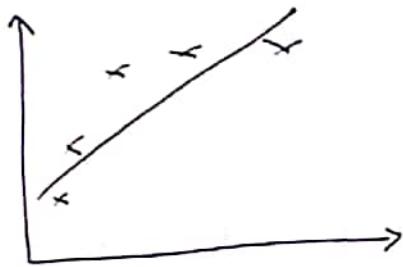
$$\max_i h_\theta^{(i)}(x)$$
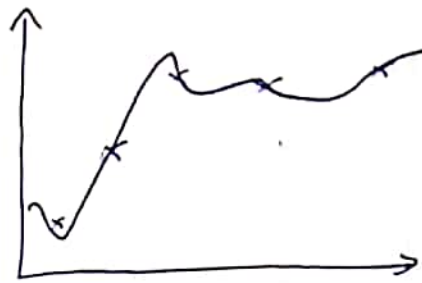
# CHAPTER 5

# Regularization

# Regularization

## 7.1) Problem of Overfitting



$\theta_0 + \theta_1 x$

$\Rightarrow$ "Under fit"

$\Rightarrow$ It has <u>high bias</u>

$\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$

$\Rightarrow$ "Overfit"

$\Rightarrow$ It has <u>high variance</u>

## Overfitting

$\rightarrow$ If we have too many features, the learned hypothesis may fit the training set very well, but fails to generalize to example.

$\Rightarrow$ Overfitting can also be found in Classification.

* Addressing Overfitting

1. Reduce number of features
   → Manually select which features to keep.
   → Model Selection algorithm

2. Regularization
   → Keep all the features, but reduce magnitude of parameters $\theta_j$.

7.2) Cost function

⟹ Small values for parameters $\theta_0, \theta_1, \theta_2 \cdots \theta_n$
   → "Simpler" hypothesis.
   → Less prone to overfitting.

⟹ As we don't know which parameter is less important, so we will penalise all the parameters:-

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

regularization Parameter

{Objective :: fit the training data well}

{Objective : Keep the parameter small}

{It controls the trade off between the two objectives}

## 7.3) Regularized - Normal equation

$$X = \begin{bmatrix} (x^{(1)})^T \\ \vdots \\ (x^{(m)})^T \end{bmatrix} \qquad y = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(m)} \end{bmatrix}$$

$$\underset{m \times (n+1)}{} \qquad \qquad \underset{m \times 1}{}$$

$$\Theta = \left( X^T X + \lambda \begin{bmatrix} 0 & & & \\ & 1 & & \\ & & 1 & \\ & & & \ddots & \\ & & & & 1 \end{bmatrix} \right)^{-1} X^T y$$

$$\underset{(n+1) \times (n+1)}{}$$

## 7.4) Regularized - logistic regression

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

$$+ \frac{\lambda}{2m} \sum_{j=1}^{n} \theta_j^2$$

———✕——— ⇒ ✕———

Scanned by CamScanner