

Threshold classifier output $h_0(x) \approx 0.5$
 If $h_0(x) \geq 0.5$, predict "y=1"
 If $h_0(x) < 0.5$, predict "y=0"

Logistic Regression

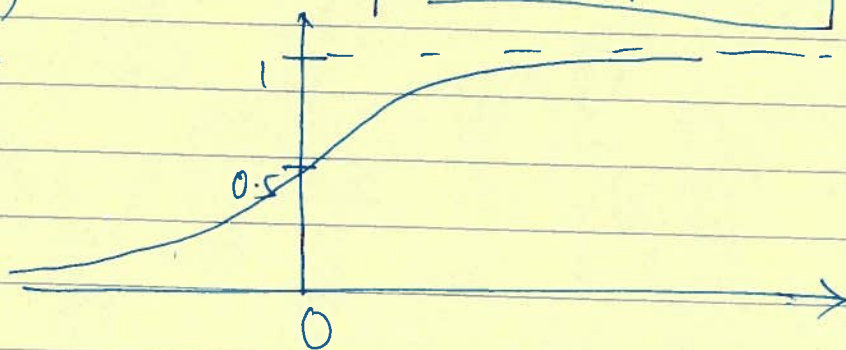
Hypothesis representation

want $0 \leq h_0(x) \leq 1$

$$h_0(x) = g(\theta^T x)$$

Sigmoid or Logistic function $g(z) = \frac{1}{1+e^{-z}}$

$$h_0(x) = \frac{1}{1+e^{-\theta^T x}}$$



Parameters θ

$h_0(x)$ = estimated probability that $y=1$ on input x
 If $x = \begin{bmatrix} x_0 \\ x_1 \end{bmatrix} = \begin{bmatrix} 1 \\ \text{tumor size} \end{bmatrix}$

$h_0(x) = 0.7$ 70% chance.

$h_0(x) = P(y=1 | x; \theta)$ $y=1$, give x parameterised by θ

$$P(y=0 | x; \theta) + P(y=1 | x; \theta) = 1$$

$$P(y=0 | x; \theta) = 1 - P(y=1 | x; \theta)$$

$$h_{\theta}(x) = g(\theta^T x) = P(y=1|x; \theta)$$

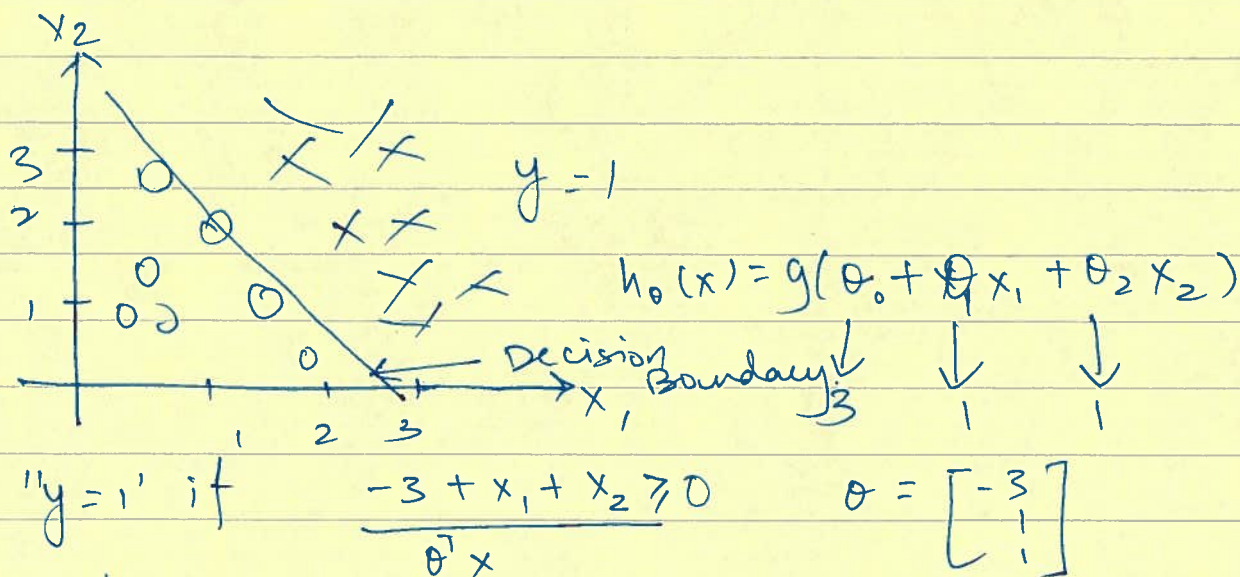
$$g(z) = \frac{1}{1 + e^{-z}}$$

$$\begin{aligned} y=1 & \text{ if } h_{\theta}(x) \geq 0.5 \\ y=0 & \text{ if } h_{\theta}(x) < 0.5 \end{aligned}$$

$$g(z) < 0.5$$

$$\begin{aligned} h_{\theta}(x) &= g(\theta^T x) \\ \theta^T x &< 0 \end{aligned}$$

$$\left. \begin{aligned} g(z) &\geq 0.5 \\ \text{when } z &\geq 0 \\ h_{\theta}(x) &= g(\theta^T x) \geq 0.5 \\ \text{wherever} \\ \theta^T x &\geq 0 \end{aligned} \right\}$$



$$\begin{aligned} &\rightarrow x_1 + x_2 \geq 3 \\ &\quad \quad \quad y=1 \end{aligned}$$

$$\begin{aligned} &\rightarrow x_1 + x_2 < 3 \\ &\quad \quad \quad y=0 \end{aligned}$$

$$\boxed{x_1 + x_2 = 3} \rightarrow h_{\theta}(x) = 0.5$$

$$\theta = \begin{bmatrix} 5 \\ -1 \\ 0 \end{bmatrix}$$

$$h_{\theta}(x) = g(5 - x_1)$$

$$\begin{aligned} y=1 & \text{ if } 5 + (-1)x_1 + 0(x_2) \geq 0 \\ & \quad \quad \quad 5 - x_1 \geq 0 \\ & \quad \quad \quad -x_1 \geq -5 \quad x \leq 5 \end{aligned}$$

$$\begin{aligned} 5 - x_1 &\geq 0 \\ x_1 &\leq 5 \end{aligned}$$

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

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

$$y=1 \text{ if } -1 + x_1^2 + x_2^2 \geq 0$$

$$x_1^2 + x_2^2 \geq 1 \quad \text{Decision Boundary.}$$

Unit circle.

Inside circle $y=0$

Outside circle $y=1$

Decision Boundary \rightarrow is a property of hypothesis and parameters not the training set.

Cost function

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

m examples $x \in \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix} \in \mathbb{R}^{n+1}$
 $x_0 = 1, y \in \{0, 1\}$

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

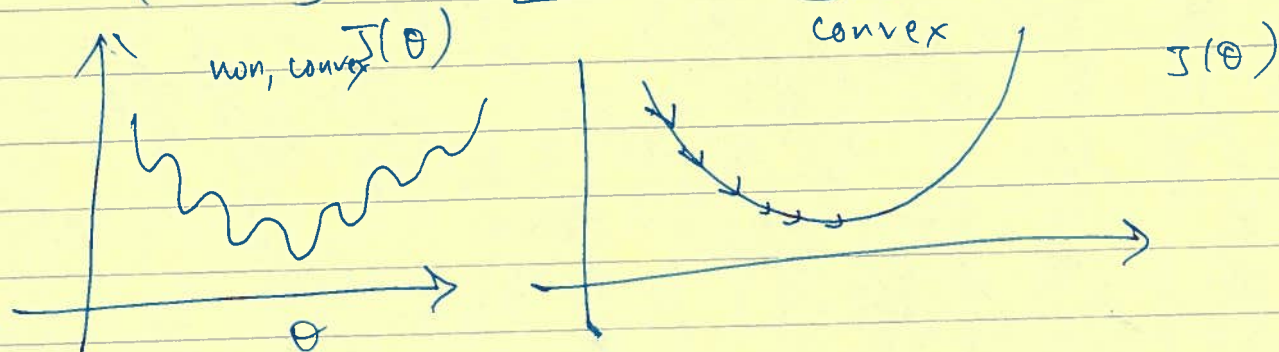
linear regression
↓
logistic.

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

cost($h_{\theta}(x^{(i)}, y)$)

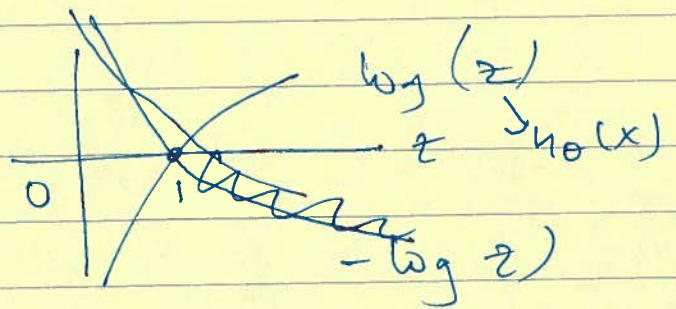
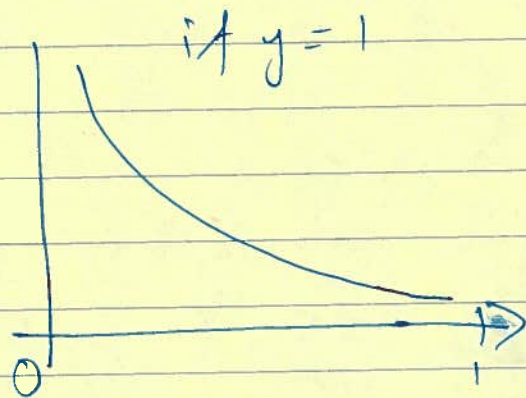
$$\text{cost}(h_{\theta}(x^{(i)}, y^{(i)})) = \frac{1}{2} (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

$$\text{cost}(h_{\theta}(x, y)) = \frac{1}{2} (h_{\theta}(x, y))^2$$



Logistic Regression Cost Function

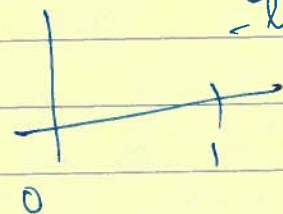
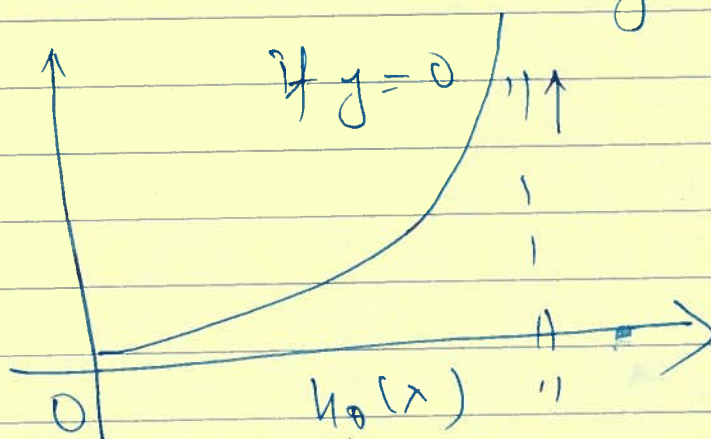
$$\text{Cost}(h_0(x), y) = \begin{cases} -\log(h_0(x)) & \text{if } y = 1 \\ -\log(1 - h_0(x)) & \text{if } y = 0 \end{cases}$$



Cost = 0 if $y = 1, h_0(x) = 1$
 But as $h_0(x) \rightarrow 0$
 Cost $\rightarrow \infty$

$h_0(x) = 0$
 Predicts

$y = 1 \mid x; 0$ by $y = 1$
 $-\log(1 - z)$



Simplified cost function and gradient descent

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

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

Note: $y=0$ or 1 always.

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

$$\text{if } y=1: \text{cost}(x) = -\log h_{\theta}(x)$$

$$\text{if } y=0: \text{cost}(x) = -\log(1-h_{\theta}(x))$$

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

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

To fit parameters θ :

$$\min_{\theta} J(\theta)$$

Output $h_{\theta}(x) = \frac{1}{1 + e^{-\theta^T x}} \quad P(y=1|x; \theta)$

Want $\min_{\theta} J(\theta)$.

Repeat $\{$

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

for $i=0$ to n .

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix}$$

$$\frac{\partial}{\partial \theta_j} J(\theta) = \frac{1}{n} \sum_{i=1}^n (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

$$h_{\theta}(x) = \theta^T x$$

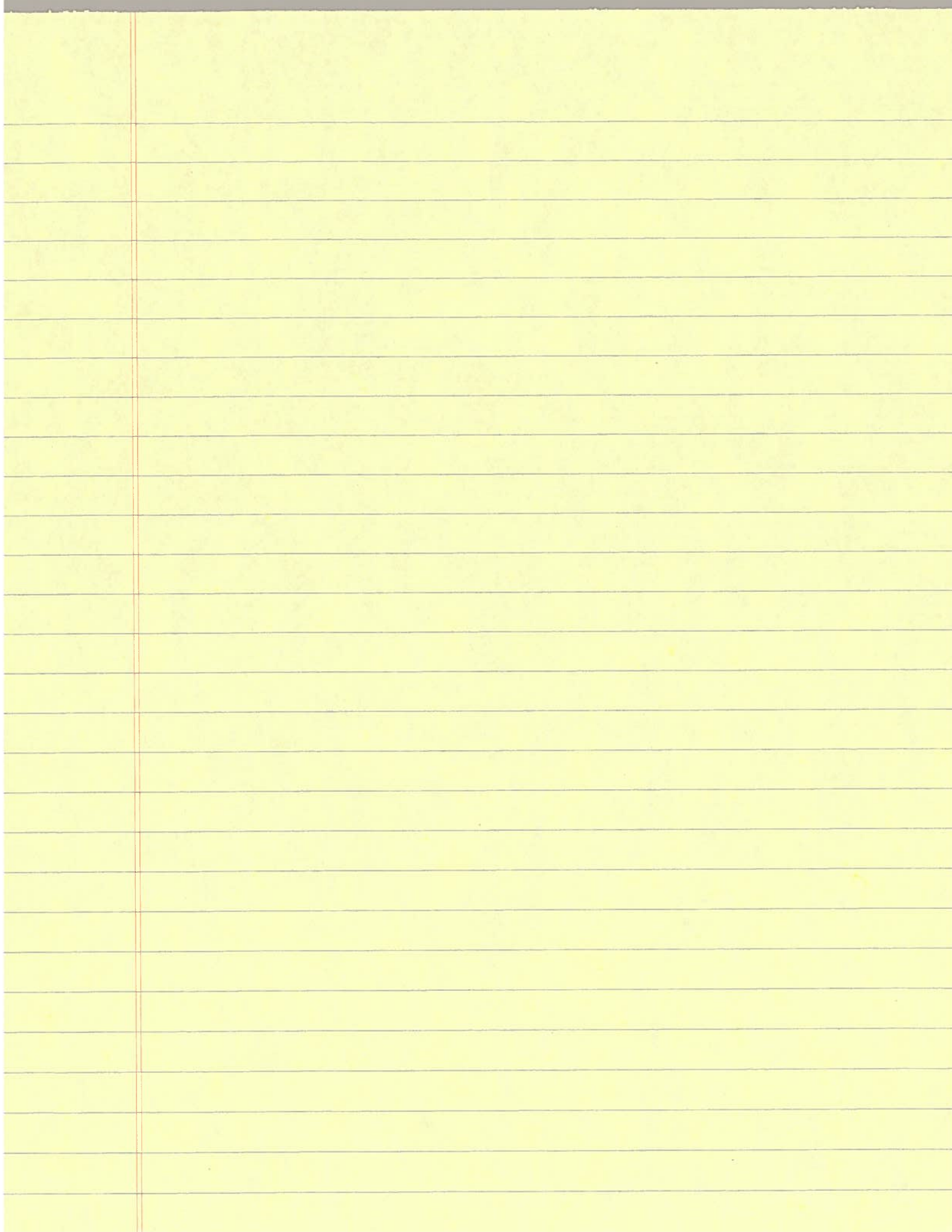
$$h_{\theta}(x) = \frac{1}{1 + e^{-\theta^T x}}$$

Vectorized implementation

$$\theta := \theta - \alpha \frac{1}{n} \sum_{i=1}^n [(h_{\theta}(x^{(i)}) - y^{(i)}) \cdot x^{(i)}]$$

vectorized

$$\theta := \theta - \frac{\alpha}{n} X^T (g(X\theta) - \bar{y})$$



Optimization Algorithm.

$$\frac{\partial J(\theta)}{\partial \theta_j} \quad (\text{for } j=0,1,\dots,n)$$

Gradient descent

{ Repeat

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

$$\theta = \begin{bmatrix} \theta_1 \\ \theta_2 \end{bmatrix}$$

$$J(\theta) = (\theta_1 - 5)^2 + (\theta_2 - 5)^2$$

$$\theta_1 = 5, \theta_2 = 5$$

$$\frac{\partial J(\theta)}{\partial \theta_1} = 2(\theta_1 - 5)$$

$$\frac{\partial J(\theta)}{\partial \theta_2} = 2(\theta_2 - 5) \quad 2$$

$$\begin{bmatrix} -6 \\ 0 \\ 1 \end{bmatrix}$$

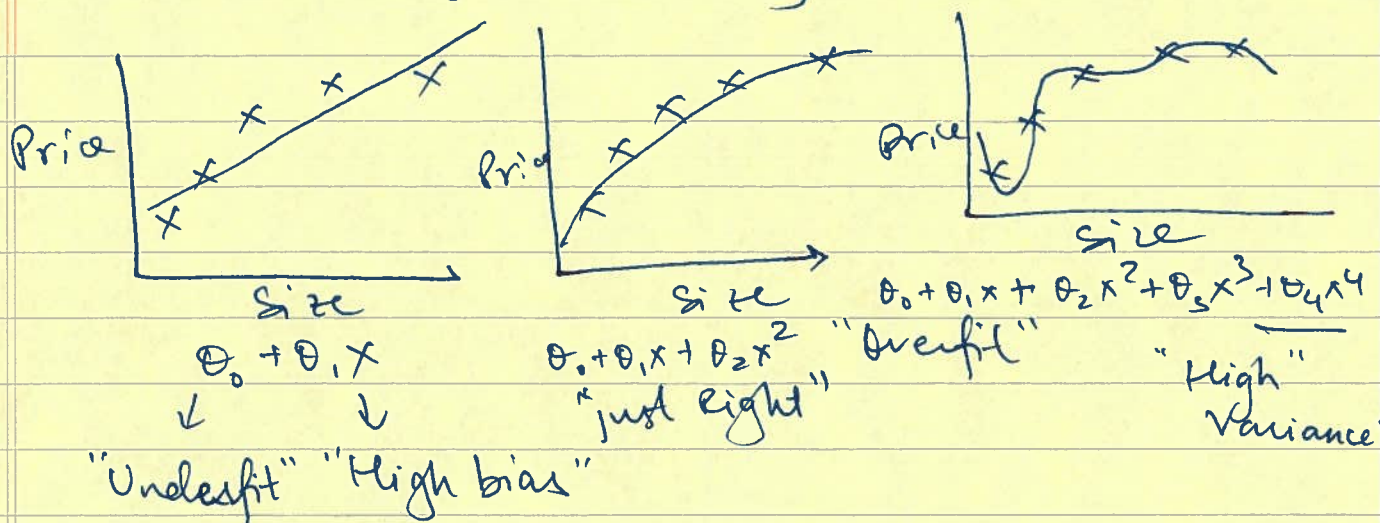
$$-6 + x_2 \geq 0$$

~~$$x_2 \leq 6$$~~

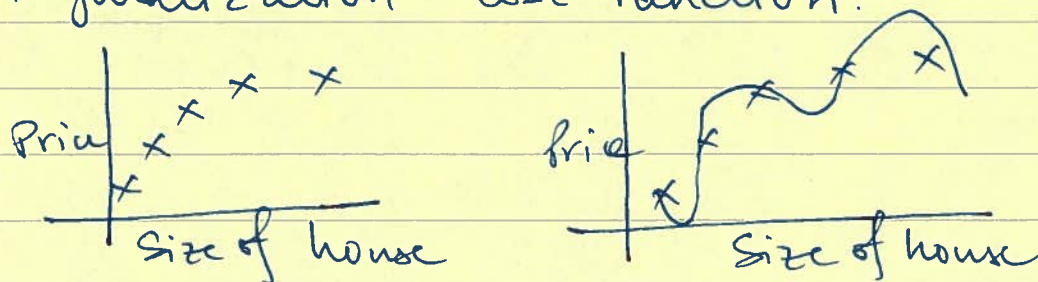
$$x_2 \geq 6.$$

Binary Classification

The Problem of Overfitting



Regularization Cost Function:



θ_3, θ_4 really small

$$\min_{\theta} \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 + 1000\theta_3^2 + 1000\theta_4^2$$

$\theta_3 \geq 0 \quad \theta_4 \geq 0$

$\theta_3, \theta_4 \geq 0$

Housing

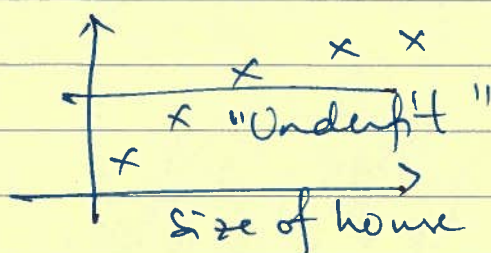
→ Features x_1, x_2, \dots, x_{100}

Parameter $\theta_0, \theta_1, \theta_2, \dots, \theta_{100}$

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

$\theta_1, \theta_2, \theta_3, \dots, \theta_n$

$\min_{\theta} J(\theta)$



Regularization
Parameter
fit Training
set well

$$\lambda = 10^{10}$$

$$\theta_1 \approx 0, \theta_2 \approx 0, \theta_3 \approx 0$$

$$h_{\theta}(x) = \theta_0$$

$$h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

Regularized Linear Regression.

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

Repeat { $\frac{\partial}{\partial \theta} J(\theta)$

$$\theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_0^{(i)}$$

$$\theta_j := \theta_j - \alpha \left[\frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)} + \frac{\lambda}{m} \theta_j \right] \quad (j = 0, 1, 2, 3, \dots, n)$$

$\frac{\partial}{\partial \theta_j} J(\theta) \rightarrow$ regularization

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

$$\frac{\partial}{\partial \theta_j} J(\theta) \stackrel{\text{set}}{=} 0$$

$$X = \begin{bmatrix} x^{(1)T} \\ \vdots \\ (x^{(m)})^T \end{bmatrix} \quad y = \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(m)} \end{bmatrix} \quad \mathbb{R}^m$$

$m \times (n+1)$

$$\min_{\theta} J(\theta) = (X^T X + \lambda \underbrace{\begin{bmatrix} 0 & 0 & 0 & 0 \\ & 1 & 0 & 0 \\ & 0 & 1 & 0 \\ & 0 & 0 & 1 \end{bmatrix}}_{(n+1) \times (n+1)})^{-1} X^T y$$

Eg $n=2$

$$\begin{bmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

Suppose $m \leq n$ \rightarrow # examples \rightarrow # features

$$\theta = \underbrace{(X^T X)^{-1}}_{\hookrightarrow \text{non invertible/singular}} X^T y$$

if $\lambda > 0$

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