

## Week - 3

⇒ classification → variable  $y$  is discrete valued

↳ we'll use logistic regression for same.

example :- email : spam | not spam?

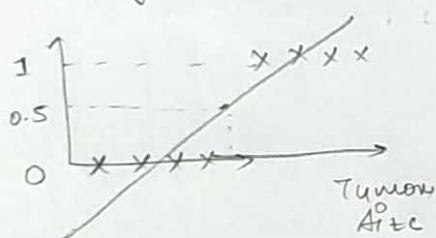
Online transactions: Fraud (Yes/No)?

Tumor malignant / Benign?

$y \in \{0, 1\}$  0: "Negative class" 1: "Positive class" Binary class.

↳ Right now, we'll deal with classification with  $y \in \{0, 1\}$  only

Training set



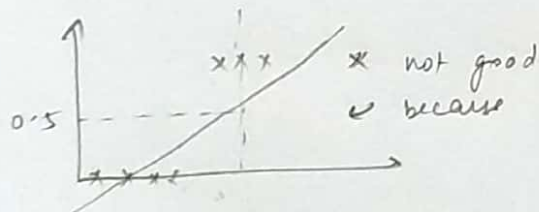
we can try to fit linear regression,

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

Threshold classifier output  $h_{\theta}(x)$  at 0.5:

If  $h_{\theta}(x) \geq 0.5$ , predict " $y=1$ "

If  $h_{\theta}(x) \leq 0.5$ , predict " $y=0$ "



classification  $y=0$  or  $1$

$h_{\theta}(x)$  can be  $> 1$  or  $< 0$

logistic regression,  $0 \leq h_{\theta}(x) \leq 1$   
classification

⇒ Hypothesis representation

logistic regression → want  $0 \leq h_{\theta}(x) \leq 1$

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

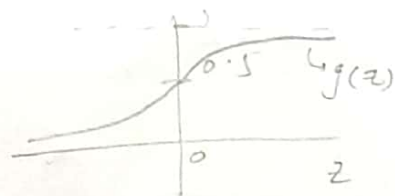
$$g(z) = \frac{1}{1 + e^{-z}}$$

} sigmoid function  
logistic function

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

$$\nearrow \frac{0.5}{}$$

Now, we need to fit parameters to our data.



Interpretations of Hypothesis output

$h_0(x)$  = estimated probability that  $y=1$  on a input  $x$ .

Example:- If  $x = \begin{bmatrix} x_0 \\ x_1 \end{bmatrix} = \begin{bmatrix} 1 \\ \text{tumor size} \end{bmatrix}$

$$h_0(x) = 0.7 \Rightarrow 0.7 \text{ chance that } y=1$$

So, it's 70% chance that tumor is ~~not~~ malignant.

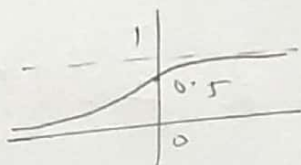
$h_0(x) = P(y=1 | x; \theta) \Rightarrow$  "Probability that  $y=1$ , given  $x$ , parametrised by  $\theta$ ".

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

$\Rightarrow$  Decision Boundary

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

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



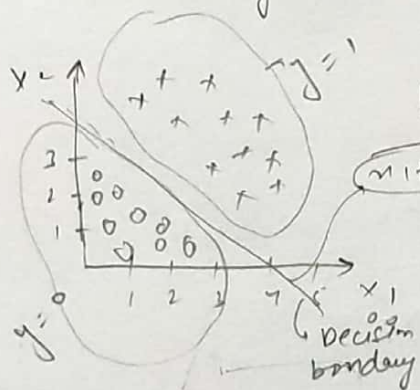
Suppose predict

" $y=1$ " if  $h_0(x) \geq 0.5$

" $y=0$ " if  $h_0(x) < 0.5$

$$\text{Now } g(z) \geq 0.5 \text{ for } z \geq 0 \Rightarrow \theta^T x \geq 0$$

Thus, " $y=1$ " when  $\theta^T x \geq 0$ , " $y=0$ ", when  $\theta^T x < 0$



$$h_0(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2)$$

$$\text{let us say } \theta_0 = -3 \quad \theta_1 = 1 \quad \theta_2 = 1$$

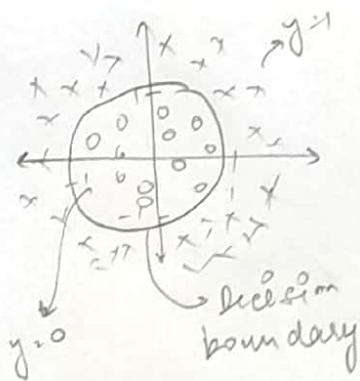
$$\theta = \begin{bmatrix} -3 \\ 1 \\ 1 \end{bmatrix}$$

$\rightarrow$  we will see how to find  $\theta$  later.

$$\text{" } y=1 \text{ " if } \underbrace{-3 + x_1 + x_2}_{\theta^T x} \geq 0 \Rightarrow x_1 + x_2 \geq 3$$

$\rightarrow$  equation of straight line

## Non-linear decision boundary



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

let's say after procedure, we have

$$\theta_0 = -1 \quad \theta_1 = 0 \quad \theta_2 = 0 \quad \theta_3 = 1 \quad \theta_4 = 1$$

$$\theta = \begin{bmatrix} -1 \\ 0 \\ 0 \\ 1 \\ 1 \end{bmatrix} \Rightarrow \text{Predict "y=1" if } -1 + x_1^2 + x_2^2 \geq 0$$

$$\Rightarrow \underbrace{x_1^2 + x_2^2 \geq 0}_{\text{equation of circle}}$$

$\Rightarrow$  Cost function

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

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

$$h_{\theta}(x) = \frac{1}{1 + e^{-\theta^T x}} \quad \text{How to choose } \theta?$$

$\rightarrow$  linear regression:  $J(\theta) = \frac{1}{m} \sum_{i=1}^m \underbrace{\frac{1}{2} (h_{\theta}(x^{(i)}) - y^{(i)})^2}_{\text{cost}(h_{\theta}(x^{(i)}), y^{(i)})}$

let's define

$$\left[ \text{cost}(h_{\theta}(x), y) = \frac{1}{2} (h_{\theta}(x) - y)^2 \right]$$

logistic regression  
cost function

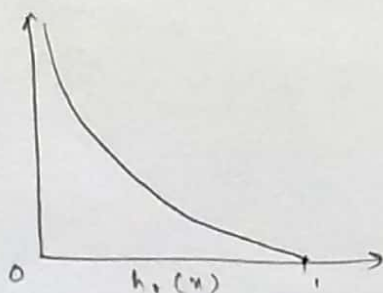
$$\text{cost}(h_{\theta}(x), y) = \begin{cases} -\log(h_{\theta}(x)) & \text{if } y = 1 \\ -\log(1 - h_{\theta}(x)) & \text{if } y = 0 \end{cases}$$

$$\text{if } y = 1$$

$$\text{cost} = 0 \quad \text{if } y = 1, h_{\theta}(x) = 1$$

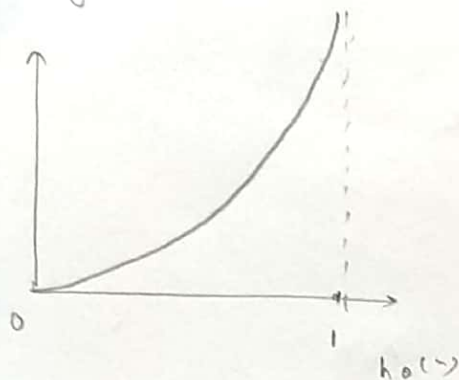
$$\text{But as } h_{\theta}(x) \rightarrow 0 \quad \text{cost} \rightarrow \infty$$

That means that if  $h_{\theta}(x) = 0$  and  $y = 1$ , we'll have a very large cost.





$$\text{If } y=0$$



→ cost function in this way guarantees that  $J(\theta)$  is convex for logistic regression.

⇒ Simplified cost function and Gradient Descent

$$J(\theta) = \frac{1}{n} \sum_{i=1}^n \text{cost}(h_0(x^{(i)}), y^{(i)})$$

$$\text{cost}(h_0(x), y) = \begin{cases} -\log(h_0(x)) & \text{if } y=1 \\ -\log(1-h_0(x)) & \text{if } y=0 \end{cases} \quad \left. \vphantom{\begin{cases} -\log(h_0(x)) \\ -\log(1-h_0(x)) \end{cases}} \right\} \begin{array}{l} \text{why? is} \\ \text{always} \\ \text{either 0 or 1} \end{array}$$

Now, let's make a combined cost fn.

$$\left[ \text{cost}(h_0(x), y) = -y \log(h_0(x)) - (1-y) \log(1-h_0(x)) \right]$$

$$\text{Thus, if } y=1 : \text{cost}(h_0(x), y) = -\log(h_0(x))$$

$$y=0 : \text{cost}(h_0(x), y) = -\log(1-h_0(x))$$

→ So, this cost fn. works

$$\text{Thus } J(\theta) = -\frac{1}{n} \sum_{i=1}^n \left[ y^{(i)} \log h_0(x^{(i)}) + (1-y^{(i)}) \log(1-h_0(x^{(i)})) \right]$$

To fit  $\theta$  such that  $J(\theta)$  is minimised.

To make a prediction given new  $x$ :

$$\text{output } h_0(x) = \frac{1}{1 + e^{-\theta^T x}} \quad \left. \vphantom{\frac{1}{1 + e^{-\theta^T x}}} \right\} \rightarrow P(y=1 | x; \theta)$$

Gradient descent

$$\text{Repeat } \left\{ \begin{array}{l} \theta_j := \theta_j - \alpha \frac{d}{d\theta_j} J(\theta) \end{array} \right.$$

$$h_0(x) = \frac{1}{1 + e^{-(\theta_0 x_0 + \theta_1 x_1)}}$$

$$J(\theta) = -\frac{1}{n} \sum_{i=1}^n y^{(i)} \log(h_0(x^{(i)})) + (1 - y^{(i)}) \log(1 - h_0(x^{(i)}))$$

$$= -\frac{1}{n} \sum_{i=1}^n y^{(i)} \log\left(\frac{1}{1 + e^{-(\theta_0 x_0 + \theta_1 x_1)}}\right) + (1 - y^{(i)}) \log\left(1 - \frac{1}{1 + e^{-(\theta_0 x_0 + \theta_1 x_1)}}\right)$$

let's take 1 term.

$$\frac{dJ(\theta)}{d\theta} = y \cdot \frac{e^{-(\theta_0 x_0 + \theta_1 x_1)}}{(1 + e^{-(\theta_0 x_0 + \theta_1 x_1)})^2} \cdot x_1$$

$$= y \frac{e^{-\theta^T x}}{1 + e^{-\theta^T x}} x_1 + (1 - y) \cdot \frac{1}{(1 + e^{-\theta^T x})^2} \cdot (-x_1)$$

$$= y \frac{e^{-\theta^T x}}{1 + e^{-\theta^T x}} x_1 + (1 - y) \frac{e^{-\theta^T x}}{(1 + e^{-\theta^T x})^2} x_1$$

$$= y x_1 \left( \frac{1}{1 + e^{-\theta^T x}} \right) - \frac{x_1}{1 + e^{-\theta^T x}}$$

$$= (y - h_0(x)) x_1$$

$$b - x_2 = 0$$

$$x_2 = b$$

so QD:

$$\theta_j := \theta_j^0 - \alpha \frac{dJ(\theta)}{d\theta_j^0}$$

$$\frac{dJ(\theta)}{d\theta_j^0} = \frac{1}{m} \sum_{i=1}^m (h_0(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

$$\theta_j^0 := \theta_j^0 - \alpha \sum_{i=1}^m (h_0(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

↳ This is exactly the same as linear regression, just  $h_0(x)$ 's definition is changed:  $\Rightarrow$  linear  $\rightarrow h_0 x = \theta^T x$   
logistic  $\rightarrow h_0 x = \frac{1}{1 + e^{-\theta^T x}}$

vectorised implementation

$$h = g(X\theta)$$

$$y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{bmatrix}_{m \times 1}$$

$$y^T = (1 \times m)$$

$$y^T = [y_1 \ y_2 \ \dots \ y_m]$$

$$X = \begin{bmatrix} x_0^1 & x_1^1 & x_2^1 & \dots & x_n^1 \\ x_0^2 & x_1^2 & x_2^2 & \dots & x_n^2 \\ \vdots & \vdots & \vdots & \ddots & \vdots \end{bmatrix}_{m \times n} \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_n \end{bmatrix}_{n \times 1} = \begin{bmatrix} x_0^1 \theta_0 + x_1^1 \theta_1 + \dots \\ x_0^2 \theta_0 + x_1^2 \theta_1 + \dots \\ \vdots \end{bmatrix}_{m \times 1}$$

$\underbrace{\quad}_{m \times 1} \underbrace{\quad}_{m \times 1} \underbrace{\quad}_{n \times 1}$

$$\log(h) = \begin{bmatrix} h_1 \\ h_2 \\ \vdots \\ h_m \end{bmatrix}_{m \times 1}$$

$$J(\theta) = \frac{1}{m} \left( -y_1 h_1 - y_2 h_2 - \dots - y_m h_m + (1-y_1)(1-h_1) + \dots + (1-y_m)(1-h_m) \right)$$

$$= \sum y h + (1-y)(h)$$

$$J(\theta) = \frac{1}{m} \left( -y^T \log(h) - (1-y)^T \log(1-h) \right)$$

$$\theta = \theta - \frac{\alpha}{m} X^T (g(X\theta) - \vec{y})$$

$$X^T \Rightarrow m \times n$$

$$\begin{bmatrix} x_0^1 & x_1^1 & x_2^1 & \dots & x_n^1 \\ x_0^2 & x_1^2 & x_2^2 & \dots & x_n^2 \\ \vdots & \vdots & \vdots & \ddots & \vdots \end{bmatrix}_{m \times n}$$

$$\begin{bmatrix} h_1 - y_1 \\ h_2 - y_2 \\ \vdots \\ h_m - y_m \end{bmatrix}_{m \times 1}$$

$$\begin{bmatrix} h_1 - y_1 \\ h_2 - y_2 \\ \vdots \\ h_m - y_m \end{bmatrix}_{m \times 1} = \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix}_{m \times 1}$$



⇒ Advanced optimisation

optimisation Algorithm : →  $\begin{bmatrix} \text{Gradient descent} \\ \text{conjugate gradient} \\ \text{BFGS, L-BFGS} \end{bmatrix}$

So, Given  $\theta$ , we have code to compute

$$J(\theta)$$

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

Now, we can supply these to various optimisation algorithms

Advantages : - No need to give  $\alpha$ , faster than gradient descent

Disadvantages :- More complex

$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix}$  → function  $[Jval, gradient] = \text{costFunction}(\theta)$

$Jval = [\text{code to compute } J(\theta)];$

$gradient[1] = [\text{code to compute } \frac{dJ(\theta)}{d\theta_0}]$

$gradient[2] = [\text{code for } \frac{dJ(\theta)}{d\theta_1}]$

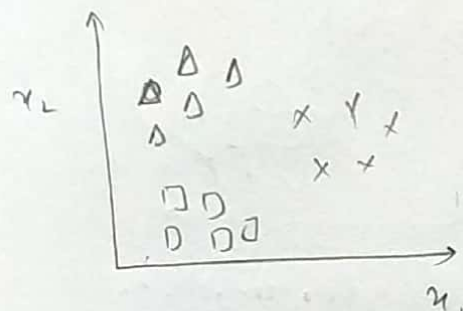
$\vdots$

⇒ Multiclass classification : One vs All.

ex. Email foldering : works, friends, family, Hobby  
 $y=1$        $y=2$        $y=3$        $y=4$

Medical diagnosis: Not ill, cold, Flu

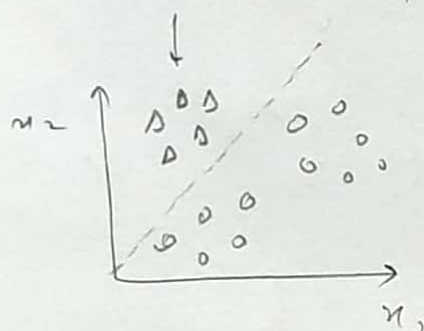
$y=1$        $y=2$        $y=2$



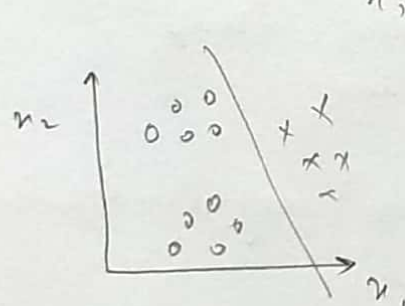
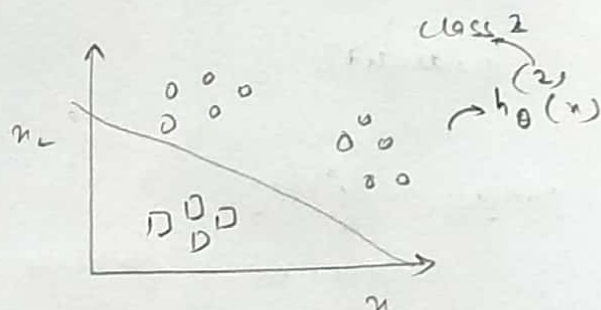
⊙ → class 1: Δ  
 class 2: □  
 class 3: X

One vs All

Now, we turn this into 3 binary classification problems



class 1  
 $\rightarrow h_{\theta}^{(1)}(x)$

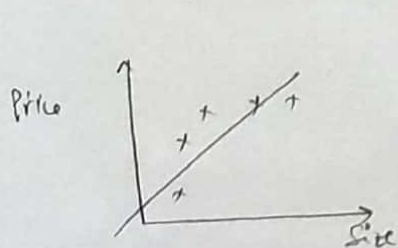


$h_{\theta}^{(2)}(x)$ , so, we fit 3 classifiers

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

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

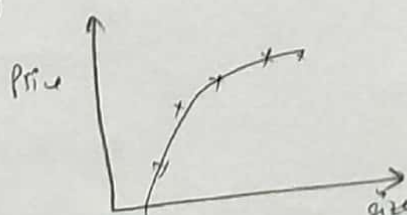
⇒ The problem of overfitting



$$\theta_0 + \theta_1 x$$

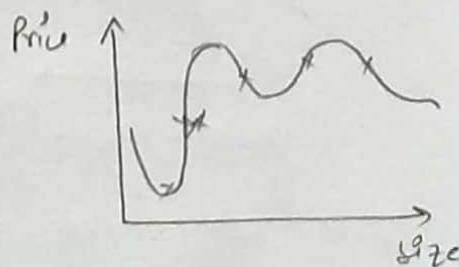
'Underfit', 'High bias'

Housing Price



$$\theta_0 + \theta_1 x + \theta_2 x^2$$

'just right'

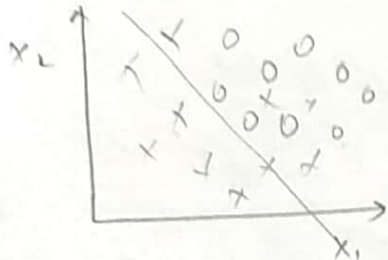


'overfit', 'High variance'



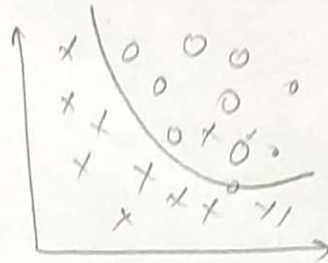
Overfitting: when we use too many features, the curve fits very well to training set,  $J(0) \approx 0$  maybe, but now, we fail to generalise over a new input  $x$ .

for logistic regression



$$h_0(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2)$$

(underfit)



$$h_0(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_2^2 + \theta_5 x_1 x_2)$$

(right)



$$h_0(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots)$$

(overfit)

→ Address overfitting

1. Reduce no. of features

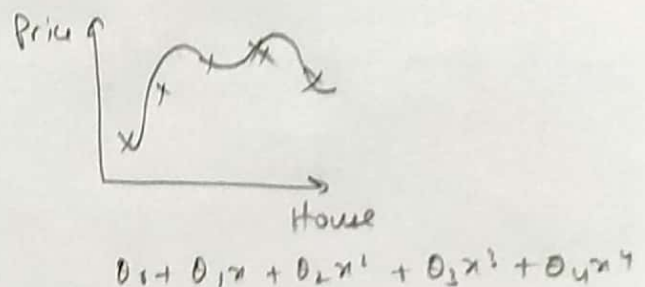
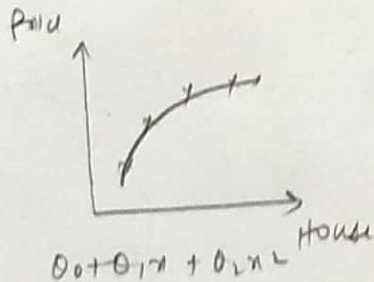
- ↳ (a) Manually select which features to keep.
- ↳ (b) Model selection algorithm (later)

2. Regularisation.

- ↳ (a) keep all the features, but reduce magnitude/values of parameters  $\theta_j$  (less assumptions)
- (b) works when we have lot of features, each contributes to predicting  $y$ .

⇒ cost function, regularisation

Intuition



Suppose, we make  $\theta_3$  and  $\theta_4$  very small, it would be as if we are getting rid of  $\theta_3$  and  $\theta_4$ .

## Regularisation

→ Small values for  $\theta_0, \theta_1, \dots, \theta_n$

- "simpler hypothesis"
- less prone to overfitting

## Housing

- Features:  $x_0, x_1, \dots, x_{100}$
- Parameters:  $\theta_0, \theta_1, \dots, \theta_{100}$

} Here, we do not know which features will give higher order terms or which features are irrelevant.

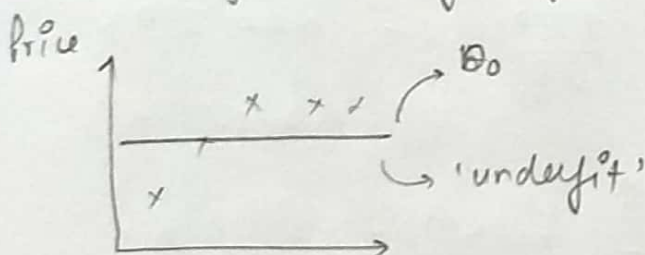
Now, we modify the cost function

$$J(\theta) = \frac{1}{2m} \left[ \sum_{i=1}^m (h_0(x^{(i)}) - y^{(i)})^2 + \underbrace{\lambda \sum_{j=1}^n \theta_j^2}_{\text{regularisation term}} \right]$$

So, in regularised linear regression, we choose  $\theta$  to minimise

$$J(\theta) = \frac{1}{2m} \left[ \sum_{i=1}^m (h_0(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^m \theta_j^2 \right]$$

What if  $\lambda$  is very large, say  $10^{10}$



$$\rightarrow \theta_1, \theta_2, \theta_3, \dots, \theta_n \approx 0$$

$$\text{So } h_0(x) \approx \theta_0$$

$$h_0(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 \dots$$

Thus, we also say that  $h_0(x)$  has high bias for  $\theta_0$ .

So,  $\lambda$  needs to be chosen wisely.

## ⇒ Regularised Linear Regression

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

↓  
Regularised cost function

We need to find  $\min_{\theta} J(\theta) \Rightarrow \theta$  corresponding to  $\min J(\theta)$ .

### Gradient Descent

Repeat-1

$$\theta_j^0 := \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

$j = (1, 2, 3, \dots, n)$

$$\theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_0^{(i)} \quad \left. \vphantom{\theta_0} \right\} \begin{array}{l} \text{writing} \\ \theta_0 \text{ case} \\ \text{separately} \end{array}$$

Now, we need to modify GD.

$$\theta_j^0 := \theta_j - \alpha \left[ \frac{1}{m} \sum (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)} + \frac{\lambda}{m} \theta_j \right]$$

↓  
because we  
don't penalise  
 $\theta_0$ .

( update rule  
now becomes.

$$\frac{dJ(\theta)}{d\theta_j} \quad \text{for new } J(\theta) \\ (j = 1, 2, 3, \dots, n)$$

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

→  $\left(1 - \frac{\alpha \lambda}{m}\right) < 1$ , this would be slightly less than 1,  
like 0.99.

$$\text{So, } \theta_j := \underbrace{0.99 \theta_j}_{\text{original update}} - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

↓  
So, basically, we are shrinking the  $\theta_j$  at every iteration so that eventually, it becomes smaller.



Normal equation

$$X = \begin{bmatrix} (x^{(1)})^T \\ (x^{(2)})^T \\ \vdots \\ (x^{(n)})^T \end{bmatrix}_{m \times n+1}$$

$$y = \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(m)} \end{bmatrix} \in \mathbb{R}^m$$

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

Now, for regularised -  $J(\theta)$

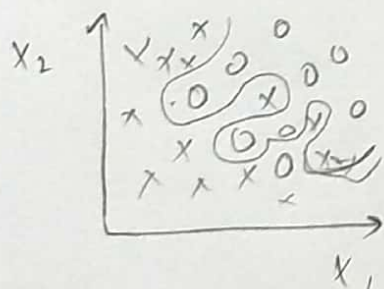
$$\theta = \left( X^T X + \lambda \begin{bmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \\ \vdots & \vdots & \vdots \\ 0 & 0 & 0 \end{bmatrix} \right)^{-1} X^T y$$

(n x 1) x (n+1)

Non invertibility issue

If  $\lambda > 0$ , then  $X^T X + \lambda \begin{bmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \\ \vdots & \vdots & \vdots \\ 0 & 0 & 0 \end{bmatrix}$  is invertible / non singular.

$\Rightarrow$  Regularised logistic Regression



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

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

So, even with a lot of features regularisation can help with overfitting problem.

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

$| \theta_0, \theta_1, \dots, \theta_n$

GD

Repeat k

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

Deal differently with  $\theta_0$ .

$$\theta_j^0 := \theta_j - \alpha \left[ \frac{1}{n} \sum_{i=1}^n (h_0(x^{(i)}) - y^{(i)}) x_j^{(i)} + \frac{d}{n} \theta_j^0 \right]$$

$(j = 1, 2, 3, \dots, n)$   
 $\theta_1, \theta_2, \dots, \theta_n$

Here

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

$$\frac{dJ(\theta)}{d\theta_j^0}$$

Advanced optimisation  $\rightarrow$  (see again)

function [jval, gradient] = costFunction(theta)

jval = [code to compute  $J(\theta)$ ]

gradient 1 = [code to compute  $\frac{dJ(\theta)}{d\theta_0}$ ]

gradient 2 = [code to compute  $\frac{dJ(\theta)}{d\theta_1}$ ]

gradient(n+1) = [code to compute  $\frac{d(J(\theta))}{d\theta_n}$ ]