

WEEK - 3

1.  $\{e^{\frac{1}{2}x}, e^{-\frac{1}{2}x}\}$

2.  $\{e^{2x}, e^{-2x}\}$

3.  $\{e^{3x}, e^{-3x}\}$

4.  $\{e^{(k+1)x}, e^{(k-1)x}\}$

WEEK - 3

## Classification Problem

(33)

Eg: Email: spam/ Not spam?

Tumor (Malignant/ Benign)

Online Transactions (Fraud/not)

$$y \in \{0, 1\}$$

Binary class

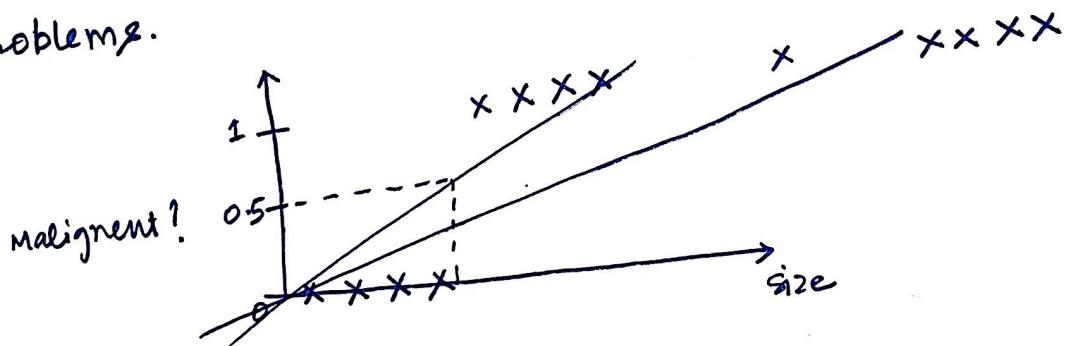
0: Negative (e.g. benign)

1: Positive (e.g. malignant)

$$y \in \{0, 1, 2, 3, \dots\}$$

Multi-class classification problem

- # Linear regression cannot be applied to classification problems.



- # Logistic Regression: classification problem.

$$0 \leq h_\theta(x) \leq 1$$

● Hypothesis : Logistic Regression Representation :

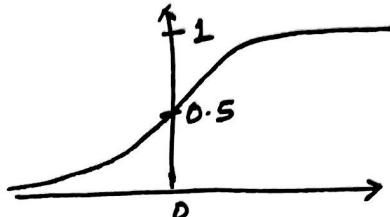
want  $0 \leq h_{\theta}(x) \leq 1$

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

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

sigmoid fn.

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



→ sigmoid fn.  
→ logistic fn.

→ Interpretation of Hypothesis Output :

$h_{\theta}(x)$  = estimated prob. that  $y=1$  on O/P.

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

$$h_{\theta}(x) = 0.7$$

Tell patient that 70% chance of tumor being malignant

$$\left. \begin{array}{l} \rightarrow h_{\theta}(x) = p(y=1 | x; \theta) \\ \rightarrow p(y=0 | x; \theta) = 1 - h_{\theta}(x) \end{array} \right\} \begin{array}{l} \text{"probability that } y=1, \\ \text{given } x, \text{ parameterized by } \theta \end{array}$$

## • Decision Boundary

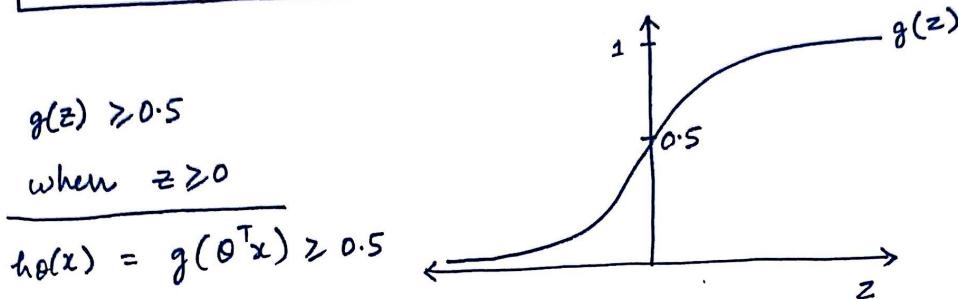
Predict

(35)

$$h_{\theta}(x) = g(\theta^T x) \quad | \quad \begin{array}{l} \text{suppose "y=1" if } h_{\theta}(x) \geq 0.5 \\ \text{"y=0" if } h_{\theta}(x) < 0.5 \end{array}$$

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

$h_{\theta}(x)$  gives a probability that  $y=1$



wherever

$$\theta^T x \geq 0$$

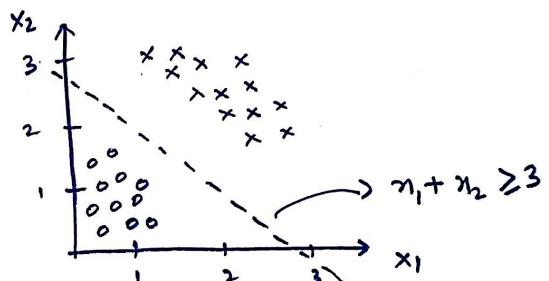
Example:

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

$$\begin{matrix} \downarrow \\ -3 \end{matrix} \quad \begin{matrix} \downarrow \\ +1 \end{matrix} \quad \begin{matrix} \downarrow \\ +1 \end{matrix}$$

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

$$\theta^T x$$



→ predict 'y=1' if  $\boxed{-3 + \theta_1 x_1 + \theta_2 x_2 \geq 0} \Rightarrow \theta_1 x_1 + \theta_2 x_2 \geq 3$

Decision boundary

is a property of

$h_{\theta}(x)$  & not of the data set

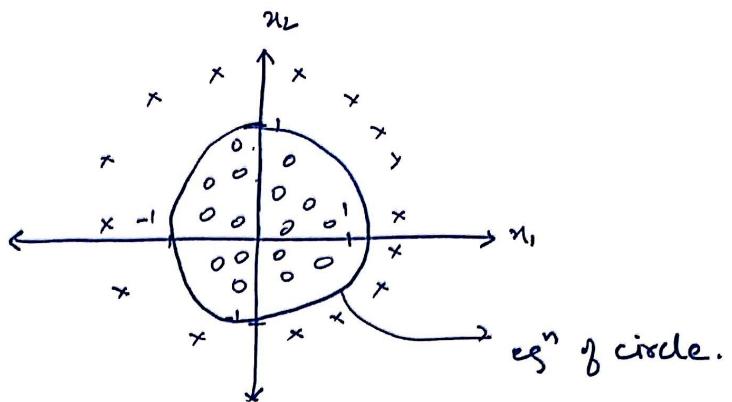
(or  
parameters)

## ● Non-linear Decision Boundaries :

(36)

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

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

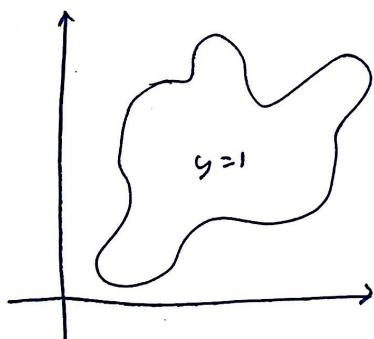


predict "y=1" if  $-1 + x_1^2 + x_2^2 \geq 0 \rightarrow x_1^2 + x_2^2 \geq 1$

\*\*

# Parameters ( $\theta$ ) defines the decision boundary

# Training set may be used to fit the parameters  $\theta$



$$h_{\theta}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_2^2 + \theta_5 x_1^2 x_2 + \theta_6 x_2^2 x_1 + \dots)$$

● COST FUNCTION:

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

m examples.

$$x \in \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix} \in \mathbb{R}^{n+1}$$

$x_0 = 1$   
 $y \in \{0, 1\}$

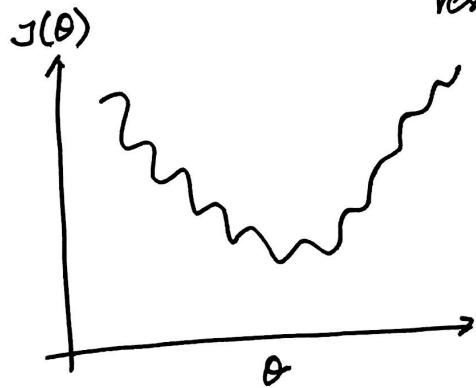
$$h_{\theta}(x) = \frac{1}{1 + e^{-\theta^T x}}$$

} → How to choose parameter  $\theta$ ?

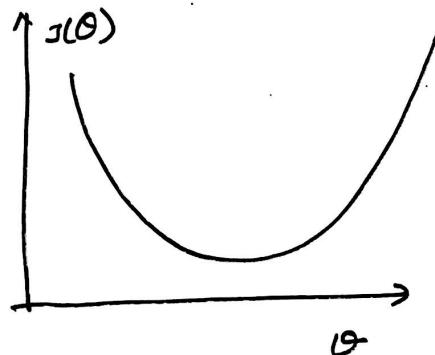
if we use the cost fn. of linear regression

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

↓ result in      ↓ want



Non-convex

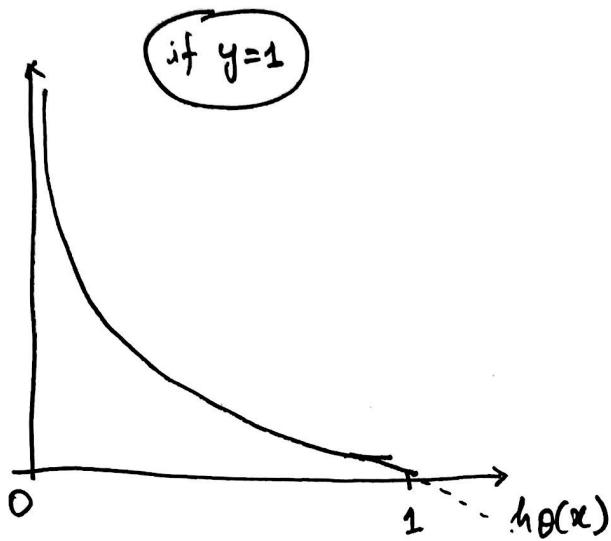


Convex

⇒ logistic regression Cost function

(38)

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

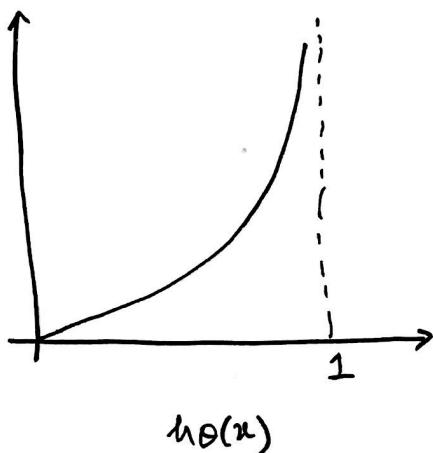


①  $\text{cost} = 0$  if  $y=1$  &  $h_{\theta}(x)=1$

② But as  $h_{\theta}(x) \rightarrow 0$  }  
 $\text{cost} \rightarrow \infty$  ]

captures intuition that if  
 $h_{\theta}(x) = 0$ , (predict  $P(y=1|x; \theta)$ ) = 0  
 but  $y=1$ , we will  
 penalize learning algorithm by  
 a very large cost.

if  $y=0$



②  $\text{cost} = 0$  if  $y=0$ ,  $h_{\theta}(x)=0$

● Simplified Cost function and Gradient descent

$$\left[ \text{cost}(\hat{h}_\theta(x), y) = \underbrace{-y \cdot \log(\hat{h}_\theta(x))}_{\text{Term 1}} - \underbrace{(1-y) \log(1-\hat{h}_\theta(x))}_{\text{Term 2}} \right]$$

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

To fit parameter  $\theta$ : → Get  $\theta$

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

$\theta$

To make a prediction given new  $x$ :

$$\text{Output: } \hat{h}_\theta(x) = \frac{1}{1 + e^{-\theta^T x}}$$

# Gradient descent:

Repeat until convergence(num-iters):

{

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

#(simultaneously  
update  
all  $\theta_j$ )

}

$$\left[ \frac{\partial \circ(\theta)}{\partial \theta_j} = \frac{1}{m} \sum_{i=1}^m (\hat{h}_\theta(x^i) - y^i) x_j^i \right] \begin{matrix} \text{derivative of } \frac{\partial J(\theta)}{\partial \theta_j} \\ \text{for logistic regression as well.} \end{matrix}$$

● Advanced Optimization Algorithms:

(40)

$$\partial J = (x^T (x^T \theta - y)) / m$$

Different Optimizations Algorithms:

- gradient descent
  - conjugate descent
  - BFGS
  - L-BFGS
- } # No need to manually pick d.  
} often faster than gradient  
descent  
} disadvantage: More complex.
- 

Example:

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

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

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

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

Cost Function in Octave

function [jval, gradient] = costFunction(theta).

$$jval = (\theta(1) - 5)^2 + (\theta(2) - 5)^2 ;$$

$$gradient = zeros(2, 1);$$

$$gradient(1) = 2 * (\theta(1) - 5)$$

$$gradient(2) = 2 * (\theta(2) - 5)$$

(41)

Data structure.

#  $\underbrace{\text{options}}_{\text{options}} = \text{optimset}(\text{'GradObj'}, \text{'on'}, \text{'MaxIter'}, \text{'100'})$ ;

$\text{initialTheta} = \text{zeros}(2, 1)$ ;

$[\text{optTheta}, \text{functionVal}, \text{exitFlag}] \dots$

$= \underbrace{\text{fminunc}}_{\text{fn. minimization}}(@\text{cost function}, \underbrace{\text{initialTheta}}_{\theta \in \mathbb{R}^d}, \text{options})$ ;

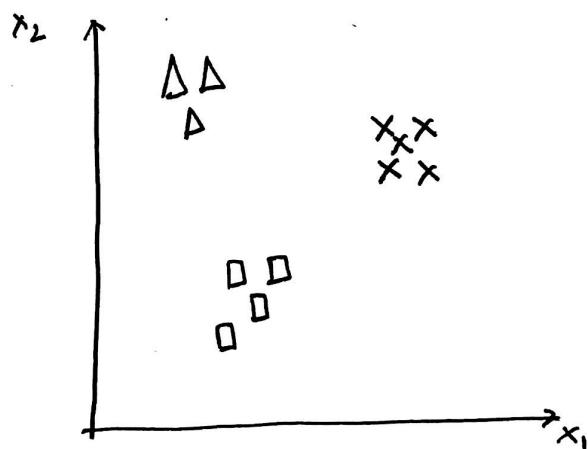
$\underbrace{\theta \in \mathbb{R}^d}_{\text{d=2}}$

fn. minimization  
unconstrained

## ● Multi-class classification Problem:

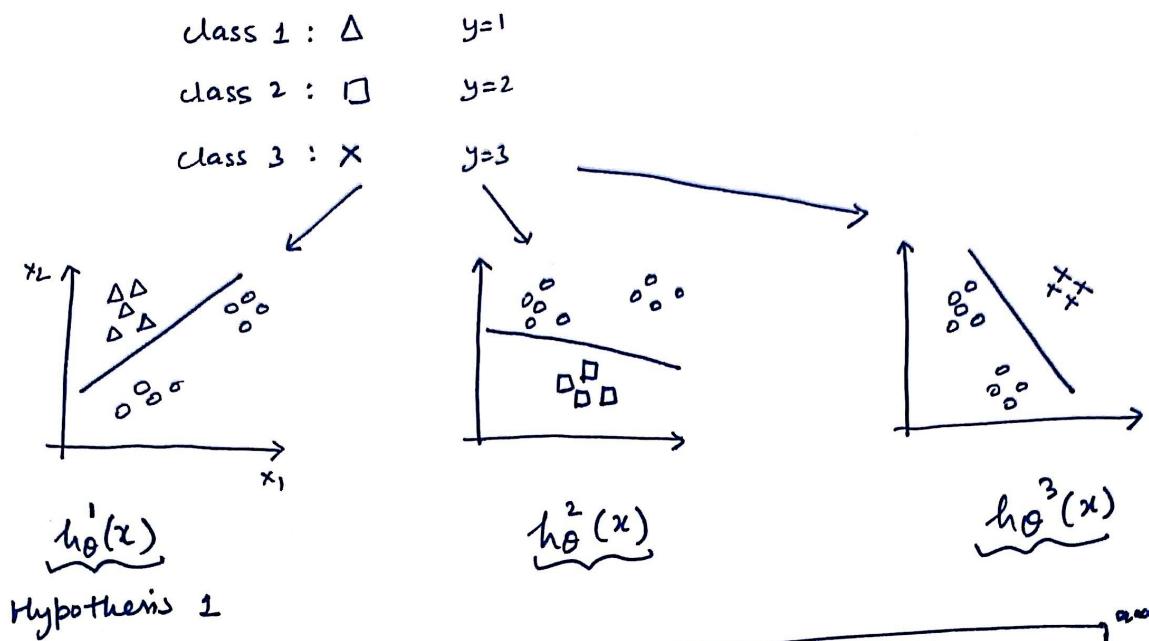
example: ① Email foldering/tagging:  $\underbrace{\text{Work}}_{y=1}$ ,  $\underbrace{\text{friends}}_{y=2}$ ,  $\underbrace{\text{family}}_{y=3}$ ,  $\underbrace{\text{Hobby}}_{y=4}$

② Medical diagrams :  $\underbrace{\text{Not ill}}_{y=1}$ ,  $\underbrace{\text{cold}}_{y=2}$ ,  $\underbrace{\text{flu.}}_{y=3}$



## # One-vs-all classification:

(42)



$$h_\theta^i(x) = P(y=i \mid x; \theta) \quad i \in \{1, 2, 3\}$$

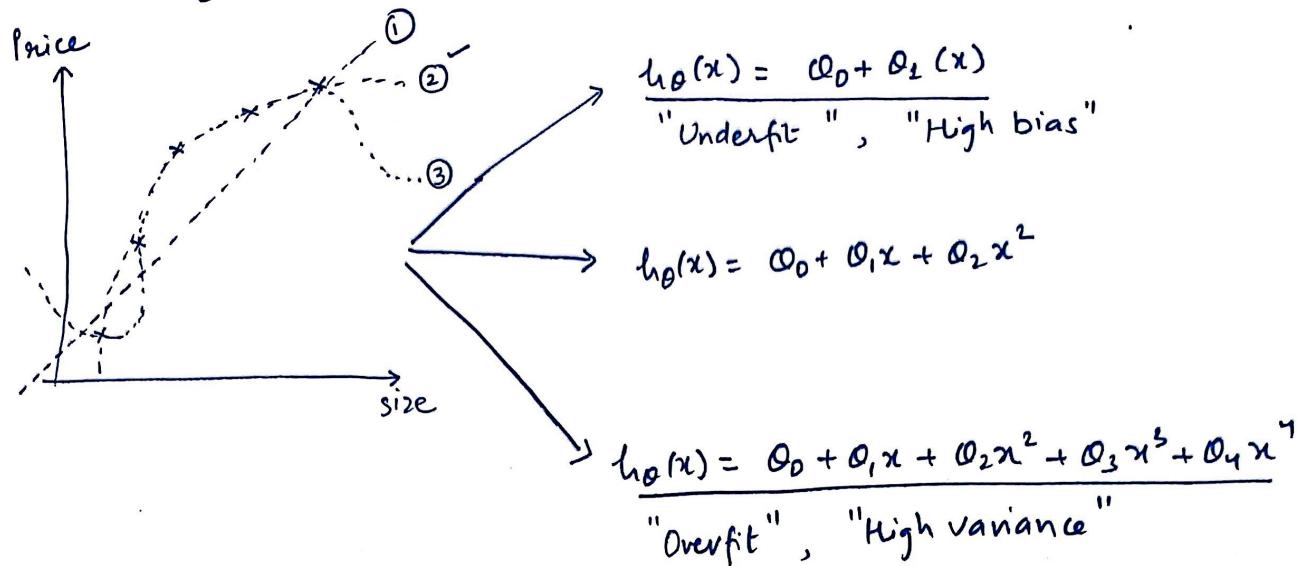
Probability that  $y=i$  (belongs to class  $i$ )  
 ↑  
 class  $i$

- Train a logistic regression classifier  $h_\theta^i(x)$  for each class  $i$  to predict the probability that  $y=i$
- On a new input  $x$ , to make a prediction, pick a class  $i$  that maximizes  $h_\theta^i(x)$ : maximizing the prob.

$$\max_i h_\theta^i(x)$$

## ● Solving the problem of Overfitting :

(43)



Overfitting: If we have too many features, the learned hypothesis may fit the training set very well

$$(J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 \approx 0), \text{ but fail to}$$

generalize to new examples (predict prices on new examples)

how well  
it new applies  
to new  
data.

[ It may work very well on training data ]  
but doesn't work on actual data )

## # Addressing Overfitting : options.

### ① Reduce number of features

- Manually select which features to keep.
- Model selection Algorithm (o)

### ② Regularization:

- Keep all the features, but reduce magnitude / values of parameters  $\theta_j$
- Works well when we have a lot of features, each of which contributes a bit to predicting  $y$ .

## \* Regularization : Cost function :

(44)

- Small values for values parameter ( $\theta_0, \dots, \theta_n$ )
  - "Simpler" hypothesis } Penalize all the
  - less prone to overfitting } parameters.

### Example :

We penalize and make  $\theta_3, \theta_4$  really small.

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

$\theta_3, \theta_4 \approx 0$

Ex.:

- Features:  $x_0, x_1, x_2, \dots, x_n$
- Parameters:  $\theta_0, \theta_1, \theta_2, \dots, \theta_{100}$

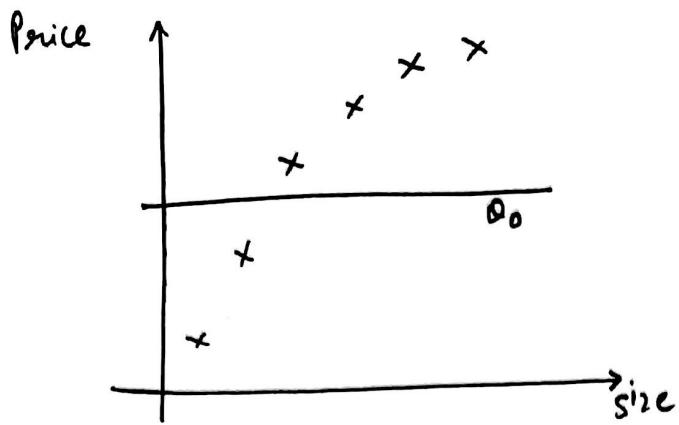
} we don't know which features/param. to penalize.

$$J(\theta) = \frac{1}{2m} \left[ \sum_{i=1}^m (h_{\theta}(x^i) - y^i)^2 + \lambda \sum_{j=1}^n \theta_j^2 \right] \rightarrow \text{This will result in } \theta_j \approx 0$$

regulation term.  
Regularization parameter [Not going to penalize  $\theta_0$ . very little difference]

\* if  $\lambda$  is set to extremely large value (for instance,  $\lambda = 10^{10}$ ) then we might have "Underfit" case.

\* while solving for  $\theta$ ,  $J(\theta)$  will become min.,  $\theta_j \approx 0$ .



our  $\min_{\theta} J(\theta)$   
will give values of  $\theta_j \approx 0$

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

we might end up  
"Underfitting our data".

## # REGULARIZED LINEAR REGRESSION ( $\min_{\theta} J(\theta)$ )

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

Repeat until convergence

①

$$\theta_0 := \theta_0 - \alpha \left[ \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^i) - y^i) \cdot x_0^i \right] \xrightarrow{\frac{\partial J}{\partial \theta_0}}$$

$$\theta_j := \theta_j - \alpha \left[ \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^i) - y^i) \cdot x_j^i + \frac{\lambda}{m} \theta_j \right] \xrightarrow{\frac{\partial J(\theta)}{\partial \theta_j}}$$

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

↓  
shrink  
the parameter  
a little bit

(2)

NORMAL EQUATION

(46)

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

design matrix

$$y = \begin{bmatrix} y^1 \\ y^2 \\ \vdots \\ y^m \end{bmatrix}_{\mathbb{R}^m}$$

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

Normal.

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

while using regularization

# Non-invertibility

when  $m \leq n$  features } too less data available to  
 training set curve fit

if  $\lambda > 0$ 

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

↳ non-invertible/singular/degenerate.

This matrix  
will be invertible.

## # REGULARIZED LOGISTIC REGRESSION :

(47)

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

$\theta \in (\theta_0, \dots, \theta_n)$

### ① Gradient descent

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

Repeat {

$$\theta_0 := \theta_0 - \alpha \left[ \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x) - y)^2 \right] x_0^{(i)}$$

$$\theta_j := \theta_j - \alpha \left[ \underbrace{\left( \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)} \right)}_{\frac{\partial J}{\partial \theta}} + \frac{\lambda}{m} \theta_j \right]_{j=1 \dots n}$$

}

## Summary of Linear/Logistic Regression

$$h_{\theta}(x) = \theta^T x = \sum_{i=0}^n \theta_i x_i$$

$\hookrightarrow$

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

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

$$h_{\theta}(x) = \underbrace{g(\theta^T x)}_{\text{sigmoid}} = \frac{1}{1 + e^{-\theta^T x}}$$

$\hookrightarrow$

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

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

$$\frac{\partial J}{\partial \theta} = \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x) - y) x^i_j \quad [ \lambda = 0 ]$$

for both linear as well as logistic regression

$$\frac{\partial J}{\partial \theta_0} = \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x) - y) x^i_0$$

for regularized linear/  
logistic regression

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

Note: for complex polynomial regression. (Map Feature)

for  $i = 1$ : degree.

for  $j = 0:i$

output  $(:, end+1) = (x1.^{(i-j)}). * (x2.^{(j)})$ ;

example:  $h_{\theta}(x) = \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_2^2 + \dots + \theta_{28} x_1^6 x_2^6$

$$\theta = 28 \times 1 \quad x = m \times 28$$

$$x = 28 \times 1 \quad y = 28 \times 1$$

## # Special: One-Vs- ALL regularized expression

$K \in \{1, \dots, 10\}$  : 10 classifiers.

↓  
classes

$\theta \in \mathbb{R}^{K \times (n+1)}$  : Each row of  $\theta$ , will contain parameters for class  $i$

### # Program

```

function [all-theta] = oneVsAll [X, y, num-labels, lambda]
{
    m = size(X, 1);
    n = size(X, 2);
    X = [ones(m, 1) X];
    all-theta = zeros(num-labels, n+1);
    for c = 1:num-labels,
        {
            initial-theta = zeros(n+1, 1)
            options = optimset ('GradObj', 'on', 'MaxIter', 50);
            all-theta(:,c) = ...
                fmincg (@t)(bcostfunction (t, X, (y==c),
                lambda)), ...
                (initial-theta, options));
        }
}

```

% ( $y == c$ ) creates a logical array for a class.

$$\begin{array}{l}
 \text{original } y = \begin{bmatrix} 1 \\ 2 \\ 3 \\ 4 \\ 5 \end{bmatrix} \\
 \text{transform. } y == 2 \Rightarrow y = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}
 \end{array}$$

function  $p = \text{predictOneVsAll}(\text{all-theta}, x)$

{

$m = \text{size}(x, 1)$

$p = \text{zeros}(m, 1)$

$x = [\text{ones}(m, 1) \ x]$ ;

$\text{predict} = \text{sigmoid}(x^* \text{all-theta}')$ ;

$[n, p] = \text{max}(\text{predict}, [ ], 2)$ ;

}