

Stanford ML

WEEK 3:

Classification : Logistic Regression

- Email: Spam / Not Spam ?
- Online Transactions: Fraudulence (Y/N ?)
- Tumour: Malignant / Benign ?

OUTPUT

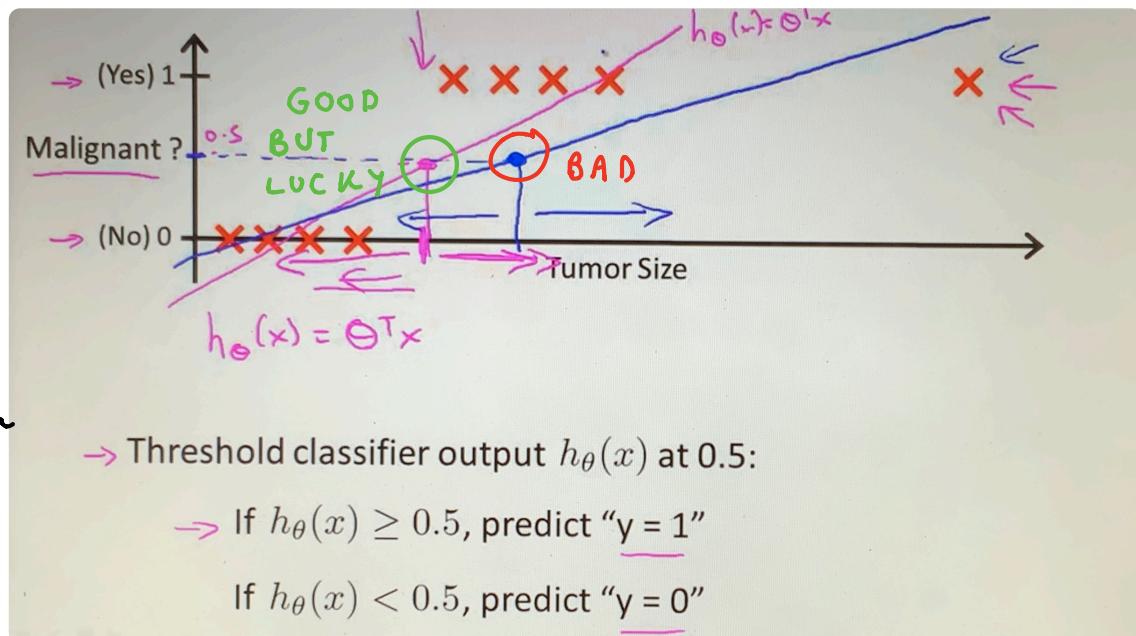
$$y \in \{0, 1\}$$

0: "Negative Class"

1: "Positive Class"

- Of course, multiclass classification is possible and common

- Example plots with linear regression and justifies need for classification algorithm



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

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

If $h_\theta(x) < 0.5$, predict "y = 0"

- For linear regression, we saw the hypothesis must take the values

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

LOGISTIC REGRESSION

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

LIN. REG.

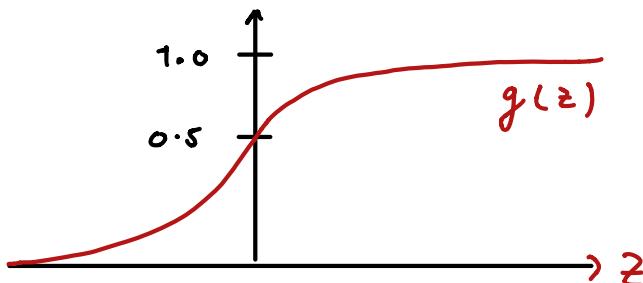
$$h_{\theta}(x) = \theta^T x$$

LOG. REG.

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

where $g(z)$ is the SIGMOID function

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



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

Interpretation of Hypothesis Output

- $h_{\theta}(x)$ = estimated probability that $y=1$ on input x

Ex:

$$\text{if } x = \begin{bmatrix} x_0 \\ x_1 \end{bmatrix} = \begin{bmatrix} 1 \\ \text{tumorsize} \end{bmatrix} \text{ and } h_{\theta}(x) = 0.7$$

∴ 70% chance of being malignant

$$\Rightarrow h_{\theta}(x) = P(y=1 | x; \theta)$$

“probability that $y=1$, given x , parameterized by θ ”

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

for BINARY CLASSIFICATION

Decision Boundary

- At what value of $h_{\theta}(x)$ do we decide that the "y = 1"

Ex:

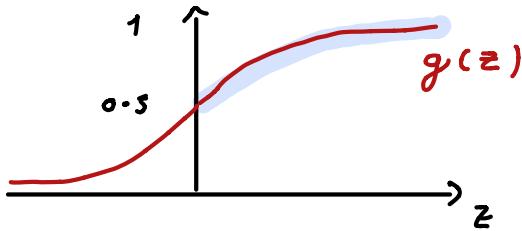
- $h_{\theta}(x) = g(\theta^T x)$
- $g(z) = \frac{1}{1 + e^{-z}}$
- Suppose predict $y=1$ if $h_{\theta}(x) \geq 0.5$
- Suppose predict $y=0$ if $h_{\theta}(x) < 0.5$

* By analogous argument

$$h_{\theta}(x) = g(\theta^T x) < 0.5$$

when $\theta^T x < 0$ so

predict $y=0$



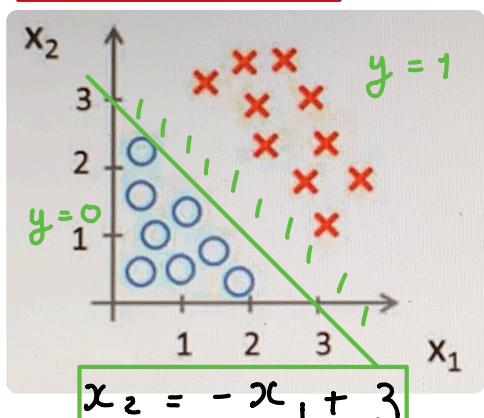
- * $g(z) \geq 0.5$ when $z \geq 0$
- * $h_{\theta}(x) = g(\theta^T x) \geq 0.5$ whenever

$$\underbrace{\theta^T x}_{z} \geq 0$$

predict $y = 1$

Ex: Constructing the Decision Boundary

TRAINING SET



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

$$\rightarrow \theta = \begin{bmatrix} -3 \\ 1 \\ 1 \end{bmatrix} \text{ chosen by some optimization procedure}$$

∴ Predict "y = 1" if:

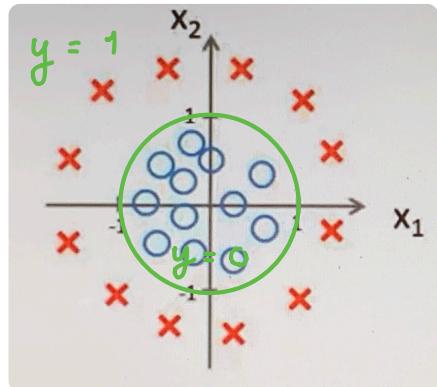
$$\underbrace{-3 + x_1 + x_2}_{\theta^T x} \geq 0$$

OR

$$x_1 + x_2 \geq 3$$

$\theta^T x$
this is the DECISION BOUNDARY

Non - Linear Decision Boundary



- In the same way that we can add higher orders of original features to conduct polynomial regression

- The same can be done for non-linear decision boundaries

EQN. OF CIRCLE

ME

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

ME

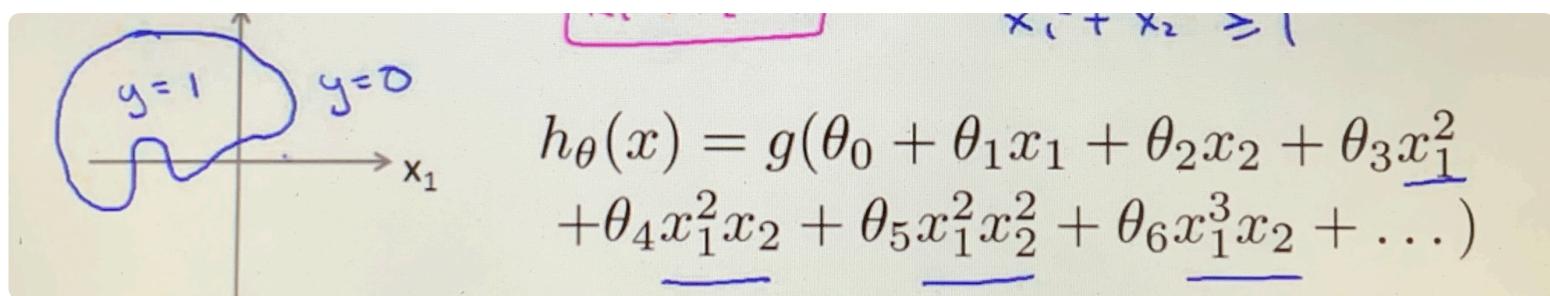
$$\theta = \begin{bmatrix} -1 \\ 0 \\ 0 \\ 1 \\ 1 \end{bmatrix} \quad \therefore \text{predict } y = 1 \quad y = -1 + x_1^2 + x_2^2$$

$x_1^2 + x_2^2 \geq -1$

CIRCLE
 $(0, 0)$
 $r^2 = 1$

NOTE :

- The decision boundary is a PROPERTY OF THE HYPOTHESIS



$$\theta^T x \geq 0 \Rightarrow y = 1$$

$$\theta^T x < 0 \Rightarrow y = 0$$

Cost Function

- Set up

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

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

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

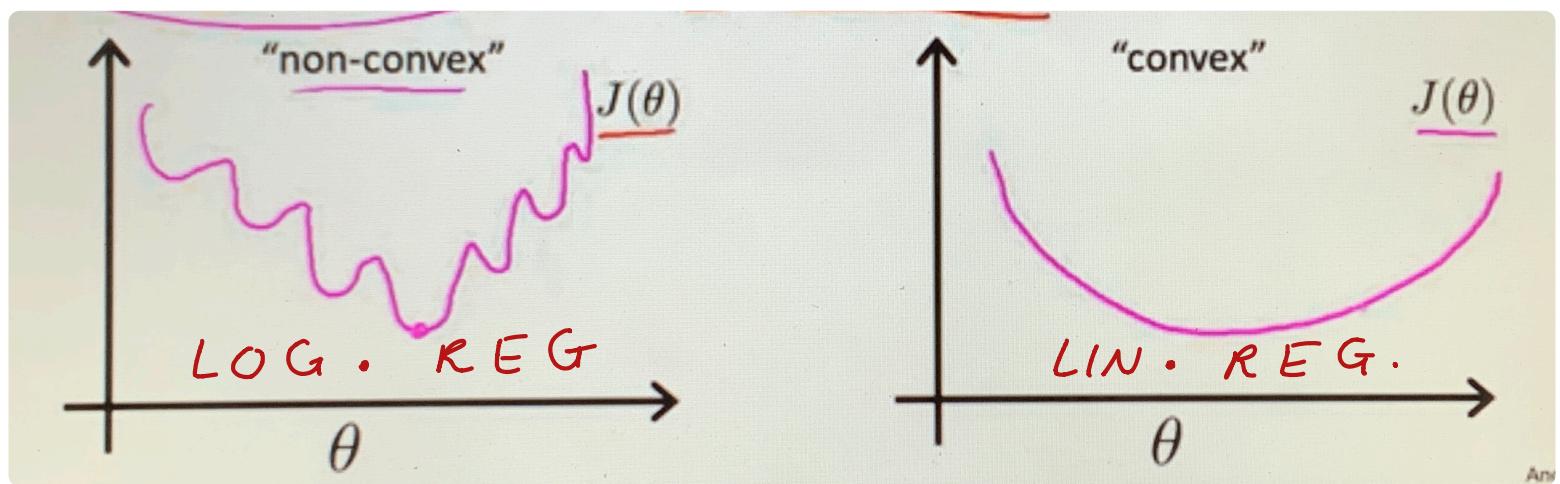
- Need to choose parameters θ

- Linear Regression

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

$$\therefore \text{Cost}(h_{\theta}(x), y) = \frac{1}{2} [h_{\theta}(x) - y]^2 \quad \forall i = 1, \dots, m$$

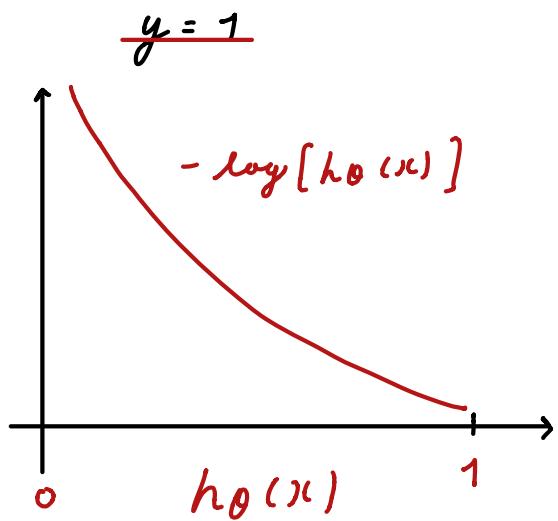
↳ For LOGISTIC REGRESSION where $h_{\theta}(x) = g(\theta^T x)$



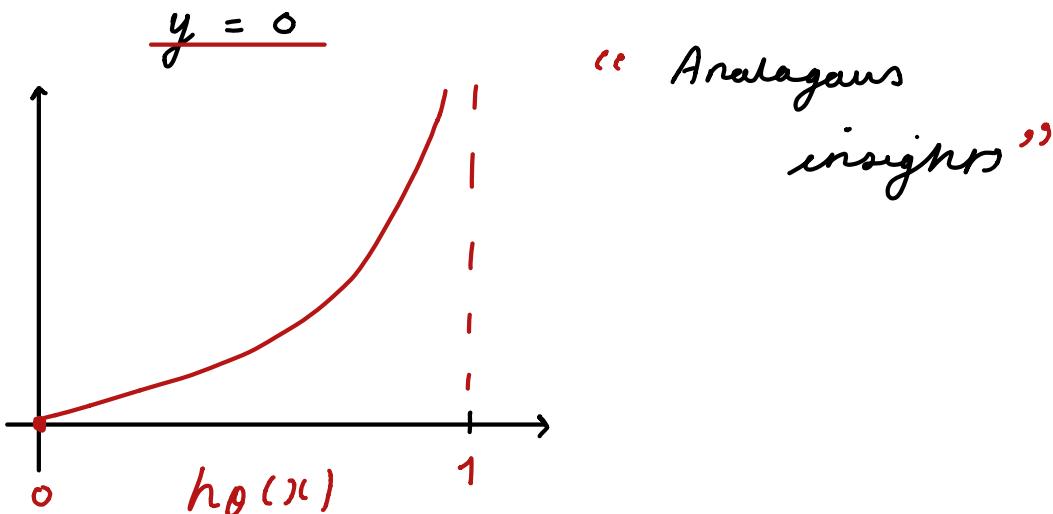
- Using the SQ. ERR cost function for LOG REG creates a "non-convex" $J(\theta)$ meaning gradient descent CANNOT be applied.

→ LOG REG Cost Function

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



- $\text{Cost} = 0$ if $y = 1$, $h_\theta(x) = 1$
- 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 penalize with a very large cost



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

$$\text{Cost}(h_\theta(x), y) \rightarrow \infty \text{ if } y=0 \text{ and } h_\theta(x) \rightarrow 1$$

$$\text{Cost}(h_\theta(x), y) \rightarrow \infty \text{ if } y=1 \text{ and } h_\theta(x) \rightarrow 0$$

Logistic regression cost function

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

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

Note: $y = 0$ or 1 always

COMBINE



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

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

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

- To fit parameters θ :

minimize $J(\theta)$, to allow us to make a prediction given new x

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

- Use GRADIENT DESCENT AGAIN

Repeat {

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta) \quad \begin{matrix} \nearrow \text{SIMULTANEOUS} \\ \nearrow \theta \text{ UPDATE} \end{matrix}$$

}

$$\frac{\partial}{\partial \theta_j} J(\theta_i) = \frac{\partial}{\partial \theta_j} \left\{ -\frac{1}{m} \sum_{i=1}^m y^{(i)} \log(\theta_j^T x) + (1-y^{(i)}) \times \log(1 - \theta_j^T x) \right\}$$

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

Repeat {

$$\theta_j := \theta_j - \alpha \sum_{i=1}^m [h_\theta(x^{(i)}) - y^{(i)}] x_j^{(i)} \quad \begin{matrix} \nearrow \text{ONLY DEF OF} \\ \nearrow h_\theta(x^{(i)}) \text{ CHANGE} \end{matrix} \quad \begin{matrix} \nearrow \text{to linear} \\ \nearrow \text{regression} \end{matrix}$$

iterations

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

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

$$\begin{aligned} \frac{\partial}{\partial \theta_j} g(\theta^T x) &= g'(\theta^T x) = - (1 + e^{-\theta_k x_k})^2 e^{-\theta_k x_k} \cdot x_j \\ &= \frac{-e^{-\theta_k x_k} x_j}{(1 + e^{-\theta_k x_k})^2} = \underbrace{g(\theta^T x)}_{\frac{1}{1 + e^{-\theta^T x}}} \underbrace{(1 - g(\theta^T x))}_{-\frac{e^{-\theta^T x}}{1 + e^{-\theta^T x}}} x_j \end{aligned}$$

$$\begin{aligned} \Rightarrow \frac{\partial}{\partial \theta_j} J(\theta) &= \left(y \frac{1}{g(\theta^T x)} - (1-y) \frac{1}{1-g(\theta^T x)} \right) \frac{\partial}{\partial \theta_j} g(\theta^T x) \\ &= \left(\frac{y}{g(\theta^T x)} - \frac{(1-y)}{1-g(\theta^T x)} \right) g(\theta^T x) (1 - g(\theta^T x)) x_j \end{aligned}$$

$$\Rightarrow \{ y (1 - g(\theta^T x)) - (1-y) g(\theta^T x) \} x_j$$

$$\therefore \frac{\partial}{\partial \theta_j} J(\theta) = \{ y - g(\theta^T x) \} x_j \cdot -\frac{1}{m} \sum$$

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

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

Gradient Descent

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

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

$$\forall j = 0, \dots, n$$

- A vectorized implementation could take the form

$$h = g(\underline{\underline{X}} \theta)$$

↗ Rows $m \times (n+1)$
 ↗ cols $(n+1) \times 1$
 $\}$ $m \times 1$ VECTOR

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

Row COL Row CCL

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

$$[\underline{\underline{X}}] = m \times (n+1) \longrightarrow [\underline{\underline{X}}^T] = (n+1) \times m$$

$$[y] = [g(\underline{\underline{X}} \theta)] = m \times 1$$

$$\Rightarrow [\theta] = (n+1) \times 1 \text{ VECTOR}$$

Advanced Optimisation

Optimization algorithm

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

Given θ , we have code that can compute

$$\begin{aligned} \rightarrow & - J(\theta) \\ \rightarrow & - \frac{\partial}{\partial \theta_j} J(\theta) \quad (\text{for } j = 0, 1, \dots, n) \end{aligned}$$

Gradient descent:

Repeat {

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

}

PROVIDE

Other Examples

① Conjugate gradient

② BFGS

③ L - BFGS

↳ quicker than standard
gradient descent

PROS:

- No need to pick α
- Faster

CONS:

- More complex

Example:

$$\begin{aligned} \rightarrow \theta &= \begin{bmatrix} \theta_1 \\ \theta_2 \end{bmatrix} \quad \min_{\theta} J(\theta) \\ \rightarrow & \theta_1 = 5, \theta_2 = 5. \\ \rightarrow & J(\theta) = (\theta_1 - 5)^2 + (\theta_2 - 5)^2 \\ \rightarrow & \frac{\partial}{\partial \theta_1} J(\theta) = 2(\theta_1 - 5) \\ \rightarrow & \frac{\partial}{\partial \theta_2} J(\theta) = 2(\theta_2 - 5) \end{aligned}$$

```

options = optimset('GradObj', 'on', 'MaxIter', '100');
initialTheta = zeros(2,1);
[optTheta, functionVal, exitFlag] ...
    = fminunc(@costFunction, initialTheta, options);

```

```

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);

```

POINTER
TO COST FUNC

$$\theta \in \mathbb{R}^d \quad d \geq 2$$

Algorithm

$$\underline{\theta} = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix} \quad \begin{array}{l} \text{---} \theta_0 \text{ ---} \theta(1) \\ \text{---} \theta_1 \text{ ---} \theta(2) \\ \vdots \\ \text{---} \theta_n \text{ ---} \theta(n+1) \end{array}$$

```
function [jVal, gradient] = costFunction(theta)
    jVal = [ code to compute  $J(\theta)$  ];
    gradient(1) = [ code to compute  $\frac{\partial}{\partial \theta_0} J(\theta)$  ];
    gradient(2) = [ code to compute  $\frac{\partial}{\partial \theta_1} J(\theta)$  ];
    :
    gradient(n+1) = [ code to compute  $\frac{\partial}{\partial \theta_n} J(\theta)$  ];
```

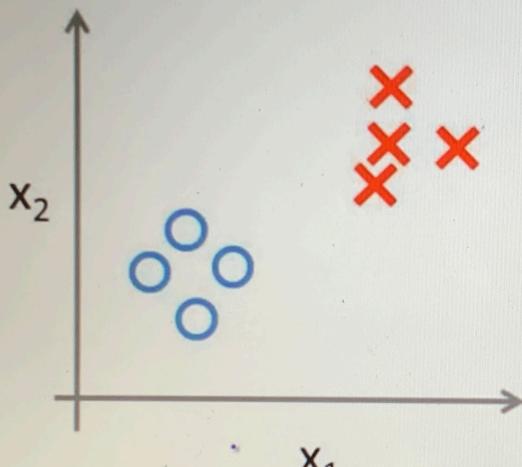
* Use for much larger problems

Multi-Class Classification

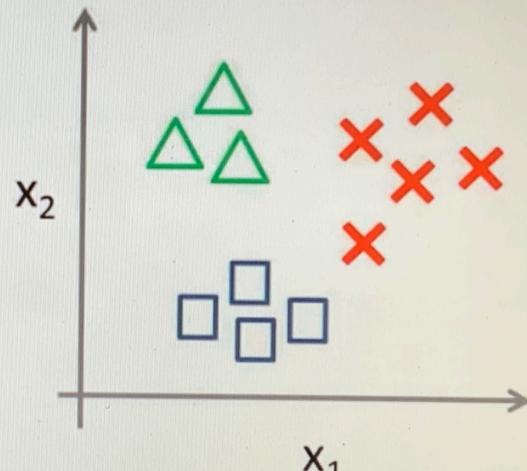
- E-mail tagging : Work, Friends, hobby
- Weather : Sunny, Cloudy, Rainy, Snow

0 1 2 3

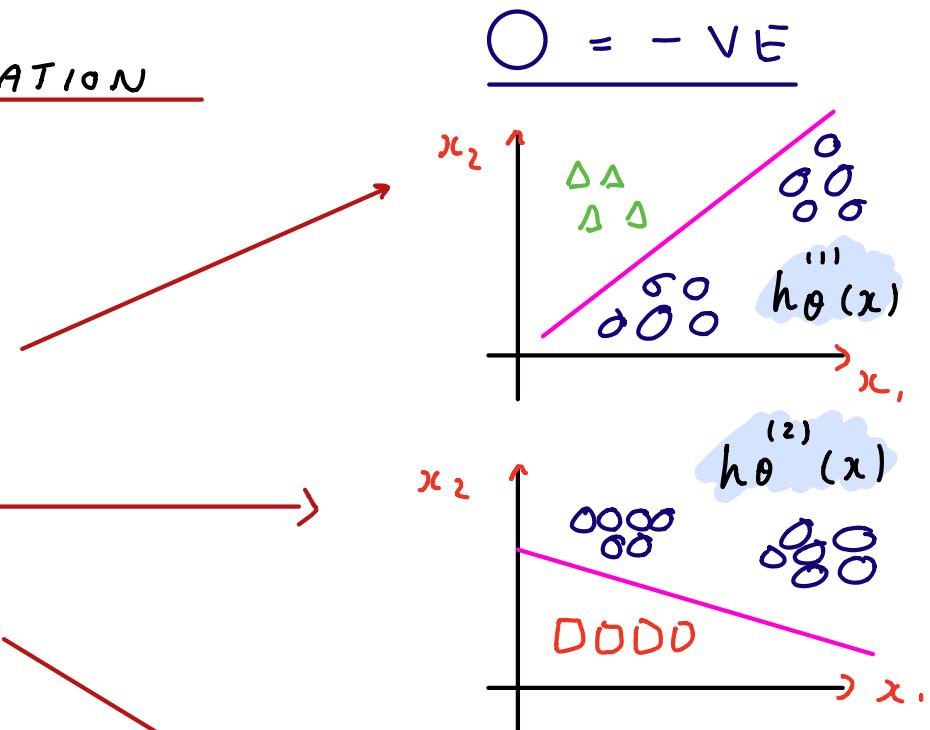
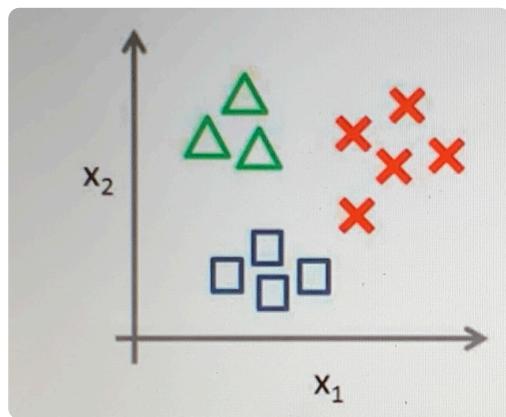
Binary classification:



Multi-class classification:



ONE VS. ALL CLASSIFICATION



- Turn into 3 separate binary problems
- Create new training set for each
- Train standard logistic regression

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

One-vs-all

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

On a new input x , to make a prediction, pick the class i that maximizes

$$\max_i \underline{h_{\theta}^{(i)}(x)}$$

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

$$h_{\theta}^{(0)}(x) = P(y=0 | x; \theta)$$

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

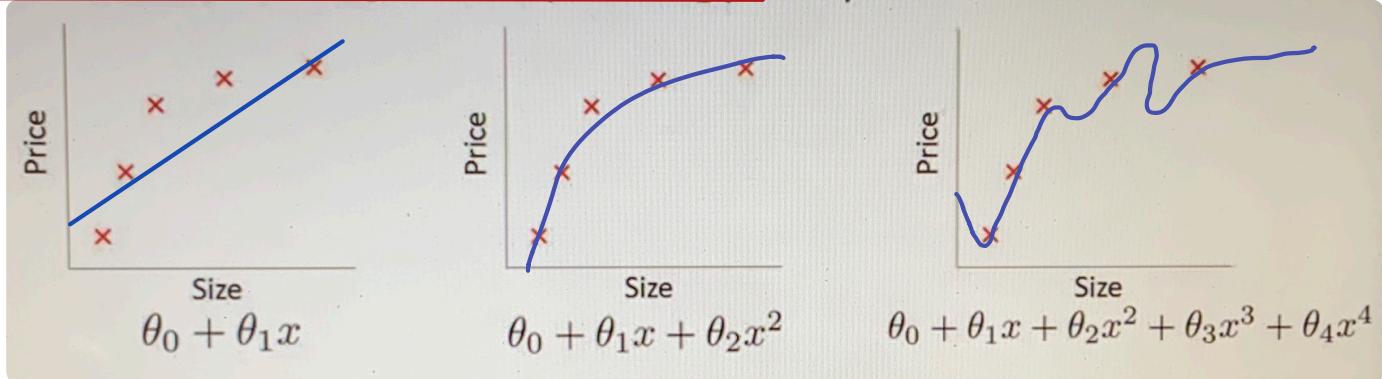
⋮

$$h_{\theta}^{(n)}(x) = P(y=n | x; \theta)$$

$$\text{prediction} = \max_i [h_{\theta}^{(i)}(x)]$$

Overshooting

LINEAR REGRESSION EXAMPLE



' UNDERFITTING '

' HIGH BIAS '

* Too few features !

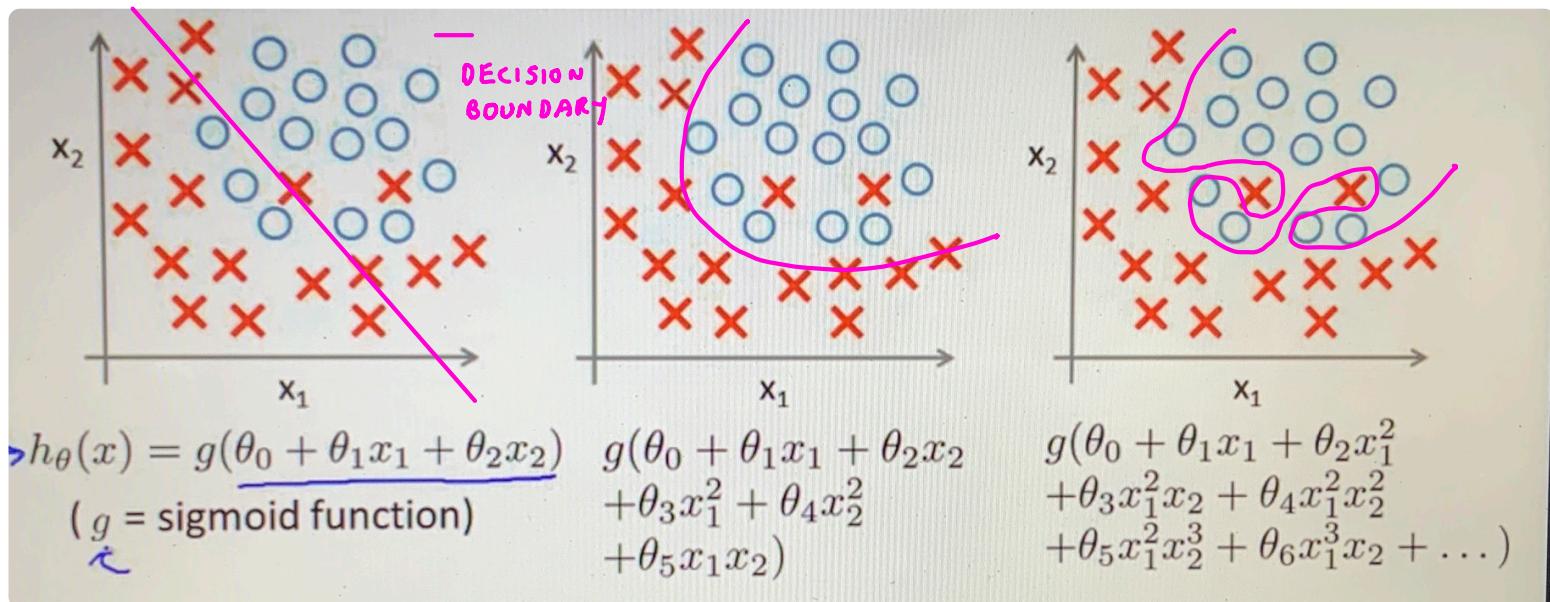
' OVERFITTING '

' HIGH VARIANCE OF HYPOTHESES '

* Too many features

Oversizing: If we have too many features, the learned hypothesis may fit the training set very well ($J(\theta) \approx 0$), but **FAILS TO GENERALIZE TO NEW CASES**

LOGISTIC REGRESSION EXAMPLE



UNDERFIT: Biased

Just Right!

OVERFIT: large variance

Addressing Oversizing:

① **Planning** hypotheses for the case of a **LOW** number of features - decide yourself

Options:

② **Reduce # of features**

- Manually select features to keep
- Model selection algorithm (**Lasso**)

③ **Regularization**

- Keep all features, $\{x_i\}$, but reduce magnitudes of $\{\theta_i\}$

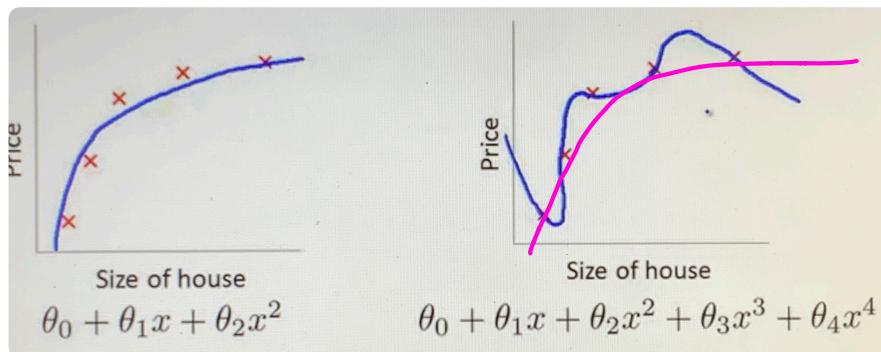
- Work well when



LOTS OF FEATURES ~ EQUIVALLY
IMPORTANT

USE REGULARIZATION

Regularization :: Cost Function



- Suppose we regularize and make θ_3 and θ_4 SMALL

PENALIZE CERTAIN FEATURES

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

- When this modified $J(\theta)$ is minimized, since $1000 \approx \text{LARGE} \dots$

$$\theta_3, \theta_4 \approx 0$$

in order for successful minimization of $J(\theta)$

In general

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

- Note that θ_0 is never penalized, we use this when we're not sure which parameters to penalize usually due to there being a large # of them.

λ - REGULARIZATION PARAMETER

- ↳ Convales the trade-off between fitting the training data well and keeping a simple $h_{\theta}(x)$ avoiding OVERFITTING

LARGE $\lambda \rightarrow$ LARGE PENALIZATION

↳ All $\{\theta_i\} \approx 0 \therefore$ UNDERFITTING

↳ $h_{\theta}(x) = \theta_0$

* λ INFLATES COST *

Regularized Linear Regression

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

$$\min \theta ; J(\theta)$$

Gradient Descent

Repeat {

$$\frac{\partial}{\partial \theta_0} J(\theta)$$

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

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

}

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

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

- Since $\alpha\lambda/m > 0$, $(1 - \alpha\lambda/m) < 1 \approx 0.99$
- \therefore We have added a term which effectively shrinks θ_j

Normal Equation

$$X = \begin{bmatrix} (x^{(1)})^\top \\ \vdots \\ (x^{(m)})^\top \end{bmatrix} ; \quad y = \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(m)} \end{bmatrix}$$

$m \times (n+1)$ $m \times 1$

$$\Rightarrow \theta = (X^\top X + \lambda \begin{smallmatrix} * & * & * \\ * & * & * \\ * & * & * \end{smallmatrix})^{-1} X^\top y$$

$(n+1) \times (n+1)$

- Note that $m < n$; $(X^\top X)^{-1} \rightarrow$ singular
- If $\lambda > 0$

$$(X^\top X + \lambda \begin{smallmatrix} * & * & * \\ * & * & * \\ * & * & * \end{smallmatrix})^{-1}$$

ALWAYS
INVERTIBLE

Regularized Logistic Regression

* LOTS OF FEATURES \longrightarrow OVERFITTING
 \longrightarrow REGULARIZATION

- New cost function:

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

L) Keeping parameters small by increasing the cost reduces the chances of overfitting

- Same gradient descent but

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

Advanced optimization

\rightarrow function [jVal, gradient] = costFunction(theta)

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

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

\rightarrow gradient(1) = [code to compute $\frac{\partial}{\partial \theta_0} J(\theta)$] ;

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

\rightarrow gradient(2) = [code to compute $\frac{\partial}{\partial \theta_1} J(\theta)$] ;

$$\left(\frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_1^{(i)} \right) + \frac{\lambda}{m} \theta_1$$

\rightarrow gradient(3) = [code to compute $\frac{\partial}{\partial \theta_2} J(\theta)$] ;

$$\vdots \quad \left(\frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_2^{(i)} \right) + \frac{\lambda}{m} \theta_2$$

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

$J(\theta)$