

Machine(Learning)

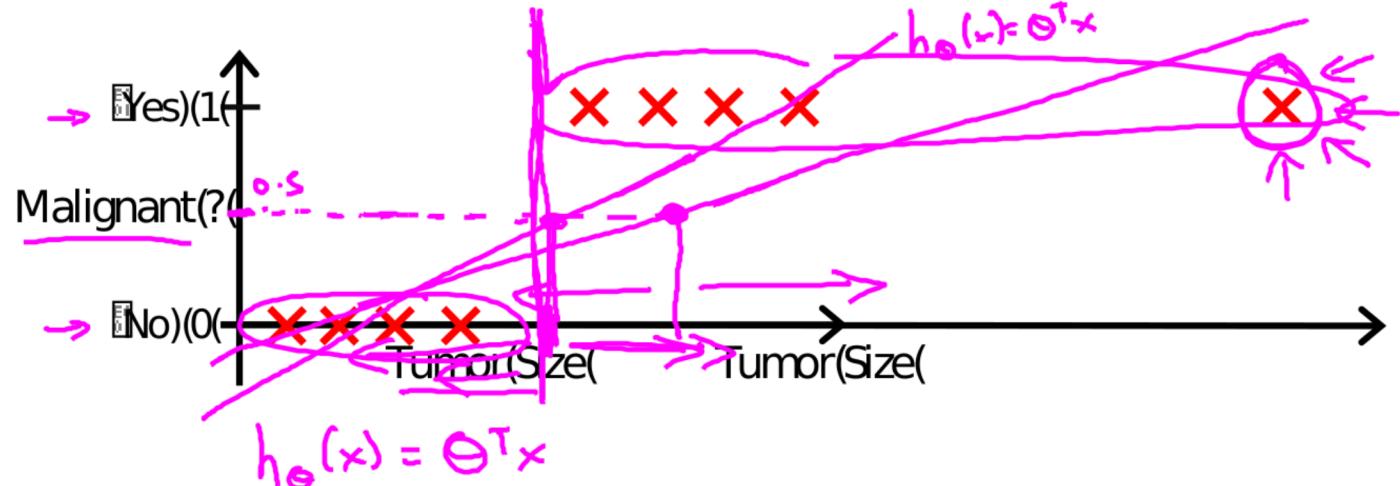
Logis~~tic~~c(

Regression(

Classification(

Classification+

- Email: (Spam / Not Spam)
 - Online Transactions: (Fraudulent / Yes / No)?
 - Tumor: (Malignant / Benign)?
- $y \in \{0, 1\}$
- 0: ("Negative Class" (e.g., benign(tumor))
 - 1: ("Positive Class" (e.g., malignant(tumor))
- $y \in \{0, 1, 2, 3\}$



→ Threshold(classifier(output)) at(0.5):

→ If($\theta_0 + \theta_1 x_1 > 0.5$, predict("y(=1")

If($\theta_0 + \theta_1 x_1 \leq 0.5$, predict("y(=0")

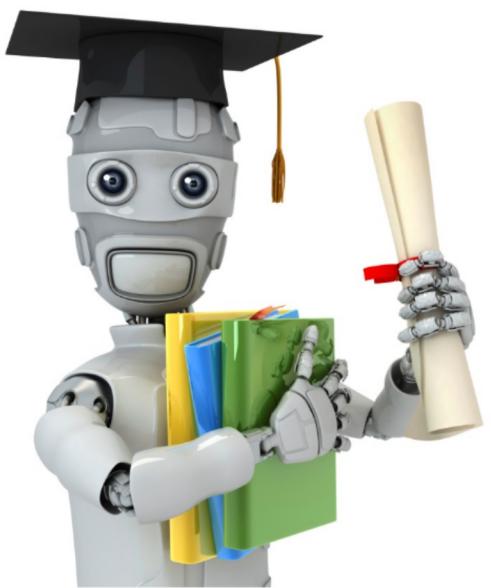
Andrew(Ng)

Classification: $y = 0$ or $y = 1$

$h_\theta(x)$ can be ≥ 1 or ≤ 0

Logistic Regression: $0 \leq h_\theta(x) \leq 1$

Classification



Machine Learning

Logistic Regression(

Hypothesis(

Representation(

Logistic Regression Model +

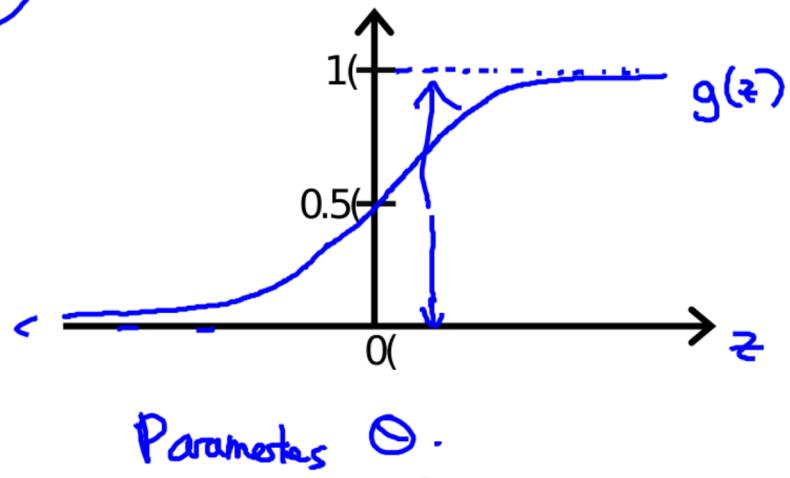
Want $0 \leq h_\theta(x) \leq 1$

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

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

↳ Sigmoid function
↳ Logistic function

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



Interpretation of Hypothesis Output +

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

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

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

$$\underline{h_{\theta}(x) = 0.7} \quad \underline{y=1}$$

Tell patient that 70% chance of tumor being malignant

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

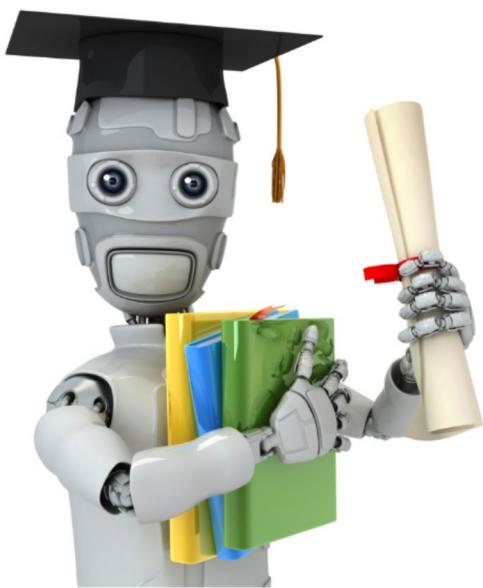
"probability that $y=1$, given x , parameterized by θ "

$$\underline{y = 0 \text{ or } 1}$$

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

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

Andrew(Ng)



Machine(Learning)

Logisc(Regression(--- Decision(boundary)

Logistic regression +

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

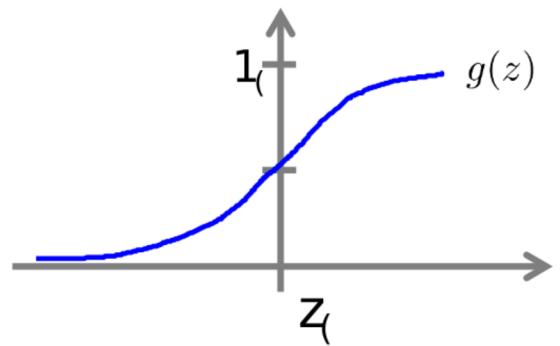
$$g(z) = \frac{1}{1+e^{-z}}$$

((Suppose predict("if($h_{\theta}(x) \geq 0.5$

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

((predict("if($h_{\theta}(x) < 0.5$

$$\theta^T x < 0$$



$$g(z) \geq 0.5$$

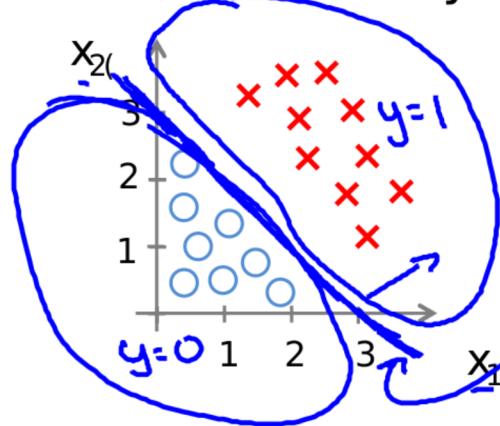
when $z \geq 0$

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

$$g(z) < 0.5$$

when $z < 0$

Decision-Boundary+



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

$$h_{\theta}(x) = g(\underline{\theta}_0 + \underline{\theta}_1 x_1 + \underline{\theta}_2 x_2)$$

$\underline{\theta}_0 = -3$
 $\underline{\theta}_1 = 1$
 $\underline{\theta}_2 = 1$

Decision boundary

Predict("if" $(-3 + x_1 + x_2 \geq 0)$)

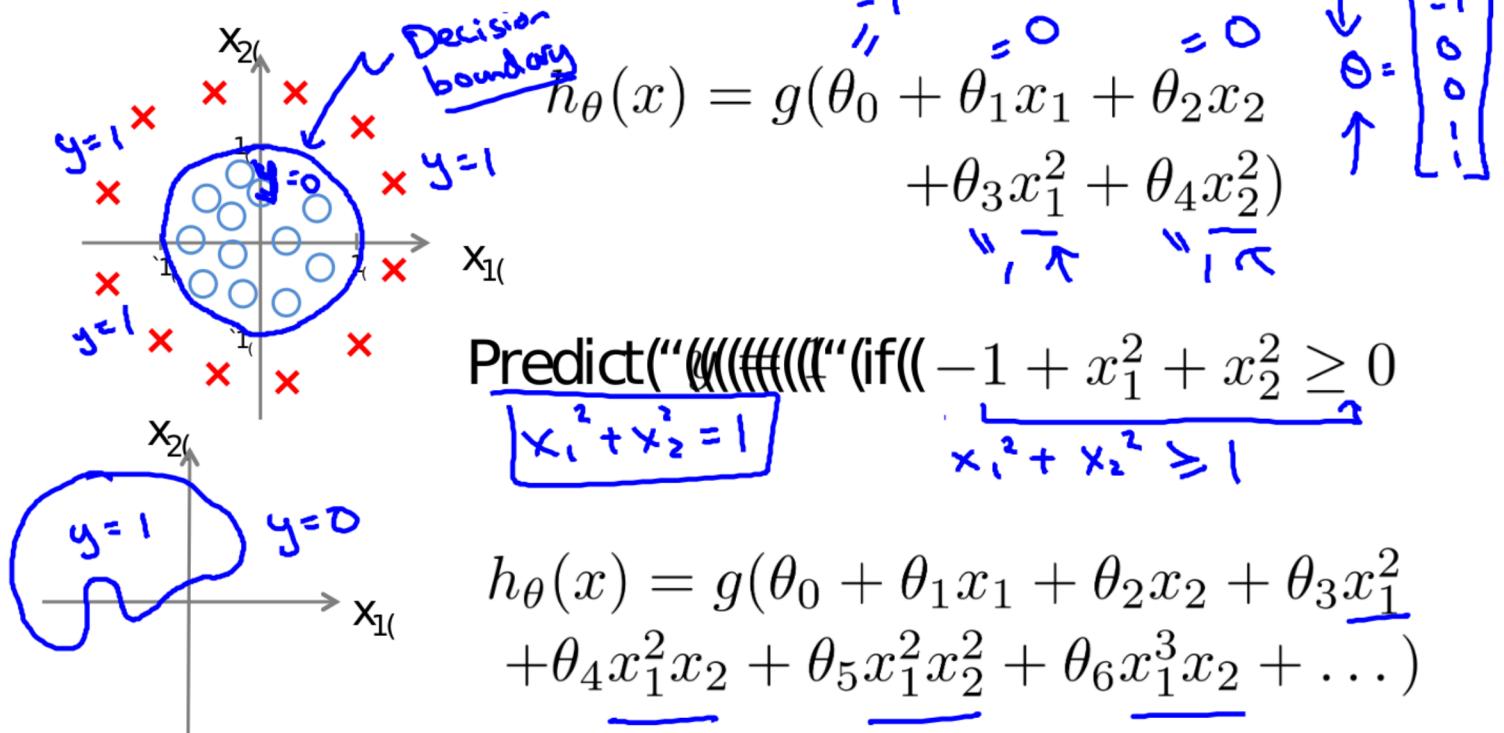
$$\Theta^T x$$

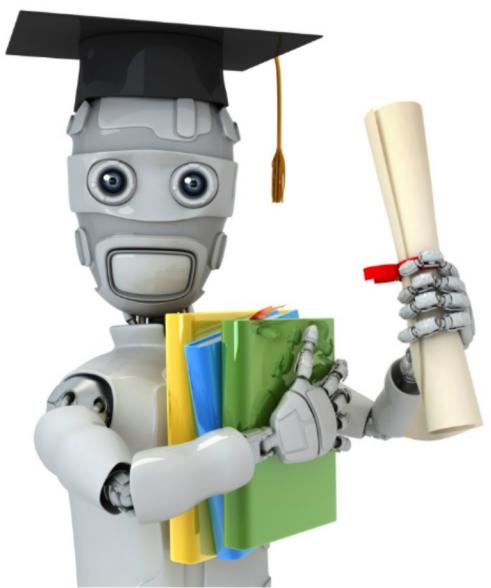
$$x_1 + x_2 \geq 3$$

$$\begin{aligned} & x_1, x_2 \\ \rightarrow h_{\theta}(x) &= 0.5 \\ & x_1 + x_2 = 3 \end{aligned}$$

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

Non-linear decision boundaries +





Machine(Learning(

Logis~~E~~c(

Regression(

Cost(fund~~E~~on(

Training set:

m(examples)

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

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

How to choose parameters ?

Cost function+

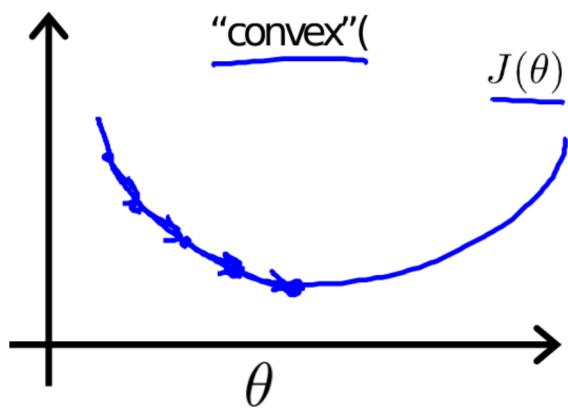
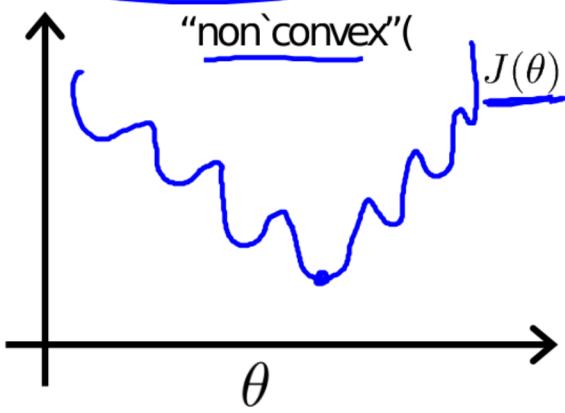
→ Linear regression: ($J(\theta) = \frac{1}{m} \sum_{i=1}^m \frac{1}{2} (h_\theta(x^{(i)}) - y^{(i)})^2$)

logistic

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

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

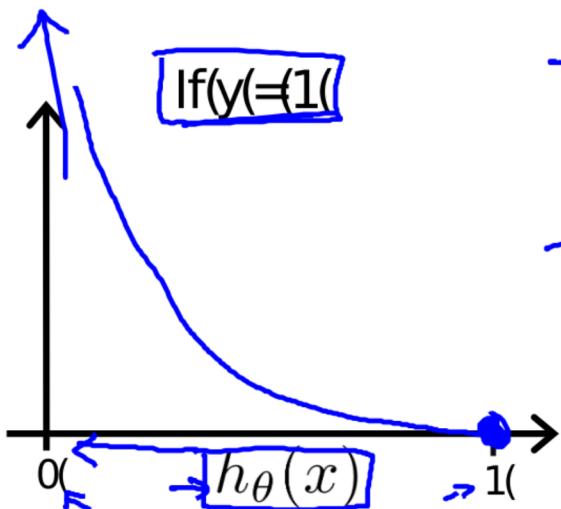
$He^{-\alpha t}$



Andrew(Ng)

Logistic regression cost function +

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

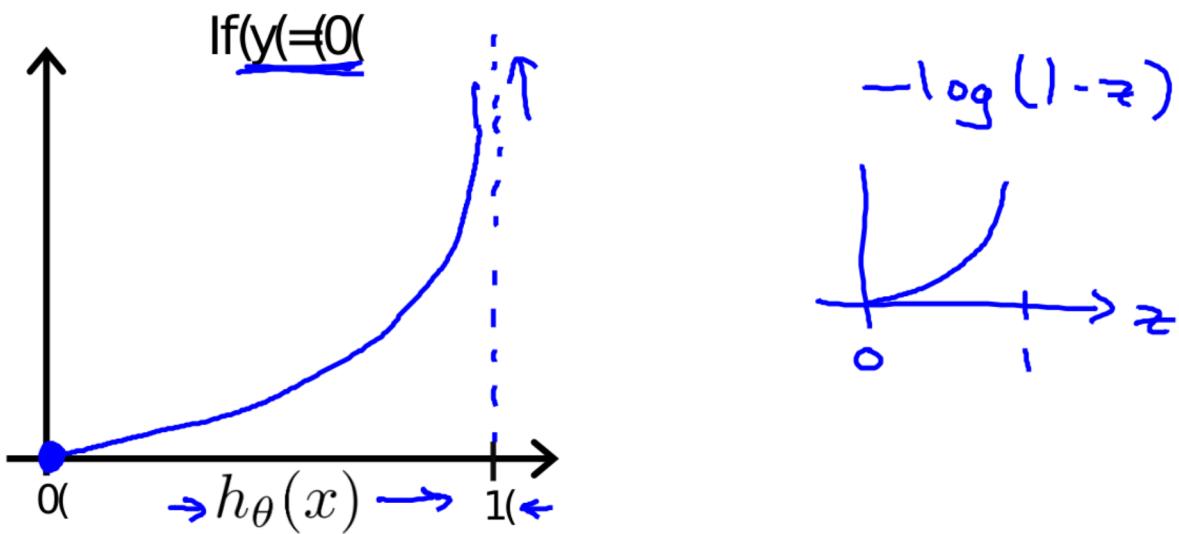


→ Cost = 0 if $y = 1, h_{\theta}(x) = 1$
 But as $h_{\theta}(x) \rightarrow 0$
 $\underline{\text{Cost}} \rightarrow \infty$

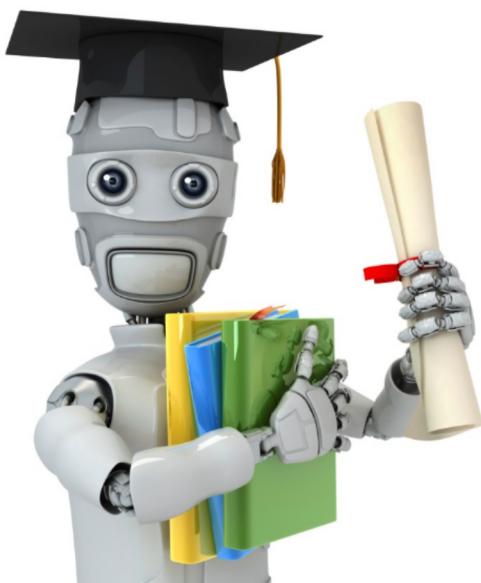
→ Captures intuition that if $h_{\theta}(x) = 0$,
 (predict $P(y = 1|x; \theta) = 0$), but $y = 1$
 we'll penalize learning algorithm by a very
 large cost.

Logistic regression cost function +

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



Andrew(Ng)



Machine Learning

Logistic Regression

Simplified cost function and gradient descent

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

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

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

$$\text{If } y=0: \text{Cost}(h_\theta(x), y) = -\log(1 - h_\theta(x))$$

Logistic regression cost function

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

To fit parameters

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

Credit Θ

To make a prediction given new

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

$$p(y=1 | x; \theta)$$

Gradient Descent +

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

Want multiple parallel gradient updates:

Repeat {

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

}

Simultaneously update all of them

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

Gradient Descent +

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

Want to minimize $J(\theta)$:

Repeat {

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

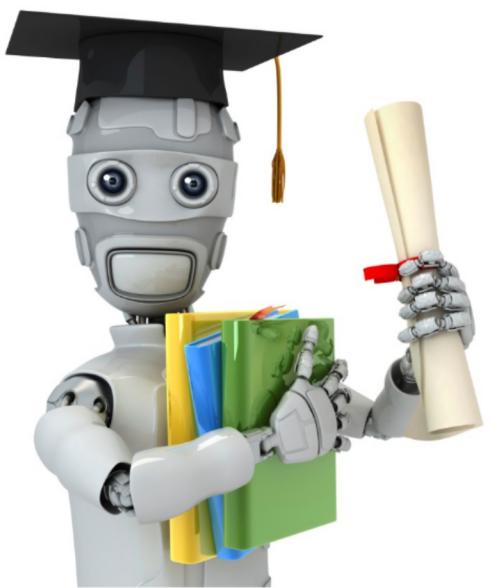
~~Simultaneously update all~~

$$\Theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_n \end{bmatrix} \quad \text{for } i=0 \text{ to } n$$

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

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

Algorithm looks identical to linear regression!



Machine Learning

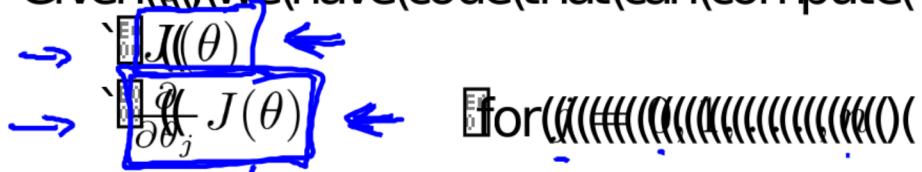
Logistic Regression

Advanced optimization

Optimization algorithm+

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

Given $J(\theta)$, we have code that can compute



Gradient descent:

Repeat {

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

}

Optimization algorithm+

Given $J(\theta)$, we have code that can compute

$$\begin{cases} J(\theta) \\ \frac{\partial}{\partial \theta_j} J(\theta) \end{cases}$$

for θ

Optimization algorithms:

- \square Gradient descent
- \square Conjugate gradient
- \square BFGS
- \square LBFGS

Advantages:

- \square No need to manually pick α
- \square Often faster than gradient descent.

Disadvantages:

- \square More complex

Example: (min $J(\theta)$)

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

```
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);
```

```
→ options = optimset('GradObj', 'on', 'MaxIter', '100');
→ initialTheta = zeros(2,1);
[optTheta, functionVal, exitFlag] ...  
= fminunc(@costFunction, initialTheta, options);
          ↑           ↑
          θ ∈ ℝd   d ≥ 2.
```

Andrew(Ng)

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

function [jVal, gradient] = costFunction(theta)

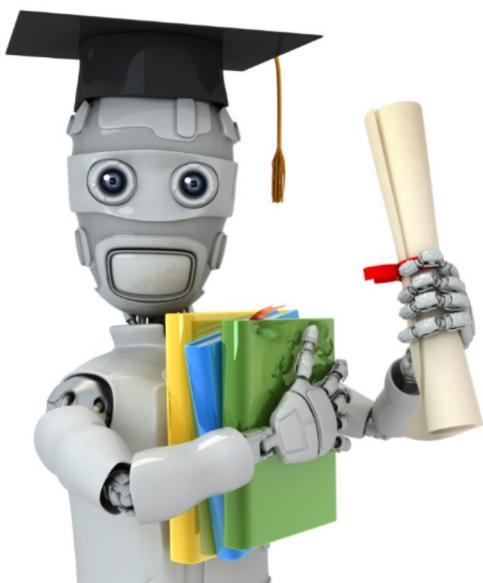
jVal = [code to compute $J(\theta)$];

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 $\frac{\partial}{\partial \theta_n} J(\theta)$];



Machine Learning

Logistic Regression

Multiclass classification: One-vs-all

Multiclass Classification +

Email (foldering/tagging): (Work, (Friends, (Family, (Hobby)

$$y=1 \quad y=2 \quad y=3 \quad y=4$$

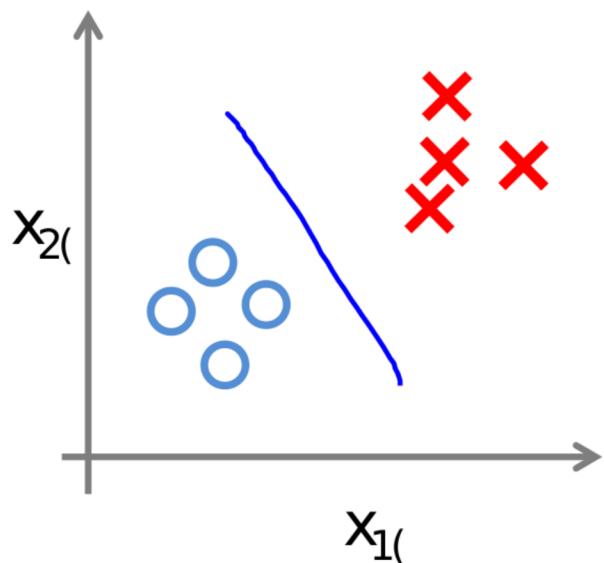
Medical (diagnoses): (Not ill, (Cold, (Flu)

$$y=1 \quad 2 \quad 3$$

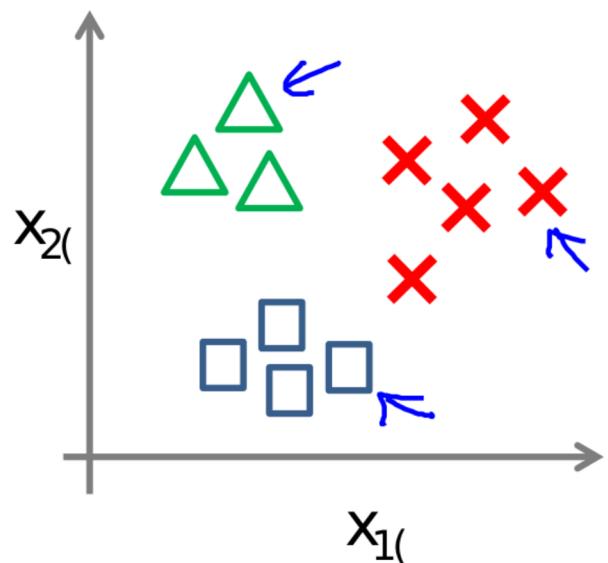
Weather: (Sunny, (Cloudy, (Rain, (Snow)

$$y=1 \quad 2 \quad 3 \quad 4 \quad \leftarrow$$

Binary(classification):

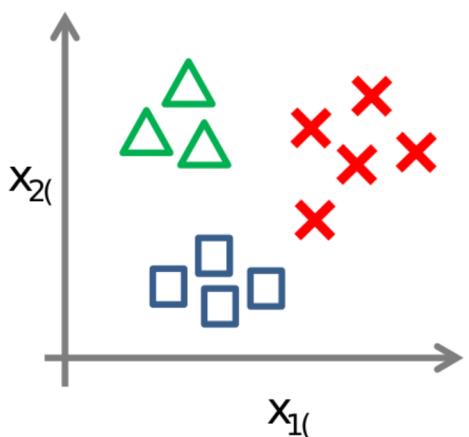


Mul[iclass(classification):



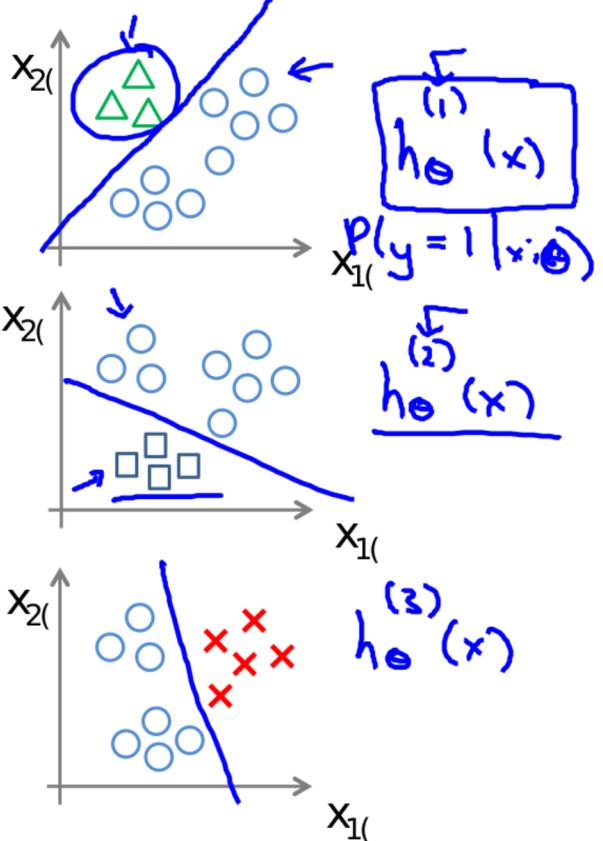
Andrew(Ng)

One-vs-all (+one-vs-rest): +



Class(1): (▲ ←
 Class(2): (□ ←
 Class(3): (✗ ←

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



Andrew(Ng)

OneVsAll +

Train(a(logistic regression(classifier)(⁽ⁱ⁾)(for(each(class(to)(predict(the(probability(that(maximizes.(

On(a(new(input(to)(make(a(prediction,(pick(the(class(that)(maximizes(

$$\max_i \underbrace{h_{\theta}^{(i)}(x)}_{\uparrow}$$