

# ANDREW NG COURSE

## • Introduction

Machine Learning gives computers the ability to learn without being explicitly programmed. (Arthur Samuel)

A computer program is said to learn from experience  $E$  with respect to some class of tasks  $T$  and performance measure  $P$ , if its performance at tasks in  $T$ , as measured by  $P$ , improves with experience  $E$

Supervised Learning: Divides into regression and classification.

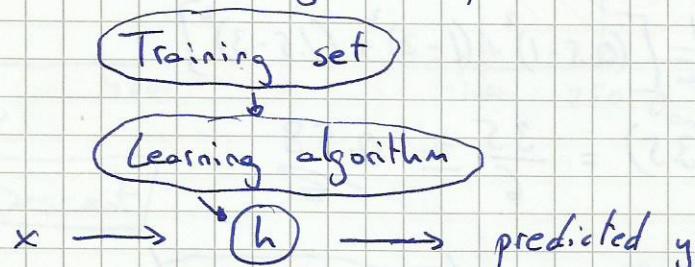
Regression  $\Rightarrow$  continuous output

Classification  $\Rightarrow$  discrete output

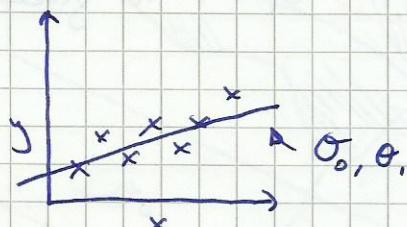
Unsupervised Learning: Allows to approach problems with little or no idea what results should look like

## • Model Representation

$x^{(i)}$  = input variables       $y^{(i)}$  = output variables

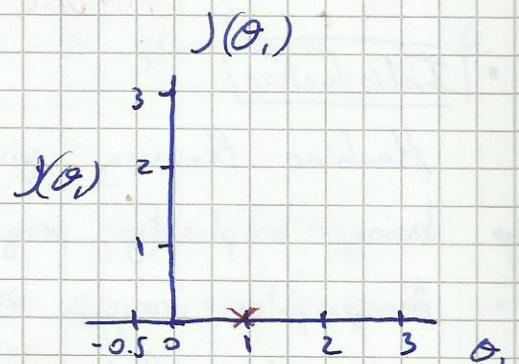
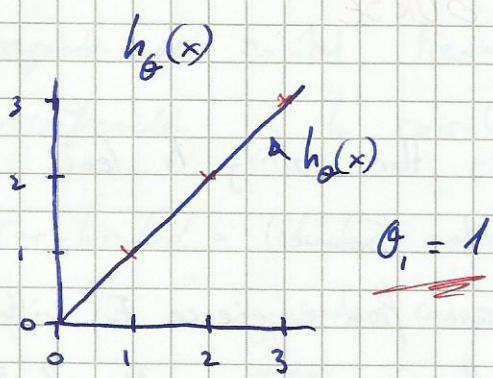


Cost function. Squared error function or Mean squared error.



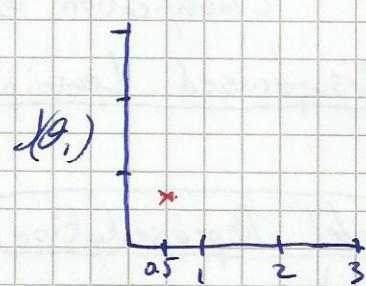
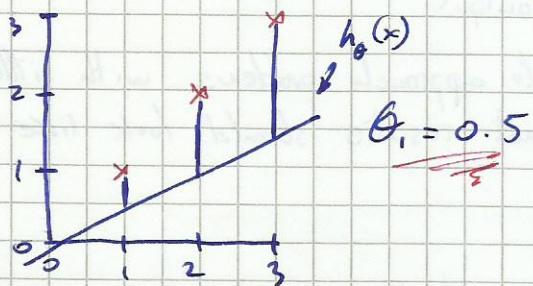
$$\text{minimize}_{\theta_0, \theta_1} \frac{1}{2n} \sum_{i=1}^n (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

Ide: Choose  $\theta_0, \theta_1$  so that  $h_{\theta}(x)$  is close to  $y$  for our training examples  $(x, y)$



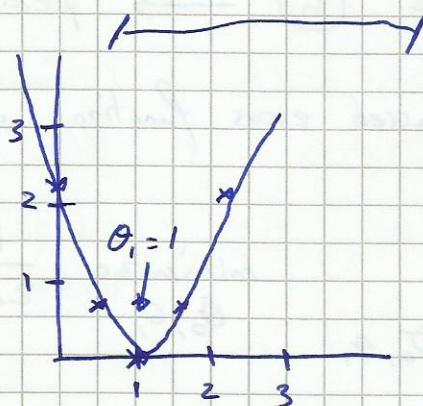
$$J(\theta_0) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta_0}(x^{(i)}) - y^{(i)})^2 = \frac{1}{2m} (0^2 + 0^2 + 0^2) = 0$$

$$\theta_0 = 0.5$$

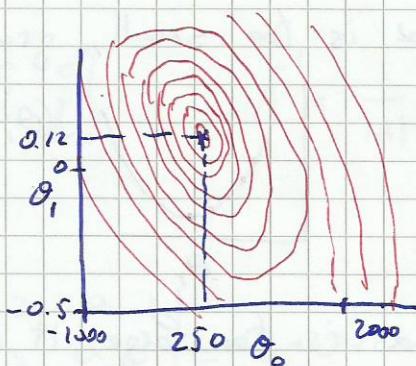
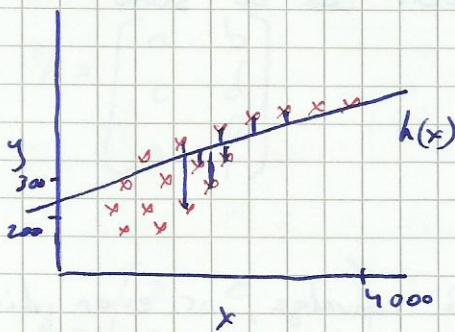
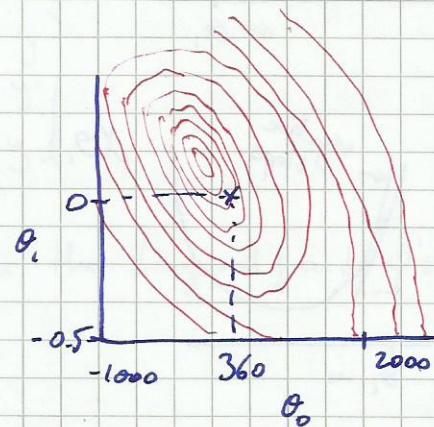
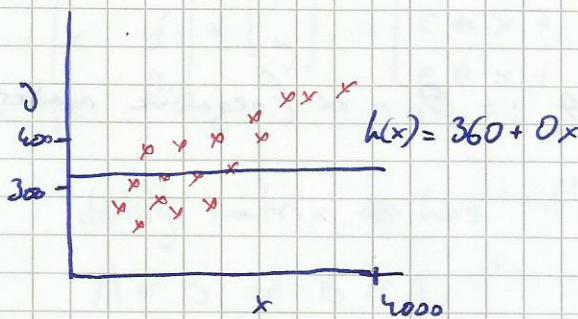
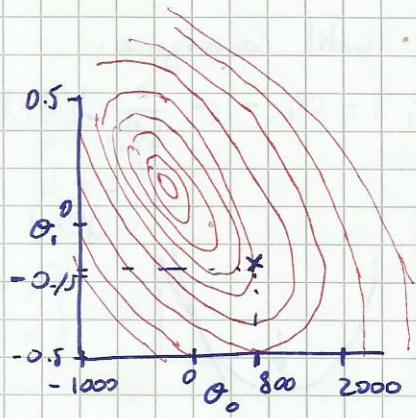
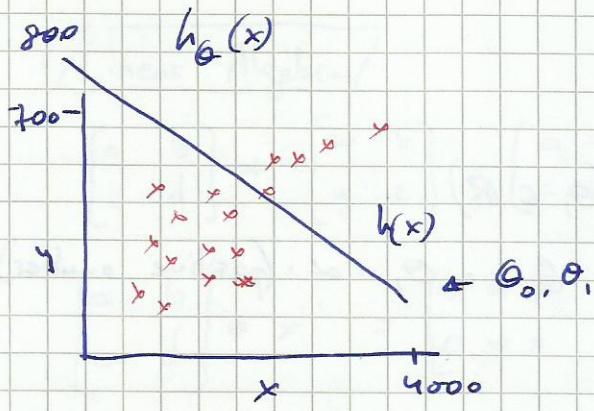


$$J(0.5) = \frac{1}{2m} [(0.5-1)^2 + (1-2)^2 + (1.5-3)^2]$$

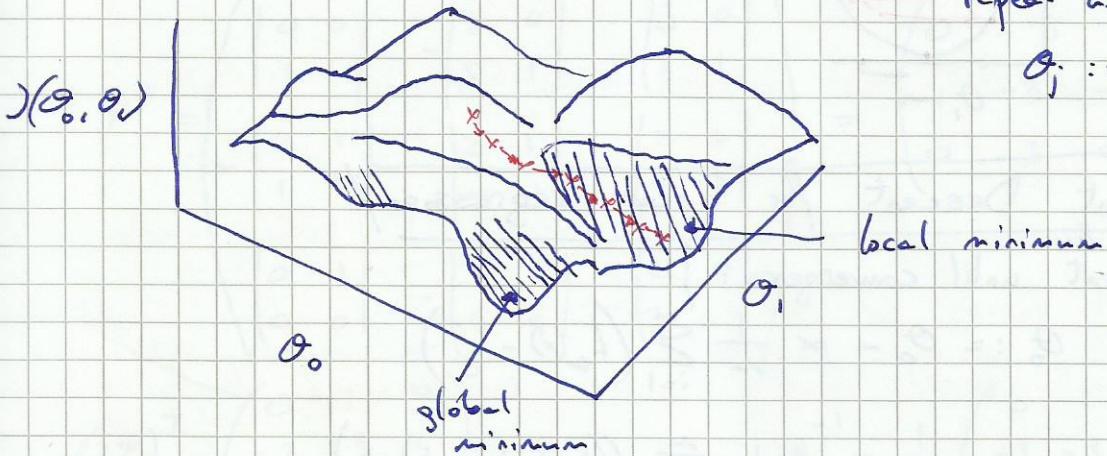
$$\frac{1}{2m} (3.5) = \frac{3.5}{6} \approx 0.58$$



$\theta_0 = 1$  is global minimum



### Gradient Descent



repeat until convergence:

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1)$$

$$\text{temp}_0 := \theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1)$$

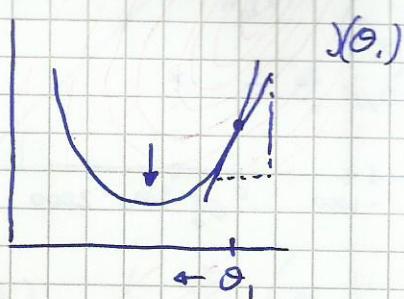
$$\text{temp}_1 := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1)$$

$$\theta_0 := \text{temp}_0$$

$$\theta_1 := \text{temp}_1$$

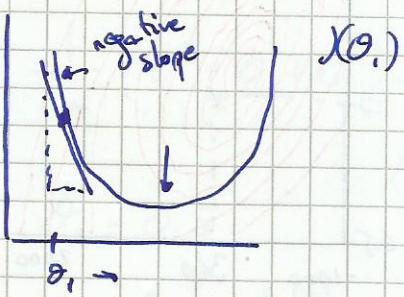
Repeat until convergence:

$$\theta_1 := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_1)$$



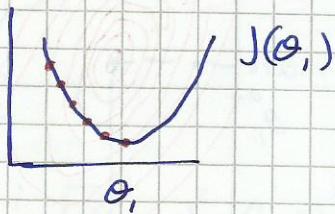
$$(\theta_1 \in \mathbb{R})$$

$$\theta_1 := \theta_1 - \alpha \cdot (\text{positive number})$$

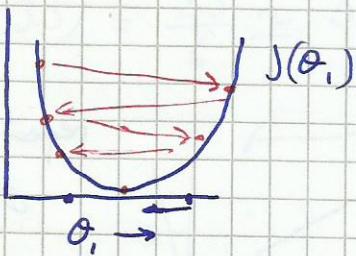


$$\theta_1 := \theta_1 - \alpha \cdot (\text{negative number})$$

If  $\alpha$  is too small, gradient descent can be slow



If  $\alpha$  is too large, it may fail to converge, or even diverge



### Gradient Descent for linear regression

repeat until convergence:

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

$$\theta_1 := \theta_1 - \alpha \sum_{i=1}^m ((h_\theta(x^i) - y^i)x^i)$$

## Linear Algebra

$$\begin{bmatrix} a & b \\ c & d \end{bmatrix} + \begin{bmatrix} w & x \\ y & z \end{bmatrix} = \begin{bmatrix} a+w & b+x \\ c+y & d+z \end{bmatrix}$$

$$\begin{bmatrix} a & b \\ c & d \end{bmatrix} * x = \begin{bmatrix} a*x & b*x \\ c*x & d*x \end{bmatrix}$$

$$\begin{bmatrix} a & b \\ c & d \\ e & f \end{bmatrix} * \begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} a*x + b*y \\ c*x + d*y \\ e*x + f*y \end{bmatrix}$$

Identity matrix has 1's on the diagonal and 0's everywhere

$$A * B \neq B * A$$

$$(A * B) * C = A * (B * C)$$

$$A = \begin{bmatrix} a & b \\ c & d \\ e & f \end{bmatrix}$$

$$A^T = \begin{bmatrix} a & c & e \\ b & d & f \end{bmatrix}$$

$$A^{-1} = \frac{1}{|A|} (A^*)^T$$

$$A = \begin{pmatrix} 2 & 0 & 1 \\ 3 & 0 & 0 \\ 5 & 1 & 1 \end{pmatrix}$$

$$A = \begin{vmatrix} 2 & 0 & 1 \\ 3 & 0 & 0 \\ 5 & 1 & 1 \end{vmatrix} = 3 + 0 + 0 - (0 + 0 + 0) = 3$$

$$A^* = \begin{pmatrix} \begin{vmatrix} 0 & 0 \\ 1 & 1 \end{vmatrix} & -\begin{vmatrix} 3 & 0 \\ 5 & 1 \end{vmatrix} & \begin{vmatrix} 3 & 0 \\ 5 & 1 \end{vmatrix} \\ -\begin{vmatrix} 0 & 1 \\ 1 & 1 \end{vmatrix} & \begin{vmatrix} 2 & 1 \\ 5 & 1 \end{vmatrix} & -\begin{vmatrix} 2 & 0 \\ 5 & 1 \end{vmatrix} \\ \begin{vmatrix} 0 & 1 \\ 0 & 0 \end{vmatrix} & -\begin{vmatrix} 2 & 1 \\ 3 & 0 \end{vmatrix} & \begin{vmatrix} 2 & 0 \\ 3 & 0 \end{vmatrix} \end{pmatrix} = \begin{pmatrix} 0 & -3 & 3 \\ 1 & -3 & -2 \\ 0 & 3 & 0 \end{pmatrix}$$

$$(A^*)^T = \begin{pmatrix} 0 & 1 & 0 \\ -3 & -3 & 3 \\ 3 & -2 & 0 \end{pmatrix}$$

$$A^{-1} = \frac{1}{3} \begin{pmatrix} 0 & 1 & 0 \\ -3 & -3 & 3 \\ 3 & -2 & 0 \end{pmatrix} = \begin{pmatrix} 0 & \frac{1}{3} & 0 \\ -1 & -1 & 1 \\ 1 & -\frac{2}{3} & 0 \end{pmatrix}$$

## • Multiple Features

Linear regression with multiple variables is also known as "multivariate linear regression"

$x_j^{(i)}$  = value of feature  $j$  in the  $i^{\text{th}}$  training example.

$\mathbf{x}^{(i)}$  = the input (features) of the  $i^{\text{th}}$  training example

$m$  = the number of training examples

$n$  = the number of features

The multivariable form of the hypothesis function:

$$h_{\theta}(\mathbf{x}) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_m x_n \quad \text{is equivalent to}$$

$$h_{\theta}(\mathbf{x}) = [\theta_0 \dots \theta_n] \begin{bmatrix} x_0 \\ \vdots \\ x_n \end{bmatrix} = \theta^T \mathbf{x}$$

## • Gradient Descent for Multiple Variables

repeat until convergence:

$$\theta_j := \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(\mathbf{x}^{(i)}) - y^{(i)}) \cdot x_{(j)}^{(i)} \quad \text{for } j := 0 \dots n$$

}

The difference between multiple variables and 1 variable is that:

- With 1 variable: ( $n=1$ )

$\theta_0$  hasn't  $x_0$  /  $\theta_1$  has  $x_1$

- With multiple variables: ( $n \geq 1$ )

$$x_0^{(i)} = 1$$

We can speed up gradient descent by having each of our input values in roughly the same range. This is because  $\theta$  will descend quickly on small ranges and slowly on large ranges. And so will oscillate inefficiently down to the optimum when the variables are very uneven.

To prevent this is to modify ranges of our input variables

$$-1 \leq x_{(i)} \leq 1 \quad \text{or} \quad -0.5 \leq x_{(i)} \leq 0.5$$

2 techniques to help are

- Feature scaling:  $x_{(i)} / \max(x_{(i)}) - \min(x_{(i)})$
- Mean normalization:  $(x_{(i)} - \text{average}(:)) / \underbrace{\max(x_{(i)}) - \min(x_{(i)})}_{\text{Can be replaced by standard deviation}}$

• Debugging Gradient Descent. Make a plot with number of iterations on the x-axis.  $J(\theta)$  over the number of iterations of gradient descent.

If  $J(\theta)$  ever increases, you need to decrease  $\alpha$ .

Declare convergence if  $J(\theta)$  decreases by less than  $10^{-3}$  in one iteration.

- For sufficiently small  $\alpha$ ,  $J(\theta)$  should decrease on every iteration
- If  $\alpha$  is too small, gradient descent can be slow to converge.

• Polynomial regression

We can combine multiple features into one. Our hypothesis function needn't be linear if that doesn't fit the data well.

To make it a square root function, we could do:

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

• Normal Equation

We will minimize  $J$  by explicitly taking its derivatives with respect to the  $\theta_j$ 's, and setting them to zero. This allows us to find the optimum theta without iteration.

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

There is no need to do feature scaling with the normal equation.

Gradient Descent

Need to choose  $\alpha$

Need many iterations

$O(Kn^2)$

With  $n$  is large

Normal equation

No need to choose  $\alpha$

No need to iterate

$O(n^3)$ , need calculate  $(X^T X)^{-1}$

With  $n$  is short

If  $X^T X$  is non invertible, the common causes might be having:

- Redundant features: 2 features are very closely related (linearly dependent)
- Too many features ( $m \leq n$ ): Delete some features or use regularization

## Classification

We want to predict take on only a small number of discrete values.

For binary classification problems,  $y \in \{0, 1\}$

### + Hypothesis representation

It doesn't make sense for  $h_\theta(x)$  to take values larger than 1 or smaller than 0 when we know that  $y \in \{0, 1\}$ .

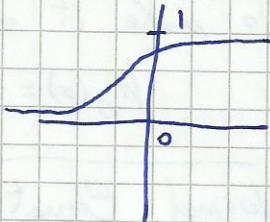
To fix it,  $0 \leq h_\theta(x) \leq 1$ .

Our new form uses the "Sigmoid Function" or also called the "Logistic Function"

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

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

$$z = \theta^T x$$



$h_\theta(x)$  give us the probability that our output is 1.

$$h_\theta(x) = P(y=1/x; \theta) \quad P(y=0/x; \theta) = 1 - P(y=1/x; \theta)$$

### + Decision Boundary

$$h_\theta(x) \geq 0.5 \rightarrow y = 1$$

$$h_\theta(x) < 0.5 \rightarrow y = 0$$

$$g(z) \geq 0.5 \quad \text{when } z \geq 0$$

\* Reminder

$$z=0, e^0=1 \Rightarrow g(0)=0.5$$

$$z \rightarrow \infty, e^{-\infty} \rightarrow 0 \Rightarrow g(\infty)=1$$

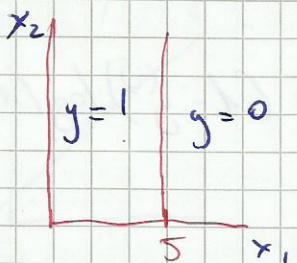
$$z \rightarrow -\infty, e^\infty \rightarrow \infty \Rightarrow g(-\infty)=0$$

If  $g = \theta^T x$ , then:  $h_{\theta}(x) = g(\theta^T x) \geq 0.5$   
when  $\theta^T x \geq 0$

From these statements, we can say:

$$\theta^T x \geq 0 \Rightarrow y = 1 \quad \text{and} \quad \theta^T x < 0 \Rightarrow y = 0$$

Example:  $\theta = \begin{bmatrix} 5 \\ -1 \\ 0 \end{bmatrix}$   $y = 1 \text{ if } 5 + (-1)x_1 + 0x_2 \geq 0$



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

The input to the sigmoid function  $g(z)$  doesn't need to be linear, and could be a function that describes a circle

### Cost Function

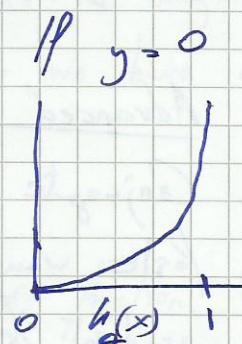
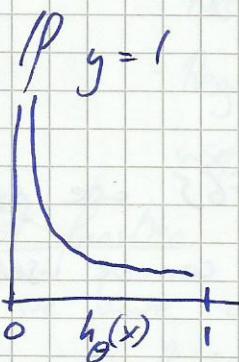
We can't use the same cost function that we use for linear regression because the Logistic Function will cause the output to be wavy. In other words, it will not be a convex function. To fix this: ↓

Our cost function looks like:

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

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

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



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

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

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

## Simplified Cost Function & Gradient Descent

We can compress our cost function's 2 conditional cases into one case:

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



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

A vectorized implementation is:

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

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

## + Gradient Descent

We can work out the derivate part using calculus to get:

Repeat {

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

}

this algorithm is identical to used in linear regression.

We still have to update all values in theta.

A vectorized implementation is:

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

## + Advanced Optimization

Conjugate gradient, BFGS, L-BFGS are more complex, faster ways to optimize  $\theta$  that can be used instead of gradient descent.

## Multiclassification: One vs all

Now  $y \in \{1, \dots, n\}$ . We divide our problem in  $n$  parts.

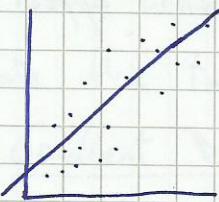
Train a logistic regression classifier  $h_{\theta}(x)$  for each class to predict

the probability that  $y = i$ . To make a prediction on a new  $x$ , pick the class that maximize  $h_\theta(x)$

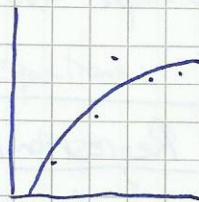
- Solving the Problem of Overfitting

It might seem that the more features we add, better. However, there is also a danger in adding too many features.

Underfitting



Correct



Overfitting



Underfitting, or high bias. It's usually caused by a function that is too simple or uses too few features.

Overfitting, or high variance. It's caused by doesn't generalize well to predict new data. It's usually caused by a complicated function that creates a lot of unnecessary curves and angles unrelated to the data.

there are two main options to address the issue of overfitting:

- 1) Reduce the number of features

- Manually select which features to keep
- Use a model selection algorithm

- 2) Regularization

- Keep all the features, but reduce the magnitude of parameters
- Regularization works well when we have a lot of slightly useful features.

- Cost function

If we have overfitting from our hypothesis function. We can modify our cost function

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

We could also regularize all of our theta parameters in a single summation as:

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

The  $\lambda$ , or lambda, is the regularization parameter. If  $\lambda$  is too large, it may smooth out the function too much and cause underfitting. If  $\lambda = 0$  or too small, doesn't affect to the function.

### • Regularized Linear Regression /

+ Gradient Descent. We don't want to penalize  $\theta_0$ .

Repeat {

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

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

$$j \in \{1, 2, \dots, n\}$$

The term  $\frac{\lambda}{m} \theta_j$  performs our regularization. With some manipulation our update rule can also be represented as:

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

$1 - \alpha \frac{\lambda}{m}$  will always be less than 1. You can see it as reducing the value of  $\theta_j$ , by some amount on every update.

### + Normal equation

$$\theta = (X^T X + \lambda L)^{-1} X^T y \quad \text{where } L = \begin{bmatrix} 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}_{(n+1) \times (n+1)}$$

Recall that if  $m < n$ , then  $X^T X$  is non-invertible.

However, when we add the term  $\lambda L$ , then  $X^T X + \lambda L$  becomes invertible.

## Regularized Logistic Regression

We can avoid overfitting, we can regularize the cost function by adding a term to the end:

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

The second sum,  $\sum_{j=1}^n \theta_j^2$  means to explicitly exclude the bias term,  $\theta_0$ . Thus, when computing the equation, we should continuously update the two following equations:

### Gradient Descent

Repeat {

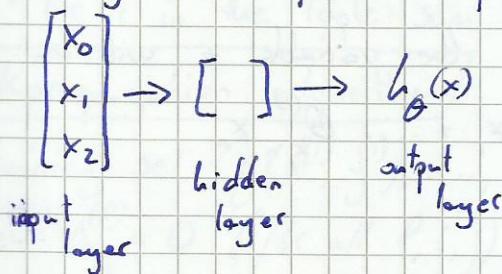
$$\theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_0^{(i)}$$

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

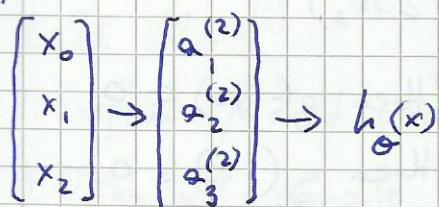
}

## Neural Networks

Visually, a simplistic representation looks like:



If we have one hidden layer:



$$a_1^{(2)} = g(\theta_{10}^{(1)} x_0 + \theta_{11}^{(1)} x_1 + \theta_{12}^{(1)} x_2 + \theta_{13}^{(1)} x_3)$$

$$a_2^{(2)} = g(\theta_{20}^{(1)} x_0 + \theta_{21}^{(1)} x_1 + \theta_{22}^{(1)} x_2 + \theta_{23}^{(1)} x_3)$$

$$a_3^{(2)} = g(\theta_{30}^{(1)} x_0 + \theta_{31}^{(1)} x_1 + \theta_{32}^{(1)} x_2 + \theta_{33}^{(1)} x_3)$$

$$h_\theta(x) = a_1^{(2)} = g(\theta_{10}^{(2)} a_0^{(2)} + \theta_{11}^{(2)} a_1^{(2)} + \theta_{12}^{(2)} a_2^{(2)} + \theta_{13}^{(2)} a_3^{(2)})$$

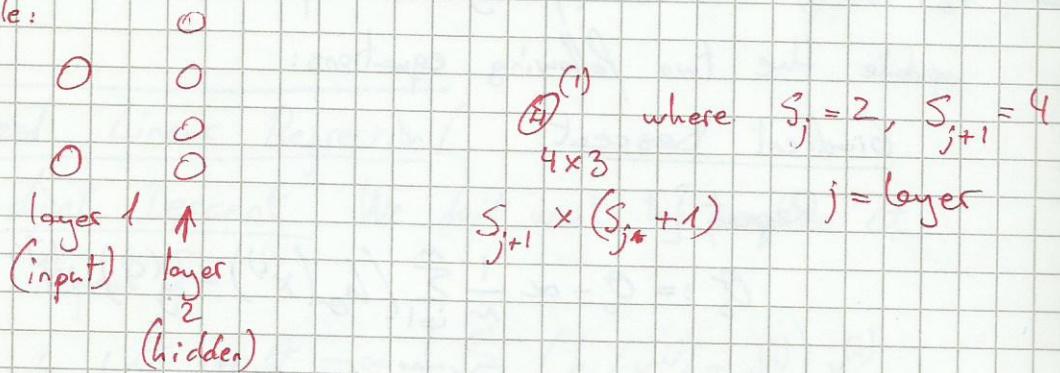
$a_i^{(j)}$  = "activation" of unit i  
in layer j

$\Theta^{(j)}$  = matrix of weights  
controlling function mapping  
from layer j to layer j+1

This is saying that we compute our activation nodes by using a  $3 \times 4$  matrix of parameters. Each layer gets its own matrix of weights,  ~~$\Theta^{(j)}$~~

If network has  $s_j$  units in layer  $j$  and  $s_{j+1}$  units in layer  $j+1$ , then  $\Theta^{(j)}$  will be of dimension  $s_{j+1} \times (s_j + 1)$ . The  $+1$  comes from the "bias nodes".

Example:



We'll do a vectorized implementation of the above functions.

We're going to define a new variable  $z_k^{(j)}$  that encompasses the parameters inside our  $g$  function.

$$a_1^{(2)} = g(z_1^{(2)}) ; \quad a_2^{(2)} = g(z_2^{(2)}) ; \quad a_3^{(2)} = g(z_3^{(2)})$$

For layer  $j=2$  and node  $k$ , the variable  $z$  will be:

$$z_k^{(2)} = \Theta_{k,0}^{(1)} x_0 + \Theta_{k,1}^{(1)} x_1 + \dots + \Theta_{k,n}^{(1)} x_n$$

• Examples

$$\Theta^{(1)} = [-30 \quad 20 \quad 20]$$

$$h_{\Theta}(x) = g(-30 + 20x_1 + 20x_2)$$

$$x_1 = 0 \text{ AND } x_2 = 0 \text{ then } g(-30) \approx 0$$

$$x_1 = 0 \text{ AND } x_2 = 1 \text{ then } g(-10) \approx 0$$

$$x_1 = 1 \text{ AND } x_2 = 0 \text{ then } g(-10) \approx 0$$

$$x_1 = 1 \text{ AND } x_2 = 1 \text{ then } g(10) \approx 1$$

## Cost Function in neural networks /

$L$  = layers in the network

$S_p$  = Number of units in layer  $p$

$K$  = Number of output units/classes

The cost function for regularized logistic regression was:

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

For neural networks: cost function

$$J(\theta) = -\frac{1}{m} \sum_{i=1}^m \sum_{k=1}^K \left[ y_k^{(i)} \log((h_\theta(x^{(i)}))_k) + (1-y_k^{(i)}) \log(1-(h_\theta(x^{(i)}))_k) \right] + \underbrace{\frac{\lambda}{2m} \sum_{l=1}^{L-1} \sum_{i=1}^{S_l} \sum_{j=1}^{S_{l+1}} (\theta_{j,i}^{(l)})^2}_{\text{cost function}}$$

- the double sum simply adds up the logistic regression costs calculated for each cell in the output layer.
- The triple sum simply adds up the squares of all the individual  $\theta$ 's in the entire network
- the  $i$  in the triple sum doesn't refer to training example  $i$

## Backpropagation algorithm /

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

Set  $\Delta_{ij}^{(l)} = 0$  (for all  $l, i, j$ )

\*Used to compute  $\frac{\partial}{\partial \theta_{ij}^{(l)}} J(\theta)$

For  $i=1$  to  $m$

1. Set  $a^{(1)} = x^{(i)}$

2. Perform forward propagation to compute  $a^{(l)}$  for  $l=2, 3, \dots, L$

3. Using  $y^{(i)}$ , compute  $\delta^{(L)} = a^{(L)} - y^{(i)}$

4. Compute  $\delta^{(L-1)}, \delta^{(L-2)}, \dots, \delta^{(2)}$

5.  $\Delta_{ij}^{(l)} := \Delta_{ij}^{(l)} + a_j^{(l)} \delta_i^{(l+1)}$



$$\Delta_{ij}^{(l)} := \frac{1}{m} \Delta_{ij}^{(l)} + \lambda \theta_{ij}^{(l)} \quad \text{if } j \neq 0$$

$$\Delta_{ij}^{(l)} := \frac{1}{m} \Delta_{ij}^{(l)} \quad \text{if } j = 0$$

$$\frac{\partial}{\partial \theta_{ij}^{(l)}} J(\theta) = \Delta_{ij}^{(l)}$$

2. Given one training example  $(x, y)$ :

Forward propagation:

$$a^{(1)} = x$$

$$z^{(2)} = \Theta^{(1)} a^{(1)}$$

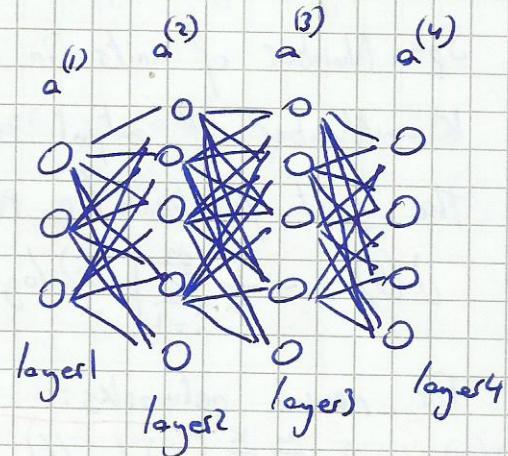
$$a^{(2)} = g(z^{(2)}) \quad (\text{add } a_0)$$

$$z^{(3)} = \Theta^{(2)} a^{(2)}$$

$$a^{(3)} = g(z^{(3)}) \quad (\text{add } a_0)$$

$$z^{(4)} = \Theta^{(3)} a^{(3)}$$

$$a^{(4)} = h_{\Theta}(x) = g(z^{(4)})$$



3. Using  $y^{(t)}$ , compute  $\delta^{(l)} = a^{(l)} - y^{(t)}$

$L$  = Total number of layers

$a^{(l)}$  = Vector of outputs of the activation units for the last layer.

Our "error values" for the last layer are simply the differences of our actual results in the last layer and the correct outputs in  $y$ .

4. Compute  $\delta^{(L-1)}, \delta^{(L-2)}, \dots, \delta^{(2)}$  using:

$$\delta^{(p)} = ((\Theta^{(p)})^T \delta^{(p+1)}) \cdot \underbrace{a^{(p)} * (1 - a^{(p)})}_{g'(z^{(p)})}$$

5. A vectorized form is  $\Delta^{(p)} := \Delta^{(p)} + \delta^{(p+1)} / a^{(p)T}$

### Unrolling Parameters

#### Learning algorithm

Have initial parameters  $\Theta^{(1)}, \Theta^{(2)}, \Theta^{(3)}$

Unroll to get initialTheta to pass to

`fminuc(@costFunction, initialTheta, options)`

function [Jval, gradientVec] = costFunction(thetaVec)

From thetaVec, get  $\Theta^{(1)}, \Theta^{(2)}, \Theta^{(3)}$

Use forward prop/back prop to compute  $\Delta^{(1)}, \Delta^{(2)}, \Delta^{(3)}$  and  $J$

Unroll  $\Delta^{(1)}, \Delta^{(2)}, \Delta^{(3)}$  to get gradientVec

## • Gradient Checking

We can approximate the derivative of our cost function with:

$$\frac{\partial}{\partial \theta_j} J(\theta) \approx \frac{J(\theta + \epsilon) - J(\theta - \epsilon)}{2\epsilon}$$

A small value for  $\epsilon$  such as  $\epsilon = 10^{-4}$ , guarantees that the math works out properly. If the value for  $\epsilon$  is too small, we can end up with numerical problems.

## • Random Initialization

Initializing all theta weights to zero doesn't work with neural networks.

### Symmetry breaking

Initialize each  $\theta_{ij}^{(l)}$  to a random value in  $[-\epsilon, \epsilon]$

## • Putting it together

- Number of input units = dimension of features  $x^{(i)}$
- Number of output units = number of classes
- Number of hidden units per layer = usually more the better
- Defaults: 1 hidden layer. If you have more than 1 hidden layer, you have the same number of units in every hidden layer

## Training a Neural Network

1. Randomly initialize the weights
2. Implement forward propagation to get  $h_\theta(x^{(i)})$  for any  $x^{(i)}$
3. Implement the cost function
4. Implement backpropagation to compute partial derivatives
5. Use gradient checking to confirm that your backpropagation works. Then disable gradient checking.
6. Use gradient descent or a built-in optimization function to minimize the cost function with the weights in theta

## Evaluating a hypothesis

Once we have done some trouble shooting for errors in our predictions by:

- Getting more training examples
- Trying smaller sets of features
- Trying additional features
- Trying polynomial features
- Increasing or decreasing  $\lambda$

We can move on to evaluate our new hypothesis

Training set: 70%      Test set: 30%

The new procedure using these two sets is then:

1. Learn  $\Theta$  and minimize  $J_{\text{train}}(\Theta)$  using the training set
2. Compute the ~~sets~~ test set error  $J_{\text{test}}(\Theta)$   
the test set error

1. For linear regression:  $J_{\text{test}}(\Theta) = \frac{1}{m_{\text{test}}} \sum_{i=1}^{m_{\text{test}}} (h_{\Theta}(x^{(i)}_{\text{test}}) - y^{(i)})^2$

2. For classification ~ Misclassification

$$\text{err}(h_{\Theta}(x), y) = \begin{cases} 1 & \text{if } h_{\Theta}(x) \geq 0.5 \text{ \& } y = 0 \\ & \text{or} \\ & h_{\Theta}(x) < 0.5 \text{ \& } y = 1 \\ 0 & \text{otherwise} \end{cases}$$

the average ~~sets~~ test error for the test set is

$$\text{Test Error} = \frac{1}{m_{\text{test}}} \sum_{i=1}^{m_{\text{test}}} \text{err}(h_{\Theta}(x^{(i)}_{\text{test}}), y^{(i)}_{\text{test}})$$

One way to break down our dataset into the ~~three~~ 3 sets:

- Training set 60%
- Cross validation set 20%
- Test set 20%

We can now calculate three separate error values for the 3 different sets using the following method:

1. Optimize the parameters in  $\Theta$  using the training set for each polynomial degree.
2. Find the polynomial degree  $d_*$  with the least error using the cross validation set
3. Estimate the generalization error using the test set with  $J_{\text{test}}(\Theta^{(d)})$  ( $d = \text{theta from polynomial with lower error}$ )

The degree of the polynomial  $d_*$  hasn't been trained using the test set.

### Diagnosing Bias vs Variance

High bias  $\rightarrow$  Underfitting

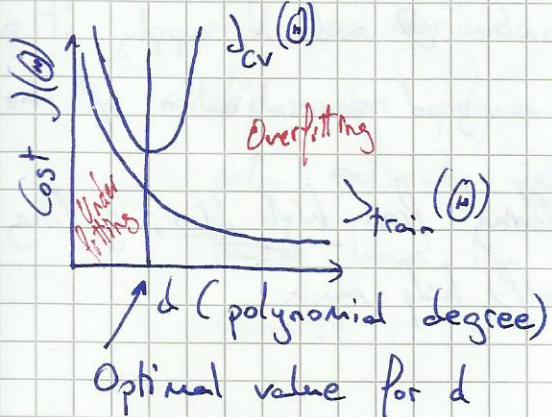
High variance  $\rightarrow$  Overfitting

The training error will tend to decrease as we increase the degree  $d$  of the polynomial.

At the same time, the cross validation error will tend to decrease as we increase  $d$  up to a point, and then it will increase as  $d$  is increased, forming a convex curve.

High bias (underfitting):  $J_{\text{cv}}(\Theta)$  and  $J_{\text{train}}(\Theta)$  will be high

High variance (overfitting):  $J_{\text{train}}(\Theta)$  will be low and  $J_{\text{cv}}(\Theta)$  will be much greater than  $J_{\text{train}}(\Theta)$



## • /Regularization and Bias/Variance/

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

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

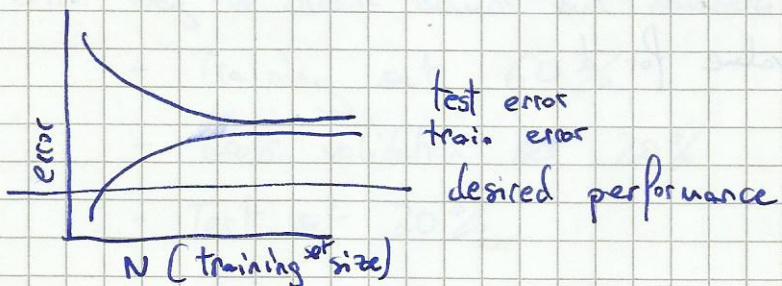
Large  $\lambda$  = high bias (underfit)

Small  $\lambda$  = high variance (overfit)

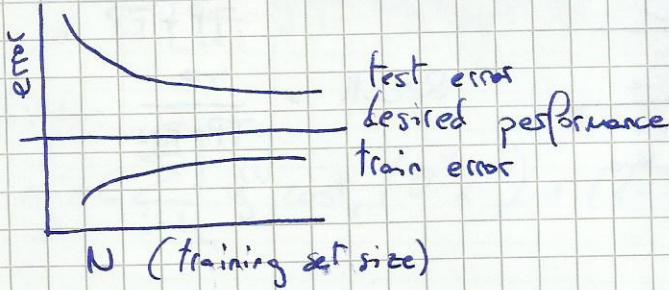
To choose the model and the regularization term  $\lambda$ , we need to:

1. Create a list of lambdas  
(i.e.  $\lambda \in \{0, 0.01, 0.02, 0.04, 0.08, 0.16, 0.32, 0.64, 1.28, 2.56, 5.12, 10.24\}$ )
2. Create a set of models with different degrees or any other variants
3. Iterate through the  $\lambda$ s and for each  $\lambda$  go through all the models to learn some  $\theta$ .
4. Compute the cross validation error using the learned  $\theta$  (computed with  $\lambda$ ) on the  $J_{cv}(\theta)$  without regularization or  $\lambda=0$
5. Select the best combo that produces the lowest error on the cross validation set
6. Using the best combo  $\theta$  and  $\lambda$ , apply it on  $J_{test}(\theta)$  to see if it has a good generalization of the problem

If a learning algorithm is suffering from high bias, getting more training data will not (by itself) help much.



If a learning algorithm is suffering from high variance, getting more training data is likely to help



### Deciding what to do next revisited

Our decision process can be broken down as follows:

- Getting more training examples. Fixes high variance
- Trying smaller sets of features. Fixes high variance
- Adding features. Fixes high bias
- Adding polynomial features. Fixes high bias
- Decreasing  $\lambda$ . Fixes high bias
- Increasing  $\lambda$ . Fixes high variance

A neural network with fewer parameters is prone to underfitting.

It's also computationally cheaper

A large neural network with more parameters is prone to overfitting.

### Model complexity effects:

- Lower-order polynomials (low model complexity) have high bias and low variance. The model fits poorly consistently
- High-order polynomials (high model complexity) fit the training data extremely well and the test data extremely poorly. These have low bias on the training data, but very high variance

## Precision / Recall /

		Actual class	
		1	0
Predicted class	1	True Positive False positive	False negative
	0	False negative	True Negative

$$\text{Precision} = \frac{TP}{TP + FP}$$

$$\text{Recall} = \frac{TP}{TP + FN}$$

In regression:

$$\text{Precision} = \frac{TP}{\text{n}^{\circ} \text{ of predicted positives}}$$

$$\text{Recall} = \frac{TP}{\text{n}^{\circ} \text{ of actual positive}}$$

To compose precision and recall is necessary create

- score model (no average)

$$F_1 \text{ score} = 2 \frac{P \cdot R}{P + R}$$

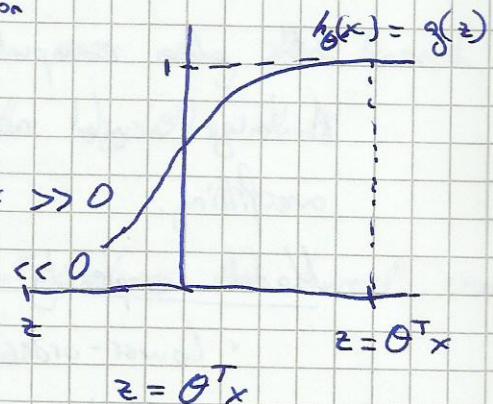
## Support Vector Machines

Alternative view of logistic regression

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

• If  $y=1$ , we want  $h_{\theta}(x) \approx 1$ ,  $\theta^T x \gg 0$

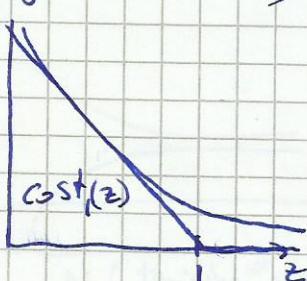
• If  $y=0$ , we want  $h_{\theta}(x) \approx 0$ ,  $\theta^T x \ll 0$



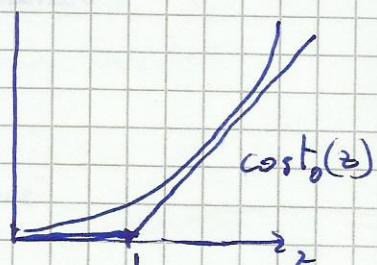
Cost of example:  $-(y \log h_{\theta}(x) + (1-y) \log (1 - h_{\theta}(x))) =$

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

If  $y=1 (\theta^T x \gg 0)$



If  $y=0 (\theta^T x \ll 0)$



Logistic regression:

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

SVM:

$$\min_{\theta} C \sum_{i=1}^m \left[ y^{(i)} \text{cost}_1(\theta^T x^{(i)}) + (1-y^{(i)}) \text{cost}_0(\theta^T x^{(i)}) + \frac{1}{2} \sum_{j=1}^n \theta_j^2 \right]$$

### Kernel

Given  $x$ ,  $f_i = \text{similarity}(x, \ell^{(i)}) = \exp \left( - \frac{\|x - \ell^{(i)}\|^2}{2\sigma^2} \right)$   
Kernel Gaussian Kernels

If  $x \approx \ell^{(i)}$ ,  $f_i \approx 1$ ; If  $x$  is far from  $\ell^{(i)}$ ,  $f_i \approx 0$

$C = \frac{1}{\lambda}$ . Large  $C$ : lower bias, high variance  
Small  $C$ : higher bias, low variance

### Logistic regression vs SVMs

$n$  = number of features ( $x \in \mathbb{R}^{n+1}$ )

$m$  = number of training examples

If  $n$  is large (relative to  $m$ ):

Use logistic regression, or SVMs without a kernel ("linear kernel")

If  $n$  is small,  $m$  is intermediate:

Use SVM with Gaussian Kernel

If  $n$  is small,  $m$  is large:

Create/add more features, then use logistic regression or SVM without a kernel

### Clustering

#### K-means

Randomly initialize  $K$  cluster centroids  $\mu_1, \mu_2, \dots, \mu_K \in \mathbb{R}^n$

Repeat {

for  $i=1$  to  $m$

$c^{(i)} :=$  index of cluster centroid closest to  $x^{(i)}$

for  $k=1$  to  $K$

$\mu_k :=$  average (mean) of points assigned to cluster  $k$

## Random initialization

For  $i=1$  to  $100\}$

Randomly initialize k-means

Run k-means. Get  $c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_k$

Compute cost function (distortion)

$$J(c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_k)$$

}

To choose the number of clusters, we could use the "elbow method"

## Data Compression

Reduce from 2-dimension to 1-dimension:

Find a direction (a vector  $a^{(1)} \in \mathbb{R}^n$ ) onto which to project the data so as to minimize the projection error

Typically, choose  $K$  to be smallest value so that

$$\frac{1}{m} \sum_{i=1}^m \|x^{(i)} - x_{approx}^{(i)}\|^2$$
$$\frac{1}{m} \sum_{i=1}^m \|x^{(i)}\|^2$$

99% variance is retained