

# Classification

Exp:

Email: Spam / Not Spam?

Online Transaction: Fraudulent (Yes/No)?

Tumor: Malignant / Benign?

binary classification

2 classes  $\Rightarrow y \in \{0, 1\}$

absent  
or  
something

present  
or  
something

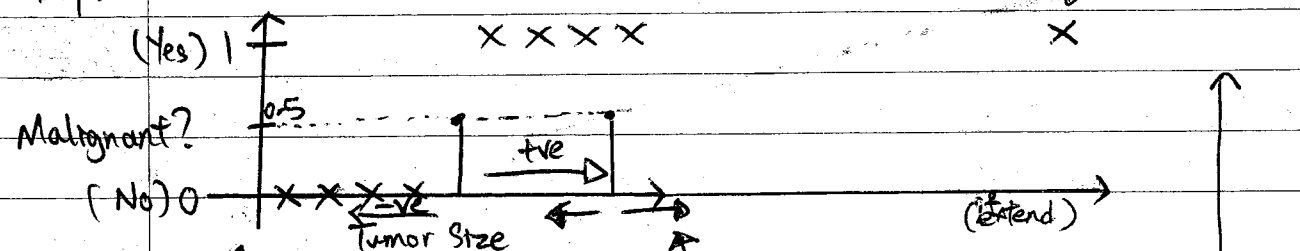
0: "Negative class" (e.g., benign tumor)

1: "Positive class" (e.g., malignant tumor)

Multiclass

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

Exp:



Use linear regression to draw out graph of line of best fit and plot threshold point to find out x.

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

magnitude that must be exceeded for certain result

[Threshold] classifier output  $h_0(x)$  at 0.5:

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

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

But when an extra example come out far away from training set, linear regression will come out with different point and different threshold.

Conclusion:

So apply linear regression on classification problem often not a good idea.

$\therefore$  Don't use linear regression on classification problem

Classification:  $y = 0$  or  $1$

Linear regression  $h_0(x)$  can be  $> 1$  or  $< 0$

Logistic Regression:  $0 \leq h_0(x) \leq 1$

$\hookrightarrow$  classification algorithm that we apply to settings where label  $y$  is discrete value (1/0)

Things Function we going to use to represent over hypothesis in classification problem.

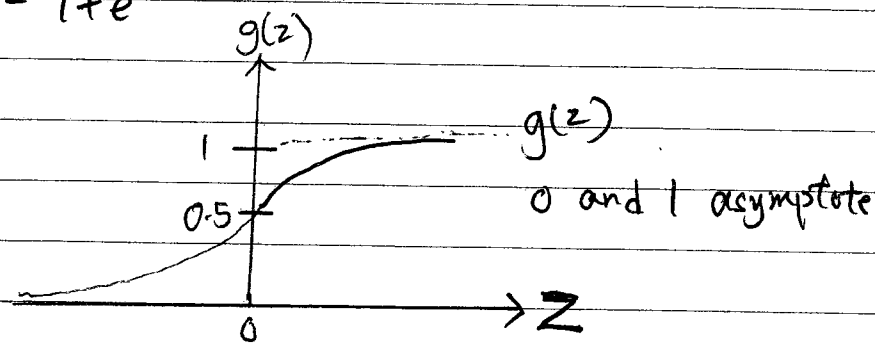
Logistic Regression Model

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

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

$g(z) = \frac{1}{1+e^{-z}}$   $\leftarrow$  sigmoid function / logistic function

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



Interpretation of Hypothesis Output

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

tumor classification example:

If

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

$h_\theta(x) = 0.7$  (patient with features  $x$ , probability  $y=1$  is 0.7).

Tell patient that 70% chance of tumor being malignant

2)  $h_\theta(x) = p(y=1|x;\theta)$  "probability that  $y=1$ , given  $x$ , parameterized by  $\theta$ "  
 $y = 0$  or  $1$

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

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

70% chance of tumor being malignant

30% being benign

To understand when it makes predictions  $y=1$  and when it makes prediction  $y=0$

Logistic regression

$$h_\theta(x) = g(\theta^T x) = P(y=1|x;\theta)$$

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

$$\begin{aligned} z=0, e^0=1, g(z)=\frac{1}{2} \\ z \rightarrow \infty, e^z \rightarrow \infty, g(z)=1 \\ z \rightarrow -\infty, e^z \rightarrow 0, g(z)=0 \end{aligned}$$

Suppose predict " $y=1$ " if  $h_\theta(x) \geq 0.5$  (estimate probability of  $y=1 \geq 0.5$ )

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

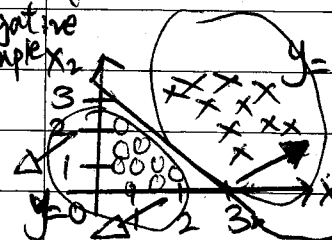
predict " $y=0$ " if  $h_\theta(x) < 0.5$  (estimate probability of  $y=1 < 0.5$ )  
 $\theta^T x < 0$

In the sigmoid function,  
 $g(z) \geq 0.5$   
when  $z \geq 0$

$$h_\theta(x) = g(\theta^T x) \geq 0.5$$

Whenever  $\theta^T x \geq 0$

$x \rightarrow$  positive example  
 $0 \rightarrow$  negative example



$$h_\theta(x) = 0.5$$

$$x_1 + x_2 = 3$$

$[y=1] x_1 + x_2 > 3$  is upper side of straight line.

$[y=0] x_1 + x_2 < 3$  is down side of s-line.

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

Let's say 2

examples

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

use this,  $y$  怎样最后  $= 1 / = 0$ .  
try  $\rightarrow \theta^T x$

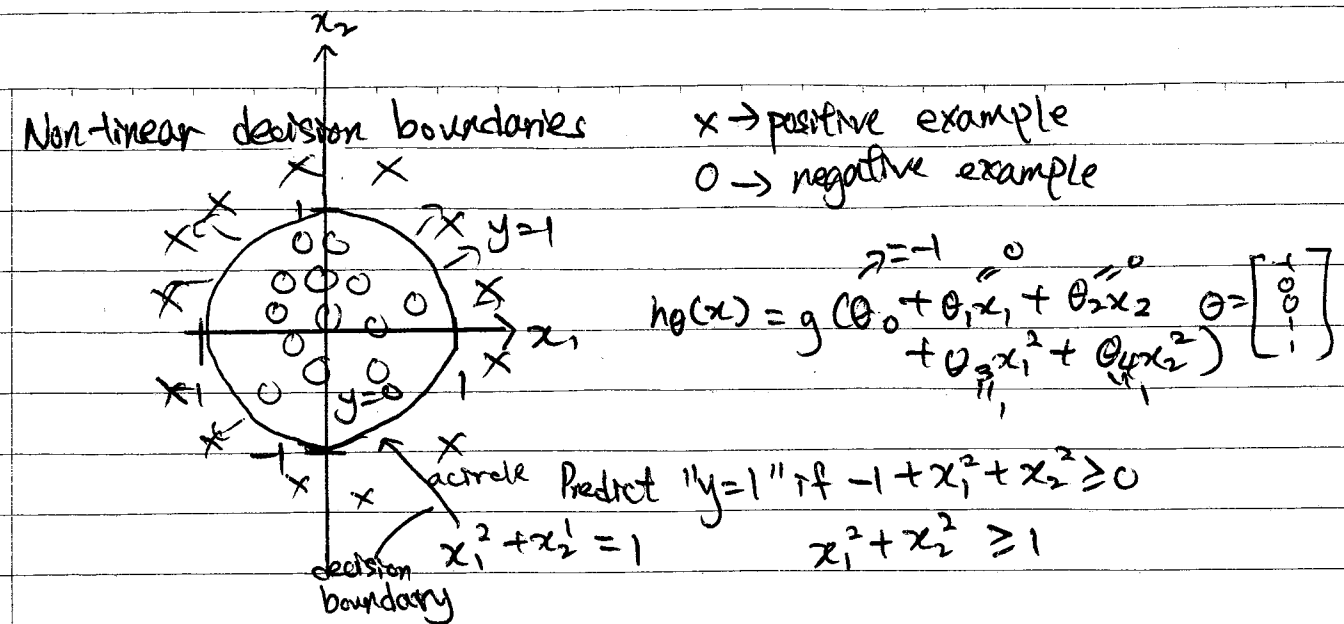
Predict " $y=1$ " if  $(-3 + x_1 + x_2) \geq 0$

$x_1 + x_2 \geq 3$  (这时  $y=1$ , 再  $y=0$ )

Def

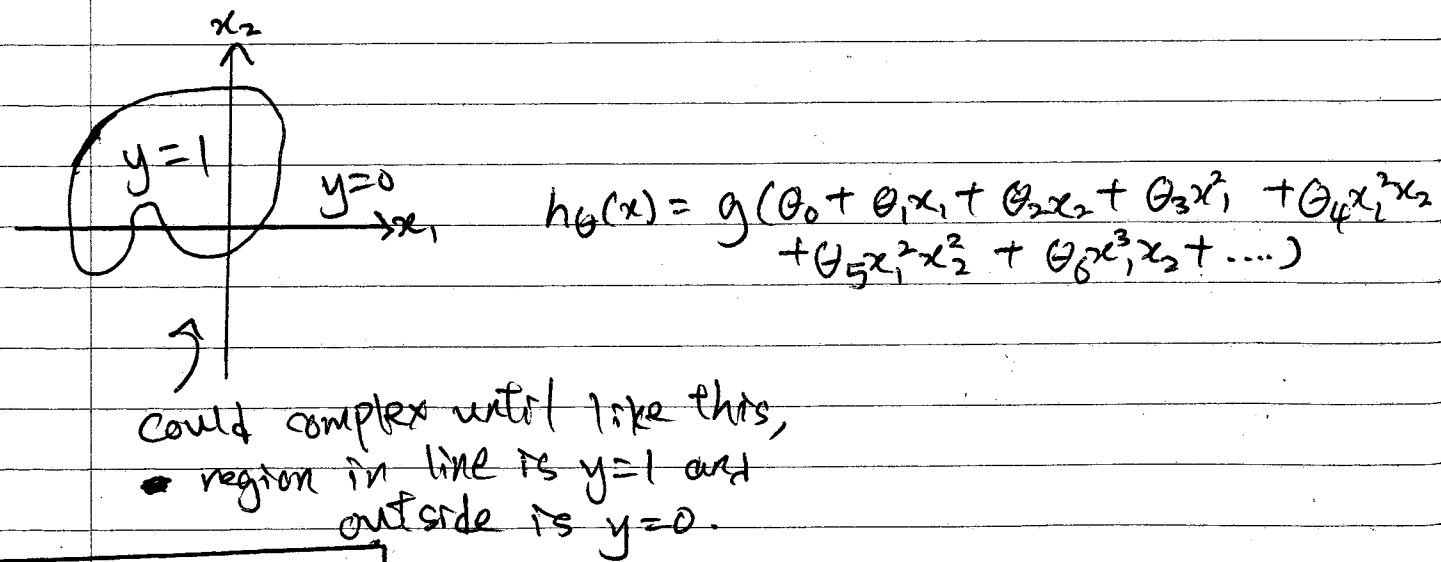
Decision boundary line = line that separate  $y=1$  and  $y=0$  region

We no need plot training set to plot decision boundary



Decision boundary is a property, not of the training set but of the hypothesis under parameters.

(Training set)  $\rightarrow$  not to define decision boundary  
may be used to fit (parameter  $\theta$ ).  $\rightarrow$  define decision boundary



Example:  $\theta = \begin{bmatrix} 5 \\ -1 \\ 0 \end{bmatrix}$

$y=1$  if  $5 + (-1)x_1 + 0x_2 \geq 0$

$5 - x_1 \geq 0 \Rightarrow x_1 \leq 5$

DB is a vertical line on graph where  $x_1 = 5$ , everything to left is  $y=1$  and to right denotes  $y=0$ .

(Higher polynomials) more complex decision boundaries  
(Higher polynomial features)

## Logistic Regression

### Cost function

To fit parameters of  $\theta$  for logistic regression.

Example: Supervised learning problem of fitting logistic regression model  
Training set:  $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})\}$

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

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

How to choose parameters  $\theta$ ?

### Cost Function

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

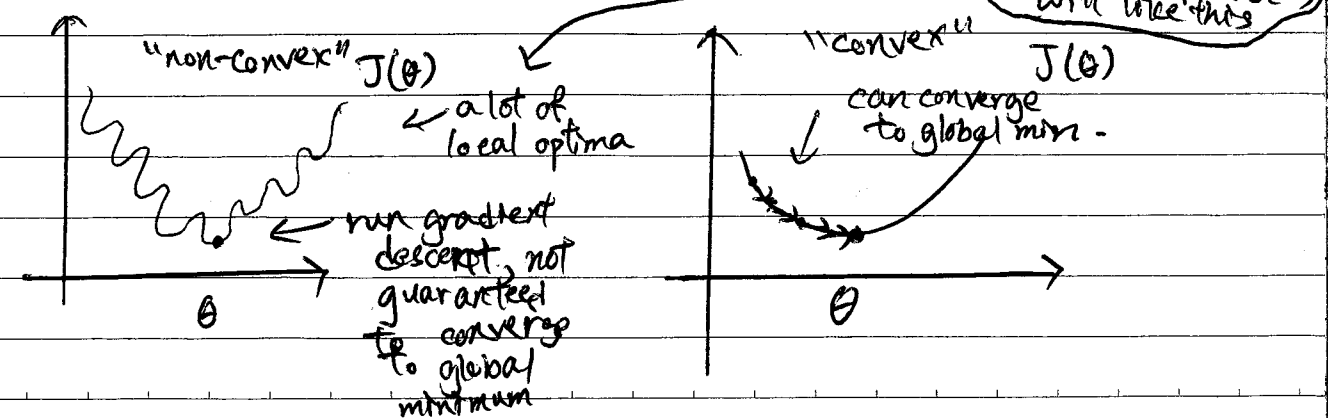
$\rightarrow \text{Cost}(h_{\theta}(x), y) = \frac{1}{2} (h_{\theta}(x) - y)^2$

this works well in linear regression,  $\rightarrow$  this learning algorithm have to pay if it outputs that value, if prediction is  $h_{\theta}(x)$  and actual label is  $y$ .

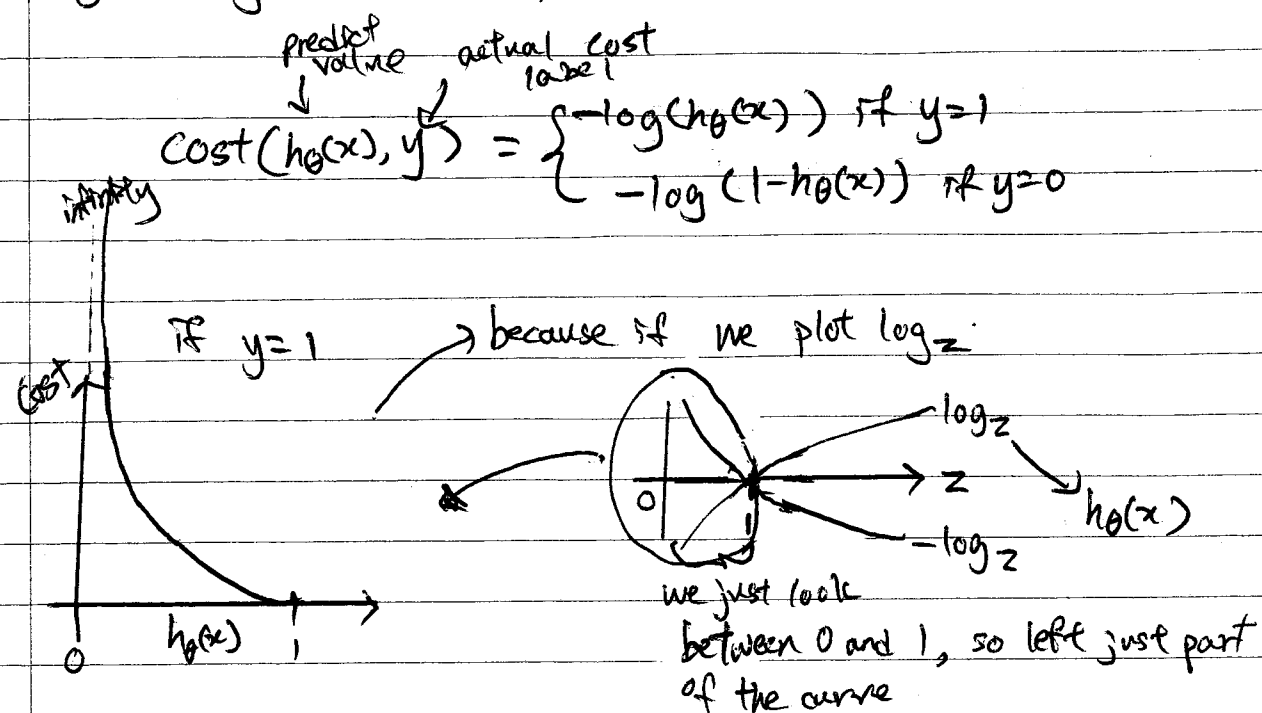
what predicted  $\rightarrow$  actual value we have,  $0$  for  $y$ .

For logistic regression  $\frac{1}{1 + e^{-\theta^T x}}$  (complicated non-linear function)

But in logistic regression, it would be a non-convex function of parameter's data



## Logistic regression cost function



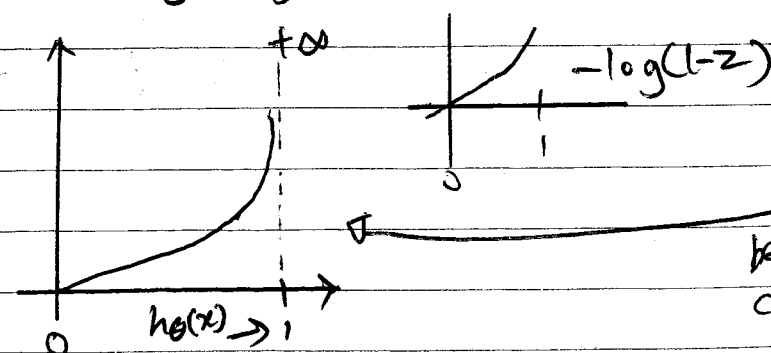
Cost = 0 if  $y=1, h_\theta(x)=1$

But as  $h_\theta(x) \rightarrow 0$   
cost  $\rightarrow \infty$

captures intuition that if  $h_\theta(x)=0$ ,  
choice of  $y=1$  is 0.  $\rightarrow$  Predict  $P(y=1|x;\theta)=0$ , but  $y=1$ ,  
we'll penalize learning algorithm by a  
very large cost

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

if  $y=0$ , but we predict  $y=1$  with almost certainty, probably 1,  
we end up paying a very large cost.



because, when near to 1,  
cost nearly infinity,  
so if it really 1,  
cost will be very large.

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

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

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

- If correct answer  $y$  is 0, then the cost function will be 0 if our hypothesis function also outputs 0. If our hypothesis approaches 1, cost function will approach infinity.
- If correct answer  $y$  is 1, then the cost function will be 0 if our hypothesis function also outputs 1. If our hypothesis approaches 0, then the cost function will approach infinity.
- Writing cost function this way guarantees  $J(\theta)$  is convex for logistic regression.

Logistic regression — Simplified cost function and gradient descent  
cost function

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

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

[Classification problem] Note:  $y=0$  or  $1$  always

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

Only 2 possible cases,  $y=1$  or  $y=0$

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

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

whole eqn = 0  
1-1=0  
1-0=1  
Same

Logistic regression cost function

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

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

vectorized implementation is

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

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

To fit parameters  $\theta$ :

min  $J(\theta)$  → Get  $\theta$

$\theta$  that fit our training set

To make a prediction given new  $x$ :

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

and output our prediction as this

→ output of hypothesis

Interpret as  $P(y=1|x;\theta)$

estimating the probability that  $y=1$

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 min  $J(\theta)$ :

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

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

so is like this

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix}$$

Repeat  $\left\{ \right.$

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

$\left. \right\}$

(simultaneously update all  $\theta_j$ )  
simultaneously update all of your values of  $\theta$

Algorithm looks identical to linear regression!

Different is that definition for this hypothesis has changed

Linear:  $R$

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

Logistic:  $R$

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

So not the same thing as gradient descent for linear regression.

[Some method] monitor linear regression to make sure it is converging. same with logistic regression.

You are running gradient descent to fit a logistic regression model with parameter  $\theta \in \mathbb{R}^{n+1}$ . Which of the following is a reasonable way to make sure learning rate  $\alpha$  is set properly and that gradient descent is running correctly?

- Plot  $J(\theta) = -\frac{1}{m} \sum_{i=1}^m [y^{(i)} \log h_{\theta}(x^{(i)}) + (1-y^{(i)}) \log(1-h_{\theta}(x^{(i)}))]$  as a function of number of iterations and make sure  $J(\theta)$  is decreasing on every iteration.

Apply feature scaling can also make gradient descent run faster for logistic regression.

Advanced Optimization - algorithm concept  $\Rightarrow$  run faster than g-descent  
Very large training set no problem

Given  $\theta$ , we have code that can compute

$$-J(\theta)$$

$$-\frac{\partial}{\partial \theta_j} J(\theta) \text{ (for } j=0,1,\dots,n)$$

to compute

Use  $\alpha$  for loop to update. Ideally, we use a vectorized implementation

vectorized implementation:

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

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

Optimization Algorithms!

- Gradient descent

- Conjugate gradient

- BFGS

- L-BFGS

Use libraries instead of writing out these sophisticated algorithms unless you are expert in numerical computing

Advantages

- No need manually pick  $\alpha$
- Often faster than g-descent

Disadvantages:

- more complex

have a default inner-loop called line search algorithm that automatically tries out diff. value for  $\alpha$  and pick good  $\alpha$ , can also pick a diff.  $\alpha$  for every iteration.

Date as a function

How to use algorithms in the previous page

Example:

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

To min  $J(\theta)$

$$\theta_1 = 5, \theta_2 = 5$$

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

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

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

function [jVal, gradient]

= costFunction(theta)

$$jVal = (\text{theta}(1) - 5)^2 + (\text{theta}(2) - 5)^2;$$

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

$$\text{gradient}(1) = 2 * (\text{theta}(1) - 5);$$

$$\text{gradient}(2) = 2 * (\text{theta}(2) - 5);$$

To apply one of advanced optimization algorithm to minimize cost function J.

Use something like g descent but more advanced. Can implement octave function

(mean you going to set gradient descent parameter to on)

then called the advanced optimization function

maximum of 100 iterations should be int, not string. thing

an options as a data structure stores options you want

Give initial guess for  $\theta$

Options = optimset('GradObj', 'on', 'MaxIter', 100);

initialTheta = zeros(2, 1);  
[optTheta, functionVal, exitFlag] ...

= fminunc(@costFunction, initialTheta, options);

Stand for function minimization unconstrained in Octave

@ symbol represent a pointer to the cost function that we defined up there and use it

$\theta \in \mathbb{R}^d$   $d \geq 2$  or fminunc can't work. (use 'help' to check details)

$$\text{theta} = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix} \leftarrow \begin{matrix} \text{theta}(1) \\ \text{theta}(2) \\ \vdots \\ \text{theta}(n+1) \end{matrix}$$

function [jVal, gradient] = costFunction(theta)  $\Rightarrow$  function that return cost function and return the gradient

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

$$\text{gradient}(1) = [\text{code to compute } \frac{\partial}{\partial \theta_1} J(\theta)];$$

$$\text{gradient}(2) = [\text{code to compute } \frac{\partial}{\partial \theta_2} J(\theta)];$$

$$\vdots$$

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

if want to use this optimization to linear regressions, just must put the appropriate code to compute in these.

## Logistic Regression

Multi-class classification: One-vs-all / One-vs-rest

When only 2 categories of data, then  $y = \{0, 1\}$ , now more than 2, then expand  $y = \{0, 1, \dots, n\}$

Example:

Email filtering / tagging: Work, Friends, Family, Hobby  
 $y=1$   $y=2$   $y=3$   $y=4$

Medical diagnosis: Not ill, Cold, Flu  
 $y=1$   $y=2$   $y=3$

Weather: Sunny, Cloudy, Rain, Snow

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

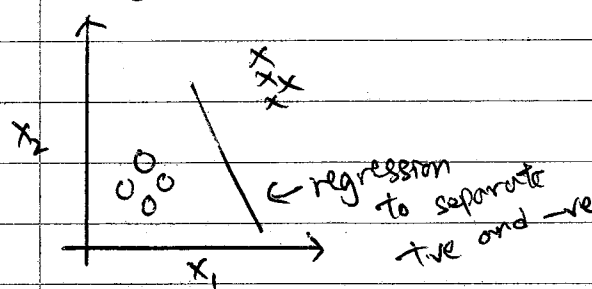
$$h_{\theta}^{(0)}(x) = P(y=0|x;\theta)$$

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

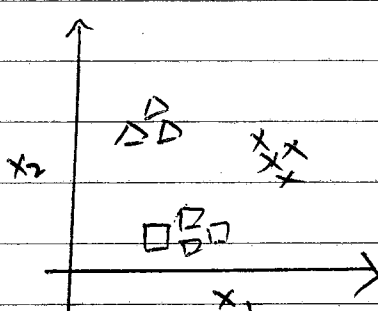
$$h_{\theta}^{(n)}(x) = P(y=n|x;\theta)$$

$$\text{prediction} = \max(h_{\theta}^{(i)}(x))$$

Binary classification:

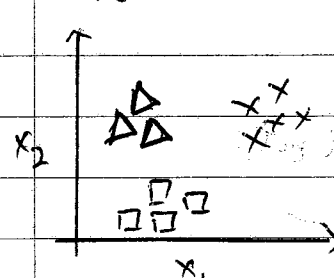


Multi-class classification:

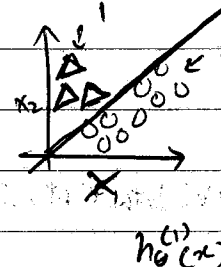


Example:

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



Class 1:  $\Delta$   
 Class 2:  $\square$   
 Class 3:  $\times$



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

$$h_{\theta}^{(2)}(x)$$

$$P(y=2|x;\theta)$$

$$h_{\theta}^{(3)}(x)$$

$$P(y=3|x;\theta)$$

estimate probability that  $y = \text{class } i$  given  $x$  and parameterized by  $\theta$ .

One-vs-all  $\Rightarrow$  Train logistic regression classifier  $h_{\theta}^{(i)}(x)$  for each class  $i$  to predict probability that  $y=i$ .

On a new input  $x$ , to make prediction, pick class  $i$  that maximizes  $\max h_{\theta}^{(i)}(x)$   $\hookrightarrow$  run 3 classifiers and pick class  $i$  that maximizes the three (value of  $i$  that gives us the highest probability)



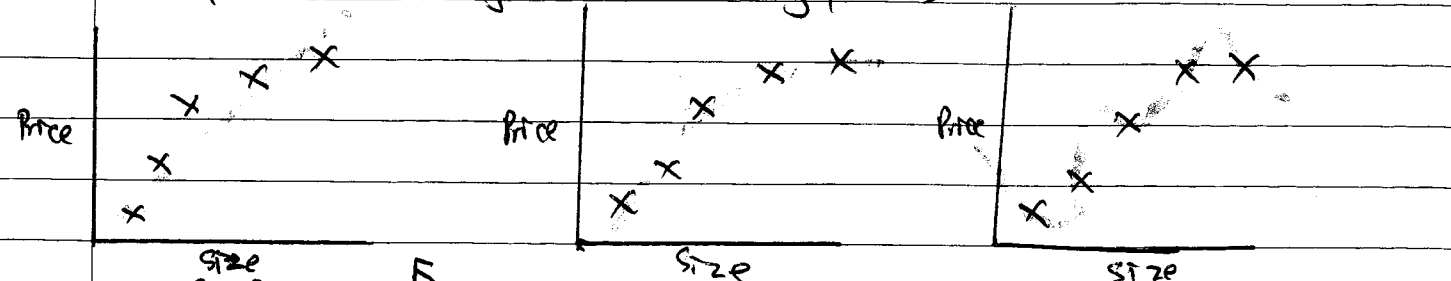
# Regularization

The problem of overfitting.

cause learning algorithm to perform ~~very~~ poorly

Date

Example: Linear regression (housing prices)



Size  $\theta_0 + \theta_1 x$   
if use linear regression, the housing price should move to right instead of keep up, so this algo, not very well suit training set.

Problem:

"Underfitting", "High bias"

=> poor fit to data

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

=> Works pretty well (Just right for fitting this data)

Size

$\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$   
=> a curve that passes through all 5 examples, but very wiggly curve, not really a good model for predicting housing prices.

Problem:

"Overfitting", "High variance"

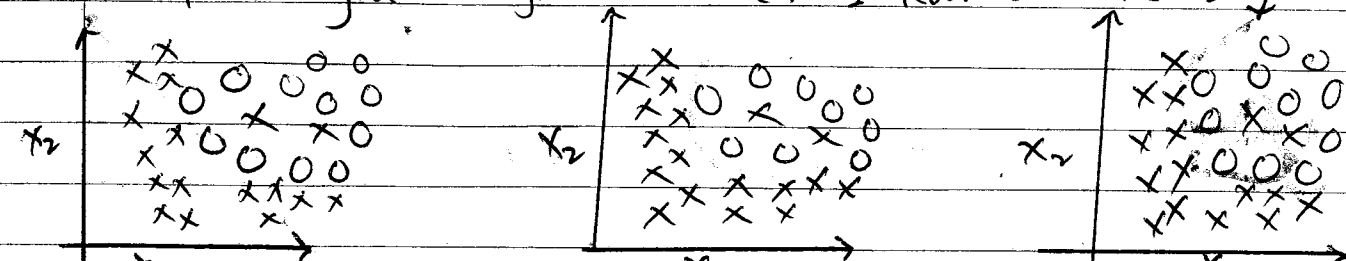
if we fit such high order polynomial, hypothesis can fit almost any function, face of possible hypothesis too large, too variable and we don't have enough data to constrain it to get good hypothesis, and so called "overfitting"

Overfitting: Problem of overfitting comes when if we have too many features, learned hypothesis will fit the training set very well

$(J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 \approx 0)$ , but end up fails too hard to fit training set, in the end fail to generalize to new examples (predict prices on new examples).

how well a hypothesis applies even to new example (data that not in training set yet)

Example: Logistic regression (with 2 features  $x_1$  and  $x_2$ )



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

$g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_2^2 + \theta_5 x_1 x_2)$   
=> we'll get a boundary like this.

$g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 x_2 + \theta_4 x_1^2 x_2^2 + \theta_5 x_1^3 x_2^2 + \theta_6 x_1^3 x_2^3 + \dots)$   
Logistic regression may contour itself, try hard to find a d-boundary that fit every example well.

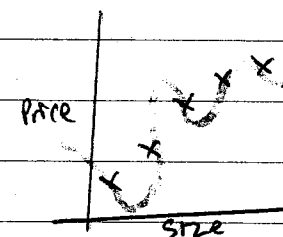
"overfitting"

(P. 对旧数据有A, 新的数据来就不合很准)

if we think overfitting is occurring, do what to address it.

Addressing overfitting:

$x_1$  = size of house  
 $x_2$  = no. of bedrooms  
 $x_3$  = no. of floors  
 $x_4$  = age of house  
 $x_5$  = average income in neighborhood  
 $x_6$  = kitchen size  
 $\vdots$   
 $x_{100}$



The way Use figure to select an appropriate degree polynomial to fit training set is not really work with many features, because it becomes very hard to plot data, much harder to visualize it and decide what features to keep or not.

Options:

1. Reduce number of features.

- Manually select which features to keep.

- Model selection algorithm (later in course) -> auto-decide which features to keep/throw (But when throwing away features just like throw away information about the problem)

2. Regularization:

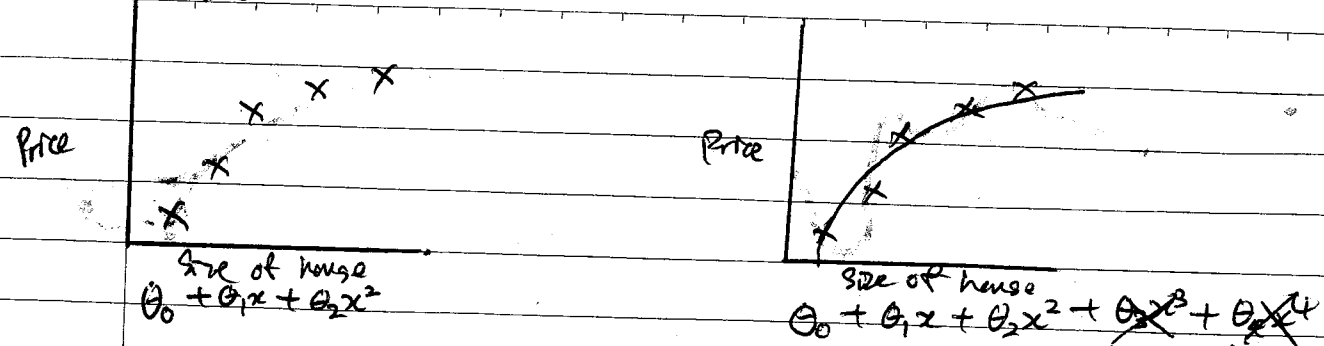
- Keep all the features, but reduce magnitude / values of parameters  $\theta$ .  
- Works well when we have a lot of features, each of which contributes a bit to predicting  $y$ .

(大数法则, 所以不想 throw)

# Regularization

## Cost function

### Intuition



Suppose we penalize and make  $\theta_3, \theta_4$  really small.

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

Example:

Give huge num to  $\theta_3$  and  $\theta_4$

$\Rightarrow$  Now the only way to make new cost function small is if  $\theta_3$  and  $\theta_4$  are small.

End up  $\theta_3 \approx 0$   $\theta_4 \approx 0$

Will make the graph back to (Good) Quadratic function, + tiny contribution from small term,  $\theta_3$  and  $\theta_4$

### Regularization

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

- "Simpler" hypothesis

- Less prone to overfitting

Small value of parameter make "simpler" hypothesis

### Housing:

Features:  $x_1, x_2, \dots, x_{100}$

Parameters:  $\theta_0, \theta_1, \theta_2, \dots, \theta_{100}$

too many do not know which one to pick, to shrink it.

If we had like 100 features, it's hard to pick which are the ones that less likely to be relevant.

### Example

Now, take linear regression cost function

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 \leftarrow \text{modify this and shrink all parameter because don't know pick which to shrink}$$

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

it starts from 1, so not going to penalize  $\theta_0$  being large.

Add extra regularization term at the end.

### Regularization

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

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

regularization term

determine how much the costs of our theta parameters are inflated.

regularization parameter

$\Rightarrow$  control a trade-off b/w 2 diff. goals.

train to fit training data well

keep parameters small

High order polynomial may have wiggly, curvy functions. If you still fit high order polynomial with all polynomial features in there but need make sure use regularized objective (2nd) to get curve more smooth.

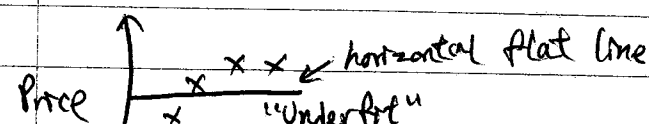
therefore keeping hypothesis relatively simple to avoid overfitting.

If  $\lambda$  set to an extremely large value (exp:  $\lambda = 10^{10}$ ),

$\Rightarrow$  Algorithm results in underfitting (fails to fit even the training set)

$\Rightarrow$  end up penalizing  $\theta_1, \theta_2, \theta_3, \theta_4 \dots$  very highly  $\Rightarrow$  which mean end up

$\theta_1 \approx 0, \theta_2 \approx 0, \theta_3 \approx 0, \theta_4 \approx 0$



Size of house

$$\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

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

make only fit horizontal straight line to data

For regularization to work well, some care should be taken to choose the regularization parameter,  $\lambda$ .



# Regularization

## Regularized linear regression

No. \_\_\_\_\_

Date \_\_\_\_\_

No. \_\_\_\_\_

Date \_\_\_\_\_

- Generalize 2 algorithms (g-descent, normal eqn) to the case of regularized linear regression

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

Find parameter theta that minimizes this cost function

Separate G-descent algorithm into 2

repeat {

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

$$\theta_j := \theta_j - \alpha \left[ \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)} + \frac{\lambda}{m} \theta_j \right] \quad (j = 1, 2, 3, \dots, n)$$

This square bracket mess is just  $\frac{\partial}{\partial \theta_j} J(\theta)$  (regularized)

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

perform update  
 $< 1$  but +ve so to shrink  $\theta_j$

Because for our regularized linear regression, we penalize  $\theta_1, \theta_2, \dots, \theta_n$  but not  $\theta_0$

## Normal equation

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

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

$$\frac{\partial}{\partial \theta_j} J(\theta) = 0 \quad \text{--- doing bunch of math, and get this eqn that minimize cost function.}$$

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

Example:  $n=2$ ,  $(3 \times 3)$  matrix as it is  $(n+1) \times (n+1)$  matrix. diagonal is one "0" and other all "1" and part apart from diagonal is all "0"

Non-invertibility (Optional / advanced)

Suppose  $m \leq n$  (#features)  $\Rightarrow$   $X$  is non-invertible if  $m < n$  may be non-invertible if  $m = n$

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

$\rightarrow$  will be non-invertible / singular / matrix will be degenerate  
 pinv, will not give a good hypothesis inv, will not work

if  $\lambda > 0$ ,

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

because this  $> 0$   
 will not be singular and will be invertible

$\therefore$  regularization take care of any non-invertibility issues of  $X^T X$  matrix as well.

# Regularization

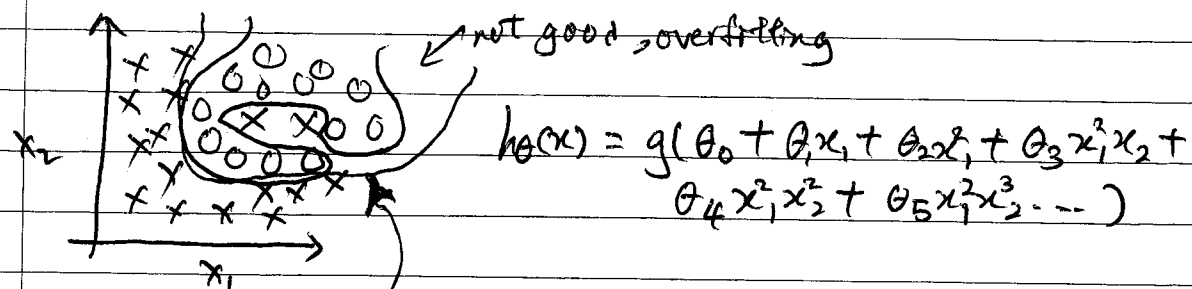
## Logistic Regression

No. \_\_\_\_\_  
Date \_\_\_\_\_

$\theta_0$  in Octave is  $\text{theta}(1)$   
 $\theta_1$  in Octave is  $\text{theta}(2)$   
 $\theta_n$  in Octave is  $\text{theta}(n+1)$

Adapt gradient descent and advanced optimization to work for regularized logistic regression.

Regularized logistic regression



When penalizing and parameter become smaller, graph become like this

Cost function:

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

Penalize  $\theta_1, \theta_2, \dots, \theta_n$  from being too large

Gradient descent

Repeat {

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

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

}

$\frac{\partial}{\partial \theta_j} J(\theta)$   $(j = 1, 2, 3, \dots, n)$

Cost function used regularization

When using regularized logistic regression, Plot

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

as a function of the number of iterations and make sure it's decreasing to monitor whether g-descent is working correctly or not.

Advanced optimization

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

function [jVal, gradient] = costFunction(theta)

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

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

include in regularization

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

$$\frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_0^{(i)} \quad \leftarrow \text{don't change}$$

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

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

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

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

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

So if implement this cost function and pass into fminunc or other advanced optimization technique, will minimize new regularized cost function  $J$  of  $\theta$  ( $J(\theta)$ )