## Logistic  Regression  for  Classification  Problems

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

$$g(z) = \frac{1}{1+e^{-z}} \longrightarrow \text{Sigmoid  function}$$

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

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

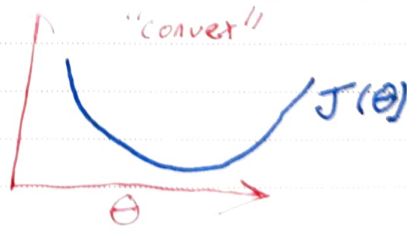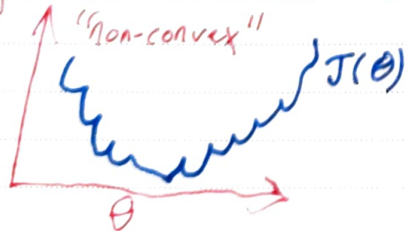$h_\theta(x)$  is  going  to be  conceptually  represented  as  the  probability  that  $y=1$
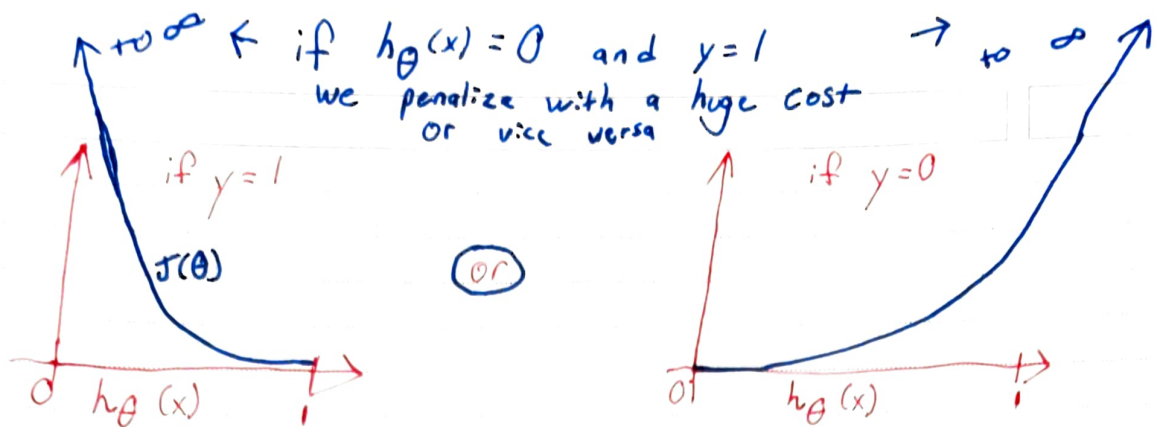
$$y \in \{0, 1\}$$

## Cost  Function :

Linear R

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

if  using  the  linear  reg........   ........tion  with  the
sigmoid  definition  of  $h_\theta(x)$.......   ....ld  be  "non-convex"

"non-convex"                    "convex"



Logistic  Regression:

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

to ∞ ← if $h_\theta(x) = 0$ and $y = 1$ → to ∞

we penalize with a huge cost
or vice versa

if $y = 1$

$J(\theta)$

(or)

if $y = 0$

$0 \quad h_\theta(x)$

$0 \quad h_\theta(x)$

## Logistic Regression Cost Function:

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

## Gradient descent:    Uni Variate

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

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

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

# Week 3 - Overfitting

Regularization helps with overfitting by making theta smaller.

Gradient Descent:

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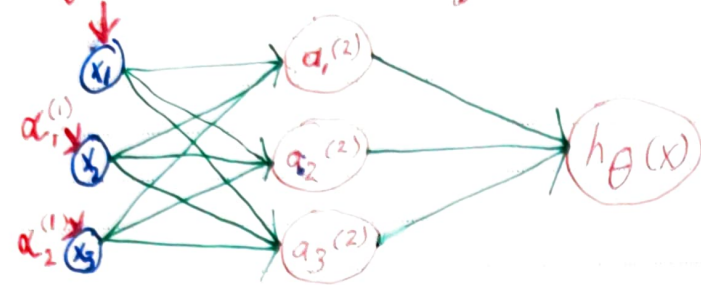
Cost Function:

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

We can find performance issues ... ie. low cost but not the lowest cost through gradient checking

Gradient Checking:

# Week 4 - Neural Networks

$a_0^{(1)} = x$    * can be thought of as the 1st activation layer



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

$\theta_j$ = matrix of weights controlling function mapping from layer $j$ to $j+1$

Layer⁽¹⁾ (input)    Layer⁽²⁾ (hidden)    Layer⁽³⁾ (output)

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

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

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

$$h_\theta(x) = a_1^{(3)} = g\left(\theta_{10}^{(2)} a_0^{(2)} + \theta_{11}^{(2)} a_1^{(2)} + \theta_{12}^{(2)} a_2^{(2)} + \theta_{13}^{(2)} a_3^{(2)}\right)$$

## Vectorized Approach

Previously    $g(z) = g(\theta^T x)$    therefore:

now    $$z_1^{(2)} = \theta_{10}^{(1)} x_0 + \theta_{11}^{(1)} x_1 + \theta_{12}^{(1)} x_2 + \theta_{13}^{(1)} x_3$$

$$z_2^{(2)} = \theta_{20}^{(1)} x_0 + \theta_{21}^{(1)} x_1 + \theta_{22}^{(1)} x_2 + \theta_{23}^{(1)} x_3$$

$$z_3^{(2)} = \theta_{30}^{(1)} x_0 + \theta_{31}^{(1)} x_1 + \theta_{32}^{(1)} x_2 + \theta_{33}^{(1)} x_3$$

$z^{(2)} = \theta^{(1)} x$  or  $z^{(2)} = \theta^{(1)} a^{(1)}$

$a^{(2)} = g(z^{(2)})$ ← element wise sigmoid of z

*we also need to add a bias - $a_0^{(2)}$ *

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

$$\boxed{h_\theta(x) = a^{(3)} = g(z^3)}$$

# Week 5 - Neural Network Backpropagation

## Cost Function

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

*not ith training example

$\rightarrow (\theta_{j,i}^{(l)})^2$

$L$ = total # of layers
$S_l$ = # of units (not counting bias) in layer $l$
$K$ = # of output units / classes
$j$ = node in layer $l$
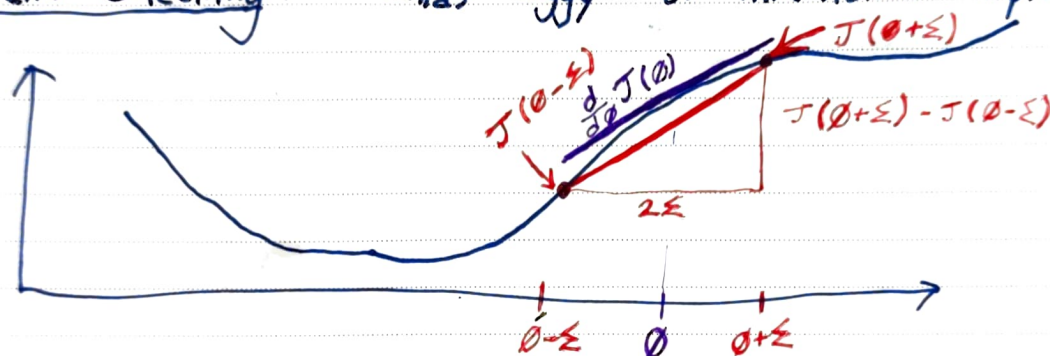
$$\boxed{g'(z^{(3)}) = a^{(3)} .* (1-a^{(3)})}$$

To calculate error for any given node:

$$\delta_j^{(4)} = a_j^{(4)} - y_j \quad \text{or} \quad \underset{\text{\# of output units}}{\underline{\delta^{(4)} = a^{(4)} - y}} \rightarrow \delta^{(3)} = (\theta^{(3)})^T \delta^{(4)} .* g'(z^{(3)})$$

$\boxed{\text{NO } \delta^{(1)}}$

## Gradient Checking: finds "buggy" or innefficient implementation



$\varepsilon \approx 10^{-4}$

$$Q \in R: \quad \frac{d}{d\theta} J\theta \approx \frac{J(\theta+\varepsilon) - J(\theta-\varepsilon)}{2\varepsilon}$$

Vectorized:
$$\frac{d}{d\theta_1} J(\theta) \approx \frac{J(\theta_1+\varepsilon, \theta_2, \theta_3, ..., \theta_n) - J(\theta_1-\varepsilon, \theta_2, \theta_3, ... \theta_n)}{2\varepsilon}$$

$$\frac{d}{d\theta_2} J(\theta) \approx \frac{J(\theta_1, \theta_2+\varepsilon, \theta_3, ..., \theta_n) - J(\theta_1, \theta_2-\varepsilon, \theta_3, ... \theta_n)}{2\varepsilon}$$

$$\vdots$$

$$\frac{d}{d\theta_n}$$

## Gradient Checking Implementation:

theta is a vector (unrolled).
for $i = (1 \, \text{to} \, n)$;
    thetaPlus = theta
    thetaPlus (i) += Epsillon
    thetaMinus = theta
    thetaMinus (i) -= Epsillon
    gradApprox (i) = (J(thetaPlus) - J(thetaMinus)) / (2 * Epsillon)
end;

1. Implement backprop to compute DVec (unrolled deltas)
2. Implement Gradient Checking to compute gradApprox
3. Turn off gradient checking before training

## Random Initialization:      How can we initialize theta?

* We cannot initialize all thetas to zero or all neurons will be the same.

aka Symmetry Breaking

Set each $\theta$ to a random value between $-\varepsilon$ and $\varepsilon$

            ↙ 10x11 matrix between 0 and 1
theta1 = rand(10, 11) * (2 * epsilon) - epsilon*
            * not the same as epsilon in gradient checking

# Week 5 - Overview for Neural Networks

## 1. Network Architecture

- input
- output
- hidden          (usually 1 layer or same # of units per layer)
  - more  units  the better  (usually more than # of inputs)

## 2. Training

1. Randomly initialize weights for $\theta$
2. Implement foward propagation to get $h_\theta(x^{(i)})$ or $a$
3. Implement cost function to get $J(\theta)$
4. Implement backward propagation to get $\frac{\partial}{\partial \theta_{jk}^{(i)}} J(\theta)$
   for $i = 1:m$
     forward prop ();   returns $a$   (activations)
     backward prop ();  returns $\frac{\partial}{\partial \theta_{jk}^{(i)}} J(\theta)$  (derivative)
5. Use gradient checking to verify backprop, then disable it.
6. Use a optimization function like gradient descent to minimize $J(\theta)$

### English Translations
1. Where do we start on the hill?
2. What is our guess?  hypothesis?
3. How wrong was our guess?
4. Which direction to the bottom of the hill? .(compass)
5. Is our compass accurate?
6. Use all the info above to make better guesses...each iteration.

$$\frac{\partial}{\partial \theta_{jk}^{(i)}} J\theta = a_j^{(\ell)} \delta_i^{(\ell+1)} *$$          *ignoring $\lambda$ / if $\lambda = 0$*

## Forward Prop:

```
y_bin = eye (num_labels) [y.T]    # create an identiy matrix for y

bias = np.ones ((m, 1))    # create a bias vector
noBias1 = np.delete (Theta1, 0, 1)
noBias2 = np.delete (Theta2, 0, 2)

a1 = concat ([bias, X], axis=1)
z2 = a1 @ Theta1.T
:     # work forward
```



## Cost Function:

$$J = (-1/m) * np.sum (y\_bin * log (h) + (1 - y\_bin) * log(1-h))$$

## Regularization:

$$reg = (lambda\_/(2*m)) * (np.sum (noBias1 **2) + np.sum (noBias2 **2))$$

$$J += reg$$

## Backprop:

```
d3 = a3 - y_bin
d2 = d3 @ noBias2 * sigmoid gradient (z2)
d1 = "no such thing as error for layer 1"
```

```
Delta1 = d2.T @ a1      # use Delta1 as an "accumulator" for our deltas
Delta2 = d3.T @ a2
```

```
Theta1 _grad = (1/m) * Delta1      # Adjust our respective theta
Theta2 _grad = (1/m) * Delta2
# Add our regularization term
Theta1_grad [:, 1:] = Theta1_grad[:, 1:] + 1 * lambda_/m * Theta1[:, 1:]
Theta2_grad [:, 1:] = Theta2_grad[:, 1:] + 1 * lambda_/m * Theta1[:, 1:]

grad = np.concatenate ([Theta1_grad.ravel(), Theta2_grad.ravel()])
```
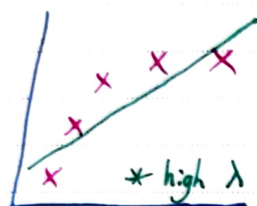
# Week 6 - ML Diagnostics

 * Don't randomly adjust parameters, or gather more data! *

Training Set: 60%
Cross Validation Set: 20%
Test Set: 20%

Plot! Plot! Plot!

1. Optimize $\theta$ w/ the training set for each polynomial
2. Find the polynomial with the lowest cost using the CV set.
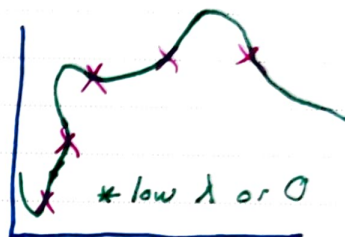3. Estimate the general error w/ the test set using the best polynomial

* The most common problems will be under or over-fitting *
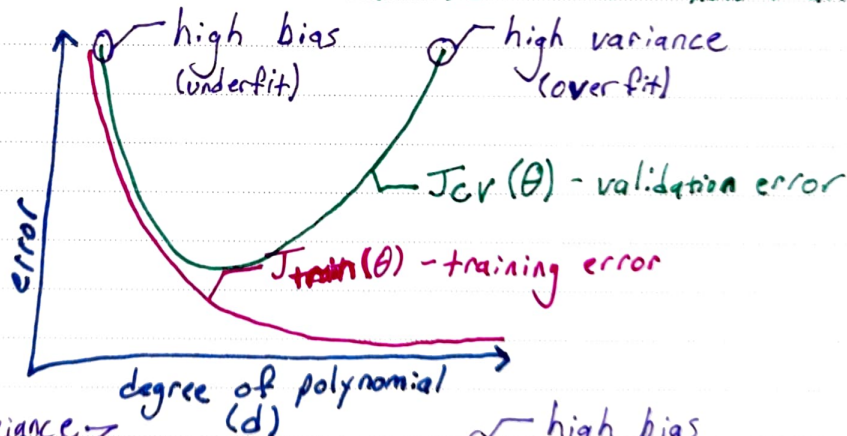


$\theta_0 + \theta_1 x$
"underfit"
$d = 1$

* high $\lambda$

$\theta_0 + \theta_1 x + \theta_2 x^2$
"Goldilocks"
$d = 2$

* good $\lambda$

$\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$
"overfit"
$d = 4$

* low $\lambda$ or $0$

* By plotting our error for Train and CV we can see how d affects J *
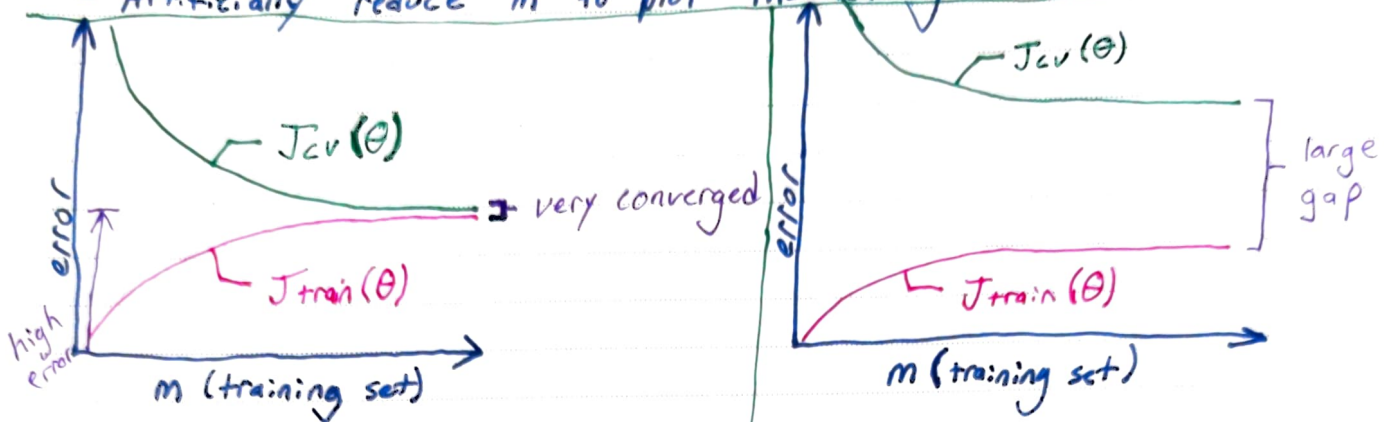
how does the # of polynomials affect our cost?



high bias (underfit)

high variance (overfit)

$J_{CV}(\theta)$ - validation error

$J_{train}(\theta)$ - training error

error

degree of polynomial (d)

how does $\lambda$ affect our cost?

high variance (overfit)

high bias (underfit)

$J_{CV}(\theta)$

$J_{train}(\theta)$

error

$\lambda$

# Week 6 - ML Diagnostics

\* Learning curves can help diagnose under or over-fitting \*
- Artificially reduce $m$ to plot the learning curves



High Bias

- more training data won't help
- lower $\lambda$ may help
- more polynomials may help
- more features may help

High Variance

- more training data may help
- higher $\lambda$ may help
- less polynomials may help
- less features may help

## Implementation of learning curves

```
for i in range (1, m+1):
    subX = X[:i, :]
    subY = y[:i]

    theta = utils.trainLinearReg (linearRegCostFunction, subX, subY, lambda_)
    error_train[i-1], _ = linearRegCostFunction (subX, subY, theta, lambda_=0)
    error_val[i-1], _ = linearRegCostFunction (Xval, yval, theta, lambd_=0)
return error_train, error_val
```

# Week 6 - Building a spam classifier

\* It is difficult to determine which features to spend time on \*

<u>Error Analysis</u> : A systematic approach to improving an algorithm

1. Start simple : Create a model in > 24 hours & test it on CV
2. Plot learning curves: decide on more data, more features, etc
3. Error Analysis : Manually review examples where the model was wrong
   \* look for error trends / patterns

<u>Error Analysis in practice:</u>

1. Categorize the reviewed errors and count them
   ie. 12: Pharma , 4: Replica /fake (54: phishing),3other, etc.
   \*how can we improve detection here →

2. What features would have helped these errors?
                 quickly
\* Our goal is to find out which scenarios are the most difficult \*
to predict. This will allow us to focus our efforts on that.

3. Numerical Evaluation : implement Accuracy % on your CV
   \* this allows you to test features that can't be evaluated by \*
   using error analysis. ie. natural language stemming

<u>Skewed Classes</u> : When your ratio of classes are skewed accuracy won't help
                 ie. 99% accuracy, but 0.5% of y=1 | 99.5% of y=0

<u>Precision/Recall</u> : A better way to measure accuracy.

Predicted ($h_\theta(x)$)

(Y)
Actual

|   | 1 | 0 |
|---|---|---|
| 1 | True Positive | False positive |
| 0 | False Negative | True Negative |

$$Precision = \frac{True\ Positives}{True\ Positives + False\ positives}$$

$$Recall = \frac{True\ Positives}{True\ Positives + False\ Negatives}$$

$*$ We can vary our threshold for predicting $y = $ <some class> $*$
ie.

- if we want $y = 1$ (cancer) only when really confident :
  - $y = 1$ if $h_\theta(x) \geq 0.7$ or even $0.9$
    $\longrightarrow$ high precision, low recall

- if we want to avoid missing positive cases :
  - $y = 1$ if $h_\theta(x) \geq 0.3$
    $\longrightarrow$ high recall, low precision

<u>$F_1$ Score</u> : a better way to measure precision/recall

$F_1$ Score $= \boxed{2\frac{PR}{P+R}}$

ie.

| P | R | $F_1$ |
|------|-----|-------|
| 0.5 | 0.4 | 0.444 |
| 0.7 | 0.1 | 0.175 |
| 0.02 | 1.0 | 0.392 |

$\leftarrow$ $F_1$ Score winner

examples


threshold
Precision
0.5
Recall