

# Classification Regression

Experiment :

Date \_\_\_\_\_

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

✓ Supervised learning  
- right ans. given

Regression

- predict continuous valued o/p

(Find structure in data)  
✓ Unsupervised learning

↳ Clustering ~ Social n/w analysis Have: some  $f^n J(\theta_0, \theta_1)$   
Want:  $\min_{\theta_0, \theta_1} J(\theta_0, \theta_1)$

\* Svd

✓ linear Regression

h-hypothesis  
(maps from x's to y's)

linear Reg (One var.)

$$h_\theta(x) = \theta_0 + \theta_1 x$$

(Cost  $f^n$ ): Idea is to  
choose  $\theta_0, \theta_1$  so that

$h_\theta(x)$  is close to  $y$  for  
over training eg.  $(x, y)$

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

$$J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2$$

Sq error  $f^n$

Gradient Descent

→ To minimize  $J(\theta_0, \theta_1)$

Outline: • Start with some  $\theta_0, \theta_1$ ,  
• Keep changing  $\theta_0, \theta_1$  to  
reduce  $J(\theta_0, \theta_1)$   
until we end up @ mi

Gradient Descent Algo:

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

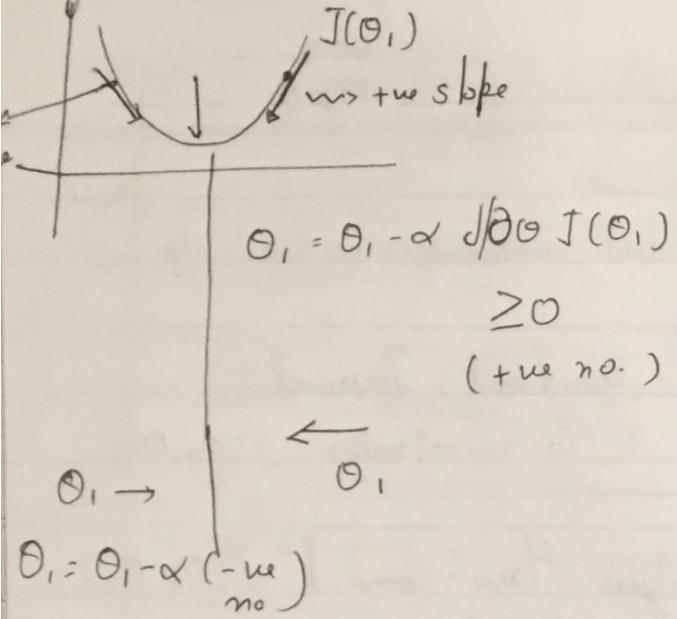
$y$  ↓ learning rate  $\alpha \frac{\partial \theta_j}{\partial \theta_j}$  (for  $j=0, 1$ )

Assignment operator

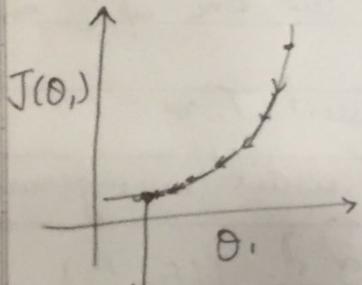
$\alpha \downarrow$  - baby steps

$\alpha \uparrow$  - giant steps

# Gradient Descent Intuition



As we approach local min., gradient descent will automatically take smaller steps. So, no need to decrease  $\alpha$  over time



gradient Descent Algo

repeat until convergence {

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

(for  $j=1 \& j=0$ )

}

$$\begin{aligned} \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1) &= \frac{\partial}{\partial \theta_j} \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2 \\ &= \frac{\partial}{\partial \theta_j} \frac{1}{2m} \sum_{i=1}^m (\underline{\theta_0 + \theta_1 x^{(i)}} - \underline{y^{(i)}})^2 \end{aligned}$$

repeat until convergence {

$$\theta_0 := \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1) = \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - \underline{y^{(i)}})$$

$$\theta_1 := \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1) = \frac{1}{m} \sum_{i=1}^m h_\theta(x^{(i)}) - \underline{y^{(i)}}$$

}

"Batch" Gradient Descent  
→ Each step uses all train eg

linear Regression Model

$$h_\theta(x) = \theta_0 + \theta_1 x$$

$$\begin{aligned} J(\theta_0, \theta_1) &= \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2 \end{aligned}$$

Scanned with CamScanner

# Matrix

Experiment:  $n \times m$ 

Date \_\_\_\_\_

$$A \cdot I = I \cdot A = A$$

$I_{n \times n}$

$\hookrightarrow m \times n$

$m \times n$        $m \times n$

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

Note:

$$AB \neq BA$$

$$A(A^{-1}) = I$$

$$AI = IA$$

## Matrix Transpose

$$A = \begin{bmatrix} 1 & 2 & 0 \\ 3 & 5 & 9 \end{bmatrix}$$

$$A^T = \begin{bmatrix} 1 & 3 \\ 2 & 5 \\ 0 & 9 \end{bmatrix}$$

Linear Regression with multiple variables

 $n$  = no. of features $x^{(i)}$  = i<sup>th</sup> features of i<sup>th</sup> training eg $x_j^{(i)}$  = value of feature j in i<sup>th</sup> "

rowno.

## New hypothesis

$$h_0(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 + \theta_4 x_4 + \dots + \theta_n x_n$$

 $(\theta_0 = 1 \text{ for } \text{constant})$ 

$$X = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \in \mathbb{R}^{n+1}$$

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_n \end{bmatrix} \in \mathbb{R}^{n+1}$$

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

$$[\theta_0, \theta_1, \theta_2, \dots, \theta_n]$$

$$\begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix}$$

 $\theta^T$  $(n+1) \times 1$   
matrix

↳ feature vector

Parameter vector

ASS TIME

# Gradient Descent

4.

Previously ( $n=1$ )

Repeat {

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

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

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

(Simultaneously update  
 $\theta_0, \theta_1$ )

New algo:

Repeat {

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

(Simultaneous update  $\theta_j$   
for  $j = 0, \dots, n$ )

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

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

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

## Feature scaling

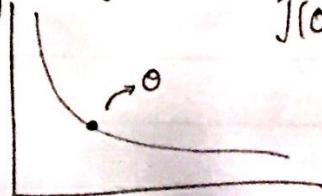
Get every feature into  $-1 \leq x_i \leq 1$  range

### Mean normalization

$$x_i \rightarrow \frac{x_i - \bar{x}_i}{s_i} \quad \begin{matrix} \rightarrow \text{to make zero mean} \\ \text{avg. value of } x_i \text{ in train set} \end{matrix}$$

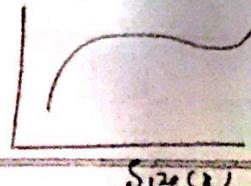
Make sure gradient descent works

$\min_{\theta} J(\theta)$   $J(\theta)$  should  $\downarrow$  with every iter



## Polynomial regression

Pred.  
 $y$



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

Normal eqn

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

What if  $X^T X$  is non-inv? (singular/degenerate)

→ Redundant features

$x_1 = \text{Size in ft}^2$

$x_2 = \dots m^2$

→ Too many features ( $m \leq n$ )

- Delete some features / use regularization

'm' training eg., 'n' features

Gradient Descent

• Need to choose  $\alpha$

• Needs many iterations

• Works well for  $n \rightarrow \infty$

Normal eqn

• No need

• Don't need to iterate

• Slow for  $n \rightarrow \infty$

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

probability that  $y=1$ , given  $x$ , parameterized by  $\theta$

$$h_\theta(u) = 0.7$$

probability that patient has tumor

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

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

Logistic Regression

Classification

Email: Spam / Not Spam?

Online Transaction: Fraudulent (Y/N)?

Tumor: Malignant / Benign?  $0 \leq h_\theta(u) \leq 1$

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

↓  
-ve class  
+ve class

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

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

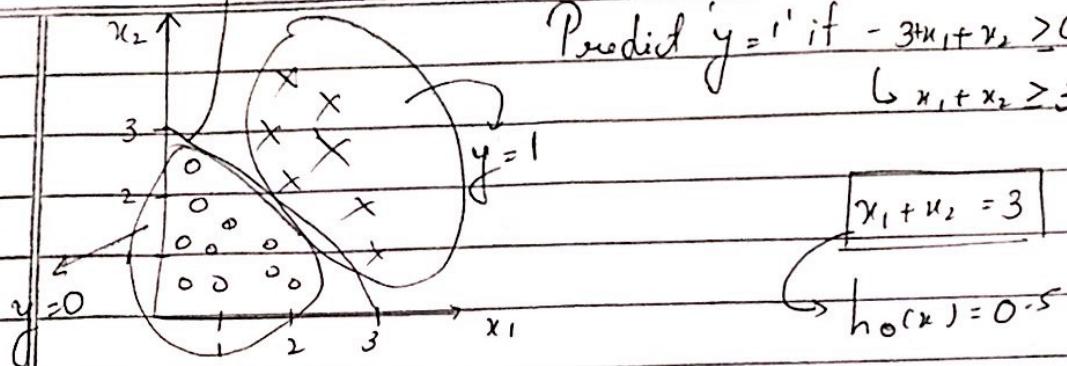
$$h_\theta(x) = \frac{1}{1 + e^{-\theta^T x}}$$

# Decision Boundary

Experiment:

Date \_\_\_\_\_

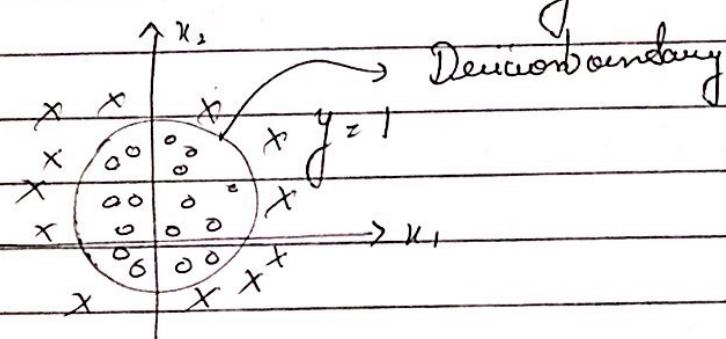
Page No. \_\_\_\_\_



$$h_0(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2)$$

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

## Non-linear decision boundary



$$h_0(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_2^2)$$

$$\begin{bmatrix} -1 \\ 0 \\ 0 \\ 1 \\ 1 \end{bmatrix}$$

Predict  $y = 1$  if  $-1 + x_1^2 + x_2^2 \geq 0$   
 $x_1^2 + x_2^2 \geq 1$

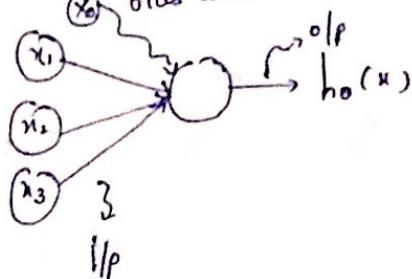
$$\frac{1}{1 + e^{-\theta^T x}}$$

Cost  $f^n$ :

$$\text{Cost}(h_0(x), y) = \frac{1}{2} ((h_0(x)) - y)^2 \quad h_0(x) \rightarrow 0 \quad (\text{cost} \rightarrow \infty)$$

$$\text{Cost}(h_0(x), y) = \begin{cases} -\log(h_0(x)) & -y = 1 \\ -\log(1 - h_0(x)) & y = 0 \end{cases}$$

Neuron model: Logistic unit if  $n/w$  has  $s_j$  units in layer  $j$ ,  
 bias unit  $x_0$   
 $s_{j+1}$  units in layer  $j+1$ , then  
 $\Theta^{(j)} - \text{dim} = s_{j+1} \times (s_j + 1)$

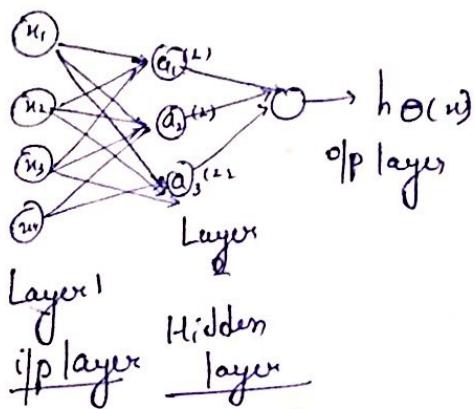


$$h_0(u) = \frac{1}{1 + e^{-\theta u}}$$

Sigmoid (logistic) activation  $f^h$

$$x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \end{bmatrix} \quad \Theta = \begin{bmatrix} \Theta_0 \\ \Theta_1 \\ \Theta_2 \\ \Theta_3 \end{bmatrix}$$

↳ "weights"



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

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

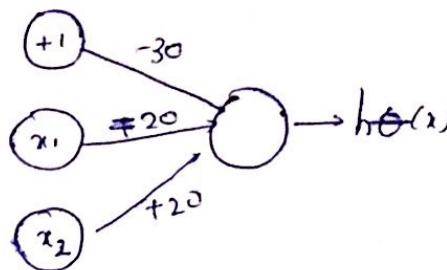
$$a_1^{(2)} = g(\Theta_{10}^{(1)} x_0 + \Theta_{11}^{(1)} x_1 + \Theta_{12}^{(1)} x_2 + \Theta_{13}^{(1)} x_3)$$

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

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

$$h_\theta(x) = o_1^{(3)} = g(\Theta_{10}^{(2)} a_0^{(2)} + \Theta_{11}^{(2)} a_1^{(2)} + \Theta_{12}^{(2)} a_2^{(2)} + \Theta_{13}^{(2)} a_3^{(2)})$$

AND  
 $x_1, x_2 \in \{0, 1\}$   
 $y = x_1 \wedge x_2$



$$h_\theta(x) = g(-30 + 20x_1 + 20x_2)$$

| $x_1$ | $x_2$ | $h_\theta(u)$      |
|-------|-------|--------------------|
| 0     | 0     | $g(-30) \approx 0$ |
| 0     | 1     | $g(-10) \approx 0$ |
| 1     | 0     | $g(10) \approx 0$  |
| 1     | 1     | $g(10) \approx 1$  |

| $x_1$ | $h_\theta(u)$      |
|-------|--------------------|
| 0     | $g(10) \approx 1$  |
| 1     | $g(-10) \approx 0$ |

Multiple o/p units

$$h_\theta(u) \approx \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} - \text{Fridge}$$

$$h_\theta(u) \approx \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} - \text{Car}$$

$$h_\theta(u) \approx \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} - \text{Motorcycle}$$

Cost  $f^h$  for logistic Reg

Vs Neural Net w/o

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

$$\frac{1}{k} \sum_{i=1}^k \left[ \log \left( \frac{1}{1 + e^{-\Theta^\top x^{(i)}}} \right) + (1 - \frac{1}{1 + e^{-\Theta^\top x^{(i)}}}) \log \left( 1 - \frac{1}{1 + e^{-\Theta^\top x^{(i)}}} \right) \right]$$

6. Artificial

Logistic regression cost fn:

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

Cost( $h_\theta(x)$ ,  $y$ )

$$= -y \log(h_\theta(x)) - (1-y) \log(1-h_\theta(x))$$

Put  $y = 1$ ,  $y = 0$

to get respective cost

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

Want: minimize  $\bar{J}(\theta)$

Repeat {

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

↓ [Simultaneous update]

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

Optimization algo

- Gradient descent
- Conjugate gradient
- BFGS
- L-BFGS

Adv. - No need to pick  $\alpha$  manually

- Often faster than GD

Disadv

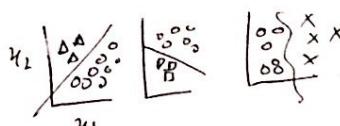
- complex

Multiclass Classification

Email: Work, Friends, Farm, Hobby

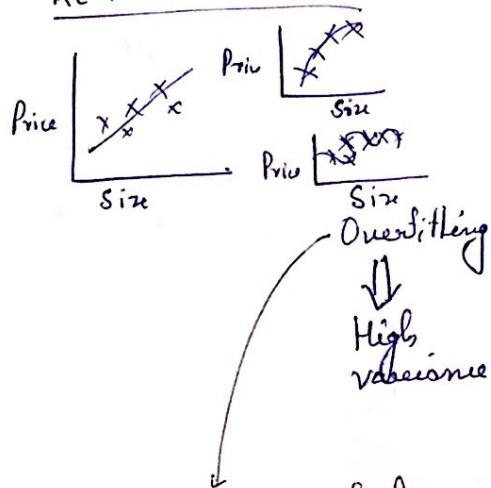
Weather: Sunny, Cloudy, Rainy

One Vs all



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

REGULARIZATION



Too many features  
→ fails to generalize new eg

Soln: to Reduce features

- Manually select specific one
- Model selection algo.

2. Regularization

- Reduce mag. of param  $\theta_j$   
(Keep all features)  $\theta_1, \theta_2, \theta_3, \dots$  if all  $\theta_j \neq 0$

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

$\theta_1, \theta_2, \theta_3, \dots$   
lost f<sup>n</sup> after regularize

Gradient Descent - linear Reg.

Repeat {

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

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

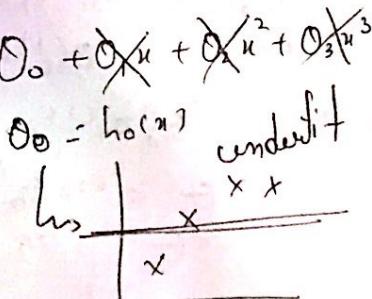
$$\theta_j := \theta_j \left( 1 - \alpha \frac{1}{m} \right) - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

→ 0.99

logistic Reg.

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

State of Art



## Precision/Recall

| Actual          | Class 0    |           |
|-----------------|------------|-----------|
| Predicted Class | True +ve   | False +ve |
| 0 - ve          | False - ve | True - ve |
| True +ve        | False - ve | True - ve |

$$\text{Precision} = \frac{\text{True} + \text{ve}}{\text{True} + \text{ve} + \text{False} + \text{ve}}$$

$$\text{Recall} = \frac{\text{True} + \text{ve}}{\text{True} + \text{ve} + \text{False} - \text{ve}}$$

of all, what fraction actually has cancer.

→ of all having cancer,  
how many correctly detected?

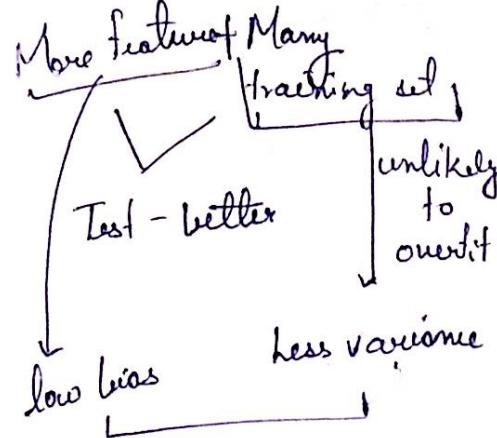
| Algo (P) | (R) | Avg. |
|----------|-----|------|
| 0.5      | 0.4 | 0.45 |
| 0.7      | 0.1 | 0.40 |
| 0.02     | 1.0 | 0.51 |

$$\text{Avg.} = \frac{P+R}{2}$$

$\bar{F}_1$  Score

$$2 \frac{PR}{P+R}$$

$D \leq F \leq 1$  gives compactness thing



$$\text{cost}_1(Z)$$

$$-\log \frac{1}{1+e^{-Z}}$$

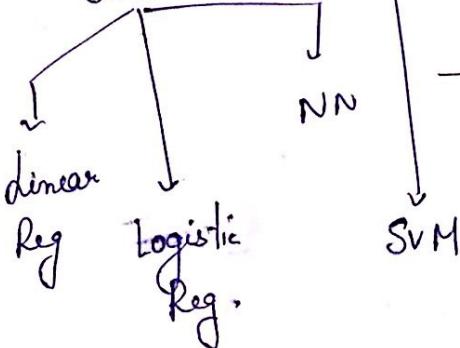


$$\text{cost}_0(Z)$$

$$-\log \left(1 - \frac{1}{1+e^{-Z}}\right)$$

Sliding window detection.  
Increase data set  
Distortions introduced!

## SUPERVISED



## Logistic Reg.

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

So, SVMs

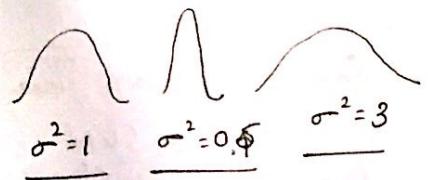
$$\min_{\theta} \frac{1}{m} \left[ \sum_{i=1}^m y^{(i)} (\text{cost}_1(\theta^T x^{(i)}) + (1-y^{(i)}) \text{cost}_0(\theta^T x^{(i)}) \right] + \frac{1}{2} \sum_{j=1}^n \theta_j^2$$

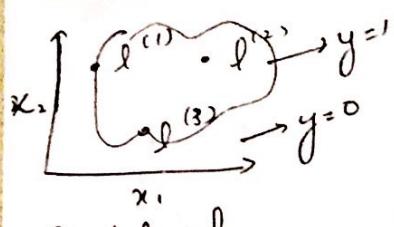
If  $y=1$ ,  $\theta^T x \geq 1$

(not just  $\geq 0$ )

If  $y=0$ ,  $\theta^T x \leq -1$  (not just  $\leq 0$ )

Kernels, ' $\sigma$ '





Predict when

$$\theta_0 + \theta_1 f_1 + \theta_2 f_2 + \theta_3 f_3 \geq 0$$

$$\theta_0 = -0.5$$

$$\theta_1 = 1$$

$$\theta_2 = 1$$

$$\theta_3 = 0$$

$$\begin{aligned} & \rightarrow \theta_0 + \theta_1 x_1 + \theta_2 x_0 + \theta_3 x_0 \\ & = -0.5 + 1 = 0.5 \geq 0 \end{aligned}$$

$$f_1, f_2, f_3 \approx 0$$

$$\theta_0 + \theta_1 f_1 + \dots \leq 0$$

Where to get  $l^{(1)}, l^{(2)}, l^{(3)}, \dots$ ?

Given  $x$ :

$$\begin{aligned} f_i &= \text{similarity}(x, l^{(i)}) \\ &= \exp\left(-\frac{\|x - l^{(i)}\|^2}{2\sigma^2}\right) \end{aligned}$$

$\sigma^2 \nearrow$ , bias  $\nearrow$ , variance  $\nearrow$   
 $\sigma^2 \searrow$ , bias  $\searrow$ , var  $\searrow$

$\left(\frac{1}{\sigma^2}\right)$   $\rightarrow$  Large C  
 $\left(\frac{1}{\sigma^2}\right)$   $\rightarrow$  Small C

Perform feature scaling b4 Gaussian Kernel

Satisfy Mercer's Theo  
to make some SVM pkgs optimize better.

SVM with kernels

Given:  $(x^{(1)}, y^{(1)})$

$(x^{(2)}, y^{(2)})$ ,

...

$(x^{(m)}, y^{(m)})$

(choose:  $l^{(1)} = x^{(1)}$

$l^{(2)} = x^{(2)}, \dots,$

$l^{(m)} = x^{(m)}$

Given eg  $x$ :

$f_1 = \text{similarity}(x, l^{(1)})$

$f_2 = \text{similarity}(x, l^{(2)})$

$\vdots$   
 $f_m$

Given  $x$ , compute features

$f \in \mathbb{R}^{m+1}$

Predict "y=1" if  $\theta^T f \geq 0$

BT PCA, data biplotting

Off-the-shelf kernels -

1.) Polynomial kernel

2.) More exotic:

Steering kernel

Chi-sq. kernel

Histogram intersection kernel

if  $n \gg m$  train eg.  
 $10K \geq 1-1K$   
 use Logistic Reg or  
 SVM w/o kernel

if  $n$  is small  
 $(1-1000)$   
 $m$  intermediate  
 $(10-10000)$

Use SVM with Gaussian

$(1-1K)$  kernel  
 $(50K)$

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

Create/add more features, then use logistic reg / linear sum

NN will work well, slow to train though.

Linear kernel = No kernel

Gaussian kernel

$$f_i = \exp\left(-\frac{\|x - l^{(i)}\|^2}{2\sigma^2}\right)$$

where  $l^{(i)} = x^{(i)}$

# UNSUPERVISED LEARNING

Find patterns in data

## K-means algo

Input:

- K (no. of clusters)

- Training set

$$\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$$

$$x^{(i)} \in \mathbb{R}^n \text{ (drop } x_0 = 1)$$

[unlabeled] convention

- Randomly initialize K.

cluster centroids  $\mu_1, \mu_2, \dots, \mu_k \in \mathbb{R}^n$  Minimize  $J(\dots)$  wrt  $c^{(1)}, c^{(2)}, \dots, c^{(n)}$

- Repeat {

for  $i = 1$  to  $m$

$c^{(i)} :=$  index (from 1 to  $k$ ) of cluster centroid closest to  $x^{(i)}$

for  $K=1$  to  $k$

$\mu_k :=$  avg (mean) of pts assigned to cluster  $k$

}

## Random initialization

For  $i = 1$  to  $100$  {  $\rightarrow 50 - 1000$

Randomly initialize k-means.

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

(compute cost  $J$  (distortion))

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

Pick clustering that gave lowest cost  $J(c^{(1)}, \dots, \mu_1, \dots, \mu_k)$

Choose value of  $k$  :-

Elbow method

See the purpose

Batch GD vs Stochastic GD

DIMENSIONALITY REDUCTION

$$2D \rightarrow 1D$$

$$3D \rightarrow 1D$$

$$1000D \rightarrow 100D$$

PCA ≠ linear Reg.

BT PCA, data preprocessing is must!

Reduce data from  $n$ -dim to  $k$ -dim.

Compute "covariance matrix"

$$\Sigma = \frac{1}{m} \sum_{i=1}^m (x^{(i)}) (x^{(i)})^T$$

(br. mat.)  $\underbrace{(n \times 1) \quad (1 \times n)}_{n \times n}$

Train set:  $x^{(1)}, x^{(2)}, \dots, x^{(n)}$

Preprocessing (feature scaling / mean normalization)

$$\mu_j = \frac{1}{m} \sum_{i=1}^m x_j^{(i)}$$

Replace each  $x_j^{(i)}$  with  $x_j^{(i)} - \mu_j$

If diff. features on diff. scales, scale features to have comparable range of values.

Compute eigen vectors of mat  $\Sigma$ :

$$[U, S, V] = SVD(Sigma)$$

$n \times n$  matrix

new centroid  
 $\downarrow$  Minimise  $J(\dots)$  wrt  $\mu_1, \dots, \mu_k$   
(holding  $c^{(1)}, \dots, c^{(n)}$  fixed)

Linear reg

vertical line

PCA

line

Choosing  $K$  in PCA  
(no. of principal components)

Reconstruction from compressed rep.

Avg. squared proj. error:

Total variation in data:

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

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

Choose  $k$  such that -

$$\frac{1}{m} \sum_{i=1}^m \|x^{(i)} - x_{\text{approx}}^{(i)}\|^2 \leq 0.01$$

$$\frac{1}{m} \sum_{i=1}^m \|x^{(i)}\|^2 \leq 0.05$$

gg. % of variance retained

g S.f. < " "

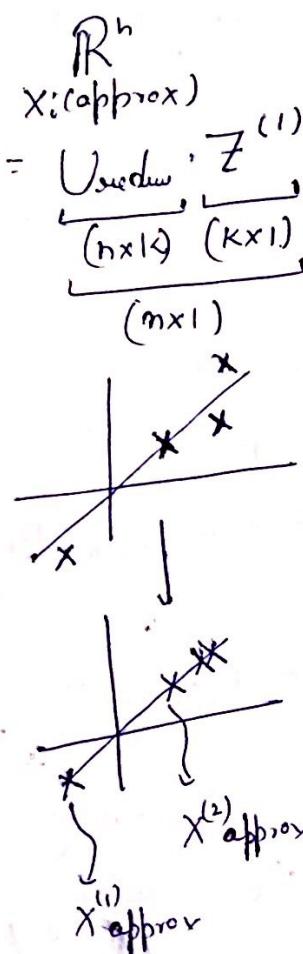
Algo:

Try PCA with  $k=1, 3, 4, \dots, 17$

Compute  $U_{\text{reduce}}, z^{(1)}, z^{(2)}, \dots, z^{(m)}, x_{\text{approx}}^{(1)}, \dots, x_{\text{approx}}^{(m)}$

Check if -

$$\left( \frac{1}{m} \sum_{i=1}^m \|x^{(i)} - x_{\text{approx}}^{(i)}\|^2 \leq 0.01? \right) \rightarrow \frac{1 - \frac{\sum_{i=1}^K S_{ii}}{\sum_{i=1}^n S_{ii}}}{\frac{\sum_{i=1}^n S_{ii}}{\sum_{i=1}^n S_{ii}}} \leq 0.01$$



$$[U, S, V] = \text{SVD}(\text{Sigma})$$

$$S = \begin{bmatrix} S_{11} & & \\ & S_{22} & \\ & & S_{33} \\ & & & \ddots \\ & & & & S_{nn} \end{bmatrix} \xrightarrow{k=3}$$

For given  $k$ ,

$$\frac{1 - \frac{\sum_{i=1}^K S_{ii}}{\sum_{i=1}^n S_{ii}}}{\frac{\sum_{i=1}^n S_{ii}}{\sum_{i=1}^n S_{ii}}} \geq 0.99$$

Underfitting Mapping

$x^{(i)} \rightarrow z^{(i)}$  can't be applied in cv & test sets,

only on training set

Bad use of PCA

- To prevent overfitting

Use  $z^{(i)}$  replacing  $x^{(i)}$  to reduce no. of features  $k < n$ ,

fewer features, less likely to overfit,

↓  
Use regularization instead

ANOMALY DETECTION:-

• Fraud Detection

→  $x^{(i)}$  - features of user i's activities

→ Model  $p(x)$  from data

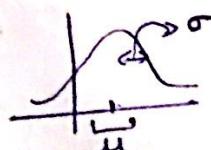
→ Identify unusual cases by checking  $p(x) \in \epsilon$

GAUSSIAN DIST.

If  $x$  is a dist. Gaussian with mean  $\mu$ , variance  $\sigma^2$

$$x \sim \mathcal{N}(\mu, \sigma^2)$$

dist. as



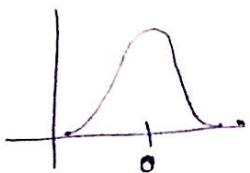
## Gaussian Dist

$$\frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$

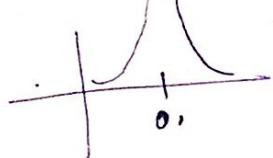
$$p(x; \mu, \sigma^2)$$

$$\sigma = \text{SD.}$$

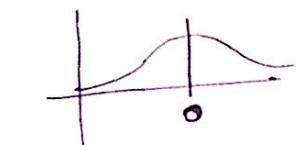
$$\mu=0, \sigma=1$$



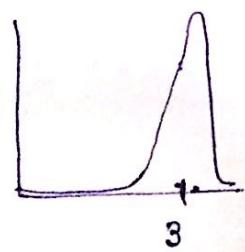
$$\mu=0, \sigma=0.5$$



$$\mu=0, \sigma=2$$



$$\mu=3, \sigma=0.5$$



$$\mu = \frac{1}{m}$$

Anomaly if  $p(x) < \epsilon$

$$\mu = \frac{1}{m} \sum_{i=1}^m x^{(i)}$$

↓ avg

$$\sigma^2 = \frac{1}{m} \sum_{i=1}^m (x^{(i)} - \mu)^2$$

↓ SD

Dataset:  $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$

## Density estimation

Training Set:

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

$$p(x) \quad \downarrow \text{feature}$$

$$= p(x_1; \mu_1, \sigma_1^2)$$

$$p(x_2; \mu_2, \sigma_2^2)$$

$$p(x_3; \mu_3, \sigma_3^2)$$

⋮

$$p(x_n; \mu_n, \sigma_n^2)$$

$$= \prod_{j=1}^n p(x_j; \mu_j, \sigma_j^2)$$

## Anomaly detection algo:

1. Choose features  $x_i$ , indicating anomalous eg  $\rightarrow$  Fraud detection

2. Fit param  $\mu_1, \mu_2, \dots, \mu_n, \sigma_1^2, \dots, \sigma_n^2 \rightarrow$  Mfg  $\rightarrow$  Monitor machines in data center

$$\mu_j = \frac{1}{m} \sum_{i=1}^m x_j^{(i)}$$

$$\sigma_j^2 = \frac{1}{m} \sum_{i=1}^m (x_j^{(i)} - \mu_j)^2$$

3. Give new eg  $x$ , compute

$$p(x) = \prod_{j=1}^n p(x_j; \mu_j, \sigma_j^2)$$

$$= \prod_{j=1}^n \frac{1}{\sqrt{2\pi\sigma_j^2}} \exp\left(-\frac{(x_j - \mu_j)^2}{2\sigma_j^2}\right)$$

## Algo evaluation:

→ fit model  $p(x)$  on training set  $\{x^{(1)}, \dots, x^{(m)}\}$   
→ On a cv / test eg  $x$ , predict

$$y = \begin{cases} 1 & \text{if } p(x) < \epsilon \\ 0 & \text{if } p(x) \geq \epsilon \end{cases}$$

anomaly ↴

normal

- Possible evaluation metrics:

- True +ve, False +ve, False -ve, True -ve

- Precision | Recall

- F1-score

## Anomaly Vs Supervised Det.

Learning

→ very small no. of true eg ( $y=1$ )  $\rightarrow$  large no. of false -ve eg.

→ large no. of -ve ( $y=0$ ) eg.

⇒ Email spam classification

⇒ Weather pred

⇒ Cancer classification

## Multivariate Gaussian Dist

Vary  $\mu, \Sigma$ , as per your case

Parameters  $\mu, \Sigma$

$$p(x; \mu, \Sigma)$$

$$= \frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}}$$

$$\exp\left(-\frac{1}{2}(x-\mu)^T \Sigma^{-1} (x-\mu)\right)$$

Fit param:

Given training set  
 $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$

$$\mu = \frac{1}{m} \sum_{i=1}^m x^{(i)}$$

$$\Sigma = \frac{1}{m} \sum_{i=1}^m (x^{(i)} - \mu)(x^{(i)} - \mu)^T$$

Algo -

1. Fit model  $p(x)$  by setting  $\mu, \Sigma$  as above

2. Given a new eg  $x$ ,

compute

$$p(x) = \frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2}(x-\mu)^T \Sigma^{-1} (x-\mu)\right)$$

Flag anomaly if  $p(x) < \epsilon$

## Original Model Vs Multivariate Gaussian

$$p(x_1; \mu_1, \sigma_1^2) \times \dots \times p(x_n; \mu_n, \sigma_n^2)$$

$$p(x; \mu, \Sigma) = \frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}}$$

$$\exp\left(-\frac{1}{2}(x-\mu)^T \Sigma^{-1} (x-\mu)\right)$$

Manually create features to capture anomalies

Computationally cheaper

Automatically does so

Computationally expensive

## Recommender Sys

Content based

Collaborative