

ML (weak AI)  
ML (weak AI)

Supervised  
eg: Face Detection  
Identification  
of numbers, animals  
etc.

Unsupervised  
eg: classify photos  
w/o focus.  
without explicit labels  
E-commerce - classify  
customers in groups

Reinforcement  
Learning to perform  
actions in an environment  
eg: Self-driving cars,  
Computer Games,  
Robotics.

### Measuring AI - Turing Test

ML initially started as a quest for AI  
& mimicking brain dynamics,  
but has now evolved to solving practical  
problems using algos which have  
little or no resemblance to brains.

Debate: AI ⊂ ML or ML ⊂ AI  
or AI & ML just have an overlap.

Goal of ML → Generalisation of learning to new &  
unknown datasets & scenarios (AI).  
Bias-Variance Tradeoff / Decomposition → pre-trained  
best way to estimate generalisation error.  
after requires domain knowledge

Bias: wrong assumption about training data (underfitting)  
Variance: High sensitivity to small fluctuations &  
random noise.  
(~~too~~ overfitting)

$$y = f(x) + \epsilon$$

$f$ : Learned function  
 $\epsilon$ : Actual function

$$\text{Bias} = E_{\theta}[f] - f$$

$$\text{Variance} = E_{\theta}\left[\left(E_{\theta}[f] - f\right)^2\right]$$

$\epsilon$ : zero mean  
 $\sigma^2$  variance.

Expected Error on unseen sample:

$$\text{Bias}^2 + \text{Var} + \underbrace{\text{Var of } \epsilon}_{\text{Var of } y}$$

# Random Variables & Probability Distributions

Variable : whose value can change.

Deterministic  
[governed by deterministic equations]

Random  
[Value is not pre-determined]  
can change in unpredictable ways but has a probability distribution.

Discrete

Discrete R.V:

Coin Toss

$$x = \begin{cases} 1 & \text{Heads} \\ 0 & \text{Tails} \end{cases}$$

$$P(X=0) = q = 1 - p$$

Bernoulli Distribution  $\leftarrow \begin{cases} P(X=1) = p \\ 0 \leq p \leq 1 \end{cases}$

Dice Roll

Binomial Distribution : N coin tosses

WRITE CODE  
to plot this  
for large 'n'.

$$P_r(X=k) = f(k, n, p) = {}^n C_k p^k (1-p)^{n-k}$$

Poisson Distribution :  $f(k; \lambda) =$

$$E[X] = \lambda = \text{Var}[X]$$

$\frac{\lambda^k e^{-\lambda}}{k!}$   
 $= P_r(X=k)$  Number of events in a fixed time interval

$$\Omega_x = \{x_1, x_2, x_3, \dots, x_n\}$$

$$0 \leq P(X=x_i) \leq 1$$

$$\sum_{i=1}^n P(X=x_i) = 1$$

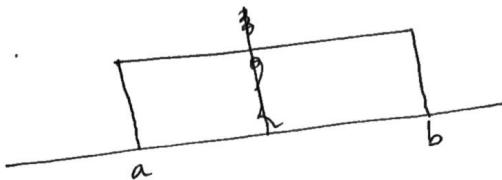
## Continuous R.V.

Any measuring device has finite precision & range

when outcome can be any real number [uncountable].

probability of  $X$  taking any one value is zero.  
Practical: when no. of possibilities is very large

### Uniform Distribution:



Distribution  $f(x) = \frac{1}{b-a}$

$$P(X=x) = 0$$

$$P(x \leq X \leq x + \Delta x) = f(x) \Delta x, \text{ if } \Delta x \approx 0$$

$$\int_{-\infty}^{\infty} f(x) dx = 1$$

$$f(x) \geq 0$$

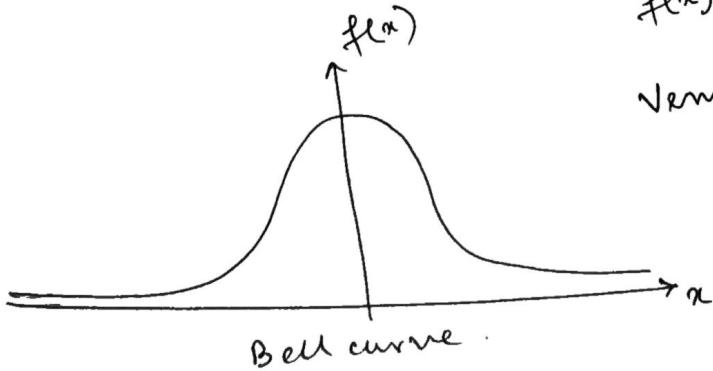
$\Rightarrow f(x) \geq 0$   
 $f(x)$  can be greater than 1 over a finite range of  $x$  values.

### Gaussian:

Velocity of gas molecules  
Noise in communication channels.  
Height of a large population.

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

Very crucial in ML too!!



## Expectation, Moments & CLT

Expectation:  $E[g(x)] = \int_{-\infty}^{\infty} g(x) f(x) dx$

Mean:  $E[x] = \int_{-\infty}^{\infty} x f(x) dx = \mu$

Variance:  $\sigma^2 = E[(x - \mu)^2] = E[x^2] - E[x]^2$

Moments:  $m_n = E[x^n] = \int_{-\infty}^{\infty} x^n f(x) dx$

Characteristic function:  $\Phi_x(\omega) = \int_{-\infty}^{\infty} f(x) e^{j\omega x} dx = E[e^{j\omega x}]$

Generating function:  $M_x(s) = \int_{-\infty}^{\infty} f(x) e^{sx} dx = E[e^{sx}]$

Moment Generating Function (MGF):  $M_x^{(n)}(0) = E[x^n] = m_n$

If all moments are finite & the following series converges absolutely, then  $M_x(s) = \sum_{n=0}^{\infty} \frac{m_n}{n!} s^n$  [does not work for log-normal distribution]

If  $X$  &  $Y$  are independent random variables,

$$\begin{aligned}\Phi_{X+Y}(\omega) &= E[e^{j\omega(X+Y)}] \\ &= E[e^{j\omega X}] E[e^{j\omega Y}] = \Phi_X(\omega) \Phi_Y(\omega)\end{aligned}$$

## MGF of Normal Distribution

$$f(x) = \frac{1}{\sqrt{2\pi}} e^{-x^2/2} \quad \int_{-\infty}^{\infty} f(x) dx = 1$$

$$\begin{aligned} M(s) &= \mathbb{E}[e^{sx}] = \int_{-\infty}^{\infty} e^{sx} \cdot f(x) dx \\ &= \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} e^{sx} e^{-x^2/2} dx \\ &= \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} e^{-(x-s)^2/2} e^{s^2/2} dx \\ &= e^{s^2/2} \cdot \left[ \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} e^{-(x-s)^2/2} dx \right] = e^{s^2/2} \end{aligned}$$

## Central Limit Theorem

$$\bar{x}_n = \frac{x_1 + x_2 + \dots + x_n}{n} \quad \{x_i\} \text{ i.i.d. rv}$$

$$\text{Let } y_i = x_i - \mu \quad \& \quad z_n = \frac{y_1 + y_2 + \dots + y_n}{\sqrt{n}\sigma^2}$$

$$M_y(s) = 1 + m_1 s + \frac{m_2}{2!} s^2 + \dots = 1 + \frac{\sigma^2}{2} s^2 + \dots$$

$$\begin{aligned} M_{z_n}(s) &= \left[ M_y\left(\frac{s}{\sqrt{n}\sigma^2}\right) \right]^n = \left[ 1 + \frac{\sigma^2}{2} \left( \frac{s}{\sqrt{n}\sigma^2} \right)^2 + \dots \right]^n \\ &= \left[ 1 + \frac{s^2}{2n} + \dots \right]^n \sim e^{s^2/2} \quad \text{for large } n \end{aligned}$$

## BAYES' THEOREM

$V = [A_1, A_2, \dots, A_n]$  is a partition of  $S$ , meaning  $A_i \neq A_j$   
 are mutually exclusive  
 &  $\sum_i p(A_i) = 1$

Let  $B$  be an arbitrary event:

$$\begin{aligned} B &= B \cap S = B \cap [A_1 \cup A_2 \cup A_3 \dots \cup A_n] \\ &= (B \cap A_1) \cup (B \cap A_2) \dots \cup (B \cap A_n) \end{aligned}$$

$B \cap A_i$  &  $B \cap A_j$  are  
 mutually exclusive  $\therefore P(B) = P(B \cap A_1) + P(B \cap A_2) + \dots + P(B \cap A_n)$

Total probability theorem:  $P(B \cap A_i) = P(B|A_i) P(A_i)$   
 $= P(A_i|B) P(B)$

$$\Rightarrow P(A_i|B) = \frac{P(B|A_i) P(A_i)}{P(B)}$$

$$= \frac{P(B|A_i) P(A_i)}{\sum_i P(B|A_i) P(A_i)}$$

BAYES' THEOREM

Example: A test for ~~Cancer~~ is known to be ~~95%~~  $95\%$  accurate.  
 A person submits to test & is found positive. Let the person  
 belong to a population of 1L people out of which 2000  
 actually have the disease. What's the prob. of the  
 person being actually infected?

$$P(T_+|c) = 0.95 \quad P(T_-|H) = 0.05$$

$$P(H) = 0.98 \quad P(c) = 0.02$$

$$P(c|T_+) = \frac{P(T_+|c) P(c)}{P(T_+)} = \underline{\underline{0.278}}$$

$$P(T_+) = P(T_+|c) P(c) + P(T_+|H) P(H)$$

Example :  $P(D|B_1) = \frac{100}{2000} = 0.05$   $P(D|B_2) = \frac{200}{500} = 0.4$

Papoulis  
page 3A

$$P(D|B_3) = \frac{100}{1000} = 0.1 \quad P(D|B_4) = \frac{100}{1000} = 0.1$$

$$P(B_1) = P(B_2) = P(B_3) = P(B_4) = 1/4$$

$$P(D) = 0.1625$$

$$P(B_2|D) = \frac{0.4 \times 0.25}{0.1625} = 0.615$$

Example  
Papoulis  
page 104

in coin tosses  $k$  heads  
Prob of head in  $(n+1)$ th toss?

$$\text{Ans: } \frac{k+1}{n+2}$$

# Naive Bayes & Gaussian Naive Bayes

## Bayer's Theorem

$\cup = [A_1, A_2, \dots, A_n]$  is a partition of  $S$   
 &  $B$  is an arbitrary event.

$$P(A_i/B) = \frac{P(B/A_i) P(A_i)}{\sum_i P(B/A_i) P(A_i)} = \frac{P(B/A_i) P(A_i)}{P(B)}$$

Dataset of 1L emails with 70% Normal & 30% Spam.  
 find total count of each word occurring in each category  
 [BAG OF WORDS model  
 $w_1 = \text{work}$

# words Two ways to calculate  $P(w_i/N) = \frac{22}{100} = 0.22$   
 # docs to calculate

$$P(w_2/N) = \frac{26}{100} = 0.26$$

$$P(w_3/N) = \frac{19}{100} = 0.19$$

$$P(w_4/N) = \frac{13}{100} = 0.13$$

$$P(w_5/N) = \frac{8}{100} = 0.08$$

$$P(w_6/N) = \frac{12}{100} = 0.12$$

$$P(w_1/S) = \frac{11}{100} = 0.11$$

$$P(w_2/S) = \frac{4}{100} = 0.04$$

$$P(w_3/S) = \frac{14}{100} = 0.14$$

$$P(w_4/S) = \frac{31}{100} = 0.31$$

$$P(w_5/S) = \frac{21}{100} = 0.21$$

$$P(w_6/S) = \frac{19}{100} = 0.19$$

$$w_2 = \text{deadline}$$

$$w_3 = \text{hand}$$

$$w_4 = \text{bonus}$$

$$w_5 = \text{luxury}$$

$$w_6 = \text{Money}$$

$$\text{Total } 100$$

$$\text{words in each category}$$

New email received with words  $w_3, w_5$  &  $w_6$  [in any order]  
Naive  $\equiv$  [we also assume word occurrence to be independent events]

$$P(N/w_3 \cap w_5 \cap w_6)$$

High bias - ignores word order  
 Low variance  
 $\hookrightarrow$  works well

$$= \frac{P(w_3 \cap w_5 \cap w_6 / N) P(N)}{P(w_3 \cap w_5 \cap w_6)}$$

$$= \frac{P(w_3/N) P(w_5/N) P(w_6/N) P(N)}{P(w_3 \cap w_5 \cap w_6)}$$

$$\propto 0.19 \times 0.08 \times 0.12 \times \frac{0.7}{0.3} = \frac{0.0012768}{0.000912}$$

$$P(S/w_3 \cap w_5 \cap w_6) = \frac{P(w_3/S) P(w_5/S) P(w_6/S) P(S)}{P(w_3 \cap w_5 \cap w_6)}$$

$$\propto 0.14 \times 0.21 \times 0.19 \times \frac{0.3}{0.5} = \frac{0.0016758}{0.0002793}$$

[Take log since numbers get very small under flow]

## Gaussian Naive Bayes

Estimating  $p(x_i | c_j)$  using Gaussian Distribution  
used when ' $x_i$ ' values can be any real number  
in a range.

Eg: Boy vs. Girl classification using

height, weight & feet size.

Use mean & variance of available data  
to fit Gaussian.

fit a Gaussian  
using mean & SD.

Bernoulli Naive Bayes : use Bernoulli distribution

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

Categorical  
Naive Bayes

word present or absent

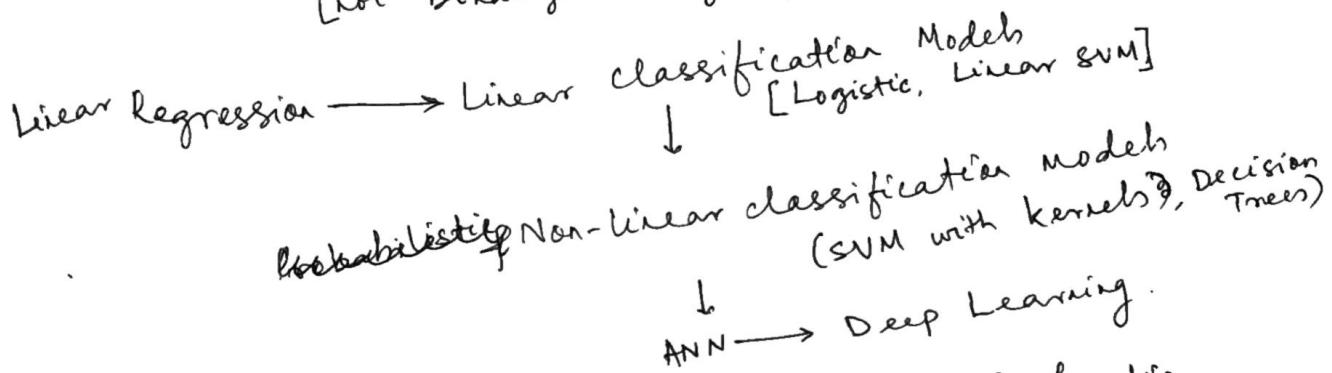
Multinomial Naive Bayes :  $x_i \in \{0, 1, 2, \dots, K\}$   
word count taken into account.

## Bayes Optimal Classifier

Model parameter dependencies leading to the  
most general classifier. Abstract notion of hard  
to implement in practice.

## Supervised Learning : Regression vs. Classification

ML is another form of ~~cost~~ curve fitting, or function approximation, but where the functional form of the curve is not explicitly obvious.  
 [not binary but gray scale]



Model encompassing a broad family of functions & changing parameters helps in navigating this set of functions.

$$\text{Learning} = \underbrace{\text{Training} + \text{Testing}}_{\text{Generalisation}}$$

Underfitting: High Bias  
 + Low Variance

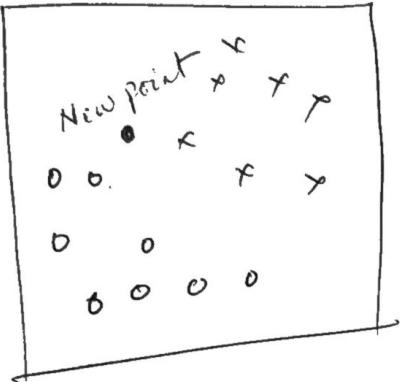
Overfitting: Low Bias + High Variance

→ reduce model complexity  
 use regularisation

	Training Error High	Testing Error High	Testing Error Low
Training Error High	Underfitting	Fluke	
Training Error Low	overfitting		Ideal Good Fit

- Naive Bayes & Bayes' optimal classifier
- Maximum Likelihood Estimator & Maximum Aposteriori Probability
- Linear Regression
- Logistic Regression
- Support Vector Machines
- Decision Trees
- K-nearest Neighbours
- ANNs.

## K- Nearest Neighbours



Measure distance from  
K-nearest neighbours  
& classify based on  
majority rule.

## Curse of Dimensionality

Consider a D-dimensional unit cube.

If we want to do our KNN estimate based on a fraction,  $f$ , of data points, we need to take a cube around our test point of edge size,

$$e_D(f) = f^{1/D}$$

$$\text{if } f = 1\% = 0.01, \quad e_D(f) = 0.215$$

$$D=3 \Rightarrow e_D(f) = 0.63 \rightarrow \text{very large cube!!}$$

$$D=10 \Rightarrow e_D(f) = 0.96!!$$

$$D=100 \Rightarrow e_D(f) = 0.96!!$$

## Blessing of Dimensionality

If Signal to Noise ratio is high.

High dimensional data can have very interesting & geometric properties which could be exploited for improving accuracy.

## KNN - Algo - pseudo- code

Accuracy -

$$\frac{TN + TP}{Total}$$

	Predicted	
	0	1
Actual	0	TN      FP
	1	FN      TP

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

$$F_1 = \frac{2}{\gamma_{\text{Precision}} + \gamma_{\text{Recall}}}$$

$$\text{Recall} = \frac{TP}{TP + FN} : \text{True Positive Rate}$$

(Sensitivity)

$$\text{Specificity} = \frac{TN}{TN + FP} : \text{True Negative Rate}$$

sensitivity

# BAYESIAN PARAMETER ESTIMATION

Event A : K Heads in n tosses

$p$ : probability of getting Head in one toss

Event A': Head in  $(n+1)$ th toss

Since we do not know the value of ' $p$ ' & are trying to estimate it, we need to consider it as a random variable.

Before event A, let's say we have no information about ' $p$ '

$$\text{So, assume } f(p) = \begin{cases} 1, & 0 \leq p \leq 1 \\ 0, & \text{otherwise} \end{cases}$$

$$P(A|p) = p^k (1-p)^{n-k}$$

$$f(p|A) = \frac{P(A|p) f(p)}{\int_0^1 P(A|p) f(p) dp} = \frac{P(A|p) f(p)}{\int_0^1 p^k (1-p)^{n-k} \cdot \frac{(n+1)!}{(n-k)! k!} dp}$$

$$= p^k (1-p)^{n-k} \cdot \frac{(n+1)!}{(n-k)! k!} \quad 0 \leq p \leq 1$$

$\sim \beta(\alpha, \beta)$   
Beta-distributed  
 $\alpha = n$   
 $\beta = k$  here.

$$P(A'|A) = \int_0^1 \underbrace{P(A'|p)}_{p} f(p|A) dp$$

$$= \frac{k+1}{n+2}$$

This is generally very complicated to estimate in most problems. Hence, we instead estimate the value of ' $\phi$ ' that best explains the observed data.

There are two ways of doing it.

MLE & MAP

MLE: Assume that you have no prior info about the parameters [uniform distribution]

MAP: Assume availability of prior info based on reasonable considerations.

also leads to "regularisation"  
which helps in preventing overfitting.

$(x, y)$  Training Data  $(\hat{x}, \hat{y})^+$

$\hat{x}$  new input data  
 $\hat{y}$  predicted output data

$$p(\hat{y} | x, y, \hat{x}) = \underbrace{\int p(\hat{y} | \phi, \hat{x})}_{\text{model}} \frac{p(\phi | x, y) d\phi}{p(y | x, \phi) p(\phi)}$$

$$\begin{aligned} p(\phi | x, y) &= \frac{p(y | x, \phi) p(\phi)}{p(y | x)} \\ &= \frac{p(y | x, \phi) p(\phi)}{\int p(y | x, \phi) p(\phi) d\phi} \end{aligned}$$

BAYESIAN PARAMETER ESTIMATION.

$$\text{MLE \& MAP : } p(\hat{y} | x, y, \hat{x}) = p(\hat{y} | \phi_E, \hat{x})$$

$$\phi_E = \arg \max_{\phi} p(\phi | x, y)$$

## MLE & Linear Regression

Let  $y_i \text{ follows } y_i = g(x_i, \theta) + \epsilon_i$  curve fitting

$$\therefore \theta_{MLE} = \arg \max \sum_i \log p(y_i | x_i, \theta)$$

$$= \arg \max \sum_i \log p(\epsilon_i)$$

$$= \arg \max \sum_i \log \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\epsilon_i^2/2\sigma^2} : \begin{matrix} \text{Assuming} \\ \text{Gaussian} \\ \text{Error/} \\ \text{Noise} \end{matrix}$$

$$= \arg \max \left[ \sum_i \log \frac{1}{\sqrt{2\pi\sigma^2}} - \frac{1}{2\sigma^2} \sum_i \epsilon_i^2 \right]$$

$$= \arg \max \arg \min \sum_i \epsilon_i^2$$

$$= \arg \min \sum_i [y_i - g(x_i, \theta)]^2$$

Least Squares Error  
Cost Function

## MLE & Classification

Predict output  $y \in \{0, 1\} \rightarrow$

using a probabilistic function

$$g(x, \theta) \text{ or } g(\theta) = P[y|g(x)]$$

Disease / Healthy  
Bank default / pay  
Fake news / factual  
could be Logistic Regression  
or ANN, etc.

$$\theta_{MLE} = \arg \max p(y_i | x_i, \theta)$$

$$= \arg \max \prod_i g_i^{y_i} (1-g_i)^{1-y_i}$$

$$= \arg \max \sum_i [y_i \log g_i + (1-y_i) \log (1-g_i)]$$

$$= \arg \min \underbrace{\sum_i y_i \log g_i + (1-y_i) \log (1-g_i)}_{\text{Binary Cross-Entropy}}$$

used in  
Logistic  
Regression

# MLE & MAP Maximum Likelihood Estimate

Naïve Bayes assumes independence of parameters which is generally not true in real world.

So, need to build models & estimate parameters.

$$f_{\theta}(\theta|x,y) = \frac{p(y|x,\theta)f_{\theta}(\theta)}{p(y|x)} \xrightarrow{\text{Likelihood}} \text{prior } p(\theta|x) = p(\theta)$$

posterior estimate

MLE  
Assume  $f_{\theta}(\theta)$  to have uniform distribution  
i.e. no prior knowledge of parameters available.

$$\begin{aligned} \therefore \theta_{MLE} &= \arg \max f_{\theta}(\theta|x,y) \\ &= \arg \max p(y|x,\theta) \xrightarrow{\text{L}} \\ &= \arg \max_i \prod p(y_i|x_i, \theta) \\ &= \arg \max_i \sum \log p(y_i|x_i, \theta) \xrightarrow{\text{Log likelihood}} \end{aligned}$$

Coin toss experiment  
 $\hat{p}$  = prob. of heads (1)      tails = 0

$$p(y_i=y_i|\hat{p}) = \hat{p}^{y_i}(1-\hat{p})^{1-y_i} \quad y_i \in \{1, 0\}$$

$\hookrightarrow$  Bernoulli dist.

$$\begin{aligned} LL &= \log p(y|\hat{p}) \\ &= \sum_{i=1}^n \log p(y_i|\hat{p}) = \sum_{i=1}^n \log \hat{p}^{y_i}(1-\hat{p})^{1-y_i} \\ &= \log \hat{p} \sum y_i + \log(1-\hat{p}) \sum (1-y_i) \\ &= n^{(1)} \log \hat{p} + n^{(0)} \log(1-\hat{p}). \end{aligned}$$

$$\frac{\partial LL}{\partial \hat{p}} = 0 \Rightarrow \frac{n^{(1)}}{\hat{p}} - \frac{n^{(0)}}{1-\hat{p}} = 0 \Rightarrow \frac{1}{\hat{p}} - 1 = \frac{n^{(0)}}{n^{(1)}} \Rightarrow \hat{p}_{MLE} = \frac{n^{(1)}}{n^{(1)} + n^{(0)}}$$

[Intuitive Estimate in MLE]

## MAP Estimation - Maximum a posteriori

$$f_p(\theta | y, x) = \frac{p(y|x, \theta) f(\theta|x)}{p(y|x)}$$

$$\begin{aligned}\theta_{\text{MAP}} &= \underset{\theta}{\operatorname{argmax}} \, f_p(\theta | y, x) \\ &= \underset{\theta}{\operatorname{argmax}} \, p(y|x, \theta) f(\theta) \\ &= \underset{\theta}{\operatorname{argmax}} \left[ \sum_i \log p(y_i|x_i, \theta) + f(\theta) \right]\end{aligned}$$

Usually clear info of  $p(\theta)$  is not available.  
so use  $f_p(\theta)$  to enforce certain constraints.  
Eg: choose  $f(\theta)$  such that there is low probability  
for  $\theta$  taking high values.

Now let's say we wish to make a prediction  
for new data point

$$p(\tilde{y}|\tilde{x}, \theta) = p(\tilde{y}|\tilde{x}, y, x) = p(\tilde{y}|\tilde{x}, \theta) f(\theta|y, x) d\theta$$

but actually,  $p(\tilde{y}|\tilde{x}, y, x) = \int p(\tilde{y}|\tilde{x}, \theta) f(\theta|y, x) d\theta$   
usually we estimate this just using  $\theta_{\text{ML}}$  or  $\theta_{\text{MAP}}$   
which is not same as RHS.  
But estimating RHS is very hard since  
it requires  $p(x)$  too.

$$\theta_{\text{MAP}} = \underset{\theta}{\operatorname{argmax}} \sum_{i=1}^n \log p$$

## MAP for Coin Toss

$$\hat{p}_{\text{MAP}} = \underset{\hat{p}}{\operatorname{argmax}} \sum_{i=1}^N \log p(y_i | \hat{p}) + \log p(\hat{p})$$

$$f = p(y_i = 1)$$

$$\begin{aligned} n^{(1)} &= \text{No. of } 1s \\ n^{(0)} &= \text{No. of } 0s \\ n^{(1)} + n^{(0)} &= n \end{aligned}$$

$$\begin{aligned} &= \underset{\hat{p}}{\operatorname{argmax}} \left[ n^{(1)} \log \hat{p} + n^{(0)} \log (1 - \hat{p}) \right. \\ &\quad \left. + \log \frac{\hat{p}^{\alpha-1} (1-\hat{p})^{\beta-1}}{B(\alpha, \beta)} \right] \\ &\quad \text{Beta-distribution} \\ &\quad \alpha, \beta > 0 \end{aligned}$$

$$= \underset{\hat{p}}{\operatorname{argmax}} \left[ n^{(1)} \log \hat{p} + n^{(0)} \log (1 - \hat{p}) \right]$$

$$+ n \log \alpha (\alpha-1) \log \hat{p} + \alpha (\beta-1) \log (1 - \hat{p}) \right]$$

$$= \underset{\hat{p}}{\operatorname{argmax}} \left[ \left( \frac{n^{(1)} + \alpha - 1}{n + \alpha + \beta - 2} \right) \log \hat{p} + \left( \frac{n^{(0)} + \beta - 1}{n + \alpha + \beta - 2} \right) \log (1 - \hat{p}) \right]$$

$$\alpha = 1, \beta = 1$$

$$\Rightarrow \hat{p}_{\text{MAP}} = \hat{p}_{\text{MLE}}$$

$$\hat{p} = \frac{n^{(1)} + \alpha - 1}{n + \alpha + \beta - 2}$$

$$\begin{aligned} \frac{\partial L}{\partial \hat{p}} &= 0 \Rightarrow \frac{n^{(1)} + \alpha - 1}{\hat{p}} - \frac{n^{(0)} + \beta - 1}{1 - \hat{p}} = 0 \\ &\Rightarrow \hat{p}^{-1} = \frac{n^{(0)} + \beta - 1}{n^{(1)} + \alpha - 1} \\ &\Rightarrow \hat{p}_{\text{MAP}} = \frac{n^{(0)} + \beta - 1}{n^{(1)} + \alpha - 1 + n^{(0)} + \beta - 2} \end{aligned}$$

$$p(\tilde{y} | y) = \frac{p(y | \tilde{y}) p(\tilde{y})}{\int_0^1 p(y | \tilde{y}) p(\tilde{y}) d\tilde{y}}$$

$$p(\tilde{y} | y) = \frac{p(y | \tilde{y}) p(\tilde{y})}{p(y)} = \frac{p(y | \tilde{y}) R(\tilde{y})}{\int_0^1 p(y | \tilde{y}) p(\tilde{y}) d\tilde{y}}$$

~~Bayesian Estimation~~

$$\begin{aligned} &= \frac{\left[ \prod_i p(y_i | \tilde{y}) \right] p(\tilde{y})}{\int_0^1 \left[ \prod_i p(y_i | \tilde{y}) \right] p(\tilde{y}) d\tilde{y}} \\ &\hookrightarrow \text{Assume Beta distribution} \\ &= \hat{p}^{n^{(1)}} (1 - \hat{p})^{n^{(0)}} \frac{\hat{p}^{\alpha-1} (1-\hat{p})^{\beta-1}}{B(n^{(1)} + \alpha, n^{(0)} + \beta)} \end{aligned}$$

# Supervised Learning: Math Details

[Gradient Descent  
+ Regularisation]

Data:  $(x_i, y_i)$

→ Gaussian Noise

$$y_i = f(x_i, \theta) + \epsilon_i$$

↑ parameters  
↓ features

$$\theta_{MLE} = \arg \min \sum_i [y_i - f(x_i, \theta)]^2 \rightarrow \text{Least Squares Error fn.}$$

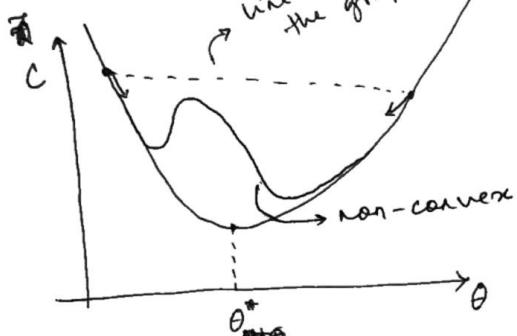
$C$ : cost function

$$\frac{\partial C}{\partial \theta} = 0 \quad \text{May not be analytically tractable.}$$

line is above the graph b/w the points

Robot star

$C(\theta)$  has to be convex.



Gradient Descent:

1. Pick a starting value of  $\theta$
2. Move on the graph in direction opposite to the gradient.

$$\theta \rightarrow \theta - \eta \frac{\partial C}{\partial \theta}$$

Learning rate

$\eta$  should neither be too high nor too low. → no convergence  
→ very slow convergence

why not  $\theta \rightarrow \theta - \eta \operatorname{sgn}\left\{\frac{\partial C}{\partial \theta}\right\}$ ?

if needed

$\frac{\partial C}{\partial \theta}$  can be estimated using finite difference or Automatic Differentiation (i.e. backpropagation in ANN)

As we approach  $\theta^*$ , we want change in  $\theta$  to become smaller & smaller.

if  $y_i, x_i \neq \theta$  are ~~orthogonal~~ vectors.

$$C = \sum_{i=1}^n [y_i - f(x_i, \theta)]^2 + \frac{\lambda}{2} \|\theta\|^2 ; \quad \|\cdot\| : \text{Norm of vector}$$

(see next page)

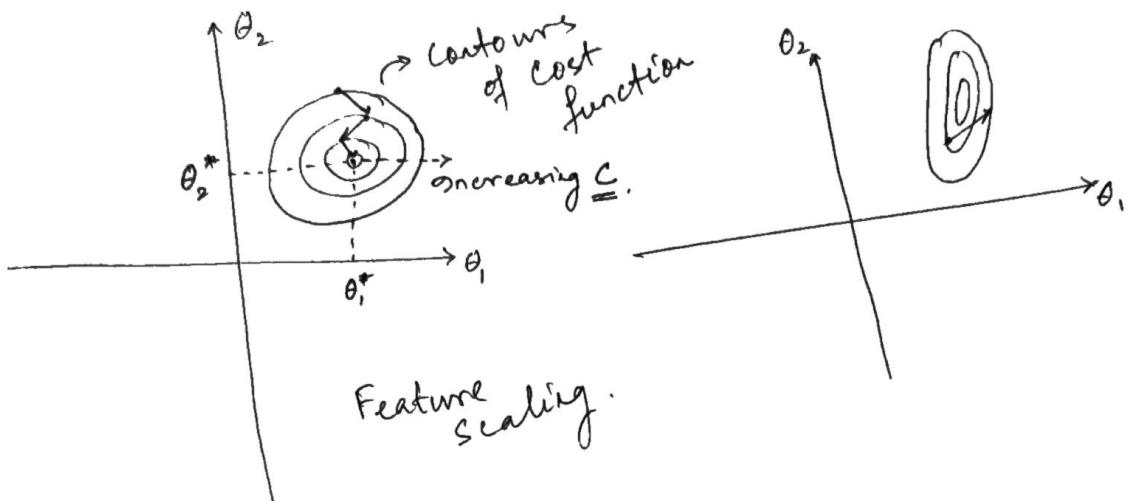
$\lambda = 0 \Rightarrow$  overfitting  
 $\lambda \approx 0 \Rightarrow$  underfitting

[Example in scikit-learn website]

Regularisation prior information prevents parameters from becoming too large.

[Ridge regression]

Lasso: L1 norm



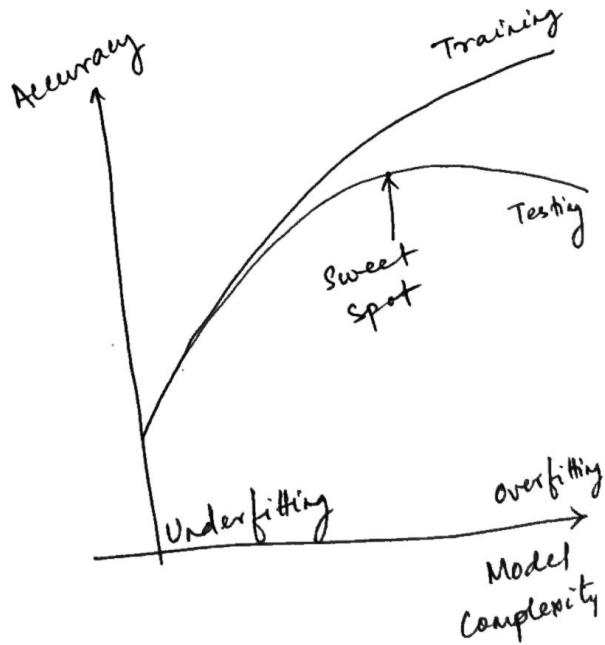
### Basic Steps of Supervised Learning :

- ① choose appropriate model, i.e.  $f(x_i, \theta)$  &  $C(\theta)$
- ② Divide data into training & testing set.
- ③ Train model parameters using gradient descent or some other optimisation method.
- ④ Stopping criteria for g.d.
- ⑤ Find training & testing accuracy.
- ⑥ If accuracy not adequate, go to step ①.

- Fix a desired precision
 
$$\|\theta^{k+1} - \theta^k\| < \epsilon$$

$$\|C^{k+1} - C^k\| < \epsilon$$

$$\|\nabla_C\| < \epsilon$$



# Linear Regression with Single Variable Feature

$$y_i = mx_i + c + \epsilon_i$$

MLE.

$$S(m, c) = \sum_{i=1}^N \epsilon_i^2$$

LSE: Least Squares Error  
 [Gives most optimal solution]  
 if  $\epsilon_i$  is Gaussian.  
 if there are outliers,  
 non-Gaussian dist. may be  
 required.

\* find  $m$  &  $c$  which minimise  $S$ .

$$\frac{\partial S}{\partial m} = 0 \Rightarrow -2 \sum_{i=1}^N x_i(y_i - mx_i - c) = 0$$

$$\Rightarrow m \sum_{i=1}^N x_i^2 + c \sum_{i=1}^N x_i = \sum_{i=1}^N x_i y_i$$

$$\frac{\partial S}{\partial c} = 0 \Rightarrow -2 \sum_{i=1}^N (y_i - mx_i - c) = 0$$

$$\Rightarrow m \sum_{i=1}^N x_i + c N = \sum_{i=1}^N y_i$$

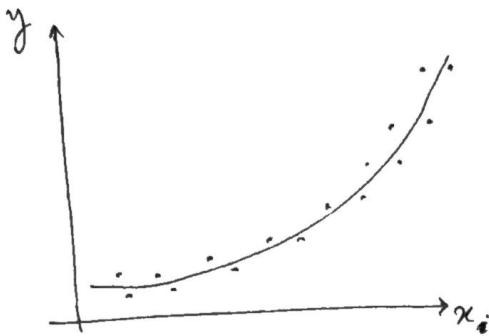
Solving the simultaneous equation, we get

$$m = \frac{N \sum x_i y_i - \sum x_i \sum y_i}{N \sum x_i^2 - (\sum x_i)^2}$$

$$c = \frac{\sum y_i - m \sum x_i}{N}$$

Linear Reg is one  
 of the most popular  
 & effective methods  
 used in science  
 engineering.

"Work Horse"



$$y_i = \beta_0 + \beta_1 x_i + \beta_2 x_i^2$$

Can still be solved using  
Linear regression methods.

$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2}$$

$$\text{where } x_{i1} = x_i$$

$$x_{i2} = x_i^2$$

In general,

$$f(x) = \beta_0 + \sum_{j=1}^m x_j \beta_j ; \quad x = [x_1, x_2, \dots, x_m]$$

$$= \sum_{j=0}^m x_j \beta_j \quad \text{where } x_0 = 1.$$

Multiple Linear Regression  
or Multivariable  
Multivariate (General)  
↳ multiple output.

$$y = w_0 + w_1 x_1 + w_2 x_2 + \epsilon$$

$$\text{Five data points : } y_i = w_0 + w_1 x_{1i} + w_2 x_{2i}; \quad 1 \leq i \leq 5$$

$$y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_5 \end{bmatrix}$$

$$x = \begin{bmatrix} 1 & x_{11} & x_{21} \\ 1 & x_{12} & x_{22} \\ 1 & x_{13} & x_{23} \\ 1 & x_{14} & x_{24} \\ 1 & x_{15} & x_{25} \end{bmatrix}$$

$$w = \begin{bmatrix} w_0 \\ w_1 \\ w_2 \end{bmatrix}$$

# MULTIPLE AND MULTIVARIATE LINEAR REGRESSION

Linear regression with single feature: MLE with Gaussian error.

$$y_i = m x_i + c + \epsilon_i$$

$$m, c \text{ obtained by minimizing } S = \sum_i (y_i - mx_i - c)^2$$

LSE.

Multiple LinReg:

$$y_i = w_0 + w_1 x_{1i} + w_2 x_{2i} + \dots + w_m x_{mi} + \epsilon_i$$

$$= w^T x_i + \epsilon_i$$

w &  $x_i$  are vectors instead of scalars.  
but  $y$  is scalar.

Y must be a vector.  
(implied case  $x_i$  is a vector)

$$S = \sum_i (y_i - w^T x_i)^2$$

~~So~~ ~~all data points~~ ~~are vectors~~

~~so~~ ~~the~~ ~~generalized~~ ~~to~~ ~~class of vectors~~

$$= \sum_i \text{vector of all data} \times \text{vector of } x_i$$

$$= (y - Xw)^T (y - Xw)$$

$$= (y^T - w^T X^T)(y - Xw)$$

$$= y^T y - 2y^T Xw + w^T X^T Xw$$

$$\nabla_w S = 2X^T Xw - 2X^T y = 0$$

$$\Rightarrow \boxed{X^T Xw = X^T y} \equiv \text{NORMAL EQUATION}$$

ORDINARY LEAST SQUARES

$$\Rightarrow w_{OLS} = (X^T X)^{-1} X^T y$$

~~very hard to compute~~

If model complexity is large, parameter values can be large leading to unstable solution, i.e.  $w_{OLS}$  can change a lot for small changes in data points. Also leads to over-fitting.

USE REGULARISATION.



## RIDGE REGRESSION

Encourage parameter values to be small [MAP]

use zero-mean Gaussian PRIOR

$$S = \frac{1}{N} \sum_{i=1}^N (y_i - w^T x_i)^2 + \lambda \underbrace{\|w\|_2^2}_{w^T w} \xrightarrow{\text{L}_2 \text{ Regularisation term or weight decay}}$$

$$w_{\text{RIDGE}} = (\lambda I_D + X^T X)^{-1} X^T y$$

$\lambda = 0 \Rightarrow$  over-fitting

$\lambda \approx 1 \Rightarrow$  under-fitting.

$$\|w\|_2 = \sqrt{w_0^2 + w_1^2 + w_2^2 + \dots}$$

$\hookrightarrow$  usually omitted from regularisation

## LASSO

Use ~~L1~~ L1 regularisation

$$S = \frac{1}{N} \sum_{i=1}^N (y_i - w^T x_i)^2 + \lambda \underbrace{\|w\|_1}_{L1 \text{ regular norm}}$$

useful when we have outliers  
(non-Gaussian errors)

$$\|w\|_1 = |w_0| + |w_1| + |w_2| + \dots$$

## SUBSET SELECTION

Pick K best features using some thresholds.

## BAYESIAN Linear Regression

Estimate errors in estimation of weights.

## MULTIVARIATE (GENERAL) LINEAR REGRESSION

multiple vector output

Same expression for Normal Eqn. works  
but  $y$  is now a matrix instead of a vector.

## LOGISTIC REGRESSION - I

Classification :  
Spam vs. Normal Email  
Malignant vs. Healthy Tissue  
Fake vs. Real News.  
Human vs. Animal image.

One of  
the most  
popular  
NL algos.

No classifier is perfect.  
 So, we are looking for a probability estimate for  $p(y_i|x_i;\theta)$

$$\hat{y}_i = h(x_i; \theta) \quad \xrightarrow{\text{hypothesis}} \quad \text{if } 0 \leq h \leq 1 \quad \text{else } x_i; \theta$$

Decision boundary: vector  
 (Linear) or HYPERPLANE

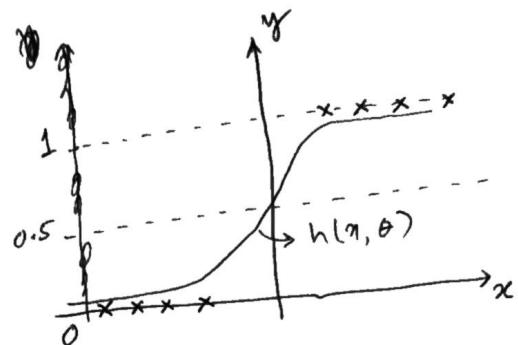
$$\theta_0 + \theta^T x_i > 0 \Rightarrow$$

$$\theta_0 + \theta^T x_i \leq 0 \Rightarrow$$

Sigmoid / Logistic function:  $h^{(m)}(x) = \frac{1}{1 + e^{-\theta^T x}}$   
or Logit

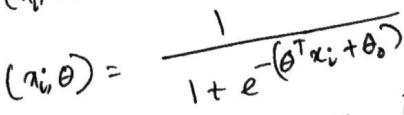
$$\text{log} \frac{h(x_i; \theta)}{1 - h(x_i; \theta)} = \theta_0 + \theta^T x_i : \text{Linear Model}$$

odds



Clearly  $y$  vs.  $x$  is  
NOT a straight line.  
So, makes no sense to  
~~fit~~ use  
linear regression

$$h(x_i|\theta) = \frac{1}{1 + e^{-(\theta^T x_i + \theta_0)}}$$

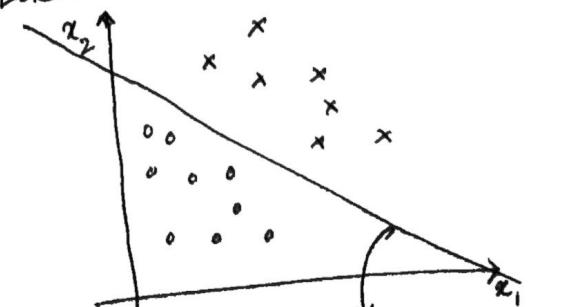


Sigmoid / Logistic function:  $h(x_i; \theta) = \frac{1}{1 + e^{-\theta_0 - \theta_1 x_i}}$

or Logit

many other options available but this one is easier to handle & interpret.

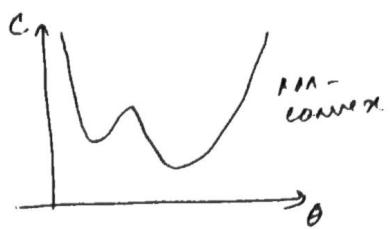
↳  $y_i = \theta_0 + \theta_1 x_i$ : Linear Model



Linear  
Decision boundary

$$\text{LSE} : C_{\text{LSE}}(\theta) = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2$$

But this not convex.



for a function,  $f(\theta)$ , to be convex.

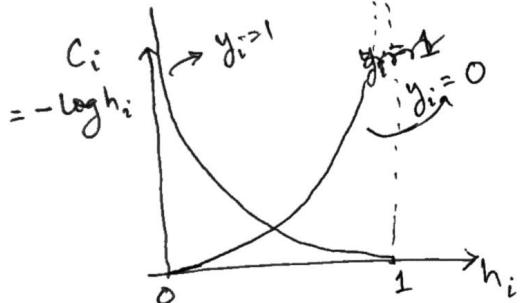
$$\text{we must have } \frac{\partial^2 f}{\partial \theta^2} > 0 \quad \forall \theta$$

for binary classification,

$$\theta_{\text{MLE}} = \underset{\theta}{\operatorname{argmax}} \left[ \sum_{i=1}^N [-y_i \log h_i - (1-y_i) \log(1-h_i)] \right]$$

$$\therefore C(\theta) = \frac{1}{N} \sum_{i=1}^N [-y_i \log h_i - (1-y_i) \log(1-h_i)] ; y_i \in \{0, 1\}$$

$\hookrightarrow$  Binary Cross Entropy.



$$\begin{aligned} \text{if } y_i = 1, \quad C_i(\theta) &= -\log h_i \\ \text{if } y_i = 0, \quad C_i(\theta) &= -\log(1-h_i) \end{aligned}$$

Gradient Descent:

$$\theta_j \rightarrow \theta_j - \eta \frac{\partial C}{\partial \theta_j}$$

$\hookrightarrow$  Learning Rate

# LOGISTIC REGRESSION - II

## MULTI-CLASS CLASSIFICATION

$S = \{\text{Spam, important, Not important}\}$

$$y = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ y^{(3)} \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} \quad \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} \quad \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}$$

ONE-HOT  
ENCODING

$$\hat{y} = \begin{bmatrix} \hat{y}^{(1)} \\ \hat{y}^{(2)} \\ \hat{y}^{(3)} \end{bmatrix}$$

$$\begin{aligned} \hat{y}^{(1)} &= P[y^{(1)} = 1] \\ \hat{y}^{(2)} &= P[y^{(2)} = 1] \\ \hat{y}^{(3)} &= P[y^{(3)} = 1] \end{aligned}$$

$$0 \leq \hat{y}^{(k)} \leq 1$$

$$\sum_k \hat{y}^{(k)} = 1$$

~~$$z^{(k)} = \theta^T x + b$$~~

$$z^{(k)} = e^{x^T \theta^{(k)} + b_0}$$

SOFTMAX function:

$$\hat{y}^{(1)} =$$

~~$$\frac{e^{z^{(1)}}}{e^{z^{(1)}} + e^{z^{(2)}} + e^{z^{(3)}}}$$~~

$$\hat{y}^{(2)} =$$

~~$$\frac{e^{z^{(2)}}}{e^{z^{(1)}} + e^{z^{(2)}} + e^{z^{(3)}}}$$~~

$$\hat{y}^{(2)} = \frac{e^{z^{(2)}}}{e^{z^{(1)}} + e^{z^{(2)}} + e^{z^{(3)}}}$$

COST / LOSS FUNCTION

for 2-classes :

$$\begin{aligned} C &= \frac{1}{N} \sum_{i=1}^N [-y_i \log \hat{y}_i - (1-y_i) \log (1-\hat{y}_i)] \\ &= \frac{1}{N} \sum_{i=1}^N [-y_i^{(1)} \log \hat{y}_i^{(1)} - y_i^{(2)} \log \hat{y}_i^{(2)}] \end{aligned}$$

$$\text{for } K\text{-classes : } C = \frac{1}{N} \sum_{i=1}^N \sum_{k=1}^K -y_i^{(k)} \log \hat{y}_i^{(k)}$$

CATEGORICAL  
CROSS ENTROPY

# LOGISTIC REGRESSION - I (contd.)

Classification :  $y_i \in \{0, 1\}$

$$\hat{y}_i = h(x_i, \theta) = \frac{1}{1 + e^{-(\theta^T x_i + \theta_0)}} \rightarrow \text{SIGMOID, LOGISTIC, LOGIT fn.}$$

$$C(\theta) = \frac{1}{N} \sum_{i=1}^N \left[ -y_i \log \hat{h}_i - (1-y_i) \log(1-\hat{h}_i) \right] \rightarrow \text{BINARY CROSS ENTROPY.}$$

Gradient Descent :  $\theta_j \rightarrow \theta_j - \eta \frac{\partial C}{\partial \theta_j}$  Learning Rate.

$$\frac{\partial C}{\partial \theta_j} = \frac{1}{N} \sum_{i=1}^N \left[ -\frac{y_i}{\hat{h}_i} \frac{\partial \hat{h}_i}{\partial \theta_j} + \frac{1-y_i}{1-\hat{h}_i} \frac{\partial \hat{h}_i}{\partial \theta_j} \right]$$

$$= \frac{1}{N} \sum_{i=1}^N \frac{(h_i - y_i)}{h_i(1-h_i)} \frac{\partial h_i}{\partial \theta_j}$$

$$= \frac{1}{N} \sum_{i=1}^N (h_i - y_i) x_{ij}$$

form Exactly same as  
Linear Regression  
just with a different expression  
for  $h_i$

$$\begin{aligned} \frac{\partial h_i}{\partial \theta_j} &= \frac{-1}{(1 + e^{-\theta^T x_i})^2} \cdot \frac{\partial}{\partial \theta_j} e^{-\theta^T x_i} \\ &= \frac{e^{-\theta^T x_i - \theta_0}}{(1 + e^{-\theta^T x_i})^2} \cdot \frac{\partial \theta^T x_i}{\partial \theta_j} \\ &= \frac{1}{1 + e^{-\theta^T x_i}} \frac{e^{-\theta^T x_i - \theta_0}}{1 + e^{-\theta^T x_i}} x_{ij} \\ &= h_i(1-h_i)x_{ij} \end{aligned}$$

The features  $x_{ij}$  used in Logit can actually be nonlinear fn. (eg. polynomial)  
of other features thereby allowing the decision boundary  
to be effectively nonlinear.

~~REGULARISATION~~

$$C(\theta) = \frac{1}{N} \sum_{i=1}^N \left[ -y_i \log \hat{h}_i - (1-y_i) \log(1-\hat{h}_i) \right] + \frac{\lambda}{N} \|\theta\|_2^2$$

on  $N \rightarrow \infty$ , this term drops out.

$$\therefore \text{approx for } \frac{\partial C}{\partial \theta_j} = \frac{1}{N} \sum_{i=1}^N (h_i - y_i) x_{ij} + \frac{2\lambda}{N} \theta_j$$

## MULTI-CLASS LOGISTIC REGRESSION

Email: Spam, important, Not-important  
Fake identification

ONE vs. ALL [or One vs. Rest]

Let  $y_i \in \{S, I, NI\} = K$ -categories

Divide into three problems

$$h^{(1)}(x_i, \theta) = \frac{1}{S} \text{ vs. } \{\overline{I}, \overline{NI}\}$$

$$h^{(2)}(x_i, \theta) = I \text{ vs. } \{S, \overline{NI}\}$$

$$h^{(3)}(x_i, \theta) = NI \text{ vs. } \{S, I\}$$

$$h(x_i, \theta) = \max \{h^{(1)}, h^{(2)}, h^{(3)}\}$$

K-binary classifications

ONE vs. ONE

S vs. I

$K_{C_2}$  binary classifications.

I vs. NI

S vs. NI

Pick category which has most classification.

These methods do not give probability estimates  
since  $h^{(1)} + h^{(2)} + h^{(3)} \neq 1$

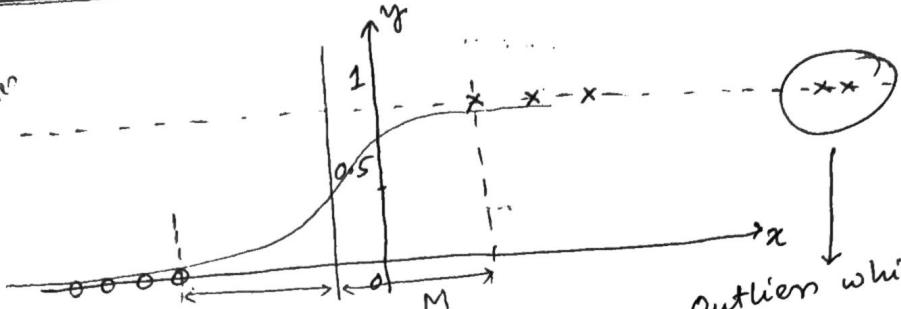
since  $h^{(1)} + h^{(2)} + h^{(3)} \neq 1$

One vs. One is also ~~leads to~~ computationally expensive.

Better to have a method which can learn  
multi-class parameters at one go.

# SUPPORT VECTOR MACHINES - I

faster than logistic regression  
extends linear boundary to non-linear boundaries



SVM: Give up probabilistic estimates for a better decision boundary  
only worry about points closest to decision boundary.

Outliers which significantly distort Logistic Regression.

Classification Rule :

$$G(x) = \text{sgn}[\theta^T x]$$

$$h(x_i, \theta) = \theta_0 + \theta^T x_i > 0 \Rightarrow \hat{y}_i = +1$$

$$h(x_i, \theta) = \theta_0 + \theta^T x_i < 0 \Rightarrow \hat{y}_i = -1$$

find  $\theta$  such that

$M$ : Margin  
 $y \in \{-1, +1\}$

$M$  is maximised for the closest points closest to Decision Boundary hyperplane.

$$\vec{n} = (\text{Normal vector}) (\theta^{(1)}, \theta^{(2)})$$

(Normal to the decision boundary)

$$\vec{n} = \vec{x}_i^\perp + y_i \frac{\vec{\theta}}{\|\vec{\theta}\|}$$

$$h(x_i, \theta) = \theta_0 + \theta^T x_i = \theta_0 + \vec{\theta} \cdot \vec{x}_i$$

$$= \theta_0 + \vec{\theta} \cdot \left[ \vec{x}_i^\perp + y_i \frac{\vec{\theta}}{\|\vec{\theta}\|} \right]$$

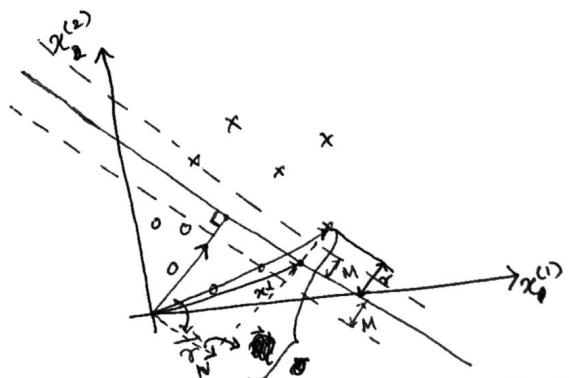
$$= (\theta_0 + \vec{\theta} \cdot \vec{x}_i^\perp) + y_i \frac{\|\vec{\theta}\|^2}{\|\vec{\theta}\|} = y_i \|\vec{\theta}\|$$

$$\Rightarrow \vec{n}_i = \frac{h(x_i, \theta)}{y_i \|\vec{\theta}\|}$$

$$= \frac{\theta_0 + \theta^T x_i}{y_i \|\vec{\theta}\|}$$

Objective : Maximize  $\vec{n}_i$  for the closest support point [Support Vector]

while ensuring  $y_i h(x_i, \theta) > 0$



$$\theta^{(0)} + \theta^{(1)} x^{(1)} + \theta^{(2)} x^{(2)} = 0$$

$$\Rightarrow x^{(2)} = -\frac{\theta^{(1)}}{\theta^{(2)}} x^{(1)} - \frac{\theta^{(0)}}{\theta^{(2)}}$$

## OBJECTIVE

$$\max_{\theta_0, \theta_1} \min_{x_i} \frac{y_i(\theta_0 + \theta_1 x_i)}{\|\theta\|}$$

Note: Scaling  $\theta, \theta_0$  does not change this function.  
 So, scale them so that  $y_i h_i = 1$  for point closest to decision boundary.

$$\max_{\theta_0, \theta_1} \frac{1}{\|\theta\|} \quad \text{s.t.} \quad y_i (\theta_0 + \theta^T x_i) \geq 1 \quad \forall i$$

QUADRATIC PROBLEM  
with N linear constraints

Quadratic program  
with  $N$  linear constraints  
[Convex]  
is regularisation.  
⇒ good generalisation

$$= \min_{\theta, \theta_0} \frac{1}{2} \|\theta\|^2 \quad \text{s.t.} \\ \rightarrow \text{some}$$

$$= \left( I_n + A^T x_i \right)^{-1}$$

Lagrange function

$$\text{Convert } \rightarrow L_D = \frac{1}{2} \| \theta \|^2 - \sum_{i=1}^N \lambda_i [y_i(\theta_0 + \theta^T x_i) - 1] \quad \begin{matrix} \text{Lagrange multipliers} \\ \lambda_i > 0 \quad x_i \\ \text{support} \end{matrix}$$

$$\text{Max} \underset{\theta}{\text{min}} \quad L_D$$

$$\frac{\partial L_D}{\partial \theta_0} = - \sum_{i=1}^N x_i y_i = 0$$

$$\nabla \theta_0 = \theta_i - \sum_{i=1}^N x_i y_i$$

$$g(x) = \sum_{i=1}^N \alpha_i y_i x_i$$

Support Vectors

Substituting this in L<sub>0</sub>, we get

$$L_0 = \sum_{i=1}^N x_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N x_i x_j y_i y_j x_i^T x_j$$

negative function

$$D = \sum_{i=1}^n x_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n c_{ij}$$

↳ lower bound on objective function  
..... simpler

lower bound on obj  
maximising this is a simpler  
quadratic programming problem.

Maximising this is a convex quadratic programming problem  
LAGRANGE DUAL OBJECTIVE FUNCTION

convex quadratic fit  
LAGRANGE DUAL OBJECTIVE FUNCTION  
(WOLFE) 11

min. of original  $L_0$  is  $- \infty$  !!

## SUPPORT VECTOR MACHINES - II

what if the classes are not linearly separable?

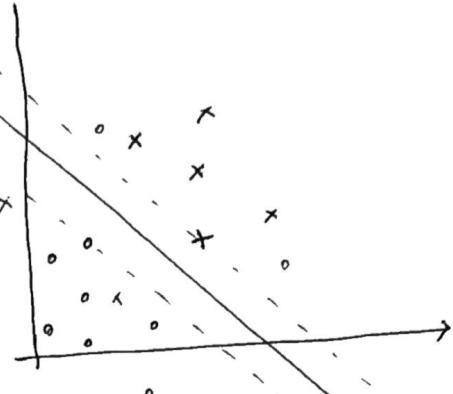
if overlap is small, use SLACK VARIABLES

$$\xi_i \geq 0, i=1,2,\dots,N$$

Margin by which points

are allowed to be on ~~the~~ wrong side.

$\xi_i \geq 0, \sum \xi_i \leq \text{constant}$  [sort of regularization]



$$\min_{\theta, \theta_0} \frac{1}{2} \|\theta\|^2 \quad \text{s.t.} \quad y_i(\theta_0 + \theta^T x_i) \geq 1 - \xi_i$$

trade-off between large margin & small training error

$$L_D = \frac{1}{2} \|\theta\|^2 + C \sum_i \xi_i - \sum_i \alpha_i [y_i(\theta_0 + \theta^T x_i) - (1 - \xi_i)] - \sum_i \lambda_i \xi_i$$

Lagrange (primal) function [minimise w.r.t.  $\theta_0, \theta, \xi_i$ ]

$$L_D = \sum_i \alpha_i - \frac{1}{2} \sum_i \sum_{i'} \alpha_i \alpha_{i'} y_i y_{i'} \boxed{x_i^T x_{i'}}$$

Lagrange (Wolfe) Dual Objective fn.

Optimisation depends only on dot product of input features.

\* What if overlap is large?

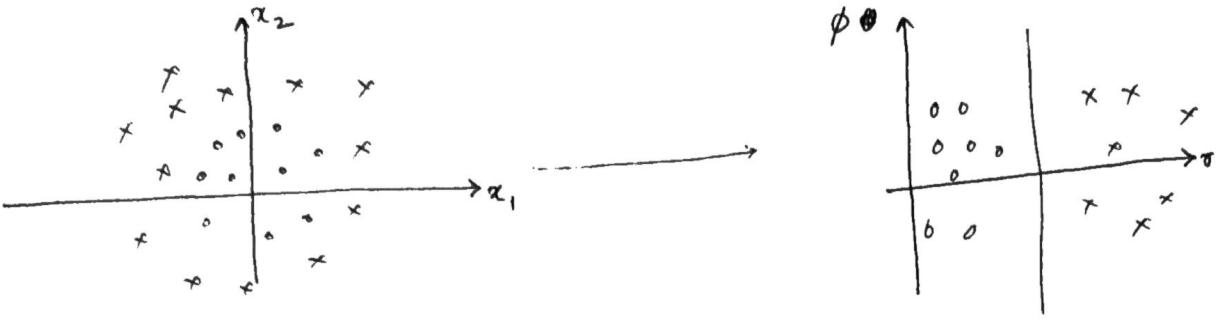
Need to transform variables, so that in the new space, the points are linearly separable.

Since dual  $L_D$  only depends on dot products, we don't need the actual transformation expression but only the new dot-product.

$\Rightarrow$  Kernel trick.

$$L_D = \sum_i \alpha_i - \sum_i \sum_{i'} \alpha_i \alpha_{i'} y_i y_{i'} K(x_i, x_{i'})$$

$\rightarrow$  positive semi-definite & symmetric.



Popular choices for Kernel.

d-th degree polynomial:

$$K(x_i, x_i) = (1 + \langle x_i, x_i \rangle)^d$$

usually best choice.  $\rightarrow$  Radial Basis (RBF): or Gaussian  $K(x_i, x_i) = \exp\left(-\gamma \frac{\text{kernel bandwidth}}{\|x_i - x_i\|^2}\right)$  [Most widely used based on Euclidean distance.]

Neural Network:  
or Sigmoidal  $K(x_i, x_i) = \tanh(\gamma_1 \langle x_i, x_i \rangle + \gamma_2)$

on original  $L_0$ , if  $C$  is large, it penalises any  $\xi_i$  & leads to a very wiggly boundary (overfit).

$$h_i = \theta_0 + \theta^T x_i$$

$$\nabla_{\theta} L_0 \frac{\partial}{\partial \theta_j} = \theta_j - \sum_i \alpha_i y_i x_i^T = 0 \Rightarrow \theta_j = \sum_i \alpha_i y_i x_i^T$$

$$\Rightarrow h_i = \theta_0 + \sum_i \alpha_i y_i x_i^T x_i'$$

$$= \theta_0 + \sum_i \alpha_i y_i K(x_i, x_i)$$

RBF: large  $\gamma$  leads to overfitting since even small distances are penalised.

Curse of Dimensionality

## SVM-II (contd.)

$$L_D = \frac{1}{2} \|\theta\|^2 + C \sum_i \xi_i - \sum_i \alpha_i [y_i(\theta_0 + \theta^T x_i) - (1 - \xi_i)] - \sum \mu_i \xi_i$$

minimize w.r.t.  $\theta, \theta_0, \xi_i$  vector

$$\frac{\partial L_D}{\partial \theta_0} = 0 \Rightarrow \sum_i \alpha_i y_i = 0$$

$$\nabla_{\theta} L_D = 0 \Rightarrow \theta = \sum_i \alpha_i x_i y_i$$

$$\frac{\partial L_D}{\partial \xi_i} = 0 \Rightarrow \alpha_i = C - \mu_i$$

Substituting all these in  $L_D$ , we obtain,

$$L_D = \sum_i \alpha_i - \frac{1}{2} \sum_i \sum_{i'} \alpha_i \alpha_{i'} y_i y_{i'} x_i^T x_{i'}$$

↳ Lagrange Dual objective function

Maximize using Quadratic Programming  
subject to  $0 \leq \alpha_i \leq C$  &  $\sum \alpha_i y_i = 0$

To get optimum solution, use Karush-Kuhn-Tucker conditions,

$$\alpha_i [y_i(\theta_0 + \theta^T x_i) - (1 - \xi_i)] = 0 \rightarrow \alpha_i \neq 0 \text{ only for support vectors}$$

$$\mu_i \xi_i = 0 \rightarrow \mu_i \neq 0 \text{ only when } \xi_i \neq 0$$

$$y_i(\theta_0 + \theta^T x_i) - (1 - \xi_i) \geq 0$$

## SVM - Regularisation

$$\min_{\theta} \frac{1}{2} \|\theta\|^2 + C \sum \xi_i$$

↑  
of Kernel

SVM - Polynomial & Sigmoid Kernel application. [when does RBF not work]

RBF is slower than polynomial kernel; especially for large data or high dimension

Polynomial kernel is popular in NLP with  $d=2$  [ $d \uparrow \rightarrow$  overfitting & numerical instability]

Refers to No Free Lunch Theorem

## SVM Multi-class

one vs. one or one vs. rest.

Lagrange Multiplier  
[Could be maximal margin]

Support vectors  
are not necessarily  
the points closest  
to decision boundary  
for when slack  
variables used.

## high $\gamma$ RBF overfitting

Logreg → find prob & change threshold.  
+ve  $\theta_j \rightarrow$  feature favours  $y=1$   
Missing values → take mean/mode or remove from data  
NOT scale well to large training sets

\* ~~Kernel - Does NOT scale well to large feature vectors~~

~~RBF is slower~~  
Hyperplane - Decision Boundary

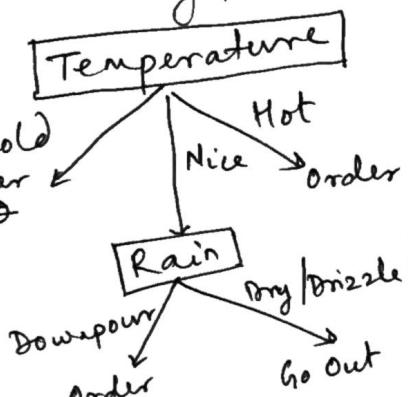
Hyper parameters vs. Parameter  
(internal model)  
(user tuned)

Generative  
vs.  
Discriminative

# DECISION TREES

Represents hierarchical decision making process.  
Dinner: Go Out or Order?

Each node represents an attribute or feature  
or value for that attribute (can be non-numerical)  
represents possible order or out



How to decide hierarchy?

Nodes with lowest entropy (best classification) is at the top/higher level.

A: Attribute

S: Total collection of examples

SA<sub>v</sub>: Subset of S for which Attribute A has value 'v'.

H<sub>S</sub>: Entropy of S

H<sub>SA<sub>v</sub></sub>: Entropy of SA<sub>v</sub>

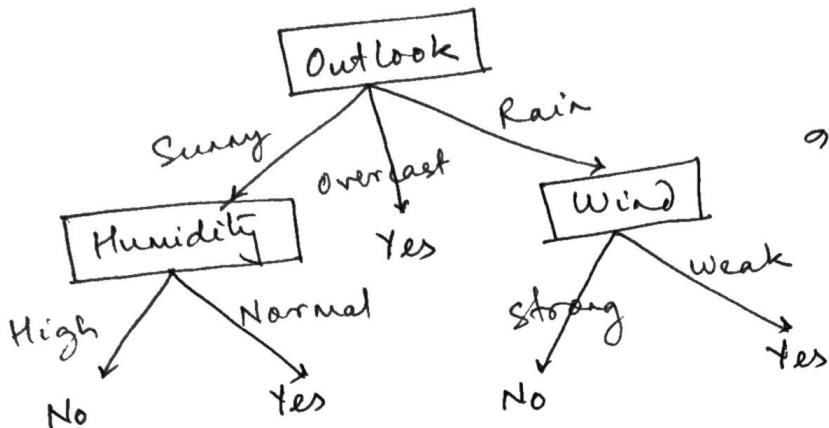
H<sub>A</sub>: Entropy of A

$$H_{SA_v} = \sum -p_i \log_2 p_i$$

$$H_A = \sum_{v \in A} \frac{|SA_v|}{|S|} H_{SA_v}$$

$$\text{Information Gain: } G(S, A) = H_S - H_A$$

ML by Tom Mitchell. [Play Tennis]



ID3:  
Iterative  
Dichotomiser 3

Gini Index / Impurity:

$$G_{AS_v} = 1 - \sum p_i^2$$

$$G_A = \sum \frac{|SA_v|}{|S|} G_{AS_v}$$

Works almost similar to entropy measure but faster to calculate.

Overfitting

- Pre-pruning: stop split beyond a significance threshold
- Post-pruning: based on validation set  
[better than pre-prune since hard to decide threshold]

Missing Features

Assign most common value or using some probability measure.

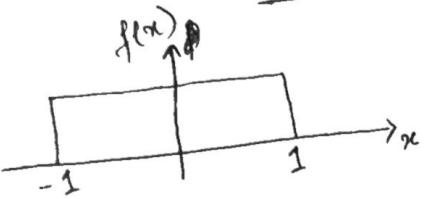
ADVANTAGE OF TREES

- Easy to interpret
- Easily handle mix data [symbols / numbers]
- Robust to outliers
- Scales well to large datasets
- ...

- DISADV.
- Low bias & high variance
  - Unstable due to hierarchical nature  
[small change in data can lead to large change in tree]

Addressed by Random Forest.

# RANDOM FOREST



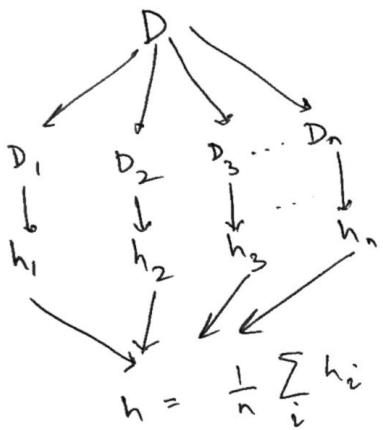
Law of Large Numbers

$$-1 \leq x \leq 1 \quad E[x] = 0$$

$x_1, x_2, x_3, \dots, x_n$  iid rv.

$$y = \frac{x_1 + x_2 + \dots + x_n}{n} \rightarrow E[x] \text{ as } n \rightarrow \infty$$

standard method of reducing variance in prediction.



RANDOM FOREST : Bagging + Decision Trees

bootstraps  
choose features randomly.

Very resistant to "curse of dimensionality"  
gets slower but still works very well.

only two hyper-parameters  
 $n$  &  $\sqrt{D}$

But loses interpretability.

Each decision tree is a weak learner

Bagging makes it a strong learner.

for ML methods with ~~high~~ bias, use BOOSTING.

Bagging runs in parallel

Boosting is serial or sequential  
[use training accuracy to assign higher weight to misclassified data & re-train. Do in loop for few times]

Bootstrap Aggregating

BAGGING : very effective way to reduce variance.

Each  $D_i$  has same size as  $D$   
[drawn randomly, can have same data point multiple times]

for each split, use a subset of features.

ENSEMBLE LEARNING

~~less~~

STACKING

## Logistic Regression vs. SVM vs. Random Forest

LR	SVM	RF
① Works best with linearly separable data	Can handle nonlinearity by using kernels.	Inherently nonlinear
② Easier to interpret results. Also gives probability estimates	Hard to interpret results no probs.	<del>sooner as SVM</del> <del>hard to interpret</del> but gives prob.
③ Usually fast	faster than LR	Usually very slow
④ Works better with balanced data		can easily handle unbalanced data
⑤ Has problems with high dimensional data (lots of features)	same as LR	Easily handles curse of dimensionality (feature subset used)
⑥ Can't handle outliers	Handles outliers better	Can handle noisy data
⑦ better for 2-class	better for 2-class	Better for multi-class classification

BEST PRACTICE — Try all 3, tune hyper-parameters & choose the one that works best.

## DECISION TREES

- Data fragmentation: Continuous partitioning leaves less data for lower-level nodes leading to weak statistical support.  
→ Can lead to multiple mis-classifications.  
→ General problem in rule based learning.

- NP-Complete:  
A problem  $h^P$  in NP is NP-complete if every other problem in NP can be transformed (or reduced) to  $h^P$  in polynomial time.  
obviously,  $h^P$  is also NP-hard.  
but NP-hard does not imply NP-complete  
(e.g. halting problem)
- Decision Trees are NP-complete.  
If you can solve Decision Tree problem in polynomial time, it would prove  $P = NP$ .
- ID3 is a greedy, heuristic algorithm  
does best fit search for locally optimal entropy values.  
(no looking back) approximates the globally optimal solution.
- Decision Tree: No real learning? Should it be called ML/AI?  
purely rule based
- In RF, choose subset of features, to avoid correlations in each tree between trees.
- ID3 was originally for categorical values.  
C4.5 allowed numerical values.  
CART has binary classification tree & allows regression.

$\text{Gini}$  represent the expected error resulting from labelling instances in the leaf randomly.

$$p(1-p) \rightarrow \text{coincidence variance}.$$

### Thresholding & Discretization

Unsupervised thresholding: mean/median

Supervised " : lower convex hull

need to know number of 'bins' a priori

requires careful data analysis  
(mean/median may not be good classifier)

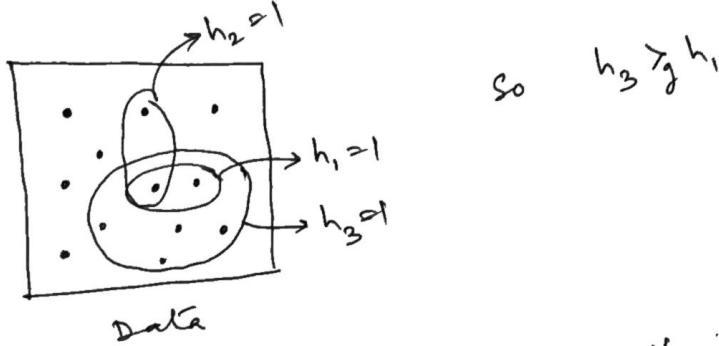
## ANN

- More general than

Let  $h_j(x)$  &  $h_k(x)$  be boolean-valued functions over  $X$ .

Then  $h_j$  is more-general-than-or-equal-to  $h_k$  ( $h_j \geq h_k$ )  
iff  $(\forall x \in X) [h_k(x) = 1 \Rightarrow h_j(x) = 1]$

if  $h_j$  is more-general-than  $h_k$   
then  $\&$   $h_k$  is more-specific-than  $h_j$

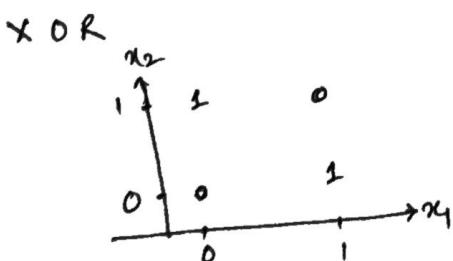


Perception = 1 neuron in ANN.  
(node of hidden layer)  
or output

Regularisation for ANN - L1, L2, Dropout.

If SGD is not guaranteed to converge.

can use adaptive learning rate  
[not necessarily better]



# Learning XOR

$$X = \{[0,0], [0,1], [1,0], [1,1]\}$$

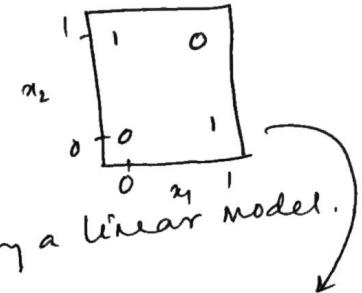
MSE loss function.

$$J(\theta) = \frac{1}{4} \sum_{x \in X} (f^*(x) - f(x; \theta))^2$$

$$f(x; w, b) = x^T w + b$$

Solving gives  $w=0$  &  $b=\gamma_2$

$\Rightarrow$  XOR can't be implemented by a linear model.



ML non-linear function

$\equiv$  Affine Transformation

+ Activation function

(Rectified Linear Unit,

ReLU)



$$h_i = g(x^T w_{:,i} + c_i)$$

$$f(x; w, c, w, b)$$

$$= w^T \max\{0, w^T x + c\} + b$$

$$w = \begin{bmatrix} 1 \\ 1 \end{bmatrix}, c = \begin{bmatrix} 0 \\ -1 \end{bmatrix}, w = \begin{bmatrix} 1 \\ -2 \end{bmatrix}, b = 0.$$

Universal Approximation Theorem

- May not be able to find parameter values corresponding to the desired function
- overfitting: May converge to the wrong function

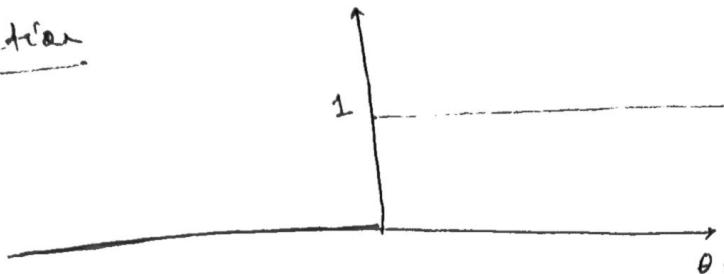
No Free Lunch Theorem

- There is no universally superior ML algo.

A feedforward network with single hidden layer is sufficient to represent any arbitrary function.

## ANN ACTIVATION FUNCTIONS

### Step function



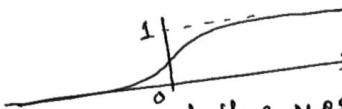
- can be used in output layer
- derivative is always zero, so can't help in gradient descent & so not for hidden layers.

### Linear function



- Gradient is a constant, so does not help in gradient descent
- Leads to overall linear classifier.

### Sigmoid function



- good classifier, one of the most widely used.
  - gives probability estimates
  - has changing gradient
  - but gradient saturates quickly leading to slow learning & can get stuck in local minima
  - can use hyperbolic tangent fn. instead.
- Vanishing gradient problem

### ReLU [Rectified Linear Unit] function



- faster than sigmoid/bath due to linearity & zero nature.
- But no learning in neg. area due to zero value.
- uses Leaky ReLU instead



## Softmax Function

- Similar to sigmoid but helps in multiclass classification.
- usually used in output layer.

Usually, use ReLU in hidden layer &  
sigmoid/softmax in output layers.

~~K - nearest neighbor~~

## ANN Optimization Methods

### BATCH GRADIENT DESCENT

- usual gradient descent
- uses whole dataset for each update
- very slow & resource consuming
- can't use new data on the fly (no online update)
- usually converges to local minima for non-convex loss functions
- done over many epochs

### STOCHASTIC GRADIENT DESCENT

- uses one data point at a time
- frequent updates with high variance
- random jumps allow convergence to global minima  
or potentially better minima.  
may not always be global since SGD can overshoot
- faster convergence
- shuffle data for each epoch

### MINI-BATCH GRADIENT

- SGD using small batch of dataset
- faster convergence than vanilla & less fluctuations than SGD with one data point
- usually called SGD.
- usually gets stuck at saddle points.

## ① SGD + Momentum

- Add update vector of past time step
- helps avoiding being caught in ravines.
- like pushing a ball down hill.
- gradient terms add up if they point in same direction

## NESTEROV SGD. [NAG - Nesterov Accelerated Gradient]

- on SGD + momentum, need to slow down

before gradient changes sign again.

- So estimate future parameter value using momentum term & then use gradient at that parameter value

$$v_t = \gamma v_{t-1} + \eta \nabla_{\theta} C(\theta - \gamma v_{t-1})$$

$\theta = \theta - v_t$  usually 0.9

## ADAGRAD

ADAM - faster  
SGD - better  
generalised

- Perform larger updates for frequent parameters
- Helps in dealing with sparsity.
- Basically use different learning rate ( $\eta$ ) for each parameter.
- $\eta$  is modified for each parameter based on past gradients.
- But  $\eta$  can keep decreasing to zero in this method. Resolved by Adadelta.

## ADAM [Adaptive Moment Estimation]

Another method to compute adaptive learning rate.

# Solving ODEs & PDEs using ANN

$$\frac{dy}{dx} = -2xy$$

$$y(0) = 1$$

$$\Rightarrow y(x) = e^{-x^2}$$

Analytical solution

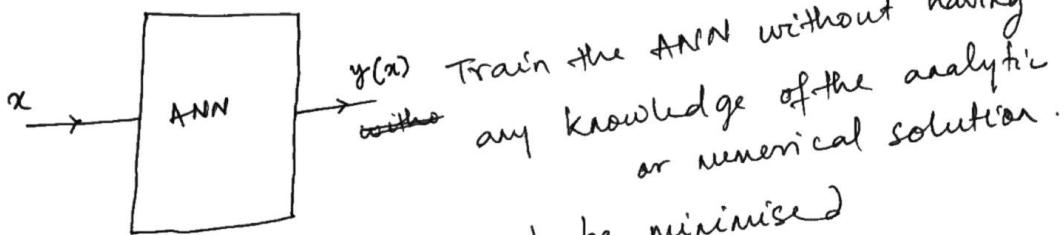
Not always possible

Numerical Simulation

- Cannot give output beyond the range on which simulation performed
- Requires human analysis of results.

ANN

- Possibility of learning the actual function [Universal Approx. Thm.]



Need a (loss) function to be minimised

$$C = \left[ \frac{dy}{dx} + 2xy \right]^2 + [y(0) - 1]^2$$

↳ Automatic Differentiation [AutoDiff] in TensorFlow

Different from both symbolic diff. & numerical differentiation.

Symbolic Diff - inefficient code (e.g. Mathematica)  
Numerical Diff - round-off errors ] unsuitable for higher & partial derivatives with many inputs needed in gradient descent

AutoDiff is NOT numerical differencer

It is both efficient & numerically stable.  
(partly symbolic & partly numerical)

# DEEP LEARNING

$DL = ANN + \text{Representation Learning}$

✓ Automatically discover the representations (abstract)  
needed for classification

✗ NO explicit feature selection

Mainly inspired by NLP, speech processing & Computer Vision.

Deep: Generally representation learning requires  
many layers of ANN.

CAP [Credit Assignment Path]: describes potentially  
causal connections between input & output.

DL usually has high CAP. [no specific threshold]

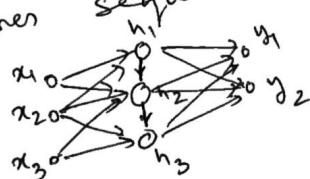
DL is computationally intensive & became popular  
largely due to GPUs.

NLP  
↳ Symbolic [1950s - 90s]  
↳ Statistical [1990s - ]  
↳ Neural [2010s - ]

DL Algorithms:

CNN - captures spatial information of images.

RNN - captures sequence information in NLP



vanishing & exploding  
gradients

LSTM - Additional structure in RNNs to  
decide which info to keep & which to  
discard.

Transformers - LSTMs are hard to train since  
it is sequential & can't be fully parallelised

Uses the "Attention Mechanism" only  
Significant improvement in language  
translation accuracy.

## Limitations or Criticisms of DL

- Lack of Theory & Explainability

DL is essentially a blackbox

We neither know the function learnt nor the features used for classification or [especially an issue in medical diagnosis] other prediction.

- Biases in prediction

Algo can learn to classify based on "problematic" features

Bias towards women & ~~other~~ minority groups (eg. automatic CV suggestions).  
for jobs & admissions.

- Cyber Threat

Easy to trick DL based surveillance systems

- Reliance on human microwork

Labelling done by humans

Continued demand for labelled data for calibration & updation of ANN.

Deep Learning is still Weak AI.

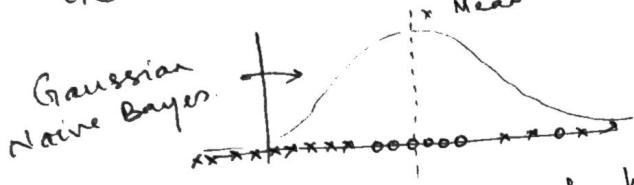
Awareness does NOT arise out of ignorance.

# GENERATIVE VS. DISCRIMINATIVE MODELS

Bayesian Reasoning  
(BPE)

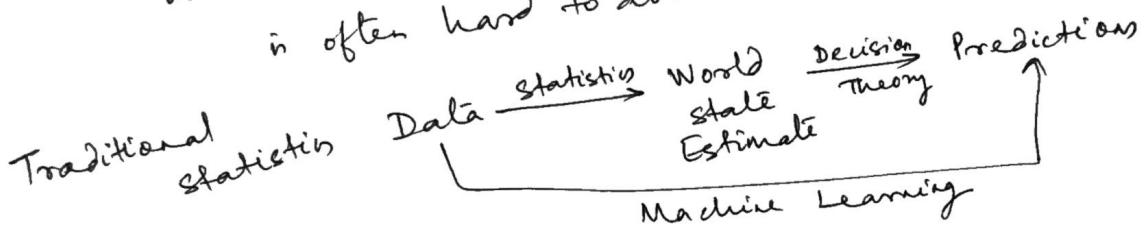
function Approximation  
(MLE & MAP)  
ML

Generative models can go wrong when there are outliers.



SVM, for example, can handle this with ease.

Generative may work here if we remove outliers or make better choice of distribution, but that is often hard to do.



## Generative

Goal: To explain data  
(concept learning)

More accurate when they use correct models

Explainable

Generate realistic data samples  
(artwork, super-resolution, colorization, etc.)

simulation & planning  
for reinforcement learning

Generative Adversarial Networks  
(GANs)

## Discriminative

Goal: To make predictions  
(shortcut coaching)

More robust against outliers & bad models

Blackbox as we use more complex models

## GANs

### Generators

Generate new data (fake) to fool the discriminator into believing that samples ~~data~~ generated is real

### vs. Discriminators

Identify fake (new generated) samples & force the generator to improve.

Applications: Art AI

Imaginary fashion models

Model dark matter distribution & predict gravitational lensing

DeepFakes!!

MNIST handwritten digits

Bayesian Networks - Help in identifying how causes ~~for~~ generate outcomes / events.  
[work both ways bottom up & top down]

# PRINCIPAL COMPONENT ANALYSIS [PCA] - I

Main objective is dimensionality reduction, i.e. to reduce the number of variables while preserving max info. Some loss of accuracy but leads to major simplifications in data analysis.

## Step 1 : Feature Scaling

Important since other features with wider range will dominate.

## Step 2 : Covariance Matrix Computation

Variables can be correlated.

Covariance helps in finding these relationships.

## Step 3 : Find eigen values & eigen vectors of covariance matrix

- Eigen vectors are linear combinations of data points & chosen such that they are independent of each other. 10-dimensional data will give 10 eigen vectors.
- Represent directions that explain maximum amount of data.
- First principal component accounts for max variance & then the 2nd one & so on.

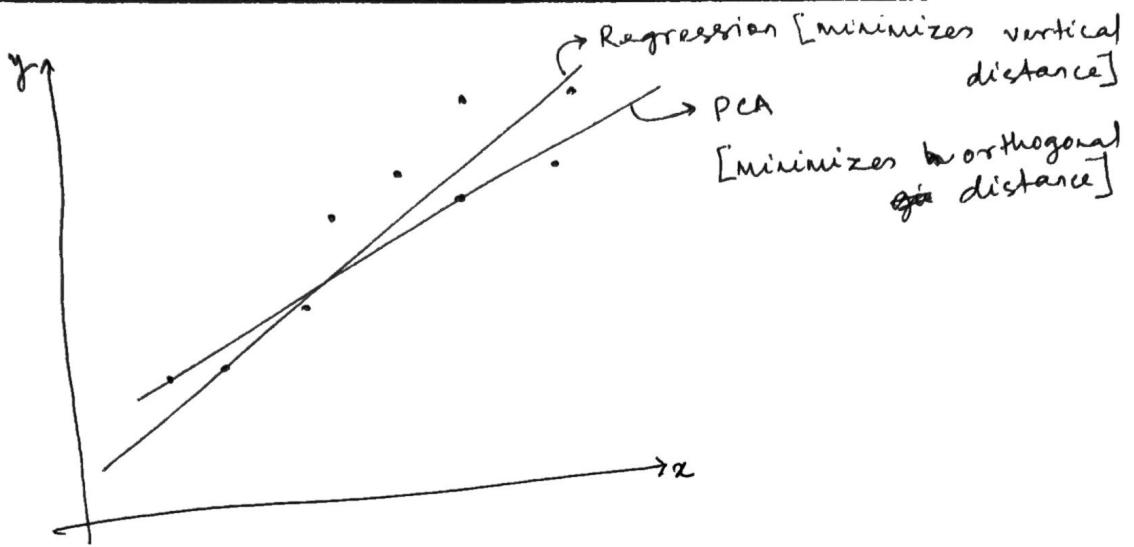
- find eigen vectors  $\{\vec{v}_i\}$  & eigen values  $\{\lambda_i\}$   
Arrange  $\{\lambda_i\}$  in decreasing order.  $\lambda_1 > \lambda_2 > \dots$

so  $\vec{v}_1$  corresponds to first component  
 $\vec{v}_2$  to second component, & so on.

- Keep a subset of the eigen vectors & discard the rest. This leads to dimensionality reduction since all the data points are now projected onto these eigen vectors through dot product.

Fraction of variance in data explained/captured

$$\text{by } \vec{v}_i = \frac{\lambda_i}{\sum \lambda_i}$$



$\{x_n\}$  Data points  $n=1, 2, \dots, N$

↓  
column vectors  
of dimensionality D.

to be projected to dimensions M < D.

PCA - II

Let M=1 for simplicity

Need to find direction  $u_1$  of projected data.

Projected data points  $\{y_n\} = \{u_1^T x_n\}$

Mean:  $\bar{y} = u_1^T \bar{x}$ , where  $\bar{x} = \frac{1}{N} \sum x_n$

Variance:  $S_y = \frac{1}{N} \sum_n \{u_1^T x_n - u_1^T \bar{x}\}^2$

$$= \frac{1}{N} \sum_n (u_1^T x_n - u_1^T \bar{x})(u_1^T x_n - u_1^T \bar{x})^T$$

$$= \frac{1}{N} \sum_n (u_1^T x_n - u_1^T \bar{x})(x_n^T u_1 - \bar{x}^T u_1)$$

$$= \frac{1}{N} \sum_n (u_1^T x_n - u_1^T \bar{x})(x_n^T u_1 - \bar{x}^T u_1) u_1$$

$$= \frac{1}{N} \sum_n u_1^T (x_n - \bar{x})(x_n^T - \bar{x}^T) u_1$$

$$= u_1^T S_x u_1, \text{ where } S_x = \frac{1}{N} \sum_n (x_n - \bar{x})(x_n - \bar{x})^T$$

Objective: Maximize  $S_y$  keeping  $\|u_1\|=1$  fixed.

∴ find  $u_1$  which maximizes

$$L = u_1^T S_x u_1 + \lambda_1 (1 - u_1^T u_1)$$

→ Lagrange multiplier.

$\nabla_{u_1} L = 0 \Rightarrow S_x u_1 = \lambda_1 u_1 \Rightarrow u_1$  is the eigen vector  
of  $S_x$  &  $\lambda_1$  is the eigen-value.

$$\Rightarrow S_y = u_1^T S_x u_1 = \lambda_1$$

⇒ choose  $u_1$  to be the eigen-vector  
with max eigen-value  
so that  $S_y$  is maximised.

Same way  $u_2, u_3, \dots$  can be chosen.

If  $M=D$ , PCA still works & is a rotation of the coordinate axis  
to maximise variance along the new directions  
or components.

What if  $D \gg N$ ? Very high dimensional data  
(small dataset of images)

finding eigen vectors of  $D \times D$  matrix  
has computational cost  $\Theta(D^3)$ .

Define  $X : (N \times D)$  dimensional centred matrix  
 $n^{th}$  row given by  $(x_n - \bar{x})^T$

Covariance Matrix :  $S = \frac{1}{N} \underbrace{X^T X}_{\hookrightarrow D \times D \text{ dimensional}}$

$$\frac{1}{N} X^T X u_i = \lambda_i u_i$$

$$\Rightarrow \frac{1}{N} X X^T X u_i = \lambda_i \underbrace{X_i u_i}_{v_i} : \text{Multiplying both sides by } X$$

$$\Rightarrow \frac{1}{N} \underbrace{X X^T}_{\hookrightarrow N \times N \text{ Dimensional}} v_i = \lambda_i v_i$$

$\Theta(N^3)$  instead of  $\Theta(D^3)$

has the same eigen values as  $S$   
[ $D-N+1$  eigen values of  $S$  are zero]

One limitation of PCA is that it is limited to  
linear transformations of data.

Can use kernel PCA to do non-linear transformations

$$x_n = \begin{pmatrix} x_n^{(1)} \\ x_n^{(2)} \end{pmatrix} \quad n=1, 2, \dots, N$$

$$D=2$$

To be projected to M=1

$$y_n = u^T x_n \quad \text{to } (u^{(1)}, u^{(2)}) \quad u = \begin{pmatrix} u^{(1)} \\ u^{(2)} \end{pmatrix}$$

$$\begin{aligned} \underline{y\text{-Variance}}: \quad S_y &= \frac{1}{N} \sum_n (y_n - \bar{y})^2 \\ &= \frac{1}{N} \sum_n (u^T x_n - u^T \bar{x})^2 \\ &= \frac{1}{N} \sum_n (u^T x_n - u^T \bar{x})(u^T x_n - u^T \bar{x})^T \\ &= \frac{1}{N} \sum_n (u^T x_n - u^T \bar{x})(x_n^T u - \bar{x}^T u) \\ &= \frac{1}{N} \sum_n (u^T x_n x_n^T u - u^T x_n \bar{x}^T u - u^T \bar{x} x_n^T u + u^T \bar{x} \bar{x}^T u) \\ &= \frac{1}{N} \sum_n u^T (x_n - \bar{x})(x_n - \bar{x})^T u \end{aligned}$$

$$= \underbrace{u^T S_x u}_{2 \times 2}$$

Maximize  $S_y$  while keeping  $\|u\|=1$ , constant.

$$\Rightarrow S_x u = \lambda u$$

$$\Rightarrow S_y = u^T \lambda u = \lambda$$

for  $D > 2$ , choose eigenvectors corresponding to max eigen-value of the co-variance matrix.

## K-MEANS CLUSTERING

Scheme to divide  $N$  data points into  $K$  clusters  
 $\{x_j\}$                                      $\{c_i\}$

$$\tau_{ji} \text{ if } x_j \text{ is assigned to } c_i \\ 0, \text{ if } x_j \text{ is assigned to other clusters}$$

Distortion Measure:  $J = \sum_{j=1}^N \sum_{i=1}^K \tau_{ji} \| \vec{x}_j - \vec{\mu}_i \|^2$   
 (Loss function)

Find  $\tau_{ji}$  &  $\vec{\mu}_i$  which minimize  $J$ .

### Iterative Process

Step 0 - choose  $\vec{\mu}_i$  randomly from given data points.

Step 1 -  $\tau_{ji} = 1$  for  $(j, i)$  such that  $\vec{\mu}_i$  is the closest mean vector for given point  $\vec{x}_j$   
 else  $\tau_{ji} = 0$

Step 2 - for  $\tau_{ji}$  estimated in step 1, update  $\vec{\mu}_i$   
 $\nabla_{\mu_i} J = 0 \Rightarrow 2 \sum_{j=1}^N \tau_{ji} (\vec{x}_j - \vec{\mu}_i) = 0$   
 $\Rightarrow \vec{\mu}_i = \frac{\sum_j \tau_{ji} \vec{x}_j}{\sum_j \tau_{ji}}$

Iterate till convergence.  
 (guaranteed)  
 could be local minima

Number of  
points in cluster  
 $c_i$

Convergence is slow & several variations available.

### Lossy Data Compression

for each of the  $N$  points, only store the identity of the cluster to which it belongs.  
 & values of the  $K$  cluster centers,  $\vec{\mu}_i$ .  
 uses significantly less data if  $K \ll N$ .  
 also called vector quantization  
 &  $\vec{\mu}_i$  code-book vectors

## Example of data compression

Each image has  $N$  pixels of  $\{R, G, B\}$   
 values stored with 8-bit precision.  
 so each image needs  $24N$  bits

$K$ -classes needs  $\log_2 K$  bits

code-book vectors need  $24K$  bits

$$\therefore \text{Total no. of bits} = 24K + N \log_2 K.$$

$$\therefore \text{Compression Ratio} = \frac{24K + N \log_2 K}{24N} \times 100\%$$

$$= \left( \frac{K}{N} + \frac{1}{24} \log_2 K \right) \times 100\%$$

$$= 13.94\% \quad \text{for } K = 10 \quad \text{and } N = 100 \times 100$$

### K-Means

Unsupervised

$K = \text{no. of clusters}$   
 into which the given  
 data needs to be  
 classified.

### KNN

Supervised

$K = \text{no. of nearest neighbors}$   
 to use for classification  
 of a new data point

1

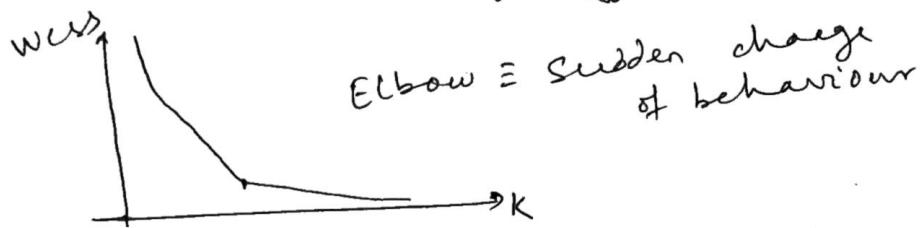
# K-Means & K-NN

## Qn K-Means

- What will you do if number of clusters is not known?

WCSS [within cluster sum of squares]

$$\sum_{i=1}^K \sum_{j=1}^N (x_j - \bar{x}_{ij} \mu_i)^2$$



- What can we do image clustering using this method?
  - ~~similar~~ Similar images may have different orientation
  - How to handle very large dataset?  
→ ~~large size~~ - use minibatches.
  - How to classify a new data point?

## K-NN

~~when k is large output happens to be noisy~~

small k - more noisy prediction  
large k - more stable prediction

- How to use k-NN for regression?

## K-Means

- Unsupervised
- K needs to be determined by elbow method
- Does not work well with clusters of different sizes
- Used for clustering
- Eager learner

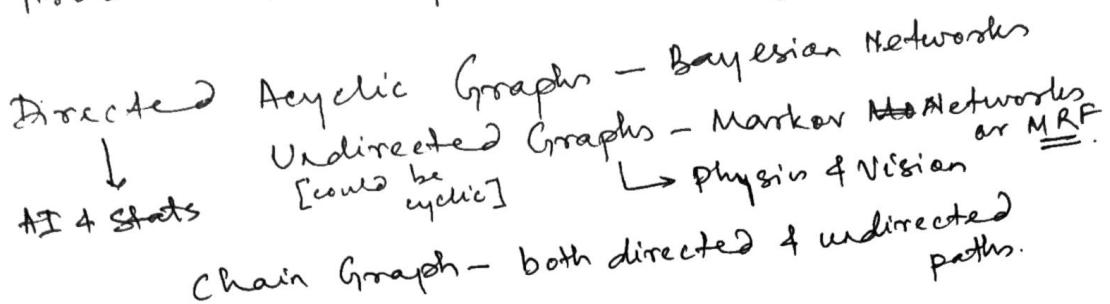
## K-NN

- Supervised
- $K \sim \sqrt{N}$
- Size difference is not a problem
- Used for classification or regression
- Lazy learner

# PROBABILISTIC GRAPH MODELS

[Bayesian Networks & Markov Network]

## Probability Theory + Graph Theory

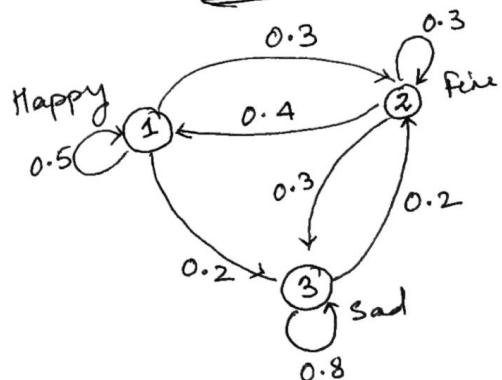


DAGs - Encode causal relationships  
Easier to learn from data  
But cannot represent cyclic dependencies

PGM ≠ Markov chain  
↓  
Nodes are elements of state space  
& edges denote transition probabilities  
Nodes are random variables & edges denote conditional dependencies.

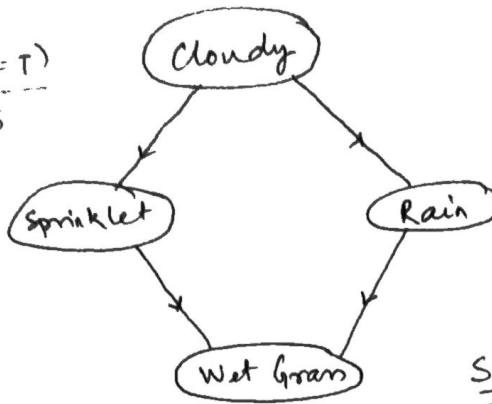
Markov: Future is independent  
of the past given the present.

Markov chain



	$P(C=F)$	$P(C=T)$
0.5		0.5

C	$P(S=F)$	$P(S=T)$		$P(R=F)$	$P(R=T)$
F	0.5	0.5		0.8	0.2
T	0.9	0.1		0.2	0.8



C	$P(R=F)$	$P(R=T)$
F	0.8	0.2
T	0.2	0.8

pfs,

$$P(C, S, R, W) = P(C)P(S|C)P(R|C, S)P(W|R, S)$$

$$= P(C)P(S|C)P(R|C)P(W|R, S)$$

[Simplification of conditional probabilities]

Most common task: probabilistic inference.

$$P(S=T|W=T) = \frac{P(S=T \wedge W=T)}{P(W=T)}$$

Bottom up ↑

$$\begin{aligned} &= \sum_{R, C \in T/F} P(S=T \wedge W=T \wedge R=F \wedge C=T/F) \\ &\quad \curvearrowleft P(W=T) \\ &\quad \curvearrowright \sum_{R, S, C \in T/F} P(S=T/F \wedge W=T \wedge R=T/F \wedge C=T/F) \end{aligned}$$

$$\begin{aligned} &\text{Top Down} \\ &\curvearrowleft P(W=T|C=T) \\ &P(W=T|C=T) \end{aligned}$$

$$= \frac{0.2781}{0.6471} = 0.430$$

∴ if grass is wet, it is more likely to be due to rain!

$$P(R=T|W=T) = \frac{0.4581}{0.6471} = 0.708$$

"Explaining Away"

$$P(S=T|W=T \wedge R=F) = 0.1945$$

Berkson's Paradox

S & R actually are independent  
Also called "selection bias"

# EXPECTATION MAXIMIZATION

## Bayesian Networks

Structure	Observability	Method
Known	Full	MLE
Known	Partial	Expectation Maximization (or Gradient Descent)
Unknown	Full	Search through Model space [NP-Hard]
Unknown	Partial	EM + Search Model space

Two coins with biases  $\theta_A$  &  $\theta_B$

Pick a coin at random & toss it 10 times  
Repeat this 5 times

Record  $X = (x_1, x_2, \dots, x_5)$  &  $(z_1, z_2, z_3, z_4, z_5) = Z$

$x_i \in \{0, 1, 2, \dots, 10\}$  = No. of heads observed

$z_i \in \{A, B\}$  = Coin chosen

Estimating  $\theta_A$  &  $\theta_B$  is easy if this full info is available  
But what if  $z_i$  are hidden variables??

Simple Approach -

Assume some initial values for  $\theta_A^{(0)}, \theta_B^{(0)}$   
Then estimate  $\{z_i\}$  based on which coin is more likely to generate the corresponding  $\{x_i\}$   
Estimate  $\theta_A^{(1)}, \theta_B^{(1)}$  using this completed info  
Repeat the above till convergence reached.

EM Approach - Compute probability of each possible completion of table.  
(create weighted training set of all these possibilities)  
Then use a modified form of MLE.

$$x = \{5, 9, 8, 4, 7\} \quad ?$$

## Simple Approach [MLE]

$$\text{Assume } \theta_A^{(0)} = 0.60 \quad \theta_B^{(0)} = \cancel{0.5} \quad 0.5$$

$$\theta_A^{(1)} = \frac{9+8+7}{30} = 0.8$$

## Expectation Maximization

another way  
to estimate  
MLE / MAP

<u>Maximization</u>						
$x_i$	$p(x_i/A)$	$p(x_i/B)$	$w_{Ai}$	$w_{Bi}$	A	B
5	0.2	0.25	0.45	0.55		
9	0.04	0.0098	0.80	0.20		
8	0.12	0.0044	0.73	0.27		
4	0.11	0.205	0.36	0.64		
7	0.22	0.123	0.65	0.35		

$$\hat{\theta}_A^{(1)} = \frac{\sum x_i w_{Ai}}{10 \sum w_{Ai}} = 0.712$$

$$\theta_B^{(1)} = \frac{\sum x_i w_{Bi}}{10 \sum w_{Bi}} = 0.583$$

$$= \frac{\sum x_i (1 - w_{Ai})}{10 \sum (1 - w_{Ai})} = \frac{\sum x_i - \sum w_{Ai}}{10(5 - \sum w_{Ai})}$$

After 10 steps converges to

$$\theta^{(10)} \approx 0.80$$

$$\theta_B^{(10)} \approx 0.52$$

May converge to local ~~maximum~~ for non-concave functions.  
but convergence is guaranteed.

# GAUSSIAN MIXTURE MODELS

Feature Extraction from speech data

Multi-object tracking

essentially, used when data is multi-modal, i.e. there are many "peaks" in the distribution.

$$P(\vec{x}) = \sum_{i=1}^K \phi_i \mathcal{N}(\vec{x} / \vec{\mu}_i, \Sigma_i)$$

Mean vector  
Covariance Matrix  
Multi-variate Gaussian Distribution.

$$\mathcal{N}(\vec{x} / \vec{\mu}_i, \Sigma_i) = \frac{1}{\sqrt{(2\pi)^K |\Sigma_i|}} \exp \left[ -\frac{1}{2} (\vec{x}_i - \vec{\mu}_i)^T \Sigma_i^{-1} (\vec{x}_i - \vec{\mu}_i) \right]$$

$$\sum \phi_i = 1$$

weights  
[ratio of points  
in each  
category]

Parameters estimated

using Expectation Maximization  
since analytic solution is  
usually not possible.

EM is guaranteed to converge

EM algo.

E-step : using  $\phi_i, \vec{\mu}_i & \Sigma_i$

estimate  $p(c_i | x_j)$  for each data point

↳ K classes/components

M-step : using assigned  $c_i$  for each  $x_j$ , [weighted]

update  $\phi_i, \vec{\mu}_i & \Sigma_i$

Initialisation :

Randomly pick  $\vec{\mu}_i$  from the given unlabelled data,  $x_j$

Set variance estimator to variance of whole data  
OR, use K-means to find initial parameter values

set  $\phi_i = \frac{1}{K} \cdot 1$

(speeds up considerably)

MNIST - K=10 components

Each image of size  $N \times N$  is a vector of size  $N^2$  of Bernoulli distributions (one per pixel) - PRML book

E-step

$$\gamma_{ji} = \frac{\phi_i N(\vec{x}_j | \vec{\mu}_i, \Sigma_i)}{\sum_{i=1}^K \phi_i N(\vec{x}_j | \vec{\mu}_i, \Sigma_i)} = p(c_i | \vec{x}_j)$$

M-step

$$\vec{\mu}_i = \frac{1}{N_i} \sum_{j=1}^N \gamma_{ji} \vec{x}_j$$

$$\vec{\Sigma}_i = \frac{1}{N_i} \sum_{j=1}^N \gamma_{ji} (\vec{x}_j - \vec{\mu}_i) (\vec{x}_j - \vec{\mu}_i)^T$$

updated value

$$\phi_i = \frac{N_i}{N}, \quad N_i = \sum_{j=1}^N \gamma_{ji}$$

: Effective no. of points in category  $i$

Log-Likelihood

$$\ln p(x | \mu, \Sigma, \phi) = \sum_{j=1}^N \ln \left\{ \sum_{i=1}^K \phi_i N(\vec{x}_j | \vec{\mu}_i, \Sigma_i) \right\}$$

Repeat the process till LL or parameters converge.

EM required since analytical solution of parameters which maximise LL not available.

K-means

hard assignment of  $\vec{x}_j$  to  $c_i$ ,  $\gamma_{ji} \in \{0, 1\}$   
 fast convergence even for high dimensional data  
 Assumes clusters to be spherical  
 Easier to interpret

GMM

Soft assignment  $0 \leq \gamma_{ji} \leq 1$

EM is slow

can handle arbitrary shapes more complex

Difficult to interpret