

# Machine Learning

(Gopakumar) T2.

(After T1 Module 2)

## Discriminant function

A classifier can be stated as :

$$c(x) = \arg \max_c g_c(x)$$

where for each class  $c$ , a discriminant function  $g_c(\cdot)$  is used to measure up to which degree objects belong to  $c$ .

## Discriminant Function analysis (DFA)

it is a statistical procedure that classifies unknown individuals and the probability of their classification into a certain group.

## Generative vs Discriminant approach

- Generative approaches estimate discriminant function by first estimating the probability distribution of data belonging to each class.
- Discriminative approaches estimate the discriminant function explicitly, without assuming probability distribution.

## Linear Discriminants

(case of 2 categories)

- Linear discriminant analysis (LDA) or normal discriminant analysis, or discriminant function analysis is generalization of Fisher's linear discriminant, a method used to find a linear combination of features that characterizes or separates two or more classes of objects or events.

$$LD = \alpha_1 X_1 + \alpha_2 X_2$$

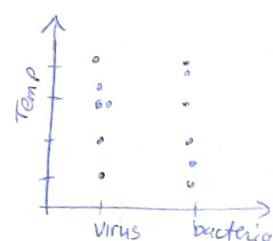
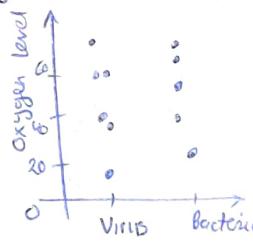
$\alpha_1, \alpha_2$  = weights

$X_1, X_2$  are features

e.g. a person can have either bacterial or viral infection we have his oxygen level & temp.

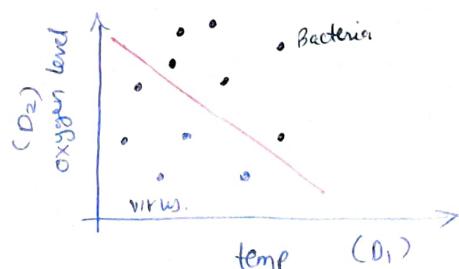
→ also we have old dataset for this

Now

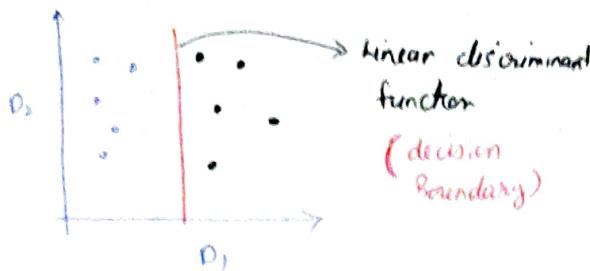


alone both features can't be used for analysis

so LD combines them



we straighten, rotate curve to get



So only  $D_1$  can be used

now on basis of  $D_1$  we find discriminant scores

Let got oxygen

$$LD = 0.11 \text{ CRP} + 0.7 \text{ temp}$$

to get score with center at 0.

$$LD = 0.11 (\text{CRP} - \overline{\text{CRP}}) + 0.70 (\text{Temp} - \overline{\text{Temp}})$$

$\Rightarrow$  So for viral  $LD < 0$

for bacteria  $LD > 0$

\* we need to find  $\alpha_1, \alpha_2$ .

$\Rightarrow$  class definition

A linear discriminant has following form,

$$g(x) = w^t x + w_0 \\ = \sum_{i=1}^d w_i x_i + w_0$$

[Decide  $w_1$  if  $g(x) > 0$  &  $w_2$  if  $g(x) < 0$ ]

if  $g(x) = 0$ , then  $x$  lies on decision boundary & can be assigned to either class

So end we need to find

$$w^t, w_0 \\ (\text{weight vector}), (\text{mean})$$

## Best Boundary

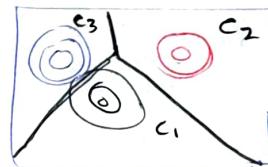
$\rightarrow$  Assume 2 classes can be separated by a linear boundary  $I(\theta)$  with some unknown parameters  $\theta$ .

$\rightarrow$  fit the best boundary to data by optimizing over parameters  $\theta$ .

### 2 methods

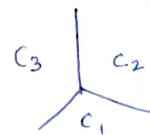
#### 1. Parametric methods

- Assume shape of density for classes is known  $p_1(x|\theta_1), p_2(x|\theta_2)$  ...
- Estimate  $\theta_1, \theta_2, \dots$  from data
- Use Bayesian classifier to find decision regions



#### 2) Discriminant Functions

- Assume discriminant functions are or known shape  $I(\theta_1), I(\theta_2)$  with parameters  $\theta_1, \theta_2, \dots$
- Estimate  $\theta_1, \theta_2, \dots$  from data.
- Use discriminant functions for classification



$\Rightarrow$  Estimating accurate density functions is much harder than estimating accurate discriminant functions

$\Rightarrow$  Discriminant functions can be more general than linear.

$\Rightarrow$  Linear discriminant functions are optimal for Gaussian distributions with equal covariance

## Decision boundary

$g(\mathbf{x}) = \omega^T \mathbf{x} + w_0 \rightarrow 0$  is a hyperplane  
 → set of vectors  $\mathbf{x}$  for which for some scalars  $\alpha_0, \dots, \alpha_d$  satisfy  $\alpha_0 + \alpha_1 \mathbf{x}^1 + \dots + \alpha_d \mathbf{x}^d = 0$

- A hyperplane is
  - a point in 1D
  - a line in 2D
  - a plane in 3D

Orientation of hyperplane is determined by  $\omega^*$  & its location by  $w_0$ .

→  $\omega$  is normal to hyperplane.

→ If  $w_0 = 0$ , it passes through origin.

- use "learning" algorithms to estimate  $\omega$  &  $w_0$  from training data  $x_k$ .
- The solution can be found by minimizing an error function, eg the "training error" or "empirical risk".

## LDF (Many classes)

- Suppose we have  $m$  classes.
- Define  $m$  linear discriminant functions  

$$g_i(\mathbf{x}) = \omega_i^T \mathbf{x} + w_0 \quad i=1, \dots, m$$
- Given  $\mathbf{x}$ , assign class  $c_i$  if  

$$g_i(\mathbf{x}) \geq g_j(\mathbf{x})$$

Such classifier is called linear machine

- A linear machine divides the feature space into  $m$  decision regions, with  $g_i(\mathbf{x})$  being the largest discriminant if  $\mathbf{x}$  is in the region  $R_i$ .

## Estimation

Let us suppose that

- 1) true class label (given to us)

$$z_k = \begin{cases} +1 & \text{if } x_k \in w_1 \\ -1 & \text{if } x_k \in w_2 \end{cases}$$

$w_1$  = class 1       $w_2$  = class 2

- 2) predicted class label

$$\hat{z}_k = \begin{cases} +1 & \text{if } g(x_k) > 0 \\ -1 & \text{if } g(x_k) \leq 0 \end{cases}$$

The solution can be found by minimizing an error function, eg "training error" or "empirical risk".

$$J(\omega, w_0) = \frac{1}{n} \sum_{k=1}^n [z_k - \hat{z}_k]^2$$

## geometric Interpretation

Using vector algebra.

$\mathbf{x}$  can be expressed as,

$$\mathbf{x} = \mathbf{x}_p + r \frac{\omega}{\|\omega\|}$$

$$g(\mathbf{x}) = \omega^T \mathbf{x} + w_0$$

$$= \omega^T (\mathbf{x}_p + r \frac{\omega}{\|\omega\|}) + w_0$$

$$= \omega^T \mathbf{x}_p + r \frac{\omega^T \omega}{\|\omega\|} + w_0$$

$$\Rightarrow \omega^T \omega = \|\omega\|^2 \quad \& \quad \omega^T \mathbf{x}_p + w_0 = 0$$

(projection)

So get

$$g(\mathbf{x}) = r \|\omega\|$$

⇒ distance of  $\mathbf{x}$  from hyperplane is given by

$$r = \frac{g(\mathbf{x})}{\|\omega\|}$$

• distance of plane from origin

$$n = 0$$

$$r = \frac{w_0}{\|w\|}$$

• Feature augmentation  
notation

Q Assume training set of samples as :-

x	class
0	1
1	0

for a bias of value 1, what will be weight of linear discriminant function

if it becomes

$$g(x) = y \begin{bmatrix} 0 \\ 1 \end{bmatrix} + 1$$

$$\text{if } y = -2$$

$$\text{for class 1. } -2 \cdot 0 + 1 = 1.$$

$$\text{for class 0. } -2 \cdot 1 + 1 = -1$$

both are different side of line with 0  
as center  
so choose  $y = -2$

### Notes

- LDF are optimal for gaussian distributions with equal covariance.
- Knowledge of class densities is not required when using linear discriminant functions. it is non-parametric approach.

### old problem

$$g(x) = w^T x + w_0$$

$$\begin{bmatrix} x_1 \\ \vdots \\ x_d \end{bmatrix}$$

### new problem

$$g(y) = a^T y$$
$$\begin{bmatrix} 1 \\ x_1 \\ \vdots \\ x_d \end{bmatrix}$$

## g2 Analysis for higher dimensions

we can consider line as

$$w_0 n_0 + w_1 n_1 + \dots + w_d n_d = 0$$

& a line it becomes

$$n_d = -\frac{w_1}{w_d} n_1 - \frac{w_0}{w_d}$$

$$y = m n + c.$$

& for more dimensions

$$w_0 + w_1 n_1 + w_2 n_2 + \dots + w_d n_d = 0$$

$$\text{or } \sum_{i=0}^d w_i n_i = 0$$

$$\Rightarrow [w_0 \ w_1 \ w_2 \ \dots \ w_d] \begin{bmatrix} 1 \\ n_1 \\ n_2 \\ \vdots \\ n_d \end{bmatrix} = 0$$

$$g(n) = \omega^T n$$

$$\text{eg } g = \omega^T x = 0.$$

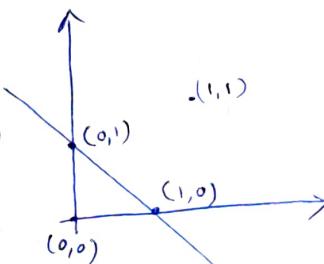
• if we got  $\omega^T = [1 \ 1 \ 1]$

$$\text{so } g(n) = [-1 \ 1 \ 1] \begin{bmatrix} 1 \\ n_1 \\ n_2 \end{bmatrix} = n_2 + n_1 - 1$$

for line  $n_2 + n_1 - 1 = 0$

$x_2 = 0$  point  $\rightarrow (1, 0)$   
gets  $x_1 = 1$

$x_1 = 0$  gets  $n_2 = 1$   $(0, 1)$



$$\text{check for } (1, 1), \ g(n) = 1 + 1 - 1 = 1 > 0$$

$$\text{for } (0, 0) = 0 + 0 - 1 < 0$$

$$\text{on line } (1, 0) = 1 + 0 - 1 = 0.$$

## Gradient descent

Gradient descent is an optimization algorithm used to find values of parameters (coefficients) of function ( $f$ ) that minimizes a cost function (cost).

Gradient descent is best used when parameters cannot be calculated analytically (eg using linear algebra) and must be searched for by an optimization algorithm

→ basically used to minimize some error function  $J(\alpha)$  with respect to  $\alpha$

$$\text{eg } J(\alpha) = \frac{1}{N} \sum_{k=1}^N [z_k - \hat{z}_{ik}]^2.$$

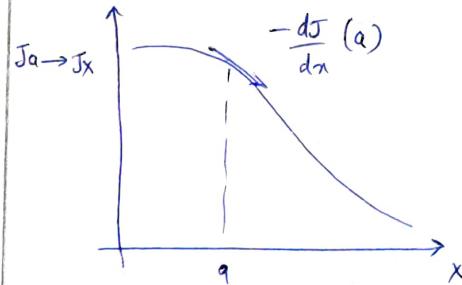
Gradient descent →

$$\alpha(k+1) = \alpha(k) + \eta(k) p_k$$

$p_k$  = search direction

$\eta(k)$  = learning rate

what is  $p_k$



$-\nabla J$  points in the direction of steepest decrease.

$$\text{so } p_k = -\nabla J(\alpha(k))$$

$$\text{so } \alpha(k+1) = \alpha(k) - \eta(k) \nabla J(\alpha(k))$$

## Learning rate

- if  $\alpha(k)$  is too small, it takes too many iterations
  - if it is too big, it might overshoot the solution (and never find it) possibly leading to oscillations (no convergence)
- ⇒ Gradient descent is very popular due to its simplicity but can get stuck in local minima.

## Difference

Gradient descent maximizes a function using knowledge of its derivative. Newton's method, a root finding algorithm, maximizes a function using knowledge of its second derivative. That can be faster when second derivative is known and easy to compute.

⇒ Newton's method is much faster than gradient descent.

## Newton descent

here

$$\alpha(k+1) = \alpha(k) + \eta(k) p_k$$

$$p_k = -H^{-1} \nabla J(\alpha(k))$$



$$\alpha(k+1) = \alpha(k) - \eta(k) H^{-1} \nabla J.$$

$H$  = Hessian matrix

$$H = \begin{pmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_1 \partial x_n} \\ \frac{\partial^2 f}{\partial x_2 \partial x_1} & \frac{\partial^2 f}{\partial x_2^2} & \cdots & \vdots \\ \vdots & \vdots & \ddots & \frac{\partial^2 f}{\partial x_n \partial x_1} \\ \frac{\partial^2 f}{\partial x_n \partial x_2} & \cdots & \cdots & \frac{\partial^2 f}{\partial x_n^2} \end{pmatrix}$$

second derivative

⇒ Gradient descent can be seen as a special case of Newton's method assuming

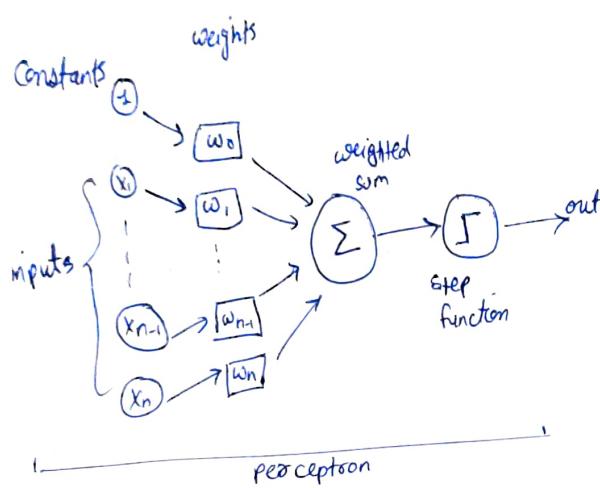
$$H = I$$

⇒ If  $J(\alpha)$  is quadratic  
Newton's method converges in one iteration

## Perceptron rule

Perceptron  $\rightarrow$  it is a single layer neural network, and a multi-layer perceptron is called a neural network.

Perception is a linear classifier (binary). Also is used in supervised learning. It helps to classify given input data.



- Perception rule minimizes following error function:

$$J_p(\alpha) = \sum_{y \in Y(\alpha)} (-\alpha^T y)$$

where  $Y(\alpha)$  is set of samples misclassified by  $\alpha$

- if  $Y(\alpha)$  is empty  $J_p(\alpha) = 0$   
otherwise  $J_p(\alpha) > 0$ .

- $J_p(\alpha)$  is  $\|\alpha\|$  times sum of distances of misclassified examples to decision boundary.

- $J_p(\alpha)$  is piecewise linear and thus suitable for gradient descent.

can be used in 2 ways

- Batch perception

here  $\alpha \leftarrow \alpha + \eta(k) \sum_{y \in Y(\alpha)} y$

Batch of misclassified examples.

$\rightarrow$  keep updating the orientation of hyperplane until all training samples are on positive side.

- Fixed increment single-sample perception.

if  $y_k$  is misclassified by  $\alpha$ .

$$\alpha \leftarrow \alpha + \eta(k) y_k$$

$\rightarrow$  update is done using one misclassified example at a time.

### Notes

- Batch algorithm leads to smoother trajectory in solution space.
- Single sample is easier to analyze and also concentrates more than necessary on any isolated "noisy" training samples.

## Perceptron Convergence theorem

→ if training samples are linearly separable  
then perceptron algorithm will terminate  
at a solution vector in finite number of  
steps.

# Machine learning

( Gopal Kumar)

- Machine learning algo's work by recognising patterns.

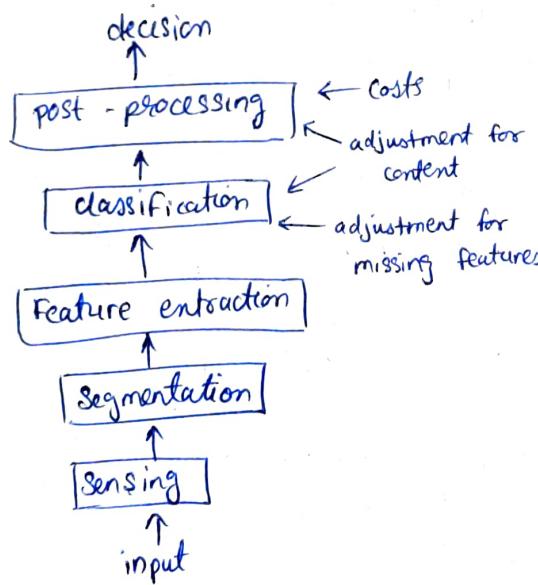
## Pattern recognition systems

### 1) Sensing

- Use of transducer (camera or microphone)
- PR system depends on bandwidth, resolution, sensitivity, distortion of transducer.

### 2) Segmentation & grouping

- Patterns should be well separated and should not overlap.



### 3) Feature extraction

- Discriminative features
- Invariant features wrt translation, rotation & scale.

### 4) Classification

Use a feature vector provided by feature extractor to assign the object to a category.

### 5) Post - processing

Exploit context, input dependent information other than from target pattern itself to improve performance.

## Bayesian Decision Theory

e.g.: catching sea bass and salmon fish

- state of nature, prior  
→ state of nature is random variable ( $\omega$ )  
→ The catch of salmon & sea bass is equiprobable

$$P(\omega_1) = P(\omega_2) \quad (\text{uniform priors})$$

$$P(\omega_1) + P(\omega_2) = 1 \quad (\text{exclusivity \& exhaustivity})$$

→ Decision rule with only prior information  
Decide  $\omega_1$  if  $P(\omega_1) > P(\omega_2)$  otherwise decide  $\omega_2$

### Use of class-conditional information

- $P(X|\omega_1)$  and  $P(X|\omega_2)$  describe diff. in lightness b/w populations of sea bass and salmon.

$X \rightarrow$  parameter / variable  
here it is lightness

→ Posterior, likelihood, evidence

$$P(\omega_j | X) = P(X | \omega_j) P(\omega_j) / P(X) \quad (\text{Bayes Rule})$$

- where in case of 2 categories

$$P(X) = \sum_{j=1}^{j=2} P(X | \omega_j) P(\omega_j)$$

## Loss function

- A loss function states exactly how costly each action is.

- As earlier, we have classes  $\{\omega_1, \omega_2, \dots, \omega_c\}$
- we also have possible actions  $\{a_1, \dots, a_d\}$
- The loss function  $\lambda(a_i/\omega_j)$  is loss function incurred for taking actions  $a_i$  when class is  $\omega_j$ .

e.g.: Zero-One loss function

$$\lambda(a_i/\omega_j) = \begin{cases} 0 & i=j \\ 1 & i \neq j \end{cases}$$

It assigns no loss to correct decision & uniform unit loss to an incorrect decision.

- expected loss (conditional risk)

$$R(a_i/x) = \sum_{j=1}^c \lambda(a_i/\omega_j) P(\omega_j/x)$$

Zero-one conditional risk

$$\begin{aligned} R(a_i/x) &= \sum_{j \neq i} P(\omega_j/x) \\ &= 1 - P(\omega_i/x) \end{aligned}$$

Hence for an observation  $x$ , we can minimize expected loss by selecting action that minimize conditional risk.

- overall risk is expected loss associated with given decision rule.

$$R = \int R(a(x)/x) p(x) dx$$

• arg prob of error

$$P(error) = \int_{-\infty}^{\infty} P(error/x) p(x) dx$$

## Pattern Classifiers

### version 1

#### Discriminant functions

- Discriminant functions are useful ways of representing pattern classifiers

- Let  $g_i(x)$  is discriminant function for  $i$ th class

This classifier will assign class  $w_i$  to feature vector  $x$  if

$$g_i(x) > g_j(x)$$

or equivalently

$$i^* = \arg \max_i g_i(x) \text{ decide } w_i^*$$

General case with risk

$$\begin{aligned} g_i(x) &= -R(d_i/x) \\ &= -\sum_{j=1}^k \lambda(\alpha_i/\omega_j) P(w_j/x) \end{aligned}$$

$\Rightarrow$  for zero-one loss fn, Bayes discriminant can be further simplified:

$$g_i(x) = P(w_i/x)$$

$$g_i(x) = P(w_i/x) = \frac{P(x/w_i) P(w_i)}{\sum_j P(x/w_j) P(w_j)}$$

Decision boundaries separate regions, they are ties among discriminant functions

2 category discriminant

$$g(x) = g_1(x) - g_2(x)$$

decide  $w_1$  if  $g(x) > 0$ , otherwise decide  $w_2$

various manipulations of discriminant

$$g(x) = P(w_1/x) - P(w_2/x)$$

$$g(x) = \ln \frac{P(w_1/x)}{P(x/w_2)} + \ln \frac{P(w_2)}{P(w_1)}$$

#### Normal density

- Continuous univariate normal, or Gaussian density

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

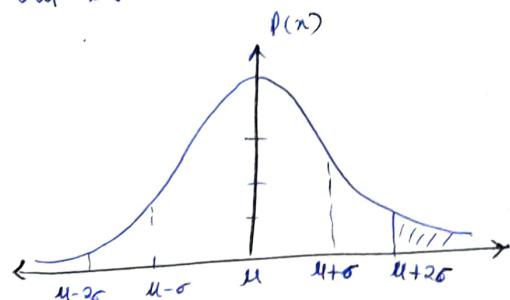
- mean is expected value of  $x$  is

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

- variance is expected squared deviation

$$\sigma^2 = E[(x-\mu)^2] = \int_{-\infty}^{\infty} (x-\mu)^2 p(x) dx$$

$\Rightarrow$  Samples from normal density tend to cluster around mean and be spread out based on variance



$\Rightarrow$  Normal density is completely specified by mean & variance. These 2 are its sufficient statistics.

$$p(x) \sim N(\mu, \sigma^2)$$

# MLE

① Pearson's method

MLE provides the parameter value(s) that makes the observed sample the most likely sample among all possible samples.

$\hat{\theta}$  is maximum likelihood estimate if  $L(\hat{\theta}) > L(\theta_0)$

for all values of  $\theta_0$  in parameter space

We want to find peak, and its position

To find MLE, we need to differentiate the log-likelihood  $\ln \ell$  & set it equal to 0.

$\frac{d}{d\theta} \ln \ell(\theta) = 0$  is called score eq<sup>n</sup>

our confidence in MLE is quantified by the pointedness of log-likelihood.

$I_\theta(\theta) = -\frac{d^2}{d\theta^2} \ln \ell(\theta)$  This is called observed information

Taking the expected value gives us expected information

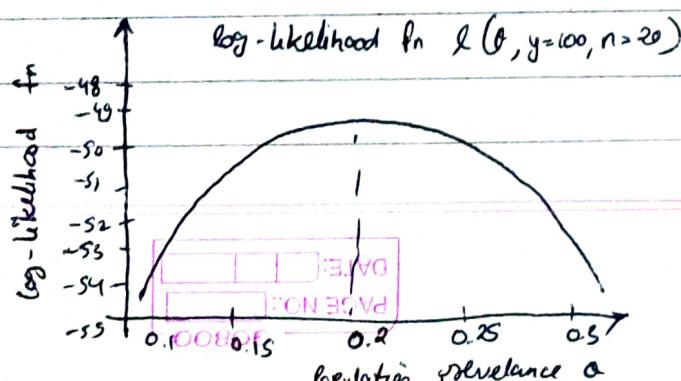
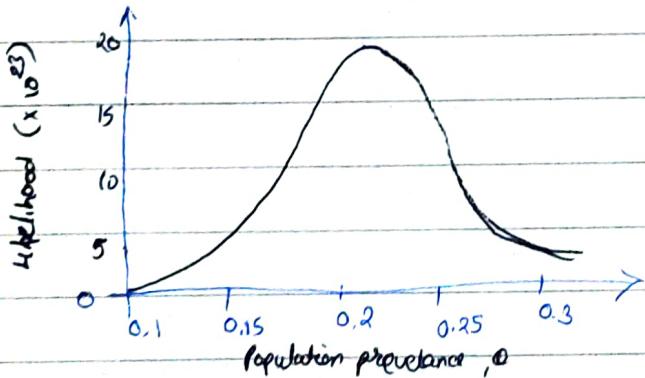
$$I(\theta) = E[I_\theta(\theta; y)]$$

variance of estimates

$$V(\hat{\theta}) \approx I(\theta)^{-1}$$

$$\ln \ell(\theta) = n \ln \theta + (y-n) \ln(1-\theta)$$

likelihood function  $\ell(\theta, y=100, n=20)$



Q Find MLE with  $y=100, n=20$  & find 95% confidence interval.

$$\begin{aligned} \frac{d}{d\theta} \ln \ell(\theta) &= n - \frac{y-n}{\theta} = 0 \Rightarrow \theta = \frac{y-n}{n} = \frac{100-20}{20} = 0.2 \\ &= \frac{y}{n} = \frac{100}{20} = 0.2 \end{aligned}$$

$$I_\theta(\theta; y) = -\frac{d^2}{d\theta^2} \ln \ell(\theta) = -\frac{n}{\theta^2(1-\theta)} = \frac{20}{0.2^2(1-0.2)} = 250$$

$$I(\theta) = \frac{n}{\theta^2(1-\theta)} \quad \text{95% CI} \quad \theta \pm 1.96 \sqrt{\frac{I(\theta)}{n}}$$

$$= 0.2 \pm 1.96 \sqrt{\frac{20}{0.2^2(1-0.2)}} = 0.2 \pm 1.96 \sqrt{\frac{20}{0.2^2(0.8)}} = 0.2 \pm 1.96 \sqrt{\frac{20}{0.16(0.8)}} = 0.2 \pm 1.96 \sqrt{\frac{20}{0.128}} = 0.2 \pm 1.96 \sqrt{156.25} = 0.2 \pm 1.96 \cdot 12.5 = 0.2 \pm 24.75 = [-24.5, 24.9]$$

## Properties of MLE

### Explicit vs Implicit MLE

- If it is possible to solve the score equation to get an expression for MLE, the MLE is explicit.
- If software algorithms are required to generate MLE expression, the MLE is implicit.

### Biasedness & consistency

The MLE is not an unbiased estimator necessarily.

MLE is asymptotically biased & consistent.  
The asymptotic distribution is normal

$$\hat{\theta} \stackrel{d}{\approx} N(\theta, I(\hat{\theta})^{-1})$$

### Asymptotic efficiency

MLE is asymptotically efficient. MLE is optimal in all but a handful of unique situations.

$$n - n = \bar{x}$$

### Invariance

$$\text{MLE for } g(\theta) = g(\hat{\theta})$$

$$P(B|A) = \frac{P(A \cap B)}{P(A)}, P(A) > 0$$

CR < 0.1  $\rightarrow$  calculation is safe

$$CR - CI / RI$$

defining most optimal solution for a problem with given constraints  
linear programming is good  
soft. power is good

### Decision Making

Single Criterion & Multiple Criteria

Adversarial search

Decision Tree + Taboo + M.L.

T  $\rightarrow$

systolic blood pressure

A sample of 5 pregnant women have their SBP taken, which is considered to be normally distributed.

$$\text{Sample: } \{135, 123, 120, 102, 110\}$$

find maximum likelihood estimate for  $\mu$  &  $\sigma^2$ .

$$L(\mu, \sigma^2) = \prod_{i=1}^5 f_i(y_i; \mu, \sigma^2) \text{ where } f(y_i; \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y_i - \mu)^2}{2\sigma^2}}$$

$$L(\mu, \sigma^2) = \left( \frac{1}{2\pi\sigma^2} \right)^{\frac{n}{2}} \exp \left( -\sum_{i=1}^n \frac{(y_i - \mu)^2}{2\sigma^2} \right)$$

$$l(\mu, \sigma^2) = \dots$$

$$\frac{dl}{d\mu} = \frac{2(\sum y - n\mu)}{2\sigma^2} = \frac{\sum y - n\mu}{\sigma^2} = 0$$

$$\left[ \hat{\mu} = \frac{\sum y}{n} \right]$$

$$\frac{dl}{d\sigma^2} = \frac{-n}{2\sigma^2} + \frac{1}{2(\sigma^2)^2} (\sum y^2 - 2\mu \sum y + n\mu^2) = 0$$

$$\hat{\sigma}^2 = \left( \frac{1}{n} \right) \left( \sum y^2 - 2\hat{\mu} \sum y + n\hat{\mu}^2 \right)$$

$$= \frac{\sum y^2}{n} - \frac{2\sum y \hat{\mu}}{n} + \frac{(\sum y)^2}{n^2}$$

$$= \frac{\sum y^2}{n} - \left( \frac{\sum y}{n} \right)^2$$

$$\frac{\sum (y - \hat{\mu})^2}{n}$$

PAT. NO. [ ] CODE [ ] DATE [ ] its position  
the

## Parametric ML algorithms

- > 1) Regression, 2) LDA
- 3) Perception

### Advantage

Simple

fast

less data

### Disadvantage

Constrained

Limited complexity

Poor fit

## Non Parametric Algo

- 1) Decision tree
- 2) Naive Bayes
- 3) SVM
- 4) Neural network

### Advantage

1) flexibility

2) power

3) performance

### Disadvantage

More Data

slower

overfitting

In parametric model, no. of parameters is fixed with respect to sample size.

In nonparametric model, the (effective) number of parameters can grow with sample size.