

16 September 2014

Probability

$$P(A) = 0.2$$

$$P(B) = 0.3$$

$$P(A \cap B) = P(A \cap B)$$

if A and B are independent,

$$P(A \cap B) = P(A) \cdot P(B)$$

→ conditional probability

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

$$\text{so } P(A \cap B) = P(A|B) \cdot P(B)$$

$$= P(B|A) \cdot P(A)$$

$$P(A|B) = P(B|A) \cdot P(A)$$

$$\text{and } P(B|A) = \frac{P(A \cap B)}{P(A)}$$

Baye's Theorem

$$P(M|H) = \frac{P(H|M)P(M)}{P(H|M)P(M) + P(H|F)P(F)}$$

evidence

class labeled

feature value

posterior prob.

prior prob.

likelihood

height

height

$$\begin{aligned} P(F | H=64) &= 1 - P(M | H=64) \\ &= \frac{P(H=64 | F) P(F)}{P(H=64 | M) P(M) + P(H=64 | F) P(F)} \end{aligned}$$

↑
feature vector

$$P(w | x) = \frac{P(x | w) p(w)}{p(x)}$$

discrete = p
continuous = p unconditional density of x

19 September 2014

Test for a rare disease (% million)

$$\begin{array}{lll} \text{Sensitivity} = 100\% & \text{False-ve (miss)} = 0\% & \text{Recall} \\ \text{specificity} = 99.9\% & \text{False +ve} = 0.1\% & \text{Precision} \end{array}$$

if you do not have the disease the test will say -ve 99.9% of the time
 if you have a disease you will be caught

$$\begin{array}{cccc} \cancel{D} & \cancel{ND} & 100\% & \\ \cancel{D} & D & 0 & \\ \cancel{ND} & \cancel{D} & 0.1 & \\ \cancel{ND} & ND & 99.9 & \end{array} \quad \begin{array}{c} 0.1 \\ 100 \\ 100 - 100 + 0.1 \end{array} \quad \begin{array}{l} P(F | D) = 1 \\ P(+ | ND) = 0.001 \end{array}$$

$$P(D | F) = \frac{P(+ | D) P(D)}{P(+)} = \frac{P(+ | D) P(D)}{P(+ | D) P(D) + P(+ | ND) P(ND)}$$

$$\begin{aligned} &= \frac{P(+ | D) P(D)}{P(+ | D) P(D) + P(+ | ND) P(ND)} = \frac{10^{-3}}{10^{-3} + 10^{-6}(1-10^{-6})} \\ &= 0.001 \end{aligned}$$

$$P(\text{error}) = ?$$

$$P(\text{error})$$

$$= \int_{-\infty}^{\infty} P(\text{error}, x) dx$$

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

if d-dimensional vector,
integrate d times of
one over each dimension

Bayes classifier

$\xrightarrow[2 \text{ classifiers}]{\text{classifier}} P(w_1|x) = 0.7$

$$P(w_2|x) = 0.3$$

we picked w_1 ($> \text{prob}$)

$$P(\text{error}|x) = 0.3$$

$$P(\text{error}|x) = \min(P(w_1|x), P(w_2|x))$$

minimize total error,

minimize the integral,

minimize $P(\text{error}|x)$

Bayes classifier is the best in
the world.

→ Nearest Neighbour Classifier

$$P(\omega_1 | x) = \frac{P(x | \omega_1) P(\omega_1)}{P(x)}$$

$$> \frac{P(x | \omega_2) P(\omega_2)}{P(x)}$$

assign it there.

$P(x | \omega_1) P(\omega_1) \rightarrow P(x | \omega_2) P(\omega_2)$

(likelihood prior)
 Baye's classifier
 probability in the world.
 (impractical)

N parameters, $10N$ samples.

Professor Pernar * ; Baye's error rate.

23 September 2014

I was Absent!
 → Student statistics.

conditional Ground truth

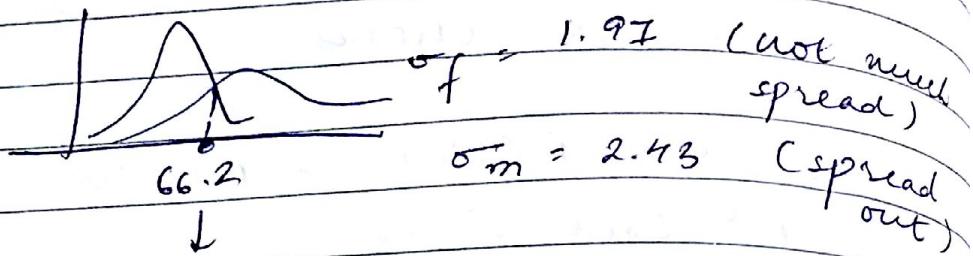
Maximum likelihood Estimate

posterior probability $P(h|w)P(w) / P(h)$

Gaussian

Now we have prob, we want pdf, how to find it?

lets assume they are Gaussian curves



decision point

$$n < 66.2 \quad F$$

$$\text{else } n \quad -\frac{1}{2} \frac{(x-\mu)^2}{\sigma^2}$$

$$\text{Gaussian: } \frac{1}{\sqrt{2\pi}\sigma} e$$

for the two classes, μ and σ are different

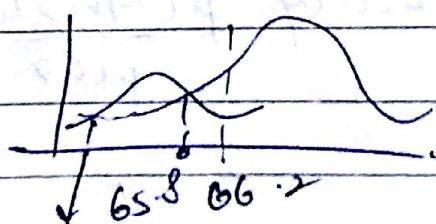
Prior scaled likelihood

$$P(W_1|W) P(W)$$

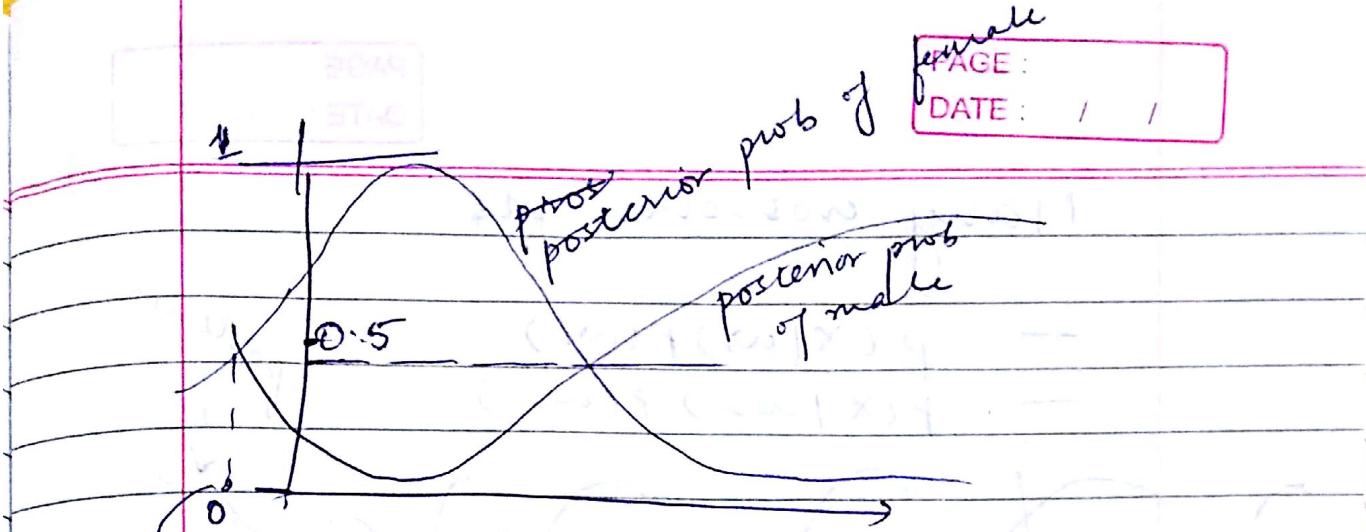
case (w_1 by 0.35 $\rightarrow P_f = 0.8$
 w_2 by 0.65 $\rightarrow P_m = 0.6$)

will the pt. of intersection change

Yes



another crossover pt. @ 45.7 , mate
 $n < 3 \times 10^{11}$

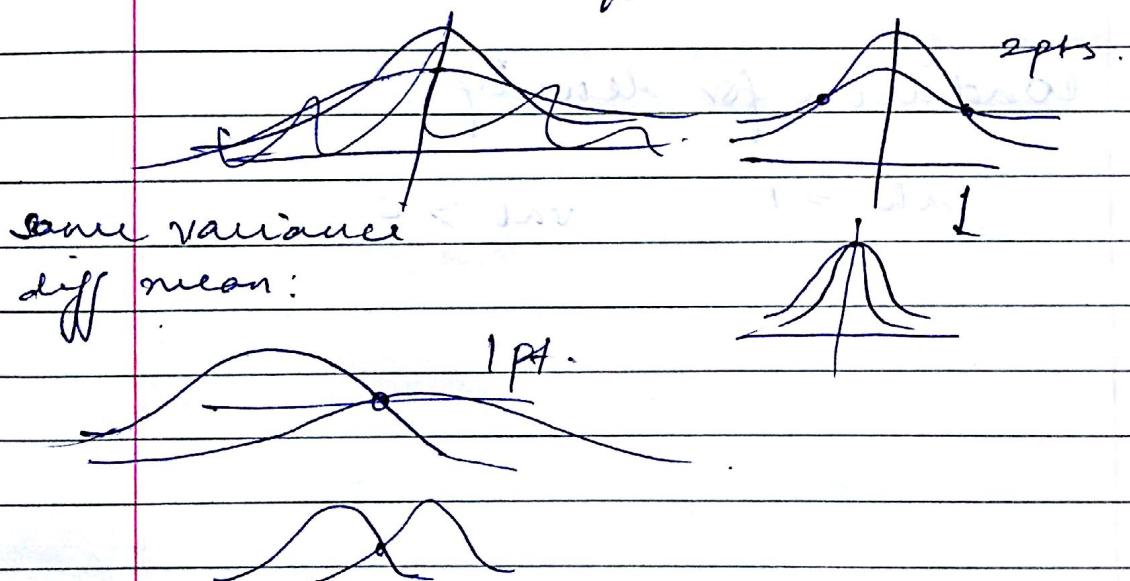


→ @ crossover pt. the posterior
prob = 0.5

1D - 2 pt of separation: non linear
2D: linear

Bayesian densities can be any
densities

Condition for only one pt. of separation
same mean, diff variance

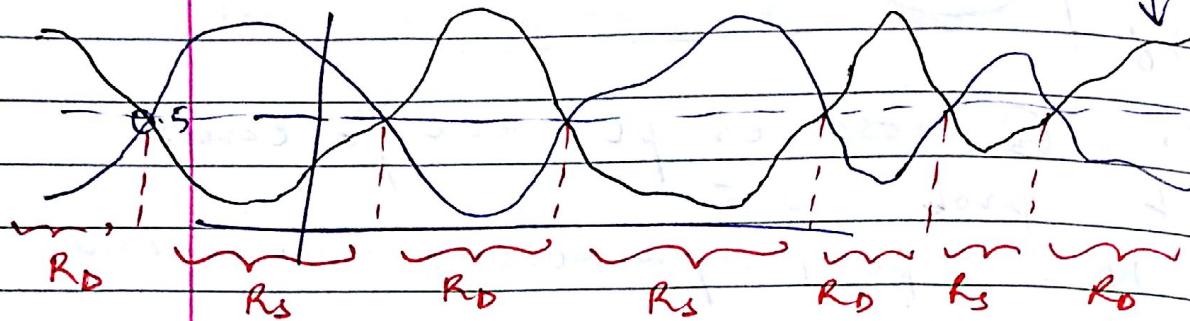


Many crossover pts.

$$- : p(x|w_1) p(w_1)$$

$$- : p(x|w_2) p(w_2)$$

posterior

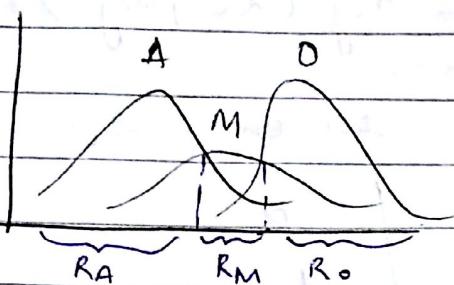


At crossover, value is 0.5
if the curve is going up.
it will continue to go up

condition for density:

$$\int dh = 1 \quad val > 0$$

26 September 2014



$$P(\omega_i | x) = \frac{P(x|\omega_i) P(\omega_i)}{P(x|\omega_1) P(\omega_1) + P(x|\omega_2) P(\omega_2) + P(x|\omega_3) P(\omega_3)}$$

* Probability theory can be extended to multiple classes.

$$P(\omega_i | x) = \frac{\sum_{i=1}^c P(x|\omega_i) P(\omega_i)}{\sum_{i=1}^c P(x|\omega_i) P(\omega_i)}$$

We are learning each class independently.

→ feature vector : n dimensional
decision boundary : n-1 dimensional
using in n dimensional space

Equidense surfaces.

→ giving credit card : loss func.

| ω_i | $c\omega$ | x | d_i | $\lambda_{ij} (d_i \omega_j)$ |
|---------------|-----------|-----|-------|-------------------------------------|
| | | 0 | 1 | |
| credit worthy | new | -1 | 0 | loss func. → zero one loss func. |

minimizing loss : $\lambda_{ij} P(\omega_j | x)$

$R(d_i | x)$ ←
↑ expected
risk (total loss)
of taking action d_i)

$$R(x_i|x) = \sum_j \gamma_{ij} (x_i|w_j) P(w_j|x)$$

action x_i

| | | A | B |
|-----------------------|---|---|---|
| | | 0 | 1 |
| ground truth w_j | A | 1 | 0 |
| | B | 0 | 1 |

$$R(x_1|x) = \lambda_{11} P(w_1|x) + \lambda_{12} P(w_2|x)$$

$$R(x_2|x) = \lambda_{21} P(w_1|x) + \lambda_{22} P(w_2|x)$$

take action that action which minimizes the risk.

$$x_1 \quad \lambda_{11} P(w_1|x) + \lambda_{12} P(w_2|x) < \lambda_{21} P(w_1|x) + \lambda_{22} P(w_2|x)$$

$$x_1 \quad (\lambda_{11} - \lambda_{21}) P(w_1|x) < (\lambda_{22} - \lambda_{12}) P(w_2|x)$$

$$\frac{P(w_1|x)}{P(w_2|x)} < \frac{\lambda_{22} - \lambda_{12}}{\lambda_{11} - \lambda_{21}}$$

$\Leftrightarrow 1$

likelihood ratio

$$\frac{P(x|w_1)}{P(x|w_2)} > \frac{\lambda_{12} - \lambda_{22}}{\lambda_{21} - \lambda_{11}} \frac{P(w_2)}{P(w_1)}$$

Gaussian Multivariate Density

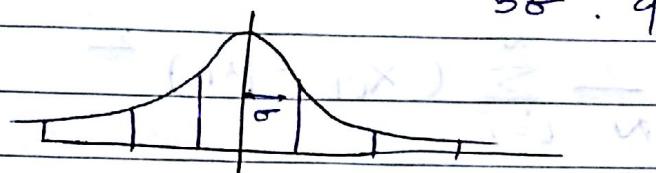
$$P(\omega|x) = p(x|\omega) P(\omega)$$

$p(x)$

$$= \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

$$N(\mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

$\sigma : 99.99\%$



x : vector of two dimensions

μ : & vector too!

$$\mu = \begin{bmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_d \end{bmatrix}$$

$$\Sigma = \begin{bmatrix} \sigma_1^2 & & & \\ & \sigma_2^2 & & \\ & & \ddots & \\ & & & \sigma_d^2 \end{bmatrix}$$

mean height 160 cm 15 cm σ
70 kg 15 kg

$$p(x|\omega) = \frac{1}{(2\pi)^d} \frac{1}{|\Sigma|^{1/2}} e^{-\frac{1}{2} (x-\mu)^\top \Sigma^{-1} (x-\mu)}$$

$$= \frac{1}{(2\pi)^d} \frac{1}{|\Sigma|^{1/2}} e^{-\frac{1}{2} (x-\mu)^\top \Sigma^{-1} (x-\mu)}$$

weight

long.

160 mm

height

$$\text{variance} = \frac{1}{N} \sum_{i=1}^N (x_{i1} - \mu_1)^2$$

$$\sigma_{11}^2 = \frac{1}{N} \sum_{i=1}^N (x_{i1} - \mu_1)^2$$

variance of the quantity x_{i1}
 how does the value of x varies
 around μ .

→ covariance

when 1 quantity changes how
 does the second quant change.

$$\sigma_{12} = \frac{1}{N} \sum_{i=1}^N (x_{i1} - \mu_1)(x_{i2} - \mu_2)$$

+ve correlated

or correlated

indep.

$$\Sigma = \begin{bmatrix} \sigma_{11} & \sigma_{12} \\ \sigma_{21} & \sigma_{22} \end{bmatrix} \quad \text{covariance matrix.}$$

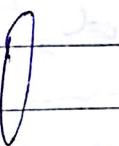
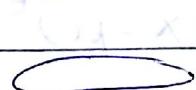
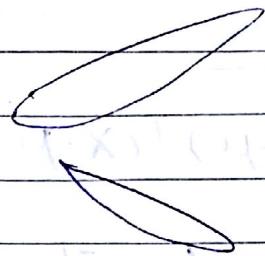
$$[\sigma_{ij}]_{d \times d}$$

Identify covariance matrix:
circle

+ve

-ve

diff diag
off diag: 0



27 September 2014

Anoop ??

→ unimodal density

→ Minimum Risk: 2-category

$$\Sigma = \begin{bmatrix} \sigma_{11} & \sigma_{12} & \dots & \sigma_{1d} \\ \vdots & & & \vdots \\ \sigma_{d1} & \dots & \dots & \sigma_{dd} \end{bmatrix}_{d \times d}$$

$$\mu = \begin{bmatrix} \mu_1 \\ \vdots \\ \mu_d \end{bmatrix}_{1 \times d}$$

$$N(\mu, \Sigma) = \frac{1}{\sqrt{2\pi}^d} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

↑ univariate
in multivariate case?

$\mu \rightarrow x$

Euclidean distance $d^2 = (x - \mu)^T (x - \mu)$

$$d_2^2 = (x - \mu)^T \Sigma^{-1} (x - \mu)$$

scaled variant
of the distance

PC Mahalanobis

$$-y_2^T (x - \mu) \Sigma^{-1} (x - \mu)$$

$$n(\mu, \Sigma) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} e^{-y_2^T (x - \mu) \Sigma^{-1} (x - \mu)}$$

Σ : symmetric
positive definite

I can do it

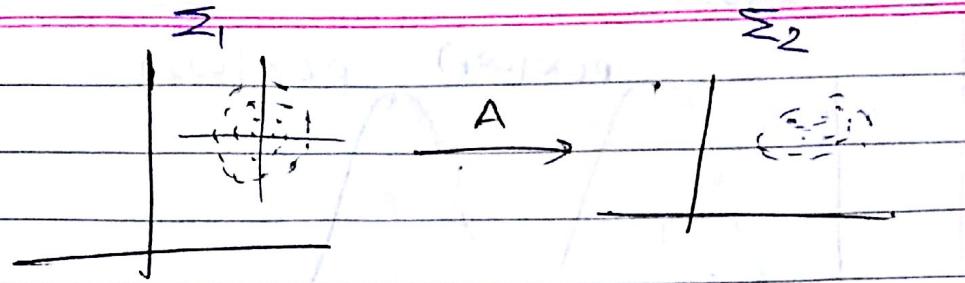
Eigen

$x^T A x > 0$

$d \times d$

1×1

transformed x



$$\Sigma_2 = A^T \Sigma_1 A$$

$$j \quad \Sigma_1 = I$$

$$\Sigma_2 = A^T A$$

Cholesky Decomposition
for positive definite symmetric
matrix

2.6. discriminant func. for
one Normal density

log / sqrt / any monofunc
func.

for normal dist.

$\text{mean} = 0$ $\sigma^2 = 1$

$\text{mean} = 1$ $\sigma^2 = 1$

$\text{mean} = 0$ $\sigma^2 = 2$

$\text{mean} = 1$ $\sigma^2 = 2$

$\text{mean} = 0$ $\sigma^2 = 3$

$\text{mean} = 1$ $\sigma^2 = 3$

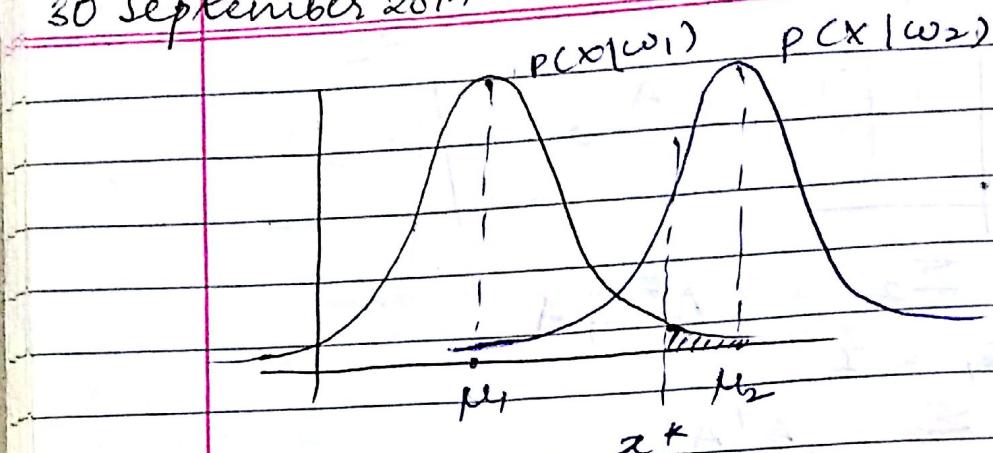
$\text{mean} = 0$ $\sigma^2 = 4$

$\text{mean} = 1$ $\sigma^2 = 4$

$\text{mean} = 0$ $\sigma^2 = 5$

$\text{mean} = 1$ $\sigma^2 = 5$

30 September 2014



Because of a noisy channel,
we send μ_1 but we get $\mu_1 +$
noise.

x^* : error/noise

what should unlock cannot $1 \rightarrow 0$
what shouldn't unlock can $0 \rightarrow 1$

x^* : noisy decision boundary

① $x > x^*, w_2$: hit $2+1$

② $x < x^*, w_2$: miss

③ $x > x^*, w_1$: false alarm $2+1$

④ $x < x^*, w_1$: correct rejection

$w_2 \rightarrow 0$ $w_2 \neq 1$) $0 \ 0 \rightarrow$ correct rejection

$w_2 \neq 0$ $w_2 \rightarrow 1$) $0 \ 1 \rightarrow$ false alarm

$w_1 \rightarrow 0$ $w_1 \neq 1$) $1 \ 0 \rightarrow$ miss

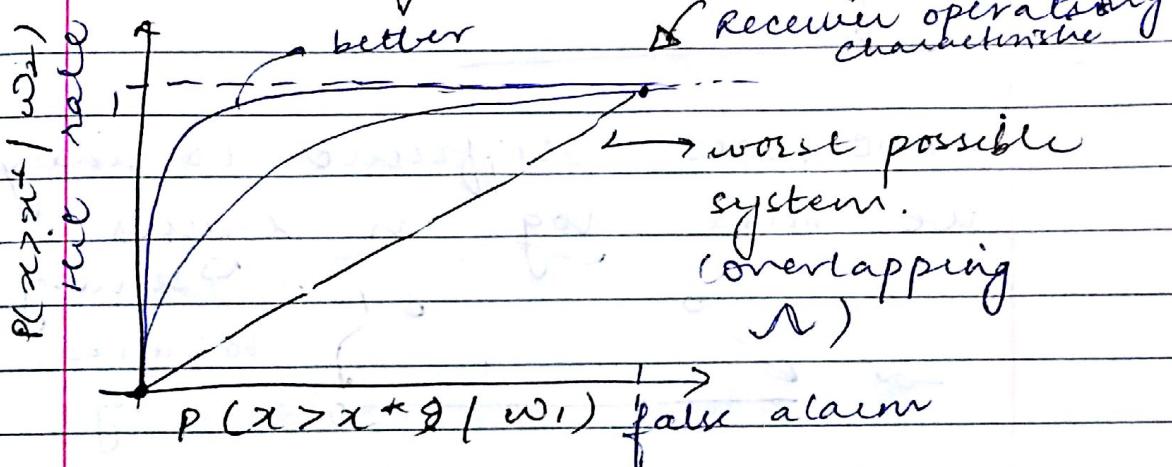
$w_1 \neq 0$ $w_1 \rightarrow 1$) $1 \ 1 \rightarrow$ hit

what is the prob. of a false
alarm, given threshold is x^*

prob. when density is given:
area under the curve.

To understand the classifier,
we need one of hit; miss
and one of false alarm and
correct rejection.

ROC plot



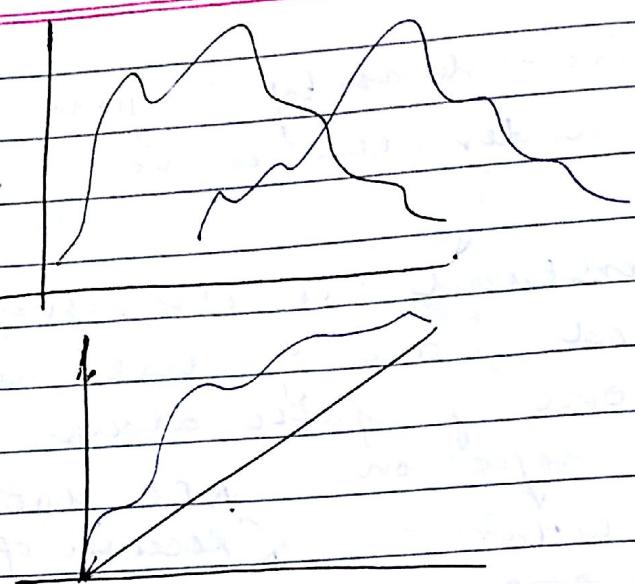
operating point of the system
dependent on the scenario.

$$d' = \frac{|\mu_2 - \mu_1|}{\sigma} \rightarrow \text{or } \frac{\sigma_1 + \sigma_2}{2}$$

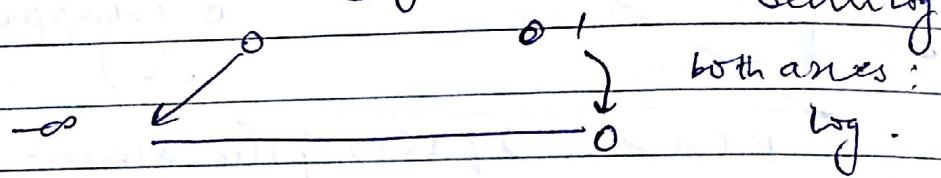
more the d' : better the curve

The same idea works even if
the distribution is not normal

Probability Density Function



• ROC plots: difficult to analyze
we take log in x-axis



98%

Equal error rate

confusion matrix

ground truth

$w_1, w_2, w_3,$

off

w_1

w_2

w_3

best: diag ≈ 1 off diag = 0

Practically impossible to learn
a complex density function.

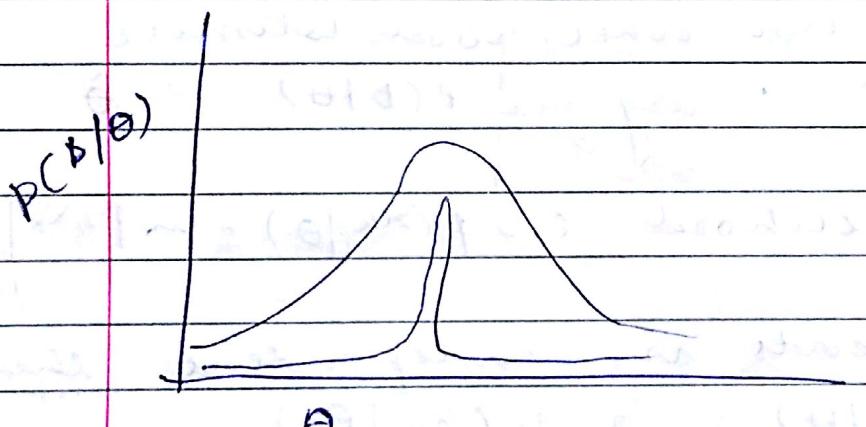
We know what dist is normal,
 $\sigma = 10$.

$\mu = ?$

$$P(D | \theta) = \prod_{k=1}^n p(x_k | \theta)$$

complete dataset

Sampling is iid
independent & identically
distributed



Likelihood of theta generally the
sample.

$$\hat{\theta} = \underset{\theta}{\operatorname{argmax}} \phi(D|\theta)$$

estimated $\hat{\theta}$ via

$\hat{\theta}$ (MLE)
Maximum estimate
likelihood discriminant

for uniform dist



$$a = \min \{D\}$$

$$b = \max \{D\}$$

7 October 2014

Maximum likelihood Estimate

$$(MLE) : \underset{\theta}{\operatorname{argmax}} \phi(D|\theta) = \hat{\theta}$$

$$\log \text{likelihood} \ln p(x_k | \theta) = \ln p(x_k | \mu, \Sigma)$$

If events are independent. Then,

$$P(D|\theta) = \prod_{i=1}^n p(x_i | \theta)$$

$$p(x_k | \mu, \Sigma) = \frac{e^{-\frac{1}{2}(x_k - \mu)^T \Sigma^{-1} (x_k - \mu)}}{(2\pi)^{d/2} |\Sigma|^{1/2}}$$

$$\ln p(x_k | \mu, \Sigma) = -\frac{1}{2} \ln [(2\pi)^d |\Sigma|] - \frac{1}{2} (x_k - \mu)^T \Sigma^{-1} (x_k - \mu)$$

For MLE :

Differentiate and put to '0'

$$\nabla_{\mu} (\ln(p(x_k | \mu))) = -\frac{1}{2} \times \Sigma^{-1} \times -2(x_k - \mu)$$

$$\begin{bmatrix} (x_1 - \mu_1) \\ (x_2 - \mu_2) \\ \vdots \\ (x_d - \mu_d) \end{bmatrix}^T \begin{bmatrix} (x_1 - \mu_1) \\ \vdots \\ (x_d - \mu_d) \end{bmatrix} = (x_1 - \mu_1)^2 + (x_2 - \mu_2)^2 + (x_3 - \mu_3)^2 + \dots + (x_d - \mu_d)^2$$

$$\nabla_{\mu} (\cdot) = \begin{bmatrix} -2(x_1 - \mu_1) \\ -2(x_2 - \mu_2) \\ \vdots \\ -2(x_d - \mu_d) \end{bmatrix} = -2(X - \mu)$$

$$\rightarrow = \Sigma^{-1} (X - \mu)$$

$$\ln p(D | \theta) = \sum_{k=1}^n \left(\frac{1}{2} \ln [(2\pi)^d |\Sigma|] - \frac{1}{2} (x_k - \mu)^T \Sigma^{-1} (x_k - \mu) \right)$$

$$\therefore \nabla_{\mu} (\ln p(D | \theta)) = \sum_{k=1}^n -2 \Sigma^{-1} (x_k - \hat{\mu}) = 0$$

↓ ↓
summation covariance
matrix

$$\Rightarrow \sum_{k=1}^n (x_k - \bar{\mu}) = 0$$

$$\bar{\mu} \sum_{k=1}^n 1 = \sum_{k=1}^n x_k$$

$$\boxed{\bar{\mu} = \frac{1}{n} \sum_{k=1}^n x_k}$$

case 2: μ and σ are unknown [1-D]

$$\theta_1 = \mu \quad \theta_2 = \sigma^2$$

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

$$\ln p(x_k | \theta) = \frac{-1}{2} \ln 2\pi \theta_2$$

$$= \frac{-1}{2} \ln 2\pi \theta_2 - \frac{1}{2} \frac{(x_k - \mu)^2}{\sigma^2}$$

$$= \frac{-1}{2} \ln 2\pi \theta_2 - \frac{1}{2} \frac{\theta_1 (x_k - \theta_1)^2}{\theta_2}$$

$$\nabla_{\theta} \ln(p(x_k | \theta)) = \begin{bmatrix} \frac{x_k - \theta_1}{\theta_2} \\ \frac{-1}{2\theta_2} + \frac{1}{2} \frac{(x_k - \theta_1)^2}{\theta_2^2} \end{bmatrix}$$

$$= \frac{1}{2} \ln \theta_2 - \frac{1}{2} \ln \theta_2$$

$$= \frac{1}{2} \times \frac{1}{\theta_2}$$

$$= \frac{1}{2} x^{-1}$$

$$= -1 x^{-2}$$

at $\hat{\theta}$

$$\textcircled{1} \quad \sum_{k=1}^n \frac{1}{\hat{\theta}_2} (x_k - \hat{\theta}_1) = 0$$

$$[\hat{\theta}_2 \neq 0]$$

$$\boxed{\hat{\theta}_1 = \frac{1}{n} \sum_{k=1}^n x_k = \bar{\mu}}$$

and

$$\textcircled{2} \quad \sum_{k=1}^n \frac{-1}{2\hat{\theta}_2} + \frac{1}{2\hat{\theta}_2^2} \sum_{k=1}^n (x_k - \hat{\theta}_1)^2 = 0$$

$$-n\hat{\theta}_2 + \sum_{k=1}^n (x_k - \hat{\theta}_1)^2 = 0$$

$$\hat{\theta}_2 = \frac{1}{n} \sum_{k=1}^n (x_k - \hat{\theta}_1)^2$$

$$\boxed{\sigma^2 = \frac{1}{n} \sum_{k=1}^n (x_k - \hat{\theta}_1)^2}$$

Multivariate MLE for $N(\mu, \Sigma)$

$$\hat{\mu} = \frac{1}{n} \sum_{k=1}^n x_k$$

$$\hat{\Sigma} = \frac{1}{n} \sum_{k=1}^n (x_k - \hat{\mu})(x_k - \hat{\mu})^t$$

unbiased

$$\boxed{\mathbb{E} \left[\frac{1}{n} \sum_{k=1}^n x_k \right] = \mu}$$

$$\mathbb{E} \left[\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 \right] = \sigma^2$$

$$= \frac{(n-1)}{n} \sigma^2$$

biased estimate

$$\sum_{n=1}^N \sum_{i=1}^n (x_i - \bar{x})^2 = \sigma^2$$

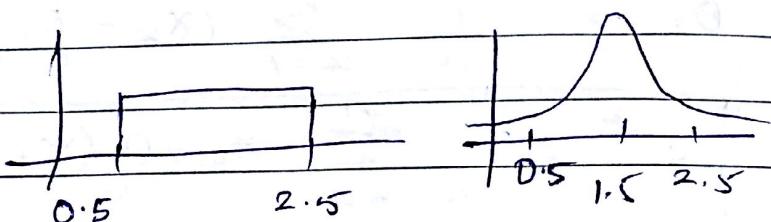
↑
unbiased MLE
(most)

10 October 2014

MLE

1. $p(x|\theta)$ is known, θ is unknown {MAP}

2. $p(\theta)$! all prior knowledge about θ
(all prior knowledge about θ
is encoded into a density func.)



$$P(\omega_i | x, D) = \frac{p(x|\omega_i, D) p(\omega_i|D)}{\sum_{i=1}^n (\dots)}$$

\downarrow

$p(x|D)$ → learn this part

$$\begin{aligned} p(x|D) &= \int p(x, \theta | D) d\theta \\ &= \int p(x|\theta, D) p(\theta|D) d\theta \\ &= \underbrace{\int p(x|\theta) p(\theta|D) d\theta}_{\text{known}} \end{aligned}$$

$$p(x|D) = \int p(x|\theta) p(\theta|D) d\theta$$

for ME :

$p(\theta|D)$ is zero at all θ
except $\hat{\theta}$
so,

$$p(x|D) = p(x|\hat{\theta})$$

(all samples $\in w_i$)
 $= p(x|w_i)$

do it for each class

$\rightarrow -$

prior

$$p(\theta|D) = \frac{p(D|\theta) p(\theta)}{p(D)}$$

$$= \frac{p(D|\theta) p(\theta)}{p(\theta)}$$

$$\mathcal{L} = \int p(\theta|D) p(\theta) d\theta$$

case 1 $\theta = \mu$ $\sigma = \text{known}$

prior of μ is modelled as
 $p(\mu) \sim N(\mu_0, \sigma_0^2)$

$$p(\theta|\mu) p(\mu) = \left[\prod_{k=1}^n p(x_k|\mu) \right] p(\mu)$$

indep identically

$$= \left[\prod_{k=1}^n \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2} \left(\frac{x_k-\mu}{\sigma} \right)^2} \right] \cdot \frac{1}{\sqrt{2\pi}\sigma_0} e^{-\frac{1}{2} \left(\frac{\mu-\mu_0}{\sigma_0} \right)^2}$$

$$\begin{aligned}
 & \text{d' } e \\
 & \rightarrow q_1 \left[\sum_{k=1}^n \left(\frac{x_k - \mu}{\sigma} \right)^2 \right] + \left(\frac{\mu - \mu_0}{\sigma_0} \right)^2 \\
 & \text{d'' } e \\
 & \rightarrow q_2 \left[\left(\frac{n}{\sigma^2} + \frac{1}{\sigma_0^2} \right) \mu^2 - 2 \left(\frac{\mu_0}{\sigma_0^2} + \frac{1}{\sigma^2} \sum_{k=1}^n x_k \right) \mu \right] \\
 & \text{d'''} e \\
 & = p(\mu | D) \\
 & \sim N(\mu_n, \sigma_n)
 \end{aligned}$$

when prior was normal density,
 $p(\mu_0)$ is also a normal density.
The shape of our belief remains
the same.

$$\times \left[\mu_n = \frac{n}{\sigma^2} + \frac{1}{\sigma_0^2} \right]$$

$$(3) p(x|D) = \int p(x|\theta) p(\theta|D) d\theta$$

$$\mu_n = \frac{n \sigma_0^2}{n \sigma_0^2 + \sigma^2} \bar{x}_n + \frac{\sigma^2}{n \sigma_0^2 + \sigma^2} \mu_0$$

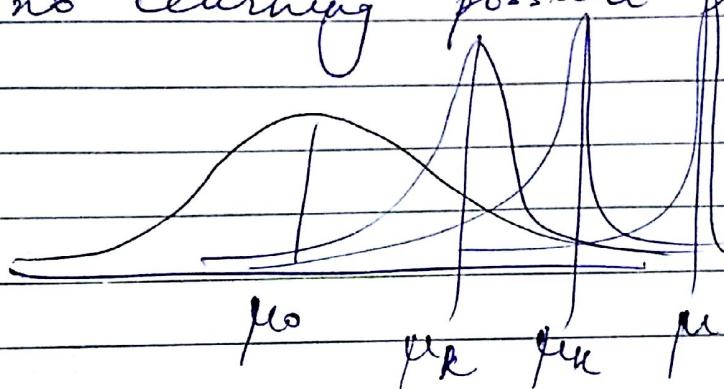
$$\sigma_n^2 = \frac{\sigma_0^2 \sigma^2}{n \sigma_0^2 + \sigma^2} = \frac{1}{\left(\frac{n}{\sigma^2} + \frac{1}{\sigma_0^2} \right)}$$

mean of the samples

$n = \text{large}$ (prior not imp.)

$\sigma^2 = \text{large}$ (unable to learn)

$\sigma_0 = 0$: fundamentalist
(believe in something ~~that~~ it's
completely certain)
(no learning possible)



Bayes & Bayesian Learning

- prior - Normal: Conjugate Prior
- posterior - Normal
- reproducing density