

Chapter 2 PDF Estimation

Parametric Estimation [Maximum Likelihood Estimation
Bayesian Estimation
Non-Parametric Estimation [KNN Estimation
Parzen Window Estimation
Expectation-Maximization Algorithm (EM)

Probability Lingo

$P(A)$ = probability of event a

$P(A|B)$ = probability of event a given event b occurred

$P(A, B|C)$ = probability of event a and b given event c occurred

$E\{f(y)\} = \int f(x)p(y)dy$

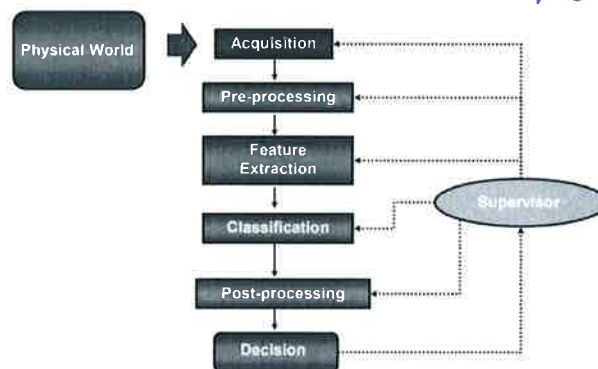
$$P(X, Y|Z) = P(Y|X, Z)P(X|Z) \text{ or } P(X|Y, Z)P(Y|Z)$$

The probability that X and Y happen if we know that Z happens is the same as the probability that X happens when we know that Z happens and that then Y happens when we know that X and Z happen.

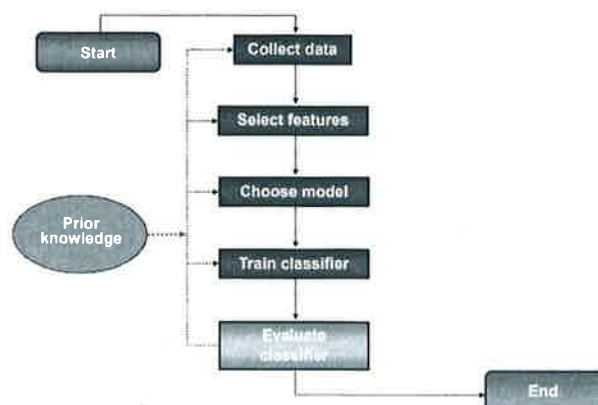
$$p(x|\theta) = p(x|m, \Sigma) = \frac{1}{\sqrt{(2\pi)^n |\Sigma|}} \exp\left[-\frac{1}{2\Sigma}(x-m)^t(x-m)\right]$$

When the covariance matrix has only the variance in the diagonal means there are no covariance between features, i.e. they are i.i.d.

ML Procedures.



Preprocessing = filtering / histogram manipulation / Segmentation - Enhancing the quality of the data Transforms can be part of pre-processing or feat. extraction



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

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

Parametric Estimation - we assume we know the model (gaussian, poisson, ...)

$$x \sim N(\mu, \sigma^2) \mid X = \{x_1, x_2, x_3 \dots x_N\} \mid \theta = [\mu, \sigma^2] \mid p(x|\theta) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

We consider that our $p(X|\theta) = p(x_1, x_2, \dots|\theta)$ and that it follows a normal distribution with variance equal to 1 and we want to discover our mean. So graphically we would have an infinite amount of normal curves for all the infinite values of μ . We get the following equation for the likelihood function:

$$\prod_{k=1}^N p(x_k|\theta) = \left(\frac{1}{\sqrt{2\pi}}\right)^N e^{-\frac{1}{2}((x_k^2 - 2x_k\mu + \mu^2) + (x_{k+1}^2 - \dots))} \quad (1)$$

$$\Rightarrow \approx \left(\frac{1}{\sqrt{2\pi}}\right)^N e^{-\frac{(\mu - \mu_0)^2}{\alpha}} \quad (2)$$

- Estimation goodness - The perfect estimation has no bias and no variance \Rightarrow but if it's good asymptotically

The bias is the shift of the estimated mean

For the variance we use:

* Cramer-Rao bound

Compute the sharpness of the estimation, this is done with the Fisher information, and then it's compared to the variance of the estimator:

$$I(\theta)_{i,j} = E\left\{\frac{\partial}{\partial \theta_i} [\ln(p(X|\theta))] \cdot \frac{\partial}{\partial \theta_j} [\ln(p(X|\theta))]\right\} \quad (3)$$

$$I(\theta)_{ij} = E\left\{\frac{\partial \ln p(X|\theta)}{\partial \theta_i} \cdot \frac{\partial \ln p(X|\theta)}{\partial \theta_j}\right\} \geq \frac{1}{I(\theta)} \quad (4)$$

2.1

* Asymptotic properties

Consistent estimator = asymptotically efficient $\frac{\text{var}(\hat{\theta})}{I^{-1}(\theta)}$ & unbiased $E\{\hat{\theta}\} = \theta$

$N \rightarrow \infty$ means large sets of observations
 μ_0 is the mean of all training samples.

$$\lim_{N \rightarrow \infty} \frac{\text{var}(\hat{\theta})}{I^{-1}(\theta)} = 1$$

$$\lim_{N \rightarrow \infty} E\{\hat{\theta}\} = \theta$$

- **Maximum Likelihood** - maximize the likelihood function

Properties \Rightarrow

- ① asymptotically unbiased
- ② asymptotically efficient
- ③ consistent

③ consistent

$$\lim_{N \rightarrow \infty} P\{\|\hat{\theta} - \theta\| < \delta\} = 1 \quad \forall \delta > 0$$

$$\hat{\theta} = \arg\max_{\theta} p(X|\theta) \Rightarrow \text{log-likelihood} \quad \hat{\theta} = \arg\max_{\theta} \ln p(X|\theta)$$

mono-dim Example here for distribution

variance σ known
with observations $\Rightarrow \mu$?

1. Calculate the likelihood function:

$$p(X|\theta) = \prod_{i=1}^N p(x_i|\theta) \quad \ln \prod_{i=1}^N \frac{1}{\sigma \sqrt{2\pi}} \exp\left(-\frac{(x_i - \mu)^2}{2\sigma^2}\right) \quad (5) \quad (\text{prove back})$$

$$\propto \ln\left(\prod_{i=1}^N p(x_i|\theta)\right) = \sum_{i=1}^N \ln(p(x_i|\theta)) \quad (6)$$

$$= \sum_{i=1}^N \ln\left((2\pi)^{-1/2} \sigma^{-1}\right) - \frac{1}{2} \frac{(x_i - \mu)^2}{\sigma^2} \quad (7)$$

$$L(\theta) = \sum_{i=1}^N -\frac{1}{2} \ln(2\pi) - \ln(\sigma) - \frac{1}{2} \frac{(x_i - \mu)^2}{\sigma^2} \quad (8)$$

2. Maximize by derivating by μ and equalising to zero:

$$\frac{\partial L(\theta)}{\partial \mu} = \sum_{i=1}^N -\frac{1}{2} \frac{2(x_i - \mu)(-1)}{\sigma^2} = 0 \quad (9)$$

$$\sum_{i=1}^N (x_i - \mu) = 0 \quad (10)$$

$$\sum_{i=1}^N x_i - \mu N = 0 \quad (11)$$

$$\mu = \frac{\sum_{i=1}^N x_i}{N} \quad (12)$$

3. Maximize by derivating by σ and equalising to zero:

$$\frac{\partial L(\theta)}{\partial \sigma} = \sum_{i=1}^N \left(-\frac{1}{\sigma} - \frac{1}{2} (x_i - \mu)^2 (-2)(\sigma^{-3})\right) = 0 \quad (13)$$

$$\frac{1}{\sigma} \sum_{i=1}^N \left(\frac{(x_i - \mu)^2}{\sigma^2} - 1\right) = 0 \quad (14)$$

$$\sum_{i=1}^N \frac{(x_i - \mu)^2}{\sigma^2} = N \quad (15)$$

$$\frac{1}{\sigma^2} \sum_{i=1}^N (x_i - \mu)^2 = N \quad (16)$$

$$\text{mean vector } \sigma^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \mu)^2 \quad (17)$$

covariance matrix

For multivariate gauss. it's replacing μ with \hat{m} and σ^2 for \sum and for the derivative of σ it's $(x_i - \hat{m})(x_i - \hat{m})^t$

Tricky point $\mu \rightarrow$ mean of the observations (estimation)
 $\sigma \rightarrow$ variance of the observations (estimation)

In ML Estimation θ fixed

Bayesian θ to be random variable (may have prior density for distribution of samples θ)

$$p(x|X) = \int_0^1 p(x, \theta|X) d\theta = \int_0^1 p(x|\theta) p(\theta|X) d\theta$$

Bayesian estimation - Weighted average of all possible models
Given observation we have our **PRIOR DENSITY** $p(\theta)$ with a $\max(p(\theta)) = M_1$
Running the X_N samples we get our **POSTERIOR DENSITY** $p(\theta|X)$ with a $\max(p(\theta|X)) = M_2$

Using Bayesian Formula

We want to find the estimation of the real distribution: $\hat{p}(x|X) = \int_0^1 p(x, \theta|X) d\theta = \int_0^1 p(x|\theta) p(\theta|X) d\theta$
 X doesn't add any value on the first prob

So $\int_0^1 p(x|\theta) p(\theta|X) d\theta$ which is basically an infinite weighted average of the probability of x given theta weighted by the probability that a given distribution from the samples (?) has theta

Bayes formula: $p(\theta|X) = \frac{p(X|\theta)p(\theta)}{p(X)} = \frac{p(X|\theta)p(\theta)}{\int p(X|\theta)p(\theta)d\theta}$

From here it's considered that $\theta = [\mu, \sigma^2]$ where only μ is unknown so we basically replace θ for μ and also that $p(\mu) = N(\mu_0, \sigma_0^2)$ prior knowledge for μ

simplified version of calculation

$p(\mu|X) = \frac{p(X|\mu)p(\mu)}{\int p(X|\mu)p(\mu)d\mu}$
Normalization Factor

$$p(\mu|X) = \frac{p(X|\mu)p(\mu)}{p(X)} = \frac{(\prod_{i=1}^N p(x_i|\mu))p(\mu)}{p(X)} \quad (18)$$

$$p(X) = \int \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(x_i - \mu)^2} d\mu \quad (19)$$

$p(x_i|\mu) \sim N(\mu, \sigma^2)$
 $p(\mu) \sim N(\mu_0, \sigma_0^2)$
minus $2(\frac{1}{\sigma^2} \sum x_i + \frac{\mu_0}{\sigma_0^2})\mu$

$$p(\mu|X) = \frac{1}{p(X)} \left(\frac{1}{\sqrt{2\pi\sigma^2}} \right)^N \frac{1}{\sqrt{2\pi\sigma_0^2}} e^{-\frac{1}{2\sigma^2} \sum_{i=1}^N (x_i - \mu)^2 - \frac{1}{2\sigma_0^2} (\mu - \mu_0)^2} \quad (20)$$

However, Remember the answer back (next page)

$$p(\mu|X) = \frac{1}{p(X)} \exp \left\{ -\frac{1}{2} \left[\left(\frac{1}{\sigma^2} + \frac{1}{\sigma_0^2} \right) \mu^2 - \frac{1}{\sigma^2} \sum_{i=1}^N x_i \mu + \frac{1}{\sigma_0^2} \mu \mu_0 \right] \right\} \quad (21)$$

$$= -\frac{1}{2} \left(\frac{\sum_{i=1}^N x_i^2}{\sigma^2} - \frac{\sum_{i=1}^N 2x_i \mu}{\sigma^2} + \frac{N\mu^2}{\sigma^2} + \left(\frac{\mu^2}{\sigma_0^2} - \frac{2\mu\mu_0}{\sigma_0^2} + \frac{\mu_0^2}{\sigma_0^2} \right) \right) \quad (22)$$

$$= -\frac{1}{2} \left(\mu^2 \left(\frac{N}{\sigma^2} + \frac{1}{\sigma_0^2} \right) - 2\mu \left(\frac{\sum_{i=1}^N x_i}{\sigma^2} + \frac{\mu_0}{\sigma_0^2} \right) + \frac{\sum_{i=1}^N x_i^2}{\sigma^2} + \frac{\mu_0^2}{\sigma_0^2} \right) \quad (23)$$

Since $p(\mu|X)$ is a normal distribution $N - (\mu_N, \sigma_N^2)$, we can realte its exponential parts:

$$\frac{1}{\sqrt{2\pi\sigma_N^2}} e^{-\frac{1}{2\sigma_N^2}(\mu - \mu_N)^2} = \frac{1}{\sqrt{2\pi\sigma_N^2}} e^{-\frac{1}{2} \left(\frac{\mu^2}{\sigma_N^2} - \frac{2\mu\mu_N}{\sigma_N^2} + \frac{\mu_N^2}{\sigma_N^2} \right)} \quad (24)$$

where:

$$\mu^2 \left(\frac{N}{\sigma^2} + \frac{1}{\sigma_0^2} \right) = \frac{\mu^2}{\sigma_N^2} \rightarrow \sigma_N^2 = \frac{\sigma_0^2 \sigma^2}{N\sigma_0^2 + \sigma^2} \quad (25)$$

and

$$-2\mu \left(\frac{\sum_{i=1}^N x_i}{\sigma^2} + \frac{\mu_0}{\sigma_0^2} \right) = -\frac{2\mu\mu_N}{\sigma_N^2} \quad (26)$$

$$\mu_N = \frac{\sigma_N^2 \sum x_i}{\sigma^2} + \frac{\sigma_N^2 \mu_0}{\sigma_0^2} = \frac{\sigma_0^2 \sigma^2}{N\sigma_0^2 + \sigma^2} \frac{\sum x_i}{\sigma^2} + \frac{\sigma_0^2 \sigma^2}{N\sigma_0^2 + \sigma^2} \frac{\mu_0}{\sigma_0^2} \quad (27)$$

$$= \frac{\sigma_0^2 \sum x_i}{N\sigma_0^2 + \sigma^2} + \frac{\sigma^2 \mu_0}{N\sigma_0^2 + \sigma^2} \mid \frac{\sum x_i}{N} = \bar{X} \quad (28)$$

$$= \frac{\sigma_0^2 N \bar{X}}{N\sigma_0^2 + \sigma^2} + \frac{\sigma^2 \mu_0}{N\sigma_0^2 + \sigma^2} \quad (29)$$

From here, we go back to the start: $\hat{p}(x|X) = \int_{\theta} p(x, \theta|X) d\theta = \int_{\theta} [p(x|\theta) - N(\mu, \sigma^2)] [p(\theta|X) - N(\mu_N, \sigma_N^2)] d\theta$ so we get that $\hat{p}(x|X) = N(\mu_N, \sigma_N^2)$

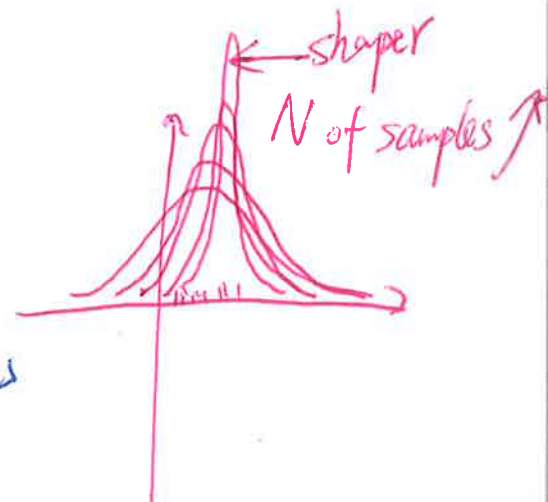
ML vs. Bayesian

ML is a consistent estimator, given high number of training samples and it's fairly simple to implement. Bayesian estimation is useful if we have prior information, which would be discarded in ML. (normally used with low number of samples)

Answer $\left\{ \begin{aligned} \mu_N &= \left(\frac{N\bar{x}}{N\bar{\sigma}_0^2 + \sigma^2} \right) \bar{X} + \left(\frac{\bar{\sigma}_0^2}{N\bar{\sigma}_0^2 + \sigma^2} \right) \mu_0 \\ \sigma_N^2 &= \frac{\bar{\sigma}_0^2 \sigma^2}{N\bar{\sigma}_0^2 + \sigma^2} \end{aligned} \right.$ properties

$\lim_{N \rightarrow \infty} \mu_N = \bar{X}$ $\sigma_N^2 = 0$
means distribution obey the distribution of samples

$\lim_{N \rightarrow 0} \mu_N = \mu_0$ $\sigma_N = \sigma_0$
when there is no posterior knowledge use the prior knowledge.



ML > Bayesian in practical
Lower computational complexity
Easier interpretability

chosen of a function of N i.e. $K_N = \sqrt{N}$

$\hat{p}(x^*) = \frac{K}{NV_K(x^*)}$ $\lim_{N \rightarrow \infty} K_N = \infty$ $\lim_{N \rightarrow \infty} \frac{K_N}{N} = 0$

$KNN \rightarrow K$ fixed $\rightarrow V$ fixed.

Parzen Window $\rightarrow V$ fixed.

Non-Parametric Estimation

We have no knowledge of the pdf.

Considering x^* is a ~~random~~ ^{generic} sample that belongs to the region R , and R has an infinitesimally small volume V :

$$P_R = \int_R p(x^*) dx^* \rightarrow p(x^*) V \quad \text{Volume of } R \quad (30)$$

The estimation simply counts the number of samples from the training set that are inside this volume, so we can write this as:

Law of Large Numbers (大数定律)

$\lim_{N \rightarrow \infty} P\{|\hat{P}_R - P_R| < \epsilon\} = 1$ $\hat{P}_R = \frac{K}{N} = p(x^*) V$ $p(x^*) V = \frac{K}{N} \rightarrow p(x^*) = \frac{K}{NV}$ $\frac{K}{N} = p(x^*) V$ $p(x^*) = \frac{K}{NV}$ (31)

Where N is the total number of samples and K the number of samples inside the volume

• K-Nearest Neighbour

Fixes K to a constant, and expands the volume until a number of K samples are inside

• Parzen windows

Fixes V to a constant, and counts the number of K samples inside.

For this we need a window function (can have different shapes):
or called kernel

$$\gamma(x) = \{1; 0\} \rightarrow \gamma\left[\frac{x_k - x^*}{h}\right] \quad (32)$$

$$p(x^*) = \frac{\sum_{k=1}^N \gamma\left[\frac{x_k - x^*}{h}\right]}{N(h^n)}$$

K is varying

h is the determined value

Where $V = h^n$ since it has sides length h in n dimensions

Developing the equation of the Bias, i.e. $E\{\hat{p}(x)\} = p(x)$:

$$\frac{1}{NV(h)} \sum_{i=1}^N E\left\{\gamma\left[\frac{x - x_i}{h}\right]\right\} \quad (34)$$

$$\frac{N}{NV(h)} E\left\{\gamma\left[\frac{x - y}{h}\right]\right\} = \frac{1}{V(h)} \int \gamma\left[\frac{x - y}{h}\right] p(y) dy = \frac{1}{V(h)} p(x) * \gamma\left[\frac{x}{h}\right] \quad (35)$$

Since the expectation of the window is the same regardless of the sample, we use y and remove the sum.

$$\hat{p}(x) = \frac{1}{N} \sum_{k=1}^N \frac{1}{V(h)} \left(\frac{x - x_k}{h} \right) \quad \boxed{y(x) \geq 0 \mid \forall x \in \mathbb{R}^n \int_{\mathbb{R}^n} y(x) d(x) = 1}$$

For the variance it's used an upper bound, and it's unbiased when

$$\text{var}(\hat{x}) \leq \frac{\dots}{Nh^n} \quad (36)$$

$$\lim_{N \rightarrow +\infty} h_N = 0 \quad (37)$$

$$\lim_{N \rightarrow +\infty} Nh_N^n = +\infty \rightarrow \text{var}(\hat{x}) \leq 0 \quad (38)$$

This means that the window acts as a filter, the smaller h the closer it is to an impulse and the lower the bias, thus why h^n should tend to 0

$$\text{Ex: } h_N = \frac{1}{\sqrt[n]{N}}$$

K and h increase = underfitting/oversmoothing/blurring

K and h decrease = overfitting

window size \uparrow blurring \uparrow approximation
window size \downarrow overfitting \downarrow performance



Estimation with incomplete data

Considering

$$p(x) = \sum_{i=1}^M P_i p(x|m_i, \Sigma_i) \quad (39)$$

as a mixture of gaussian func. where P_i is the prior probability, m_i the mean vector and Σ_i the covariance matrix.

Objective is to estimate the distribution of each gaussian. We consider $Z = (X, Y)$ where X are the samples and Y is the missing information

- **Expectation maximization** - more elaborate technique than ML

Expectation step (k iteration index):

$$E\{\ln(p(X, Y|\theta)) | X, \theta^k\} = \int \ln[p(X, Y|\theta)] p(Y|X, \theta^k) dY = Q(\theta, \theta^k) \quad (40)$$

Read as assuming X samples we have and theta at iteration k

This is "expectation" step and now the maximization step:

$$\theta^{k+1} = \arg \max_{\theta} Q(\theta, \theta^k) \quad (42)$$

Ex: To estimate θ we should compute the maximum likelihood on Z and to "get" Y we average ($E\{\}$) on what we have (X)

Considering $\theta = [m, \Sigma]$; $m = [\mu_1, \mu_2]$ and $\Sigma = [[\sigma_1^2, 0], [0, \sigma_2^2]]$

$$Z = \{[0, 2], [1, 0], [2, 2], [x_{4,1}, 4]\} \mid x_{4,1} = Y \quad (43)$$

$$Q(\theta, \theta^0) = \int \ln(p(x_1, x_2, x_3, x_4|\theta)) p(x_{4,1}|X, \theta^0) dx_{4,1} \quad (44)$$

$$\int \ln\left[\prod_{i=1}^3 p(x_i|\theta) p(x_4|\theta)\right] p(x_{4,1}|\theta^0) dx_{4,1} \quad (45)$$

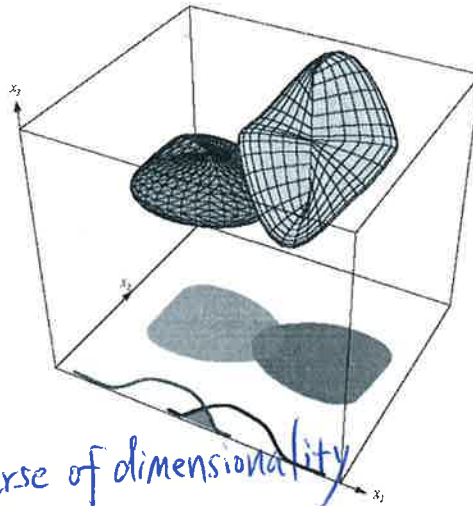
$$\sum_{i=1}^3 \left[\int \ln(p(x_i|\theta)) p(x_{4,1}|\theta^0) dx_{4,1} \right] + \int \ln(p(x_4|\theta)) p(x_{4,1}|\theta^0) dx_{4,1} \quad (46)$$

$$\sum_{i=1}^3 \ln(p(x_i|\theta)) + \int \ln(p(x_{4,1}, x_{4,2}|\theta)) p(x_{4,1}|\theta^0) \quad (47)$$

Feature reduction Chapter3

Adding features can help separate data:

There is a non-zero Bayes error in the 1-D x_1 space or the 2-D x_1, x_2 space. However, the Bayes error vanishes in the 3-D x_1, x_2, x_3 space because of non-overlapping densities.



Hughes effect: adding more features to separate data becomes a disadvantage after a specific point Why? *classification accuracy affected*

- Adding features changes the distribution type of data, i.e. no longer gaussian
- Overfitting

So we reduce features to not get this problems and reduce computational complexity

The reason behind this is that using more dimensions requires more training data, since in each dimension increases volume. Consider a circle inside a square:

- 1D the circle and the square are a line of length $= 2R = L$
- 2D the circle area is now $A_c = \pi R^2$ and the square is L^2 which gives a difference of $(1 - \frac{\pi}{4})$
- 3D the sphere inside the cube gives a difference of

**Feature selection**

Consider $F = \{X_1, X_2, X_3 \dots X_n\}$ with n Features and f_k and $f_{\bar{k}}$ being k selected elements and \bar{k} eliminated ones thus $k + \bar{k} = n$

We want to pick the features that SEPARATE the best possible the classes so we want F^* :

$$F^* = \operatorname{argmax}_{F'} \{J(F')\} \quad \text{opportunity function} \quad (48)$$

F' is the subset of F with the selected features

Types of $J()$ s**Divergence:**

We use the likelihood ratio: $\frac{p(x|w_1)}{p(x|w_2)}$

If the distributions are overlapped the ratio will be $=1$, else if it's very big or small it means there is low overlap

But this is not symmetric ($1/2 \neq 2/1$), so we use 2 ratios:

$$D_{ij}(F') = E\{L'_{ij}(x)\} + E\{L'_{ji}(x)\} = E\{\ln p(x|w_i) - \ln p(x|w_j)\} + E\{\ln p(x|w_j) - \ln p(x|w_i)\}$$

$$D_{1,2} = E\{\ln \frac{p(x|w_1)}{p(x|w_2)} | w_1\} + E\{\ln \frac{p(x|w_2)}{p(x|w_1)} | w_2\} \quad (49)$$

$$\Rightarrow \int \left\{ \frac{p(x|w_1)}{p(x|w_2)} \ln \frac{p(x|w_1)}{p(x|w_2)} p(x|w_2) dx + \int \frac{p(x|w_2)}{p(x|w_1)} \ln \frac{p(x|w_2)}{p(x|w_1)} p(x|w_1) dx \right\} \quad (50)$$

$$\Rightarrow \int \left\{ \frac{p(x|w_1)}{p(x|w_2)} \ln \frac{p(x|w_1)}{p(x|w_2)} p(x|w_2) dx + \int \frac{p(x|w_2)}{p(x|w_1)} \ln \frac{p(x|w_2)}{p(x|w_1)} p(x|w_1) dx \right\} \quad (51)$$

Bhattacharyya:

Minimize the bayess error using the area of overlap between two distributions with a crossing point in T :

$$p(\text{err}|w_1)P(w_1) + p(\text{err}|w_2)P(w_2) \quad (52)$$

$$p(\text{err}|w_2) = \int_T^{+\infty} p(x|w_1) dx \quad (53)$$

$$p(\text{err}|w_1) = \int_{-\infty}^T p(x|w_2) dx \quad (54)$$

$$\int_{-\infty}^{+\infty} \min\{p(x|w_1)P(w_1), p(x|w_2)P(w_2)\} dx \quad (55)$$

To solve this we would need the real densities and probabilities, and we only have the estimates, and solving the integral is complex so we use the upper bound ε_u i.e. the Chernoff bound. The smallest the value the better separability

Find an upper bound $\min\{a, b\} \Rightarrow \text{inequality } \min\{a, b\} \leq a^s b^{1-s}, 0 \leq s \leq 1$
Chernoff bound

$$\min\{a, b\} \leq a^s b^{1-s} ; 0 \leq s \leq 1 \quad (56)$$

$$\int [p(x|w_1)P(w_1)]^s [p(x|w_2)P(w_2)]^{1-s} dx \quad \text{constant} \quad (57)$$

$$= P(w_1)^s P(w_2)^{1-s} \int p(x|w_1)^s p(x|w_2)^{1-s} dx \quad (58)$$

$$= P(w_1)^s P(w_2)^{1-s} \exp[-\mu_{1,2}(s)] \quad (59)$$

Where for a gaussian distribution the chernoff distance $\mu_{1,2}(s) =$

$$\frac{s(1-s)}{2} (m_1 - m_2)^2 \left\{ s \sum_1 + (1-s) \sum_j \right\}^{-1} (m_1 - m_2) + \frac{1}{2} \ln \left(\frac{|s \sum_1 + (1-s) \sum_j|}{|\sum_1|^s |\sum_2|^{1-s}} \right) \quad (60)$$

A specific case is where $s = 0.5$ where we get the battacharyya bound, and for gaussians the equation above with 0.5 is called battacharyya distance $B_{1,2}$ and it has the following properties:

$$B_{1,2} = \mu_{1,2}(0.5) \quad (61)$$

$$w_1 = w_2 \rightarrow B_{1,2} = 0 \quad (62)$$

$$w_1 \neq w_2 \rightarrow B_{1,2} > 0 \quad (63)$$

$$B_{1,2} = B_{2,1} \quad (64)$$

And it has no saturating behaviour just like divergence

- Jeffries-Matusita: Basically find one that saturates, so it's the MSE of the densities:

JM is an ave. distance between two density

$$JM = \left[\int \left(\sqrt{p(x|w_1)} - \sqrt{p(x|w_2)} \right)^2 dx \right]^{\frac{1}{2}}$$

$$JM_{1,2}^2 = \int (\sqrt{p(x|w_1)} - \sqrt{p(x|w_2)})^2 dx \quad (65)$$

$$\int p(x|w_1) dx - 2 \int \sqrt{p(x|w_1)p(x|w_2)} dx + \int p(x|w_2) dx \quad (66)$$

$$1 + 1 - 2 \int \sqrt{p(x|w_1)p(x|w_2)} dx \quad (67)$$

This is written as a function of the Battacharyya distance (because):

$$\int \sqrt{p(x|w_1)p(x|w_2)} dx = e^{-B_{1,2}} \quad (68)$$

$$JM_{1,2} = \sqrt{2(1 - e^{-B_{1,2}})} \quad (69)$$

JM saturates at $\sqrt{2}$

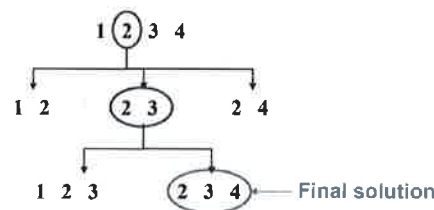
Search Strategies

Used to find the best subset of features that optimize the criterion

- SFS - Sequential forward selection

1. Pick the best mean distance in a 1D/1 feature plain, meaning along each feature calculate the distance between each class and calculate the mean. Then compare all the means from all the features as pick the highest one.
2. Pick the best mean distance using the previously stored one, now in 2D/2 features, meaning the selected feature is now the one of the axis and the other is the other feature you're testing
3. keep repeating

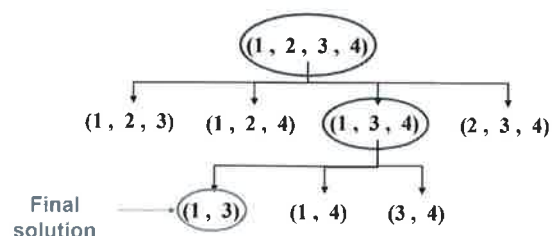
● Example: $m=3$; $n=4$



- SBS

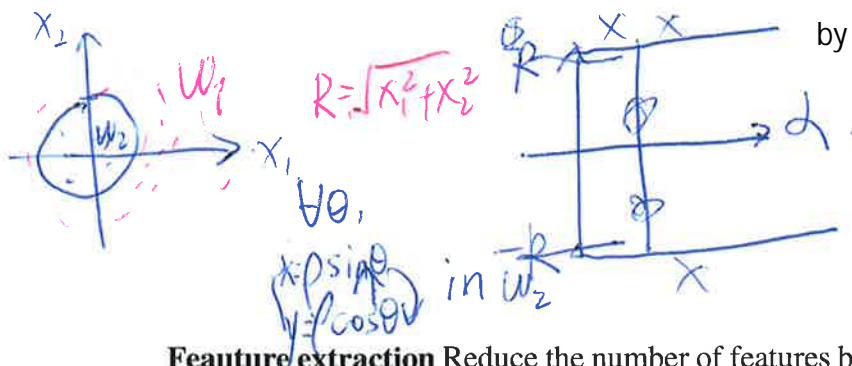
Pick a group of features and progressively remove one till you find the best one

● Example: $m=2$; $n=4$



Consider n features and m subsets of features, the best search depends on these numbers

SFS is better if you want to go up to $m = n/2$, else SBS is better, since it's faster



Feature extraction Reduce the number of features by using transformations, not eliminating any possible relevant data

- PCA - Principal Component Analysis

Lingo: *best represents the data in least square sense*

$$\det(\sum -\lambda I) = 0 \rightarrow \text{find } \underline{\text{eigenvalues}} \text{ aka char. equation} \quad (70)$$

特征值

Best represent the data in less features. It is unsupervised i.e. no labels are needed

Apply a transformation to the data:

① Reconstruct Samples

② Minimize the error of the distance *between the reconstructed data with real data*

$$x = \sum_{i=1}^n y_i \phi_i \quad (71)$$

x_i ← reconstructed

phi_i ← bases for the

y_i ← x_i subspace

$$x' = \sum_{i=1}^m y_i \Phi_i \quad |m < n$$

$$\text{error}(J_m) = \sum_{i=1}^N \|x_i - x'\|^2 \quad (72) \quad (73)$$

reconstruction

"I will project x along each ϕ_i , and each projection will give me each value y_i . These values together will provide me the new coordinates for x "

What are the ϕ_i s so that this is minimized? we pick the eigenvectors, i.e. the directions where the variance is the highest, i.e. where the entropy is highest

1. Calculate the baricenter (mean of all data) *重心(均值)*

2. Use this to move the data to the center of the space by $x' = x - m$ (Data Shifting)

3. Calculate the covariance matrix *(协方差矩阵)*

- ★ 4. Calculate the eigen values and with each of them calculate their eigen vector *(特征值/特征向量)*

5. Calculate $y = \Phi^t x$ ex: *(变换)*

$$\Phi = [\phi_1, \phi_2] \quad x = [x_1, x_2]$$

$$\phi_1 = [n_1, n_2] \quad \phi_2 = [n_3, n_4]$$

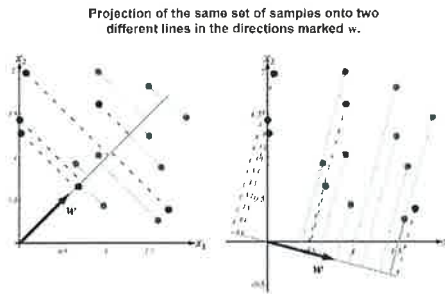
$$y = [y_1, y_2] \quad y_1 = x_1 n_1 + x_2 n_2 \quad y_2 = x_3 n_3 + x_4 n_4$$

In the case of reducing to one dimension we would only use ϕ_1 and thus the output will only be y_1

- LDA - Linear Discriminant Analysis

Best separates the data in less features

Goal: find a projection to a line w where the points of the classes are separated, where $y = W^t x$



To determine how separated the classes are we use the **Fisher's linear discriminant** and we want to MAXIMIZE IT WITH RESPECT TO w :

均値 $y \in W_i$ 方差 \rightarrow criterion function

$$\frac{\|m'_1 - m'_2\|^2}{\sigma_1'^2 + \sigma_2'^2} \rightarrow J = \frac{\|m'_1 - m'_2\|^2}{S'_1 + S'_2} \quad (74)$$

$$m' = W^t m \quad (75)$$

$$W^t S_b W = W^t (m_1 - m_2)(m_1 - m_2)^t W \quad (76)$$

$$(W^t m_1 - W^t m_2)(W^t (m_1 - m_2)) = (m'_1 - m'_2)^2 \quad (77)$$

$$\text{where } W^t m_1 = m'_1 \text{ \& } W^t m_2 = m'_2 \quad (78)$$

Scatter Matrix $S_i = \sum_{x \in X_i} (x - m_i)(x - m_i)^t$

Where $S_w = S_1 + S_2$ is the scatter matrix within class and $S_b = (m_1 - m_2)(m_1 - m_2)^t$ between classes.

The same done in eq 53-55 can be done for S_w :

$$W^t S_w W = W^t (S_1 + S_2) W = W^t S_1 W + W^t S_2 W \quad (79)$$

$$\text{ex for } S_1 \rightarrow W^t \left(\sum_{x \in \text{class}_1} (x - m_1)(x - m_1)^t \right) W \quad (80)$$

$$= \sum_{x \in \text{class}_1} (W^t x - W^t m_1)(W^t (x - m_1)) \quad (81)$$

$$= \sum_{x \in \text{class}_1} (W^t x - W^t m_1)^2 = \sum_{x \in \text{class}_1} (y - m'_1)^2 = S'_1 \quad (82)$$

and we can write $J(w)$:

$$J(w) = \frac{w^t S_b w}{w^t S_w w} \quad w = S_w^{-1} (m_1 - m_2) \quad (83)$$

And maximized (i.e. the w direction that maximizes the separability) is

$$w = S_w^{-1}(m_1 - m_2)$$

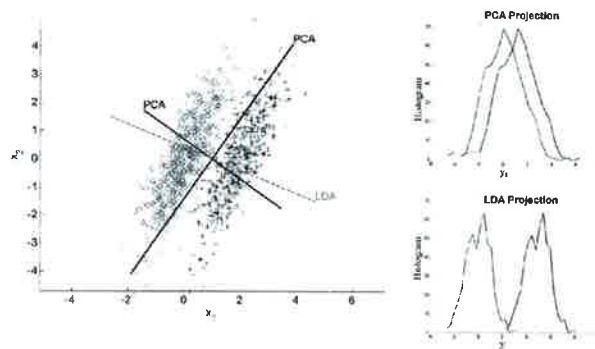
*For more classes/dimension it's $|J(w)| \rightarrow$ maximized are the eigen vectors of $w = S_w^{-1}S_B$

$$y = W^t x \mid W = [w_1, w_2, w_3, \dots, w_{c-1}] \quad (84)$$

$$S_b = n_1 S_{b1} + n_2 S_{b2} + \dots n_c S_{bc} \mid S_w = S_1 + S_2 + S_3 + \dots S_c \quad (85)$$

Where n 's are the number of samples per class, so S_b is a weighted avg.

$J(W)$ is the same formula only with W being the vector of w 's and to calculate it as an escalar we use the determinant of the parts of the numerator and the denominator



Multiclass Case

C classes $\Rightarrow C-1$ discriminant functions

$$S_w = \sum_{i=1}^C S_i$$

$$S_B = \sum_{i=1}^C [\#(X_i)] (m_i - m)(m_i - m)^t$$

$$m = \frac{1}{N} \sum_{x \in X} x$$

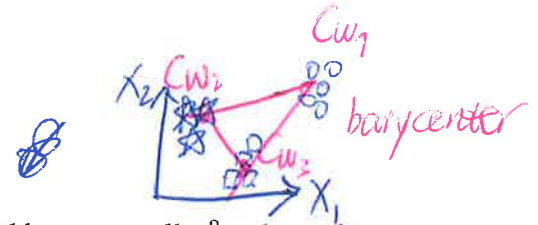
Global mean vector

Intro. design of a supervised classifier depends on

1. Available set of features
2. training samples
3. classification model
4. cost function to optimize

Supervised Classification

Bayesian Classification
Minimum Risk Theory
Discriminant Functions
Decision Trees
Accuracy Evaluation



• Parametric

– MDM - minimal distance to means Classes should have a small σ^2 or have the same shape

1. Calculate the mean of each class (aka baricenter aka class prototype)
2. Given a new point calculate the distance to each of the class prototypes
3. Assign it to the closest one

– Box

1. Calculate the baricenter of each class
2. Calculate the σ^2 of the data in the direction of each axis
3. Create boxes with size of $\alpha\sigma^2$
4. Given a new point assign it in the box it falls



*The boxes don't necessarily have equal length and width

*If a new point is outside of all boxes it can be assigned to the closest one or left unclassified

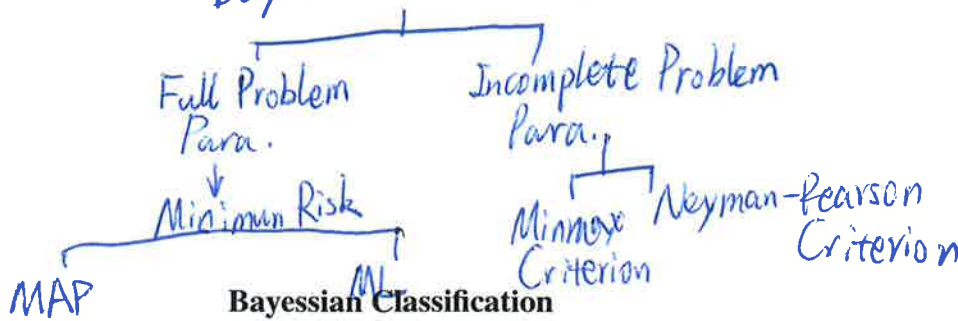
*Boxes can't overlap

– Bayesian Use all available data for classification

SEE NEXT PAGE

• Non Parametric

- Bayesian
- DT
- Neural Nets
- SVM



- Full problem parametrization

- Minimum risk (explained after)

- * MAP Criterion - Maximum a posteriori aka Minimum error crit.

- It's called both minimum and maximum because they can be seen both ways, in the case of the maximum a posteriori probability the equation comes from:

aka.

$$x \in W_i \Leftrightarrow P(w_i|x) \geq P(w_j|x) \quad \forall i=1,2,\dots,C$$

$$\hat{w} = \arg \max_{w_i} \{P(w_i|x)\} \quad (86)$$

Bayesian Formula

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

rewrite the posterior probability ($P(w_i|x)$) as the prior ($P(w_i)$)

$$P(w_i|x) = \frac{P(w_i)p(x|w_i)}{p(x)} \quad (87)$$

$$\hat{w} = \arg \max_{w_i} \{p(x|w_i)P(w_i)\} \quad (88)$$

$$\rightarrow p(x|w_i)P(w_i) \geq p(x|w_j)P(w_j) \quad (89)$$

ave. probability of error

$P(x \rightarrow w_2|w_1)$ denote error

$P_e(x|w_1)$

$$P_e = \sum_{i=1}^C P_e(x|w_i)P(w_i)$$

$$= \sum_{i=1}^C p(x|R_i|w_i)P(w_i)$$

1. Now with the other definition, the minimum error criterion starts by using the total probability theorem:

$$P_{err} = P(x \rightarrow w_2|w_1)P(w_1) + P(x \rightarrow w_1|w_2)P(w_2) \quad (90)$$

$$= \int_T^{+\infty} p(x|w_1)P(w_1)dx + \int_{-\infty}^T p(x|w_2)P(w_2)dx \quad (91)$$

$$= \int_T^{+\infty} P(w_1|x)p(x)dx + \int_{-\infty}^T P(w_2|x)p(x)dx \quad (92)$$

in eq 58-59 we change from posterior to prior

2. Defining the probability of error:

$$P_{err} = E\{P(err|x)\} = \int_{-\infty}^{+\infty} P(err|x)p(x)dx \quad (93)$$

$$= \int_{-\infty}^{+T} P(err|x)p(x)dx + \int_T^{+\infty} P(err|x)p(x)dx \quad (94)$$

Equalizing this equation to the one in step 1 we reach to the conclusion that:

3. The probability of error can only be given by $P(w_1|x)$ or $P(w_2|x)$, so since for a given point one will be higher than the other we minimize it:

$$P(err|x) = \min\{P(w_1|x), P(w_2|x)\} \quad (95)$$

4. final equation:

$$P_{errmin} = \int_{-\infty}^{+\infty} \min\{P(w_1|x), P(w_2|x)\}p(x)dx \quad (96)$$

* ML Criterion - Maximum likelihood criterion

1. Usign the same definition as MAP:

$$\hat{w} = \arg \max_{w_i} \{P(w_i|x)\} = \arg \max_{w_i} \{p(x|w_i)\} \quad (97)$$

$$\rightarrow p(x|w_i) \geq p(x|w_j) \quad (98)$$

but the difference is that we assume both distributions have the same prior probability $P(w_i)$

Minimum risk theory:

Depending on the application there can be a related cost to assigning incorrectly a new point of data. With this we can integrate this information:

$$A = \{\alpha_1, \alpha_2, \dots, \alpha_R\} \quad (99)$$

*Elements in the diagonal are correctly assigned and the data in the matrix should

no-call
call fireman

1000 euros

0

fire no fire

It is thus possible to define a $R \times C$ cost matrix Λ :

$$\Lambda = \begin{bmatrix} \lambda(\alpha_1|w_1) & \lambda(\alpha_1|w_2) & \dots & \lambda(\alpha_1|w_C) \\ \lambda(\alpha_2|w_1) & \lambda(\alpha_2|w_2) & \dots & \lambda(\alpha_2|w_C) \\ \vdots & \vdots & \ddots & \vdots \\ \lambda(\alpha_R|w_1) & \lambda(\alpha_R|w_2) & \dots & \lambda(\alpha_R|w_C) \end{bmatrix}$$

row 1 take action 1 given different situations

$\lambda_{ij} = \lambda(\alpha_i|w_j)$ ($\lambda_{ij} \geq 0$) is the cost (loss) incurred for taking action α_i given the class is w_j .

be "provided by an expert in the domain"

1. Calculate the conditional risk/expectation:

given observation $x \Rightarrow$ taking action

$$R(\alpha_i|x) = \sum_{j=1}^C \lambda(\alpha_i|w_j) P(w_j|x) \quad (100)$$

$$(101)$$

★ 2. Compare which one is lower (ex with two options):

$$R(\alpha_1|x) \leq R(\alpha_2|x) \quad (102)$$

$$\lambda_{1,1}P(w_1|x) + \lambda_{1,2}P(w_2|x) \leq \lambda_{2,1}P(w_1|x) + \lambda_{2,2}P(w_2|x) \quad (103)$$

$$\lambda_{1,1}p(x|w_1)P(w_1) + \lambda_{1,2}p(x|w_2)P(w_2) \leq \lambda_{2,1}p(x|w_1)P(w_1) + \lambda_{2,2}p(x|w_2)P(w_2) \quad (104)$$

$$\frac{p(x|w_1)}{p(x|w_2)} P(w_1) (\lambda_{1,1} - \lambda_{2,1}) \leq P(w_2) (\lambda_{2,2} - \lambda_{1,2}) \quad (105)$$

$$\frac{p(x|w_1)}{p(x|w_2)} > || < \frac{P(w_2) (\lambda_{2,2} - \lambda_{1,2})}{P(w_1) (\lambda_{1,1} - \lambda_{2,1})} \quad (106)$$

If $>$ assign it to w_1 , if $<$ assign it to w_2 .

*If $\lambda_{1,2}$ and $\lambda_{2,1}$ are equal we get the MAP formula

$$\lambda_{1,2} = \lambda_{2,1} \text{ MinRisk} = \text{MAP}$$

Minimum risk criterion

if $\lambda_{2,2} = \lambda_{1,1} = 0$ (no punishment for correct d)

$$\Delta_1 \Lambda(x) = \frac{p(x|w_1)}{p(x|w_2)} > \frac{P(w_2) \lambda_{1,2}}{P(w_1) \lambda_{2,1}}$$

Discriminant Functions

A classifier assigns a feature vector \mathbf{x} to a class that has a higher discriminant function value, the discriminant functions ($g_i(x)$) we've seen:

- Minimize conditional risk: $g_i(x) = -R(\alpha_i|x)$
- Minimize error: $g_i(x) = P(w_i|x) \rightarrow \ln(p(x|w_i)P(w_i))$
- From here, we'll consider **MAP** as our discriminant function, and throw a $\ln()$ cause it's gaussian

remember

$$p(\mathbf{x}|\omega_i) = \frac{1}{(2\pi)^{n/2} |\Sigma_i|^{1/2}} \exp\left[-\frac{1}{2}(\mathbf{x} - \mathbf{m}_i)^T \Sigma_i^{-1} (\mathbf{x} - \mathbf{m}_i)\right] = N(\mathbf{m}_i, \Sigma_i)$$

with $i = 1, \dots, C$

- In this case, the discriminant function is written as

$$g_i(\mathbf{x}) = \ln p(\mathbf{x}|\omega_i) + \ln P(\omega_i)$$



$$g_i(\mathbf{x}) = -\frac{1}{2}(\mathbf{x} - \mathbf{m}_i)^T \Sigma_i^{-1} (\mathbf{x} - \mathbf{m}_i) - \frac{n}{2} \ln 2\pi - \frac{1}{2} \ln |\Sigma_i| + \ln P(\omega_i)$$

*If $g_1(x) = g_2(x)$ it means that the point is on the decision boundary and thus must be assigned randomly

1. $\sum = \sigma^2 I$ (aka isolevels take a spherical shape)

With this we know:

$$\sum_{i=1}^n = \frac{1}{\sigma^2} I \quad (107)$$

$$|\sum| = \sigma^{2n} \quad (108)$$

So simplifying the og equation:

$$g_i(x) = -\frac{1}{2}(x - m_i)^t \frac{I}{\sigma^2} (x - m_i) + \ln(P(w_i)) \quad (109)$$

$$= \frac{-1}{2\sigma^2} (x^t x - x^t m_i - m_i^t x - m_i^t m_i) + \ln(P(w_i)) \quad | \quad x^t m_i = m_i^t x \quad (110)$$

$$= \frac{m_i^t x}{\sigma^2} - \frac{m_i^t m_i}{2\sigma^2} + \ln(P(w_i)) \quad | \quad x^t x \rightarrow \text{ignored} \quad (111)$$

$$= w_i^t x + bias_i \quad (112)$$

where $w_i = \frac{m_i}{\sigma^2}$ and $bias_i = \frac{-m_i^t m_i}{2\sigma^2} + \ln(P(w_i))$

Now to get the equation to the decision boundary we calculate $g_i(x) = g_j(x)$ and

- In this case, they can be written as

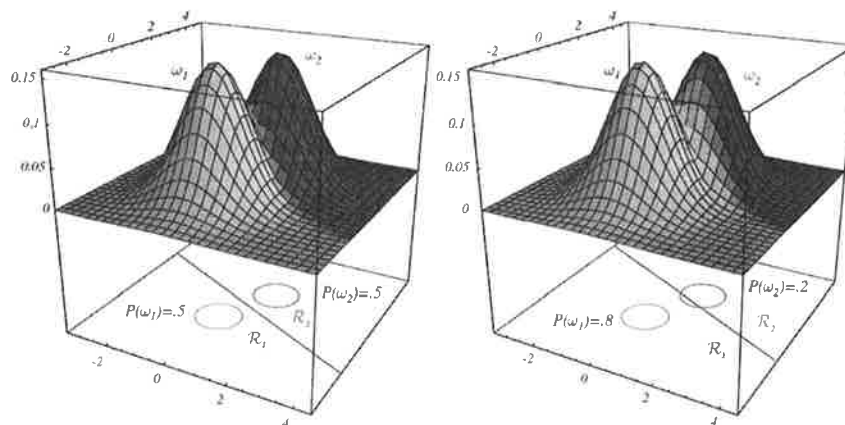
$$\mathbf{w}'(\mathbf{x} - \mathbf{x}_0) = 0$$

where

$$\begin{cases} \mathbf{w} = \mathbf{m}_i - \mathbf{m}_j \\ \mathbf{x}_0 = \frac{1}{2}(\mathbf{m}_i + \mathbf{m}_j) - \frac{\sigma^2}{\|\mathbf{m}_i - \mathbf{m}_j\|^2} \ln \frac{P(\omega_i)}{P(\omega_j)} (\mathbf{m}_i - \mathbf{m}_j) \end{cases}$$

Where $\frac{1}{2}(\mathbf{m}_i + \mathbf{m}_j)$ is the midpoint between the means and the rest is a bias "correction term mostly dependant on the prior probabilities" that gives more area to the biggest prob.

Examples with 2-D distributions



2. $\sum_i = \sum$ (aka all distributions don't have a standard shape?)
Same demonstrations as before but this time

- Decision boundaries are

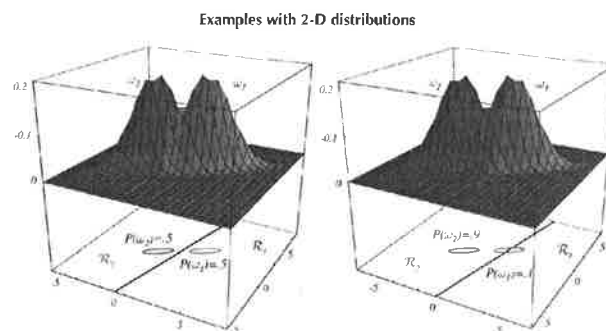
$$\mathbf{w}'(\mathbf{x} - \mathbf{x}_0) = 0$$

where

$$\begin{cases} \mathbf{w} = \Sigma^{-1}(\mathbf{m}_i - \mathbf{m}_j) \\ \mathbf{x}_0 = \frac{1}{2}(\mathbf{m}_i + \mathbf{m}_j) - \frac{\ln[P(\omega_i)/P(\omega_j)](\mathbf{m}_i - \mathbf{m}_j)}{(\mathbf{m}_i - \mathbf{m}_j)' \Sigma^{-1}(\mathbf{m}_i - \mathbf{m}_j)} \end{cases}$$

- Hyperplane passes through \mathbf{x}_0 but is not necessarily orthogonal to the line between the means.

w is now "rotated" by the covariance matrix, to take into account the shapes of the distributions. (Meaning the decision boundary line is no longer perpendicular to w , now it can be tilted)



3. \sum_i is arbitrary

- Discriminant functions are

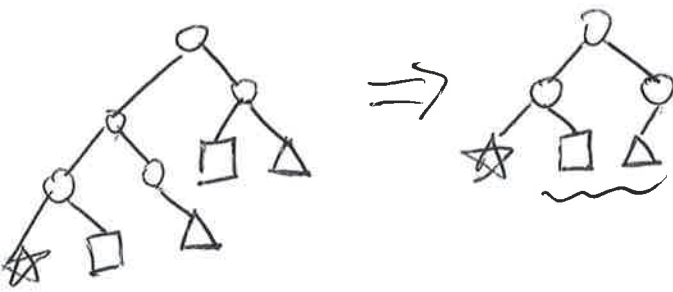
$$g_i(\mathbf{x}) = \mathbf{x}' \mathbf{W}_i \mathbf{x} + \mathbf{w}_i' \mathbf{x} + w_{i0} \quad \text{Quadratic discriminant}$$

where

$$\begin{cases} \mathbf{W}_i = -\frac{1}{2} \Sigma_i^{-1} \\ \mathbf{w}_i = \Sigma_i^{-1} \mathbf{m}_i \\ w_{i0} = -\frac{1}{2} \mathbf{m}_i' \Sigma_i^{-1} \mathbf{m}_i - \frac{1}{2} \ln |\Sigma_i| + \ln P(\omega_i) \end{cases}$$

- Decision boundaries are hyperquadrics.

Basically means the decision boundary no longer needs to be a line, it can take any shape



we prefer rules ^{that} making decision tree
looks like right

Decision trees

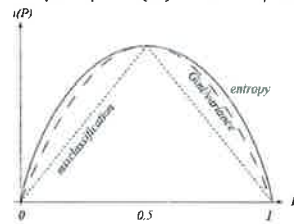
Used to classify non numerical data. We create a tree of questions that lead to the final class.

- CART - Classification and regression trees

Prefered to use a branching facto of 2, each time the data is split it's basically going into a smaller subset of data

Impurity A pure partition is when a question splits the data into perfectly pure subsets, i.e. they only contain one type of class

- Entropy impurity: $i(N) = - \sum_j P(\omega_j) \log_2 (P(\omega_j))$
- Variance impurity: $i(N) = P(\omega_1)P(\omega_2)$
- Gini impurity: $i(N) = \sum_{i \neq j} P(\omega_i)P(\omega_j) = 1 - \sum_j P^2(\omega_j)$
- Misclassification impurity: $i(N) = 1 - \max P(\omega_j)$



We use one of these impurities to calculate the drop in impurity:

$$\Delta i(N) = i(N) - P_L i(N_L) - (1 - P_L) i(N_R) \quad (113)$$

L and R being left and right node. This has a problem that the splits are local, i.e. they don't guarantee that the next one will give good results

✓ - Twoing Criterion What happens when I have more than two classes before a split??

1. We create superclasses, in the case of 3 classes we create:

$$C1 = [w_1] \quad C2 = [w_2, w_3] \quad C1 = [w_2] \quad C2 = [w_1, w_3] \quad C1 = [w_3] \quad C2 = [w_1, w_2] \quad (114)$$

2. Then calculate the impurity of each case, for each superclass
3. pick the best one

- When to stop splitting? Splitting too much leads to overfitting to the training data

Solutions:

* Stopped splitting

- Validation dataset, with accuracy measure aka cross-validation
- Impurity threshold - meaning stop splitting if Δi goes lower than a specified value

- Complexity-accuracy tradeoff criterion

$$\alpha(\text{tree size}) + \sum_{j=1}^N i(j) \quad (115)$$

Evaluate using the number of branches (ex) and the sum of the impurity of each end node at the current depth level. Bigger tree = overfit, unless there still is high impurity

- Hypothesis testing - Check if the new split is "statistically significant" i.e. is the imp. drop close to 0?

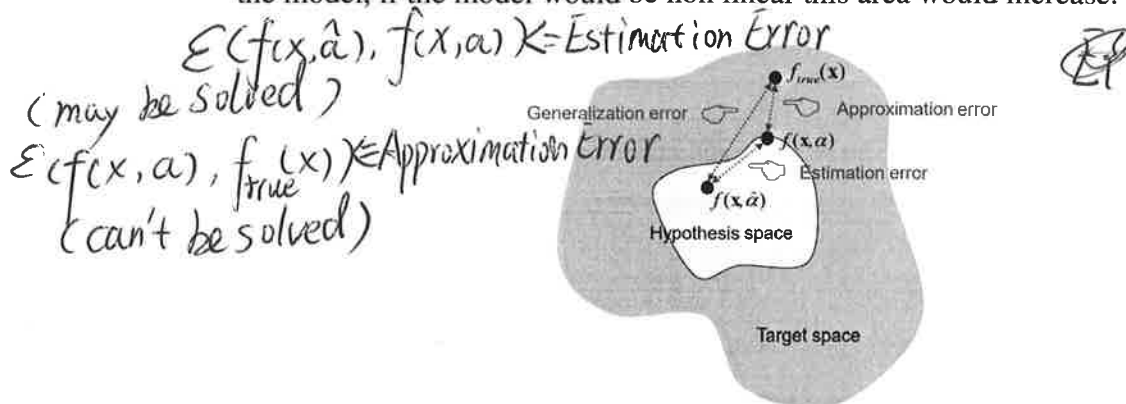
* Prunning

Let the tree reach overfit / impurity = 0 and remove branches that have a low impurity drop by "merging" the classes into one / merging the decision areas into one

It gets more computationally demanding if the training set is big

- Generalization error

When we train a linear model, there will be an error between the approximated model and the ground truth. Possibly caused by the limited dataset size, but even with an infinite amount of data we won't get the perfect model, and this is caused because it's linear, and it can't fit the real model perfectly. Here the hypothesis is given by the linear limitation of the model, if the model would be non linear this area would increase.



Estimation error caused by data limitations and **Approximation error** is caused by the model linear limitation

Gen. error:

$$E\{R(\alpha|x, y)\} = \int_{(x,y)} R(\alpha|x, y)p(x, y)dx dy \quad (116)$$

"The error is the average (E) of the cost (R) for any pair of points (x, y) belonging to the input and output spaces"

Where alpha in the case of a linear model is $\alpha = [w, b]$ equal to the weight and the bias $y' = wx + b$. And R is the loss function

The problem is that $p(x, y)$ is the "true joint pdf" which we don't have. So we estimate using the **empirical risk** i.e. the error computed on the training samples:

$$E\{R(\alpha|x, y)\} \approx \frac{1}{N} \sum_{i=1}^N L(y_i, f(x_i, \alpha)) \quad (117)$$

Accuracy Estimation

- Exhaustive cross-validation
 - Leave one out:
 - Remove one sample from training data, and use it as test. Then start from scratch and pick another sample as test, and keep repeating
- Non-exhaustive
 - Hold out
 - Split dataset in training and test (commonly by half)
 - K-fold
 - Make K clusters and apply same logic of leave one out
 - Monte Carlo
 - Repeat hold out various times
- Nested

Confusion matrix

Confusion Matrix and Common Derivations

	Decision (D_0)	Decision (D_1)
True (H_0)	TN	FP
True (H_1)	FN	TP

$$Accuracy = \frac{TN+TP}{TN+TP+FN+FP}$$

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

$$Sensitivity = \frac{TP}{TP+FN}$$

↑
Also called
"Recall"

$$Specificity = \frac{TN}{TN+FP}$$

$$F_\beta = \frac{(1+\beta^2)(Precision \cdot Sensitivity)}{\beta^2 Precision + Sensitivity}$$

$$F_1 = \frac{2 \cdot Precision \cdot Sensitivity}{Precision + Sensitivity}$$

Where the accuracy is the "overall accuracy" and the specificity is the "accuracy"

ROC curve Receiver Operating Characteristics

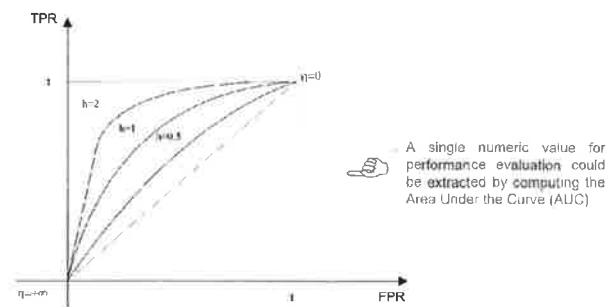
TPR - true positive rate $\int_T^{+\infty} p(x|w_2)dx$

FPR - false positive rate $\int_T^{+\infty} p(x|w_1)dx$

If the classes are not overlapped the curve is squared from 00 to 10 to 11

If the classes overlap completely it's TPR = FPR

Partial overlap are the cases in between



Comparing Classifiers

If I have two classifiers and want to know which one works better, how can I do this? And if the difference is really small, can I tell if it's just noise? We consider H_0 as the hypothesis that there is a difference between classif. and the opposite for H_1

- T Test**

$$T = \frac{\frac{\Delta a}{k}}{\sqrt{\frac{Var(\Delta)}{k}}} \quad (118)$$

Where $\frac{\Delta a}{k}$ is the average difference between the accuracy of the two classifiers given k folds And the $\sqrt{\frac{Var(\Delta)}{k}}$ the avg of the std deviations??

Then if T is lower than a given Tcritical value the difference is not statistically significant

$$T \leq T_c \rightarrow \text{reject } H_0 \mid T > T_c \rightarrow \text{accept } H_0 \quad (119)$$

- McNemar's Test**

Compare on each sample the two classifiers and checks if they agree and disagree in the same way using a contingency table:

Sample	C_1	C_2
#1	correct	correct
#2	correct	wrong
...
#i	wrong	wrong
...
#N	wrong	correct

	C_2 correct	C_2 wrong
C_1 correct	N_{11}	N_{10}
C_1 wrong	N_{01}	N_{00}

χ^2 distribution

It has it's own formula and a similar threshold to T test:

$$Z^2 = \frac{(N_{10} - N_{01})^2}{N_{10} + N_{01}} \quad (120)$$