

# Maximum-Likelihood and Bayesian Parameter Estimation

- Sufficient Statistics
- Common Probability Distributions
- Problems of Dimensionality

# Problems of Dimensionality

Problems involving 50 or 100 features  
(binary valued)

- Classification accuracy depends upon the dimensionality and the amount of training data
- Case of two classes multivariate normal with the same covariance

$$P(\text{error}) = \frac{1}{\sqrt{2\pi}} \int_{r/2}^{\infty} e^{-\frac{u^2}{2}} du$$

$$\text{where : } r^2 = (\mu_1 - \mu_2)^t \Sigma^{-1} (\mu_1 - \mu_2)$$

$$\lim_{r \rightarrow \infty} P(\text{error}) = 0$$

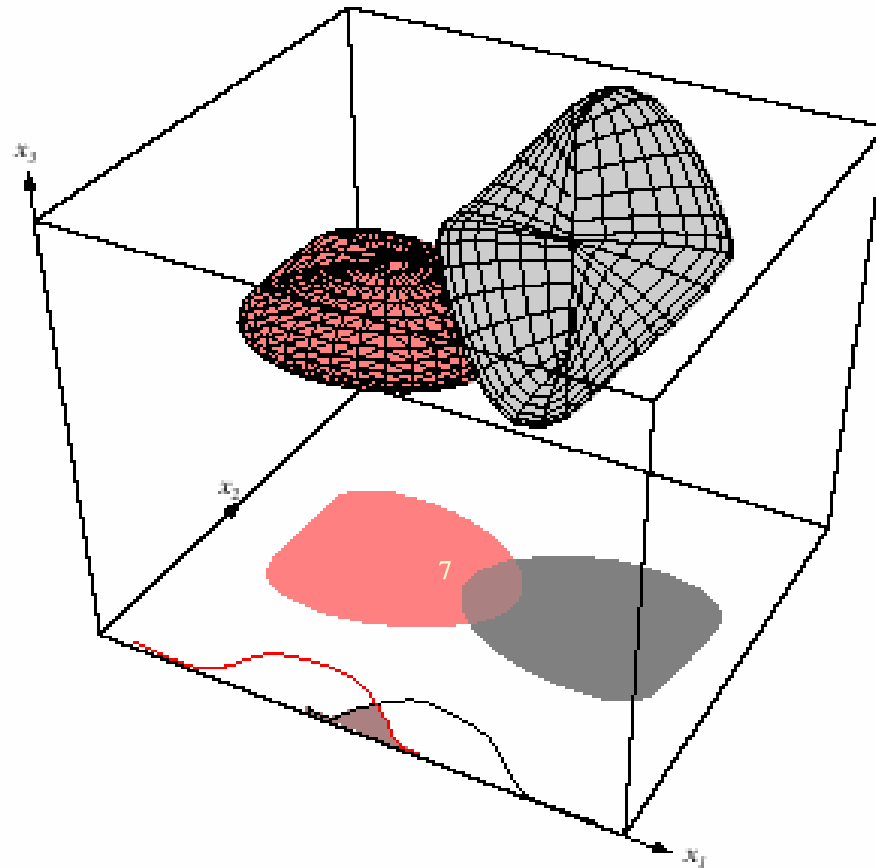
# Error Rate and Dimensionality

- If features are independent then:

$$\Sigma = \text{diag}(\sigma_1^2, \sigma_2^2, \dots, \sigma_d^2)$$
$$r^2 = \sum_{i=1}^d \left( \frac{\mu_{i1} - \mu_{i2}}{\sigma_i} \right)^2$$

- Most useful features are the ones for which the difference between the means is large relative to the standard deviation
- It has frequently been observed in practice that, beyond a certain point, the inclusion of additional features leads to worse rather than better performance.

# Decrease in error rate with features



Non-overlapping distributions in 3-dimensions where Bayes error vanishes. When projected to  $x_1$ - $x_2$  sub-space or  $x_1$  sub-space there is a greater overlap and hence greater Bayes error

# Computational Complexity

Design methodology is affected by computational difficulty

“big oh” notation  $f(x) = O(h(x))$  “big oh of  $h(x)$ ”

If:

$$\exists \mathbf{c}, \mathbf{x}_0 \ni |\mathbf{f}(\mathbf{x})| \leq \mathbf{c}|\mathbf{h}(\mathbf{x})|$$
$$\forall \mathbf{x} > \mathbf{x}_0$$

(An upper bound on  $f(x)$  grows no worse than  $h(x)$  for sufficiently large  $x$ )

Example:  $f(x) = 2 + 3x + 4x^2$

$h(x) = x^2$  and  $c$  can be appropriately chosen

Therefore  $f(x) = O(x^2)$

# Big Theta Notation

“Big oh” is not unique

e.g.,  $f(x) = O(x^2)$ ;  $f(x) = O(x^3)$ ;  $f(x) = O(x^4)$

Thus we introduce “big theta” notation

$f(x) = \theta(h(x))$  if:

$$\exists \mathbf{x}_0, \mathbf{c}_1, \mathbf{c}_2 \in \forall \mathbf{x} > \mathbf{x}_0$$

$$0 \leq \mathbf{c}_1 \mathbf{h}(\mathbf{x}) \leq \mathbf{f}(\mathbf{x}) \leq \mathbf{c}_2 \mathbf{h}(\mathbf{x})$$

Therefore  $f(x) = \theta(x^2)$  but  $f(x) \neq \theta(x^3)$

# Complexity of ML Estimation

Gaussian priors in  $d$  dimensions classifier with  $n$  training samples for each of  $c$  classes

- For each category, we have to compute the discriminant function

$$g(\mathbf{x}) = -\frac{1}{2}(\mathbf{x} - \hat{\boldsymbol{\mu}})^t \overbrace{\boldsymbol{\Sigma}^{-1}}^{O(n \cdot d^2)} (\mathbf{x} - \hat{\boldsymbol{\mu}}) - \overbrace{\frac{d}{2} \ln 2\pi}^{O(1)} - \underbrace{\frac{1}{2} \ln |\hat{\boldsymbol{\Sigma}}|}_{O(d^2 \cdot n)} + \underbrace{\ln P(\omega)}_{O(n)}$$

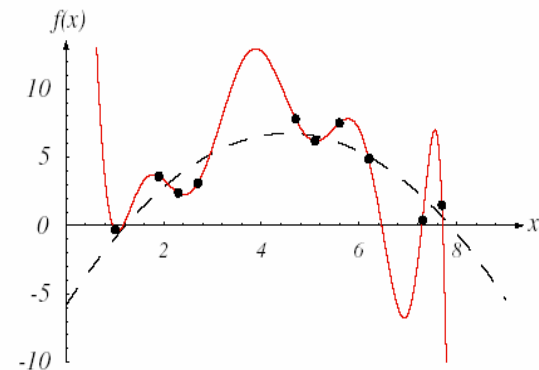
Total =  $O(d^2 \cdot n)$

Total for  $c$  classes =  $O(cd^2 \cdot n) \cong O(d^2 \cdot n)$

- Cost increase when  $d$  and  $n$  are large!

# Overfitting

- Frequently, no of available samples usually inadequate, how to proceed?
- Solutions:
  - **Reduce dimensionality**: redesign feature extractor, or use a subset of features
  - **Assume all  $c$  classes share same covariance matrix**: pool available data
  - **Assume statistical independence**: all off-diagonal elements of covariance matrix are zero
- Classifier performs better than if we overfit the data, why?



Training Data from quadratic function plus Gaussian noise  
 $f(x)=ax^2+bx+c+\varepsilon$  where  $p(\varepsilon) \sim N(0, \sigma^2)$ . Tenth order fits perfectly but second order performs better for new samples



# Capacity of a Separating Plane

Overdetermined solution is significant for classification as it is for estimation

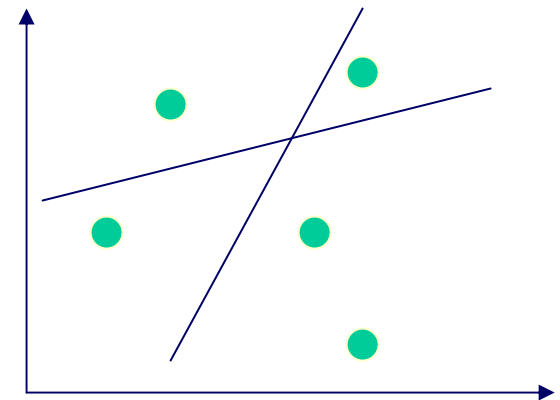
Task of partitioning a  $d$ -dimensional feature space by a hyperplane  $n$  samples in general position labelled either  $\omega_1$  or  $\omega_2$ .

Of the  $2^n$  possible dichotomies a fraction  $f(n,d)$  are linear dichotomies

$$f(n,d) = \begin{cases} 1 & n \leq d+1 \\ \frac{2}{2^n} \sum_{i=0}^d \binom{n-1}{i} & n > d+1 \end{cases}$$

Fraction of dichotomies of  $n$  points in  $d$  dimensions that are linear

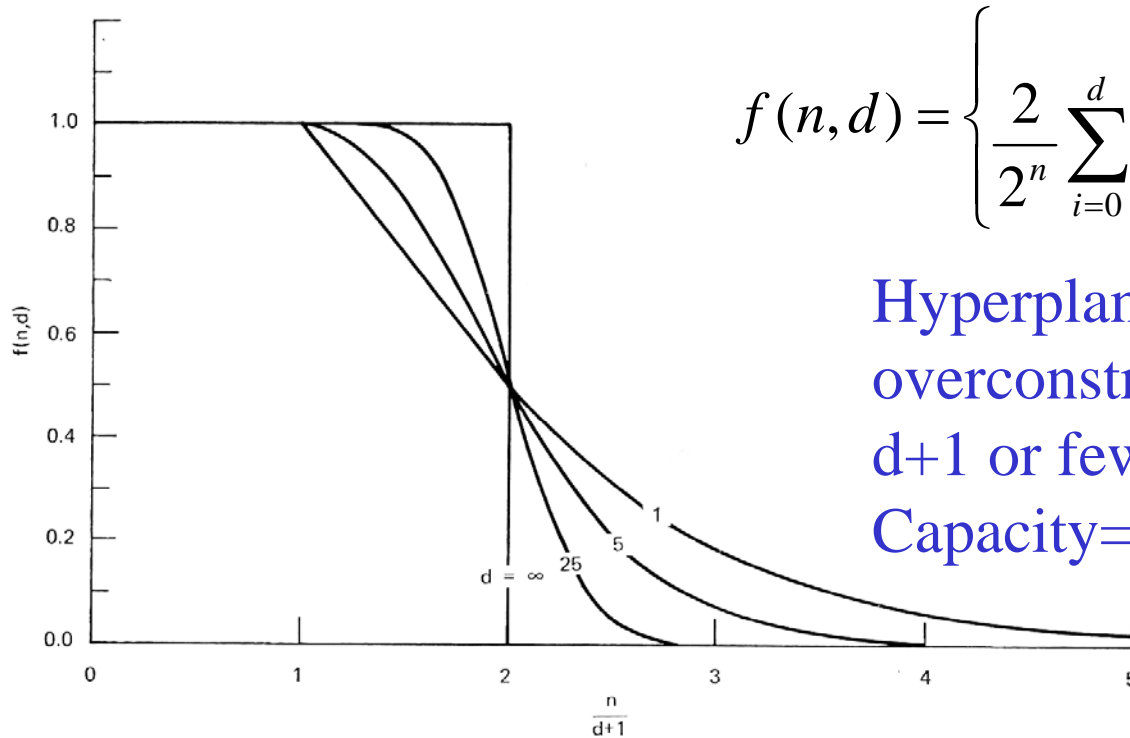
$$f(5,2) = 11/16$$



5 points can be labelled as either class in 32 ways

# Capacity of a Separating Plane

Of the  $2^n$  possible dichotomies a fraction  $f(n, d)$  are linear dichotomies



$$f(n, d) = \begin{cases} 1 & n \leq d + 1 \\ \frac{2}{2^n} \sum_{i=0}^d \binom{n-1}{i} & n > d + 1 \end{cases}$$

Hyperplane is not  
overconstrained in classifying  
 $d+1$  or fewer points  
Capacity =  $2(d+1)$

Fraction of dichotomies of  $n$  points in  $d$  dimensions that are linear