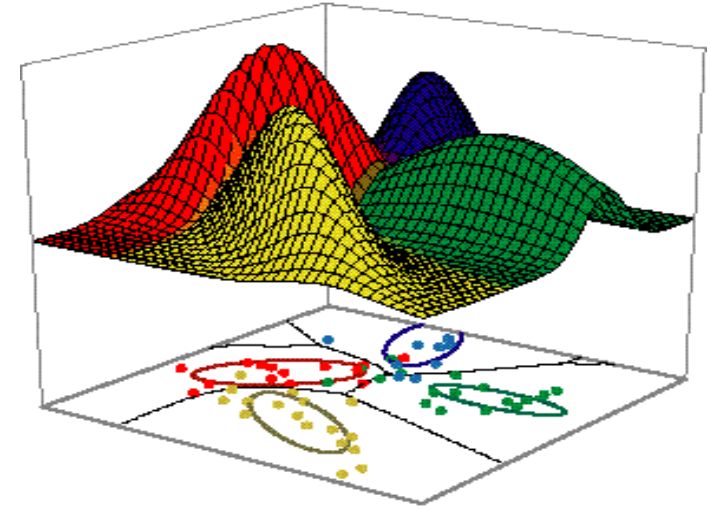


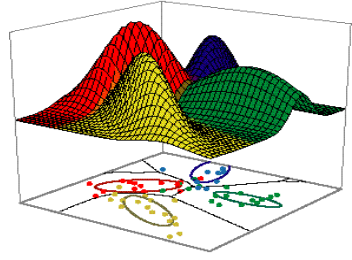
Part 6: Parameter Estimation



Parametric vs. Nonparametric
Maximum Likelihood Parameter Estimation
Bayesian Parameter Estimation

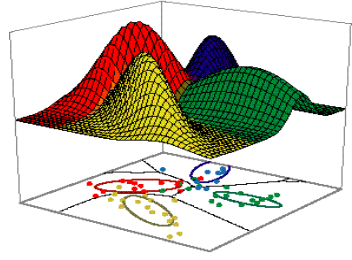
Some materials in these slides were taken from [Pattern Classification](#) (2nd ed) by R. O. Duda, P. E. Hart and D. G. Stork, John Wiley & Sons, 2000, **Chapter 3.1-3.5.**

Introduction



- Until now, have assumed that we know the distribution of all class-conditional density functions & priors
 - e.g., know that $p(x|\omega_1) \sim N(0.5, 6.2)$
 - If we know this, we can create optimal classifiers.
- In reality, we rarely know this information!
 - How do we estimate these density functions from the data?
 - Samples are often too small for direct class-conditional estimation (large dimension of feature space!)

Introduction



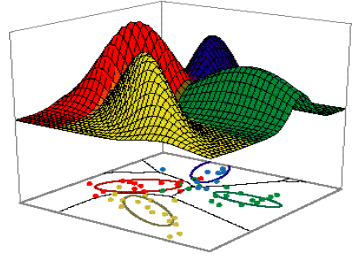
- 2 Fundamental Approaches:
 - Parameter Estimation:
 - We may know (or assume) form of distribution; estimate parameters of that distribution from the data
→ requires prior knowledge
 - e.g., assume that $p(x|\omega_1) \sim N(\mu, \sigma)$, now estimate parameters μ & σ from the data
 - Now only have to estimate 2 parameters instead of entire distribution.
 - Nonparametric Density Estimation:
 - Don't assume any particular form of distribution
 - Let the data determine the density function directly.
 - May not result in analytical solution

Parameter Estimation



- Typically discuss 2 approaches to estimating $p(x|\omega_1)$
 - Results are nearly identical, but the approaches are different
 - 1) Maximum Likelihood Parameter Estimation
 - Assume parameters, θ , have fixed but unknown values
 - Estimate the parameter values that make the data most likely to have been observed given the value of θ
 - 2) Bayesian Parameter Estimation
 - Assume that parameters, θ , are random variables with some known prior distribution.
 - Observations convert this to a *posteriori* density
 - Additional observations sharpen the *a posteriori* density and cause it to peak near actual values of θ
- In either approach, use $P(\omega_i | x)$ for classification rule!

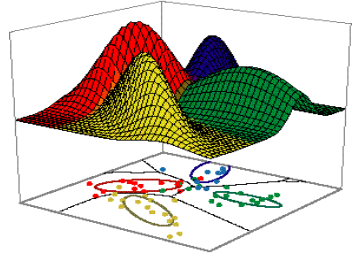
Maximum Likelihood Estimation



- Has good convergence properties as the sample size increases
- Simpler than any alternative techniques
- General principle
 - Assume we have c classes and
$$P(x | \omega_j) \sim N(\mu_j, \Sigma_j)$$
$$P(x | \omega_j) \equiv P(x | \omega_j, \theta_j) \text{ where:}$$

$$\theta_j = (\mu_j, \Sigma_j) = (\mu_j^1, \mu_j^2, \dots, \mu_j^d, \sigma_j^{11}, \sigma_j^{22}, \dots, \sigma_j^{dd}, \text{cov}(x_j^k, x_j^l) \dots)$$

Maximum Likelihood Estimation



- Use the information provided by the training samples to estimate $\theta = (\theta_1, \theta_2, \dots, \theta_c)$, each θ_i ($i = 1, 2, \dots, c$) is associated with a category
- Suppose that D contains n samples, x_1, x_2, \dots, x_n **drawn independently**

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

$P(D | \theta)$ is called the likelihood of θ w.r.t. the set of samples

- ML estimate of θ is, by definition the value that maximizes $P(D | \theta)$. Referred to as: $\hat{\theta}$

“It is the value of θ that best agrees with the actually observed training sample”

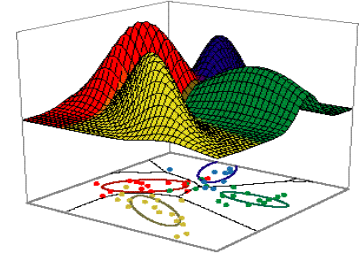
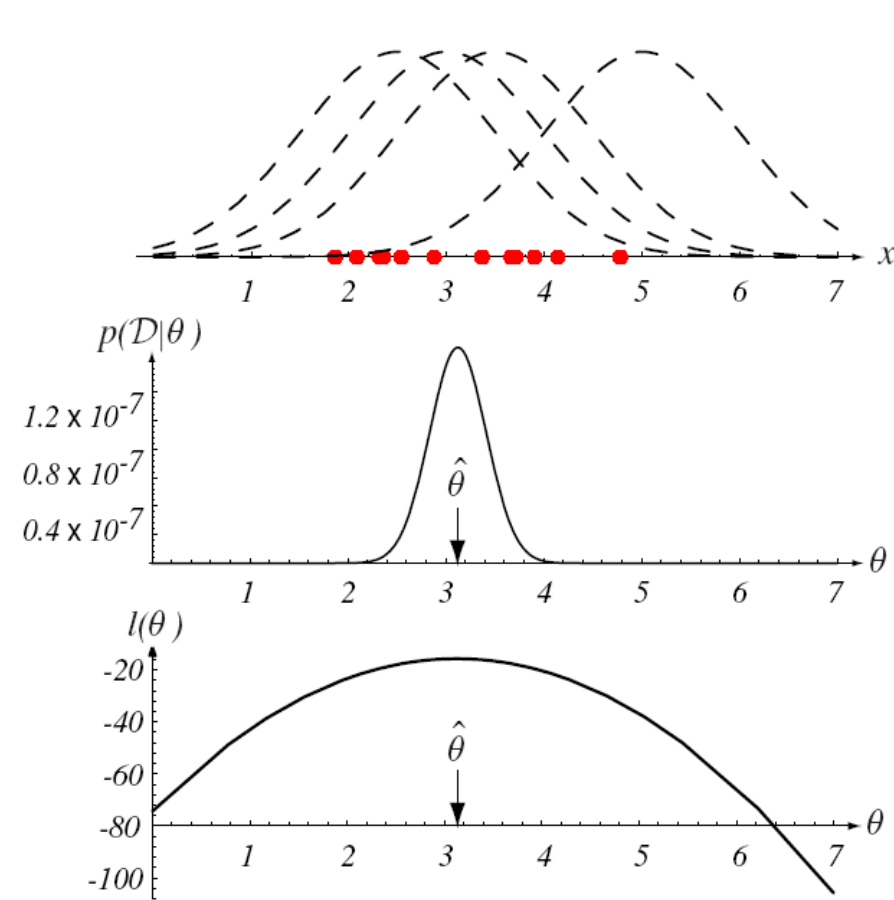
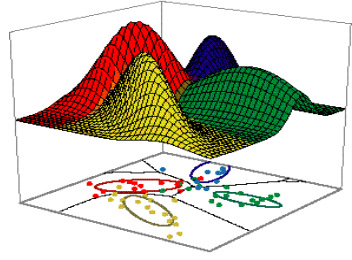


FIGURE 3.1. The top graph shows several training points in one dimension, known or assumed to be drawn from a Gaussian of a particular variance, but unknown mean. Four of the infinite number of candidate source distributions are shown in dashed lines. The middle figure shows the likelihood $p(\mathcal{D}|\theta)$ as a function of the mean. If we had a very large number of training points, this likelihood would be very narrow. The value that maximizes the likelihood is marked $\hat{\theta}$; it also maximizes the logarithm of the likelihood—that is, the log-likelihood $l(\theta)$, shown at the bottom. Note that even though they look similar, the likelihood $p(\mathcal{D}|\theta)$ is shown as a function of θ whereas the conditional density $p(x|\theta)$ is shown as a function of x . Furthermore, as a function of θ , the likelihood $p(\mathcal{D}|\theta)$ is not a probability density function and its area has no significance. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.

Maximum Likelihood Estimation



- Optimal estimation
 - Let $\theta = (\theta_1, \theta_2, \dots, \theta_p)^t$ and let ∇_{θ} be the gradient operator

$$\nabla_{\theta} = \left[\frac{\partial}{\partial \theta_1}, \frac{\partial}{\partial \theta_2}, \dots, \frac{\partial}{\partial \theta_p} \right]^t$$

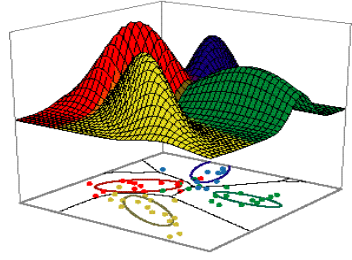
- We define $l(\theta)$ as the log-likelihood function

$$l(\theta) = \ln P(D | \theta)$$

- New problem statement:
determine θ that maximizes the log-likelihood

$$\hat{\theta} = \arg \max_{\theta} l(\theta) \quad \text{implicit dependence on } D$$

Maximum Likelihood Estimation



Set of necessary conditions for an optimum is:

$$\nabla_{\theta} l = \sum_{k=1}^n \nabla_{\theta} \ln P(x_k | \theta)$$

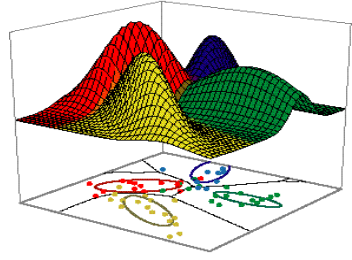
Since samples drawn IID

$$\nabla_{\theta} l = 0$$

Take the gradient, and set to zero

A solution could represent global min/max, local min/max, or point of inflection

Maximum Likelihood Estimation



- Example of a specific case: unknown μ
 - $p(x_i | \mu) \sim N(\mu, \Sigma)$
(Samples are drawn from a multivariate normal population)

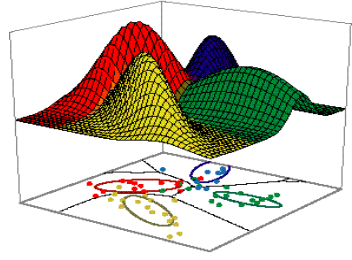
$$\ln p(x_k | \mu) = -\frac{1}{2} \ln[(2\pi)^d |\Sigma|] - \frac{1}{2} (x_k - \mu)^t \Sigma^{-1} (x_k - \mu)$$

$$\theta = \mu \text{ therefore: } \nabla_{\mu} \ln p(x_k | \mu) = \Sigma^{-1} (x_k - \mu)$$

- The ML estimate for μ must satisfy:

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

Maximum Likelihood Estimation



- Multiplying by Σ and rearranging, we obtain:

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

Simply the arithmetic average of the samples of the training samples!

Conclusion:

If $p(x_k | \omega_j)$ ($j = 1, 2, \dots, c$) is supposed to be Gaussian in a d -dimensional feature space; then we can estimate the vector $\theta = (\theta_1, \theta_2, \dots, \theta_c)^t$ and perform an optimal classification!

Maximum Likelihood Estimation

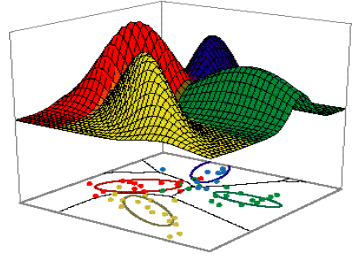


- ML Estimation:
 - 1D Gaussian Case: *unknown* μ and σ
 $\theta = (\theta_1, \theta_2) = (\mu, \sigma^2)$
 - Consider for a single point x_k :

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

$$\nabla_{\theta} l = \begin{pmatrix} \frac{\partial}{\partial \theta_1} (\ln P(x_k | \theta)) \\ \frac{\partial}{\partial \theta_2} (\ln P(x_k | \theta)) \end{pmatrix} = 0 \rightarrow \begin{cases} \frac{1}{\theta_2} (x_k - \theta_1) = 0 \\ -\frac{1}{2\theta_2} + \frac{(x_k - \theta_1)^2}{2\theta_2^2} = 0 \end{cases}$$

Maximum Likelihood Estimation



Summation over all x_k :

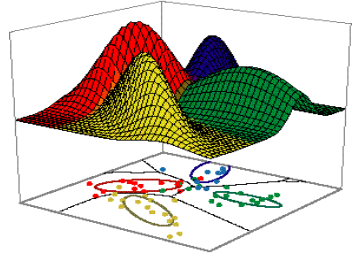
$$\begin{cases} \sum_{k=1}^n \frac{1}{\hat{\theta}_2} (x_k - \hat{\theta}_1) = 0 & (1) \\ -\sum_{k=1}^n \frac{1}{\hat{\theta}_2} + \sum_{k=1}^n \frac{(x_k - \hat{\theta}_1)^2}{\hat{\theta}_2^2} = 0 & (2) \end{cases}$$

Combining (1) and (2), one obtains:

$$\hat{\mu} = \sum_{k=1}^n \frac{x_k}{n} \quad \hat{\sigma}^2 = \frac{\sum_{k=1}^n (x_k - \hat{\mu})^2}{n}$$

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

Maximum Likelihood Estimation



- Bias
 - Bias of an estimate is the systematic/deterministic error
 - i.e., diff between true value, θ , and expected value of its estimate, $\hat{\theta}$.
 - ML estimate for σ^2 is biased

$$E \left[\frac{1}{n} \sum (x_i - \bar{x})^2 \right] = \frac{n-1}{n} \cdot \sigma^2 \neq \sigma^2$$

- An elementary unbiased estimator for Σ is:

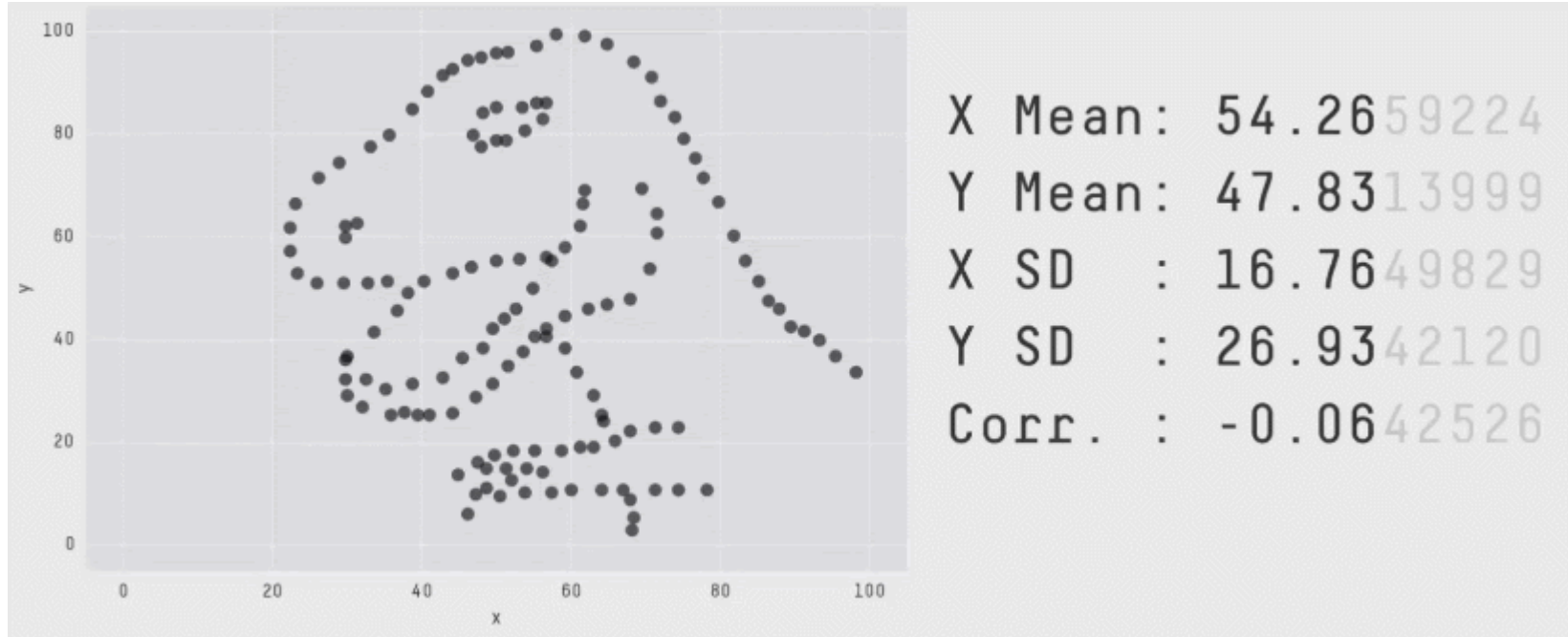
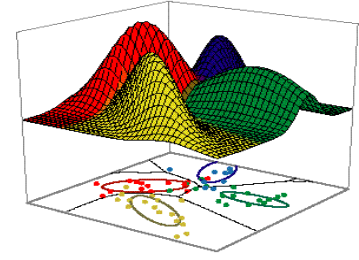
$$C = \frac{1}{n-1} \sum_{k=1}^{k=n} (x_k - \bar{x})^2$$

Sample covariance

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

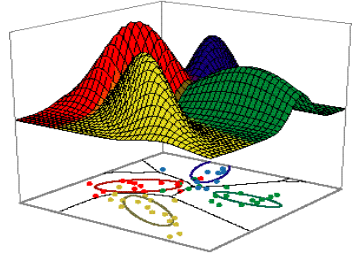
Sample covariance matrix (multivariate)

Warning: Visualize Data Before Assuming Distribution



Justin Matejka and George Fitzmaurice. **Same stats, different graphs: Generating datasets with varied appearance and identical statistics through simulated annealing.** In Proceedings of the 2017 CHI Conference on Human Factors in Computing Systems, pages 1290–1294. ACM, 2017.

Bayesian Estimation

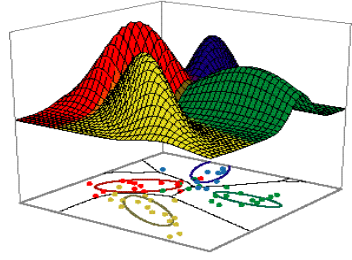


- (aka Bayesian learning when used for pattern classification)
- Assume form of density is known, parameter θ is not
 - In MLE, θ was assumed to be fixed, but unknown
 - In BE θ is a random variable, training data allows us to convert dist on this variable into a posterior probability density $P(\theta | x)$
 - Either way, $p(x)$ is not known, but $p(x|\theta)$ is completely known
- Ultimate goal: compute $P(\omega_i | x, D)$
 - Given the sample D , Bayes' formula can be written

$$P(\omega_i | x, D) = \frac{P(x | \omega_i, D)P(\omega_i | D)}{\sum_{j=1}^c P(x | \omega_j, D)P(\omega_j | D)}$$

- Use information from D to help determine class-conditional and prior probabilities. Will henceforth assume priors are known or easily estimated (simplify notation...)

Bayesian Estimation

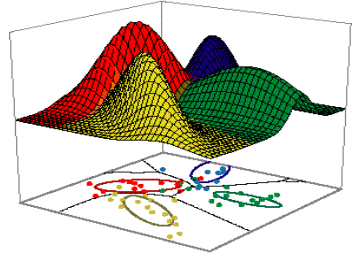


- Assume that samples in D_i have no influence on $p(x|w_j, D)$ if $i \neq j$, so:

$$P(\omega_i | x, D) = \frac{P(x | \omega_i, D_i)P(\omega_i)}{\sum_{j=1}^c P(x | \omega_j, D_j)P(\omega_j)}$$

- Allows us to work with each class separately

Bayesian Estimation



- Gaussian Case

- Goal: Estimate θ using the *a posteriori* density $P(\theta \mid D)$
(recall that ML maximized $P(\mathcal{D}|\theta)$ instead)
- The univariate case: $P(\mu \mid D)$
 μ is the only unknown parameter

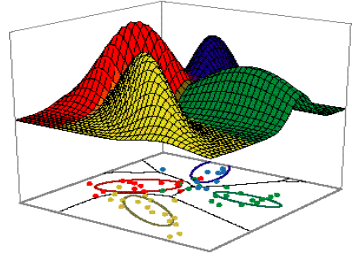
$$p(x \mid \mu) \sim N(\mu, \sigma^2)$$

$$p(\mu) \sim N(\mu_0, \sigma_0^2) \longleftarrow \text{Prior distribution of the unknown mean } \mu.$$

(μ_0 and σ_0 are known!)

Imagine value is drawn for μ from a population governed by $p(\mu)$. Once drawn, it becomes the true value for μ and completely determines density of x . Now n Samples $D = \{x_1, \dots, x_n\}$ are drawn independently from resulting population.

Bayesian Estimation



$$p(\mu | \mathbf{D}) = \frac{P(\mathbf{D} | \mu) p(\mu)}{\int P(\mathbf{D} | \mu) p(\mu) d\mu} \quad (1)$$

$$= \alpha \prod_{k=1}^{k=n} P(x_k | \mu) p(\mu) \quad \begin{array}{l} \alpha \text{ depends on } D, \text{ but is} \\ \text{indep of } \mu \end{array}$$

- Reproducing density $p(\mu|D)$ is the product of two Normal distributions. Results in a new Normal:

$$P(\mu | \mathbf{D}) \sim N(\mu_n, \sigma_n^2) \quad (2)$$

- Equating (1) and (2) yields *(skipping math on p93):*

$$\mu_n = \left(\frac{n\sigma_0^2}{n\sigma_0^2 + \sigma^2} \right) \hat{\mu}_n + \frac{\sigma^2}{n\sigma_0^2 + \sigma^2} \cdot \mu_0 \quad \begin{array}{l} \mu_n \text{ is best guess for } \mu \text{ after } n \\ \text{points. Weighted combination of} \\ \text{sample mean, and prior.} \end{array}$$

$$\text{and } \sigma_n^2 = \frac{\sigma_0^2 \sigma^2}{n\sigma_0^2 + \sigma^2} \quad \begin{array}{l} \sigma_n \text{ measures uncertainty about } \mu. \text{ as } n \uparrow, \sigma_n \downarrow. \\ \text{Leads to a dirac dist about true } \mu \rightarrow \text{Bayesian Learning} \end{array}$$

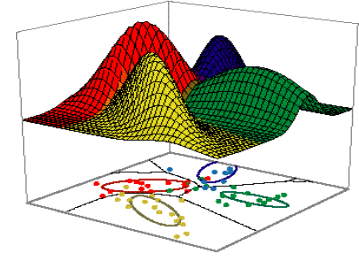
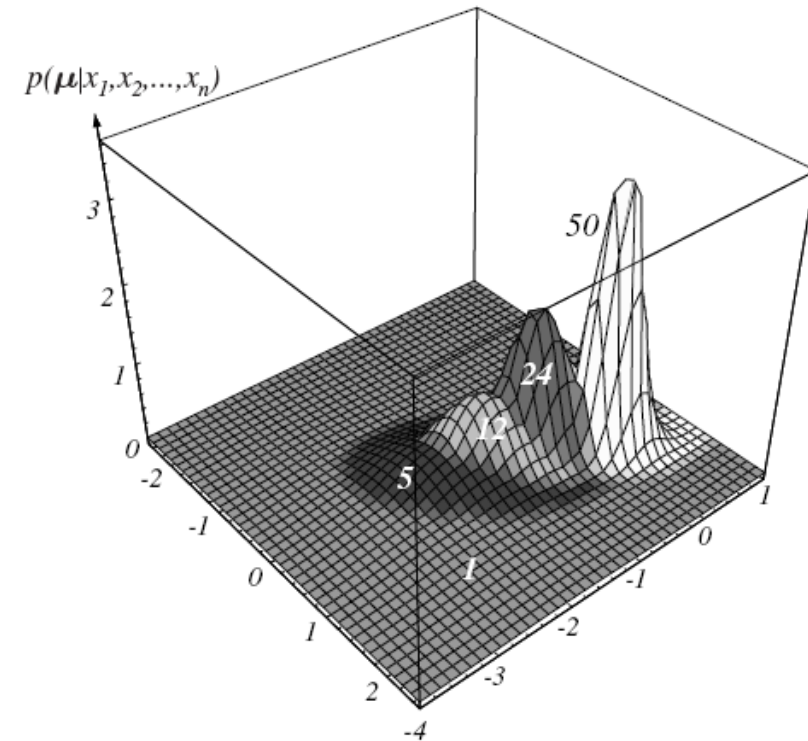
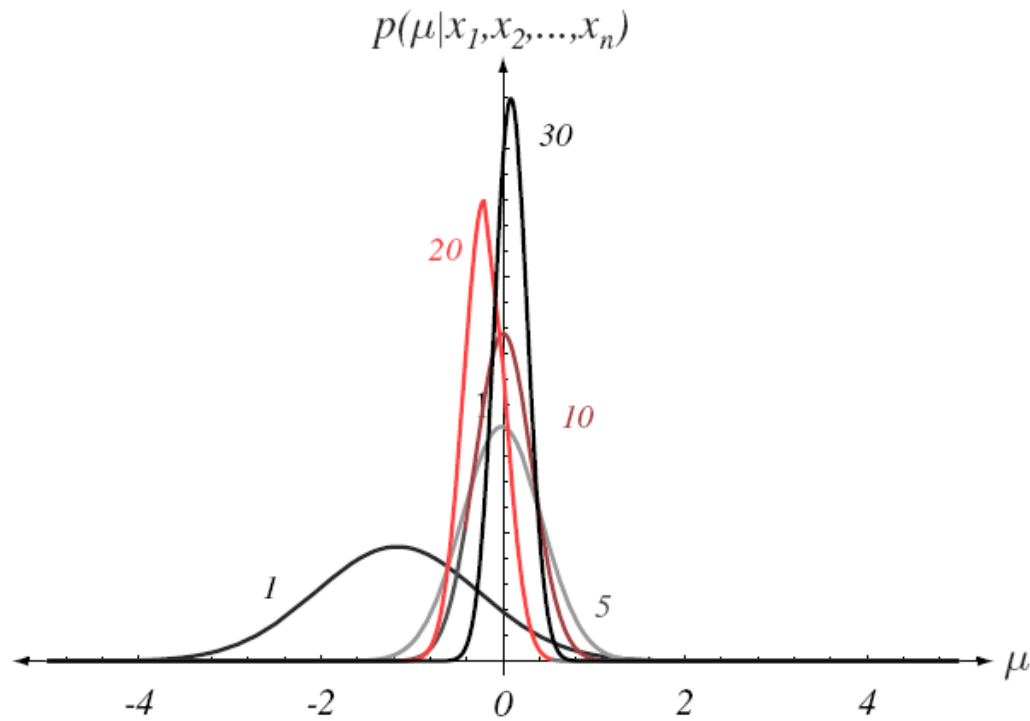
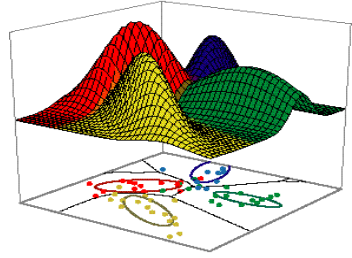


FIGURE 3.2. Bayesian learning of the mean of normal distributions in one and two dimensions. The posterior distribution estimates are labeled by the number of training samples used in the estimation. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.

Bayesian Estimation



- The univariate case $P(x | \mathcal{D})$
 - $P(\mu | D)$ computed
 - $P(x | D)$ remains to be computed!

$$p(x | \mathbf{D}) = \int p(x | \mu) p(\mu | \mathbf{D}) d\mu \text{ is Gaussian}$$

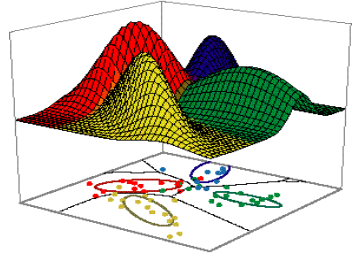
It provides: $p(x | \mathbf{D}) \sim N(\mu_n, \sigma^2 + \sigma_n^2)$ *Variance increased due to uncertainty about mean.*

(This is desired class-conditional density $P(x | \omega_j, D_j)$)

Finally, using Bayes formula to combine $P(x | \omega_j, D_j)$ together with $P(\omega_j)$, we obtain the Bayesian classification rule:

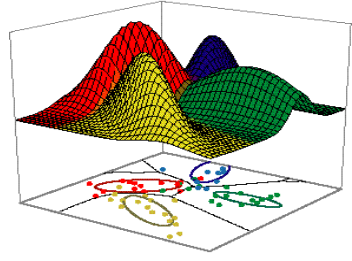
$$\underset{\omega_j}{Max} [P(\omega_j | x, \mathbf{D})] \equiv \underset{\omega_j}{Max} [P(x | \omega_j, \mathbf{D}_j) \cdot P(\omega_j)]$$

Bayesian Estimation : General Theory



- $p(x \mid \mathcal{D})$ computation can be applied to any situation in which the unknown density can be parameterized (not just univariate Gaussian)
- The basic assumptions are:
 - The form of $p(x \mid \theta)$ is assumed known, but the value of θ is not known exactly
 - Our knowledge about θ is assumed to be contained in a known prior density $p(\theta)$
 - The rest of our knowledge θ is contained in a set D of n random variables x_1, x_2, \dots, x_n that follows $p(x)$

Bayesian Estimation : General Theory



The basic problem is:

“1) Compute the posterior density $P(\theta | \mathcal{D})$ ”

Using Bayes formula, we have:

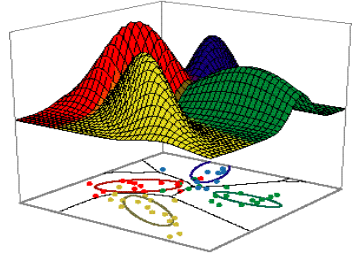
$$p(\theta | \mathbf{D}) = \frac{p(\mathbf{D} | \theta) p(\theta)}{\int p(\mathbf{D} | \theta) p(\theta) d\theta}$$

And by independence assumption: $P(\mathbf{D} | \theta) = \prod_{k=1}^{k=n} P(x_k | \theta)$

then “2) Derive $P(\mathbf{x} | \mathcal{D})$ ”

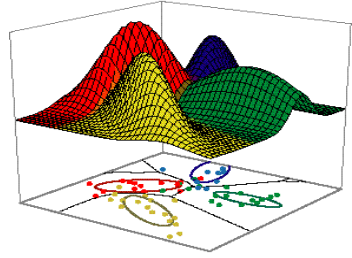
$$p(x | \mathbf{D}) = \int p(\mathbf{x} | \theta) p(\theta | \mathbf{D}) d\theta$$

Recursive Bayesian Approach



- Given:
$$P(\mathbf{D} \mid \theta) = \prod_{k=1}^{k=n} P(x_k \mid \theta)$$
- We get:
$$p(\mathbf{D}^n \mid \theta) = p(x_n \mid \theta) p(D^{n-1} \mid \theta)$$
- And:
$$p(\theta \mid \mathbf{D}^n) = \frac{p(\mathbf{x}_n \mid \theta) p(\theta \mid D^{n-1})}{\int p(\mathbf{x}_n \mid \theta) p(\theta \mid D^{n-1}) d\theta} \quad p(\theta \mid \mathbf{D}^0) = p(\theta)$$
- Recursive relationship.
 - Example of *incremental* or *on-line* learning
 - Learn as data is collected.
 - Many methods require all training data must be present before learning takes place.

ML vs. Bayes Methods



- For reasonable priors that include true solution, ML & Bayes equivalent as $n \rightarrow \infty$
 - However, we are normally faced with limited data
 - If $p(\theta|D)$ is broad or asymmetric around ML estimate for θ , resulting $p(x|D)$ will differ.
- Several criteria for comparing:
 - Computational complexity \rightarrow ML preferred
 - Compute gradient rather than multidimensional integration.
 - Interpretability \rightarrow ML often preferred.
 - Returns single best model, not weighted distribution of possible models.
 - However, Bayes makes bias/variance trade-off more explicit.
 - Prior information \rightarrow Bayes preferred
 - If dependable prior information is available, Bayes can use it more completely.
 - If prior does not contain true model, ML will never find it. Bayes can (not strictly constrained by assumed parametric model)
 - Flat or uninformative priors make methods similar.
- Recall 3 sources of error in pattern classification
 - 1) Bayes or Indistinguishability Error: Can never be eliminated.
 - 2) Model Error: Model based on prior information. Generally not dependent on parameter estimation approach that will be taken.
 - 3) Estimation error: Need more training data \rightarrow independent of P.E. method.