

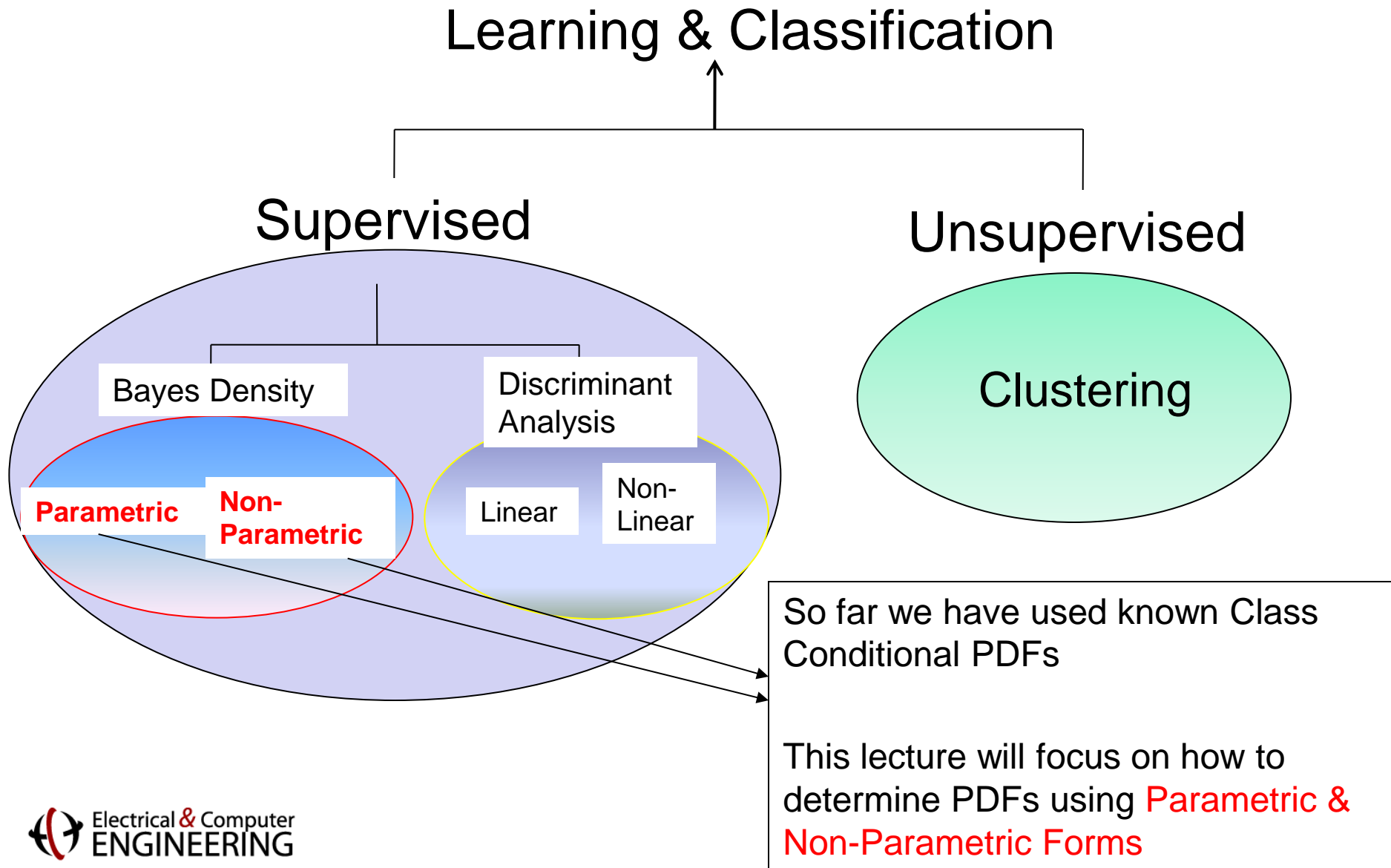
Prof. Marios Savvides

# Pattern Recognition Theory

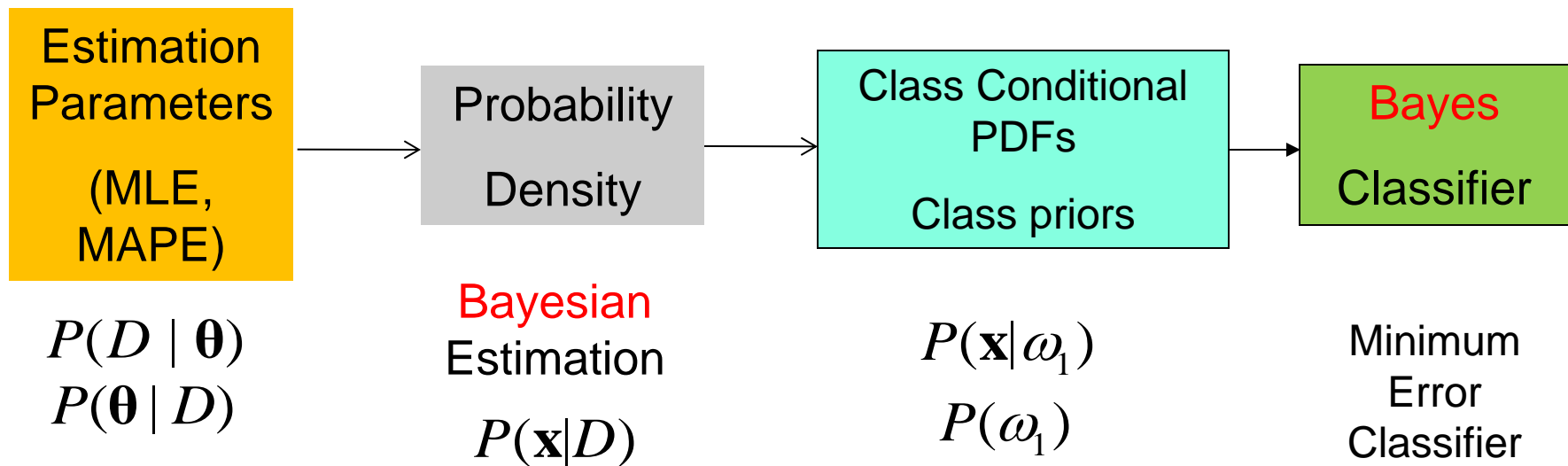
## Lecture 7: Parametric & Non Parametric Density Estimation

All graphics from Pattern Classification, Duda, Hart and Stork, Copyright © John Wiley and Sons, 2001

# Overview Of Pattern Recognition



# Big Picture



# Overview

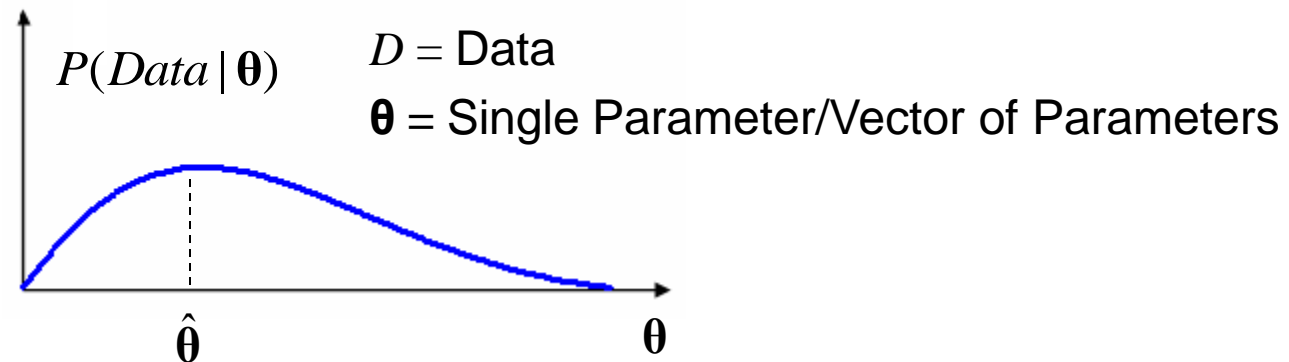
- Previous lectures have shown how to develop classifiers when the underlying statistical structure is known
- Parametric Estimation
  - This method assumes a **particular form** of a PDF (e.g. Gaussian) is known so that we only need to determine the **parameters** (e.g. Mean & Variance)
    - Maximum Likelihood Estimation (MLE)
    - Maximum A Posteriori (Bayesian) Estimation (MAPE)
- Non-Parametric Density Estimation
  - This method **does not assume ANY knowledge** about the density
    - Parzen Windows
    - Kernel Density Estimation
    - K-Nearest Neighbor Rule

# ML Estimation (MLE)

- Maximum Likelihood Estimation

- Assume  $P(\mathbf{x}|\omega)$  has a known parametric form uniquely determined by the parameter vector  $\theta$
- The parameters are assumed to be **FIXED ( i.e. NON RANDOM)** but unknown
- Suppose we have a dataset  $D$  with the samples in  $D$  having been drawn **independently** according to the probability law  $P(\mathbf{x}|\omega)$
- The MLE is the value of  $\theta$  that best explains the data and **once we know this value, we know  $P(\mathbf{x}|\omega)$**

$$\hat{\theta} = \arg \max_{\theta} \{P(D | \theta)\}$$



**“Choose the value of  $\theta$  that is the most likely to give rise to the data we observe”**

# MLE

$$D = \{x_1, x_2, \dots, x_N\} \quad \text{N independent observations}$$

$$P(D | \theta) = P(x_1, x_2, \dots, x_N | \theta) = \prod_{k=1}^N P(x_k | \theta)$$



The likelihood of observing  
a particular pattern (random variable)

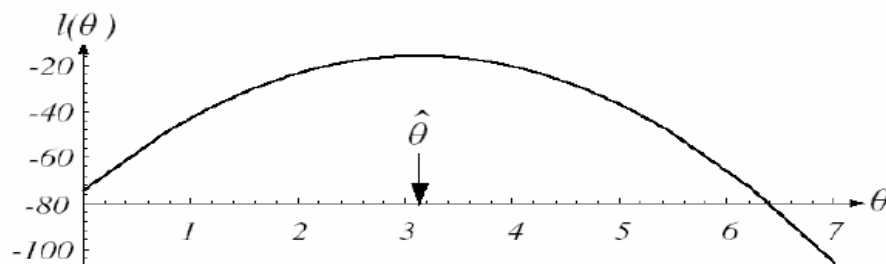
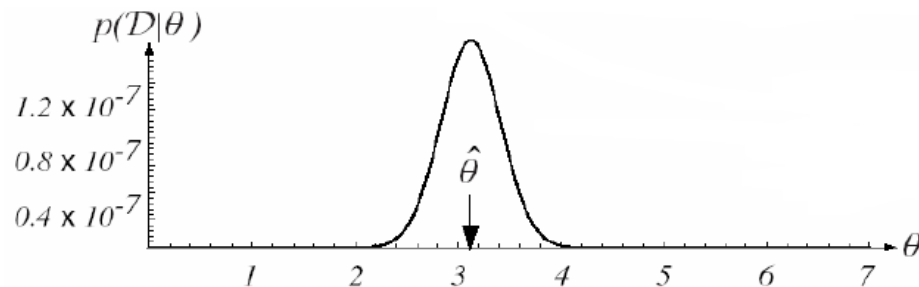
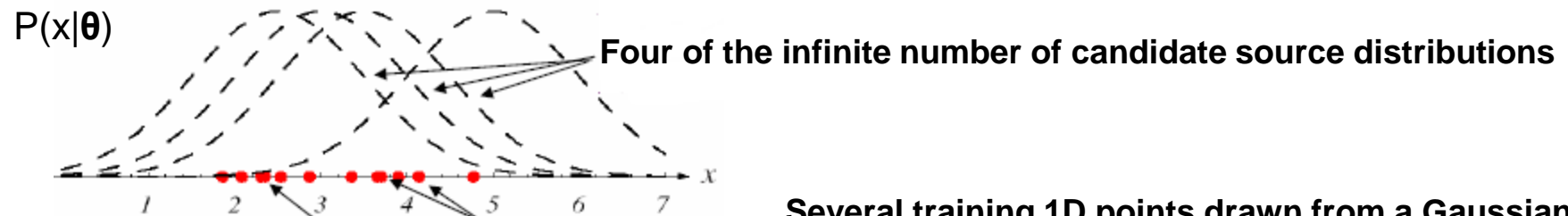
$$\hat{\theta} = \arg \max_{\theta} \{P(Data | \theta)\}$$

“Choose the value of  $\theta$  that is the most likely to give rise to the data we observe”

# MLE contd..

- It is convenient to work with the [log of the likelihood](#)

$$\hat{\theta} = \arg \max_{\theta} \{P(D|\theta)\} = \arg \max_{\theta} \{\log(P(D|\theta))\}$$



# How To Solve For The ML Estimate?

- Let  $\boldsymbol{\theta}$  be the p-component parameter vector  $\boldsymbol{\theta} = [\theta_1, \theta_2, \dots, \theta_p]^T$
- Let this be the gradient operator  $\nabla_{\boldsymbol{\theta}} = \left[ \frac{\partial}{\partial \theta_1}, \frac{\partial}{\partial \theta_2}, \dots, \frac{\partial}{\partial \theta_p} \right]^T$
- We have  $P(D | \boldsymbol{\theta}) = \prod_{k=1}^n P(x_k | \boldsymbol{\theta})$
- We define  $l(\boldsymbol{\theta})$  the log-likelihood of the function

$$l(\boldsymbol{\theta}) = \log(P(D | \boldsymbol{\theta})) = \sum_{k=1}^n \log(P(x_k | \boldsymbol{\theta}))$$

- And

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

- A set of necessary condition for the ML estimate can be obtained from the set of  $p$  equations:

$$\nabla_{\boldsymbol{\theta}} l(\boldsymbol{\theta}) = \mathbf{0}$$



# MLE Example: Univariate Gaussian

- Now assume **neither the mean nor the covariance** matrix are known
- First consider univariate case:

$$\boldsymbol{\theta} = \begin{bmatrix} \theta_1 = \mu \\ \theta_2 = \sigma^2 \end{bmatrix} \quad P(D | \boldsymbol{\theta}) = \prod_{k=1}^n P(x_k | \boldsymbol{\theta}) \quad \log P(x_k | \boldsymbol{\theta}) = -\frac{1}{2} \log(2\pi\theta_2) - \frac{1}{2\theta_2} (x_k - \theta_1)^2$$

- Its derivative is:

$$\nabla_{\boldsymbol{\theta}} l = \nabla_{\boldsymbol{\theta}} \log(P(x_k | \boldsymbol{\theta})) = \begin{bmatrix} \frac{1}{\theta_2} (x_k - \theta_1) \\ -\frac{1}{2\theta_2} + \frac{(x_k - \theta_1)^2}{2\theta_2^2} \end{bmatrix} = 0$$

- Setting it to zero leads to:

$$\sum_{k=1}^n \frac{1}{\theta_2} (x_k - \theta_1) = 0 \quad \text{and} \quad -\sum_{k=1}^n \frac{1}{\theta_2} + \sum_{k=1}^n \frac{(x_k - \theta_1)^2}{\theta_2^2} = 0$$

- Rearranging:

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

**Sample Mean**

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

**Sample Variance**

**ML Estimate**

# MLE Example: Multivariate Gaussian

$$P(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{d/2} |\boldsymbol{\Sigma}|^{1/2}} \exp\left[-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^t \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right]$$

$$l(\theta) = \log(P(D | \theta)) = \sum_{k=1}^n \log(P(\mathbf{x}_k | \theta))$$

Consider **only the mean is unknown**:

$$\log P(\mathbf{x}_k | \boldsymbol{\mu}) = -\frac{1}{2} \log\left((2\pi)^d |\boldsymbol{\Sigma}|\right) - \frac{1}{2}(\mathbf{x}_k - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x}_k - \boldsymbol{\mu})$$

Derivative of log likelihood must be set to 0 to obtain the MLE

$$\nabla_{\boldsymbol{\mu}} \log(P(D | \boldsymbol{\mu})) = \sum_{k=1}^n \boldsymbol{\Sigma}^{-1}(\mathbf{x}_k - \boldsymbol{\mu}) = \mathbf{0}$$

The ML estimate must satisfy:

$$\hat{\boldsymbol{\mu}} = \frac{1}{n} \sum_{k=1}^n \mathbf{x}_k$$

**Sample Mean -> ML Estimate**

# MLE Example: Multivariate Gaussian

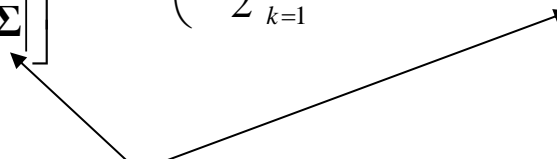
- Neither the mean nor the covariance matrix are known

$$\theta = \begin{bmatrix} \theta_1 = \boldsymbol{\mu} \\ \theta_2 = \boldsymbol{\Sigma} \end{bmatrix} \quad \log P(\mathbf{x}_k | \boldsymbol{\theta}) = -\frac{1}{2} \log \left( (2\pi)^d |\boldsymbol{\Sigma}| \right) - \frac{1}{2} (\mathbf{x}_k - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}_k - \boldsymbol{\mu})$$

- Derivative of log likelihood is:

$$\nabla_{\boldsymbol{\theta}} l = \nabla_{\boldsymbol{\theta}} \log(P(\mathbf{x}_k | \boldsymbol{\theta})) = \begin{bmatrix} \boldsymbol{\Sigma}^{-1} (\mathbf{x}_k - \boldsymbol{\mu}) \\ ? \end{bmatrix}$$

- How to take the gradient of a determinant of a matrix?

$$\begin{aligned} P(\mathbf{x} | \boldsymbol{\Sigma}) &= \prod_{k=1}^n P(\mathbf{x}_k | \boldsymbol{\Sigma}) = \prod_{k=1}^n \left\{ \frac{1}{\sqrt{(2\pi)^d |\boldsymbol{\Sigma}|}} \exp \left( -\frac{1}{2} (\mathbf{x}_k - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}_k - \boldsymbol{\mu}) \right) \right\} \\ &= \frac{1}{\left[ (2\pi)^d |\boldsymbol{\Sigma}| \right]^{n/2}} \exp \left( -\frac{1}{2} \sum_{k=1}^n (\mathbf{x}_k - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}_k - \boldsymbol{\mu}) \right) \end{aligned}$$


Need to take gradient with respect to  $\boldsymbol{\Sigma}$

# MLE Example: Multivariate Gaussian

$$P(\mathbf{x} | \Sigma) = \prod_{k=1}^n P(\mathbf{x}_k | \Sigma) = \prod_{k=1}^n \left\{ \frac{1}{\sqrt{(2\pi)^d |\Sigma|}} \exp \left( -\frac{1}{2} (\mathbf{x}_k - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{x}_k - \boldsymbol{\mu}) \right) \right\}$$

$$= \frac{1}{[(2\pi)^d |\Sigma|]^{n/2}} \exp \left( -\frac{1}{2} \sum_{k=1}^n (\mathbf{x}_k - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{x}_k - \boldsymbol{\mu}) \right)$$

Scalar



$$\mathbf{b}^T \mathbf{B} \mathbf{b} = \text{trace}(\mathbf{b}^T \mathbf{B} \mathbf{b}) = \text{trace}(\mathbf{B} \mathbf{b} \mathbf{b}^T)$$

$$\text{trace}(\mathbf{A} + \mathbf{B}) = \text{trace}(\mathbf{A}) + \text{trace}(\mathbf{B})$$

$$\text{trace}(\mathbf{C}(\mathbf{A} + \mathbf{B})) = \text{trace}(\mathbf{C}\mathbf{A}) + \text{trace}(\mathbf{C}\mathbf{B})$$

$$\sum_{k=1}^n (\mathbf{x}_k - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{x}_k - \boldsymbol{\mu})$$

$$= (\mathbf{x}_1 - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{x}_1 - \boldsymbol{\mu}) + \dots + (\mathbf{x}_n - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{x}_n - \boldsymbol{\mu})$$

$$= \text{trace} \left( \Sigma^{-1} (\mathbf{x}_1 - \boldsymbol{\mu}) (\mathbf{x}_1 - \boldsymbol{\mu})^T \right) + \dots + \text{trace} \left( \Sigma^{-1} (\mathbf{x}_n - \boldsymbol{\mu}) (\mathbf{x}_n - \boldsymbol{\mu})^T \right)$$

$$= \text{trace} \left( \Sigma^{-1} (\mathbf{x}_1 - \boldsymbol{\mu}) (\mathbf{x}_1 - \boldsymbol{\mu})^T + \dots + \Sigma^{-1} (\mathbf{x}_n - \boldsymbol{\mu}) (\mathbf{x}_n - \boldsymbol{\mu})^T \right)$$

$$= \text{trace} \left( \Sigma^{-1} \sum_{k=1}^n (\mathbf{x}_k - \boldsymbol{\mu}) (\mathbf{x}_k - \boldsymbol{\mu})^T \right)$$



$$\mathbf{A} = \frac{1}{n} \sum_{k=1}^n (\mathbf{x}_k - \boldsymbol{\mu}) (\mathbf{x}_k - \boldsymbol{\mu})^T$$

# MLE Example: Multivariate Gaussian

We can now rewrite:

$$\mathbf{A} = \frac{1}{n} \sum_{k=1}^n (\mathbf{x}_k - \boldsymbol{\mu})(\mathbf{x}_k - \boldsymbol{\mu})^T$$

$$\sum_{k=1}^n (\mathbf{x}_k - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}_k - \boldsymbol{\mu}) = \text{trace} \left( \boldsymbol{\Sigma}^{-1} \sum_{k=1}^n (\mathbf{x}_k - \boldsymbol{\mu})(\mathbf{x}_k - \boldsymbol{\mu})^T \right) = \text{trace} (n \boldsymbol{\Sigma}^{-1} \mathbf{A}) = n \cdot \text{trace} (\boldsymbol{\Sigma}^{-1} \mathbf{A})$$

$$\begin{aligned} p(\mathbf{x} | \boldsymbol{\Sigma}) &= \prod_{k=1}^n p(\mathbf{x}_k | \boldsymbol{\Sigma}) = \frac{1}{[(2\pi)^d |\boldsymbol{\Sigma}|]^{n/2}} \exp \left[ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right] \\ &= \frac{1}{[(2\pi)^{dn/2} |\boldsymbol{\Sigma}|]^{n/2}} \exp \left[ -\frac{n}{2} \text{trace}(\boldsymbol{\Sigma}^{-1} \mathbf{A}) \right] \end{aligned}$$

- Now define  $\mathbf{B} = \boldsymbol{\Sigma}^{-1} \mathbf{A}$  and let  $\lambda_1, \lambda_2, \dots, \lambda_d$  be its eigenvalues

$$\begin{aligned} p(\mathbf{x} | \boldsymbol{\Sigma}) &= \frac{1}{[(2\pi)^{dn/2} |\boldsymbol{\Sigma}|]^{n/2}} \exp \left[ -\frac{n}{2} \text{trace}(\mathbf{B}) \right] \\ &= \frac{1}{(2\pi)^{dn/2}} \left[ \frac{|\mathbf{B}|}{|\mathbf{A}|} \right]^{n/2} \exp \left[ -\frac{n}{2} \text{trace}(\mathbf{B}) \right] \\ &= \frac{1}{(2\pi)^{dn/2}} |\mathbf{A}|^{-n/2} \left( \prod_{i=1}^d \lambda_i \right)^{n/2} \exp \left[ -\frac{n}{2} \sum_{i=1}^d \lambda_i \right] \end{aligned}$$

$$\det(\mathbf{AB}) = \det(\mathbf{A}) \det(\mathbf{B})$$

$$\det(\mathbf{A}^{-1}) = \det(\mathbf{A})^{-1}$$

$\text{trace}(\mathbf{A})$  **Sum of eigenvalues**

$\det(\mathbf{A})$  **Product of eigenvalues**

# MLE Example: Multivariate Gaussian

Now consider the log likelihood

$$P(\mathbf{x} | \mathbf{\Sigma}) = (2\pi)^{-dn/2} |\mathbf{A}|^{-n/2} \left( \prod_{i=1}^d \lambda_i \right)^{n/2} \exp \left( -\frac{n}{2} \sum_{i=1}^d \lambda_i \right)$$

$$\log P(\mathbf{x} | \mathbf{\Sigma}) = -\frac{dn}{2} \log(2\pi) - \frac{n}{2} \log |\mathbf{A}| + \frac{n}{2} \sum_{i=1}^d \log(\lambda_i) - \frac{n}{2} \sum_{i=1}^d \lambda_i$$

Remember that  $\mathbf{A}$  is fixed, and taking the gradient with respect to  $\mathbf{\Sigma}$  is now equivalent to taking the derivatives with respect to the eigenvalues of  $\mathbf{B}$

$$\frac{\partial}{\partial \lambda_i} \log P(\mathbf{x} | \mathbf{\Sigma}) = \frac{n}{2\lambda_i} - \frac{n}{2} = 0 \Rightarrow \lambda_i = 1 \text{ for } i = 1, \dots, d$$

Thus  $\mathbf{B}$  must have all eigenvalues of 1 and is equivalent to  $\mathbf{I}$ . This means that the likelihood is maximized if  $\mathbf{\Sigma}^{-1}\mathbf{A}=\mathbf{I}$  or  $\mathbf{\Sigma}=\mathbf{A}$

$$\hat{\mathbf{\Sigma}}_{\text{ML}} = \frac{1}{n} \sum_{k=1}^n (\mathbf{x}_k - \hat{\mathbf{\mu}})(\mathbf{x}_k - \hat{\mathbf{\mu}})^T$$

**MLE - Sample Covariance**

# MLE vs The Bayesian Approach

- **Maximum Likelihood Estimation (MLE)**

- The parameters are assumed to be **FIXED** ( i.e. **NON RANDOM**) but unknown
- The ML seeks the solutions that best explains the data

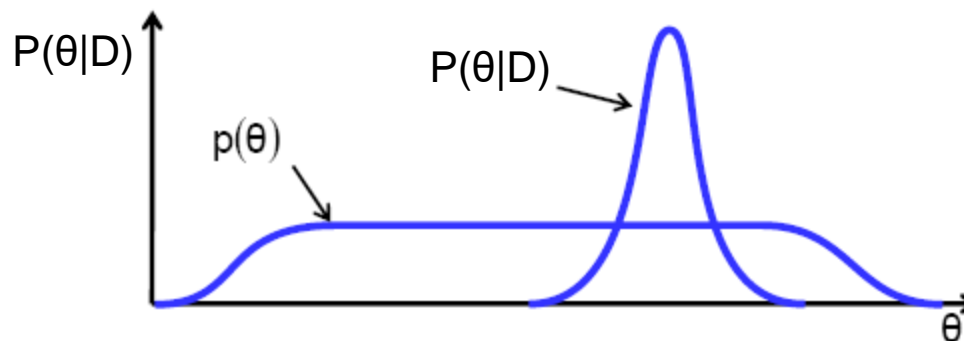
$$\hat{\theta} = \arg \max_{\theta} \{P(Data | \theta)\}$$

- **Bayesian Estimation (BE)**

- The parameters are assumed to be **RANDOM VARIABLES** with some known **PRIORI DISTRIBUTION**
- Bayesian approach aims at estimating the posterior density  $P(\theta | Data)$
- **The MAPE (Maximum A Posteriori Estimate)** of  $\theta$  is the value of  $\theta$  that maximizes the posterior density (i.e. it is the mode of the posterior)

# Bayesian Estimation

- In the Bayesian approach, our uncertainty about the parameters is represented by a PDF
- The parameters are described by a prior density  $P(\theta)$  which indicates which parameters are more likely than others
- We make use of Bayes theorem to find the posterior  $P(\theta|D)$
- Ideally, we want the training data to “sharpen” the posterior  $P(\theta|D)$  or reduce our uncertainty about the parameters





# Bayesian Estimation

- We ideally want to estimate a PDF -  $P(x)$ . Best we can do is estimate it by observing the training data to obtain  $P(x|D)$ . We also assume that it has a known parametric form. So  $P(x|\theta)$  is completely known. But  $\theta$  is random (unlike in MLE) and has its own PDF.

$$P(x|D) = \int P(x, \theta | D) d\theta \quad \leftarrow \text{Using the theorem of total probability}$$

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

Known

Unknown

This integration is typically very hard

$\theta$  is random and has its own PDF

Applying Bayes rule:

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

Posterior

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

# MAP Estimation (MAPE)

- MLE aims at determining the value  $\theta = \theta_{\text{MLE}}$  that maximizes the likelihood

$$\theta_{\text{MLE}} = \arg \max_{\theta} \{P(D | \theta)\}$$

- MAPE assumes that  $\theta$  is not a fixed, but is an RV with a PDF given by  $P(\theta)$  and it aims at maximizing the posterior

$$\theta_{\text{MAP}} = \arg \max_{\theta} \{P(\theta | D)\}$$

$$\text{where } P(\theta | D) = \frac{P(D | \theta)P(\theta)}{P(D)} = \frac{P(D | \theta)P(\theta)}{\int P(D | \theta)P(\theta)d\theta} = \frac{P(\theta) \prod_{k=1}^N P(x_k | \theta)}{\int P(D | \theta)P(\theta)d\theta}$$

Since  $P(D)$  is constant we can also view MAP estimate as

$$\theta_{\text{MAP}} = \arg \max_{\theta} \{P(D | \theta)P(\theta)\}$$

- Thus the MAPE of  $\theta$  is simply the mode of the posterior  $P(\theta | D)$  and MAPE differs from MLE as it determines a value of  $\theta$  which maximizes the posterior instead of the likelihood

# MAPE Example: Univariate Gaussian

- Compute  $P(\theta|D)$  and the desired PDF  $P(x|D)$  where  $p(\mu) = N(\mu_0, \sigma_0^2)$
- $\sigma$  (the uncertainty about  $\mu$ ) is known and fixed.
- Assume the known prior knowledge about the mean can be expressed by a *known* prior density assumed to be normal:

$$p(\mu) = N(\mu_0, \sigma_0^2) \quad \text{Known!}$$

$$P(\mu | D) = \frac{P(D | \mu)P(\mu)}{\int P(D | \mu)P(\mu)d\mu} = \alpha \prod_{k=1}^n P(x_k | \mu)P(\mu)$$

$$P(\mu | D) = \alpha \prod_{k=1}^n \overbrace{\frac{1}{\sqrt{2\pi}\sigma} \exp\left[-\frac{1}{2}\left(\frac{x_k - \mu}{\sigma}\right)^2\right]}^{P(x_k|\mu)} \overbrace{\frac{1}{\sqrt{2\pi}\sigma_0} \exp\left[-\frac{1}{2}\left(\frac{\mu - \mu_0}{\sigma_0}\right)^2\right]}^{P(\mu)}$$

$$= \alpha' \exp\left[-\frac{1}{2}\left(\sum_{k=1}^n \left(\frac{x_k - \mu}{\sigma}\right)^2 + \left(\frac{\mu - \mu_0}{\sigma_0}\right)^2\right)\right]$$

$$= \alpha'' \exp\left[-\frac{1}{2}\left[\left(\frac{n}{\sigma^2} + \frac{1}{\sigma_0^2}\right)\mu^2 - 2\left(\frac{1}{\sigma^2}\sum_{k=1}^n x_k + \frac{\mu_0}{\sigma_0^2}\right)\mu\right]\right] \quad \text{GAUSSIAN!}$$

# MAPE Example: Univariate Gaussian

Gaussian Likelihood && Gaussian Prior  $\rightarrow$  Gaussian Posterior

*compare*

$$P(\mu | D) = \alpha \exp \left[ -\frac{1}{2} \left[ \underbrace{\left( \frac{n}{\sigma^2} + \frac{1}{\sigma_0^2} \right)}_{\text{precision}} \mu^2 - 2 \underbrace{\left( \frac{1}{\sigma^2} \sum_{k=1}^n x_k + \frac{\mu_0}{\sigma_0^2} \right)}_{\text{precision} \times \text{mean}} \mu \right] \right]$$

$$P(\mu | D) = \frac{1}{\sqrt{2\pi\sigma_n^2}} \exp \left[ -\frac{1}{2} \left( \frac{\mu - \mu_n}{\sigma_n} \right)^2 \right] = \alpha \exp \left[ -\frac{1}{2} \left[ \mu^2 \left( \frac{1}{\sigma_n^2} \right) - 2\mu \left( \frac{\mu_n}{\sigma_n^2} \right) + \left( \frac{\mu_n}{\sigma_n} \right)^2 \right] \right]$$

We get  $\frac{1}{\sigma_n^2} = \frac{n}{\sigma^2} + \frac{1}{\sigma_0^2}$  and  $\frac{\mu_n}{\sigma_n^2} = \frac{1}{\sigma^2} \sum_{k=1}^n x_k + \frac{\mu_0}{\sigma_0^2} = \frac{n}{\sigma^2} \hat{\mu}_n + \frac{\mu_0}{\sigma_0^2}$

where  $\hat{\mu}_n = \frac{1}{n} \sum_{k=1}^n x_k$  and is the sample mean

So

$$\sigma_n^2 = \frac{\sigma^2 \sigma_0^2}{n\sigma_0^2 + \sigma^2}$$

$$\mu_n = \left( \frac{n\sigma_0^2}{n\sigma_0^2 + \sigma^2} \right) \hat{\mu}_n + \left( \frac{\sigma^2}{n\sigma_0^2 + \sigma^2} \right) \mu_0$$

# MAPE Example: Univariate Gaussian

$$\mu_n = \left( \frac{n\sigma_0^2}{n\sigma_0^2 + \sigma^2} \right) \hat{\mu}_n + \left( \frac{\sigma^2}{n\sigma_0^2 + \sigma^2} \right) \mu_0$$

$$\sigma_n^2 = \frac{\sigma^2 \sigma_0^2}{n\sigma_0^2 + \sigma^2}$$

Best guess for  $\mu$

Uncertainty for  $\mu$

- $\mu_n$  is a linear combination of  $\hat{\mu}_n$  and  $\mu_0$  (always lies between them)
  - For small samples -> More weights on prior ( $\hat{\mu}_0$ )
  - For large samples -> More weights on observation ( $\hat{\mu}_n$ )
- If  $\sigma_n \neq 0$ , then  $\mu_n \rightarrow \hat{\mu}_n$  as  $n \rightarrow \infty$  (MLE = MAPE)
- If  $\sigma_0 = 0$ , then  $\mu_n = \mu_0$
- If  $\sigma_0 \gg \sigma$  then  $\mu_n = \hat{\mu}_n$  (MLE = MAPE)

# MAPE Example: Univariate Gaussian

- $\mu_n$  represents our best guess for  $\mu$  after observing  $n$  samples, and  $\sigma_n^2$  measures our uncertainty about this guess.

$$\mu_n = \underbrace{\left( \frac{n\sigma_0^2}{n\sigma_0^2 + \sigma^2} \right)}_{ML} \hat{\mu}_n + \underbrace{\left( \frac{\sigma^2}{n\sigma_0^2 + \sigma^2} \right)}_{PRIOR\_KNOWLEDGE} \mu_0$$

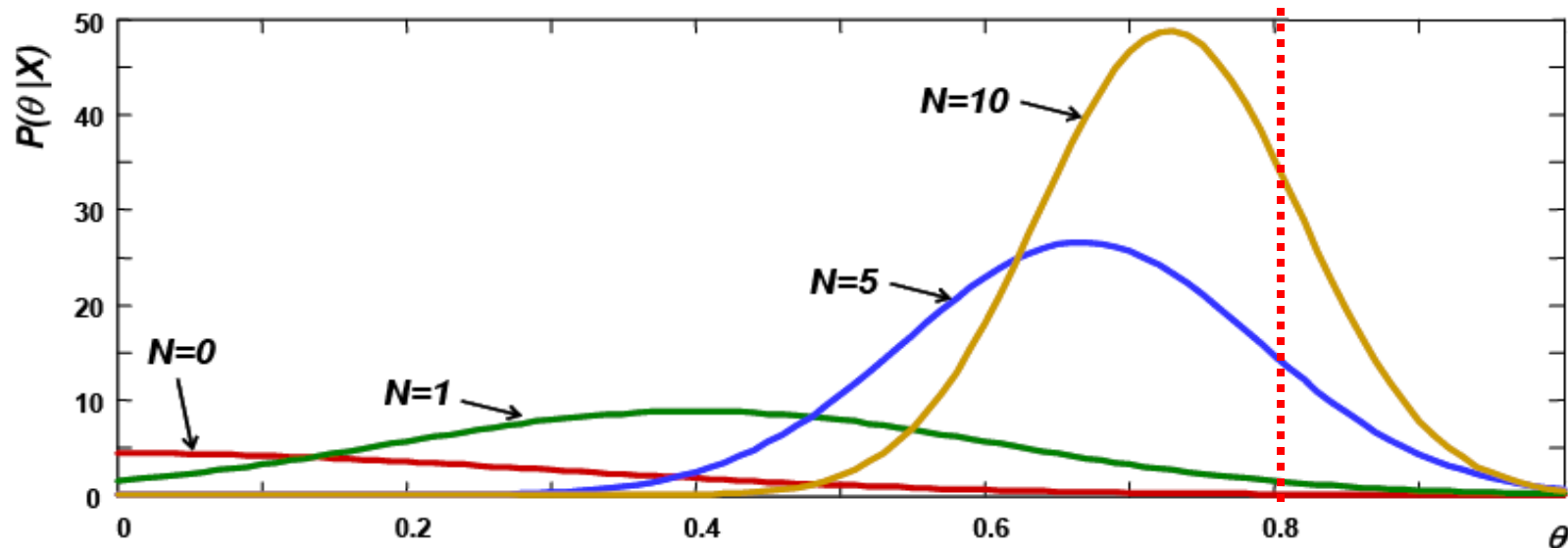
$$\sigma_n^2 = \frac{\sigma^2 \sigma_0^2}{n\sigma_0^2 + \sigma^2}$$

Decreases as  $n \rightarrow \infty$

- Each additional observation decreases our uncertainty and  $P(\mu|D)$  becomes narrower and sharply peaked around the true value of  $\mu$  and becomes a Dirac delta function. This is called **Bayesian Learning**.
- A class of PDF  $P(\theta)$  is said to be **conjugate** to a class of likelihood functions  $P(x|\theta)$  if the resulting posterior distributions  $P(\theta|x)$  are in the **same family as  $P(\theta)$** .
  - The Gaussian family is conjugate to itself if the likelihood function is Gaussian, choosing a Gaussian prior will ensure that the posterior distribution is also Gaussian

# MAPE Example: Gaussians

- Assume that the true mean of the  $P(x)$  is  $N(0.8, 0.09)$ .
  - In reality this is something we cannot know
- We generate a number of examples from this distribution
- We don't know where the mean will be, so we assume a Gaussian prior
  - $P_0(\mu) = N(0, 0.09)$
- As the number of training examples increases, the estimates  $\mu_N$  approaches its true value of 0.8 and the spread decreases



# MAPE Example: Multivariate Gaussian

- Compute  $P(\boldsymbol{\mu}|D)$  and the desired PDF  $P(\mathbf{x}|D)$  where  $P(\mathbf{x} | \boldsymbol{\mu}) = N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$
- Assume the known prior knowledge about the mean can be expressed by a *known* prior density assumed to be normal:

$$P(\boldsymbol{\mu}) = N(\boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0) \quad \text{Known!}$$

$$P(\boldsymbol{\mu} | D) = \frac{P(D | \boldsymbol{\mu})P(\boldsymbol{\mu})}{\int P(D | \boldsymbol{\mu})P(\boldsymbol{\mu})d\boldsymbol{\mu}} = \alpha \prod_{k=1}^n P(\mathbf{x}_k | \boldsymbol{\mu})P(\boldsymbol{\mu})$$

$$P(\boldsymbol{\mu} | D) = \alpha \prod_{k=1}^n P(\mathbf{x}_k | \boldsymbol{\mu})P(\boldsymbol{\mu})$$

$$= \alpha' \exp \left[ -\frac{1}{2} \left( \boldsymbol{\mu}^T (n\boldsymbol{\Sigma}^{-1} + \boldsymbol{\Sigma}_0^{-1}) \boldsymbol{\mu} - 2\boldsymbol{\mu}^T \left( \boldsymbol{\Sigma}^{-1} \sum_{k=1}^n \mathbf{x}_k + \boldsymbol{\Sigma}_0^{-1} \boldsymbol{\mu}_0 \right) \right) \right]$$

Which has a Gaussian Form

$$= \alpha'' \exp \left[ -\frac{1}{2} (\boldsymbol{\mu} - \boldsymbol{\mu}_n)^T \boldsymbol{\Sigma}_n^{-1} (\boldsymbol{\mu} - \boldsymbol{\mu}_n) \right]$$



# MAPE Example: Multivariate Gaussian

- Thus  $P(\boldsymbol{\mu}|D)$  is  $N(\boldsymbol{\mu}_n, \boldsymbol{\Sigma}_n)$ . Equating coefficients, we obtain:

$$\boldsymbol{\Sigma}_n^{-1} = n\boldsymbol{\Sigma}^{-1} + \boldsymbol{\Sigma}_0^{-1} \quad \text{and} \quad \boldsymbol{\Sigma}_n^{-1}\boldsymbol{\mu}_n = n\boldsymbol{\Sigma}^{-1}\hat{\boldsymbol{\mu}}_n + \boldsymbol{\Sigma}_0^{-1}\boldsymbol{\mu}_0$$

where  $\hat{\boldsymbol{\mu}}_n = \frac{1}{n} \sum_{k=1}^n \mathbf{x}_k$  is the sample mean

After some manipulation (which we don't prove) we get:

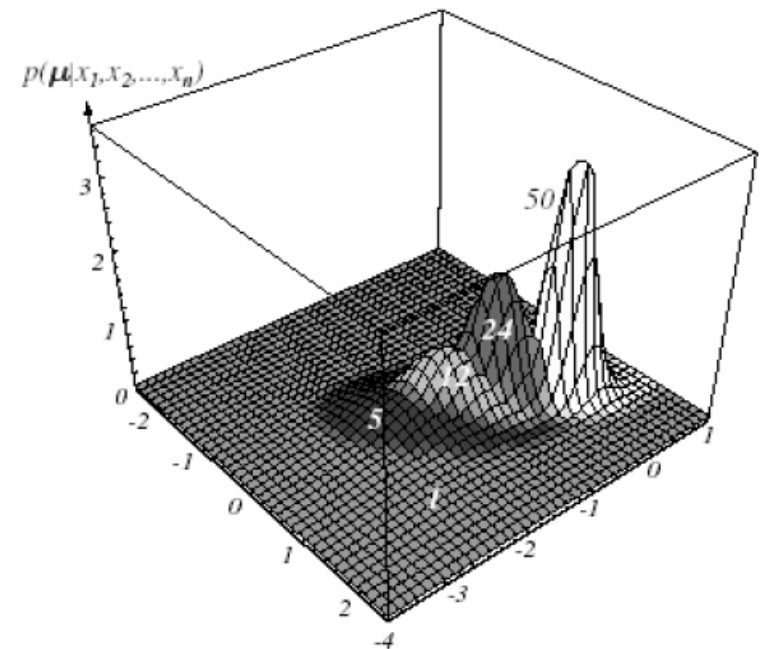
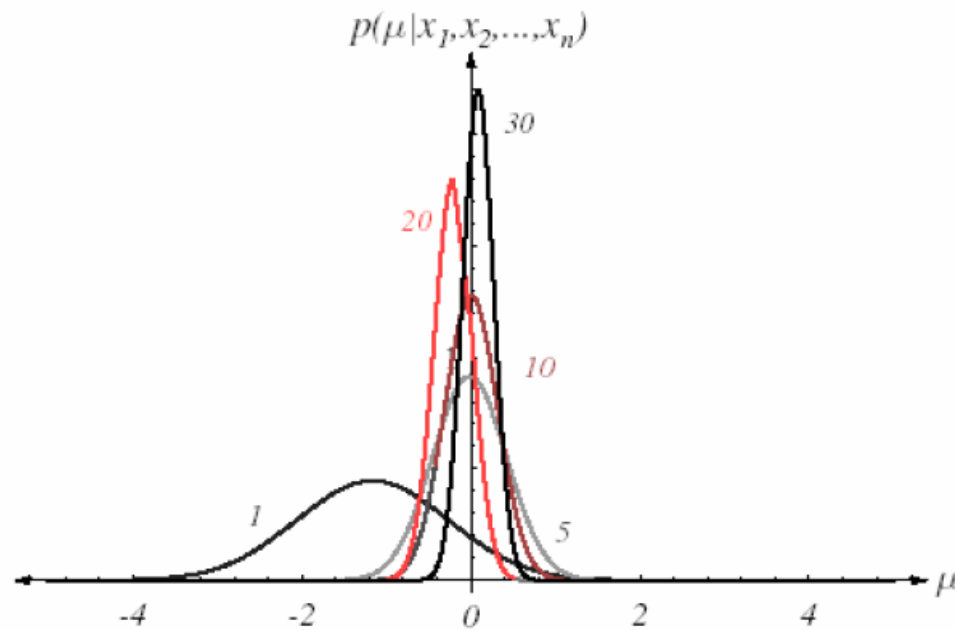
$$\boldsymbol{\mu}_n = \boldsymbol{\Sigma}_0 \left( \boldsymbol{\Sigma}_0 + \frac{1}{n} \boldsymbol{\Sigma} \right)^{-1} \hat{\boldsymbol{\mu}}_n + \frac{1}{n} \boldsymbol{\Sigma} \left( \boldsymbol{\Sigma}_0 + \frac{1}{n} \boldsymbol{\Sigma} \right)^{-1} \boldsymbol{\mu}_0$$

$$\boldsymbol{\Sigma}_n = \boldsymbol{\Sigma}_0 \left( \boldsymbol{\Sigma}_0 + \frac{1}{n} \boldsymbol{\Sigma} \right)^{-1} \frac{1}{n} \boldsymbol{\Sigma}$$

- Finally, it can be shown that

$$P(\mathbf{x} | D) = N(\boldsymbol{\mu}_n, \boldsymbol{\Sigma} + \boldsymbol{\Sigma}_n)$$

# MAPE Example: Gaussians



**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.

# Recap: ML vs Bayesian Estimation

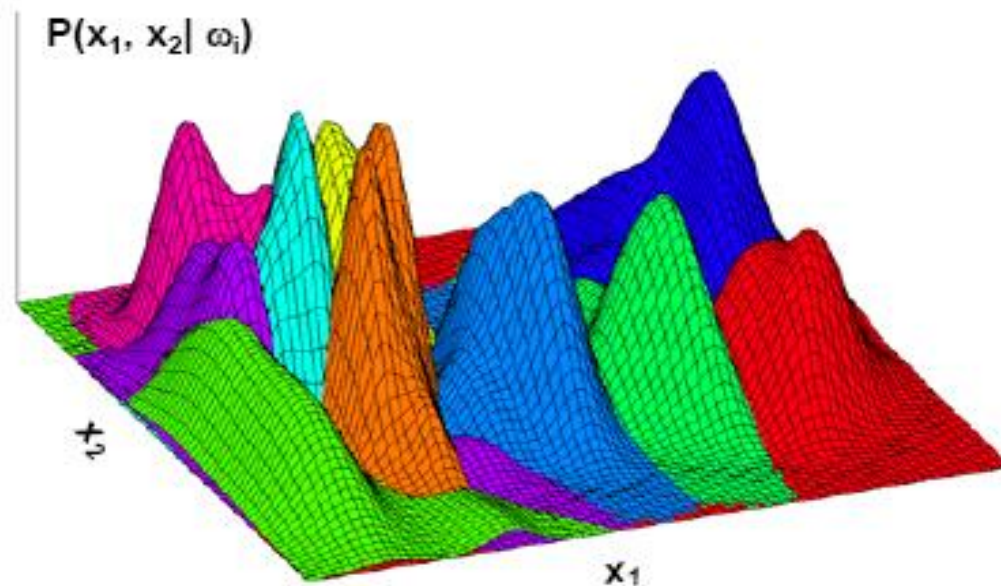
- BE assumes that the parameters come from a distribution with known priors
- BE provides a distribution for  $\theta$  rather than point values. BE provides more information, but is in general much more difficult to compute (integration)
- For most times, if the assumptions are correct, MLE gives good enough results

# Recap: ML vs Bayesian Estimation

- Amount of Training data
  - The two methods are equivalent assuming infinite training samples.
  - They differ for smaller training data
- Computational Complexity
  - ML uses **differential calculus** or gradient search
  - Bayesian estimation needs complex multidimensional **integration techniques**
- Solution Complexity
  - ML solution is easy to interpret
  - A Bayes Estimation solution **might not be** of the parametric form assumed
- Prior distribution
  - If the prior  $P(\theta)$  is **uniform (flat)**, BE solutions are equivalent to ML solutions

# Non Parametric Density Estimation

- Problems with parametric techniques
  - Most forms of distributions are unimodal
  - Most of the time we assumed that features are independent
- In non-parametric approach we don't need these assumptions
  - We have to attempt to estimate the density directly from the data without making any parametric assumptions about the true density

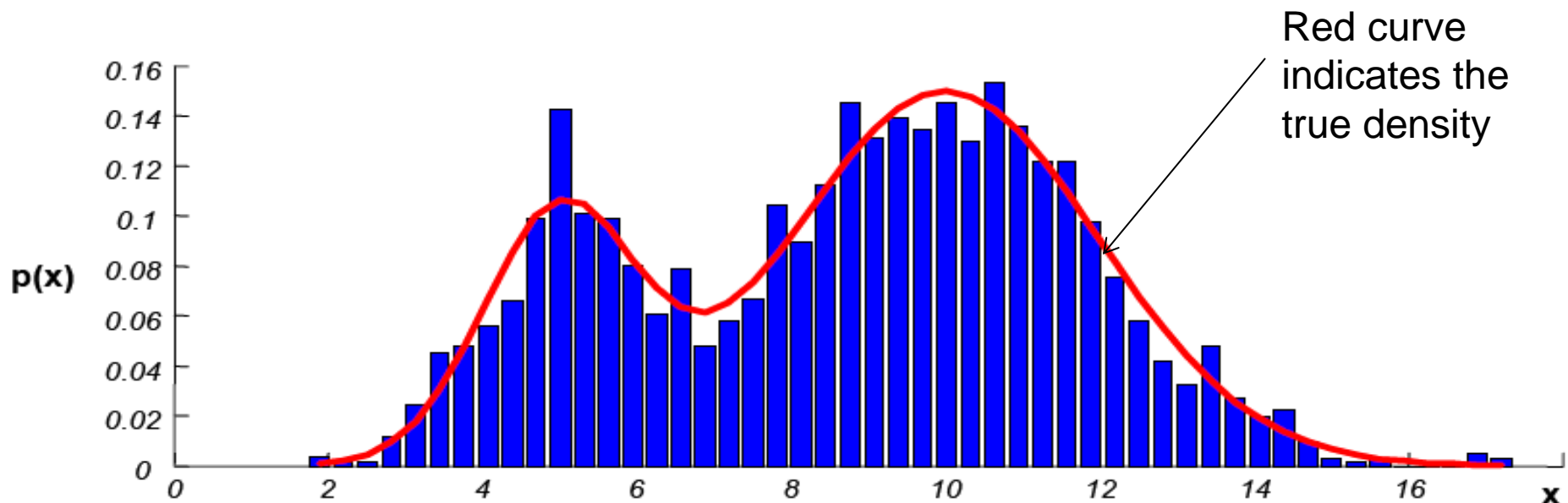


# The Histogram

- The simplest form of non-parametric density estimation
  - Divide the sample space into a number of bins and approximate the density at the center of each bin using the fraction of training samples that fall into the corresponding bin.

$$P_H(x) = \frac{1}{n} \frac{[\text{number of samples in the same bin as } x]}{[\text{width of bin containing } x]}$$

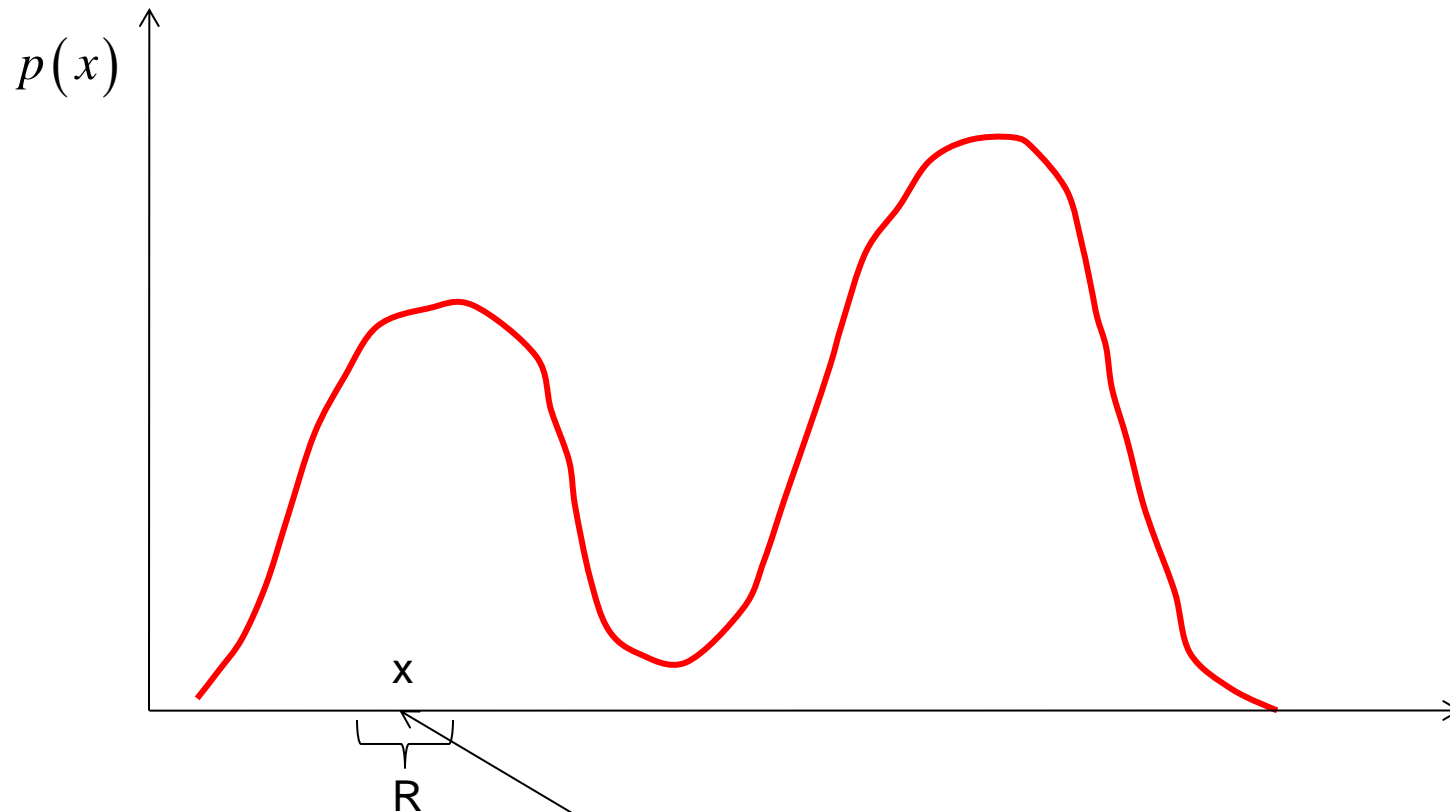
- The histogram requires two tuning options to be defined:
  - Bin width
  - Starting position of the first bin



# The Histogram

- The histogram is simple, but has flaws
  - The bin location can introduce discontinuity artifacts that are not due to the true underlying density
  - The histogram approach becomes impractical when the number of dimensions increases. The number of bins needed grows exponentially.
- Histograms are useful for visualizations in one or two dimensions

# Non Parametric Density Estimation



Let's estimate distribution at this point  $x$

Consider a region  $R$  around the point  $x \rightarrow$  analogous to the bin in a histogram



# Non Parametric Density Estimation

- From basic probability, we know that
  - The probability that a vector  $x$ , drawn from some distribution  $p(x)$ , will fall in a given region  $R$  of the sample space is:

$$P = \int_R p(x') dx' \quad (\text{i.e. a space averaged value})$$

- Suppose now that  $n$  vectors  $\{x_1, x_2, \dots, x_n\}$  are drawn from the distribution. The probability that  $k$  of these  $n$  vectors fall in  $R$  is given by the binomial distribution
- From the properties of the binomial pmf, we have that

$$P_k = \binom{n}{k} P^k (1-P)^{n-k}$$

$$E[k] = nP \quad \text{Var}[k] = nP(1-P)$$

$$E[k/n] = P \quad \text{Var}[k/n] = P(1-P)/n$$

# Non Parametric Density Estimation

$$E[k / n] = P \quad \text{Var}[k / n] = P(1 - P) / n$$

- As  $n \rightarrow \infty$ , the distribution becomes sharper and hence a good estimate of the probability  $P$  can be obtained from the mean fraction of the points that fall within  $R$

$$P = \int_R p(x') dx' \cong \frac{k}{n}$$

- If we now assume that  $R$  is so small that  $p(x)$  is almost constant within it, then

$$\int_R p(x') dx' \cong p(x)V$$

*where  $V$  is the volume enclosed by the region  $R$*

For this one dimensional example,  $V$  will just be the length of the region

# Non Parametric Density Estimation

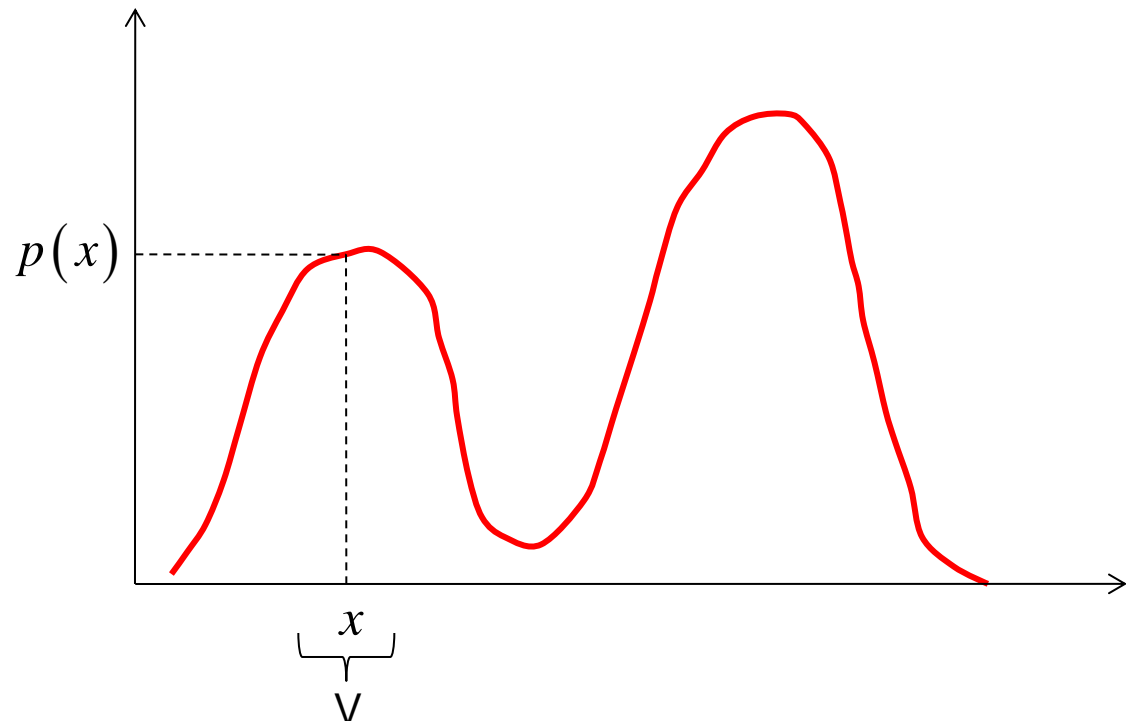
$$P = \int_R p(x') dx' \cong \frac{k}{n}$$

$$\int_R p(x') dx' \cong p(x)V$$

– Merging the two results we obtain:

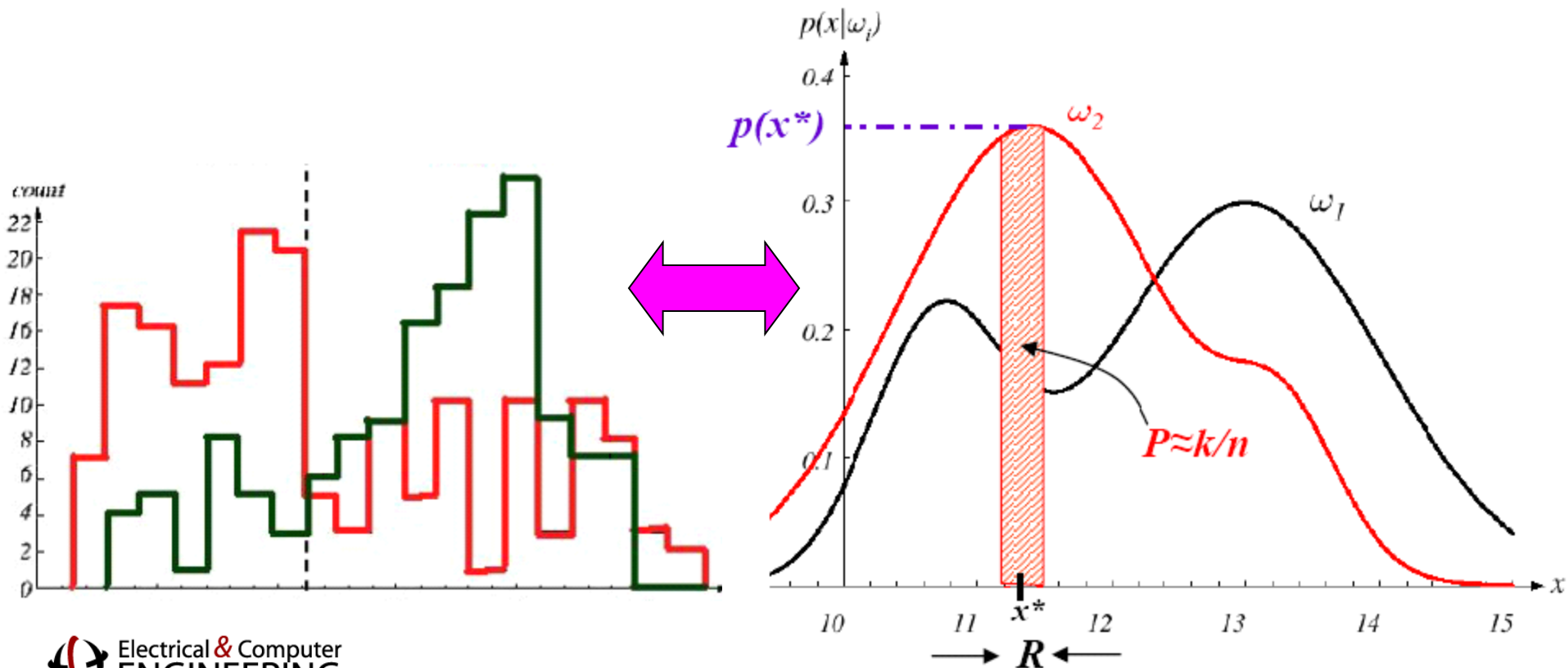
$$p(x) = \frac{k}{nV}$$

What we get is a  
space averaged  
value depending  
on  $V$



# Non Parametric Density Estimation

- Example:
  - If we observe a large number of fish, count those whose length fall within the range defined by  $R$ ,  $k/n$  can be used as an estimate of  $P$  when  $n \rightarrow \infty$



# What are the Assumptions Here?

$$p(x) \cong \frac{k}{nV}$$

1.  $V$  is so small that  $p(x)$  is constant in  $R$ .
  - But if  $V$  is really small, then it may not contain any samples! Then  $(k/n)=0$
2.  $V$  is sufficiently large so that we get enough number of samples ( $k$ ) such that
 
$$P \cong \frac{k}{n}$$
  - Two contradictory assumptions! So we'll have to accept
    - A certain amount of variance in  $(k/n)$
    - A space averaged value for  $p(x)$  (since  $V$  will not be very small)

# Approaches

$$p(x) \cong \frac{k}{nV}$$

- Two basic approaches can be used:
  - Choose a fixed value of  $V$  around the point  $x$  and determine  $k$  from the data. Then determine  $p(x)$  from above. This is called Kernel Density Estimation (KDE).
  - Choose a fixed value of  $k$  and determine the corresponding volume  $V$  around  $x$  that includes  $k$  samples. This is called  $k$  Nearest Neighbor (kNN) approach.

# Approaches

$$p(x) \cong \frac{k}{nV}$$

- Two basic approaches can be used:
  - Choose a fixed value of  $V$  around the point  $x$  and determine  $k$  from the data. Then determine  $p(x)$  from above. This is called Kernel Density Estimation (KDE).
  - Choose a fixed value of  $k$  and determine the corresponding volume  $V$  around  $x$  that includes  $k$  samples. This is called  $k$  Nearest Neighbor (kNN) approach.

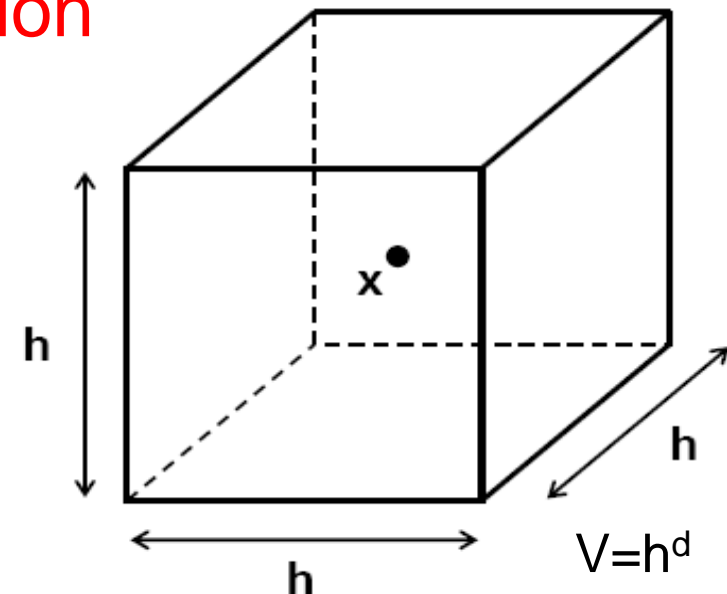
# Parzen Windows

- Suppose that the region  $R$  that encloses  $k$  examples is a  $d$ -dimensional hypercube with sides of length  $h$  centered at the estimation point  $x$ .
- To find the number of instances that fall within this region we define a **Kernel Function**

$$\varphi(u_j) = \begin{cases} 1 & \text{if } |u_j| < 1/2 \quad \forall j = 1, \dots, d \\ 0 & \text{otherwise} \end{cases}$$



Emanuel Parzen



$$\varphi\left(\frac{(x - x_i)}{h}\right) = \begin{cases} 1 & \text{if } x_i \text{ inside hypercube of side } h \text{ centered on } x \\ 0 & \text{otherwise} \end{cases}$$



# Parzen Windows

- The total number of points inside the hypercube is then

$$k = \sum_{i=1}^n \varphi\left(\frac{(x - x_i)}{h}\right)$$

- Substitute back into the expression for the density estimate:

$$p_{KDE}(x) = \frac{1}{nh^d} \sum_{i=1}^n \varphi\left(\frac{(x - x_i)}{h}\right)$$

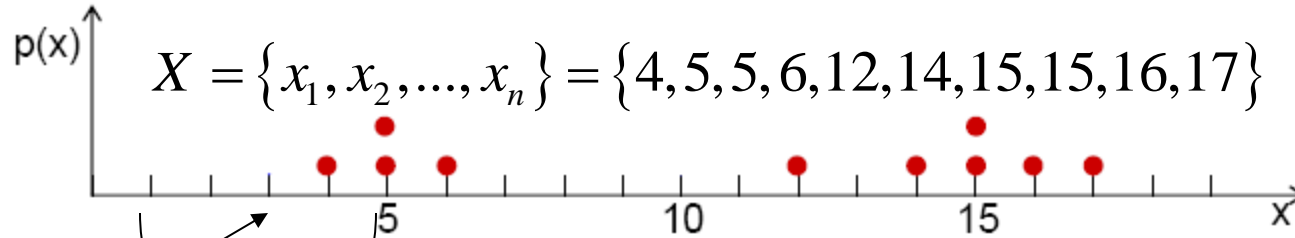
Number of  
samples falling  
into R

Volume covered by  
the window function

Width of the window function

# Parzen Windows Example

- Given the dataset below, use Parzen windows of bandwidth  $h=4$  to estimate the density  $p(x)$



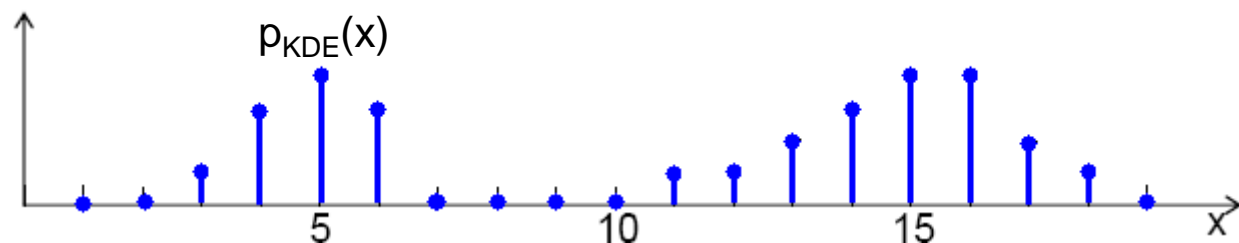
$$p_{KDE}(x=3) = \frac{1}{nh^d} \sum_{n=1}^n \varphi\left(\frac{x-x_n}{h}\right) = \frac{1}{10.4^1} \left[ \underbrace{\varphi\left(\frac{3-4}{4}\right)}_{-1/4} + \underbrace{\varphi\left(\frac{3-5}{4}\right)}_{-1/2} + \underbrace{\varphi\left(\frac{3-5}{4}\right)}_{-1/2} + \underbrace{\varphi\left(\frac{3-6}{4}\right)}_{-1} + \dots + \underbrace{\varphi\left(\frac{3-17}{4}\right)}_{-14/4} \right]$$

$$p_{KDE}(x=5) = \frac{1}{10.4^1} [1+1+1+1+0+0+0+0+0+0] = 0.10$$

$$= \frac{1}{10.4^1} [1+0+0+0+\dots+0] = 0.025$$

$$p_{KDE}(x=10) = \frac{1}{10.4^1} [0+0+0+0+0+0+0+0+0+0] = 0$$

$$p_{KDE}(x=15) = \frac{1}{10.4^1} [0+0+0+0+0+0+1+1+1+1+0] = 0.10$$



# Smooth Kernels

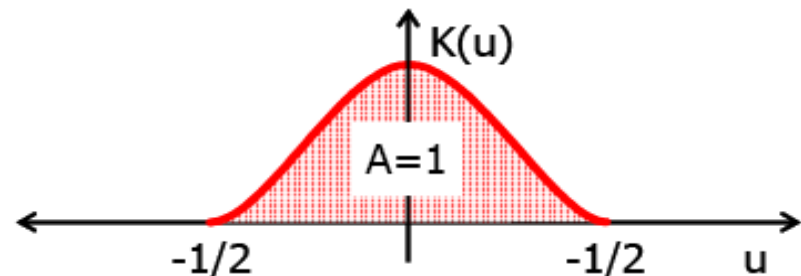
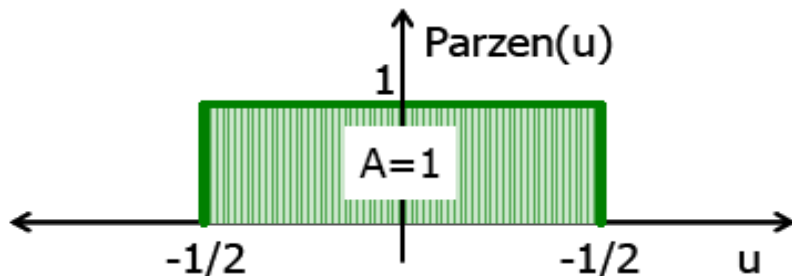
- The Parzen window has several issues
  - Errors due to discontinuities at window edges (same problem as with histograms)
  - Weights all points equally regardless of their distance to the estimation point  $x$
- We can do better by generalizing to a smooth kernel function  $\phi$  such that

$$\int_{\mathbb{R}^d} \phi(x) dx = 1 \quad \text{and} \quad \phi(x) \geq 0$$

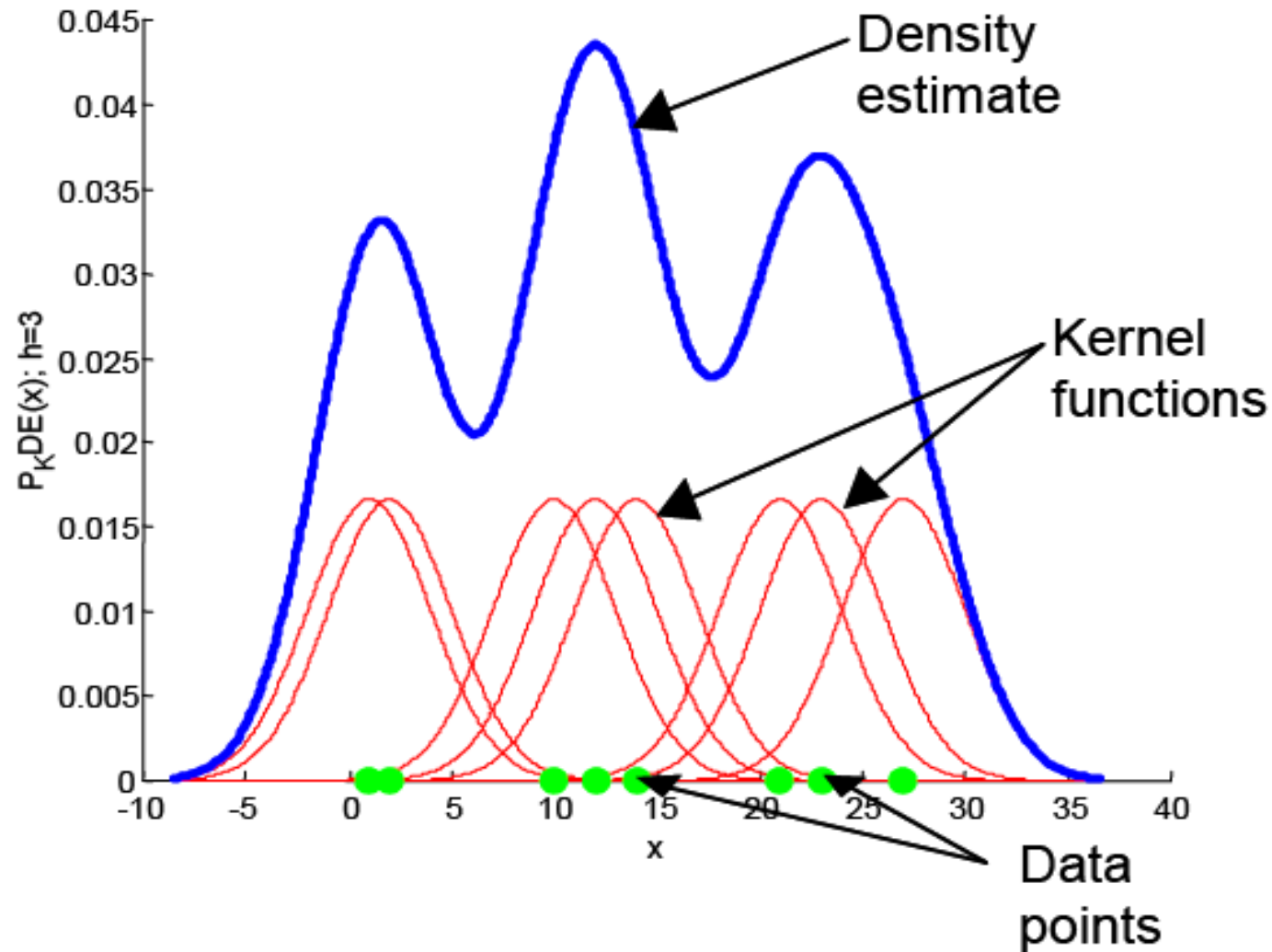
- Usually but not necessarily,  $\phi(u)$  will be a radially symmetric and unimodal PDF, such as the multivariate Gaussian

$$p_{KDE}(x) = \frac{1}{nh^d} \sum_{n=1}^n \phi\left(\frac{(x - x_n)}{h}\right) = \frac{1}{nh^d} \sum_{n=1}^n \frac{1}{(2\pi)^{d/2}} \exp\left\{-\frac{\|x - x_n\|^2}{2h^2}\right\}$$

- The expression of the density estimate remains the same as with Parzen windows

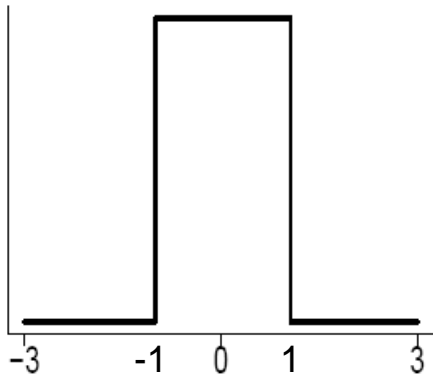


# Kernel Density Estimation



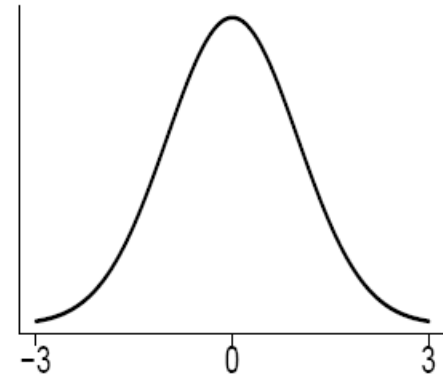
# Commonly Used Kernels

**Boxcar:**  $K(x) = \frac{1}{2} I(x)$

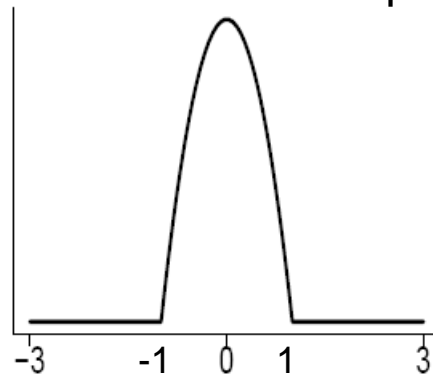


$$I(x) = \begin{cases} 1 & \text{if } |x| \leq 1 \\ 0 & \text{else} \end{cases}$$

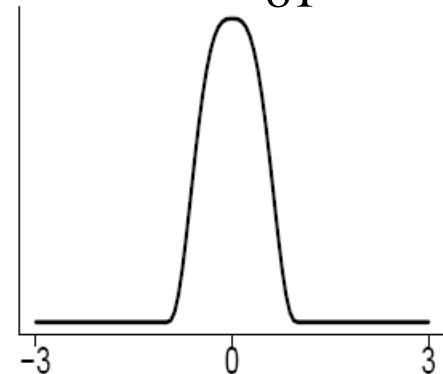
**Gaussian:**  $K(x) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{x^2}{2}\right)$



**Epanechnikov:**  $K(x) = \frac{3}{4} (1 - x^2) I(x)$

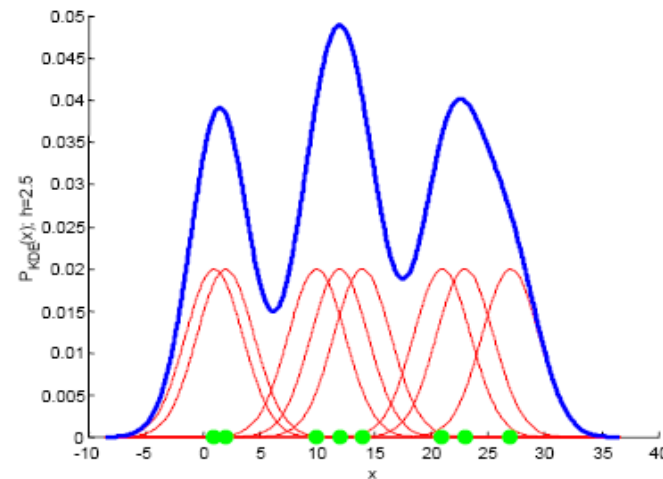
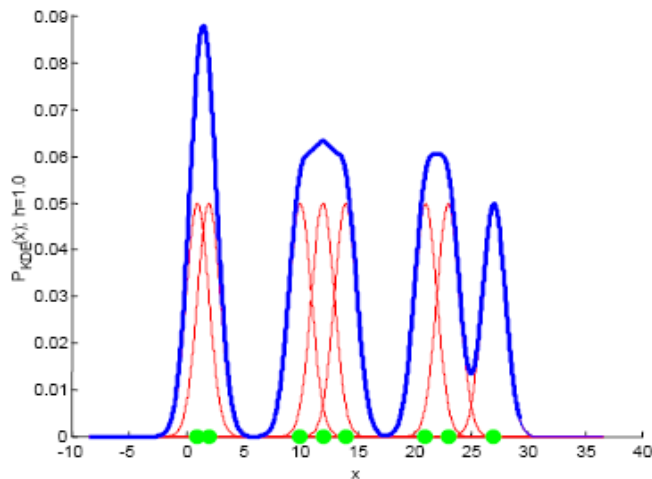


**Tricube:**  $K(x) = \frac{70}{81} (1 - |x|^3)^3 I(x)$

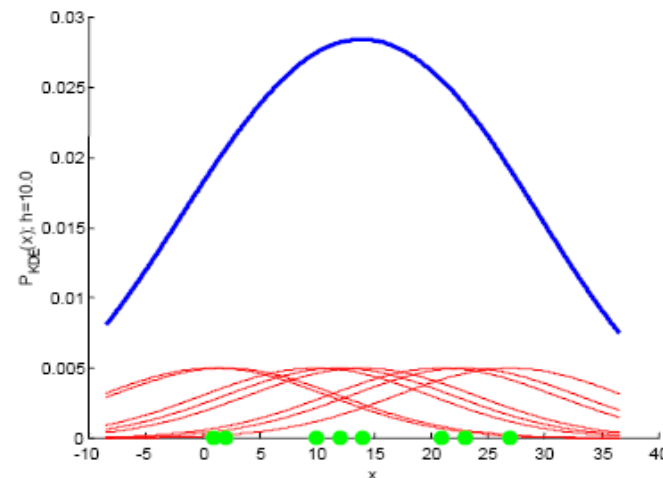
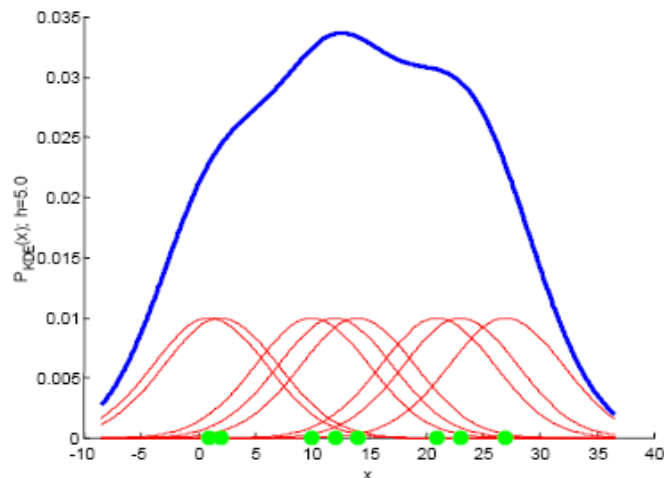


# Importance of Bandwidth

- A good compromise is needed for the bandwidth parameter
  - A large bandwidth will over smooth the density and erode the true structure of the data
  - A small bandwidth will yield a spiky density estimate

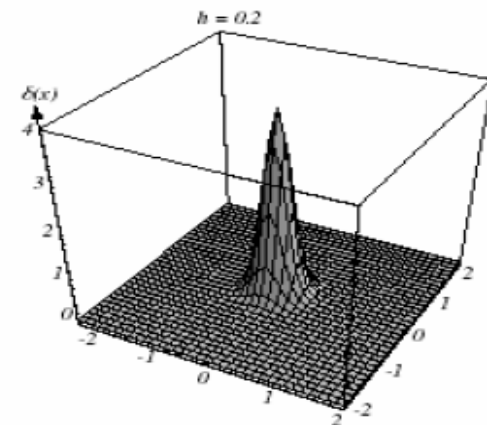
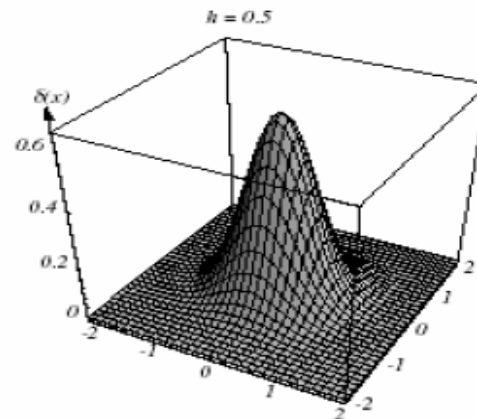
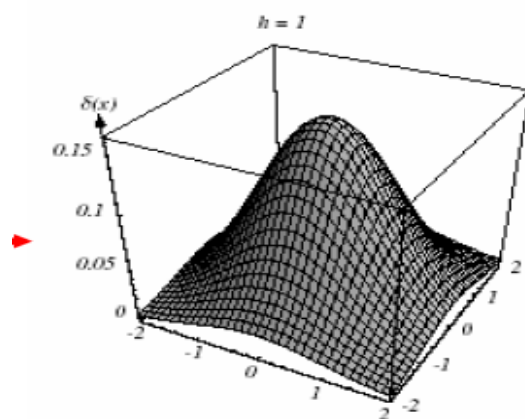


Small BW



Large BW

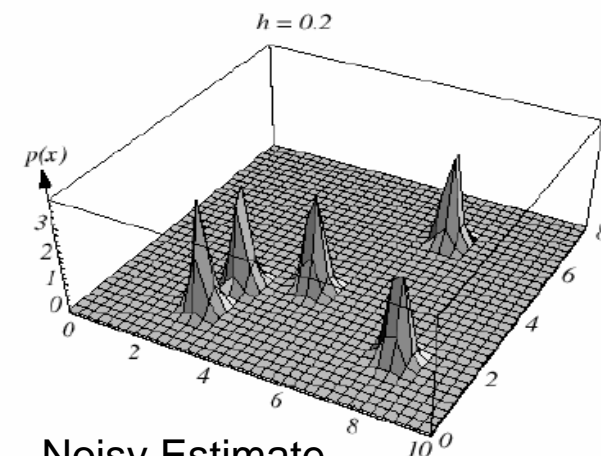
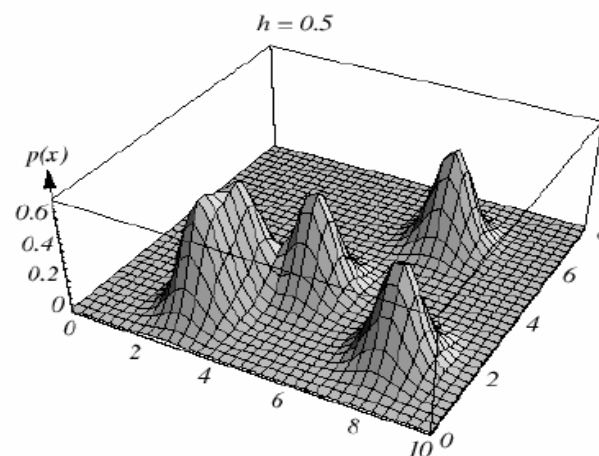
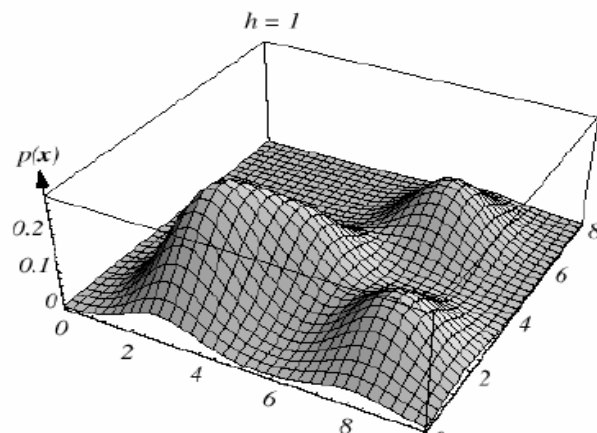
# Importance of Bandwidth



2-D

Large  $V_n \rightarrow$  too little resolution

Small  $V_n \rightarrow$  too much statistical variability



Out Of Focus Estimate

Noisy Estimate

# Choosing the Bandwidth

- We would like to get a value of the smoothing parameter that minimizes the error between the estimated density and the true density
  - A natural measure is the MSE at the estimation point  $x$ , given by:

$$MSE(p_{KDE}) = E\left[(p_{KDE}(x) - p(x))^2\right] = \left\{E[p_{KDE}(x) - p(x)]\right\}^2 + \text{var}(p_{KDE}(x))$$

$$\therefore \text{var}(p_{KDE}(x)) = E\left[(p_{KDE}(x) - p(x))^2\right] - \left\{E[p_{KDE}(x) - p(x)]\right\}^2$$



# Bias-Variance Tradeoff

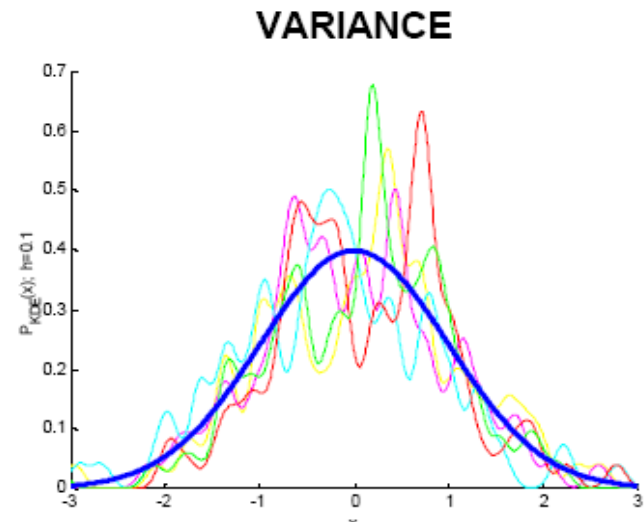
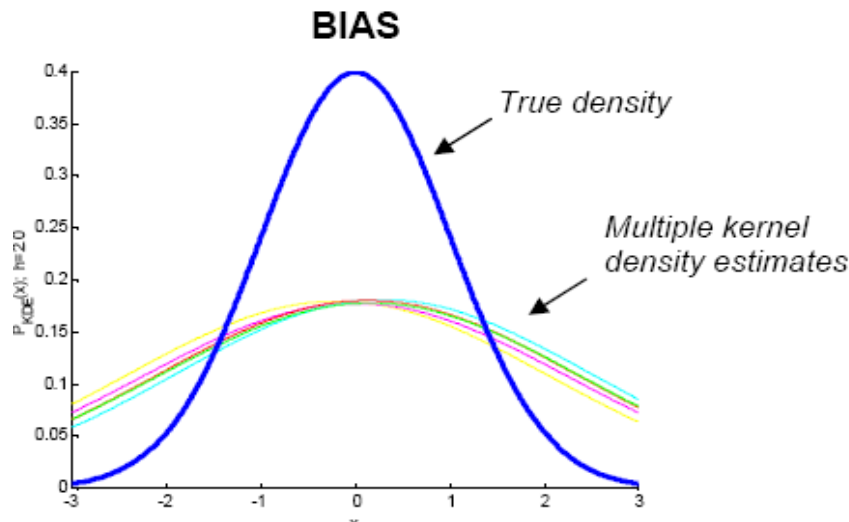
$$MSE \left( p_{KDE} \right) = \left\{ E \left[ p_{KDE} (x) - p(x) \right] \right\}^2 + \text{var} \left( p_{KDE} (x) \right)$$

$$MSE \left( p_{KDE} \right) = (\text{Bias})^2 + \text{Variance}$$

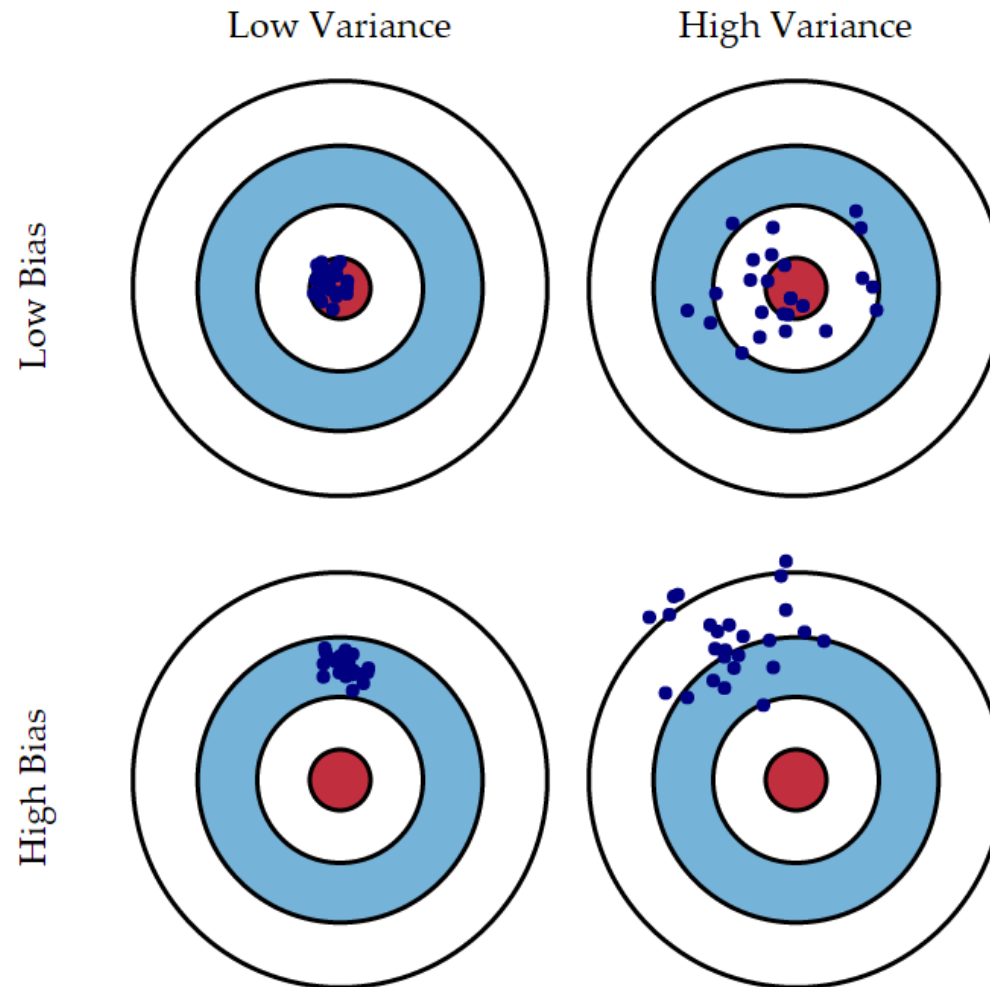
- Bias can be reduced at the expense of variance and vice versa
  - The bias of an estimate is the systematic error incurred in the estimation
  - The variance of an estimate is the random error incurred in the estimation

# Bias-Variance Tradeoff

- The bias-variance dilemma of the bandwidth means that:
  - A large bandwidth will reduce the differences among the estimates of  $p_{\text{KDE}}(x)$  for different datasets (i.e. the variance) but will increase the bias of  $p_{\text{KDE}}(x)$  with respect to the true density  $p(x)$
  - A small bandwidth will reduce the bias of  $p_{\text{KDE}}(x)$  at the expense of a large variance in the estimates of  $p_{\text{KDE}}(x)$

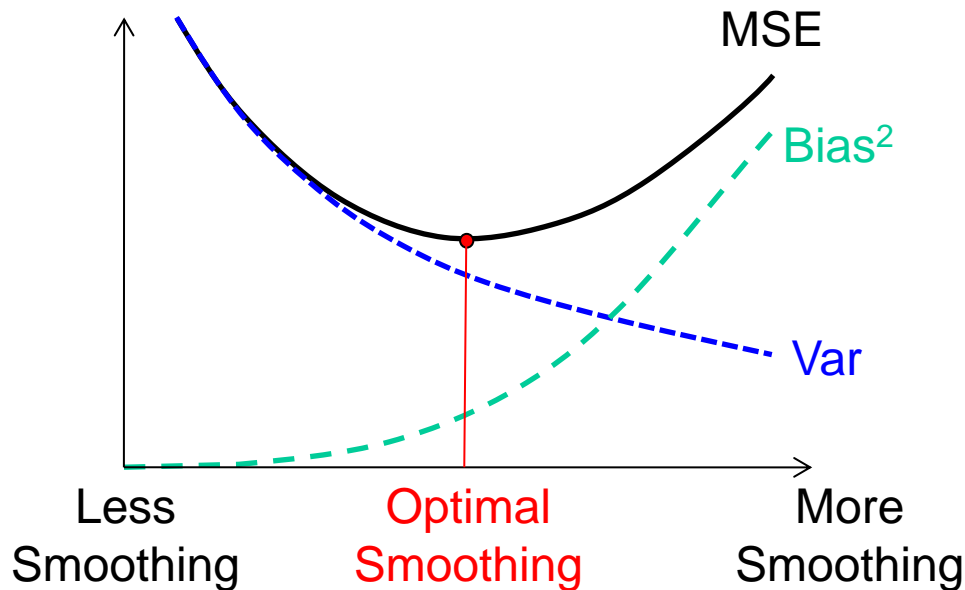


# Bias-Variance Tradeoff



# Bias-Variance Tradeoff

The bias increases and the variance decreases with the amount of smoothing



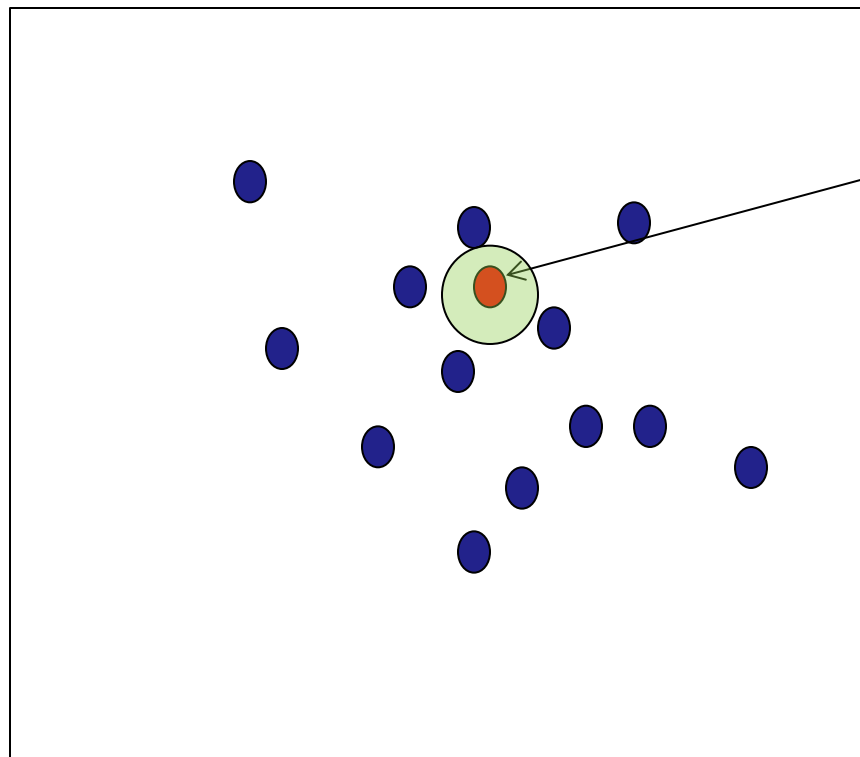
# Density Estimation

$$p(x) \cong \frac{k}{nV}$$

- Two basic approaches can be used:
  - Choose a fixed value of  $V$  around the point  $x$  and determine  $k$  from the data. Then determine  $p(x)$  from above. This is called Kernel Density Estimation (KDE).
  - Choose a fixed value of  $k$  and determine the corresponding volume  $V$  around  $x$  that includes  $k$  samples. This is called  $k$  Nearest Neighbor (kNN) approach.

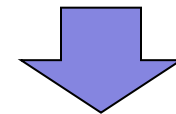
# kNN Density Estimation

- In the kNN method we grow the volume surrounding the estimation point  $x$  until it encloses a total of  $k$  data points.



Estimate the  
density at this  
point

Volume  $V$  enclose  
 $k=4$  samples



$$p(x) \cong \frac{k / n}{V}$$

# kNN Estimate of Bayes Probabilities

- kNN can be used to estimate the prior, posterior, and likelihoods.
  - Assume that we have  $n$  total samples, and  $c$  total classes, and  $n_i$  samples inside each class such that  $n = n_1 + n_2 + \dots + n_c$
  - To classify an unknown data point  $\mathbf{x}$ , we draw a hyper-sphere  $V$  around  $\mathbf{x}$ . Assume that this volume contains a total of  $k$  samples,  $k_i$  of which come from class  $\omega_i$ , then,

$$p(\mathbf{x} | \omega_i) \cong \frac{k_i}{n_i V} \quad p(\mathbf{x}) \cong \frac{k}{nV} \quad p(\omega_i) \cong \frac{n_i}{n}$$

- Put them all together using Bayes formula

$$P(\omega_i | \mathbf{x}) = \frac{p(\mathbf{x} | \omega_i) p(\omega_i)}{p(\mathbf{x})} \cong \frac{\frac{k_i}{n_i V} \cdot \frac{n_i}{n}}{\frac{k}{nV}} = \frac{k_i}{k}$$

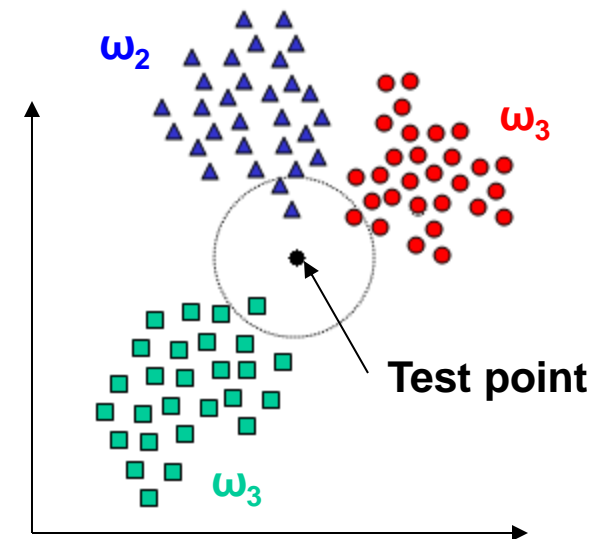
- For  $k = 1$ , this is called the Nearest Neighbor rule

# kNN Classification Rule

- kNN Rule is a simple and intuitive method that classifies unlabeled instances based on their proximity to examples in the training set.

$$p(\omega_i | \mathbf{x}) = \frac{k_i}{k}$$

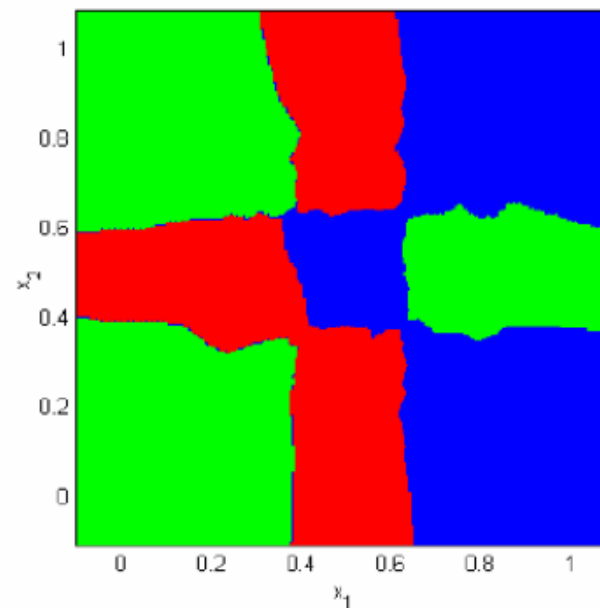
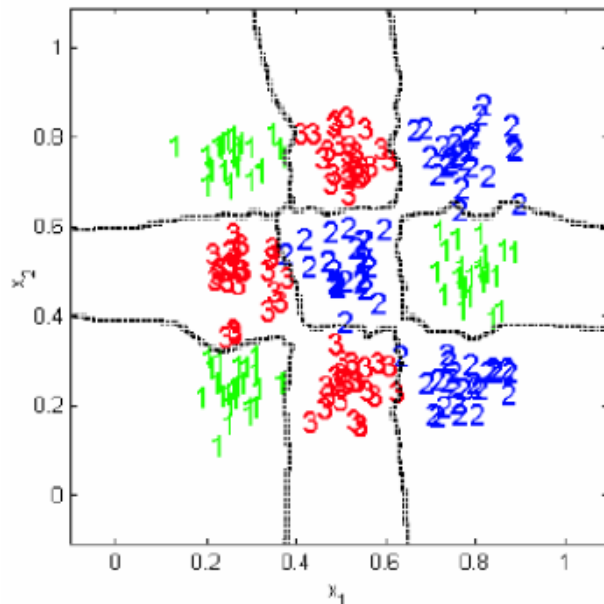
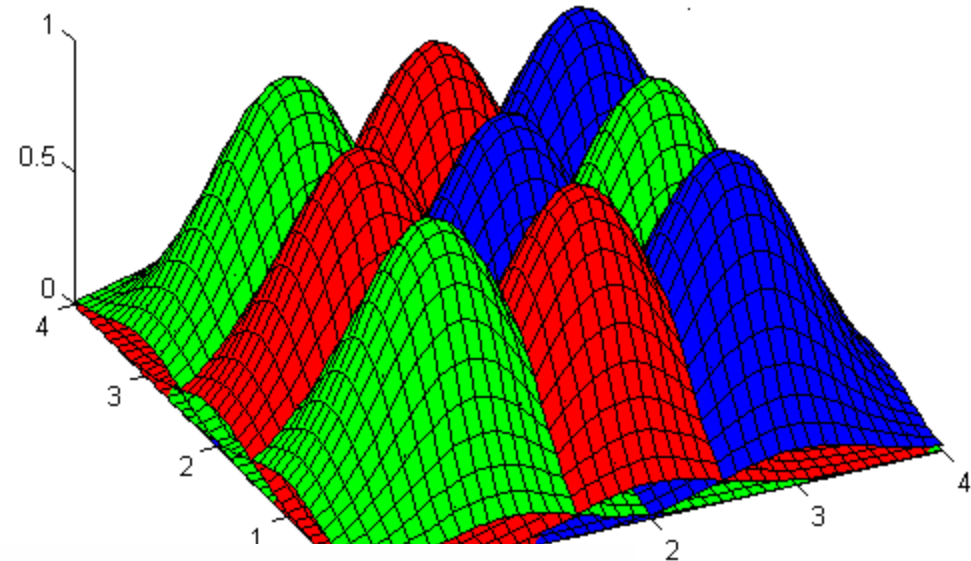
- Given an unlabeled test sample, find the k “closest” labeled examples in the training set and **assign it to the class that appears the most frequently** in this k-subset.
- kNN rule requires:
  - An integer k
  - A set of labeled training instances
  - A metric to measure proximity.
- In this example:
  - Use k = 4
  - Euclidian distance
  - The test point is assigned to class 2





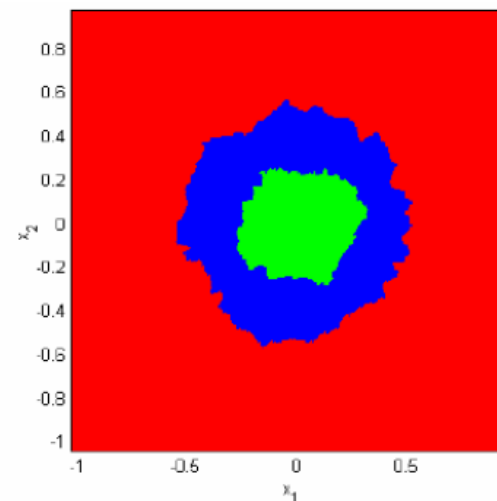
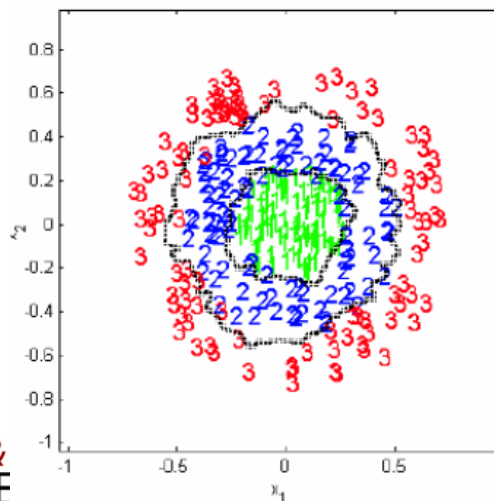
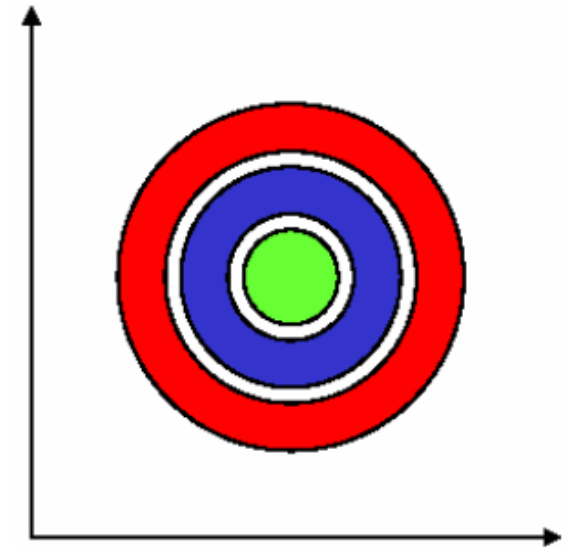
# kNN Classification Rule

- 3-class problem unimodal non-linearly separable.
- Use  $k=5$  and euclidian distance
- The resulting boundaries and decision regions are shown below



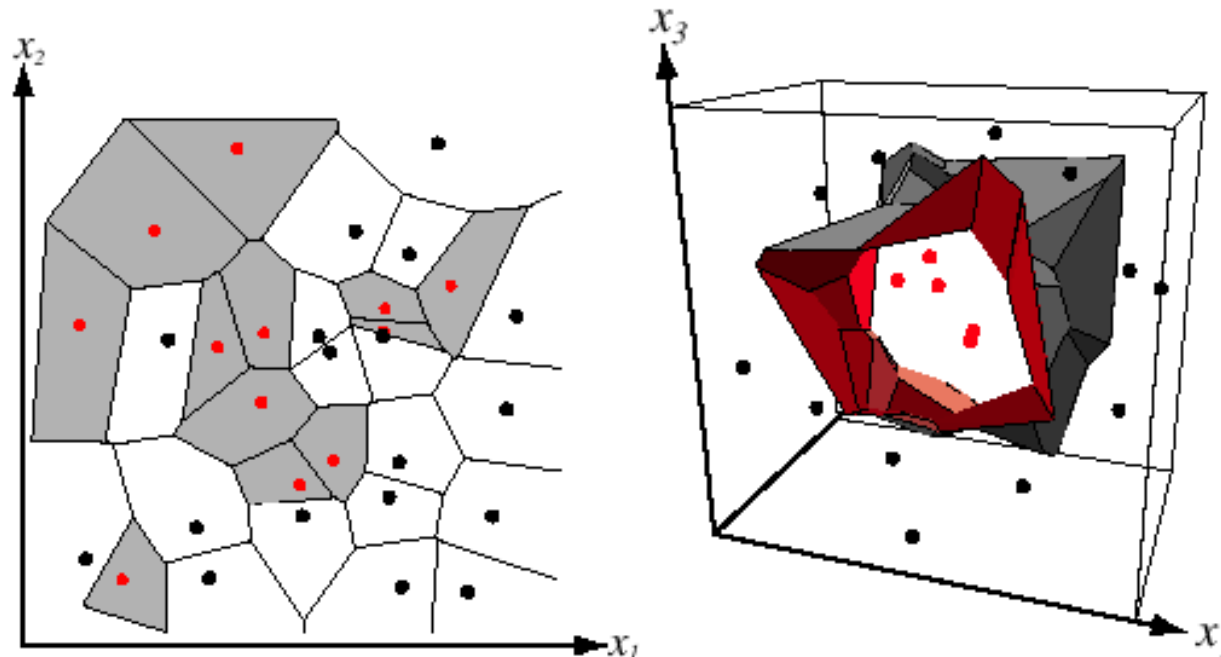
# kNN Classification Rule

- 3-class problem unimodal non-linearly separable.
- Use  $k=5$  and euclidian distance
- The resulting boundaries and decision regions are shown below



# kNN Decision Boundary

The nearest neighbor rule leads to a Voronoi Tessellation



**FIGURE 4.13.** In two dimensions, the nearest-neighbor algorithm leads to a partitioning of the input space into Voronoi cells, each labeled by the category of the training point it contains. In three dimensions, the cells are three-dimensional, and the decision boundary resembles the surface of a crystal. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.



**Georgy Voronoy**

1868-1908



**Johann Peter  
Gustav Lejeune  
Dirichlet**

1805-1859

# Recap

- Non-Parametric Density Estimation
  - Problem Formulation
  - Different Approaches
- Histograms
- Kernels and Kernel Density Estimation
  - Parzen Windows
- kNN Density Estimation and Classification