

Module 1

Bayes Decision Theory

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

$g(w_i|x) \rightarrow$ omit $P(x)$

Discriminant function

Decision Surface (D2D) (NDF)

Gaussian fn for NDF

Discrimin fn for NDF

$\frac{-1}{2} \frac{(x-u)^2}{\sigma^2}$

$$1D \quad P(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2} \frac{(x-u)^2}{\sigma^2}}$$

Module 2

Ref.

"chapter 2 - f Dnd"

Parameter Estimation Methods

$$P(x) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} e^{-\frac{1}{2} (x-u)^T \Sigma^{-1} (x-u)}$$

Σ = covariance matrix.

✓ Supervised scenario \rightarrow PARAMETER ESTIMATION:
 $(\mu, \Sigma)!!$

$$P(w_i | X_t) = \frac{P(w_i) P(x|w_i)}{P(x)}.$$

$P(w_i)$? \rightarrow estimate

$\hookrightarrow \frac{P(x|w_i)}{\text{, I don't know what is}}$ \rightarrow estimate \rightarrow Gaussian Distrib^n in multivariable

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

Estimate what the $\hat{\mu}, \hat{\Sigma}$ ^{is} _{^n (parameters)} distribution follows.

Parameter Estimation

- optimal classifier \checkmark -> if we knew the prior $P(\omega_i)$ and the class-
conditional densities $p(x|\omega_i)$. \checkmark
- But is it always possible to have this information?
- What we have is a number of design samples or training data —
particular representatives of the patterns we want to classify.

Parameter Estimation

- Can we use training data to design the classifier?
- Consider Supervised Learning scenario (labelled data, class labels are available)
Where the ~~data~~ samples and the class to which they belong is known.
Then we can use the [✓] samples to estimate the unknown probabilities and probability densities. $P(x|w_i)$?

Parameter Estimation

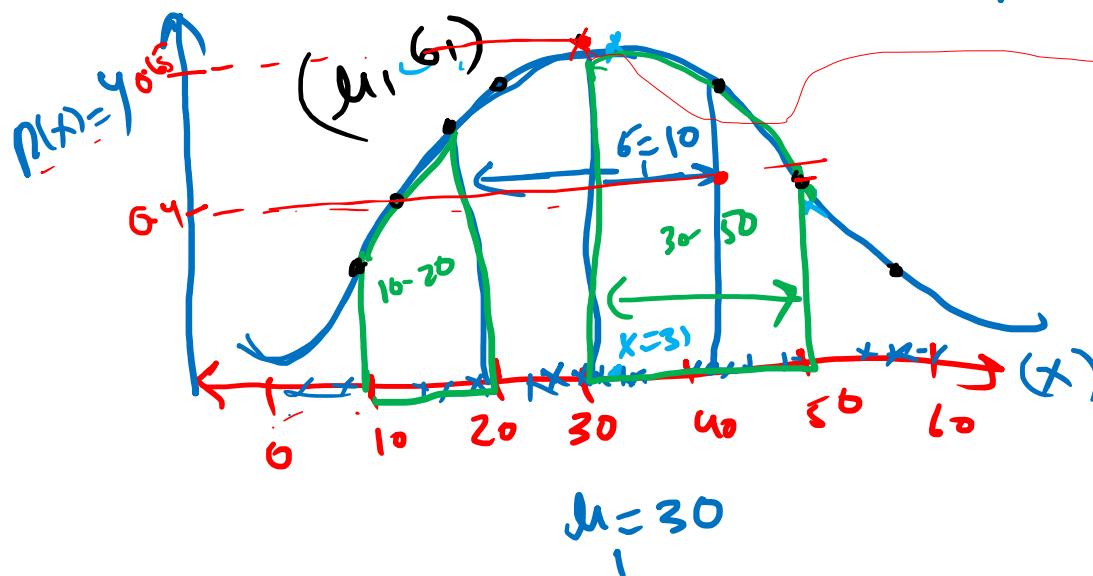
- We assume that $p(x|\omega_i)$ is a normal density with mean μ_i and covariance matrix Σ_i ,

$$p(x|\omega_i) = \mathcal{N}(\mu_i, \Sigma_i)$$

- Instead of estimating an unknown function $p(x|\omega_i)$ -> now we have to estimating the parameters μ_i and Σ_i

Maximum Likelihood Estimation Probability

- 1D. NDF
- Likelihood vs Probability



$$(\mu, \sigma) \\ \mu = 30, \sigma = 10$$

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

$$P(X=30) = 0.65 \\ P(X=50) = 0.40$$

$$P(X=30 \text{ to } 50) \quad | \quad \mu_1 = 30, \sigma_1 = 10$$

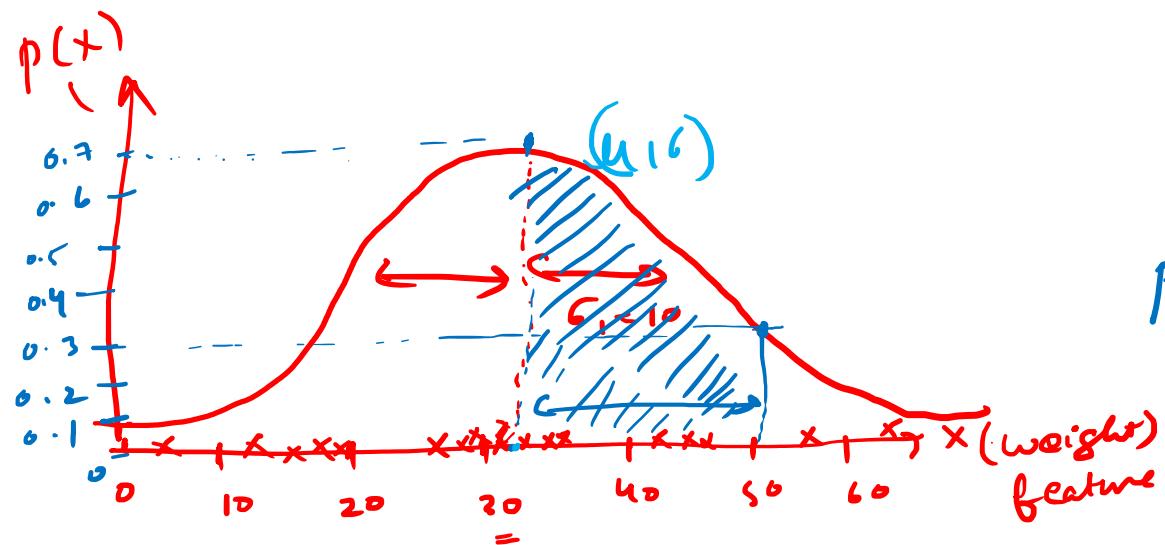
$$P(X=10 \text{ to } 20) \quad | \quad \mu_1 = 30, \sigma_1 = 10$$

prob.

Maximum Likelihood Estimation

- Likelihood vs Probability

Example of fishes



$$\begin{aligned}\mu &= 30 \\ \sigma &= 10\end{aligned}$$

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

$$\mu = 30, \sigma = 10$$

$$p(x=30) = 0.7$$

(given that $\mu=30, \sigma=10$)

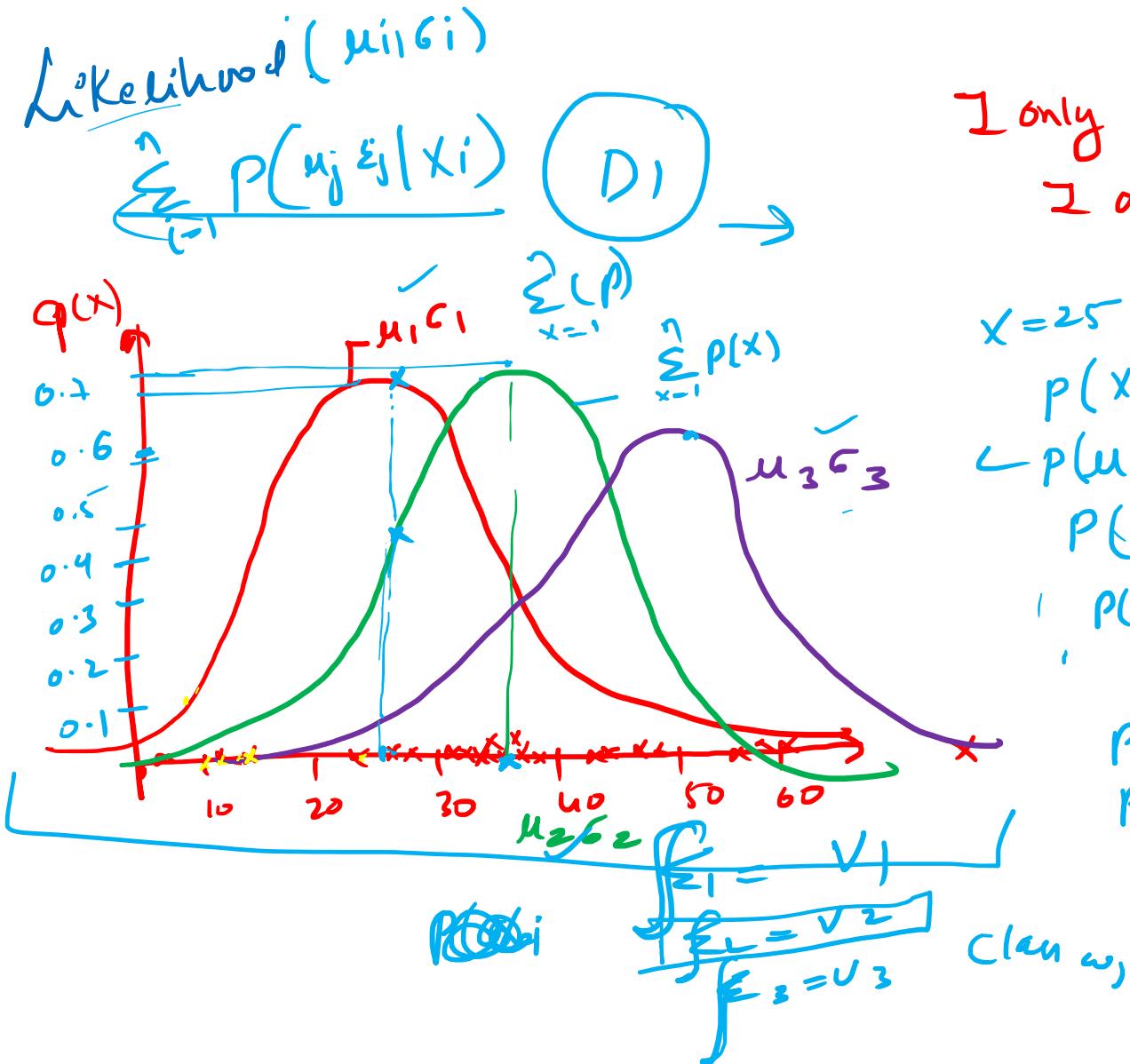
$$p(x=50) = 0.3$$

($\mu=30, \sigma=10$)

$$p(x=30 \text{ to } 50 \mid \mu=30, \sigma=10)$$

What is $p(x)$. given μ, σ .

Max. Likelihood Estimation.



I only have samples points (x) \rightarrow Training data.
I don't know what is the μ_i & σ_i ?

$$x = 25$$

$$p(x = 25) ?$$

$$p(\mu_1 | c_1 | x = 25) = 0.7 \quad (\checkmark)$$

$$p(\mu_2 | c_2 | x = 25) = 0.2 \quad \text{decided}$$

$$p(\mu_3 | c_3 | x = 25) = 0.1$$

$$p(x = 35)$$

$$p(\mu_1 | c_1 | 35) = 0.6$$

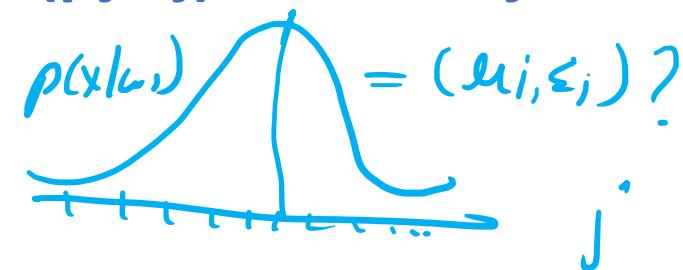
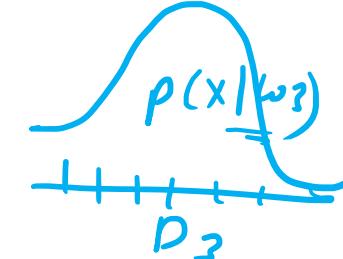
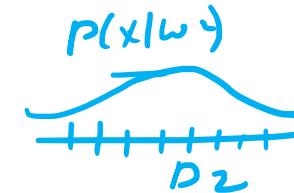
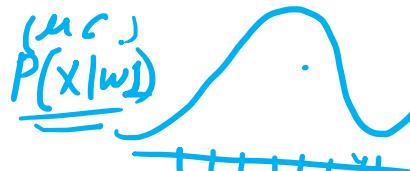
$$p(\mu_2 | c_2 | 35) = 0.25 \rightarrow \checkmark (\mu_2 | c_2)$$

$$p(\mu_3 | c_3 | 35) = 0.15$$

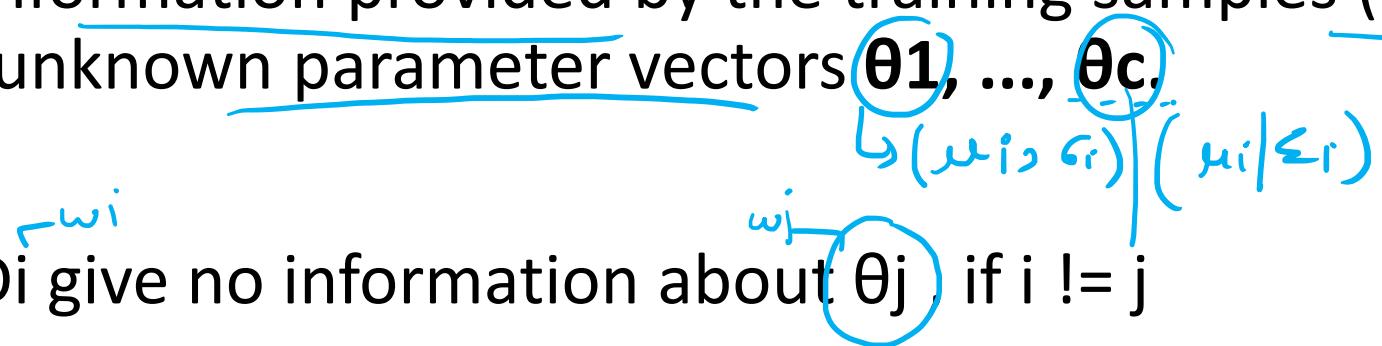
Maximum Likelihood Estimation

- Let say we have 'c' classes.
- We also have 'c' sets, D_1, \dots, D_c , of data samples
- Samples in D_j having been drawn independently according to i.i.d. the probability law $p(x|\omega_j)$.
- The distribution $(x|\omega_j)$ has a known parametric form -> determined uniquely by parameter vector θ_j .
- **example, we might have $p(x|\omega_j) \sim N(\mu_j, \Sigma_j)$, where θ_j consists of the components of μ_j and Σ_j .**

$$p(x|\omega_i) p(\omega_i) \Rightarrow g_j(x_i)$$



Maximum Likelihood Estimation

- Objective : information provided by the training samples (D_1, \dots, D_c) to estimate unknown parameter vectors $\theta_1, \dots, \theta_c$.


Annotations: w_i points to w_i , (w_i, c_i) points to (w_i, c_i) , $(x_i | c_i)$ points to $(x_i | c_i)$, w_j points to w_j , θ_j points to θ_j , θ_i points to θ_i .
- samples in D_i give no information about θ_j if $i \neq j$
- that is, we shall assume that the parameters for the different classes are functionally independent.
- This permits us to work with each class separately, and to simplify our notation by deleting indications of class distinctions.
- $p(x|\omega, \theta) \rightarrow p(x|\theta)$
 $p(x|w, \theta) \rightarrow p(x|\theta)$

Maximum Likelihood Estimation

~~Assume $p(x|w)$ → Gaussian Dist'n.~~

- Use a set D of training samples drawn independently from the probability density $p(x|\theta)$ to estimate unknown parameter vector θ .
- Suppose that D contains n number of samples: x_1, \dots, x_n .
- Since the samples were drawn independently, we have

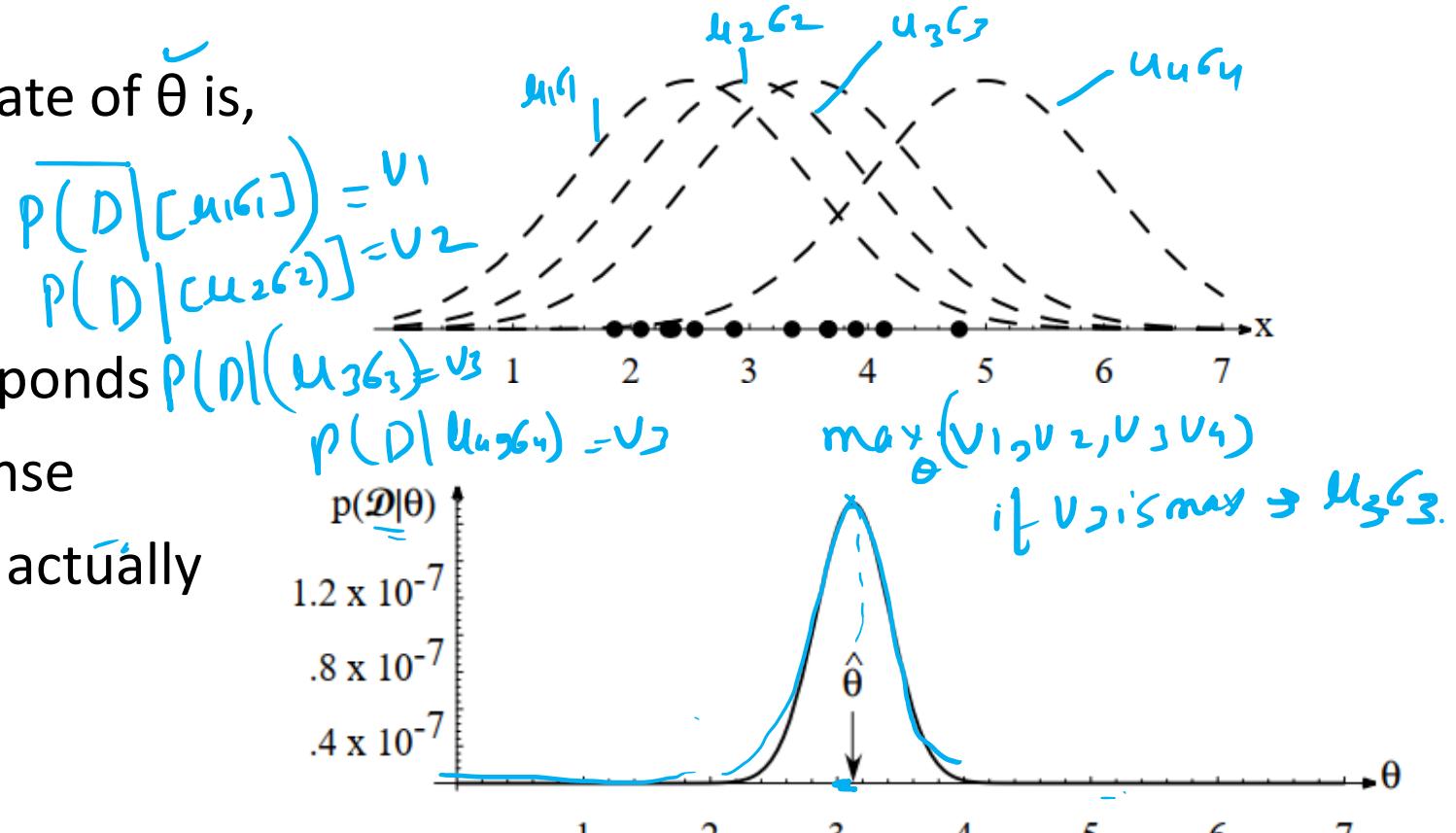
$$p(D|\theta) = \prod_{k=1}^n p(x_k|\theta) = \prod_{i=1}^n p(x_i|\mu_i, \sigma_i)$$
$$\theta = \begin{bmatrix} \mu_i \\ \sigma_i \end{bmatrix}$$

Maximum Likelihood Estimation

- $p(D|\theta)$ is called the likelihood of θ with respect to the set of samples\

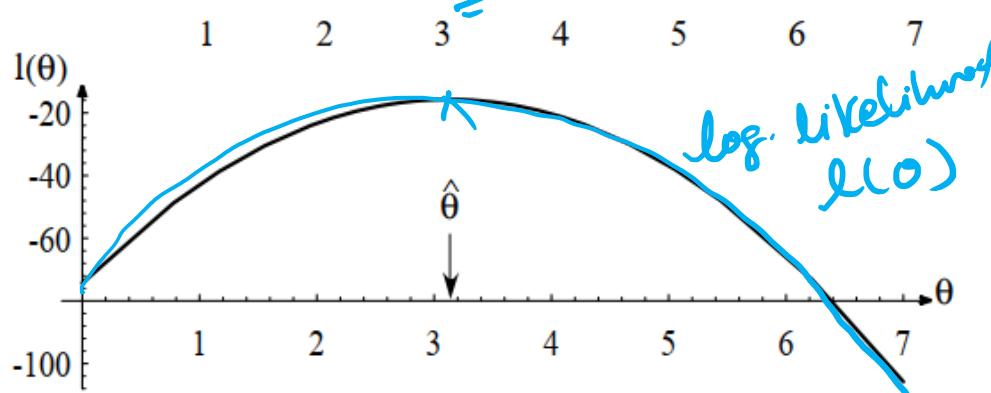
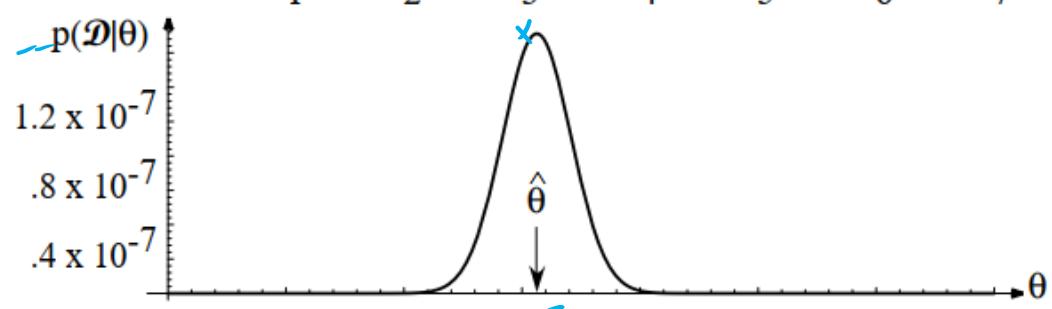
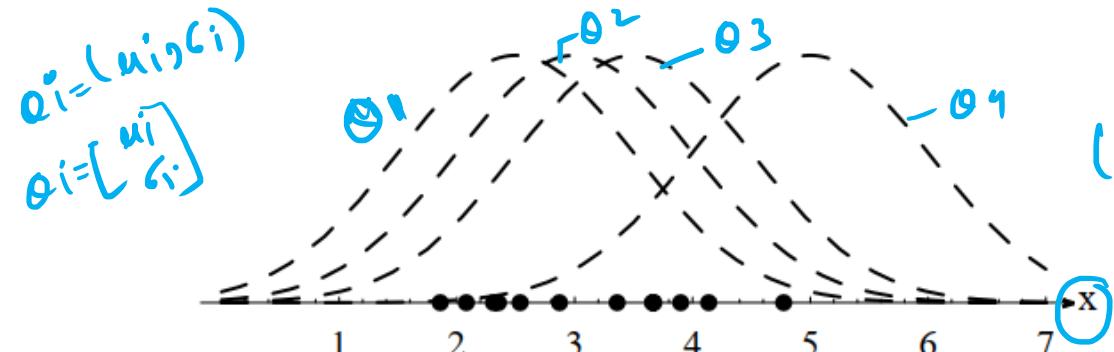
- The maximum likelihood estimate of θ is,

value $\hat{\theta}$ that maximizes $p(D|\theta)$.
 $\hat{\theta}$ $\approx \mu_3$



- Intuitively, this estimate corresponds to the value of θ that in some sense best agrees with or supports the actually Observed training samples

Maximum Likelihood Estimation



$$P(D|\theta) = \text{is max if a part of } \hat{\theta}$$

$$P(D|\theta) = \prod_{k=1}^n P(x_k|\theta_i) = \text{max. likelihood.}$$

$$\log P(D|\theta) = \log \text{of max. likelihood}$$

$$\log P(D|\theta).$$

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 $p(D|\theta)$ 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 — i.e., the log-likelihood $l(\theta)$, shown at the bottom. Note especially that the likelihood lies in a different space from $p(x|\hat{\theta})$, and the two can have different functional forms.

Maximum Likelihood Estimation

It is usually easier to work with the logarithm of the likelihood than with the likelihood itself

- Lets define $\underline{l}(\theta)$ as the log-likelihood function

$$l(\theta) \equiv \ln p(\mathcal{D}|\theta).$$

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

$\underline{l}(\theta) = \ln p(\mathcal{D}|\theta)$

$\ln(p(\mathcal{D}|\theta)) = \sum_{k=1}^n \ln p(x_k|\theta)$

$\ln(p(\mathcal{D}|\theta)) = \sum_{k=1}^n \ln p(x_k|\theta)$

$\ln(p(\mathcal{D}|\theta)) = \ln(p(a) \cdot p(b) \cdots h) = \ln p(a) + \ln p(b) + \cdots$

$\ln(p(\mathcal{D}|\theta)) = \sum_{k=1}^n \ln p(x_k|\theta)$

$\ln(p(\mathcal{D}|\theta)) = \sum_{k=1}^n \ln p(x_k|\theta)$

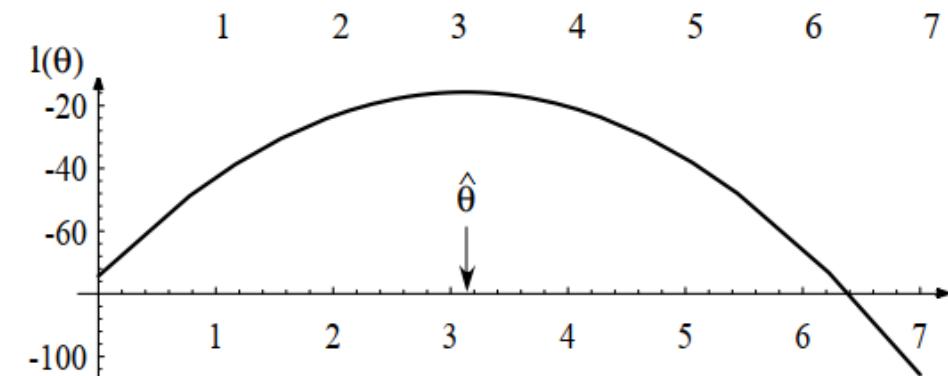
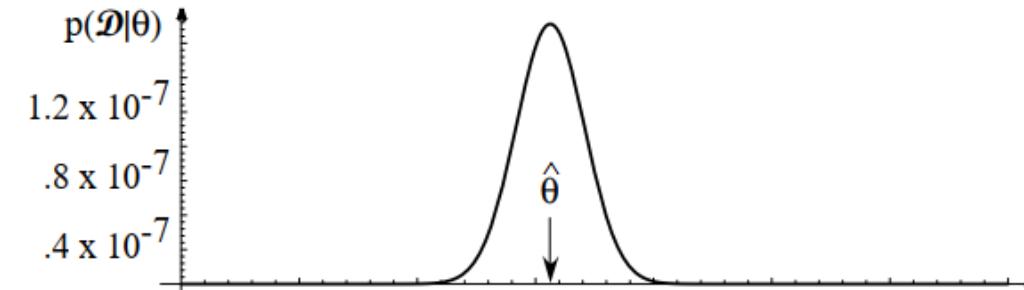
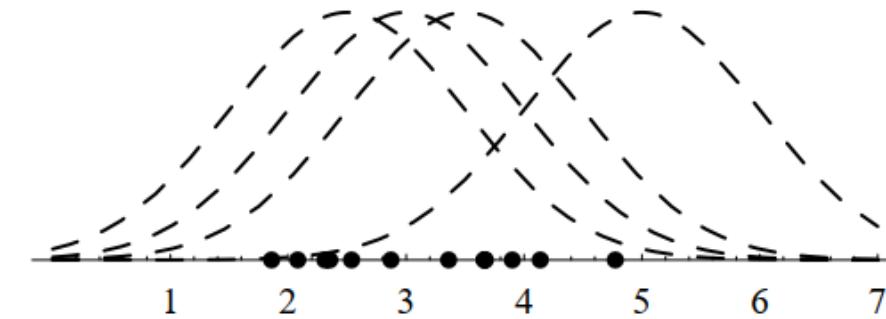
Maximum Likelihood Estimation

we don't know^{it} what is θ !

- We can then write our solution formally as the argument $\hat{\theta}$ that maximizes the loglikelihood, i.e.,

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

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



Taking $\partial \theta$ & $p_{\text{mod}} = \text{diff}$ (P&S)

Maximum Likelihood Estimation

- We can then write our solution formally as the argument θ that maximizes the loglikelihood, i.e.,

Likelihood

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

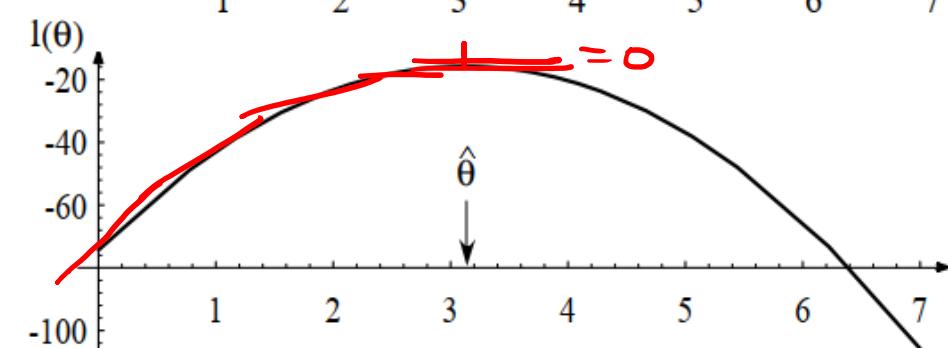
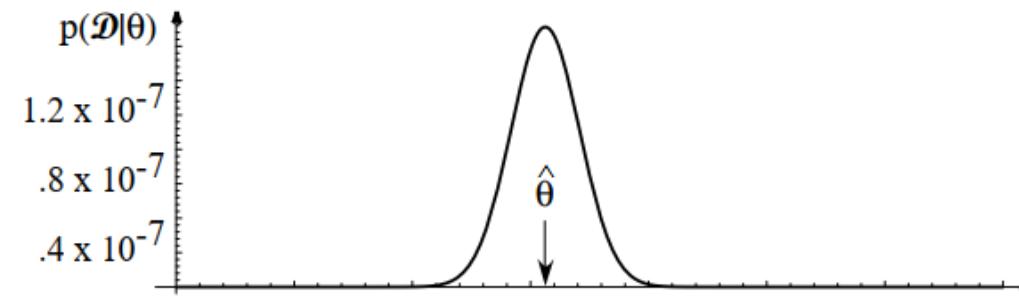
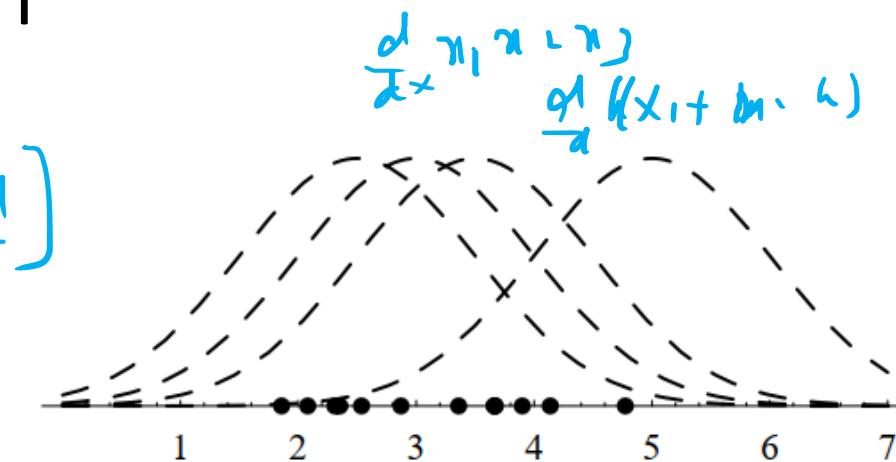
✓ Take first derivative:

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

Equate to zero to get a maxima:

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

eq for max for $\theta \rightarrow l(\theta)$ is maximum



Maximum Likelihood Estimation

- We can then write our solution formally as the argument $\hat{\theta}$ that maximizes the loglikelihood, i.e.,

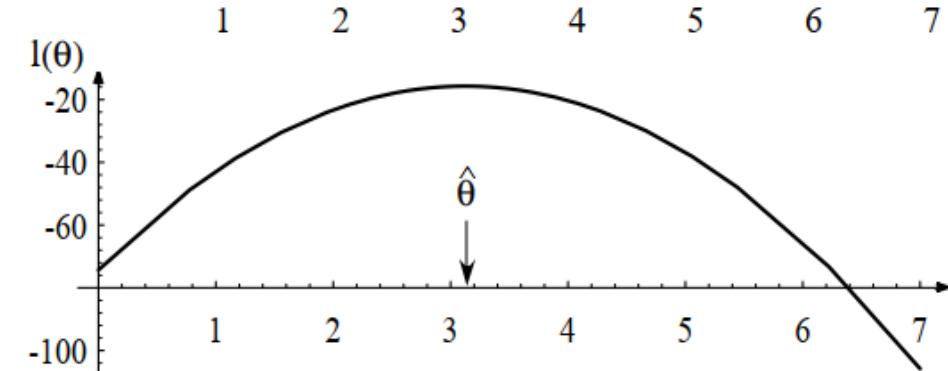
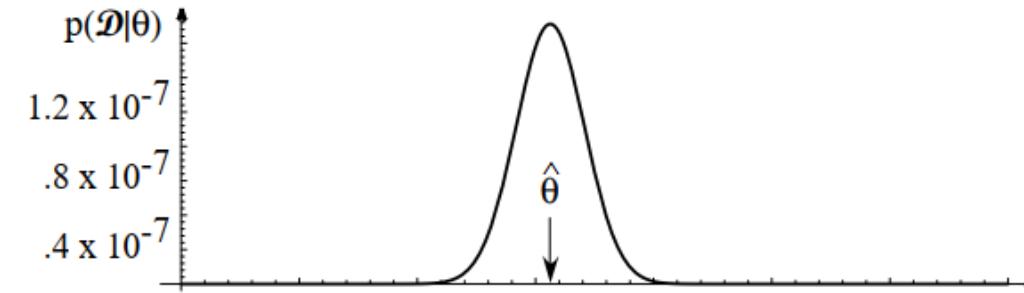
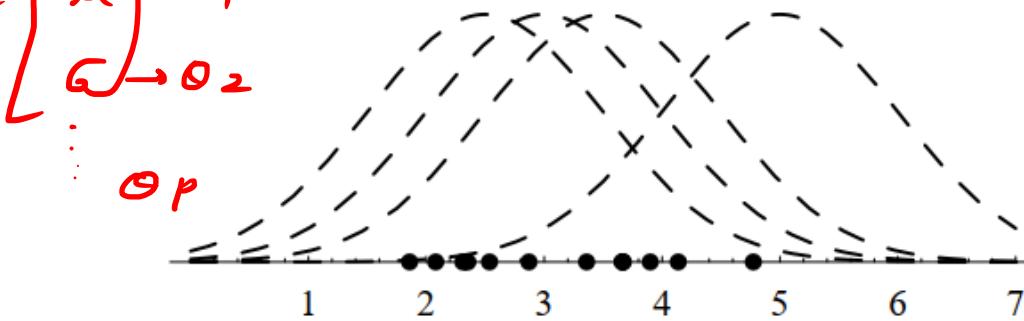
p -component vector $\underline{\theta} = (\theta_1, \dots, \theta_p)^t$, and $\nabla_{\underline{\theta}}$ be the gradient operator

$$\nabla_{\underline{\theta}} \equiv \begin{bmatrix} \frac{\partial}{\partial \theta_1} \\ \vdots \\ \frac{\partial}{\partial \theta_p} \end{bmatrix}.$$

$$\nabla_{\underline{\theta}} l = \sum_{k=1}^n \nabla_{\underline{\theta}} \ln p(\mathbf{x}_k | \underline{\theta}).$$

$$\boxed{\nabla_{\underline{\theta}} l = 0.} =$$

$$\underline{\theta} = \begin{bmatrix} u \\ \omega \\ \vdots \\ \theta_p \end{bmatrix}$$



The Gaussian Case: Unknown μ

- For simplicity, consider first the case where only the mean is unknown.

$$\boldsymbol{\theta} = \begin{bmatrix} \theta_1 = \mu \\ \vdots \\ \theta_r \end{bmatrix}$$

$$\underline{\theta_1} = \underline{\mu}$$

The Gaussian Case: Unknown μ

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

- For simplicity, consider first the case where only the mean is unknown.
- Under this condition, we consider a sample point x_k and find

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

The Gaussian Case: Unknown μ

- For simplicity, consider first the case where only the mean is unknown.
- Under this condition, we consider a sample point x_k and find

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

- and

$$\nabla_{\theta} \ln p(x_k | \mu) = \left(\Sigma^{-1} (x_k - \mu) \right) \rightarrow = 0.$$

The Gaussian Case: Unknown μ

- that the maximum likelihood estimate for μ must satisfy

$$\sum_{k=1}^n \Sigma^{-1}(\mathbf{x}_k - \hat{\boldsymbol{\mu}}) = 0,$$

Equating the derivative to zero

- Multiplying by Σ and rearranging, we obtain

$$\boxed{\hat{\boldsymbol{\mu}} = \frac{1}{n} \sum_{k=1}^n \mathbf{x}_k.} = ① \text{ Proved over here with the help of the equations.}$$

This is a very satisfying result. It says that the maximum likelihood estimate for the unknown population mean is just the arithmetic average of the training samples — the sample mean

The Gaussian Case: Unknown μ and Σ

- Here neither the mean μ nor the covariance matrix Σ is known.
- Thus, these unknown parameters constitute the components of the parameter vector θ .
- Consider first the univariate case with $\theta_1 = \mu$ and $\theta_2 = \sigma^2$.
- Here the log-likelihood of a single point is

$$\ln p(x|w_1) = \frac{1}{\sqrt{2\pi\theta_2}} e^{-\frac{1}{2} \left(\frac{x-\mu}{\theta_2}\right)^2}$$

$$\begin{aligned}\mu_1 &= \theta_1 \\ \sigma &= \theta_2\end{aligned}$$

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

The Gaussian Case: Unknown μ and Σ

- Here neither the mean μ nor the covariance matrix Σ is known.
- Thus, these unknown parameters constitute the components of the parameter vector θ .
- Consider first the univariate case with $\theta_1 = \mu$ and $\theta_2 = \sigma^2$.
- Here the log-likelihood of a single point is

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

- and its derivative is

$$\nabla_{\theta} l = \nabla_{\theta} \ln p(x_k|\theta) = \left[\begin{array}{c} \frac{1}{\theta_2}(x_k - \theta_1) \\ -\frac{1}{2\theta_2} + \frac{(x_k - \theta_1)^2}{2\theta_2^2} \end{array} \right].$$

$\frac{\partial}{\partial \theta_1} = 0$
 $\frac{\partial}{\partial \theta_2} = 0$

The Gaussian Case: Unknown μ and Σ

- follow book.
- Taking derivative and equating to zero we have: $\sum_{k=1}^n \frac{1}{\hat{\theta}_2} (x_k - \hat{\theta}_1) = 0$ — //w.r.t θ_1 th We have proved by taking MLE & eq 1st derivative to zero.

$$-\sum_{k=1}^n \frac{1}{\hat{\theta}_2} + \sum_{k=1}^n \frac{(x_k - \hat{\theta}_1)^2}{\hat{\theta}_2^2} = 0, \quad \text{— //w.r.t } \theta_2$$

We get

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

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

Result 2. of MLE.
sample mean \rightarrow mean of the
given pop set (\bar{x})
& sample var \rightarrow variance of
the given pop set (s^2)

The Gaussian Case: Unknown μ and Σ

Today's class

- the maximum likelihood estimates for μ and Σ are given by

*Find your
classmate
and do
MLC."*

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

$$\hat{\Sigma} = \frac{1}{n} \sum_{k=1}^n (\mathbf{x}_k - \hat{\mu})(\mathbf{x}_k - \hat{\mu})^t.$$

(1) Prob (2) Likelihood
(3) log likelihood.

$$\begin{aligned} \text{Max likelihood} & P(D | \theta) \\ & = \prod_{k=1}^n P(\mathbf{x}_k | \theta) \\ & \theta = \underline{(\theta_1, \theta_2, \dots)} \end{aligned}$$

Thus, once again we find that the maximum likelihood estimate for the mean vector is the sample mean. The maximum likelihood estimate for the covariance matrix is the arithmetic average of the n matrices $(\mathbf{x}_k - \hat{\mu})(\mathbf{x}_k - \hat{\mu})^t$. Since the true covariance matrix is the expected value of the matrix $(\mathbf{x} - \hat{\mu})(\mathbf{x} - \hat{\mu})^t$, this is also a very satisfying result.

$\theta_1 = \mu_1 = \text{sample mean}$

$\theta_2 - \sigma^2 = \text{sample variance}$

Bayesian estimation

- Maximum likelihood method considers the true parameter vector, θ , to be a fixed (constant) value,
- Bayesian learning considers θ to be a random variable,
- We don't know the exact value of parameter, but we know that the parameter falls in a certain range of values.
- We capture our lack of knowledge about the value of parameter -> through a probability density over the parameter space.

Bayesian estimation

- We have a **prior density $p(\theta)$ the parameter** (prior meaning without having any look at the training data).
- prior density $p(\theta)$ captures information about the parameters.
- We can then have access to our data **D** (the training samples).
- We then see how our data **D** is useful in mapping from **prior density $p(\theta)$ into posterior density $p(\theta|D)$ for the parameter θ .**

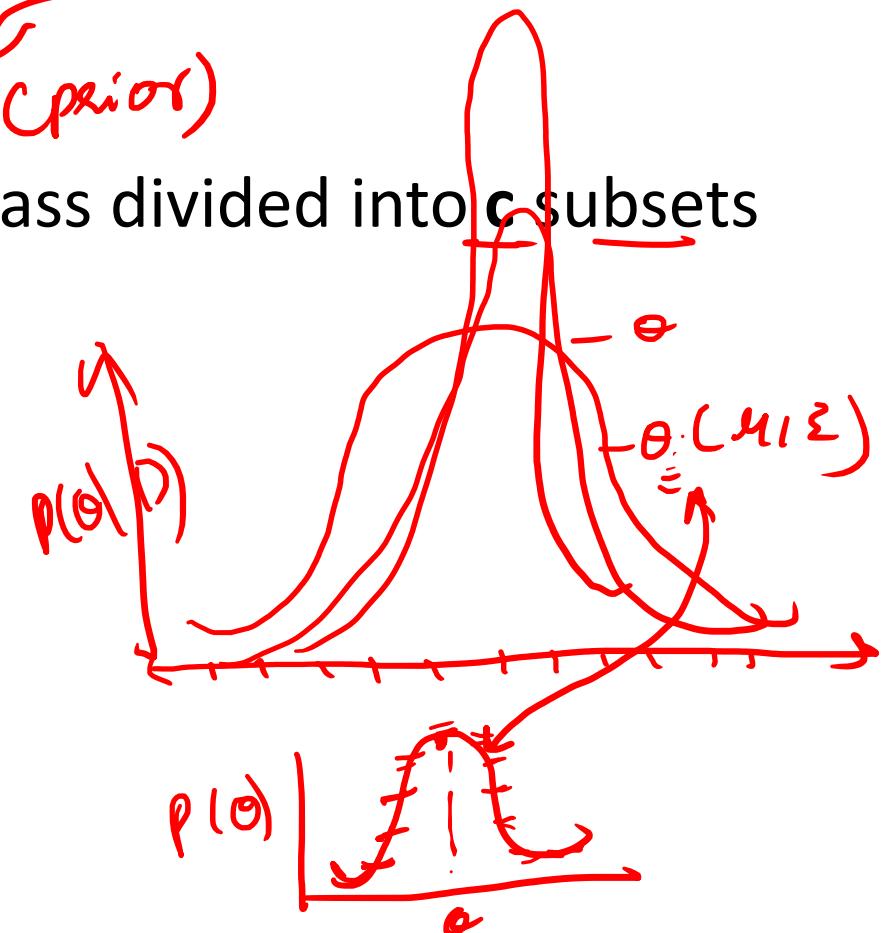
Notations

- The parameter is denoted by θ
- The prior density of parameter -> $p(\theta)$
- We have the training samples for ever class divided into c subsets
 - D_1, \dots, D_c ,
 - with the samples in D_i belonging to class ω_i
- We have posterior density $p(\theta | D)$

Notations

- The parameter is denoted by θ *= fixed - ω is random variable*
- The prior density of parameter $\rightarrow p(\theta)$ *(prior)*
- We have the training samples for every class divided into c subsets
 - $D_1, \dots, D_c = []$
 - with the samples in D_i belonging to class ω_i
- We have posterior density $p(\theta | D)$

$$P(\omega | D)$$



Recap Bayesian Rule

- Let D denote the set of samples, then we can emphasize the role of the samples by saying that our
- goal is to compute the posterior probabilities $P(\omega_i | x, D)$

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

- Given ‘ c ’ classes we have c separate parameter estimation problems $p(x | D)$ corresponding to each class ω_i .

Recap Bayesian Rule

- Let D denote the set of samples, then we can emphasize the role of the samples by saying that our
- goal is to compute the posterior probabilities $P(\omega_i | x, D)$

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

freakin' aw $\frac{p(x | \omega_j, D)}{P(\omega_i | D)}$ $\leftarrow r(x)$

- Given ' c ' classes we have c separate parameter estimation problems $p(x | D)$ corresponding to each class ω_i .

The Parameter Distribution

- Although desired pdf $p(x)$ is unknown, we assume that it has a known parametric form.
- Lets say parametric form is a gaussian distribution.
- The only thing assumed unknown is the value of a parameter vector θ .

The Parameter Distribution

- our basic goal is to compute $p(x|D)$
- We do this by integrating the joint density $p(x, \theta|D)$ over θ .
- That is

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

- we can write $p(x, \theta|D)$ as the product $p(x|\theta, D) \cdot p(\theta|D)$
- Since the selection of x and that of the training samples in D is done independently, the first factor $p(x|\theta, D)$ is merely $p(x|\theta)$

The Parameter Distribution

- That is, the distribution of x is known completely once we know the value of the parameter vector θ

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

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

prior $\underline{\theta} - p(\theta)$

data seen $\mathcal{D} -$

$p(\theta|\mathcal{D})$ = posterior prob.

- This key equation links the desired class-conditional density $p(x|\mathcal{D})$ to the **posterior density $p(\theta|\mathcal{D})$** for the unknown parameter vector $\underline{\theta}$

The Parameter Distribution

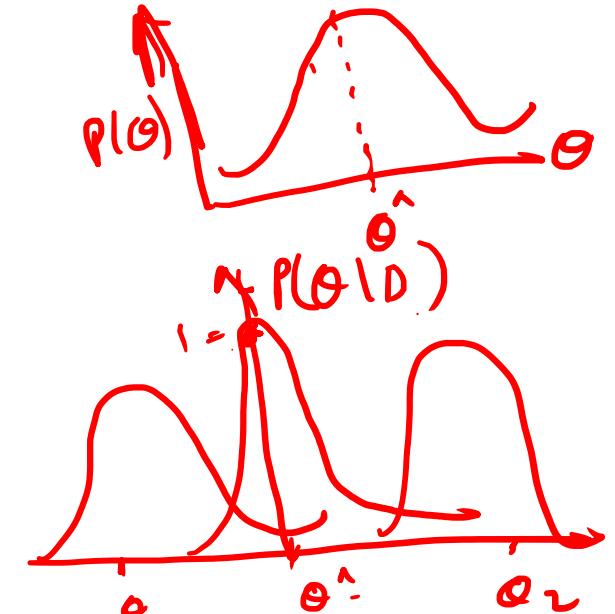
$$p(\mathbf{x}|\mathcal{D}) = \int p(\mathbf{x}|\boldsymbol{\theta})p(\boldsymbol{\theta}|\mathcal{D}) d\boldsymbol{\theta}.$$

- Lets say we have a given range of parameter $\boldsymbol{\theta}$.
- If $p(\boldsymbol{\theta}|\mathcal{D})$ peaks very sharply at some value $\boldsymbol{\theta}^*$, we obtain $p(\mathbf{x}|\mathcal{D})=p(\mathbf{x}|\boldsymbol{\theta}^*)$,
- i.e., the result we would obtain by substituting the estimate $\boldsymbol{\theta}^*$ for the true parameter vector.
- when the unknown densities have a known parametric form, the samples exert their influence on $p(\mathbf{x}|\mathcal{D})$ through the posterior density $p(\boldsymbol{\theta}|\mathcal{D})$.

The Parameter Distribution

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

p(x|D) circled in red. Red annotations: "unknown densities" above the integral, "scattered" near the second term, and a checkmark to the right.



- Lets say we have a given range of parameter θ .
- If $p(\theta|D)$ peaks very sharply at some value $\hat{\theta}$, we obtain $p(x|D) \approx p(x|\hat{\theta})$,
- i.e., the result we would obtain by substituting the estimate $\hat{\theta}$ for the true parameter vector.
- when the unknown densities have a known parametric form, the samples exert their influence on $p(x|D)$ through the posterior density $p(\theta|D)$.

Bayesian Parameter Estimation: Gaussian Case

3.2

- Bayesian estimation techniques to calculate the a posteriori density $p(\theta|D)$ and the desired probability density $p(x|D)$ for the case where
- $p(x|\mu) \sim N(\mu, \Sigma)$.

$$p(x|\theta) = p(x|\mu) \sim N(\mu, \Sigma)$$

The Univariate Case: $p(\mu | D)$

- the case where μ is the only unknown parameter.

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

- where the only unknown quantity is the mean μ .
- We assume that whatever prior knowledge we might have about μ can be expressed by a known **prior density $p(\mu)$** .

$$p(\mu) \sim N(\mu_0, \sigma_0^2),$$

- where both μ_0 and σ_0^2 are known
- And μ_0 represents our best a priori guess for μ , and σ_0^2 measures our uncertainty about this guess.

The Univariate Case: $p(\mu | D)$

- the case where μ is the only unknown parameter.

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

- where the only unknown quantity is the mean μ .

- We assume that whatever prior knowledge we might have about μ can be expressed by a known prior density $p(\mu)$.

Prior dens $p(\mu)$
 $\sigma = \mu$ $p(\mu)$

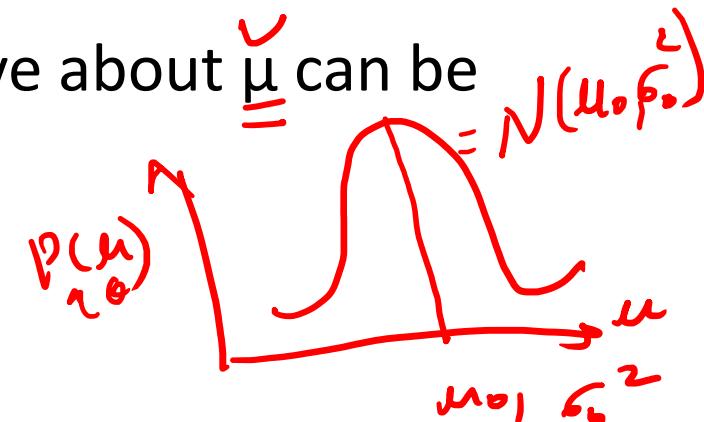
$$p(\mu) \sim N(\mu_0, \sigma_0^2),$$

- where both μ_0 and σ_0^2 are known
- And μ_0 represents our best a priori guess for μ , and σ_0^2 measures our uncertainty about this guess.

$$p(\theta | D)$$

$$p(\mu | D)$$

$$p(x|\theta) \Rightarrow p(x|\mu) = N(\mu, \sigma^2)$$



The Univariate Case: $p(\mu | D)$

- Having selected the a priori density for μ , we can view the situation as follows.
- Imagine that a value is drawn for μ from a population governed by the probability law $p(\mu)$.

The Univariate Case: $p(\mu | D)$

- Having selected the a priori density for μ , we can view the situation as follows.
- Imagine that a value is drawn for μ from a population governed by the probability law $p(\mu)$.
- Once this value is drawn, it becomes the true value of μ and completely determines the density for x .
- Suppose now that n samples x_1, \dots, x_n are independently drawn from the resulting population.



The Univariate Case: $p(\mu | D)$

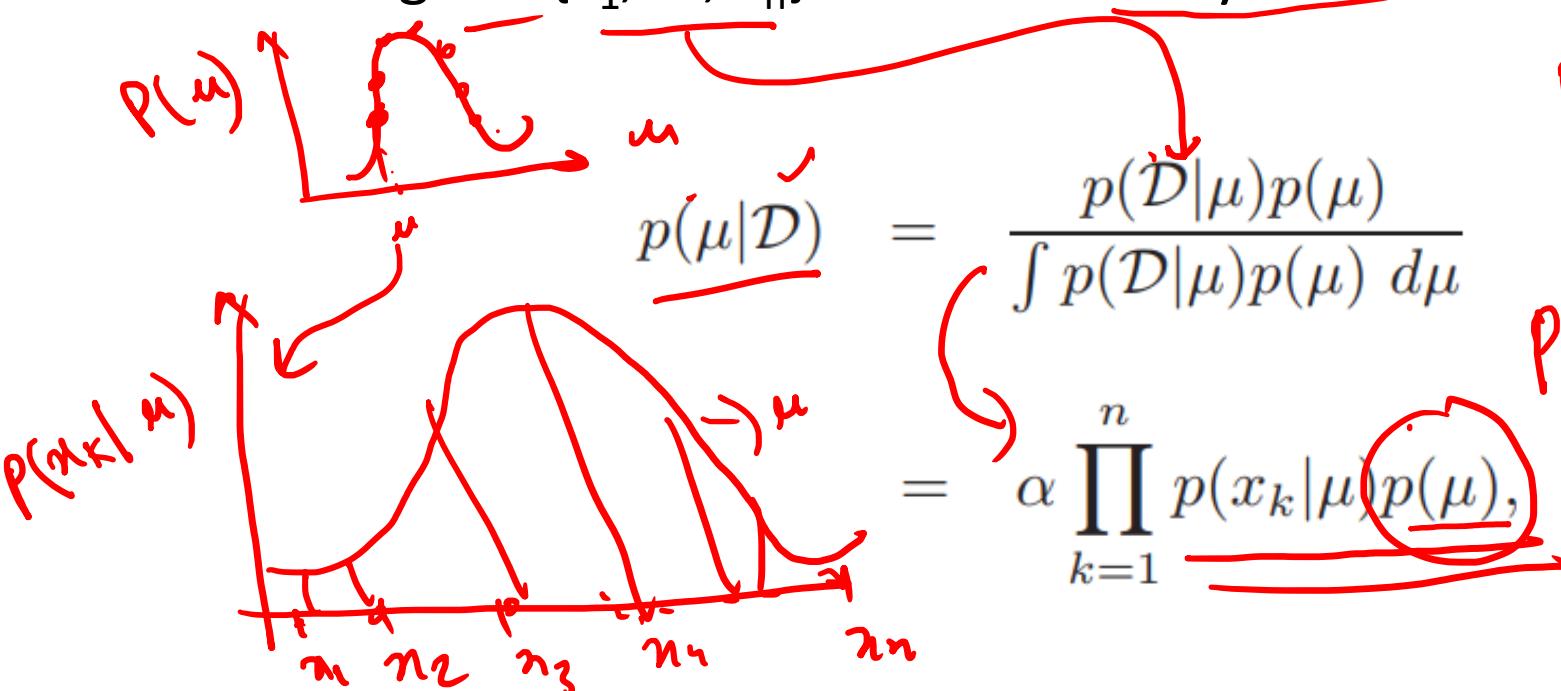
- Letting $D = \{x_1, \dots, x_n\}$ and we use Bayes' formula

$$\begin{aligned} p(\mu | D) &= \frac{p(D|\mu)p(\mu)}{\int p(D|\mu)p(\mu) d\mu} \\ &= \alpha \prod_{k=1}^n p(x_k|\mu)p(\mu), \end{aligned}$$

- where α is a normalization factor that depends on D but is independent of μ
- it relates the prior density $p(\mu)$ to an a posteriori density $p(\mu|D)$

The Univariate Case: $p(\mu | D)$

- Letting $D = \{x_1, \dots, x_n\}$ and we use Bayes' formula



- where α is a normalization factor that depends on D but is independent of μ
- it relates the prior density $p(\mu)$ to an posterior density $p(\mu|D)$

$p(\mu | D) = \frac{p(D|\mu) \cdot p(\mu)}{\int p(D|\mu) \cdot p(\mu) d\mu}$

$p(\mu | D)$

The Univariate Case: $p(\mu | D)$

$$p(\mu) = N(\mu_0, \sigma_0^2)$$
$$p(x_k | \mu) = N(\mu, \sigma^2)$$

- Since

$$p(\mu | \mathcal{D}) = \alpha \prod_{k=1}^n p(x_k | \mu) p(\mu),$$

- And $p(x_k | \mu) \sim N(\mu, \sigma^2)$ and $p(\mu) \sim N(\mu_0, \sigma_0^2)$, we have

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

The Univariate Case: $p(\mu | D)$

- Since

$$p(\mu | \mathcal{D}) = \alpha \prod_{k=1}^n p(x_k | \mu) p(\mu),$$

- And $p(x_k | \mu) \sim N(\mu, \sigma^2)$ and $p(\mu) \sim N(\mu_0, \sigma_0^2)$, we have

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

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

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

The Univariate Case: $p(\mu | D)$

- On simplifications

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

- Thus, $p(\mu | D)$ is an exponential function of a quadratic function of μ , i.e., is again a normal density.
- Since this is true for any number of training samples,
- $p(\mu | D)$ remains normal as the number n of samples is increased.
- $p(\mu | D)$ is said to be a **reproducing density**

The Univariate Case: $p(\mu | D)$

- If we write $p(\mu|D) \sim N(\mu_n, \sigma_n^2)$, then μ_n and σ_n^2 can be found by equating coefficients

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

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

The Univariate Case: $p(\mu | D)$

- If we write $p(\mu | \mathcal{D}) \sim N(\mu_n, \sigma_n^2)$, then μ_n and σ_n^2 can be found by equating coefficients

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

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

- Which gives

$$\frac{1}{\sigma_n^2} = \frac{n}{\sigma^2} + \frac{1}{\sigma_0^2} \quad \text{and} \quad \frac{\mu_n}{\sigma_n^2} = \frac{n}{\sigma^2} \bar{x}_n + \frac{\mu_0}{\sigma_0^2}, \quad \text{and} \quad \bar{x}_n = \frac{1}{n} \sum_{k=1}^n x_k.$$

where \bar{x}_n is the sample mean

The Univariate Case: $p(\mu | D)$

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

We solve explicitly for μ_n and σ_n^2 and obtain

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

Sample *estimated μ_0 of $P(\mu)$*

and

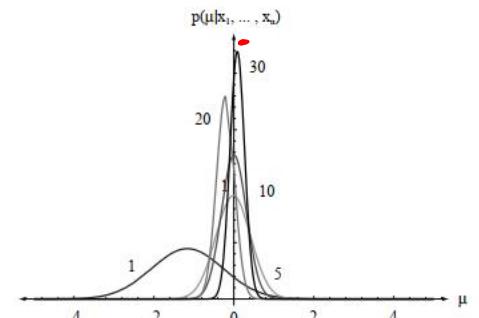
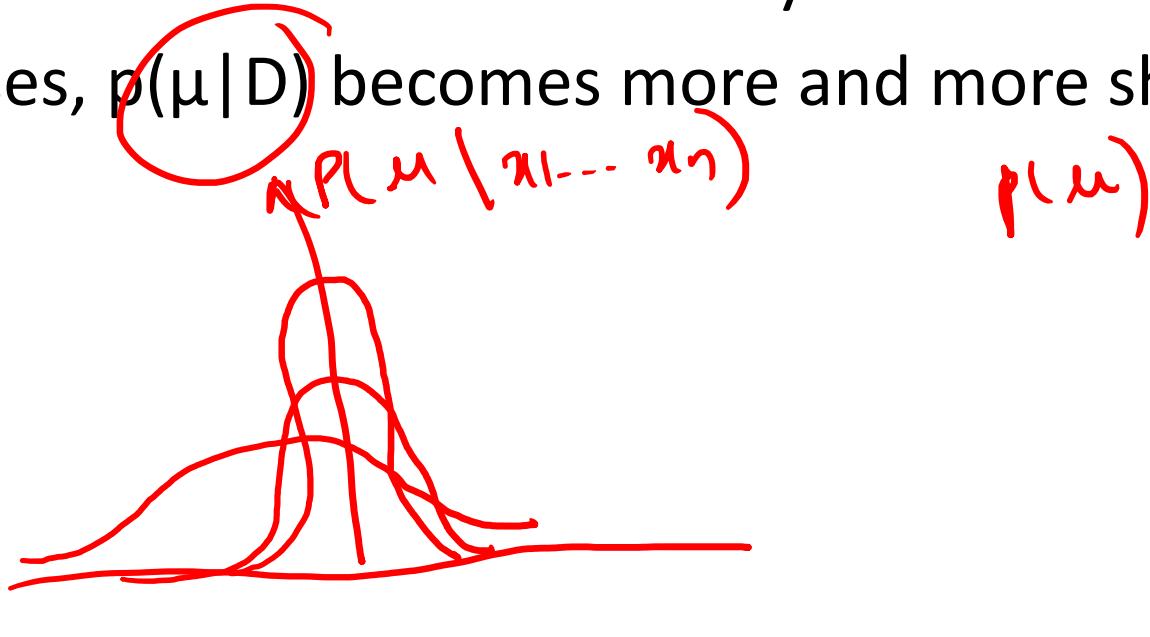
$$\sigma_n^2 = \frac{\sigma_0^2 \sigma^2}{n\sigma_0^2 + \sigma^2}. \quad p(\mu)$$

- These equations show how the prior information is combined with the empirical information in the samples to obtain the a posteriori density $p(\mu | D)$

$$p(\mu | D)$$

The Univariate Case: $p(\mu | D)$

- Roughly speaking, μ_n represents our best guess for μ after observing n samples, and σ_n^2 measures our uncertainty about this guess.
- Since σ_n^2 decreases monotonically with n — each additional observation decreases our uncertainty about the true value of μ .
- As n increases, $p(\mu | D)$ becomes more and more sharply peaked,



Computing $p(x | D)$

- $p(\mu | D)$ is computed
- $\underline{p(x | D)}$ remains

$$\begin{aligned} p(x|D) &= \int p(x|\mu)p(\mu|D) d\mu \\ &= \int \frac{1}{\sqrt{2\pi}\sigma} \exp \left[-\frac{1}{2} \left(\frac{x-\mu}{\sigma} \right)^2 \right] \frac{1}{\sqrt{2\pi}\sigma_n} \exp \left[-\frac{1}{2} \left(\frac{\mu-\mu_n}{\sigma_n} \right)^2 \right] d\mu \\ &= \frac{1}{2\pi\sigma\sigma_n} \exp \left[-\frac{1}{2} \frac{(x-\mu_n)^2}{\sigma^2 + \sigma_n^2} \right] f(\sigma, \sigma_n), \end{aligned} \tag{37}$$

$$p(x|D) \sim N(\mu_n, \sigma^2 + \sigma_n^2)$$

(Desired class-conditional density $\underline{p(x | D_j, \omega_j)}$)

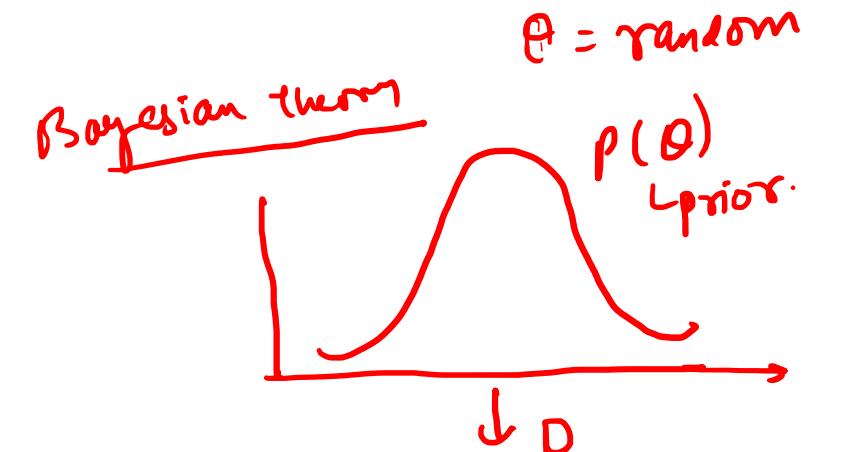
The Multivariate Case

- Univariate case : $p(x|\mathcal{D}) \sim N(\mu_n, \sigma^2 + \sigma_n^2).$
- Multivariate case: $p(\mathbf{x}|\mathcal{D}) \sim N(\boldsymbol{\mu}_n, \boldsymbol{\Sigma} + \boldsymbol{\Sigma}_n)$

The Multivariate Case

- Univariate case : $p(x|D) \sim N(\mu_n, \sigma^2 + \sigma_n^2)$.

- Multivariate case: $p(x|D) \sim N(\mu_n, \Sigma + \Sigma_n)$



$$p(\theta|D) = \text{posterior}$$

Please go & read chapter 3!!!.

Bayesian Parameter Estimation: General Theory

- The form of the density $p(x|\theta)$ is assumed to be known, but the value of the parameter vector θ is not known exactly.
- Our initial knowledge about θ is assumed to be contained in a known a priori density $p(\theta)$.
- The rest of our knowledge about θ is contained in a set D of n samples x_1, \dots, x_n drawn independently according to the unknown probability density $p(x)$.

$p(x|\theta) = \text{normal}$



$p(\theta|D)$

Bayesian Parameter Estimation: General Theory

- Compute posterior density $p(\theta | D)$

By independence assumption:

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

$$p(\mathcal{D}|\theta) = \prod_{k=1}^n p(\mathbf{x}_k|\theta).$$

- then $p(x | D)$ using

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

Bayesian Parameter Estimation: General Theory

- Compute posterior density $p(\theta | D)$

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

- then $p(x | D)$ using

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

