

***Parametric Density Estimation:***

***Bayesian Estimation***

**C7**

# *Bayesian Parameter Estimation*

---

- Suppose we have some idea of the range where parameters  $\theta$  should be
  - Shouldn't we formalize such prior knowledge in hopes that it will lead to better parameter estimation?
- Let  $\theta$  be a random variable with prior distribution  $\mathbf{P}(\theta)$ 
  - This is the key difference between ML and Bayesian parameter estimation
  - This key assumption allows us to fully exploit the information provided by the data

# ***Bayesian Parameter Estimation***

---

- $\theta$  is a random variable with prior  $p(\theta)$ 
  - Unlike MLE case,  $p(x|\theta)$  is a conditional density
- The training data  $D$  allow us to convert  $p(\theta)$  to a posterior probability density  $p(\theta|D)$  .
  - After we observe the data  $D$ , using Bayes rule we can compute the posterior  $p(\theta|D)$
- But  $\theta$  is not our final goal, our final goal is the unknown  $p(\mathbf{x})$
- Therefore a better thing to do is to maximize  $p(\mathbf{x}/D)$ , this is as close as we can come to the unknown  $p(\mathbf{x})$  !

## Bayesian Estimation: Formula for $p(\mathbf{x}|D)$

- From the definition of joint distribution:

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

- Using the definition of conditional probability:

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

- But  $p(\mathbf{x}|\theta, D) = p(\mathbf{x}|\theta)$  since  $p(\mathbf{x}|\theta)$  is completely specified by  $\theta$

$$p(\mathbf{x} | D) = \int \overset{\text{known}}{p(\mathbf{x} | \theta)} \overset{\text{unknown}}{p(\theta | D)} d\theta$$

- Using Bayes formula,

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

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

# ***Bayesian Estimation vs. MLE***

---

- So in principle  $p(\mathbf{x}/D)$  can be computed
  - In practice, it may be hard to do integration analytically, may have to resort to numerical methods

$$p(\mathbf{x} / D) = \int p(\mathbf{x} / \theta) \frac{\prod_{k=1}^n p(\mathbf{x}_k / \theta) p(\theta)}{\int \prod_{k=1}^n p(\mathbf{x}_k / \theta) p(\theta) d\theta} d\theta$$

- Contrast this with the MLE solution which requires differentiation of likelihood to get  $p(\mathbf{x} / \hat{\theta})$ 
  - Differentiation is easy and can always be done analytically

# *Bayesian Estimation vs. MLE*

---

*support  $\theta$  receives  
from the data*

$$p(\mathbf{x} | D) = \int \underbrace{p(\mathbf{x} | \theta)}_{\text{proposed model with certain } \theta} \underbrace{p(\theta | D)}_{\text{support } \theta \text{ receives from the data}} d\theta$$

*proposed model  
with certain  $\theta$*

- The above equation implies that if we are less certain about the exact value of  $\theta$ , we should consider a weighted average of  $p(\mathbf{x}|\theta)$  over the possible values of  $\theta$ .
- Contrast this with the MLE solution which always gives us a single model:

$$p(\mathbf{x} | \hat{\theta})$$

# Bayesian Estimation for Gaussian with unknown $\mu$

- Let  $p(\mathbf{x} | \mu)$  be  $\mathbf{N}(\mu, \sigma^2)$  that is  $\sigma^2$  is known, but  $\mu$  is unknown and needs to be estimated, so  $\theta = \mu$
- Assume a prior over  $\mu$ :  $p(\mu) \sim N(\mu_0, \sigma_0^2)$
- $\mu_0$  encodes some prior knowledge about the true mean  $\mu$ , while  $\sigma_0^2$  measures our prior uncertainty.
- The posterior distribution is:

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

# Bayesian Estimation for Gaussian with unknown $\mu$

- Where factors that do not depend on  $\mu$  have been absorbed into the constants  $\alpha'$  and  $\alpha''$
- $p(\mu | D)$  is an exponent of a quadratic function of  $\mu$  i.e. it is a normal density.
- $p(\mu | D)$  remains normal for any number of training samples.

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

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

# Bayesian Estimation for Gaussian with unknown $\mu$

- then identifying the coefficients, we get

$$\frac{1}{\sigma_n^2} = \frac{n}{\sigma^2} + \frac{1}{\sigma_0^2} \qquad \frac{\mu_n}{\sigma_n^2} = \frac{n}{\sigma^2} \hat{\mu}_n + \frac{\mu_0}{\sigma_0^2}$$

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

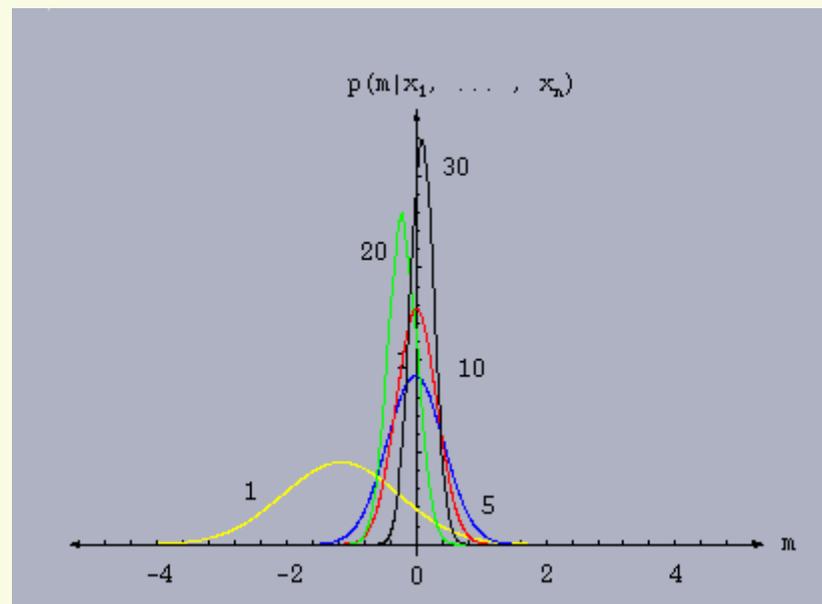
- Solving explicitly for  $\mu_n$  and  $\sigma_n^2$  we obtain:

$$\mu_n = \left( \frac{n\sigma_0^2}{n\sigma_0^2 + \sigma^2} \right) \hat{\mu}_n + \frac{\sigma^2}{n\sigma_0^2 + \sigma^2} \mu_0 \quad \text{our best guess after observing } n \text{ samples}$$

$$\sigma_n^2 = \frac{\sigma_0^2 \sigma^2}{n\sigma_0^2 + \sigma^2} \quad \text{uncertainty about the guess, decreases monotonically with } n$$

# Bayesian Estimation for Gaussian with unknown $\mu$

- Each additional observation decreases our uncertainty about the true value of  $\mu$ .
- As  $n$  increases,  $p(\mu | D)$  becomes more and more sharply peaked, approaching a Dirac delta function as  $n$  approaches infinity. This behavior is known as Bayesian Learning.



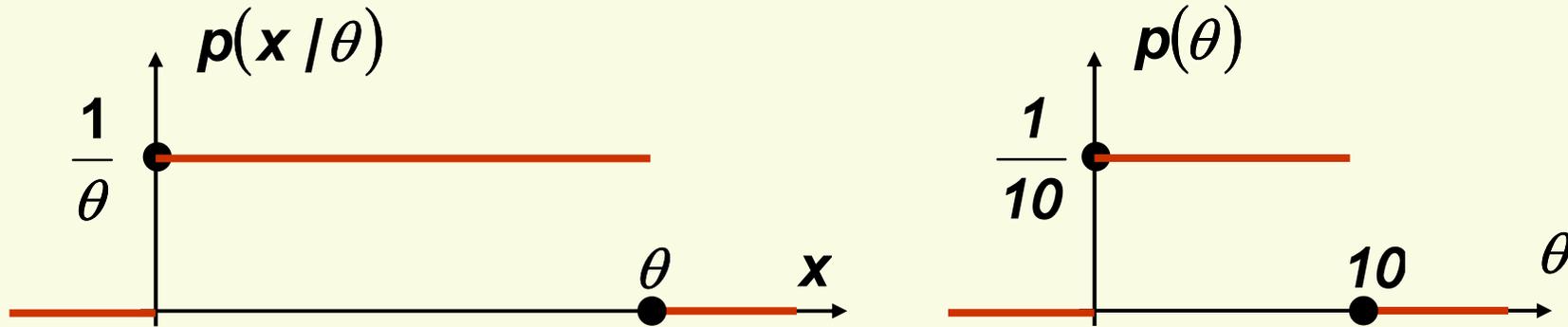
# Bayesian Estimation for Gaussian with unknown $\mu$

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

- In general,  $\mu_n$  is a linear combination of  $\hat{\mu}_n$  and  $\mu_0$ , with coefficients that are non-negative and sum to 1.
- Thus  $\mu_n$  lies somewhere between  $\hat{\mu}_n$  and  $\mu_0$ .
- If  $\sigma_0 \neq 0$ ,  $\mu_n \rightarrow \hat{\mu}_n$  as  $n \rightarrow \infty$
- If  $\sigma_0 = 0$ , our a priori certainty that  $\mu = \mu_0$  is so strong that no number of observations can change our opinion.
- If  $\sigma_0 \approx \sigma$ , a priori guess is very uncertain, and we take  $\mu_n = \hat{\mu}_n$

## Bayesian Estimation: Example for $U[0, \theta]$

- Let  $X$  be  $U[0, \theta]$ . Recall  $p(x|\theta) = 1/\theta$  inside  $[0, \theta]$ , else 0



- Suppose we assume a  $U[0, 10]$  prior on  $\theta$ 
  - good prior to use if we just know the range of  $\theta$  but don't know anything else

## Bayesian Estimation: Example for $U[0, \theta]$

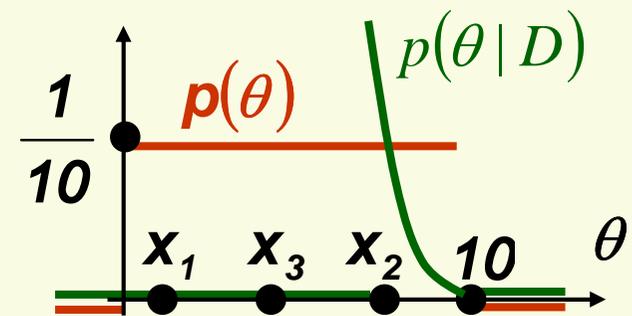
- We need to compute  $p(\mathbf{x} | \mathbf{D}) = \int p(\mathbf{x} | \theta) p(\theta | \mathbf{D}) d\theta$
- using  $p(\theta | \mathbf{D}) = \frac{p(\mathbf{D} | \theta) p(\theta)}{\int p(\mathbf{D} | \theta) p(\theta) d\theta}$  and  $p(\mathbf{D} | \theta) = \prod_{k=1}^n p(x_k | \theta)$

- When computing MLE of  $\theta$ , we had

$$p(\mathbf{D} | \theta) = \begin{cases} \frac{1}{\theta^n} & \text{for } \theta \geq \max\{x_1, \dots, x_n\} \\ 0 & \text{otherwise} \end{cases}$$

- Thus

$$p(\theta | \mathbf{D}) = \begin{cases} c \frac{1}{\theta^n} & \text{for } \max\{x_1, \dots, x_n\} \leq \theta \leq 10 \\ 0 & \text{otherwise} \end{cases}$$



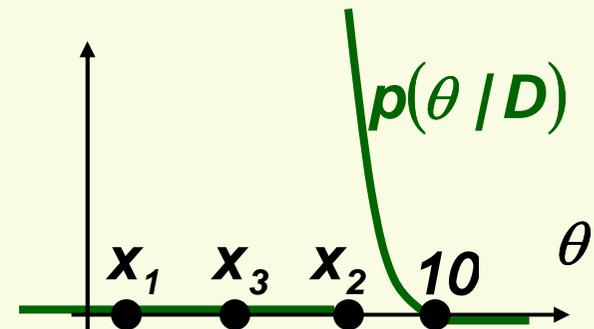
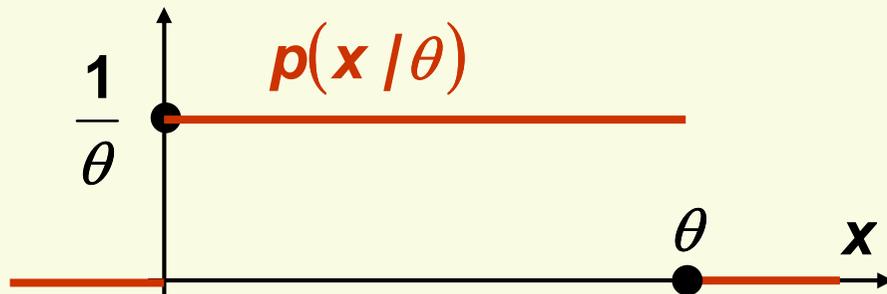
- where  $c$  is the normalizing constant, i.e.

$$c = \frac{1}{\int_{\max\{x_1, \dots, x_n\}}^{10} \frac{d\theta}{\theta^n}}$$

# Bayesian Estimation: Example for $U[0, \theta]$

- We need to compute  $p(\mathbf{x} | D) = \int p(\mathbf{x} | \theta) p(\theta | D) d\theta$

$$p(\theta | D) = \begin{cases} c \frac{1}{\theta^n} & \text{for } \max\{x_1, \dots, x_n\} \leq \theta \leq 10 \\ 0 & \text{otherwise} \end{cases}$$



- We have 2 cases:
- case  $x < \max\{x_1, x_2, \dots, x_n\}$

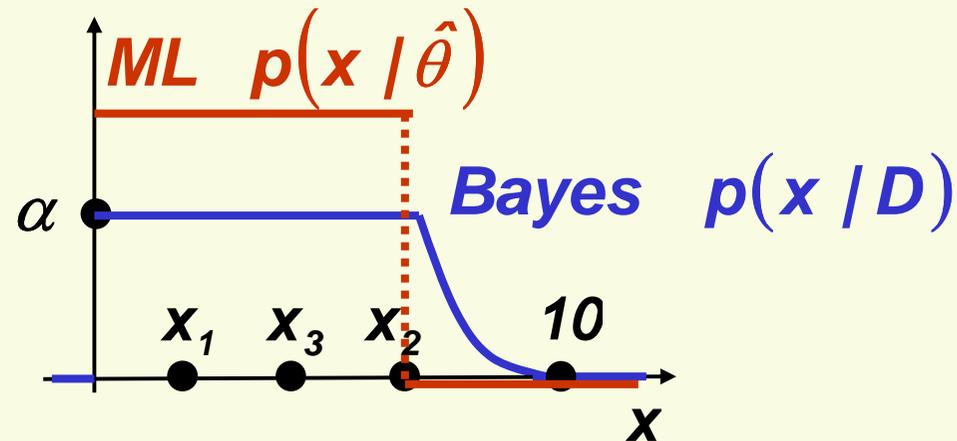
$$p(\mathbf{x} | D) = \int_{\max\{x_1, \dots, x_n\}}^{10} c \frac{1}{\theta^{n+1}} d\theta = \boxed{\alpha}$$

*constant independent of x*

- case  $x > \max\{x_1, x_2, \dots, x_n\}$

$$p(\mathbf{x} | D) = \int_x^{10} c \frac{1}{\theta^{n+1}} d\theta = \left. -\frac{c}{n\theta^n} \right|_x^{10} = \boxed{\frac{c}{nx^n}} - \frac{c}{n10^n}$$

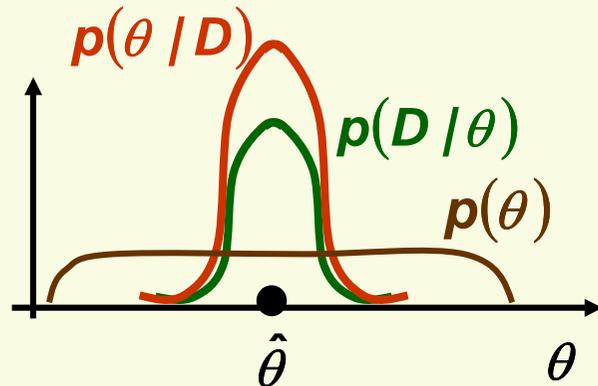
## Bayesian Estimation: Example for $U[0, \theta]$



- Note that even after  $x > \max \{x_1, x_2, \dots, x_n\}$ , Bayes density is not zero, which makes sense
- curious fact: Bayes density is not uniform, i.e. does not have the functional form that we have assumed!

## ML vs. Bayesian Estimation with Broad Prior

- Suppose  $p(\theta)$  is flat and broad (close to uniform prior)
- $p(\theta|D)$  tends to sharpen if there is a lot of data



- Thus  $p(D|\theta) \propto p(\theta|D)p(\theta)$  will have the same sharp peak as  $p(\theta|D)$
- But by definition, peak of  $p(D|\theta)$  is the ML estimate  $\hat{\theta}$
- The integral is dominated by the peak:  
$$p(x|D) = \int p(x|\theta)p(\theta|D)d\theta \approx p(x|\hat{\theta}) \int p(\theta|D)d\theta = p(x|\hat{\theta})$$
- Thus as  $n$  goes to infinity, Bayesian estimate will approach the density corresponding to the MLE!

# *ML vs. Bayesian Estimation*

- **Number of training data**
  - The two methods are equivalent assuming infinite number of training data (and prior distributions that do not exclude the true solution).
  - For small training data sets, they give different results in most cases.
- **Computational complexity**
  - ML uses differential calculus or gradient search for maximizing the likelihood.
  - Bayesian estimation requires complex multidimensional integration techniques.

# *ML vs. Bayesian Estimation*

- Solution complexity
  - Easier to interpret ML solutions (i.e., must be of the assumed parametric form).
  - A Bayesian estimation solution might not be of the parametric form assumed. Hard to interpret, returns weighted average of models.
- Prior distribution
  - If the prior distribution  $p(\theta)$  is uniform, Bayesian estimation solutions are equivalent to ML solutions.

# *ML vs. Bayesian Estimation*

- Broad or asymmetric  $p(\theta/D)$ 
  - In this case, the two methods will give different solutions.
  - Bayesian methods will explicitly exploit such information.
- General comments
  - There are strong theoretical and methodological arguments supporting Bayesian estimation.
  - In practice, ML estimation is simpler and can lead to comparable performance.