

Recall from MAP estimation example:

$$f(x|\mu) = \mathcal{N}(\mu, \sigma^2)$$

$\uparrow$   
known

$D = (x_1, \dots, x_n)$  observed samples

$$f(\mu) = \mathcal{N}(\mu_0, \sigma_0^2)$$

$\uparrow \quad \nearrow$   
known

$$\underbrace{f(\mu|D)}_{\text{posterior}} = \frac{f(D|\mu) f(\mu)}{f(D)} \propto f(D|\mu) f(\mu)$$

$$= \prod_{k=1}^n f(x_k|\mu) f(\mu)$$
$$= \prod_{k=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x_k - \mu)^2}{2\sigma^2}} \cdot \frac{1}{\sqrt{2\pi\sigma_0^2}} e^{-\frac{(\mu - \mu_0)^2}{2\sigma_0^2}}$$

$$\Rightarrow f(\mu|D) \propto e^{-\frac{1}{2\sigma^2} \sum_{k=1}^n (x_k - \mu)^2 - \frac{1}{2\sigma_0^2} (\mu - \mu_0)^2}$$

$$\propto e^{-\frac{1}{2\sigma^2} \sum_k (x_k^2 - 2x_k\mu + \mu^2) - \frac{1}{2\sigma_0^2} (\mu^2 - 2\mu\mu_0 + \mu_0^2)}$$

$$\propto e^{-\frac{1}{2} \left\{ \underbrace{\left( \frac{n}{\sigma^2} + \frac{1}{\sigma_0^2} \right)}_{\frac{1}{\sigma_n^2}} \mu^2 - 2 \underbrace{\left( \frac{1}{\sigma^2} \sum_{k=1}^n x_k + \frac{\mu_0}{\sigma_0^2} \right)}_{\frac{\mu_n}{\sigma_n^2}} \mu \right\}}$$
$$\frac{\mu_n}{\sigma_n^2} = \frac{n}{\sigma^2} \hat{\mu}_n + \frac{\mu_0}{\sigma_0^2}$$

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

sample mean

$$\left\{ \begin{array}{l} \frac{1}{\sigma_n^2} = \frac{n}{\sigma^2} + \frac{1}{\sigma_0^2} \quad \text{or} \quad \sigma_n^2 = \frac{\sigma_0^2 \sigma^2}{\sigma^2 + \sigma_0^2 n} \\ \frac{\mu_n}{\sigma_n^2} = \frac{n}{\sigma^2} \hat{\mu}_n + \frac{\mu_0}{\sigma_0^2} \quad \text{or} \quad \mu_n = \left( \frac{n\sigma_0^2}{n\sigma_0^2 + \sigma^2} \right) \hat{\mu}_n + \frac{\sigma^2}{n\sigma_0^2 + \sigma^2} \mu_0 \end{array} \right.$$

note  $f(\mu|D)$  is exponential function of a quadratic

function of  $\mu$ : normal density

$f(\mu|D)$  is called reproducing density

$f(\mu)$  is called conjugate prior.

Write  $f(\mu|D) = \mathcal{N}(\mu_n, \sigma_n^2)$

$$= \frac{1}{\sqrt{2\pi}\sigma_n} \exp \left[ -\frac{1}{2} \left( \frac{\mu - \mu_0}{\sigma_n} \right)^2 \right]$$

i.e.  $f(\mu|D) = \mathcal{N}(\hat{\mu}_{\text{MAP}}, \sigma_n^2)$

- $\mu_n = \hat{\mu}_{\text{MAP}}$  represents best guess for  $\mu$  after observing  $n$  samples
- $\sigma_n^2$  measures uncertainty about this guess
- $\sigma_n^2$  decreases monotonically with  $n$ ,  $\rightarrow 0$  as  $\sigma^2/n$   
each additional observation decreases uncertainty about true value of  $\mu$
- As  $n$  increases,  $f(\mu|D)$  becomes more and more sharply peaked (see figure).

### Bayesian parameter estimation

- Want to estimate density  $f(\underline{x}|D)$
- Assume  $f(\underline{x}|\underline{\theta})$  known  
prior distribution on  $\underline{\theta}$ :  $f(\underline{\theta})$   
a posteriori dist. on  $\underline{\theta}$ :  $f(\underline{\theta}|D)$
- $$f(\underline{x}|D) = \int f(\underline{x}, \underline{\theta}|D) d\underline{\theta}$$
$$= \int f(\underline{x}|\underline{\theta}, D) f(\underline{\theta}|D) d\underline{\theta}$$

- Since  $\underline{x}$  and  $\underline{\theta}$  selected independently,

$$f(\underline{x} | \underline{\theta}, \underline{D}) = f(\underline{x} | \underline{\theta})$$

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

- If  $f(\underline{\theta} | \underline{D})$  peaks sharply about  $\hat{\underline{\theta}}$ , then  $f(\underline{x} | \underline{D}) \sim f(\underline{x} | \hat{\underline{\theta}})$

- But there is usual uncertainty about exact value of  $\underline{\theta}$ , so average  $f(\underline{x} | \underline{\theta})$  over  $\underline{\theta}$ .

Back to Gaussian 1-D example:

Assume  $f(x|\mu) = \mathcal{N}(\mu, \sigma^2)$   
└ known

Assume  $\mu$  is a r.v.  $\sim f(\mu) = \mathcal{N}(\mu_0, \sigma_0^2)$   
┆ ┆  
best prior guess for  $\mu$  uncertainty about prior guess.

Observed samples  $\underline{D} = \{x_1, \dots, x_n\}$

$$f(x|\underline{D}) = \int f(x|\mu) \underbrace{f(\mu|\underline{D})}_{\text{posterior}} d\mu$$

$$= \int \mathcal{N}(x|\mu, \sigma^2) \mathcal{N}(\mu|\mu_n, \sigma_n^2) d\mu$$

$$= \mathcal{N}(x|\mu_n, \sigma^2 + \sigma_n^2)$$

$$\text{i.e. } f(x|\underline{D}) = \int \frac{1}{\sqrt{2\pi}\sigma^2} \exp\left[-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2\right] \cdot$$

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

$$\propto \frac{1}{2\pi\sigma\sigma_n} \exp\left[-\frac{1}{2} \frac{(x-\mu_n)^2}{\sigma^2+\sigma_n^2}\right]$$

$$\Rightarrow f(x|D) \stackrel{=}{=} \mathcal{N}(\mu_n, \sigma^2+\sigma_n^2)$$

- $f(x|\mu) = \mathcal{N}(\mu, \sigma^2) \longrightarrow f(x|D) = \mathcal{N}(\mu_n, \sigma^2+\sigma_n^2)$
- Conditional mean  $\mu_n$  treated as true mean
- Variance increased to account for additional uncertainty in  $x$  from lack of exact knowledge about  $\mu$ . ( $\sigma^2 \rightarrow \sigma^2+\sigma_n^2$ )
- Contrast Bayesian approach with ML methods which only estimate  $\hat{\mu}$  and  $\hat{\sigma}^2$ , rather than estimate a distribution  $p(x|D)$ .

### Multivariate case

$$f(\underline{x}|\underline{\mu}) = \underbrace{\mathcal{N}(\underline{\mu}, \underline{\Sigma})}_{\substack{\text{unknown} \\ \text{known}}} = \frac{1}{(2\pi)^{d/2} |\underline{\Sigma}|^{1/2}} e^{-\frac{1}{2}(\underline{x}-\underline{\mu})^T \underline{\Sigma}^{-1}(\underline{x}-\underline{\mu})}$$

$$f(\underline{\mu}) = \underbrace{\mathcal{N}(\underline{\mu}_0, \underline{\Sigma}_0)}_{\text{known}} = \frac{1}{(2\pi)^{d/2} |\underline{\Sigma}_0|^{1/2}} e^{-\frac{1}{2}(\underline{\mu}-\underline{\mu}_0)^T \underline{\Sigma}_0^{-1}(\underline{\mu}-\underline{\mu}_0)}$$

$$D = \{\underline{x}_1, \dots, \underline{x}_n\}$$

$$f(\underline{\mu}|D) \propto \prod_{k=1}^n f(\underline{x}_k|\underline{\mu}) f(\underline{\mu})$$

$$\propto \exp\left[-\frac{1}{2}(\underline{\mu}^T (n\underline{\Sigma}^{-1} + \underline{\Sigma}_0^{-1})\underline{\mu} - 2\underline{\mu}^T (\underline{\Sigma}^{-1} \sum_{k=1}^n \underline{x}_k + \underline{\Sigma}_0^{-1} \underline{\mu}_0))\right]$$

$$\propto \exp \left[ -\frac{1}{2} (\underline{\mu} - \underline{\mu}_n)^T \Sigma_n^{-1} (\underline{\mu} - \underline{\mu}_n) \right]$$

$$= n(\underline{\mu}_n, \Sigma_n)$$

equating coefficients,

$$\Sigma_n^{-1} = n \Sigma^{-1} + \Sigma_0^{-1}$$

$$\Sigma_n^{-1} \underline{\mu}_n = n \Sigma^{-1} \hat{\underline{\mu}}_n - \Sigma_0^{-1} \underline{\mu}_0$$

where  $\hat{\underline{\mu}}_n = \frac{1}{n} \sum_{k=1}^n \underline{x}_k$  is sample mean =  $\hat{\underline{\mu}}_{ML}$

Now  $(A^{-1} + B^{-1})^{-1} = A(A+B)^{-1}B = B(A+B)^{-1}A$

for any nonsingular  $d \times d$  matrices  $A, B$ .

$$\Sigma_n = \Sigma_0 (\Sigma_0 + \frac{1}{n} \Sigma)^{-1} \frac{1}{n} \Sigma$$

$$\underline{\mu}_n = \Sigma_0 (\Sigma_0 + \frac{1}{n} \Sigma)^{-1} \hat{\underline{\mu}}_n + \frac{1}{n} \Sigma (\Sigma_0 + \frac{1}{n} \Sigma)^{-1} \underline{\mu}_0$$

linear combination of  $\hat{\underline{\mu}}_n$  and  $\underline{\mu}_0$

Now  $f(\underline{x} | D) = \int f(\underline{x} | \underline{\mu}) f(\underline{\mu} | D) d\underline{\mu}$

observe  $\underline{x} = \underline{\mu} + \underline{z}$   
(given  $D$ )

where  $\underline{\mu} \sim n(\underline{\mu}_n, \Sigma_n)$

$\underline{z} \sim n(\underline{0}, \Sigma)$  indep of  $\underline{\mu}$

Thus  $f(\underline{x} | D) = n(\underline{\mu}_n, \Sigma_n + \Sigma)$

## Bernoulli Distribution

66

$X \in \{0, 1\}$   $X=1$  heads,  $X=0$  tails

$$p(X=1|\mu) = \mu \quad 0 \leq \mu \leq 1$$

$$p(X=0|\mu) = 1 - \mu$$

$$E[X] = \mu, \quad \text{Var}(X) = \mu(1 - \mu)$$

~~...~~

$$p(x_1, \dots, x_n | \mu) = \mu^{\sum x_i} (1 - \mu)^{n - \sum x_i}$$

Given  $D = \{x_1, \dots, x_n\}$

$x_k$  iid  $\sim \text{Bern}(\mu)$

we showed

$$\hat{\mu}_{ML} = \frac{1}{n} \sum_{k=1}^n x_k \equiv m$$

## Beta distribution

normalization

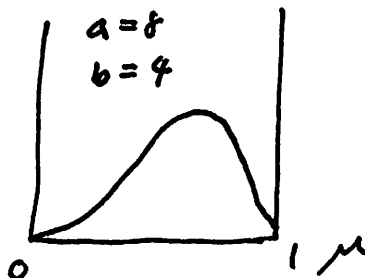
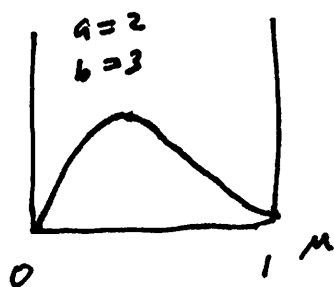
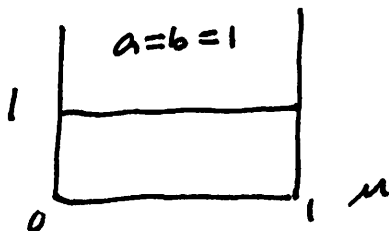
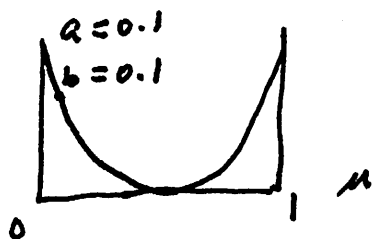
$$\text{Beta}(\mu | a, b) = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \mu^{a-1} (1-\mu)^{b-1} \quad \mu \text{ a r.v.}$$

hyperparameters  $\mu \in [0, 1]$

where  $\Gamma(x) = \int_0^\infty u^{x-1} e^{-u} du$  gamma function

$$E[\mu] = \frac{a}{a+b}$$

$$\text{Var}(\mu) = \frac{ab}{(a+b)^2(a+b+1)}$$



$$f(\mu|D) \propto p(D|\mu)f(\mu)$$

$$\propto \left( \prod_{k=1}^n \mu^{x_k} (1-\mu)^{1-x_k} \right) \text{Beta}(\mu|a, b)$$

$$\propto \mu^{\sum x_k} (1-\mu)^{n - \sum x_k} \mu^{a-1} (1-\mu)^{b-1}$$

$$\propto \mu^m (1-\mu)^{n-m} \mu^{a-1} (1-\mu)^{b-1} \quad m = \sum_{k=1}^n x_k$$

$$\propto \mu^{m+a-1} (1-\mu)^{n-m+b-1}$$

$$\Rightarrow f(\mu|D) = \text{Beta}(\mu|m+a, n-m+b)$$

a: pseudo count for  $x=1$

b: " for  $x=0$

$$\text{Now } P(X=1|D) = \int_0^1 P(X=1|\mu) f(\mu|D) d\mu$$

$$= \int_0^1 \mu f(\mu|D) d\mu$$

$$= E[\mu|D] \quad \text{where } f(\mu|D)$$

$$= \text{Beta}(\mu|m+a, n-m+b)$$

$$= \frac{m+a}{m+a+n-m+b}$$

$$= \frac{m+a}{n+a+b}$$

$$P(X=0|D) = \frac{n-m+b}{n+a+b}$$

$$\text{Given } D, \quad X \sim \text{Bernoulli}\left(\frac{m+a}{n+a+b}\right)$$

~~or a priori~~,

Recap:

$$D = \{x_1, \dots, x_n\}$$

(68)

ML learning:  $\underline{\theta}$  fixed. no prior over  $\underline{\theta}$

Find parameter which maximizes likelihood of data:

$$\hat{\underline{\theta}}_{ML} = \underset{\underline{\theta}}{\operatorname{argmax}} f(D|\underline{\theta}) \quad \left( f(D|\underline{\theta}) = \prod_{k=1}^n f(x_k|\underline{\theta}) \right)$$

MAP Learning:

by independence

$\underline{\theta}$  is a r.v.  $\sim f(\underline{\theta})$ .

Find parameter to maximize posterior.

$$\begin{aligned} \hat{\underline{\theta}}_{MAP} &= \underset{\underline{\theta}}{\operatorname{argmax}} f(\underline{\theta}|D) \\ &= \underset{\underline{\theta}}{\operatorname{argmax}} f(D|\underline{\theta}) f(\underline{\theta}) \end{aligned}$$

Bayesian learning

$\underline{\theta}$  is a r.v.  $\sim f(\underline{\theta})$

compute posterior dist<sup>n</sup> of  $\underline{\theta} = f(\underline{\theta}|D)$

$$\begin{aligned} \text{Then estimate } f(x|D) &= \int f(x, \underline{\theta}|D) d\underline{\theta} \\ &= \int f(x|\underline{\theta}, D) f(\underline{\theta}|D) d\underline{\theta} \\ &= \int f(x|\underline{\theta}) f(\underline{\theta}|D) d\underline{\theta} \end{aligned}$$

Recursive Bayes incremental learning

Let  $D^n = \{x_1, \dots, x_n\}$  if  $n > 1$ ,

$$\underbrace{f(D^n|\underline{\theta})}_{\text{likelihood}} = \prod_{k=1}^n f(x_k|\underline{\theta}) = f(x_n|\underline{\theta}) \prod_{k=1}^{n-1} f(x_k|\underline{\theta})$$
$$= f(x_n|\underline{\theta}) f(D^{n-1}|\underline{\theta})$$



where  $D^{n-1} = \{x_1, \dots, x_{n-1}\}$

$$\begin{aligned}
 f(\theta | D^n) &= \frac{f(D^n | \theta) f(\theta)}{\int f(D^n | \theta) f(\theta) d\theta} = \frac{f(x_n | \theta) f(D^{n-1} | \theta) f(\theta)}{\int f(x_n | \theta) f(D^{n-1} | \theta) f(\theta) d\theta} \\
 &= \frac{f(x_n | \theta) f(\theta | D^{n-1}) f(D^{n-1})}{\int f(x_n | \theta) f(\theta | D^{n-1}) f(D^{n-1}) d\theta} \quad (1)
 \end{aligned}$$

- Repeated use gives sequence:

$$f(\theta | D^0) \equiv f(\theta), f(\theta | x_1), f(\theta | x_1, x_2), \dots$$

- Given  $f(\theta | D^{n-1})$ , obtain  $f(\theta | D^n)$  with  $f(x_n | \theta)$  and (1)

- Recursive approach.

- on-line learning: incremental, learning goes on as data is collected

- Batch: non-incremental, all training data must be present before learning takes place. Uses all training data for learning.

Example: Assume

70

$$f(x|\theta) = u(0, \theta) = \begin{cases} 1/\theta, & 0 \leq x \leq \theta \\ 0, & \text{o/w.} \end{cases}$$

and assume  $f(\theta|D^0) = f(\theta) = u(0, 10)$

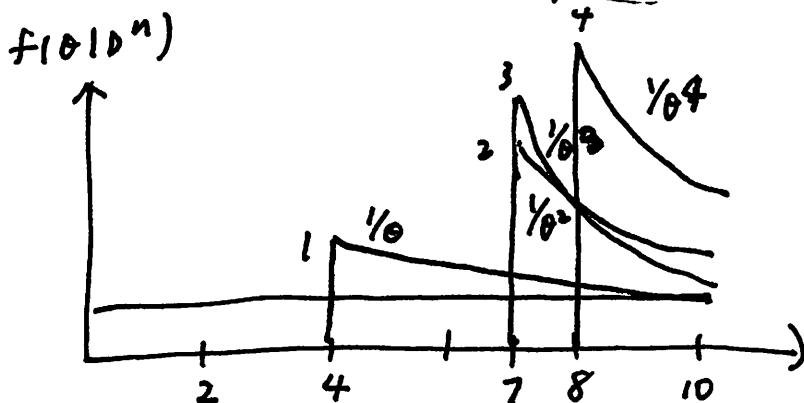
$$D = \{4, 7, 2, 8\}$$

$$f(\theta|D^1) \propto f(x_1|\theta) f(\theta|D^0) \\ \propto \begin{cases} \frac{1}{\theta}, & 4 \leq \theta \leq 10 \\ 0, & \text{o/w.} \end{cases}$$

$$f(\theta|D^2) \propto \underbrace{f(x_2|\theta)}_{\begin{cases} \frac{1}{\theta}, & 7 \leq \theta \leq 10 \\ 0, & \text{o/w} \end{cases}} \underbrace{f(\theta|D^1)}_{\begin{cases} \frac{1}{\theta}, & 4 \leq \theta \leq 10 \\ 0, & \text{o/w} \end{cases}} \\ \propto \begin{cases} \frac{1}{\theta^2}, & 7 \leq \theta \leq 10 \\ 0, & \text{o/w.} \end{cases}$$

:

$$f(\theta|D^n) \propto \begin{cases} \frac{1}{\theta^n}, & \max\{D^n\} \leq \theta \leq 10 \\ 0, & \text{o/w.} \end{cases}$$

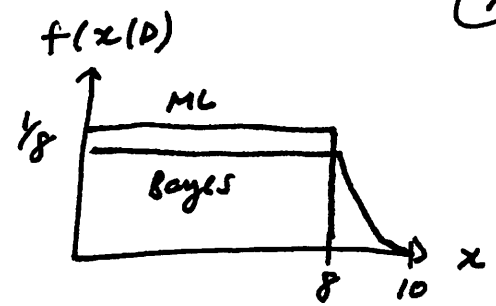


$$f(\theta|D^4) \propto \begin{cases} \frac{1}{\theta^4}, & 8 \leq \theta \leq 10 \\ 0, & \text{o/w.} \end{cases}$$

Given  $D = \{4, 7, 2, 8\}$

$$\hat{\theta}_{ML} = \arg \max_{\theta} f(D|\theta) = 8$$

$$\Rightarrow f(x|D) \sim u(0, 8)$$



Bayesian methodology:

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

has tail above  $x=8$ : prior influence remains