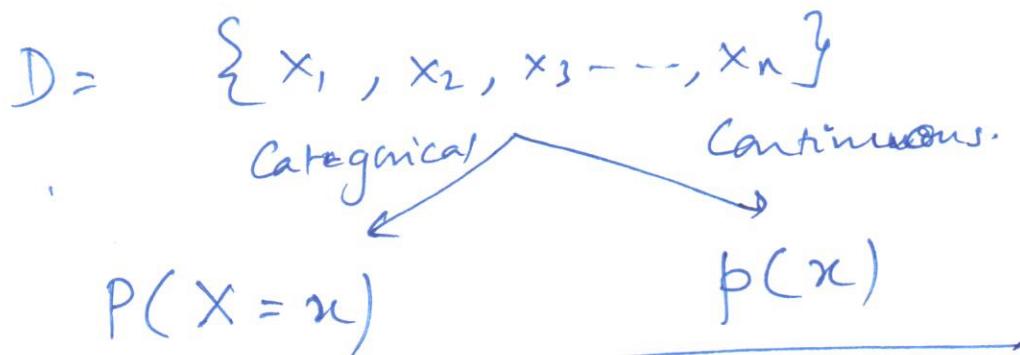


Likelihood



let x be continuous

$f(x_1) \rightarrow$ pdf value at x_1

$f(x_2) \rightarrow$ pdf $\dots x_2$

\vdots

$f(x_n) \rightarrow \dots$

$$p(x_1, x_2, x_3, \dots, x_n) = \frac{p(x_1) p(x_2) \dots p(x_n)}{\text{assuming independence}}$$

i.i.d

independent & identically distributed

$$= \prod_{i=1}^n p(x_i)$$

$$p(D) = \prod_{i=1}^n p(x_i)$$

Likelihood.

$$\sum_{i=1}^n \log(p(x_i))$$

$$\log p(D) =$$

log-likelihood

$$D = \{1, 4, 16, 64\}$$

What is likelihood of D given h .

$$\begin{aligned} l(D|h) &= \prod_{i=1}^4 p(x_i|h) \\ &= p(1|h) p(4|h) p(16|h) p(64|h) \end{aligned}$$

$h \rightarrow$ all even numbers.

$$p(1| \text{even numbers}) = 0 \quad \text{∅} \quad l(D|h) = 0$$

$h \rightarrow$ odd numbers. $\therefore l(D|h) = 0$

$h \rightarrow$ squares.

$$p(1|\text{squares}) = \frac{1}{10}$$

$$p(4|\text{squares}) = \frac{1}{10}$$

$$p(16|\text{squares}) = \frac{1}{10}$$

$$p(64|\text{squares}) = \frac{1}{10}$$

$1, 4, 16, 25, \dots, 100$
10 possible "squares"

$$l(D|\text{squares}) = \frac{1}{10} * \frac{1}{10} * \frac{1}{10} * \frac{1}{10} = 10^{-4}$$

$h \rightarrow$ Powers of 4

$1, 4, 16, 64$

$$p(1|\text{powers of 4}) = \frac{1}{4}$$

$$p(64|\text{powers of 4}) = \frac{1}{4}$$

$$l(D|\text{powers of 4}) = \frac{1}{4} * \frac{1}{4} * \frac{1}{4} * \frac{1}{4} = \underline{\underline{3.906 \cdot 10^{-3}}}$$

$h \rightarrow$ all numbers between 1 and 100

$$p(1 | \text{all numbers}) = \frac{1}{100}$$

$$l(D | \text{all numbers}) = \frac{1}{100} * \frac{1}{100} * \frac{1}{100} * \frac{1}{100} = 10^{-8}$$

Prior $p(h)$

Posterior $p(h|D) = \frac{p(D|h) p(h)}{\sum_{\text{all } h'} p(D|h') p(h')}$

$$\frac{p(\text{powers of } 4 | D)}{0 * 0.075 + 0 * 0.3 + \dots} = \frac{3.906 * 10^{-3} * 0.075}{\dots}$$

Inference D — Training data. $p(h)$ — prior.

new data example x^*

What is the prob that x^* was also generated by the same hypothesis as D

option 0: $p(x^* | \text{prior})$

$$x^* = 4$$
$$p(4 | h) = \frac{1}{100}$$

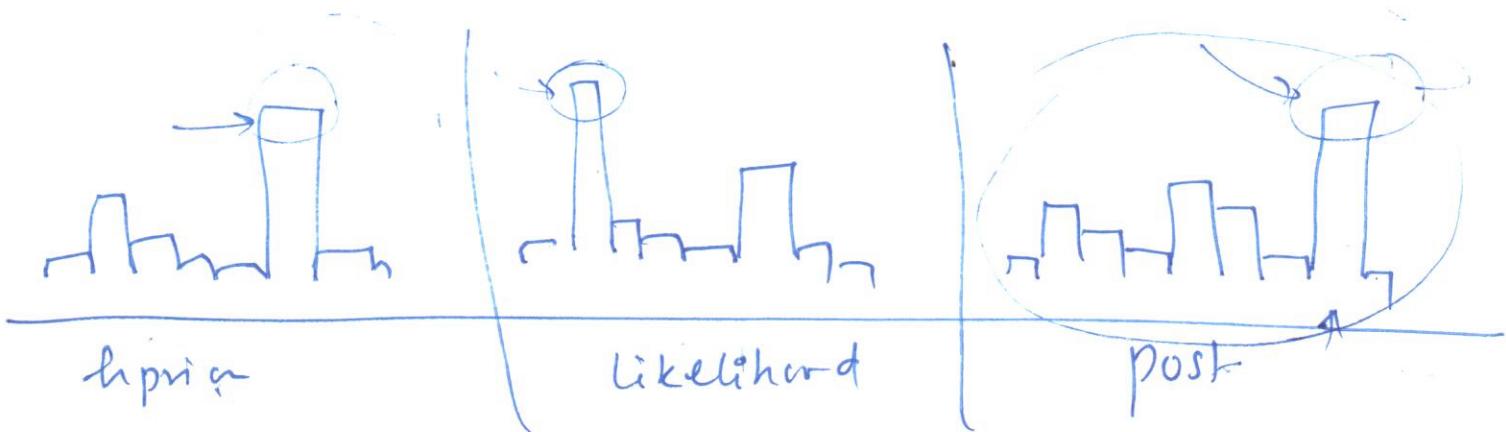
example
let $p(h \rightarrow \text{numbers between } 1 \text{ & } 100)$ be the highest prior.

Option 1 $p(x^* | h_{MLE})$

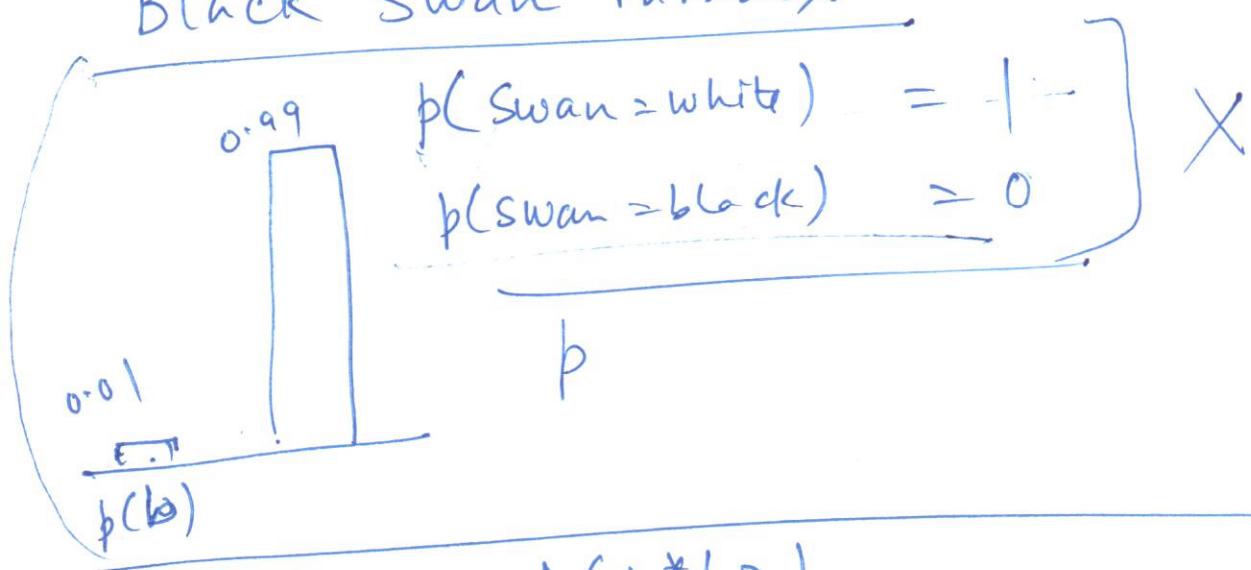
let $h \rightarrow$ powers of 4
have the high likelihood

$$p(x^* | D) = p(x^* | h_{MLE}) = \frac{1}{4}$$

Option 2 $p(x^* | h_{MAP})$

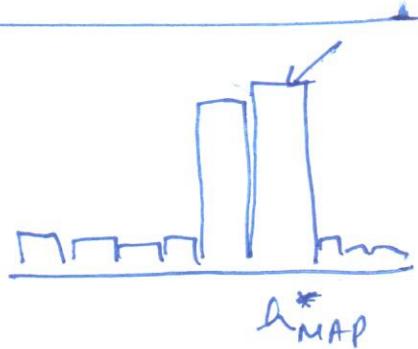


Black Swan Paradox



Option 3: $\underline{p(x^* | D)}$

Bayesian Averaging



Daniel Salgo (F)
Dan Goings (M)

$$\{1, 4, 16, 64\}$$

h_1
 h_2
 h_3
 h_4
 h_5
 \vdots
 h_{10}

$$\{x^*\}$$

coins:
 $\{H, H, H, T, H\} \rightarrow D$



$$P(X^* = \text{heads} | D)$$

$$P(X^* = \text{tails} | D)$$

$$D = \{H, H, H, T, T, H\}$$

$\theta \rightarrow$ Prob. of X to be H

$$\begin{aligned} \text{likelikwd of } D &= P(x_1 | \theta) P(x_2 | \theta) \dots \\ &= \theta * \theta * \theta * (1-\theta) * (1-\theta) * \theta \\ &= \theta^4 (1-\theta)^2 \end{aligned}$$

if D has N_1 heads and N_0 tails.

$$\text{lik. } l(D|\theta) = \theta^{N_1} (1-\theta)^{N_0}$$

Taking logs:

$$\underline{l(l(D|\theta)) = N_1 \log \theta + N_0 \log (1-\theta)}$$

Find θ that maximizes $l(l(D|\theta))$

$$\theta_{\max} \text{ will be : } \frac{d}{d\theta} \ell \ell(D|\theta) = 0$$

$$\begin{aligned} & \frac{N_1}{\theta} + \frac{N_0}{(1-\theta)} (-1) \\ &= \frac{N_1}{\theta} - \frac{N_0}{(1-\theta)} \end{aligned}$$

Setting to 0

$$\frac{N_1}{\theta} = \frac{N_0}{1-\theta} \Rightarrow N_1 - N_1 \theta = N_0 \theta$$

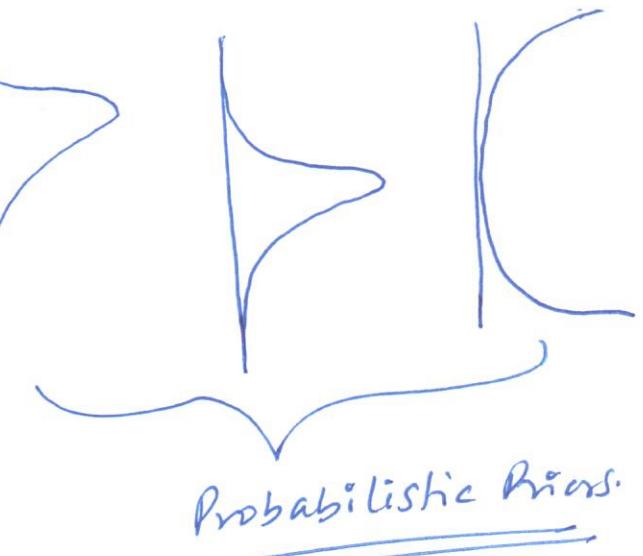
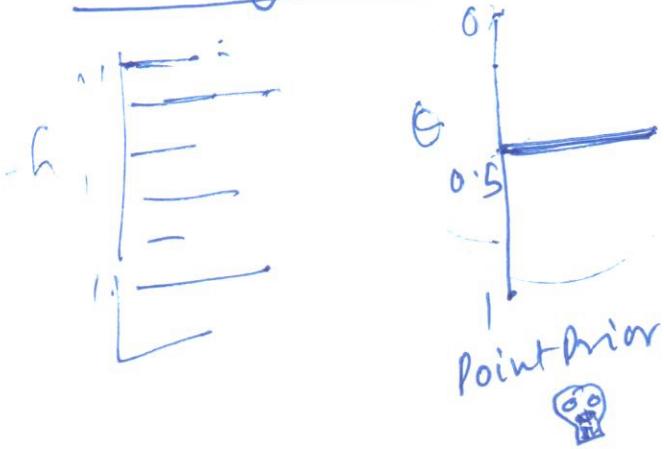
$\theta = \frac{N_1}{N_0 + N_1}$

MLE for θ

$N = \text{Total Size of } D : N_0 + N_1 = N$

$$P(X^* = \text{heads} | D) = \theta_{\text{MLE}} = \frac{N_1}{N_0 + N_1}$$

How do you define a prior?



Use a probability distribution to represent prior.

$\theta \rightarrow$ also a random variable
 $\underline{\theta \in (0,1)}$

Beta Distribution

$a, b.$

$$p(\theta=0.7 | a, b) = \theta^{a-1} (1-\theta)^{b-1}$$

$$\text{If } \theta = 0.7 \quad a = 3 \quad b = 2$$

$$p(0.7 | 3, 2) = (0.7)^2 (1-0.7)^1 \\ = \underline{\hspace{2cm}}$$

$$E[\theta] = \int \theta p(\theta | a, b) d\theta$$

$$= \int \theta \theta^{a-1} (1-\theta)^{b-1} d\theta$$

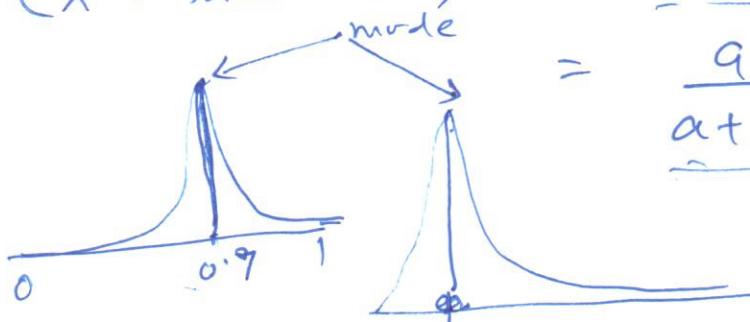
skipping calculus

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

$$\text{Var}[\theta] = E[(\theta - E(\theta))^2] = \frac{ab}{(a+b)^2(a+b+1)}$$

If prior for θ is a Beta(a, b)

Ignoring D:

$$P(X^* = \text{heads}) = \frac{\theta \text{ that maximizes pdf.}}{\frac{a-1}{a+b-2}} \quad \left[\begin{array}{l} \text{Same as the} \\ \text{expected value of } \theta \end{array} \right]$$


We have likelihood $p(D|\theta) = \theta^{N_1} (1-\theta)^{N_0}$
prior $p(\theta|a,b) \propto \theta^{a-1} (1-\theta)^{b-1}$

We need posterior $p(\theta|D) = \frac{p(D|\theta) p(\theta)}{\int_{\theta \in [0,1]} p(D|\theta) p(\theta) d\theta}$

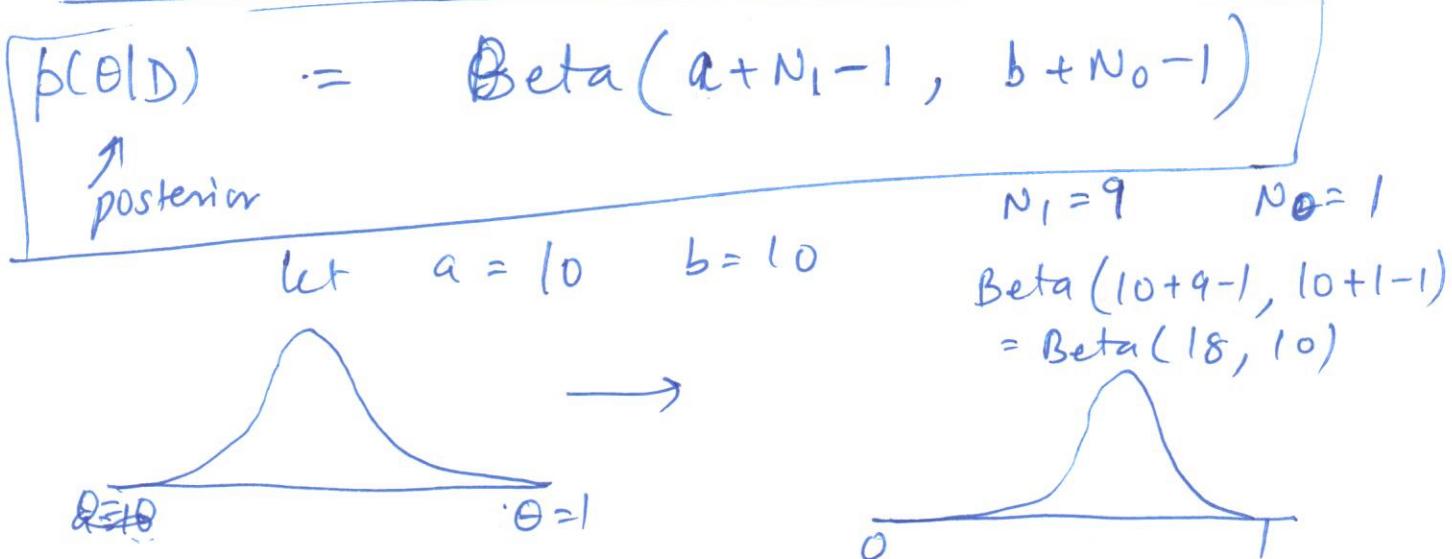
Bayes Rule.

for some pairs of $p(D|\theta)$ and $p(\theta)$
 \downarrow lik. \uparrow prior

Computing $\int p(D|\theta) p(\theta) d\theta$
is easy.

Conjugate pairs.

If $p(\theta)$ and $p(D|\theta)$ are conjugates
then $p(\theta|D)$ will have the same form as $p(\theta)$

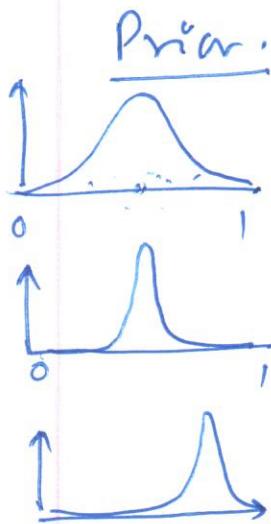


Learn parameters of a Bernoulli distribution,
given some samples drawn from it.

$$0 \leq \theta \leq 1$$

$$D = \{H, H, H, T, T, H\}$$

N_0 - # Tails in D
 N_1 - # Heads in D



Beta(a, b)

$$\hat{\theta}_{\text{Prior}} = \underset{\theta}{\operatorname{arg\,max}} \text{pdf}(\theta | a, b)$$

$$= \frac{a-1}{a+b-2}$$

$$\hat{\theta}_{\text{MLE}} = \frac{N_1}{N_0 + N_1}$$

①

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

↑
posterior

Conjugate pairs.

$\hat{p}(\theta | D)$ is also a Beta distribution (a', b')

$$\begin{aligned} a' &= a + N_1 \\ b' &= b + N_0 \end{aligned}$$

$$\underset{\theta}{\operatorname{arg\,max}} \text{pdf}(\theta | a', b')$$

$$\hat{\theta}_{\text{MAP}} = \frac{a + N_1 - 1}{a + b + N_0 + N_1 - 2}$$

Given a new instance x^*
What is $P(x^* = 1 | D)$

$\hat{\theta}_{\text{Prior}}$

$\hat{\theta}_{\text{MLE}}$

$\hat{\theta}_{\text{MAP}}$

Frequentist:

Bayesian Averaging.

Bayesian:

~~Side discussion~~ Beta ($\theta | a, b$) $\propto \theta^{a-1} (1-\theta)^{b-1}$

$$\begin{aligned} p(\theta | D) &= \frac{\theta^{N_1} (1-\theta)^{N_0} \theta^{a-1} (1-\theta)^{b-1}}{\int_{\theta} \theta^{N_1} (1-\theta)^{N_0} \theta^{a-1} (1-\theta)^{b-1} d\theta} \\ &= \frac{\theta^{N_1+a-1} (1-\theta)^{N_0+b-1}}{\int \theta^{N_1+a-1} (1-\theta)^{N_0+b-1} d\theta} \\ &= \frac{1}{Z} \theta^{N_1+a-1} (1-\theta)^{N_0+b-1} \\ &\quad \xrightarrow{\text{constant}} \\ p(\theta | D) &\propto \theta^{N_1+a-1} (1-\theta)^{N_0+b-1} \end{aligned}$$

Gaussian Distribution

$$x \sim \mathcal{N}(x | \mu, \Sigma)$$

$d \times 1$ $d \times 1$ $d \times d$

Let $d = 1$

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

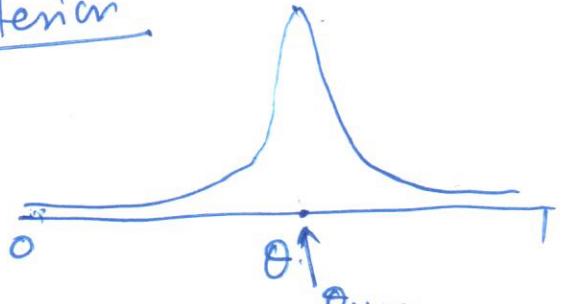
scalar scalar

$D = (x_1, x_2, x_3, \dots, x_N) \leftarrow \text{data.}$

Likelihood of D , $\ell(D) = p(x_1 | \mu, \sigma^2) * p(x_2 | \mu, \sigma^2) \dots$

$$\text{log-lik: } ll(D) = \sum_{i=1}^N \log p(x_i | \mu, \sigma^2)$$

Posterior



$\theta = 0.1$
 0.2
 0.3
 0.4
 0.5
 0.6
 0.7
 0.8
 0.9
 1.0

$$\begin{aligned}
 & \cancel{0.1 * p(\theta=0.1|D)} \\
 & 0.2 * p(\theta=0.2|D) \\
 & 0.3 * p(\theta=0.3|D) \\
 & \vdots \\
 & \frac{1}{\sum \theta p(\theta|D)}
 \end{aligned}$$

x^*

This is the same expression as $E[\theta]$ posterior.

$$\boxed{
 \int_{-\infty}^{\theta_{MAP}} \theta p(\theta|D) d\theta
 }$$

$$P(x^*=1|D) = \frac{a + N_1}{a + b + N_0 + N_1}$$

$$\ell\ell(D) = \sum_{i=1}^N \log \left[\frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{1}{2\sigma^2}(x_i - \mu)^2\right) \right]$$

$$= \sum_{i=1}^N \left[\log \frac{1}{\sqrt{2\pi}\sigma} + \left(-\frac{1}{2\sigma^2}(x_i - \mu)^2 \right) \right]$$

$$= \sum_{i=1}^N \left[-\frac{1}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2}(x_i - \mu)^2 \right]$$

$$\frac{\partial \ell\ell(D)}{\partial \mu} = \sum_{i=1}^N \left[\frac{\partial}{\partial \mu} \left(-\frac{1}{2} \log(2\pi\sigma^2) \right) - \frac{1}{2\sigma^2} \frac{\partial}{\partial \mu} (x_i - \mu)^2 \right]$$

$$= -\sum_{i=1}^N \frac{1}{2\sigma^2} 2(x_i - \mu) (-1)$$

$$= \frac{1}{\sigma^2} \sum_{i=1}^N (x_i - \mu)$$

$$\frac{\partial \ell\ell(D)}{\partial \mu} = 0 \Rightarrow \frac{1}{\sigma^2} \sum_{i=1}^N (x_i - \mu) = 0$$

$$\mu = \frac{1}{N} \sum_{i=1}^N x_i$$
Sample mean

MLE for Gaussian Distributions.

$$x \in \mathbb{R} \quad d=1$$

$$\text{pdf}(x) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left[-\frac{(x-\mu)^2}{2\sigma^2}\right]$$

$$D = \{x_1, x_2, \dots, x_N\}$$

$$\ell(D|\mu, \sigma^2) = \prod_{i=1}^N \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x_i-\mu)^2}{2\sigma^2}\right)$$

$$\ell\ell(D|\mu, \sigma^2) = \sum_{i=1}^N \log \left[\frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x_i-\mu)^2}{2\sigma^2}\right) \right]$$

To get MLE for μ, σ :

$$\frac{\partial}{\partial \mu} \ell\ell(D|\mu, \sigma^2) = 0 \quad \boxed{\mu_{\text{MLE}} = \frac{1}{N} \sum_{i=1}^N x_i}$$

$$\frac{\partial}{\partial \sigma} \ell\ell(D|\mu, \sigma^2) = \sum_{i=1}^N \left[\frac{\partial}{\partial \sigma} \left[\log \frac{1}{\sqrt{2\pi}\sigma} - \frac{(x_i-\mu)^2}{2\sigma^2} \right] \right]$$

$$= \sum_{i=1}^N \left[-\frac{\partial}{\partial \sigma} (\log \sqrt{2\pi}\sigma) - \frac{\partial}{\partial \sigma} \left[\frac{(x_i-\mu)^2}{2\sigma^2} \right] \right]$$

$$= -\frac{\partial}{\partial \sigma} \left[\frac{N}{2} \log \sigma + \frac{1}{2\sigma} \sum_{i=1}^N (x_i-\mu)^2 \right] - \frac{(x_i-\mu)^2}{2\sigma^3}$$

$$\begin{aligned} & \log(\sqrt{2\pi}\sigma) \\ & = \log \sqrt{2\pi} + \log \sigma \end{aligned}$$

Setting

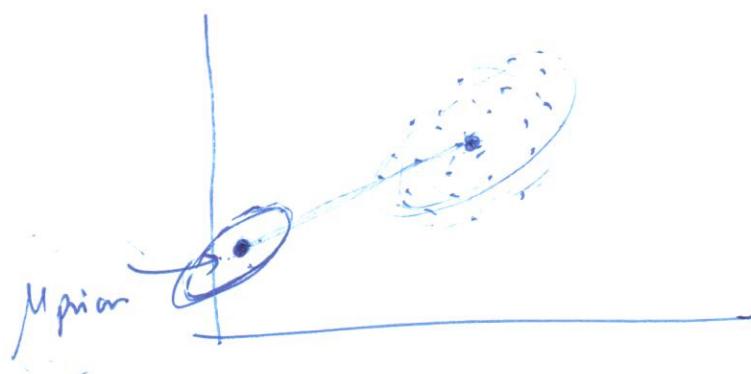
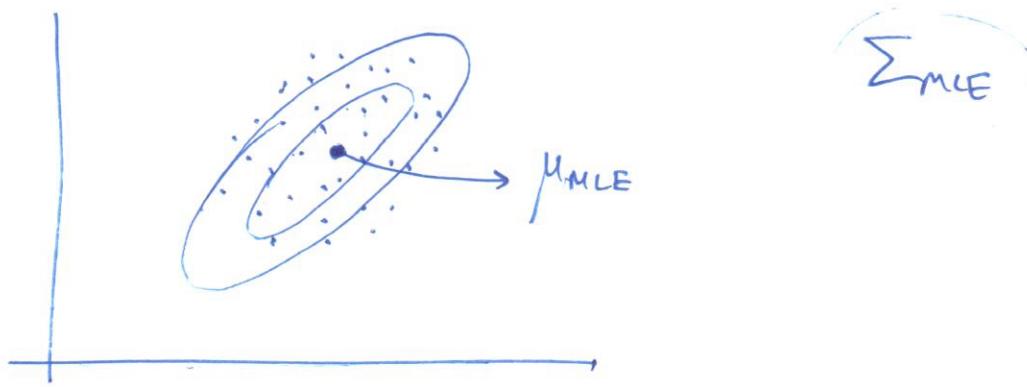
$$\frac{\partial}{\partial \sigma} \ell\ell(D|\mu, \sigma^2) = 0 \quad \boxed{\sigma_{\text{MLE}}^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \mu_{\text{MLE}})^2}$$

$$\text{if } d \geq 1 \quad \text{pdf}(x|\mu, \Sigma) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \exp\left[-\frac{1}{2} (x-\mu)^T \Sigma^{-1} (x-\mu)\right]$$

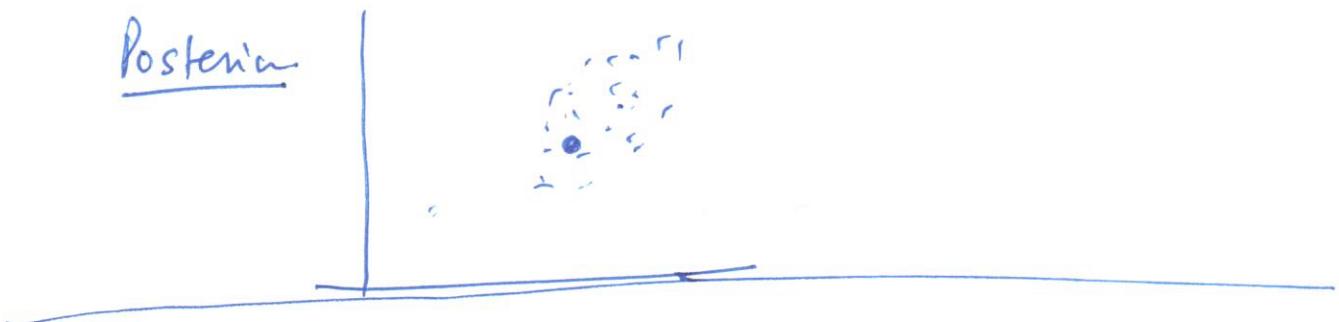
$$\mu_{\text{MLE}} = \frac{1}{N} \sum_{i=1}^N x_i \quad \boxed{\sum_{\text{MLE}} = \frac{1}{N} \sum_{i=1}^N \underbrace{(x_i - \mu_{\text{MLE}})}_{D \times 1} \underbrace{(x_i - \mu_{\text{MLE}})^T}_{1 \times D}}$$

Matrix
Cookbook
Handouts

$d=2$



Prior



$$D = \begin{bmatrix} x_1 & [3, 2] \\ x_2 & [1, 3] \\ x_3 & [2, 1] \end{bmatrix} \quad \mu = \frac{1}{N} \sum x_i$$

$$\mu = \frac{3+1+2}{3}, \frac{2+3+1}{3}$$

$$\Sigma_{MLE} = \frac{1}{N} \sum_{i=1}^N (x_i - \mu)(x_i - \mu)^T$$

$$x_1 = \begin{bmatrix} 3 \\ 2 \end{bmatrix} \quad x_2 = \begin{bmatrix} 1 \\ 3 \end{bmatrix} \quad x_3 = \begin{bmatrix} 2 \\ 1 \end{bmatrix} \quad \mu = \begin{bmatrix} 2 \\ 2 \end{bmatrix}$$

$$\Sigma_{MLE} = \frac{1}{3} \left[(x_1 - \mu)(x_1 - \mu)^T + (x_2 - \mu)(x_2 - \mu)^T + (x_3 - \mu)(x_3 - \mu)^T \right]$$

$$x_1 - \mu = \begin{bmatrix} 1 \\ 0 \end{bmatrix} \quad (x_1 - \mu)^T = \begin{bmatrix} 1 & 0 \end{bmatrix}$$

$$(x_1 - \mu)(x_1 - \mu)^T = \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix}$$

$$(x_2 - \mu)(x_2 - \mu)^T = \begin{bmatrix} -1 \\ 1 \end{bmatrix} \begin{bmatrix} -1 & 1 \end{bmatrix} = \begin{bmatrix} 1 & -1 \\ -1 & 1 \end{bmatrix}$$

$$(x_3 - \mu)(x_3 - \mu)^T = \begin{bmatrix} 0 \\ -1 \end{bmatrix} \begin{bmatrix} 0 & -1 \end{bmatrix} = \begin{bmatrix} 0 & 0 \\ 0 & 1 \end{bmatrix}$$

$$\Sigma_{MLE} = \frac{1}{3} \begin{bmatrix} 2 & -1 \\ 0 & 2 \end{bmatrix} = \begin{bmatrix} 2/3 & -1/3 \\ 0 & 2/3 \end{bmatrix}$$

$$x = N \times D$$

\oplus

$$\text{np.mean}(x)$$

$$\text{np.mean}(x, axis=0)$$

$$\begin{cases} \text{np.cov}(x) - N \times N \\ \text{np.cov}(x.T) \end{cases}$$