Assignment #5

ABHINAV KUMAR
180050003

BHASKAR GUPTA
180050022

**Q.1**

MLE estimate $\hat{\mu}_{MLE} = \Sigma x_i / n$

i) Gaussian prior

$$P(\mu|data) = \frac{P(data|\mu)P(\mu; 10.5, 1)}{\int P(data|\mu)P(\mu; 10.5, 1)\,d\mu}$$

MAP Estimate: $\frac{\partial}{\partial\mu} P(\mu|data) = 0 \equiv \frac{\partial}{\partial\mu}\log\left[P(\mu|data)\right] = 0$

$\log P(\mu|data) \propto \cancel{\sum_i e^{-(x_i-\mu)^2/2\sigma^2}}$

$$\sum_i \frac{(x_i-\mu)^2}{2\sigma_{true}^2} + \frac{(\mu-10.5)^2}{2(1)}$$

$$\frac{\partial}{\partial\mu}\log P(\mu|data) = \frac{\sum(x_i-\mu)}{\sigma_{true}^2} + \frac{(10.5-\mu)}{(1)} = 0$$

$$(\qquad)$$

$$\Rightarrow \frac{n\mu}{\sigma_{true}^2} + \mu = 10.5 + \frac{\sum x_i}{\sigma_{true}^2}$$

$$\mu = \frac{10.5 + \dfrac{\sum x_i}{\sigma_{true}^2}}{\dfrac{n}{\sigma_{true}^2} + 1}$$

ii) Uniform prior

$$P(\mu|data) = \frac{P(data|\mu)\times(1/2)}{\int_\mu P(data|\mu)\times\frac{1}{2}\,d\mu} \longrightarrow e^{-\Sigma(x_i-\mu)^2/2\sigma_{true}^2}$$

$\Rightarrow \frac{\partial}{\partial\mu}\log P(\mu|data) = 0 \Rightarrow \mu = \Sigma x_i/n$

But since $\mu$ has a uniform prior in $(9.5, 11.5)$, and Likelihood function is parabolic,

if a. $\bar{x} < 9.5$, MAP estimate in $(9.5, 11.5)$ is $9.5$

if $\bar{x} > 11.5$, MAP estimate in $(9.5, 11.5)$ is $11.5$

else MAP estimate is $\bar{x}$.

b) As N increases, the relative error for all the three estimates of $\mu$ converges to that of MLE estimate. And, Gaussian prior boxplot has the minimum variance amongst all three estimate as The posterior $\wedge$ which is is product of two Gaussians,
                     mean

has $\sigma^{2*}$ (variance) $= \dfrac{(\sigma^2)(\sigma_0^2/n)}{\sigma^2 + \sigma_0^2/n} < \left(\dfrac{\sigma_0^2}{n}\right)$ $\begin{bmatrix} \text{ML estimator} \\ \text{Variance} \end{bmatrix}$

Hence, Gaussian prior is preferable.

Q2) $y = -\frac{1}{\lambda} \log x$.

$g(x) = -\frac{1}{\lambda} \log x$.

$x = e^{-\lambda y}$

$\therefore g^{-1}(y) = e^{-\lambda y}$

$p(y) = \left| q_x(g^{-1}(y)) \frac{d}{dy} (g^{-1}(y)) \right|$

$p(y) = \lambda e^{-\lambda y}$

$P(Y/\lambda) = \lambda^N e^{-\lambda \sum\limits_{i=1}^{N} y_i}$

To find $\hat{\lambda}_{ML}$ we will differentiate $P(Y/\lambda)$ w.r.t. $\lambda$.

$\frac{d}{d\lambda} P(Y/\lambda) = 0$.

$N\lambda^{N-1} e^{-\lambda \sum y_i} - \left(\sum y_i\right)\lambda^N e^{-\lambda \sum y_i} = 0$.

$\hat{\lambda}_{ML} = \frac{N}{\sum y_i}$

$$P(\lambda|Y) = \frac{P(Y|\lambda)P(\lambda)}{\int_0^\infty P(Y|\lambda)P(\lambda)d\lambda}$$

$$= \frac{\lambda^N e^{-\lambda \Sigma y_i} \lambda^{\alpha-1} e^{-\beta\lambda}}{\int_0^\infty \lambda^N e^{-\lambda \Sigma y_i} \lambda^{\alpha-1} e^{-\beta\lambda} d\lambda} \qquad \left(\begin{array}{l}\text{ignored the}\\ \text{constant as}\\ \text{they will}\\ \text{be cancelled}\end{array}\right)$$

$$= \frac{\lambda^{N+\alpha-1} e^{-\lambda(\beta+\Sigma y_i)}}{\frac{\Gamma(N+\alpha)}{(\Sigma y_i + \beta)^{N+\alpha}}}$$

$$E(\lambda) = \int_0^\infty \lambda \cdot \frac{\lambda^{N+\alpha-1} e^{-\lambda(\beta+\Sigma y_i)}}{\frac{\Gamma(N+\alpha)}{(\Sigma y_i + \beta)^{N+\alpha}}} d\lambda$$

$$= \frac{N+\alpha}{\Sigma y_i + \beta}$$

We can find it by using the property

$$\Gamma(x+1) = x\Gamma(x)$$

$$\boxed{\begin{array}{l} \lambda_{posterior~mean} = \dfrac{N+\alpha}{\Sigma y_i + \beta} \\[2em] \hat{\lambda}_{ML} = \dfrac{N}{\Sigma y_i} \end{array}}$$

Posterior mean for small N. M.L.E for larger N.

The posterior mean estimate is more accurate for small N. As N increases MLE becomes more efficient.

P.M.E. Tends to M.L.E as $N \to \infty$ and has less variance.

Posterior is more valuable for limited data, but as N increases weight of the biased prior causes M.L.E. to be better.

(Q3) Points are of the form
$$(r_0 \cos\theta, \; r_0 \sin\theta)$$

$$E[r_0 \cos\theta, \; r_0 \sin\theta] = [0, 0].$$

$$C = E\left( \begin{bmatrix} r_0 \cos\theta \\ r_0 \sin\theta \end{bmatrix} \begin{bmatrix} r_0 \cos\theta & r_0 \sin\theta \end{bmatrix} \right)$$

$$= E \begin{bmatrix} r_0^2 \cos^2\theta & r_0^2 \sin\theta\cos\theta \\ r_0^2 \sin\theta\cos\theta & r_0^2 \sin^2\theta \end{bmatrix}$$

$$C = \begin{bmatrix} r_0^2/2 & 0 \\ 0 & r_0^2/2 \end{bmatrix}$$

The theoretical and empirical mean and covariance matrix are closely correlated.

The modal point is theoretically $(0, 0)$ but no data point is present there.

We can even find the modal radius:

$$P(r, \theta) = \frac{1}{2\pi \left(\frac{r_o^2}{2}\right)} e^{\frac{-r^2}{r_o^2}} \, r \, dr \, d\theta$$

$$P(r) = \frac{1}{\left(\frac{r_o^2}{2}\right)} e^{-r^2/r_o^2} \, r \, dr.$$

$$\frac{dP(r)}{dr} = 0$$

$$\Rightarrow e^{-r^2/r_o^2} \left[ 1 - \frac{2r^2}{r_o^2} e^{-r^2/r_o^2} \right] = 0$$

$$r = \frac{r_o}{\sqrt{2}} \quad (\text{modal radius})$$

The Gaussian model is not satisfactory as the modal radius is at $\frac{1}{\sqrt{2}}$ times the distance of the data

$$\text{Likelihood } (L) = \prod_{i=1}^{N} p(x_i ; \mu ; C)$$

$$L = \frac{1}{\left(2\pi \sqrt{|C|}\right)^N} e^{-\frac{1}{2} \sum_i \left(x_i - \mu\right)^T C^{-1} (x - \mu)}$$

Taking log and differentiating w.r.t. $\mu$ we get.

$$-\frac{1}{2} \sum_i \frac{d}{d\mu} \left[ (x_i - \mu)^T C^{-1} (x_i - \mu) \right] = 0$$

$$\sum_i C^{-1} (x_i - \mu) = 0$$

$$\hat{\mu} = \frac{\sum_i x_i}{N}$$

$$\frac{\partial}{\partial C} \left( \log |C| + \frac{1}{2} \sum_i (x_i - \mu)^T C^{-1} (x_i - \mu) \right) = 0$$

$$\hat{C} = \frac{\sum_i (x_i - \hat{\mu})(x_i - \hat{\mu})^T}{N}$$

$$\hat{C} = \frac{\sum_i x_i x_i^T}{N} - \hat{\mu} \hat{\mu}^T$$

Q.3 c) Yes. They match the theoretical values

$$\hat{Cov} = \begin{bmatrix} .5 & -.0007 \\ -.0007 & .5 \end{bmatrix} \qquad C_{the} = \begin{bmatrix} .5 & 0 \\ 0 & .5 \end{bmatrix}$$

$$\hat{\mu} = \begin{bmatrix} .0016 \\ .0008 \end{bmatrix} \qquad \mu_{theo} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$N = 100000$$

**Q.4**

**a)** $\hat{\Theta}_{ML}$

   likelihood function $f(x;\Theta) = \dfrac{1}{\Theta^n} \longrightarrow$ decreasing function,

$\Theta \geqslant \max \{x_1 \ldots x_n\}$

$\Rightarrow$ Max-likelihood estimate is $\boxed{\hat{\Theta}_{ML} = \max \{x_1 \ldots x_n\}}$

$P(\Theta|x) \propto P(x|\Theta) \times P(\Theta)$

$P(x|\Theta) \propto \dfrac{1}{\Theta^n}$ ,   ⊙ $\Theta \geqslant \max\{x_1 \ldots x_n\}$
                          otherwise $0$.

$P(\Theta) \propto \left(\dfrac{\Theta_m}{\Theta}\right)^{\alpha}$   for $\Theta \geqslant \Theta_m$

              $0$       otherwise

$\Rightarrow P(x|\Theta) P(\Theta) \propto \dfrac{1}{\Theta^{n+\alpha}}$ ,   $\Theta \geqslant \underset{max}{(\max\{x_1 \ldots x_n\}, \Theta)}$

                 $0$       otherwise

$\Rightarrow$ Posterior is also Pareto with parameters
$\left(n+\alpha, \max\{x_1 \ldots x_n, \Theta\}\right)$

MAP estimate $\boxed{\hat{\Theta}_{MP} = \max \{x_1 \ldots x_n, \Theta\}}$

$\Big[$Same reasoning as MLE case,   $P(\Theta|x)$ is decreasing$\Big]$.

b)  $\hat{\Theta}_{ML} = \max \{X_1 \text{---} X_n\}$

 $\hat{\Theta}_{MP} = \max \{X_1 \text{---} X_n, \Theta_m\}$

As sample size tends to infinity & by the prior we know that $\hat{\Theta}_1 > \Theta_m$, there will be $X_i$'s $> \hat{\Theta}_m$ & hence (more chances of finding)

 $\hat{\Theta}_{ML} = \max \{X_1 \text{---} X_n\} = \max \{X_1 \text{---} X_n, \Theta_m\}$

(since , $X_i > \Theta_m$ , both equalities are the same)

$\Rightarrow \hat{\Theta}_{MAP} \rightarrow \hat{\Theta}_{ML}$ or rather $\hat{\Theta}_{ML} = \hat{\Theta}_{MAP}$ if samples are perfectly random.

And above is desirable as MLE estimate has asymptotically lowest MSE &.

c)  Posterior $P(\Theta|x)$ is also proportional to

  $\text{Pareto}(n+\alpha, \underbrace{\max \{X_1 \text{---} X_n, \Theta\}}_{\Theta_m'})$

$\Rightarrow E_{P(\Theta|x)}[\hat{\Theta}] \propto \displaystyle\int_{\Theta_m'} \frac{(\Theta_m')^{n+\alpha}}{(\Theta)^{n+\alpha+1}} \Theta \, d\Theta$

   $\propto (\Theta_m')^{n+\alpha} \cdot \dfrac{1}{(n+\alpha-2)(\Theta_m')^{n+\alpha-2}} \propto \dfrac{(\Theta_m')^2}{n+\alpha-2}$

More rigorous $\rightarrow$

$P(\Theta|x) = c \cdot \dfrac{\dfrac{1}{\Theta^n} \times \left(\dfrac{\Theta_m}{\Theta}\right)^\alpha}{\displaystyle\int \dfrac{1}{\Theta^n} \times \dfrac{(\Theta_m)^\alpha}{\Theta^\alpha} d\Theta}$    $\Theta \geq \underbrace{\max \{X_1 \text{---} X_n, \Theta\}}_{\Theta_m'}$

  $= \dfrac{c}{\Theta^{n+\alpha}} \times \dfrac{1}{\displaystyle\int_{\Theta_m}^{\infty} \dfrac{d\Theta}{\Theta^{n+\alpha}}} = \dfrac{c/\Theta^{n+\alpha}}{\dfrac{1}{(n+\alpha-1)(\Theta_m')^{n+\alpha-1}}}$

Scanned by CamScanner

$$\Rightarrow E_{P(\theta|x)}[\hat{\theta}] = \int_{\theta_m'} \frac{c}{\theta^{n+\alpha}} \times \frac{(\theta_m')^{n+\alpha-1}}{9} \times (n+\alpha-1)\,\theta\,d\theta$$

$$= \frac{(c)(\theta_m')^{n+\alpha-1} \ (n+\alpha-1)}{(\theta_m')^{n+\alpha-2} \ (n+\alpha-2)}$$

$$\boxed{= \frac{(c)(\theta_m')(n+\alpha-1)}{(n+\alpha-2)}}$$

$$\theta_m' = \max \{\theta_m, x_1, \cdots x_n\}$$

d) @

$$\lim_{n\to\infty} (c)\,\theta_m'\,\frac{(n+\alpha-1)}{(n+\alpha-2)} = c\,\theta_m'$$

The exact value of $\hat{\theta}_{pm}$ depends on coefficient in Pareto but asymptotically, $\hat{\theta}_m'$ also tends to $\hat{\theta}_{ML}$ by the same reasoning as in (b). This is desirable as we want our estimate to have minimum MSE & MLE has minimum estimate.