

1. (a). We have

$$\begin{aligned} \text{MSE}(\theta, \delta) &= \mathbb{E}_\theta \|\delta(x) - \mathbb{E}_\theta \delta(x) + \mathbb{E}_\theta \delta(x) - \theta\|_2^2 \\ &= \mathbb{E}_\theta \|\delta(x) - \mathbb{E}_\theta \delta(x)\|_2^2 + 2\mathbb{E}_\theta \langle \delta(x) - \mathbb{E}_\theta \delta(x), \mathbb{E}_\theta \delta(x) - \theta \rangle + \mathbb{E}_\theta \|\mathbb{E}_\theta \delta(x) - \theta\|_2^2 \\ &= \mathbb{E}_\theta \|\delta(x) - \mathbb{E}_\theta \delta(x)\|_2^2 + \|\mathbb{E}_\theta \delta(x) - \theta\|_2^2 \end{aligned}$$

Where the last step comes from $\mathbb{E}_\theta \langle \delta(x) - \mathbb{E}_\theta \delta(x), \mathbb{E}_\theta \delta(x) - \theta \rangle = 0$

Notice that $\mathbb{E}_\theta \|\delta(x) - \mathbb{E}_\theta \delta(x)\|_2^2 = \text{tr}[\text{Var}_\theta(\delta)]$

and $\|\mathbb{E}_\theta \delta(x) - \theta\|_2^2 = \|\text{Bias}_\theta(\delta)\|_2^2$

There we have $\text{MSE}(\theta, \delta) = \text{tr}[\text{Var}_\theta(\delta)] + \|\text{Bias}_\theta(\delta)\|_2^2$.

(b). Using (1), we have

$$\begin{aligned} \text{MSE}(\theta, \delta_X) &= \text{Bias}_\theta(\delta_X)^2 + \text{Var}_\theta(\delta_X) \\ &= (\gamma(\theta_0 - \theta))^2 + (1-\gamma)^2 n\theta(1-\theta)/n^2 \\ &= \gamma^2(\theta_0 - \theta)^2 + (1-\gamma)^2 \theta(1-\theta)/n \end{aligned} \quad \square$$

(c). Notice that $\text{MSE}(\theta, \delta_X)$ is a convex function of γ .

$$\begin{aligned} \text{Thus } 0 &= \frac{d}{d\gamma} \text{MSE}(\theta, \delta_X) \Big|_{\gamma=\gamma^*} = 2\gamma^*(\theta_0 - \theta)^2 + 2(1-\gamma^*)\theta(1-\theta)/n \\ \gamma^* &= \theta(1-\theta) / (n(\theta_0 - \theta)^2 + \theta(1-\theta)). \end{aligned}$$

If $\theta \rightarrow \theta_0$ and n is fixed, then $\gamma^* \rightarrow 1$.

If $n \rightarrow \infty$ and θ is fixed, then $\gamma^* \rightarrow 0$.

These limits make sense because in the former case,

θ_0 gives a better estimate, so δ_X should put more weight on θ_0 , i.e., $\gamma \rightarrow 1$; in the latter case, the variance of X/n is smaller ~~than~~ than the bias of θ_0 , so δ_X should move towards X/n , i.e., $\gamma \rightarrow 0$.

(d). δ_0 is not inadmissible.

In fact, for any $\gamma \in (0, 1)$, there exists $\tilde{\theta} \xrightarrow{\gamma \rightarrow 0 \text{ or } \gamma \rightarrow 1} \theta_0$, such that $\text{MSE}(\tilde{\theta}, \delta_X) > \text{MSE}(\tilde{\theta}, \delta_0)$ for $\epsilon > 0$. Thus " δ_0 is dominated by any δ_X " is invalid. \square

2. (a). Let $g_1 = c^{-1}$, $g_2 = (1-c)^{-1}$, and $f_i(x) = \left(e^{\eta_i^T T(x)} h_i(x) \right)^{\frac{1}{g_i}}$ $i=1,2$.

Hölder's inequality implies:

$$\int \exp(c\eta_1 + (1-c)\eta_2)^T T(x) h(x) d\mu(x) \leq \left(\int \exp(c\eta_1^T T(x))^{\frac{1}{c}} h(x) d\mu(x) \right)^c \left(\int \exp((1-c)\eta_2^T T(x))^{\frac{1}{1-c}} h(x) d\mu(x) \right)^{1-c}$$

Taking logarithm on both sides, notice that the left hand side becomes $A(c\eta_1 + (1-c)\eta_2)$ and the right hand side becomes $cA(\eta_1) + (1-c)A(\eta_2)$. Thus

$$A(c\eta_1 + (1-c)\eta_2) \leq cA(\eta_1) + (1-c)A(\eta_2).$$

□

(b). $\forall \eta_1, \eta_2 \in \Xi_1$ and $\lambda \in (0,1)$, we have

$$A(\lambda\eta_1 + (1-\lambda)\eta_2) \leq \lambda A(\eta_1) + (1-\lambda)A(\eta_2) < \infty$$

This means $\lambda\eta_1 + (1-\lambda)\eta_2 \in \Xi_1$, so Ξ_1 is convex. □

3. (a) Introduce Lagrangian multipliers $\lambda \in \mathbb{R}$, $\gamma \in \mathbb{R}^S$, we obtain Lagrangian

$$L(p; \lambda, \gamma) = -\sum_x p(x) \log p(x) + \lambda(1 - \sum_x p(x)) + \gamma^T (\alpha - \sum_x p(x) T(x))$$

Since the objective and feasible set are all convex, the solution of the primal problem can be written as:

$$p(x) = \exp(-\gamma^T T(x) - \lambda - 1)$$

which belongs to the s -parameter exponential family. □

(b) Introduce Lagrangian multipliers $\lambda_1, \lambda_2, \lambda_3 \in \mathbb{R}$ and let

$$L(p; \lambda_1, \lambda_2, \lambda_3) = \int_{\mathbb{R}} J(p, x; \lambda_1, \lambda_2, \lambda_3) dx + \lambda_1 + \lambda_2 + \lambda_3$$

$$\text{where } J(p, x; \lambda_1, \lambda_2, \lambda_3) = -p \log p - \lambda_1 p - \lambda_2 x p - \lambda_3 x^2 p,$$

calculus of variation gives

$$p = \frac{\delta J}{\delta p} = -(1 + \log p - \lambda_1 - \lambda_2 x - \lambda_3 x^2), \text{ i.e. } p(x) = \exp(-\lambda_3 x^2 - \lambda_2 x - \lambda_1 - 1).$$

Thus p must be density of Gaussian or Exponential distribution.

If $\mu^2 \neq \sigma^2$, then p is Gaussian, i.e. $X \sim N(\mu, \sigma^2)$.

If $\mu^2 = \sigma^2$, then we need to compare the entropy.

$$\text{Ent}(X \sim N(\mu, \sigma^2)) = \frac{1}{2}(1 + \log(2\pi) + \log \mu^2)$$

$$\text{Ent}(X \sim \text{Exp}(\mu)) = 1 + \log \mu < \text{Ent}(X \sim N(\mu, \sigma^2))$$

Thus in the case, X is still Gaussian, i.e. $X \sim N(\mu, \sigma^2)$. □

4. (a). We can write the density function of $T(k, \theta)$ as

$$p_{k, \theta}(x) = \exp(\langle \eta, T(x) \rangle - A(\eta)) \text{ where}$$

natural parameter is $\eta = (-\theta^{-1}, k-1)$

sufficient statistic is $T(x) = (x, \log x)$

carrier density is $h(x) = 1$

log-partition function is $A(\eta) = \log T(k) + k \log \theta$.

The canonical form reads

$$p_{k, \theta}(x) = \exp(\langle (-\theta^{-1}, k-1), (x, \log x) \rangle - (\log T(k) + k \log \theta)) \quad \square$$

(b). It follows that:

$$E[X] = \frac{\partial}{\partial \eta_1} A(\eta) = k\theta,$$

$$\text{Var}(X) = \frac{\partial^2}{\partial \eta_1^2} A(\eta) = k\theta^2. \quad \square$$

(c) We have $M_X(u)$

$$= \int \frac{x^{k-1} e^{-\theta^{-1}x + ux}}{T(k) \theta^k} dx$$

$$= (1 - \theta u)^{-k}.$$

$$\text{Thus } M_{X_+}(u) = \prod_{i=1}^n M_{X_i}(u) = (1 - \theta u)^{\sum_{i=1}^n k_i}$$

It follows that $X_+ \sim T(\sum_{i=1}^n k_i, \theta)$. \square

5 (a). The density function of $Y = (X_{(1)}, \dots, X_{(n)})$ is given by

$$p_Y(y_1, \dots, y_n) = n! \exp\left(-\sum_{i=1}^n y_i\right) \mathbb{1}_{(y_1 \leq \dots \leq y_n)}$$

Now consider linear transform $Z = AY$ where

$$A = \begin{pmatrix} 1 & 0 & \dots & 0 \\ -1 & 1 & \dots & 0 \\ 0 & -1 & 1 & \dots & 0 \\ \vdots & \vdots & \ddots & \ddots & \vdots \\ 0 & \dots & 0 & -1 & 1 \end{pmatrix} \in \mathbb{R}^{n \times n}, \quad |J_A(Y)| = n!$$

$$\begin{aligned} \text{Then } p_Z(z_1, \dots, z_n) &= p_Y(A^{-1}Z) \cdot |J_{A^{-1}}(Z)| \\ &= n! \cdot \exp\left(-\sum_{i=1}^n z_i\right) \cdot (n!)^{-1} \cdot \mathbb{1}_{(z_i \geq 0, \forall i)} \\ &= \exp\left(-\sum_{i=1}^n z_i\right) \end{aligned}$$

It follows that $z_i \stackrel{\text{iid}}{\sim} \text{Exp}(1)$. \square

(b). Since F^{-1} is monotone increasing, we have

$$\begin{aligned} &(F^{-1}(U_{(1)}), \dots, F^{-1}(U_{(n)})) \\ &= (F^{-1}(U_{(1)}))_{(1)}, \dots, (F^{-1}(U_{(n)}))_{(n)} \end{aligned}$$

where RHS is the order statistics of $(F^{-1}(U_1), \dots, F^{-1}(U_n))$.

Notice that $F^{-1}(U_i) \stackrel{d}{=} W_i$, thus

$$(W_{(1)}), \dots, W_{(n)} \stackrel{d}{=} (F^{-1}(U_{(1)}))_{(1)}, \dots, (F^{-1}(U_{(n)}))_{(n)}$$

$$\stackrel{d}{=} (F^{-1}(U))_{(1)}, \dots, (F^{-1}(U))_{(n)}$$

$$= (F^{-1}(W_{(1)}))_{(1)}, \dots, (F^{-1}(W_{(n)}))_{(n)}$$

\square

(c). Notice that $\frac{1}{n}Z_1 + \dots + \frac{1}{n+r+1}Z_r = X_{(1)} + X_{(2)} - X_{(1)} + \dots + X_{(n)} - X_{(n-r)}$

$$\text{Let } G(w) = 1 - e^{-w} = X_{(1)}$$

Then G is the CDF of $\text{Exp}(1)$. We thus have

$$\begin{aligned} &\stackrel{d}{=} (F^{-1}(G(X_{(1)})))_{(1)}, \dots, (F^{-1}(G(X_{(n)})))_{(n)} \\ &\stackrel{d}{=} F^{-1}(U_{(1)}), \dots, F^{-1}(U_{(n)}) \quad (\text{since } (G(X_{(i)}))_{i=1}^n = (U_{(i)})_{i=1}^n) \end{aligned}$$

$$\stackrel{d}{=} (W_{(1)}), \dots, W_{(n)} \quad (\text{due to (b)})$$

In the penultimate step, we used the fact that

$G(X) \stackrel{d}{=} U$ where $U = (U_1, \dots, U_n)$ is iid uniform random variables. \square

6.(a). We will show that for arbitrary $\theta_1 \neq \theta_2$,

$$p_{\theta_1}(x)/p_{\theta_1}(T(x)) = p_{\theta_2}(x)/p_{\theta_2}(T(x)) \quad (*)$$

Therefore $p_{\theta}(x)/p_{\theta}(T(x))$ does not depend on θ ,
it follows that T is sufficient.

To show this, consider $q = \gamma \delta_{\theta_1} + (1-\gamma) \delta_{\theta_2}$,
where δ_{θ} is the Dirac-Delta measure at θ .

Thus ~~$\gamma p_{\theta_1}(x) + (1-\gamma) p_{\theta_2}(x)$~~

$$\gamma p_{\theta_1}(x) / (\gamma p_{\theta_1}(x) + (1-\gamma) p_{\theta_2}(x))$$

$$= \mathbb{P}_{\text{post}}(\theta|x)$$

$$= \mathbb{P}_{\text{post}}(\theta|T(x))$$

$$= \gamma p_{\theta_1}(T(x)) / (\gamma p_{\theta_1}(T(x)) + (1-\gamma) p_{\theta_2}(T(x)))$$

Rearranging, we obtain (*). \square

(b). If T is sufficient, by Factorization Thm, $p_{\theta}(x) = h(x) g_{\theta}(T(\theta))$.

Thus $\mathbb{P}_{\text{post}}(\theta|x)$

$$= p_{\theta}(x) g_{\theta}(\theta) / \int p_{\theta}(x) g_{\theta}(\theta') d\theta'$$

(cancelling out $h(x)$)

$$= g_{\theta}(T(\theta)) g_{\theta}(\theta) / \int g_{\theta}(T(\theta)) g_{\theta}(\theta') d\theta'$$

Now in the factorization, we can choose

$g_{\theta}(T(\theta))$ to be the density function of $T(x)$ given θ .

Therefore, $\mathbb{P}_{\text{post}}(\theta|x)$

$$= p_{\theta}(T) g_{\theta}(\theta) / \int p_{\theta}(T) g_{\theta}(\theta') d\theta'$$

$$= \mathbb{P}_{\text{post}}(\theta|T(x)).$$

This means that the posterior depends on x only
through $T(x)$. \square

7. The density function of $(X_{(1)}, \dots, X_{(r)})$ is given by

$$\begin{aligned}
 & p_{X_{(1)}, \dots, X_{(r)}}(y_1, \dots, y_r) \\
 &= \int_{\mathbb{R}^{n-r}} n! \mu^{-n} \exp\left(-\sum_{i=1}^n (y_i - \sigma) \mu^{-1}\right) \mathbb{1}_{(\sigma < y_1 < \dots < y_n)} dy_{n-1} \dots dy_{r+1} \\
 &= \int_{\mathbb{R}^{n-r-1}} n! \mu^{-(n-1)} \exp\left(-\sum_{i=1}^{n-1} (y_i - \sigma) \mu^{-1} - (y_{n-1} - \sigma) \mu^{-1}\right) \mathbb{1}_{(\sigma < y_1 < \dots < y_{n-1})} \\
 &\quad dy_{n-1} \dots dy_{r+1} \\
 &\quad \text{(integrating over } y_n)
 \end{aligned}$$

$$\dots$$

$$= \int_{\mathbb{R}} h(n-1, \dots, n-r+1) \mu^{-r} \exp\left(-\sum_{i=1}^r (y_i - \sigma) \mu^{-1} - r(y_r - \sigma) \mu^{-1}\right) \mathbb{1}_{(\sigma < y_1 < \dots < y_r)} dy$$

Let $g_{\mu, \sigma}(u, v) = h(n-1, \dots, n-r+1) \cdot \mu^{-r} \cdot \exp\left(-\sum_{i=1}^r (u_i - \sigma) \mu^{-1} - r(u_r - \sigma) \mu^{-1}\right) \cdot \mathbb{1}_{(\sigma < u)}$

and $h(x) = \mathbb{1}_{(x_1 < x_2 < \dots < x_r)}$ where $x = (x_1, \dots, x_r)$,

then we can write $p_{X_{(1)}, \dots, X_{(r)}}$ as

$$\begin{aligned}
 & p_{X_{(1)}, \dots, X_{(r)}}(y_1, \dots, y_r) \\
 &= g_{\mu, \sigma}\left(y_1, \sum_{i=1}^r y_i + r y_r\right) \cdot h(y_1, \dots, y_r) \\
 &= g_{\mu, \sigma}(Y) \cdot h(y_1, \dots, y_r)
 \end{aligned}$$

Thus Factorization Theorem asserts that Y is a sufficient statistics of $(X_{(1)}, \dots, X_{(r)})$. \square