

1: Properties of exponential families

(a) Want to show that if θ_1 and θ_2 are in Θ_N , then for any $\lambda \in [0, 1]$, $\lambda\theta_1 + (1 - \lambda)\theta_2 \in \Theta_N$.

Observe that

$$\begin{aligned} & \int \exp\{\langle \lambda\theta_1 + (1 - \lambda)\theta_2, \mathbf{T}(x) \rangle\} h(x) d\mathbf{x} \\ &= \int \exp\{\langle \lambda\theta_1 + (1 - \lambda)\theta_2, \mathbf{T}(x) \rangle\} h(x)^\lambda h(x)^{1-\lambda} d\mathbf{x} \\ &= \int (\exp\{\langle \theta_1, \mathbf{T}(x) \rangle\} h(x))^\lambda (\exp\{\langle \theta_2, \mathbf{T}(x) \rangle\} h(x))^{1-\lambda} d\mathbf{x} \\ &\leq \left(\int (\exp\{\langle \theta_1, \mathbf{T}(x) \rangle\} h(x)) d\mathbf{x} \right)^\lambda \left(\int (\exp\{\langle \theta_2, \mathbf{T}(x) \rangle\} h(x)) d\mathbf{x} \right)^{1-\lambda} \quad [By \text{ Holder's inequality}] \\ &< \infty \end{aligned}$$

That ϕ is a convex function follows from exactly the same line of arguments as above.

(b)

$$\frac{\partial \phi}{\partial \theta_i}(\theta) = \frac{1}{Z(\theta)} \int T_i(x) \exp\{\langle \theta, \mathbf{T}(x) \rangle\} h(x) d\mathbf{x} = \mathbf{E}_\theta\{T_i(\mathbf{X})\}$$

Now,

$$Z(\theta) \frac{\partial \phi}{\partial \theta_i}(\theta) = \int T_i(x) \exp\{\langle \theta, \mathbf{T}(x) \rangle\} h(x) d\mathbf{x} = \mathbf{E}_\theta\{T_i(\mathbf{X})\}$$

Differentiating the above equation w.r.t. θ_j and then dividing both sides by $Z(\theta)$, we get

$$\frac{\partial \phi}{\partial \theta_i}(\theta) \frac{\partial \phi}{\partial \theta_j}(\theta) + \frac{\partial^2 \phi}{\partial \theta_i \partial \theta_j}(\theta) = \frac{1}{Z(\theta)} \int T_i(x) T_j(x) \exp\{\langle \theta, \mathbf{T}(x) \rangle\} h(x) d\mathbf{x}$$

Hence proved.

(c)

$$\begin{aligned} & \mathbf{E}_\theta \left\{ \left[\frac{1}{h(\mathbf{x})} \frac{\partial h}{\partial x_i}(\mathbf{x}) + \langle \theta, \frac{\partial \mathbf{T}}{\partial x_i}(\mathbf{x}) \rangle \right] g(\mathbf{x}) \right\} \\ &= \int \frac{\partial \mathbf{P}_\theta(\mathbf{x})}{\partial x_i} g(\mathbf{x}) d\mathbf{x} \\ &= - \int \mathbf{P}_\theta(\mathbf{x}) \frac{\partial g}{\partial x_i}(\mathbf{x}) d\mathbf{x} \quad [By \text{ Integration by parts}] \\ &= -\mathbf{E}_\theta \left\{ \frac{\partial g}{\partial x_i}(\mathbf{x}) \right\} \end{aligned}$$

Hence proved.

(d) If $X \sim N(\mu, I)$, then trivially we have

$$\mathbb{E}\{(\mathbf{x} - \boldsymbol{\mu}) g(\mathbf{x})\} = \mathbb{E}\{\nabla g(\mathbf{x})\}.$$

where we applied Stein's identity component wise.

Now, let $X \sim N(\mu, \Sigma)$, and whiten X .

Take $Y = AX$ where A^{-1} is the square root of Σ . Let $h(\mathbf{x}) = g(A^{-1}\mathbf{x})$.

Using Stein's identity for Y and the differentiable function h ., we get

$$\mathbb{E}\{(\mathbf{y} - A\boldsymbol{\mu}) h(\mathbf{x})\} = \mathbb{E}\{\nabla h(\mathbf{x})\}.$$

$$\text{Also, } \nabla h(\mathbf{x}) = A^{-1}\nabla g(\mathbf{x})$$

Plugging in $Y = AX$, we get,

$$\mathbb{E}\{(A\mathbf{x} - A\boldsymbol{\mu}) g(\mathbf{x})\} = A^{-1}\mathbb{E}\{\nabla g(\mathbf{x})\}.$$

Hence,

$$\mathbb{E}\{(\mathbf{x} - \boldsymbol{\mu}) g(\mathbf{x})\} = \Sigma\mathbb{E}\{\nabla g(\mathbf{x})\}.$$

2: Exercises on sufficient statistics

(a) We claim that $T(x) = \{\frac{p_2}{p_1}, \frac{p_3}{p_1}, \dots, \frac{p_k}{p_1}\}$ is a sufficient statistic. Observe that for any j , such that $1 \leq j \leq k$, we have,

$$p_j = \frac{p_j}{p_1} p_1$$

Let $g(T(x), j) = \frac{p_j}{p_1}$ and $h(x) = p_1$, then by Factorization theorem, we have that $T(x)$ is a sufficient statistic.

(b) The joint density of X_1, \dots, X_n is

$$\begin{aligned} f(x_1, \dots, x_n) &= \frac{1}{\theta_2 - \theta_1} \prod_{i=1}^k \mathbf{1}_{\{\theta_1 \leq x_i \leq \theta_2\}} \\ &= \frac{1}{\theta_2 - \theta_1} \mathbf{1}_{\{x_{\min} \geq \theta_1\}} \mathbf{1}_{\{x_{\max} \leq \theta_2\}} \end{aligned}$$

By Factorization theorem, we have the result.

(c)

$$\begin{aligned} p_\theta &= \frac{1}{(\sigma\sqrt{2\pi})^n} \exp\left\{-\frac{(x - A\boldsymbol{\theta})^T(x - A\boldsymbol{\theta})}{2\sigma^2}\right\} \\ &= \frac{1}{(\sigma\sqrt{2\pi})^n} \exp\left\{-\frac{1}{2\sigma^2}(x^T x + \boldsymbol{\theta}^T A^T A\boldsymbol{\theta} - 2x^T A\boldsymbol{\theta})\right\} \end{aligned}$$

Hence, again by Factorization theorem, $A^T x$ is a sufficient statistic.

3: Optimal linear estimation in heteroscedastic Gaussian model

Assume $\sigma_1, \dots, \sigma_d > 0$ to be known, and consider the statistical model $P_\theta = \mathbf{N}(\theta \mathbf{1}, \Sigma)$, where $\Sigma = \text{diag}(\sigma_1^2, \dots, \sigma_d^2)$, and $\theta \in \Theta = \mathbb{R}$ (with $\mathbf{1}$ denoting the all-ones vector). In other words, $X_i = \theta + \sigma_i G_i$ where $(G_i)_{i \leq d} \sim \text{iid } \mathbf{N}(0, 1)$. (Here $\langle \mathbf{u}, \mathbf{v} \rangle = \sum_{i=1}^m u_i v_i$ denotes the usual scalar product of $\mathbf{u}, \mathbf{v} \in \mathbb{R}^m$.)

(a)

$$\begin{aligned} p_\theta &= \frac{1}{(\prod_{i=1}^d \sigma_i \sqrt{2\pi})^n} \exp\left\{-\sum_{i=1}^d \left(\frac{x_i - \theta}{\sigma_i}\right)^2\right\} \\ &= \frac{1}{(\prod_{i=1}^d \sigma_i \sqrt{2\pi})^n} \exp\left\{\sum_{i=1}^d \left(\frac{x_i^2}{\sigma_i^2} + \frac{\theta^2}{\sigma_i^2} - \frac{2\theta x_i}{\sigma_i^2}\right)\right\} \\ &= \frac{1}{(\prod_{i=1}^d \sigma_i \sqrt{2\pi})^n} \exp\left\{\sum_{i=1}^d \frac{x_i^2}{\sigma_i^2}\right\} \exp\left\{\sum_{i=1}^d \frac{\theta^2}{\sigma_i^2}\right\} \exp\left\{-2\theta \sum_{i=1}^d \frac{x_i}{\sigma_i^2}\right\} \end{aligned}$$

So, by Factorization theorem, $l(\mathbf{x}) = \langle \mathbf{c}, \mathbf{x} \rangle$ is a sufficient statistic where $\mathbf{c} = \left(\frac{1}{\sigma_1^2}, \dots, \frac{1}{\sigma_d^2}\right)$

(b) We use Rao-Blackwell theorem to do this problem.

First observe that $\langle \mathbf{a}, \mathbf{1} \rangle = 1$, as

$$\begin{aligned} \mathbb{E}(\langle \mathbf{a}, \mathbf{x} \rangle - \theta)^2 &= \text{Var}(\langle \mathbf{a}, \mathbf{x} \rangle) + (\text{Bias}(\langle \mathbf{a}, \mathbf{x} \rangle))^2 \\ &= \mathbf{a}^T \Sigma \mathbf{a} + \theta^2 (\langle \mathbf{a}, \mathbf{1} \rangle - 1)^2 \end{aligned}$$

So, if $\langle \mathbf{a}, \mathbf{1} \rangle \neq 1$, then we can not avoid having the corresponding minimum value as ∞ .

Now, for any \mathbf{a} such that $\langle \mathbf{a}, \mathbf{1} \rangle = 1$, consider the Rao-Blackwell estimator $\mathbb{E}(\langle \mathbf{a}, \mathbf{x} \rangle | \langle \mathbf{c}, \mathbf{x} \rangle)$ corresponding to the estimator $\langle \mathbf{a}, \mathbf{x} \rangle$, where we have used the sufficiency of $\langle \mathbf{c}, \mathbf{x} \rangle$.

By Rao-Blackwell theorem, the corresponding Rao-Blackwell estimator has lower risk.

So, it suffices to consider only the Rao-Blackwell estimators, i.e., we only to optimize among all Rao-Blackwell estimators such that $\langle \mathbf{a}, \mathbf{1} \rangle = 1$.

Let us find the explicit form of $\mathbb{E}(\langle \mathbf{a}, \mathbf{x} \rangle | \langle \mathbf{c}, \mathbf{x} \rangle)$.

Note,

$$\begin{pmatrix} \langle \mathbf{a}, \mathbf{x} \rangle \\ \langle \mathbf{c}, \mathbf{x} \rangle \end{pmatrix} \sim N \left[\begin{pmatrix} \theta \langle \mathbf{a}, \mathbf{1} \rangle \\ \theta \langle \mathbf{c}, \mathbf{1} \rangle \end{pmatrix}, \begin{pmatrix} \mathbf{a}^T \Sigma \mathbf{a} & \mathbf{a}^T \Sigma \mathbf{c} \\ \mathbf{a}^T \Sigma \mathbf{c} & \mathbf{c}^T \Sigma \mathbf{c} \end{pmatrix} \right]$$

Also note,

$$\begin{aligned} \mathbf{a}^T \Sigma \mathbf{c} &= \langle \mathbf{a}, \mathbf{1} \rangle = 1 \\ \mathbf{c}^T \Sigma \mathbf{c} &= \langle \mathbf{c}, \mathbf{1} \rangle \end{aligned}$$

Using the expression for conditional mean for bivariate Normal, we have,

$$\mathbb{E}(\langle \mathbf{a}, \mathbf{x} \rangle | \langle \mathbf{c}, \mathbf{x} \rangle) = \frac{\langle \mathbf{c}, \mathbf{x} \rangle}{\langle \mathbf{c}, \mathbf{1} \rangle}$$

which is free of \mathbf{a} , so it is the required optimal linear estimator.

[Remark: You can also go about the problem by first invoking Rao-Blackwell Theorem and then observing that $\langle \mathbf{a}, \mathbf{1} \rangle = 1$ or simply use Lagrangian Method for optimization].

(c) Let $X \sim N(\theta \mathbf{1}, \Sigma)$, take $Y = AX$ where A^{-1} is the square root of Σ , so $Y \sim N(\theta A \mathbf{1}, I)$.

Let a_i be the i th row sum of the matrix A .

It can be easily verified that the sufficient statistic is again of the form $l(\mathbf{y}) = \langle \mathbf{c}, \mathbf{y} \rangle$, where $\mathbf{c} = (a_1, \dots, a_d)$.

Again following the same argument as in part b), we have the optimal linear estimator is $\frac{\langle \mathbf{c}, \mathbf{y} \rangle}{\langle \mathbf{c}, A \mathbf{1} \rangle} = \frac{\langle \mathbf{c}, A \mathbf{x} \rangle}{\langle \mathbf{c}, A \mathbf{1} \rangle}$.