

$$1. \quad f(x) = \frac{1}{\sqrt{(2\pi)^k |\Sigma|}} e^{-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)},$$

where $x, \mu \in \mathbb{R}^k$, Σ is a k by k positive definite matrix and $|\Sigma|$ is its determinant.

Show that $\int_{\mathbb{R}^k} f(x) dx = 1$.

Since Σ is a positive-definite matrix.

There exist a orthogonal matrix $P \in M_{k \times k}(\mathbb{R})$ and diagonal matrix $D \in M_{k \times k}(\mathbb{R})$

$$\text{s.t. } \Sigma = P D P^* = P D P^T = P D^{\frac{1}{2}} D^{\frac{1}{2}} P^T = A A^T, \quad \text{where } A = P D^{\frac{1}{2}}, \quad D_{ij}^{\frac{1}{2}} = \sqrt{D_{ij}}$$

$$|\Sigma| = |A A^T|, \quad \text{since } A \text{ is invertible, } |\Sigma| = |A| |A^T| = |A|^2$$

$$\text{Let } y = A^{-1}(x - \mu) \Rightarrow x = \mu + Ay, \quad dx = |A| dy$$

$$\text{and } (x - \mu)^T \Sigma^{-1} (x - \mu) = (Ay)^T (A A^T)^{-1} (Ay) = y^T A^T (A^T)^{-1} A^{-1} Ay = y^T y$$

$$\begin{aligned} \text{Hence, } \int_{\mathbb{R}^k} f(x) dx &= \frac{1}{\sqrt{(2\pi)^k |A|^2}} \int_{\mathbb{R}^k} e^{-\frac{1}{2} y^T y} |A| dy \\ &= \frac{1}{\sqrt{(2\pi)^k |A|^2}} \int_{\mathbb{R}^k} e^{-\frac{1}{2} (y_1^2 + y_2^2 + \dots + y_k^2)} |A| dy_1 dy_2 \dots dy_k \\ &= \frac{1}{\sqrt{(2\pi)^k |A|^2}} |A| \prod_{i=1}^k \int_{\mathbb{R}} e^{-\frac{1}{2} y_i^2} dy_i \\ &= \prod_{i=1}^k \frac{1}{\sqrt{2\pi}} e^{-\frac{y_i^2}{2}} dy_i \\ &= 1 \end{aligned}$$

2.

Let A, B be $n \times n$ matrices and x be $n \times 1$ vector.

(a) Show that $\frac{\partial}{\partial A} \text{trace}(AB) = B^T$

$$\text{pf: } \frac{\partial}{\partial A} \text{trace}(AB) = \frac{\frac{\partial}{\partial a_{ij}} \sum_{i=1}^n \sum_{j=1}^n a_{ij} b_{ji}}{\frac{\partial}{\partial a_{ij}} \sum_{i=1}^n \sum_{j=1}^n a_{ij}} = \sum_{i=1}^n \sum_{j=1}^n b_{ji} = B^T$$

(b) Show that $x^T A x = \text{trace}(x x^T A)$

$$\begin{aligned} x^T A x &= [x_1 \dots x_n] \begin{bmatrix} a_{11} & \dots & a_{1n} \\ \vdots & \ddots & \vdots \\ a_{n1} & \dots & a_{nn} \end{bmatrix} \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} \\ &= \begin{bmatrix} \langle x, A^{(1)} \rangle & \dots & \langle x, A^{(n)} \rangle \end{bmatrix} \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} \\ &= [x_1 \langle x, A^{(1)} \rangle + \dots + x_n \langle x, A^{(n)} \rangle] \\ &= x_1 (a_{11} x_1 + a_{12} x_2 + \dots + a_{1n} x_n) + \dots + x_n (a_{n1} x_1 + a_{n2} x_2 + \dots + a_{nn} x_n) \\ &= \sum_{i=1}^n x_i \left(\sum_{j=1}^n a_{ji} x_j \right) \\ &= \sum_{i=1}^n \sum_{j=1}^n x_i x_j a_{ji} \end{aligned}$$

$$\begin{aligned} \text{trace}(x x^T A) &= \text{trace} \left(\begin{bmatrix} x_1 x_1 & x_1 x_2 & \dots & x_1 x_n \\ x_2 x_1 & x_2 x_2 & \dots & x_2 x_n \\ \vdots & \vdots & \ddots & \vdots \\ x_n x_1 & x_n x_2 & \dots & x_n x_n \end{bmatrix} \begin{bmatrix} a_{11} & \dots & a_{1n} \\ \vdots & \ddots & \vdots \\ a_{n1} & \dots & a_{nn} \end{bmatrix} \right) \\ &= \sum_{i=1}^n \sum_{j=1}^n (x_i x_j) (a_{ji}) \end{aligned}$$

$$\text{Hence, } x^T A x = \sum_{i=1}^n \sum_{j=1}^n x_i x_j a_{ji} = \text{trace}(x x^T A)$$

(c)

Derive the maximum likelihood estimator for multivariate Gaussian.

Let $x_1, \dots, x_n \stackrel{\text{iid}}{\sim} f(x)$

$$L(x_1, \dots, x_n | \mu, \Sigma) = \prod_{i=1}^n (2\pi)^{-\frac{n}{2}} |\Sigma|^{-\frac{1}{2}} e^{-\frac{1}{2} (x_i - \mu)^T \Sigma^{-1} (x_i - \mu)}$$

$$\log L = \ell(\mu, \Sigma) = -\frac{nk}{2} \log(2\pi) - \frac{n}{2} \log |\Sigma| - \frac{1}{2} \sum_{i=1}^n (x_i - \mu)^T \Sigma^{-1} (x_i - \mu)$$

$$\frac{\partial \ell}{\partial \mu} = \frac{\partial}{\partial \mu} \left(-\frac{1}{2} \sum_{i=1}^n (x_i - \mu)^T \Sigma^{-1} (x_i - \mu) \right) = -\frac{1}{2} \sum_{i=1}^n \frac{\partial}{\partial \mu} \left(x_i^T \Sigma^{-1} x_i - x_i^T \Sigma^{-1} \mu - \mu^T \Sigma^{-1} x_i + \mu^T \Sigma^{-1} \mu \right)$$

$$\begin{aligned} &= -\frac{1}{2} \sum_{i=1}^n -(\Sigma^{-1})^T x_i - (\Sigma^{-1}) x_i + \frac{\partial}{\partial \mu} (\text{trace}(\mu \mu^T \Sigma^{-1})) \\ &= -\frac{1}{2} \sum_{i=1}^n -2 \Sigma^{-1} x_i + 2 \mu \Sigma^{-1} \\ &= \sum_{i=1}^n \Sigma^{-1} (x_i - \mu) \\ &= \Sigma^{-1} \left(\sum_{i=1}^n x_i - n \mu \right) = 0 \Rightarrow \hat{\mu}_{MLE} = \frac{1}{n} \sum_{i=1}^n x_i \end{aligned}$$

$$\begin{aligned} \frac{\partial \ell}{\partial \Sigma} &= \frac{\partial}{\partial \Sigma} \left(-\frac{n}{2} \log |\Sigma| - \frac{1}{2} \sum_{i=1}^n (x_i - \mu)^T \Sigma^{-1} (x_i - \mu) \right) \\ &= -\frac{n}{2} (\Sigma^{-1})^T - \frac{1}{2} \sum_{i=1}^n \frac{\partial}{\partial \Sigma} \text{trace}(\Sigma^{-1} (x_i - \mu) (x_i - \mu)^T) \\ &= -\frac{n}{2} (\Sigma^{-1})^T - \frac{1}{2} \sum_{i=1}^n -(\Sigma^{-1} (x_i - \mu) (x_i - \mu)^T \Sigma^{-1})^T = 0 \\ &\Rightarrow \sum_{i=1}^n (\Sigma^{-1})^T (x_i - \mu) (x_i - \mu)^T (\Sigma^{-1})^T = n (\Sigma^{-1})^T \\ &\Rightarrow \frac{1}{n} \sum_{i=1}^n (x_i - \hat{\mu}_{MLE}) (x_i - \hat{\mu}_{MLE})^T = \hat{\Sigma}_{MLE} \end{aligned}$$

3. Unanswer Questions

上課有提到 exponential family 的 canonical form 求 $E(T(\eta))$

2.1 Densities and Parameters

Let μ be a measure on \mathbb{R}^n , let $h: \mathbb{R}^n \rightarrow \mathbb{R}$ be a nonnegative function, and let T_1, \dots, T_s be measurable functions from \mathbb{R}^n to \mathbb{R} . For $\eta \in \mathbb{R}^s$, define

$$A(\eta) = \log \int \exp \left[\sum_{i=1}^s \eta_i T_i(x) \right] h(x) d\mu(x). \quad \begin{array}{l} \text{h.d.} = \text{d.v.} \\ \text{p.e.e.} = \text{v.} \\ \text{v.d.e.} = \text{h.} \end{array} \quad (2.1)$$

Whenever $A(\eta) < \infty$, the function p_η given by

$$p_\eta(x) = \exp \left[\sum_{i=1}^s \eta_i T_i(x) - A(\eta) \right] h(x), \quad x \in \mathbb{R}^n, \quad (2.2)$$

integrates to one; that is, $\int p_\eta d\mu = 1$. So, this construction gives a family of probability densities indexed by η . The set

$$\Xi = \{ \eta : A(\eta) < \infty \}$$

is called the natural parameter space, and the family of densities $\{p_\eta : \eta \in \Xi\}$ is called an s-parameter exponential family in canonical form.

2.2 Differential Identities

In canonical exponential families it is possible to relate moments and cumulants for the statistics T_1, \dots, T_s to derivatives of A . The following theorem plays a central role.

Theorem 2.4. Let Ξ_f be the set of values for $\eta \in \mathbb{R}^s$ where

$$\int |f(x)| \exp \left[\sum_{i=1}^s \eta_i T_i(x) \right] h(x) d\mu(x) < \infty.$$

Then the function $E_{p_\eta} \{ f(T) \}$

$$g(\eta) = \int f(x) \exp \left[\sum_{i=1}^s \eta_i T_i(x) \right] h(x) d\mu(x)$$

is continuous and has continuous partial derivatives of all orders for $\eta \in \Xi_f$ (the interior of Ξ_f). Furthermore, these derivatives can be computed by differentiation under the integral sign.

A proof of this result is given in Brown (1986), a monograph on exponential families with statistical applications. Although the proof is omitted here, key ideas from it are of independent interest and are presented in the next section. As an application of this result, if $f = 1$, then $\Xi_f = \Xi$, and, by (2.1),

$$g(\eta) = e^{A(\eta)} = \int \exp \left[\sum_{i=1}^s \eta_i T_i(x) \right] h(x) d\mu(x).$$

Differentiating this expression with respect to η_j , which can be done under the integral if $\eta \in \Xi^o$, gives

$$\begin{aligned} e^{A(\eta)} \frac{\partial A(\eta)}{\partial \eta_j} &= \int \frac{\partial}{\partial \eta_j} \exp \left[\sum_{i=1}^s \eta_i T_i(x) \right] h(x) d\mu(x) \\ &= \int T_j(x) \exp \left[\sum_{i=1}^s \eta_i T_i(x) \right] h(x) d\mu(x). \end{aligned}$$

Using the definition (2.2) of p_η , division by $e^{A(\eta)}$ gives

$$\frac{\partial A(\eta)}{\partial \eta_j} = \int T_j(x) p_\eta(x) d\mu(x).$$

This shows that if data X has density p_η with respect to μ , then

$$E_{p_\eta} T_j(X) = \frac{\partial A(\eta)}{\partial \eta_j} \quad (2.4)$$

for any $\eta \in \Xi^o$.

$$\frac{\partial A(\eta)}{\partial \eta_j} = \lim_{\epsilon \rightarrow 0} \frac{A(\eta + \epsilon e_j) - A(\eta)}{\epsilon}, \text{ where } e_j = \begin{bmatrix} 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix} \text{ with } 1 \text{ at } j\text{-th position}$$

Derivatives at the origin give cumulants, and thus cumulants for the sum T will equal the sum of the corresponding cumulants of Y_1, \dots, Y_n . This is a well-known result for the mean and variance.

If X has density from a canonical exponential family (2.2), and if $T = T(X)$, then T has moment generating function

$$\begin{aligned} E_\eta e^{u \cdot T(X)} &= \int e^{u \cdot T(x)} e^{\eta \cdot T(x) - A(\eta)} h(x) d\mu(x) \\ &= e^{A(u+\eta) - A(\eta)} \int e^{(u+\eta) \cdot T(x) - A(u+\eta)} h(x) d\mu(x), \end{aligned}$$

provided $u + \eta \in \Xi$. The final integrand is $p_{u+\eta}$, which integrates to one. So, the moment generating function is $e^{A(u+\eta) - A(\eta)}$, and the cumulant generating function is

$$K_T(u) = A(u + \eta) - A(\eta).$$

Taking derivatives, the cumulants for T are

$$\kappa_{r_1, \dots, r_s} = \frac{\partial^{r_1}}{\partial \eta_1^{r_1}} \cdots \frac{\partial^{r_s}}{\partial \eta_s^{r_s}} A(\eta).$$

Example 2.10. If X has the Poisson distribution with mean λ , then

$$P(X = x) = \frac{\lambda^x e^{-\lambda}}{x!} = \frac{1}{x!} e^x \log \lambda - \lambda, \quad x = 0, 1, \dots$$

The mass functions for X form an exponential family, but the family is not in canonical form. The canonical parameter here is $\eta = \log \lambda$. The mass function expressed using η is

$$P(X = x) = \frac{1}{x!} \exp[\eta x - e^\eta], \quad x = 0, 1, \dots,$$

and so $A(\eta) = e^\eta$. Taking derivatives, all of the cumulants of $T = X$ are $e^\eta = \lambda$.

Example 2.11. The class of normal densities formed by varying μ with σ^2 fixed can be written as

$$p_\mu(x) = \exp \left[\frac{\mu x}{\sigma^2} - \frac{\mu^2}{2\sigma^2} \right] \frac{e^{-x^2/(2\sigma^2)}}{\sqrt{2\pi\sigma^2}}.$$

These densities form an exponential family with $T(x) = x$, canonical parameter $\eta = \mu/\sigma^2$, and $A(\eta) = \sigma^2 \eta^2/2$. The first two cumulants are $\kappa_1 = A'(\eta) = \sigma^2 \eta = \mu$ and $\kappa_2 = A''(\eta) = \sigma^2$. Because A is quadratic, all higher-order cumulants, $\kappa_3, \kappa_4, \dots$, are zero.