

1. Given

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

where $x, \mu \in \mathbb{R}^k$, Σ is a k -by- k positive definite matrix and $|\Sigma|$ is its determinant.
Show that $\int_{\mathbb{R}^k} f(x) dx = 1$.

pf: $\because \Sigma$ is positive definite.

\therefore By Cholesky decomposition, \exists ! 對角線都嚴格大於 0 的下三角矩陣 L s.t.
 $\Sigma = LL^T$. (L^{-1} exists)

Let $y = L^{-1}(x-\mu)$, then $x = \mu + Ly$, and $dx = |\det L| dy = |\Sigma|^{\frac{1}{2}} dy$

$$\Rightarrow -\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu) = -\frac{1}{2}(Ly)^T (LL^T)^{-1}(Ly) = -\frac{1}{2}y^T L^T (L^T)^{-1} L^{-1} Ly = -\frac{1}{2}y^T y = -\frac{1}{2}\|y\|^2$$

$$\begin{aligned} \therefore \int_{\mathbb{R}^k} \frac{1}{\sqrt{(2\pi)^k |\Sigma|}} e^{-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)} dx &= \int_{\mathbb{R}^k} \frac{1}{\sqrt{(2\pi)^k |\Sigma|}} e^{-\frac{1}{2}\|y\|^2} |\Sigma|^{\frac{1}{2}} dy \\ &= \int_{\mathbb{R}^k} \frac{1}{\sqrt{(2\pi)^k}} e^{-\frac{1}{2}\|y\|^2} dy \end{aligned}$$

Consider $\int_{\mathbb{R}^k} e^{-\|z\|^2} dz$

By Fubini's Theorem, $\int_{\mathbb{R}^k} e^{-\|z\|^2} dz = \pi^{\frac{k}{2}}$

Let $z = \frac{1}{\sqrt{2}}y$, then $\|z\| = \frac{1}{\sqrt{2}}\|y\| \Rightarrow \|z\|^2 = \frac{1}{2}\|y\|^2$ and $dz = \left(\frac{1}{\sqrt{2}}\right)^k dy$

$$\begin{aligned} \therefore \int_{\mathbb{R}^k} \frac{1}{\sqrt{(2\pi)^k}} e^{-\frac{1}{2}\|y\|^2} dy &= \int_{\mathbb{R}^k} \frac{1}{\sqrt{(2\pi)^k}} e^{-\|z\|^2} \left(\frac{1}{\sqrt{2}}\right)^k dz \\ &= \frac{1}{\sqrt{\pi^k}} \int_{\mathbb{R}^k} e^{-\|z\|^2} dz \\ &= \frac{1}{\sqrt{\pi^k}} \cdot \pi^{\frac{k}{2}} \\ &= 1 \end{aligned}$$

2. Let A, B be n -by- n matrices and x be a n -by-1 vector.

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

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

(c) ~~(b)~~ Derive the maximum likelihood estimators for a multivariate Gaussian.

$$(a) \quad \text{Def: } \left[\frac{\partial}{\partial A} C \right]_{ij} = \frac{\partial}{\partial A_{ij}} C \quad \text{i.e.,} \quad \frac{\partial}{\partial A} C = \begin{bmatrix} \frac{\partial C}{\partial A_{11}} & \frac{\partial C}{\partial A_{12}} & \cdots & \frac{\partial C}{\partial A_{1n}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial C}{\partial A_{n1}} & \frac{\partial C}{\partial A_{n2}} & \cdots & \frac{\partial C}{\partial A_{nn}} \end{bmatrix}$$

$$\text{trace}(AB) = \sum_{p=1}^n (AB)_{pp} = \sum_{p=1}^n \left(\sum_{q=1}^n A_{pq} B_{qp} \right)$$

$$\left[\frac{\partial}{\partial A} \text{trace}(AB) \right]_{ij} = \frac{\partial}{\partial A_{ij}} \text{trace}(AB) = \frac{\partial}{\partial A_{ij}} \left[\sum_{p=1}^n \sum_{q=1}^n A_{pq} B_{qp} \right]$$

$$= \sum_{p=1}^n \sum_{q=1}^n \frac{\partial}{\partial A_{ij}} A_{pq} B_{qp}$$

$$= \sum_{p=1}^n \sum_{q=1}^n \delta_{ip} \delta_{jq} B_{qp} = B_{ji}$$

$$\Rightarrow \frac{\partial}{\partial A} \text{trace}(AB) = \begin{bmatrix} B_{11} & B_{21} & \cdots & B_{n1} \\ B_{12} & B_{22} & \cdots & B_{n2} \\ \vdots & \vdots & \ddots & \vdots \\ B_{1n} & B_{2n} & \cdots & B_{nn} \end{bmatrix} = B^T$$

$$(b) \quad x^T A x = [x_1 \ x_2 \ \cdots \ x_n] \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n1} & a_{n2} & \cdots & a_{nn} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$

$$= x_1 \sum_{p=1}^n x_p A_{p1} + x_2 \sum_{p=1}^n x_p A_{p2} + \cdots + x_n \sum_{p=1}^n x_p A_{pn}$$

$$= \sum_{q=1}^n \left(x_q \sum_{p=1}^n x_p A_{pq} \right) = \sum_{q=1}^n \sum_{p=1}^n x_p x_q A_{pq}$$

$$x x^T = [x_1 \ x_2 \ \cdots \ x_n] \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} = \begin{bmatrix} x_1^2 & x_1 x_2 & \cdots & x_1 x_n \\ x_1 x_2 & x_2^2 & \cdots & x_2 x_n \\ \vdots & \vdots & \ddots & \vdots \\ x_1 x_n & x_2 x_n & \cdots & x_n^2 \end{bmatrix} =: \bar{X}$$

$$\begin{aligned}
 \text{trace}(\chi \chi^T A) &= \sum_{g=1}^n (\bar{\Sigma} A)_{gg} = \sum_{g=1}^n \sum_{p=1}^n \bar{\Sigma}_{gp} A_{pg} \\
 &= \sum_{g=1}^n \sum_{p=1}^n \chi_g \chi_p A_{pg} \quad \because \bar{\Sigma}_{ij} = \chi_i \chi_j \\
 &= \sum_{g=1}^n \sum_{p=1}^n \chi_p \chi_g A_{pg}
 \end{aligned}$$

$$\therefore \chi^T A \chi = \text{trace}(\chi \chi^T A)$$

(c) Let $\chi^{(1)}, \chi^{(2)}, \dots, \chi^{(N)} \in \mathbb{R}^n$ and $\chi^{(k)} \sim N(\mu, \Sigma)$

$$\begin{aligned}
 \text{Likelihood function } L(\mu, \Sigma) &= \prod_{i=1}^N P(\chi^{(i)}) \\
 &= \prod_{i=1}^N \frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma|^{\frac{1}{2}}} e^{-\frac{1}{2}(\chi^{(i)} - \mu)^T \Sigma^{-1} (\chi^{(i)} - \mu)}
 \end{aligned}$$

We want to find μ, Σ s.t. $L(\mu, \Sigma)$ is max.

Define $\ell(\mu, \Sigma) := \ln L(\mu, \Sigma)$

$$\begin{aligned}
 &= \ln \left\{ \left[\frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma|^{\frac{1}{2}}} \right]^N e^{-\frac{1}{2}(\chi^{(1)} - \mu)^T \Sigma^{-1} (\chi^{(1)} - \mu)} \cdots e^{-\frac{1}{2}(\chi^{(N)} - \mu)^T \Sigma^{-1} (\chi^{(N)} - \mu)} \right\} \\
 &= -N \ln \left[(2\pi)^{\frac{n}{2}} |\Sigma|^{\frac{1}{2}} \right] + \sum_{i=1}^N -\frac{1}{2} (\chi^{(i)} - \mu)^T \Sigma^{-1} (\chi^{(i)} - \mu) \\
 &= -\frac{nN}{2} \ln(2\pi) - \frac{N}{2} \ln |\Sigma| - \frac{1}{2} \sum_{i=1}^N (\chi^{(i)} - \mu)^T \Sigma^{-1} (\chi^{(i)} - \mu)
 \end{aligned}$$

$$\because \frac{\partial}{\partial \mu} \left[(\chi - \mu)^T \Sigma^{-1} (\chi - \mu) \right] = -2 \Sigma^{-1} (\chi - \mu)$$

$$\therefore \frac{\partial}{\partial \mu} \ell(\mu, \Sigma) = -\frac{1}{2} \sum_{i=1}^N \left[-2 \Sigma^{-1} (\chi^{(i)} - \mu) \right]$$

$$\text{Let } \frac{\partial}{\partial \mu} \ell(\mu, \Sigma) = 0, \text{ i.e., } \sum_{i=1}^N \Sigma^{-1} (\chi^{(i)} - \mu) = 0$$

$$\Rightarrow \mu N = \sum_{i=1}^N \chi^{(i)} \Rightarrow \mu = \frac{1}{N} \sum_{i=1}^N \chi^{(i)}$$

$$\begin{aligned} \therefore \frac{\partial}{\partial \Sigma} \ln |\Sigma| &= \Sigma^{-1} \quad \text{and} \quad \frac{\partial}{\partial \Sigma} (x-\mu)^T \Sigma^{-1} (x-\mu) = \frac{\partial}{\partial \Sigma} [\text{trace}((x-\mu)(x-\mu)^T \Sigma^{-1})] \\ &= \frac{\partial}{\partial \Sigma} [\text{trace}(\Sigma^{-1} (x-\mu)(x-\mu)^T)] \\ &= -\Sigma^{-1} (x-\mu)(x-\mu)^T \Sigma^{-1} \end{aligned}$$

$$\begin{aligned} \therefore \frac{\partial}{\partial \Sigma} l(\mu, \Sigma) &= -\frac{N}{2} \Sigma^{-1} - \frac{1}{2} \left\{ \sum_{i=1}^N [-\Sigma^{-1} (x^{(i)} - \mu)(x^{(i)} - \mu)^T \Sigma^{-1}] \right\} \\ &= -\frac{N}{2} \Sigma^{-1} + \frac{1}{2} \sum_{i=1}^N \Sigma^{-1} (x^{(i)} - \mu)(x^{(i)} - \mu)^T \Sigma^{-1} \end{aligned}$$

$$\text{Let } \frac{\partial}{\partial \Sigma} l(\mu, \Sigma) = 0, \text{ then } \sum_{i=1}^N \Sigma^{-1} (x^{(i)} - \mu)(x^{(i)} - \mu)^T = N$$

$$\Rightarrow \Sigma = \frac{1}{N} \sum_{i=1}^N (x^{(i)} - \mu)(x^{(i)} - \mu)^T$$

3. Unanswered Questions

There are unanswered questions from the lecture, and there are likely more questions we haven't covered.

- Take a moment to think about these questions.
- Write down the ones you find important, confusing, or interesting.
- You do **not** need to answer them—just state them clearly.

The exponentially family 的形式是 $p(y; \eta) = b(y) \exp(\eta^T y - a(\eta))$
(可用來說明在一些 distribution 下, hypothesis function 為什麼是如此假設)

Q: 為什麼形式是這樣?

從哪裡推導來的?

背後還有什麼更深的理論支持嗎?