

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

$$\text{Let } y = x - \mu \Rightarrow dy = dx$$

$$\text{Then } \int_{\mathbb{R}^k} \frac{1}{\sqrt{(2\pi)^k |\Sigma|}} e^{-\frac{1}{2} y^T \Sigma^{-1} y} dy$$

$\because \Sigma$ is a k -by- k positive definite matrix

\therefore By Cholesky decomposition, we can get $\Sigma = LL^T$, where L is an invertible lower triangular matrix.

$$\text{Then } |\Sigma| = |LL^T| = |L| |L^T| = |L|^2 \Rightarrow |L| = |\Sigma|^{\frac{1}{2}}$$

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

$$\begin{aligned} \text{And } (x - \mu)^T \Sigma^{-1} (x - \mu) &= (x - \mu)^T (LL^T)^{-1} (x - \mu) \\ &= (x - \mu)^T (L^{-T} L^{-1}) (x - \mu) \\ &= (L^{-1}(x - \mu))^T (L^{-1}(x - \mu)) \\ &= y^T y = \|y\|^2 \end{aligned}$$

$$\begin{aligned} \text{Hence, } \int_{\mathbb{R}^k} \frac{1}{\sqrt{(2\pi)^k |\Sigma|}} e^{-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)} dx &= \int_{\mathbb{R}^k} \frac{1}{\sqrt{(2\pi)^k |\Sigma|}} e^{-\frac{1}{2}\|y\|^2} |L| dy \\ &= \frac{1}{\sqrt{(2\pi)^k |\Sigma|}} \cdot |\Sigma|^{\frac{1}{2}} \int_{\mathbb{R}^k} e^{-\frac{1}{2}\|y\|^2} dy \\ &= \frac{1}{\sqrt{(2\pi)^k}} \prod_{i=1}^k \left(\int_{-\infty}^{\infty} e^{-\frac{1}{2}y_i^2} dy_i \right) \quad \text{by Gauss integral, } \int_{-\infty}^{\infty} e^{-\frac{1}{2}x^2} dx = \sqrt{2\pi} \\ &= \frac{1}{\sqrt{(2\pi)^k}} (\sqrt{2\pi})^k = 1. \end{aligned}$$

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

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

(a) Suppose $(AB)_{ij} = \sum_{k=1}^n A_{ik} B_{kj}$

$$\text{trace}(AB) = \sum_{i=1}^n (AB)_{ii} = \sum_{i=1}^n \sum_{k=1}^n A_{ik} B_{ki}$$

$$\frac{\partial}{\partial A_{pq}} \text{trace}(AB) = \frac{\partial}{\partial A_{pq}} \left(\sum_{i=1}^n \sum_{k=1}^n A_{ik} B_{ki} \right), \quad \because \frac{\partial A_{ik}}{\partial A_{pq}} = \begin{cases} 1 & \text{if } i=p \text{ and } k=q \\ 0 & \text{o.w.} \end{cases}$$

$$= 0 \cdot \sum_{i=1}^n \sum_{k \neq q} B_{ki} + 1 \cdot B_{qp}$$

$$= B_{qp}$$

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

(b) $\because x^T A x$ is a scalar, $\therefore x^T A x = \text{trace}(x^T A x)$

Suppose A, B, C are $n \times n$ matrix.

Note that $\text{trace}(ABC) = \text{trace}(BCA) = \text{trace}(CAB)$

Hence, $\text{trace}(x^T A x) = \text{trace}(x x^T A)$.

(c) Suppose $X_i \stackrel{\text{iid}}{\sim} N_n(\mu, \Sigma)$ for $i=1, 2, \dots, m$, where $\mu \in \mathbb{R}^n$ and

Σ is a $n \times n$ matrix.

$$f(x; \mu, \Sigma) = \frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma|^{\frac{1}{2}}} \exp\left(-\frac{1}{2} (x-\mu)^T \Sigma^{-1} (x-\mu)\right)$$

$$L(\mu, \Sigma) = \prod_{i=1}^m f(x_i; \mu, \Sigma) = \frac{1}{(2\pi)^{\frac{mn}{2}} |\Sigma|^{\frac{m}{2}}} \prod_{i=1}^m \exp\left(-\frac{1}{2} (x_i-\mu)^T \Sigma^{-1} (x_i-\mu)\right)$$

$$\ell(\mu, \Sigma) = \ln L(\mu, \Sigma) = -\frac{mn}{2} \ln(2\pi) - \frac{m}{2} \ln|\Sigma| - \frac{1}{2} \sum_{i=1}^m (x_i-\mu)^T \Sigma^{-1} (x_i-\mu)$$

$$\frac{\partial}{\partial \mu} \ell(\mu, \Sigma) = -\frac{1}{2} \sum_{i=1}^m \frac{\partial}{\partial \mu} ((x_i-\mu)^T \Sigma^{-1} (x_i-\mu))$$

$$= -\frac{1}{2} \sum_{i=1}^m (-2 \Sigma^{-1} (x_i-\mu))$$

$$= \Sigma^{-1} \sum_{i=1}^m (x_i-\mu)$$

Let $\frac{\partial}{\partial \mu} \ell(\mu, \Sigma) = 0 \Rightarrow \sum_{i=1}^m (x_i-\mu) = 0 \Rightarrow \hat{\mu} = \frac{\sum_{i=1}^m x_i}{m}$

$\because \sum_{i=1}^m (x_i-\mu)^T \Sigma^{-1} (x_i-\mu)$ is a scalar, \therefore by (b), $\sum_{i=1}^m (x_i-\mu)^T \Sigma^{-1} (x_i-\mu) = \text{trace}\left(\sum_{i=1}^m (x_i-\mu)(x_i-\mu)^T \Sigma^{-1}\right)$

Hence, $\ell(\mu, \Sigma) = -\frac{mn}{2} \ln(2\pi) - \frac{m}{2} \ln|\Sigma| - \frac{1}{2} \sum_{i=1}^m (x_i-\mu)^T \Sigma^{-1} (x_i-\mu)$

$$= -\frac{mn}{2} \ln(2\pi) - \frac{m}{2} \ln|\Sigma| - \frac{1}{2} \text{trace}\left(\sum_{i=1}^m (x_i-\mu)(x_i-\mu)^T \Sigma^{-1}\right)$$

Let $S = \sum_{i=1}^m (x_i-\mu)(x_i-\mu)^T$

$$\frac{\partial}{\partial \Sigma} l(u, \Sigma) = \frac{m}{2} \Sigma - \frac{1}{2} S$$

$$\text{Let } \frac{\partial}{\partial \Sigma} l(u, \Sigma) = 0 \Rightarrow \frac{m}{2} \Sigma - \frac{1}{2} S = 0 \Rightarrow \hat{\Sigma} = \frac{S}{m} = \frac{1}{m} \sum_{i=1}^n (x_i - u)(x_i - u)^T$$

#3.

1. softmax 的對稱性使不同參數組合給出相同輸出，這會如何影響可辨識性與正則化設計？
2. cross-entropy 對於標註噪聲 (label noise) 敏感度如何？什麼情況需要改用其他損失或加權？