

#1

Singular values are square root of eigenvalues of $A^T A$.

Let $A = U \Sigma V^T$ be the SVD of A . For $A \in M_{m \times n}$, we have

$U \in M_{m \times m}$, $V \in M_{n \times n}$ unitary and $\Sigma \in M_{m \times n}$ rectangular diagonal with diagonal elements $\sigma_1 \geq \sigma_2 \geq \dots \geq 0$. Let $\| \cdot \|$ denote $\| \cdot \|_2$.

Claim: $\|A\|_2 = \sigma_1$

(z) Let $V = [v_1 \dots v_n]$, $U = [u_1 \dots u_m]$. As V, U are unitary, $\|v_i\|_2 = \|u_i\|_2 = 1$.

$$\begin{aligned} A v_i &= U \Sigma V^T v_i = U \Sigma e_i, \text{ where } e_i = (1, 0, \dots, 0) \in \mathbb{R}^n \\ &= U \sigma_i \bar{e}_i, \text{ where } \bar{e}_i = (1, 0, \dots, 0) \in \mathbb{R}^m \\ &= \sigma_i (U \bar{e}_i) = \sigma_i u_i. \end{aligned}$$

Thus, $\|A\|_2 = \max_{\|x\|_2=1} \|Ax\|_2 \geq \|Av_i\|_2 = \|\sigma_i u_i\|_2 = \sigma_i \|u_i\|_2 = \sigma_i$.

(\leq) Now, for any $x \in \mathbb{R}^n$ with $\|x\|_2=1$, let $y = V^T x$. Then, $Vy = VV^T x = x$.

As V^T is unitary, $\|y\|_2 = \|V^T x\|_2 = \|x\|_2 = 1$

$Ax = U \Sigma V^T x = U \Sigma y$. As U is unitary, $\|Ax\|_2 = \|U \Sigma y\|_2 = \|\Sigma y\|_2$. $\cdots \textcircled{1}$

$$\|\Sigma y\|_2^2 = \sum_{i=1}^m (\sigma_i y_i)^2 \quad (\sigma_i \text{ might be } 0 \text{ for some } i)$$

$$\leq \sum_{i=1}^m (\sigma_i y_i)^2 = \sigma_1^2 \|y\|_2^2 = \sigma_1^2. \quad \textcircled{2}$$

$\textcircled{1}, \textcircled{2} \Rightarrow \|Ax\|_2 = \|\Sigma y\|_2 \leq \sigma_1$, and we are done. (By taking max. over $\|x\|_2=1$)

#2

$$(a) KL(p || q) := \int p(x) \log \left(\frac{p(x)}{q(x)} \right) dx$$

By convention $\partial \log(\frac{0}{0}) = \partial \log(\frac{0}{a}) = a \log(\frac{a}{0}) = \infty$ from lecture slide,

i) If $[a, b] \subseteq [c, d]$, then integrand $= 0 \log(\frac{0}{?}) = 0$ outside $[a, b]$.

$$\text{Thus, } KL(U_{[a,b]} || U_{[c,d]}) = \int_a^b \frac{1}{b-a} \log \frac{d-c}{b-a} dx = \log \frac{d-c}{b-a}.$$

ii) Otherwise, $a < c$ or $d < b$. Then, integrand $= \frac{1}{b-a} \log(\frac{1}{0})$ on (a, c) or (d, b) .

Integration of f(x) with value ∞ at positive measure set is ∞ .

$$\text{Thus, } KL(U_{[a,b]} || U_{[c,d]}) = \infty.$$

$$\begin{aligned}
 (b) \text{KL}(N(\mu_0, \sigma_0^2) \| N(\mu_1, \sigma_1^2)) &= \int \frac{1}{\sigma_0 \sqrt{2\pi}} e^{-\frac{1}{2\sigma_0^2}(x-\mu_0)^2} \log \left(\frac{\sigma_1}{\sigma_0} e^{\frac{1}{2\sigma_1^2}(x-\mu_1)^2 - \frac{1}{2\sigma_0^2}(x-\mu_0)^2} \right) dx \\
 &= \underbrace{\log \left(\frac{\sigma_1}{\sigma_0} \right)}_{①} \int \frac{1}{\sigma_0 \sqrt{2\pi}} e^{-\frac{1}{2\sigma_0^2}(x-\mu_0)^2} dx + \underbrace{\int \left(\frac{(x-\mu_1)^2}{2\sigma_1^2} - \frac{(x-\mu_0)^2}{2\sigma_0^2} \right) \frac{1}{\sigma_0 \sqrt{2\pi}} e^{-\frac{1}{2\sigma_0^2}(x-\mu_0)^2} dx}_{②}
 \end{aligned}$$

Let's denote $p(x; \mu, \sigma) = \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{1}{2\sigma^2}(x-\mu)^2}$.

$$\text{Note that } ① \int p(x; \mu, \sigma) dx = E[1] = 1$$

$$② \int x p(x; \mu, \sigma) dx = E[X] = \mu \quad \text{for } X \sim N(\mu, \sigma^2)$$

$$③ \int (x-\mu)^2 p(x; \mu, \sigma) dx = E[(X-\mu)^2] = \sigma^2 \quad \text{for } X \sim N(\mu, \sigma^2).$$

By (1), ① = $\log \left(\frac{\sigma_1}{\sigma_0} \right)$.

$$\begin{aligned}
 ② &= \int \left(\frac{(x-\mu_1)^2}{2\sigma_1^2} - \frac{(x-\mu_0)^2}{2\sigma_0^2} \right) p(x; \mu_0, \sigma_0) dx \\
 &= \underbrace{\int \frac{(x-\mu_1)^2}{2\sigma_1^2} p(x; \mu_0, \sigma_0) dx}_{③} - \underbrace{\int \frac{(x-\mu_0)^2}{2\sigma_0^2} p(x; \mu_0, \sigma_0) dx}_{④}.
 \end{aligned}$$

$$\text{By (3), } ④ = \frac{\sigma_0^2}{2\sigma_1^2} = \frac{1}{2}.$$

$$\begin{aligned}
 \text{Finally, } ③ &= \int \frac{1}{2\sigma_1^2} (x-\mu_0 + \mu_0 - \mu_1)^2 p(x; \mu_0, \sigma_0) dx \\
 &= \int \frac{1}{2\sigma_1^2} \left[(x-\mu_0)^2 + 2(\mu_0 - \mu_1)(x-\mu_0) + (\mu_0 - \mu_1)^2 \right] p(x; \mu_0, \sigma_0) dx \\
 &= \frac{1}{2\sigma_1^2} \left[\underbrace{\sigma_0^2}_{\text{By (3)}} + 2(\mu_0 - \mu_1)x_0 + \underbrace{(\mu_0 - \mu_1)^2}_{\text{By (2)}} \times 1 \right] \\
 &= \frac{1}{2\sigma_1^2} [\sigma_0^2 + (\mu_0 - \mu_1)^2]
 \end{aligned}$$

$$\text{Thus, } \text{KL}(N(\mu_0, \sigma_0^2) || N(\mu_1, \sigma_1^2)) = \log\left(\frac{\sigma_1}{\sigma_0}\right) + \frac{1}{2\sigma_1^2} [\sigma_0^2 + (\mu_0 - \mu_1)^2] - \frac{1}{2}$$

$$= \log\left(\frac{\sigma_1}{\sigma_0}\right) + \frac{1}{2\sigma_1^2} [\sigma_0^2 - \sigma_1^2 + (\mu_0 - \mu_1)^2]$$

#3

Notes

When handling multivariate calculus, $x \in \mathbb{R}^n$, $y = f(x) \in \mathbb{R}^m$, we have two choices of representation $\frac{dy}{dx} \in M_{m \times n}$ or $\frac{dy}{x} \in M_{n \times m}$.

For example, $y = Ax$, $A \in M_{m \times n}$ gives $\frac{dy}{dx} = A$ in first representation
 $\frac{dy}{x} = A^T$ in the 2nd representation

This matters when applying chain rule. $x \in \mathbb{R}^n$, $y = f(x) \in \mathbb{R}^m$, $z = g(y) \in \mathbb{R}^l$.

$$\frac{dz}{dx} = \frac{dz}{dy} \times \frac{dy}{dx} \in M_{l \times n}, \text{ with } \frac{dz}{dy} \in M_{l \times m}, \frac{dy}{dx} \in M_{m \times n} \text{ in 1st rep.}$$

$$\frac{dz}{dx} = \frac{dy}{dx} \times \frac{dz}{dy} \in M_{n \times l}, \text{ with } \frac{dz}{dy} \in M_{m \times l}, \frac{dy}{dx} \in M_{m \times n} \text{ in 2nd rep.}$$

Let's be consistent with the 1st representation.

For example, $\frac{d}{dt} t^T t = 2t^T \in M_{1 \times n}$ for $t \in \mathbb{R}^n$.

(a) By chain rule, $\frac{d}{dt} \sin(\log(t^T t)) = \cos(\log(t^T t)) \cdot \frac{1}{t^T t} \cdot 2t^T$

(b) For simplicity, for $M \in M_{m \times n}$, $y = f(M) \in \mathbb{R}$, let $\frac{dy}{dM} \in M_{m \times n}$
with $\left(\frac{df}{dM}\right)_{ij} = \frac{df}{dM_{ij}}$.

$$\begin{aligned} \text{tr}(AXB) &= \sum_{i=1}^D (AXB)_{ii} \\ &= \sum_{i=1}^D \sum_{j=1}^E A_{ij} (XB)_{ji} \\ &= \sum_{i=1}^D \sum_{j=1}^E \sum_{k=1}^F A_{ij} X_{jk} B_{ki} \\ &= \sum_{j=1}^E \sum_{k=1}^F X_{jk} \left(\sum_{i=1}^D B_{ki} A_{ij} \right) = \sum_{j=1}^E \sum_{k=1}^F X_{jk} (BA)_{kj} \end{aligned}$$

Hence, $\frac{d}{dx_{jk}} \text{tr}(AXB) = (BA)_{kj}$. i.e. $\frac{d}{dx} \text{tr}(AXB) = (BA)^T$

Joint distribution

#4

(a) Given Samples x_1, \dots, x_n of X_1, \dots, X_n ,

$$\begin{aligned} l &:= \log P(X_1=x_1, \dots, X_n=x_n) = \log \prod_{i=1}^n P(X_i=x_i) \\ &= \sum_{i=1}^n \log P(X_i=x_i) = \sum_{i=1}^n \log \frac{1}{(2\pi)^{k/2} |\Sigma^{1/2}|} e^{-\frac{1}{2}(x_i-\mu)^T \Sigma^{-1} (x_i-\mu)} \\ &= nC - \sum_{i=1}^n \frac{1}{2} (x_i-\mu)^T \Sigma^{-1} (x_i-\mu) \end{aligned}$$

$$\frac{\partial l}{\partial \mu} = \sum_{i=1}^n (\bar{x}_i - \mu) = \sum_{i=1}^n (\bar{x}_i - n\bar{\mu})$$

$$= n(\bar{x} - \mu), \text{ where } \bar{x} = \frac{1}{n} \sum_{i=1}^n x_i.$$

Thus, $\frac{\partial l}{\partial \mu} = 0$ gives $\hat{\mu} = \bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$. (MLE)

$$(b) \text{ Posterior } p(\mu | x_1, \dots, x_n) = \frac{p(x_1, \dots, x_n | \mu) p(\mu)}{\int p(x_1, \dots, x_n | \mu) p(\mu) d\mu}$$

$$\begin{aligned} \bar{l} &= \frac{\prod_{i=1}^n N(x_i | \mu, \Sigma) N(\mu | \mu_0, \Sigma_0)}{\int \prod_{i=1}^n N(x_i | \mu, \Sigma) d\mu} \\ &= \frac{\frac{1}{(2\pi)^{nk/2} |\Sigma|^{n/2}} e^{-\frac{1}{2} \sum_{i=1}^n (x_i-\mu)^T \Sigma^{-1} (x_i-\mu)}}{(2\pi)^{n/2} |\Sigma_0|^{n/2}} \frac{1}{(2\pi)^{k/2} |\Sigma_0|^{k/2}} e^{-\frac{1}{2} (\mu-\mu_0)^T \Sigma_0^{-1} (\mu-\mu_0)} \\ &\quad (\text{constant wrt } \mu) \end{aligned}$$

As denominator of (*) is indep. of μ ,

$$\begin{aligned} \frac{d}{d\mu} \log p(\mu | x_1, \dots, x_n) &= \frac{d}{d\mu} \left[(\text{constant wrt } \mu) - \frac{1}{2} \sum_{i=1}^n (x_i-\mu)^T \Sigma^{-1} (x_i-\mu) \right. \\ &\quad \left. - \frac{1}{2} (\mu-\mu_0)^T \Sigma_0^{-1} (\mu-\mu_0) \right] \\ &= -n\Sigma^{-1}(\bar{x}-\mu) - \Sigma_0^{-1}(\mu-\mu_0) \end{aligned}$$

Hence, $\frac{\partial \bar{l}}{\partial \mu} = 0$ gives

$$n\Sigma^{-1}(\bar{x}-\mu) = \Sigma_0^{-1}(\mu-\mu_0)$$

$$\Rightarrow (n\Sigma^{-1} + \Sigma_0^{-1})\mu = n\Sigma^{-1}\bar{x} + \Sigma_0^{-1}\mu_0$$

$$\Rightarrow \mu = (n\Sigma^{-1} + \Sigma_0^{-1})^{-1} (n\Sigma^{-1}\bar{x} + \Sigma_0^{-1}\mu_0)$$

i.e. MAP is $\hat{\mu} = (n\Sigma^{-1} + \Sigma_0^{-1})^{-1} (n\Sigma^{-1}\bar{x} + \Sigma_0^{-1}\mu_0)$

#5

Note that if $L: \mathbb{R}^n \rightarrow \mathbb{R}$ is linear, $\varphi: \mathbb{R} \rightarrow \mathbb{R}$ is convex, then $\varphi \circ L$ is convex.
 $\because \varphi \circ L(\lambda x + (1-\lambda)y) = \varphi(\lambda L(x) + (1-\lambda)L(y)) \leq \lambda \varphi(L(x)) + (1-\lambda)\varphi(L(y)).$

① First we show $g: \mathbb{R} \rightarrow \mathbb{R}$ by $g(t) = -\log \sigma(t)$ is convex.

$$\therefore \sigma(t) = \frac{1}{1+e^{-t}}, \quad \sigma'(t) = \frac{e^{-t}}{(1+e^{-t})^2} = e^{-t} \sigma(t)^2.$$

$$\Rightarrow g'(t) = -\frac{\sigma'(t)}{\sigma(t)} = -e^{-t} \sigma(t)$$

$$\Rightarrow g''(t) = e^{-t} \sigma(t) - e^{-t} \sigma'(t)$$

$$= e^{-t} (\sigma(t) - \sigma'(t))$$

$$= e^{-t} (\sigma(t) - e^{-t} \sigma(t)^2)$$

$$= e^{-t} \sigma(t)^2 \left(\frac{1}{\sigma(t)} - e^{-t} \right)$$

$$= e^{-t} \sigma(t)^2 (1+e^{-t} - e^{-t}) = e^{-t} \sigma(t)^2 > 0.$$

Thus, g is convex.

② $h: \mathbb{R} \rightarrow \mathbb{R}$ by $h(t) = -\log(1-\sigma(t))$ is convex.

$$\therefore h(t) = -\log(\sigma(-t)) = g(-t), \text{ where } g \text{ is convex, and } t \mapsto -t \text{ is linear.}$$

$$\sigma(t) + \sigma(-t) = 1.$$

③ f is convex.

$L_i(\theta)$ by $L_i(\theta) = x_i^\top \theta$ is linear, g, h are convex.

$$f = \sum_{i=1}^N y_i g(L_i(\theta)) + (1-y_i) h(L_i(\theta)).$$

$\theta \mapsto g(L_i(\theta))$ is convex, $\theta \mapsto h(L_i(\theta))$ is convex.

$y_i \geq 0, 1-y_i \geq 0$. Thus, f is a positive linear combination of convex functions and so convex.