Notes on "Pattern Recognition and Machine Learning"

Zhen Wang

July 18, 2013

Contents

The more evenly a distribution spreads, the higher its entropy is. Besides, the *noiseless coding theorem* states that:

Conclusion 1.1.1. The entropy is a lower bound on the average number of bits needed to transmit the state of a random variable. We can arrive the lower bound by choosing an efficient coding scheme.

The maximum entropy configuration can be found by maximizing the following equation which includes the normalization constraint on the probability distribution:

$$\tilde{H} = -\sum_{i=1}^{M} p_i \ln(p_i) + \lambda (\sum_{i=1}^{M} p_i - 1)$$
(3)

where the λ is a Lagrangian multiplier and doesn't need to be determined. To find stationary points, we check its derivatives:

$$\frac{\partial \tilde{H}}{\partial p_i} = -(\ln(p_i) + 1) + \lambda = 0 \text{ for } i = 1, \dots, M$$

$$\sum_{i=1}^{M} p_i - 1 = 0$$

We have stationary $p_i = \frac{1}{M}$, i = 1, ..., M and its second derivative is:

$$\frac{\partial \tilde{H}}{\partial p_i \partial p_j} = -I_{ij} \frac{1}{p_i} = -I_{ij} M \tag{4}$$

Since M > 0, the second derivative is negative definite (diagonal matrix's eigenvalues are elements on its diagonal), thus the stationary point is indeed a maximum.

The entropy for continuous variable can't be elegantly defined and we use the concept—differential entropy:

Definition 1.1.3.

$$H[X] = -\int \Pr(x) \ln \Pr(x) dx$$
 (5)

We find the maximum of H[x] using Lagrange multipliers with normalization constraint of probability distribution and contraints over the first and second moments of Pr(X):

$$\lambda_{1}\left(\int_{-\infty}^{\infty} \Pr(x) dx - 1\right)$$

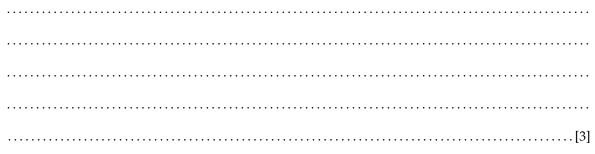
$$\lambda_{2}\left(\int_{-\infty}^{\infty} x \Pr(x) dx - \mu\right)$$

$$\lambda_{3}\left(\int_{-\infty}^{\infty} (x - \mu)^{2} \Pr(x) dx - \sigma^{2}\right)$$

Ising the calculus of variations and set the derivative of the functional to zero giving us:	
	[2]
$Pr(X) = \exp\{-1 + \lambda_1 + \lambda_2 X + \lambda_3 (X - \mu)^2\}$	(6)

......

back substitution of (6) into the three constraint equations gives us:



$$\Pr(X) = \frac{1}{(2\pi\sigma^2)^{\frac{1}{2}}} \exp\left\{-\frac{(X-\mu)^2}{2\sigma^2}\right\}$$
 (7)

so we have:

Conclusion 1.1.2. The maximum entropy configuration for continuous random variable is the Guassian.

Then we evaluate its differential entropy. $\ln \Pr(X) = \frac{-1}{2} \ln (2\pi\sigma^2) - \frac{(X-\mu)^2}{2\sigma^2}$ and $-\int_{-\infty}^{\infty} \Pr(X) \frac{-1}{2} \ln (2\pi\sigma^2) dX = \frac{1}{2} \ln (2\pi\sigma^2) \ln (2\pi\sigma^2)$ $\frac{1}{2}\ln(2\pi\sigma^2)$. So the difficulty is to calculate

$$\frac{1}{(2\pi\sigma^2)^{\frac{1}{2}}} \int_{-\infty}^{\infty} \exp\left\{\frac{-(X-\mu)^2}{2\sigma^2}\right\} \frac{(X-\mu)^2}{2\sigma^2} dX$$

$$= \frac{(2\sigma^2)^{\frac{1}{2}}}{(2\pi\sigma^2)^{\frac{1}{2}}} \int_{-\infty}^{\infty} \exp\left\{\frac{-(X-\mu)^2}{2\sigma^2}\right\} \frac{(X-\mu)^2}{2\sigma^2} d\frac{(X-\mu)}{(2\sigma^2)^{\frac{1}{2}}}$$
(8)

By subsection integral method, we know that:

$$\int_{-\infty}^{\infty} x^{2} \exp(-x^{2}) dx = \int_{-\infty}^{\infty} \frac{d \exp(-x^{2})}{dx} \frac{-x}{2} dx$$

$$= \exp(-x^{2}) \frac{-x}{2} \Big|_{-\infty}^{\infty} - \int_{-\infty}^{\infty} \exp(-x^{2}) \frac{d \frac{-x}{2}}{dx} dx$$

$$= 0 + \frac{1}{2} \int_{-\infty}^{\infty} \exp(-x^{2}) dx$$

$$= \frac{\pi^{\frac{1}{2}}}{2}$$
(9)

Thus, based on (9) and (8), we have $H[X] = \frac{1}{2}\{1 + \ln(2\pi\sigma^2)\}$. Suppose we have a joint distribution $\Pr(X,Y)$. If we observe a specific value X=x, then the additional information needed to specify the corresponding value of Y = y is given by $-\ln \Pr(y|x)$. Intuitively, h(x, y) = h(x) + h(y|x), by the production rule of probability, $-\ln \Pr(x, y) = -\ln \{\Pr(x) \Pr(y|x)\} = -\ln \{\Pr(x, y) = -\ln \{\Pr(x) \Pr(y|x)\} = -\ln \{\Pr(x) \Pr(x)\} = -\ln \{\Pr(x)\} = -\ln \{\Pr(x)\}$ $-\ln \Pr(x) - \ln \Pr(y|x)$, which comfirms that our definition makes sense.

Then we define the *conditional entropy* of Y given X as the average additional information needed to specify Y:

Definition 1.1.4.

$$H[Y|X] = -\int \int Pr(X,Y) \ln Pr(Y|X) dY dX$$
 (10)

Then we confirm that H[X,Y] = H[Y|X] + H[X]:

$$H[X,Y] = -\int \int \Pr(X,Y) \ln \Pr(X,Y) dXdY$$

$$= -\int \int \Pr(X,Y) \ln \Pr(Y|X) + \Pr(X,Y) \ln \Pr(X) dXdY$$

$$= H[Y|X] - \int \ln \Pr(X) dX \int \Pr(X,Y) dY$$

$$= H[Y|X] - \int \Pr(X) \ln \Pr(X) dX$$

$$= H[Y|X] + H[X]$$
(11)

1.1.1 Relative entropy and mutual information

Definition 1.1.5. *Convex* functions are functions that satisfy the following inequality:

$$f(\lambda a + (1 - \lambda)b) \le \lambda f(a) + (1 - \lambda)f(b) \tag{12}$$

where $0 \le \lambda \le 1$.

A convex function satisfies the Jensen's inequality:

$$f(\sum_{i=1}^{M} \lambda_i x_i) \le \sum_{i=1}^{M} \lambda_i f(x_i)$$
(13)

where $\lambda_i \geq 0$ and $\sum_{i=1}^{M} \lambda_i = 1$. We can prove the above inequality by induction:

Proof.

$$\sum_{i=1}^{N+1} \lambda_i f(x_i) = (1 - \lambda_{N+1}) \sum_{i=1}^{N} \frac{\lambda_i}{1 - \lambda_{N+1}} f(x_i) + \lambda_{N+1} f(x_{N+1})$$
(14)

: (13) is valid for previous N cases : (14) $\geq (1 - \lambda_{N+1}) f(\sum_{i=1}^{N} \frac{\lambda_i}{1 - \lambda_{N+1}} x_i) + \lambda_{N+1} f(x_{N+1})$

 $f(\cdot)$ is a convex function $f(\cdot)$ we have:

$$(1 - \lambda_{N+1}) f(\sum_{i=1}^{N} \frac{\lambda_i}{1 - \lambda_{N+1}} x_i) + \lambda_{N+1} f(x_{N+1}) \ge f(\sum_{i=1}^{N+1} \lambda_i x_i)$$
(15)

Thus, we derive the N + 1 case from previous case.

We can interpret λ_i as the probability over $X = x_i$, then (13) can be written as:

$$f(E[X]) \le E[f(X)] \tag{16}$$

which is also valid for continuous case.

Suppose a random variable X is transmitted and its distribution p(X) is unknown. We use distribution q(X) to approximate it. Thus the average *additional* amount of information required to specify X = x is:

Definition 1.1.6.

$$KL(p||q) = -\int p(X) \ln q(X) dX - \left(-\int p(X) \ln p(X) dX\right)$$

$$= -\int p(X) \ln \left\{\frac{q(X)}{p(X)}\right\} dX$$
(17)

Obviously, *Kullback-Leibler divergence* is not symmetric. Since $(-\ln x)'' = \frac{1}{x^2} > 0$, $-\ln x$ is a convex function. Thus we can apply Jensen inequality to it:

$$\int -\ln\left\{\frac{q(X)}{p(X)}\right\}p(X)dX \ge -\ln\int p(X)\frac{q(X)}{p(X)}dX = 0$$
(18)

Obviously, KL(p||q) = 0 only when q(X) = p(X). Besides, we have:

Conclusion 1.1.3. KL divergence is the *lower bound* for average additional information that must be transmitted. We can arrive the bound by choosing an efficient coding scheme.

Both data compression and density estimation are aimed at modelling an unknown probability distribution. Suppose q(X) is characterized by a bunch of parameters θ and we estimate θ by minimizing KL divergence. However, p(X) is unknown, we only have its samples x_i , i = 1, ..., N and we use this traning set to approximate p(X):

$$KL(p||q) \simeq \sum_{i=1}^{N} \{-\ln q(x_i|\theta) + \ln p(x_i)\}$$
 (19)

Note that the second term on the right-hand side of (19) is indepedent of θ , and the first term is the negative log likelihood function for θ . Thus we have:

Conclusion 1.1.4. Minimizing KL divergence ⇔ Maximizing the likelihood function.

Given two random variables X, Y, they are independent if Pr(X,Y) = Pr(X) Pr(Y). We can measure how independent they are by:

Definition 1.1.7.

$$I[X,Y] \equiv KL(\Pr(X,Y) \| \Pr(X) \Pr(Y))$$

$$= -\int \int \Pr(X,Y) \ln\left(\frac{\Pr(X) \Pr(Y)}{\Pr(X,Y)}\right) dXdY$$
(20)

which is called the *mutual information* between *X* and *Y*.

Mutual information is related to conditional entropy:

$$I[X,Y] = -\int \int \Pr(X,Y) \ln\left(\frac{\Pr(X)\Pr(Y)}{\Pr(X,Y)}\right) dXdY$$

$$= H[X] + H[Y] - \left(-\int \int \Pr(X,Y) \ln \Pr(X,Y) dXdY\right)$$

$$= H[X] + H[Y] - H[X,Y]$$

$$= H[X] - H[X|Y] = H[Y] - H[Y|X] \quad \text{(by (11))}$$
(21)

Thus, the mutual information can also be viewed as the residue in the uncertainty about *X* after being told the value of *Y* (or vice versa).

2 Probability Distribution

2.1 The Gaussian Distribution

Preliminaries:

Conclusion 2.1.1. For quadratic form $\mathbf{x}^{T}A\mathbf{x}$, matrix A can be taken to be symmetric without loss of generality.

Proof.
$$\mathbf{x}^{\mathrm{T}}A\mathbf{x} = \sum_{i=1}^{D} \sum_{j=1}^{D} x_{i}x_{j}A_{ij}$$
 $\mathbf{x}^{\mathrm{T}}A\mathbf{x} = \sum_{i=1}^{D} \sum_{j=1}^{D} x_{i}x_{j}A_{ij}$ $\mathbf{x}^{\mathrm{T}}A\mathbf{x} = \sum_{i=1}^{D} \sum_{j=1}^{D} x_{i}A_{ij}$

Conclusion 2.1.2. Matrix *A* is symmetric $\Leftrightarrow A\mathbf{x} \cdot \mathbf{y} = \mathbf{x} \cdot A\mathbf{y}$.

Proof.
$$(\Rightarrow) : A\mathbf{x} \cdot \mathbf{y} = (A\mathbf{x})^{\mathrm{T}} \mathbf{y} = \mathbf{x}^{\mathrm{T}} A \mathbf{y} = \mathbf{x} \cdot A \mathbf{y}$$

 $(\Leftarrow) A \mathbf{e}_i \cdot \mathbf{e}_j = A_{ji} = \mathbf{e}_i \cdot A \mathbf{e}_j = A_{ij}$
 $: \forall (i,j) : A_{ij} = A_{ji} : A = A^{\mathrm{T}}$

Conclusion 2.1.3. Real symmetric matrix $A_{n \times n}$ has n real eigenvalues (counting multiplicities).

Proof. Suppose
$$A\mathbf{x} = \lambda \mathbf{x}$$
 where $\mathbf{x} \in \mathbb{C}^n$. Let $q = \bar{\mathbf{x}}^T A \mathbf{x}$
 $\because \bar{q} = \mathbf{x}^T A \bar{\mathbf{x}} = \bar{q}^T = \bar{\mathbf{x}}^T A \mathbf{x} = q \quad \therefore q \text{ is actually real.}$
 $\because q = \bar{\mathbf{x}}^T A \mathbf{x} = \bar{\mathbf{x}}^T \lambda \mathbf{x} = \lambda \|\mathbf{x}\|^2 \quad \therefore \lambda \text{ is also real.}$

Conclusion 2.1.4. Symmetric matrix $A_{n \times n}$ has n eigenvectors forming an orthonormal set.

Proof. (Assumption) for $i = 1, ..., n-1 \quad (n > 1)$, symmetric matrix $A_{i \times i}$ has i orthonormal eigenvectors and $A_{n \times n}$ is symmetric but doesn't have n orthonormal eigenvectors.

Pick λ_1 , one of $A_{n \times n}$'s n eigenvalues and pick a corresponding unit eigenvector \mathbf{u}_1 (i.e., $A\mathbf{u}_1 = \lambda_1\mathbf{u}_1$). Let $W = \text{Nul }\mathbf{u}_1^T$ be the null space of \mathbf{u}_1 . Obviously, it is a n-1 dimensional subspace of \mathbb{R}^n . Then we could choose the basis for W consisting of n-1 orthonormal vectors $\mathbb{B} = \{\mathbf{u}_2, \dots, \mathbf{u}_n\}$ and let $P_{\mathbb{B}} = [\mathbf{u}_2, \dots, \mathbf{u}_n]$. Although $P_{\mathbb{B}}$ is not square, its rank is n-1 and thus gives an one-to-one and onto mapping from \mathbb{R}^{n-1} to W which implies the existence of its *left inverse* $P_{\mathbb{B}}^{-1}$.

$$W = \text{Nul } \mathbf{u_1}^{\text{T}} \qquad \qquad \Re^{n-1}$$

$$\mathbf{w} \qquad \qquad \stackrel{P_{\mathcal{B}}^{-1}}{\longleftarrow} \qquad [\mathbf{w}]_{\mathcal{B}}$$

Then we prove that multiplication by A defines a linear transformation $T: W \to W$: $\therefore A$ is symmetric. $\therefore \forall \mathbf{w} \in W : \mathbf{u}_1 \cdot A\mathbf{w} = A\mathbf{u}_1 \cdot \mathbf{w} = \lambda_1 \mathbf{u}_1 \cdot \mathbf{w} = 0$ Note that the codomain of mapping T is W but the range may not.

$$W = \{ \mathbf{w} \in \Re^n | \mathbf{u_1}^T \mathbf{w} = 0 \} \xrightarrow{T : A} V = \{ A \mathbf{w} | \mathbf{w} \in W \} \subseteq W$$

$$\downarrow P_{\mathcal{B}}^{-1} \qquad \qquad \downarrow P_{\mathcal{B}}^{-1}$$

$$\Re^{n-1} = \{ [\mathbf{w}]_{\mathcal{B}} | \mathbf{w} \in W \} \xrightarrow{T' : M} M[\mathbf{w}]_{\mathcal{B}} = [A \mathbf{w}]_{\mathcal{B}} \in \Re^{n-1}$$

Now our task is to prove the existence of matrix *M* in the above relationships:

Suppose $\mathbf{w} = r_2 \mathbf{u}_2 + \ldots + r_n \mathbf{u}_n$ and since T is a linear transformation which preserves addition and scalar multiplication, we have:

$$T(\mathbf{w}) = r_2 T(\mathbf{u}_2) + \ldots + r_n T(\mathbf{u}_n)$$
(22)

Rewrite (22) in basis B we have:

$$[T(\mathbf{w})]_{\mathbb{B}} = r_2[T(\mathbf{u}_2)]_{\mathbb{B}} + \ldots + r_n[T(\mathbf{u}_n)]_{\mathbb{B}}$$
(23)

Note that $[\mathbf{w}]_{\mathbb{B}} = [r_2 \quad \cdots \quad r_n]^{\mathrm{T}}$ and thus:

$$[T(\mathbf{w})]_{\mathbb{B}} = M[\mathbf{w}]_{\mathbb{B}} \tag{24}$$

where $M = [[T(\mathbf{u}_2)]_{\mathbb{B}} \cdots [T(\mathbf{u}_n)]_{\mathbb{B}}]$

M is related to *A* in terms of eigenvalues and eigenvectors:

Suppose $M[\mathbf{x}]_{\mathbb{B}} = \lambda[\mathbf{x}]_{\mathbb{B}}$, according to (24), $[A\mathbf{x}]_{\mathbb{B}} = \lambda[\mathbf{x}]_{\mathbb{B}}$. \therefore $[\cdot]_{\mathbb{B}}$ is actually multiplication by $P_{\mathbb{B}}^{-1}$: $W \to \mathbb{R}^{n-1}$ which is a linear transformation with one-to-one and onto properties. \therefore \mathbf{x} and $[\mathbf{x}]_{\mathbb{B}}$ possess uniquely correspondence and \mathbf{x} is eigenvector of A. Besides, since we choose the inverse mapping $P_{\mathbb{B}}$ to be constructed by basis of W, $\mathbf{x} \in W$.

If there is an orthonormal basis of \mathbb{R}^{n-1} consisting of n-1 eigenvectors of $M: [\mathbf{x}_2]_{\mathbb{B}}, \ldots, [\mathbf{x}_n]_{\mathbb{B}}$, then $\mathbf{x}_2, \ldots, \mathbf{x}_n$ are eigenvectors of $A: \forall i=2,\ldots,n: \mathbf{x}_i \in W \quad \therefore \mathbf{u}_1 \cdot \mathbf{x}_i = 0$. Now the remaining task is to prove $\mathbf{x}_2, \ldots, \mathbf{x}_n$ are orthonormal. Remember that $[\mathbf{x}_2]_{\mathbb{B}}, \ldots, [\mathbf{x}_n]_{\mathbb{B}}$ are orthonormal so that if multiplication by $P_{\mathbb{B}}$ preserves dot products (in turn length), then their correspondences are also orthonormal. So the next steps are:

- multiplication by $P_{\mathbb{B}}$ preserves dot products.
- M has n-1 orthonormal eigenvectors.

Preserving dot products means $\forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^{n-1} : \mathbf{x} \cdot \mathbf{y} = (P_{\mathbb{B}}\mathbf{x}) \cdot (P_{\mathbb{B}}\mathbf{y})$. It is straightforward by noticing that $P_{\mathbb{B}}^{\mathsf{T}}P_{\mathbb{B}} = I_{(n-1)\times(n-1)}$. Besides, equality is symmetric and we can say $[\cdot]_{\mathbb{B}}$ preserves dot products (i.e., both $P_{\mathbb{B}}$ and $P_{\mathbb{B}}^{-1}$).

To show that M has n-1 orthonormal eigenvectors seems difficult, however, note that by induction, a symmetric matrix $M_{(n-1)\times(n-1)}$ has n-1 orthonormal eigenvectors. Thus, we try to prove M is symmetric.

 \therefore *A* is symmetric. \therefore $A\mathbf{x} \cdot \mathbf{y} = \mathbf{x} \cdot A\mathbf{y} \cdot \therefore$ []_B is an isomorphism and preserves dot products. \therefore according to (24), $M[\mathbf{x}]_{\mathbb{B}} \cdot [\mathbf{y}]_{\mathbb{B}} = [\mathbf{x}]_{\mathbb{B}} \cdot M[\mathbf{y}]_{\mathbb{B}}$ which implies that M is symmetric.

Up to now, we showed that A has n orthonormal eigenvectors $\{\mathbf{u}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$ which contradicts our initial assumption. Thus, there doesn't exist such a n and $\forall n$: symmetric matrix $A_{n \times n}$ has n orthonormal eigenvectors.

Conclusion 2.1.5. Symmetric matrix A is orthonormally diagonalizable, i.e., $A = U\Lambda U^{T}$ where U is orthogonal (i.e., columns are unit orthogonal vectors and rows are unit orthogonal vectors) and Λ is diagonal matrix.

Proof. We have proved that symmetric matrix $A_{n \times n}$ has n orthonormal eigenvectors $\{\mathbf{u}_1, \dots, \mathbf{u}_n\}$ $\{A_{i} = \mathbf{u}_i\}$

$$\lambda_i \mathbf{u}_i$$
). Let $U = [\mathbf{u}_1 \quad \cdots \quad \mathbf{u}_n]$ and $\Lambda = \begin{bmatrix} \lambda_1 \\ & \ddots \\ & & \lambda_n \end{bmatrix}$, we show that $A = U \Lambda U^T$: $\therefore U$ consists of

unit orthonormal columns. \therefore right multiplied by $U, AU = [\lambda_1 \mathbf{u}_1 \cdots \lambda_n \mathbf{u}_n]$ and $U \wedge U^T U = U \wedge = [\lambda_1 \mathbf{u}_1 \cdots \lambda_n \mathbf{u}_n]$

The multivariate Gaussian distribution takes the form:

Definition 2.1.1.

$$\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{\frac{D}{2}}} \frac{1}{|\boldsymbol{\Sigma}|^{\frac{1}{2}}} \exp\left\{-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^{\mathrm{T}} \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})\right\}$$
(25)

where μ is the mean vector, Σ is a $D \times D$ covariance matrix with its determinant $|\Sigma|$. The functional dependence of the Gaussian on \mathbf{x} is through the quadratic form:

$$\Delta^2 = (\mathbf{x} - \boldsymbol{\mu})^{\mathrm{T}} \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu})$$
 (26)

The quantity Δ is called the *Mahalanobis distance* from μ to x.

(Based on above propositions) We can choose Σ^{-1} to be symmetric without loss of generality. Actually, we firstly choose the covariance matrix Σ to be symmetric, then Σ is orthogonal diagonalizable:

$$\Sigma = U\Lambda U^{T}$$

$$= \begin{bmatrix} \mathbf{u}_{1} & \cdots & \mathbf{u}_{D} \end{bmatrix} \begin{bmatrix} \lambda_{1} & & \\ & \ddots & \\ & \lambda_{D} \end{bmatrix} \begin{bmatrix} \mathbf{u}_{1}^{T} \\ \vdots \\ \mathbf{u}_{D}^{T} \end{bmatrix}$$

$$= \begin{bmatrix} \lambda_{1}\mathbf{u}_{1} & \cdots & \lambda_{D}\mathbf{u}_{D} \end{bmatrix} \begin{bmatrix} \mathbf{u}_{1}^{T} \\ \vdots \\ \mathbf{u}_{D}^{T} \end{bmatrix}$$

$$= \sum_{i=1}^{D} \lambda_{i}\mathbf{u}_{i}\mathbf{u}_{i}^{T}$$

$$(27)$$

Because Λ is diagonal matrix, its inverse is $\Lambda^{-1} = \begin{bmatrix} \frac{1}{\lambda_1} & & \\ & \ddots & \\ & & \frac{1}{\lambda_D} \end{bmatrix}$. In addition, U is orthogonal ($UU^T = \frac{1}{\lambda_1} + \frac{1}{\lambda_2} = \frac{1}{\lambda_2} = \frac{1}{\lambda_2} = \frac{1}{\lambda_2} + \frac{1}{\lambda_2} = \frac{1}{\lambda_2$

I). Thus, we have:

$$\Sigma^{-1} = \sum_{i=1}^{D} \frac{1}{\lambda_i} \mathbf{u}_i \mathbf{u}_i^{\mathrm{T}}$$
 (28)

This equation also implicitly proves that

Conclusion 2.1.6. Inverse of symmetric matrix is also symmetric.

Substitute (28) into (26), the quadratic form becomes:

$$\Delta^2 = \sum_{i=1}^D \frac{y_i^2}{\lambda_i} \tag{29}$$

where $y_i = \mathbf{u}_i^{\mathrm{T}}(\mathbf{x} - \boldsymbol{\mu})$. We can interpret $\{y_i\}$ as a new coordinate system by projecting $(\mathbf{x} - \boldsymbol{\mu})$ to $\{\mathbf{u}_i\}$. The functional dependence of $\mathbf{y} = \begin{bmatrix} y_1 & \cdot & y_D \end{bmatrix}$ on \mathbf{x} is captured by the following equation:

$$\mathbf{y} = U^{\mathrm{T}}(\mathbf{x} - \boldsymbol{\mu}) \tag{30}$$

To be properly normalized, Σ should be <i>positive definite</i> . Otherwise, see "Chap12":														
					• •									
					• •									
					• •									

Integral by substitution (\mathbf{x} replaced by \mathbf{y}) requires calculating *Jacobi determinant* where $J_{ij} = \frac{\partial x_i}{\partial y_j} = .$ From (30) we know that $\mathbf{x} = U\mathbf{y} + \boldsymbol{\mu}$ and thus $J_{ij} = \frac{\partial x_i}{\partial y_j} = U_{ji}^T$. Obviously, this implies that J = U. Then we have:

$$|J|^{2} = |U|^{2}$$

$$= |U^{T}||U| \quad (\because |A| = |A^{T}|)$$

$$= |U^{T}U| \quad (\because |A||B| = |AB|)$$

$$= |I| = 1$$
(31)

Use these properties of determinant in above derivation, we also have:

$$|\Sigma| = |U\Lambda U^{T}|$$

$$= |U||\Lambda||U^{T}|$$

$$= |U^{T}U||\Lambda|$$

$$= \prod_{i=1}^{D} \lambda_{i}$$
(32)

Up to now, we know the substitution results:

$$\Pr(\mathbf{y}) = \Pr(\mathbf{x})|J|$$

$$= \frac{1}{(2\pi)^{\frac{D}{2}}} \frac{1}{\prod_{i=1}^{D} \lambda_{i}^{\frac{1}{2}}} \exp\left\{-\frac{1}{2} \sum_{i=1}^{D} \frac{y_{i}^{2}}{\lambda_{i}}\right\}$$

$$= \prod_{i=1}^{D} \frac{1}{(2\pi\lambda_{i})^{\frac{1}{2}}} \exp\left\{-\frac{y_{i}^{2}}{2\lambda_{i}}\right\}$$
(33)

From the perspective of y coordinate system, it is straightforward to show that:

$$\int \Pr(\mathbf{y}) d\mathbf{y} = \prod_{i=1}^{D} \int_{-\infty}^{\infty} \frac{1}{(2\pi\lambda_i)^{\frac{1}{2}}} \exp\left\{-\frac{y_i^2}{2\lambda_i}\right\} dy_i = 1$$
 (34)

We now check the moment and second moment of Gaussian. By substitution $\mathbf{z} = \mathbf{x} - \boldsymbol{\mu}$ we have $\mathbb{E}[\mathbf{x}] = \frac{1}{(2\pi)^{\frac{D}{2}}} \frac{1}{|\Sigma|^{\frac{1}{2}}} \int \exp{\{-\frac{1}{2}\mathbf{z}^T\Sigma^{-1}\mathbf{z}\}}(\mathbf{z} + \boldsymbol{\mu})\mathrm{d}\mathbf{z}$. Note that $\exp(\cdot)$ is an even function and the integral region is the whole \mathbb{R}^D . Thus the term of exponential function multiplied by z will vanish and $\mathbb{E}[\mathbf{x}] = \boldsymbol{\mu}$. Univariate case we consider $\mathbb{E}[x^2]$, while in the multivariate case (D dimensional) we have D^2 pairs $x_i x_j$. Thus we consider:

$$\mathbb{E}[\mathbf{x}\mathbf{x}^{\mathrm{T}}] = \frac{1}{(2\pi)^{\frac{D}{2}}} \frac{1}{|\Sigma|^{\frac{1}{2}}} \int \exp\left\{-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^{\mathrm{T}} \Sigma^{-1}(\mathbf{x} - \boldsymbol{\mu})\right\} \mathbf{x}\mathbf{x}^{\mathrm{T}} d\mathbf{x}$$

$$= \frac{1}{(2\pi)^{\frac{D}{2}}} \frac{1}{|\Sigma|^{\frac{1}{2}}} \int \exp\left\{-\frac{1}{2}\mathbf{z}^{\mathrm{T}} \Sigma^{-1} \mathbf{z}\right\} (\mathbf{z} + \boldsymbol{\mu}) (\mathbf{z} + \boldsymbol{\mu})^{\mathrm{T}} d\mathbf{z}$$
(35)

Consider $(\mathbf{z} + \boldsymbol{\mu})(\mathbf{z} + \boldsymbol{\mu})^T$ we firstly find that $\mathbf{z}\boldsymbol{\mu}^T$ and $\boldsymbol{\mu}\mathbf{z}^T$ will vanish because of the same reason leveraged in computing moment of Gaussian above. Secondly, $\boldsymbol{\mu}\boldsymbol{\mu}^T$ doesn't contain integral variable \mathbf{z} and thus contributes $\boldsymbol{\mu}\boldsymbol{\mu}^T$ to final integral result since $\Pr(\cdot)$ is normalized. Remember that from (30) we have

z = Uy and substitution results:

$$\frac{1}{(2\pi)^{\frac{D}{2}}} \frac{1}{|\Sigma|^{\frac{1}{2}}} \int \exp\left\{-\frac{1}{2}\mathbf{z}^{\mathsf{T}}\Sigma^{-1}\mathbf{z}\right\} \mathbf{z}\mathbf{z}^{\mathsf{T}} d\mathbf{z}$$

$$= \frac{1}{(2\pi)^{\frac{D}{2}}} \frac{1}{|\Sigma|^{\frac{1}{2}}} \int \exp\left\{-\sum_{i=1}^{D} \frac{y_{i}^{2}}{2\lambda_{i}}\right\} (U\mathbf{y})(U\mathbf{y})^{\mathsf{T}} d\mathbf{y}$$

$$= \frac{1}{(2\pi)^{\frac{D}{2}}} \frac{1}{|\Sigma|^{\frac{1}{2}}} \int \exp\left\{-\sum_{i=1}^{D} \frac{y_{i}^{2}}{2\lambda_{i}}\right\} (\sum_{j=1}^{D} y_{j}\mathbf{u}_{j}) (\sum_{k=1}^{D} y_{k}\mathbf{u}_{k}^{\mathsf{T}}) d\mathbf{y}$$

$$= \frac{1}{(2\pi)^{\frac{D}{2}}} \frac{1}{|\Sigma|^{\frac{1}{2}}} \sum_{j=1}^{D} \sum_{k=1}^{D} \mathbf{u}_{j}\mathbf{u}_{k}^{\mathsf{T}} \int \exp\left\{-\sum_{i=1}^{D} \frac{y_{i}^{2}}{2\lambda_{i}}\right\} y_{j} y_{k} d\mathbf{y}$$
(36)

Since $y_j y_k$ is odd function of y_j , y_k respectively, the term with different j, k will vanish. Besides, $\int \mathcal{N}(x|\mu,\sigma)x^2 dx = \sigma$. Thus, we can further simplify (36) to:

$$\frac{1}{(2\pi)^{\frac{D}{2}}} \frac{1}{|\Sigma|^{\frac{1}{2}}} \sum_{j=1}^{D} \mathbf{u}_{j} \mathbf{u}_{j}^{\mathsf{T}} \int \exp\left\{-\sum_{i=1}^{D} \frac{y_{i}^{2}}{2\lambda_{i}}\right\} y_{j}^{2} d\mathbf{y} = \sum_{j=1}^{D} \mathbf{u}_{j} \mathbf{u}_{j}^{\mathsf{T}} \lambda_{j} = \Sigma$$
(37)

Based on above analysis and calculation we have:

$$cov[\mathbf{x}] = \mathbb{E}[(\mathbf{x} - \mathbb{E}[\mathbf{x}])(\mathbf{x} - \mathbb{E}[\mathbf{x}])^{\mathrm{T}}]$$

$$= \mathbb{E}[\mathbf{x}\mathbf{x}^{\mathrm{T}}] - \mathbb{E}[\mathbf{x}]\mathbb{E}[\mathbf{x}]^{\mathrm{T}}$$

$$= \Sigma + \mu\mu^{\mathrm{T}} - \mu\mu^{\mathrm{T}}$$

$$= \Sigma.$$
(38)

Because Σ governs the covariance, it is called *covariance matrix* and Σ^{-1} is called the *precision matrix*.

2.1.1 Conditional Gaussian distribution

(preliminaries)

Conclusion 2.1.7.

$$\begin{bmatrix} A & B \\ C & D \end{bmatrix}^{-1} = \begin{bmatrix} M & -MBD^{-1} \\ -D^{-1}CM & D^{-1} + D^{-1}CMBD^{-1} \end{bmatrix}$$
(39)

where A, B, D, D are respectively $p \times p$, $p \times q$, $q \times p$, $q \times q$ matrices and $M = (A - BD^{-1}C)^{-1}$. Note that D is required to be *invertible*. M is known as the *Schur complement* of the matrix on the left-hand side of (39) w.r.t. the submatrix D.

Proof. Multiply
$$\begin{bmatrix} M & -MBD^{-1} \\ -D^{-1}CM & D^{-1} + D^{-1}CMBD^{-1} \end{bmatrix}$$
 by $\begin{bmatrix} A & B \\ C & D \end{bmatrix}$ and check whether it results $I_{(p+q)\times(p+q)} = \begin{bmatrix} I_{p\times p} & 0 \\ 0 & I_{q\times q} \end{bmatrix}$.

Firstly, $MA - MBD^{-1}C = M(A - BD^{-1}C) = I_{p \times p}$.

Secondly, $MB - MBD^{-1}D = MB - MB = 0$.

Thirdly,
$$-D^{-1}CMA + D^{-1}C + D^{-1}CMBD^{-1}C = D^{-1}C(-MA + I + MBD^{-1}C) = D^{-1}CM(-A + M^{-1} + BD^{-1}C) = 0.$$

Finally,
$$-D^{-1}CMB + D^{-1}D + D^{-1}CMBD^{-1}D = I_{q \times q}$$
.

Suppose \mathbf{x} is a D-dimensional vector with Gaussian distribution $\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \Sigma)$ and we partition \mathbf{x} into two disjoint subsets \mathbf{x}_a and \mathbf{x}_b . We can always put \mathbf{x}_a on the first M components of \mathbf{x} by *permutation* over $\boldsymbol{\mu}$ and Σ without loss of generality.

•
$$\mathbf{x} = \begin{bmatrix} \mathbf{x}_a \\ \mathbf{x}_b \end{bmatrix}$$

$$\bullet \ \mu = \begin{bmatrix} \mu_a \\ \mu_b \end{bmatrix}$$

•
$$\Sigma = \begin{bmatrix} \Sigma_{aa} & \Sigma_{ab} \\ \Sigma_{ba} & \Sigma_{bb} \end{bmatrix}$$
 and the following properties result from $\Sigma^{T} = \Sigma$

- $\Sigma_{aa} = \Sigma_{aa}^{T}$ and $\Sigma_{bb} = \Sigma_{bb}^{T}$

- $\Sigma_{ab}^{T} = \Sigma_{ba}$

$$\bullet \ \Lambda = \Sigma^{-1} = \begin{bmatrix} \Lambda_{aa} & \Lambda_{ab} \\ \Lambda_{ba} & \Lambda_{bb} \end{bmatrix}$$

Using this partitioning we firstly obtain

$$\begin{split} &-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu})^{\mathsf{T}}\boldsymbol{\Sigma}^{-1}(\mathbf{x}-\boldsymbol{\mu}) \\ &= -\frac{1}{2}\left[(\mathbf{x}_{a}-\boldsymbol{\mu}_{a})^{\mathsf{T}} \quad (\mathbf{x}_{b}-\boldsymbol{\mu}_{b})^{\mathsf{T}}\right] \begin{bmatrix} \boldsymbol{\Lambda}_{aa} & \boldsymbol{\Lambda}_{ab} \\ \boldsymbol{\Lambda}_{ba} & \boldsymbol{\Lambda}_{bb} \end{bmatrix} \begin{bmatrix} \mathbf{x}_{a}-\boldsymbol{\mu}_{a} \\ \mathbf{x}_{b}-\boldsymbol{\mu}_{b} \end{bmatrix} \\ &= -\frac{1}{2}\left[(\mathbf{x}_{a}-\boldsymbol{\mu}_{a})^{\mathsf{T}}\boldsymbol{\Lambda}_{aa} + (\mathbf{x}_{b}-\boldsymbol{\mu}_{b})^{\mathsf{T}}\boldsymbol{\Lambda}_{ba} \quad (\mathbf{x}_{a}-\boldsymbol{\mu}_{a})^{\mathsf{T}}\boldsymbol{\Lambda}_{ab} + (\mathbf{x}_{b}-\boldsymbol{\mu}_{b})^{\mathsf{T}}\boldsymbol{\Lambda}_{bb}\right] \begin{bmatrix} \mathbf{x}_{a}-\boldsymbol{\mu}_{a} \\ \mathbf{x}_{b}-\boldsymbol{\mu}_{b} \end{bmatrix} \\ &= -\frac{1}{2}\left\{(\mathbf{x}_{a}-\boldsymbol{\mu}_{a})^{\mathsf{T}}\boldsymbol{\Lambda}_{aa}(\mathbf{x}_{a}-\boldsymbol{\mu}_{a}) + (\mathbf{x}_{b}-\boldsymbol{\mu}_{b})^{\mathsf{T}}\boldsymbol{\Lambda}_{ba}(\mathbf{x}_{a}-\boldsymbol{\mu}_{a}) + (\mathbf{x}_{a}-\boldsymbol{\mu}_{a})^{\mathsf{T}}\boldsymbol{\Lambda}_{ab}(\mathbf{x}_{b}-\boldsymbol{\mu}_{b}) + (\mathbf{x}_{b}-\boldsymbol{\mu}_{b})^{\mathsf{T}}\boldsymbol{\Lambda}_{bb}(\mathbf{x}_{b}-\boldsymbol{\mu}_{b})\right\} \end{split}$$

We find that as a function of x_a , this is again a quadratic form.

Conclusion 2.1.8. If two sets of random variables are jointly Gaussian, then the conditional probability distribution of one set conditioned on the other is agian Gaussian.

To determine the mean and covariance of $\Pr(\mathbf{x}_a|\mathbf{x}_b)$, we should firstly note that the exponent in a general $\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}\boldsymbol{\Sigma})$ can be written as:

$$-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^{\mathrm{T}} \Sigma (\mathbf{x} - \boldsymbol{\mu}) = -\frac{1}{2} \mathbf{x}^{\mathrm{T}} \Sigma^{-1} \mathbf{x} + \mathbf{x}^{\mathrm{T}} \Sigma^{-1} \boldsymbol{\mu} - \frac{1}{2} \boldsymbol{\mu}^{\mathrm{T}} \Sigma^{-1} \boldsymbol{\mu}$$
(41)

We pick out all terms that are second order in \mathbf{x}_a from (40): $-\frac{1}{2}\mathbf{x}_a^{\mathrm{T}}\Lambda_{aa}\mathbf{x}_a$. Thus the covariance of $\Pr(\mathbf{x}_a|\mathbf{x}_b)$ is Λ_{aa}^{-1} .

We pick out all terms that are linear in \mathbf{x}_a from (40):

$$\frac{1}{2}\mathbf{x}_{a}^{\mathsf{T}}\Lambda_{aa}\boldsymbol{\mu}_{a} + \frac{1}{2}\boldsymbol{\mu}_{a}^{\mathsf{T}}\Lambda_{aa}\mathbf{x}_{a} - \frac{1}{2}\mathbf{x}_{a}^{\mathsf{T}}\Lambda_{ab}(\mathbf{x}_{b} - \boldsymbol{\mu}_{b}) - \frac{1}{2}(\mathbf{x}_{b} - \boldsymbol{\mu}_{b})^{\mathsf{T}}\Lambda_{ba}\mathbf{x}_{a}
= \mathbf{x}_{a}^{\mathsf{T}}\Lambda_{aa}\boldsymbol{\mu}_{a} - \mathbf{x}_{a}^{\mathsf{T}}\Lambda_{ab}(\mathbf{x}_{b} - \boldsymbol{\mu}_{b}) \quad (\because \mathbf{x}^{\mathsf{T}}A\mathbf{y} = \mathbf{y}^{\mathsf{T}}A^{\mathsf{T}}\mathbf{x} \text{ and } \Lambda_{ba}^{\mathsf{T}} = \Lambda_{ab})
= \mathbf{x}_{a}^{\mathsf{T}}\{\Lambda_{aa}\boldsymbol{\mu}_{a} - \Lambda_{ab}(\mathbf{x}_{b} - \boldsymbol{\mu}_{b})\}$$
(42)

Thus the mean of $\Pr(\mathbf{x}_a|\mathbf{x}_b)$ is $\Lambda_{aa}^{-1}\{\Lambda_{aa}\boldsymbol{\mu}_a - \Lambda_{ab}(\mathbf{x}_b - \boldsymbol{\mu}_b)\} = \boldsymbol{\mu}_a - \Lambda_{aa}^{-1}\Lambda_{ab}(\mathbf{x}_b - \boldsymbol{\mu}_b)$. If we can ensure Σ_{bb} to be invertible (how?),

٠.			٠.		•			•			٠.				•		•						٠.		•	٠.	•	•		•			•					•				•						•		•	٠.			٠.					•		•			
																									•			•																																	•			. .
• •	• •	• •	٠.	• •	•	• •	• •	•	• •	٠.	• •	•	• •	• •	•	• •	•	• •	• •	•	•	•	• •	• •	•	٠.	•	•	• •	•	٠.	• •	•	٠.	•	•	•	•	• •	•	٠.	•	• •	•	• •	•	• •	•	٠.	•	٠.	•	• •	• •	•	• •	• •	• •	•	• •	•	•	٠.	•

......[5] we can make use of (39) to express Λ_{aa} as $(\Sigma_{aa} - \Sigma_{ab}\Sigma_{bb}^{-1}\Sigma_{ba})^{-1}$ and hence:

$$\Sigma_{a|b} = \Lambda_{aa}^{-1} = (\Sigma_{aa} - \Sigma_{ab} \Sigma_{bb}^{-1} \Sigma_{ba}) \tag{43}$$

we can also make use of (39) to express Λ_{ab} as $-\Lambda_{aa}\Sigma_{ab}\Sigma_{bb}^{-1}$ and hence:

$$\mu_{a|b} = \mu_{a} - \Lambda_{aa}^{-1} \Lambda_{ab} (\mathbf{x}_{b} - \mu_{b})$$

$$= \mu_{a} + \Lambda_{aa}^{-1} \Lambda_{aa} \Sigma_{ab} \Sigma_{bb}^{-1} (\mathbf{x}_{b} - \mu_{b})$$

$$= \mu_{a} + \Sigma_{ab} \Sigma_{bb}^{-1} (\mathbf{x}_{b} - \mu_{b})$$
(44)

2.1.2 Marginal Gaussian distributions

(Preliminary)

Conclusion 2.1.9. Suppose A, B, C, D are matrices of dimension $p \times p$, $p \times q$, $q \times p$, $q \times q$ respectively. Besides, D is *invertible*. Then we have:

$$\det\begin{pmatrix} \begin{bmatrix} A & B \\ C & D \end{bmatrix} = \det(D)\det(A - BD^{-1}C)$$
(45)

We calculate $Pr(\mathbf{x}_a)$ by integrated out \mathbf{x}_b :

$$Pr(\mathbf{x}_a) = \int Pr(\mathbf{x}_a, \mathbf{x}_b) d\mathbf{x}_b$$
 (46)

Because we want to integrate out x_b , we should firstly consider terms that are functional to x_b in (40). Since b is symmetric to a in (40), we can directly use (41), (43) and (44) to derive:

$$-\frac{1}{2}\mathbf{x}_{b}^{\mathsf{T}}\Lambda_{bb}\mathbf{x}_{b} + \mathbf{x}_{b}^{\mathsf{T}}\{\Lambda_{bb}\boldsymbol{\mu}_{b} - \Lambda_{ba}(\mathbf{x}_{a} - \boldsymbol{\mu}_{a})\}$$

$$= -\frac{1}{2}\mathbf{x}_{b}^{\mathsf{T}}\Lambda_{bb}\mathbf{x}_{b} + \mathbf{x}_{b}^{\mathsf{T}}\mathbf{m} \quad \text{(where } \mathbf{m} = \Lambda_{bb}\boldsymbol{\mu}_{b} - \Lambda_{ba}(\mathbf{x}_{a} - \boldsymbol{\mu}_{a}))$$

$$= -\frac{1}{2}\mathbf{x}_{b}^{\mathsf{T}}\Lambda_{bb}\mathbf{x}_{b} + \mathbf{x}_{b}^{\mathsf{T}}\Lambda_{bb}(\Lambda_{bb}^{-1}\mathbf{m}) - \frac{1}{2}(\Lambda_{bb}^{-1}\mathbf{m})^{\mathsf{T}}\Lambda_{bb}(\Lambda_{bb}^{-1}\mathbf{m}) + \frac{1}{2}(\Lambda_{bb}^{-1}\mathbf{m})^{\mathsf{T}}\Lambda_{bb}(\Lambda_{bb}^{-1}\mathbf{m})$$

$$= -\frac{1}{2}(\mathbf{x}_{b} - \Lambda_{bb}^{-1}\mathbf{m})^{\mathsf{T}}\Lambda_{bb}(\mathbf{x}_{b} - \Lambda_{bb}^{-1}\mathbf{m}) + \frac{1}{2}\mathbf{m}^{\mathsf{T}}\Lambda_{bb}^{-1}\mathbf{m}$$

$$(47)$$

Note that **m** doesn't contain x_b and the first term in the right-hand side of (47) is the standard quadratic form. Assuming x_a is M dimensional, integration results:

$$\int \exp\{-\frac{1}{2}(\mathbf{x}_b - \Lambda_{bb}^{-1}\mathbf{m})^{\mathrm{T}}\Lambda_{bb}(\mathbf{x}_b - \Lambda_{bb}^{-1}\mathbf{m})\}d\mathbf{x}_b = (2\pi)^{\frac{D-M}{2}}|\Lambda_{bb}^{-1}|^{\frac{1}{2}}$$
(48)

while $\frac{1}{2}\mathbf{m}^{\mathrm{T}}\Lambda_{hh}^{-1}\mathbf{m}$ remains as exponent. We firstly expand it as:

$$\frac{1}{2}\mathbf{m}^{T}\Lambda_{bb}^{-1}\mathbf{m} = \frac{1}{2}(\Lambda_{bb}\mu_{b} - \Lambda_{ba}(\mathbf{x}_{a} - \mu_{a}))^{T}\Lambda_{bb}^{-1}(\Lambda_{bb}\mu_{b} - \Lambda_{ba}(\mathbf{x}_{a} - \mu_{a}))$$

$$= \frac{1}{2}\{\mu_{b}^{T}\Lambda_{bb}\mu_{b} - \mu_{b}^{T}\Lambda_{ba}(\mathbf{x}_{a} - \mu_{a}) - (\mathbf{x}_{a} - \mu_{a})^{T}\Lambda_{ab}\mu_{b} + (\mathbf{x}_{a} - \mu_{a})^{T}\Lambda_{ab}\Lambda_{bb}^{-1}\Lambda_{ba}(\mathbf{x}_{a} - \mu_{a})\}$$

$$= \frac{1}{2}(\mathbf{x}_{a} - \mu_{a})^{T}\Lambda_{ab}\Lambda_{bb}^{-1}\Lambda_{ba}(\mathbf{x}_{a} - \mu_{a}) - (\mathbf{x}_{a} - \mu_{a})^{T}\Lambda_{ab}\mu_{b} + \frac{1}{2}\mu_{b}^{T}\Lambda_{bb}\mu_{b}$$
(49)

Then we consider terms that are not functional to x_b in (40):

$$-\frac{1}{2}(\mathbf{x}_a - \boldsymbol{\mu}_a)^{\mathrm{T}} \Lambda_{aa}(\mathbf{x}_a - \boldsymbol{\mu}_a) + (\mathbf{x}_a - \boldsymbol{\mu}_a)^{\mathrm{T}} \Lambda_{ab} \boldsymbol{\mu}_b - \frac{1}{2} \boldsymbol{\mu}_b^{\mathrm{T}} \Lambda_{bb} \boldsymbol{\mu}_b$$
 (50)

Combining (49) and (50) we get:

$$\exp\left\{-\frac{1}{2}(\mathbf{x}_a - \boldsymbol{\mu}_a)^{\mathrm{T}}(\Lambda_{aa} - \Lambda_{ab}\Lambda_{bb}^{-1}\Lambda_{ba})(\mathbf{x}_a - \boldsymbol{\mu}_a)\right\}$$
(51)

Then we can prove it is normalized by transforming original normalization term through (45) and substituting (48) into it:

$$\frac{1}{(2\pi)^{\frac{D}{2}}} \frac{1}{|\Sigma|^{\frac{1}{2}}} (2\pi)^{\frac{D-M}{2}} |\Lambda_{bb}^{-1}|^{\frac{1}{2}} \\
= \frac{1}{(2\pi)^{\frac{M}{2}}} |\Lambda|^{\frac{1}{2}} \frac{1}{|\Lambda_{bb}|^{\frac{1}{2}}} \quad (\because AA^{-1} = I, |A||A^{-1}| = 1) \\
= \frac{1}{(2\pi)^{\frac{M}{2}}} \frac{(|\Lambda_{bb}||\Lambda_{aa} - \Lambda_{ab}\Lambda_{bb}^{-1}\Lambda_{ba}|)^{\frac{1}{2}}}{|\Lambda_{bb}|^{\frac{1}{2}}} \\
= \frac{1}{(2\pi)^{\frac{M}{2}}} \frac{1}{|(\Lambda_{aa} - \Lambda_{ab}\Lambda_{bb}^{-1}\Lambda_{ba})^{-1}|^{\frac{1}{2}}} \tag{52}$$

Conclusion 2.1.10. If two sets of random variables are jointly Gaussian, then the marginal probability distribution of one set integrated out the other is again Gaussian.

$$Pr(\mathbf{x}_a) = \int Pr(\mathbf{x}_a, \mathbf{x}_b) d\mathbf{x}_b = \mathcal{N}(\mathbf{x}_a | \boldsymbol{\mu}_a, \boldsymbol{\Sigma}_{aa})$$
 (53)

2.1.3 Bayes' theorem for Gaussian variables

Up to now, the whole story is given a quadratic form $(\mathbf{x} - \boldsymbol{\mu})^{\mathrm{T}} \Lambda(\mathbf{x} - \boldsymbol{\mu})$, we could partition $\mathbf{x}, \boldsymbol{\mu}, \Lambda$ with consistent size and transform it into two quadratic forms:

$$-\frac{1}{2}\mathbf{x}^{T}\Lambda\mathbf{x} = -\frac{1}{2}(\mathbf{x}_{b} - \Lambda_{bb}^{-1}\mathbf{m})^{T}\Lambda_{bb}(\mathbf{x}_{b} - \Lambda_{bb}^{-1}\mathbf{m}) - \frac{1}{2}(\mathbf{x}_{a} - \mu_{a})^{T}(\Lambda_{aa} - \Lambda_{ab}\Lambda_{bb}^{-1}\Lambda_{ba})(\mathbf{x}_{a} - \mu_{a})$$
(54)

where $\mathbf{m} = \Lambda_{bb}\mu_b - \Lambda_{ba}(\mathbf{x}_a - \mu_a)$ and thus " μ " for the first term on the right-hand side of (54) is $\mu_b - \Lambda_{bb}^{-1}\Lambda_{ba}(\mathbf{x}_a - \mu_a)$ which is a linear function to \mathbf{x}_a .

Hence, the inverse operation can be applied to $\Pr(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \Lambda^{-1})$, $\Pr(\mathbf{y}|\mathbf{x}) = \mathcal{N}(\mathbf{y}|A\mathbf{x} + \mathbf{b}, L^{-1})$ resulting their joint Gaussian distribution $\Pr(\mathbf{z}) = \Pr(\mathbf{x}, \mathbf{y}) = \Pr(\mathbf{y}|\mathbf{x}) \Pr(\mathbf{x})$.

Suppose the quadratic form is $(\mathbf{z} - \mathbb{E}[\mathbf{z}])^T R(\mathbf{z} - \mathbb{E}[\mathbf{z}])$ and $R = \begin{bmatrix} R_{xx} & R_{xy} \\ R_{yx} & R_{yy} \end{bmatrix}$, $\mathbb{E}[\mathbf{z}] = \begin{bmatrix} \mathbb{E}[\mathbf{z}]_x \\ \mathbb{E}[\mathbf{z}]_y \end{bmatrix}$, then we calculate them according to (54):

- Obviously, $R_{yy} = L$, $\mathbb{E}[\mathbf{z}]_x = \mu$.
- :: A**x** + **b** = $\mathbb{E}[\mathbf{z}]_y R_{yy}^{-1} R_{yx} (\mathbf{x} \mathbb{E}[\mathbf{z}]_x)$:: $A = -R_{yy}^{-1} R_{yx} = -L^{-1} R_{yx}$, i.e., $R_{yx} = -LA$.
- Now $A\mathbf{x} + \mathbf{b} = \mathbb{E}[\mathbf{z}]_{\nu} + A(\mathbf{x} \mu)$, thus, $\mathbb{E}[\mathbf{z}]_{\nu} = A\mu + \mathbf{b}$.
- Since *R* is required to be symmetric, $R_{xy} = R_{yx}^{T} = -A^{T}L$.
- $\therefore R_{xx} R_{xy}R_{yy}^{-1}R_{yx} = R_{xx} A^{T}LA = \Lambda \therefore R_{xx} = \Lambda + A^{T}LA.$

Use (39), we know that:

$$\Pr(\mathbf{z}) = \mathcal{N}(\mathbf{z} \begin{vmatrix} \boldsymbol{\mu} \\ A\boldsymbol{\mu} + \mathbf{b} \end{vmatrix}, \begin{bmatrix} \Lambda + A^{\mathrm{T}}LA & -A^{\mathrm{T}}L \\ -LA & L \end{bmatrix}^{-1})$$

$$= \mathcal{N}(\mathbf{z} \begin{vmatrix} \boldsymbol{\mu} \\ A\boldsymbol{\mu} + \mathbf{b} \end{vmatrix}, \begin{bmatrix} \Lambda^{-1} & \Lambda^{-1}A^{\mathrm{T}} \\ A\Lambda^{-1} & L^{-1} + A\Lambda^{-1}A^{\mathrm{T}} \end{bmatrix})$$
(55)

2.1.4 Maximum likelihood for the Gaussian

(Preliminary)

Definition 2.1.2. The trace of a square matrix is the sum of the elements on its diagonal, i.e., $\text{Tr}(A_{n \times n}) = \sum_{i=1}^{n} A_{ii}$.

Conclusion 2.1.11.

$$Tr(AB) = Tr(BA) \tag{56}$$

By which we can generalize to *cyclic* property of the trace operator: Tr(ABC) = Tr(CAB) = Tr(BCA).

Proof. Suppose A, B is of size $n \times m$ and $m \times n$.

 $\forall i \in 1, \ldots, n : (AB)_{ii} = \sum_{k=1}^m A_{ik} B_{ki}$ Thus $\text{Tr}(AB) = \sum_{i=1}^n \sum_{k=1}^m A_{ik} B_{ki}$. The other, $\text{Tr}(BA) = \sum_{i=1}^m \sum_{k=1}^n B_{ik} A_{ki}$. Obviously, in these two cases, each element of A is matched to a corresponding element of B in the same way.

Conclusion 2.1.12.

$$Tr(\mathbf{x}^{\mathrm{T}}A\mathbf{x}) = Tr(A\mathbf{x}\mathbf{x}^{\mathrm{T}}) \tag{57}$$

Proof. Obviously, the left-hand side equals $\sum_{i=1}^{M} \sum_{j=1}^{M} \mathbf{x}_i \mathbf{x}_j A_{ij}$

$$\mathbf{x}\mathbf{x}^{\mathrm{T}} = \begin{bmatrix} \mathbf{x}_{1}\mathbf{x}_{1} & \cdots & \mathbf{x}_{1}\mathbf{x}_{M} \\ \vdots & \ddots & \vdots \\ \mathbf{x}_{M}\mathbf{x}_{1} & \cdots & \mathbf{x}_{M}\mathbf{x}_{M} \end{bmatrix}$$
 and suppose $A\mathbf{x}\mathbf{x}^{\mathrm{T}} = R$, then $R_{ii} = \sum_{j=1}^{M} A_{ij}\mathbf{x}_{j}\mathbf{x}_{i}$. Trace is the sum of

elements in diagonal, thus it equals $\sum_{i=1}^{M} \sum_{j=1}^{M} A_{ij} \mathbf{x}_{i} \mathbf{x}_{i}$

Conclusion 2.1.13.

$$\frac{\partial \text{Tr}(AB)}{\partial x} = \text{Tr}(\frac{\partial A}{\partial x}B) \tag{58}$$

when B is fixed or is not functional to x.

Proof. Just write out the subscripts.

Conclusion 2.1.14.

$$\frac{\partial}{\partial r}(A^{-1}) = -A^{-1}\frac{\partial A}{\partial r}A^{-1} \tag{59}$$

Proof.
$$: 0 = \frac{\partial I}{\partial x} = \frac{\partial AA^{-1}}{\partial x} = \frac{\partial A}{\partial x}A^{-1} + A\frac{\partial A^{-1}}{\partial x}$$

Conclusion 2.1.15. For a symmetric, strictly positive definite matrix $A_{M\times M}$, we have:

$$\frac{\partial}{\partial x} \ln|A| = \text{Tr}(A^{-1} \frac{\partial A}{\partial x}) \tag{60}$$

Proof. : $A = A^{T}$: by (32) we have $|A| = \prod_{i=1}^{M} \lambda_{i}$

Then the left-hand side can be expressed as:

$$\sum_{i=1}^{M} \frac{1}{\lambda_i} \frac{\partial \lambda_i}{\partial x} \tag{61}$$

Use (27) and (28), we express $A^{-1} \frac{\partial A}{\partial x}$ as:

$$A^{-1}\frac{\partial A}{\partial x} = \sum_{i=1}^{M} \frac{1}{\lambda_{i}} \mathbf{u}_{i} \mathbf{u}_{i}^{T} \sum_{j=1}^{M} \frac{\partial \lambda_{j} \mathbf{u}_{j} \mathbf{u}_{j}^{T}}{\partial x}$$

$$= \sum_{i=1}^{M} \sum_{j=1}^{M} \frac{1}{\lambda_{i}} \mathbf{u}_{i} \mathbf{u}_{i}^{T} \frac{\partial \lambda_{j} \mathbf{u}_{j} \mathbf{u}_{j}^{T}}{\partial x}$$

$$= \sum_{i=1}^{M} \sum_{j=1}^{M} \frac{1}{\lambda_{i}} \mathbf{u}_{i} \mathbf{u}_{i}^{T} (\frac{\partial \lambda_{j}}{\partial x} \mathbf{u}_{j} \mathbf{u}_{j}^{T} + \lambda_{j} \frac{\partial \mathbf{u}_{j} \mathbf{u}_{j}^{T}}{\partial x})$$

$$= \sum_{i=1}^{M} \frac{1}{\lambda_{i}} \frac{\partial \lambda_{i}}{\partial x} \mathbf{u}_{i} \mathbf{u}_{i}^{T} + \sum_{i=1}^{M} \sum_{j=1}^{M} \frac{\lambda_{j}}{\lambda_{i}} \mathbf{u}_{i} \mathbf{u}_{i}^{T} \frac{\partial \mathbf{u}_{j} \mathbf{u}_{j}^{T}}{\partial x} \quad (\because \mathbf{u}_{i}^{T} \mathbf{u}_{j} = I_{ij})$$

$$(62)$$

 $\because \forall i \in 1, \dots, M : \text{Tr}(\mathbf{u}_i \mathbf{u}_i^{\text{T}}) = \|\mathbf{u}_i\|^2 = 1 \therefore \text{Tr}(\sum_{i=1}^M \frac{1}{\lambda_i} \frac{\partial \lambda_i}{\partial x} \mathbf{u}_i \mathbf{u}_i^{\text{T}}) = \sum_{i=1}^M \frac{1}{\lambda_i} \frac{\partial \lambda_i}{\partial x} \text{Tr}(\mathbf{u}_i \mathbf{u}_i^{\text{T}}) = (61). \text{ Hence, next step is to prove the remaining term in the right-hand side of (62) vanishes.}$

$$\operatorname{Tr}\left\{\sum_{i=1}^{M} \sum_{j=1}^{M} \frac{\lambda_{j}}{\lambda_{i}} \mathbf{u}_{i} \mathbf{u}_{i}^{\mathrm{T}} \left(\frac{\partial \mathbf{u}_{j}}{\partial x} \mathbf{u}_{j}^{\mathrm{T}} + \mathbf{u}_{j} \frac{\partial \mathbf{u}_{j}^{\mathrm{T}}}{\partial x}\right)\right\} \quad (\because (AB)' = A'B + AB')$$

$$= \operatorname{Tr}\left\{\sum_{i=1}^{M} \sum_{j=1}^{M} \frac{\lambda_{j}}{\lambda_{i}} \mathbf{u}_{i} \mathbf{u}_{i}^{\mathrm{T}} \frac{\partial \mathbf{u}_{j}}{\partial x} \mathbf{u}_{j}^{\mathrm{T}}\right\} + \operatorname{Tr}\left(\sum_{i=1}^{M} \mathbf{u}_{i} \frac{\partial \mathbf{u}_{i}^{\mathrm{T}}}{\partial x}\right) \quad (\because \mathbf{u}_{i}^{\mathrm{T}} \mathbf{u}_{j} = I_{ij})$$

$$= \sum_{i=1}^{M} \sum_{j=1}^{M} \frac{\lambda_{j}}{\lambda_{i}} \operatorname{Tr}\left(\mathbf{u}_{i} \mathbf{u}_{i}^{\mathrm{T}} \frac{\partial \mathbf{u}_{j}}{\partial x} \mathbf{u}_{j}^{\mathrm{T}}\right) + \operatorname{Tr}\left(\sum_{i=1}^{M} \mathbf{u}_{i} \frac{\partial \mathbf{u}_{i}^{\mathrm{T}}}{\partial x}\right)$$

$$= \sum_{i=1}^{M} \sum_{j=1}^{M} \frac{\lambda_{j}}{\lambda_{i}} \operatorname{Tr}\left(\mathbf{u}_{j}^{\mathrm{T}} \mathbf{u}_{i} \mathbf{u}_{i}^{\mathrm{T}} \frac{\partial \mathbf{u}_{j}}{\partial x}\right) + \operatorname{Tr}\left(\sum_{i=1}^{M} \mathbf{u}_{i} \frac{\partial \mathbf{u}_{i}^{\mathrm{T}}}{\partial x}\right) \quad (\because \operatorname{Tr}(AB) = \operatorname{Tr}(BA))$$

$$= \operatorname{Tr}\left(\sum_{i=1}^{M} \sum_{j=1}^{M} \frac{\lambda_{j}}{\lambda_{i}} \mathbf{u}_{j}^{\mathrm{T}} \mathbf{u}_{i} \mathbf{u}_{i}^{\mathrm{T}} \frac{\partial \mathbf{u}_{j}}{\partial x}\right) + \operatorname{Tr}\left(\sum_{i=1}^{M} \mathbf{u}_{i} \frac{\partial \mathbf{u}_{i}^{\mathrm{T}}}{\partial x}\right)$$

$$= \operatorname{Tr}\left(\sum_{i=1}^{M} \sum_{j=1}^{M} \frac{\lambda_{j}}{\lambda_{i}} \mathbf{u}_{i}^{\mathrm{T}}\right) + \operatorname{Tr}\left(\sum_{i=1}^{M} \mathbf{u}_{i} \frac{\partial \mathbf{u}_{i}^{\mathrm{T}}}{\partial x}\right)$$

$$= \operatorname{Tr}\left(\sum_{i=1}^{M} \frac{\partial \mathbf{u}_{i}}{\partial x} \mathbf{u}_{i}^{\mathrm{T}}\right) + \operatorname{Tr}\left(\sum_{i=1}^{M} \mathbf{u}_{i} \frac{\partial \mathbf{u}_{i}^{\mathrm{T}}}{\partial x}\right)$$

$$= \operatorname{Tr}\left(\sum_{i=1}^{M} \frac{\partial \mathbf{u}_{i}}{\partial x} \mathbf{u}_{i}^{\mathrm{T}}\right) + \operatorname{Tr}\left(\sum_{i=1}^{M} \mathbf{u}_{i} \frac{\partial \mathbf{u}_{i}^{\mathrm{T}}}{\partial x}\right)$$

$$= \operatorname{Tr}\left(\sum_{i=1}^{M} \frac{\partial \mathbf{u}_{i} \mathbf{u}_{i}^{\mathrm{T}}}{\partial x}\right) = \operatorname{Tr}\left(\frac{\partial \operatorname{U}^{\mathrm{U}^{\mathrm{T}}}}{\partial x}\right) = \operatorname{Tr}\left(\frac{\partial \operatorname{U}^{\mathrm{U}^{\mathrm{T}}}}{\partial x}\right) = \operatorname{Tr}\left(\frac{\partial \operatorname{U}^{\mathrm{U}^{\mathrm{T}}}}{\partial x}\right) = \operatorname{Tr}\left(\frac{\partial \operatorname{U}^{\mathrm{U}^{\mathrm{T}}}}{\partial x}\right)$$

Given $X = \{x_1, ..., x_N\}$ where each x_i is drawn *independently* from a *D*-dimensional Gaussian distribution. We estimate the parameters of the unknown distribution by maximizing the log likelihood function:

$$\ln \Pr(\mathbf{X}|\mu, \Sigma) = -\frac{ND}{2} \ln(2\pi) - \frac{N}{2} \ln|\Sigma| - \frac{1}{2} \sum_{i=1}^{N} (\mathbf{x}_i - \mu)^{\mathrm{T}} \Sigma^{-1} (\mathbf{x}_i - \mu)$$
 (64)

Firstly, we calculate the derivative of (64) w.r.t μ . Suppose resulting vector is \mathbf{r} , $\forall i \in 1, \ldots, D: \mathbf{r}_i = \frac{\partial \ln \Pr(\mathbf{X}|\boldsymbol{\mu},\boldsymbol{\Sigma})}{\partial \mathbf{r}_i}$. Consider a quadratic form $\Delta^2 = (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) = \sum_{i=1}^D \sum_{j=1}^D (\mathbf{x}_i - \boldsymbol{\mu}_i) (\mathbf{x}_j - \boldsymbol{\mu}_j) \boldsymbol{\Sigma}_{ij}^{-1} = \sum_{i=1}^D \sum_{j=1}^D (\mathbf{x}_i \mathbf{x}_j - \mathbf{x}_j \boldsymbol{\mu}_i - \mathbf{x}_i \boldsymbol{\mu}_j + \boldsymbol{\mu}_i \boldsymbol{\mu}_j) \boldsymbol{\Sigma}_{ij}^{-1}$.

$$\therefore \frac{\partial \Delta^2}{\partial \mu_i} = 2 \sum_{j=1}^{D} (\mu_j - \mathbf{x}_j) \Sigma_{ij}^{-1}
\therefore \frac{\partial \Delta^2}{\partial \mu} = 2 \Sigma^{-1} (\mu - \mathbf{x})$$
(65)

According to (65), we know that $\frac{\partial \ln \Pr(\mathbf{X}|\mu,\Sigma)}{\partial \mu} = \sum_{i=1}^{N} \Sigma^{-1}(\mathbf{x}_n - \mu)$ and set it to zero so that:

$$\mu_{\text{ML}} = \frac{1}{N} \sum_{i=1}^{N} \mathbf{x}_n \tag{66}$$

Then we fix μ in (64) to $\mu_{\rm ML}$ and calculate derivative of (64) w.r.t. Σ . Suppose the answer is matrix R

where $R_{ij} = \frac{\partial \ln \Pr(\mathbf{X}|\mu,\Sigma)}{\partial \Sigma_{ij}}$. Hence, we can consider the derivative of (64) w.r.t. a certain x:

$$\frac{\partial \ln \Pr(\mathbf{X}|\boldsymbol{\mu},\boldsymbol{\Sigma})}{\partial \boldsymbol{x}} = -\frac{N}{2} \frac{\ln |\boldsymbol{\Sigma}|}{\partial \boldsymbol{x}} - \frac{1}{2} \frac{\partial \sum_{n=1}^{N} (\mathbf{x}_{n} - \boldsymbol{\mu})^{\mathrm{T}} \boldsymbol{\Sigma}^{-1} (\mathbf{x}_{n} - \boldsymbol{\mu})}{\partial \boldsymbol{x}} \\
= -\frac{N}{2} \operatorname{Tr}(\boldsymbol{\Sigma}^{-1} \frac{\partial \boldsymbol{\Sigma}}{\partial \boldsymbol{x}}) - \frac{1}{2} \sum_{n=1}^{N} \frac{\partial \operatorname{Tr}((\mathbf{x}_{n} - \boldsymbol{\mu})^{\mathrm{T}} \boldsymbol{\Sigma}^{-1} (\mathbf{x}_{n} - \boldsymbol{\mu}))}{\partial \boldsymbol{x}} \quad (\because (60)) \\
= -\frac{N}{2} \operatorname{Tr}(\boldsymbol{\Sigma}^{-1} \frac{\partial \boldsymbol{\Sigma}}{\partial \boldsymbol{x}}) - \frac{1}{2} \sum_{n=1}^{N} \frac{\partial \operatorname{Tr}(\boldsymbol{\Sigma}^{-1} (\mathbf{x}_{n} - \boldsymbol{\mu}) (\mathbf{x}_{n} - \boldsymbol{\mu})^{\mathrm{T}})}{\partial \boldsymbol{x}} \quad (\because (57)) \\
= -\frac{N}{2} \operatorname{Tr}(\boldsymbol{\Sigma}^{-1} \frac{\partial \boldsymbol{\Sigma}}{\partial \boldsymbol{x}}) - \frac{1}{2} \frac{\partial \operatorname{Tr}(\boldsymbol{\Sigma}^{-1} \sum_{n=1}^{N} (\mathbf{x}_{n} - \boldsymbol{\mu}) (\mathbf{x}_{n} - \boldsymbol{\mu})^{\mathrm{T}})}{\partial \boldsymbol{x}} \quad (\because (58)) \\
= -\frac{N}{2} \operatorname{Tr}(\boldsymbol{\Sigma}^{-1} \frac{\partial \boldsymbol{\Sigma}}{\partial \boldsymbol{x}}) - \frac{1}{2} \operatorname{Tr}(\frac{\partial \boldsymbol{\Sigma}^{-1}}{\partial \boldsymbol{x}} \sum_{n=1}^{N} (\mathbf{x}_{n} - \boldsymbol{\mu}) (\mathbf{x}_{n} - \boldsymbol{\mu})^{\mathrm{T}}) \quad (\because (59)) \\
= -\frac{N}{2} \operatorname{Tr}(\frac{\partial \boldsymbol{\Sigma}}{\partial \boldsymbol{x}} \boldsymbol{\Sigma}^{-1} I) + \frac{1}{2} \operatorname{Tr}(\frac{\partial \boldsymbol{\Sigma}}{\partial \boldsymbol{x}} \boldsymbol{\Sigma}^{-1} \sum_{n=1}^{N} (\mathbf{x}_{n} - \boldsymbol{\mu}) (\mathbf{x}_{n} - \boldsymbol{\mu})^{\mathrm{T}} \boldsymbol{\Sigma}^{-1}) \quad (\because (56)) \\
= -\frac{N}{2} \operatorname{Tr}(\frac{\partial \boldsymbol{\Sigma}}{\partial \boldsymbol{x}} \boldsymbol{\Sigma}^{-1} \boldsymbol{\Sigma}^{-1}) + \frac{1}{2} \operatorname{Tr}(\frac{\partial \boldsymbol{\Sigma}}{\partial \boldsymbol{x}} \boldsymbol{\Sigma}^{-1} \sum_{n=1}^{N} (\mathbf{x}_{n} - \boldsymbol{\mu}) (\mathbf{x}_{n} - \boldsymbol{\mu})^{\mathrm{T}} \boldsymbol{\Sigma}^{-1}) \\
= \frac{1}{2} \operatorname{Tr}(\frac{\partial \boldsymbol{\Sigma}}{\partial \boldsymbol{x}} \boldsymbol{\Sigma}^{-1} \{\sum_{n=1}^{N} (\mathbf{x}_{n} - \boldsymbol{\mu}) (\mathbf{x}_{n} - \boldsymbol{\mu}) (\mathbf{x}_{n} - \boldsymbol{\mu})^{\mathrm{T}} \boldsymbol{\Sigma}^{-1}) \\
= \frac{1}{2} \operatorname{Tr}(\frac{\partial \boldsymbol{\Sigma}}{\partial \boldsymbol{x}} \boldsymbol{\Sigma}^{-1} \{\sum_{n=1}^{N} (\mathbf{x}_{n} - \boldsymbol{\mu}) (\mathbf{x}_{n} - \boldsymbol{\mu})^{\mathrm{T}} - N \boldsymbol{\Sigma} \} \boldsymbol{\Sigma}^{-1}) \right)$$

Thus, we require that $(\mathbf{x}_n - \boldsymbol{\mu})(\mathbf{x}_n - \boldsymbol{\mu})^{\mathrm{T}} - N\Sigma = 0$, i.e.,

$$\Sigma = \frac{1}{N} \sum_{i=1}^{N} (\mathbf{x}_n - \boldsymbol{\mu}) (\mathbf{x}_n - \boldsymbol{\mu})^{\mathrm{T}}$$
(68)