

Probabilistic and Unsupervised Learning

Phu Sakulwongtana

1 Probability Basis

1.1 Statistics Introduction

Remark 1. Before we start anything, let's recall the definition of all distributions (that would be used):

Distribution	Sample Space	Probability Density Function	$\mathbb{E}[\mathbf{x}]$	$\text{var}(\mathbf{x})$	Parameter
Bernoulli	$\{0, 1\}$	$\theta^x(1 - \theta)^{1-x}$	θ	$\theta(1 - \theta)$	$\theta \in [0, 1]$
Bernoulli(+)	$\left\{ \begin{array}{l} \mathbf{x} \in \{0, 1\}^D \\ \sum_{i=1}^D x_i = 1 \end{array} \right.$	$\prod_{i=1}^D \theta_i^{x_i}$	θ_i	$\text{var}[x_i] = \theta_i(1 - \theta_i)$	$\left\{ \begin{array}{l} 0 \leq \boldsymbol{\theta} \leq 1 \\ \sum_{i=1}^D \theta_i = 1 \end{array} \right.$
Binomial	$[N]$	$\binom{N}{x} \theta^x(1 - \theta)^{1-x}$	$N\theta$	$N\theta(1 - \theta)$	$\theta \in [0, 1]$
Multinomial	$\left\{ \begin{array}{l} \mathbf{x} \in [N]^D \\ \sum_{i=1}^D x_i = N \end{array} \right.$	$\frac{N!}{x_1 x_2 \cdots x_D} \prod_{i=1}^D \theta_i^{x_i}$	$N\theta_i$	$N\theta_i(1 - \theta_i)$	$\left\{ \begin{array}{l} 0 \leq \boldsymbol{\theta} \leq 1 \\ \sum_{i=1}^D \theta_i = 1 \end{array} \right.$
Gaussian	\mathbb{R}	$\frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{1}{2\sigma^2}(x - \mu)^2\right\}$	μ	σ^2	$\left\{ \begin{array}{l} \mu \in \mathbb{R} \\ \sigma \in \mathbb{R}_{\geq 0} \end{array} \right.$
Multinormal	\mathbb{R}^D	$\frac{1}{\sqrt{ 2\pi\boldsymbol{\Sigma} }} \exp\left\{-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right\}$	$\boldsymbol{\mu}$	$\boldsymbol{\Sigma}$	$\left\{ \begin{array}{l} \boldsymbol{\mu} \in \mathbb{R}^D \\ \boldsymbol{\Sigma} \in \mathbb{S}_+^{D \times D} \end{array} \right.$
Beta	$[0, 1]$	$\frac{\Gamma}{\Gamma(a) + \Gamma(b)} x^{a-1}(1 - x)^{b-1}$	$\frac{a}{a + b}$	$\frac{ab}{(a + b)^2(a + b + 1)}$	$\left\{ \begin{array}{l} a > 0 \\ b > 0 \end{array} \right.$
Dirichlet	$\left\{ \begin{array}{l} \mathbf{x} \in \mathbb{R}^D \\ 0 \leq \mathbf{x} \leq 1 \\ \sum_{i=1}^D x_i = 1 \end{array} \right.$	$\frac{\Gamma(\hat{\theta})}{\Gamma(\theta_1) \cdots \Gamma(\theta_D)} \prod_{i=1}^D x_i^{\theta_i - 1}$	$\theta_i / \hat{\theta}$	$\left\{ \begin{array}{l} \text{var}[x_i] = \frac{\theta_i(\hat{\theta} - \theta_i)}{\hat{\theta}^2(\hat{\theta} + 1)} \\ \text{cov}[x_i x_j] = \frac{-\theta_i \theta_j}{\hat{\theta}^2(\hat{\theta} + 1)} \end{array} \right.$	$\theta_i > 0$
Gamma	$x > 0$	$\frac{1}{\Gamma(a)} b^a x^{a-1} \exp(-bx)$	$\frac{a}{b}$	$\frac{a}{b^2}$	$\left\{ \begin{array}{l} a > 0 \\ b > 0 \end{array} \right.$
Wishart*	$\boldsymbol{\Lambda}^{-1} \in \mathbb{S}_+^{D \times D}$	$B(\mathbf{W}, \nu) \boldsymbol{\Lambda} ^{(\nu - D - 1)/2} \exp\left(-\frac{1}{2} \text{Tr}(\mathbf{W}^{-1} \boldsymbol{\Lambda})\right)$	$\nu \mathbf{W}$	—	$\left\{ \begin{array}{l} \mathbf{W} \in \mathbb{S}_+^{D \times D} \\ \nu > D - 1 \end{array} \right.$
Poisson	\mathbb{N}_0	$\frac{\lambda^x e^{-\lambda}}{x!}$	λ	λ	$\lambda > 0$

Table 1: *1D Wishart is Gamma with $a = \nu/2$ and $1/2W$

where $\mathbb{S}_+^{D \times D}$ is the set of positive definite matrix of size $D \times D$, also we have the following addition definitions:

$$\Gamma(z) = \int_0^\infty x^{z-1} \exp(-x) dx \quad \hat{\theta} = \sum_{i=1}^D \theta_i$$

$$B(\mathbf{W}, \nu) = |\mathbf{W}|^{-\nu/2} \left(2^{\nu D/2} \pi^{D(D-1)/4} \prod_{i=1}^D \Gamma\left(\frac{\nu+1-i}{2}\right) \right)^{-1}$$

Remark 2. (Basic Quantities) Now we will consider the following probability facts, which would be useful in the future:

$$\begin{aligned} \mathbb{E}[f(\mathbf{x})] &= \int p(\mathbf{x})f(\mathbf{x}) d\mathbf{x} & \text{cov}(\mathbf{x}, \mathbf{y}) &= \mathbb{E}_{\mathbf{x}, \mathbf{y}}[(\mathbf{x} - \mathbb{E}[\mathbf{x}])(\mathbf{y} - \mathbb{E}[\mathbf{y}])^T] \\ & & &= \mathbb{E}_{\mathbf{x}, \mathbf{y}}[\mathbf{x}\mathbf{y}^T] - \mathbb{E}[\mathbf{x}]\mathbb{E}[\mathbf{y}^T] \\ \text{var}[f(\mathbf{x})] &= \mathbb{E}[(f(\mathbf{x}) - \mathbb{E}[f(\mathbf{x})])^2] \\ &= \mathbb{E}[f(\mathbf{x})^2] - \mathbb{E}[f(\mathbf{x})]^2 \end{aligned}$$

Remark 3. (Additional Quantities) We have the following equality:

$$\mathbb{E}_{\mathbf{x}}[f(\mathbf{x})] = \mathbb{E}_{\mathbf{z}} [\mathbb{E}_{\mathbf{x}|\mathbf{z}}[f(\mathbf{x})]] \quad \mathbb{V}_{\mathbf{x}}[\mathbf{x}] = \mathbb{E}_{\mathbf{z}} [\mathbb{V}[\mathbf{x}|\mathbf{z}]] + \mathbb{V}_{\mathbf{z}}[\mathbb{E}[\mathbf{x}|\mathbf{z}]]$$

1.2 Linear Algebra

Proposition 1.1. (Woodbury Identity) This following identity can helps the computation as follows:

$$(\mathbf{A} + \mathbf{B}\mathbf{D}^{-1}\mathbf{C})^{-1} = \mathbf{A}^{-1} - \mathbf{A}^{-1}\mathbf{B}(\mathbf{D} + \mathbf{C}\mathbf{A}^{-1}\mathbf{B})^{-1}\mathbf{C}\mathbf{A}^{-1}$$

This identity is useful when $\mathbf{A} \in \mathbb{R}^{a \times a}$ is large and diagonal (easy to invert), while $\mathbf{B} \in \mathbb{R}^{a \times b}$ has many rows but few columns ($a > b$) conversely for $\mathbf{C} \in \mathbb{R}^{b \times a}$. The RHS is simpler than LHS.

Proof. This can be proven easily as:

$$\begin{aligned} & \left[\mathbf{A}^{-1} - \mathbf{A}^{-1}\mathbf{B}(\mathbf{D} + \mathbf{C}\mathbf{A}^{-1}\mathbf{B})^{-1}\mathbf{C}\mathbf{A}^{-1} \right] (\mathbf{A} + \mathbf{B}\mathbf{D}^{-1}\mathbf{C}) \\ &= \mathbf{A}^{-1}(\mathbf{A} + \mathbf{B}\mathbf{D}^{-1}\mathbf{C}) - \mathbf{A}^{-1}\mathbf{B}(\mathbf{D} + \mathbf{C}\mathbf{A}^{-1}\mathbf{B})^{-1}\mathbf{C}\mathbf{A}^{-1}(\mathbf{A} + \mathbf{B}\mathbf{D}^{-1}\mathbf{C}) \\ &= \mathbf{A}^{-1}\mathbf{A} + \mathbf{A}^{-1}\mathbf{B}\mathbf{D}^{-1}\mathbf{C} - \mathbf{A}^{-1}\mathbf{B}(\mathbf{D} + \mathbf{C}\mathbf{A}^{-1}\mathbf{B})^{-1}\mathbf{C}\mathbf{A}^{-1}\mathbf{A} \\ & \quad - \mathbf{A}^{-1}\mathbf{B}(\mathbf{D} + \mathbf{C}\mathbf{A}^{-1}\mathbf{B})^{-1}\mathbf{C}\mathbf{A}^{-1}\mathbf{B}\mathbf{D}^{-1}\mathbf{C} \\ &= \mathbf{I} + \mathbf{A}^{-1}\mathbf{B}\mathbf{D}^{-1}\mathbf{C} - \mathbf{A}^{-1}\mathbf{B}(\mathbf{D} + \mathbf{C}\mathbf{A}^{-1}\mathbf{B})^{-1}\mathbf{D}\mathbf{D}^{-1}\mathbf{C} \\ & \quad - \mathbf{A}^{-1}\mathbf{B}(\mathbf{D} + \mathbf{C}\mathbf{A}^{-1}\mathbf{B})^{-1}\mathbf{C}\mathbf{A}^{-1}\mathbf{B}\mathbf{D}^{-1}\mathbf{C} \\ &= \mathbf{I} + \mathbf{A}^{-1}\mathbf{B}\mathbf{D}^{-1}\mathbf{C} - \mathbf{A}^{-1}\mathbf{B} \left[(\mathbf{D} + \mathbf{C}\mathbf{A}^{-1}\mathbf{B})^{-1}\mathbf{D} + (\mathbf{D} + \mathbf{C}\mathbf{A}^{-1}\mathbf{B})^{-1}\mathbf{C}\mathbf{A}^{-1}\mathbf{B} \right] \mathbf{D}^{-1}\mathbf{C} \\ &= \mathbf{I} + \mathbf{A}^{-1}\mathbf{B}\mathbf{D}^{-1}\mathbf{C} - \mathbf{A}^{-1}\mathbf{B} \left[(\mathbf{D} + \mathbf{C}\mathbf{A}^{-1}\mathbf{B})^{-1}(\mathbf{D} + \mathbf{C}\mathbf{A}^{-1}\mathbf{B}) \right] \mathbf{D}^{-1}\mathbf{C} \\ &= \mathbf{I} + \mathbf{A}^{-1}\mathbf{B}\mathbf{D}^{-1}\mathbf{C} - \mathbf{A}^{-1}\mathbf{B}\mathbf{D}^{-1}\mathbf{C} = \mathbf{I} \end{aligned}$$

□

Proposition 1.2. (Another Identity) Another useful identity can be stated as:

$$(\mathbf{P}^{-1} + \mathbf{B}^T \mathbf{R}^{-1} \mathbf{B})^{-1} \mathbf{B}^T \mathbf{R}^{-1} = \mathbf{P} \mathbf{B}^T (\mathbf{B} \mathbf{P} \mathbf{B}^T + \mathbf{R})^{-1}$$

Proof. This can be proven easily as:

$$\begin{aligned}
& (\mathbf{P}^{-1} + \mathbf{B}^T \mathbf{R}^{-1} \mathbf{B})^{-1} \mathbf{B}^T \mathbf{R}^{-1} (\mathbf{B} \mathbf{P} \mathbf{B}^T + \mathbf{R}) \\
&= (\mathbf{P}^{-1} + \mathbf{B}^T \mathbf{R}^{-1} \mathbf{B})^{-1} \mathbf{B}^T \mathbf{R}^{-1} \mathbf{B} \mathbf{P} \mathbf{B}^T + (\mathbf{P}^{-1} + \mathbf{B}^T \mathbf{R}^{-1} \mathbf{B})^{-1} \mathbf{B}^T \mathbf{R}^{-1} \mathbf{R} \\
&= (\mathbf{P}^{-1} + \mathbf{B}^T \mathbf{R}^{-1} \mathbf{B})^{-1} \mathbf{B}^T \mathbf{R}^{-1} \mathbf{B} \mathbf{P} \mathbf{B}^T + (\mathbf{P}^{-1} + \mathbf{B}^T \mathbf{R}^{-1} \mathbf{B})^{-1} \mathbf{B}^T \\
&= (\mathbf{P}^{-1} + \mathbf{B}^T \mathbf{R}^{-1} \mathbf{B})^{-1} \left[\mathbf{B}^T \mathbf{R}^{-1} \mathbf{B} \mathbf{P} + \mathbf{P}^{-1} \mathbf{P} \right] \mathbf{B}^T \\
&= (\mathbf{P}^{-1} + \mathbf{B}^T \mathbf{R}^{-1} \mathbf{B})^{-1} \left[\mathbf{B}^T \mathbf{R}^{-1} \mathbf{B} + \mathbf{P}^{-1} \right] \mathbf{P} \mathbf{B}^T = \mathbf{P} \mathbf{B}^T
\end{aligned}$$

□

Remark 4. (More Matrix Identities) We have the following:

$$\text{Tr}(\mathbf{A} \mathbf{B} \mathbf{C}) = \text{Tr}(\mathbf{C} \mathbf{A} \mathbf{B}) = \text{Tr}(\mathbf{B} \mathbf{C} \mathbf{A}) \quad |\mathbf{A}^{-1}| = \frac{1}{|\mathbf{A}|} \quad |\mathbf{A} \mathbf{B}| = |\mathbf{A}| |\mathbf{B}|$$

Remark 5. (Additional Identity)

$$\left| \mathbf{I}_N + \mathbf{A} \mathbf{B}^T \right| = \left| \mathbf{I}_M + \mathbf{A}^T \mathbf{B} \right|$$

This also implies that $\left| \mathbf{I}_N + \mathbf{a} \mathbf{b}^T \right| = 1 + \mathbf{a}^T \mathbf{b}$

Remark 6. (Matrix Derivative - Basics) We still use the following facts:

$$\frac{\partial}{\partial \mathbf{x}} (\mathbf{x}^T \mathbf{a}) = \frac{\partial}{\partial \mathbf{x}} (\mathbf{a}^T \mathbf{x}) = \mathbf{a} \quad \frac{\partial}{\partial \mathbf{x}} \mathbf{A} \mathbf{B} = \frac{\partial \mathbf{A}}{\partial \mathbf{x}} \mathbf{B} + \frac{\partial \mathbf{B}}{\partial \mathbf{x}} \mathbf{A} \quad \frac{\partial}{\partial \mathbf{x}} \mathbf{x}^T \mathbf{A} \mathbf{x} = (\mathbf{A} + \mathbf{A}^T) \mathbf{x}$$

Proposition 1.3. *The derivative of the inverse matrix is given by:*

$$\frac{\partial}{\partial \mathbf{x}} (\mathbf{A}^{-1}) = -\mathbf{A}^{-1} \frac{\partial \mathbf{A}}{\partial \mathbf{x}} \mathbf{A}^{-1}$$

Proof. We consider differentiate the following equation $\mathbf{A}^{-1} \mathbf{A} = \mathbf{I}$ as we have:

$$\mathbf{0} = \left(\frac{\partial}{\partial \mathbf{x}} \mathbf{I} \right) \mathbf{A}^{-1} = \left(\frac{\partial}{\partial \mathbf{x}} \mathbf{A}^{-1} \mathbf{A} \right) \mathbf{A}^{-1} = \left(\frac{\partial \mathbf{A}^{-1}}{\partial \mathbf{x}} \mathbf{A} + \mathbf{A}^{-1} \frac{\partial \mathbf{A}}{\partial \mathbf{x}} \right) \mathbf{A}^{-1} = \frac{\partial \mathbf{A}^{-1}}{\partial \mathbf{x}} + \mathbf{A}^{-1} \frac{\partial \mathbf{A}}{\partial \mathbf{x}} \mathbf{A}^{-1}$$

With algebraic manipulation the proposition is proven. □

Remark 7. (Additional Matrix Derivative)

$$\begin{aligned}
\frac{\partial}{\partial \mathbf{A}} \text{Tr}(\mathbf{A} \mathbf{B}) &= \mathbf{B}^T & \frac{\partial}{\partial \mathbf{A}} \text{Tr}(\mathbf{A}^T \mathbf{B}) &= \mathbf{B} \\
\frac{\partial}{\partial \mathbf{A}} \text{Tr}(\mathbf{A}) &= \mathbf{I} & \frac{\partial}{\partial \mathbf{A}} \text{Tr}(\mathbf{A} \mathbf{B} \mathbf{A}^T) &= \mathbf{A} (\mathbf{B} + \mathbf{B}^T)
\end{aligned}$$

Proposition 1.4. *The matrix can be diagonalization as:*

$$\mathbf{A} = \mathbf{U} \mathbf{\Lambda} \mathbf{U}^T = \sum_{i=1}^N \lambda_i \mathbf{u}_i \mathbf{u}_i^T$$

\mathbf{U} to be the matrix constructed that has the column as the eigenvectors \mathbf{u}_i . The matrix $\mathbf{\Lambda}$ is the diagonal matrix, whose diagonal element is the eigenvalue λ_i .

Proof. Recall that the square matrix $\mathbf{A} \in \mathbb{R}^{M \times M}$'s eigenvalue and eigenvector, which are given by:

$$\mathbf{A} \mathbf{u}_i = \lambda_i \mathbf{u}_i$$

For $i = 1, \dots, M$ where \mathbf{u}_i is the eigenvector and λ_i is the corresponding eigenvalue. This means that $\mathbf{A} \mathbf{U} = \mathbf{U} \mathbf{\Lambda}$; furthermore, $\mathbf{U} \mathbf{U}^T = \mathbf{U}^T \mathbf{U} = \mathbf{I}$ as we have $|\mathbf{U}| = 1$. The identity follows by right multiplying with \mathbf{U}^T . □

Proposition 1.5. We can show that the determinant and trace of the matrix to be:

$$|\mathbf{A}| = \prod_{i=1}^M \lambda_i \quad \text{Tr}(\mathbf{A}) = \sum_{i=1}^M \lambda_i$$

This follows from the identity of determinant and the cyclic properties of trace.

Proposition 1.6. We now going to show very useful identity

$$\frac{\partial}{\partial x} \ln |\mathbf{A}| = \text{Tr} \left(\mathbf{A}^{-1} \frac{\partial \mathbf{A}}{\partial x} \right)$$

Proof. We consider the LHS first, as we have:

$$\frac{\partial}{\partial x} \ln |\mathbf{A}| = \frac{\partial}{\partial x} \ln \left(\prod_{i=1}^M \lambda_i \right) = \frac{\partial}{\partial x} \sum_{i=1}^M \ln \lambda_i = \sum_{i=1}^M \frac{1}{\lambda_i} \frac{\partial \lambda_i}{\partial x}$$

Now, consider the RHS, which we have:

$$\begin{aligned} \text{Tr} \left(\left[\sum_{i=1}^M \frac{1}{\lambda_i} \mathbf{u}_i \mathbf{u}_i^T \right] \left[\sum_{i=1}^M \frac{\partial \lambda_i}{\partial x} \mathbf{u}_i \mathbf{u}_i^T \right] \right) &= \text{Tr} \left(\sum_{i=1}^M \sum_{j=1}^M \frac{1}{\lambda_i} \mathbf{u}_i \mathbf{u}_i^T \frac{\partial \lambda_j}{\partial x} \mathbf{u}_j \mathbf{u}_j^T \right) \\ &= \sum_{i=1}^M \sum_{j=1}^M \frac{1}{\lambda_i} \frac{\partial \lambda_j}{\partial x} \text{Tr} (\mathbf{u}_i \mathbf{u}_i^T \mathbf{u}_j \mathbf{u}_j^T) \\ &= \sum_{i=1}^M \sum_{j=1}^M \frac{1}{\lambda_i} \frac{\partial \lambda_j}{\partial x} \text{Tr} (\mathbf{u}_i^T \mathbf{u}_j \mathbf{u}_j^T \mathbf{u}_i) = \sum_{i=1}^M \frac{1}{\lambda_i} \frac{\partial \lambda_i}{\partial x} \end{aligned}$$

The equality is proven. \square

Corollary 1.1. The above proposition implies that:

$$\frac{\partial}{\partial \mathbf{A}} \ln |\mathbf{A}| = (\mathbf{A}^{-1})^T$$

Proof. Consider the following partial derivative, we have:

$$\frac{\partial}{\partial a_{cd}} \ln |\mathbf{A}| = \text{Tr} \left(\mathbf{A}^{-1} \frac{\partial \mathbf{A}}{\partial a_{cd}} \right) = a_{dc}^{-1}$$

where we have a_{ij}^{-1} be the (i, j) -th element of the matrix \mathbf{A}^{-1} . Thus we have proven the equality. \square

Proposition 1.7. We can show that:

$$\frac{\partial}{\partial \mathbf{A}} \text{Tr}[\mathbf{A}^T \mathbf{B} \mathbf{A} \mathbf{C}] = \mathbf{B} \mathbf{A} \mathbf{C} + \mathbf{B}^T \mathbf{A} \mathbf{C}^T$$

Proof. We will use the identity mapping $F_1(\cdot)$ and $F_2(\cdot)$ to make the differetiation easier:

$$\begin{aligned} \frac{\partial}{\partial \mathbf{A}} \text{Tr}[\mathbf{A}^T \mathbf{B} \mathbf{A} \mathbf{C}] &= \frac{\partial}{\partial \mathbf{A}} \text{Tr}[F_1(\mathbf{A})^T \mathbf{B} F_2(\mathbf{A}) \mathbf{C}] \\ &= \frac{\partial}{\partial \mathbf{F}_1} \text{Tr}[\mathbf{F}_1^T \mathbf{B} \mathbf{F}_2 \mathbf{C}] \frac{\partial F_1}{\partial \mathbf{A}} + \frac{\partial}{\partial \mathbf{F}_2} \text{Tr}[\mathbf{F}_1^T \mathbf{B} \mathbf{F}_2 \mathbf{C}] \frac{\partial F_2}{\partial \mathbf{A}} \\ &= \frac{\partial}{\partial \mathbf{F}_1} \text{Tr}[\mathbf{F}_1^T \mathbf{B} \mathbf{F}_2 \mathbf{C}] \frac{\partial F_1}{\partial \mathbf{A}} + \frac{\partial}{\partial \mathbf{F}_2} \text{Tr}[\mathbf{C} \mathbf{F}_1^T \mathbf{B} \mathbf{F}_2] \frac{\partial F_2}{\partial \mathbf{A}} \\ &= \frac{\partial}{\partial \mathbf{F}_1} \text{Tr}[\mathbf{F}_1^T \mathbf{B} \mathbf{F}_2 \mathbf{C}] \frac{\partial F_1}{\partial \mathbf{A}} + \frac{\partial}{\partial \mathbf{F}_2} \text{Tr}[\mathbf{F}_2^T \mathbf{B}^T \mathbf{F}_1 \mathbf{C}^T] \frac{\partial F_2}{\partial \mathbf{A}} \\ &= \mathbf{B} \mathbf{F}_2 \mathbf{C} + \mathbf{B}^T \mathbf{F}_1 \mathbf{C}^T = \mathbf{B} \mathbf{A} \mathbf{C} + \mathbf{B}^T \mathbf{A} \mathbf{C}^T \end{aligned}$$

where we have $F_1(\mathbf{A}) = F_2(\mathbf{A}) = \mathbf{F}_1 = \mathbf{F}_2 = \mathbf{A}$ \square

Remark 8. (Notes on Symmetric Matrix) We can construct the symmetric matrix from any kind of matrix \mathbf{A} using the following formula:

$$\mathbf{M} = \frac{\mathbf{A} + \mathbf{A}^T}{2}$$

One can show that its eigenvalue is real (given real symmetric matrix), as we consider the complex of, assuming \mathbf{x} is an eigenvector of \mathbf{M} with its eigenvalue to be λ :

$$\langle \mathbf{M}\mathbf{x}, \mathbf{M}\mathbf{x} \rangle = \mathbf{x}^* \mathbf{M}^* \mathbf{M}\mathbf{x} = \mathbf{x}^* \mathbf{M}\mathbf{M}\mathbf{x} = \mathbf{x}^* \lambda^2 \mathbf{x} = \lambda^2 \|\mathbf{x}\|^2$$

Where $\mathbf{M}^* = \bar{\mathbf{M}}^T$, and so λ^2 is real a non-negative, thus being a real number.

Remark 9. (Notes on Square-Root of Matrix) Given positive semi-definite matrix \mathbf{A} , one can show that there a matrix \mathbf{B} such that $\mathbf{A} = \mathbf{B}\mathbf{B}$ (or $\mathbf{B}^T\mathbf{B}$ as \mathbf{B} is symmetric as we will shown later). Given the eigendecomposition of \mathbf{A} to be $\mathbf{U}\mathbf{\Lambda}\mathbf{U}^T$, matrix \mathbf{B} is $\mathbf{U}\sqrt{\mathbf{\Lambda}}\mathbf{U}^T$, where $\sqrt{\mathbf{\Lambda}}$ is the matrix contains square root of the diagonal of $\mathbf{\Lambda}$:

$$\mathbf{B}\mathbf{B} = (\mathbf{U}\sqrt{\mathbf{\Lambda}}\mathbf{U}^T)(\mathbf{U}\sqrt{\mathbf{\Lambda}}\mathbf{U}^T) = \mathbf{U}\sqrt{\mathbf{\Lambda}}\sqrt{\mathbf{\Lambda}}\mathbf{U}^T = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^T = \mathbf{A}$$

Please note that, since \mathbf{A} is positive semi-definite, the eigenvalue is non-negative implies that \mathbf{B} has real value eigenvalue (and non-negative), and so \mathbf{B} is symmetric. Finally, if \mathbf{A} is positive define, then there is a *unique* \mathbf{B} .

Proposition 1.8. (Partition Matrix) The block matrix can be inversed as:

$$\begin{pmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{pmatrix} = \begin{pmatrix} \mathbf{M} & -\mathbf{M}\mathbf{B}\mathbf{D}^{-1} \\ -\mathbf{D}^{-1}\mathbf{C}\mathbf{M} & \mathbf{D}^{-1} + \mathbf{D}^{-1}\mathbf{C}\mathbf{M}\mathbf{B}\mathbf{D}^{-1} \end{pmatrix}$$

where $\mathbf{M} = (\mathbf{A} - \mathbf{B}\mathbf{D}^{-1}\mathbf{C})^{-1}$

1.3 Optimization

Definition 1.1. (Constraint Optimization) The constraint optimization is the optimization problem is in the form of:

$$\begin{aligned} \min_{\mathbf{x}} \quad & f_0(\mathbf{x}) \\ \text{s.t} \quad & f_i(\mathbf{x}) \leq 0 \quad i \in \mathcal{I} = [m] \\ & h_i(\mathbf{x}) = 0 \quad i \in \mathcal{E} = [p] \end{aligned}$$

Definition 1.2. (KKT Condition) The constraint optimization problem given above can be solved using the KKT condition. Before that, we consider the Lagragian to be defined as:

$$\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}) = f_0(\mathbf{x}) + \sum_{i \in \mathcal{I}} \lambda_i f_i(\mathbf{x}) + \sum_{i \in \mathcal{E}} \mu_i h_i(\mathbf{x})$$

The KKT condition is given by:

$$\begin{aligned} \nabla \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}) &= 0 \\ f_i(\mathbf{x}) &\leq 0 \quad \text{for } i \in \mathcal{I} \\ h_i(\mathbf{x}) &= 0 \quad \text{for } i \in \mathcal{E} \\ \lambda_i &\geq 0 \quad \text{for } i \in \mathcal{I} \\ \lambda_i f_i(\mathbf{x}) &= 0 \quad \text{for } i \in \mathcal{I} \end{aligned}$$

1.4 What are we going to do ?

Theorem 1.1. (Bayes' Theorem) One can show that:

$$p(A|B) = \frac{p(B|A)p(A)}{p(B)}$$

Remark 10. There are many ways to learn the parameter given the dataset $\mathcal{D} = \{\mathbf{x}_i\}_{i=1}^N$, as we have:

- *Maximum Likelihood:* Find the parameter $\boldsymbol{\theta}_{\text{ML}}$ such that it maximizes the log-likelihood as we have:

$$\boldsymbol{\theta}_{\text{ML}} = \arg \max_{\boldsymbol{\theta}} \log p(\mathcal{D}|\boldsymbol{\theta})$$

- *Bayesian Inference:* Find the distribution over the parameter $\boldsymbol{\theta}$ using Bayes' Theorem:

$$p(\boldsymbol{\theta}|\mathcal{D}) = \frac{p(\boldsymbol{\theta})p(\mathcal{D}|\boldsymbol{\theta})}{p(\mathcal{D})}$$

- *Maximum A Posteriori:* Find the mode of the posterior distribution over parameter

$$\boldsymbol{\theta}_{\text{MAP}} = \arg \max_{\boldsymbol{\theta}} \log P(\boldsymbol{\theta}|\mathcal{D})$$

The main problems/solutions of this works is simply trying to get better estimate of $p(\boldsymbol{\theta}|\mathcal{D})$ as it maybe intractable to calculate.

Definition 1.3. (Bayesian Model) The model is $\mathcal{M} = \{P(\cdot|\theta) : \theta \in \mathcal{T}\}$, where they are the distribution of a single random variable $X \sim P(\cdot|\theta)$. Given the prior, we also have a prior π on parameter space \mathcal{T} . The data is generated by the following process:

$$\Theta \sim \pi \quad X_1, \dots, X_n | \Theta \sim_{\text{iid}} p(\cdot|\Theta)$$

The tuple (\mathcal{M}, π) is the Bayesian model.

Remark 11. (Model Selection) Given various kinds of model $\mathcal{M}_1, \mathcal{M}_2, \dots$. The following set of likelihood associated with \mathcal{M}_i is

$$\{p(\mathbf{x}|\boldsymbol{\theta}_i, \mathcal{M}_i) : \boldsymbol{\theta}_i \in \mathcal{T}_i\}$$

We are interested in selecting the \mathcal{M}_i . Starting with the prior $p(\mathcal{M}_i)$ and the prior probability of parameter, given the model \mathcal{M}_i , is $p(\boldsymbol{\theta}_i|\mathcal{M}_i)$. Finally, the data probability is, where we assume the iid of the dataset:

$$p(\mathcal{D}|\boldsymbol{\theta}_i, \mathcal{M}_i) = \prod_{i=1}^N p(\mathbf{x}_i|\boldsymbol{\theta}_i, \mathcal{M}_i)$$

Now, we can find the posterior of the parameter given the model \mathcal{M}_i together with dataset evidence:

$$p(\boldsymbol{\theta}_i|\mathcal{D}, \mathcal{M}_i) = \frac{p(\mathcal{D}|\boldsymbol{\theta}_i, \mathcal{M}_i)P(\boldsymbol{\theta}_i|\mathcal{M}_i)}{p(\mathcal{D}|\mathcal{M}_i)} \quad p(\mathcal{D}|\mathcal{M}_i) = \int p(\mathcal{D}|\boldsymbol{\theta}_i, \mathcal{M}_i)p(\boldsymbol{\theta}_i|\mathcal{M}_i) d\boldsymbol{\theta}_i$$

We can perform Bayesian inference over the model \mathcal{M}_i

$$p(\mathcal{M}_i|\mathcal{D}) = \frac{p(\mathcal{D}|\mathcal{M}_i)p(\mathcal{M}_i)}{p(\mathcal{D})}$$

Now, we have the distribution over possible models.

1.5 Exponential Family and Friends

Definition 1.4. (Exponential Family) The set of probability distribution $\{p(\cdot|\theta) : \theta \in \mathcal{T}\}$, where \mathcal{T} is the parameter space, is exponential family if we have the distribution of the form:

$$p(\mathbf{x}|\boldsymbol{\theta}) = f(\mathbf{x})g(\boldsymbol{\theta}) \exp\left(\boldsymbol{\phi}(\boldsymbol{\theta})^T \mathbf{T}(\mathbf{x})\right)$$

where each components are given (and named) as:

- Sufficient Statistics: $\mathbf{T} : \mathcal{X} \rightarrow \mathbb{R}^m$
- Natural Parameter: $\phi : \mathcal{T} \rightarrow \mathbb{R}^m$
- Auxilliary Functions $f : \mathcal{X} \rightarrow \mathbb{R}_{\geq 0}$ and $g : \mathcal{T} \rightarrow \mathbb{R}_{> 0}$ (normalizing factor)

Please note that the function g has the following properties:

$$g(\boldsymbol{\theta}) \int f(\mathbf{x}) \exp\left(\boldsymbol{\phi}(\boldsymbol{\theta})^T \mathbf{T}(\mathbf{x})\right) d\mathbf{x} = 1$$

Remark 12. Let's consider the example of exponential families as:

Distribution	$\boldsymbol{\phi}(\boldsymbol{\theta})$	$\mathbf{T}(\mathbf{x})$
Bernoulli/Binomial	$\ln\left(\frac{\theta}{1-\theta}\right)$	x
Bernoulli(+)/Multinomial	$[\ln \theta_1, \ln \theta_2, \dots, \ln \theta_D]$	$[x_1, x_2, \dots, x_D]$
Gaussian	$[\boldsymbol{\mu}/\sigma^2, -1/2\sigma^2]$	$[x, x^2]$
Multinomial	$\left[-\frac{1}{2} \text{Vec}(\boldsymbol{\Sigma}^{-1}), \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}\right]$	$[\text{Vec}(\mathbf{x}\mathbf{x}^T), \mathbf{x}]$
Beta	$[a-1, b-1]$	$[\ln x, \ln(1-x)]$
Dirichlet	$[a_1-1, a_2-1, \dots, a_D-1]$	$[\ln x_1, \ln x_2, \dots, \ln x_D]$
Gamma	$[a-1, -b]$	$[\ln x, x]$
Poisson	$\ln \lambda$	x

where we have the following operation $\text{Vec} : \mathbb{R}^{n \times m} \rightarrow \mathbb{R}^{n \cdot m \times 1}$ is defined as:

$$\text{Vec}(\mathbf{X}) = [X_{11}, \dots, X_{n1}, \dots, X_{1m}, \dots, X_{nm}]^T$$

Proposition 1.9. *The normal distribution can be written as:*

$$\begin{aligned} & \frac{1}{\sqrt{|2\pi\boldsymbol{\Sigma}|}} \exp\left\{-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x}-\boldsymbol{\mu})\right\} \\ &= \frac{1}{\sqrt{|2\pi\boldsymbol{\Sigma}|}} \exp\left\{-\frac{1}{2}\boldsymbol{\mu}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}\right\} \exp\left\{\left[-\frac{1}{2} \text{Vec}(\boldsymbol{\Sigma}^{-1}), \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}\right]^T [\text{Vec}(\mathbf{x}\mathbf{x}^T), \mathbf{x}]\right\} \end{aligned}$$

Please note that $\text{Tr}(\mathbf{A}\mathbf{B}) = \text{vec}(\mathbf{A})^T \text{vec}(\mathbf{B})$, which is proven by the expanding the equation.

Proof. Now, we expand the normal distribution:

$$\begin{aligned} & \frac{1}{\sqrt{|2\pi\boldsymbol{\Sigma}|}} \exp\left\{-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x}-\boldsymbol{\mu})\right\} \\ &= \frac{1}{\sqrt{|2\pi\boldsymbol{\Sigma}|}} \exp\left\{-\frac{1}{2}(\mathbf{x}^T \boldsymbol{\Sigma}^{-1} - \boldsymbol{\mu}^T \boldsymbol{\Sigma}^{-1})(\mathbf{x}-\boldsymbol{\mu})\right\} \\ &= \frac{1}{\sqrt{|2\pi\boldsymbol{\Sigma}|}} \exp\left\{-\frac{1}{2}(\mathbf{x}^T \boldsymbol{\Sigma}^{-1} \mathbf{x} - 2\boldsymbol{\mu}^T \boldsymbol{\Sigma}^{-1} \mathbf{x} + \boldsymbol{\mu}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu})\right\} \\ &= \frac{1}{\sqrt{|2\pi\boldsymbol{\Sigma}|}} \exp\left\{-\frac{1}{2}\boldsymbol{\mu}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}\right\} \exp\left\{-\frac{1}{2}\mathbf{x}^T \boldsymbol{\Sigma}^{-1} \mathbf{x} + \boldsymbol{\mu}^T \boldsymbol{\Sigma}^{-1} \mathbf{x}\right\} \end{aligned}$$

Now, we consider the quadratic $\mathbf{x}^T \Sigma^{-1} \mathbf{x}$, as we have:

$$\text{Tr}(\mathbf{x}^T \Sigma^{-1} \mathbf{x}) = \text{Tr}(\Sigma^{-1} \mathbf{x} \mathbf{x}^T) = \text{Vec}(\Sigma^{-1})^T \text{Vec}(\mathbf{x} \mathbf{x}^T)$$

And, so we have the following:

$$\begin{aligned} & \frac{1}{\sqrt{|2\pi\Sigma|}} \exp\left\{-\frac{1}{2}\boldsymbol{\mu}^T \Sigma^{-1} \boldsymbol{\mu}\right\} \exp\left\{-\frac{1}{2}\mathbf{x}^T \Sigma^{-1} \mathbf{x} + \boldsymbol{\mu}^T \Sigma^{-1} \mathbf{x}\right\} \\ &= \frac{1}{\sqrt{|2\pi\Sigma|}} \exp\left\{-\frac{1}{2}\boldsymbol{\mu}^T \Sigma^{-1} \boldsymbol{\mu}\right\} \exp\left\{-\frac{1}{2}\text{Vec}(\Sigma^{-1})^T \text{Vec}(\mathbf{x} \mathbf{x}^T) + \boldsymbol{\mu}^T \Sigma^{-1} \mathbf{x}\right\} \end{aligned}$$

and so we have proven the proposition. \square

Remark 13. Now, we can consider the iid observations $\{\mathbf{x}_i\}_{i=1}^N$ of exponential family and we have:

$$\prod_{i=1}^N \left[f(\mathbf{x}_i) g(\boldsymbol{\theta}) \exp\left(\boldsymbol{\phi}(\boldsymbol{\theta})^T \mathbf{T}(\mathbf{x}_i)\right) \right] = g(\boldsymbol{\theta})^N \left(\prod_{i=1}^N f(\mathbf{x}_i) \right) \exp\left(\sum_{i=1}^N \boldsymbol{\phi}(\boldsymbol{\theta})^T \mathbf{T}(\mathbf{x}_i) \right)$$

Definition 1.5. (Conjugate Prior) The conjugate prior of the exponential family is the probability distribution of the form of:

$$p(\boldsymbol{\theta}|\boldsymbol{\tau}, \nu) = f(\boldsymbol{\tau}, \nu) g(\boldsymbol{\theta})^\nu \exp\left(\boldsymbol{\phi}(\boldsymbol{\theta})^T \boldsymbol{\tau}\right)$$

where $\nu > 0$ and $\boldsymbol{\tau} \in \mathbb{R}^m$. It is designed so that the posterior given this prior will have the same distribution as the conjugate prior.

Remark 14. The conjugate prior allow use to find the posterior with ease as we don't have to find the normalization of the Bayes' theorem:

$$\begin{aligned} p(\boldsymbol{\theta}|\mathcal{D}) &\propto p(\mathcal{D}|\boldsymbol{\theta}) p(\boldsymbol{\theta}|\boldsymbol{\tau}, \nu) \propto g(\boldsymbol{\theta})^{N+\nu} \exp\left(\boldsymbol{\phi}(\boldsymbol{\theta})^T \left(\boldsymbol{\tau} + \sum_{i=1}^N \mathbf{T}(\mathbf{x}_i)\right)\right) \\ &= F\left(\boldsymbol{\tau} + \sum_{i=1}^N \mathbf{T}(\mathbf{x}_i), N + \nu\right) g(\boldsymbol{\theta})^{N+\nu} \exp\left(\boldsymbol{\phi}(\boldsymbol{\theta})^T \left(\boldsymbol{\tau} + \sum_{i=1}^N \mathbf{T}(\mathbf{x}_i)\right)\right) \end{aligned}$$

And so we have $\boldsymbol{\tau}$ quantify the pseudo-observations via the sufficient statistics. Furthermore, ν is the pseudo-count of the pseudo-observation (can also be seen as the weight of prior belief).

Remark 15. Now, we have the following list of conjugate prior, where it is shown in the table below:

Distribution	Conjugate	Distribution (2)	Conjugate (2)
Bernoulli	Beta	Multinormal (unknown $\boldsymbol{\mu}$)	Multinormal
Poisson	Gamma	Multinormal (unknown $\boldsymbol{\Lambda} = \Sigma^{-1}$)	Wishart
Multinormal	Dirichlet	Multinormal (unknown $\boldsymbol{\Lambda}$)	Normal-Wishart
Normal (unknown $\boldsymbol{\mu}$)	Normal		
Normal (unknown $\tau = \sigma^{-1}$)	Gamma		
Normal (unknown $\boldsymbol{\mu}, \tau$)	Normal-Gamma		

Definition 1.6. (Normal-Gamma Distribution) The Normal-Gamma is defined as:

$$p(x, \tau|\boldsymbol{\mu}, \lambda, \alpha, \beta) = \frac{\beta^\alpha \sqrt{\lambda}}{\Gamma(\alpha) \sqrt{2\pi}} \tau^{\alpha-1/2} \exp(-\beta\tau) \exp\left(-\frac{\lambda\tau(x - \boldsymbol{\mu})^2}{2}\right)$$

similarly the Normal-Wishart distribution is given by $\mathcal{N}(\boldsymbol{\mu}|\boldsymbol{\mu}_0, (\lambda\boldsymbol{\Lambda})^{-1})\mathcal{W}(\boldsymbol{\Lambda}|\mathbf{W}, \nu)$, where \mathcal{W} is the Wishart distribution, where we consider $(\boldsymbol{\mu}, \boldsymbol{\Sigma}) \sim \text{NW}(\boldsymbol{\mu}_0, \lambda, \mathbf{W}, \nu)$, given as:

$$\begin{aligned} & \mathcal{N}(\boldsymbol{\mu}|\boldsymbol{\mu}_0, (\lambda\boldsymbol{\Lambda})^{-1})\mathcal{W}(\boldsymbol{\Lambda}|\mathbf{W}, \nu) \\ &= B(\mathbf{W}, \nu) |\boldsymbol{\Lambda}|^{(\nu-D-1)/2} \exp\left(-\frac{1}{2}\text{Tr}(\mathbf{W}^{-1}\boldsymbol{\Lambda})\right) \frac{1}{\sqrt{|2\pi(\lambda\boldsymbol{\Lambda})^{-1}|}} \exp\left\{-\frac{\lambda}{2}(\boldsymbol{\mu} - \boldsymbol{\mu}_0)^T \boldsymbol{\Lambda}(\boldsymbol{\mu} - \boldsymbol{\mu}_0)\right\} \\ &= \frac{B(\mathbf{W}, \nu) |\boldsymbol{\Lambda}|^{(\nu-D-1)/2}}{\sqrt{|2\pi(\lambda\boldsymbol{\Lambda})^{-1}|}} \exp\left(-\frac{1}{2}\text{Tr}[\mathbf{W}^{-1}\boldsymbol{\Lambda}] - \frac{\lambda}{2}(\boldsymbol{\mu} - \boldsymbol{\mu}_0)^T \boldsymbol{\Lambda}(\boldsymbol{\mu} - \boldsymbol{\mu}_0)^T\right) \end{aligned}$$

1.6 Everything You Always Wanted to Know About Gaussian But Were Afraid to Ask..

Proposition 1.10. (Gaussian Integration) We can show that:

$$\int_{-\infty}^{\infty} \exp(-x^2) dx = \sqrt{\pi}$$

Remark 16. (Shape of Gaussian) Recalling the definition of multivariate Gaussian distribution:

$$\frac{1}{\sqrt{|2\pi\boldsymbol{\Sigma}|}} \exp\left\{-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right\}$$

where we will define the Mahalanobis distance of relating to the Gaussian to be $\Delta^2 = (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})$. Now, it is clear that the shape of the Gaussian depends on this value. Now, using the eigendecomposition on the covariance function, we have:

$$\boldsymbol{\Sigma}^{-1} = \sum_{i=1}^D \frac{1}{\lambda_i} \mathbf{u}_i \mathbf{u}_i^T \implies \Delta^2 = \sum_{i=1}^D \frac{y_i^2}{\lambda_i}$$

where $(\lambda_i, \mathbf{u}_i)$ is the eigenvalue/eigenvector pair of $\boldsymbol{\Sigma}$, and $y_i = \mathbf{u}_i^T(\mathbf{x} - \boldsymbol{\mu})$. Now for the full vector \mathbf{y} is equal to $U(\mathbf{x} - \boldsymbol{\mu})$. Therefore, the shape of Gaussian is characterized as:

- Ellipsoids with the center at $\boldsymbol{\mu}$
- Axis is in the direction of eigenvector \mathbf{u}_i .
- Scaling of each direction is the eigenvector λ_i associated with \mathbf{u}_i

Remark 17. (Normalizing Factor of Gaussian) To show that the multivariate Gaussian indeed is normalized, we consider a change of coordinate to the eigen-basis consider above. To do this, we consider Jacobina matrix \mathbf{J} as:

$$J_{ij} = \frac{\partial x_i}{\partial y_j} = U_{ji}$$

Since \mathbf{U} is orthonormal, we have: $|\mathbf{J}|^2 = |\mathbf{U}^T|^2 = |\mathbf{U}^T| |\mathbf{U}| = |\mathbf{U}^T \mathbf{U}| = |\mathbf{I}| = 1$. Now, using the definition of determinant: $|\boldsymbol{\Sigma}|^{1/2} = \prod_{i=1}^D |\lambda_j|^{1/2}$. And, so one can perform the integration over as:

$$\int p(\mathbf{x}) d\mathbf{x} = \int p(\mathbf{y}) |\mathbf{J}| d\mathbf{y} = \int \prod_{i=1}^D \frac{1}{\sqrt{2\pi\lambda_i}} \exp\left\{-\frac{y_i^2}{2\lambda_i}\right\} d\mathbf{y} = \prod_{i=1}^D \int \frac{1}{\sqrt{2\pi\lambda_i}} \exp\left\{-\frac{y_i^2}{2\lambda_i}\right\} dy_i = 1$$

The final equality integration follows from the Gaussian integration.

Proposition 1.11. Consider the following multivariate Gaussian distribution:

$$\begin{pmatrix} \mathbf{x}_a \\ \mathbf{x}_b \end{pmatrix} \sim \mathcal{N} \left(\begin{pmatrix} \boldsymbol{\mu}_a \\ \boldsymbol{\mu}_b \end{pmatrix}, \begin{pmatrix} \boldsymbol{\Sigma}_{aa} & \boldsymbol{\Sigma}_{ab} \\ \boldsymbol{\Sigma}_{ba} & \boldsymbol{\Sigma}_{bb} \end{pmatrix} \right)$$

where $\boldsymbol{\Lambda}^{-1} = \boldsymbol{\Sigma}$, then we can show that:

$$p(\mathbf{x}_a | \mathbf{x}_b) = \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_{a|b}, \boldsymbol{\Lambda}_{aa}^{-1}) \quad \text{where} \quad \boldsymbol{\mu}_{a|b} = \boldsymbol{\mu}_a - \boldsymbol{\Lambda}_{aa}^{-1} \boldsymbol{\Lambda}_{ab} (\mathbf{x}_b - \boldsymbol{\mu}_b)$$

One can simply the equation by consider the value $\mathbf{K} = \boldsymbol{\Sigma}_{ab} \boldsymbol{\Sigma}_{bb}^{-1}$ and since $\boldsymbol{\Sigma}_{ab} = \boldsymbol{\Sigma}_{ba}$, then we have:

$$\begin{aligned} \boldsymbol{\mu}_{a|b} &= \boldsymbol{\mu}_a + \mathbf{K}(\mathbf{x}_b - \boldsymbol{\mu}_b) \\ \boldsymbol{\Sigma}_{a|b} &= \boldsymbol{\Sigma}_{aa} - \mathbf{K} \boldsymbol{\Sigma}_{bb} \mathbf{K}^T \\ &= \boldsymbol{\Sigma}_{aa} - \boldsymbol{\Sigma}_{ab} \boldsymbol{\Sigma}_{bb}^{-1} \boldsymbol{\Sigma}_{ba} \end{aligned}$$

The above equation follows from the block-matrix inverse.

Proof. We will consider only the quadratic term inside Gaussian, as we have:

$$\begin{aligned} -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}) &= \\ -\frac{1}{2}(\mathbf{x}_a - \boldsymbol{\mu}_a)^T \boldsymbol{\Lambda}_{aa}(\mathbf{x}_a - \boldsymbol{\mu}_a) &- \frac{1}{2}(\mathbf{x}_a - \boldsymbol{\mu}_a)^T \boldsymbol{\Lambda}_{ab}(\mathbf{x}_b - \boldsymbol{\mu}_b) \\ -\frac{1}{2}(\mathbf{x}_b - \boldsymbol{\mu}_b)^T \boldsymbol{\Lambda}_{ba}(\mathbf{x}_a - \boldsymbol{\mu}_a) &- \frac{1}{2}(\mathbf{x}_b - \boldsymbol{\mu}_b)^T \boldsymbol{\Lambda}_{bb}(\mathbf{x}_b - \boldsymbol{\mu}_b) \end{aligned}$$

Now, we are interested to find the condition distribution as it will have the form of:

$$-\frac{1}{2}(\mathbf{x}_a - \boldsymbol{\mu}_{a|b})^T \boldsymbol{\Sigma}_{a|b}^{-1}(\mathbf{x}_a - \boldsymbol{\mu}_{a|b}) = \underbrace{-\frac{1}{2} \mathbf{x}_a^T \boldsymbol{\Sigma}_{a|b}^{-1} \mathbf{x}_a}_{\textcircled{1}} + \underbrace{\mathbf{x}_a^T \boldsymbol{\Sigma}_{a|b}^{-1} \boldsymbol{\mu}_{a|b}}_{\textcircled{2}} + \text{const}$$

Now, let's consider the term for full Gaussian that have the form that matches the condition distribution:

1. The first one is simple as we have: $-\frac{1}{2} \mathbf{x}_a^T \boldsymbol{\Lambda}_{aa} \mathbf{x}_a$, and so, one can conclude that $\boldsymbol{\Sigma}_{a|b}^{-1} = \boldsymbol{\Lambda}_{aa}$.
2. For the second term, we consider equation with $\mathbf{x}_a^T(\dots)$, as we have (from term 2 and 3):

$$\begin{aligned} \mathbf{x}_a^T \boldsymbol{\Lambda}_{aa} \boldsymbol{\mu}_a - \frac{1}{2} \mathbf{x}_a^T \boldsymbol{\Lambda}_{ab} \mathbf{x}_b + \frac{1}{2} \mathbf{x}_a^T \boldsymbol{\Lambda}_{ab} \boldsymbol{\mu}_b + \frac{1}{2} \boldsymbol{\mu}_b^T \boldsymbol{\Lambda}_{ba} \mathbf{x}_a - \frac{1}{2} \boldsymbol{\mu}_b^T \boldsymbol{\Lambda}_{ba} \boldsymbol{\mu}_a \\ = \mathbf{x}_a^T \left[\boldsymbol{\Lambda}_{aa} \boldsymbol{\mu}_a + \boldsymbol{\Lambda}_{ab} (\mathbf{x}_b - \boldsymbol{\mu}_b) \right] \end{aligned}$$

For the second equality, we use the fact that $\boldsymbol{\Lambda}_{ba}^T = \boldsymbol{\Lambda}_{ab}$. Now we simply apply the inverse of $\boldsymbol{\Lambda}_{aa}$ to get the mean, which means:

$$\begin{aligned} \mathbf{x}_a^T \boldsymbol{\Lambda}_{aa} \boldsymbol{\mu}_{a|b} &= \mathbf{x}_a^T \left[\boldsymbol{\Lambda}_{aa} \boldsymbol{\mu}_a + \boldsymbol{\Lambda}_{ab} (\mathbf{x}_b - \boldsymbol{\mu}_b) \right] \\ \implies \boldsymbol{\mu}_{a|b} &= \boldsymbol{\mu}_a + \boldsymbol{\Lambda}_{aa}^{-1} \boldsymbol{\Lambda}_{ab} (\mathbf{x}_b - \boldsymbol{\mu}_b) \end{aligned}$$

as required.

Please note that one can use the block inverse above to calculate $\boldsymbol{\Lambda}_{aa}$ and $\boldsymbol{\Lambda}_{ab}$ in terms of $\boldsymbol{\Sigma}$. □

Proposition 1.12. Consider the following multivariate Gaussian distribution:

$$\begin{pmatrix} \mathbf{x}_a \\ \mathbf{x}_b \end{pmatrix} \sim \mathcal{N} \left(\begin{pmatrix} \boldsymbol{\mu}_a \\ \boldsymbol{\mu}_b \end{pmatrix}, \begin{pmatrix} \boldsymbol{\Sigma}_{aa} & \boldsymbol{\Sigma}_{ab} \\ \boldsymbol{\Sigma}_{ba} & \boldsymbol{\Sigma}_{bb} \end{pmatrix} \right)$$

where $\boldsymbol{\Lambda}^{-1} = \boldsymbol{\Sigma}$, then we can show that: $p(\mathbf{x}_a) = \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_a, \boldsymbol{\Sigma}_{aa})$

Proof. We will use the full expansion of Gaussian like above proof. However, we will consider the term with \mathbf{x}_b , first as we have (similar to the conditional case):

$$\begin{aligned}
& -\frac{1}{2}\mathbf{x}_b^T \mathbf{\Lambda}_{bb} \mathbf{x}_b + \mathbf{x}_b^T \mathbf{\Lambda}_{bb} \boldsymbol{\mu}_b - \frac{1}{2}\mathbf{x}_b^T \mathbf{\Lambda}_{ba} \mathbf{x}_a + \frac{1}{2}\mathbf{x}_b^T \mathbf{\Lambda}_{ba} \boldsymbol{\mu}_a - \frac{1}{2}\mathbf{x}_a^T \mathbf{\Lambda}_{ab} \mathbf{x}_b + \frac{1}{2}\boldsymbol{\mu}_a^T \mathbf{\Lambda}_{ab} \mathbf{x}_b \\
& = -\frac{1}{2}\mathbf{x}_b^T \mathbf{\Lambda}_{bb} \mathbf{x}_b + \mathbf{x}_b^T \mathbf{\Lambda}_{bb} \boldsymbol{\mu}_b - \mathbf{x}_b^T \mathbf{\Lambda}_{ba} \mathbf{x}_a + \mathbf{x}_b^T \mathbf{\Lambda}_{ba} \boldsymbol{\mu}_a \\
& = -\frac{1}{2} \left[\mathbf{x}_b^T \mathbf{\Lambda}_{bb} \mathbf{x}_b + 2\mathbf{x}_b^T \mathbf{\Lambda}_{bb} \mathbf{\Lambda}_{bb}^{-1} \underbrace{(\mathbf{\Lambda}_{bb} \boldsymbol{\mu}_b - \mathbf{\Lambda}_{ba} (\mathbf{x}_a - \boldsymbol{\mu}_a))}_{\mathbf{m}} \right] \\
& = -\frac{1}{2} \left[\mathbf{x}_b^T \mathbf{\Lambda}_{bb} \mathbf{x}_b - 2\mathbf{x}_b^T \mathbf{\Lambda}_{bb} \mathbf{\Lambda}_{bb}^{-1} \mathbf{m} + (\mathbf{\Lambda}_{bb}^{-1} \mathbf{m})^T \mathbf{\Lambda}_{bb} (\mathbf{\Lambda}_{bb}^{-1} \mathbf{m}) \right] + \frac{1}{2} \mathbf{m}^T \mathbf{\Lambda}_{bb}^{-1} \mathbf{m} \\
& = -\frac{1}{2} (\mathbf{x}_b - \mathbf{\Lambda}_{bb}^{-1} \mathbf{m})^T \mathbf{\Lambda}_{bb} (\mathbf{x}_b - \mathbf{\Lambda}_{bb}^{-1} \mathbf{m}) + \frac{1}{2} \mathbf{m}^T \mathbf{\Lambda}_{bb}^{-1} \mathbf{m}
\end{aligned}$$

Now, we can integrate out the quantities (that depends on \mathbf{x}_b) i.e:

$$\int \exp \left\{ -\frac{1}{2} (\mathbf{x}_b - \mathbf{\Lambda}_{bb}^{-1} \mathbf{m})^T \mathbf{\Lambda}_{bb} (\mathbf{x}_b - \mathbf{\Lambda}_{bb}^{-1} \mathbf{m}) \right\} d\mathbf{x}_b$$

Since it is a Gaussian integration, we didn't have to perform any thing further as this would yields similar value for normalization factor. Now, consider the vales related to \mathbf{x}_a (without \mathbf{x}_b) and the leftout value from above:

$$\begin{aligned}
& \frac{1}{2} \mathbf{m}^T \mathbf{\Lambda}_{bb}^{-1} \mathbf{m} - \frac{1}{2} \mathbf{x}_a^T \mathbf{\Lambda}_{aa} \mathbf{x}_a + \mathbf{x}_a^T \mathbf{\Lambda}_{aa} \boldsymbol{\mu}_a + \frac{1}{2} \mathbf{x}_a^T \mathbf{\Lambda}_{ab} \boldsymbol{\mu}_b + \frac{1}{2} \boldsymbol{\mu}_b^T \mathbf{\Lambda}_{ba} \mathbf{x}_a \\
& = \frac{1}{2} \left[\mathbf{\Lambda}_{bb} \boldsymbol{\mu}_b - \mathbf{\Lambda}_{ba} (\mathbf{x}_a - \boldsymbol{\mu}_a) \right]^T \mathbf{\Lambda}_{bb}^{-1} \left[\mathbf{\Lambda}_{bb} \boldsymbol{\mu}_b - \mathbf{\Lambda}_{ba} (\mathbf{x}_a - \boldsymbol{\mu}_a) \right] \\
& \quad - \frac{1}{2} \mathbf{x}_a^T \mathbf{\Lambda}_{aa} \mathbf{x}_a + \mathbf{x}_a^T \mathbf{\Lambda}_{aa} \boldsymbol{\mu}_a + \mathbf{x}_a^T \mathbf{\Lambda}_{ab} \boldsymbol{\mu}_b \\
& = \frac{1}{2} \left[\boldsymbol{\mu}_b^T \mathbf{\Lambda}_{bb} \boldsymbol{\mu}_b - (\mathbf{x}_a - \boldsymbol{\mu}_a)^T \mathbf{\Lambda}_{ba}^T \boldsymbol{\mu}_b \right. \\
& \quad \left. - \boldsymbol{\mu}_b^T \mathbf{\Lambda}_{ba} (\mathbf{x}_a - \boldsymbol{\mu}_a) + (\mathbf{x}_a - \boldsymbol{\mu}_a)^T \mathbf{\Lambda}_{ba}^T \mathbf{\Lambda}_{bb}^{-1} \mathbf{\Lambda}_{ba} (\mathbf{x}_a - \boldsymbol{\mu}_a) \right] \\
& \quad - \frac{1}{2} \mathbf{x}_a^T \mathbf{\Lambda}_{aa} \mathbf{x}_a + \mathbf{x}_a^T \mathbf{\Lambda}_{aa} \boldsymbol{\mu}_a + \mathbf{x}_a^T \mathbf{\Lambda}_{ab} \boldsymbol{\mu}_b \\
& = \frac{1}{2} \left[\boldsymbol{\mu}_b^T \mathbf{\Lambda}_{bb} \boldsymbol{\mu}_b - \mathbf{x}_a^T \mathbf{\Lambda}_{ba}^T \boldsymbol{\mu}_b + \boldsymbol{\mu}_a^T \mathbf{\Lambda}_{ba}^T \boldsymbol{\mu}_b - \boldsymbol{\mu}_b^T \mathbf{\Lambda}_{ba} \mathbf{x}_a + \boldsymbol{\mu}_b^T \mathbf{\Lambda}_{ba} \boldsymbol{\mu}_a \right. \\
& \quad \left. + \mathbf{x}_a^T \mathbf{\Lambda}_{ba}^T \mathbf{\Lambda}_{bb}^{-1} \mathbf{\Lambda}_{ba} \mathbf{x}_a - \boldsymbol{\mu}_a^T \mathbf{\Lambda}_{ba}^T \mathbf{\Lambda}_{bb}^{-1} \mathbf{\Lambda}_{ba} \mathbf{x}_a - \mathbf{x}_a^T \mathbf{\Lambda}_{ba}^T \mathbf{\Lambda}_{bb}^{-1} \mathbf{\Lambda}_{ba} \boldsymbol{\mu}_a + \boldsymbol{\mu}_a^T \mathbf{\Lambda}_{ba}^T \mathbf{\Lambda}_{bb}^{-1} \mathbf{\Lambda}_{ba} \boldsymbol{\mu}_a \right] \\
& \quad - \frac{1}{2} \mathbf{x}_a^T \mathbf{\Lambda}_{aa} \mathbf{x}_a + \mathbf{x}_a^T \mathbf{\Lambda}_{aa} \boldsymbol{\mu}_a + \mathbf{x}_a^T \mathbf{\Lambda}_{ab} \boldsymbol{\mu}_b \\
& = \frac{1}{2} \left[-2\mathbf{x}_a^T \mathbf{\Lambda}_{ba}^T \boldsymbol{\mu}_b + \mathbf{x}_a^T \mathbf{\Lambda}_{ba}^T \mathbf{\Lambda}_{bb}^{-1} \mathbf{\Lambda}_{ba} \mathbf{x}_a - 2\mathbf{x}_a^T \mathbf{\Lambda}_{ba}^T \mathbf{\Lambda}_{bb}^{-1} \mathbf{\Lambda}_{ba} \boldsymbol{\mu}_a \right] \\
& \quad - \frac{1}{2} \mathbf{x}_a^T \mathbf{\Lambda}_{aa} \mathbf{x}_a + \mathbf{x}_a^T \mathbf{\Lambda}_{aa} \boldsymbol{\mu}_a + \mathbf{x}_a^T \mathbf{\Lambda}_{ab} \boldsymbol{\mu}_b + \text{const} \\
& = \frac{1}{2} \mathbf{x}_a^T \mathbf{\Lambda}_{ba}^T \mathbf{\Lambda}_{bb}^{-1} \mathbf{\Lambda}_{ba} \mathbf{x}_a - \mathbf{x}_a^T \mathbf{\Lambda}_{ba}^T \mathbf{\Lambda}_{bb}^{-1} \mathbf{\Lambda}_{ba} \boldsymbol{\mu}_a - \frac{1}{2} \mathbf{x}_a^T \mathbf{\Lambda}_{aa} \mathbf{x}_a + \mathbf{x}_a^T \mathbf{\Lambda}_{aa} \boldsymbol{\mu}_a + \text{const} \\
& = -\frac{1}{2} \mathbf{x}_a^T \left[\mathbf{\Lambda}_{aa} - \mathbf{\Lambda}_{ba}^T \mathbf{\Lambda}_{bb}^{-1} \mathbf{\Lambda}_{ba} \right] \mathbf{x}_a + \mathbf{x}_a^T \left[\mathbf{\Lambda}_{aa} - \mathbf{\Lambda}_{ba}^T \mathbf{\Lambda}_{bb}^{-1} \mathbf{\Lambda}_{ba} \right] \boldsymbol{\mu}_a + \text{const}
\end{aligned}$$

Now, we can compare this to the form:

$$-\frac{1}{2} (\mathbf{x}_a - \boldsymbol{\mu}_a^*)^T \boldsymbol{\Sigma}_a^* (\mathbf{x}_a - \boldsymbol{\mu}_a^*) = -\frac{1}{2} \mathbf{x}_a^T \boldsymbol{\Sigma}_a^* \mathbf{x}_a + \mathbf{x}_a^T \boldsymbol{\Sigma}_a^* \boldsymbol{\mu}_a^* + \text{const}$$

and we have $\boldsymbol{\Sigma}_a^*$ (marginalized) is equal to $(\mathbf{\Lambda}_{aa} - \mathbf{\Lambda}_{ba}^T \mathbf{\Lambda}_{bb}^{-1} \mathbf{\Lambda}_{ba})^{-1}$, and so we have $\boldsymbol{\mu}_a^* = \boldsymbol{\mu}_a$. Furthermore, from the partition inverse of $(\mathbf{\Lambda}_{aa} - \mathbf{\Lambda}_{ba}^T \mathbf{\Lambda}_{bb}^{-1} \mathbf{\Lambda}_{ba})^{-1} = \boldsymbol{\Sigma}_{aa}$. Thus, we completed the proof. \square

Proposition 1.13. (Linear Gaussian Model) Consider the following distributions:

$$p(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Lambda}^{-1}) \quad p(\mathbf{y}|\mathbf{x}) = \mathcal{N}(\mathbf{y}|\mathbf{A}\mathbf{x} + \mathbf{b}, \mathbf{L}^{-1})$$

Now, we can show that the marginal distribution and the conditional distribution of \mathbf{x} given \mathbf{y} is given by:

$$p(\mathbf{y}) = \mathcal{N}(\mathbf{y}|\mathbf{A}\boldsymbol{\mu} + \mathbf{b}, \mathbf{L}^{-1} + \mathbf{A}\boldsymbol{\Lambda}^{-1}\mathbf{A}^T) \quad p(\mathbf{x}|\mathbf{y}) = \mathcal{N}\left(\mathbf{x} \mid \boldsymbol{\Sigma} \left\{ \mathbf{A}^T \mathbf{L}(\mathbf{y} - \mathbf{b}) + \boldsymbol{\Lambda}\boldsymbol{\mu} \right\}, \boldsymbol{\Sigma}\right)$$

where $\boldsymbol{\Sigma} = (\boldsymbol{\Lambda} + \mathbf{A}^T \mathbf{L} \mathbf{A})^{-1}$

Proof. Let's consider the joint distribution first, where we denote $\mathbf{z} = \begin{pmatrix} \mathbf{x} \\ \mathbf{y} \end{pmatrix}$, and consider the inside of exponential as we have:

$$\begin{aligned} & -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Lambda}(\mathbf{x} - \boldsymbol{\mu}) - \frac{1}{2}(\mathbf{y} - \mathbf{A}\mathbf{x} - \mathbf{b})^T \mathbf{L}(\mathbf{y} - \mathbf{A}\mathbf{x} - \mathbf{b}) + \text{const} \\ &= -\frac{1}{2} \left[\mathbf{x}^T \boldsymbol{\Lambda} \mathbf{x} - 2\mathbf{x}^T \boldsymbol{\Lambda} \boldsymbol{\mu} + \boldsymbol{\mu}^T \boldsymbol{\Lambda} \boldsymbol{\mu} \right] \\ & \quad - \frac{1}{2} \left[\mathbf{y}^T \mathbf{L} \mathbf{y} - \mathbf{x}^T \mathbf{A}^T \mathbf{L} \mathbf{y} - \mathbf{b}^T \mathbf{L} \mathbf{y} - \mathbf{y}^T \mathbf{L} \mathbf{A} \mathbf{x} + \mathbf{x}^T \mathbf{A}^T \mathbf{L} \mathbf{A} \mathbf{x} + \mathbf{b}^T \mathbf{L} \mathbf{A} \mathbf{x} \right. \\ & \quad \left. - \mathbf{y}^T \mathbf{L} \mathbf{b} + \mathbf{x}^T \mathbf{A}^T \mathbf{L} \mathbf{b} + \mathbf{b}^T \mathbf{L} \mathbf{b} \right] + \text{const} \\ &= -\frac{1}{2} \left[\mathbf{x}^T \boldsymbol{\Lambda} \mathbf{x} - 2\mathbf{x}^T \boldsymbol{\Lambda} \boldsymbol{\mu} + \mathbf{y}^T \mathbf{L} \mathbf{y} - \mathbf{x}^T \mathbf{A}^T \mathbf{L} \mathbf{y} - \mathbf{b}^T \mathbf{L} \mathbf{y} \right. \\ & \quad \left. - \mathbf{y}^T \mathbf{L} \mathbf{A} \mathbf{x} + \mathbf{x}^T \mathbf{A}^T \mathbf{L} \mathbf{A} \mathbf{x} + \mathbf{b}^T \mathbf{L} \mathbf{A} \mathbf{x} - \mathbf{y}^T \mathbf{L} \mathbf{b} + \mathbf{x}^T \mathbf{A}^T \mathbf{L} \mathbf{b} \right] + \text{const} \\ &= -\frac{1}{2} \left[\mathbf{x}^T (\boldsymbol{\Lambda} + \mathbf{A}^T \mathbf{L} \mathbf{A}) \mathbf{x} + \mathbf{y}^T \mathbf{L} \mathbf{y} - \mathbf{x}^T \mathbf{A}^T \mathbf{L} \mathbf{y} - \mathbf{y}^T \mathbf{L} \mathbf{A} \mathbf{x} \right. \\ & \quad \left. - 2\mathbf{x}^T \boldsymbol{\Lambda} \boldsymbol{\mu} - 2\mathbf{b}^T \mathbf{L} \mathbf{y} + 2\mathbf{b}^T \mathbf{L} \mathbf{A} \mathbf{x} \right] + \text{const} \\ &= -\frac{1}{2} \mathbf{x}^T \left(\boldsymbol{\Lambda} - \frac{1}{2} \mathbf{A}^T \mathbf{L} \mathbf{A} \right) \mathbf{x} - \frac{1}{2} \mathbf{y}^T \mathbf{L} \mathbf{y} + \frac{1}{2} \mathbf{x}^T \mathbf{A}^T \mathbf{L} \mathbf{y} + \frac{1}{2} \mathbf{y}^T \mathbf{L} \mathbf{A} \mathbf{x} \\ & \quad + \mathbf{x}^T \boldsymbol{\Lambda} \boldsymbol{\mu} + \mathbf{b}^T \mathbf{L} \mathbf{y} - \mathbf{b}^T \mathbf{L} \mathbf{A} \mathbf{x} + \text{const} \\ &= -\frac{1}{2} \begin{pmatrix} \mathbf{x} \\ \mathbf{y} \end{pmatrix}^T \begin{pmatrix} \boldsymbol{\Lambda} + \mathbf{A}^T \mathbf{L} \mathbf{A} & -\mathbf{A}^T \mathbf{L} \\ -\mathbf{L} \mathbf{A} & \mathbf{L} \end{pmatrix} \begin{pmatrix} \mathbf{x} \\ \mathbf{y} \end{pmatrix} + \begin{pmatrix} \mathbf{x} \\ \mathbf{y} \end{pmatrix}^T \begin{pmatrix} \boldsymbol{\Lambda} \boldsymbol{\mu} - \mathbf{A}^T \mathbf{L} \mathbf{b} \\ \mathbf{L} \mathbf{b} \end{pmatrix} + \text{const} \end{aligned}$$

Now, using the same pattern matching, we can see that the covariance of \mathbf{z} is equal to:

$$\begin{pmatrix} \boldsymbol{\Lambda} + \mathbf{A}^T \mathbf{L} \mathbf{A} & -\mathbf{A}^T \mathbf{L} \\ -\mathbf{L} \mathbf{A} & \mathbf{L} \end{pmatrix}^{-1} = \begin{pmatrix} \boldsymbol{\Lambda}^{-1} & \boldsymbol{\Lambda}^{-1} \mathbf{A}^T \\ \mathbf{A} \boldsymbol{\Lambda}^{-1} & \mathbf{L}^{-1} + \mathbf{A} \boldsymbol{\Lambda}^{-1} \mathbf{A}^T \end{pmatrix}$$

Similarly, the mean is equal to:

$$\begin{pmatrix} \boldsymbol{\Lambda}^{-1} & \boldsymbol{\Lambda}^{-1} \mathbf{A}^T \\ \mathbf{A} \boldsymbol{\Lambda}^{-1} & \mathbf{L}^{-1} + \mathbf{A} \boldsymbol{\Lambda}^{-1} \mathbf{A}^T \end{pmatrix} \begin{pmatrix} \boldsymbol{\Lambda} \boldsymbol{\mu} - \mathbf{A}^T \mathbf{L} \mathbf{b} \\ \mathbf{L} \mathbf{b} \end{pmatrix} = \begin{pmatrix} \boldsymbol{\mu} \\ \mathbf{A} \boldsymbol{\mu} + \mathbf{b} \end{pmatrix}$$

Now, we can use results above (conditional and marginalized) to get the result. \square

Remark 18. (Conjugate Prior of Multinormal) The proof follows from [here](#). We are now consider the likelihood of Multinormal distribution given the dataset:

$$\begin{aligned} & \prod_{i=1}^N \frac{1}{\sqrt{|2\pi\boldsymbol{\Lambda}^{-1}|}} \exp \left\{ -\frac{1}{2}(\mathbf{x}_i - \boldsymbol{\mu})^T \boldsymbol{\Lambda}(\mathbf{x}_i - \boldsymbol{\mu}) \right\} \\ &= \frac{1}{|2\pi\boldsymbol{\Lambda}^{-1}|^{N/2}} \exp \left\{ -\frac{1}{2} \sum_{i=1}^N (\mathbf{x}_i - \boldsymbol{\mu})^T \boldsymbol{\Lambda}(\mathbf{x}_i - \boldsymbol{\mu}) \right\} \\ &\propto |\boldsymbol{\Lambda}|^{N/2} \exp \left\{ -\frac{1}{2} \sum_{i=1}^N (\mathbf{x}_i - \boldsymbol{\mu})^T \boldsymbol{\Lambda}(\mathbf{x}_i - \boldsymbol{\mu}) \right\} \end{aligned}$$

Now, we will consider the Normal-Wishart distribution in the similar form as:

$$\begin{aligned} \frac{B(\mathbf{W}, \nu) |\mathbf{\Lambda}|^{(\nu-D-1)/2}}{\sqrt{|2\pi(\lambda\mathbf{\Lambda})^{-1}|}} \exp\left(-\frac{1}{2} \text{Tr} \left[\mathbf{W}^{-1} \mathbf{\Lambda} \right] - \frac{\lambda}{2} (\boldsymbol{\mu} - \boldsymbol{\mu}_0)^T \mathbf{\Lambda} (\boldsymbol{\mu} - \boldsymbol{\mu}_0)^T\right) \\ \propto |\mathbf{\Lambda}|^{(\nu-D)/2} \exp\left(-\frac{1}{2} \text{Tr} \left[\mathbf{W}^{-1} \mathbf{\Lambda} \right]\right) \exp\left(-\frac{\lambda}{2} (\boldsymbol{\mu} - \boldsymbol{\mu}_0)^T \mathbf{\Lambda} (\boldsymbol{\mu} - \boldsymbol{\mu}_0)^T\right) \end{aligned}$$

Now, we simply have to multiply the conjugate prior and the normal distribution, which gives us:

$$\begin{aligned} & |\mathbf{\Lambda}|^{N/2} \exp\left\{-\frac{1}{2} \sum_{i=1}^N (\mathbf{x}_i - \boldsymbol{\mu})^T \mathbf{\Lambda} (\mathbf{x}_i - \boldsymbol{\mu})\right\} \\ & \quad |\mathbf{\Lambda}|^{(\nu-D)/2} \exp\left(-\frac{1}{2} \text{Tr} \left[\mathbf{W}^{-1} \mathbf{\Lambda} \right]\right) \exp\left(-\frac{\lambda}{2} (\boldsymbol{\mu} - \boldsymbol{\mu}_0)^T \mathbf{\Lambda} (\boldsymbol{\mu} - \boldsymbol{\mu}_0)^T\right) \\ & = |\mathbf{\Lambda}|^{(\nu-D+N)/2} \exp\left(-\frac{1}{2} \text{Tr} \left[\mathbf{W}^{-1} \mathbf{\Lambda} \right]\right) \\ & \quad \exp\left\{-\frac{1}{2} \left[\sum_{i=1}^N (\mathbf{x}_i^T \mathbf{\Lambda} \mathbf{x}_i) + N \boldsymbol{\mu}^T \mathbf{\Lambda} \boldsymbol{\mu} - 2N \bar{\mathbf{x}}^T \mathbf{\Lambda} \boldsymbol{\mu} \right]\right\} \\ & \quad \exp\left\{-\frac{\lambda}{2} \left[\boldsymbol{\mu}^T \mathbf{\Lambda} \boldsymbol{\mu} - 2\boldsymbol{\mu}_0^T \mathbf{\Lambda} \boldsymbol{\mu} + \boldsymbol{\mu}_0^T \mathbf{\Lambda} \boldsymbol{\mu}_0 \right]\right\} \\ & = |\mathbf{\Lambda}|^{(\nu-D+N)/2} \exp\left(-\frac{1}{2} \left[\text{Tr} \left[\mathbf{W}^{-1} \mathbf{\Lambda} \right] + \sum_{i=1}^N (\mathbf{x}_i^T \mathbf{\Lambda} \mathbf{x}_i) + (N + \lambda) \boldsymbol{\mu}^T \mathbf{\Lambda} \boldsymbol{\mu} \right. \right. \\ & \quad \left. \left. - 2(N \bar{\mathbf{x}}^T + \lambda \boldsymbol{\mu}_0^T) \mathbf{\Lambda} \boldsymbol{\mu} + \lambda \boldsymbol{\mu}_0^T \mathbf{\Lambda} \boldsymbol{\mu}_0 \right]\right) \end{aligned}$$

Let's consider exclusively for the expression in the square bracket:

$$\begin{aligned} & \text{Tr} \left[\mathbf{W}^{-1} \mathbf{\Lambda} \right] + \sum_{i=1}^N (\mathbf{x}_i^T \mathbf{\Lambda} \mathbf{x}_i) + \lambda \boldsymbol{\mu}_0^T \mathbf{\Lambda} \boldsymbol{\mu}_0 + (N + \lambda) \boldsymbol{\mu}^T \mathbf{\Lambda} \boldsymbol{\mu} - 2(N \bar{\mathbf{x}}^T + \lambda \boldsymbol{\mu}_0^T) \mathbf{\Lambda} \boldsymbol{\mu} \\ & = \text{Tr} \left[\mathbf{W}^{-1} \mathbf{\Lambda} \right] + \sum_{i=1}^N (\mathbf{x}_i^T \mathbf{\Lambda} \mathbf{x}_i) + \lambda \boldsymbol{\mu}_0^T \mathbf{\Lambda} \boldsymbol{\mu}_0 - \frac{1}{\lambda + N} (\lambda \boldsymbol{\mu}_0 + N \bar{\mathbf{x}})^T \mathbf{\Lambda} (\lambda \boldsymbol{\mu}_0 + N \bar{\mathbf{x}}) \\ & \quad + (N + \lambda) \boldsymbol{\mu}^T \mathbf{\Lambda} \boldsymbol{\mu} - 2\boldsymbol{\mu}^T \mathbf{\Lambda} (N \bar{\mathbf{x}} + \lambda \boldsymbol{\mu}_0) + \frac{1}{\lambda + N} (\lambda \boldsymbol{\mu}_0 + N \bar{\mathbf{x}})^T \mathbf{\Lambda} (\lambda \boldsymbol{\mu}_0 + N \bar{\mathbf{x}}) \end{aligned}$$

Please note that the blue part, which is equal to:

$$(N + \lambda) \left(\boldsymbol{\mu} - \frac{\lambda \boldsymbol{\mu}_0 + N \bar{\mathbf{x}}}{\lambda + N} \right)^T \mathbf{\Lambda} \left(\boldsymbol{\mu} - \frac{\lambda \boldsymbol{\mu}_0 + N \bar{\mathbf{x}}}{\lambda + N} \right)$$

Now for the red part, we consider adding and subtracting $2N \bar{\mathbf{x}}^T \mathbf{\Lambda} \bar{\mathbf{x}}$, which we have:

$$\begin{aligned} & \sum_{i=1}^N (\mathbf{x}_i^T \mathbf{\Lambda} \mathbf{x}_i) - 2 \sum_{i=1}^N \mathbf{x}_i^T \mathbf{\Lambda} \bar{\mathbf{x}} + \sum_{i=1}^N \bar{\mathbf{x}} \mathbf{\Lambda} \bar{\mathbf{x}} + N \bar{\mathbf{x}}^T \mathbf{\Lambda} \bar{\mathbf{x}} \\ & = \sum_{i=1}^N (\mathbf{x}_i - \bar{\mathbf{x}})^T \mathbf{\Lambda} (\mathbf{x}_i - \bar{\mathbf{x}}) + N \bar{\mathbf{x}}^T \mathbf{\Lambda} \bar{\mathbf{x}} \end{aligned}$$

For the **brown** equation (together with $N\bar{\mathbf{x}}^T\Lambda\bar{\mathbf{x}}$), we have:

$$\begin{aligned}
& \lambda\boldsymbol{\mu}_0^T\Lambda\boldsymbol{\mu}_0 - \frac{1}{\lambda+N}(\lambda\boldsymbol{\mu}_0 + N\bar{\mathbf{x}})^T\Lambda(\lambda\boldsymbol{\mu}_0 + N\bar{\mathbf{x}}) + N\bar{\mathbf{x}}^T\Lambda\bar{\mathbf{x}} \\
&= \lambda\boldsymbol{\mu}_0^T\Lambda\boldsymbol{\mu}_0 + N\bar{\mathbf{x}}^T\Lambda\bar{\mathbf{x}} - \frac{1}{\lambda+N}\left[\lambda^2\boldsymbol{\mu}_0^T\Lambda\boldsymbol{\mu}_0 + 2N\lambda\bar{\mathbf{x}}^T\Lambda\boldsymbol{\mu}_0 + N^2\bar{\mathbf{x}}^T\Lambda\bar{\mathbf{x}}\right] \\
&= \frac{\lambda+N}{\lambda+N}\lambda\boldsymbol{\mu}_0^T\Lambda\boldsymbol{\mu}_0 + \frac{\lambda+N}{\lambda+N}N\bar{\mathbf{x}}^T\Lambda\bar{\mathbf{x}} - \frac{1}{\lambda+N}\left[\lambda^2\boldsymbol{\mu}_0^T\Lambda\boldsymbol{\mu}_0 + 2N\lambda\bar{\mathbf{x}}^T\Lambda\boldsymbol{\mu}_0 + N^2\bar{\mathbf{x}}^T\Lambda\bar{\mathbf{x}}\right] \\
&= \frac{1}{\lambda+N}\left[\lambda^2\boldsymbol{\mu}_0^T\Lambda\boldsymbol{\mu}_0 + N\lambda\boldsymbol{\mu}_0^T\Lambda\boldsymbol{\mu}_0 + N\lambda\bar{\mathbf{x}}^T\Lambda\bar{\mathbf{x}} + N^2\bar{\mathbf{x}}^T\Lambda\bar{\mathbf{x}}\right] \\
&\quad - \frac{1}{\lambda+N}\left[\lambda^2\boldsymbol{\mu}_0^T\Lambda\boldsymbol{\mu}_0 + 2N\lambda\bar{\mathbf{x}}^T\Lambda\boldsymbol{\mu}_0 + N^2\bar{\mathbf{x}}^T\Lambda\bar{\mathbf{x}}\right] \\
&= \frac{1}{\lambda+N}\left[N\lambda\boldsymbol{\mu}_0^T\Lambda\boldsymbol{\mu}_0 + N\lambda\bar{\mathbf{x}}^T\Lambda\bar{\mathbf{x}} + 2N\lambda\bar{\mathbf{x}}^T\Lambda\boldsymbol{\mu}_0\right] \\
&= \frac{N\lambda}{\lambda+N}(\bar{\mathbf{x}} - \boldsymbol{\mu}_0)^T\Lambda(\bar{\mathbf{x}} - \boldsymbol{\mu}_0)
\end{aligned}$$

Combining the first part (black, **red**, and **brown**), and together with the trace trick:

$$\text{Tr}\left(\mathbf{W}^{-1}\Lambda + \sum_{i=1}^N(\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^T\Lambda + \frac{N\lambda}{\lambda+N}(\bar{\mathbf{x}} - \boldsymbol{\mu}_0)(\bar{\mathbf{x}} - \boldsymbol{\mu}_0)^T\Lambda\right)$$

Now, the distribution becomes:

$$\begin{aligned}
|\Lambda|^{(\nu-D+N)/2} \exp\left(-\frac{1}{2}\text{Tr}\left(\left[\mathbf{W}^{-1} + \sum_{i=1}^N(\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^T + \frac{N\lambda}{\lambda+N}(\bar{\mathbf{x}} - \boldsymbol{\mu}_0)(\bar{\mathbf{x}} - \boldsymbol{\mu}_0)^T\right]\Lambda\right)\right. \\
\left.- \frac{N+\lambda}{2}\left(\boldsymbol{\mu} - \frac{\lambda\boldsymbol{\mu}_0 + N\bar{\mathbf{x}}}{\lambda+N}\right)^T\Lambda\left(\boldsymbol{\mu} - \frac{\lambda\boldsymbol{\mu}_0 + N\bar{\mathbf{x}}}{\lambda+N}\right)\right)
\end{aligned}$$

Which is simply a Normal-Wishart Distribution as required.

Proposition 1.14. (Maximum Likelihood of Mean) We can show that for Gaussian distribution, the maximum likelihood of $\boldsymbol{\mu}$, given the dataset $\{\mathbf{x}_i\}_{i=1}^N$ is

$$\hat{\boldsymbol{\mu}}_{ML} = \frac{1}{N}\sum_{i=1}^N\mathbf{x}_i$$

Proof. Starting with the log-likelihood to be:

$$l(\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \log\prod_{i=1}^N\mathcal{N}(\mathbf{x}_i|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = -\frac{N}{2}\log|2\pi\boldsymbol{\Sigma}| - \frac{1}{2}\sum_{i=1}^N(\mathbf{x}_i - \boldsymbol{\mu})^T\boldsymbol{\Sigma}^{-1}(\mathbf{x}_i - \boldsymbol{\mu})$$

And, so we can consider the derivative over $\boldsymbol{\mu}$ as:

$$\begin{aligned}
\frac{\partial(-l)}{\partial\boldsymbol{\mu}} &= \frac{\partial}{\partial\boldsymbol{\mu}}\left[\frac{N}{2}\log|2\pi\boldsymbol{\Sigma}| + \frac{1}{2}\sum_{i=1}^N(\mathbf{x}_i - \boldsymbol{\mu})^T\boldsymbol{\Sigma}^{-1}(\mathbf{x}_i - \boldsymbol{\mu})\right] \\
&= \frac{1}{2}\sum_{i=1}^N\frac{\partial}{\partial\boldsymbol{\mu}}(\mathbf{x}_i - \boldsymbol{\mu})^T\boldsymbol{\Sigma}^{-1}(\mathbf{x}_i - \boldsymbol{\mu}) \\
&= \frac{1}{2}\sum_{i=1}^N\left(\frac{\partial}{\partial\boldsymbol{\mu}}[\boldsymbol{\mu}^T\boldsymbol{\Sigma}^{-1}\boldsymbol{\mu}] - 2\frac{\partial}{\partial\boldsymbol{\mu}}[\boldsymbol{\mu}^T\boldsymbol{\Sigma}^{-1}\mathbf{x}_i]\right) = N\boldsymbol{\Sigma}^{-1}\boldsymbol{\mu} - \boldsymbol{\Sigma}^{-1}\sum_{i=1}^N\mathbf{x}_i
\end{aligned}$$

Setting the derivative to zero and we yields the result. \square

Proposition 1.15. (Maximum Likelihood of Covariance) We can show that for Gaussian distribution, the maximum likelihood of Σ , given the dataset $\{\mathbf{x}_i\}_{i=1}^N$ is

$$\hat{\Sigma}_{ML} = \frac{1}{N} \sum_{i=1}^N (\mathbf{x}_i - \boldsymbol{\mu})(\mathbf{x}_i - \boldsymbol{\mu})^T$$

Proof. We consider the derivative over Σ^{-1} as (please note that we have to constraint of the covariance to be positive definite but it turn out we don't have to as the result already satisfies the constraint):

$$\begin{aligned} \frac{\partial(-l)}{\partial \Sigma^{-1}} &= \frac{\partial}{\partial \Sigma^{-1}} \left[\frac{N}{2} \log |2\pi \Sigma| + \frac{1}{2} \sum_{i=1}^N (\mathbf{x}_i - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{x}_i - \boldsymbol{\mu}) \right] \\ &= -\frac{\partial}{\partial \Sigma^{-1}} \left[\frac{N}{2} \log |\Sigma^{-1}| \right] + \frac{1}{2} \sum_{i=1}^N \frac{\partial}{\partial \Sigma^{-1}} [(\mathbf{x}_i - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{x}_i - \boldsymbol{\mu})] \\ &= -\frac{N}{2} \Sigma^T + \frac{1}{2} \sum_{i=1}^N (\mathbf{x}_i - \boldsymbol{\mu})(\mathbf{x}_i - \boldsymbol{\mu})^T \end{aligned}$$

Setting the derivative to zero and we yields the result. \square

Definition 1.7. (Linear Regression) Now, given the data: $\mathcal{D} = \{(\mathbf{x}_i, \mathbf{y})\}_{i=1}^N$. We have \mathbf{y}_i is conditionally independent given \mathbf{x}_i . Now, we consider the supervised learning as we have a linear function \mathbf{x} together with Gaussian noise:

$$p(\mathbf{y}|\mathbf{x}, \mathbf{W}, \Sigma_y) = \frac{1}{\sqrt{|2\pi \Sigma_y|}} \exp \left\{ -\frac{1}{2} (\mathbf{y} - \mathbf{W}\mathbf{x})^T \Sigma_y^{-1} (\mathbf{y} - \mathbf{W}\mathbf{x}) \right\}$$

Proposition 1.16. (Maximum Likelihood of Linear Regression) The maximum likelihood is given by the following:

$$\hat{\mathbf{W}} = \sum_{i=1}^N \mathbf{y}_i \mathbf{x}_i^T \left(\sum_{i=1}^N \mathbf{x}_i \mathbf{x}_i^T \right)^{-1}$$

Proof. Consider the log-likelihood of the linear regression:

$$\begin{aligned} l &= \log \prod_{i=1}^N \frac{1}{\sqrt{|2\pi \Sigma_y|}} \exp \left\{ -\frac{1}{2} (\mathbf{y}_i - \mathbf{W}\mathbf{x}_i)^T \Sigma_y^{-1} (\mathbf{y}_i - \mathbf{W}\mathbf{x}_i) \right\} \\ &= \sum_{i=1}^N -\frac{1}{2} \log |2\pi \Sigma_y| - \frac{1}{2} (\mathbf{y}_i - \mathbf{W}\mathbf{x}_i)^T \Sigma_y^{-1} (\mathbf{y}_i - \mathbf{W}\mathbf{x}_i) \\ &= -\frac{N}{2} \log |2\pi \Sigma_y| - \frac{1}{2} \sum_{i=1}^N (\mathbf{y}_i - \mathbf{W}\mathbf{x}_i)^T \Sigma_y^{-1} (\mathbf{y}_i - \mathbf{W}\mathbf{x}_i) \end{aligned}$$

We consider the derivative, as we have:

$$\begin{aligned} \frac{\partial(-l)}{\partial \mathbf{W}} &= \frac{\partial}{\partial \mathbf{W}} \left[\frac{N}{2} \log |2\pi \Sigma_y| + \frac{1}{2} \sum_{i=1}^N (\mathbf{y}_i - \mathbf{W}\mathbf{x}_i)^T \Sigma_y^{-1} (\mathbf{y}_i - \mathbf{W}\mathbf{x}_i) \right] \\ &= \frac{1}{2} \sum_{i=1}^N \frac{\partial}{\partial \mathbf{W}} [(\mathbf{y}_i - \mathbf{W}\mathbf{x}_i)^T \Sigma_y^{-1} (\mathbf{y}_i - \mathbf{W}\mathbf{x}_i)] \\ &= \frac{1}{2} \sum_{i=1}^N \frac{\partial}{\partial \mathbf{W}} [\mathbf{y}_i^T \Sigma_y^{-1} \mathbf{y}_i + \mathbf{x}_i^T \mathbf{W}^T \Sigma_y^{-1} \mathbf{W} \mathbf{x}_i - 2\mathbf{x}_i^T \mathbf{W}^T \Sigma_y^{-1} \mathbf{y}_i] \end{aligned}$$

$$\begin{aligned}
&= \frac{1}{2} \sum_{i=1}^N \left[\frac{\partial}{\partial \mathbf{W}} \text{Tr}[\mathbf{x}_i^T \mathbf{W}^T \Sigma_y^{-1} \mathbf{W} \mathbf{x}_i] - 2 \frac{\partial}{\partial \mathbf{W}} \text{Tr}[\mathbf{x}_i^T \mathbf{W}^T \Sigma_y^{-1} \mathbf{y}_i] \right] \\
&= \frac{1}{2} \sum_{i=1}^N \left[\frac{\partial}{\partial \mathbf{W}} \text{Tr}[\mathbf{W}^T \Sigma_y^{-1} \mathbf{W} \mathbf{x}_i \mathbf{x}_i^T] - 2 \frac{\partial}{\partial \mathbf{W}} \text{Tr}[\mathbf{W}^T \Sigma_y^{-1} \mathbf{y}_i \mathbf{x}_i^T] \right] \\
&= \frac{1}{2} \sum_{i=1}^N [2 \Sigma_y^{-1} \mathbf{W} \mathbf{x}_i \mathbf{x}_i^T - 2 \Sigma_y^{-1} \mathbf{y}_i \mathbf{x}_i^T]
\end{aligned}$$

Please noted that we use the derivative of the trace above, and setting the derivative to zero and we have the answer as we wanted. \square

Proposition 1.17. *The MAP estimate given prior over the weight $p(\mathbf{w}) = \mathcal{N}(\mathbf{0}, \mathbf{A}^{-1})$ is:*

$$\mathbf{w}^{MAP} = \underbrace{\left(\mathbf{A} + \frac{\sum_{i=1}^N \mathbf{x}_i \mathbf{x}_i^T}{\sigma_y^2} \right)^{-1}}_{\Sigma_w} \frac{\sum_{i=1}^N \mathbf{y}_i \mathbf{x}_i}{\sigma_y^2} = \left(A \sigma_y^2 + \sum_{i=1}^N \mathbf{x}_i \mathbf{x}_i^T \right)^{-1} \sum_{i=1}^N \mathbf{y}_i \mathbf{x}_i$$

We will denote the a Please note that we consider the prediction in one-dimensional, as we have

$$p(\mathcal{D} | \{\mathbf{x}_i\}_{i=1}^N, \mathbf{w}, \sigma_y) = \frac{1}{\sqrt{2\pi\sigma_y^2}} \exp \left\{ -\frac{1}{2\sigma_y^2} \sum_{i=1}^N (y_i - \mathbf{w}^T \mathbf{x}_i)^2 \right\}$$

Proof. We start with the log-posterior on \mathbf{w} is given as:

$$\begin{aligned}
\log p(\mathbf{w} | \mathcal{D}, A, \sigma_y) &= \log p(\mathcal{D} | \{\mathbf{x}_i\}_{i=1}^N, \mathbf{w}, \sigma_y) + \log p(\mathbf{w} | A) + \text{const} \\
&= -\frac{1}{2} \mathbf{w}^T \mathbf{A} \mathbf{w} - \frac{1}{2\sigma_y^2} \sum_{i=1}^N (y_i - \mathbf{w}^T \mathbf{x}_i)^2 + \text{const} \\
&= -\frac{1}{2} \mathbf{w}^T \mathbf{A} \mathbf{w} - \frac{1}{2\sigma_y^2} \sum_{i=1}^N [y_i^2 - 2y_i \mathbf{w}^T \mathbf{x}_i + (\mathbf{w}^T \mathbf{x}_i)^2] + \text{const} \\
&= -\frac{1}{2} \mathbf{w}^T \mathbf{A} \mathbf{w} - \frac{1}{2\sigma_y^2} \sum_{i=1}^N 2y_i \mathbf{w}^T \mathbf{x}_i - \frac{1}{2} \mathbf{w}^T \left(\sigma_y^{-2} \sum_{i=1}^N \mathbf{x}_i \mathbf{x}_i^T \right) \mathbf{w} + \text{const} \\
&= -\frac{1}{2} \mathbf{w}^T \left(\mathbf{A} + \sigma_y^{-2} \sum_{i=1}^N \mathbf{x}_i \mathbf{x}_i^T \right) \mathbf{w} - \mathbf{w}^T \sum_{i=1}^N (y_i \mathbf{x}_i) \sigma_y^{-2} + \text{const} \\
&= -\frac{1}{2} \mathbf{w}^T \Sigma_w^{-1} \mathbf{w} - \mathbf{w}^T \Sigma_w^{-1} \Sigma_w \sum_{i=1}^N (y_i \mathbf{x}_i) \sigma_y^{-2} + \text{const} \\
&= \log \mathcal{N} \left(\Sigma_w \sum_{i=1}^N (y_i \mathbf{x}_i) \sigma_y^{-2}, \Sigma_w \right)
\end{aligned}$$

As the Gaussian gives the maximum probability at its mean; therefore, we have completed the prove. \square

2 Latent Variable Model

Definition 2.1. (Latent Variable Model) The latent variable model can be seen as:

$$\begin{aligned} \mathbf{z} &\sim p(\boldsymbol{\theta}_z) \\ \mathbf{x}|\mathbf{z} &\sim p(\boldsymbol{\theta}_z) \\ p(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta}_x, \boldsymbol{\theta}_z) &= p(\mathbf{x}|\mathbf{z}; \boldsymbol{\theta}_x)p(\mathbf{z}; \boldsymbol{\theta}_z) \\ p(\mathbf{x}; \boldsymbol{\theta}_x, \boldsymbol{\theta}_z) &= \int p(\mathbf{x}|\mathbf{z}; \boldsymbol{\theta}_x)p(\mathbf{z}; \boldsymbol{\theta}_z) d\mathbf{z} \end{aligned}$$

Note that $p(\mathbf{z})$, $p(\mathbf{x}|\mathbf{z})$ and $p(\mathbf{x}, \mathbf{z})$ are exponential family but $p(\mathbf{x})$ doesn't have to be an exponential family.

2.1 PCA Formulation

Remark 19. We will consider the family of PCA formulation. We will start with PCA definition, which can be formulated in 2 ways: Maximal Variance and Average Projection Cost. Then, we will consider the PPCA a probabilistic version of PCA.

Definition 2.2. (Maximal Variance) Consider the dataset $\{\mathbf{x}_i\}_{i=1}^N$ where $\mathbf{x}_i \in \mathbb{R}^D$. We want to project the data onto space with dimension $M < D$. To do this, we want to find a subspace orthonormal basis \mathbf{u}_i for $i = 1, \dots, M$, such that $\mathbf{u}_i^T \mathbf{u}_j = \delta_{ij}$. So that the empirical projected variance, given as:

$$\frac{1}{N} \sum_{j=1}^N (\mathbf{u}_i^T \mathbf{x}_j - \mathbf{u}_i^T \bar{\mathbf{x}})^2$$

is maximized for $i = 1, \dots, M$ where \mathbf{u}_1 gives the highest variance, where $\bar{\mathbf{x}} = 1/N \sum_{i=1}^N \mathbf{x}_i$

Proposition 2.1. *The maximum projected variance direction $\mathbf{u}_1, \dots, \mathbf{u}_M$ is the M eigenvectors of the data-covariance matrix:*

$$\mathbf{S} = \frac{1}{N} \sum_{i=1}^N (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^T$$

associated with the following eigenvalues $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_M \geq \dots \geq \lambda_D$.

Proof. Let's start with the first direction \mathbf{u}_1 as we can show that the empirical projected variance:

$$\frac{1}{N} \sum_{j=1}^N (\mathbf{u}_1^T \mathbf{x}_j - \mathbf{u}_1^T \bar{\mathbf{x}})^2 = \mathbf{u}_1^T \mathbf{S} \mathbf{u}_1$$

We will consider the following constraint optimization problem as we have the following Lagrange multiplier and set the derivative to be 0:

$$\frac{\partial}{\partial \mathbf{u}} \left[\mathbf{u}_1^T \mathbf{S} \mathbf{u}_1 + \lambda_1 (1 - \mathbf{u}_1^T \mathbf{u}_1) \right] = 2\mathbf{S} \mathbf{u}_1 - 2\lambda_1 \mathbf{u}_1 = 0$$

Please note that the matrix \mathbf{S} is symmetric. This leads us to the following, equation, which is:

$$\mathbf{S} \mathbf{u}_1 = \lambda_1 \mathbf{u}_1$$

Furthermore, if we multiply \mathbf{u}_1^T on the LHS together with the constraint $\mathbf{u}_1^T \mathbf{u}_1 = 1$, then we have $\mathbf{u}_1^T \mathbf{S} \mathbf{u}_1 = \lambda_1$. This means that first maximal variance projection direction is the eigenvector \mathbf{u}_1 that has the highest associated eigenvalue. Furthermore, the orthogonal comes from the properties of eigenvectors. \square

Definition 2.3. (Minimum-Error Formulation) We get the formulation PCA where we have orthonormal set D -dimensional basis vector $\{\mathbf{u}_i\}$ where $i = 1, \dots, D$ and $\mathbf{u}_i^T \mathbf{u}_j = \delta_{ij}$. This means that the data point can be represented by the basis:

$$\mathbf{x}_i = \sum_{j=1}^D \alpha_{ij} \mathbf{u}_j = \sum_{j=1}^D (\mathbf{x}_i^T \mathbf{u}_j) \mathbf{u}_j$$

The value of α_{ij} can be found by the inner product $\alpha_{ij} = \mathbf{x}_i^T \mathbf{u}_j$ as above (by the orthogonal properties). Now, we will consider the approximation using the projection over linear subspace $M < D$:

$$\tilde{\mathbf{x}}_i = \sum_{j=1}^M z_{ij} \mathbf{u}_j + \sum_{j=M+1}^D b_j \mathbf{u}_j$$

where the $\{b_j\}$ are component that is same for all data points. Now, we are free to find the $\{b_j\}$, $\{z_{ij}\}$ and $\{\mathbf{u}_j\}$ given the following objective:

$$\frac{1}{N} \sum_{i=1}^N \|\mathbf{x}_i - \tilde{\mathbf{x}}_i\|^2$$

Proposition 2.2. *The solution to the minimum error formulation is the same as the maximal variance formulation. This gives us the difference interpretation of the PCA.*

Proof. Let's start with finding the value $\{z_{ij}\}$, first. As we want to find derivative of z_{ij} to be zero:

$$\begin{aligned} & \frac{1}{N} \sum_{i=1}^N \left\| \mathbf{x}_i - \sum_{j=1}^M z_{ij} \mathbf{u}_j - \sum_{j=M+1}^D b_j \mathbf{u}_j \right\|^2 \\ &= \frac{1}{N} \sum_{i=1}^N \left[\left(\mathbf{x}_i - \sum_{j=1}^M z_{ij} \mathbf{u}_j - \sum_{j=M+1}^D b_j \mathbf{u}_j \right)^T \left(\mathbf{x}_i - \sum_{j=1}^M z_{ij} \mathbf{u}_j - \sum_{j=M+1}^D b_j \mathbf{u}_j \right) \right] \\ &= \frac{1}{N} \sum_{i=1}^N \left[\mathbf{x}_i^T \mathbf{x}_i - \sum_{j=1}^M z_{ij} \mathbf{u}_j^T \mathbf{x}_i - \sum_{j=M+1}^D b_j \mathbf{u}_j^T \mathbf{x}_i - \sum_{j=1}^M z_{ij} \mathbf{x}_i^T \mathbf{u}_j \right. \\ & \quad + \left(\sum_{j=1}^M z_{ij} \mathbf{u}_j^T \right) \left(\sum_{j=1}^M z_{ij} \mathbf{u}_j \right) + \left(\sum_{j=M+1}^D b_j \mathbf{u}_j^T \right) \left(\sum_{j=1}^M z_{ij} \mathbf{u}_j \right) \\ & \quad \left. - \sum_{j=M+1}^D b_j \mathbf{x}_i^T \mathbf{u}_j + \left(\sum_{j=1}^M z_{ij} \mathbf{u}_j^T \right) \left(\sum_{j=M+1}^D b_j \mathbf{u}_j \right) + \left(\sum_{j=M+1}^D b_j \mathbf{u}_j^T \right) \left(\sum_{j=M+1}^D b_j \mathbf{u}_j \right) \right] \\ &= \frac{1}{N} \sum_{i=1}^N \left[\mathbf{x}_i^T \mathbf{x}_i - \sum_{j=1}^M z_{ij} \mathbf{u}_j^T \mathbf{x}_i - \sum_{j=M+1}^D b_j \mathbf{u}_j^T \mathbf{x}_i - \sum_{j=1}^M z_{ij} \mathbf{x}_i^T \mathbf{u}_j - \sum_{j=M+1}^D b_j \mathbf{x}_i^T \mathbf{u}_j \right. \\ & \quad + \left(\sum_{a=1}^M z_{ia} \mathbf{u}_a^T \right) \left(\sum_{b=1}^M z_{ib} \mathbf{u}_b \right) + \left(\sum_{a=M+1}^D b_a \mathbf{u}_a^T \right) \left(\sum_{b=1}^M z_{ib} \mathbf{u}_b \right) \\ & \quad \left. + \left(\sum_{a=1}^M z_{ia} \mathbf{u}_a^T \right) \left(\sum_{b=M+1}^D b_b \mathbf{u}_b \right) + \left(\sum_{a=M+1}^D b_a \mathbf{u}_a^T \right) \left(\sum_{b=M+1}^D b_b \mathbf{u}_b \right) \right] \\ &= \frac{1}{N} \sum_{i=1}^N \left[\mathbf{x}_i^T \mathbf{x}_i - 2 \sum_{j=1}^M z_{ij} \mathbf{u}_j^T \mathbf{x}_i - 2 \sum_{j=M+1}^D b_j \mathbf{u}_j^T \mathbf{x}_i + \sum_{a=1}^M \sum_{b=1}^M z_{ia} z_{ib} \mathbf{u}_a^T \mathbf{u}_b \right. \\ & \quad \left. + 2 \sum_{a=M+1}^D \sum_{b=1}^M b_a z_{ib} \mathbf{u}_a^T \mathbf{u}_b + \sum_{a=M+1}^D \sum_{b=M+1}^D b_a b_b \mathbf{u}_a^T \mathbf{u}_b \right] \end{aligned}$$

$$\begin{aligned}
&= \frac{1}{N} \sum_{i=1}^N \left[\mathbf{x}_i^T \mathbf{x}_i - 2 \sum_{j=1}^M z_{ij} \mathbf{u}_j^T \mathbf{x}_i - 2 \sum_{j=M+1}^D b_j \mathbf{u}_j^T \mathbf{x}_i + \sum_{j=1}^M z_{ij}^2 + \sum_{j=M+1}^D b_j^2 + 2 \sum_{a=M+1}^D \sum_{b=1}^M b_a z_{ib} \mathbf{u}_a^T \mathbf{u}_b \right] \\
&= \frac{1}{N} \sum_{i=1}^N \left[\mathbf{x}_i^T \mathbf{x}_i - 2 \sum_{j=1}^M z_{ij} \mathbf{u}_j^T \mathbf{x}_i - 2 \sum_{j=M+1}^D b_j \mathbf{u}_j^T \mathbf{x}_i + \sum_{j=1}^M z_{ij}^2 + \sum_{j=M+1}^D b_j^2 \right]
\end{aligned}$$

Now, let's consider the derivative with respected to z_{ab} as we now have:

$$\begin{aligned}
&\frac{\partial}{\partial z_{ab}} \frac{1}{N} \sum_{i=1}^N \left[\mathbf{x}_i^T \mathbf{x}_i - 2 \sum_{j=1}^M z_{ij} \mathbf{u}_j^T \mathbf{x}_i - 2 \sum_{j=M+1}^D b_j \mathbf{u}_j^T \mathbf{x}_i + \sum_{j=1}^M z_{ij}^2 + \sum_{j=M+1}^D b_j^2 \right] \\
&= -\frac{\partial}{\partial z_{ab}} \frac{2}{N} \sum_{i=1}^N \sum_{j=1}^M z_{ij} \mathbf{u}_j^T \mathbf{x}_i + \frac{\partial}{\partial z_{ab}} \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^M z_{ij}^2 \\
&= -\frac{2}{N} \mathbf{u}_b^T \mathbf{x}_a + \frac{2}{N} z_{ab} = 0
\end{aligned}$$

And so, we have the $z_{ij} = \mathbf{x}_i^T \mathbf{u}_j$ for $j = 1, \dots, M$. Now, we consider the derivative of the with respected to b_a as we have:

$$\begin{aligned}
&\frac{\partial}{\partial b_a} \frac{1}{N} \sum_{i=1}^N \left[\mathbf{x}_i^T \mathbf{x}_i - 2 \sum_{j=1}^M z_{ij} \mathbf{u}_j^T \mathbf{x}_i - 2 \sum_{j=M+1}^D b_j \mathbf{u}_j^T \mathbf{x}_i + \sum_{j=1}^M z_{ij}^2 + \sum_{j=M+1}^D b_j^2 \right] \\
&= -\frac{\partial}{\partial b_a} \frac{2}{N} \sum_{i=1}^N \sum_{j=M+1}^D b_j \mathbf{u}_j^T \mathbf{x}_i + \frac{\partial}{\partial b_a} \frac{1}{N} \sum_{i=1}^N \sum_{j=M+1}^D b_j^2 \\
&= -\frac{\partial}{\partial b_a} 2 \sum_{j=M+1}^D b_j \mathbf{u}_j^T \left(\frac{1}{N} \sum_{i=1}^N \mathbf{x}_i \right) + \frac{\partial}{\partial b_a} \sum_{j=M+1}^D b_j^2 \\
&= -\frac{\partial}{\partial b_a} 2 \sum_{j=M+1}^D b_j \mathbf{u}_j^T \left(\frac{1}{N} \sum_{i=1}^N \mathbf{x}_i \right) + \frac{\partial}{\partial b_a} \sum_{j=M+1}^D b_j^2 \\
&= -2 \mathbf{u}_a^T \bar{\mathbf{x}} + 2 b_a
\end{aligned}$$

And so, we have $b_j = \bar{\mathbf{x}}^T \mathbf{u}_j$ for $j = M+1, \dots, D$. To find the \mathbf{u}_i , we have the following:

$$\begin{aligned}
\mathbf{x}_i - \tilde{\mathbf{x}}_i &= \mathbf{x}_i - \sum_{j=1}^M z_{ij} \mathbf{u}_j - \sum_{j=M+1}^D b_j \mathbf{u}_j \\
&= \sum_{j=1}^M (\mathbf{x}_i^T \mathbf{u}_j) \mathbf{u}_j + \sum_{j=M+1}^D (\mathbf{x}_i^T \mathbf{u}_j) \mathbf{u}_j - \sum_{j=1}^M (\mathbf{x}_i^T \mathbf{u}_j) \mathbf{u}_j - \sum_{j=M+1}^D (\bar{\mathbf{x}}^T \mathbf{u}_j) \mathbf{u}_j \\
&= \sum_{j=M+1}^D (\mathbf{x}_i^T \mathbf{u}_j) \mathbf{u}_j - \sum_{j=M+1}^D (\bar{\mathbf{x}}^T \mathbf{u}_j) \mathbf{u}_j \\
&= \sum_{j=M+1}^D \{(\mathbf{x}_i - \bar{\mathbf{x}})^T \mathbf{u}_j\} \mathbf{u}_j
\end{aligned}$$

And, so we now have, the following objective:

$$\frac{1}{N} \sum_{i=1}^N \left[\left(\sum_{j=M+1}^D \{(\mathbf{x}_i - \bar{\mathbf{x}})^T \mathbf{u}_j\} \mathbf{u}_j \right)^T \left(\sum_{j=M+1}^D \{(\mathbf{x}_i - \bar{\mathbf{x}})^T \mathbf{u}_j\} \mathbf{u}_j \right) \right]$$

$$\begin{aligned}
&= \frac{1}{N} \sum_{i=1}^N \left[\left(\sum_{a=M+1}^D \mathbf{u}_a^T \{ \mathbf{u}_a^T (\mathbf{x}_i - \bar{\mathbf{x}}) \} \right) \left(\sum_{b=M+1}^D \{ (\mathbf{x}_i - \bar{\mathbf{x}})^T \mathbf{u}_b \} \mathbf{u}_b \right) \right] \\
&= \frac{1}{N} \sum_{i=1}^N \left[\left(\sum_{a=M+1}^D \sum_{b=M+1}^D \{ \mathbf{u}_a^T (\mathbf{x}_i - \bar{\mathbf{x}}) (\mathbf{x}_i - \bar{\mathbf{x}})^T \mathbf{u}_b \} \mathbf{u}_a^T \mathbf{u}_b \right) \right] \\
&= \frac{1}{N} \sum_{i=1}^N \sum_{a=M+1}^D \mathbf{u}_a^T (\mathbf{x}_i - \bar{\mathbf{x}}) (\mathbf{x}_i - \bar{\mathbf{x}})^T \mathbf{u}_a = \sum_{a=M+1}^D \mathbf{u}_a^T \mathbf{S} \mathbf{u}_a
\end{aligned}$$

Now, we have the same objective to the maximal variance formulation. And the proposition is proven. \square

2.2 Probabilistic PCA

Definition 2.4. (PPCA) We consider the following system of probabilities:

$$p(\mathbf{z}) = \mathcal{N}(\mathbf{z} | \mathbf{0}, \mathbf{I}) \quad p(\mathbf{x} | \mathbf{z}) = \mathcal{N}(\mathbf{x} | \mathbf{W}\mathbf{z} + \boldsymbol{\mu}, \Psi)$$

which we can consider $p(\mathbf{z})$ to be the PCA projection, while $p(\mathbf{x} | \mathbf{z})$ is the reconstruction of the PCA. This means that we can perform PPCA of it.

Proposition 2.3. We consider the marginalization of Gaussian to be:

$$p(\mathbf{x}) = \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}, \mathbf{C}) \quad \text{where} \quad \mathbf{C} = \Psi + \mathbf{W}\mathbf{W}^T$$

Proof. We consider the linear Gaussian model in this case, where we use the marginalization (see above) to get the value of $p(\mathbf{x})$ \square

Proposition 2.4. We consider the inference of the latent and we have:

$$p(\mathbf{z} | \mathbf{x}) = \mathcal{N}(\boldsymbol{\Sigma}^{-1} \mathbf{W}^T \Psi^{-1} (\mathbf{x} - \boldsymbol{\mu}), \boldsymbol{\Sigma}^{-1})$$

where we have $\boldsymbol{\Sigma} = \mathbf{I} + \mathbf{W}^T \Psi^{-1} \mathbf{W}$.

Remark 20. We have the projection to be:

$$\hat{\mathbf{x}}_i = \mathbf{W} \boldsymbol{\Sigma}^{-1} \mathbf{W}^T \Psi^{-1} (\mathbf{x}_i - \boldsymbol{\mu})$$

As we have the PCA projection that also take noise into consideration. Furthermore, if $\Psi = \psi^2 \mathbf{I}$ and $\psi \rightarrow 0$, then it leads to the PCA estimation (given the correct \mathbf{W} , which we will explore later).

Remark 21. (Likelihood of PPCA) Now, we are left to find the actual value of \mathbf{W} and we will assume that we are aware of $\boldsymbol{\mu}$ (which is usually $\mathbf{0}$), while the covariance matrix is assumed to be $\Psi = \psi^2 \mathbf{I}$. The log-likelihood of this PPCA is (using the marginalized):

$$\begin{aligned}
l = \log p(\{\mathbf{x}_i\}_{i=1}^N | \boldsymbol{\mu}, \mathbf{C}) &= \log \prod_{i=1}^N \frac{1}{\sqrt{|2\pi\mathbf{C}|}} \exp \left\{ -\frac{1}{2} (\mathbf{x}_i - \boldsymbol{\mu})^T \mathbf{C}^{-1} (\mathbf{x}_i - \boldsymbol{\mu}) \right\} \\
&= -\frac{N}{2} \log |\mathbf{C}| - \frac{1}{2} \sum_{i=1}^N (\mathbf{x}_i - \boldsymbol{\mu})^T \mathbf{C}^{-1} (\mathbf{x}_i - \boldsymbol{\mu}) + \text{const.} \\
&= -\frac{N}{2} \log |\mathbf{C}| - \frac{N}{2} \text{Tr} \left(\frac{1}{N} \sum_{i=1}^N (\mathbf{x}_i - \boldsymbol{\mu})^T \mathbf{C}^{-1} (\mathbf{x}_i - \boldsymbol{\mu}) \right) + \text{const.} \\
&= -\frac{N}{2} \log |\mathbf{C}| - \frac{N}{2} \text{Tr} \left(\mathbf{C}^{-1} \frac{1}{N} \sum_{i=1}^N (\mathbf{x}_i - \boldsymbol{\mu}) (\mathbf{x}_i - \boldsymbol{\mu})^T \right) + \text{const.} \\
&= -\frac{N}{2} \log |\mathbf{C}| - \frac{N}{2} \text{Tr} (\mathbf{C}^{-1} \mathbf{S}) + \text{const.}
\end{aligned}$$

Now $\mathbf{C} = \mathbf{W}\mathbf{W}^T + \psi^2 \mathbf{I}$ is given above.

Proposition 2.5. *Non-Trivial Solution of Maximum Likelihood estimate of \mathbf{W} is equal to:*

$$\mathbf{W}_{ML} = \mathbf{U}(\mathbf{\Lambda} - \psi^2 \mathbf{I})^{1/2} \mathbf{V}^T$$

where $\mathbf{U} \in \mathbb{R}^{D \times M}$ is the first M eigenvector of \mathbf{S} the empirical variance matrix and $\mathbf{\Lambda} \in \mathbb{R}^{M \times M}$ be the matrix, which is the eigenvalues of \mathbf{S} . Finally, $\mathbf{V} \in \mathbb{R}^{M \times M}$ is an arbitrary orthogonal matrix.

Proof. Let's consider the derivative of the log-likelihood with respect to \mathbf{W} as we now have:

$$\begin{aligned} \frac{\partial l}{\partial \mathbf{W}} &= \frac{\partial}{\partial \mathbf{W}} \left[-\frac{N}{2} \log |\mathbf{C}| - \frac{N}{2} \text{Tr}(\mathbf{C}^{-1} \mathbf{S}) \right] \\ &= -N \left[-\mathbf{C}^{-1} \mathbf{W} + \mathbf{C}^{-1} \mathbf{S} \mathbf{C}^{-1} \mathbf{W} \right] \end{aligned}$$

Setting this to $\mathbf{0}$, and we can see that, we have the following equation: $\mathbf{S} \mathbf{C}^{-1} \mathbf{W} = \mathbf{W}$, there are 2 outcome to the solution: $\mathbf{W} = \mathbf{0}$, which can be shown to be minimum. Or, Consider \mathbf{W} in the SVD form i.e $\mathbf{W} = \mathbf{U} \mathbf{L} \mathbf{V}^T$ for orthogonal matrix \mathbf{U} and \mathbf{V} , while \mathbf{L} is diagonal matrix. This would entail:

$$\begin{aligned} \mathbf{S}(\mathbf{U} \mathbf{L} \mathbf{V}^T \mathbf{V} \mathbf{L}^T \mathbf{U}^T + \psi^2 \mathbf{I})^{-1} \mathbf{U} \mathbf{L} \mathbf{V}^T &= \mathbf{U} \mathbf{L} \mathbf{V}^T \\ \implies \mathbf{S}(\mathbf{U} \mathbf{L}^2 \mathbf{U}^T + \psi^2 \mathbf{I})^{-1} \mathbf{U} &= \mathbf{U} \\ \implies \mathbf{S} \mathbf{U} (\mathbf{L}^2 + \psi^2 \mathbf{I})^{-1} &= \mathbf{U} \\ \implies \mathbf{S} \mathbf{U} = \mathbf{U} (\mathbf{L}^2 + \psi^2 \mathbf{I}) \end{aligned}$$

For the second implication, we have:

$$\mathbf{U} (\mathbf{L}^2 + \psi^2 \mathbf{I}) = (\mathbf{U} \mathbf{L}^2 \mathbf{U}^T + \psi^2 \mathbf{I}) \mathbf{U} \implies (\mathbf{U} \mathbf{L}^2 \mathbf{U}^T + \psi^2 \mathbf{I})^{-1} \mathbf{U} = \mathbf{U} (\mathbf{L}^2 + \psi^2 \mathbf{I})^{-1}$$

We can see that the \mathbf{U} is the eigenvector of \mathbf{S} , where the corresponding eigenvalues are $\lambda_i = l_i^2 + \psi^2$, and so we can rewrite the weight to be:

$$\mathbf{W} = \mathbf{U}(\mathbf{\Lambda} - \psi^2 \mathbf{I})^{1/2} \mathbf{V}^T$$

where $\mathbf{\Lambda}$ be the matrix, which is the eigenvalues of \mathbf{S} . □

2.3 Other Related Models

Definition 2.5. (Factor Analysis) The factor analysis is PPCA:

$$p(\mathbf{z}) = \mathcal{N}(\mathbf{z} | \mathbf{0}, \mathbf{I}) \quad p(\mathbf{x} | \mathbf{z}) = \mathcal{N}(\mathbf{x} | \mathbf{W} \mathbf{z} + \boldsymbol{\mu}, \boldsymbol{\Psi})$$

but we consider the matrix $\boldsymbol{\Psi}$ to be $D \times D$ diagonal matrix. This means that the inferences still holds. However, the training is much harder now as we may not find the \mathbf{W} from the data in closed form.

Definition 2.6. (Canonical Correlation Analysis) The data vector $\mathcal{D} = \{(\mathbf{u}_1, \mathbf{v}_1), (\mathbf{u}_2, \mathbf{v}_2), \dots\}$ where $\mathbf{u}_i \in \mathcal{U}$ and $\mathbf{v}_i \in \mathcal{V}$. We want to find the correlation:

- We find te unti vector $\mathbf{a} \in \mathcal{U}$ and $\mathbf{b} \in \mathcal{V}$ such that the correlation $\mathbf{u}_i^T \mathbf{a}$ and $\mathbf{v}_i^T \mathbf{b}$ is the maximum the covariance between them.
- This also requires some in the orthogonal subspace.

Now, the probabilistic CCA is the generative model with latent $\mathbf{z}_i \in \mathbb{R}^K$ such that:

$$\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \quad \mathbf{u} \sim \mathcal{N}(\boldsymbol{\Upsilon} \mathbf{z}, \boldsymbol{\Psi}_u) \quad \mathbf{v} \sim \mathcal{N}(\boldsymbol{\Phi} \mathbf{z}, \boldsymbol{\Psi}_v)$$

where we have $\boldsymbol{\Psi}_u \succeq 0$ and $\boldsymbol{\Psi}_v \succeq 0$. This is block diagonal noise. There are certain restriction of Gaussian FA and PCA as it is modelled the distribution that is too restrictive.

Definition 2.7. (Mixture Distribution) The mixture distribution has simple discrete latent variable:

$$s_i \sim \text{Discrete}[\boldsymbol{\pi}] \quad \mathbf{x}_i | s_i \sim P_{s_i}[\boldsymbol{\theta}_{s_i}]$$

The mixture can be seen as a mixture of multiple sources of data. The probability density of the single data point is given as:

$$p(\mathbf{x}_i) = \sum_{i=1}^k p(\mathbf{x}_i | s_i = m) p(s_i = m) = \sum_{i=1}^k \pi_m p(s_i = m)$$

The most notable mixture distribution is the mixture of Gaussian distribution.

Remark 22. Please note that once can perform a Bayesian inference to infer the probability that particular point \mathbf{x} belongs to the cluster m of the mixture distribution:

$$p(s_i = m | \mathbf{x}) = \frac{p_m(\mathbf{x}) \pi_m}{\sum_{i=1}^k p_i(\mathbf{x}) \pi_i}$$

Remark 23. (Mixture of Gaussian) Let's consider the mixture of Gaussian, where we have the following mixture distribution:

$$p(\{\mathbf{x}_i\}_{i=1}^N | \{\boldsymbol{\mu}_i\}_{i=1}^k, \{\boldsymbol{\Sigma}_i\}_{i=1}^k, \boldsymbol{\pi}) = \prod_{i=1}^n \sum_{m=1}^k \pi_m \frac{1}{\sqrt{|2\pi\boldsymbol{\Sigma}_m|}} \exp\left\{-\frac{1}{2}(\mathbf{x}_i - \boldsymbol{\mu}) \boldsymbol{\Sigma}_m^{-1} (\mathbf{x}_i - \boldsymbol{\mu})\right\}$$

Again it is hard to find the solution to the problem, and so we will consider the method to solve such the problem, which is called Expectation-Maximization (EM).

Remark 24. (Mixture of Factor Analyzers) Now, we consider the clustering and dimensionality reduction:

$$p(\mathbf{x} | \boldsymbol{\theta}) = \sum_{m=1}^k \pi_m \mathcal{N}(\boldsymbol{\mu}_k | \mathbf{W}_k \mathbf{W}_k^T + \boldsymbol{\Psi})$$

where π_k is the mixing proportion, while the parameter are $\boldsymbol{\theta} = \left\{ \{\pi_m, \boldsymbol{\mu}_m, \mathbf{W}_m\}_{m=1}^k, \boldsymbol{\Psi} \right\}$. Please note that this model has 2 kinds of latent variables, which are:

- Cluster indicator variable π_m for $m \in \{1, \dots, k\}$
- Continuous factor $\mathbf{z}_{im} \in \mathbb{R}^M$

Together giving us the following data generating distribution:

$$p(\mathbf{x} | \boldsymbol{\theta}) = \sum_{m=1}^k p(\pi_m) \int p(\mathbf{z}) p(\mathbf{x}_m | \mathbf{z}, \boldsymbol{\theta}) \, d\mathbf{z}$$

We can use EM to perform an optimization over it.

3 Expectation-Maximization

3.1 Methods

Remark 25. (General Form of The Objective) Now, we shall consider the latent variable model to be:

$$p(\mathbf{x} | \mathbf{z}, \boldsymbol{\theta}_x) = \frac{f_x(\mathbf{x}) \exp(\boldsymbol{\phi}(\boldsymbol{\theta}_x, \mathbf{z})^T \mathbf{T}_x(\mathbf{x}))}{Z_x(\boldsymbol{\phi}(\boldsymbol{\theta}_x, \mathbf{z}))} \quad p(\mathbf{z} | \boldsymbol{\theta}_x) = \frac{f_z(\mathbf{z}) \exp(\boldsymbol{\theta}_z^T \mathbf{T}_z(\mathbf{z}))}{Z_z(\boldsymbol{\theta}_z)}$$

and so the marginalization of the latent variable model is equal to:

$$p(\mathbf{x}|\boldsymbol{\theta}_x, \boldsymbol{\theta}_z) = \int \frac{f_x(\mathbf{x}) \exp(\boldsymbol{\phi}(\boldsymbol{\theta}_x, \mathbf{z})^T \mathbf{T}_x(\mathbf{x}))}{Z_x(\boldsymbol{\phi}(\boldsymbol{\theta}_x, \mathbf{z}))} \cdot \frac{f_z(\mathbf{z}) \exp(\boldsymbol{\theta}_z^T \mathbf{T}_z(\mathbf{z}))}{Z_z(\boldsymbol{\theta}_z)} d\mathbf{z}$$

And so, the log-likelihood over the dataset $\{\mathbf{x}_i\}_{i=1}^N$ is equal to:

$$l(\boldsymbol{\theta}_x, \boldsymbol{\theta}_z) = \sum_{i=1}^N \log \int \frac{f_x(\mathbf{x}_i) \exp(\boldsymbol{\phi}(\boldsymbol{\theta}_x, \mathbf{z})^T \mathbf{T}_x(\mathbf{x}_i))}{Z_x(\boldsymbol{\phi}(\boldsymbol{\theta}_x, \mathbf{z}))} \cdot \frac{f_z(\mathbf{z}) \exp(\boldsymbol{\theta}_z^T \mathbf{T}_z(\mathbf{z}))}{Z_z(\boldsymbol{\theta}_z)} d\mathbf{z}$$

Theorem 3.1. (Jensen's Inequality) Getting the set of weights $\{\alpha_i\}_{i=1}^N$ where $\sum \alpha_i = 1$ and $\{x_i > 0\}_{i=1}^N$, then we can show that:

$$\log \left(\sum_{i=1}^N \alpha_i x_i \right) \geq \sum_{i=1}^N \alpha_i \log(x_i)$$

for the concave f (and α_i is the probability measure). This also implies that: $f(\mathbb{E}_\alpha[x]) \geq \mathbb{E}_\alpha[f(x)]$ with the equality iff $f(x)$ is almost surely constant or linear support of α .

Definition 3.1. (Free Energy) Consider the latent variable model again: Given the observed data $\mathcal{X} = \{\mathbf{x}_i\}$ and the set of latent variables $\mathcal{Z} = \{\mathbf{z}_i\}$ with the parameter $\boldsymbol{\theta} = \{\boldsymbol{\theta}_x, \boldsymbol{\theta}_z\}$:

$$l(\boldsymbol{\theta}) = \log p(\mathcal{X}|\boldsymbol{\theta}) = \log \int p(\mathcal{Z}, \mathcal{X}|\boldsymbol{\theta}) d\mathcal{Z}$$

We will consider the Jensen's inequality for arbitrary distribution $q(\mathcal{Z})$. We can find the lower bound:

$$l(\boldsymbol{\theta}) = \log \int q(\mathcal{Z}) \frac{p(\mathcal{Z}, \mathcal{X}|\boldsymbol{\theta})}{q(\mathcal{Z})} d\mathcal{Z} \geq \int q(\mathcal{Z}) \log \frac{p(\mathcal{Z}, \mathcal{X}|\boldsymbol{\theta})}{q(\mathcal{Z})} d\mathcal{Z} = F(q, \boldsymbol{\theta})$$

We denote $F(q, \boldsymbol{\theta})$ to be a free energy.

Remark 26. (Other Form of Free Energy) We can consider the free energy to be:

$$\begin{aligned} \int q(\mathcal{Z}) \log \frac{p(\mathcal{Z}, \mathcal{X}|\boldsymbol{\theta})}{q(\mathcal{Z})} d\mathcal{Z} &= \int q(\mathcal{Z}) \log p(\mathcal{Z}, \mathcal{X}|\boldsymbol{\theta}) d\mathcal{Z} - \int q(\mathcal{Z}) \log q(\mathcal{Z}) d\mathcal{Z} \\ &= \int q(\mathcal{Z}) \log p(\mathcal{Z}, \mathcal{X}|\boldsymbol{\theta}) d\mathcal{Z} + H[q] \end{aligned}$$

and so the free-energy is equal to $F(q, \boldsymbol{\theta}) = \langle \log p(\mathcal{Z}, \mathcal{X}|\boldsymbol{\theta}) \rangle_{q(\mathcal{Z})} + H[q]$

Definition 3.2. (Expectation Maximization) The EM-Step follows the following step to be:

- *E-Step:* Optimize $F(q, \boldsymbol{\theta})$ with respected to $q(\mathcal{Z})$ as:

$$q^{(k)}(\mathcal{Z}) = \arg \max_{q(\mathcal{Z})} F(q(\mathcal{Z}), \boldsymbol{\theta}^{(k-1)})$$

- *M-Step:* Optimize $F(q, \boldsymbol{\theta})$ with respected to $\boldsymbol{\theta}$ as:

$$\boldsymbol{\theta}^{(k)} = \arg \max_{\boldsymbol{\theta}} F(q^{(k)}(\mathcal{Z}), \boldsymbol{\theta}) = \arg \max_{\boldsymbol{\theta}} \langle \log P(\mathcal{Z}, \mathcal{X}|\boldsymbol{\theta}) \rangle_{q^{(k)}(\mathcal{Z})}$$

Remark 27. (Simplification of E-Step) To consider the E-step, we have the following free-energy principle:

$$\begin{aligned} F(q, \boldsymbol{\theta}) &= \int q(\mathcal{Z}) \frac{p(\mathcal{Z}, \mathcal{X}|\boldsymbol{\theta})}{q(\mathcal{Z})} d\mathcal{Z} \\ &= \int q(\mathcal{Z}) \frac{p(\mathcal{Z}|\mathcal{X}, \boldsymbol{\theta}) p(\mathcal{X}|\boldsymbol{\theta})}{q(\mathcal{Z})} d\mathcal{Z} \\ &= \int q(\mathcal{Z}) \log p(\mathcal{X}|\boldsymbol{\theta}) d\mathcal{Z} + \int q(\mathcal{Z}) \log \frac{p(\mathcal{Z}|\mathcal{X}, \boldsymbol{\theta})}{q(\mathcal{Z})} d\mathcal{Z} \\ &= l(\boldsymbol{\theta}) - \text{KL}[q(\mathcal{Z})||p(\mathcal{Z}|\mathcal{X}, \boldsymbol{\theta})] \end{aligned}$$

where $l(\boldsymbol{\theta})$ is the log-likelihood and to minimize the $q(\mathcal{Z})$ by making KL-divergence equal to 0, when setting:

$$q^{(k)}(\mathcal{Z}) = p(\mathcal{Z}|\mathcal{X}, \boldsymbol{\theta}^{(k-1)})$$

Furthermore, we can see that after E-step the free energy is equal to likelihood.

Theorem 3.2. (Improvement of EM) *The results θ given by EM step never decrease in the likelihood.*

Proof. We have the following chain of inequalities:

$$l(\boldsymbol{\theta}^{(k-1)}) =_1 F(q^{(k)}, \boldsymbol{\theta}^{(k-1)}) \leq_2 F(q^{(k)}, \boldsymbol{\theta}^{(k)}) \leq_3 l(\boldsymbol{\theta}^{(k)})$$

Let's consider each step of the EM as we have:

1. The E-step gives free energy to the likelihood.
2. The M-step maximizes the free energy with respect to $\boldsymbol{\theta}$
3. $F \leq l$ comes from the lower-bound property of the free energy.

Thus the theorem is proved □

Theorem 3.3. (Convergence of EM to Optimum) *The fixed point of the EM algorithm is the stationary point of log-likelihood $l(\boldsymbol{\theta})$ and it is maxima.*

Proof. Please note that the fixed point of the EM is when (M-step with complete E-step):

$$\left. \frac{\partial}{\partial \boldsymbol{\theta}} \langle \log p(\mathcal{X}, \mathcal{Z}|\boldsymbol{\theta}) \rangle_{p(\mathcal{Z}|\mathcal{X}, \boldsymbol{\theta}^*)} \right|_{\boldsymbol{\theta}^*} = \mathbf{0}$$

Consider the log-likelihood, which we have:

$$\begin{aligned} l(\boldsymbol{\theta}) &= \log p(\mathcal{X}|\boldsymbol{\theta}) = \langle \log p(\mathcal{X}|\boldsymbol{\theta}) \rangle_{p(\mathcal{Z}|\mathcal{X}, \boldsymbol{\theta}^*)} \\ &= \left\langle \log \frac{p(\mathcal{Z}, \mathcal{X}|\boldsymbol{\theta})}{p(\mathcal{Z}|\mathcal{X}, \boldsymbol{\theta})} \right\rangle_{p(\mathcal{Z}|\mathcal{X}, \boldsymbol{\theta}^*)} \\ &= \langle p(\mathcal{Z}, \mathcal{X}|\boldsymbol{\theta}) \rangle_{p(\mathcal{Z}|\mathcal{X}, \boldsymbol{\theta}^*)} - \langle p(\mathcal{Z}|\mathcal{X}, \boldsymbol{\theta}) \rangle_{p(\mathcal{Z}|\mathcal{X}, \boldsymbol{\theta}^*)} \end{aligned}$$

Consider the derivative with respect to $\boldsymbol{\theta}$ of the log-likelihood evaluated at $\boldsymbol{\theta}$ and we have:

$$\begin{aligned} \left. \frac{\partial}{\partial \boldsymbol{\theta}} l(\boldsymbol{\theta}) \right|_{\boldsymbol{\theta}^*} &= \left. \frac{\partial}{\partial \boldsymbol{\theta}} \langle p(\mathcal{Z}, \mathcal{X}|\boldsymbol{\theta}) \rangle_{p(\mathcal{Z}|\mathcal{X}, \boldsymbol{\theta}^*)} \right|_{\boldsymbol{\theta}^*} - \left. \frac{\partial}{\partial \boldsymbol{\theta}} \langle p(\mathcal{Z}|\mathcal{X}, \boldsymbol{\theta}) \rangle_{p(\mathcal{Z}|\mathcal{X}, \boldsymbol{\theta}^*)} \right|_{\boldsymbol{\theta}^*} \\ &= \mathbf{0} - \left. \frac{\partial}{\partial \boldsymbol{\theta}} \langle p(\mathcal{Z}|\mathcal{X}, \boldsymbol{\theta}) \rangle_{p(\mathcal{Z}|\mathcal{X}, \boldsymbol{\theta}^*)} \right|_{\boldsymbol{\theta}^*} = \mathbf{0} \end{aligned}$$

The second term is equal to $\mathbf{0}$ as the KL-divergence is zero. Now, we are left to show that the second derivative is “negative”, where we have:

$$\left. \frac{\partial^2}{\partial^2 \boldsymbol{\theta}} l(\boldsymbol{\theta}) \right|_{\boldsymbol{\theta}^*} = \left. \frac{\partial^2}{\partial^2 \boldsymbol{\theta}} \langle p(\mathcal{Z}, \mathcal{X}|\boldsymbol{\theta}) \rangle_{p(\mathcal{Z}|\mathcal{X}, \boldsymbol{\theta}^*)} \right|_{\boldsymbol{\theta}^*} - \left. \frac{\partial^2}{\partial^2 \boldsymbol{\theta}} \langle p(\mathcal{Z}|\mathcal{X}, \boldsymbol{\theta}) \rangle_{p(\mathcal{Z}|\mathcal{X}, \boldsymbol{\theta}^*)} \right|_{\boldsymbol{\theta}^*}$$

The first term is negative due to the improvement of EM, now the second term is positive because it is the minimum of KL-divergence. And, the point $\boldsymbol{\theta}^*$ is minima. □

Definition 3.3. (Partial M-Step and E-Step) Let's start with the *M-Step* first, as long as we increase the $F(q, \boldsymbol{\theta})$ in the M-step, the theorem above still holds. And so, for partial M-step, after E-step, we have the following gradient update:

$$\left. \frac{\partial}{\partial \boldsymbol{\theta}} \right|_{\boldsymbol{\theta}^{(k-1)}} \langle \log p(\mathcal{Z}, \mathcal{X}|\boldsymbol{\theta}) \rangle_{p(\mathcal{Z}|\mathcal{X}, \boldsymbol{\theta})} = \left. \frac{\partial}{\partial \boldsymbol{\theta}} \right|_{\boldsymbol{\theta}^{(k-1)}} \log p(\mathcal{X}|\boldsymbol{\theta})$$

For *E-Step*, we can use the typical gradient based scheme on $q(\mathcal{Z})$ but the theorem above may not hold.

3.2 Practical

Remark 28. (Useful Derivative Results) Now, we are going to consider EM applied to the mixture of Gaussian that we have described above. There are some derivatives that would be useful in the future. Starting with the log-likelihood:

$$\begin{aligned} l(\{\boldsymbol{\mu}_i\}_{i=1}^k, \{\boldsymbol{\Sigma}_i\}_{i=1}^k, \boldsymbol{\pi}) &= \sum_{i=1}^N \log \sum_{m=1}^k \pi_m \underbrace{\frac{1}{\sqrt{|2\pi\boldsymbol{\Sigma}_m|}} \exp\left\{-\frac{1}{2}(\mathbf{x}_i - \boldsymbol{\mu})\boldsymbol{\Sigma}_m^{-1}(\mathbf{x}_i - \boldsymbol{\mu})\right\}}_{p_m(\mathbf{x}_i; \boldsymbol{\theta}_m)} \\ &= \sum_{i=1}^N \log \sum_{m=1}^k \pi_m p_m(\mathbf{x}_i; \boldsymbol{\theta}_m) \end{aligned}$$

where we denote $\boldsymbol{\theta}_m = \{\{\boldsymbol{\mu}_i\}_{i=1}^k, \{\boldsymbol{\Sigma}_i\}_{i=1}^k\}$. Now consider the derivatives of the log-likelihood with respect to both $\boldsymbol{\mu}_m$ and $\boldsymbol{\Sigma}_m$: Starting with outer derivatives, we have:

$$\begin{aligned} \frac{\partial}{\partial \boldsymbol{\theta}_m} \sum_{i=1}^N \log \sum_{m=1}^k \pi_m p_m(\mathbf{x}_i; \boldsymbol{\theta}_m) &= \sum_{i=1}^N \frac{\pi_m}{\sum_{m=1}^k \pi_m p_m(\mathbf{x}_i; \boldsymbol{\theta}_m)} \frac{\partial p_m(\mathbf{x}_i; \boldsymbol{\theta}_m)}{\partial \boldsymbol{\theta}_m} \\ &= \sum_{i=1}^N \underbrace{\frac{\pi_m p_m(\mathbf{x}_i; \boldsymbol{\theta}_m)}{\sum_{m=1}^k \pi_m p_m(\mathbf{x}_i; \boldsymbol{\theta}_m)}}_{\gamma_{im}} \frac{\partial \log p_m(\mathbf{x}_i; \boldsymbol{\theta}_m)}{\partial \boldsymbol{\theta}_m} \end{aligned}$$

Please remember γ_{im} as it is going to be useful later on. We are left to find the derivatives on the RHS, starting with the mean $\boldsymbol{\mu}_m$, which we can use the results from before:

$$\frac{\partial \log p_m(\mathbf{x}_i; \boldsymbol{\theta}_m)}{\partial \boldsymbol{\mu}_m} = \boldsymbol{\Sigma}_m^{-1}(\mathbf{x}_i - \boldsymbol{\mu}_m)$$

Similarly for $\boldsymbol{\Sigma}^{-1}$ we have:

$$\frac{\partial \log p_m(\mathbf{x}_i; \boldsymbol{\theta}_m)}{\partial \boldsymbol{\Sigma}_m} = \frac{1}{2} [\boldsymbol{\Sigma}_m - (\mathbf{x}_i - \boldsymbol{\mu}_m)(\mathbf{x}_i - \boldsymbol{\mu}_m)^T]$$

We can plug both into get the full derivative.

Proposition 3.1. (E-Step for mixture of Gaussian) *The E-step for mixture of Gaussian is done by setting:*

$$\gamma_{im}^{(t)} = \frac{\pi_m p_m(\mathbf{x}_i; \boldsymbol{\theta}_m)}{\sum_{m=1}^k \pi_m p_m(\mathbf{x}_i; \boldsymbol{\theta}_m)}$$

This follows directly from the Bayes' rule.

Proposition 3.2. (M-Step for mixture of Gaussian) *The M-step for mixture of Gaussian is done by setting:*

$$\begin{aligned} \boldsymbol{\mu}_m^{(t)} &= \frac{1}{N_m} \sum_{i=1}^N \gamma_{im}^{(t)} \mathbf{x}_i \\ \boldsymbol{\Sigma}_m^{(t)} &= \frac{1}{N_m} \sum_{i=1}^N \gamma_{im}^{(t)} (\mathbf{x}_i - \boldsymbol{\mu}_m)(\mathbf{x}_i - \boldsymbol{\mu}_m)^T \\ \pi_m^{(t)} &= \frac{N_m}{N} \end{aligned}$$

where $N_m = \sum_{i=1}^N \gamma_{im}^{(t)}$.

Proof. Setting the derivative of log-likelihood to zero for $\boldsymbol{\mu}_m$ and $\boldsymbol{\Sigma}_m$, which are:

$$\frac{\partial l}{\partial \boldsymbol{\mu}_m} = \sum_{i=1}^N \gamma_{im} \boldsymbol{\Sigma}_m^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_m) \quad \frac{\partial l}{\partial \boldsymbol{\Sigma}_m} = \frac{1}{2} \sum_{i=1}^N \gamma_{im} \left[\boldsymbol{\Sigma}_m - (\mathbf{x}_i - \boldsymbol{\mu}_m)(\mathbf{x}_i - \boldsymbol{\mu}_m)^T \right]$$

with simple rearrangement, we have the required update. We are left with $\pi_m^{(t)}$ as we required the constraint so that $\sum_{m=1}^k \pi_m^{(t)} = 1$. We have the following Lagrangian:

$$\begin{aligned} \frac{\partial}{\partial \pi_m} \left(\sum_{i=1}^N \log \sum_{m=1}^k \pi_m p_m(\mathbf{x}_i; \boldsymbol{\theta}_m) - \lambda \left[\sum_{m=1}^k \pi_m - 1 \right] \right) \\ = \sum_{i=1}^N \frac{p_m(\mathbf{x}_i; \boldsymbol{\theta}_m)}{\sum_{m=1}^k \pi_m p_m(\mathbf{x}_i; \boldsymbol{\theta}_m)} + \lambda \end{aligned}$$

Setting this to zero, and we have:

$$\begin{aligned} \sum_{i=1}^N \frac{p_m(\mathbf{x}_i; \boldsymbol{\theta}_m)}{\sum_{m=1}^k \pi_m p_m(\mathbf{x}_i; \boldsymbol{\theta}_m)} + \lambda &= 0 \\ \Rightarrow \sum_{m=1}^k \sum_{i=1}^N \frac{p_m(\mathbf{x}_i; \boldsymbol{\theta}_m) \pi_m}{\sum_{m=1}^k \pi_m p_m(\mathbf{x}_i; \boldsymbol{\theta}_m)} + \lambda \sum_{m=1}^k \pi_m &= 0 \\ \Rightarrow N + \lambda &= 0 \end{aligned}$$

and so, we have $\lambda = -N$. Now we simply times π_k on both side of the original derivative and we get the result. \square

Remark 29. (Connection between K-Mean and Gaussian Mixture Model) If we consider $\phi_m = 1/k$ and $\boldsymbol{\Sigma}_m = \sigma^2 \mathbf{I}$ where $\sigma^2 \rightarrow 0$, the leads us to the responsibility:

$$r_{im} = \delta \left(m, \arg \min_l \|\mathbf{x}_i - \boldsymbol{\mu}_l\|^2 \right)$$

We find the data in which it is closest to the mean, which we allocate the weight to be 1 and 0 for all others. Now, we simply have to update the mean (M-step), which is given by:

$$\boldsymbol{\mu}_m = \frac{\sum_i \gamma_{im} \mathbf{x}_i}{\sum_i \gamma_{im}}$$

This is exactly the K-mean algorithm.

Remark 30. (EM for Factor Analysis) Now, we are interested in training the factor analysis using the EM algorithm. Recall that factor analysis consists of the following model:

$$p(\mathbf{z}) = \mathcal{N}(\mathbf{z} | \mathbf{0}, \mathbf{I}) \quad p(\mathbf{x} | \mathbf{z}) = \mathcal{N}(\mathbf{x} | \mathbf{W}\mathbf{z}, \boldsymbol{\Psi})$$

Furthermore, we are going to assume that $\boldsymbol{\Psi} \in \mathbb{R}^{D \times D}$ where it is a diagonal matrix with values $\{\psi_{ii}\}_{i=1}^D$. And so recall that the marginal distribution, which is:

$$p(\mathbf{x} | \boldsymbol{\theta}) = \mathcal{N}(\mathbf{x} | \mathbf{0}, \boldsymbol{\Psi} + \mathbf{W}\mathbf{W}^T)$$

We can see that the model parameter is $\boldsymbol{\theta} = \{\mathbf{W}, \boldsymbol{\Psi}\}$. Now let's recall the EM algorithm, and see what we have to do:

- E-Step: We have to find the posterior distribution of each latent variable that corresponds to each data point in dataset $\{\mathbf{x}_i\}_{i=1}^N$:

$$q_i^{(t)}(\mathbf{z}_i) = p(\mathbf{z}_i | \mathbf{x}_i, \boldsymbol{\theta}^{(t-1)})$$

- M-Step: We have to find $\boldsymbol{\theta}^{(t)}$ as we have:

$$\arg \max_{\boldsymbol{\theta}} F(q_i^{(t)}, \boldsymbol{\theta}) = \arg \max_{\boldsymbol{\theta}} \sum_{i=1}^N \int q_i^{(t)}(\mathbf{z}_i) \left[\log p(\mathbf{z}_i | \boldsymbol{\theta}) + \log p(\mathbf{x}_i | \mathbf{z}_i, \boldsymbol{\theta}) \right] d\mathbf{z}_i$$

Proposition 3.3. (E-Step for Factor Analysis) Finding the posterior of \mathbf{z}_i for the dataset $\{\mathbf{x}_i\}_{i=1}^N$, is:

$$p(\mathbf{z}_i | \mathbf{x}_i) = \mathcal{N}(\mathbf{z}_i | \boldsymbol{\Sigma}^{-1} \mathbf{W}^T \boldsymbol{\Psi}^{-1}(\mathbf{x}_i, \boldsymbol{\Sigma}^{-1}))$$

where $\boldsymbol{\Sigma} = \mathbf{I} + \mathbf{W}^T \boldsymbol{\Psi}^{-1} \mathbf{W}$. This is the old result from PPCA that we have derived earlier. We denote the mean to be $\boldsymbol{\mu}_i = \boldsymbol{\Sigma}^{-1} \mathbf{W}^T \boldsymbol{\Psi}^{-1}(\mathbf{x}_i)$.

Proposition 3.4. (M-Step for Factor Analysis) The M-Step update of the factor analysis is given by:

$$\mathbf{W} = \left(\sum_{i=1}^N \mathbf{x}_i \boldsymbol{\mu}_i^T \right) \left(\sum_{i=1}^N \boldsymbol{\mu}_i \boldsymbol{\mu}_i^T + N \boldsymbol{\Sigma} \right)^{-1} \quad \boldsymbol{\Psi} = \mathbf{W} \boldsymbol{\Sigma} \mathbf{W}^T + \frac{1}{N} \sum_{i=1}^N (\mathbf{x}_i - \mathbf{W} \boldsymbol{\mu}_i)(\mathbf{x}_i - \mathbf{W} \boldsymbol{\mu}_i)^T$$

Proof. M-step of the Factor analysis, we are interested into find the variables $\boldsymbol{\theta}$ such that

$$\arg \max_{\boldsymbol{\theta}} F(q, \boldsymbol{\theta}) = \sum_{i=1}^N \langle \log p(\mathbf{z}_i | \boldsymbol{\theta}) + \log p(\mathbf{x}_i | \mathbf{z}_i, \boldsymbol{\theta}) \rangle_{q_i(\mathbf{z}_i)}$$

Now let's consider the sum of the log, and so we have:

$$\begin{aligned} \log p(\mathbf{z}_i | \boldsymbol{\theta}) + \log p(\mathbf{x}_i | \mathbf{z}_i, \boldsymbol{\theta}) &= -\frac{1}{2} \log |\boldsymbol{\Psi}| - \frac{1}{2} (\mathbf{x}_i - \mathbf{W} \mathbf{z}_i)^T \boldsymbol{\Psi}^{-1} (\mathbf{x}_i - \mathbf{W} \mathbf{z}_i) + \text{const.} \\ &= -\frac{1}{2} \log |\boldsymbol{\Psi}| - \frac{1}{2} \left[\mathbf{x}_i^T \boldsymbol{\Psi}^{-1} \mathbf{x}_i - 2 \mathbf{x}_i^T \boldsymbol{\Psi}^{-1} \mathbf{W} \mathbf{z}_i + \mathbf{z}_i^T \mathbf{W}^T \boldsymbol{\Psi}^{-1} \mathbf{W} \mathbf{z}_i \right] + \text{const.} \\ &= -\frac{1}{2} \log |\boldsymbol{\Psi}| - \frac{1}{2} \left[\mathbf{x}_i^T \boldsymbol{\Psi}^{-1} \mathbf{x}_i - 2 \mathbf{x}_i^T \boldsymbol{\Psi}^{-1} \mathbf{W} \mathbf{z}_i + \text{Tr}[\mathbf{W}^T \boldsymbol{\Psi}^{-1} \mathbf{W} \mathbf{z}_i \mathbf{z}_i^T] \right] + \text{const.} \end{aligned}$$

Now, we consider the expectation over $q_i^{(t)}(\mathbf{z}_i)$ as we now have:

$$\begin{aligned} &\langle \log p(\mathbf{z}_i | \boldsymbol{\theta}) + \log p(\mathbf{x}_i | \mathbf{z}_i, \boldsymbol{\theta}) \rangle_{q_i^{(t)}(\mathbf{z}_i)} \\ &= \left\langle -\frac{1}{2} \log |\boldsymbol{\Psi}| - \frac{1}{2} \left[\mathbf{x}_i^T \boldsymbol{\Psi}^{-1} \mathbf{x}_i - 2 \mathbf{x}_i^T \boldsymbol{\Psi}^{-1} \mathbf{W} \mathbf{z}_i + \text{Tr}[\mathbf{W}^T \boldsymbol{\Psi}^{-1} \mathbf{W} \mathbf{z}_i \mathbf{z}_i^T] \right] \right\rangle_{q_i^{(t)}(\mathbf{z}_i)} \\ &= -\frac{1}{2} \log |\boldsymbol{\Psi}| - \frac{1}{2} \left[\mathbf{x}_i^T \boldsymbol{\Psi}^{-1} \mathbf{x}_i - 2 \mathbf{x}_i^T \boldsymbol{\Psi}^{-1} \mathbf{W} \langle \mathbf{z}_i \rangle_{q_i^{(t)}(\mathbf{z}_i)} + \text{Tr}[\mathbf{W}^T \boldsymbol{\Psi}^{-1} \mathbf{W} \langle \mathbf{z}_i \mathbf{z}_i^T \rangle_{q_i^{(t)}(\mathbf{z}_i)}] \right] \\ &= -\frac{1}{2} \log |\boldsymbol{\Psi}| - \frac{1}{2} \left[\mathbf{x}_i^T \boldsymbol{\Psi}^{-1} \mathbf{x}_i - 2 \mathbf{x}_i^T \boldsymbol{\Psi}^{-1} \mathbf{W} \boldsymbol{\mu}_i + \text{Tr}[\mathbf{W}^T \boldsymbol{\Psi}^{-1} \mathbf{W} (\boldsymbol{\mu}_i \boldsymbol{\mu}_i^T + \boldsymbol{\Sigma})] \right] \end{aligned}$$

Using the expectation over the Gaussian distribution to be $\langle \mathbf{z}_i \rangle = \boldsymbol{\mu}_i$ and $\langle \mathbf{z}_i \mathbf{z}_i^T \rangle = \boldsymbol{\mu}_i \boldsymbol{\mu}_i^T + \boldsymbol{\Sigma}$. Now, let's consider the derivative with respect to \mathbf{W} and $\boldsymbol{\Psi}^{-1}$:

$$\begin{aligned} &\frac{1}{2} \frac{\partial}{\partial \mathbf{W}} \sum_{i=1}^N \left[2 \mathbf{x}_i^T \boldsymbol{\Psi}^{-1} \mathbf{W} \boldsymbol{\mu}_i - \text{Tr}[\mathbf{W}^T \boldsymbol{\Psi}^{-1} \mathbf{W} (\boldsymbol{\mu}_i \boldsymbol{\mu}_i^T + \boldsymbol{\Sigma})] \right] \\ &= \frac{1}{2} \sum_{i=1}^N \left[2 \frac{\partial}{\partial \mathbf{W}} \text{Tr}[\mathbf{W} \boldsymbol{\mu}_i \mathbf{x}_i^T \boldsymbol{\Psi}^{-1}] - \frac{\partial}{\partial \mathbf{W}} \text{Tr}[\mathbf{W}^T \boldsymbol{\Psi}^{-1} \mathbf{W} (\boldsymbol{\mu}_i \boldsymbol{\mu}_i^T + \boldsymbol{\Sigma})] \right] \\ &= \frac{1}{2} \sum_{i=1}^N \left[2 \boldsymbol{\Psi}^{-1} \mathbf{x}_i \boldsymbol{\mu}_i^T - 2 (\boldsymbol{\Psi}^{-1} \mathbf{W} (\boldsymbol{\mu}_i \boldsymbol{\mu}_i^T + \boldsymbol{\Sigma})) \right] \\ &= \boldsymbol{\Psi}^{-1} \left[\sum_{i=1}^N \mathbf{x}_i \boldsymbol{\mu}_i^T - \mathbf{W} \left(\sum_{i=1}^N \boldsymbol{\mu}_i \boldsymbol{\mu}_i^T + N \boldsymbol{\Sigma} \right) \right] \end{aligned}$$

Setting the derivative to zero and we yields:

$$\mathbf{W} = \left(\sum_{i=1}^N \mathbf{x}_i \boldsymbol{\mu}_i^T \right) \left(\sum_{i=1}^N \boldsymbol{\mu}_i \boldsymbol{\mu}_i^T + N \boldsymbol{\Sigma} \right)^{-1}$$

Let's consider the case for $\boldsymbol{\Psi}^{-1}$, which we have:

$$\begin{aligned} \frac{\partial}{\partial \boldsymbol{\Psi}^{-1}} \sum_{i=1}^N \left(-\frac{1}{2} \log |\boldsymbol{\Psi}| - \frac{1}{2} \left[\mathbf{x}_i^T \boldsymbol{\Psi}^{-1} \mathbf{x}_i - 2 \mathbf{x}_i^T \boldsymbol{\Psi}^{-1} \mathbf{W} \boldsymbol{\mu}_i + \text{Tr}[\mathbf{W}^T \boldsymbol{\Psi}^{-1} \mathbf{W} (\boldsymbol{\mu}_i \boldsymbol{\mu}_i^T + \boldsymbol{\Sigma})] \right] \right) \\ = \frac{1}{2} \sum_{i=1}^N \left(-\frac{\partial}{\partial \boldsymbol{\Psi}^{-1}} \log |\boldsymbol{\Psi}| - \frac{\partial}{\partial \boldsymbol{\Psi}^{-1}} \mathbf{x}_i^T \boldsymbol{\Psi}^{-1} \mathbf{x}_i + 2 \frac{\partial}{\partial \boldsymbol{\Psi}^{-1}} \mathbf{x}_i^T \boldsymbol{\Psi}^{-1} \mathbf{W} \boldsymbol{\mu}_i \right. \\ \left. - \frac{\partial}{\partial \boldsymbol{\Psi}^{-1}} \text{Tr}[\mathbf{W}^T \boldsymbol{\Psi}^{-1} \mathbf{W} (\boldsymbol{\mu}_i \boldsymbol{\mu}_i^T + \boldsymbol{\Sigma})] \right) \\ = \frac{1}{2} \sum_{i=1}^N \left(\boldsymbol{\Psi} - \mathbf{x}_i \mathbf{x}_i^T + 2 \mathbf{x}_i \boldsymbol{\mu}_i^T \mathbf{W}^T - \mathbf{W} \boldsymbol{\mu}_i \boldsymbol{\mu}_i^T \mathbf{W}^T - \mathbf{W} \boldsymbol{\Sigma} \mathbf{W}^T \right) \end{aligned}$$

Setting the derivative to zero, and we have:

$$\begin{aligned} \frac{1}{2} \sum_{i=1}^N \left(\boldsymbol{\Psi} - \mathbf{x}_i \mathbf{x}_i^T + 2 \mathbf{x}_i \boldsymbol{\mu}_i^T \mathbf{W}^T - \mathbf{W} \boldsymbol{\mu}_i \boldsymbol{\mu}_i^T \mathbf{W}^T - \mathbf{W} \boldsymbol{\Sigma} \mathbf{W}^T \right) &= 0 \\ \implies N \boldsymbol{\Psi} - N \mathbf{W} \boldsymbol{\Sigma} \mathbf{W}^T - \sum_{i=1}^N \left[\mathbf{x}_i \mathbf{x}_i^T - 2 \mathbf{x}_i \boldsymbol{\mu}_i^T \mathbf{W}^T + \mathbf{W} \boldsymbol{\mu}_i \boldsymbol{\mu}_i^T \mathbf{W}^T \right] &= 0 \\ \implies \boldsymbol{\Psi} = \mathbf{W} \boldsymbol{\Sigma} \mathbf{W}^T + \frac{1}{N} \sum_{i=1}^N (\mathbf{x}_i - \mathbf{W} \boldsymbol{\mu}_i) (\mathbf{x}_i - \mathbf{W} \boldsymbol{\mu}_i)^T \end{aligned}$$

Which is the result as required. □

3.3 Additional EM Methods

Remark 31. This is going to be more generalized version of EM as we may dealing with the exponential family, which occurs for both latent variable and the observable.

Proposition 3.5. *Given the exponential family of the form:*

$$p(\mathbf{x}|\boldsymbol{\theta}) = \frac{f(\mathbf{x}) \exp(\boldsymbol{\theta}^T \mathbf{T}(\mathbf{x}))}{Z(\boldsymbol{\theta})}$$

We can show that the

$$\frac{\partial}{\partial \boldsymbol{\theta}} \log Z(\boldsymbol{\theta}) = \langle \mathbf{T}(\mathbf{x}) \rangle_{p(\mathbf{x}|\boldsymbol{\theta})}$$

Proof.

$$\begin{aligned} \frac{\partial}{\partial \boldsymbol{\phi}} \log Z(\boldsymbol{\phi}) &= \frac{1}{Z(\boldsymbol{\phi})} \frac{\partial}{\partial \boldsymbol{\phi}} Z(\boldsymbol{\phi}) \\ &= \frac{1}{Z(\boldsymbol{\phi})} \frac{\partial}{\partial \boldsymbol{\phi}} \iint f(\mathbf{x}, \mathbf{z}) \exp \left\{ \boldsymbol{\theta}^T \mathbf{T}(\mathbf{x}) \right\} d\mathbf{x} d\mathbf{z} \\ &= \iint \frac{1}{Z(\boldsymbol{\phi})} f(\mathbf{x}, \mathbf{z}) \frac{\partial}{\partial \boldsymbol{\phi}} \exp \left\{ \boldsymbol{\theta}^T \mathbf{T}(\mathbf{x}) \right\} d\mathbf{x} d\mathbf{z} \\ &= \iint \frac{1}{Z(\boldsymbol{\phi})} f(\mathbf{x}, \mathbf{z}) \exp \left\{ \boldsymbol{\theta}^T \mathbf{T}(\mathbf{x}) \right\} \mathbf{T}(\mathbf{x}) d\mathbf{x} d\mathbf{z} \\ &= \langle \mathbf{T}(\mathbf{x}) \rangle_{p(\mathbf{x}|\boldsymbol{\theta})} \end{aligned}$$

□

Definition 3.4. (EM For Exponential Family) We now consider the *joint* probability distribution over \mathbf{x} and \mathbf{z} to be an exponential family, now we have:

$$p(\mathbf{x}, \mathbf{z} | \phi) = \frac{f(\mathbf{x}, \mathbf{z}) \exp\left(\phi(\boldsymbol{\theta}_x, \boldsymbol{\theta}_z)^T \mathbf{T}(\mathbf{x}, \mathbf{z})\right)}{Z(\phi)}$$

Let's consider the free-energy (consider only terms that depends on ϕ), using the previous result:

$$\begin{aligned} F(q, \phi) &= \int q(\mathbf{z}) \log P(\mathbf{x}, \mathbf{z} | \phi) \, d\mathbf{z} + H[q] \\ &= \int q(\mathbf{z}) \left[\phi(\boldsymbol{\theta}_x, \boldsymbol{\theta}_z)^T \mathbf{T}(\mathbf{x}, \mathbf{z}) - \log Z(\phi) \right] \, d\mathbf{z} + \text{const.} \\ &= \phi(\boldsymbol{\theta}_x, \boldsymbol{\theta}_z)^T \langle \mathbf{T}(\mathbf{x}, \mathbf{z}) \rangle_{q(\mathbf{z})} - \log Z(\phi) \end{aligned}$$

For the EM algorithm, we have the following:

- For *E-Step*, we only need to compute the sufficient statistics under the distribution $q(\mathbf{z})$.
- For *M-step*, we solve for the following equation, where we use the derivative from above:

$$\frac{\partial F}{\partial \phi} = \langle \mathbf{T}(\mathbf{x}, \mathbf{z}) \rangle_{q(\mathbf{z})} - \langle \mathbf{T}(\mathbf{x}, \mathbf{z}) \rangle_{p(\mathbf{x}, \mathbf{z} | \phi)} = \mathbf{0}$$

Remark 32. (EM For Exponential Family Mixture) We consider a short-hand version where we consider the 1-hot vector i.e for the mixture component m , we have:

$$\mathbf{s}_i = m \iff \mathbf{s}_i = [0, 0, \dots, \underbrace{1}_{m\text{-th position}}, \dots, 0]^T$$

To compute the components, we consider the “list” of parameters, which is a matrix $\Theta = [\boldsymbol{\theta}_m]$ (has the size of $\mathbb{R}^{D \times M}$). The log-likelihood is the given by:

$$\begin{aligned} \log p\left(\{\mathbf{x}_i\}_{i=1}^N, \{\mathbf{s}_i\}_{i=1}^N\right) &= \sum_{i=1}^N \log \left(\pi_{i s_i} \left[\frac{f(\mathbf{x}_i) \exp\left(\boldsymbol{\theta}_{s_i}^T \mathbf{T}(\mathbf{x}_i)\right)}{Z(\boldsymbol{\theta}_{s_i})} \right] \right) \\ &= \sum_{i=1}^N \left[\log \pi_{i s_i} + \boldsymbol{\theta}_{s_i}^T \mathbf{T}(\mathbf{x}_i) - \log Z(\boldsymbol{\theta}_{s_i}) \right] + \text{const} \\ &= \sum_{i=1}^N \left[(\log \boldsymbol{\pi})^T \mathbf{s}_i + (\Theta \mathbf{s}_i)^T \mathbf{T}(\mathbf{x}_i) - \mathbf{s}_i^T \log Z(\Theta) \right] + \text{const} \end{aligned}$$

where $Z(\Theta)$ is the normalizing factor of each component, so it has the size of M . If we were to consider the EM-algorithm for this, we have:

- E-step, we calculate the expectation over the latent variable, which gives us:

$$\gamma_{im} = \sum_{i=1}^N \langle \mathbf{s}_i \rangle_p \quad \sum_{i=1}^N \mathbf{T}(\mathbf{x}_i) \langle \mathbf{s}_i^T \rangle_p$$

- M-step, maximizing the log-joint probability, we have the following solutions (to the derivative):

$$\boldsymbol{\pi}^{(k+1)} \propto \sum_{i=1}^N \langle \mathbf{s}_i \rangle_p \quad \left\langle \mathbf{T}(\mathbf{x}_i) \middle| \boldsymbol{\theta}_m^{(k+1)} \right\rangle = \frac{\sum_{i=1}^N \mathbf{T}(\mathbf{x}_i) \langle [\mathbf{s}_i]_m \rangle_p}{\sum_{i=1}^N \langle [\mathbf{s}_i]_m \rangle_p}$$

Thinking of this as the Gaussian mixture model, where we also have similar form.

Remark 33. (EM for MAP) Recall the probability distribution over the exponential model; however, we also added the prior over the parameters:

$$p(\mathbf{x}, \mathbf{z} | \phi) = \frac{f(\mathbf{x}, \mathbf{z}) \exp\left(\phi(\boldsymbol{\theta}_x, \boldsymbol{\theta}_z)^T \mathbf{T}(\mathbf{x}, \mathbf{z})\right)}{Z(\phi)} \quad p(\phi | \nu, \boldsymbol{\tau}) = \frac{F(\nu, \boldsymbol{\tau}) \exp\left(\phi(\boldsymbol{\theta}_x, \boldsymbol{\theta}_z)^T \boldsymbol{\tau}\right)}{Z(\phi)^\nu}$$

To consider the free-energy over the dataset \mathcal{X} and set of latent variables \mathcal{Z} , we have:

$$\begin{aligned} F(q, \phi) &= \int q(\mathbf{z}) \log P(\mathcal{X}, \mathcal{Z}, \phi) \, d\mathbf{z} + H[q] \\ &= \int q(\mathcal{Z}) \left[\phi(\boldsymbol{\theta}_x, \boldsymbol{\theta}_z)^T \left(\boldsymbol{\tau} + \sum_{i=1}^N \mathbf{T}(\mathbf{x}_i, \mathbf{z}_i) \right) - (N + \nu) \log Z(\phi) \right] \, d\mathcal{Z} + \text{const.} \\ &= \phi(\boldsymbol{\theta}_x, \boldsymbol{\theta}_z)^T \left(\boldsymbol{\tau} + \sum_{i=1}^N \langle \mathbf{T}(\mathbf{x}_i, \mathbf{z}_i) \rangle_{q(\mathbf{z}_i)} \right) - (N + \nu) \log Z(\phi) + \text{const.} \end{aligned}$$

Now, we carry on with the EM step as usual.

4 Latent Variable Models for Time Series

4.1 Problems

Definition 4.1. (Time Series with Latent Variable) The latent variable model for time series (with first order markov model) is given by:

$$p(\mathbf{z}_{1:T}, \mathbf{x}_{1:T}) = p(\mathbf{z}_1) p(\mathbf{x}_1 | \mathbf{z}_1) \prod_{i=2}^T p(\mathbf{z}_i | \mathbf{z}_{i-1}) p(\mathbf{x}_i | \mathbf{z}_i)$$

Definition 4.2. (Discrete Model) The discrete model is described using the following set of parameters as there are K number of latent discrete values:

- Initial State Probability $p(\mathbf{z}_1)$: is given as $\pi_z = p(\mathbf{z}_1 = j)$
- Transition Matrix $p(\mathbf{z}_{t+1} | \mathbf{z}_t)$: is given as $\Phi_{ij} = p(\mathbf{z}_{t+1} = j | \mathbf{z}_t = i)$
- Emission Distribution $p(\mathbf{x}_t | \mathbf{z}_t)$:
 - Continuous variable: $\mathbf{A}_j(\mathbf{x}) = p(\mathbf{x}_t = \mathbf{x} | \mathbf{z}_t = j)$
 - Discrete variable: $\mathbf{A}_{jk}(\mathbf{x}) = p(\mathbf{x}_t = k | \mathbf{z}_t = j)$

And we have the following process for a discrete model:

$$\mathbf{z}_1 \sim \boldsymbol{\pi} \quad \mathbf{z}_{t+1} | \mathbf{z}_t \sim \boldsymbol{\Phi}_{\mathbf{z}_t} \quad \mathbf{x}_t | \mathbf{z}_t \sim \mathbf{A}_{\mathbf{z}_t}$$

Definition 4.3. (Continuous Variable) We consider the continuous variable model, which we model the time series with Gaussian and Linear transformation, which we have the following process:

$$\mathbf{z}_1 \sim \mathcal{N}(\boldsymbol{\mu}_0, \mathbf{Q}_0) \quad \mathbf{z}_t | \mathbf{z}_{t-1} \sim \mathcal{N}(\mathbf{A} \mathbf{z}_{t-1}, \mathbf{Q}) \quad \mathbf{x}_t | \mathbf{z}_t \sim \mathcal{N}(\mathbf{C} \mathbf{z}_t, \mathbf{R})$$

We still use the same joint distribution as above. This has the other name called Linear Gaussian State Space Model (LGSSM).

Definition 4.4. (Inference Problems) We are interesting to find the following joint distributions:

$$p(\mathbf{z}_t|\mathbf{x}_{1:t}) \quad p(\mathbf{z}_t|\mathbf{x}_{1:T})$$

The first one is called filtering distribution, while the second one is called smoothing distribution.

Remark 34. (Problem for the Problem) It is clear that we can integrate out the posterior to get the answer for filtering (and smoothing), but it is very hard:

$$p(\mathbf{z}|\mathbf{x}_1, \dots, \mathbf{x}_t) = \int \cdots \int p(\mathbf{z}_1, \dots, \mathbf{z}_t|\mathbf{x}_1, \dots, \mathbf{x}_t) d\mathbf{z}_1 \cdots d\mathbf{z}_t$$

For the discrete variable, we can have sum instead of integration.

4.2 Solving Inference Problem

Lemma 4.1. (Bayesian Filtering) One can compute the filtering recursively as:

$$p(\mathbf{z}_t|\mathbf{x}_{1:t}) \propto \int p(\mathbf{x}_t|\mathbf{z}_t)p(\mathbf{z}_t|\mathbf{z}_{t-1})p(\mathbf{z}_{t-1}|\mathbf{x}_{1:t-1}) d\mathbf{z}_{t-1}$$

We have the same treatment for discrete variable.

Proof. We have the following:

$$\begin{aligned} p(\mathbf{z}_t|\mathbf{x}_{1:t}) &= \int p(\mathbf{z}_t, \mathbf{z}_{t-1}|\mathbf{x}_t, \mathbf{x}_{1:t-1}) d\mathbf{z}_{t-1} \\ &= \int \frac{p(\mathbf{z}_t, \mathbf{z}_{t-1}|\mathbf{x}_{1:t-1})}{p(\mathbf{x}_t|\mathbf{x}_{1:t-1})} d\mathbf{z}_{t-1} \\ &\propto \int p(\mathbf{x}_t|\mathbf{z}_t, \mathbf{z}_{t-1}, \mathbf{x}_{1:t-1})p(\mathbf{z}_t|\mathbf{z}_{t-1}, \mathbf{x}_{1:t-1})p(\mathbf{z}_{t-1}|\mathbf{x}_{1:t-1}) d\mathbf{z}_{t-1} \\ &= \int p(\mathbf{x}_t|\mathbf{z}_t)p(\mathbf{z}_t|\mathbf{z}_{t-1})p(\mathbf{z}_{t-1}|\mathbf{x}_{1:t-1}) d\mathbf{z}_{t-1} \end{aligned}$$

And we have the values as required. □

Definition 4.5. (Forward Message) We consider the quantity $\alpha_t(i)$ for which it is the joint probability of observation and the current latent variable $\mathbf{z}_t = i$ i.e:

$$\alpha_t(i) = p(\mathbf{x}_1, \dots, \mathbf{x}_t, \mathbf{z}_t = i|\boldsymbol{\theta})$$

We call this a forward message.

Proposition 4.1. The forward message can be calculated recursively as:

$$\alpha_{t+1}(i) = \left(\sum_{j=1}^K \alpha_t(j)\Phi_{ij} \right) \mathbf{A}_i(\mathbf{x}_{t+1})$$

where at the start: $\alpha_1(i) = \pi_i \mathbf{A}_i(\mathbf{x})$

Proof. We consider the following:

$$\begin{aligned} \alpha_{t+1}(i) &= \sum_{j=1}^K p(\mathbf{x}_1, \dots, \mathbf{x}_t, \mathbf{x}_{t+1}, s_t = j | s_{t+1} = i) \\ &= \sum_{j=1}^K p(\mathbf{x}_1, \dots, \mathbf{x}_t | s_t = j) p(s_{t+1} = i | s_t = j) p(\mathbf{x}_{t+1} | s_{t+1} = i) \\ &= \left(\sum_{j=1}^K \alpha_t(j)\Phi_{ij} \right) \mathbf{A}_i(\mathbf{x}_{t+1}) \end{aligned}$$

□

Remark 35. (Filtering For Discrete Model) One can solve the filtering problem, as we have:

$$p(\mathbf{s}_t = i | \mathbf{x}_1, \dots, \mathbf{x}_t, \boldsymbol{\theta}) = \frac{\alpha_t(i)}{\sum_{j=1}^K \alpha_t(j)}$$

Remark 36. (Usefulness of Filtering Method) Given the values of $\alpha_T(i)$, we can calculate the likelihood of the parameter $\boldsymbol{\theta}$ in $\mathcal{O}(TK^2)$ time instead of exponential time via the marginalization:

$$p(\mathbf{x}_1, \dots, \mathbf{x}_T | \boldsymbol{\theta}) = \sum_{k=1}^K \alpha_T(k)$$

Proposition 4.2. (First Step Filter) We can show that the $p(\mathbf{z}_1 | \mathbf{x}_{1:1})$ is equal to:

$$p(\mathbf{z}_1 | \mathbf{x}_{1:1}) = \mathcal{N}\left(\mathbf{z}_1 \mid \underbrace{\hat{\mathbf{z}}_1^0 + \mathbf{K}_1(\mathbf{x}_1 - \mathbf{C}\hat{\mathbf{z}}_1^0)}_{\mathbf{z}_1^1}, \underbrace{\hat{\mathbf{V}}_1^0 - \mathbf{K}_1\mathbf{C}\hat{\mathbf{V}}_1^0}_{\hat{\mathbf{V}}_1^1}\right)$$

where we set $\mathbf{K}_1 = \hat{\mathbf{V}}_1^0 \mathbf{C}^T [\mathbf{R} + \mathbf{C}\hat{\mathbf{V}}_1^0 \mathbf{C}^T]^{-1}$.

Proof. For the state space model, we start with $p(\mathbf{z}_1 | \mathbf{x}_1)$. We denote $\hat{\mathbf{z}}_1^0 = \boldsymbol{\mu}_0$ and $\hat{\mathbf{V}}_1^0 = \mathbf{Q}_0$. Now, we apply linear Gaussian model above to get the inference. Recall that

$$p(\mathbf{z}_1) = \mathcal{N}(\mathbf{z}_1 | \hat{\mathbf{z}}_1^0, \hat{\mathbf{V}}_1^0) \quad p(\mathbf{x}_1 | \mathbf{z}_1) = \mathcal{N}(\mathbf{x}_1 | \mathbf{C}\mathbf{z}_1, \mathbf{R})$$

Now, we have the linear Gaussian model as

$$p(\mathbf{z}_1 | \mathbf{x}_1) = \mathcal{N}\left(\mathbf{z}_1 \mid \boldsymbol{\Sigma} \left[\mathbf{C}^T \mathbf{R}^{-1} \mathbf{x}_1 + (\hat{\mathbf{V}}_1^0)^{-1} \hat{\mathbf{z}}_1^0 \right], \boldsymbol{\Sigma} \right)$$

where $\boldsymbol{\Sigma}$ is simplified using Woodbury identity:

$$\begin{aligned} \boldsymbol{\Sigma} &= \left[(\hat{\mathbf{V}}_1^0)^{-1} + \mathbf{C}^T \mathbf{R}^{-1} \mathbf{C} \right]^{-1} \\ &= \hat{\mathbf{V}}_1^0 - \hat{\mathbf{V}}_1^0 \mathbf{C}^T \left[\mathbf{R} + \mathbf{C}\hat{\mathbf{V}}_1^0 \mathbf{C}^T \right]^{-1} \mathbf{C}\hat{\mathbf{V}}_1^0 \\ &= \hat{\mathbf{V}}_1^0 - \mathbf{K}_1 \mathbf{C}\hat{\mathbf{V}}_1^0 \end{aligned}$$

Now, we consider the mean, as we now have:

$$\begin{aligned} &\left[\hat{\mathbf{V}}_1^0 - \mathbf{K}_1 \mathbf{C}\hat{\mathbf{V}}_1^0 \right] \left[\mathbf{C}^T \mathbf{R}^{-1} \mathbf{x}_1 + (\hat{\mathbf{V}}_1^0)^{-1} \hat{\mathbf{z}}_1^0 \right] \\ &= \hat{\mathbf{V}}_1^0 \mathbf{C}^T \mathbf{R}^{-1} \mathbf{x}_1 + \hat{\mathbf{z}}_1^0 - \mathbf{K}_1 \mathbf{C}\hat{\mathbf{V}}_1^0 \mathbf{C}^T \mathbf{R}^{-1} \mathbf{x}_1 - \mathbf{K}_1 \mathbf{C}\hat{\mathbf{z}}_1^0 \\ &= \hat{\mathbf{z}}_1^0 - \mathbf{K}_1 \mathbf{C}\hat{\mathbf{z}}_1^0 + (\hat{\mathbf{V}}_1^0 \mathbf{C}^T \mathbf{R}^{-1} - \mathbf{K}_1 \mathbf{C}\hat{\mathbf{V}}_1^0 \mathbf{C}^T \mathbf{R}^{-1}) \mathbf{x}_1 \end{aligned}$$

Let's consider the value, which we can show that it is indeed \mathbf{K}_1 , where we will consider $\mathbf{B} = \mathbf{C}\hat{\mathbf{V}}_1^0 \mathbf{C}^T$:

$$\begin{aligned} &\hat{\mathbf{V}}_1^0 \mathbf{C}^T \mathbf{R}^{-1} - \mathbf{K}_1 \mathbf{B} \mathbf{R}^{-1} = \mathbf{K}_1 \\ \Leftrightarrow &\hat{\mathbf{V}}_1^0 \mathbf{C}^T \mathbf{R}^{-1} = \mathbf{K}_1 + \mathbf{K}_1 \mathbf{B} \mathbf{R}^{-1} \\ \Leftrightarrow &\hat{\mathbf{V}}_1^0 \mathbf{C}^T \mathbf{R}^{-1} = \mathbf{K}_1 (\mathbf{I} + \mathbf{B} \mathbf{R}^{-1}) \\ \Leftrightarrow &\hat{\mathbf{V}}_1^0 \mathbf{C}^T \mathbf{R}^{-1} (\mathbf{I} + \mathbf{B} \mathbf{R}^{-1})^{-1} = \mathbf{K}_1 \\ \Leftrightarrow &\hat{\mathbf{V}}_1^0 \mathbf{C}^T \mathbf{R}^{-1} (\mathbf{I} + \mathbf{B} \mathbf{R}^{-1})^{-1} = \hat{\mathbf{V}}_1^0 \mathbf{C}^T [\mathbf{R} + \mathbf{B}]^{-1} \\ \Leftrightarrow &\mathbf{R}^{-1} (\mathbf{I} + \mathbf{B} \mathbf{R}^{-1})^{-1} = [\mathbf{R} + \mathbf{B}]^{-1} \\ \Leftrightarrow &[\mathbf{R} (\mathbf{I} + \mathbf{B} \mathbf{R}^{-1})]^{-1} = [\mathbf{R} + \mathbf{B}]^{-1} \end{aligned}$$

And so equation the is proven. \square

Proposition 4.3. (General Time Filtering) We want to find the value $p(\mathbf{z}_t|\mathbf{x}_{1:t})$, which we can show that it is equal to, given $p(\mathbf{z}_{t-1}|\mathbf{x}_{1:t-1}) = \mathcal{N}(\mathbf{z}_{t-1}|\hat{\mathbf{z}}_{t-1}^{t-1}, \hat{\mathbf{V}}_{t-1}^{t-1})$:

$$\mathcal{N}\left(\mathbf{z}_t \left| \hat{\mathbf{z}}_{t-1}^t + \mathbf{K}_t(\mathbf{x}_t - \mathbf{C}\hat{\mathbf{z}}_{t-1}^t), \hat{\mathbf{V}}_t^{t-1} - \mathbf{K}_t\mathbf{C}\hat{\mathbf{V}}_t^{t-1}\right.\right)$$

where $\hat{\mathbf{z}}_{t-1}^t = \mathbf{A}\hat{\mathbf{z}}_{t-1}^{t-1}$ and $\hat{\mathbf{V}}_t^{t-1} = \mathbf{Q} + \mathbf{A}\hat{\mathbf{V}}_{t-1}^{t-1}\mathbf{A}^T$ and $\mathbf{K}_t = \hat{\mathbf{V}}_t^{t-1}\mathbf{C}^T [\mathbf{R} + \mathbf{C}\hat{\mathbf{V}}_t^{t-1}\mathbf{C}^T]^{-1}$

Proof. Now, we want to find the probability $p(\mathbf{z}_t|\mathbf{x}_{1:t-1})$, first, which we can use the marginal distribution from the linear Gaussian model between $p(\mathbf{z}_t|\mathbf{z}_{t-1})$ and $p(\mathbf{z}_{t-1}|\mathbf{x}_{1:t-1})$, which we marginalize out \mathbf{z}_{t-1} , which we have:

$$\begin{aligned} & \int \mathcal{N}(\mathbf{z}_{t-1}|\hat{\mathbf{z}}_{t-1}^{t-1}, \hat{\mathbf{V}}_{t-1}^{t-1})\mathcal{N}(\mathbf{z}_t|\mathbf{A}\mathbf{z}_{t-1}, \mathbf{Q}) d\mathbf{z}_{t-1} \\ &= \mathcal{N}\left(\mathbf{z}_t \left| \mathbf{A}\hat{\mathbf{z}}_{t-1}^{t-1}, \mathbf{Q} + \mathbf{A}\hat{\mathbf{V}}_{t-1}^{t-1}\mathbf{A}^T\right.\right) \\ &= \mathcal{N}\left(\mathbf{z}_t \left| \hat{\mathbf{z}}_{t-1}^t, \hat{\mathbf{V}}_t^{t-1}\right.\right) \end{aligned}$$

This follows directly from linear Gaussian model. Now, we follows the same method above (as we now need to find distribution of \mathbf{z}_t given *addiitonal* information of \mathbf{x}_t), and so we have:

$$\mathcal{N}\left(\mathbf{z}_t \left| \hat{\mathbf{z}}_{t-1}^t + \mathbf{K}_t(\mathbf{x}_t - \mathbf{C}\hat{\mathbf{z}}_{t-1}^t), \hat{\mathbf{V}}_t^{t-1} - \mathbf{K}_t\mathbf{C}\hat{\mathbf{V}}_t^{t-1}\right.\right)$$

where $\mathbf{K}_t = \hat{\mathbf{V}}_t^{t-1}\mathbf{C}^T [\mathbf{R} + \mathbf{C}\hat{\mathbf{V}}_t^{t-1}\mathbf{C}^T]^{-1}$. Thus complete the prove. \square

Remark 37. (Bayesian Smoothing - For Discrete Model) To find the smoothing, we consider calculating marginal posterior as:

$$p(\mathbf{z}_t|\mathbf{x}_{1:T}) = \frac{p(\mathbf{z}_t, \mathbf{x}_{t+1:T}|\mathbf{x}_{1:t})}{p(\mathbf{x}_{t+1:T}|\mathbf{x}_{1:t})} = \frac{p(\mathbf{x}_{t+1:T}|\mathbf{z}_t)p(\mathbf{z}_t|\mathbf{x}_{1:t})}{p(\mathbf{x}_{t+1:T}|\mathbf{x}_{1:t})} = \frac{p(\mathbf{x}_{t+1:T}|\mathbf{z}_t)p(\mathbf{z}_t, \mathbf{x}_{1:t})}{p(\mathbf{z}_{1:T})}$$

We can see that the probability distribution in **blue** are the forward passing, while the probability distribution in **red** is the backward distribution.

Definition 4.6. (Backward Message) Consider the value

$$\beta_t(i) = p(\mathbf{x}_{t+1:T}|\mathbf{z}_t = i)$$

We call it a backward message.

Proposition 4.4. Backward message can be calculated in recursive manner:

$$\beta_t(i) = \sum_{j=1}^K \Phi_{ij} \mathbf{A}_j(\mathbf{x}_{t+1})\beta_{t+1}(j)$$

where $\beta_T(i) = 1/K$, we consider a uniform distribution over the state.

Proof. We consider the following

$$\begin{aligned} \beta_t(i) &= \sum_{j=1}^K p(s_{t+1} = j, \mathbf{x}_{t+1}, \mathbf{x}_{t+2:T}|\mathbf{z}_t = i) \\ &= \sum_{j=1}^K p(s_{t+1} = j|\mathbf{z}_t = i)p(\mathbf{x}_{t+1}|\mathbf{z}_t = i)p(\mathbf{x}_{t+2:T}|\mathbf{z}_{t+1} = j) \\ &= \sum_{j=1}^K \Phi_{ij} \mathbf{A}_j(\mathbf{x}_{t+1})\beta_{t+1}(j) \end{aligned}$$

And so it proposition is proven. \square

Definition 4.7. (Forward-Backward Algorithm) The smooth distribution can be calculated as:

$$\gamma_t(i) = p(s_t = i | \mathbf{x}_{1:T}) = \frac{p(s_t = i, \mathbf{x}_{1:t})p(\mathbf{x}_{t+1:T} | s_t = i)}{p(\mathbf{x}_{1:T})} = \frac{\alpha_t(i)\beta_t(i)}{\sum_{j=1}^k \alpha_t(j)\beta_t(j)}$$

Remark 38. (Problem with Choosing State) We consider the value of $\gamma_t(i)$ is computed using forward and backward algorithm. Choosing the state i_t^* with the largest $\gamma_t(i)$ is the best way since the path has the maximum expected number of correct state. However, this may leads to the path, which is impossible.

Definition 4.8. (Viterbi Decoding) The best path algorithm is called Viterbi decoding as we can use the Bellman's dynamics programming. This is done by compute the most probable state sequence (instead of expected state):

$$\arg \max_{\mathbf{z}_{1:T}} p(\mathbf{z}_{1:T} | \mathbf{x}_{1:T}, \boldsymbol{\theta})$$

And so, we use the same recursion max instead of \sum , just like the Bellman equation algorithm.

Proposition 4.5. (Smoothing for LGSSM) We receives the value $p(\mathbf{z}_{t+1} | \mathbf{x}_{1:T}) = \mathcal{N}(\mathbf{z}_{t+1} | \hat{\mathbf{z}}_{t+1}^T, \hat{\mathbf{V}}_{t+1}^T)$ from the previous time step. Then we can show the following recursive calculation:

$$p(\mathbf{z}_t | \mathbf{x}_{1:T}) = \mathcal{N}\left(\mathbf{z}_t \mid \hat{\mathbf{z}}_t + \mathbf{J}_t(\mathbf{z}_{t+1} - \mathbf{A}\hat{\mathbf{z}}_t^t), \hat{\mathbf{V}}_t^t + \mathbf{J}_t(\hat{\mathbf{V}}_{t+1}^T - \hat{\mathbf{V}}_{t+1}^t)\mathbf{J}_t^T\right)$$

where $\mathbf{J}_t = \hat{\mathbf{V}}_t^t \mathbf{A}^T (\hat{\mathbf{V}}_{t+1}^t)^{-1}$, and the first time step can be calculated from result of filtering.

Proof. We consider the following steps:

$$\begin{aligned} p(\mathbf{z}_t | \mathbf{x}_{1:T}) &= \int p(\mathbf{z}_t, \mathbf{z}_{t+1} | \mathbf{x}_{1:T}) d\mathbf{z}_{t+1} \\ &= \int p(\mathbf{z}_t | \mathbf{z}_{t+1}, \mathbf{x}_{1:T}) p(\mathbf{z}_{t+1} | \mathbf{x}_{1:T}) d\mathbf{z}_{t+1} \\ &= \int p(\mathbf{z}_t | \mathbf{z}_{t+1}, \mathbf{x}_{1:t}) p(\mathbf{z}_{t+1} | \mathbf{x}_{1:T}) d\mathbf{z}_{t+1} \\ &= \int \left[\frac{p(\mathbf{z}_{t+1}, \mathbf{z}_t | \mathbf{x}_{1:t})}{p(\mathbf{z}_{t+1} | \mathbf{x}_{1:t})} \right] p(\mathbf{z}_{t+1} | \mathbf{x}_{1:T}) d\mathbf{z}_{t+1} \end{aligned}$$

The third equality comes from the markov property. Let's consider finding the value $p(\mathbf{z}_t | \mathbf{z}_{t+1}, \mathbf{x}_{1:t})$ using linear Gaussian model as we have proven above: Noted that the value $p(\mathbf{z}_t | \mathbf{x}_{1:t}) = \mathcal{N}(\mathbf{z}_t | \hat{\mathbf{z}}_t^t, \hat{\mathbf{V}}_t^t)$. We consider the mean vector $[\mathbf{z}_t, \mathbf{z}_{t+1}]^T$, we will have to consider the covariance:

$$\begin{aligned} \text{Cov}(\mathbf{z}_t, \mathbf{z}_{t+1} | \mathbf{x}_{1:t}) &= \mathbb{E} \left[(\mathbf{z}_t - \hat{\mathbf{z}}_t^t) (\mathbf{z}_{t+1} - \hat{\mathbf{z}}_{t+1}^t)^T \mid \mathbf{x}_{1:t} \right] \\ &= \mathbb{E} \left[(\mathbf{z}_t - \hat{\mathbf{z}}_t^t) (\mathbf{A}\mathbf{z}_t + \mathbf{w} - \mathbf{A}\hat{\mathbf{z}}_t^t)^T \mid \mathbf{x}_{1:t} \right] \\ &= \mathbb{E} \left[(\mathbf{z}_t - \hat{\mathbf{z}}_t^t) (\mathbf{A}\mathbf{z}_t - \mathbf{A}\hat{\mathbf{z}}_t^t)^T \mid \mathbf{x}_{1:t} \right] + \mathbb{E} \left[(\mathbf{z}_t - \hat{\mathbf{z}}_t^t) \mid \mathbf{x}_{1:t} \right] \mathbb{E}[\mathbf{w}^T] \\ &= \mathbb{E} \left[(\mathbf{z}_t - \hat{\mathbf{z}}_t^t) (\mathbf{z}_t - \hat{\mathbf{z}}_t^t)^T \mid \mathbf{x}_{1:t} \right] \mathbf{A}^T = \hat{\mathbf{V}}_t^t \mathbf{A}^T \end{aligned}$$

where $p(\mathbf{w}) = \mathcal{N}(\mathbf{0}, \mathbf{Q})$. Now, we have the following normal distribution for the joint :

$$\mathcal{N} \left(\begin{bmatrix} \mathbf{z}_t \\ \mathbf{z}_{t+1} \end{bmatrix} \mid \begin{bmatrix} \hat{\mathbf{z}}_t^t \\ \hat{\mathbf{z}}_{t+1}^t \end{bmatrix}, \begin{bmatrix} \hat{\mathbf{V}}_t^t & \hat{\mathbf{V}}_t^t \mathbf{A}^T \\ \mathbf{A} \hat{\mathbf{V}}_t^t & \hat{\mathbf{V}}_{t+1}^t \end{bmatrix} \right)$$

Please note that the value $\hat{\mathbf{V}}_{t+1}^t$ can be found in the intermediate value of the filtering. Now, we use conditional Gaussian result from above, and we have:

$$\begin{aligned} p(\mathbf{z}_t | \mathbf{z}_{t+1}, \mathbf{x}_{1:t}) &= \mathcal{N} \left(\mathbf{z}_t \mid \hat{\mathbf{z}}_t + \mathbf{J}_t(\mathbf{z}_{t+1} - \hat{\mathbf{z}}_{t+1}^t), \hat{\mathbf{V}}_t^t - \mathbf{J}_t \hat{\mathbf{V}}_{t+1}^t \mathbf{J}_t^T \right) \\ &= \mathcal{N} \left(\mathbf{z}_t \mid \hat{\mathbf{z}}_t + \mathbf{J}_t(\mathbf{z}_{t+1} - \mathbf{A}\hat{\mathbf{z}}_t^t), \hat{\mathbf{V}}_t^t - \mathbf{J}_t \hat{\mathbf{V}}_{t+1}^t \mathbf{J}_t^T \right) \end{aligned}$$

where $\mathbf{J}_t = \hat{\mathbf{V}}_t^t \mathbf{A}^T (\hat{\mathbf{V}}_{t+1}^t)^{-1}$. Now consider the marginalization using a normal Gaussian linear model together with $p(\mathbf{z}_{t+1} | \hat{\mathbf{z}}_{t+1}^T, \hat{\mathbf{V}}_{t+1}^T)$, as we now have:

$$\begin{aligned} p(\mathbf{z}_t | \mathbf{x}_{1:T}) &= \mathcal{N} \left(\mathbf{z}_t \mid \hat{\mathbf{z}}_t + \mathbf{J}_t (\mathbf{z}_{t+1} - \mathbf{A} \hat{\mathbf{z}}_t^t), \hat{\mathbf{V}}_t^t - \mathbf{J}_t \hat{\mathbf{V}}_{t+1}^t \mathbf{J}_t^T + \mathbf{J}_t \hat{\mathbf{V}}_{t+1}^T \mathbf{J}_t \right) \\ &= \mathcal{N} \left(\mathbf{z}_t \mid \hat{\mathbf{z}}_t + \mathbf{J}_t (\mathbf{z}_{t+1} - \mathbf{A} \hat{\mathbf{z}}_t^t), \hat{\mathbf{V}}_t^t + \mathbf{J}_t (\hat{\mathbf{V}}_{t+1}^T - \hat{\mathbf{V}}_{t+1}^t) \mathbf{J}_t^T \right) \end{aligned}$$

Thus the prove is now complete. \square

4.3 Learning Parameter

Definition 4.9. (EM-Algorithm for SSM) We have the following free-energy:

$$F(q, \theta) = \int q(\mathbf{z}_{1:T}) \left[\log p(\mathbf{x}_{1:T}, \mathbf{z}_{1:T} | \theta) - \log q(\mathbf{z}_{1:T}) \right] d\mathbf{z}_{1:T}$$

Now, we consider the following EM-step, as we have:

- *E-Step*: We find $q^*(\mathbf{z}_{1:T}) = p(\mathbf{z}_{1:T} | \mathbf{x}_{1:T}, \theta)$, which is already done in the Kalman smoothing.
- *M-Step*: We want to find θ by maximizing the $\langle \log p(\mathbf{z}_{1:T}, \mathbf{x}_{1:T} | \theta) \rangle_{q(\mathbf{z}_{1:T} | \mathbf{x}_{1:T})}$, where we have:

$$p(\mathbf{z}_{1:T}, \mathbf{x}_{1:T} | \theta) = p(\mathbf{z}_1) p(\mathbf{x}_1 | \mathbf{z}_1) \prod_{t=2}^T p(\mathbf{z}_t | \mathbf{z}_{t-1}) p(\mathbf{x}_t | \mathbf{z}_t)$$

which are all Gaussian, this leads to least square problem as we will see. This also works with discrete case.

Proposition 4.6. (M-Step for C) The update for \mathbf{C} is:

$$\mathbf{C} = \mathbf{R}^{-1} \sum_{t=1}^T \langle \mathbf{z}_t \rangle_q \mathbf{x}_t^T \left(\sum_{t=1}^T \langle \mathbf{z}_t \mathbf{z}_t^T \rangle_q \right)^{-1}$$

Proof. We have the following log-likelihood (with removed unnecessary terms):

$$\begin{aligned} \frac{\partial}{\partial \mathbf{C}} \log p(\mathbf{z}_{1:T}, \mathbf{x}_{1:T} | \theta) &= \frac{\partial}{\partial \mathbf{C}} \left[\sum_{t=1}^T \langle \mathbf{P}(\mathbf{x}_t | \mathbf{z}_t) \rangle_q \right] \\ &= \frac{\partial}{\partial \mathbf{C}} \left[-\frac{1}{2} \sum_{t=1}^T \langle (\mathbf{x}_t - \mathbf{C} \mathbf{z}_t)^T \mathbf{R}^{-1} (\mathbf{x}_t - \mathbf{C} \mathbf{z}_t) \rangle_q \right] \\ &= \frac{\partial}{\partial \mathbf{C}} \left[-\frac{1}{2} \sum_{t=1}^T \langle \mathbf{x}_t^T \mathbf{R}^{-1} \mathbf{x}_t^T - 2 \mathbf{x}_t^T \mathbf{R}^{-1} \mathbf{C} \mathbf{z}_t + \mathbf{z}_t^T \mathbf{C}^T \mathbf{R}^{-1} \mathbf{C} \mathbf{z}_t \rangle_q \right] \\ &= \frac{\partial}{\partial \mathbf{C}} \left[\text{Tr} \left[\mathbf{C} \sum_{t=1}^T \langle \mathbf{z}_t \rangle_q \mathbf{x}_t^T \mathbf{R}^{-1} \right] - \frac{1}{2} \text{Tr} \left[\mathbf{C}^T \mathbf{R}^{-1} \mathbf{C} \left\langle \sum_{t=1}^T \mathbf{z}_t \mathbf{z}_t^T \right\rangle_q \right] \right] \\ &= \mathbf{R}^{-1} \sum_{t=1}^T \langle \mathbf{z}_t \rangle_q \mathbf{x}_t^T - \mathbf{R}^{-1} \mathbf{C} \left\langle \sum_{t=1}^T \mathbf{z}_t \mathbf{z}_t^T \right\rangle_q \end{aligned}$$

Setting this to 0, and we have the update. \square

Proposition 4.7. (M-Step for A) The update for \mathbf{A} is:

$$\mathbf{A} = \mathbf{R}^{-1} \sum_{t=1}^T \langle \mathbf{z}_t \mathbf{z}_{t+1}^T \rangle_q \left(\sum_{t=1}^T \langle \mathbf{z}_t \mathbf{z}_t^T \rangle_q \right)^{-1}$$

Proof. We see that the we have the same situation as the \mathbf{C} , and so the update is the same. Please note that we replace \mathbf{x}_t with \mathbf{z}_{t+} . \square

Proposition 4.8. (M-Step for HMM) We can take the free-energy with respected to the parameter θ , and we have:

- Initial State Distribution $\pi_i = \gamma_i(i)$, which is expected number of times in state i at the start.
- Define the expected transition from state $i \rightarrow j$ to be:

$$\xi_t(i \rightarrow j) = p(z_t = i, z_{t+1} = j | \mathbf{x}_{1:T}) = \frac{\alpha_t(i) \Phi_{ij} \mathbf{A}_j(\mathbf{x}_{t+1}) \beta_{t+1}(j)}{p(\mathbf{x}_{1:T})}$$

- We have the transition probability and output probability to be:

$$\hat{\Phi}_{ij} = \sum_{t=1}^{T-1} \xi_t(i \rightarrow j) \Big/ \sum_{t=1}^T \gamma_t(i) \quad \hat{A}_{ik} = \sum_{t:\mathbf{x}_t=k} \gamma_t(i) \Big/ \sum_{t=1}^T \gamma_t(i)$$

Proof. As we have computed the E-step, now we consider the free energy, as M-Step trying to maximize the following quantity

$$\begin{aligned} & \langle \log p(z_{1:T}, \mathbf{x}_{1:T}) \rangle_{q(z_{1:T} | \mathbf{x}_{1:T}, \theta)} \\ &= \mathbb{E}_{q(z_{1:T} | \mathbf{x}_{1:T})} \left[\log p(z_1 | \boldsymbol{\pi}) + \sum_{t=2}^T \log p(z_t | z_{t-1}, \Phi_{ij}) + \sum_{t=1}^T \log p(\mathbf{x}_t | z_t, A_{z_t}(\mathbf{x}_t)) \right] \\ &= \mathbb{E}_{q(z_1 | \mathbf{x}_{1:T})} [\log p(z_1 | \boldsymbol{\pi})] + \sum_{t=2}^T \mathbb{E}_{q(z_t, z_{t-1} | \mathbf{x}_{1:T})} [\log p(z_t | z_{t-1}, \Phi_{ij})] + \sum_{t=1}^T \mathbb{E}_{q(z_t)} [\log p(\mathbf{x}_t | z_t, A_{z_t}(\mathbf{x}_t))] \\ &= \sum_{i=1}^L \gamma_1(i) \pi_i + \sum_{t=2}^T \sum_{i=1}^L \sum_{j=1}^L \xi(z_i = i, z_{t-1} = j) \Phi_{ji} + \sum_{t=1}^T \sum_{i=1}^L \gamma_t(i) A_{i\mathbf{x}_i} \end{aligned}$$

Recall that the number of state to be L . We will denote the emission matrix $\mathbf{A} \in \mathbb{R}^{L \times O}$ where O is the number of possible observable. Furthermore, we recall that we have the following constraint to be:

$$\sum_j \Phi_{ij} = 1 \quad \sum_{\mathbf{x} \in \mathcal{X}} A_i(\mathbf{x}) = 1 \quad \sum_i \pi_i = 1$$

for all $i = 1, \dots, L$ in the first and second constraint. Now, we have the following Lagrange multiplier:

$$\begin{aligned} \mathcal{L} &= \sum_{i=1}^L \gamma_1(i) \pi_i + \sum_{t=2}^T \sum_{i=1}^L \sum_{j=1}^L \xi(z_i = i, z_{t-1} = j) \Phi_{ji} + \sum_{t=1}^T \sum_{i=1}^L \gamma_t(i) A_{i\mathbf{x}_i} \\ &\quad - \sum_{i=1}^L \lambda_i \left(\sum_{j=1}^L \Phi_{ij} - 1 \right) - \sum_{i=1}^L \nu_i \left(\sum_{j=1}^L A_{ij} - 1 \right) - \eta \left(\sum_{i=1}^L \pi_i - 1 \right) \end{aligned}$$

We have the following derivative for π_a :

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial \pi_a} &= \frac{\partial}{\partial \pi_a} \sum_{i=1}^L \gamma_1(i) \log \pi_i - \frac{\partial}{\partial \pi_a} \eta \left(\sum_{i=1}^L \pi_i - 1 \right) = \frac{\gamma_1(a)}{\pi_a} - \eta = 0 \\ &\iff \gamma_1(a) = \nu \pi_a \\ &\iff \sum_{a=1}^L \gamma_1(a) = \nu \sum_{a=1}^L \pi_a = \nu \end{aligned}$$

Solving algebra yields to solution for π_a . Let's consider the value of Φ_{ab} , as we have:

$$\begin{aligned}
\frac{\partial \mathcal{L}}{\partial \Phi_{ab}} &= \frac{\partial \mathcal{L}}{\partial \Phi_{ab}} \sum_{t=2}^T \sum_{i=1}^L \sum_{j=1}^L \xi(z_t = i, z_{t-1} = j) \Phi_{ji} - \frac{\partial \mathcal{L}}{\partial \Phi_{ab}} \sum_{i=1}^L \lambda_i \left(\sum_{j=1}^L \Phi_{ij} - 1 \right) \\
&= \sum_{t=1}^T \xi(z_t = b, z_{t-1} = b) \frac{1}{\Phi_{ab}} - \lambda_a = 0 \\
\iff \lambda_a \Phi_{ab} &= \sum_{t=1}^T \xi(z_t = b, z_{t-1} = a) \\
\iff \sum_{j=1}^L \lambda_a \Phi_{aj} &= \sum_{j=1}^L \sum_{t=1}^T \xi(z_t = j, z_{t-1} = a) \\
\iff \lambda_a &= \sum_{j=1}^L \sum_{t=1}^T \xi(z_t = j, z_{t-1} = a) = \sum_{t=1}^T \gamma_{t-1}(a)
\end{aligned}$$

And, so the equality is proven. Note that there are difference in time notation but they are equivalent. The proof for A_{ij} is the same, so we will not go into details. \square

Remark 39. (Practical Numerical Method) The message passing algorithm can be implode i.e $\alpha_t(i) = p(\mathbf{x}_{1:t}, \mathbf{z}_t = i) \rightarrow 0$ to fix this, we do the following rescale:

$$\bar{\alpha}_t(i) = \mathbf{A}_i(\mathbf{x}_t) \sum_{j=1}^L \tilde{\alpha}_{t-1}(j) \Phi_{ij} \quad \rho_t = \sum_{i=1}^L \bar{\alpha}_t(i) \quad \tilde{\alpha}(i) = \bar{\alpha}_t(i) / \rho_t$$

Proposition 4.9. *We can show that: $\rho_t = p(\mathbf{x}_t | \mathbf{x}_{1:t-1}, \boldsymbol{\theta})$*

Proof. We will consider the quantity $p(z_t | \mathbf{x}_{1:t})$ to be equal to:

$$\tilde{\alpha}_t(z) = p(z_t | \mathbf{x}_1, \dots, \mathbf{x}_t) = \frac{\alpha_t(i)}{p(\mathbf{x}_{1:t})}$$

As, we will denote $\rho_t = p(\mathbf{x}_t | \mathbf{x}_{1:t-1})$, and from the product rule, we have $p(\mathbf{x}_{1:t}) = \prod_{i=1}^t \rho_i$, and so, we can recover:

$$\alpha_t(i) = \left(\prod_{i=1}^t \rho_i \right) \tilde{\alpha}_t(i) = \left(\prod_{i=1}^{t-1} \rho_i \right) \rho_t \tilde{\alpha}_t(i)$$

Let's consider the recursive property of the α , recall its formula:

$$\begin{aligned}
\alpha_{t+1}(i) &= \left(\sum_{j=1}^L \alpha_t(j) \Phi_{ij} \right) \mathbf{A}_i(\mathbf{x}_{t+1}) \\
\iff \left(\prod_{i=1}^n \rho_i \right) \rho_{n+1} \tilde{\alpha}_{t+1}(i) &= \left(\sum_{j=1}^L \left[\left(\prod_{i=1}^n \rho_i \right) \tilde{\alpha}(z_n) \right] \Phi_{ij} \right) \mathbf{A}_i(\mathbf{x}_{t+1}) \\
\iff \rho_{n+1} \tilde{\alpha}_{t+1}(i) &= \left(\sum_{j=1}^L \tilde{\alpha}(z_n) \Phi_{ij} \right) \mathbf{A}_i(\mathbf{x}_{t+1})
\end{aligned}$$

Thus, we have proven the identity. For $\beta_t(i)$. we can consider the fact that $\hat{\beta}_t(i)$ is:

$$\beta_t(i) = \left(\prod_{i=t+1}^N c_i \right) \hat{\beta}_t(i)$$

\square

Remark 40. (Decoding/Derivation of Viterbi) We are interested in finding the best hidden state path given the observations $\mathbf{o}_{1:t}$, together with the following quantity $\delta_t(i)$:

$$\begin{aligned}\{z_1^*, \dots, z_t^*\} &= \arg \max_{z_1, \dots, z_t} p(z_1, \dots, z_t, \mathbf{x}_{1:t} = \mathbf{o}_{1:t} | \boldsymbol{\theta}) \\ \delta_t(i) &= p(z_1, \dots, z_{t-1}, z_t = i, \mathbf{x}_{1:t} = \mathbf{o}_{1:t} | \boldsymbol{\theta})\end{aligned}$$

We can see that we can use the induction step to find this value:

$$\delta_{t+1}(j) = \left[\max_i \delta_t(i) \Phi_{ij} \right] A_i(\mathbf{o}_{t+1})$$

This recursion gives us the Viterbi algorithm.

Definition 4.10. (Viterbi Algorithm) The algorithm to find the best path is given by, where the sequence of hidden state is stored in ψ :

- Initialization:

$$\delta_1(i) = \pi_i b_i(\mathbf{o}_1) \quad \psi_1(i) = 0$$

- Induction Step:

$$\delta_t(j) = \left[\max_{i \in [L]} \delta_{t-1}(i) a_{ij} \right] b_j(\mathbf{o}_t) \quad \psi_t(j) = \arg \max_{i \in [L]} \delta_{t-1}(i) a_{ij}$$

- Termination:

$$P^* = \max_{i \in [L]} [\delta_T(i)] \quad q_T^* = \arg \max_{i \in [L]} [\delta_T(i)]$$

Remark 41. (Multiple sequences) We now consider the update for HMM but with multiple sequences $l \in [L]$, which we can show that the batch update is:

$$\pi_i = \frac{1}{L} \sum_{l=1}^L \gamma_1^{(l)}(i) \quad \Phi_{ij} = \frac{\sum_{l=1}^L \sum_{t=1}^{T^{(i)}-1} \xi_t^{(l)}(i \rightarrow j)}{\sum_{l=1}^L \sum_{t=1}^{T^{(i)}} \gamma_t^{(l)}(i)} \quad \mathbf{A}_{ik} = \frac{\sum_{l=1}^L \sum_{t=1}^{T^{(i)}-1} \delta(\mathbf{x}_t = k) \gamma_t^{(l)}(i \rightarrow j)}{\sum_{l=1}^L \sum_{t=1}^{T^{(i)}} \gamma_t^{(l)}(i)}$$

5 Markov Chain (Monte Carlo)

5.1 Markov Chain

Definition 5.1. (Transition Matrix) We consider the probability of state $i \in [d]$ will move to state $j \in [d]$ i.e $P_{i \rightarrow j}$, where d is the number of states. This information for the markov chain is stored in $d \times d$ transtion matrix:

$$\mathbf{P} = \begin{pmatrix} P_{1 \rightarrow 1} & \cdots & P_{1 \rightarrow d} \\ \vdots & \ddots & \vdots \\ P_{d \rightarrow 1} & \cdots & P_{d \rightarrow d} \end{pmatrix}$$

Definition 5.2. (Initial Distribution) At the start, we consider the distribution over the states:

$$\mathbf{p}_{\text{init}} = \left(p(x_0 = 1), \dots, p(x_0 = d) \right)^T$$

Lemma 5.1. Given the state distribution p_t , we can find the probability distribution over the next state after transtion \mathbf{p}_{t+1} is:

$$\mathbf{p}_{t+1} = \mathbf{P} \mathbf{p}_t$$

Note that $\mathbf{p}_t \in \mathbb{R}^{d \times 1}$ for any $t \in \mathbb{N}$

Proof. Consider the distribution of the first state:

$$\mathbf{p}_{t+1}(s) = p(x_{t+1} = s) = \sum_{i=1}^d p(x_{t+1} = s | x_t = i) p(x_t = i) = \sum_{i=1}^d P_{si} \mathbf{p}_t[i]$$

And so it is proven. \square

Definition 5.3. (Stationary) Stationary markov chain is when the transtion matrix \mathbf{P} doesn't depends on time step.

Definition 5.4. (Equilibrium Distribution) Let \mathbf{P} be transtion matrix, then a distribution \mathbf{p} such that: $\mathbf{P}\mathbf{p} = \mathbf{p}$ is called equilibrium distribution.

Remark 42. Please note that if we consider applying transtion matrix to any starting distribution:

$$\mathbf{p}_\infty = \lim_{n \rightarrow \infty} \mathbf{P}^n \mathbf{p}_{\text{init}}$$

assuming that it exists, then we can show that it gives rise to (somewhat) equilibrium distribution:

$$\mathbf{P}\mathbf{p}_\infty = \mathbf{P} \lim_{n \rightarrow \infty} \mathbf{P}^n \mathbf{p}_{\text{init}} = \lim_{n \rightarrow \infty} \mathbf{P}^{n+1} \mathbf{p}_{\text{init}} = \mathbf{p}_\infty$$

Definition 5.5. (Aperiodic MC) A stationary MC is aperiodic if: $P_{ii} = 0$ for $i \in [d]$.

Definition 5.6. (Irreducible) MC is irreducible if there is a path from each state to every other state in the transtion.

Theorem 5.1. *If first-order (depends only last time step), stationary MC with finite state space is a periodic and irreducible, then:*

- *Limit distribution \mathbf{p}_∞ exists.*
- *Limit distribution is also the equilibrium distribution.*
- *Equilibrium distribution is unique.*
- *Equilibrium distribution doesn't depends on initial distribution.*

Definition 5.7. (Aperiodic State) The state i is a aperiodic if there exits a time t such that for all $t \geq t_0$ with positive integer t_0 , the state comes back to itself with positive probability.

Definition 5.8. (Positive Recurrence State) The state i is positive recurrence if the expected number of times that it will return to itself is finite.

Definition 5.9. (Ergodicity) There exists a positive integer t_0 , such that for all pair of states i, j , if markov chain starts at time 0 at state i , then for all $t \geq t_0$, the probability of being in state j at time t is more than 0. Note that this implies all states are positive recurrence and aperiodic.

Theorem 5.2. (Result of Ergodicity) *The ergodicity implies the existence of unique stationary distribution \mathbf{p}_∞ given any initialized probability distribution.*

Definition 5.10. (Detailed Balance wrt. \mathbf{p}^*) The markov chain is detailed balance with respect to a stationary distribution \mathbf{p}^* if for all states i, j :

$$p_i^* P_{ij} = p_j^* P_{ji}$$

This also means that the markov chain is *reversible* under this stationary distribution.

Remark 43. (Solving Detailed Balance) Detailed balance implies that we have to solve $|\mathcal{S}|^2$ equation, which is more than the simple stationary equation, which require only $|\mathcal{S}|$ equations.

Proposition 5.1. *The stationary distribution of detailed balance markov chain with respect to stationary distribution \mathbf{p}^* is \mathbf{p}^**

Proof. We consider the summation over the detailed balance equation, as we have:

$$\sum_{i=1}^{|\mathcal{S}|} p_i^* P_{ij} = \sum_{i=1}^{|\mathcal{S}|} p_j^* P_{ji} = p_j^* \sum_{i=1}^{|\mathcal{S}|} P_{ji} = p_j^*$$

This complete the proof. □

Remark 44. (Ergodicity and Detailed Balance) Eventhough the detailed balance with respect to stationary distribution \mathbf{p}^* , will have the stationary distribution to be \mathbf{p}^* , this doesn't means that it is unique. For the uniqueness to hold, ergodicity is required.

Remark 45. (Finding Equilibrium) There are 2 ways we can find the equilibrium: power method, and eigenvalue. Both of them is based on finding the underlying eigenvalue and eigenvectors.

Definition 5.11. (Random Walk) Given a directed graph G with d vertices, a random walk generates a sequence of nodes: x_0, x_1, \dots by first select the vertex x_0 at random and, at time t , we uniformly at random select the children of x_{t-1} to get to vertex x_t .

Remark 46. Random walk can be defined as a markov chain, where:

$$\mathbf{p}_{\text{init}} = \left(\frac{1}{p}, \dots, \frac{1}{p} \right) \quad P_{ij} = \begin{cases} \frac{1}{\text{number of edges}} & \text{if } i \text{ links to } j \\ 0 & \text{otherwise} \end{cases}$$

Remark 47. PageRank algorithm applies random walk to a graph, in order to make the original markov chain defined by a transtion matrix \mathbf{T} to be irreducible and aperiodic:

$$P_{ij} = (1 - \alpha)\mathbf{T}_{ij} + \frac{\alpha}{d}$$

where $\alpha \in (0, 1)$

5.2 Sampling Algorithm

Definition 5.12. (Sampling Algorithm) Given a distribution p , the sampling algorithm samples to get the random point $\mathbf{x}_1, \dots, \mathbf{x}_n$ such that the marginal distribution is p . Ideally, the drawing are independent.

Remark 48. (Usefulness of Sampling Algorithm) There are many use-cases of sampling algorithm notably:

- Compute expectation. Given the output $\mathbf{x}_1, \dots, \mathbf{x}_N$ to be independent from distribution p , the expectation is:

$$\mathbb{E}[f(\mathbf{x})] \approx \frac{1}{N} \sum_{i=1}^N f(\mathbf{x}_i)$$

- We can approximate the distribution using the sample to perform a density estimation.

Definition 5.13. (Sampling in Bayesian Model) Sometimes the posterior over parameter $\hat{q}(\boldsymbol{\theta}|\mathbf{x}_{1:n})$ can't be computed analytically, we can still sample from it to get $\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_m$ given the posterior \hat{q} . We can now use it for computing expectation like above or perform a prediction:

$$p(\mathbf{x}_{n+1}|\mathbf{x}_{1:n}) = \int_{\Theta} p(\mathbf{x}_{n+1}|\boldsymbol{\theta})\hat{q}(\boldsymbol{\theta}|\mathbf{x}_{1:n}) d\boldsymbol{\theta} \approx \frac{1}{m} \sum_{i=1}^m p(\mathbf{x}_{n+1}|\boldsymbol{\theta}_i)$$

Definition 5.14. (Boxed Rejection Sampling) We perform the following procedure, as we want to sample from distribution \tilde{p} :

- We generate the $X_i \sim \mathcal{U}[a, b]$, $Y_i \sim \mathcal{U}[0, c]$ uniformly to define the box.
- If $Y_i \leq \tilde{p}(X_i)$, then we keep the sample X_i .
- Otherwise, we reject the sample and repeat the process.

Remark 49. (Invariance to Scale) Please note that rejection sampling still works if we have $Y_i \sim \mathcal{U}[0, kc]$ given the distribution $kp(\cdot)$ for $k > 0$. This means that rejection sampling only requires us to know the shape of the distribution, as long as kc is big enough to cover the whole distribution p . This also means that, in Bayesian inference, we don't have to know the evidence in order to sample from the posterior:

$$p(\boldsymbol{\theta}|\mathcal{X}) = \frac{p(\mathcal{X}|\boldsymbol{\theta})p(\boldsymbol{\theta})}{Z} \quad \text{where} \quad Z = \int p(\mathcal{X}|\boldsymbol{\theta})p(\boldsymbol{\theta}) \, d\boldsymbol{\theta}$$

We only need $\tilde{p} = p(\mathcal{X}|\boldsymbol{\theta})p(\boldsymbol{\theta})$ to do rejection sampling.

Definition 5.15. (Rejection Sampling) One doesn't have to use the uniform distribution to sample proposed points. We simply have to find the distribution r such that $\tilde{p} < r$ everywhere:

- Sample $X_i \sim r$
- Sample $Y_i|X_i \sim \mathcal{U}[0, r(X_i)]$
- If $Y_i \leq \tilde{p}(X_i)$, then we keep the sample X_i .
- Otherwise, we reject the sample and repeat the process.

The scaling properties still hold. It will generate iid samples $\theta_1, \dots, \theta_m$, where the samples $1/m \sum_{i=1}^m f(x_i)$ is an unbiased estimate of $\mathbb{E}_p[f(x)]$.

Remark 50. Given the height of r , $|A|$ is too high, and the height of \tilde{p} , $|B|$. The rejection sampling will accept a sample with probability $|B|/|A|$, so if $|A|$ is too high, the algorithm may be inefficient.

Definition 5.16. (Importance Sampling) Importance sampling is given distribution $p = 1/Z_p \tilde{p}$ that we want to sample from and an arbitrary proposal distribution $q = 1/Z_q \tilde{q}$ (that we can sample), then:

- Draw $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m$ iid from proposal q .
- Expectation $\mathbb{E}_p[f(\mathbf{x})]$ is approximated as:

$$\frac{1}{m} \sum_{i=1}^m \frac{f(\mathbf{x}_i)[\tilde{p}(\mathbf{x}_i)/\tilde{q}(\mathbf{x}_i)]}{\sum_{j=1}^m \tilde{p}(\mathbf{x}_j)/\tilde{q}(\mathbf{x}_j)}$$

Proposition 5.2. *The importance sampling gives unbiased estimate of $\mathbb{E}_p[f(\mathbf{x})]$.*

Proof. We have the following sampling:

$$\mathbb{E}_p[f(\mathbf{x})] = \int f(\mathbf{x})p(\mathbf{x}) \, d\mathbf{x} = \int f(\mathbf{x})\frac{p(\mathbf{x})}{q(\mathbf{x})}q(\mathbf{x}) \, d\mathbf{x} = \mathbb{E} \left[f(\mathbf{x})\frac{p(\mathbf{x})}{q(\mathbf{x})} \right] \approx \frac{1}{m} \sum_{i=1}^m \frac{f(\mathbf{x}_i)p(\mathbf{x}_i)}{q(\mathbf{x}_i)}$$

where $\mathbf{x}_1, \dots, \mathbf{x}_m \sim q$. Now, consider the ratio of normalized factors:

$$\frac{Z_p}{Z_q} = \frac{\int \tilde{p}(\mathbf{x}) \, d\mathbf{x}}{Z_q} = \frac{\int \tilde{p}(\mathbf{x})\frac{q(\mathbf{x})}{q(\mathbf{x})} \, d\mathbf{x}}{Z_q} = \int \tilde{p}(\mathbf{x})\frac{q(\mathbf{x})}{Z_q q(\mathbf{x})} \, d\mathbf{x} = \mathbb{E}_q \left[\frac{\tilde{p}(\mathbf{x})}{\tilde{q}(\mathbf{x})} \right] \approx \frac{1}{m} \sum_{i=1}^m \frac{\tilde{p}(\mathbf{x}_i)}{\tilde{q}(\mathbf{x}_i)}$$

Now, the estimator of f can be found as

$$\mathbb{E}_p[f(\mathbf{x})] \approx \frac{1}{m} \sum_{i=1}^m \frac{f(\mathbf{x}_i)p(\mathbf{x}_i)}{q(\mathbf{x}_i)} = \frac{1}{m} \sum_{i=1}^m \frac{f(\mathbf{x}_i)p(\mathbf{x}_i)}{q(\mathbf{x}_i)} \frac{\tilde{p}(\mathbf{x}_i)}{\tilde{q}(\mathbf{x}_i)} \frac{Z_q}{Z_p} = \frac{1}{m} \sum_{i=1}^m \frac{f(\mathbf{x}_i)[\tilde{p}(\mathbf{x}_i)/\tilde{q}(\mathbf{x}_i)]}{\sum_{j=1}^m \tilde{p}(\mathbf{x}_j)/\tilde{q}(\mathbf{x}_j)}$$

□

5.3 More Probabilistic Models

Definition 5.17. (Random Field) Given a weighted undirected graph $\mathcal{N} = (v_{\mathcal{N}}, w_{\mathcal{N}})$ where $v_{\mathcal{N}}$ is the vertex set, and $w_{\mathcal{N}}$ is the set of edge weights. The edge weights are scalar $w_{ij} \in \mathbb{R}$. An edge weight w_{ij} means to edge between i and j . Each vertex v_i is associated with random variable Θ_i . The neighbourhood of vertex v_i is the set:

$$\partial(i) = \{j : w_{ij} \neq 0\}$$

The set $\{\Theta_j : j \in \partial(i)\}$ of random variable associated with neighbourhood is called Markov blanket of Θ_i

Definition 5.18. (Markov Property) The Markov properties is when:

$$p(\theta_i | \theta_j, j \neq i) = p(\theta_i | \theta_j, j \in \partial(i))$$

Each Θ is conditionally independent of remaining field given its Markov blanket. A Markov random field is a random field that is Markov.

Definition 5.19. (Energy Function) Any (strictly positive) probability or density can be rewritten as:

$$p(x) = \frac{1}{Z} \exp(-H(x)) \quad \text{where} \quad H : \mathcal{X} \rightarrow \mathbb{R}_+ \text{ and } Z = \int \exp(-H(x)) dx$$

where H is called potential.

Remark 51. The Markov random field (MRF) density for random variable $\Theta_{1:n}$ can be written as:

$$p(\theta_1, \dots, \theta_n) = \frac{1}{Z} \exp(-H(\theta_1, \dots, \theta_n))$$

Definition 5.20. (Potts Model) Let \mathcal{N} be a neighbourhood of graph with weight w_{ij} and $\beta > 0$. The MRF, where we have:

$$p(\theta_1, \dots, \theta_n) = \frac{1}{Z(\beta)} \exp\left(\beta \sum_{ij} w_{ij} \mathbb{I}\{\theta_i = \theta_j\}\right)$$

Note that the energy is additive over pairs. Positive weight encourage smoothness:

- $w_{ij} > 0$, this means that $\theta_i = \theta_j$ increases probability.
- $w_{ij} < 0$, this means that $\theta_i = \theta_j$ decreases probability.
- $w_{ij} = 0$, this means that no interaction between θ_i and θ_j .

Definition 5.21. (Ising Model) We have the following distribution:

$$p(\theta_{1:n}) = \frac{1}{Z(\beta)} \exp\left(\sum_{(i,j) \in \text{edge}} \beta \mathbb{I}[\theta_i = \theta_j]\right)$$

Given rejection sampling, we can sample the distribution without the normalizing factor $Z(\beta)$, where $\theta_i \in \{-1, +1\}$ and $w_{ij} \in \{0, 1\}$. This model is on the d -dimensional grid.

Remark 52. (Usage of Markov Random Field) Given the problem with observation x_i for each i location on a grid. If we want to model the observation with a distribution $p(x_i | \Theta_i)$ for each location, with its own parameter Θ_i . We can use MRF as a prior distribution, while we represent $p(x_i | \Theta_i)$ as the emission probability:

- Define a joint $(\Theta_1, \dots, \Theta_n)$ as an MRF on the grid graph.
- Given a positive weight, the MRF will encourage the model to explain a neighbourhood of x_i by similar parameter value. This leads to smoothing in the result.

Definition 5.22. (Bayesian Mixture Model) The model of the form:

$$\pi(\mathbf{x}) = \sum_{k=1}^K c_k p(\mathbf{x}|\Theta_k)$$

is called Bayesian model if $p(\mathbf{x}|\theta)$ is an exponential model and:

- $\Theta_1, \dots, \Theta_K \sim q$ where q is a prior over Θ
- (c_1, \dots, c_K) is sampled from K -dimensional dirichlet distribution.

Definition 5.23. (Posterior Bayesian Mixture Model) Another intractable model that may need sampling is:

$$\hat{q}_n(c_{1:k}, \boldsymbol{\theta}|\mathbf{x}_{1:n}) \propto \prod_{i=1}^n \left(\sum_{k=1}^K c_k p(\mathbf{x}_i|\boldsymbol{\theta}_k) \right) \left(\prod_{k=1}^K q(\boldsymbol{\theta}_k) \right) q_{\text{dir}}(c_{1:K})$$

individual evaluate of non-normalized \tilde{q} is numerically unstable but given specific value of c, \mathbf{x} or $\boldsymbol{\theta}$, this collapse to $\sum_{k=1}^K c_k p(\mathbf{x}_i|\boldsymbol{\theta}_k)$ making it tractable. Furthermore, please note that: we can multiple the Bayesian Mixture model $\prod_{k=1}^K q(\boldsymbol{\theta}_k)$ with MRF prior to get smoothing effect.

Remark 53. (Problems and Solution) If MRF is used as prior, we have to compute or approximate the posterior distribution. The solution of MRF distribution on grids are not analytically tractable, so we have to perform sampling and inference using Markov chain sampling algorithm.

5.4 Markov Chain Monte Carlo

Definition 5.24. (MCMC) We want to sample from distribution with density p . Suppose, we can define Markov chain with invariance distribution i.e $P_{\text{inv}} \equiv P$. If we sample $\mathbf{x}_1, \mathbf{x}_2, \dots$ from the chain, then once it has converged and we obtain the sample.

Definition 5.25. (Continuous MC) A continuous Markov chain is defined by an initial p_{init} and conditional probability $t(y|x)$ is transition probability or transition kernel. For example, markov chain on \mathbb{R}^2 , we can define Markov chain by:

$$x_{i+1}|x_i = x_i \sim g(\cdot|x_i, \sigma_i^2)$$

where g is the spherical Gaussian with fixed variance. Suppose the state \mathcal{X} is uncountable, so the transition matrix is substituted by conditional probability t . The distribution p_{inv} with density p_{inv} is invariance if:

$$\int_{\mathcal{X}} t(y|x)p_{\text{inv}}(x) dx = p_{\text{inv}}(y)$$

Remark 54. (Several Problems) There are several problems that we have to solve. We have to construct MC with invariance distribution p . We can't actually start sampling with $x_1 \sim p$ as if we know how to sample, our method would be pointless. Furthermore, each point x_i is marginally distribution as $x_i \sim p$ but the points are not iid.

Definition 5.26. (Metropolis-Hasting) We define the conditional probability $q(y|x)$ on \mathcal{X} . We then define a rejection kernel A as we have:

$$A(\mathbf{x}_{i+1}|\mathbf{x}_i) = \min \left\{ 1, \frac{q(\mathbf{x}_i|\mathbf{x}_{i+1})p(\mathbf{x}_{i+1})}{q(\mathbf{x}_{i+1}|\mathbf{x}_i)p(\mathbf{x}_i)} \right\}$$

Knowing \tilde{p} (unnormalized) is enough for the rejection kernel $A(\cdot|\cdot)$ as the normalizing factor is cancelled. We define the transition probability of the chain as:

$$t(\mathbf{x}_{i+1}|\mathbf{x}_i) = q(\mathbf{x}_{i+1}|\mathbf{x}_i)A(\mathbf{x}_{i+1}|\mathbf{x}_i) + \delta_x(\mathbf{x}_{i+1})c(\mathbf{x}_i) \quad \text{where} \quad c(\mathbf{x}_i) = \int q(\mathbf{y}|\mathbf{x}_i)(1 - A(\mathbf{y}|\mathbf{x}_i)) d\mathbf{y}$$

To sample from the MH algorithm, at each step $i + 1$, generate a proposal $\mathbf{x}^* \sim q(\cdot|\mathbf{x}_i)$ and $\mathbf{u}_i \sim \mathcal{U}[0, 1]$:

- If $u_i > A(\mathbf{x}^*|\mathbf{x}_i)$ reject a proposal, and we set $\mathbf{x}_{i+1} = \mathbf{x}_i$
- If $u_i \leq A(\mathbf{x}^*|\mathbf{x}_i)$, accept proposal, and we set $\mathbf{x}_{i+1} = \mathbf{x}^*$

Remark 55. (Derivation of Metropolis-Hasting) Now, we will consider how Metropolis-Hasting can to be. Starting with the detailed balance equation, as we have:

$$p(\mathbf{x}_{i+1}|\mathbf{x}_i)p(\mathbf{x}) = p(\mathbf{x}|\mathbf{x}_{i+1})p(\mathbf{x}_{i+1}) \iff \frac{p(\mathbf{x}_{i+1}|\mathbf{x})}{p(\mathbf{x}|\mathbf{x}_{i+1})} = \frac{p(\mathbf{x}_{i+1})}{p(\mathbf{x})}$$

Now, we will separate the transition step to be:

- Proposal Distribution $q(\mathbf{x}_{i+1}|\mathbf{x}_i)$
- Acceptance Distribution $A(\mathbf{x}_{i+1}|\mathbf{x}_i)$

Now, we have:

$$\frac{A(\mathbf{x}_{i+1}|\mathbf{x}_i)}{A(\mathbf{x}_i|\mathbf{x}_{i+1})} = \frac{p(\mathbf{x}_{i+1})q(\mathbf{x}_i|\mathbf{x}_{i+1})}{p(\mathbf{x}_i)q(\mathbf{x}_{i+1}|\mathbf{x}_i)}$$

Now, if we use the Metropolis-Hasting rejection kernel, we now have:

$$A(\mathbf{x}_{i+1}|\mathbf{x}_i)p(\mathbf{x}_i)q(\mathbf{x}_{i+1}|\mathbf{x}_i) = A(\mathbf{x}_i|\mathbf{x}_{i+1})p(\mathbf{x}_{i+1})q(\mathbf{x}_i|\mathbf{x}_{i+1}) \\ \iff \min \left\{ P(\mathbf{x}_i)q(\mathbf{x}_{i+1}|\mathbf{x}_i), p(\mathbf{x}_{i+1})q(\mathbf{x}_i|\mathbf{x}_{i+1}) \right\} = \min \left\{ q(\mathbf{x}_{i+1}|\mathbf{x}_i)p(\mathbf{x}_i), p(\mathbf{x}_{i+1})q(\mathbf{x}_i|\mathbf{x}_{i+1}) \right\}$$

Thus the detailed balance is satisfied.

Remark 56. There are several observations on the Metropolis-Hasting as we have:

- We accept if the second term is larger than 1: $q(\mathbf{x}_i|\mathbf{x}_{i+1})p(\mathbf{x}_{i+1}) > q(\mathbf{x}_{i+1}|\mathbf{x}_i)p(\mathbf{x}_i)$. We accept if the proposal increases the probability under p .
- If it decreases the probability, we still accept with a probability which depends on the difference to the current probability.
- We can see this as noisy hill-climbing as it tends to move to the probability under p .
- However, there are some probability that the sampling can move down-hill with certain probability. Finally, we can show that it won't stuck at the local maxima.

Definition 5.27. (Burn-in and Mixing-Time) The first m samples are called burn-in phase. The first m samples are discarded. And, we have:

$$\mathbf{x}_1, \dots, \mathbf{x}_{m-1}, \mathbf{x}_m, \mathbf{x}_{m+1}, \dots$$

We don't know the m but we can use a certain heuristic called convergence diagnostic.

Definition 5.28. (Sequential Dependence) Even after burn-in MC are not iid, we can use the following strategy as we have to consider the when the sample can be used:

- Estimate empirically how many steps L are needed for \mathbf{x}_i and \mathbf{x}_{i+L} to be independent. We keep every L -th sample and discard in the between.
- The most common method uses is called auto-correlation function as we have:

$$\text{Auto}(\mathbf{x}_i, \mathbf{x}_j) = \frac{\mathbb{E}[(\mathbf{x}_i - \boldsymbol{\mu}_i)(\mathbf{x}_j - \boldsymbol{\mu}_j)]}{\sigma_i \sigma_j}$$

where $\boldsymbol{\mu}_i$ is the mean and σ_i is the standard deviation of \mathbf{x}_i . We can use this auto-correlation to calculate the value L .

Definition 5.29. (Gelman-Rubin Criterion) We can also start several chain at random for chain k the sample x_i^k has the marginal of p_i^i :

- The distribution has to converge to all $p_i = p_{inv}$, which are all idetical.
- We can use hypothesis testing to compare p_i^k for difference k .
- Once the test doesn't reject anymore, assume that the chain has passed the burn-in phase

Remark 57. (Rules for Selecting a Proposal Distribution) Selecting a proposal distribution, we have to be aware of the tradeoff, where if $\text{var}(q)$ is too large will overstep p and leads to rejection. If $\text{var}(q)$ is too small as many steps will be needed to achives a good converge of the domain.

- If p is unimodal and can be roughly approximate by Gaussian, $\text{var}(q)$ should be choosen to be smaller than the covariance of p .
- With complex posterior, choosing q is difficult but it is important to convergence speed. If we know some information about the posterior, this might help choosing q .

There are mmany other ways to sample for example mixture of proposal.

Definition 5.30. (Gibbs Sampling) Suppose $p(\mathbf{x})$ is a distribution on \mathbb{R}^d and so $\mathbf{x} = (\mathbf{x}_1, \dots, \mathbf{x}_D)$. The full conditional probability of the entry \mathbf{x}_d is given by the other entries to be:

$$p(\mathbf{x}_d | \mathbf{x}_1, \dots, \mathbf{x}_{d-1}, \mathbf{x}_{d+1}, \dots, \mathbf{x}_D)$$

The Gibbs sampler is the special case of MH-algorithm where the proposal of \mathbf{x}_d is the full conditional over \mathbf{x}_d

Remark 58. (Connection to HM) Suppose p is a distribution on \mathbb{R}^D so each sampler is of the form $\mathbf{x}_i = (\mathbf{x}_{i,1}, \dots, \mathbf{x}_{i,D})$. We generate a proposal \mathbf{x}_{i+1} as:

$$\begin{aligned} \mathbf{x}_{i+1,1} &\sim p(\cdot | \mathbf{x}_{i,2}, \dots, \mathbf{x}_{i,D}) \\ &\vdots \\ \mathbf{x}_{i+1,d} &\sim p(\cdot | \mathbf{x}_{i+1,1}, \dots, \mathbf{x}_{i+1,d-1}, \mathbf{x}_{i,d+1}, \dots, \mathbf{x}_{i,D}) \\ &\vdots \\ \mathbf{x}_{i+1,D} &\sim p(\cdot | \mathbf{x}_{i,1}, \dots, \mathbf{x}_{i,D-1}) \end{aligned}$$

This is like Metropolis-Hasting as we use the proposal distribution and accepting probability to be 1.

Remark 59. (MRF with Gibbs Sampling) Consider D nodes, each dimension d and markov property:

$$\begin{aligned} \tilde{p}(\theta_d | \theta_1, \dots, \theta_{d-1}, \theta_{d+1}) \\ = \exp \left(\beta (\mathbb{I} \{ \theta_d = \theta_{\text{left}} \} + \mathbb{I} \{ \theta_d = \theta_{\text{right}} \}), \mathbb{I} \{ \theta_d = \theta_{\text{up}} \}, \mathbb{I} \{ \theta_d = \theta_{\text{down}} \} \right) \end{aligned}$$

And so we can sample it.