

# Conditional Gaussian and Marginal Gaussian

Xinghu Yao

September 22, 2018

## 1 Question

Given a joint Gaussian distribution  $\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma})$  with  $\boldsymbol{\Lambda} \equiv \boldsymbol{\Sigma}^{-1}$  and

$$\mathbf{x} = \begin{pmatrix} \mathbf{x}_a \\ \mathbf{x}_b \end{pmatrix}, \quad \boldsymbol{\mu} = \begin{pmatrix} \boldsymbol{\mu}_a \\ \boldsymbol{\mu}_b \end{pmatrix} \quad (1)$$

$$\boldsymbol{\Sigma} = \begin{pmatrix} \boldsymbol{\Sigma}_{aa} & \boldsymbol{\Sigma}_{ab} \\ \boldsymbol{\Sigma}_{ba} & \boldsymbol{\Sigma}_{bb} \end{pmatrix}, \quad \boldsymbol{\Lambda} = \begin{pmatrix} \boldsymbol{\Lambda}_{aa} & \boldsymbol{\Lambda}_{ab} \\ \boldsymbol{\Lambda}_{ba} & \boldsymbol{\Lambda}_{bb} \end{pmatrix} \quad (2)$$

Try to derive the following conditional distribution and marginal distribution:

$$p(\mathbf{x}_a|\mathbf{x}_b) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_{a|b}, \boldsymbol{\Lambda}_{aa}^{-1}) \quad (3)$$

$$\boldsymbol{\mu}_{a|b} = \boldsymbol{\mu}_a - \boldsymbol{\Lambda}_{aa}^{-1} \boldsymbol{\Lambda}_{ab} \mathbf{x}_b - \boldsymbol{\mu}_b. \quad (4)$$

$$p(\mathbf{x}_b) = \mathcal{N}(\mathbf{x}_b|\boldsymbol{\mu}_b, \boldsymbol{\Sigma}_{bb}). \quad (5)$$

## 2 Conditional Gaussian

According to the definition of condition distribution, we have  $p(\mathbf{x}_a|\mathbf{x}_b) = \frac{p(\mathbf{x}_a, \mathbf{x}_b)}{p(\mathbf{x}_b)}$ . Thus, through fixing  $\mathbf{x}_b$  to the observed value and normalizing the resulting expression with  $p(\mathbf{x}_b)$ , we can get the conditional distribution  $p(\mathbf{x}_a|\mathbf{x}_b)$ .

We firstly consider the quadratic form in the exponent of the Gaussian distribution give by Eq. (1) and Eq. (2). In fact, we have

$$\begin{aligned} & -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}) \\ &= -\frac{1}{2} \begin{pmatrix} \mathbf{x}_a - \boldsymbol{\mu}_a \\ \mathbf{x}_b - \boldsymbol{\mu}_b \end{pmatrix}^T \begin{pmatrix} \boldsymbol{\Lambda}_{aa} & \boldsymbol{\Lambda}_{ab} \\ \boldsymbol{\Lambda}_{ba} & \boldsymbol{\Lambda}_{bb} \end{pmatrix} \begin{pmatrix} \mathbf{x}_a - \boldsymbol{\mu}_a \\ \mathbf{x}_b - \boldsymbol{\mu}_b \end{pmatrix} \\ &= -\frac{1}{2} [(\mathbf{x}_a - \boldsymbol{\mu}_a)^T, (\mathbf{x}_b - \boldsymbol{\mu}_b)^T] \begin{pmatrix} \boldsymbol{\Lambda}_{aa} & \boldsymbol{\Lambda}_{ab} \\ \boldsymbol{\Lambda}_{ba} & \boldsymbol{\Lambda}_{bb} \end{pmatrix} \begin{pmatrix} \mathbf{x}_a - \boldsymbol{\mu}_a \\ \mathbf{x}_b - \boldsymbol{\mu}_b \end{pmatrix} \\ &= -\frac{1}{2}(\mathbf{x}_a - \boldsymbol{\mu}_a)^T \boldsymbol{\Lambda}_{aa}(\mathbf{x}_a - \boldsymbol{\mu}_a) - \frac{1}{2}(\mathbf{x}_a - \boldsymbol{\mu}_a)^T \boldsymbol{\Lambda}_{ab}(\mathbf{x}_b - \boldsymbol{\mu}_b) \\ &\quad - \frac{1}{2}(\mathbf{x}_b - \boldsymbol{\mu}_b)^T \boldsymbol{\Lambda}_{ba}(\mathbf{x}_a - \boldsymbol{\mu}_a) - \frac{1}{2}(\mathbf{x}_b - \boldsymbol{\mu}_b)^T \boldsymbol{\Lambda}_{bb}(\mathbf{x}_b - \boldsymbol{\mu}_b). \end{aligned} \quad (6)$$

We see that as a function of  $\mathbf{x}_a$ , this is a quadratic form, thus the conditional distribution  $p(\mathbf{x}_a|\mathbf{x}_b)$  is a Gaussian distribution. Noticing that the exponent of a general Gaussian distribution  $\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma})$  can be written

$$-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}) = -\frac{1}{2}\mathbf{x}^T \boldsymbol{\Sigma}^{-1}\mathbf{x} + \mathbf{x}^T \boldsymbol{\Sigma}^{-1}\boldsymbol{\mu} + \text{const} \quad (7)$$

where 'const' denotes terms which are independent of  $\mathbf{x}$ . Consider the functional dependence of Eq. (6) on  $\mathbf{x}_b$  in which  $\mathbf{x}_a$  is regarded as a constant. The second order of  $\mathbf{x}_a$  can be written

$$-\frac{1}{2}\mathbf{x}_a^T \boldsymbol{\Lambda}_{aa}\mathbf{x}_a \quad (8)$$

Comparing Eq. (7) and Eq. (8), we can immediately get the covariance matrix of  $p(\mathbf{x}_a|\mathbf{x}_b)$  is given by

$$\Sigma_{a|b} = \Lambda_{aa}^{-1}. \quad (9)$$

Now consider the linear form of  $\mathbf{x}_a$  in Eq. (6)

$$\mathbf{x}_a^T \{ \Lambda_{aa} \boldsymbol{\mu}_a - \Lambda_{ab} (\mathbf{x}_b - \boldsymbol{\mu}_b) \} \quad (10)$$

Comparing Eq. (7) Eq. (10), we can get  $\Sigma^{-1} \boldsymbol{\mu} = \mathbf{x}_a^T \{ \Lambda_{aa} \boldsymbol{\mu}_a - \Lambda_{ab} (\mathbf{x}_b - \boldsymbol{\mu}_b) \}$ . Thus, we have

$$\begin{aligned} \boldsymbol{\mu}_{a|b} &= \Sigma_{a|b} \{ \Lambda_{aa} \boldsymbol{\mu}_a - \Lambda_{ab} (\mathbf{x}_b - \boldsymbol{\mu}_b) \} \\ &= \boldsymbol{\mu}_a - \Lambda_{ab}^{-1} \Lambda_{ab} (\mathbf{x}_b - \boldsymbol{\mu}_b) \end{aligned} \quad (11)$$

Combing Eq. (9) and Eq. (11), we can get  $p(\mathbf{x}_a|\mathbf{x}_b) = \mathcal{N}(\mathbf{x}_a|\boldsymbol{\mu}_{a|b}, \Lambda_{aa}^{-1})$  where  $\boldsymbol{\mu}_{a|b} = \boldsymbol{\mu}_a - \Lambda_{ab}^{-1} \Lambda_{ab} \mathbf{x}_b + \boldsymbol{\mu}_b$ .

### 3 Marginal Gaussian

In fact, the Marginal Gaussian distribution can be written as

$$p(\mathbf{x}_a) = \int p(\mathbf{x}_a, \mathbf{x}_b) d\mathbf{x}_b \quad (12)$$

Picking out those terms only involve  $\mathbf{x}_b$  in the joint distribution  $p(\mathbf{x}_a, \mathbf{x}_b)$ , we have

$$-\frac{1}{2} \mathbf{x}_b^T \Lambda_{bb} \mathbf{x}_b + \mathbf{x}_b^T \mathbf{m} = -\frac{1}{2} (\mathbf{x}_b - \Lambda_{bb}^{-1} \mathbf{m})^T \Lambda_{bb} (\mathbf{x}_b - \Lambda_{bb}^{-1} \mathbf{m}) + \frac{1}{2} \mathbf{m}^T \Lambda_{bb}^{-1} \mathbf{m} \quad (13)$$

where  $\mathbf{m} = \Lambda_{bb} \boldsymbol{\mu}_b - \Lambda_{ba} (\mathbf{x}_a - \boldsymbol{\mu}_a)$  Now we turn to consider the general Gaussian distribution, which is

$$\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \Sigma) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma|^{1/2}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right\} \quad (14)$$

From Eq. (14) we can see the coefficient of Gaussian distribution is independent of the mean and only governed by the determinant of the covariance matrix. Back to Eq. 13, we can see the integration over  $\mathbf{x}$  is as follows

$$\int \exp \left\{ -\frac{1}{2} (\mathbf{x}_b - \Lambda_{bb}^{-1} \mathbf{m})^T \Lambda_{bb} (\mathbf{x}_b - \Lambda_{bb}^{-1} \mathbf{m}) \right\} d\mathbf{x}_b. \quad (15)$$

This integration is irrelevant with the mean so we can margin out  $\mathbf{x}_b$  after the integration. Combing the  $\mathbf{x}_b^T \mathbf{m}$  in Eq. (13) with the remaining terms from Eq. (6), we have

$$\begin{aligned} &\frac{1}{2} [\Lambda_{bb} \boldsymbol{\mu}_b - \Lambda_{ba} (\mathbf{x}_a - \boldsymbol{\mu}_a)]^T \Lambda_{bb}^{-1} [\Lambda_{bb} \boldsymbol{\mu}_b - \Lambda_{ba} (\mathbf{x}_a - \boldsymbol{\mu}_a)] - \frac{1}{2} \mathbf{x}_a^T \Lambda_{aa} \mathbf{x}_a + \mathbf{x}_a^T (\Lambda_{aa} \boldsymbol{\mu}_a + \Lambda_{ab} \boldsymbol{\mu}_b) + \text{const} \\ &= \frac{1}{2} \mathbf{x}_a^T (\Lambda_{aa} - \Lambda_{ab} \Lambda_{bb}^{-1} \Lambda_{ba}) \mathbf{x}_a + \mathbf{x}_a^T (\Lambda_{aa} - \Lambda_{ab} \Lambda_{bb}^{-1} \Lambda_{ba}) \boldsymbol{\mu}_a + \text{const} \end{aligned} \quad (16)$$

where 'const' denotes quantities independent of  $\mathbf{x}_a$ . Comparing Eq. (16) with Eq. (7), we can get the covariance of  $p(\mathbf{x}_a)$  is

$$\Sigma_a = (\Lambda_{aa} - \Lambda_{ab} \Lambda_{bb}^{-1} \Lambda_{ba})^{-1} = \Sigma_{aa} \quad (17)$$

and the mean of  $p(\mathbf{x}_a)$  is given by

$$\boldsymbol{\mu}_a = \Sigma_a (\Lambda_{aa} - \Lambda_{ab} \Lambda_{bb}^{-1} \Lambda_{ba}) \boldsymbol{\mu}_a \quad (18)$$

We can summarized Eq. (17) and Eq. (18) as

$$p(\mathbf{x}_a) = \mathcal{N}(\mathbf{x}_a|\boldsymbol{\mu}_a, \Sigma_{aa}) \quad (19)$$

# Bayes' theorem for Gaussian variables

Xinghu Yao

September 22, 2018

## 1 Question

Given a marginal Gaussian distribution for  $\mathbf{x}$  and a conditional Gaussian distribution for  $\mathbf{y}$  given  $\mathbf{x}$  in the form

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

$$p(\mathbf{y}|\mathbf{x}) = \mathcal{N}(\mathbf{y}|\mathbf{A}\mathbf{x} + \mathbf{b}, \mathbf{L}^{-1}) \quad (2)$$

the marginal distribution of  $\mathbf{y}$  and the conditional distribution of  $\mathbf{x}$  given  $\mathbf{y}$  are given by

$$p(\mathbf{y}) = \mathcal{N}(\mathbf{y}|\mathbf{A}\boldsymbol{\mu} + \mathbf{b}, \mathbf{A}\boldsymbol{\Lambda}^{-1}\mathbf{A}^T) \quad (3)$$

$$p(\mathbf{x}|\mathbf{y}) = \mathcal{N}(\mathbf{x}|\boldsymbol{\Sigma}\{\mathbf{A}^T\mathbf{L}\mathbf{y} - \mathbf{b} + \boldsymbol{\Lambda}\boldsymbol{\mu}\}, \boldsymbol{\Sigma}) \quad (4)$$

where

$$\boldsymbol{\Sigma} = (\boldsymbol{\Lambda} + \mathbf{A}^T\mathbf{L}\mathbf{A})^{-1} \quad (5)$$

## 2 Solution

First, it is easy (just follow the definition of matrix multiplication) to prove the following equation for the inverse of a partitioned matrix

$$\begin{pmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{pmatrix} = \begin{pmatrix} \mathbf{M} & -\mathbf{M}\mathbf{B}\mathbf{D}^{-1} \\ -\mathbf{D}^{-1}\mathbf{C}\mathbf{M} & \mathbf{D}^{-1} + \mathbf{D}^{-1}\mathbf{C}\mathbf{M}\mathbf{B}\mathbf{D}^{-1} \end{pmatrix} \quad (6)$$

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

We define  $\mathbf{z} = (\mathbf{x}, \mathbf{y})^T$  as the joint distribution over  $\mathbf{x}$  and  $\mathbf{y}$  and consider the log of  $\mathbf{z}$

$$\begin{aligned} \ln p(\mathbf{z}) &= \ln p(\mathbf{x}) + \ln p(\mathbf{y}|\mathbf{x}) \\ &= -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Lambda}(\mathbf{x} - \boldsymbol{\mu}) - \frac{1}{2}(\mathbf{y} - \mathbf{A}\mathbf{x} - \mathbf{b})^T \mathbf{L}(\mathbf{y} - \mathbf{A}\mathbf{x} - \mathbf{b}) + \text{const} \end{aligned} \quad (7)$$

The second order terms in Eq. 7 can be written as

$$\begin{aligned} & -\frac{1}{2}\mathbf{x}^T(\boldsymbol{\Lambda} + \mathbf{A}^T\mathbf{L}\mathbf{A})\mathbf{x} - \frac{1}{2}\mathbf{y}^T\mathbf{L}\mathbf{y} + \frac{1}{2}\mathbf{y}^T\mathbf{L}\mathbf{A}\mathbf{x} + \frac{1}{2}\mathbf{x}^T\mathbf{A}^T\mathbf{L}\mathbf{y} \\ &= -\frac{1}{2}\begin{pmatrix} \mathbf{x} \\ \mathbf{y} \end{pmatrix}^T \begin{pmatrix} \boldsymbol{\Lambda} + \mathbf{A}^T\mathbf{L}\mathbf{A} & -\mathbf{A}^T\mathbf{L} \\ -\mathbf{L}\mathbf{A} & \mathbf{L} \end{pmatrix} \begin{pmatrix} \mathbf{x} \\ \mathbf{y} \end{pmatrix} \\ &= -\frac{1}{2}\mathbf{z}^T \mathbf{R} \mathbf{z} \end{aligned} \quad (8)$$

The linear terms in Eq. 7 can be written as

$$\mathbf{x}^T \boldsymbol{\Lambda} \boldsymbol{\mu} - \mathbf{x}^T \mathbf{A}^T \mathbf{L} \mathbf{b} + \mathbf{y}^T \mathbf{L} \mathbf{b} = \begin{pmatrix} \mathbf{x} \\ \mathbf{y} \end{pmatrix}^T \begin{pmatrix} \boldsymbol{\Lambda} \boldsymbol{\mu} - \mathbf{A}^T \mathbf{L} \mathbf{b} \\ \mathbf{L} \mathbf{b} \end{pmatrix}. \quad (9)$$

Thus, the inverse matrix of covariance matrix can be written as

$$\mathbf{R} = \begin{pmatrix} \boldsymbol{\Lambda} + \mathbf{A}^T\mathbf{L}\mathbf{A} & -\mathbf{A}^T\mathbf{L} \\ -\mathbf{L}\mathbf{A} & \mathbf{L} \end{pmatrix} \quad (10)$$

And using the Eq. 6, we have

$$\text{cov}[\mathbf{z}] = \mathbf{R}^{-1} = \begin{pmatrix} \mathbf{\Lambda}^{-1} & \mathbf{\Lambda}^{-1}\mathbf{A}^T \\ \mathbf{A}\mathbf{\Lambda}^{-1} & \mathbf{L}^{-1} + \mathbf{A}\mathbf{\Lambda}^{-1}\mathbf{A}^T \end{pmatrix}. \quad (11)$$

Similar to the former report, we have

$$\mathbb{E}[\mathbf{z}] = \mathbf{R}^{-1} \begin{pmatrix} \mathbf{A}\boldsymbol{\mu} - \mathbf{A}^T\mathbf{L}\mathbf{b} \\ \mathbf{L}\mathbf{b} \end{pmatrix} = \begin{pmatrix} \boldsymbol{\mu} \\ \mathbf{A}\boldsymbol{\mu} + \mathbf{b} \end{pmatrix}. \quad (12)$$

Making use of Eq. 11 and Eq. 12, we can get the mean and covariance of the marginal distribution  $p(\mathbf{y})$ , which is

$$\mathbb{E}[\mathbf{y}] = \mathbf{A}\boldsymbol{\mu} + \mathbf{b} \quad (13)$$

$$\text{cov}[\mathbf{y}] = \mathbf{L}^{-1} + \mathbf{A}\mathbf{\Lambda}^{-1}\mathbf{A}^T. \quad (14)$$

This means  $p(\mathbf{y}) = \mathcal{N}(\mathbf{y}|\mathbf{A}\boldsymbol{\mu} + \mathbf{b}, \mathbf{L}^{-1} + \mathbf{A}\mathbf{\Lambda}^{-1}\mathbf{A}^T)$ . By using the last reports results, we have

$$\mathbb{E}[\mathbf{x}|\mathbf{y}] = (\mathbf{\Lambda} + \mathbf{A}^T\mathbf{L}\mathbf{A}^{-1}) \{ \mathbf{A}^T\mathbf{L}(\mathbf{y} - \mathbf{b}) + \mathbf{\Lambda}\boldsymbol{\mu} \} \quad (15)$$

$$\text{cov}[\mathbf{x}|\mathbf{y}] = \mathbf{\Lambda} + \mathbf{A}^T\mathbf{L}\mathbf{A}^{-1}. \quad (16)$$

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

# Gradient and EM algorithms for PPCA

Xinghu Yao

September 24, 2018

## 1 Gradient method for PPCA

The log likelihood function of PPCA can be written as

$$\begin{aligned}\ln p(\mathbf{X}|\boldsymbol{\mu}, \mathbf{W}, \sigma^2) &= \sum_{n=1}^N \ln p(\mathbf{x}_n|\mathbf{W}, \boldsymbol{\mu}, \sigma^2) \\ &= -\frac{ND}{2} \ln(2\pi) - \frac{N}{2} \ln|\mathbf{C}| - \frac{1}{2} \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu})^T \mathbf{C}^{-1} (\mathbf{x}_n - \boldsymbol{\mu}).\end{aligned}\quad (1)$$

where the  $D \times D$  covariance matrix  $\mathbf{C}$  is defined by

$$\mathbf{C} = \mathbf{W}\mathbf{W}^T + \sigma^2 \mathbf{L}. \quad (2)$$

Setting the derivation w.r.t.  $\boldsymbol{\mu}$  equal to zero gives:

$$\boldsymbol{\mu} = -\frac{1}{N} \sum_{n=1}^N \mathbf{x}_n = \bar{\mathbf{x}}. \quad (3)$$

The log-likelihood is then simplified as:

$$\ln p(\mathbf{X}|\boldsymbol{\mu}, \mathbf{W}, \sigma^2) = -\frac{ND}{2} \ln(2\pi) - \frac{N}{2} \ln|\mathbf{C}| - \frac{1}{2} \sum_{n=1}^N (\mathbf{x}_n - \bar{\mathbf{x}})^T \mathbf{C}^{-1} (\mathbf{x}_n - \bar{\mathbf{x}}) \quad (4)$$

or can be written as:

$$\mathbf{L} \triangleq \ln p(\mathbf{X}|\boldsymbol{\mu}, \mathbf{W}, \sigma^2) = -\frac{N}{2} \{D \ln(2\pi) + \ln|\mathbf{C}| + \text{Tr}(\mathbf{C}^{-1} \mathbf{S})\} \quad (5)$$

where  $\mathbf{S}$  is the data covariance matrix defined by

$$\mathbf{S} = \frac{1}{N} \sum_{n=1}^N (\mathbf{x}_n - \bar{\mathbf{x}}) (\mathbf{x}_n - \bar{\mathbf{x}})^T \quad (6)$$

The gradient of the log-likelihood with respect to  $\mathbf{W}$  may be obtained from standard matrix differentiation results:

$$\frac{\partial \mathbf{L}}{\partial \mathbf{W}} = N (\mathbf{C}^{-1} \mathbf{S} \mathbf{C}^{-1} \mathbf{W} - \mathbf{C}^{-1} \mathbf{W}). \quad (7)$$

At the stationary points:

$$\mathbf{S} \mathbf{C}^{-1} \mathbf{W} = \mathbf{W} \quad (8)$$

By using SVD method and some interesting tricks, we can solve this problem and all solutions of  $\mathbf{W}$  can be written as

$$\mathbf{W}_{ML} = \mathbf{U}_M (\mathbf{L}_M - \sigma^2 \mathbf{I})^{1/2} \mathbf{R} \quad (9)$$

where  $\mathbf{U}_M$  is a  $D \times M$  matrix whose columns are given by any subset of the eigenvectors of the data covariance matrix  $\mathbf{S}$ , the  $M \times M$  diagonal matrix  $\mathbf{L}_M$  has elements given by the corresponding eigenvalues  $\lambda_i$ , and  $\mathbf{R}$  is an arbitrary  $M \times M$  orthogonal matrix. In fact, when the  $M$  largest eigenvalues are chosen, the maximum of

the likelihood function is obtained. In this case, the columns of  $\mathbf{W}$  define the principle subspace of standard PCA and the corresponding maximum likelihood solution for  $\sigma^2$  is then given by

$$\sigma_{ML}^2 = \frac{1}{D-M} \sum_{i=M+1}^D \lambda_i \quad (10)$$

which is the average of the discarded eigenvalues.

## 2 EM algorithms for PPCA

We first take the expectation of the complete-data log-likelihood w.r.t. the posterior distribution of the latent distribution evaluated using 'old' parameter values. Maximization of this expected complete data log-likelihood then yields the 'new' parameter values. The complete-data log-likelihood function takes the form

$$\ln p(\mathbf{X}, \mathbf{Z} | \boldsymbol{\mu}, \mathbf{W}, \sigma^2) = \sum_{n=1}^N \{\ln p(\mathbf{x}_n | \mathbf{z}_n) + \ln p(\mathbf{z}_n)\} \quad (11)$$

where the  $n^{\text{th}}$  row of the matrix  $\mathbf{Z}$  is given by  $\mathbf{z}_n$ . Recall that  $p(\mathbf{z}) = \mathcal{N}(\mathbf{z} | \mathbf{0}, \mathbf{I})$ ,  $p(\mathbf{x} | \mathbf{z}) = \mathcal{N}(\mathbf{x} | \mathbf{W}\mathbf{z} + \boldsymbol{\mu}, \sigma^2 \mathbf{I})$ . Thus, the expectation w.r.t. the posterior distribution over the latent variables can be written as

$$\mathbb{E} [\ln p(\mathbf{X}, \mathbf{Z} | \boldsymbol{\mu}, \mathbf{W}, \sigma^2)] = - \sum_{n=1}^N \left\{ \begin{aligned} &\frac{D}{2} \ln(2\pi\sigma^2) + \frac{1}{2} \text{Tr}(\mathbb{E}[\mathbf{z}_n \mathbf{z}_n^T]) \\ &+ \frac{1}{2\sigma^2} \|\mathbf{x}_n - \boldsymbol{\mu}\|^2 - \frac{1}{\sigma^2} \mathbb{E}[\mathbf{z}_n]^T \mathbf{W}^T (\mathbf{x}_n - \boldsymbol{\mu}) \\ &+ \frac{1}{2\sigma^2} \text{Tr}(\mathbb{E}[\mathbf{z}_n \mathbf{z}_n^T] \mathbf{W}^T \mathbf{W}) + \frac{M}{2 \ln(2\pi)} \end{aligned} \right\} \quad (12)$$

**E-Step:** We use the old parameter to evaluate

$$\mathbb{E}[\mathbf{z}_n] = \mathbf{M}^{-1} \mathbf{W}^T (\mathbf{x}_n - \bar{\mathbf{x}}) \quad (13)$$

$$\mathbb{E}[\mathbf{z}_n \mathbf{z}_n^T] = \sigma^2 \mathbf{M}^{-1} + [\mathbf{z}_n][\mathbf{z}_n]^T \quad (14)$$

This follows directly from

$$p(\mathbf{z} | \mathbf{x}) = \mathcal{N}(\mathbf{z} | \mathbf{M}^{-1} \mathbf{W}^T (\mathbf{x} - \boldsymbol{\mu}), \sigma^2 \mathbf{M}^{-1}), \mathbf{M} = \mathbf{W}^T \mathbf{W} + \sigma^2 \mathbf{I} \quad (15)$$

together with the standard result

$$\mathbb{E}[\mathbf{z}_n \mathbf{z}_n^T] = \text{cov}[\mathbf{z}] + \mathbb{E}[\mathbf{z}_n] \mathbb{E}[\mathbf{z}_n]^T \quad (16)$$

Substituting Eq. (15) and Eq. (16) into Eq. (12), we can compute the expectation result.

**M-Step:** We can get the two M-equations by setting the derivatives w.r.t  $\mathbf{W}$  and  $\sigma^2$  to zero, which is

$$\frac{\partial \mathbb{E} [\ln p(\mathbf{X}, \mathbf{Z} | \boldsymbol{\mu}, \mathbf{W}, \sigma^2)]}{\partial \mathbf{W}} = \sum_{n=1}^N \left\{ \frac{1}{\sigma^2} (\mathbf{x}_n - \boldsymbol{\mu}) \mathbb{E}[\mathbf{z}_n]^T - \frac{1}{\sigma^2} \mathbf{W} \mathbb{E}[\mathbf{z}_n \mathbf{z}_n^T] \right\} = 0 \quad (17)$$

$$\begin{aligned} \frac{\partial \mathbb{E} [\ln p(\mathbf{X}, \mathbf{Z} | \boldsymbol{\mu}, \mathbf{W}, \sigma^2)]}{\partial \sigma^2} &= \sum_{n=1}^N \left\{ -\frac{D}{2\sigma^2} - \frac{1}{\sigma^4} \mathbb{E}[\mathbf{z}_n]^T \mathbf{W}^T (\mathbf{x}_n - \boldsymbol{\mu}) + \right. \\ &\quad \left. \frac{1}{2\sigma^4} \|\mathbf{x}_n - \boldsymbol{\mu}\|^2 + \frac{1}{2\sigma^4} \text{Tr}(\mathbb{E}[\mathbf{z}_n \mathbf{z}_n^T] \mathbf{W}^T \mathbf{W}) \right\} = 0 \end{aligned} \quad (18)$$

It is worth to say that we used  $\frac{\partial}{\partial \mathbf{A}} \text{Tr}(\mathbf{A} \mathbf{B} \mathbf{A}^T) = \mathbf{A}(\mathbf{B} + \mathbf{B}^T)$ ,  $\frac{\partial}{\partial \mathbf{A}} \text{Tr}(\mathbf{A} \mathbf{B}) = \mathbf{B}^T$ . So, it is clear that we can get the following equations

$$\mathbf{W}_{\text{new}} = \left[ \sum_{n=1}^N (\mathbf{x}_n - \bar{\mathbf{x}}) \mathbb{E}[\mathbf{z}_n]^T \right] \left[ \sum_{n=1}^N \mathbb{E}[\mathbf{z}_n \mathbf{z}_n^T] \right]^{-1} \quad (19)$$

$$\sigma_{\text{new}}^2 = \frac{1}{ND} \sum_{n=1}^N \left\{ \|\mathbf{x}_n - \bar{\mathbf{x}}\|^2 - 2 \mathbb{E}[\mathbf{z}_n]^T \mathbf{W}_{\text{new}}^T (\mathbf{x}_n - \bar{\mathbf{x}}) + \text{Tr}(\mathbb{E}[\mathbf{z}_n \mathbf{z}_n^T] \mathbf{W}_{\text{new}}^T \mathbf{W}_{\text{new}}) \right\}. \quad (20)$$

# Bayesian Model Comparison

Xinghu Yao

September 25, 2018

## 1 Question

The Bayesian view of model comparison simply involves the use of probabilities to represent uncertainty in the choice of model, along with a consistent application of the sum and product rules of probability. Consider a data set  $\mathcal{D}$  and a set of models  $\{\mathcal{M}_i\}$  having parameters  $\{\theta_i\}$ . For each model we define a likelihood function  $p(\mathcal{D}|\theta_i, \mathcal{M}_i)$ . If we introduce a prior  $p(\theta_i|\mathcal{M}_i)$  for the various models. Try to approximate the distribution as follows using the Laplace Approximation.

$$p(\mathcal{D}) = \int p(\mathcal{D}|\theta)p(\theta)d\theta \quad (1)$$

Note that in Eq. (1) each  $\mathcal{M}_i$  is omitted to keep the notation uncluttered.

## 2 The Laplace Approximation

Laplace approximation is a simple but widely used framework which aims to find a Gaussian approximations to a probability density defined over a set of continuous variables. For a given distribution  $p(\mathbf{z}) = f(\mathbf{z})/Z$  defined over an  $M$ -dimensional space  $\mathbf{z}$ . At a stationary point  $\mathbf{z}_0$  the gradient  $\nabla f(\mathbf{z})$  will vanish. We therefor consider a Taylor expansion of  $\ln f(\mathbf{z})$  centred on the mode  $\mathbf{z}_0$  so that

$$\ln f(\mathbf{z}) \simeq \ln f(\mathbf{z}_0) - \frac{1}{2}(\mathbf{z} - \mathbf{z}_0)^T \mathbf{A}(\mathbf{z} - \mathbf{z}_0) \quad (2)$$

Where the  $M \times M$  Hessian matrix  $A$  is defined by

$$\mathbf{A} = -\nabla \nabla \ln f(\mathbf{z})|_{\mathbf{z}=\mathbf{z}_0} \quad (3)$$

and  $\nabla$  is the gradient operator. Taking the exponential of both sides we obtain

$$f(\mathbf{z}) \simeq f(\mathbf{z}_0) \exp \left\{ -\frac{1}{2}(\mathbf{z} - \mathbf{z}_0)^T \mathbf{A}(\mathbf{z} - \mathbf{z}_0) \right\} \quad (4)$$

The distribution  $q(\mathbf{z})$  is proportional to  $f(\mathbf{z})$  and the appropriate normalization coefficient can be found by using the standard result for a normalized multivariate Gaussian, giving

$$q(\mathbf{z}) = \frac{|\mathbf{A}|^{1/2}}{(2\pi)^{M/2}} \exp \left\{ -\frac{1}{2}(\mathbf{z} - \mathbf{z}_0)^T \mathbf{A}(\mathbf{z} - \mathbf{z}_0) \right\} = \mathcal{N}(\mathbf{z}|\mathbf{z}_0, \mathbf{A}^{-1}) \quad (5)$$

where  $|\mathbf{A}|$  denotes the determinant of  $\mathbf{A}$ . This Gaussian distribution will be well defined provided its precision matrix, given by  $A$ , is positive definite, which implies that the stationary point  $\mathbf{z}_0$  must be a local maximum, not a minimum of a saddle point. In fact, we can also obtain an approximation to the normalization  $Z$ . Because we have

$$q(\mathbf{z}) \simeq p(\mathbf{z}) = \frac{1}{Z} f(\mathbf{z}) \quad (6)$$

Thus, the approximation to the normalization constant  $Z$  have the following form

$$Z \simeq \frac{f(\mathbf{z})}{q(\mathbf{z})} = f(\mathbf{z}_0) \frac{(2\pi)^{M/2}}{|\mathbf{A}|^{1/2}} \quad (7)$$

### 3 Solution

Identifying  $f(\boldsymbol{\theta}) = p(\mathcal{D}|\boldsymbol{\theta})p(\boldsymbol{\theta})$  and  $Z = p(\mathcal{D})$ . According to Eq. (7), we have

$$\begin{aligned}\ln p(\mathcal{D}) &= \ln(Z) \simeq \ln f(\mathbf{z}_0) \frac{(2\pi)^{M/2}}{|\mathbf{A}|^{1/2}} \\ &= \ln p(\mathcal{D}|\boldsymbol{\theta}_{\text{MAP}}) + \ln p(\boldsymbol{\theta}_{\text{MAP}}) + \frac{M}{2}\ln(2\pi) - \frac{1}{2}\ln|\mathbf{A}|\end{aligned}\quad (8)$$

where  $\boldsymbol{\theta}_{\text{MAP}}$  is the value of  $\boldsymbol{\theta}$  at the mode of the posterior distribution, and  $\mathbf{A}$  is the Hessian matrix of second derivatives of the negative log posterior

$$\mathbf{A} = -\nabla\nabla\ln p(\mathcal{D}|\boldsymbol{\theta}_{\text{MAP}})p(\boldsymbol{\theta}_{\text{MAP}}) = -\nabla\nabla\ln p(\boldsymbol{\theta}_{\text{MAP}}|\mathcal{D}) \quad (9)$$

From Eq. (9), we have

$$\begin{aligned}\mathbf{A} &= -\nabla\nabla\ln p(\mathcal{D}|\boldsymbol{\theta}_{\text{MAP}})p(\boldsymbol{\theta}_{\text{MAP}}) \\ &= \mathbf{H} - \nabla\nabla\ln p(\boldsymbol{\theta}_{\text{MAP}})\end{aligned}\quad (10)$$

and if  $p(\boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\theta}|\mathbf{m}, \mathbf{V}_0)$ , this becomes

$$\mathbf{A} = \mathbf{H} + \mathbf{V}_0^{-1}. \quad (11)$$

If we assume that the prior is broad or equivalently that the number of data points is large, we can neglect the term  $\mathbf{V}_0^{-1}$  compared to  $\mathbf{H}$ . Using this result, Eq. (8) can be rewritten in the form

$$\ln p(\mathcal{D}) \simeq \ln p(\mathcal{D}|\boldsymbol{\theta}_{\text{MAP}}) - \frac{1}{2}(\boldsymbol{\theta}_{\text{MAP}} - \mathbf{m})^T \mathbf{V}_0^{-1}(\boldsymbol{\theta}_{\text{MAP}} - \mathbf{m}) - \frac{1}{2}\ln|\mathbf{H}| + \text{const} \quad (12)$$

We now again invoke the broad prior assumption, allowing us to neglect the second term on the right hand side of Eq. (12). Since we assume i.i.d data,  $\mathbf{H} = -\nabla\nabla\ln p(\mathcal{D}|\boldsymbol{\theta}_{\text{MAP}})$  consists of a sum of terms, one term for each datum, and we can consider the following approximation:

$$\mathbf{H} = \sum_{n=1}^N \mathbf{H}_n = N\hat{\mathbf{H}} \quad (13)$$

where  $\mathbf{H}_n$  is the contribution from the  $n^{\text{th}}$  data point and

$$\hat{\mathbf{H}} = \frac{1}{N} \sum_{n=1}^N \mathbf{H}_n \quad (14)$$

Combining this with the properties of the determinant, we have

$$\ln|\mathbf{H}| = \ln|N\hat{\mathbf{H}}| = \ln\left(N^M|\hat{\mathbf{H}}|\right) = M\ln N + \ln|\hat{\mathbf{H}}| \quad (15)$$

Where  $M$  is the dimensionality of  $\boldsymbol{\theta}$ . Note that we are assuming that  $\hat{\mathbf{H}}$  has full rank  $M$ . Finally, using this result together Eq. (12), we obtain

$$\ln p(\mathcal{D}) \simeq \ln p(\mathcal{D}|\boldsymbol{\theta}_{\text{MAP}}) - \frac{1}{2}M\ln N \quad (16)$$



# Auto-Encoding Variational Bayes

Xinghu Yao

September 26, 2018

## 1 Problem Establishment

Using auto-encoding variational Bayes we can perform efficient approximate inference and learning with directed probabilistic models whose continuous latent variables and/or parameters have intractable posterior distributions.

Let us consider some dataset  $\mathbf{X} = \{\mathbf{x}^{(i)}\}_{i=1}^N$  consisting of  $N$  i.i.d samples of some continuous or discrete variable  $\mathbf{x}$ . We assume that the data are generated by some random process, involving an unobserved continuous random variable  $\mathbf{z}$ . The process consists of two steps: (1) a value  $\mathbf{z}^{(i)}$  is generated from some prior distribution  $p_{\theta}(\mathbf{z})$ . (2) a value  $\mathbf{x}^{(i)}$  is generated from some conditional distribution  $p_{\theta}(\mathbf{x}|\mathbf{z})$ .

We can use two networks to train a generative model and the two parts of our networks have the following relationship.

$$\mathbf{x} \xrightarrow{F_{\theta}} (\mathbf{z}|\mathbf{x}) \xrightarrow{G_{\theta}} \hat{\mathbf{x}} \quad (1)$$

where  $F_{\theta}$  is a network to cover the latent hidden variable  $\mathbf{z}$  and  $G_{\theta}$  is another network to decode  $\mathbf{x}$  using the hidden variable  $\mathbf{z}$ . Thus, the marginal likelihood of this structure can be written as:

$$\begin{aligned} \ln p(\mathbf{x}) &= \ln \int p(\mathbf{x}|\mathbf{z})p(\mathbf{z})d\mathbf{z} \\ &= \ln \int \frac{p(\mathbf{x}|\mathbf{z})}{q(\mathbf{z}|\mathbf{x})}p(\mathbf{z})q(\mathbf{z}|\mathbf{x})d\mathbf{z} \\ &= \ln \mathbb{E}_{\mathbf{z} \sim q(\mathbf{z}|\mathbf{x})} \left( \frac{p(\mathbf{x}|\mathbf{z})}{q(\mathbf{z}|\mathbf{x})}p(\mathbf{z}) \right) \\ &\geq \mathbb{E}_{\mathbf{z} \sim q(\mathbf{z}|\mathbf{x})} \ln \left( \frac{p(\mathbf{x}|\mathbf{z})}{q(\mathbf{z}|\mathbf{x})}p(\mathbf{z}) \right) \quad (\text{Jensen Inequality}) \\ &= \mathbb{E}_{\mathbf{z} \sim q(\mathbf{z}|\mathbf{x})} \ln \frac{p(\mathbf{x}|\mathbf{z})}{q(\mathbf{z}|\mathbf{x})} + \mathbb{E}_{\mathbf{z} \sim q(\mathbf{z}|\mathbf{x})} \ln p(\mathbf{z}) \\ &= \int q(\mathbf{z}|\mathbf{x}) \ln p(\mathbf{x}|\mathbf{z})d\mathbf{z} - \text{KL} [q(\mathbf{z}|\mathbf{x})||p(\mathbf{z})] \end{aligned} \quad (2)$$

where  $q(\mathbf{z}|\mathbf{x})$  is the probabilistic encoder, since given a datapoint  $\mathbf{x}$  it produce a distribution over the possible values of the code  $\mathbf{z}$  from which the datapoint  $\mathbf{x}$  could have been generated. In a similar vein we will refer to  $p(\mathbf{x}|\mathbf{z})$  as a probabilistic decoder, since given a code  $\mathbf{z}$  it produce a distribution over the possible corresponding values of  $\mathbf{x}$ . Thus, we can get the following optimization problem

$$\max_{q(\mathbf{z}|\mathbf{x})} \mathbb{E}_{\mathbf{z} \sim q(\mathbf{z}|\mathbf{x})} [\ln p(\mathbf{x}|\mathbf{z})] - \text{KL} [q(\mathbf{z}|\mathbf{x})||p(\mathbf{z})] \quad (3)$$

## 2 The reparameterization trick

In order to solve our problem we invoke an alternative method for generating samples from  $q(\mathbf{z}|\mathbf{x})$ . The essential parameterization trick is quite simple. Let  $z \sim p(z|x) = \mathcal{N}(\mu, \sigma^2)$ . In this case, a valid reparameterization is  $z = \mu + \sigma\epsilon$ , where  $\epsilon$  is an auxiliary noise variable  $\epsilon \sim \mathcal{N}(0, 1)$