

2. Gaussian Distribution

(1) Introduction and Background

$$N(x|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{(x-\mu)^2}{2\sigma^2}\right\}$$

$$N(\underline{x}|\underline{\mu}, \underline{\Sigma}) = \frac{1}{\sqrt{2\pi^D|\underline{\Sigma}|}} \exp\left\{-\frac{1}{2}(\underline{x}-\underline{\mu})^T \underline{\Sigma}^{-1} (\underline{x}-\underline{\mu})\right\}$$

D is the dimension of \underline{x}

Because of very good properties for Gaussians, say Central Limit Theorem, we will focus greatly on this distribution.

the value

$$\Delta^2 = (\underline{x}-\underline{\mu})^T \underline{\Sigma}^{-1} (\underline{x}-\underline{\mu})$$

is called Mahalanobis distance

Here we will state a very useful hypothesis that

Matrix $\underline{\Sigma}$ can be viewed as symmetric

Actually, this assumption is not limited, We know that what really matters is the value $(\underline{x}-\underline{\mu})^T \underline{\Sigma}^{-1} (\underline{x}-\underline{\mu})$, which is a scalar. Also, $(\underline{x}-\underline{\mu})^T (\underline{x}-\underline{\mu})$ is a scalar. So we can find a magic number λ such that $(\underline{x}-\underline{\mu})^T \underline{\lambda \cdot I} (\underline{x}-\underline{\mu})$. We find another matrix λI which could act the same effect as $\underline{\Sigma}^{-1}$. So, why not choose a symmetric matrix?

Now, $\underline{\Sigma}$ is a real symmetric matrix. According to Spectrum Theorem, $\underline{\Sigma}$'s eigenvalues are real. $\underline{\Sigma}$ can be decomposed as

$$\underline{\Sigma} = \sum_{i=1}^r \lambda_i \mathbf{q}_i \mathbf{q}_i^T \quad \begin{array}{l} r \text{ is the rank of } \underline{\Sigma}. \text{ Here it must be full rank.} \\ \text{So } r = D \end{array}$$

Then

$$\underline{\Sigma}^{-1} = \sum_{i=1}^D \frac{1}{\lambda_i} \mathbf{q}_i \mathbf{q}_i^T$$

$$\begin{aligned} \text{So, } \Delta^2 &= (\underline{x} - \underline{\mu})^T \left(\sum_{i=1}^D \frac{1}{\lambda_i} \mathbf{q}_i \mathbf{q}_i^T \right) (\underline{x} - \underline{\mu}) \\ &= \sum_{i=1}^D \frac{y_i^2}{\lambda_i} \quad y_i = (\underline{x} - \underline{\mu})^T \cdot \mathbf{q}_i \end{aligned}$$

So if we can have a group of points which could guarantee $y_i = (\underline{x} - \underline{\mu})^T \mathbf{q}_i$ to be a constant, then the corresponding Gaussian distribution is constant. Mathematically, ^{if} all the λ_i are positive, points \underline{x} are on an ellipse. λ_i is the axis length. Actually, $\{\lambda_i\}$ are non-negative, or the Gaussian distribution will not be normalized. In other words, $\underline{\Sigma}$ is PSD.

(2) Conditional and Marginal Gaussian Distribution

First, we throw out the conclusion.

If two sets of variables are jointly Gaussian, the conditional distribution of one set on the other is also Gaussian. Marginal Distribution is again Gaussian

Now we will prove these important conclusions.

① Conditioned Distribution $P(\underline{x}_a | \underline{x}_b)$

Now \underline{x} is a D -dimensional vector with Gaussian distribution $\mathcal{N}(\underline{x} | \underline{\mu}, \underline{\Sigma})$. Let's split it into two parts $\underline{x}_a, \underline{x}_b$.

$$\underline{x} = \begin{pmatrix} \underline{x}_a \\ \underline{x}_b \end{pmatrix} \begin{matrix} M \\ D-M \end{matrix} \quad \underline{\mu} = \begin{pmatrix} \underline{\mu}_a \\ \underline{\mu}_b \end{pmatrix} \quad \underline{\Sigma} = \begin{pmatrix} \Sigma_a & \Sigma_{ab} \\ \Sigma_{ba} & \Sigma_b \end{pmatrix}$$

Note that $\underline{\Sigma}$ is symmetric, so $\Sigma_{ab} = \Sigma_{ba}^T$, Σ_a & Σ_b are symmetric.

Let's denote

$$\underline{\Lambda} = \underline{\Sigma}^{-1} \quad \underline{\Lambda} = \begin{pmatrix} \Lambda_a & \Lambda_{ab} \\ \Lambda_{ba} & \Lambda_b \end{pmatrix}$$

$\underline{\Lambda}$ is also a symmetric matrix. $\underline{\Lambda}^T = \underline{\Lambda}$

The exponent part

$$\begin{aligned} -\frac{1}{2} (\underline{x} - \underline{\mu})^T \underline{\Sigma}^{-1} (\underline{x} - \underline{\mu}) &= -\frac{1}{2} \begin{pmatrix} \underline{x}_a - \underline{\mu}_a \\ \underline{x}_b - \underline{\mu}_b \end{pmatrix}^T \begin{pmatrix} \Lambda_a & \Lambda_{ab} \\ \Lambda_{ba} & \Lambda_b \end{pmatrix} \begin{pmatrix} \underline{x}_a - \underline{\mu}_a \\ \underline{x}_b - \underline{\mu}_b \end{pmatrix} \\ &= -\frac{1}{2} \left\{ \begin{aligned} &(\underline{x}_a - \underline{\mu}_a)^T \Lambda_a (\underline{x}_a - \underline{\mu}_a) + (\underline{x}_a - \underline{\mu}_a)^T \Lambda_{ab} (\underline{x}_b - \underline{\mu}_b) \\ &+ (\underline{x}_b - \underline{\mu}_b)^T \Lambda_{ba} (\underline{x}_a - \underline{\mu}_a) + (\underline{x}_b - \underline{\mu}_b)^T \Lambda_b (\underline{x}_b - \underline{\mu}_b) \end{aligned} \right\} \quad \dots \quad (*) \end{aligned}$$

We again get quadratic form. because we get terms like $a^T \Sigma b$

So the result will also be a Gaussian distribution.

Our target now turns to find the mean μ_{ab} and Covariance Σ_{ab} .

Our method is coefficient analysis.

Because x_b is given, we can simply treat μ_b, x_b as constant.

For a quadratic form

$$-\frac{1}{2} (\underline{x} - \underline{\mu})^T \underline{\Sigma}^{-1} (\underline{x} - \underline{\mu}) = -\frac{1}{2} \underline{x}^T \underline{\Sigma}^{-1} \underline{x} + \underline{x}^T \underline{\Sigma}^{-1} \underline{\mu} - \frac{1}{2} \underline{\mu}^T \underline{\Sigma}^{-1} \underline{\mu}$$

We can find out the coefficient corresponding to $x^T x$ and x no x , constant!

In (X) expression, we have

$$-\frac{1}{2} \underline{x}_a^T \Delta_a \underline{x}_a \rightsquigarrow -\frac{1}{2} \underline{x}^T \underline{\Sigma}^{-1} \underline{x}$$

$$\underline{x}_a^T \{ \Delta_a \underline{\mu}_a - \Delta_{ab} (\underline{x}_b - \underline{\mu}_b) \} \rightsquigarrow \underline{x}^T \underline{\Sigma}^{-1} \underline{\mu}$$

So, we have

$$\Delta_{ab} = \Delta_a, \quad \mu_{ab} = \Delta_a^{-1} \{ \Delta_a \mu_a - \Delta_{ab} (\underline{x}_b - \underline{\mu}_b) \} \\ = \mu_a - \Delta_a^{-1} \Delta_{ab} (\underline{x}_b - \underline{\mu}_b)$$

For partition matrix

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

$$M = (A - BD^{-1}C)^{-1}$$

\uparrow
Schur complement r.p.t
D

$$\text{So } \Lambda_a = (\Sigma_a - \Sigma_{ab} \Sigma_b^{-1} \Sigma_{ba})^{-1}$$

$$\Lambda_{ab} = -(\Sigma_a - \Sigma_{ab} \Sigma_b^{-1} \Sigma_{ba})^{-1} \Sigma_{ab} \Sigma_b^{-1}$$

Now, the mean and variance for $\underline{x}_a | \underline{x}_b$ is

$$\underline{\mu}_{a|b} = \underline{\mu}_a + \Sigma_{ab} \Sigma_b^{-1} (\underline{x}_b - \underline{\mu}_b)$$

$$\Sigma_{a|b} = \Sigma_a - \Sigma_{ab} \Sigma_b^{-1} \Sigma_{ba}$$

\Rightarrow no \underline{x}_a , only linearly dependent on \underline{x}_b

② Marginal Gaussian Distribution

Marginal Distribution for \underline{x}_a is

$$P(\underline{x}_a) = \int_{\underline{x}_b} P(\underline{x}_a, \underline{x}_b) d\underline{x}_b$$

So we will extract all terms involving \underline{x}_b and treat the rest as constant.

According to (*) 2 pages before, terms relate to \underline{x}_b are

$$-\frac{1}{2} \underline{x}_b^T \Lambda_b \underline{x}_b + \underline{x}_b^T \underline{m} \quad \underline{m} = \Lambda_b \underline{\mu}_b - \Lambda_{ba} (\underline{x}_a - \underline{\mu}_a)$$

make it a quadratic form

$$-\frac{1}{2} (\underline{x}_b - \Lambda_b^{-1} \underline{m})^T \Lambda_b (\underline{x}_b - \Lambda_b^{-1} \underline{m}) + \underbrace{\frac{1}{2} \underline{m}^T \Lambda_b^{-1} \underline{m}}_{\text{no } \underline{x}_b \text{ terms}}$$

$$\int \underbrace{\exp\left\{-\frac{1}{2} (\underline{x}_b - \Lambda_b^{-1} \underline{m})^T \Lambda_b (\underline{x}_b - \Lambda_b^{-1} \underline{m})\right\}}_{\text{un-normalized Gaussian}} d\underline{x}_b = \sqrt{(2\pi)^{D-M} |\Sigma_b|}$$

un-normalized Gaussian

Now all the terms independent from \underline{x}_b are

$$\frac{1}{2} \underline{m}^T \underline{\Lambda}_b^{-1} \underline{m} - \frac{1}{2} \underline{x}_a^T \underline{\Lambda}_a \underline{x}_a + \underline{x}_a^T (\underline{\Lambda}_a \underline{\mu}_a + \underline{\Lambda}_{ab} \underline{\mu}_b) + \text{const}$$

no $\underline{x}_a, \underline{x}_b$

$$= \frac{1}{2} [\underline{\Lambda}_b \underline{\mu}_b - \underline{\Lambda}_{ba} (\underline{x}_a - \underline{\mu}_a)]^T \underline{\Lambda}_b^{-1} [\underline{\Lambda}_b \underline{\mu}_b - \underline{\Lambda}_{ba} (\underline{x}_a - \underline{\mu}_a)]$$

$$- \frac{1}{2} \underline{x}_a^T \underline{\Lambda}_a \underline{x}_a + \underline{x}_a^T (\underline{\Lambda}_a \underline{\mu}_a + \underline{\Lambda}_{ab} \underline{\mu}_b) + \text{const}$$

$$= -\frac{1}{2} \underline{x}_a^T (\underline{\Lambda}_a - \underline{\Lambda}_{ab} \underline{\Lambda}_b^{-1} \underline{\Lambda}_{ba}) \underline{x}_a + \underline{x}_a^T (\underline{\Lambda}_a - \underline{\Lambda}_{ab} \underline{\Lambda}_b^{-1} \underline{\Lambda}_{ba}) \underline{\mu}_a + \text{const}$$

$$= -\frac{1}{2} \underline{x}_a^T \underline{\Sigma}_a^{-1} \underline{x}_a + \underline{x}_a^T \underline{\Sigma}_a^{-1} \underline{\mu}_a + \text{const}$$

Recall that, a valid quadratic form has.

$$-\frac{1}{2} (\underline{x} - \underline{\mu})^T \underline{\Sigma}^{-1} (\underline{x} - \underline{\mu}) = -\frac{1}{2} \underline{x}^T \underline{\Sigma}^{-1} \underline{x} + \underline{x}^T \underline{\Sigma}^{-1} \underline{\mu} + \text{const}$$

So, here we have the mean and variance for marginal distribution of \underline{x}_a

$$\boxed{E(\underline{x}_a) = \underline{\mu}_a \quad \text{Cov}(\underline{x}_a) = \underline{\Sigma}_a}$$

(3) Bayes' Theorem for Gaussian Variables

In section (2), we give out the mean and variance for conditional and marginal Gaussian Variables. Which is

$$\begin{array}{ccc} \text{Given} & \begin{array}{l} P(\underline{x}, \underline{y}) \\ \underline{\mu}_x \ \underline{\mu}_y \\ \Sigma_x, \Sigma_y, \Sigma_{xy}, \Sigma_{yx} \end{array} & \Rightarrow \begin{array}{l} P(\underline{x}|\underline{y}) \ P(\underline{y}|\underline{x}) \\ P(\underline{x}) \ P(\underline{y}) \end{array} \end{array}$$

In this section, we will apply Bayes Law to do the inverse. We know that, $P(\underline{x}|\underline{y})$ is a linear function of \underline{y} . So let's use these start condition

$$\begin{array}{l} \text{Given } P(\underline{x}) = \mathcal{N}(\underline{x} | \underline{\mu}, \Lambda^{-1}) \\ P(\underline{y}|\underline{x}) = \mathcal{N}(\underline{y} | A\underline{x} + \underline{b}, L^{-1}) \end{array}$$

① find out the joint distribution

$$\text{let } \underline{z} = \begin{pmatrix} \underline{x} \\ \underline{y} \end{pmatrix}$$

$$\begin{aligned} \log P(\underline{z}) &= \ln P(\underline{x}) + \ln P(\underline{y}|\underline{x}) \\ &= -\frac{1}{2}(\underline{x} - \underline{\mu})^T \Lambda (\underline{x} - \underline{\mu}) - \frac{1}{2}(\underline{y} - A\underline{x} - \underline{b})^T L (\underline{y} - A\underline{x} - \underline{b}) + \text{const} \end{aligned}$$

Consider all the second order term (term involves two x or two y or x and y).

$$\begin{aligned} & -\frac{1}{2} \underline{x}^T (\Lambda + A^T L A) \underline{x} - \frac{1}{2} \underline{y}^T L \underline{y} + \frac{1}{2} \underline{y}^T L A \underline{x} + \frac{1}{2} \underline{x}^T A^T L \underline{y} \\ &= -\frac{1}{2} \begin{pmatrix} \underline{x} \\ \underline{y} \end{pmatrix}^T \underbrace{\begin{pmatrix} \Lambda + A^T L A & -A^T L \\ -L A & L \end{pmatrix}}_R \begin{pmatrix} \underline{x} \\ \underline{y} \end{pmatrix} = -\frac{1}{2} \underline{z}^T R \underline{z} \end{aligned}$$

Again comparing this with the expansion of quadratic form of Gaussian,

$$\text{Cov}(\underline{z}) = \mathbf{R}^{-1} = \begin{pmatrix} \Delta^{-1} & \Delta^{-1} \mathbf{A}^T \\ \mathbf{A} \Delta^{-1} & \mathbf{L}^{-1} + \mathbf{A} \Delta^{-1} \mathbf{A}^T \end{pmatrix}$$

Now inspect the linear terms (involving only one x or y)

$$\underline{x}^T \Delta \underline{\mu} - \underline{x}^T \mathbf{A}^T \mathbf{L} \underline{b} + \underline{y}^T \mathbf{L} \underline{b} = \begin{pmatrix} \underline{x} \\ \underline{y} \end{pmatrix}^T \begin{pmatrix} \Delta \underline{\mu} - \mathbf{A}^T \mathbf{L} \underline{b} \\ \mathbf{L} \underline{b} \end{pmatrix}$$

So

$$\bar{\mathbb{H}}(\underline{z}) = \mathbf{R}^{-1} \begin{pmatrix} \Delta \underline{\mu} - \mathbf{A}^T \mathbf{L} \underline{b} \\ \mathbf{L} \underline{b} \end{pmatrix} = \begin{pmatrix} \underline{\mu} \\ \mathbf{A} \underline{\mu} + \underline{b} \end{pmatrix}$$

(4) Maximum Likelihood for Gaussian.

Now our dataset is $X = (\underline{x}_1, \underline{x}_2, \dots, \underline{x}_N)^T$. The log likelihood is

$$\ln P(X|\underline{\mu}, \Sigma) = -\frac{ND}{2} \ln(2\pi) - \frac{N}{2} \ln|\Sigma| - \frac{1}{2} \sum_{n=1}^N (\underline{x}_n - \underline{\mu})^T \Sigma^{-1} (\underline{x}_n - \underline{\mu})$$

$$\frac{\partial}{\partial \underline{\mu}} \ln P(X|\underline{\mu}, \Sigma) = \sum_{n=1}^N \Sigma^{-1} (\underline{x}_n - \underline{\mu}) = 0 \Rightarrow \underline{\mu}_{ML} = \frac{1}{N} \sum_{n=1}^N \underline{x}_n$$

$$\frac{\partial}{\partial \Sigma} \ln P(X|\underline{\mu}, \Sigma) = \frac{\partial}{\partial \Sigma} \left(-\frac{N}{2} \ln|\Sigma| - \frac{1}{2} \sum_{n=1}^N (\underline{x}_n - \underline{\mu})^T \Sigma^{-1} (\underline{x}_n - \underline{\mu}) \right)$$

According to «The Matrix Cookbook»

$$\frac{\partial}{\partial \Sigma} -\frac{N}{2} \ln|\Sigma| = -\frac{N}{2} \Sigma^{-T}$$

$$\frac{\partial}{\partial \Sigma} -\frac{1}{2} \sum_{n=1}^N (\underline{x}_n - \underline{\mu})^T \Sigma^{-1} (\underline{x}_n - \underline{\mu}) = \frac{1}{2} \sum_{n=1}^N \Sigma^{-T} (\underline{x}_n - \underline{\mu}) (\underline{x}_n - \underline{\mu})^T \Sigma^{-1}$$

$$\text{So } \frac{\partial}{\partial \Sigma} \ln P(X|\underline{\mu}, \Sigma) = -\frac{N}{2} \Sigma^{-T} + \frac{1}{2} \Sigma^{-T} \left(\sum_{n=1}^N (\underline{x}_n - \underline{\mu}) (\underline{x}_n - \underline{\mu})^T \Sigma^{-1} \right) = 0$$

$$\Rightarrow N = \Sigma^{-T} \sum_{n=1}^N (\underline{x}_n - \underline{\mu}) (\underline{x}_n - \underline{\mu})^T \Rightarrow \Sigma_{ML} = \frac{1}{N} \sum_{n=1}^N (\underline{x}_n - \underline{\mu}_{ML}) (\underline{x}_n - \underline{\mu}_{ML})^T$$

But, when we check the expectation of our MLE, we have

$$\mathbb{E}(\underline{\mu}_{ML}) = \frac{1}{N} \sum_{n=1}^N \mathbb{E} \underline{x}_n = \frac{1}{N} \cdot N \cdot \underline{\mu} = \underline{\mu}$$

$$\begin{aligned} \mathbb{E}(\Sigma_{ML}) &= \frac{1}{N} \sum_{n=1}^N \mathbb{E} \left\{ (\underline{x}_n - \frac{1}{N} \sum_a \underline{x}_a) (\underline{x}_n - \frac{1}{N} \sum_b \underline{x}_b)^T \right\} \\ &= \frac{1}{N} \sum_{n=1}^N \mathbb{E} \left\{ \underline{x}_n \underline{x}_n^T - \frac{2}{N} \underline{x}_n \sum_a \underline{x}_a^T + \frac{1}{N^2} \sum_a \sum_b \underline{x}_a \underline{x}_b^T \right\} \end{aligned}$$

$$= \frac{1}{N} \left\{ \underline{\mu} \underline{\mu}^T + \Sigma - 2 \left(\underline{\mu} \underline{\mu}^T + \frac{1}{N} \Sigma \right) + \underline{\mu} \underline{\mu}^T + \frac{1}{N} \Sigma \right\} = \frac{N-1}{N} \Sigma$$

Biased Estimation.

We can construct an unbiased estimator

$$\hat{\Sigma} = \frac{1}{N-1} \sum_{n=1}^N (\underline{x}_n - \underline{\mu}_{ML}) (\underline{x}_n - \underline{\mu}_{ML})^T$$

(5) Bayesian Inference.

Like Multinomial Distribution, we also want to find a conjugate prior for μ .

① Let's assume σ^2 's known. The likelihood function $p(\underline{X}|\mu)$ is

$$p(\underline{X}|\mu) = \prod_{n=1}^N \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{1}{2\sigma^2} (X_n - \mu)^2\right\}$$

If we treat $p(\underline{X}|\mu)$ as a function of μ , it seems like it is a Gaussian.

So let's choose the prior to be another Gaussian

$$P(\mu) \sim \mathcal{N}(\mu|\mu_0, \sigma_0^2).$$

Then

$$P(\mu|\underline{X}) \propto p(\underline{X}|\mu) \cdot P(\mu)$$

The exponent part for $p(\underline{X}|\mu) \cdot P(\mu)$ is

$$\begin{aligned} & -\frac{1}{2\sigma^2} \cdot \sum_{n=1}^N (X_n - \mu)^2 - \frac{1}{2\sigma_0^2} (\mu - \mu_0)^2 \\ = & -\underbrace{\left(\frac{N}{2\sigma^2} + \frac{1}{2\sigma_0^2}\right)}_{\frac{1}{2\sigma_N^2}} \mu^2 + \underbrace{\left(\frac{1}{\sigma^2} \sum_{n=1}^N X_n + \frac{1}{\sigma_0^2} \mu_0\right)}_{\frac{\mu_N}{\sigma_N^2}} \mu + \text{const} \end{aligned}$$

$$\boxed{\mu_{ML} = \frac{1}{N} \sum_{n=1}^N X_n}$$

$$\hookrightarrow \frac{1}{\sigma_N^2} = \frac{N}{\sigma^2} + \frac{1}{\sigma_0^2}$$

$$\begin{aligned} \hookrightarrow \mu_N &= \left(\frac{N}{\sigma^2} + \frac{1}{\sigma_0^2}\right)^{-1} \left(\frac{\mu_0}{\sigma_0^2} + \frac{1}{\sigma^2} \sum_{n=1}^N X_n\right) \\ &= \frac{\sigma^2}{N\sigma_0^2 + \sigma^2} \mu_0 + \frac{N\sigma_0^2}{N\sigma_0^2 + \sigma^2} \mu_{ML} \end{aligned}$$

So, we find the posterior's μ_N & σ_N^2

② Let's assume the mean is known. we want to know the posterior for precision $\lambda = \frac{1}{\sigma^2}$

$$P(X|\lambda) = \prod_{n=1}^N \frac{1}{\sqrt{2\pi}} \cdot \sqrt{\lambda} \cdot \exp\left\{-\frac{\lambda}{2} (X_n - \mu)^2\right\}$$

$$\propto \lambda^{N/2} \cdot \exp\left\{-\frac{\lambda}{2} \sum_{n=1}^N (X_n - \mu)^2\right\}$$

The corresponding prior is Gamma distribution (not gamma function)

$$\text{Gamma}(\lambda|a, b) = \frac{1}{\Gamma(a)} b^a \lambda^{a-1} \exp(-b\lambda)$$

$$\mathbb{E}(\lambda) = \frac{a}{b} \quad \text{Var}(\lambda) = \frac{a}{b^2}$$

So, the posterior will be (Given the prior to be $\text{Gamma}(\lambda|a_0, b_0)$)

$$P(\lambda|X) \propto \lambda^{a_0-1} \lambda^{N/2} \exp\left\{-b_0\lambda - \frac{\lambda}{2} \sum_n (X_n - \mu)^2\right\}$$

$$\propto \text{Gamma}(\lambda|a_N, b_N)$$

$$a_N = a_0 + \frac{N}{2}$$

$$b_N = b_0 + \frac{1}{2} \sum (X_n - \mu)^2 = b_0 + \frac{N}{2} \sigma_{ML}^2$$

similar