

PROBABILITY DISTRIBUTIONS

Chia-Ping Chen

Professor

National Sun Yat-sen University

Department of Computer Science and Engineering

Machine Learning

We introduce the distributions of a few random variables.
Then we learn how to estimate such distributions from data.

- Binary Variables
- K -ary Variables
- Gaussian Variables
- Exponential Family
- Non-parametric Methods for Density Estimation

Binary Variables

Definition. Let X be a random variable. X is a **binary random variable** if it has exactly 2 possible values. In particular, X is a **Bernoulli random variable** if the possible values are 0 and 1. The distribution of a Bernoulli random variable is a **binary distribution function** given by

$$p(x) = \text{Bern}(x|\mu) = \mu^x(1 - \mu)^{1-x}, \quad x = 0, 1$$

It is clear that

$$p(1) = \mu, \quad p(0) = 1 - \mu$$

and

$$\begin{aligned}\mathbb{E}(X) &= \mu \\ \text{var}[X] &= \mu(1 - \mu)\end{aligned}$$

MAXIMUM LIKELIHOOD ESTIMATION

Let X be a Bernoulli random variable with $p(x) = \text{Bern}(x|\mu)$. What is the **maximum-likelihood estimate** of μ based on $\mathcal{D} = \{x_1, \dots, x_N\}$?

The likelihood of \mathcal{D} as a function of μ is

$$p(\mathcal{D}|\mu) = \prod_{n=1}^N \text{Bern}(x_n|\mu) = \prod_{n=1}^N \mu^{x_n} (1 - \mu)^{1-x_n}$$
$$\Rightarrow \log p(\mathcal{D}|\mu) = \sum_{n=1}^N [x_n \log \mu + (1 - x_n) \log(1 - \mu)]$$

Setting the derivative with respect to μ to 0, we get

$$\mu_{\text{ML}} = \frac{1}{N} \sum_{n=1}^N x_n$$

In **Bayesian learning**, we treat μ as a random variable.

- 1 We assume a prior distribution $p(\mu)$ of μ .
- 2 We update the distribution to $p(\mu|\mathcal{D})$ with data \mathcal{D} .
- 3 We obtain approximate $p(x)$ via $p(\mu|\mathcal{D})$ (point estimates or integration).

Benefit: Consider the density estimation of a Bernoulli random variable X based on $\mathcal{D} = \{x_1, \dots, x_N\}$. When N is not large enough, the maximum likelihood estimation of μ is prone to over-fitting. The issue can be alleviated by Bayesian learning.

Definition. Consider Bayesian learning of a Bernoulli random variable X with $p(x) = \text{Bern}(x|\mu)$ based on \mathcal{D} . A **conjugate prior** $p(\mu)$ of μ makes the posterior $p(\mu|\mathcal{D})$ have the same functional form as $p(\mu)$.

According to the Bayes' rule

$$p(\mu|\mathcal{D}) \propto p(\mu)p(\mathcal{D}|\mu)$$

Hence, the problem we face here is to make $p(\mu)$ and $p(\mu|\mathcal{D})$ belong to the same family, given data likelihood $p(\mathcal{D}|\mu)$.

DECIDING PRIOR FROM LIKELIHOOD FUNCTION

We can re-write the likelihood function as

$$\begin{aligned} p(\mathcal{D}|\mu) &= \prod_{n=1}^N \mu^{x_n} (1 - \mu)^{1-x_n} \\ &= \mu^{(\sum_{n=1}^N x_n)} (1 - \mu)^{(\sum_{n=1}^N (1-x_n))} \\ &= \mu^m (1 - \mu)^{N-m} \end{aligned}$$

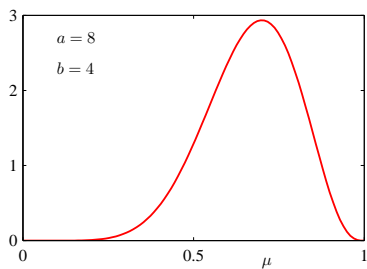
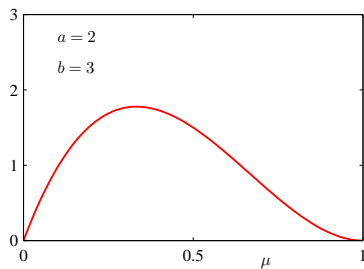
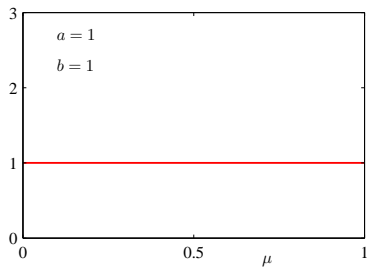
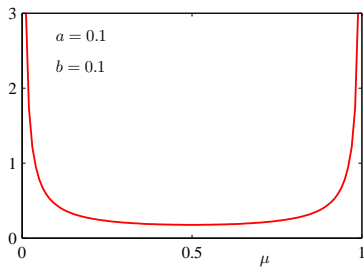
where $m = \sum_{n=1}^N x_n$ is the number of 1s in \mathcal{D} .

As a function of μ , the likelihood $p(\mathcal{D}|\mu)$ is proportional to powers of μ and $(1 - \mu)$. Thus, if we let the prior $p(\mu)$ be proportional to powers of μ and $(1 - \mu)$, then the posterior $p(\mu|\mathcal{D}) \propto p(\mu)p(\mathcal{D}|\mu)$ will also be proportional to powers of μ and $(1 - \mu)$.

Definition. A beta distribution is

$$p(\mu) = \text{Beta}(\mu|a, b) = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \mu^{a-1} (1-\mu)^{b-1}, \quad 0 \leq \mu \leq 1$$

The parameters in the distribution of a parameter, here a and b of $p(\mu)$, are **hyperparameters**.

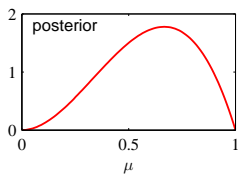
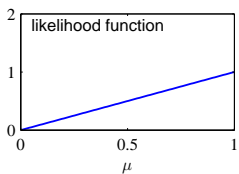
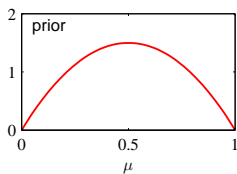


BETA PRIOR AND POSTERIOR

Let X be a Bernoulli random variable with distribution $p(x) = \text{Bern}(x|\mu)$. Then a beta distribution $p(\mu) = \text{Beta}(\mu|a, b)$ is a conjugate prior for μ .

Let m be the number of 1s and $l = N - m$ be the number of 0s in the data set $\mathcal{D} = \{x_1, \dots, x_N\}$. Then

$$\begin{aligned} p(\mu|\mathcal{D}) &\propto p(\mu)p(\mathcal{D}|\mu) \\ &= \text{Beta}(\mu|a, b)p(\mathcal{D}|\mu) \\ &\propto \mu^{a-1}(1-\mu)^{b-1}\mu^m(1-\mu)^l \\ &\propto \mu^{a+m-1}(1-\mu)^{b+l-1} \\ \Rightarrow p(\mu|\mathcal{D}) &= \frac{\Gamma(a+m+b+l)}{\Gamma(a+m)\Gamma(b+l)}\mu^{a+m-1}(1-\mu)^{b+l-1} \\ &= \text{Beta}(\mu|a+m, b+l) \end{aligned}$$



One step of sequential Bayesian inference. $\mu \sim \text{Beta}(\mu|a, b)$ from $a = 2, b = 2$ to $a = 3, b = 2$ with a single observation of $x = 1$.

INTERPRETATION OF HYPERPARAMETERS

Let X be a Bernoulli random variable with $p(x) = \text{Bern}(x|\mu)$. We see

$$p(\mu) = \text{Beta}(\mu|a, b)$$

↓

↓ Bayesian learning

↓

$$p(\mu|\mathcal{D}) = \text{Beta}(\mu|a + m, b + l)$$

The hyperparameters a and b in the prior distribution $p(\mu)$ can be interpreted as the **effective numbers of observations** for $x = 1$ and $x = 0$ prior to any observation.

Let X be a Bernoulli random variable with $p(x) = \text{Bern}(x|\mu)$. When μ is treated as a random variable, the distribution of X is the integration over the distribution of μ (the sum rule).

$$\begin{aligned}P(X = 1|\mathcal{D}) &= \int_0^1 P(X = 1, \mu|\mathcal{D})d\mu \\&= \int_0^1 P(X = 1|\mu)p(\mu|\mathcal{D})d\mu \\&= \int_0^1 \mu p(\mu|\mathcal{D})d\mu \\&= \mathbb{E}[\mu|\mathcal{D}] \\&= \frac{a + m}{a + m + b + l}\end{aligned}$$

On average, data observation reduces parameter uncertainty.

By the total variance theorem of probability theory

$$\text{var}_{\theta}(\boldsymbol{\theta}) = \mathbb{E}_{\mathcal{D}}[\text{var}_{\theta}(\boldsymbol{\theta}|\mathcal{D})] + \text{var}_{\mathcal{D}}(\mathbb{E}_{\theta}[\boldsymbol{\theta}|\mathcal{D}])$$

It follows from $\text{var}_{\mathcal{D}}(\mathbb{E}_{\theta}[\boldsymbol{\theta}|\mathcal{D}]) \geq 0$ that

$$\text{var}_{\theta}(\boldsymbol{\theta}) \geq \mathbb{E}_{\mathcal{D}}[\text{var}_{\theta}(\boldsymbol{\theta}|\mathcal{D})]$$

K-ary Variables

Definition. Let X be a discrete random variable. X is a **K -ary random variable** if it has exactly K possible values.

A value (a.k.a. **state**) of X can be represented by a vector of size K , i.e. $\mathbf{x} = (x_1, \dots, x_K)^T$, called 1-of- K (a.k.a. 1-hot) representation: one component is 1 for state identity, and the other components are 0. The set of possible values are $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(K)}\}$ with

$$x_j^{(k)} = \delta_{kj}$$

For example

$$\mathbf{x}^{(3)} = (0, 0, 1, 0, 0, 0)^T$$

We also denote a K -ary random variable by \mathbf{x} since the values are represented by vectors.

Definition. Let \mathbf{x} be a K -ary random variable. The distribution of \mathbf{x} is a **K -ary distribution**

$$P(\mathbf{x} = \mathbf{x}^{(j)}) = p(\mathbf{x}^{(j)}|\boldsymbol{\mu}) = \mu_j, \quad j = 1, \dots, K$$

where $\boldsymbol{\mu} = (\mu_1, \dots, \mu_K)^T$ is a vector of parameters.

The parameters $\boldsymbol{\mu}$ must satisfy

$$\mu_k \geq 0, \quad \sum_{k=1}^K \mu_k = 1$$

The expectation of \mathbf{x} is

$$\mathbb{E}[\mathbf{x}|\boldsymbol{\mu}] = \sum_{k=1}^K \mathbf{x}^{(k)} p(\mathbf{x}^{(k)}|\boldsymbol{\mu}) = (\mu_1, \dots, \mu_K)^T = \boldsymbol{\mu}$$

Let \mathbf{x} be a K -ary random variable with K -ary distribution $p(\mathbf{x}^{(j)}|\boldsymbol{\mu})$, and $\mathcal{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ be observations of \mathbf{x} . Using 1-of- K representation, the likelihood of a data point \mathbf{x}_n can be written as

$$p(\mathbf{x}_n|\boldsymbol{\mu}) = \prod_{k=1}^K \mu_k^{x_{nk}}$$

So the likelihood of \mathcal{D} is

$$\begin{aligned} p(\mathcal{D}|\boldsymbol{\mu}) &= \prod_{n=1}^N p(\mathbf{x}_n|\boldsymbol{\mu}) = \prod_{n=1}^N \prod_{k=1}^K \mu_k^{x_{nk}} = \prod_{k=1}^K \mu_k^{(\sum_{n=1}^N x_{nk})} \\ &= \prod_{k=1}^K \mu_k^{m_k} \end{aligned}$$

where m_k is the number of points in \mathcal{D} with $\mathbf{x}_n = \mathbf{x}^{(k)}$ or $x_{nk} = 1$.

MAXIMUM LIKELIHOOD ESTIMATION

The log data likelihood is

$$\log p(\mathcal{D}|\boldsymbol{\mu}) = \sum_{k=1}^K m_k \log \mu_k$$

Here the parameters μ_1, \dots, μ_K are not independent, so we need to maximize the **Lagrangian**

$$L(\boldsymbol{\mu}, \lambda) = \sum_{k=1}^K m_k \log \mu_k + \lambda \left(\sum_{k=1}^K \mu_k - 1 \right)$$

Setting the derivative of $L(\boldsymbol{\mu}, \lambda)$ with respect to $\boldsymbol{\mu}$ to 0, we get

$$\mu_k^{\text{ML}} = \frac{m_k}{N}$$

In Bayesian learning, we start with a distribution over the parameters and update the distribution with data. Again, we use conjugate prior so the posterior is an update of the hyperparameters with data.

This dependency of the data likelihood function on $\boldsymbol{\mu}$ decides the conjugate prior of $\boldsymbol{\mu}$. In the case of a K -ary random variable, the data likelihood function depends on powers of μ_k

$$p(\mathcal{D}|\boldsymbol{\mu}) = \prod_{k=1}^K \mu_k^{m_k}$$

Definition. A Dirichlet distribution is

$$\text{Dir}(\boldsymbol{\mu}|\boldsymbol{\alpha}) = \frac{\Gamma(\alpha_0)}{\Gamma(\alpha_1)\dots\Gamma(\alpha_K)} \prod_{k=1}^K \mu_k^{\alpha_k-1}$$

where

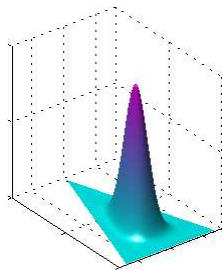
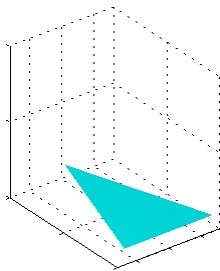
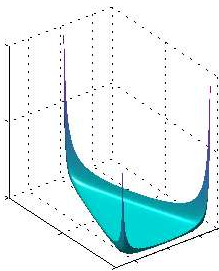
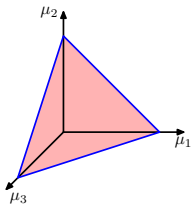
$$0 \leq \mu_k \leq 1, \quad \sum_{k=1}^K \mu_k = 1$$

Here $\boldsymbol{\alpha} = (\alpha_1, \dots, \alpha_K)^T$ are hyperparameters, and

$$\alpha_0 = \sum_{k=1}^K \alpha_k$$

Note $\text{Dir}(\boldsymbol{\mu}|\boldsymbol{\alpha})$ depends on powers of μ_k .

DIRICHLET: SUPPORT AND DENSITY



Dirichlet distributions for $\alpha_k = 0.1, 1, 10$, respectively

Let \mathbf{x} be a K -ary random variable with $p(\mathbf{x}^{(j)}|\boldsymbol{\mu})$. Then a Dirichlet distribution $\text{Dir}(\boldsymbol{\mu}|\boldsymbol{\alpha})$ is a conjugate prior for $\boldsymbol{\mu}$.

Based on \mathcal{D} , the posterior distribution of $\boldsymbol{\mu}$ is

$$p(\boldsymbol{\mu}|\mathcal{D}) \propto p(\boldsymbol{\mu})p(\mathcal{D}|\boldsymbol{\mu}) = \text{Dir}(\boldsymbol{\mu}|\boldsymbol{\alpha})p(\mathcal{D}|\boldsymbol{\mu}) \propto \prod_{k=1}^K \mu_k^{\alpha_k + m_k - 1}$$

By normalization

$$p(\boldsymbol{\mu}|\mathcal{D}) = \frac{\Gamma(\alpha_0 + N)}{\Gamma(\alpha_1 + m_1) \dots \Gamma(\alpha_K + m_K)} \prod_{k=1}^K \mu_k^{\alpha_k + m_k - 1} = \text{Dir}(\boldsymbol{\mu}|\boldsymbol{\alpha}')$$

where $\alpha'_k = \alpha_k + m_k$.

Gaussian Variables

Definition. A **Gaussian distribution** or **Gaussian PDF** is

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

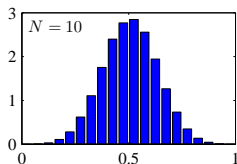
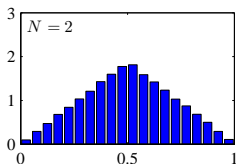
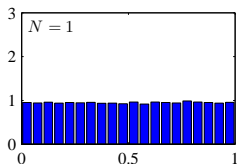
It is governed by parameters μ and σ^2 .

A random variable with a Gaussian distribution is a Gaussian random variable. Let X be a Gaussian random variable with distribution $p(x) = \mathcal{N}(x|\mu, \sigma^2)$. Then it can be shown that

$$\mu = \mathbb{E}[X], \quad \sigma^2 = \text{var}(X)$$

CENTRAL LIMIT THEOREM

The mean (or the sum) of a set of **i.i.d.** random variables has a distribution that becomes increasingly Gaussian as the number of terms in the sum increases.



Histogram plots of the mean of N uniform random variables in $[0, 1]$ for $N = 1, 2, 10$ respectively.

Definition. A multi-variate Gaussian distribution is

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

where D is the dimension of \mathbf{x} . It is governed by parameter vector $\boldsymbol{\mu}$ and parameter matrix $\boldsymbol{\Sigma}$.

A random vector with a multi-variate Gaussian distribution is a Gaussian random vector. Let \mathbf{x} be a Gaussian random vector with distribution $p(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma})$. Then it can be shown that

$$\int p(\mathbf{x})d\mathbf{x} = 1, \mathbb{E}[\mathbf{x}] = \boldsymbol{\mu}, \mathbb{E}[\mathbf{x}\mathbf{x}^T] = \boldsymbol{\mu}\boldsymbol{\mu}^T + \boldsymbol{\Sigma}, \text{var}(\mathbf{x}) = \boldsymbol{\Sigma}$$

MAHALANOBIS DISTANCE

Definition. Let \mathbf{x} be a Gaussian random vector of dimension D with $p(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma})$. The dependence of the PDF on \mathbf{x} is through quadratic form

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

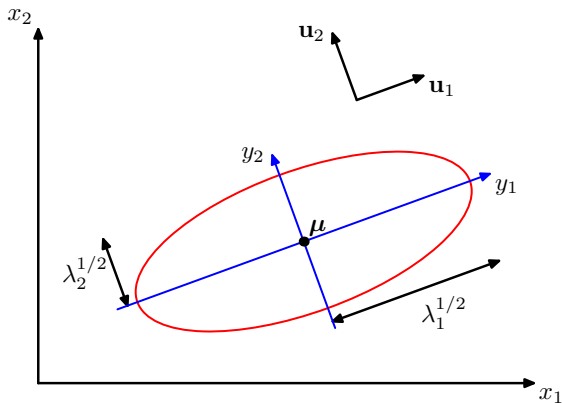
Δ is the **Mahalanobis distance** from $\boldsymbol{\mu}$ to \mathbf{x} .

Let $\mathbf{u}_1, \dots, \mathbf{u}_D$ be orthonormal eigenvectors of $\boldsymbol{\Sigma}$ with eigenvalues $\lambda_1 \leq \dots \leq \lambda_D$. Then $\boldsymbol{\Sigma} = \sum_{i=1}^D \lambda_i \mathbf{u}_i \mathbf{u}_i^T$, $\boldsymbol{\Sigma}^{-1} = \sum_{i=1}^D \lambda_i^{-1} \mathbf{u}_i \mathbf{u}_i^T$, and

$$\Delta^2 = (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) = \sum_{i=1}^D \frac{y_i^2}{\lambda_i}$$

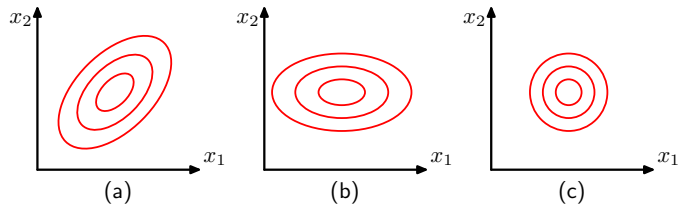
where $y_i = \mathbf{u}_i^T (\mathbf{x} - \boldsymbol{\mu})$ and $(\mathbf{x} - \boldsymbol{\mu}) = \sum_{i=1}^D y_i \mathbf{u}_i$.

CONTOUR OF A GAUSSIAN PDF



The contour of a 2-D Gaussian PDF, on which $\Delta^2 = 1$
so the density is $e^{-1/2}$ of the value at $\mathbf{x} = \boldsymbol{\mu}$.

EXAMPLES



Examples 2-D Gaussian PDF with general, diagonal, and isotropic covariance matrix, respectively.

Let \mathbf{x} be a Gaussian random vector with $p(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma})$.
Then the expectation of \mathbf{x} is $\boldsymbol{\mu}$.

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

where $\mathbf{z} = \mathbf{x} - \boldsymbol{\mu}$.

SECOND MOMENT

Let \mathbf{x} be a Gaussian random vector with $p(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma})$.
Then the expectation of $\mathbf{x}\mathbf{x}^T$ is $\boldsymbol{\mu}\boldsymbol{\mu}^T + \boldsymbol{\Sigma}$.

$$\begin{aligned}\mathbb{E}[\mathbf{x}\mathbf{x}^T] &= \frac{1}{(2\pi)^{\frac{D}{2}} |\boldsymbol{\Sigma}|^{\frac{1}{2}}} \int \exp\left\{-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right\} \mathbf{x}\mathbf{x}^T d\mathbf{x} \\ &= \frac{1}{(2\pi)^{\frac{D}{2}}} \frac{1}{|\boldsymbol{\Sigma}|^{\frac{1}{2}}} \int \exp\left\{-\frac{1}{2}\mathbf{z}^T \boldsymbol{\Sigma}^{-1}\mathbf{z}\right\} (\mathbf{z} + \boldsymbol{\mu})(\mathbf{z} + \boldsymbol{\mu})^T d\mathbf{z} \\ &= \frac{1}{(2\pi)^{\frac{D}{2}}} \frac{1}{|\boldsymbol{\Sigma}|^{\frac{1}{2}}} \int \exp\left\{-\frac{1}{2}\mathbf{z}^T \boldsymbol{\Sigma}^{-1}\mathbf{z}\right\} (\mathbf{z}\mathbf{z}^T + \boldsymbol{\mu}\boldsymbol{\mu}^T) d\mathbf{z} \\ &= \frac{1}{(2\pi)^{\frac{D}{2}}} \frac{1}{|\boldsymbol{\Sigma}|^{\frac{1}{2}}} \int \exp\left\{-\frac{1}{2}\mathbf{z}^T \boldsymbol{\Sigma}^{-1}\mathbf{z}\right\} \mathbf{z}\mathbf{z}^T d\mathbf{z} + \boldsymbol{\mu}\boldsymbol{\mu}^T \\ &= \sum_{i=1}^D \mathbf{u}_i \mathbf{u}_i^T \lambda_i + \boldsymbol{\mu}\boldsymbol{\mu}^T = \boldsymbol{\Sigma} + \boldsymbol{\mu}\boldsymbol{\mu}^T\end{aligned}$$

Recall that $\mathbf{z} = \mathbf{x} - \boldsymbol{\mu} = \sum y_i \mathbf{u}_i$.

$$\begin{aligned}
 & \frac{1}{(2\pi)^{\frac{D}{2}}} \frac{1}{|\boldsymbol{\Sigma}|^{\frac{1}{2}}} \int \exp \left\{ -\frac{1}{2} \mathbf{z}^T \boldsymbol{\Sigma}^{-1} \mathbf{z} \right\} \mathbf{z} \mathbf{z}^T d\mathbf{z} \\
 &= \sum_{i=1}^D \sum_{j=1}^D \int \frac{1}{(2\pi)^{\frac{D}{2}}} \frac{1}{|\boldsymbol{\Sigma}|^{\frac{1}{2}}} \exp \left\{ -\frac{1}{2} \sum_{k=1}^D \frac{y_k^2}{\lambda_k} \right\} y_i y_j \mathbf{u}_i \mathbf{u}_j^T d\mathbf{y} \\
 &= \sum_{i=1}^D \sum_{j=1}^D \mathbf{u}_i \mathbf{u}_j^T \int \frac{1}{(2\pi)^{\frac{D}{2}}} \frac{1}{|\boldsymbol{\Sigma}|^{\frac{1}{2}}} \exp \left\{ -\frac{1}{2} \sum_{k=1}^D \frac{y_k^2}{\lambda_k} \right\} y_i y_j d\mathbf{y} \\
 &= \sum_{i=1}^D \lambda_i \mathbf{u}_i \mathbf{u}_i^T = \boldsymbol{\Sigma}
 \end{aligned}$$

where

$$\int \frac{1}{(2\pi)^{\frac{D}{2}}} \frac{1}{|\boldsymbol{\Sigma}|^{\frac{1}{2}}} \exp \left\{ -\frac{1}{2} \sum_{k=1}^D \frac{y_k^2}{\lambda_k} \right\} y_i y_j d\mathbf{y} = \begin{cases} 0, & j \neq i \\ \lambda_i, & j = i \end{cases}$$

COVARIANCE MATRIX

Let \mathbf{x} be a Gaussian random vector with $p(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma})$.
Then the covariance matrix of \mathbf{x} is $\boldsymbol{\Sigma}$.

$$\begin{aligned}\text{cov}[\mathbf{x}] &= \mathbb{E} \left[(\mathbf{x} - \mathbb{E}[\mathbf{x}])(\mathbf{x} - \mathbb{E}[\mathbf{x}])^T \right] \\ &= \mathbb{E} \left[\mathbf{x}\mathbf{x}^T \right] - \mathbb{E}[\mathbf{x}] \mathbb{E}[\mathbf{x}^T] \\ &= \boldsymbol{\Sigma}\end{aligned}$$

FROM JOINT GAUSSIAN TO CONDITIONAL GAUSSIAN

Let random vectors \mathbf{x}_a and \mathbf{x}_b be joint Gaussian. Then \mathbf{x}_a is conditional Gaussian given \mathbf{x}_b .

Let the PDF of $\mathbf{x} = (\mathbf{x}_a, \mathbf{x}_b)^T$ be $p(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma})$. Partition \mathbf{x} , $\boldsymbol{\mu}$, $\boldsymbol{\Sigma}$ and the precision matrix $\boldsymbol{\Lambda} = \boldsymbol{\Sigma}^{-1}$ as follows

$$\mathbf{x} = \begin{pmatrix} \mathbf{x}_a \\ \mathbf{x}_b \end{pmatrix}, \boldsymbol{\mu} = \begin{pmatrix} \boldsymbol{\mu}_a \\ \boldsymbol{\mu}_b \end{pmatrix}, \boldsymbol{\Sigma} = \begin{pmatrix} \boldsymbol{\Sigma}_{aa} & \boldsymbol{\Sigma}_{ab} \\ \boldsymbol{\Sigma}_{ba} & \boldsymbol{\Sigma}_{bb} \end{pmatrix}, \boldsymbol{\Lambda} = \begin{pmatrix} \boldsymbol{\Lambda}_{aa} & \boldsymbol{\Lambda}_{ab} \\ \boldsymbol{\Lambda}_{ba} & \boldsymbol{\Lambda}_{bb} \end{pmatrix}$$

It can be shown that the conditional distribution of \mathbf{x}_a given \mathbf{x}_b is $p(\mathbf{x}_a|\mathbf{x}_b) = \mathcal{N}(\mathbf{x}_a|\boldsymbol{\mu}_{a|b}, \boldsymbol{\Sigma}_{a|b})$ where

$$\boldsymbol{\mu}_{a|b} = \boldsymbol{\mu}_a - \boldsymbol{\Lambda}_{aa}^{-1} \boldsymbol{\Lambda}_{ab} (\mathbf{x}_b - \boldsymbol{\mu}_b) = \boldsymbol{\mu}_a + \boldsymbol{\Sigma}_{ab} \boldsymbol{\Sigma}_{bb}^{-1} (\mathbf{x}_b - \boldsymbol{\mu}_b)$$

$$\boldsymbol{\Sigma}_{a|b} = \boldsymbol{\Lambda}_{aa}^{-1} = \boldsymbol{\Sigma}_{aa} - \boldsymbol{\Sigma}_{ab} \boldsymbol{\Sigma}_{bb}^{-1} \boldsymbol{\Sigma}_{ba}$$

Note $\boldsymbol{\Lambda}_{a|b} = \boldsymbol{\Sigma}_{a|b}^{-1} = \boldsymbol{\Lambda}_{aa}$.

Let random vectors \mathbf{x}_a and \mathbf{x}_b be joint Gaussian. Then \mathbf{x}_a (and \mathbf{x}_b) is Gaussian.

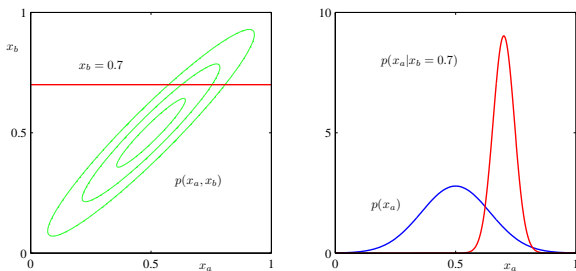
Let the PDF of $\mathbf{x} = (\mathbf{x}_a, \mathbf{x}_b)^T$ be $p(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma})$. Partition $\mathbf{x}, \boldsymbol{\mu}, \boldsymbol{\Sigma}$ and the precision matrix $\boldsymbol{\Lambda} = \boldsymbol{\Sigma}^{-1}$ as follows

$$\mathbf{x} = \begin{pmatrix} \mathbf{x}_a \\ \mathbf{x}_b \end{pmatrix}, \boldsymbol{\mu} = \begin{pmatrix} \boldsymbol{\mu}_a \\ \boldsymbol{\mu}_b \end{pmatrix}, \boldsymbol{\Sigma} = \begin{pmatrix} \boldsymbol{\Sigma}_{aa} & \boldsymbol{\Sigma}_{ab} \\ \boldsymbol{\Sigma}_{ba} & \boldsymbol{\Sigma}_{bb} \end{pmatrix}, \boldsymbol{\Lambda} = \begin{pmatrix} \boldsymbol{\Lambda}_{aa} & \boldsymbol{\Lambda}_{ab} \\ \boldsymbol{\Lambda}_{ba} & \boldsymbol{\Lambda}_{bb} \end{pmatrix}$$

It can be shown that the distribution of \mathbf{x}_a is $p(\mathbf{x}_a) = \mathcal{N}(\mathbf{x}_a|\boldsymbol{\mu}_{a|\emptyset}, \boldsymbol{\Sigma}_a)$ where

$$\begin{aligned} \boldsymbol{\mu}_{a|\emptyset} &= \boldsymbol{\mu}_a \\ \boldsymbol{\Sigma}_a &= \boldsymbol{\Sigma}_{aa} = (\boldsymbol{\Lambda}_{aa} - \boldsymbol{\Lambda}_{ab}\boldsymbol{\Lambda}_{bb}^{-1}\boldsymbol{\Lambda}_{ba})^{-1} \end{aligned}$$

CONDITIONAL GAUSSIAN



Left: Contours of joint Gaussian PDF.

Right: A marginal and a conditional Gaussian PDF.

Let \mathbf{y} be Gaussian and \mathbf{z} be conditional Gaussian given \mathbf{y} .
Then \mathbf{y} and \mathbf{z} are joint Gaussian.

Let the PDF of \mathbf{y} be $p(\mathbf{y}) = \mathcal{N}(\mathbf{y}|\boldsymbol{\mu}, \boldsymbol{\Lambda}^{-1})$ and the conditional PDF of \mathbf{z} given \mathbf{y} be $p(\mathbf{z}|\mathbf{y}) = \mathcal{N}(\mathbf{z}|\mathbf{A}\mathbf{y} + \mathbf{b}, \mathbf{L}^{-1})$. Then the joint PDF of $\mathbf{x} = (\mathbf{y}, \mathbf{z})^T$ is

$$\begin{aligned} p(\mathbf{x}) &= p(\mathbf{y})p(\mathbf{z}|\mathbf{y}) \\ &= \mathcal{N} \left(\begin{pmatrix} \mathbf{y} \\ \mathbf{z} \end{pmatrix} \middle| \begin{pmatrix} \boldsymbol{\mu} \\ \mathbf{A}\boldsymbol{\mu} + \mathbf{b} \end{pmatrix}, \begin{pmatrix} \boldsymbol{\Lambda} + \mathbf{A}^T \mathbf{L} \mathbf{A} & -\mathbf{A}^T \mathbf{L} \\ -\mathbf{L} \mathbf{A} & \mathbf{L} \end{pmatrix}^{-1} \right) \end{aligned}$$

Let \mathbf{y} be Gaussian and \mathbf{z} be conditional Gaussian given \mathbf{y} .
Since \mathbf{y} and \mathbf{z} are joint Gaussian, we have

- \mathbf{z} is Gaussian with PDF
- \mathbf{y} is conditional Gaussian given \mathbf{z}

The PDF of \mathbf{z} is

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

The conditional PDF of \mathbf{y} given \mathbf{z} is

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

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

Let \mathbf{x} be Gaussian with PDF $p(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma})$. In the **density estimation** of \mathbf{x} , we use a data set to estimate $p(\mathbf{x})$.

- Maximum likelihood estimate
- Bayesian learning

Let $\mathcal{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ be a data set of a Gaussian random vector \mathbf{x} . For data likelihood, we have

$$\begin{aligned} p(\mathbf{x}_n | \boldsymbol{\mu}, \boldsymbol{\Sigma}) &= \frac{1}{(2\pi)^{\frac{D}{2}} |\boldsymbol{\Sigma}|^{\frac{1}{2}}} \exp \left\{ -\frac{1}{2} (\mathbf{x}_n - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}_n - \boldsymbol{\mu}) \right\} \\ \Rightarrow \log p(\mathbf{x}_n | \boldsymbol{\mu}, \boldsymbol{\Sigma}) &= -\frac{D}{2} \log 2\pi - \frac{1}{2} \log |\boldsymbol{\Sigma}| - \frac{1}{2} (\mathbf{x}_n - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}_n - \boldsymbol{\mu}) \\ \Rightarrow \log p(\mathcal{D} | \boldsymbol{\mu}, \boldsymbol{\Sigma}) &= \sum_{n=1}^N \log p(\mathbf{x}_i | \boldsymbol{\mu}, \boldsymbol{\Sigma}) \\ &= -\frac{ND}{2} \log 2\pi - \frac{N}{2} \log |\boldsymbol{\Sigma}| \\ &\quad - \frac{1}{2} \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}_n - \boldsymbol{\mu}) \end{aligned}$$

Maximum likelihood estimate maximizes data likelihood.

- For $\boldsymbol{\mu}$

$$\begin{aligned}\nabla_{\boldsymbol{\mu}} \log p(\mathcal{D}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = 0 &\Rightarrow \sum_{n=1}^N \boldsymbol{\Sigma}^{-1}(\mathbf{x}_n - \boldsymbol{\mu}) = 0 \\ &\Rightarrow \boldsymbol{\mu}_{\text{ML}} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n\end{aligned}$$

- For $\boldsymbol{\Sigma}$, it can be shown that

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

Let x be a Gaussian random variable with PDF $p(x) = \mathcal{N}(x|\mu, \sigma^2)$, and $\mathcal{D} = \{x_1, \dots, x_N\}$ be a data set. We consider the following scenarios of Bayesian learning of Gaussian distribution.

- Given σ^2 and a conjugate prior of μ
- Given μ and a conjugate prior of precision $\lambda = 1/\sigma^2$
- Given a conjugate prior of μ and λ

CONJUGATE PRIOR OF MEAN

The likelihood of \mathcal{D} as a function of μ is

$$\begin{aligned} p(\mathcal{D}|\mu) &= \prod_{n=1}^N p(x_n|\mu) \\ &= \frac{1}{(2\pi\sigma^2)^{N/2}} \exp \left\{ -\frac{1}{2\sigma^2} \sum_{n=1}^N (x_n - \mu)^2 \right\} \end{aligned}$$

It is a log quadratic function of μ , so a conjugate prior of μ is log quadratic, i.e. Gaussian

$$p(\mu) = \mathcal{N}(\mu|\mu_0, \sigma_0^2)$$

According to the Bayes' rule

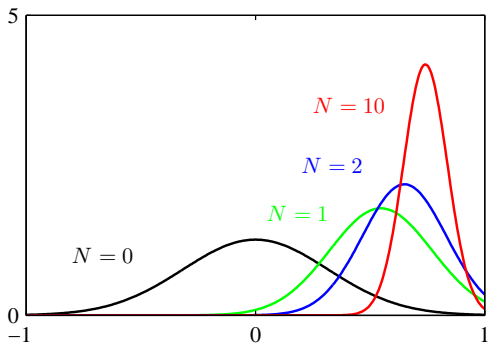
$$p(\mu|\mathcal{D}) \propto p(\mu)p(\mathcal{D}|\mu)$$

we have

$$p(\mu|\mathcal{D}) = \mathcal{N}(\mu|\mu_N, \sigma_N^2)$$

where

$$\mu_N = \frac{\sigma^2}{N\sigma_0^2 + \sigma^2}\mu_0 + \frac{N\sigma_0^2}{N\sigma_0^2 + \sigma^2}\mu_{\text{ML}}$$
$$\frac{1}{\sigma_N^2} = \frac{1}{\sigma_0^2} + \frac{N}{\sigma^2}$$



Plots of $p(\mu|\mathcal{D})$ assuming Gaussian $p(\mu)$ and $p(x|\mu) = \mathcal{N}(x|\mu, \sigma^2)$.
 The true $p(x)$ is Gaussian with mean 0.8 and variance 0.1.

The likelihood of \mathcal{D} as a function of λ is

$$p(\mathcal{D}|\lambda) = \prod_{n=1}^N p(x_n|\lambda) \\ \propto \lambda^{N/2} \exp \left\{ -\frac{\lambda}{2} \sum_{n=1}^N (x_n - \mu)^2 \right\}$$

So a conjugate prior of λ is **gamma distribution**

$$p(\lambda) = \text{Gam}(\lambda|a, b) = \frac{1}{\Gamma(a)} b^a \lambda^{a-1} e^{-b\lambda}, \quad \lambda \geq 0$$

BAYESIAN LEARNING OF PRECISION

According to the Bayes' rule

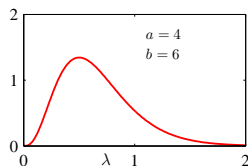
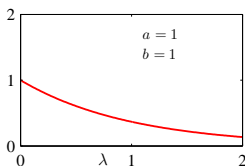
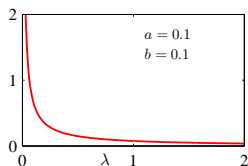
$$p(\lambda|\mathcal{D}) \propto p(\lambda)p(\mathcal{D}|\lambda)$$

we have

$$p(\lambda|\mathcal{D}) = \text{Gam}(\lambda|a_N, b_N)$$

where

$$a_N = a + \frac{N}{2}, \quad b_N = b + \frac{N}{2}\sigma_{\text{ML}}^2$$



STUDENT'S t -DISTRIBUTION

The marginal distribution of x is

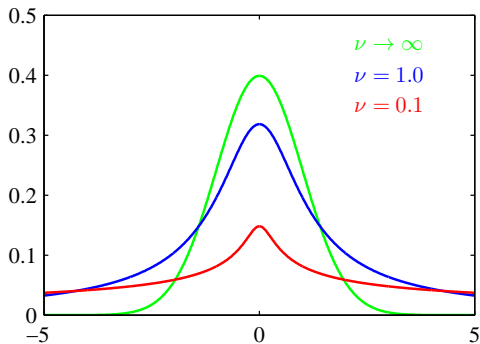
$$\begin{aligned} p(x|\mu, a, b) &= \int_0^\infty \mathcal{N}(x|\mu, \lambda^{-1}) \text{Gam}(\lambda|a, b) d\lambda \\ &= \frac{b^a}{\Gamma(a)} \left(\frac{1}{2\pi}\right)^{1/2} \left[b + \frac{(x - \mu)^2}{2} \right]^{-a-1/2} \Gamma(a + 1/2) \end{aligned}$$

This is a **Student's t -distribution**

$$\text{St}(x|\mu, \lambda, \nu) = \frac{\Gamma(\nu/2 + 1/2)}{\Gamma(\nu/2)} \left(\frac{\lambda}{\pi\nu}\right)^{1/2} \left[1 + \frac{\lambda(x - \mu)^2}{\nu} \right]^{-\nu/2-1/2}$$

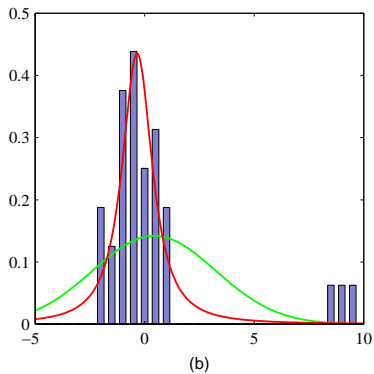
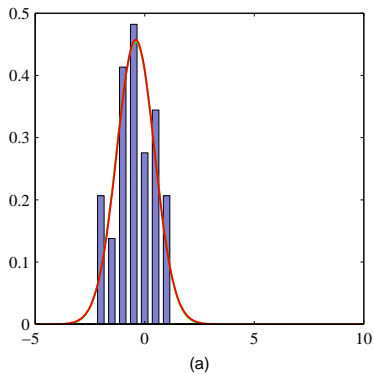
with parameters $\nu = 2a$ and $\lambda = a/b$.

LONG TAIL OF STUDENT'S t -DISTRIBUTION



Plots of Student's t -distribution $\text{St}(x|\mu, \lambda, \nu)$
with 3 ν 's, $\mu = 0$, and $\lambda = 1$.

ROBUSTNESS TO OUTLIERS



CONJUGATE PRIOR OF MEAN AND PRECISION

The likelihood of \mathcal{D} as a function of μ and λ is

$$\begin{aligned} p(\mathcal{D}|\mu, \lambda) &= \prod_{n=1}^N p(x_n|\mu, \lambda) \\ &= \prod_{n=1}^N \left(\frac{\lambda}{2\pi}\right)^{1/2} \exp\left\{-\frac{\lambda}{2}(x_n - \mu)^2\right\} \\ &\propto \left[\lambda^{1/2} \exp\left(-\frac{\lambda\mu^2}{2}\right)\right]^N \exp\left\{\lambda\mu \sum_{n=1}^N x_n - \frac{\lambda}{2} \sum_{n=1}^N x_n^2\right\} \end{aligned}$$

So a conjugate prior of μ and λ is

$$p(\mu, \lambda) \propto \left[\lambda^{1/2} \exp\left(-\frac{\lambda\mu^2}{2}\right)\right]^{\beta_0} \exp\{c\lambda\mu - d\lambda\}$$

NORMAL-GAMMA DISTRIBUTION

The conjugate prior can be further reduced

$$\begin{aligned} p(\mu, \lambda) &\propto \left[\lambda^{1/2} \exp\left(-\frac{\lambda\mu^2}{2}\right) \right]^{\beta_0} \exp\{c\lambda\mu - d\lambda\} \\ &= \exp\left\{c\lambda\mu - \frac{\beta}{2}\lambda\mu^2\right\} \lambda^{\beta_0/2} \exp\{-d\lambda\} \\ &= \exp\left\{-\frac{\beta_0\lambda}{2}(\mu - c/\beta_0)^2\right\} \lambda^{\beta_0/2} \exp\left\{-\left(d - \frac{c^2}{2\beta_0}\right)\lambda\right\} \\ &= \mathcal{N}(\mu|\alpha_0, (\beta_0\lambda)^{-1})\text{Gam}(\lambda|a_0, b_0) \end{aligned}$$

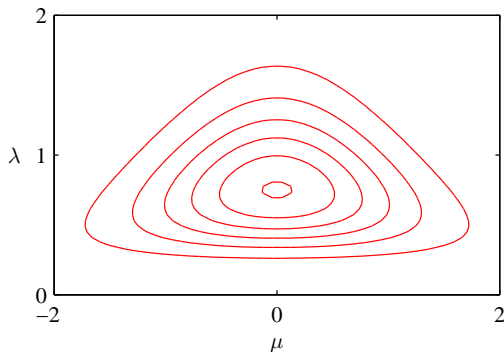
with

$$\alpha_0 = \frac{c}{\beta_0}, \quad a_0 = 1 + \frac{\beta_0}{2}, \quad b_0 = d - \frac{c^2}{2\beta_0}$$

This is a **normal-gamma distribution** defined by

$$\text{Nor-Gam}(\mu, \lambda|\alpha, \beta, a, b) = \mathcal{N}(\mu|\alpha, (\beta\lambda)^{-1})\text{Gam}(\lambda|a, b)$$

EXAMPLE



Plot of $\text{Nor-Gam}(\mu, \lambda | \alpha, \beta, a, b)$ with
 $\alpha = 0, \beta = 2, a = 5, b = 6.$

According to the Bayes' rule

$$p(\mu, \lambda | \mathcal{D}) \propto p(\mu, \lambda) p(\mathcal{D} | \mu, \lambda)$$

The posterior distribution of μ and λ is

$$p(\mu, \lambda | \mathcal{D}) = \text{Nor-Gam}(\mu, \lambda | \alpha_N, \beta_N, a_N, b_N)$$

with parameters

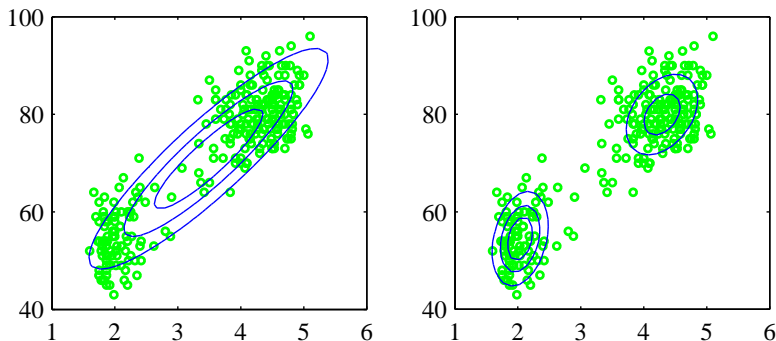
$$\beta_N = \beta_0 + N, \quad \alpha_N = \frac{c_N}{\beta_N}, \quad a_N = 1 + \frac{\beta_N}{2}, \quad b_N = d_N - \frac{c_N^2}{2\beta_N}$$

where

$$c_N = c + \sum_{n=1}^N x_n, \quad d_N = d + \frac{1}{2} \sum_{n=1}^N x_n^2$$

AN ISSUE WITH GAUSSIAN PDF

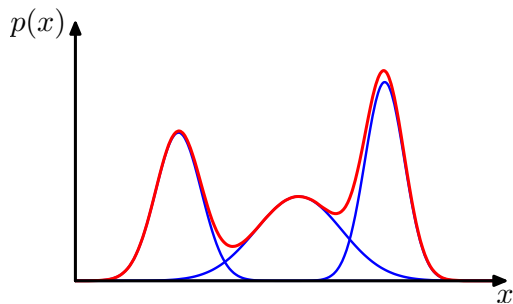
A Gaussian distribution is not good for data with multiple clusters.



Fitting the old faithful data set with a Gaussian (left) and a mixture of Gaussians (right).

SUPERPOSITION OF GAUSSIANS

With sufficient components, a linear combination of Gaussian PDFs can approximate any distribution.



Definition. A **mixture of Gaussians** is a superposition of Gaussians

$$p(\mathbf{x}) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

where $\mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$ is a **component** of the mixture, with a **mixing coefficient** π_k . The mixing coefficients of a mixture of Gaussians must satisfy

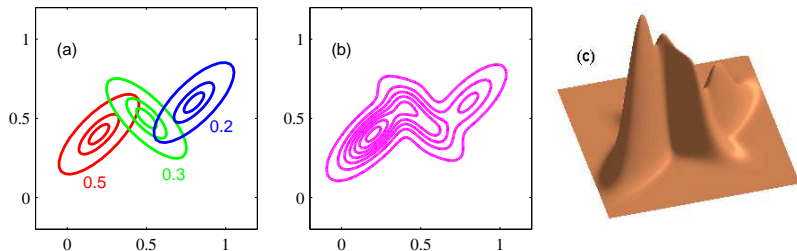
$$\pi_k \geq 0, \quad \sum_k \pi_k = 1$$

The parameters of a mixture of Gaussians are

$$\pi_k, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k \text{ for } k = 1, \dots, K$$

2-D EXAMPLE

A mixture of bi-variate Gaussians.



Left: Contours of the PDFs of 3 Gaussian components

Middle: Contours of the PDF of the mixture

Right: 3-D plot of the PDF of the mixture

Exponential Family*

Definition. Let \mathbf{x} be a random vector with a parametric distribution $p(\mathbf{x}|\boldsymbol{\eta})$. The distribution $p(\mathbf{x}|\boldsymbol{\eta})$ is in the **exponential family** if it has the following form

$$p(\mathbf{x}|\boldsymbol{\eta}) = h(\mathbf{x})g(\boldsymbol{\eta}) \exp \left\{ \boldsymbol{\eta}^T \mathbf{u}(\mathbf{x}) \right\}$$

- $\boldsymbol{\eta}$ are the natural parameters.
- $g(\boldsymbol{\eta})$ ensures normalization of $p(\mathbf{x}|\boldsymbol{\eta})$, i.e.

$$g(\boldsymbol{\eta}) \int h(\mathbf{x}) \exp \left\{ \boldsymbol{\eta}^T \mathbf{u}(\mathbf{x}) \right\} d\mathbf{x} = 1$$

A Bernoulli distribution is in the exponential family.

$$\begin{aligned}\text{Bern}(x|\mu) &= \mu^x(1 - \mu)^{1-x} \\ &= \exp\{x \log \mu + (1 - x) \log(1 - \mu)\} \\ &= (1 - \mu) \exp\left\{x \log\left(\frac{\mu}{1 - \mu}\right)\right\} \\ &= h(x)g(\eta) \exp\{\eta^T \mathbf{u}(x)\}\end{aligned}$$

with

$$\begin{aligned}\eta &= \log\left(\frac{\mu}{1 - \mu}\right) \\ \mathbf{u}(x) &= x \\ g(\eta) &= 1 - \mu = \frac{1}{1 + \exp(\eta)} \\ h(x) &= 1\end{aligned}$$

A K -ary distribution is in the exponential family.

$$\begin{aligned} p(\mathbf{x}|\boldsymbol{\mu}) &= \prod_{k=1}^K \mu_k^{x_k} = \exp \left\{ \sum_{k=1}^K x_k \log \mu_k \right\} \\ &= \exp \left\{ \sum_{k=1}^{K-1} x_k \log \mu_k + \left(1 - \sum_{k=1}^{K-1} x_k \right) \log \mu_K \right\} \\ &= \mu_K \exp \left\{ \sum_{k=1}^{K-1} x_k \log \left(\frac{\mu_k}{\mu_K} \right) \right\} \\ &= h(\mathbf{x})g(\boldsymbol{\eta}) \exp\{\boldsymbol{\eta}^T \mathbf{u}(\mathbf{x})\} \end{aligned}$$

with

$$\eta_k = \log \left(\frac{\mu_k}{\mu_K} \right), \quad g(\boldsymbol{\eta}) = \mu_K = 1 - \sum_{k=1}^{K-1} \mu_k, \quad h(\mathbf{x}) = 1, \quad \mathbf{u}(\mathbf{x}) = \mathbf{x}$$

A Gaussian PDF is in the exponential family.

$$\begin{aligned}
 \mathcal{N}(x|\mu, \sigma^2) &= \frac{1}{(2\pi\sigma^2)^{\frac{1}{2}}} \exp\left\{-\frac{1}{2\sigma^2}(x - \mu)^2\right\} \\
 &= \frac{1}{(2\pi\sigma^2)^{\frac{1}{2}}} \exp\left\{-\frac{1}{2\sigma^2}x^2 + \frac{\mu}{\sigma^2}x - \frac{1}{2\sigma^2}\mu^2\right\} \\
 &= h(x)g(\boldsymbol{\eta}) \exp\{\boldsymbol{\eta}^T \mathbf{u}(x)\}
 \end{aligned}$$

with

$$\boldsymbol{\eta} = \begin{pmatrix} \frac{\mu}{\sigma^2} \\ -\frac{1}{2\sigma^2} \end{pmatrix}, \quad \mathbf{u}(x) = \begin{pmatrix} x \\ x^2 \end{pmatrix}$$

$$h(x) = (2\pi)^{-\frac{1}{2}}, \quad g(\boldsymbol{\eta}) = (-2\eta_2)^{\frac{1}{2}} \exp\left(\frac{\eta_1^2}{4\eta_2}\right)$$

DENSITY ESTIMATION AND LIKELIHOOD

Let \mathbf{x} be a random vector with

$$p(\mathbf{x}|\boldsymbol{\eta}) = h(\mathbf{x})g(\boldsymbol{\eta}) \exp \left\{ \boldsymbol{\eta}^T \mathbf{u}(\mathbf{x}) \right\}$$

Consider the density estimation of \mathbf{x} with $\mathcal{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$.

$$p(\mathcal{D}|\boldsymbol{\eta}) = \prod_{n=1}^N p(\mathbf{x}_n|\boldsymbol{\eta})$$

$$= \prod_{n=1}^N \left(h(\mathbf{x}_n)g(\boldsymbol{\eta}) \exp \left\{ \boldsymbol{\eta}^T \mathbf{u}(\mathbf{x}_n) \right\} \right)$$

$$\Rightarrow \log p(\mathcal{D}|\boldsymbol{\eta}) = \sum_{n=1}^N \log h(\mathbf{x}_n) + N \log g(\boldsymbol{\eta}) + \boldsymbol{\eta}^T \sum_{n=1}^N \mathbf{u}(\mathbf{x}_n)$$

$$\Rightarrow \nabla \log p(\mathcal{D}|\boldsymbol{\eta}) = N \nabla \log g(\boldsymbol{\eta}) + \sum_{n=1}^N \mathbf{u}(\mathbf{x}_n)$$

Maximum likelihood estimate is a stationary point of $\log p(\mathcal{D}|\boldsymbol{\eta})$

$$\nabla \log p(\mathcal{D}|\boldsymbol{\eta}) \Big|_{\boldsymbol{\eta}=\boldsymbol{\eta}_{\text{ML}}} = 0$$

With $p(\mathbf{x}|\boldsymbol{\eta})$ in the exponential family, we have

$$-\nabla \log g(\boldsymbol{\eta}_{\text{ML}}) = \frac{1}{N} \sum_{n=1}^N \mathbf{u}(\mathbf{x}_n)$$

Note $\boldsymbol{\eta}_{\text{ML}}$ depends on the data set only through $\sum_{n=1}^N \mathbf{u}(\mathbf{x}_n)$. This is an example of **sufficient statistics**.

Let \mathbf{x} be a random vector with $p(\mathbf{x}|\boldsymbol{\eta})$ in the exponential family

$$p(\mathbf{x}|\boldsymbol{\eta}) = h(\mathbf{x})g(\boldsymbol{\eta}) \exp \left\{ \boldsymbol{\eta}^T \mathbf{u}(\mathbf{x}) \right\}$$

For a conjugate prior of $\boldsymbol{\eta}$, we match the dependency of the data likelihood function on $\boldsymbol{\eta}$, i.e. a power of $g(\boldsymbol{\eta})$ and an exponent linear in $\boldsymbol{\eta}$

$$p(\boldsymbol{\eta}|\boldsymbol{\chi}, \nu) = f(\boldsymbol{\chi}, \nu)g(\boldsymbol{\eta})^\nu \exp \left\{ \nu \boldsymbol{\eta}^T \boldsymbol{\chi} \right\}$$

According to the Bayes' rule

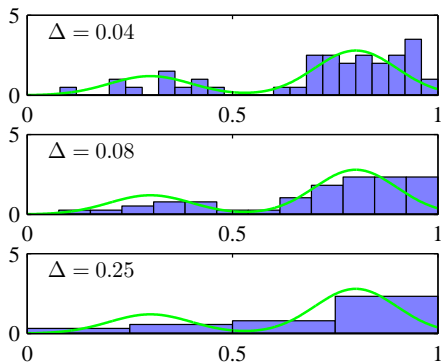
$$\begin{aligned} p(\boldsymbol{\eta}|\mathcal{D}, \boldsymbol{\chi}, \nu) &\propto p(\boldsymbol{\eta}|\boldsymbol{\chi}, \nu)p(\mathcal{D}|\boldsymbol{\eta}) \\ &= g(\boldsymbol{\eta})^{\nu+N} \exp \left\{ \boldsymbol{\eta}^T \left(\sum_{n=1}^N \mathbf{u}(\mathbf{x}_n) + \nu\boldsymbol{\chi} \right) \right\} \\ &= g(\boldsymbol{\eta})^{\nu'} \exp \left\{ \nu' \boldsymbol{\eta}^T \boldsymbol{\chi}' \right\} \\ &\propto p(\boldsymbol{\eta}|\boldsymbol{\chi}', \nu') \end{aligned}$$

where

$$\nu' = \nu + N, \quad \boldsymbol{\chi}' = \frac{1}{\nu'} \left(\nu\boldsymbol{\chi} + \sum_{n=1}^N \mathbf{u}(\mathbf{x}_n) \right)$$

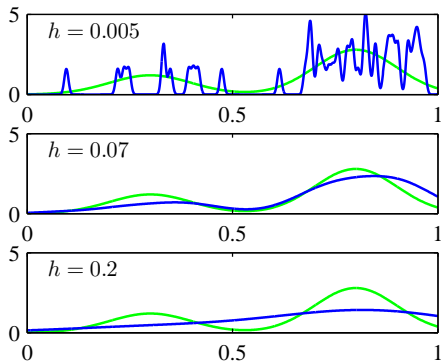
Nonparametric Methods

HISTOGRAM DENSITY ESTIMATION



3 cases of histogram density estimation with 50 data points generated from the distribution shown by the green curve.

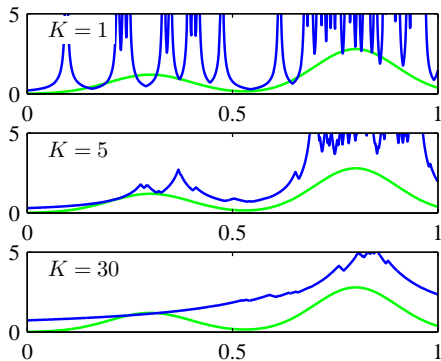
KERNEL DENSITY ESTIMATION



3 cases of kernel density estimation with the same data set

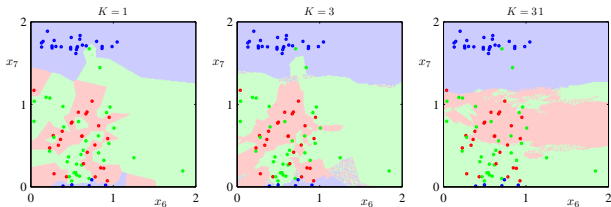
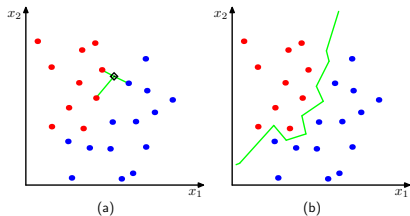
$$p(x) = \frac{1}{N} \sum_n k(x, x_n), \text{ where } k(x, x') = \frac{1}{(2\pi h^2)^{1/2}} \exp\left\{-\frac{(x - x')^2}{2h^2}\right\}$$

K -NEAREST-NEIGHBOR DENSITY ESTIMATION



3 cases of KNN density estimation with the same data set. Here $\hat{p}(x) = \frac{K}{NV(x)}$ where $V(x)$ is the volume of a sphere centered on x and containing K data points.

K -NEAREST-NEIGHBOR CLASSIFIERS



* **Joint, Marginal, and Conditional Gaussians**

For a Gaussian random vector \mathbf{x} , the log distribution

- is quadratic in \mathbf{x}
- the second-order term depends on precision/covariance
- the first-order term depends on precision and mean

$$\begin{aligned}\log \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Lambda}^{-1}) &= -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Lambda}(\mathbf{x} - \boldsymbol{\mu}) + \text{const} \\ &= -\frac{1}{2}\mathbf{x}^T \boldsymbol{\Lambda} \mathbf{x} + \mathbf{x}^T \boldsymbol{\Lambda} \boldsymbol{\mu} + \text{const}\end{aligned}$$

PARTITION BY TWO SUB-VECTORS

Let $\mathbf{x} \sim \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Lambda}^{-1})$ be a Gaussian random vector. Partition \mathbf{x} and $\boldsymbol{\mu}$ into sub-vectors

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

Partition the covariance and precision into sub-matrices

$$\boldsymbol{\Sigma} = \begin{pmatrix} \boldsymbol{\Sigma}_{aa} & \boldsymbol{\Sigma}_{ab} \\ \boldsymbol{\Sigma}_{ba} & \boldsymbol{\Sigma}_{bb} \end{pmatrix}, \quad \boldsymbol{\Lambda} = \begin{pmatrix} \boldsymbol{\Lambda}_{aa} & \boldsymbol{\Lambda}_{ab} \\ \boldsymbol{\Lambda}_{ba} & \boldsymbol{\Lambda}_{bb} \end{pmatrix} = \boldsymbol{\Sigma}^{-1}$$

CONDITIONAL GAUSSIAN PROPERTY

The conditional distribution of \mathbf{x}_a given \mathbf{x}_b is Gaussian.

The conditional distribution of \mathbf{x}_a given $\mathbf{x}_b = \mathbf{x}_b$ is Gaussian

$$\mathcal{N}(\mathbf{x}_a | \boldsymbol{\mu}_{a|b}, \boldsymbol{\Lambda}_{a|b}^{-1})$$

where the conditional mean is

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

and the conditional precision is

$$\boldsymbol{\Lambda}_{a|b} = \boldsymbol{\Lambda}_{aa}$$

The conditional distribution of \mathbf{x}_a given $\mathbf{x}_b = \mathbf{x}_b$ is

$$p(\mathbf{x}_a | \mathbf{x}_b) = \frac{p(\mathbf{x}_a, \mathbf{x}_b)}{p(\mathbf{x}_b)} \propto p(\mathbf{x}_a, \mathbf{x}_b)$$

The logarithm of $p(\mathbf{x}_a | \mathbf{x}_b)$, apart from a constant,

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

is **quadratic** in \mathbf{x}_a . Hence \mathbf{x}_a given $\mathbf{x}_b = \mathbf{x}_b$ is Gaussian.

The **second-order** term in $\log p(\mathbf{x}_a|\mathbf{x}_b)$ is

$$-\frac{1}{2}\mathbf{x}_a^T\Lambda_{aa}\mathbf{x}_a$$

Since \mathbf{x}_a given $\mathbf{x}_b = \mathbf{x}_b$ is Gaussian, this term must be

$$-\frac{1}{2}\mathbf{x}_a^T\Lambda_{a|b}\mathbf{x}_a$$

Hence

$$\Lambda_{a|b} = \Lambda_{aa}$$

Note this conditional precision is independent of \mathbf{x}_b .

The **first-order** term in $\log p(\mathbf{x}_a|\mathbf{x}_b)$ is

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

Since \mathbf{x}_a given $\mathbf{x}_b = \mathbf{x}_b$ is Gaussian, this term must be

$$\mathbf{x}_a^T \Lambda_{a|b} \boldsymbol{\mu}_{a|b}$$

So $\Lambda_{aa} \boldsymbol{\mu}_a - \Lambda_{ab} (\mathbf{x}_b - \boldsymbol{\mu}_b) = \Lambda_{a|b} \boldsymbol{\mu}_{a|b}$. Hence

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

Note this conditional mean is **linear** in \mathbf{x}_b .

INVERSE OF A PARTITIONED MATRIX

Consider $\begin{pmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{pmatrix}$ where \mathbf{A} and \mathbf{D} are invertible. Then

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

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

Note that $\mathbf{M}^{-1} = \mathbf{A} - \mathbf{B}\mathbf{D}^{-1}\mathbf{C}$ is also known as the **Schur complement** of \mathbf{D} .

INVERSE COVARIANCE MATRIX

For a covariance matrix

$$\begin{aligned} \begin{pmatrix} \Lambda_{aa} & \Lambda_{ab} \\ \Lambda_{ba} & \Lambda_{bb} \end{pmatrix} &= \begin{pmatrix} \Sigma_{aa} & \Sigma_{ab} \\ \Sigma_{ba} & \Sigma_{bb} \end{pmatrix}^{-1} \\ &= \begin{pmatrix} \left(\Sigma_{aa} - \Sigma_{ab} \Sigma_{bb}^{-1} \Sigma_{ba} \right)^{-1} & - \left(\Sigma_{aa} - \Sigma_{ab} \Sigma_{bb}^{-1} \Sigma_{ba} \right)^{-1} \Sigma_{ab} \Sigma_{bb}^{-1} \\ * & * \end{pmatrix} \end{aligned}$$

Hence

$$\begin{aligned} \Lambda_{aa} &= \left(\Sigma_{aa} - \Sigma_{ab} \Sigma_{bb}^{-1} \Sigma_{ba} \right)^{-1} \\ \Lambda_{ab} &= - \left(\Sigma_{aa} - \Sigma_{ab} \Sigma_{bb}^{-1} \Sigma_{ba} \right)^{-1} \Sigma_{ab} \Sigma_{bb}^{-1} \end{aligned}$$

Similarly

$$\begin{aligned} \Sigma_{aa} &= \left(\Lambda_{aa} - \Lambda_{ab} \Lambda_{bb}^{-1} \Lambda_{ba} \right)^{-1} \\ \Sigma_{ab} &= - \left(\Lambda_{aa} - \Lambda_{ab} \Lambda_{bb}^{-1} \Lambda_{ba} \right)^{-1} \Lambda_{ab} \Lambda_{bb}^{-1} \end{aligned}$$

ALTERNATIVE PARAMETERIZATION

We can use covariance instead of precision.

The conditional mean and covariance of \mathbf{x}_a given $\mathbf{x}_b = \mathbf{x}_b$ can be expressed by covariance as

$$\boldsymbol{\mu}_{a|b} = \boldsymbol{\mu}_a + \boldsymbol{\Sigma}_{ab}\boldsymbol{\Sigma}_{bb}^{-1}(\mathbf{x}_b - \boldsymbol{\mu}_b)$$

$$\boldsymbol{\Sigma}_{a|b} = \boldsymbol{\Sigma}_{aa} - \boldsymbol{\Sigma}_{ab}\boldsymbol{\Sigma}_{bb}^{-1}\boldsymbol{\Sigma}_{ba}$$

This follow from

$$\begin{aligned}\boldsymbol{\Lambda}_{aa}^{-1}\boldsymbol{\Lambda}_{ab} &= \left(\boldsymbol{\Sigma}_{aa} - \boldsymbol{\Sigma}_{ab}\boldsymbol{\Sigma}_{bb}^{-1}\boldsymbol{\Sigma}_{ba}\right) \left(-\left(\boldsymbol{\Sigma}_{aa} - \boldsymbol{\Sigma}_{ab}\boldsymbol{\Sigma}_{bb}^{-1}\boldsymbol{\Sigma}_{ba}\right)^{-1}\boldsymbol{\Sigma}_{ab}\boldsymbol{\Sigma}_{bb}^{-1}\right) \\ &= -\boldsymbol{\Sigma}_{ab}\boldsymbol{\Sigma}_{bb}^{-1}\end{aligned}$$

$$\boldsymbol{\Sigma}_{a|b} = \boldsymbol{\Lambda}_{a|b}^{-1} = \boldsymbol{\Lambda}_{aa}^{-1} = \boldsymbol{\Sigma}_{aa} - \boldsymbol{\Sigma}_{ab}\boldsymbol{\Sigma}_{bb}^{-1}\boldsymbol{\Sigma}_{ba}$$

\mathbf{x}_a is Gaussian with mean $\boldsymbol{\mu}_a$ and covariance $\boldsymbol{\Sigma}_{aa}$.

Proof. By marginalization

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

The integration function $p(\mathbf{x}_a, \mathbf{x}_b)$ is exponential with exponent

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

The terms involved in the integration over \mathbf{x}_b is

$$\begin{aligned} & -\frac{1}{2}\mathbf{x}_b^T \Lambda_{bb}\mathbf{x}_b - \mathbf{x}_b^T \Lambda_{ba}\mathbf{x}_a + \mathbf{x}_b^T (\Lambda_{ba}\boldsymbol{\mu}_a + \Lambda_{bb}\boldsymbol{\mu}_b) \\ & = -\frac{1}{2}\mathbf{x}_b^T \Lambda_{bb}\mathbf{x}_b + \mathbf{x}_b^T \mathbf{m} \\ & = -\frac{1}{2}(\mathbf{x}_b - \Lambda_{bb}^{-1}\mathbf{m})^T \Lambda_{bb}(\mathbf{x}_b - \Lambda_{bb}^{-1}\mathbf{m}) + \frac{1}{2}\mathbf{m}^T \Lambda_{bb}^{-1}\mathbf{m} \end{aligned}$$

where $\mathbf{m} = \Lambda_{bb}\boldsymbol{\mu}_b + \Lambda_{ba}\boldsymbol{\mu}_a - \Lambda_{ba}\mathbf{x}_a = \Lambda_{bb}\boldsymbol{\mu}_b - \Lambda_{ba}(\mathbf{x}_a - \boldsymbol{\mu}_a)$.

After the integration over \mathbf{x}_b , the remaining exponent is

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

Thus, the covariance of \mathbf{x}_a is

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

The mean of \mathbf{x}_a is

$$\boldsymbol{\mu}_{a|\emptyset} = \boldsymbol{\mu}_a$$

SUMMARY

Suppose $\mathbf{x} \sim \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma})$ and

$$\mathbf{x} = \begin{pmatrix} \mathbf{x}_a \\ \mathbf{x}_b \end{pmatrix}, \boldsymbol{\mu} = \begin{pmatrix} \boldsymbol{\mu}_a \\ \boldsymbol{\mu}_b \end{pmatrix}, \boldsymbol{\Sigma} = \begin{pmatrix} \boldsymbol{\Sigma}_{aa} & \boldsymbol{\Sigma}_{ab} \\ \boldsymbol{\Sigma}_{ba} & \boldsymbol{\Sigma}_{bb} \end{pmatrix}, \boldsymbol{\Lambda} = \begin{pmatrix} \boldsymbol{\Lambda}_{aa} & \boldsymbol{\Lambda}_{ab} \\ \boldsymbol{\Lambda}_{ba} & \boldsymbol{\Lambda}_{bb} \end{pmatrix}$$

The conditional distribution of \mathbf{x}_a given $\mathbf{x}_b = \mathbf{x}_b$ is Gaussian

$$\begin{aligned} p(\mathbf{x}_a|\mathbf{x}_b) &= \mathcal{N}(\mathbf{x}_a|\boldsymbol{\mu}_{a|b}, \boldsymbol{\Sigma}_{a|b}) \\ \boldsymbol{\mu}_{a|b} &= \boldsymbol{\mu}_a - \boldsymbol{\Lambda}_{aa}^{-1} \boldsymbol{\Lambda}_{ab}(\mathbf{x}_b - \boldsymbol{\mu}_b) \\ \boldsymbol{\Sigma}_{a|b} &= \boldsymbol{\Lambda}_{aa}^{-1} \end{aligned}$$

The marginal distribution of \mathbf{x}_a is Gaussian

$$\begin{aligned} p(\mathbf{x}_a) &= \mathcal{N}(\mathbf{x}_a|\boldsymbol{\mu}_{a|\emptyset}, \boldsymbol{\Sigma}_{a|\emptyset}) \\ \boldsymbol{\mu}_{a|\emptyset} &= \boldsymbol{\mu}_a \\ \boldsymbol{\Sigma}_{a|\emptyset} &= \boldsymbol{\Sigma}_{aa} \end{aligned}$$

A linear Gaussian model for \mathbf{y} and \mathbf{z} assumes

- \mathbf{y} is Gaussian
- \mathbf{z} is conditional Gaussian given $\mathbf{y} = \mathbf{y}$
- The conditional mean of \mathbf{z} given $\mathbf{y} = \mathbf{y}$ is linear in \mathbf{y}
- The conditional covariance of \mathbf{z} given $\mathbf{y} = \mathbf{y}$ does not depend on \mathbf{y}

That is

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

and

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

The joint distribution of a linear Gaussian model is Gaussian.

By product rule

$$p(\mathbf{y}, \mathbf{z}) = p(\mathbf{y})p(\mathbf{z}|\mathbf{y})$$

So

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

which is quadratic in $\begin{pmatrix} \mathbf{y} \\ \mathbf{z} \end{pmatrix} = \mathbf{x} = \begin{pmatrix} \mathbf{x}_a \\ \mathbf{x}_b \end{pmatrix}$.

PRECISION AND COVARIANCE

The second-order term in $\log p(\mathbf{x})$ is

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

Hence the precision and covariance of \mathbf{x} are

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

The first-order term in $\log p(\mathbf{x})$ is

$$\begin{aligned} \mathbf{x}^T \Sigma_{\mathbf{x}}^{-1} \boldsymbol{\mu}_{\mathbf{x}} &= \mathbf{y}^T \Lambda \boldsymbol{\mu} - \mathbf{y}^T \mathbf{A}^T \mathbf{L} \mathbf{b} + \mathbf{z}^T \mathbf{L} \mathbf{b} \\ &= \begin{pmatrix} \mathbf{y} \\ \mathbf{z} \end{pmatrix}^T \begin{pmatrix} \Lambda \boldsymbol{\mu} - \mathbf{A}^T \mathbf{L} \mathbf{b} \\ \mathbf{L} \mathbf{b} \end{pmatrix} \end{aligned}$$

Hence the mean of \mathbf{x} is

$$\begin{aligned} \boldsymbol{\mu}_{\mathbf{x}} &= \Sigma_{\mathbf{x}} \begin{pmatrix} \Lambda \boldsymbol{\mu} - \mathbf{A}^T \mathbf{L} \mathbf{b} \\ \mathbf{L} \mathbf{b} \end{pmatrix} \\ &= \begin{pmatrix} \boldsymbol{\mu} \\ \mathbf{A} \boldsymbol{\mu} + \mathbf{b} \end{pmatrix} \\ &= \begin{pmatrix} \boldsymbol{\mu}_a \\ \boldsymbol{\mu}_b \end{pmatrix} \end{aligned}$$

Consider a linear Gaussian model for \mathbf{y} and \mathbf{z} where \mathbf{y} is Gaussian and \mathbf{z} is conditional Gaussian given $\mathbf{y} = \mathbf{y}$. Then

- \mathbf{z} is Gaussian.
- \mathbf{y} is conditional Gaussian given $\mathbf{z} = \mathbf{z}$.