

Notes 2-Probability Distribution

Junda Du

2019-8-12

1 Binary Variables

Bernoulli Distribution Points:

1. Definition: $x \in \{0, 1\}$, $p(x = 1|\mu) = \mu$, the probability distribution over x is $Bern(x|\mu) = \mu^x(1 - \mu)^{1-x}$, which is known as *Bernoulli* distribution. It's normalized and easy to get $E[x] = \mu$, $var[x] = \mu(1 - \mu)$.

2. ML to estimate parameters given data set $D(\text{default} : \{x_1, \dots, x_N\} \text{ iid})$:

- $p(D|\mu) = \prod_{n=1}^N p(x_n|\mu) = \prod_{n=1}^N \mu^{x_n}(1 - \mu)^{1-x_n}$
- $\ln p(D|\mu) = \sum_{n=1}^N \ln p(x_n|\mu) = \sum_{n=1}^N \{x_n \ln \mu + (1 - x_n) \ln(1 - \mu)\}$
- Set the derivative of $\ln p(D|\mu)$ with respect to μ equal to zero, then obtain ML estimator $\mu_{ML} = \frac{1}{N} \sum_{n=1}^N x_n$

3. Bayesian treatment for the over-fitting problem given small data sets.

- Introduce a conjugate prior distribution $p(\mu) \propto \mu^a(1 - \mu)^b$ according to $p(D|\mu) = \prod_{n=1}^N p(x_n|\mu) = \prod_{n=1}^N \mu^{x_n}(1 - \mu)^{1-x_n}$ [Same form for prior variable, const a, b, c, \dots for $f(x_n)$ or non-prior variables], so that the posterior probability $p(\mu|D)$ will have the same functional form as the prior probability.
- Choose a prior $p(\mu) = \text{Beta}(\mu|a, b) = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \mu^{a-1}(1-\mu)^{b-1}$, which is normalized by $\frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)}$. Its $E[\mu] = \frac{a}{a+b}$, $var[\mu] = \frac{ab(a+b)^2(a+b+1)}{(a+b)^2(a+b+1)}$.
- The posterior probability $p(\mu|m, l, a, b) \propto \mu^{m+a-1}\mu^{l+b-1}$, $l = N - m$, then normalizing it according to the form of the prior probability, get $p(\mu|m, l, a, b) = \frac{\Gamma(m+a+l+b)}{\Gamma(m+a)\Gamma(m+b)} \mu^{m+a-1}\mu^{l+b-1}$, $l = N - m$
- When getting a new observation, we can take the old prior probability as the prior probability, and multiply it by the likelihood function for the new observation to get the current prior probability. It's called the *sequential* approach and depends only on the assumption of i.i.d. data.

- Utilize the posterior probability to predict, $p(x = 1|D) \xrightarrow{\text{Bayesian's theorem}} \int_0^1 p(x = 1|\mu)p(\mu|D)d\mu = \int_0^1 \mu p(\mu|D)d\mu = E[\mu|D] \xrightarrow{\text{Beta distribution's identity}} \frac{\frac{m+a}{m+a+l+b}}$

4. Show that **generally**, as we observe more and more data, the uncertainty represented by posterior distribution **on average** will steadily decrease. First, let take the mean for example, the prior mean $E_\theta[\theta] = \int \theta p(\theta)d\theta$, the posterior mean on average $E_D[E_\theta[\theta|D]] = \int \{\int \theta p(\theta|D)d\theta\}p(D)dD$, easily get $E_\theta[\theta] = E_D[E_\theta[\theta|D]]$. Similarly, we can show that $var_\theta[\theta] = E_D[var_\theta[\theta|D]] + var_D[E_\theta[\theta|D]]$, and variance ≥ 0 .

Binomial Distribution Points:

1. Definition: Given the data set $D [p(x_i|\mu = Bern(x_i|\mu))]$, m is the number of the observations of $x = 1$. The probability distribution over m $p(m|N, \mu) \propto \mu^m(1 - \mu)^{N-m}$. After normalization,

$$p(m|N, \mu) = Bin(m|N, \mu) = \binom{N}{m} \mu^m (1 - \mu)^{N-m}, \binom{N}{m} = \frac{N!}{(N-m)!m!}$$

2. $E[m] \xrightarrow{\text{definition}} \sum_{m=0}^N m Bin(m|N, \mu) \xrightarrow{m=x_1+\dots+x_N} \sum_{n=0}^N E[x_n] = N\mu$
 $var[m] \xrightarrow{\text{definition}} E[(m - E[m])^2] = \sum_{m=0}^N (m - E[m])^2 Bin(m|N, \mu)$
 $\xrightarrow{m=x_1+\dots+x_N, iid} \sum_{n=0}^N var[x_n] = N\mu(1 - \mu)$

2 Multinomial Variables

1. Definition: K-dimensional vector \mathbf{x} only have one element $x_k = 1$, all remaining elements = 0, $p(x_k = 1) = \mu_k, \sum_k \mu_k = 1$. The probability over \mathbf{x} is $p(\mathbf{x}|\boldsymbol{\mu}) = \prod_{k=1}^K \mu_k^{x_k}, \sum_k \mu_k = 1$. Easily get $\sum_{\mathbf{x}} p(\mathbf{x}|\boldsymbol{\mu}) = \sum_{k=1}^K \mu_k = 1, E[\mathbf{x}|\boldsymbol{\mu}] = \sum_{\mathbf{x}} p(\mathbf{x}|\boldsymbol{\mu}) \mathbf{x} = (\mu_1, \dots, \mu_M)^T = \boldsymbol{\mu}$
2. ML to estimate parameters

- $p(D|\boldsymbol{\mu}) = \prod_{n=1}^N \prod_{k=1}^K \mu_k^{x_{nk}} = \prod_{k=1}^K \mu_k^{(\sum_n x_{nk})} = \prod_{k=1}^K \mu_k^{m_k}, m_k = \sum_n x_{nk}$
- $\ln p(D|\boldsymbol{\mu}) + \lambda(\sum_k \mu_k - 1) = \sum_{k=1}^K m_k \ln \mu_k + \lambda(\sum_k \mu_k - 1)$
- Set the derivative with respect to μ_k to zero, obtain $\mu_k = -\frac{m_k}{\lambda}$, substituting it into constraint $\sum_k \mu_k = 1$ to give $\lambda = -N$, so $\mu_k^{ML} = \frac{m_k}{N}$

3. Bayesian's treatment

- Conjugate prior $p(\boldsymbol{\mu}|\boldsymbol{\alpha}) \propto \prod_{k=1}^K \mu_k^{\alpha_k-1}, 0 \leq \mu_k \leq 1, \sum_k \mu_k = 1$, according to $p(D|\boldsymbol{\mu}) = \prod_{k=1}^K \mu_k^{m_k}$
 - After normalization, $p(\boldsymbol{\mu}|\boldsymbol{\alpha}) = Dir(\boldsymbol{\mu}|\boldsymbol{\alpha}) = \frac{\Gamma(\alpha_0)}{\Gamma(\alpha_1)\dots\Gamma(\alpha_K)} \prod_{k=1}^K \mu_k^{\alpha_k-1}, \alpha_0 = \sum_{k=1}^K \alpha_k$
 - $p(\boldsymbol{\mu}|D, \boldsymbol{\alpha}) \propto p(D|\boldsymbol{\mu})p(\boldsymbol{\mu}|\boldsymbol{\alpha}) \propto \prod_{k=1}^K \mu_k^{\alpha_k+m_k-1}$, then normalizing it according to the form of the prior probability, get $p(\boldsymbol{\mu}|D, \boldsymbol{\alpha}) = Dir(\boldsymbol{\mu}|\boldsymbol{\alpha} + \mathbf{m}) = \frac{\Gamma(\alpha_0+N)}{\Gamma(\alpha_1+m_1)\dots\Gamma(\alpha_K+m_K)} \prod_{k=1}^K \mu_k^{\alpha_k+m_k-1}$
4. *Multinomial* distribution
- $$Mult(m_1, \dots, m_K | \boldsymbol{\mu}, N) \propto p(D|\boldsymbol{\mu}) \xrightarrow{\text{Normalizing}} \left(\frac{N}{m_1 m_2 \dots m_K} \right) \prod_{k=1}^K \mu_k^{m_k}, \sum_k m_k = N, \left(\frac{N}{m_1 m_2 \dots m_K} \right) = \frac{N!}{m_1! \dots m_K!}$$

3 The Gaussian Distribution

Points:

- Definition of Gaussian distribution:

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

$$N(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{\frac{D}{2}} |\boldsymbol{\Sigma}|^{\frac{1}{2}}} \exp\left\{-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right\}$$
 for multivariable.
- Situations it arises:

A single variable's distribution that maximizes the entropy is the Gaussian.

The sum of a set of random variables has a Gaussian distribution.
- Manipulating Gaussian distributions:
 - Represent Gaussian in a simpler form:
 1. We focus on the part dependent on \mathbf{x} : $\Delta^2 = (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})$, which is a quadratic form. Δ is called *Mahalanobis distance* from $\boldsymbol{\mu}$ to \mathbf{x} , and reduces to Euclidean distance when $\boldsymbol{\Sigma}$ is the identity matrix.
 2. Consider eigenvector equation for $\boldsymbol{\Sigma}$, $\boldsymbol{\Sigma}\mathbf{u}_i = \lambda_i \mathbf{u}_i$, where $i = 1, \dots, D$ and eigenvector can be chosen to form an orthonormal set because $\boldsymbol{\Sigma}$ must be real and symmetric?. So $\mathbf{u}_i^T \mathbf{u}_j = I_{ij}$.

3. Then Σ, Σ^{-1} can be expressed as these form: $\Sigma = \sum_{i=1}^D \lambda_i \mathbf{u}_i \mathbf{u}_i^T, \Sigma^{-1} = \sum_{i=1}^D \frac{1}{\lambda_i} \mathbf{u}_i \mathbf{u}_i^T$.
4. So the quadratic form becomes $\Delta^2 = \sum_{i=1}^D \frac{y_i^2}{\lambda_i}, y_i = \mathbf{u}_i^T (\mathbf{x} - \boldsymbol{\mu})$.
5. In going from the \mathbf{x} to the \mathbf{y} coordinate system, we have a Jacobian matrix \mathbf{J} with elements given by $J_{ij} = \frac{\delta x_i}{\delta y_j} = U_j i$, because $\mathbf{y} = \mathbf{U}(\mathbf{x} - \boldsymbol{\mu}), \mathbf{U} = (\mathbf{u}_1^T, \dots, \mathbf{u}_D^T)^T$, thus $\mathbf{U}\mathbf{U}^T = \mathbf{I}$. So $|\mathbf{J}|^2 = |\mathbf{U}^T|^2 = |\mathbf{U}^T| |\mathbf{U}| = |\mathbf{U}^T \mathbf{U}| = |\mathbf{I}| = 1$. Also $|\Sigma|^{1/2} = \prod_{j=1}^D \lambda_j^{1/2}$.
6. Finally $p(\mathbf{y}) = p(\mathbf{x}) |\mathbf{J}| = \prod_{j=1}^D \frac{1}{(2\pi \lambda_j)^{1/2}} \exp\{-\frac{y_j^2}{2\lambda_j}\}$, which is the product of D independent univariate Gaussian distributions.
7. Provide an interpretation of the parameters $\boldsymbol{\mu}$ and Σ :
- Let's look at the moments and the covariance of the Gaussian distribution.
1. For first order moments/expectation:

$$\begin{aligned}
E[\mathbf{x}] &= \frac{1}{(2\pi)^{\frac{D}{2}} |\Sigma|^{\frac{1}{2}}} \int \exp\{-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu})\} \mathbf{x} d\mathbf{x} \\
&\stackrel{\mathbf{z}=\mathbf{x}-\boldsymbol{\mu}}{=} \frac{1}{(2\pi)^{\frac{D}{2}} |\Sigma|^{\frac{1}{2}}} \int \exp\{-\frac{1}{2}\mathbf{z}^T \Sigma^{-1} \mathbf{z}\} (\mathbf{z} + \boldsymbol{\mu}) d\mathbf{z} \\
&\stackrel{\text{even, odd, symmetry}}{=} \boldsymbol{\mu}
\end{aligned}$$

2. For second order moments:

$$\begin{aligned}
E[\mathbf{x}\mathbf{x}^T] &= \frac{1}{(2\pi)^{\frac{D}{2}} |\Sigma|^{\frac{1}{2}}} \int \exp\{-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu})\} \mathbf{x}\mathbf{x}^T d\mathbf{x} \\
&\stackrel{\mathbf{z}=\mathbf{x}-\boldsymbol{\mu}}{=} \frac{1}{(2\pi)^{\frac{D}{2}} |\Sigma|^{\frac{1}{2}}} \int \exp\{-\frac{1}{2}\mathbf{z}^T \Sigma^{-1} \mathbf{z}\} (\mathbf{z} + \boldsymbol{\mu})(\mathbf{z} + \boldsymbol{\mu})^T d\mathbf{z} \\
&\stackrel{\text{even, odd, symmetry}}{=} \frac{1}{(2\pi)^{\frac{D}{2}} |\Sigma|^{\frac{1}{2}}} \int \exp\{-\frac{1}{2}\mathbf{z}^T \Sigma^{-1} \mathbf{z}\} (\mathbf{z}\mathbf{z}^T + \boldsymbol{\mu}\boldsymbol{\mu}^T) d\mathbf{z} \\
&\stackrel{\mathbf{z}=\mathbf{U}^{-1}\mathbf{y}=\mathbf{U}^T\mathbf{y}=\sum_{j=1}^D y_j \mathbf{u}_j}{\Sigma^{-1}=\sum_{i=1}^D \frac{1}{\lambda_i} \mathbf{u}_i \mathbf{u}_i^T}{=} \boldsymbol{\mu}\boldsymbol{\mu}^T + \frac{1}{(2\pi)^{\frac{D}{2}} |\Sigma|^{\frac{1}{2}}} \sum_{i=1}^D \sum_{j=1}^D \mathbf{u}_i \mathbf{u}_j^T \int \exp\{-\sum_{k=1}^D \frac{y_k^2}{2\lambda_k}\} y_i y_j \\
&\stackrel{\mathbf{u}_i^T \mathbf{u}_j = I_{ij}, |\Sigma|^{1/2} = \prod_{j=1}^D \lambda_j^{1/2}}{\text{split } i \text{ from the integral}}{=} \boldsymbol{\mu}\boldsymbol{\mu}^T + \sum_{i=1}^D \mathbf{u}_i \mathbf{u}_i^T \lambda_i = \boldsymbol{\mu}\boldsymbol{\mu}^T + \Sigma
\end{aligned}$$

3. For covariance:

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

4. As we can see:

$$\boldsymbol{\mu} = E[\mathbf{x}], \Sigma = \text{cov}[\mathbf{x}] = E[(\mathbf{x} - E[\mathbf{x}])(\mathbf{x} - E[\mathbf{x}])^T]$$

3.1 Conditional Gaussian distribution

We want to get $p(\mathbf{x}_a|\mathbf{x}_b)$, and $p(\mathbf{x}) = p(\mathbf{x}_a, \mathbf{x}_b)$ is Gaussian. We set

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

where $\Sigma_{ab} = \Sigma_{ba}^T$ because Σ is symmetric, $\Lambda = \Sigma^{-1}$ and $\Lambda_{ab} = \Lambda_{ba}^T$ because the inverse of a symmetric matrix is also symmetric, but Λ_{aa} is not simply given by Σ_{aa}^{-1} . Next we will discuss the relationship between Λ_{aa} and Σ_{aa}^{-1} .

There is a identity for the inverse of a partitioned matrix,

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

The quantity \mathbf{M}^{-1} is known as the *Schur complement* of the left-hand side with respect to \mathbf{D} .

Using the definition

$$\begin{pmatrix} \Sigma_{aa} & \Sigma_{ab} \\ \Sigma_{ba} & \Sigma_{bb} \end{pmatrix}^{-1} = \begin{pmatrix} \Lambda_{aa} & \Lambda_{ab} \\ \Lambda_{ba} & \Lambda_{bb} \end{pmatrix},$$

we have $\Lambda_{aa} = (\Sigma_{aa} - \Sigma_{ab}\Sigma_{bb}^{-1}\Sigma_{ba})^{-1}$, $\Lambda_{ab} = -(\Sigma_{aa} - \Sigma_{ab}\Sigma_{bb}^{-1}\Sigma_{ba})^{-1}\Sigma_{ab}\Sigma_{bb}^{-1}$
 $p(\mathbf{x}_a|\mathbf{x}_b)$ is simply evaluated from the joint distribution $p(\mathbf{x}_a, \mathbf{x}_b)$ by fixing \mathbf{x}_b and then normalizing to obtain a valid probability over \mathbf{x}_a . We only consider the exponent part because normalization process is adapt coefficient to the exponent.

$$\begin{aligned} -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1}(\mathbf{x} - \boldsymbol{\mu}) &= -\frac{1}{2}(\mathbf{x}_a - \boldsymbol{\mu}_a)^T \Lambda_{aa}(\mathbf{x}_a - \boldsymbol{\mu}_a) - \frac{1}{2}(\mathbf{x}_a - \boldsymbol{\mu}_a)^T \Lambda_{ab}(\mathbf{x}_b - \boldsymbol{\mu}_b) \\ &\quad - \frac{1}{2}(\mathbf{x}_b - \boldsymbol{\mu}_b)^T \Lambda_{ba}(\mathbf{x}_a - \boldsymbol{\mu}_a) - \frac{1}{2}(\mathbf{x}_b - \boldsymbol{\mu}_b)^T \Lambda_{bb}(\mathbf{x}_b - \boldsymbol{\mu}_b) \\ &\quad \underline{\underline{\text{fix } \mathbf{x}_b, \text{complete the square}}} - \frac{1}{2}\mathbf{x}_a^T \Lambda_{aa} \mathbf{x}_a + \mathbf{x}_a^T \{\Lambda_{aa}\boldsymbol{\mu}_a - \Lambda_{ab}(\mathbf{x}_b - \boldsymbol{\mu}_b)\} + \text{const} \end{aligned}$$

There is a technique called "**completing the square**", in which we obtain the required mean and variance by comparing given exponent with a general Gaussian distribution $N(\mathbf{x}|\boldsymbol{\mu}, \Sigma)$'s exponent $-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1}(\mathbf{x} - \boldsymbol{\mu}) = -\frac{1}{2}\mathbf{x}^T \Sigma^{-1} \mathbf{x} + \mathbf{x}^T \Sigma^{-1} \boldsymbol{\mu} + \text{const}$

So we can obtain

$$\Sigma_{a|b} = \Lambda_{aa}^{-1},$$

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

we can also substitute $\Lambda_{aa} = (\Sigma_{aa} - \Sigma_{ab}\Sigma_{bb}^{-1}\Sigma_{ba})^{-1}$, $\Lambda_{ab} = -(\Sigma_{aa} - \Sigma_{ab}\Sigma_{bb}^{-1}\Sigma_{ba})^{-1}\Sigma_{ab}\Sigma_{bb}^{-1}$ to obtain a new representation.

3.2 Marginal Gaussian distribution

Usually, we use $p(\mathbf{x}_a) = \int p(\mathbf{x}_a, \mathbf{x}_b) d\mathbf{x}_b$ to get marginal probability over \mathbf{x}_a , so first we consider the terms involving \mathbf{x}_b and then completing the square to facilitate integration.

$$\begin{aligned}
-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}) &= -\frac{1}{2}(\mathbf{x}_a - \boldsymbol{\mu}_a)^T \boldsymbol{\Lambda}_{aa}(\mathbf{x}_a - \boldsymbol{\mu}_a) - \frac{1}{2}(\mathbf{x}_a - \boldsymbol{\mu}_a)^T \boldsymbol{\Lambda}_{ab}(\mathbf{x}_b - \boldsymbol{\mu}_b) \\
&\quad - \frac{1}{2}(\mathbf{x}_b - \boldsymbol{\mu}_b)^T \boldsymbol{\Lambda}_{ba}(\mathbf{x}_a - \boldsymbol{\mu}_a) - \frac{1}{2}(\mathbf{x}_b - \boldsymbol{\mu}_b)^T \boldsymbol{\Lambda}_{bb}(\mathbf{x}_b - \boldsymbol{\mu}_b) \\
&\quad \xrightarrow{\text{fix } \mathbf{x}_a, \text{complete the square}} -\frac{1}{2}\mathbf{x}_b^T \boldsymbol{\Lambda}_{bb} \mathbf{x}_b + \mathbf{x}_b^T \mathbf{m} + \mathbf{R}(\mathbf{x}_a), \mathbf{m} = \boldsymbol{\Lambda}_{bb} \boldsymbol{\mu}_b - \boldsymbol{\Lambda}_{ba}(\mathbf{x}_a - \boldsymbol{\mu}_a) \\
&\quad \xrightarrow{\text{facilitate integration}} -\frac{1}{2}(\mathbf{x}_b - \boldsymbol{\Lambda}_{bb}^{-1} \mathbf{m})^T \boldsymbol{\Lambda}_{bb}(\mathbf{x}_b - \boldsymbol{\Lambda}_{bb}^{-1} \mathbf{m}) + \frac{1}{2} \mathbf{m}^T \boldsymbol{\Lambda}_{bb}^{-1} \mathbf{m} + \mathbf{R}(\mathbf{x}_a)
\end{aligned}$$

Since $\frac{1}{2} \mathbf{m}^T \boldsymbol{\Lambda}_{bb}^{-1} \mathbf{m} + \mathbf{R}(\mathbf{x}_a)$ is independent on \mathbf{x}_b , and $\int \exp\{-\frac{1}{2}(\mathbf{x}_b - \boldsymbol{\Lambda}_{bb}^{-1} \mathbf{m})^T \boldsymbol{\Lambda}_{bb}(\mathbf{x}_b - \boldsymbol{\Lambda}_{bb}^{-1} \mathbf{m})\} d\mathbf{x}_b$ becomes the coefficient because it's a normalized Gaussian. So exponent part becomes

$$\begin{aligned}
\frac{1}{2} \mathbf{m}^T \boldsymbol{\Lambda}_{bb}^{-1} \mathbf{m} + \mathbf{R}(\mathbf{x}_a) &= \frac{1}{2} [\boldsymbol{\Lambda}_{bb} \boldsymbol{\mu}_b - \boldsymbol{\Lambda}_{ba}(\mathbf{x}_a - \boldsymbol{\mu}_a)]^T \boldsymbol{\Lambda}_{bb}^{-1} [\boldsymbol{\Lambda}_{bb} \boldsymbol{\mu}_b - \boldsymbol{\Lambda}_{ba}(\mathbf{x}_a - \boldsymbol{\mu}_a)] \\
&\quad - \frac{1}{2} \mathbf{x}_a^T \boldsymbol{\Lambda}_{aa} \mathbf{x}_a + \mathbf{x}_a^T (\boldsymbol{\Lambda}_{aa} \boldsymbol{\mu}_a + \boldsymbol{\Lambda}_{ab} \boldsymbol{\mu}_b) + \text{const} \\
&\quad \xrightarrow{\text{fix } \mathbf{x}_b, \text{complete the square}} -\frac{1}{2} \mathbf{x}_a^T (\boldsymbol{\Lambda}_{aa} - \boldsymbol{\Lambda}_{ab} \boldsymbol{\Lambda}_{bb}^{-1} \boldsymbol{\Lambda}_{ba}) \mathbf{x}_a + \\
&\quad \mathbf{x}_a^T (\boldsymbol{\Lambda}_{aa} - \boldsymbol{\Lambda}_{ab} \boldsymbol{\Lambda}_{bb}^{-1} \boldsymbol{\Lambda}_{ba}) \boldsymbol{\mu}_a + \text{const}
\end{aligned}$$

So we can obtain
 $\boldsymbol{\Sigma}_a = (\boldsymbol{\Lambda}_{aa} - \boldsymbol{\Lambda}_{ab} \boldsymbol{\Lambda}_{bb}^{-1} \boldsymbol{\Lambda}_{ba})^{-1}$,
 $\boldsymbol{\mu}_a = \boldsymbol{\Sigma}_a (\boldsymbol{\Lambda}_{aa} - \boldsymbol{\Lambda}_{ab} \boldsymbol{\Lambda}_{bb}^{-1} \boldsymbol{\Lambda}_{ba}) \boldsymbol{\mu}_a = \boldsymbol{\mu}_a$
we can also substitute $\boldsymbol{\Lambda}_{aa} = (\boldsymbol{\Sigma}_{aa} - \boldsymbol{\Sigma}_{ab} \boldsymbol{\Sigma}_{bb}^{-1} \boldsymbol{\Sigma}_{ba})^{-1}$, $\boldsymbol{\Lambda}_{ab} = -(\boldsymbol{\Sigma}_{aa} - \boldsymbol{\Sigma}_{ab} \boldsymbol{\Sigma}_{bb}^{-1} \boldsymbol{\Sigma}_{ba})^{-1} \boldsymbol{\Sigma}_{ab} \boldsymbol{\Sigma}_{bb}^{-1}$
to obtain $\boldsymbol{\Sigma}_a = \boldsymbol{\Sigma}_{aa}$.

3.3 Bayes' theorem for Gaussian distribution

We want to obtain $p(\mathbf{y}), p(\mathbf{x}|\mathbf{y})$ given $p(\mathbf{x}) = N(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Lambda}^{-1}), p(\mathbf{y}|\mathbf{x}) = N(\mathbf{y}|\mathbf{A}\mathbf{x} + \mathbf{b}, \mathbf{L}^{-1})$.

1. First, we can obtain $p(\mathbf{z}) = p(\mathbf{x}, \mathbf{y})$ by $\ln p(\mathbf{z}) = \ln p(\mathbf{x}) + \ln p(\mathbf{y}|\mathbf{x})$.
2. Second, writing $\ln p(\mathbf{z})$ as the general Gaussian form $-\frac{1}{2} \mathbf{x}^T \boldsymbol{\Sigma}^{-1} \mathbf{x} + \mathbf{x}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu} + \text{const}$ by considering first and second order item respectively.
3. Third, use conclusion of last two sections to obtain $p(\mathbf{y}), p(\mathbf{x}|\mathbf{y})$.

Ignoring the details, we can obtain
 $p(\mathbf{y}) = N(\mathbf{y}|\mathbf{A}\boldsymbol{\mu} + \mathbf{b}, \mathbf{L}^{-1} + \mathbf{A}\boldsymbol{\Lambda}^{-1}\mathbf{A}^T)$,
 $p(\mathbf{x}|\mathbf{y}) = N(\mathbf{x}|\boldsymbol{\Sigma}\{\mathbf{A}^T \mathbf{L}(\mathbf{y} - \mathbf{b}) + \boldsymbol{\Lambda}\boldsymbol{\mu}\}, \boldsymbol{\Sigma})$

3.4 Maximum likelihood for the Gaussian

Given data set $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_N)^T$, assuming \mathbf{x}_n are drawn independently from a multivariate Gaussian distribution, we want to estimate the parameters of the Gaussian distribution by maximum likelihood.

1. Write the log likelihood function.
2. Set the derivative with respect to $\boldsymbol{\mu}, \boldsymbol{\Sigma}$ to zero respectively. Obtain

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

3. Evaluate the expectation of estimated parameters $E[\boldsymbol{\mu}_{ML}] = \boldsymbol{\mu}, E[\boldsymbol{\Sigma}_{ML}] = \frac{N-1}{N} \boldsymbol{\Sigma}$, then correct the bias obtaining

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

3.5 Sequential estimation

Robbins-Monro algorithm: If a problem is to find θ satisfying $f(\theta) = 0$, and $f(\theta)$ can be written as the form $f(\theta) = E[z|\theta] = \int zp(z|\theta)dz$, the root $\theta^{(N)} = \theta^{(N-1)} + a_{N-1}z(\theta^{(N-1)})$.

Take the maximum likelihood solution θ_{ML} for example.

By definition, θ_{ML} is a stationary point of the log likelihood function,

$$\left. \frac{\partial}{\partial \theta} \left\{ \frac{1}{N} \sum_{n=1}^N \ln p(\mathbf{x}_n | \theta) \right\} \right|_{\theta_{ML}} = 0.$$

It can be written as following form

$$\frac{\partial}{\partial \theta} \left\{ \frac{1}{N} \sum_{n=1}^N \ln p(\mathbf{x}_n | \theta) \right\} = \frac{1}{N} \sum_{n=1}^N \frac{\partial}{\partial \theta} \ln p(\mathbf{x}_n | \theta) = E_x \left[\frac{\partial}{\partial \theta} \ln p(\mathbf{x} | \theta) \right] = 0$$

$z(\theta) = \frac{\partial}{\partial \theta} \ln p(\mathbf{x} | \theta)$. Thus,

$$\theta^{(N)} = \theta^{(N-1)} + a_{N-1} z(\theta^{(N-1)}) = \theta^{(N-1)} + a_{N-1} \frac{\partial}{\partial \theta^{(N-1)}} \ln p(\mathbf{x}_N | \theta^{(N-1)})$$

Take the sequential estimation of mean of a Gaussian distribution as a more specific example,

$$z(\mu^{(N-1)}) = \frac{\partial}{\partial \mu^{(N-1)}} \ln p(\mathbf{x}_N | \mu^{(N-1)}, \sigma^2) = \frac{1}{\sigma^2} (\mathbf{x}_N - \mu^{(N-1)})$$

So we can update μ by

$$\begin{aligned}\mu^{(N)} &= \mu^{(N-1)} + a_{N-1} z(\mu^{(N-1)}) \\ &= \mu^{(N-1)} + a_{N-1} \frac{1}{\sigma^2} (\mathbf{x}_N - \mu^{(N-1)}) \\ &\stackrel{a_{N-1} = \sigma^2/N}{=} \mu^{(N-1)} + \frac{1}{N} (\mathbf{x}_N - \mu^{(N-1)})\end{aligned}$$

3.6 Bayesian inference for the Gaussian

Considering a single Gaussian random variable x , given a set of N observations $\mathbf{X} = \{x_1, \dots, x_N\}$, the likelihood function is

$$p(\mathbf{X}|\mu, \sigma^2) = \prod_{n=1}^N p(x_n|\mu, \sigma^2) = \frac{1}{(2\pi\sigma^2)^{\frac{N}{2}}} \exp\left\{-\frac{1}{2\sigma^2} \sum_{n=1}^N (x_n - \mu)^2\right\}$$

- Choose a conjugate prior $p(\mu)$ assuming variance σ^2 is known
Conjugate prior $p(\mu) = N(\mu|\mu_0, \sigma_0^2)$ [Ignoring the process \propto], posterior $p(\mu|\mathbf{X}) \propto p(\mathbf{X}|\mu)p(\mu)$.
By completing the square [i.e. considering first and second order item], we can obtain $p(\mu|\mathbf{X}) = N(\mu|\mu_N, \sigma_N^2)$,

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

From sequential view of the inference problem, we can update $p(\mu|\mathbf{X})$ by

$$p(\mu|\mathbf{X}) \propto \left[p(\mu) \prod_{n=1}^{N-1} p(\mathbf{x}_n|\mu) \right] p(\mathbf{x}_N|\mu)$$

- Choose a conjugate prior $p(\lambda)$, $\lambda = \frac{1}{\sigma^2}$ assuming mean μ is known
Conjugate prior $p(\lambda) \propto \lambda^{a-1} \exp(-b\lambda)$, after normalization $p(\lambda) = \text{Gam}(\lambda|a, b) = \frac{1}{\Gamma(a)} b^a \lambda^{a-1} \exp(-b\lambda)$.
Then we can obtain $p(\lambda|\mathbf{X})$, which is also a gamma distribution.

- Choose a conjugate prior $p(\mu, \lambda)$
The conjugate prior and posterior are both *normal-gamma* or *Gaussian-gamma* distribution.

3.7 Student's t-distribution

- We can get Student's t-distribution by integrating out the precision for univariate Gaussian multiplying the conjugate prior of precision (Gamma distribution).
- It also can be obtained by adding up an infinite number of Gaussian distributions having the same mean but different precisions.
- It's more robust than Gaussian model, i.e. less sensitive to outliers.

3.8 Periodic variables

Here we'll talk about the *circular normal* or *von Mises* distribution. Let's how it's extended from two-variate Gaussian distribution.

Consider a Gaussian distribution over two variables $\mathbf{x} = (x_1, x_2)$, having the mean $\boldsymbol{\mu} = (\mu_1, \mu_2)$ and a covariance matrix $\boldsymbol{\Sigma} = \sigma^2 \mathbf{I}$, so

$$p(x_1, x_2) = \frac{1}{2\pi\sigma^2} \exp\left\{-\frac{(x_1 - \mu_1)^2 + (x_2 - \mu_2)^2}{2\sigma^2}\right\}$$

Transforming from Cartesian coordinates to polar coordinates by $x_1 = r\cos\theta, x_2 = r\sin\theta$ also for $\boldsymbol{\mu}$ mapped by $\mu_1 = r_0\cos\theta_0, \mu_2 = r_0\sin\theta_0$

Then we obtain

$$p(\theta|\theta_0, m) = \frac{1}{2\pi I_0(m)} \exp\{m \cos(\theta - \theta_0)\}, I_0(m) = \frac{1}{2\pi} \int_0^{2\pi} \exp\{m \cos\theta\} d\theta.$$

The exponent is important, and $I_0(m)$ is the normalization coefficient. We consider the maximum likelihood estimators for the parameters θ_0 and m for the Mises distribution. The likelihood function is given by $\ln p(D|\theta_0, m) = -N \ln(2\pi) - N \ln I_0(m) + m \sum_{n=1}^N \cos(\theta_n - \theta_0)$.

$$\begin{aligned} \frac{\partial \ln p(D|\theta_0, m)}{\partial \theta_0} = 0 &\Leftrightarrow \sum_{n=1}^N \sin(\theta_n - \theta_0) = 0 \\ \xleftrightarrow{\sin(A-B) = \sin A \cos B - \cos A \sin B} \theta_0^{ML} &= \tan^{-1} \left\{ \frac{\sum_n \sin \theta_n}{\sum_n \cos \theta_n} \right\} \end{aligned}$$

which is the result obtained for the mean of observations viewed in two-dimensional Cartesian space.?

$$\begin{aligned} \frac{\partial \ln p(D|\theta_0, m)}{\partial m} &= 0 \\ \xleftrightarrow{I'_0(m) = I'_1(m)} A(m) &= \frac{1}{N} \sum_{n=1}^N \cos(\theta_n - \theta_0^{ML}), A(m) = \frac{I_1(m)}{I_0(m)} \\ \Leftrightarrow A(m_{ML}) &= \left(\frac{1}{N} \sum_{n=1}^N \cos \theta_n \right) \cos \theta_0^{ML} - \left(\frac{1}{N} \sum_{n=1}^N \sin \theta_n \right) \sin \theta_0^{ML} \end{aligned}$$

3.9 Mixtures of Gaussians

4 The Exponent Family

A little hard...

5 Nonparametric Method

The nonparametric method focus mainly on frequentist method.

5.1 Histogram density

The probability density of i_{th} box is obtained by $p_i = \frac{n_i}{N\Delta_i}$, where n_i is the number of observations of x falling in bin i , N is the total number of observations, and Δ_i is the width of i_{th} bin.

- Points:
 - When Δ is very small. the resulting density model is very spiky, with a lot of structure that is not present in the underlying distribution that generated the data set. Conversely, if Δ is too large, the result is too smooth and so fails to capture the bimodal property of the green curve. So Δ should be neither too small nor too large.
 - A histogram density model is also dependent on the choice of edge location for the bins.
- Advantages:
 - It's advantageous if the data set is large because once the histogram has been computed, the data set itself can be discarded.
 - It's easily applied if the data points are arriving sequentially.
 - It's useful for obtaining a quick visualization of data in one or two dimensions.
- Disadvantages:
 - The estimated density has discontinuities that are due to the bin edges ? rather than any property of the underlying distribution that generated the data.
 - Its scaling with dimensionality. In a space of high dimensionality D , the quantity of data needed to provide meaningful estimates of local probability density and the calculation complexity would be prohibitive. (M_D)

5.2 Kernel density estimators

The probability density is obtained by $p(\mathbf{x}) = \frac{K}{NV}$, where K is the number of observations falling in volume V around \mathbf{x} , N is the total number of observations.

The kernel approach fix V and determine K . There are some examples for kernel approach.

- Take the region R to be a small hypercube centred on the point \mathbf{x} .

$$p(\mathbf{x}) = \frac{1}{N} \sum_{n=1}^N \frac{1}{h^D} k\left(\frac{\mathbf{x} - \mathbf{x}_n}{h}\right), k(\mathbf{u}) = \begin{cases} 1, & |u_i| \leq \frac{1}{2}, i = 1, \dots, D \\ 0, & \text{else} \end{cases},$$

where $h_D = V, K = \sum_{n=1}^N k\left(\frac{\mathbf{x} - \mathbf{x}_n}{h}\right)$, h is the smooth factor. It suffer from one of the same problems that the histogram method suffered from, namely the presence of artificial discontinuities.

- A smoother density model by choosing a smoother kernel function Gaussian.

$$p(\mathbf{x}) = \frac{1}{N} \sum_{n=1}^N \frac{1}{(2\pi h^2)^{1/2}} \exp\left\{-\frac{\|\mathbf{x} - \mathbf{x}_n\|^2}{2h^2}\right\}$$

- A general kernel function should satisfies:

$$k(\mathbf{u}) \geq 0, \int k(\mathbf{u}) d\mathbf{u} = 1$$

5.3 Nearest-neighbor methods

The nearest-neighbor method fix K and determine V . It can be extended to classification problem, which is specification of Bayesian's estimation.