# Machine Learning (Homework #2)

A014587 羅羿牧

## 1. Information Theory

(a) Please show that the maximum entropy distribution for a continuous variable with three constrains

$$\int_{-\infty}^{\infty} p(x)dx = 1$$

$$\int_{-\infty}^{\infty} xp(x)dx = \mu$$

$$\int_{-\infty}^{\infty} (x-\mu)^2 p(x)dx = \sigma^2$$

is a Gaussian distribution.

Ans.



(b) Gaussian distribution is given by

$$p(x) = \frac{1}{(2\pi\sigma^2)^{1/2}} \exp\left\{ -\frac{(x-\mu)^2}{2\sigma^2} \right\}$$

Please derive the corresponding entropy.

Ans.

(b)

$$H[x] = -\int p(x)\ln p(x)\,dx = -\int p(x)\left[-\frac{1}{2}\ln(2\pi\sigma^2) - \frac{(x-\mu)^2}{2\sigma^2}\right]dx$$

$$= \frac{1}{2}\left[\ln(2\pi\sigma^2) + \frac{1}{\sigma^2}\int p(x)(x-\mu)^2\,dx\right]$$

$$= \frac{1}{2}\left[\ln(2\pi\sigma^2) + 1\right]$$

## 2. Bayesian Inference for the Gaussian

We develop a Bayesian learning by introducing prior distributions to estimate Gaussian parameters $\boldsymbol{\mu}$ and $\boldsymbol{\Sigma}$. Traditionally, batch learning is performed by using the whole training set where high computational complexity is caused. If training data is sufficiently large, it is suitable to use sequential learning (on-line learning) algorithm. Please solve the following question. The file **r2.mat** contains a 1000-point sequence, which is generated by the following multivariate Gaussian distribution $\mathcal{N}(\mathbf{x}|\boldsymbol{\mu},\boldsymbol{\Sigma})$ with $\boldsymbol{\mu} = [1,\ -1]^{\mathrm{T}}$ and $\boldsymbol{\Sigma}$ ($\boldsymbol{\Sigma}$ is unknown). The sequential learning of the posterior distribution of $\boldsymbol{\Lambda}$ ($\boldsymbol{\Lambda} = \boldsymbol{\Sigma}^{-1}$) with the contribution from the final data $\mathbf{x}_N$ can be expressed as follows:

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

(a) Please derive the posterior distribution of precision matrix $\boldsymbol{\Lambda}$, $p(\boldsymbol{\Lambda}|\mathbf{X}) = \mathcal{W}(\boldsymbol{\Lambda}|\mathbf{W_\Lambda}, \mathcal{V}_\Lambda)$, in details where $\mathcal{V}_\Lambda$ is called the *degrees of freedom* of the distribution and $\mathbf{W_\Lambda}$ is a $D \times D$ symmetric matrix. Here, we apply the conjugate prior of $\boldsymbol{\Lambda}$ which is a *Wishart* distribution $p(\boldsymbol{\Lambda}) = \mathcal{W}(\boldsymbol{\Lambda}|\mathbf{W_0}, \mathcal{V}_0)$.
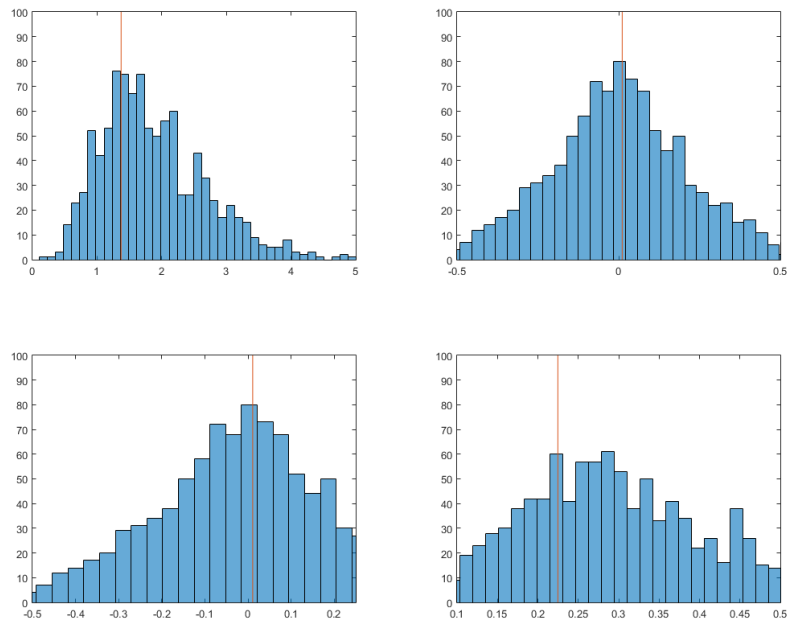
Ans.

(a)

$$p(\mu\,|\,m_1,\ m_2,\ a,\ b) \propto \mu^{m_1+a-1}(1-\mu)^{m_2+b-1}$$

$$= \frac{\Gamma(m_1+m_2+a+b)}{\Gamma(m_1+a)\,\Gamma(m_2+b)}\,\mu^{m_1+a-1}(1-\mu)^{m_2+b-1}$$

$$\mu_{MAP} = \frac{\partial \ln p(\mu|m_1,m_2,a,b)}{\partial\mu} = \frac{\partial}{\partial\mu}\ln\left[K\cdot\mu^{m_1+a-1}(1-\mu)^{m_2+b-1}\right]$$

$$\Rightarrow \frac{\partial}{\partial\mu}\left[(m_1+a-1)\ln\mu + (m_2+b-1)\ln(1-\mu)\right] = 0$$

$$\Rightarrow \frac{m_1+a-1}{\mu} = \frac{m_2+b-1}{1-\mu} \Rightarrow (m_1+a-1)(1-\mu) = \mu(m_2+b-1)$$
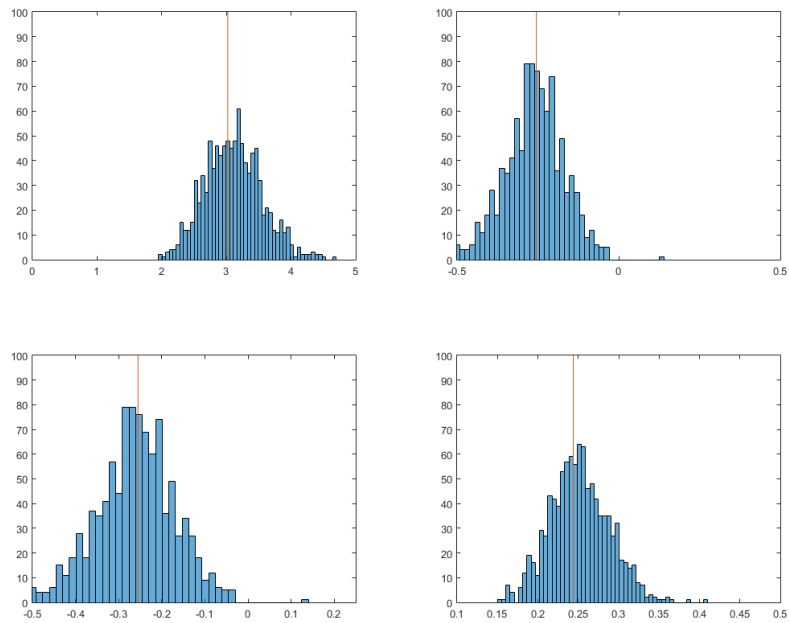
$$\Rightarrow \mu = \frac{m_1+a-1}{m_1+m_2+a+b-2}$$

(b) Please consider the *Wishart* prior $p_1(\mathbf{\Lambda}) = \mathcal{W}(\mathbf{\Lambda}|\begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}, 1)$ and find the MAP

solution of $\mathbf{\Lambda}$ (or $\mathbf{\Sigma}$) for $N = 10$, 100, and 500. ($\mathbf{\Lambda}_{\textbf{MAP}} = \text{argmax}_{\mathbf{\Lambda}} \, p(\mathbf{\Lambda}|\mathbf{X})$) You may also directly use the Matlab command 'wishrnd' to generate many samples of $\mathbf{\Lambda}$ and compare their corresponding $p(\mathbf{\Lambda})$ to obtain the approximate MAP solution.
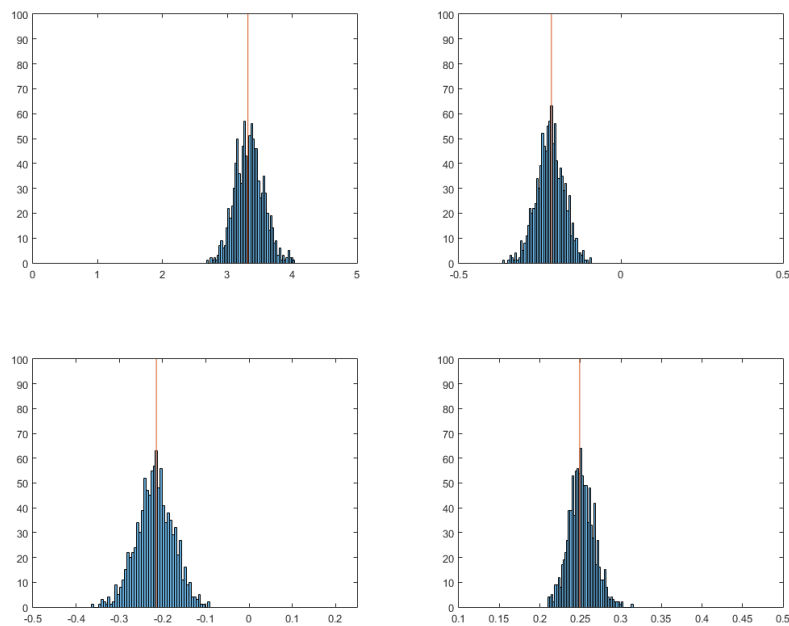
Ans.
- $N = 10$



- $N = 100$

- $N = 500$



## 3. Bayesian Inference for the Binomial

A discrete variable is given with two possible states. Suppose we draw this variable $N$ times, the outcomes of the $N$ trials are recorded as **O.mat**. Let $D = (m_1, m_2)$ denote the numbers of occurrences of two states from the draws. These draws can be represented by a binomial distribution $\text{Bin}(m|N, \mu)$ where $\mu$ denotes the

probability or parameter of the first state which satisfies $\mu \geq 0$. Please solve the following problems.

(a) Please apply the conjugate prior of $\mu$, which is a Beta distribution, $\text{Beta}(\mu|a,b) = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)}\mu^{a-1}(1-\mu)^{b-1}$, derive the posterior distribution $p(\mu|D, a, b)$, and show the derivation of MAP solution $\mu_{\text{MAP}}$ in details.

Ans.

$$P(\Lambda|X) \propto \left[ P(\Lambda) \prod_{n=1}^{N-1} P(X_n|\Lambda) \right] P(X_N|\Lambda) \quad \propto \quad P(\Lambda) \prod_{n=1}^{N} P(X_n|\Lambda)$$

$$= \mathcal{W}(\Lambda|W_o, \nu_o) \prod_{n=1}^{N} N(X_n|\mu, \Lambda^{-1})$$

$$\propto |\Lambda|^{(\nu_o-D-1)/2} \exp\left\{-\frac{1}{2}\text{Tr}(W_o^{-1}\Lambda)\right\} \cdot |\Lambda|^{N/2}\exp\left\{-\frac{1}{2}\sum_{n=1}^{N}(X_n-\mu)^T\Lambda(X_n-\mu)\right\}$$

$$= |\Lambda|^{(\nu_o+N)-D-1)/2}\exp\left\{-\frac{1}{2}\text{Tr}(W_o^{-1}\Lambda)\right\}\exp\left\{-\frac{1}{2}\text{Tr}(\Lambda\sum_{n=1}^{N}(X_n-\mu)(X_n-\mu)^T)\right\}$$

$$= |\Lambda|^{(\nu_\Lambda-D-1)/2}\exp\left\{-\frac{1}{2}\text{Tr}(W_\Lambda^{-1}\Lambda)\right\} = \mathcal{W}(\Lambda|W_\Lambda, \nu_\Lambda)$$

where $\nu_\Lambda = \nu_o + N$, $\quad W_\Lambda^{-1} = W_o^{-1} + \sum_{h=1}^{N}(X_n-\mu)(X_n-\mu)^T$

(b) **Programming**:

You can use Beta random variable for parameter [1]. Please use the recorded data **O.mat** and plot the prior and posterior distributions from 50 data samples and from the whole data samples. The parameters of the prior distribution are given as $a = b = 0.1$.

Ans.