

Homework 4, Part II: Multivariate Normal & Monte-Carlo Method

Submission Guidelines: Please briefly summarize your observation in a technical report (no more than 2 pages), compress your report and source code into one .zip file, and submit the compressed file via E3.

Problem 1 (Multivariate Normal for Regression)

(15+25=40 points)

One interesting application of multivariate normal (MVN) random variables is to solve regression tasks. In this problem, you will implement a simple MVN-based predictor that predicts the outputs of the testing queries based on the training data. Specifically, let $D_{train} = \{(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)\}$ be the training dataset and let $D_{test} = \{x_{N+1}, \dots, x_{N+M}\}$ be the testing queries. The goal is to predict $\{y_{N+1}, \dots, y_{N+M}\}$ that correspond to $\{x_{N+1}, \dots, x_{N+M}\}$.

(a) As a prep work, please show the following property that we discussed in class: Let Z_1, Z_2 be a pair of bivariate normal random variable with mean μ_1, μ_2 , variance σ_1^2, σ_2^2 , and correlation coefficient ρ . Show that conditioned on that $Z_1 = z_1$, the conditional distribution of Z_2 is normal with mean $\mu_2 + \frac{\rho\sigma_2(z_1 - \mu_1)}{\sigma_1}$ and variance $(1 - \rho^2)\sigma_2^2$.

(b) The property in (a) can be extended to the multivariate normal case. Suppose that for every $k \in \{N + 1, \dots, N + M\}$, $\{Y_1, \dots, Y_N, Y_k\}$ is multivariate normal with mean vector $\mu = [\mu_1, \dots, \mu_N, \mu_k]^\top$ and a $(N + 1) \times (N + 1)$ covariance matrix Σ , where the covariance between Y_i and Y_j (denoted by $\Sigma_{i,j}$) has the following form:

$$\Sigma_{i,j} = \sigma_f^2 \exp\left(-\frac{(x_i - x_j)^2}{2\ell^2}\right) + \sigma^2 \delta_{i,j}, \forall i, j,$$

where σ_f is a scale factor, ℓ is called the lengthscale, σ^2 is some positive constant (usually called the noise parameter), and $\delta_{i,j}$ is the delta function (i.e. $\delta_{i,j} = 1$ if $i = j$ and $\delta_{i,j} = 0$ if $i \neq j$). Given that $Y_1 = y_1, \dots, Y_N = y_N$, it can be shown that the conditional distribution of Y_k is normal with mean $K(x_k, x_{1:N})[K(x_{1:N}, x_{1:N}) + \sigma^2 I]^{-1} y_{1:N}$ and variance $K(x_k, x_k) - K(x_k, x_{1:N})[K(x_{1:N}, x_{1:N}) + \sigma^2 I]^{-1} K(x_{1:N}, x_k)$, where

- $K(x_k, x_k) = \Sigma_{k,k}$ is a scalar.
- $K(x_k, x_{1:N}) = [\Sigma_{k,1}, \dots, \Sigma_{k,N}]$ is a $1 \times N$ vector.
- $K(x_{1:N}, x_{1:N})$ is an $N \times N$ matrix with the (i, j) -th entry equal to $\Sigma_{i,j}$.
- I is an identity matrix of size $N \times N$.
- $K(x_{1:N}, x_k)$ is the transpose of $K(x_k, x_{1:N})$.
- $y_{1:N} = [y_1, \dots, y_N]^\top$ is an $N \times 1$ vector.

Based on the above conditional distribution, please write a program (e.g. in Python or MATLAB) to find the predictive distributions of the outputs of the test query points $\{x_{N+1}, \dots, x_{N+M}\}$. What is the prediction result of the testing dataset under $\sigma_f = 1, \sigma = 0.1, \ell = 0.5$? What is the prediction result of the testing dataset if ℓ is set to be 0.1 instead? How about $\ell = 2.0$?

Problem 2 (Monte-Carlo Method for Integration)

(15+15=30 points)

In probability research, one common task is to find the expected value of a function that depends on a random variable. If the random variable is continuous, then calculating such an expected value may involve a complex integral that has no simple closed-form expression. In this scenario, one work-around is to leverage the Monte-Carlo method to find an approximate answer.

(a) In this subproblem, you will be asked to tackle the following integral: Let $Z \sim \mathcal{N}(0, 1)$. Define another random variable $Y = \cos(Z) + \sin(2Z)$. Our goal is to find out the expected value of Y using Monte Carlo

method. Specifically, let z_1, z_2, \dots, z_n be n independent numbers drawn from a standard normal distribution. Then,

$$E[Y] = \int (\cos(z) + \sin(2z)) \frac{1}{\sqrt{2\pi}} e^{-\frac{z^2}{2}} dz \approx \frac{1}{n} \sum_{i=1}^n (\cos(z_i) + \sin(2z_i)).$$

Please write a short program (either in Python or MATLAB) to implement the above procedure. Suppose we set $n = 10^3$ and repeat the same estimation procedure for 20 times. What are the estimation results? What if we reconfigure $n = 10^5$ and again repeat the same estimation procedure for 20 times. Do you observe any differences in the estimation results?

(b) Define a closed region $A = \{(x, y) : (x - 0.2)^2 + 2(y + 0.3)^2 \leq 0.25\}$. Our goal is to approximately compute the area of the region A (denoted by Δ_A) using Monte Carlo integration. To begin with, we may write Δ_A as

$$\Delta_A = \int_A dx dy.$$

To apply the Monte-Carlo method, we consider two independent random variables $X \sim \text{Unif}(-1, 1)$ and $Y \sim \text{Unif}(-1, 1)$. Then, the joint PDF of X and Y (denoted by $f_{XY}(x, y)$) is $1/4$ for any $x \in [-1, 1], y \in [-1, 1]$ and is 0 elsewhere. Accordingly, we have

$$\Delta_A = 4 \cdot P((X, Y) \in A) = 4 \int_{-1}^1 \int_{-1}^1 \mathbb{I}_{\{(x, y) \in A\}} f_{XY}(x, y) dx dy \approx 4 \cdot \frac{1}{n} \sum_{i=1}^n \mathbb{I}_{\{(x_i, y_i) \in A\}},$$

where each (x_i, y_i) is an independent sample from the joint distribution f_{XY} and $\mathbb{I}_{\{\cdot\}}$ is the indicator function. Therefore, Δ_A can be approximated by counting the number of random points that fall in the region A . Please write a short program (either in Python or MATLAB) to implement the above procedure. What are your estimates of Δ_A under $n = 10^1, 10^3, 10^5$, and 10^7 ?