

2023 Spring 《数理统计》第 1 次作业

2023 年 5 月 5 日

问题：已知利用加噪的正弦函数 $y_i = 10 \sin(0.3x_i) + \epsilon_i$ 在区间 $[0, 20]$ 内生成了一组由 100 个数据对 (x_i, y_i) 构成的数据，其中 ϵ_i 表示数据噪声. 现在我们利用多项式回归模型

$$y_i = \mathbf{w}^\top \mathbf{x}_i + \epsilon_i, \quad \mathbf{x}_i = [1, x_i, \dots, x_i^m]^\top$$

来建模该组数据，其中 \mathbf{w} 表示模型参数， \mathbf{x}_i 表示模型的第 i 个输入， y_i 表示对应的观测值，而噪声 ϵ_i 服从高斯分布 $\mathcal{N}(\epsilon_i; 0, \sigma_0^2)$ ， m 表示多项式的次数. 若假设模型参数 \mathbf{w} 的先验分布为

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w}; \mathbf{0}, \sigma_w^2 \mathbf{I}),$$

试求模型参数 \mathbf{w} 的最大似然估计和贝叶斯估计.

数据下载地址：

github.com/Mephestopheles/Mathematical-Statistics-2023Spring/tree/main/Assignment1

训练数据由 data.csv 给出，其中数据形式如下表所示：

No	x	y
1	1.6302148253728963	6.241160657344439
2	0.8351081552813885	4.0151310636573365
\vdots	\vdots	\vdots
99	17.60037082165521	-9.361072876930848
100	18.97287441658695	-8.612303468512675

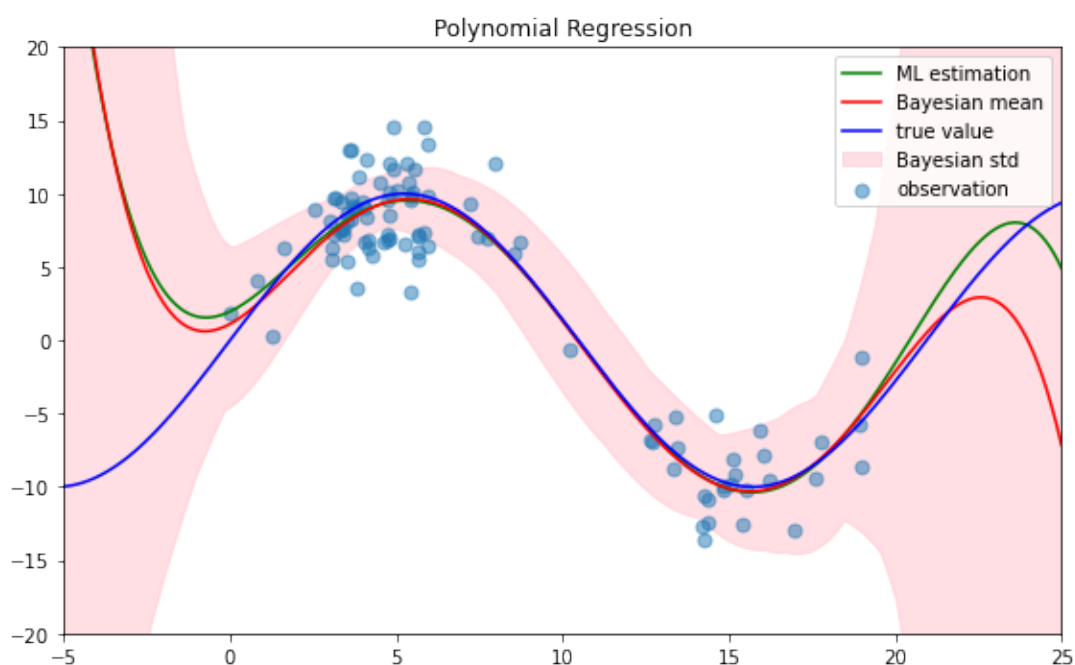
作业要求：

- 1) 分别给出模型参数 \mathbf{w} 的最大似然估计和贝叶斯估计的理论结果；

- 2) 实现最大似然估计和贝叶斯估计的代码，要求只调用 numpy 软件包中的基本运算以及必要的数据库读取工具（如 pandas 库）和绘图工具（如 matplotlib 库）；
- 3) 探索多项式次数 m 对最大似然估计和贝叶斯估计的影响；
- 4) 在贝叶斯估计的实验中，探索采用不同的 σ_0 和 σ_w 取值对实验结果的影响。

实验报告需包含（但不限于）：

- 1) 分别对单数据点和多数数据点的情形，给出模型参数 w 最大似然估计和贝叶斯估计的理论推导；
- 2) 分别为最大似然估计和贝叶斯估计绘制如下所示的曲线图（包含观测值、真实值和估计值）：



其中贝叶斯估计需要绘制出后验分布 $p(w|\mathbf{X}, \mathbf{y})$ 所表示的曲线簇（可以根据后验分布先采样得到若干个模型参数，从而进行绘制）；

- 3) 通过观察实验结果，结合理论知识，说明最大似然估计与贝叶斯估计之间的关系；
- 4) 实验结果及讨论，包括多项式次数、超参数取值和模型参数的先验分布等。