Project Report of Machine Learning (Assignment 4)

Author: 赵文浩 23020201153860 (计算机科学系)

■ Task Description

试推导 Bayesian Linear Regression 中 y_{new} 的概率分布。

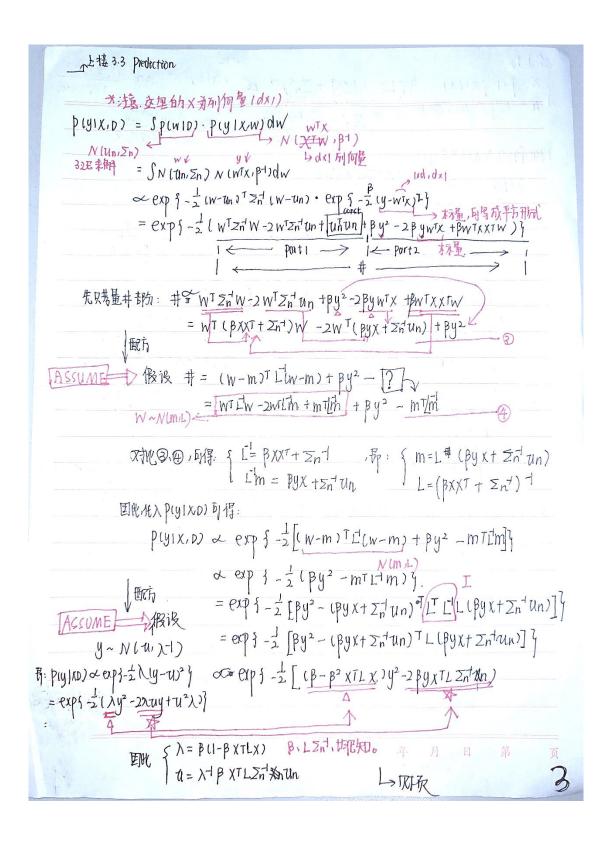
P.S. 本次任务涉及较多符号计算,因此这里采用纸质的形式展示推导过程。

```
23020201153860 赵女浩
                                                                                                                                                                                                                                          计算机科学系
                    Assignment 4 — Boyesian Linear Regression
                   I. Peview of Bayesian 特更Prior
                                                                                                                                                             献船函数估计 likelihood
                      O R叶斯红基种的: P(AB) = P(A). P(B)A)
                      O Bayesian Method 两物可聚 后的 posterior
                             Step D Inference — 推理求命验: posterior (w)
                             Step ② Prediction — 根据Inference 弥婀: スキーy*
               I. Linear Regression
                             这里对线性阿归进行足义说明,以便推导。
                            结定数据: Data → D= f2X1,y1>,2X2,y2>...2X1,y1>7 i=1.n. Xitpd, y1 tR
  对假设W也用助人 Soursian 统
                                                                                    据《ite Baffet 偶然性(噪音)
1 P(w)= N(0, 21)
                                                                                  易知 yi~N(WTXi, Bit)
       II Linear Regression with Bayesian
                   3.1 Interence — Likelihood
                             Inference的目的是求多数 w的后驱与布 PCWID)
                        Therefore the state of the st
```

```
类似地, Prior: P(W)= N(D, ot I)
                                                                                x exp g-dwTwg
           3.2 Interence - Posterior
           毎面 posterior: p(W|D) ベ p(y|W,X).P(W)

    exp (- B (y-xw)) (y-xw) 3. exp (- 2 wTw).

                                                                                              = ex p s - \frac{\beta}{2} (y - xw)^{T} (y - xw) - \frac{\alpha}{2} w Tw 
                                                                                              = exp = \frac{B}{2} (y^{T}y - 2y^{T}xw + wx^{T}xw) - \frac{A}{2} w^{T}w^{2}
                                          = exp s = B w x x x w - 2 w w + B y x x w - B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B y x x w + B x x x w + B x x x w + B x x x w + B x x x w + B x x x w + B x x x w + B x x x x w + B x x x x w + B x x x x w + B x x x x w + B x x x x w + B x x x x x w + B x x x x x w + B x x x x x x w
                                                                    P(W|D) = N(w|u_p, \Sigma_p)
                                           展刊 Sp(WID) = N(Un, En)
   = β (β X T X + Δ I) T X T Y 1 Pidge 回日 (Em)的特果相同
3.3 Prediction
            预测 Prediction的目的是从新数据样本 X*, 结结结定的数据集D, 预测Y*, 即预测矫:
             P(y|X,D) = Sp(y,w|X,D)dw w为何是 * w为文的的无关; mew = Sp(w|X,D)p(y|X,w,D)dw * yen D无关
                     =Sp(w|D), p(y|X,w)dW wTX N(知, p1) 予例何是(dxi)
```



1上接 $\begin{cases} \lambda = \beta (1 - \beta X T L X) & \sharp P L = (\beta X X T + \sum_{n=1}^{N-1})^{-1} & (AtuvT)^{-1} = A^{-1} - \underbrace{A^{\frac{1}{2}} uvTA^{\frac{1}{2}}}_{1 + vTA^{\frac{1}{2}} uv} = \sum_{n=1}^{N-1} \frac{A^{\frac{1}{2}} uvTA^{\frac{1}{2}}}{1 + \chi T \sum_{n} \beta X} = \sum_{n=1}^{N-1} \frac{\beta Z_{n} X X T Z_{n}}{1 + \beta X T Z_{n} X}$ 因此将让代入 XTLX中间得。 国此将L代入XTLX中可得: XTLX = X TENX - BX TENXXTENX = XIZnX + BXIZnX: XIZnX - BXIZnXXIZnX $\begin{array}{ll}
\lambda = \beta (1 - \beta XTLX) & u = \lambda^{-1} \beta XTLZn^{-1}Un \\
= \beta (1 - \beta XTZnX) & = (\beta^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{-1}t^{$

4