

回归分析-第一次上机

南开大学统计与数据科学学院，马东升

2024 年 10 月 29 日

1 第一问

1.1 题目

检查复共线性

1.2 源代码 (R)

```
1 #导入数据
2 library(readxl)
3 data <- read_excel("longley.xlsx")
4 head(data)
5
6 # 提取前六列
7 X <- data[, 1:6]
8
9 # 中心化和标准化
10 X_scaled <- scale(X)
11
12 # 计算  $X'X$ 
13 XTX <- t(X_scaled) %*% X_scaled
14
15 # 计算特征值
16 eig_values <- eigen(XTX)$values
17
18 # 按从大到小排序
19 sorted_eig_values <- sort(eig_values, decreasing = TRUE)
20
21 #计算条件数
22 k=sorted_eig_values[1]/sorted_eig_values[6]
23
24 k
```

1.3 统计分析结论

根据 <http://stat.ethz.ch/R-manual/R-devel/library/datasets/html/longley.html> 的数据说明，我们以前六列为自变量，第七列为因变量。

1.2 中代码输出的结果为 $X^T X$ 的条件数，即

$$k = \frac{\lambda_{max}}{\lambda_{min}} = 12220.0 > 1000$$

这说明了数据的复共线性非常严重。

2 第二问

2.1 题目

使用主成分回归解决复共线性，选择适当个数的主成分

2.2 源代码 (R)

第一段

```
1 # 提取前六列
2 X <- data[, 1:6]
3
4 # 对X进行PCA变换
5 pca_result <- prcomp(X, center = TRUE, scale. = TRUE)
6
7 # 查看PCA的结果，包括主成分得分和主成分载荷
8 summary(pca_result)
```

第二段

```
1 Y <- as.numeric(data$Employed)
2
3 # 主成分得分 (PCA变换后的数据)
4 pca_scores <- pca_result$x[, 1:2]
5
6 # 对Y进行回归分析，使用前两个主成分
7 regression_model <- lm(Y ~ pca_scores)
8
9 # 提取回归系数
10 coefficients <- regression_model$coefficients
11
12 # 恢复到原始变量空间
13 # 通过PCA载荷矩阵 (主成分载荷矩阵) 来恢复系数
14 pca_loadings <- pca_result$rotation[, 1:2]
15
16 # 计算标准化后的系数 beta
17 beta_standardized <- pca_loadings %*% coefficients[2:3]
18
19 # 提取原始数据的均值和标准差
20 X_means <- colMeans(X)
21 X_sds <- apply(X, 2, sd)
22
23 # 将系数恢复到原始变量空间
24 beta_original <- beta_standardized / X_sds
25
26 # 恢复截距项
27 original_intercept <- coefficients[1] - sum((X_means / X_sds) * beta_standardized)
28
29 # 输出恢复后的系数
30 list(
```

```

31 original_intercept = original_intercept ,
32 original_coefficients = beta_original
33 )

```

2.3 统计分析结论

第一段代码，我们得出两个主成分的累计变差贡献率：

$$\frac{\lambda_1 + \lambda_2}{\sum_{i=1}^6 \lambda_i} = 0.9631 > 0.85$$

后三者的特征根都比较小（小于 1），所以我们直接选择两个主成分。

第二段代码，我们用主成分做回归后，再恢复到原始系数，最终得到如下的经验回归方程：

$$\begin{aligned} \hat{Employed} = & -258.6257 + 0.0691GNP.deflator + 0.0075GNP \\ & + 0.0029Unemployed + 0.0090Armed.Forces \\ & + 0.1014Population + 0.1529Year \end{aligned}$$

我们发现个别自变量的系数已经比较小了，这也符合我们做 PCA 的目的。

3 第三问

3.1 题目

使用岭回归解决复共线性，并采用不同方法估计岭参数

3.2 源代码（R）

第一段

```

1 library(glmnet)
2 # 先定义 ridge_regression 函数
3 ridge_regression <- function(X, Y, k) {
4   # 确保 Y 是数值型向量
5   Y <- as.numeric(Y)
6
7   # 对 X 进行中心化和标准化
8   X_scaled <- scale(X)
9
10  # 执行岭回归，alpha=0 表示岭回归（L2 正则化）
11  ridge_model <- glmnet(X_scaled, Y, alpha = 0, lambda = k, standardize = FALSE)
12
13  # 提取标准化后的系数（包含截距）
14  beta_standardized <- coef(ridge_model)
15
16  # 将系数转换为数值向量
17  beta_standardized <- as.vector(beta_standardized)
18
19  # 计算预测值
20  Y_pred <- cbind(1, X_scaled) %*% beta_standardized # 计算预测值

```

```

21
22 # 计算残差平方和 (RSS)
23 RSS <- sum((Y - Y_pred) ^ 2)
24
25 # 恢复到原始变量空间的系数
26 X_means <- attr(X_scaled, "scaled:center")
27 X_sds <- attr(X_scaled, "scaled:scale")
28
29 beta_original <- beta_standardized[-1] / X_sds # 去掉截距项并进行恢复
30 intercept_original <- beta_standardized[1] - sum(X_means * beta_original)
31
32 # 返回结果
33 return(list(
34   beta_standardized = beta_standardized,
35   beta_original = c(intercept_original, beta_original),
36   RSS = RSS # 返回 RSS 值
37 ))
38 }

```

第二段

```

1 X <- data[, 1:6]
2
3 # 确保Y是数值型向量
4 Y <- as.numeric(data$Employed)
5
6 # 标准化 X
7 X_scaled <- scale(X)
8
9 # 对 X'X 进行特征值分解 (spectral decomposition)
10 XtX <- t(X_scaled) %*% X_scaled
11 eigen_decomp <- eigen(XtX)
12
13 # 提取特征向量矩阵 Phi 和特征值矩阵 Lambda
14 Phi <- eigen_decomp$vectors
15 Lambda <- diag(eigen_decomp$values)
16
17 # 计算新变量矩阵 Z = X * Phi
18 Z <- X_scaled %*% Phi
19
20 # 进行典则形式的回归 y = alpha_0 + Z * alpha + e
21 canonical_model <- lm(Y ~ ., data = as.data.frame(Z))
22
23 # 提取回归系数 alpha (去掉截距项)
24 alpha <- coef(canonical_model)[-1]
25
26 # 计算残差的方差 sigma^2
27 residuals <- residuals(canonical_model)
28 sigma_squared <- var(residuals)
29
30 # 选取 alpha 的最大值
31 max_alpha2 <- max(alpha**2)
32

```

```

33 # 根据 H-K 公式计算岭回归系数 k
34 k <- sigma_squared / max_alpha2
35 k

    第三段

1
2 # 定义岭迹图绘制函数
3 ridge_trace <- function(X, Y, k_values) {
4   # 创建一个矩阵以存储标准化系数
5   beta_matrix <- matrix(0, nrow = length(k_values), ncol = ncol(X) + 1)
6
7   # 计算每个 k 值下的标准化系数
8   for (i in seq_along(k_values)) {
9     k <- k_values[i]
10    ridge_result <- ridge_regression(X, Y, k)
11
12    # 将标准化系数存入矩阵
13    beta_matrix[i, ] <- ridge_result$beta_standardized
14  }
15
16  # 将 k_values 转换为矩阵 (确保是列矩阵)
17  k_values_matrix <- matrix(k_values, ncol = 1)
18
19  # 绘制岭迹图
20  matplot(k_values_matrix, beta_matrix[, -1], type = "l", lty = 1, col = 1:ncol(X),
21          xlab = "岭回归系数 k", ylab = "标准化系数",
22          main = "岭迹图", ylim = range(beta_matrix[, -1]), lwd=2)
23  legend("topright", legend = colnames(X), col = 1:ncol(X), lty = 1)
24  abline(h = 0, col = "gray", lwd = 2)
25 }
26 k_values <- seq(0.01, 10, length.out = 10000)
27 ridge_trace(X = data[, 1:6], Y = data$Employed, k_values = k_values)
28 k_values <- seq(1, 4, length.out = 3000)
29 ridge_trace(X = data[, 1:6], Y = data$Employed, k_values = k_values)

```

第四段

```

1 ridge_regression(X = data[, 1:6], Y = data$Employed, k = 0)
2 ridge_regression(X = data[, 1:6], Y = data$Employed, k = 0.00115)
3 ridge_regression(X = data[, 1:6], Y = data$Employed, k = 3)

```

3.3 统计分析结论

第一段代码是岭回归的代码，第四段代码是最终做的岭回归；第二段代码为用 Horel-Kennard 公式确定 k，第三段代码为用岭迹图确定 k。

3.3.1 用 Horel-Kennard 公式确定 k

Horel-Kennard 公式确定的 k 值为：

$$\hat{k} = \frac{\hat{\sigma}^2}{\max_i \hat{\alpha}_i^2}$$

计算得到 $\hat{k} = 0.00115$ ，这个值略小，这是因为 $\hat{\sigma}^2 = 0.0558$ 比较小，而且 $\max_i \hat{\alpha}_i^2 = 48.4885$ 也比较大。最终得到的经验回归方程为：

$$\begin{aligned}\hat{Employed} = & -2429.5210 - 0.0137GNP.deflator - 0.0027GNP \\ & - 0.0153Unemployed - 0.0088Armed.Forces \\ & - 0.1608Population + 1.2910Year\end{aligned}$$

虽然这是由 Horel-Kennard 公式确定的 k ，但是很明显，该回归方程不符合经济学意义， $Employed$ 至少应该和 GNP 及 $GNP.deflator$ 正相关（经济越好或越热，就业越好）。所以我个人认为，这个 k 不是一个特别好的 k 值。

3.3.2 用岭迹图确定 k

我们用第三段代码来画出岭迹图，注意这里是标准化后的系数。

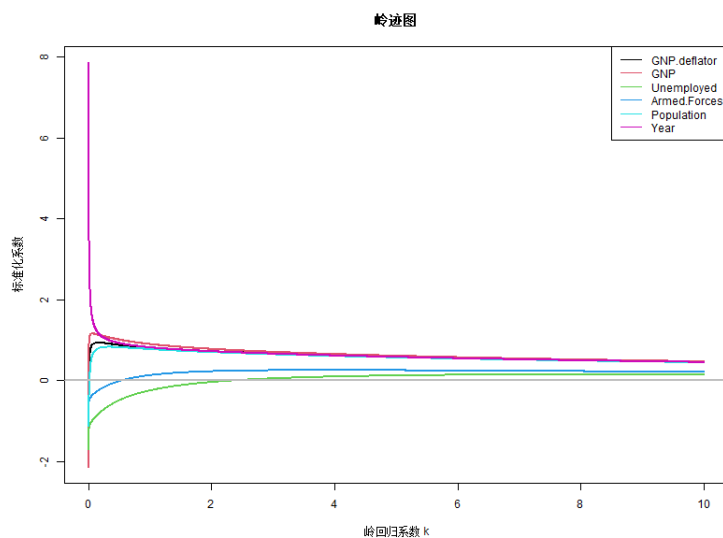


图 1: $k \in [0, 10]$ 的岭迹图

从图 1 可以看出 $k \in [0, 1]$ 这个区间确实不是一个岭回归的好区间，原因是系数尚不稳定，所以我们缩小区间到 $k \in [1, 4]$ ，得到图 2。

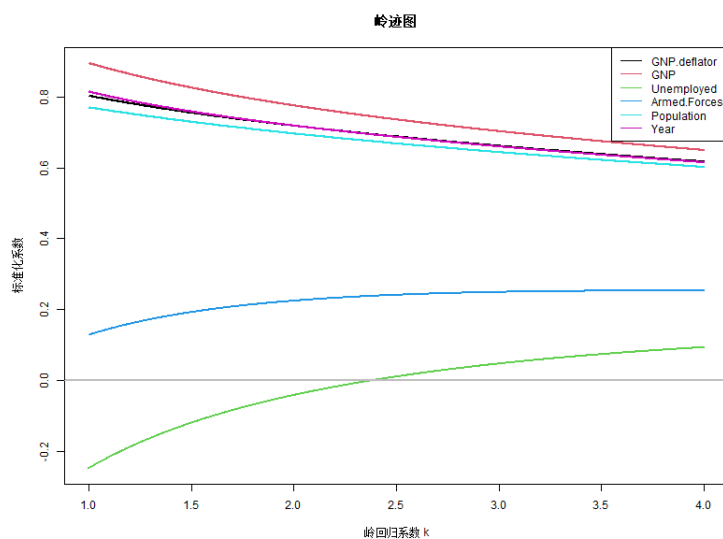


图 2: $k \in [1, 4]$ 的岭迹图

为了使系数符号稳定，虽然会使得 RSS 较大，但我们也不得不选择较大的 k ，比如 $k = 3$ 。
 $k = 3$ 时，我们得到经验回归方程：

$$\begin{aligned}\hat{Employed} = & -226.5122 + 0.0613GNP.deflator + 0.0070GNP \\ & + 0.0005Unemployed + 0.0035Armed.Forces \\ & + 0.0925Population + 0.1386Year\end{aligned}$$

我们发现这一方程和 PCA 得到的方程很相似，虽然 $RSS_{k=3} = 15.1309$ （相较于 $RSS_{k=0} = 0.8542$ ）确实较大，但是至少系数的经济意义是相对正确的。所以我个人认为， $k = 3$ 是一个较好的值。