# STA409

## Answer to Assignment 3

12112627

李乐平

## 1. Solution.

(1).



(2). The family-wise error rate (FWER) is computed by

$$FWER = 1 - (1 - 0.05)^{42} = 0.8840$$

That means we only have a chance of 11.6% that we do not make any Type I Error without any adjustment.

(3).

|  | Diabetes | Asthma | Cancer | Multiple Sclerosis | Thyroid Disease | Liver Disease | Arthritis |
|---|---|---|---|---|---|---|---|
| Leo | 0.7953 | 0.8187 | 0.9528 | 0.1092 | 0.4827 | 0.1701 | 0.9060 |
| Purple | 0.6143 | 0.2577 | 0.4849 | 0.2487 | 0.8314 | 0.7887 | 0.1878 |
| 0:00-1:00 | 0.0628 | 0.4728 | 0.1874 | 0.0088 | 0.8947 | 0.0103 | 0.5419 |
| Red Hair | 0.2353 | 0.5705 | 0.7993 | 0.2556 | 0.3873 | 0.3861 | 0.7048 |
| First Name C | 0.4391 | 0.3479 | 0.3406 | 0.7265 | 0.2134 | 0.6235 | 0.8320 |
| Summer | 0.6329 | 0.3941 | 0.8162 | 0.3259 | 0.9229 | 0.3629 | 0.3189 |

In the table above, 2 combinations of disease and risk factor are determined to be dependent, which are actually Type I Errors. That implies it makes sense to do adjustment to reduce the chance of Type I Error happens under the increment of number of experiments.

## 2. Solution.

(1). The ANOVA table and coefficients table are given by

| ANOVA Table | | | | | |
|---|---|---|---|---|---|
| Source | DF | Sum of Squares | Mean Square | F-value | p-value |
| Model | 4 | 23665352 | 5916338 | 22.9782 | $5.0715 \times 10^{-13}$ |

| | | | | |
|---|---|---|---|---|
| Error | 88 | 22657938 | 257476.5682 | | |

where MS = sum_of_squares / DF; F-value = $MS_{Model}$ / $MS_{Error}$; p-value = $Pr(F(DF_{Model}, DF_{Error}) > F)$.

| Coefficients Table | | | | |
|---|---|---|---|---|
| Parameter | Estimate | Standard Error | t-value | p-value |
| Intercept | 3526.4 | 327.7 | 10.7610 | 0 |
| Gender | 722.5 | 117.8 | 6.1333 | $2.3917 \times 10^{-8}$ |
| Education | 90.02 | 24.69 | 3.6460 | $4.5033 \times 10^{-4}$ |
| Experience | 1.2690 | 0.5877 | 2.1593 | $3.3547 \times 10^{-2}$ |
| Months | 23.406 | 5.201 | 4.5003 | $2.0765 \times 10^{-5}$ |

where t-value = Estimate / SE; p-value = $Pr(|t(DF_{Error})| > |t|)$

In this question, the model for test is given by
$$Salary = \beta_0 + \beta_1 Gender + \beta_2 Education + \beta_3 Experience + \beta_4 Months + \varepsilon$$
The overall F-test is rejected, indicating at least one of $\beta$s are significantly different from 0. And t-tests show that every predictor variable is significant.

(2). The $R^2$ and adjusted $R^2$ are computed by
$$R^2 = 1 - \frac{SSE}{SST} = 1 - \frac{23665652}{46323290} = 0.5109$$

$$Adjusted - R^2 = 1 - \frac{SSE/(n-p-1)}{SST/(n-1)} = 1 - \frac{23665652/88}{46323290/92} = 0.4886$$

(3). A partial F-test is utilized to test which of the full model and the reduced model is better.
$$F = \frac{(SST_R - SST_F)/(p-k)}{SST_F/(n-p-1)} = \frac{(38460756 - 22657938)/3}{22657938/88} = 20.4586$$

Because $Pr(F(3, 88) > F) = 3.8052 \times 10^{-10}$, hence the reduced model is rejected.

## 3. Solution.

Now we want to show that $\hat{\boldsymbol{\beta}}^{WLS} = (X^T W X)^{-1} X^T W \boldsymbol{y}$ is the best linear unbiased estimate of $\boldsymbol{\beta}$, where $W = [Var(\boldsymbol{y})]^{-1}$.

The variance of $\hat{\boldsymbol{\beta}}^{WLS}$ is
$$Var(\hat{\boldsymbol{\beta}}^{WLS}) = ((X^T W X)^{-1} X^T W) Var(\boldsymbol{y})((X^T W X)^{-1} X^T W)^T$$
$$= (X^T W X)^{-1}$$

Assume there is another linear estimation of $\boldsymbol{\beta}$: $\widetilde{\boldsymbol{\beta}} = A\boldsymbol{y}, A = (X^T W X)^{-1} X^T W + B$, then it must satisfies
$$E(\widetilde{\boldsymbol{\beta}}) = E[((X^T W X)^{-1} X^T W + B)(X\boldsymbol{\beta} + \varepsilon)] = (I + BX)\boldsymbol{\beta} = \boldsymbol{\beta},$$
i.e. $BX = 0$, or it would not be unbiased estimator of $\boldsymbol{\beta}$.

Yet
$$Var(\widetilde{\boldsymbol{\beta}}) = ((X^T W X)^{-1} X^T W + B) W^{-1} ((X^T W X)^{-1} X^T W + B)^T$$
$$= (X^T W X)^{-1} + B W^{-1} B^T$$
$$= Var(\hat{\boldsymbol{\beta}}^{WLS}) + B W^{-1} B^T$$
$$\geq Var(\hat{\boldsymbol{\beta}}^{WLS})$$

then it is clear that $\hat{\boldsymbol{\beta}}^{WLS}$ is the best linear unbiased estimate of $\boldsymbol{\beta}$.
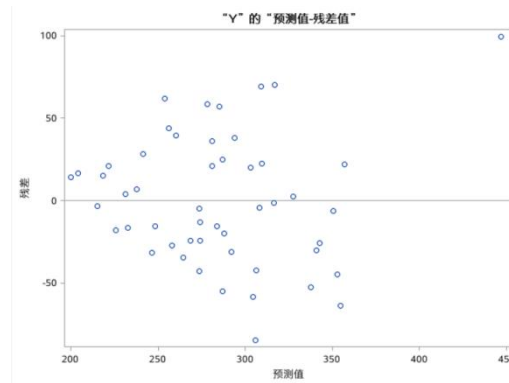
## 4. Solution.

(1). The fitted linear model is stated by

$$Y = -556.5680 + 0.0724X_1 + 1.5521X_2 - 0.0043X_3 + \varepsilon$$

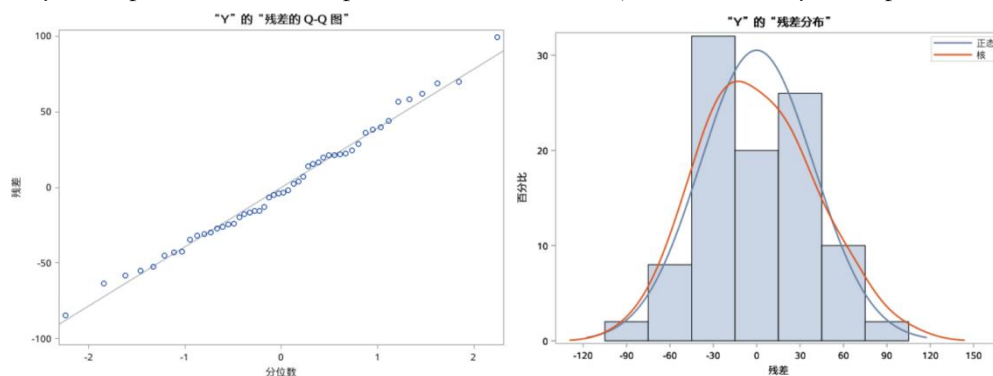And the assumptions are checked as follows.

Linearity assumption: The predicted value - residual plot shows that the mean of the residuals are roughly 0 at any place, which implies the validity of linearity assumption.



Homoscedasticity assumption: The predicted value - residual plot does not show homoscedaticity since it is narrow at left side and wide at right side. Also, the White test and Breusch-Pagan test reject the homoscedasticity assumption with p-values less than 0.01.

| 异方差性检验 | | | | | |
|---|---|---|---|---|---|
| 方程 | 检验 | 统计量 | 自由度 | Pr > 卡方 | 变量 |
| Y | White 检验 | 22.68 | 9 | 0.0070 | 所有变量的叉积 |
| | Breusch-Pagan | 15.59 | 3 | 0.0014 | 1, X1, X2, X3 |

Normality: The Q-Q plot as well as distribution of residuals does not show violation of the normality assumption. Also, the Shapiro-Wilk test does not reject the normality assumption.



| 正态性检验 | | | | |
|---|---|---|---|---|
| 检验 | | 统计量 | p 值 | |
| Shapiro-Wilk | W | 0.962619 | Pr < W | 0.1145 |
| Kolmogorov-Smirnov | D | 0.077587 | Pr > D | >0.1500 |
| Cramer-von Mises | W-Sq | 0.0368 | Pr > W-Sq | >0.2500 |
| Anderson-Darling | A-Sq | 0.295345 | Pr > A-Sq | >0.2500 |

Independence: The data are collected from 50 states independently, hence we can assume the independence of the data.

(2). Notice that the state Alaska has both high leverage, Cook's distance and studentized residual, which shows Alaska is an outlier and influential observation. Also shown in the plot.

| | State | leverage ▼ | cookD | rstudent |
|---|---|---|---|---|
| 1 | AK | 0.4419099163 | 2.1327949114 | 3.7099223829 |
| 2 | UT | 0.2941338295 | 0.0002247544 | -0.04594205 |
| 3 | NM | 0.1829094976 | 0.054304552 | 0.9847410811 |
| 4 | VT | 0.1400044576 | 0.0235191903 | 0.7566427892 |
| 5 | WV | 0.1334181929 | 0.0053186159 | 0.3682198446 |
| 6 | FL | 0.1329380512 | 0.0121187746 | 0.5580649664 |
| 7 | MS | 0.121715957 | 0.024280392 | -0.834377609 |
| 8 | LA | 0.1150154266 | 0.0134134888 | -0.638374691 |
| 9 | ND | 0.1046879429 | 0.0682663128 | -1.551362037 |
| 10 | CT | 0.101921737 | 0.01292351 | -0.670860476 |
| 11 | RI | 0.1011891856 | 0.0370314221 | 1.1510950224 |
| 12 | NJ | 0.088828501 | 0.0452504612 | -1.375748431 |
| 13 | TX | 0.0872943074 | 0.0003731912 | -0.123585963 |
| 14 | NY | 0.0864580927 | 0.077778339 | 1.8610093979 |
| 15 | CA | 0.0861878909 | 0.0078594908 | 0.57311144 |



Considering Alaska is away from the mainland of U.S., it may be dropped from the analysis.

(3). If we only refit the model using weighted GLM without dropping outlier to solve the heteroscedasticity, the refitted model is stated as

$$Y = -423.0778 + 0.0649X_1 + 1.1954X_2 + 0.0223X_3 + \varepsilon$$

| 参数 | 估计 | 标准误差 | t 值 | Pr > |t| |
|---|---|---|---|---|
| 截距 | -423.0777720 | 116.1533088 | -3.64 | 0.0007 |
| X1 | 0.0648833 | 0.0114096 | 5.69 | <.0001 |
| X2 | 1.1953871 | 0.2978225 | 4.01 | 0.0002 |
| X3 | 0.0222879 | 0.0468463 | 0.48 | 0.6365 |

However, after dropping the data of Alaska, the fitted linear model is stated as

$$Y = -277.5773 + 0.0483X_1 + 0.8869X_2 + 0.0668X_3 + \varepsilon$$

This time the homoscedasticity assumption is accepted. And the other 3 assumptions are also valid (omitted).

| | | | 参数估计 | | | | |
|---|---|---|---|---|---|---|---|
| 变量 | 标签 | 自由度 | 参数估计 | 标准误差 | t 值 | Pr > |t| | 方差膨胀 |
| Intercept | Intercept | 1 | -277.57731 | 132.42286 | -2.10 | 0.0417 | 0 |
| X1 | X1 | 1 | 0.04829 | 0.01215 | 3.98 | 0.0003 | 2.23892 |
| X2 | X2 | 1 | 0.88693 | 0.33114 | 2.68 | 0.0103 | 1.26399 |
| X3 | X3 | 1 | 0.06679 | 0.04934 | 1.35 | 0.1826 | 1.89792 |

"Y" 的 "预测值-残差值"

| 异方差性检验 | | | | | |
|---|---|---|---|---|---|
| 方程 | 检验 | 统计量 | 自由度 | Pr > 卡方 | 变量 |
| Y | White 检验 | 8.69 | 9 | 0.4665 | 所有变量的叉积 |
| | Breusch-Pagan | 5.03 | 3 | 0.1698 | 1, X1, X2, X3 |

Compared with the former fitted model,

$$Y = -556.5680 + 0.0724X_1 + 1.5521X_2 - 0.0043X_3 + \varepsilon$$

all the assumptions are accepted, and the fitted coefficients differs significantly, which indicates the Alaska row is indeed an outlier.
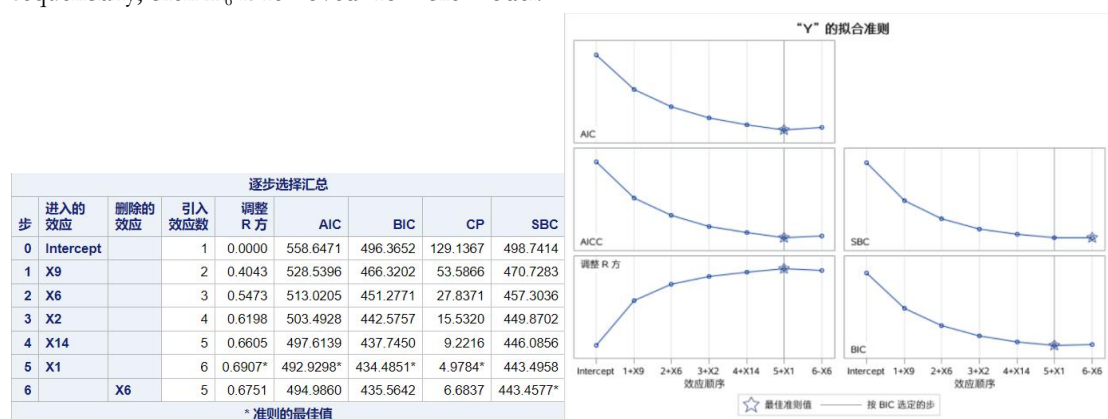
## 5. Solution.

(1). By the ranked table of pairwise Pearson correlations, we can observe that $X_{12}$ and $X_{13}$ are highly correlated, whose correlation coefficient is 0.98384.



Pearson 相关系数, N = 60
Prob > |r|, H0: Rho=0

(2). By checking the tolerances and VIFs of $X_{12}$ and $X_{13}$, with the criterion Tolerance < 0.1 and VIF > 10, multicollinarity does exist between $X_{12}$ and $X_{13}$.

| 参数估计 | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| 变量 | 标签 | 自由度 | 参数估计 | 标准误差 | t 值 | Pr > \|t\| | 容差 | 方差膨胀 |
| Intercept | Intercept | 1 | 1763.99793 | 437.33031 | 4.03 | 0.0002 | . | 0 |
| X1 | X1 | 1 | 1.90536 | 0.92374 | 2.06 | 0.0451 | 0.24308 | 4.11389 |
| X2 | X2 | 1 | -1.93762 | 1.10839 | -1.75 | 0.0874 | 0.16277 | 6.14355 |
| X3 | X3 | 1 | -3.10040 | 1.90167 | -1.63 | 0.1102 | 0.25203 | 3.96777 |
| X4 | X4 | 1 | -9.06517 | 8.48622 | -1.07 | 0.2912 | 0.13387 | 7.47004 |
| X5 | X5 | 1 | -106.83103 | 69.78007 | -1.53 | 0.1329 | 0.23215 | 4.30762 |
| X6 | X6 | 1 | -17.15689 | 11.86012 | -1.45 | 0.1551 | 0.20574 | 4.86054 |
| X7 | X7 | 1 | -0.65111 | 1.76777 | -0.37 | 0.7144 | 0.25033 | 3.99478 |
| X8 | X8 | 1 | 0.00360 | 0.00403 | 0.89 | 0.3761 | 0.60303 | 1.65828 |
| X9 | X9 | 1 | 4.45958 | 1.32721 | 3.36 | 0.0016 | 0.14750 | 6.77960 |
| X10 | X10 | 1 | -0.18715 | 1.66169 | -0.11 | 0.9108 | 0.35192 | 2.84158 |
| X11 | X11 | 1 | -0.16741 | 3.22730 | -0.05 | 0.9589 | 0.11472 | 8.71707 |
| X12 | X12 | 1 | -0.67216 | 0.49102 | -1.37 | 0.1780 | 0.01014 | 98.63993 |
| X13 | X13 | 1 | 1.34010 | 1.00559 | 1.33 | 0.1895 | 0.00953 | 104.98240 |
| X14 | X14 | 1 | 0.08626 | 0.14752 | 0.58 | 0.5617 | 0.23647 | 4.22893 |
| X15 | X15 | 1 | 0.10674 | 1.16943 | 0.09 | 0.9277 | 0.52436 | 1.90709 |

(3). The process of model selection is shown as follows. $X_9$, $X_6$, $X_2$, $X_{14}$ and $X_1$ are added to model sequentially, then $X_6$ is removed from the model.



| 逐步选择汇总 | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| 步 | 进入的效应 | 删除的效应 | 引入效应数 | 调整 R 方 | AIC | BIC | CP | SBC |
| 0 | Intercept | | 1 | 0.0000 | 558.6471 | 496.3652 | 129.1367 | 498.7414 |
| 1 | X9 | | 2 | 0.4043 | 528.5396 | 466.3202 | 53.5866 | 470.7283 |
| 2 | X6 | | 3 | 0.5473 | 513.0205 | 451.2771 | 27.8371 | 457.3036 |
| 3 | X2 | | 4 | 0.6198 | 503.4928 | 442.5757 | 15.5320 | 449.8702 |
| 4 | X14 | | 5 | 0.6605 | 497.6139 | 437.7450 | 9.2216 | 446.0856 |
| 5 | X1 | | 6 | 0.6907* | 492.9298* | 434.4851* | 4.9784* | 443.4958 |
| 6 | | X6 | 5 | 0.6751 | 494.9860 | 435.5642 | 6.6837 | 443.4577* |
| * 准则的最佳值 | | | | | | | | |

The final model is stated as

$$Y = 857.43 + 2.06\ X_1\ -1.77\ X_2 + 4.08\ X_9 + 0.33\ X_{14} + \ \varepsilon$$

| 参数估计 | | | | |
|---|---|---|---|---|
| 参数 | 自由度 | 估计 | 标准误差 | t 值 |
| Intercept | 1 | 857.431124 | 26.234966 | 32.68 |
| X1 | 1 | 2.059246 | 0.524393 | 3.93 |
| X2 | 1 | -1.771640 | 0.527540 | -3.36 |
| X9 | 1 | 4.078669 | 0.671468 | 6.07 |
| X14 | 1 | 0.330551 | 0.077299 | 4.28 |

This linear model describes the relationship between the total age-adjusted mortality rate (per 100,000) and five explanatory variables. The intercept term of 857.43 indicates that when all explanatory variables are zero, the expected total age-adjusted mortality rate is 857.43. For the coefficients of each explanatory variable, the interpretations are as follows: The coefficient of $X_1$ is 2.06, meaning that when the mean annual precipitation increases by one unit (in inches), the total age-adjusted mortality rate is expected to increase by 2.06 units. The coefficient of $X_2$ is -1.77, indicating that when the mean January temperature increases by one unit (in degrees Fahrenheit), the

total age-adjusted mortality rate is expected to decrease by 1.77 units. The coefficient of $X_9$ is 4.08, suggesting that when the percent of nonwhite population increases by one unit, the total age-adjusted mortality rate is expected to increase by 4.08 units. Lastly, the coefficient of $X_{14}$ is 0.33, signifying that for each unit increase in the relative pollution potential of sulfur dioxide, the total age-adjusted mortality rate is expected to increase by 0.33 units. ε represents the error term, which accounts for the unexplained portion of the model.
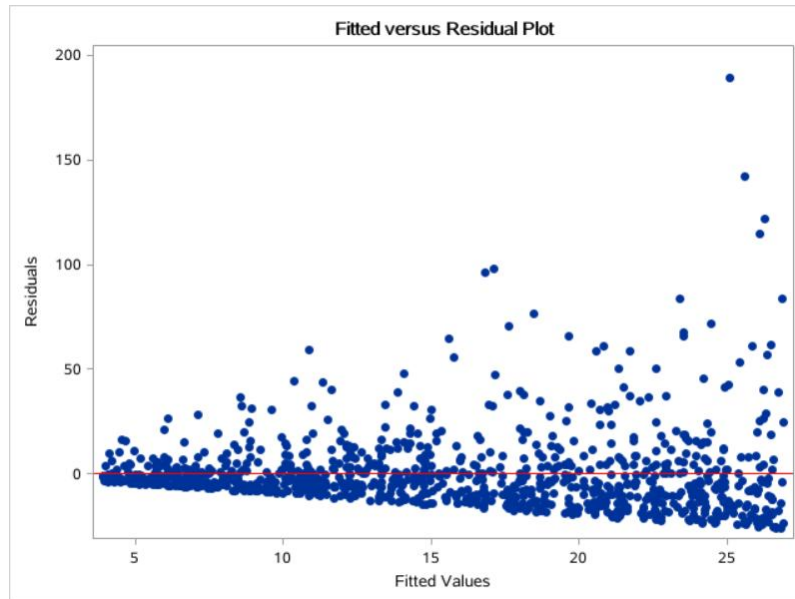
(4). The lines showing the parameter estimates of $X_1$, $X_2$, $X_6$, $X_9$, $X_{12}$, $X_{13}$ and $X_{14}$ against $\lambda$ is plotted as follows.



From the plot, it's evident that as $\lambda$ increases slightly, the coefficients of $X_{12}$ and $X_{13}$ experience rapid absolute value decreases, which aligns with their high sample correlation coefficient of 0.98 as noted in question 5.(1). These coefficients quickly converge towards zero, exhibiting a nearly symmetrical pattern about the zero line. Additionally, the effects of $X_1$, $X_2$, $X_6$, and $X_9$ seem to be initially overestimated, with their absolute values decreasing as $\lambda$ increases and stabilizing at non-zero values. Conversely, the effect of $X_{14}$ appears to be initially underestimated, with its coefficient slightly increasing as $\lambda$ increases. Notably, coefficients tend to stabilize around $\lambda = 0.2$, suggesting that estimates at this level of regularization are more suitable for assessing the explanatory variables' effects.

## 6. Solution.
(1). The Fitted versus Residual plot is shown below. Obviously it does not follows the homoscedasticity assumption.

Fitted versus Residual Plot

(2). After the transformation of $Y_i' = \log Y_i$, the transformed plot is as follows. And now it seems more like a distribution of homoscedasticity.



Fitted versus Residual Plot (Transformed)