

# Logit 模型 - 预测

宋歌 2015080086 数52

5/23/2018

## 1 研究目的

给定数据集，建立不同的模型，进行回归诊断，并对各个模型的预测能力进行评估。

## 2 实验过程及结果讨论

### 2.1 数据分析

- 已知：

$$\text{greID} = \begin{cases} \text{低} & \text{gre} \in [220, 400] \\ \text{中} & \text{gre} \in (400, 600] \\ \text{高} & \text{gre} \in (600, 800] \end{cases}$$

- 数据类型：

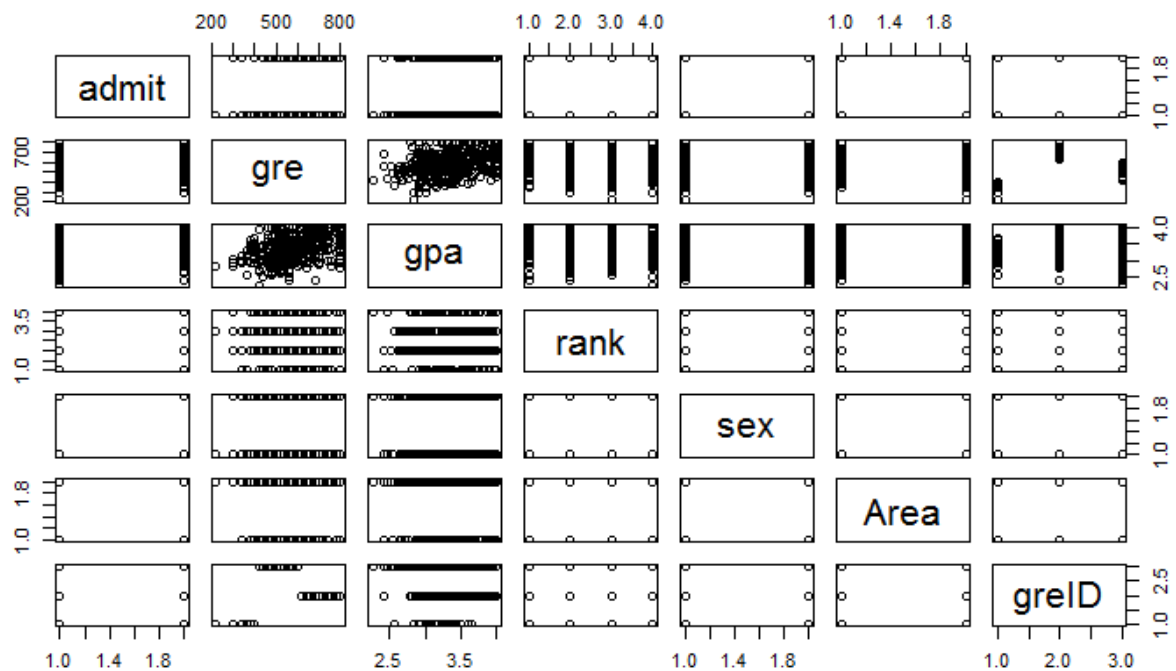
```
#读取并分析原始数据
dat <- read.table('pgBinary.txt')
summary(dat)
```

```
##      admit      gre      gpa      rank
## Min.   :0.0000 Min.   :220.0 Min.   :2.260 Min.   :1.000
## 1st Qu.:0.0000 1st Qu.:520.0 1st Qu.:3.130 1st Qu.:2.000
## Median :0.0000 Median :580.0 Median :3.395 Median :2.000
## Mean   :0.3175 Mean   :587.7 Mean   :3.390 Mean   :2.485
## 3rd Qu.:1.0000 3rd Qu.:660.0 3rd Qu.:3.670 3rd Qu.:3.000
## Max.   :1.0000 Max.   :800.0 Max.   :4.000 Max.   :4.000
##      sex      Area greID
## female:215 A:108 低: 31
## male :185 C:292 高:174
##                中:195
##
##
##
```

连续型变量：gre, gpa;

分类型变量：admit为二元分类变量，rank为四元分类变量，sex为二元分类变量，Area为二元分类变量，greID为三元分类变量；

- 初步观察数据之间的关系——数据散点图:



从散点图中也可观察出，除了gre,gpa其余变量均为分类变量；

其中gre与gpa之间有一定的线性关系；gre与greID之间呈现出分段函数的形式；

## 2.2 admit作为响应变量，gre,gpa,rank,sex,Area 作为协变量，建立预测模型并评估

### 2.2.1 研究方案

- 选取训练集与测试集：用简单随机抽样的方法选取80% 的数据作为训练集，剩下的作为测试集。
- 在R中利用glm建立logistic线性模型mylogit：二元分类变量admit作为响应变量，gre,gpa,rank,sex,Area作为协变量，在训练集上进行回归。
- 利用训练集上拟合出来的模型，在测试集上进行预测。
- 利用ROC曲线分别评估：模型在训练集上的拟合程度，模型在测试集上的预测能力。

### 2.2.2 详细方法与结果讨论

- 数据处理

```
#admit, rank, sex, Area, greID化为哑变量
dat$admit <- factor(dat$admit)
dat$rank <- factor(dat$rank)
dat$sex <- factor(dat$sex)
dat$Area <- factor(dat$Area)
dat$greID <- factor(dat$greID)
```

- 建立训练集和测试集

```
#建立训练集和测试集 - 简单随机抽样
set.seed(2015080086)
train_sub <- sample(nrow(dat), 4/5*nrow(dat))
dat_train <- dat[train_sub,]
dat_test <- dat[-train_sub,]
```

用简单随机抽样的方法，抽取数据集的80%作为训练集，剩下的20%作为测试集；

- 在训练集上进行Logistic回归

```
#admit~gre+gpa+rank+sex+Area在训练集上进行logistic回归
mylogit <- glm(admit ~ gre + gpa + rank + sex + Area, data = dat_train, binomial(link = 'logit'))
summary(mylogit)
```

```
##
## Call:
## glm(formula = admit ~ gre + gpa + rank + sex + Area, family = binomial(link = "logit"),
##      data = dat_train)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -1.6984  -0.9173  -0.6180   1.1193   2.1187
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept) -4.141696   1.295964  -3.196 0.001394 **
## gre          0.003004   0.001221   2.462 0.013835 *
## gpa          0.783863   0.368346   2.128 0.033332 *
## rank2       -0.562438   0.350944  -1.603 0.109013
## rank3       -1.301363   0.384090  -3.388 0.000704 ***
## rank4       -1.773979   0.491304  -3.611 0.000305 ***
## sexmale     -0.020431   0.252669  -0.081 0.935553
## AreaC       -0.201162   0.272216  -0.739 0.459920
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##      Null deviance: 413.11  on 319  degrees of freedom
## Residual deviance: 372.31  on 312  degrees of freedom
## AIC: 388.31
##
## Number of Fisher Scoring iterations: 4
```

从p值可以观察到，变量sex,Area,rank2均不显著；

- 评估模型的拟合程度与预测能力

由于该模型的响应变量admit是二值变量，故可以用ROC曲线来评价该“分类器”，ROC曲线下的面积，也即AUC值越大，该模型拟合程度越好；

```

mylogitfit <- mylogit$fitted.values
mylogitfit_roc <- roc(dat_train$admit, mylogitfit)

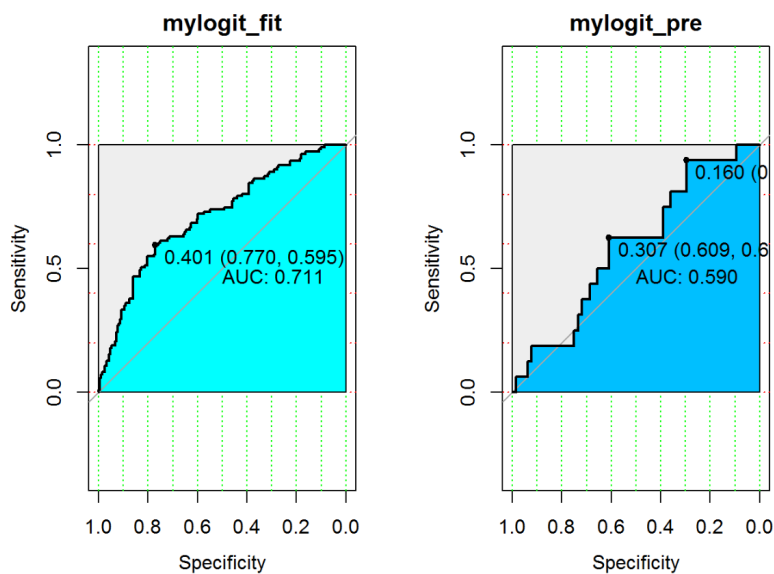
mylogitpre <- predict(mylogit, newdata = dat_test, type = "response")
mylogitpre_roc <- roc(dat_test$admit, mylogitpre)

par(mfrow = c(1, 2))

#评估mylogit在训练集上的拟合程度
plot(mylogitfit_roc, print.auc=TRUE, auc.polygon=TRUE, grid=c(0.1, 0.2),
     grid.col=c("green", "red"), max.auc.polygon=TRUE,
     auc.polygon.col="cyan", print.thres=TRUE, main = "mylogit_fit")

#评估mylogit在测试集上的预测能力
plot(mylogitpre_roc, print.auc=TRUE, auc.polygon=TRUE, grid=c(0.1, 0.2),
     grid.col=c("green", "red"), max.auc.polygon=TRUE,
     auc.polygon.col="deepskyblue", print.thres=TRUE, main = "mylogit_pre")

```



通过比较AUC值，可以从图中观察到一个很自然的结果：模型在测试集上的拟合程度优于模型在训练集上的拟合程度。

## 2.3 加入协变量greID，评估并比较预测精度

### 2.3.1 研究方案

- 在R中利用glm建立logistic线性模型logit1：二元分类变量admit作为响应变量，gre,gpa,rank,sex,Area,greID作为协变量，在之前定义的训练集上进行回归。
- 利用训练集上拟合出来的模型，在测试集上进行预测。
- 利用ROC曲线分别评估：模型在训练集上的拟合程度，模型在测试集上的预测能力。并将其与mylogit模型比较。

### 2.3.2 详细方法与结果讨论

- 加入协变量greID进行回归

```
#将greID加入协变量重新回归
logit1 <- glm(admit ~ gre + gpa + rank + sex + Area + greID, data = dat_train, binomial(link = 'logit'))
summary(logit1)
```

```
##
## Call:
## glm(formula = admit ~ gre + gpa + rank + sex + Area + greID,
##      family = binomial(link = "logit"), data = dat_train)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -1.6523  -0.9173  -0.6068   1.1072   2.1381
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept) -3.840966    1.548996  -2.480 0.013151 *
## gre           0.001626    0.002347   0.693 0.488481
## gpa           0.769780    0.368340   2.090 0.036630 *
## rank2        -0.570771    0.351340  -1.625 0.104258
## rank3        -1.303868    0.385300  -3.384 0.000714 ***
## rank4        -1.775480    0.490635  -3.619 0.000296 ***
## sexmale       -0.034026    0.253795  -0.134 0.893347
## AreaC        -0.175047    0.274900  -0.637 0.524277
## greID高        0.729463    1.005006   0.726 0.467944
## greID中        0.468243    0.761876   0.615 0.538824
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##      Null deviance: 413.11  on 319  degrees of freedom
## Residual deviance: 371.77  on 310  degrees of freedom
## AIC: 391.77
##
## Number of Fisher Scoring iterations: 4
```

- 从p值可以判断变量greID不显著；
- 将该回归结果与mylogit模型的回归结果相比较，可以看到加入greID后模型拟合的结果变化并不大，同样可以推断出变量greID并不显著，即对响应变量admit无影响。
- 比较mylogit(388.31)和logit1(391.77)模型的AIC值同样可知，logit1模型中包含了比mylogit模型更多的不显著的变量，也即greID对admit 无影响。
- 还观察到在mylogit模型中变量gre是有一定显著性的，但在logit1模型中加入了协变量greID之后，变量gre和greID在logit1模型中都是不显著的了。通过之前的数据散点图，以及之后第四部分的验证，我们知道gre和greID之间是高度线性相关的，故可以推断，加入greID之后，logit1模型的预测变量之间产生了共线性（或者共线性变严重），使得原本显著的变量gre在logit1模型中变得不再显著了。

## ● 模型的评估与比较

```

logitlfit <- logitl$fitted.values
logitlfit_roc <- roc(dat_train$admit, logitlfit)

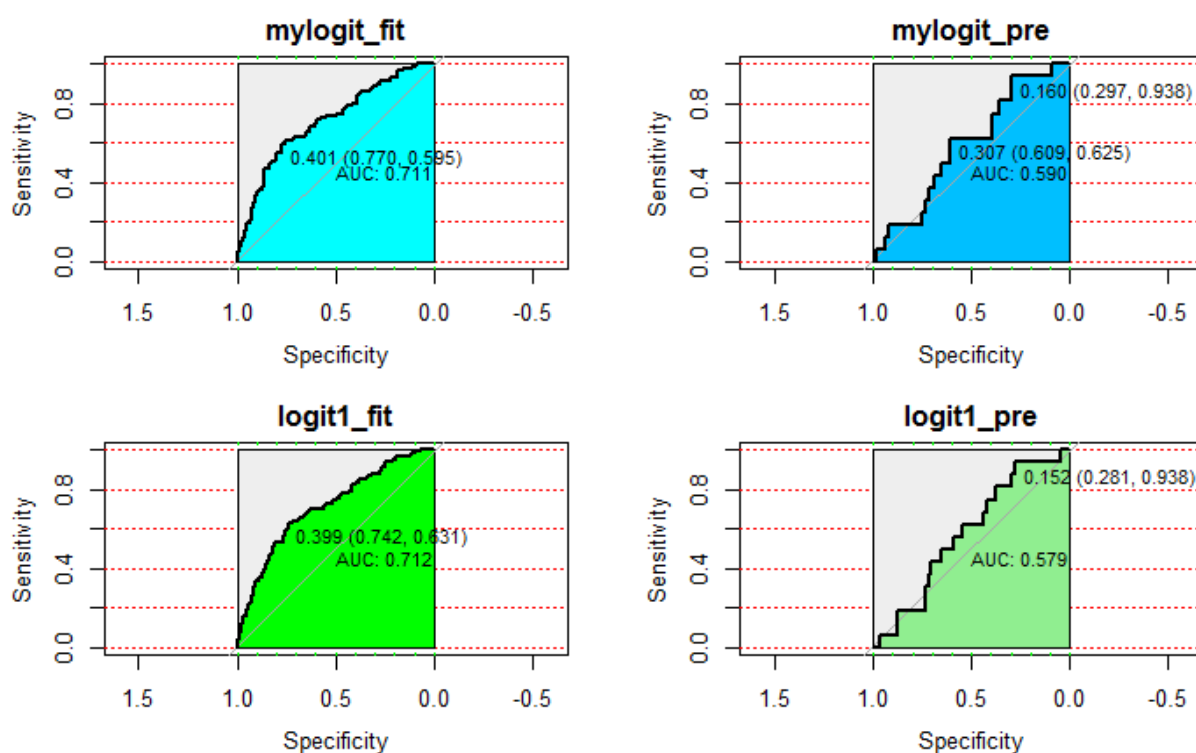
logitlpre <- predict(logitl, newdata = dat_test, type = "response")
logitlpre_roc <- roc(dat_test$admit, logitlpre)

par(mfrow = c(1,2))

#评估logitl在训练集上的拟合程度
plot(logitlfit_roc, print.auc=TRUE, auc.polygon=TRUE, grid=c(0.1, 0.2),
     grid.col=c("green", "red"), max.auc.polygon=TRUE,
     auc.polygon.col="green", print.thres=TRUE, main = "logitl_fit")

#评估logitl在测试集上的预测能力
plot(logitlpre_roc, print.auc=TRUE, auc.polygon=TRUE, grid=c(0.1, 0.2),
     grid.col=c("green", "red"), max.auc.polygon=TRUE,
     auc.polygon.col="lightgreen", print.thres=TRUE, main = "logitl_pre")

```



通过比较AUC值，可以观察到加入了协变量greID的logit1模型在训练集上的拟合程度略微好于mylogit模型，但在测试集上的预测能力低于mylogit模型；

又根据之前所得到的“变量greID不显著”该结果，我们可以推断：logit1模型比起mylogit模型，出现了过拟合的现象。即在训练集上的表现变好，很好地拟合了训练集上的数据，但在测试集上的表现变差，因为加入了无关的变量greID。

## 2.4 用greID替换gre，评估并比较拟合程度和预测能力

### 2.4.1 研究方案

- 在R中利用glm建立logistic线性模型logit2：二元分类变量admit作为响应变量，gpa,rank,sex,Area,greID作为协变量，在之前定义的训练集上进行回归。

- 利用训练集上拟合出来的模型，在测试集上进行预测。
- 利用ROC曲线分别评估：模型在训练集上的拟合程度，模型在测试集上的预测能力。并将其与mylogit、logit1模型比较。

## 2.4.2 详细方法与结果讨论

- 用greID替换gre进行回归

```
#gre用greID替代
logit2 <- glm(admit ~ greID + gpa + rank + sex + Area, data = dat_train, binomial(link = 'logit'))
summary(logit2)
```

```
##
## Call:
## glm(formula = admit ~ greID + gpa + rank + sex + Area, family = binomial(link = "logit"),
##      data = dat_train)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -1.5925  -0.9130  -0.6152   1.1026   2.1701
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept) -3.34587    1.36708  -2.447 0.014387 *
## greID高      1.24930    0.66954   1.866 0.062053 .
## greID中      0.72886    0.66226   1.101 0.271086
## gpa          0.80125    0.36511   2.195 0.028195 *
## rank2       -0.57989    0.35030  -1.655 0.097846 .
## rank3       -1.31757    0.38414  -3.430 0.000604 ***
## rank4       -1.78262    0.48934  -3.643 0.000270 ***
## sexmale     -0.04658    0.25294  -0.184 0.853885
## AreaC       -0.14684    0.27183  -0.540 0.589067
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##      Null deviance: 413.11  on 319  degrees of freedom
## Residual deviance: 372.25  on 311  degrees of freedom
## AIC: 390.25
##
## Number of Fisher Scoring iterations: 4
```

- 可以看到，在logit2模型中，变量greID有微弱的显著性，低于变量gre在mylogit模型中的显著性，说明变量gre对admit的影响强于变量greID对admit的影响；
- 通过比较mylogit(388.31),logit1(391.77),logit2(390.25)的AIC值，我们可以加强之前的结论：greID对响应变量admit的无关性强于gre对响应变量admit的无关性。

- 模型的评估与比较

```

logit2fit <- logit2$fitted.values
logit2fit_roc <- roc(dat_train$admit, logit2fit)

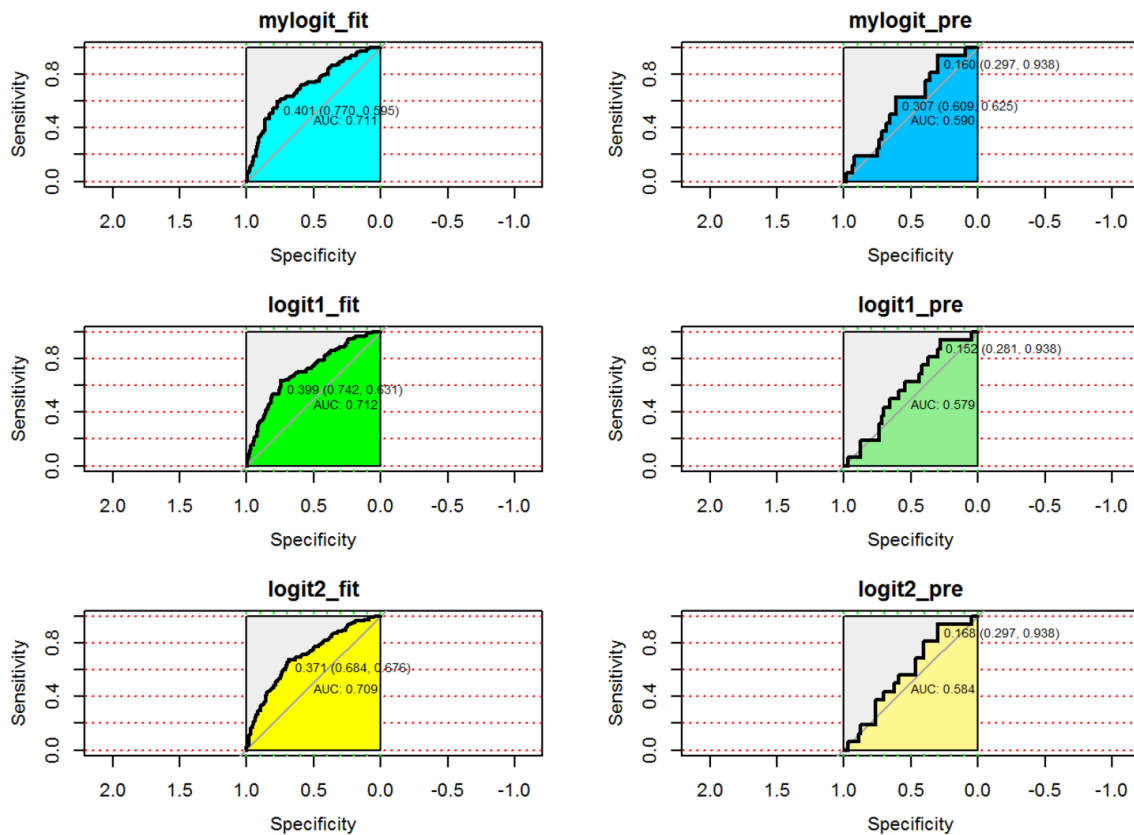
logit2pre <- predict(logit2, newdata = dat_test, type = "response")
logit2pre_roc <- roc(dat_test$admit, logit2pre)

par(mfrow = c(1,2))

#评估logit2在训练集上的拟合程度
plot(logit2fit_roc, print.auc=TRUE, auc.polygon=TRUE, grid=c(0.1, 0.2),
     grid.col=c("green", "red"), max.auc.polygon=TRUE,
     auc.polygon.col="yellow", print.thres=TRUE, main = "logit2_fit")

#评估logit2在测试集上的预测能力
plot(logit2pre_roc, print.auc=TRUE, auc.polygon=TRUE, grid=c(0.1, 0.2),
     grid.col=c("green", "red"), max.auc.polygon=TRUE,
     auc.polygon.col="khakil", print.thres=TRUE, main = "logit2_pre")

```



- 三个模型都验证了一个平凡的事实：模型在测试集上的拟合程度优于模型在训练集上的拟合程度。
- 通过比较AUC值，我们观察到训练集上的拟合程度 $\text{logit1} > \text{mylogit} > \text{logit2}$ ，测试集上的拟合程度（预测能力） $\text{mylogit} > \text{logit2} > \text{logit1}$ 。验证了我们之前得到的结论：加入协变量greID的logit1模型，比起mylogit模型和logit2模型都出现了过拟合的情况，表现为拟合程度最好但预测能力最差；变量gre比变量greID更显著，对admit的影响更大，也可以理解为gre包含了比greID更多的信息，选择gre作为协变量比选择greID作为协变量更优，从而mylogit模型从整体上都比logit2模型更优，表现为mylogit模型比logit2模型在训练集上拟合得更好，在测试集上也拟合得更好；gre与greID之间产生了共线性，也可能导致logit1模型的预测能力最差。



## 2.5 gre和greID分别作为响应变量，建立和其余变量之间的线性模型并做诊断

### 2.5.1 gre作为响应变量

#### 1) 研究方案

- gre为连续型变量，故以gre为响应变量，admit,gpa,rank,sex,Area为协变量，建立一般的线性模型，在之前定义的训练集上进行回归。
- 回归诊断：残差图，响应图，Durbin-Watson检验

#### 2) 详细方法与结果讨论

- 训练集上建立线性模型

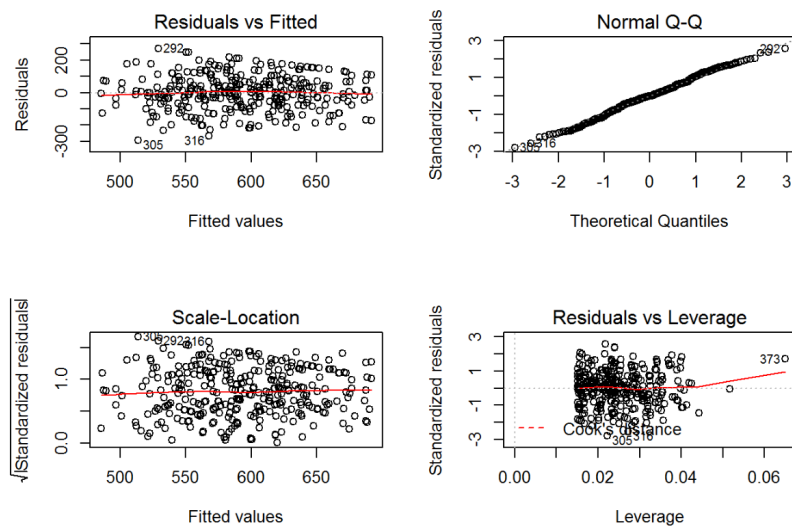
```
#以gre为响应变量在训练集上建立线性模型
l1 <- lm(gre ~ admit + gpa + rank + sex + Area, data = dat_train)
summary(l1)

##
## Call:
## lm(formula = gre ~ admit + gpa + rank + sex + Area, data = dat_train)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -293.798  -64.785    0.144   70.746  270.767
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  213.785     58.945   3.627 0.000335 ***
## admit1       32.313     13.218   2.445 0.015051 *
## gpa          106.559     16.368   6.510 3e-10 ***
## rank2         3.061     18.003   0.170 0.865100
## rank3        -16.332     19.077  -0.856 0.392585
## rank4        -12.019     22.051  -0.545 0.586112
## sexmale        1.369     12.018   0.114 0.909349
## AreaC         14.783     13.195   1.120 0.263453
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 106.6 on 312 degrees of freedom
## Multiple R-squared:  0.1622, Adjusted R-squared:  0.1434
## F-statistic: 8.631 on 7 and 312 DF,  p-value: 1.101e-09
```

变量rank,sex,Area均不显著；模型的拟合优度很低；模型是显著的；

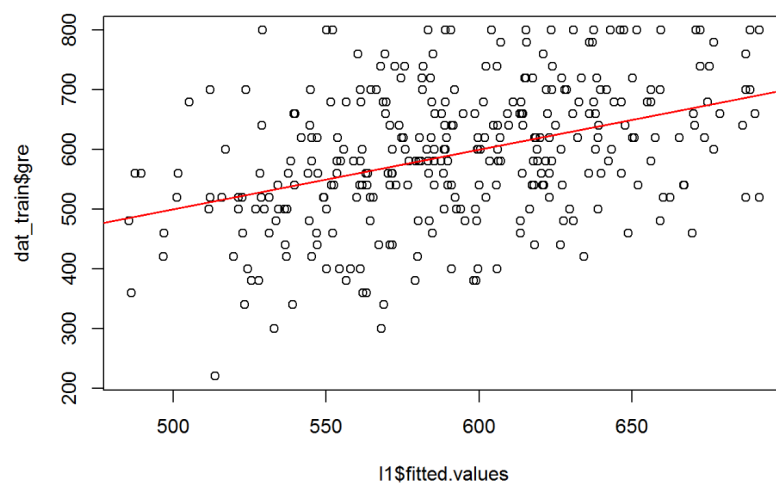
- 回归诊断：残差图与响应图

```
#残差图
par(mfrow = c(2, 2))
plot(l1)
```



残差基本符合正态性，独立性，同方差性假设；

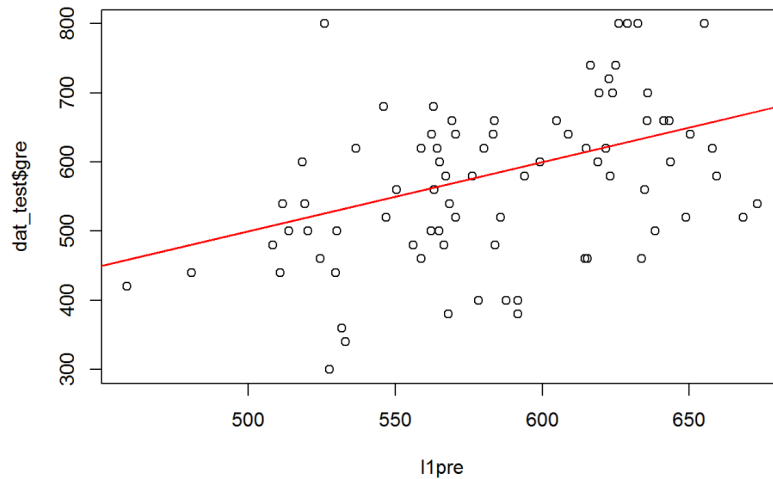
```
#响应图
plot(l1$fitted.values, dat_train$gre)
lines(dat_train$gre, dat_train$gre, col = 'red')
```



散点整体趋势符合 $y = x$ ，但过于分散，也印证了该模型的拟合优度很低；

- 回归诊断：评估在测试集上的预测能力

```
#评估l1在测试集上的预测能力
l1pre <- predict(l1, newdata = dat_test)
plot(l1pre, dat_test$gre)
lines(dat_test$gre, dat_test$gre, col = 'red')
```



```
sqrt(sum((l1pre - dat_test$gre)^2) / 80) #平均误差
```

```
## [1] 104.5817
```

散点并未很好地呈现 $y = x$ ，平均误差并非远小于变量gre自身尺度，该模型在预测集上的表现并不好。

#### ● 回归诊断：Durbin-Watson检验

```
durbinWatsonTest(l1)
```

```
## lag Autocorrelation D-W Statistic p-value
## 1 -0.009940962 2.018743 0.81
## Alternative hypothesis: rho != 0
```

p值很大，接受0假设，认为误差之间相互独立。

## 2.5.2 greID作为响应变量

### 1) 研究方案

- 在R中利用polr建立有序logistic回归线性模型l2：以三元分类变量greID为响应变量，admit,gpa,rank,sex,Area为协变量，在之前定义的训练集上进行回归。
- 利用训练集上拟合出来的模型，在测试集上进行预测。
- 利用正确率评估模型的拟合程度和预测能力。

### 2) 详细方法与结果讨论

- 数据处理：

```
#定义低中高顺序以便回归
y <- c()
for(i in 1:400){
  if(dat$greID[i] == "低"){
    y[i] <- "a"
  }
  if(dat$greID[i] == "中"){
    y[i] <- "b"
  }
  if(dat$greID[i] == "高"){
    y[i] <- "c"
  }
}

y <- factor(y)
y_train <- y[train_sub]
y_test <- y[-train_sub]

#提取数据集里的低中高为123
Y <- unclass(y)

Y_train <- Y[train_sub]
Y_test <- Y[-train_sub]
```

将“低，中，高”换为“a,b,c”，从而在回归时，R能够按照我们想要的顺序识别哑变量greID；同时将“低，中，高”也化为数值型1,2,3，便于进行回归诊断；

- 训练集上建立有序logistic线性回归模型

```
l2 <- polr(y_train ~ admit + gpa + rank + sex + Area, data = dat_train, method = "logistic")
summary(l2)
```

```
##
## Re-fitting to get Hessian
```

```
## Call:
## polr(formula = y_train ~ admit + gpa + rank + sex + Area, data = dat_train,
##       method = "logistic")
##
## Coefficients:
##              Value Std. Error t value
## admit1    0.595843   0.2533  2.35196
## gpa       1.847288   0.3284  5.62575
## rank2     0.213051   0.3450  0.61759
## rank3    -0.118604   0.3644 -0.32551
## rank4    -0.139622   0.4154 -0.33613
## sexmale   0.169110   0.2294  0.73703
## AreaC    -0.006064   0.2517 -0.02409
##
## Intercepts:
##      Value Std. Error t value
## a|b  3.7535  1.1442   3.2806
## b|c  6.8380  1.1912   5.7406
##
## Residual Deviance: 526.9958
## AIC: 544.9958
```

- 评估在训练集上的拟合程度

```
#评估l2在训练集上的拟合程度
l2fit <- c()
for(i in 1:320){
  l2fit[i] <- which.max(l2$fitted.values[i,]) #对每一个样本点取概率最大的作为拟合值
}

#考查拟合正确率
tn_train <- as.numeric((Y_train == l2fit))
tr_train <- sum(tn_train) / 320
tr_train
```

```
## [1] 0.58125
```

拟合正确率略高于0.5，拟合得并不好：

- 评估在测试集上的预测能力

```
#评估l2在测试集上的预测能力
l2pre <- predict(l2, newdata = dat_test)
l2pre <- unclass(l2pre)

#考查预测正确率
tn_test <- as.numeric((Y_test == l2pre))
tr_test <- sum(tn_test) / 80
tr_test
```

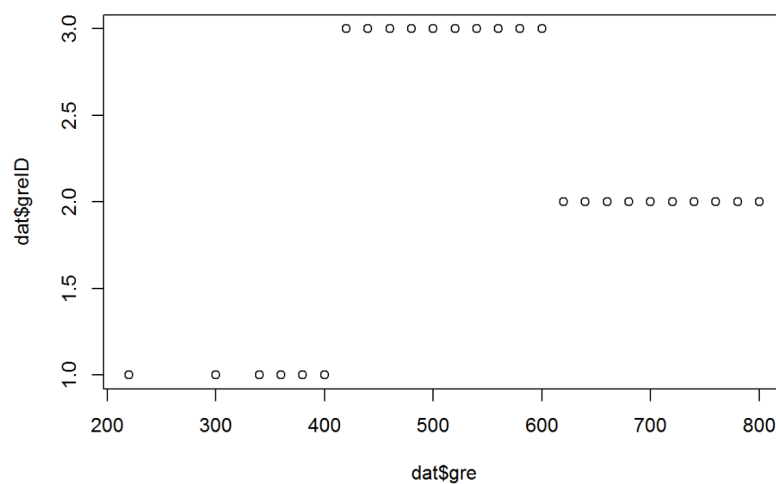
```
## [1] 0.6
```

预测正确率0.6，预测能力一般：

## 2.6 把gre和greID分别作响应变量和协变量建立回归模型，并探讨二者之间的关系

### 2.6.1 画数据散点图初步观察

```
#画图初步观察
plot(dat$gre, dat$greID)
```



可见gre与greID之间有类似于分段函数的关系；

## 2.6.2 gre作为响应变量

### 1) 研究方案

- gre作响应变量，greID作协变量，在数据集上建立简单线性模型f。
- 回归诊断：残差图，响应图，平均误差，Durbin-Watson检验

### 2) 详细方法与结果讨论

- 数据集上建立一般线性回归模型

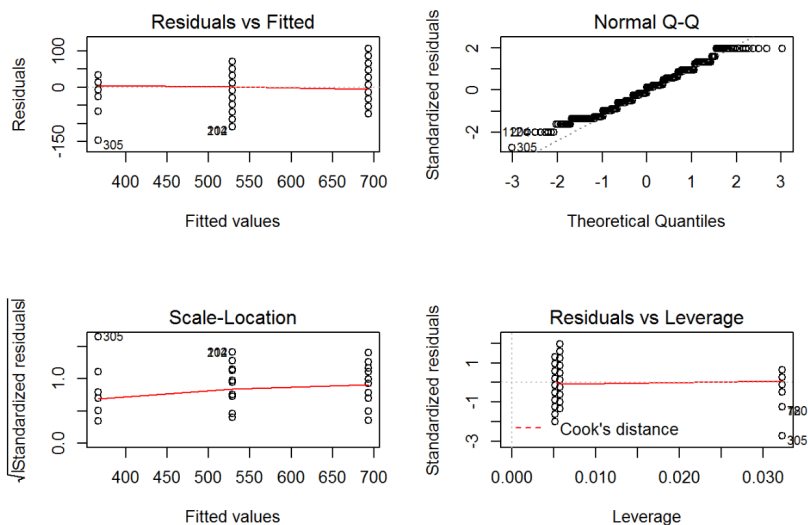
```
f <- lm(gre ~ y, data = dat)
summary(f)

##
## Call:
## lm(formula = gre ~ y, data = dat)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -146.452  -48.615    6.667   33.548  106.667
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   366.452      9.811   37.35  <2e-16 ***
## yb             162.164     10.562   15.35  <2e-16 ***
## yc             326.882     10.649   30.70  <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 54.62 on 397 degrees of freedom
## Multiple R-squared:  0.7775, Adjusted R-squared:  0.7764
## F-statistic: 693.7 on 2 and 397 DF,  p-value: < 2.2e-16
```

模型显著；变量均显著；模型拟合优度较高；

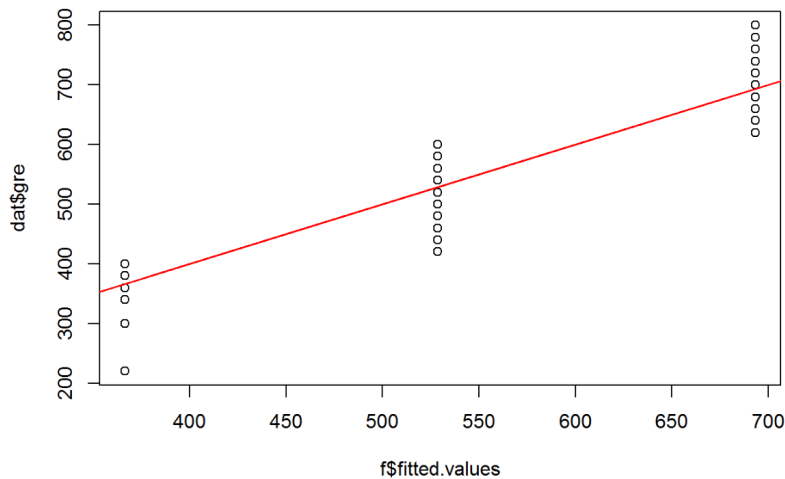
- 回归诊断

```
#残差图
par(mfrow = c(2,2))
plot(f)
```



残差基本符合正态性，独立性，同方差性假设；

```
#响应图
plot(f$fitted.values, dat$gre)
lines(dat$gre, dat$gre, col = 'red')
```



散点集中在直线 $y = x$ 两侧，拟合程度较好；

```
#平均误差
sqrt(sum((f$fitted.values - dat$gre)^2)) / 400
```

```
## [1] 2.720871
```

平均误差远小于gre自身尺度，拟合程度较好；

```
#Durbin-Watson检验
durbinWatsonTest(f)
```

```
## lag Autocorrelation D-W Statistic p-value
## 1 -0.03944415 2.074431 0.442
## Alternative hypothesis: rho != 0
```

p值很大，接受0假设，认为误差之间相互独立。

### 2.6.3 greID作为响应变量

#### 1) 研究方案

- 在R中利用polr建立有序logistic回归线性模型g：以三元分类变量greID为响应变量，gre为协变量，在数据集上进行回归。
- 回归诊断：考查正确率

#### 2) 详细方法与结果讨论

- 数据集上建立回归模型

```
#greID作响应变量
g <- polr(y ~ gre, data = dat, start = c(1, 400, 600), method = "logistic")
summary(g)

##
## Re-fitting to get Hessian

## Call:
## polr(formula = y ~ gre, data = dat, start = c(1, 400, 600), method = "logistic")
##
## Coefficients:
##      Value Std. Error t value
## gre 0.9771    0.0763   12.81
##
## Intercepts:
##      Value      Std. Error  t value
## a|b    400.7286      0.0009 459127.5775
## b|c    595.8703     51.2192   11.6337
##
## Residual Deviance: 0.008013004
## AIC: 6.008013
```

- 回归诊断

```
#回归诊断
gfit <- c()
for(i in 1:400){
  gfit[i] <- which.max(g$fitted.values[i,]) #对每一个样本点取概率最大的作为拟合值
}

tn <- as.numeric((Y == gfit)) #考查拟合正确率
tr <- sum(tn) / 400
tr
```

```
## [1] 1
```

拟合正确率为1.

## 2.6.4 gre与greID之间的关系

由以上回归结果可知

$$\begin{aligned} \text{gre} &= 366.452 + 162.164\text{中} + 326.882\text{高} \\ \text{gre} &= 366.452\text{低} + 528.616\text{中} + 693.334\text{高} \\ \text{logit}P(\text{greID} = \text{低}) &= 400.73 - 0.977\text{gre} \\ \text{logit}P(\text{greID} = \text{中}) &= 595.87 - 0.977\text{gre} \end{aligned}$$

(不知道理解得对不对……)