

《多元统计分析》课后作业

姓名： 李倩倩

学号： 2024017349

班级： 统计 24-1 班

中国石油大学（北京）克拉玛依校区文理学院数学与统计系

Friday 22nd May, 2026

作业要求

1. 可以和其他同学讨论作业当中的问题，但应当自己独立完成作业
2. 计算、证明等要有过程，要有主要步骤的说明
3. 请将计算、绘图所用的 R 代码以及生成的结果和图像一并添加在作业文件当中
4. 请使用 \LaTeX 编辑并生成 PDF 格式的文件，第 X 周作业文件命名方式：学号-姓名-X.pdf
5. 评分标准：每一问得分 $\in \{2, 1, 0\}$
 - 2: 按时完成并上交作业，且答案基本正确
 - 1: 按时完成并上交作业，且答案部分正确
 - 0: 答案完全错误，或者迟交作业(规定时间72小时之后)
6. 请将完成的 PDF 格式的作业文件发送至邮箱：xiaolei@cup.edu.cn
7. 每位同学可以有一次迟交作业的机会，但不得晚于规定时间三日之后
8. 第 10 周作业截止时间：2026年5月22日24:00

目录

Chapter 1

第 10 周作业

第 10 周作业截止时间: 2026年5月22日24:00

第 10 周作业完成时间: Friday 22nd May, 2026 08:59

1. 数据文件 `2018-Mean Expenditure of Urban Residents.csv` 包含了 2018 年中国 31 个省(市)的城镇居民人均生活费用, 各变量的含义如下: Province 省市名称, Food 食品, Cloth 衣着, Residential 居住, Expenditure 生活用品与服务, Trans-Com (交通通讯), Education (教育娱乐), Healthcare (医疗保健), and Others (其它). 作上述变量的因子分析.

(a) [2 分] 读入数据, 从相关矩阵出发作因子分析.

【解】 读入数据后, 取除 `Province` 以外的 8 个支出变量作因子分析. 以下代码默认在数据文件所在目录运行; 若在其它目录运行, 将文件名替换为题目给出的绝对路径即可. 由于各变量量纲相同但数值水平差异较大, 下面从相关矩阵出发作最大似然因子分析, 并使用方差最大正交旋转.

```
exp_file <- "2018_Mean Expenditure of Urban Residents.csv"
exp_dat <- read.csv(exp_file, stringsAsFactors = FALSE)
rownames(exp_dat) <- exp_dat$Province
X1 <- exp_dat[, -1]

R1 <- cor(X1)
round(R1, 3)

fa1 <- factanal(covmat = R1, factors = 2, n.obs = nrow(X1),
               rotation = "varimax")

L1 <- unclass(loadings(fa1))
ord1 <- order(colSums(L1^2), decreasing = TRUE)
L1 <- L1[, ord1]
for (j in seq_len(ncol(L1))) {
  if (L1[which.max(abs(L1[, j])), j] < 0) L1[, j] <- -L1[, j]
```

```

}
colnames(L1) <- c("F1", "F2")

scores1 <- as.matrix(scale(X1)) %*% solve(R1) %*% L1
colnames(scores1) <- c("F1", "F2")

```

相关矩阵为:

	Food	Cloth	Residential	Expenditure	Trans_Com	Education	Healthcare	Others
Food	1.000	0.580	0.834	0.460	0.850	0.820	0.582	0.865
Cloth	0.580	1.000	0.667	0.463	0.801	0.678	0.798	0.808
Residential	0.834	0.667	1.000	0.281	0.869	0.900	0.738	0.912
Expenditure	0.460	0.463	0.281	1.000	0.493	0.376	0.441	0.470
Trans_Com	0.850	0.801	0.869	0.493	1.000	0.861	0.776	0.914
Education	0.820	0.678	0.900	0.376	0.861	1.000	0.843	0.922
Healthcare	0.582	0.798	0.738	0.441	0.776	0.843	1.000	0.857
Others	0.865	0.808	0.912	0.470	0.914	0.922	0.857	1.000

(b) [2 分] 确定公共因子数量并给出理由。

【解】 先计算相关矩阵的特征根和累计贡献率:

```

lambda1 <- eigen(R1)$values
prop1 <- lambda1 / sum(lambda1)
eig_tab1 <- data.frame(No = 1:length(lambda1),
                      Eigenvalue = lambda1,
                      Proportion = prop1,
                      Cumulative = cumsum(prop1))
round(eig_tab1, 4)

plot(lambda1, type = "b", pch = 19, xlab = "Component",
      ylab = "Eigenvalue", main = "Urban Residents: Scree Plot")
abline(h = 1, lty = 2, col = "gray50")

```

输出结果为:

	No	Eigenvalue	Proportion	Cumulative
1	1	6.0870	0.7609	0.7609
2	2	0.8250	0.1031	0.8640
3	3	0.5497	0.0687	0.9327
4	4	0.2711	0.0339	0.9666
5	5	0.0986	0.0123	0.9789
6	6	0.0858	0.0107	0.9897

7	7	0.0562	0.0070	0.9967
8	8	0.0265	0.0033	1.0000

第 1 个特征根解释了 76.09% 的信息，第 1、2 个特征根累计解释了 86.40% 的信息；从碎石图看，第 2 个特征根之后下降趋缓，因此本文保留 2 个公共因子。其中 Kaiser 准则只会保留 1 个公共因子，但为刻画第二类支出结构差异，并与后续公共因子平面和第 2 公共因子排序相一致，取公共因子数为 $m = 2$ 。

运行上述代码即可得到城镇居民生活费用数据的碎石图。

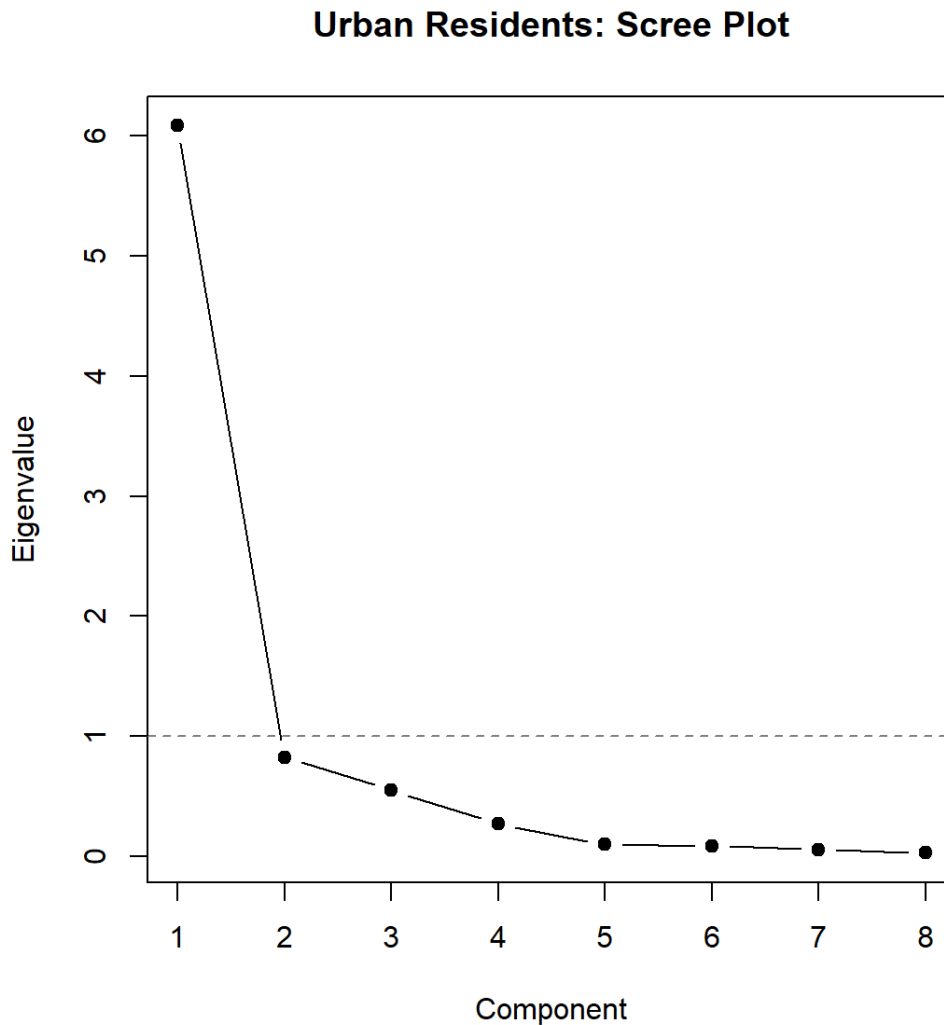


图 1.1: 城镇居民生活费用数据的碎石图

(c) [2 分] 给出因子旋转之后的因子载荷矩阵。

【解】 旋转后的因子载荷矩阵如下，其中 h^2 为共同度， u^2 为特殊方差。

```
load_tab1 <- data.frame(L1,
                        h2 = rowSums(L1^2),
```

```

                                u2 = 1 - rowSums(L1^2))
round(load_tab1, 3)
colSums(L1^2)
colSums(L1^2) / ncol(X1)

```

表 1.1: 城镇居民生活费用数据旋转后的因子载荷矩阵.

变量	F_1	F_2	h^2	u^2
Food	0.902	0.302	0.904	0.096
Cloth	0.451	0.691	0.680	0.320
Residential	0.777	0.512	0.865	0.135
Expenditure	0.315	0.357	0.227	0.773
Trans_Com	0.746	0.562	0.873	0.127
Education	0.678	0.656	0.890	0.110
Healthcare	0.330	0.941	0.995	0.005
Others	0.740	0.651	0.971	0.029

两个因子的方差贡献分别为 $3.3910/8 = 42.39\%$ 和 $3.0145/8 = 37.68\%$, 累计解释约 80.07% 的相关矩阵总变异.

(d) [2 分] 作变量在公共因子平面上的散点图, 对公共因子作出解释.

【解】

```

plot(L1[, 1], L1[, 2], pch = 19,
     xlab = "Factor 1", ylab = "Factor 2",
     main = "Urban Residents: Variable Loadings")
abline(h = 0, v = 0, col = "gray70")
text(L1[, 1], L1[, 2], labels = rownames(L1), pos = 4, cex = 0.8)

```

从载荷看, F_1 在 Food、Residential、Trans_Com、Education、Others 上均有较高正载荷, 反映城镇居民总体消费水平, 尤其与食品、居住、交通通讯、教育娱乐及其它消费支出关系较强, 可解释为“综合生活消费水平因子”. F_2 在 Healthcare、Cloth、Education、Others 上载荷较高, 反映医疗保健、衣着以及文教娱乐等改善型或服务型支出的相对水平, 可解释为“医疗衣着与发展型消费因子”.

Urban Residents: Variable Loadings

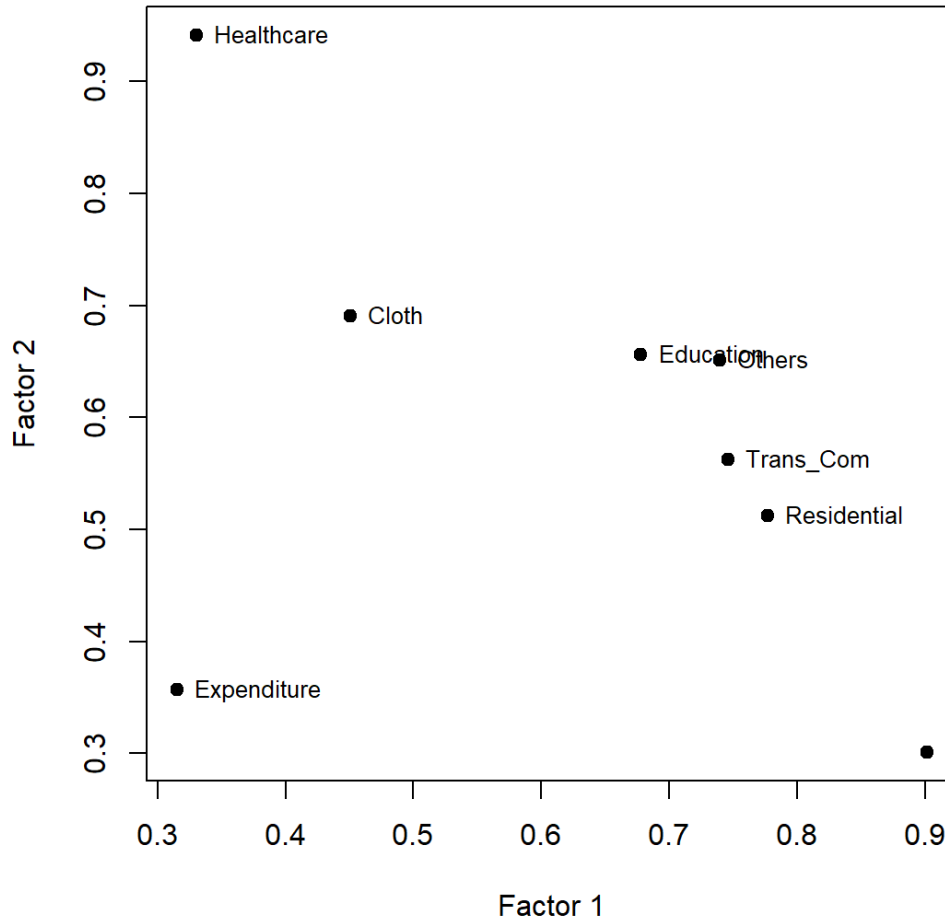


图 1.2: 变量在第 1、第 2 公共因子平面上的散点图.

(e) [2 分] 利用因子得分，作观测数据在公共因子平面上的散点图.

【解】 利用回归法估计因子得分:

$$\hat{F} = ZR^{-1}\hat{\Lambda},$$

其中 Z 为标准化后的观测矩阵， R 为相关矩阵， $\hat{\Lambda}$ 为旋转后的载荷矩阵.

```
score_dat1 <- data.frame(Province = exp_dat$Province, scores1)

plot(score_dat1$F1, score_dat1$F2, pch = 19,
      xlab = "Factor score 1", ylab = "Factor score 2",
      main = "Urban Residents: Factor Scores")
abline(h = 0, v = 0, col = "gray70")
text(score_dat1$F1, score_dat1$F2,
      labels = score_dat1$Province, pos = 4, cex = 0.65)
```

运行上述代码即可得到 31 个省 (市) 在第 1、第 2 公共因子平面上的散点图。

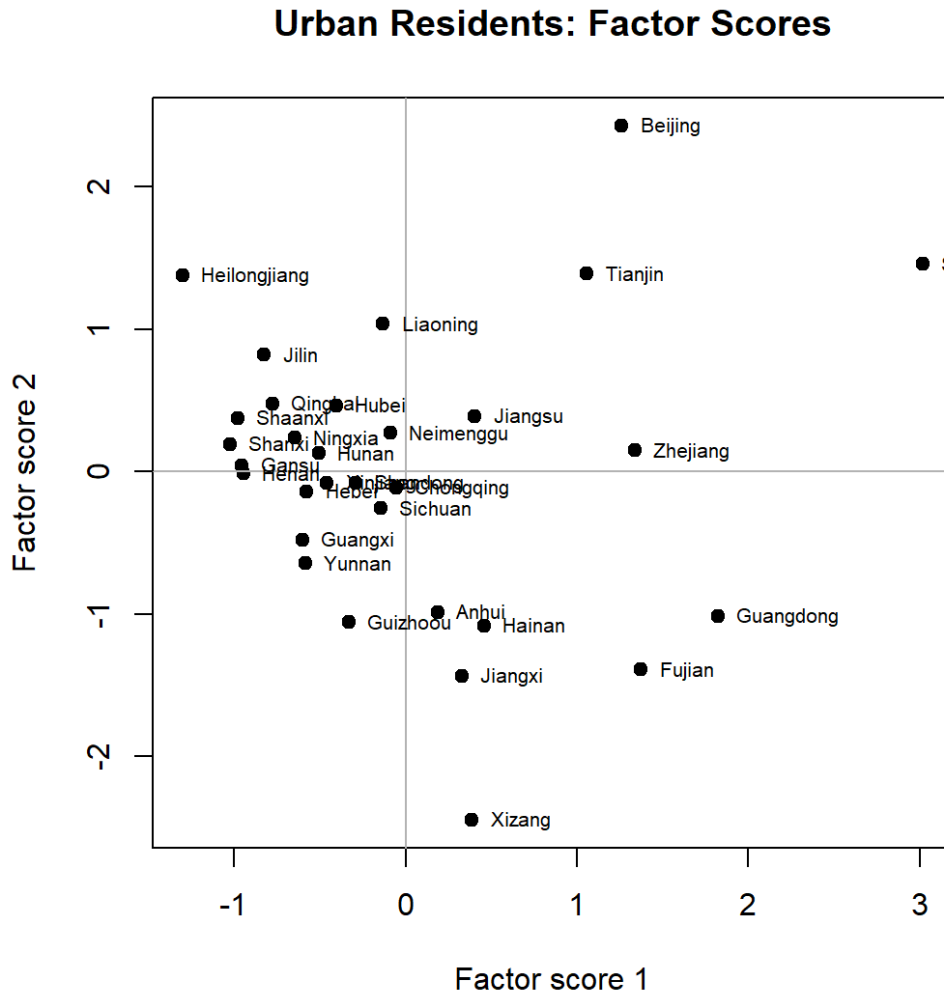


图 1.3: 31 个省 (市) 在第 1、第 2 公共因子平面上的散点图。

(f) [2 分] 利用第 1 个公共因子的得分对我国 31 个省 (市) 进行排序并作简要分析。

【解】

```
rank_F1 <- score_dat1[order(-score_dat1$F1), c("Province", "F1")]
rank_F1$Rank <- seq_len(nrow(rank_F1))
rank_F1 <- rank_F1[, c("Rank", "Province", "F1")]
rank_F1$F1 <- round(rank_F1$F1, 3)
rank_F1
```

排序结果如下:

Rank	Province	F1
------	----------	----

1	Shanghai	3.019
2	Guangdong	1.820
3	Fujian	1.375
4	Zhejiang	1.337
5	Beijing	1.261
6	Tianjin	1.059
7	Hainan	0.460
8	Jiangsu	0.405
9	Xizang	0.388
10	Jiangxi	0.328
11	Anhui	0.185
12	Chongqing	-0.056
13	Neimenggu	-0.089
14	Liaoning	-0.131
15	Sichuan	-0.145
16	Shandong	-0.293
17	Guizhoou	-0.333
18	Hubei	-0.404
19	Xinjiang	-0.458
20	Hunan	-0.506
21	Hebei	-0.576
22	Yunnan	-0.586
23	Guangxi	-0.602
24	Ningxia	-0.647
25	Qinghai	-0.779
26	Jilin	-0.824
27	Henan	-0.947
28	Gansu	-0.954
29	Shaanxi	-0.979
30	Shanxi	-1.026
31	Heilongjiang	-1.302

第 1 公共因子主要反映总体生活消费水平。上海、广东、福建、浙江、北京、天津得分较高，说明这些地区城镇居民在食品、居住、交通通讯、教育娱乐等方面的总体支出水平较高；山西、黑龙江、陕西、甘肃、河南等地得分偏低，说明综合消费水平相对较低。

(g) [2 分] 利用第 2 个公共因子的得分对我国 31 个省 (市) 进行排序并作简要分析。

【解】

```
rank_F2 <- score_dat1[order(-score_dat1$F2), c("Province", "F2")]
rank_F2$Rank <- seq_len(nrow(rank_F2))
rank_F2 <- rank_F2[, c("Rank", "Province", "F2")]
rank_F2$F2 <- round(rank_F2$F2, 3)
```

rank_F2

排序结果如下:

Rank	Province	F2
1	Beijing	2.430
2	Shanghai	1.456
3	Tianjin	1.388
4	Heilongjiang	1.377
5	Liaoning	1.036
6	Jilin	0.818
7	Qinghai	0.477
8	Hubei	0.465
9	Jiangsu	0.390
10	Shaanxi	0.374
11	Neimenggu	0.270
12	Ningxia	0.236
13	Shanxi	0.193
14	Zhejiang	0.148
15	Hunan	0.128
16	Gansu	0.041
17	Henan	-0.016
18	Shandong	-0.079
19	Xinjiang	-0.080
20	Chongqing	-0.115
21	Hebei	-0.141
22	Sichuan	-0.259
23	Guangxi	-0.477
24	Yunnan	-0.645
25	Anhui	-0.989
26	Guangdong	-1.016
27	Guizhou	-1.059
28	Hainan	-1.083
29	Fujian	-1.390
30	Jiangxi	-1.435
31	Xizang	-2.446

第 2 公共因子可解释为医疗衣着与发展型消费因子. 北京、上海、天津以及东北地区的黑龙江、辽宁、吉林得分较高, 说明这些地区在医疗保健、衣着等方面相对支出较突出; 西藏、江西、福建、海南、贵州等地区得分较低, 说明该类支出在其消费结构中的相对水平较弱.

(h) [2 分] 利用第 1、第 2 两个公共因子的得分对我国 31 个省 (市) 城镇居民人均生活费用进行综

合排序，并作简要分析。

【解】 以两个公共因子旋转后的方差贡献作为权数：

$$w_1 = \frac{3.3910}{3.3910 + 3.0145} = 0.5294, \quad w_2 = \frac{3.0145}{3.3910 + 3.0145} = 0.4706.$$

综合得分为

$$F = 0.5294F_1 + 0.4706F_2.$$

```
ss1 <- colSums(L1^2)
w1 <- ss1 / sum(ss1)
score_dat1$Composite <- as.vector(as.matrix(score_dat1[, c("F1", "F2")]) %*% w1)

rank_all1 <- score_dat1[order(-score_dat1$Composite),
                          c("Province", "F1", "F2", "Composite")]
rank_all1$Rank <- seq_len(nrow(rank_all1))
rank_all1 <- rank_all1[, c("Rank", "Province", "F1", "F2", "Composite")]
rank_all1[, c("F1", "F2", "Composite")] <-
  round(rank_all1[, c("F1", "F2", "Composite")], 3)
rank_all1
```

综合排序结果如下：

Rank	Province	F1	F2	Composite
1	Shanghai	3.019	1.456	2.283
2	Beijing	1.261	2.430	1.811
3	Tianjin	1.059	1.388	1.214
4	Zhejiang	1.337	0.148	0.777
5	Guangdong	1.820	-1.016	0.486
6	Liaoning	-0.131	1.036	0.418
7	Jiangsu	0.405	0.390	0.398
8	Neimenggu	-0.089	0.270	0.080
9	Fujian	1.375	-1.390	0.074
10	Hubei	-0.404	0.465	0.005
11	Heilongjiang	-1.302	1.377	-0.041
12	Jilin	-0.824	0.818	-0.051
13	Chongqing	-0.056	-0.115	-0.084
14	Qinghai	-0.779	0.477	-0.188
15	Shandong	-0.293	-0.079	-0.193
16	Sichuan	-0.145	-0.259	-0.199
17	Hunan	-0.506	0.128	-0.208
18	Ningxia	-0.647	0.236	-0.231
19	Hainan	0.460	-1.083	-0.266
20	Xinjiang	-0.458	-0.080	-0.280

21	Shaanxi	-0.979	0.374	-0.342
22	Anhui	0.185	-0.989	-0.367
23	Hebei	-0.576	-0.141	-0.371
24	Shanxi	-1.026	0.193	-0.452
25	Gansu	-0.954	0.041	-0.486
26	Jiangxi	0.328	-1.435	-0.502
27	Henan	-0.947	-0.016	-0.509
28	Guangxi	-0.602	-0.477	-0.543
29	Yunnan	-0.586	-0.645	-0.614
30	Guizhoou	-0.333	-1.059	-0.674
31	Xizang	0.388	-2.446	-0.945

综合排序中上海、北京、天津位居前三，说明直辖市城镇居民人均生活费用整体较高；浙江、广东、江苏等东部沿海省份也处于前列。西藏、贵州、云南、广西、河南等地区综合得分较低，说明总体消费水平和发展型消费水平相对偏低。个别地区如广东第 1 因子得分很高但第 2 因子得分偏低，因此综合名次低于北京、天津；辽宁第 1 因子不高但第 2 因子较高，综合名次有所上升。

2. 有 48 人申请到某公司就业。该公司对申请者的 15 项指标进行打分，这 15 项指标分别是：FL (求职信的形式)，APP (外貌)，AA (专业能力)，LA (讨人喜欢程度)，SC (自信心)，LC (洞察力)，HON (诚实度)，SMS (推销能力)，EXP (经验)，DRV (驾驶水平)，AMB (事业心)，GSP (理解能力)，POT (潜在能力)，KJ (社交能力)，SUIT (适应能力)。结果见数据文件 `Applicants.csv`。

(a) [2 分] 读入数据，从相关矩阵出发作因子分析。

【解】 读入 48 名申请者的 15 项指标评分，以 15 个评分变量的相关矩阵作为因子分析的出发点。以下代码默认在数据文件所在目录运行；若在其它目录运行，将文件名替换为题目给出的绝对路径即可。数据文件中 ID=2 的 AA 取值为 56，下面按原始数据记录直接参与计算。

```
app_file <- "Applicants.csv"
app_dat <- read.csv(app_file, stringsAsFactors = FALSE)
rownames(app_dat) <- app_dat$ID
X2 <- app_dat[, -1]

R2 <- cor(X2)
round(R2, 3)

fa2 <- factanal(covmat = R2, factors = 3, n.obs = nrow(X2),
               rotation = "varimax")

L2 <- unclass(loadings(fa2))
ord2 <- order(colSums(L2^2), decreasing = TRUE)
L2 <- L2[, ord2]
```

```

for (j in seq_len(ncol(L2))) {
  if (L2[which.max(abs(L2[, j])), j] < 0) L2[, j] <- -L2[, j]
}
colnames(L2) <- c("F1", "F2", "F3")

scores2 <- as.matrix(scale(X2)) %*% solve(R2) %*% L2
colnames(scores2) <- c("F1", "F2", "F3")

```

数据维数为 48×15 ，后续最大似然因子分析、因子旋转和因子得分均基于相关矩阵 R_2 完成。

(b) [2 分] 确定公共因子数量并给出理由。

【解】 计算相关矩阵 R_2 的特征根：

```

lambda2 <- eigen(R2)$values
prop2 <- lambda2 / sum(lambda2)
eig_tab2 <- data.frame(No = 1:length(lambda2),
                      Eigenvalue = lambda2,
                      Proportion = prop2,
                      Cumulative = cumsum(prop2))

round(eig_tab2, 4)

plot(lambda2, type = "b", pch = 19, xlab = "Component",
      ylab = "Eigenvalue", main = "Applicants: Scree Plot")
abline(h = 1, lty = 2, col = "gray50")

```

输出结果为：

	No	Eigenvalue	Proportion	Cumulative
1	1	7.5628	0.5042	0.5042
2	2	1.9923	0.1328	0.6370
3	3	1.3939	0.0929	0.7299
4	4	0.9721	0.0648	0.7947
5	5	0.7203	0.0480	0.8428
6	6	0.5860	0.0391	0.8818
7	7	0.4209	0.0281	0.9099
8	8	0.3445	0.0230	0.9329
9	9	0.2900	0.0193	0.9522
10	10	0.2394	0.0160	0.9681
11	11	0.1665	0.0111	0.9792
12	12	0.1128	0.0075	0.9868
13	13	0.0943	0.0063	0.9931
14	14	0.0670	0.0045	0.9975
15	15	0.0372	0.0025	1.0000

前三个特征根均大于 1，累计贡献率为 72.99%；从碎石图看，第 3 个特征根之后曲线明显趋缓，第 4 个特征根已小于 1。因此取公共因子数 $m = 3$ 。

运行上述代码即可得到申请者评分数据的碎石图。

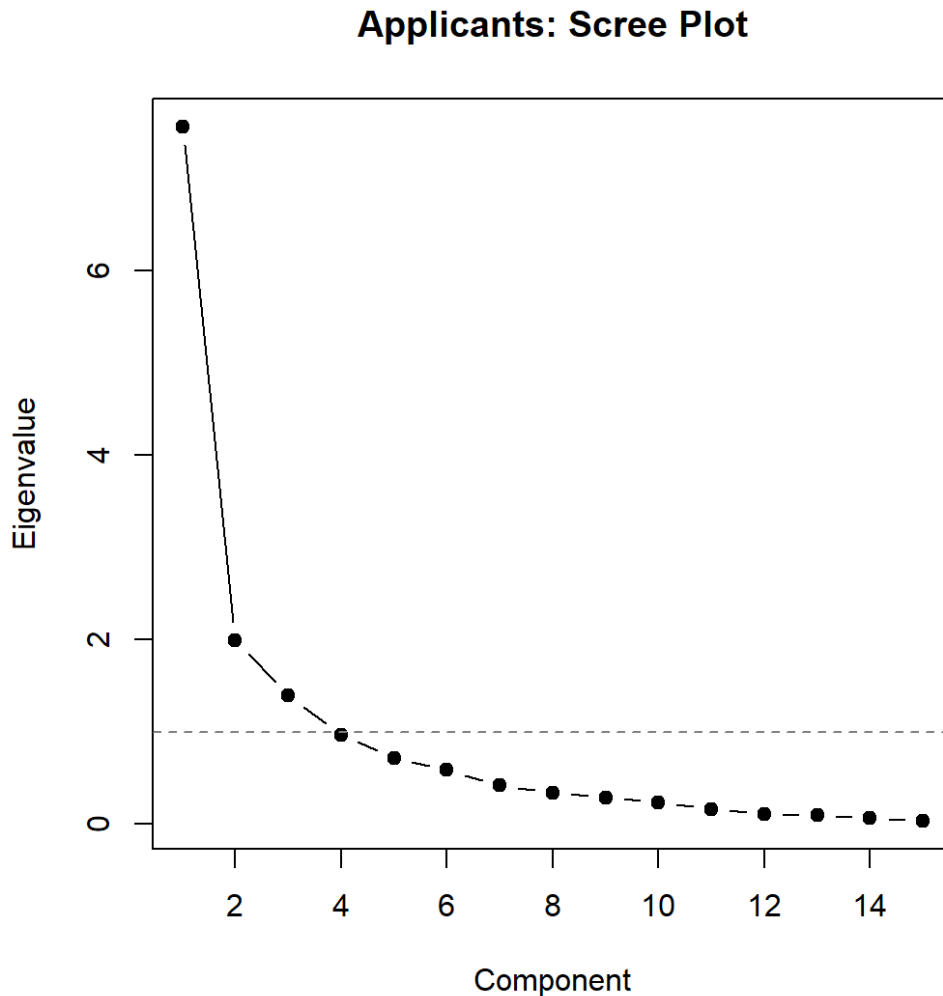


图 1.4: 申请者评分数据的碎石图。

(c) [2 分] 给出因子旋转之后的因子载荷矩阵。

【解】

```
load_tab2 <- data.frame(L2,
                        h2 = rowSums(L2^2),
                        u2 = 1 - rowSums(L2^2))
round(load_tab2, 3)
colSums(L2^2)
colSums(L2^2) / ncol(X2)
```

旋转后的因子载荷矩阵如下：

表 1.2: 申请者评分数据旋转后的因子载荷矩阵.

变量	F_1	F_2	F_3	h^2	u^2
FL	0.070	0.586	0.120	0.362	0.638
APP	0.468	0.285	0.279	0.377	0.623
AA	0.177	0.164	0.031	0.059	0.941
LA	0.175	0.224	0.956	0.995	0.005
SC	0.915	-0.101	0.172	0.877	0.123
LC	0.828	0.114	0.326	0.806	0.194
HON	0.241	-0.182	0.674	0.545	0.455
SMS	0.874	0.264	0.156	0.858	0.142
EXP	0.095	0.782	-0.053	0.623	0.377
DRV	0.765	0.380	0.183	0.763	0.237
AMB	0.890	0.159	0.162	0.843	0.157
GSP	0.790	0.276	0.318	0.801	0.199
POT	0.717	0.359	0.418	0.818	0.182
KJ	0.428	0.253	0.578	0.581	0.419
SUIT	0.350	0.845	0.080	0.842	0.158

三个因子的方差贡献分别为 36.41%、15.94%、15.33%，累计贡献约为 67.67%.

(d) [2 分] 作变量在公共因子平面上的散点图，对公共因子作出解释.

【解】

```
plot(L2[, 1], L2[, 2], pch = 19,
      xlab = "Factor 1", ylab = "Factor 2",
      main = "Applicants: Variable Loadings")
abline(h = 0, v = 0, col = "gray70")
text(L2[, 1], L2[, 2], labels = rownames(L2), pos = 4, cex = 0.8)
```

由载荷矩阵可知， F_1 在 SC、LC、SMS、DRV、AMB、GSP、POT 上载荷较高，主要反映自信心、洞察力、推销能力、驾驶水平、事业心、理解能力和潜在能力，可解释为“综合职业能力与进取性因子”。 F_2 在 SUIT、EXP、FL 上载荷较高，同时 DRV、POT 也有一定正载荷，可解释为“经验与适应表达因子”。 F_3 在 LA、HON、KJ 上载荷较高，可解释为“人际评价与品格因子”。公共因子平面图展示的是 F_1 与 F_2 ，其中 SC、SMS、AMB、LC、GSP 等变量集中在 F_1 高载荷方向，SUIT、EXP、FL 集中在 F_2 高载荷方向。

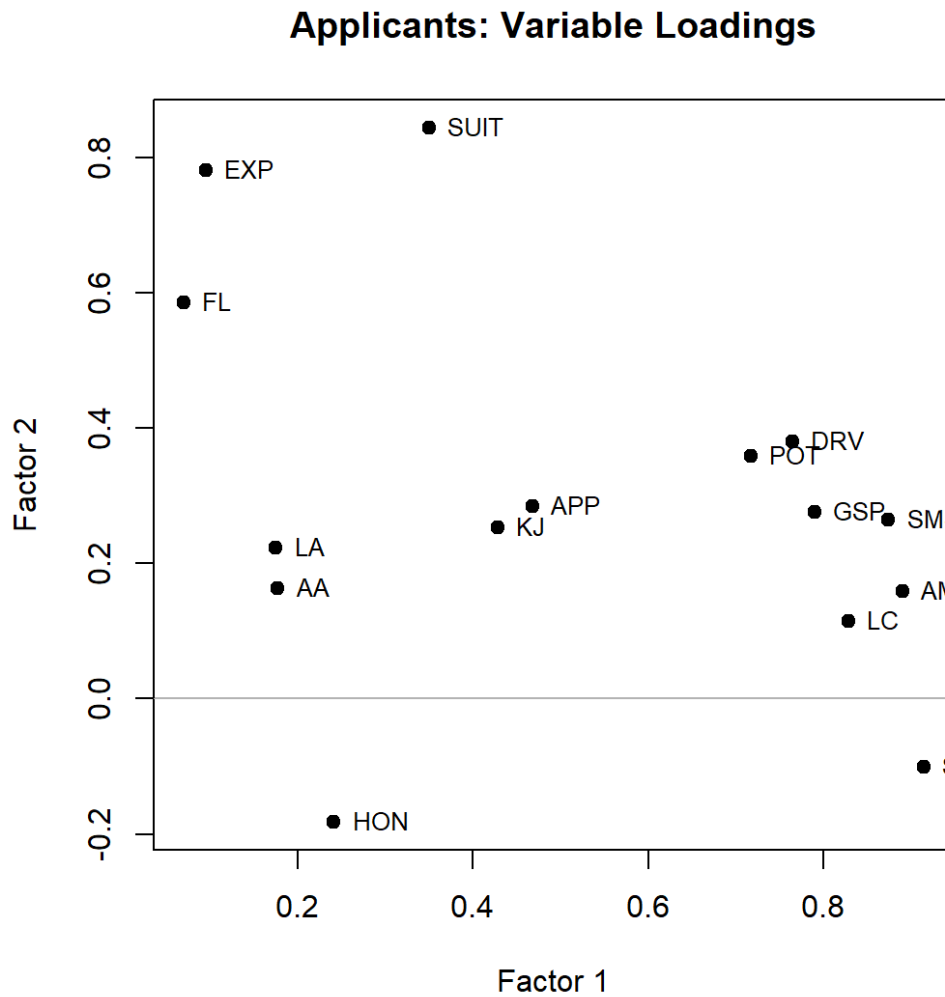


图 1.5: 申请者评分变量在第 1、第 2 公共因子平面上的散点图.

(e) [2 分] 该公司准备录用其中 6 人, 利用公共因子得分对该公司的录用结果给出建议.

【解】 以三个公共因子旋转后的方差贡献作为权数:

$$w_1 = 0.5380, \quad w_2 = 0.2355, \quad w_3 = 0.2265.$$

综合得分为

$$F = 0.5380F_1 + 0.2355F_2 + 0.2265F_3.$$

```
score_dat2 <- data.frame(ID = app_dat$ID, scores2)
ss2 <- colSums(L2^2)
w2 <- ss2 / sum(ss2)
score_dat2$Composite <- as.vector(as.matrix(score_dat2[, c("F1", "F2", "F3")]) %*% w2)

rank_app <- score_dat2[order(-score_dat2$Composite),
                        c("ID", "F1", "F2", "F3", "Composite")]
```

```

rank_app$Rank <- seq_len(nrow(rank_app))
rank_app <- rank_app[, c("Rank", "ID", "F1", "F2", "F3", "Composite")]
round(head(rank_app, 12), 3)

plot(score_dat2$F1, score_dat2$F2, pch = 19,
      xlab = "Factor score 1", ylab = "Factor score 2",
      main = "Applicants: Factor Scores")
abline(h = 0, v = 0, col = "gray70")
text(score_dat2$F1, score_dat2$F2,
      labels = score_dat2$ID, pos = 4, cex = 0.65)

```

综合得分前 12 名如下:

Rank	ID	F1	F2	F3	Composite
1	40	0.919	1.226	0.975	1.004
2	39	0.782	1.259	0.995	0.942
3	2	1.076	0.755	0.312	0.827
4	8	0.874	1.189	0.251	0.807
5	23	1.023	0.242	0.821	0.793
6	22	1.079	0.280	0.435	0.745
7	24	0.633	0.375	1.231	0.708
8	7	0.615	1.298	0.271	0.698
9	10	1.955	0.247	-1.945	0.670
10	9	0.352	1.176	0.350	0.546
11	16	0.545	0.985	-0.016	0.522
12	3	0.867	0.499	-0.324	0.510

运行上述代码即可得到 48 名申请者在第 1、第 2 公共因子平面上的散点图.

Applicants: Factor Scores

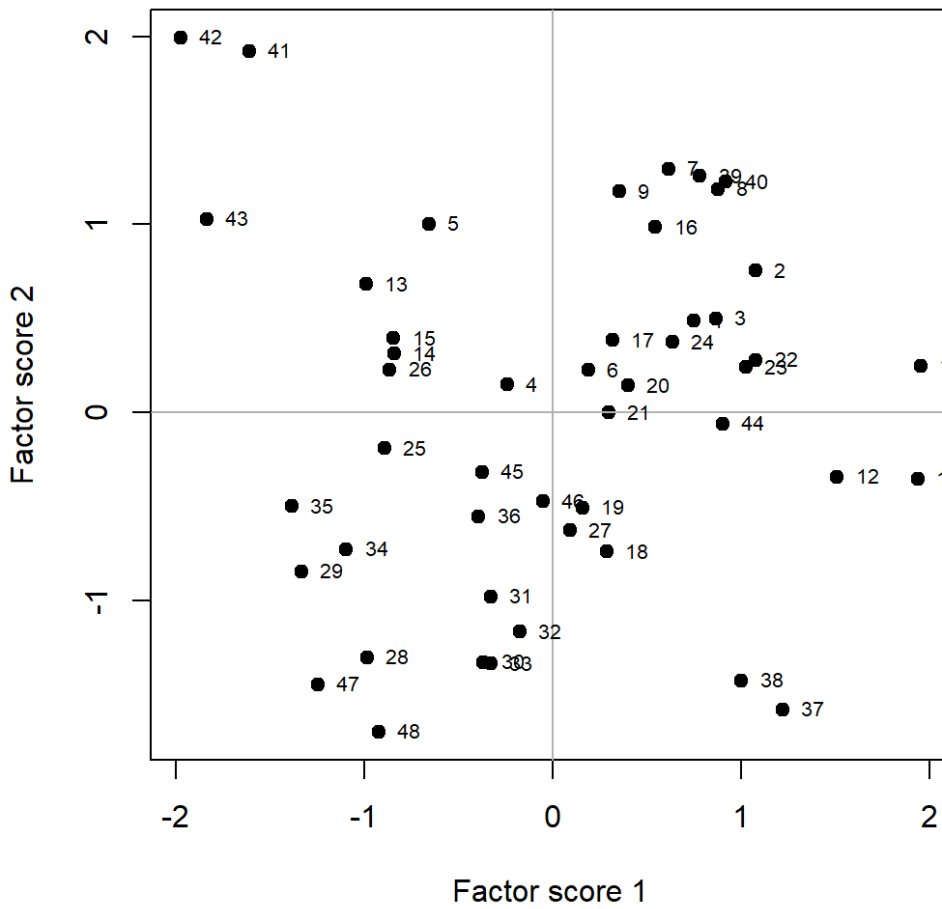


图 1.6: 48 名申请者在第 1、第 2 公共因子平面上的散点图

因此建议录用综合得分排名前 6 的申请者:

40, 39, 2, 8, 23, 22.

其中 40 号和 39 号在三个因子上表现都较好, 综合优势最明显; 2 号、8 号、23 号、22 号在综合职业能力与进取性因子上得分较高, 也具有较好的录用价值. 24 号、7 号可作为候补人选. 需要注意的是, 10 号虽然 F_1 得分很高, 但 F_3 得分较低, 说明其职业能力指标突出而人际评价与品格类指标相对不足, 综合排序未进入前 6.