

习题 4.3

a. 画 Age(y) 和 Diameter(x) 的散点图

散点图的代码和结果如下所示：

```
trees = read.csv('trees.csv')
plot(trees$diam, trees$age, pch=19, xlab='Diameter, x', ylab='Age, y')
```

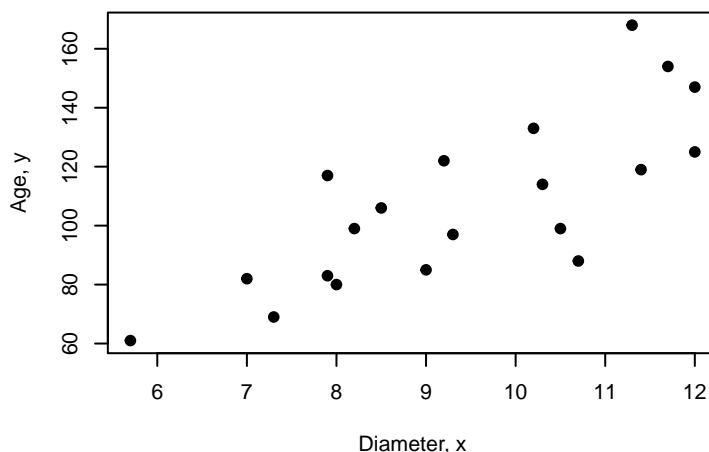


图 1: 树龄-直径散点图

b. 树龄总体均值的比估计和标准误

树龄总体均值的比估计为

$$\bar{y}_r = \hat{B}\bar{x}_U = \frac{\bar{y}}{\bar{x}}\bar{x}_U$$

估计量的标准误为

$$SE(\bar{y}_r) = \sqrt{\hat{V}(\bar{y}_r)} = \sqrt{\frac{1-f}{n} \left(\frac{\bar{x}_U}{\bar{x}} \right)^2 s_e^2}$$

其中

$$s_e^2 = s_y^2 - 2\hat{B}s_{yx} + \hat{B}^2s_x^2 = \frac{1}{n-1} \sum_{i \in \mathcal{S}} (y_i - \hat{B}x_i)^2$$

计算的代码和结果如下：

```
x_U = 10.3; n = 20; N = 1132
B_hat = mean(trees$age) / mean(trees$diam)
y_r = B_hat * x_U # 树龄总体均值的比估计
s2_e = var(trees$age - B_hat * trees$diam)
SE_ratio = sqrt((1-n/N)/n * (x_U/mean(trees$diam))^2 * s2_e) # 比估计的标准误
```

得树龄总体均值的比估计为 $\bar{y}_r = 117.620$ ，标准误为 $SE(\bar{y}_r) = 4.355$ ，其中 $s_e^2 = 321.933$ 。

c1. 树龄总体均值的回归估计和标准误

对样本数据采用最小二乘法估计系数 \hat{B}_1 ，则树龄总体均值的回归估计为

$$\hat{y}_{reg} = \bar{y} + \hat{B}_1(\bar{x}_U - \bar{x})$$

估计量的标准误为

$$SE(\hat{y}_{reg}) = \sqrt{\hat{V}(\hat{y}_{reg})} = \sqrt{\frac{1-f}{n} s_e^2}$$

其中

$$s_e^2 = \frac{1}{n-2} \sum_{i \in \mathcal{S}} e_i^2 = \frac{1}{n-2} \sum_{i \in \mathcal{S}} [y_i - (\hat{B}_0 + \hat{B}_1 x_i)^2]$$

计算的代码和结果如下：

```
trees_fit = lm(age~diam, data=trees)
B0 = trees_fit$coefficients[1]; B1 = trees_fit$coefficients[2]
y_reg = mean(trees$age) + B1 * (x_U - mean(trees$diam)) # 树龄总体均值的回归估计
e = trees_fit$residuals
s2_e_reg = sum(e^2) / (n-2)
SE_reg = sqrt((1-n/N)/n * s2_e_reg) # 回归估计的标准误
```

得树龄总体均值的回归估计为 $\bar{y}_{reg} = 118.363$ ，其中 $\hat{B}_1 = 12.250$ 。标准误为 $SE(\hat{y}_{reg}) = 4.071$ ，其中 $s_e^2 = 337.385$ 。

c2. 树龄总体均值的简单估计和标准误

计算的代码和结果如下：

```
y_bar = mean(trees$age)
SE_simple = sqrt((1-n/N)/n * var(trees$age))
```

得树龄总体均值的简单估计为 $\bar{y} = 107.40$ ，标准误为 $SE(\bar{y}) = \sqrt{\frac{1-f}{n} s_y^2} = 6.352$ 。

d. 三种估计量的比较

将三种估计量标在图 1 的散点图中，结果如图 2 所示：

```
library(latex2exp)
plot(trees$diam, trees$age, pch=20, col='red', xlab='Diameter, x', ylab='Age, y')
abline(a=B0, b=B1, col='blue')
abline(h=y_bar, lty=3); abline(h=y_r, lty=3); abline(h=y_reg, lty=3)
text(x=5.7, y=103, labels=TeX("$\\bar{y}=107.40$"), cex=0.8, adj=0)
text(x=5.7, y=113, labels=TeX("$\\bar{y}_r=117.62$"), cex=0.8, adj=0)
text(x=5.7, y=125, labels=TeX("$\\bar{y}_{reg}=118.36$"), cex=0.8, adj=0)
abline(v=mean(trees$diam), lty=3); abline(v=x_U, lty=3)
```

```
text(x=9.7, y=155, labels=TeX("$\\bar{x}=9.4$"), cex=0.8)
text(x=10.7, y=155, labels=TeX("$\\bar{x}_U=10.3$"), cex=0.8)
```

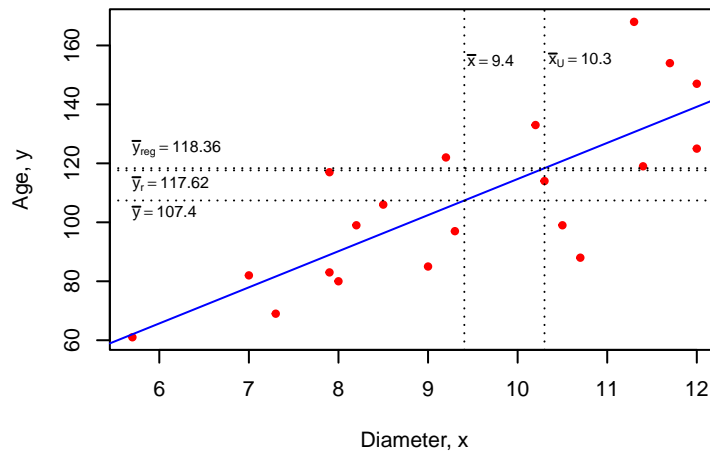


图 2: 树龄-直径散点图及估计量

结论: 由树龄-直径散点图可知, 树龄和直径间存在较好的正比例关系, 并且样本量为 20, \bar{x}_U 容易测量, 所以比估计和回归估计优于简单估计。而回归估计的标准误差略小于比估计的标准误差, 但是样本容量不太大, 回归估计的偏倚可能大于比估计的偏倚, 并且树龄和直径间的线性关系更接近正比例关系, 所以, 本题可以直接使用比估计。

习题 4.10

a. 画 volume vs. diameter 的散点图

散点图的代码和结果如下所示:

```
cherry = read.csv('cherry.csv')
diam = cherry$diameter; height = cherry$height; vol = cherry$volume
cherry_fit = lm(vol~diam)
B0 = cherry_fit$coefficients[1]; B1 = cherry_fit$coefficients[2]
plot(diam, vol, pch=20, col='red', xlab='Diameter', ylab='Volume', xlim=c(0,21), ylim=c(0,80))
abline(a=B0, b=B1, col='blue')
```

由图 3 的 Volume-Diameter 散点图可知, Volume 和 Diameter 之间有较明显的线性关系。

b. Volume 总体总和的比估计及 95% 置信区间

Volume 总体总和的比估计为

$$\hat{t}_{yr} = \hat{B}t_x = \frac{\bar{y}}{\bar{x}}t_x$$

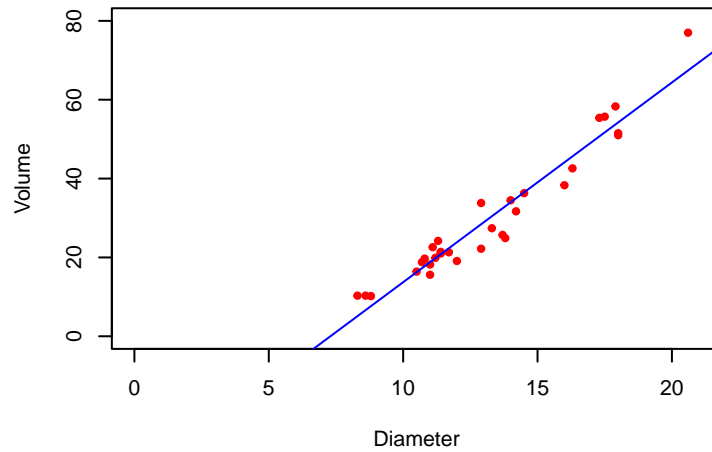


图 3: Volume-Diameter 散点图

比估计的标准误差为

$$SE(\hat{t}_{yr}) = \sqrt{\hat{V}(\hat{t}_{yr})} = \sqrt{\frac{1-f}{n} \left(\frac{t_x}{\bar{x}}\right)^2 s_e^2}$$

其中

$$s_e^2 = s_y^2 - 2\hat{B}s_{yx} + \hat{B}^2 s_x^2 = \frac{1}{n-1} \sum_{i \in \mathcal{S}} (y_i - \hat{B}x_i)^2$$

则 Volume 总体总和的 95% 置信区间为

$$[\hat{t}_{yr} - z_{\alpha/2} SE(\hat{t}_{yr}), \hat{t}_{yr} + z_{\alpha/2} SE(\hat{t}_{yr})]$$

计算的代码和结果如下所示:

```
N = 2967; n = 31; t_x = 41835; y_bar = mean(vol); x_bar = mean(diam)
B_hat = y_bar / x_bar
t_yr = B_hat * t_x # volume 总体总和的比估计
s2_e = var(vol - B_hat * diam)
SE_ratio = sqrt((1/n-1/N) * (t_x/x_bar)^2 * s2_e)
CI_lb_ratio = t_yr - qnorm(0.975) * SE_ratio # volume 总体总和的 95% 置信下限
CI_ub_ratio = t_yr + qnorm(0.975) * SE_ratio # volume 总体总和的 95% 置信上限
```

得 Volume 总体总和的比估计为 $\hat{t}_{yr} = 95272.159$, 95% 置信区间为 $[84548.344, 105995.973]$. 其中, $s_e^2 = 94.053$, $SE(\hat{t}_{yr}) = 5471.434$.

c. Volume 总体总和的回归估计及 95% 置信区间

对样本数据采用最小二乘法估计系数 \hat{B}_1 , 则 Volume 总体总和的回归估计为

$$\hat{t}_{y,reg} = N\hat{\bar{y}}_{reg} = \hat{t}_y + \hat{B}_1(t_x - \hat{t}_x)$$

估计量的标准误为

$$SE(\hat{t}_{y,reg}) = N \cdot SE(\hat{\bar{y}}_{reg}) = N \sqrt{\hat{V}(\hat{\bar{y}}_{reg})} = N \sqrt{\frac{1-f}{n} s_e^2}$$

其中

$$s_e^2 = \frac{1}{n-2} \sum_{i \in \mathcal{S}} e_i^2 = \frac{1}{n-2} \sum_{i \in \mathcal{S}} [y_i - (\hat{B}_0 + \hat{B}_1 x_i)^2]$$

则 Volume 总体总和的 95% 置信区间为

$$[\hat{t}_{y,reg} - z_{\alpha/2} SE(\hat{t}_{y,reg}), \hat{t}_{y,reg} + z_{\alpha/2} SE(\hat{t}_{y,reg})]$$

```
t_y_hat = y_bar * N; t_x_hat = x_bar * N
t_yreg = t_y_hat + B1 * (t_x - t_x_hat) # volume 总体总和的回归估计
s2_e = sum(cherry_fit$residuals^2) / (n-2)
SE_reg = N * sqrt((1/n-1/N) * s2_e)
CI_lb_reg = t_yreg - qnorm(0.975) * SE_reg # volume 总体总和的 95% 置信下限
CI_ub_reg = t_yreg + qnorm(0.975) * SE_reg # volume 总体总和的 95% 置信上限
```

计算得 Volume 总体总和的回归估计为 $\hat{t}_{y,reg} = 102318.90$ ，其中 $\hat{B}_1 = 5.066$ 。标准误为 $SE(\hat{t}_{y,reg}) = 2253.969$ ，其中 $s_e^2 = 18.079$ ，则 Volume 总体总和的 95% 置信区间为 $[97901.16, 106736.60]$ 。

d. Volume 总体总和的简单估计、95% 置信区间及三种估计的比较

```
t_y_hat = N * y_bar # volume 总体总和的简单估计
SE_simple = N * sqrt((1/n-1/N) * var(vol))
CI_lb_simple = t_y_hat - qnorm(0.975) * SE_simple # volume 总体总和的 95% 置信下限
CI_ub_simple = t_y_hat + qnorm(0.975) * SE_simple # volume 总体总和的 95% 置信上限
```

计算得 Volume 总体总和的简单估计为 $\hat{t}_y = N\bar{y} = 89517.261$ ，标准误为 $SE(\hat{t}_y) = 8713.665$ ，则 95% 置信区间为 $[72438.791, 106595.731]$ 。

三种估计量的比较结果如表 1 所示：

表 1: Volume 总体总和的三种估计量比较

Method	Point Estimate	Standard Error	95% CI Lower Bound	95% CI Upper Bound
Simple Estimation	89517.261	8713.665	72438.791	106595.731
Ratio Estimation	95272.159	5471.434	84548.344	105995.973
Regression Estimation	102318.90	2253.969	97901.16	106736.60

结论：由图 3 的 Volume-Diameter 散点图可知，这两个变量间有较明显的线性关系，且回归直线的截距明显不为 0。此外，Diameter 容易测量，其总体总和易得，并且样本容量（为 31）足够大，比较适合用回归估计。计算结果也表明回归估计的标准误明显小于比估计，且比估计的标准误明显小于简单估计。

e. 有多个辅助变量时的回归估计

对样本数据采用最小二乘法估计系数，则 Volume 总体总和的回归估计为

$$\hat{t}_{y,reg} = N\hat{y}_{reg} = N(\hat{B}_0 + \hat{B}_1\bar{x}_{1U} + \hat{B}_2\bar{x}_{2U}) = N\hat{B}_0 + \hat{B}_1t_{x_1} + \hat{B}_2t_{x_2}$$

估计量的标准误为

$$SE(\hat{t}_{y,reg}) = N \cdot SE(\hat{y}_{reg}) = N\sqrt{\hat{V}(\hat{y}_{reg})} = N\sqrt{\frac{1-f}{n}s_e^2}$$

其中

$$s_e^2 = \frac{1}{n-q-1} \sum_{i \in \mathcal{S}} e_i^2 = \frac{1}{n-q-1} \sum_{i \in \mathcal{S}} [y_i - (\hat{B}_0 + \hat{B}_1x_{i1} + \hat{B}_2x_{i2})^2], \quad q = 2$$

```
t_x1 = 41835; t_x2 = 240327
cherry_fit2 = lm(vol~diam+height); coef = cherry_fit2$coefficients
t_yreg2 = N * coef[1] + coef[2] * t_x1 + coef[3] * t_x2
s2_e = sum(cherry_fit2$residuals^2) / (n-3)
SE_reg2 = N * sqrt((1/n-1/N) * s2_e)
```

得 Volume 总体总和的回归估计为 $\hat{t}_{y,reg} = 106447.70$ ，标准误为 $SE(\hat{t}_{y,reg}) = 2057.75$ ，其中 $s_e^2 = 15.07$ 。

习题 4.11

a. Physicians 数量的直方图

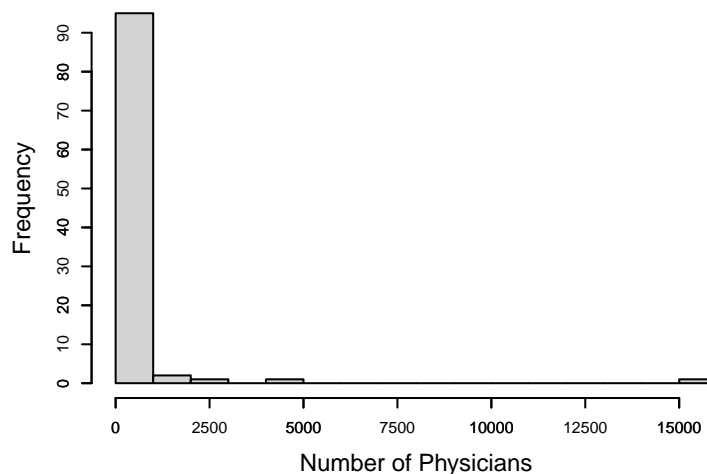


图 4: Physicians 数量直方图

b. Physicians 数量总体总和的简单估计及标准误

```

county = read.csv("counties.csv"); physi = county$physician
N = 3141; n = 100; y_bar = mean(physi)
t_y_hat = N * y_bar # Physicians 数量总体总和的简单估计
SE_simple = N * sqrt((1/n-1/N) * var(physi))

```

计算得 Physicians 数量总体总和的简单估计为 $\hat{t}_y = N\bar{y} = 933411$ ，标准误为 $SE(\hat{t}_y) = N\sqrt{\frac{1-f}{n}s_y^2} = 491982.787$ 。

c. Physicians vs. Population 散点图

```

popu = county$totpop/1000000
county_fit = lm(physi~popu)
B0 = county_fit$coefficients[1]; B1 = county_fit$coefficients[2]
plot(popu, physi, col='red', pch=20,
      xlab='Population (in 1,000,000s)', ylab='Physicians')
abline(a=B0, b=B1, col='blue')

```

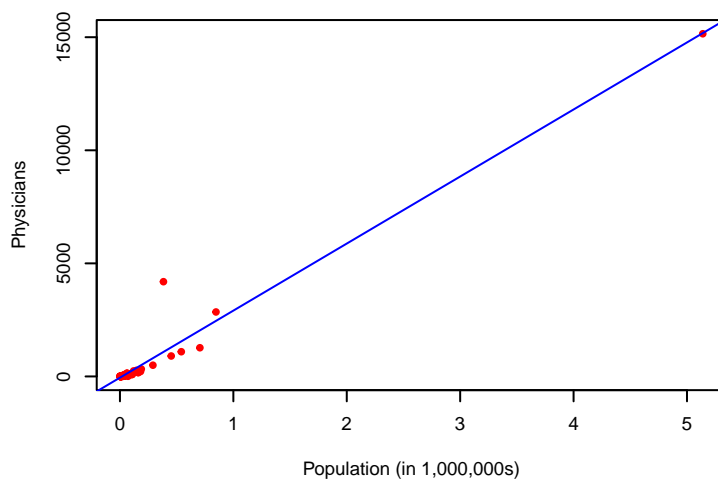


图 5: Physicians vs. Population 散点图

由图 5 可知，回归直线似乎是经过原点的，即 Physicians 和 Population 更可能是正比例关系，所以比估计更合适。

d. Physicians 数量总体总和的比估计和回归估计

d.1 比估计及其标准误

Physicians 总体总和的比估计为

$$\hat{t}_{yr} = \hat{B}t_x = \frac{\bar{y}}{\bar{x}}t_x$$

比估计的标准误为

$$SE(\hat{t}_{yr}) = \sqrt{\hat{V}(\hat{t}_{yr})} = \sqrt{\frac{1-f}{n} \left(\frac{t_x}{\bar{x}}\right)^2 s_e^2}$$

其中

$$s_e^2 = s_y^2 - 2\hat{B}s_{yx} + \hat{B}^2 s_x^2 = \frac{1}{n-1} \sum_{i \in \mathcal{S}} (y_i - \hat{B}x_i)^2$$

计算的代码和结果如下所示：

```
t_x = 255077536; popu = county$totpop; x_bar = mean(popu)
B_hat = y_bar / x_bar
t_yr = B_hat * t_x # physicians 总体总和的比估计
s2_e = var(physi - B_hat * popu)
SE_ratio = sqrt((1/n-1/N) * (t_x/x_bar)^2 * s2_e) # 比估计的标准误
```

得 Physicians 总体总和的比估计为 $\hat{t}_{yr} = 639506$ ，标准误为 $SE(\hat{t}_{yr}) = 87885.27$ ，其中 $s_e^2 = 172267.87$ 。

d.2 回归估计及其标准误

对样本数据采用最小二乘法估计系数 \hat{B}_1 ，则 Physicians 总体总和的回归估计为

$$\hat{t}_{y,reg} = N\hat{y}_{reg} = \hat{t}_y + \hat{B}_1(t_x - \hat{t}_x)$$

估计量的标准误为

$$SE(\hat{t}_{y,reg}) = N \cdot SE(\hat{y}_{reg}) = N\sqrt{\hat{V}(\hat{y}_{reg})} = N\sqrt{\frac{1-f}{n} s_e^2}$$

其中

$$s_e^2 = \frac{1}{n-2} \sum_{i \in \mathcal{S}} e_i^2 = \frac{1}{n-2} \sum_{i \in \mathcal{S}} [y_i - (\hat{B}_0 + \hat{B}_1 x_i)]^2$$

计算代码和结果如下：

```
physician_fit = lm(physi~popu)
B1 = physician_fit$coefficients[2]
t_y_hat = y_bar * N; t_x_hat = x_bar * N
t_yreg = t_y_hat + B1 * (t_x - t_x_hat) # physicians 总体总和的回归估计
s2_e = sum(physician_fit$residuals^2) / (n-2)
SE_reg = N * sqrt((1/n-1/N) * s2_e) # 回归估计的标准误
```

计算得 Physicians 总体总和的回归估计为 $\hat{t}_{y,reg} = 585871$ ，其中 $\hat{B}_1 = 0.00296$ 。标准误为 $SE(\hat{t}_{y,reg}) = 105177.418$ ，其中 $s_e^2 = 115813.892$ 。

e. 哪种估计方法更接近真值

Physicians 总体总和的比估计为 639,506，回归估计为 585,871，而真值为 532,638，则回归估计更加接近真值。

习题 4.19

1. 生成总体

令 $x_i \sim \text{Exp}(10)$, $\varepsilon_i \sim N(0, \sigma^2 x_i)$, 取 $\beta = 2$, $\sigma^2 = 1$, 有 $y_i = \beta x_i + \varepsilon_i$. 则 $\{y_i; i = 1, \dots, 1000\}$ 为总体. 生成总体的代码如下所示:

```
set.seed(123)
x_pop = rexp(1000, rate=10); beta = 2
y_pop = c() # 初始化 y_i
for (i in 1:1000){
  epsilon_i = rnorm(1, mean=0, sd=sqrt(x_pop[i]))
  y_i = beta * x_pop[i] + epsilon_i
  y_pop = append(y_pop, y_i)
}
```

可以画散点图来观察 y 围绕趋势线 $y = \beta x$ 的波动. 由下图可知, 波动程度随着 x 的增大而增大:

```
plot(x_pop, y_pop, col='red', pch=20, xlab='x', ylab='y', cex.lab=1.3)
abline(a=0, b=2, col='blue'); text(0.6, 0.9, 'y = 2x', col='blue')
```

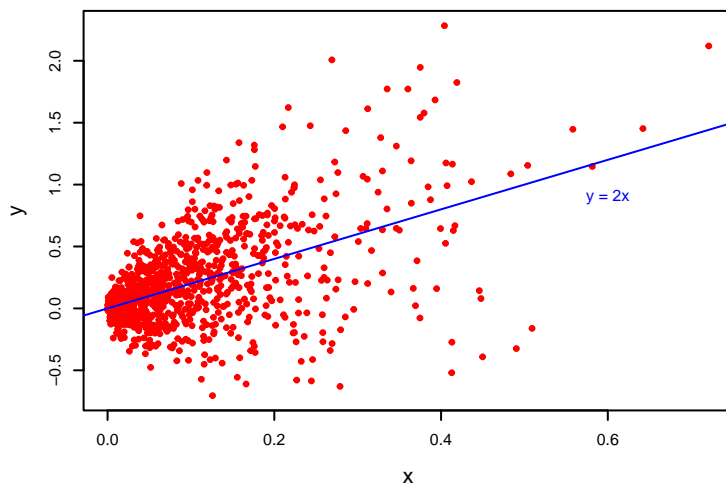


图 6: y-x 散点图

2. 比较两个方差估计量

我们通常使用 $\hat{V}_1(\hat{y}_r)$ 和 $\hat{V}_2(\hat{y}_r)$ 来估计 $V(\hat{y}_r)$, 有

$$\hat{V}_1(\hat{y}_r) = \frac{1-f}{n} \left(\frac{\bar{x}_U}{\bar{x}} \right)^2 s_e^2, \quad \hat{V}_2(\hat{y}_r) = \frac{1-f}{n} s_e^2$$

其中

$$s_e^2 = s_y^2 - 2\hat{B}s_{yx} + \hat{B}^2 s_x^2 = \frac{1}{n-1} \sum_{i \in \mathcal{S}} (y_i - \hat{B}x_i)^2$$

每次从总体中抽取容量为 50 的 SRS，共抽取 100 次，可根据上述公式计算得每次的 $\hat{V}_1(\hat{y}_r)$ 和 $\hat{V}_2(\hat{y}_r)$ ，画出两者的抽样分布图并进行比较。

```
set.seed(321); n=50; N=1000; x_U = mean(x_pop)
V1_hat = c(); V2_hat = c() # 初始化两种方差估计量
for (i in 1:100){
  S = sample(1:1000, n); x_sample = x_pop[S]; y_sample = y_pop[S]
  x_bar = mean(x_sample); y_bar = mean(y_sample); B_hat = y_bar / x_bar
  s2_e = var(y_sample - B_hat * x_sample)
  V1_temp = (1/n-1/N) * (x_U/x_bar)^2 * s2_e; V2_temp = (1/n-1/N) * s2_e
  V1_hat = append(V1_hat, V1_temp); V2_hat = append(V2_hat, V2_temp)
}

par(mfrow=c(1,2))
hist(V1_hat, freq=F, nclass=15, ylim=c(0,900),
     main=TeX("Histogram of  $\hat{V}_1(\hat{y}_r)$ ", bold=T), xlab=NA)
d1 = seq(min(V1_hat), max(V1_hat), by=0.0001)
lines(x=d1, y=dnorm(d1, mean(V1_hat), sd(V1_hat)), col='red', lty=4)
lines(density(V1_hat), col='blue', lty=4)

hist(V2_hat, freq=F, nclass=15, ylim=c(0,900),
     main=TeX("Histogram of  $\hat{V}_2(\hat{y}_r)$ ", bold=T), xlab=NA)
d2 = seq(min(V2_hat), max(V2_hat), by=0.0001)
lines(x=d2, y=dnorm(d2, mean(V2_hat), sd(V2_hat)), col='red', lty=4)
lines(density(V2_hat), col='blue', lty=4)
legend('right', c('normal', 'kdensity'), lty=c(4,4), col=c('red', 'blue'), cex=0.8)
```

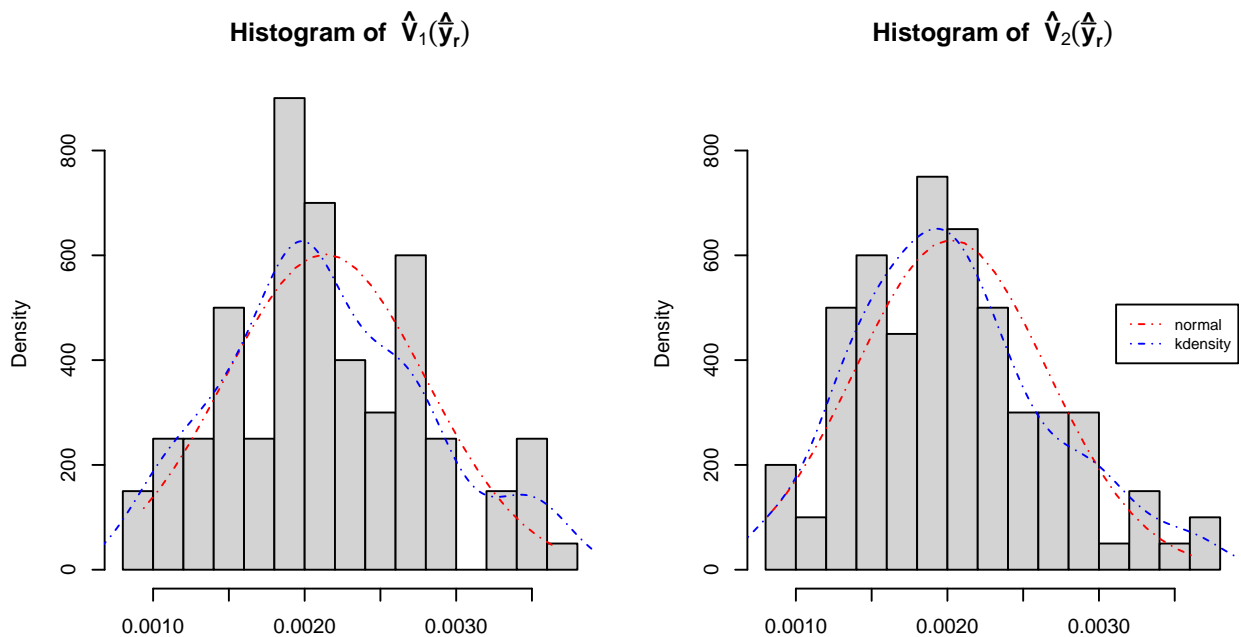


图 7: 两种方差估计量的比较

结论：由图 7 的直方图可以看出， $\hat{V}_1(\hat{y}_r)$ 的抽样分布更接近正态，且更加集中在均值附近，而 $\hat{V}_2(\hat{y}_r)$ 的抽样分布有一定的右偏。所以， $\hat{V}_1(\hat{y}_r)$ 是更好的方差估计量。

习题 4.38

1. 模拟样本 (x_i, y_i)

假设一本英语字典共 7,000 页，收录的英语单词共 28,000 个，则 $N = 7000$ ， $n = 30$ ， $t_x = 28000$ 为已知。假设抽出的样本中，每页的单词数 $x_i \sim U(20, 60)$ ，每页认识的单词数 $y_i = (\beta + \varepsilon_i)x_i$ ，其中 $\beta = 0.4$ ， $\varepsilon_i \sim N(0, 0.05)$ ， x_i, y_i 均为非负整数。用此方法生成容量为 30 的样本：

```
set.seed(1111)
N=7000; n=30; t_x=28000; beta = 0.4
x_sample = 20 + floor(41*runif(30)); epsi = rnorm(30,0,0.05)
y_sample = round((beta + epsi) * x_sample)
plot(x_sample,y_sample,pch=20,col='red',
      xlab='Words per Page',ylab='Known Words per Page')
```

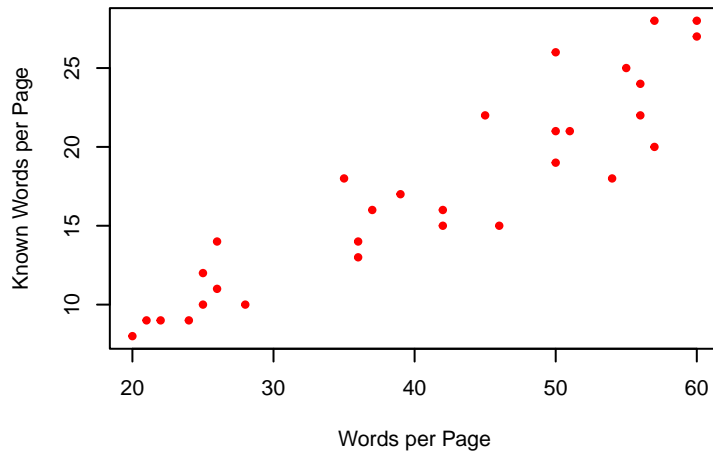


图 8: 每页认识单词数-每页单词数散点图

2. 认识单词总数的比估计

认识单词总数的比估计为

$$\hat{t}_{yr} = \hat{B}t_x = \frac{\bar{y}}{\bar{x}}t_x$$

比估计的标准误为

$$SE(\hat{t}_{yr}) = \sqrt{\hat{V}(\hat{t}_{yr})} = \sqrt{\frac{1-f}{n} \left(\frac{t_x}{\bar{x}} \right)^2 s_e^2}$$

其中

$$s_e^2 = s_y^2 - 2\hat{B}s_{yx} + \hat{B}^2 s_x^2 = \frac{1}{n-1} \sum_{i \in \mathcal{S}} (y_i - \hat{B}x_i)^2$$

计算的代码和结果如下所示:

```
x_bar = mean(x_sample); y_bar = mean(y_sample); B_hat = y_bar / x_bar
t_yr = B_hat * t_x # 认识单词总数的比估计
s2_e = var(y_sample - B_hat * x_sample)
SE_ratio = sqrt((1/n-1/N) * (t_x/x_bar)^2 * s2_e) # 比估计的标准误
```

得认识单词总数的比估计为 $\hat{t}_{yr} = 11759$, 标准误为 $SE(\hat{t}_{yr}) = 302.748$, 其中 $s_e^2 = 5.931$.

3. 认识单词比例的比估计

认识单词比例的比估计为

$$\hat{B} = \frac{\bar{y}}{\bar{x}}$$

比估计的标准误为

$$SE(\hat{B}) = \sqrt{\hat{V}_1(\hat{B})} = \sqrt{\frac{1-f}{n\bar{x}^2} s_e^2}$$

```
B_hat = y_bar / x_bar
SE_B = sqrt((1/n-1/N)/(n*x_bar^2) * s2_e)
```

计算得认识单词比例的比估计为 $\hat{B} = 0.420$, 标准误为 $SE(\hat{B}) = 0.00197$.