

Flag variable: two-sample Z-test

(1)

Multinomial variable: χ^2 -test

Continuous variable: two-sample t-test

$H_0: \mu_1 = \mu_2$ $H_a: \mu_1 \neq \mu_2$: two-sample (2 t-test

$$t_{data} = \frac{\bar{n}_1 - \bar{n}_2}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}} = \frac{0.1}{\sqrt{\frac{(5.2)^2}{2000} + \frac{(4.9)^2}{600}}} = \frac{0.1}{\sqrt{0.01352 + 0.04001}}$$

$$\frac{0.1}{\sqrt{0.0535}} = 0.43$$

$$P\text{-value} = 2P(t_{599} > 0.43) = 0.66$$

با مقدار p یا $p\text{-value}$ ($0.1 < 0.66$) ، می توان نتیجه گرفت که

H_0 قویاً برقرار است و train و test از یکدیگر مستقلند ، در نتیجه

partition مذکور valid است .

$H_0: P_{\text{mixed train}} = P_{\text{mixed test}}$
 $P_{\text{single train}} = P_{\text{single test}}$
 $P_{\text{other train}} = P_{\text{other test}}$

$H_a: \text{test is different from } H_0$
 (4)

Expected:

	Mixed	Single	Other
Trained	1040×2000 2600 13	1000×2000 2600 13	560×2000 2600 13
Test	1040×600 2600	1000×600 2600	560×600 2600 13

E:	M	S	O
Train	800	769	431
Test	240	231	129

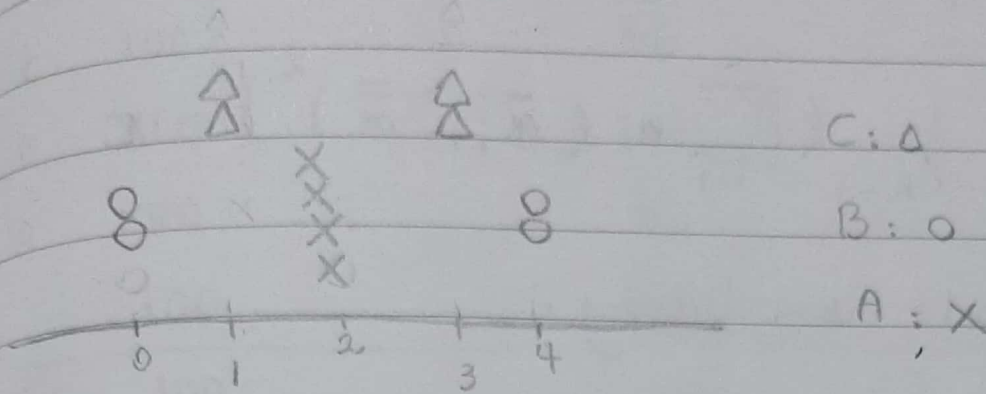
O:	M	S	O
Train	800	750	450
Test	240	250	110

$$\sum \sum \frac{(O_{ij} - E_{ij})^2}{E_{ij}} = 0 + 0 + \frac{1}{769} + \frac{(19)^2}{231} + \frac{(19)^2}{431} + \frac{(19)^2}{29} = 5.2$$

$$P\text{-value} = P(X^2_2 > 5.2) = 0.07$$

1 CS.
 valid (0.07 < α) p-value
 0.1
 sam (subset) test, train

(6) a) dotplot بر روی محور را در نظر بگیرید



همانطور که دیده می شود، میانگین هر 3 بودن با یکدیگر برابر است (2)،

اما هیچ overlap میان گروه ها نیست. عبارت دیگر، بین dotplot و تست

فرض H_0 (مستقل بودن 3 زیرگروه) صادق وجود دارد. تحت

شود که 3 زیرگروه مستقل نیستند چون توزیع ها مختلفی دارند. اگرچه میانگین

همان یکی است، اما واریانس هر subset متفاوت است.

(b) تفاوت میانگین گروه ها : Between-sample variability

واریانس داخل هر نمونه (تست)
به میانگین آن نمونه : Within-sample variability

Between-sample variability:

(c)

$$MSTR = \frac{1}{k-1} \left(\sum_{i=1}^k n_i (\bar{x}_i - \bar{\bar{x}})^2 \right)$$

Within-Sample variability:

$$MSE = \frac{1}{n_t - k} \left(\sum_{i=1}^k (n_i - 1) s_i^2 \right)$$

$$s_i^2 = \frac{1}{n_i - 1} \sum_{j=1}^{n_i} (x_{ij} - \bar{x}_i)^2$$

$$F_{data} = \frac{MSTR}{MSE}$$

(d) با توجه آماره محاسبه شده $MSTR$ و MSE می‌توانیم

از MSE بسیار بیشتر باشد، یعنی تفاوت در میانگین از تفاوت در واریانس

بسیار بیشتر باشد، آن‌گاه فرض صفر (برابری میانگین‌ها) رد می‌شود.

$H_0: \mu_{\text{credit card}} = \mu_{\text{Debit card}} = \mu_{\text{check}}$

$H_a: H_0 \text{ is wrong}$

Credit Card	Debit Card	check	(7)
100	80	50	
110	120	70	120
90	90	80	160
100	110	80	240
$\bar{x} = 100$	100	70	
$\bar{\bar{x}} = 90$			

$$MSTR = \frac{1}{k-1} \sum_{i=1}^k (n_i (\bar{x}_i - \bar{\bar{x}})^2) = \frac{2}{3} (10^2 + 10^2 + 20^2) = \boxed{1200}$$

$$MSE = \frac{1}{n_t - k} \left[\sum_{i=1}^k (n_i - 1) \left(\frac{1}{n_i - 1} \sum_{j=1}^{n_i} (x_{ij} - \bar{x}_i)^2 \right) \right]$$

$$= \frac{1}{12-3} \left[(3 \times \frac{1}{3} \times (0 + 10^2 + 10^2 + 0)) + (3 \times \frac{1}{3} \times (20^2 + 20^2 + 10^2 + 10^2)) + (3 \times \frac{1}{3} \times (20^2 + 0^2 + 10^2 + 10^2)) \right] =$$

$$\frac{1}{9} [200 + 1000 + 600] = \frac{1}{9} [1800] = \boxed{200}$$

$$\Rightarrow F = \frac{MSTR}{MSE} = 6$$

$$p\text{-value} = P(F_{(2,9)} > 6) = 0.02$$

بدین روش پویش بودن p-value از $\alpha = 0.05$ و H_0 را رد
 میکنیم (میانگین ها برابر نیستند).

$$\overline{\text{weight}}_i = -180 + 5 \times \text{height}_i$$

(10)

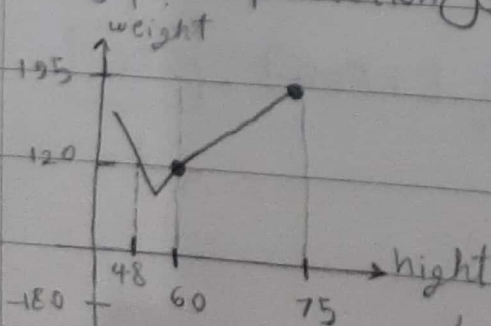
a) $5 \times 3 = 15$ pound

b) $\overline{\text{weight}}_i = -180 + 5(65) = 145$ pound

c) $\overline{\text{weight}}_i = -180 + 5(48) = 60$ pound

ما اینجا باید مراقب extrapolation باشیم، زیرا مدل در همین جا

برای height های [60 تا 75] اینجای prediction انجام میدهد.



ما واقعاً از رابطه $\text{weight} = \text{height}$

در خارج این بازه اطلاعی نداریم (مسلماً ما سه نفر نداریم)
 دارای دلالتی ندارد.

۵/۱/۱۴۰۲
(d) یا هر inch اتراس قَد، به متوسط دن 5 pound افزوده

مستود.
(دست مورد نه $60 \leq \text{height} \leq 75$)

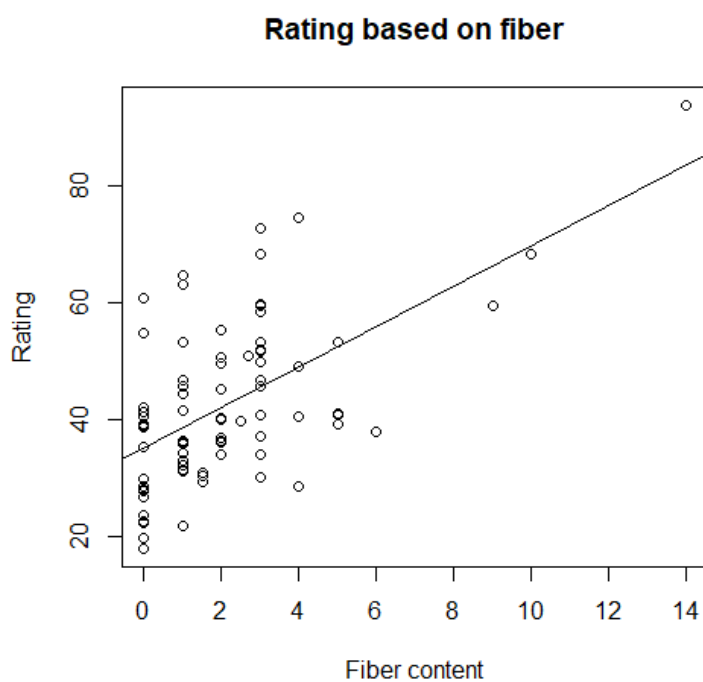
(e) عرض از مبدأ 180 - در اینجا مفروض ندارد. زیرا می توان

سقفی با $\text{height} = 0$, $\text{weight} < 0$ داشت.

با اجرای کد زیر داریم: (با توجه به coefficient ها)

```
cereals = read.csv(file.choose(), header = TRUE);
cereals[1:5, ];

fibers = cereals$fiber;
ratings = cereals$rating;
lm.out = lm(ratings~fibers)
plot(
  fibers,
  ratings,
  main = "Rating based on fiber",
  xlab = "Fiber content",
  ylab = "Rating"
)
abline(lm.out)
summary(lm.out)
```




```

Call:
lm(formula = ratings ~ fibers)

Residuals:
    Min       1Q   Median       3Q      Max
-20.436  -8.159  -2.037   6.491  27.216

Coefficients:
              Estimate Std. Error t value Pr(>|t|)
(Intercept)  35.2566     1.7674   19.948 < 2e-16 ***
fibers        3.4430     0.5524    6.233 2.45e-08 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 11.48 on 75 degrees of freedom
Multiple R-squared:  0.3412,    Adjusted R-squared:  0.3325
F-statistic: 38.85 on 1 and 75 DF,  p-value: 2.445e-08

```

The regression equation is:

$$\text{rating}[i] = 35.2566 + 3.4430 * \text{fibers}[i]$$

(13)

مقدار عرض از مبدا (در اینجا 35.2566) در این مثال معنی میدهد. معنای آن این است که اگر یک محصولی اصلاً fiber نداشته باشد، rating آن 35.2566٪ می باشد.

(15)

با توجه به مقدار R-squared در summary، که برابر 0.3 است، نمیتوان گفت که مدل ما بخوبی fit شده روی داده train.

(16)

با جاگذاری مقدار $\text{fiber} = 3$ در معادله بالا: (point estimation)
 $\text{rating}[3] = 35.2566 + 3.4430 * 3 = 45.5856$

مقدار فوق، برآورد نقطه ای میانگین rating یک محصول با 3 gram fiber است.

(17)

با اجرای کد زیر داریم:

```
> predict(lm.out, data.frame(fibers = 3), interval = "confidence")
      fit      lwr      upr
1 45.58553 42.81792 48.35314
> |
```

بازه اطمینان برای میانگین rating یک محصول با fiber = 3gram برابر [42.81792, 48.35314] می باشد. (با مرکز 45.58553، همان برآورد نقطه ای سوال قبل)

(18)

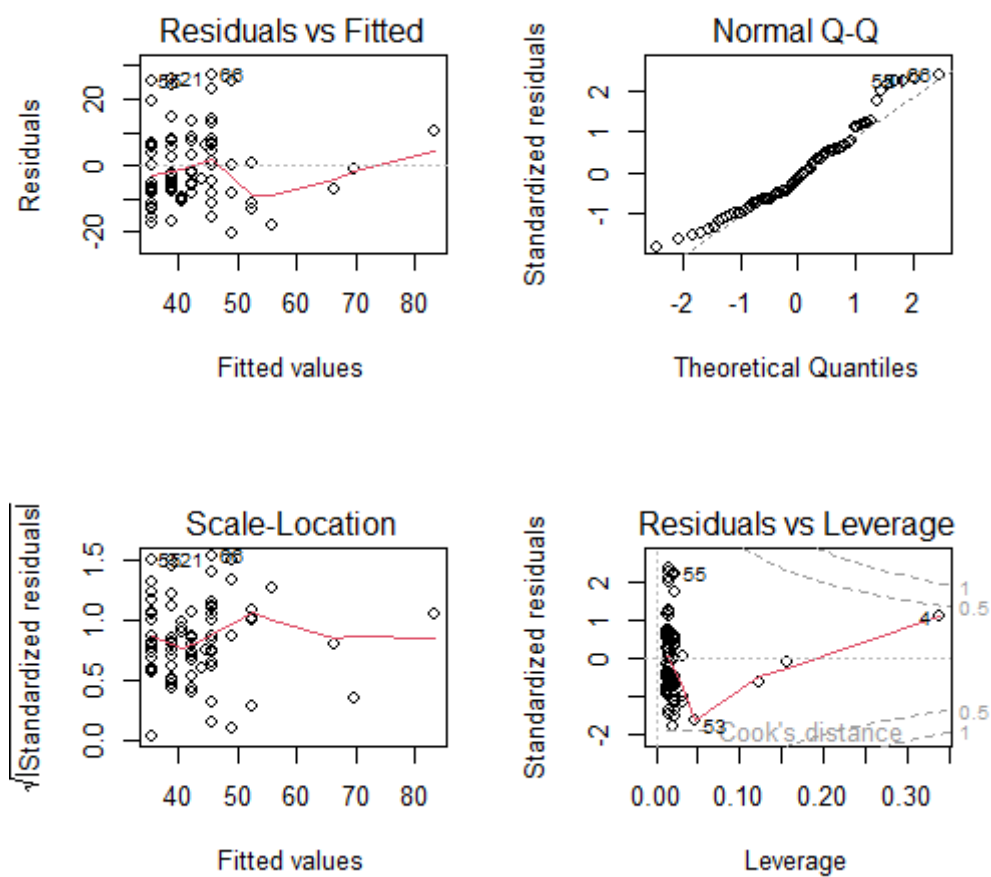
با اجرای کد زیر داریم:

```
> predict(lm.out, data.frame(fibers = 3), interval = "prediction")
      fit      lwr      upr
1 45.58553 22.55513 68.61593
> |
```

بازه پیش بینی rating یک محصول با fiber = 3gram (دقت شود که در اینجا میانگین برآورد نمیشود و مقدار متغیر هدف یک random record برآورد میشود) برابر [22.55513, 68.61593] می باشد.

(19

با توجه به نمودار های residual بر حسب x_i ، دیده میشود که متغیر هدف (rating) کمی به راست چوله است (نمودار QQ). اما میتوان از خطی بودن مدل و همگنی واریانس ها اطمینان حاصل کرد (نمودار residuals vs fitted)



(20)

با ران کردن کد زیر داریم:

```
22 sugars = cereals$sugars;
23 mreg.out = lm(ratings~sugars + fibers)|
24 summary(mreg.out)
```

24:39 (Top Level) ⌵

Console Terminal Background Jobs ×

R 4.2.2 · ~/

Call:
lm(formula = ratings ~ sugars + fibers)

Residuals:

Min	1Q	Median	3Q	Max
-12.133	-4.247	-1.031	2.620	16.398

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	51.6097	1.5463	33.376	< 2e-16 ***
sugars	-2.1837	0.1621	-13.470	< 2e-16 ***
fibers	2.8679	0.3023	9.486	2.02e-14 ***

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 6.219 on 74 degrees of freedom
Multiple R-squared: 0.8092, Adjusted R-squared: 0.804
F-statistic: 156.9 on 2 and 74 DF, p-value: < 2.2e-16

> |

The estimated multiple regression equation is:

$$\text{rating}[i] = 51.6097 + 2.8679 * \text{fibers}[i] + -2.1837 * \text{sugars}[i]$$

(21)

تفسیر ضریب fiber در معادله سوال قبل:

با ثابت نگه داشتن مقدار sugar، میانگین rating برای یک محصول با هر گرم افزایش fiber، 2.8679 درصد افزایش می یابد.

(22)

مقدار r-squared در حالت multiple regression حدود 0.8 می باشد که نسبت به حالت simple regression با مقدار 0.3 بسیار بزرگتر است. بدلیل استفاده از داده بیشتر در مدل، رگرسیون حاصل بهتر توانسته record ها را برآورد کند.