

به نام خداوند بخشنده مهربان

عنوان اسلایدها:

مروری بر رگرسیون خطی ۱

تهیه کننده: محراب عتیقی

رگرسیون چیست؟

- رگرسیون یک نوع تحلیل آماری برای بررسی وابستگی یا ارتباط بین متغیرها می باشد. و به سوالاتی در رابطه با ارتباط بین متغیرها پاسخ می دهد:

(۱) آیا بین متغیرها ارتباطی وجود دارد؟

(۲) یک متغیر خاص تحت تاثیر کدام متغیر یا متغیرهای دیگر است؟

انواع متغیر ها در تحلیل رگرسیون

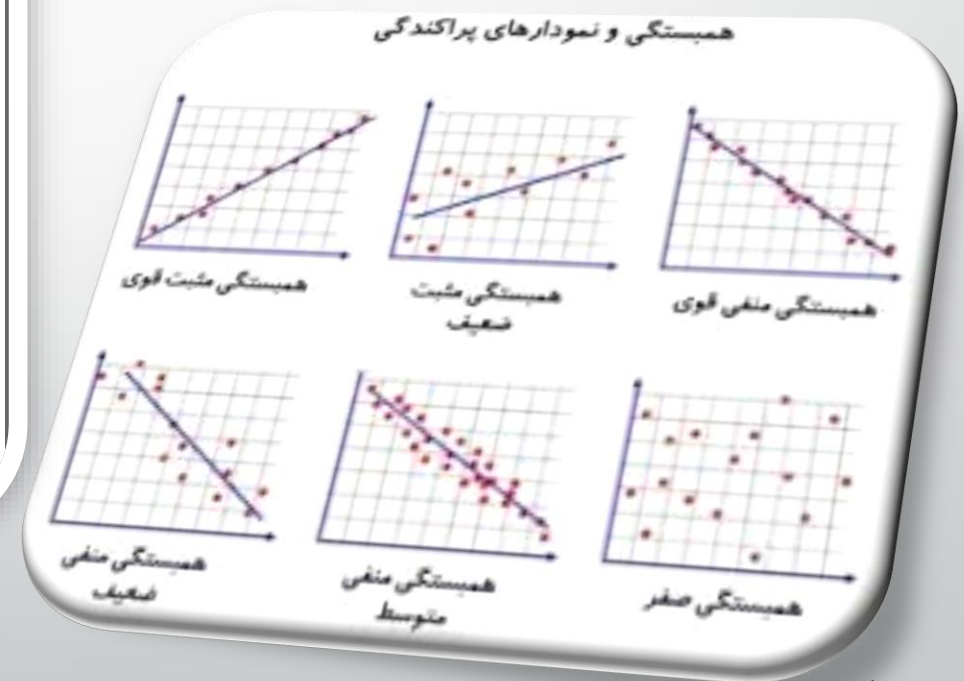
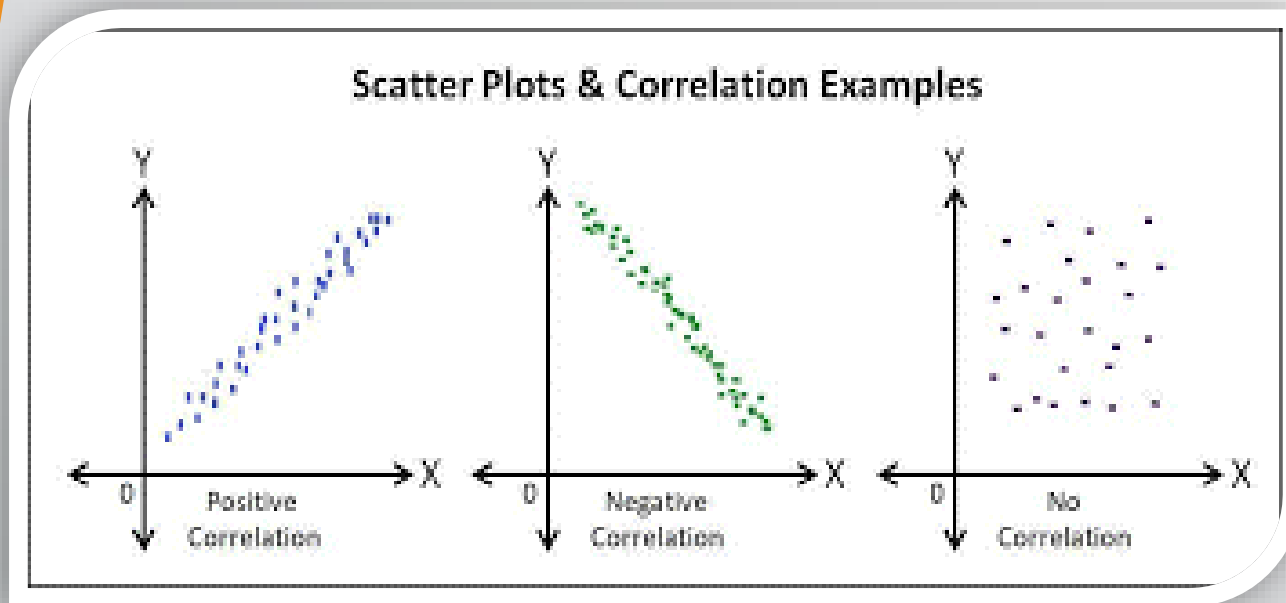
متغیر مستقل (X):

متغیری که قابل تغییر است و بر متغیر وابسته اثرگذار است، که به آن متغیر پیشگو و یا پیش‌بین نیز گفته می‌شود.

متغیر وابسته (y):

متغیری که تحت تاثیر متغیر و یا متغیرهای مستقل است و با تغییر متغیر مستقل نیز تغییر می‌کند.

انواع رابطه متغیر مستقل و متغیر وابسته به کمک نمودار پراکنش:



مدل رگرسیون خطی ساده

مدل رگرسیون خطی ساده مدلی، شامل یک متغیر مستقل و یک متغیر وابسته می باشد. این مدل به صورت زیر نوشته می شود:

$$y = \beta_0 + \beta_1 x + \varepsilon$$

متغیر مستقل (پیشگو) x →
عرض از مبدا جامعه β_0 →
شیب جامعه β_1 →
خطای تصادفی ε ←
متغیر پاسخ (وابسته) y ←

که در آن شیب و عرض از مبدا، پارامترهای ثابت ولی مجهول هستند.

تعاریف:

متغیر: صفتی است که از یک فرد به فرد دیگر تغییر می کند.

متغیر کمی: صفتی است که قابل سنجش و اندازه گیری است.

متغیر کیفی: صفتی است که قابل سنجش و اندازه گیری نیست.

پارامتر: مقادیری مربوط به جامعه که مجهول اما ثابت هستند.

برآورد: مقادیری بدست آمده از نمونه تصادفی که معلوم و تصادفی هستند.

خطا چیست؟

خطا یا باقی مانده همان فاصله نقاط پیش بینی شده تا خط رگرسیونی جامعه می باشد.

توزیع متغیر تصادفی خطا

- خطاها متغیرهای تصادفی هستند.

- $\varepsilon \sim N(0, \sigma^2)$

- خطاها ناهمبسته هستند، یعنی مقدار یک خطا به مقدار خطای دیگری وابسته نیست.

- واریانس σ^2 ، پارامتری نامعلوم است.

- $\sum \varepsilon_i = \sum e_i = 0$

توزیع متغیر تصادفی پاسخ

- متغیر پاسخ y ، یک متغیر تصادفی هست.
- پاسخ‌ها نیز ناهمبسته هستند. (چون خطاها ناهمبسته هستند).
- متغیر پاسخ y ، به ازای هر مقدار x یک توزیع احتمال نرمال دارد. (فقط میانگین توزیع متغیر پاسخ ما تغییر می‌کند)
- $y \sim N(\beta_0 + \beta_1 x, \sigma^2)$

ضرایب مدل رگرسیون خطی ساده و تفسیر آنها

- **عرض از مبدا (β_0):** برابر است با میانگین توزیع متغیر پاسخ Y ، هنگامی که $X=0$ باشد. به شرط اینکه مقدار صفر در دامنه مقادیر X وجود داشته باشد، در غیر این صورت تفسیر ندارد.
- **شیب (β_1):** به ازای یک واحد تغییر در مقدار X ، میانگین توزیع متغیر پاسخ Y ، به میزان β_1 تغییر می کند.

برآورد حداقل مربعات ضرایب رگرسیونی

• برآورد شیب:

$$\widehat{\beta}_1 = \frac{s_{xy}}{s_{xx}}$$

$$s_{xy} = \sum (x_i - \bar{x})(y_i - \bar{y}) \quad , \quad s_{xx} = \sum (x_i - \bar{x})^2 \quad \text{که در آن:}$$

• برآورد عرض از مبدا:

$$\widehat{\beta}_0 = \bar{y} - \widehat{\beta}_1 \bar{x} \quad \bullet$$

ویژگی‌های برآوردگرهای حداقل مربعات ضرایب رگرسیون

$\hat{\beta}_1, \hat{\beta}_2$ نااریب هستند و دارای کمترین واریانس هستند.

ولذا طبق قضیه گوس-مارکوف، این برآوردگرهای حداقل مربعات برای مدل ما بهترین برازش را دارند.

برآورد واریانس:

$$\widehat{\sigma^2} = \frac{SSE}{n - 2} = MSE$$

$$SSE = \sum e_i^2 \quad \text{که در آن}$$

ویژگی‌های برآوردگر واریانس:

- میانگین مربعات خطا، MSE برآوردگری ناریب از σ^2 است.
- برآوردگر σ^2 ، وابسته به مدل است.

توزیع برآوردگرهای ضرایب خط رگرسیونی

$$\widehat{\beta}_1 \sim N\left(\beta_1, \frac{\sigma^2}{s_{xx}}\right),$$
$$\widehat{\beta}_0 \sim N\left(\beta_0, \sigma^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{s_{xx}}\right)\right)$$

آزمون فرض درباره عرض از مبدا خط رگرسیونی:

$$\begin{cases} H_0 : \beta_0 = \beta_{00} \\ H_1 : \beta_0 \neq \beta_{00} \end{cases} , t_0 = \frac{\widehat{\beta}_0 - \beta_{00}}{\sqrt{MSE\left(\frac{1}{n} + \frac{\bar{x}^2}{s_{xx}}\right)}} , |t_0| > t_{(n-2, 1-\frac{\alpha}{2})}$$

آزمون فرض درباره شیب خط رگرسیونی:

$$\begin{cases} H_0 : \beta_1 = 0 \\ H_1 : \beta_1 \neq 0 \end{cases}, t_0 = \frac{\widehat{\beta}_1 - 0}{\sqrt{MSE/s_{xx}}} \quad , |t_0| > t_{(n-2, 1-\frac{\alpha}{2})}$$

$$, Z_0 = \frac{\widehat{\beta}_1 - 0}{\sqrt{\sigma^2/s_{xx}}}$$

مجموع مربعات کل:

$$SST = \sum_{i=1}^n (y_i - \bar{y})^2, df = n - 1$$

زیر یکدرجه آزادی مربوط، روی انحرافات ازدست رفته است.

مجموع مربعات رگرسیون:

$$SSR = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2, df = 1$$

مجموع مربعات خطا:

$$SSE = \sum_{i=1}^n (y_i - \hat{y}_i)^2, df = n - 2$$

$$SST = SSR + SSE$$

همیشه داریم:

آزمون آنالیز واریانس درباره شیب خط رگرسیونی:

$$\begin{cases} H_0 : \beta_1 = 0 \\ H_1 : \beta_1 \neq 0 \end{cases}, F_0 = \frac{MSR}{MSE}, F_0 > F_{(1-\alpha, 1, n-1)}$$

$$, MSR = \frac{SSR}{df_{MSR}}, MSE = \frac{SSE}{df_{MSE}}$$

فواصل اطمینان برای شیب و عرض از مبدا و واریانس:

- فاصله اطمینان شیب خط رگرسیون:

$$\frac{\widehat{\beta}_1 - \beta_1}{\sqrt{MSE/s_{xx}}} \sim t_{(n-2)} \rightarrow \left(\widehat{\beta}_1 - t_{(n-2, 1-\frac{\alpha}{2})} \sqrt{\frac{MSE}{s_{xx}}} < \beta_1 < \widehat{\beta}_1 + t_{(n-2, 1-\frac{\alpha}{2})} \sqrt{\frac{MSE}{s_{xx}}} \right)$$

- فاصله اطمینان عرض از مبدا خط رگرسیون:

$$\frac{\widehat{\beta}_0 - \beta_0}{\sqrt{MSE \left(\frac{1}{n} + \frac{\bar{x}^2}{s_{xx}} \right)}} \sim t_{(n-2)} \rightarrow \left(\widehat{\beta}_0 - t_{(n-2, 1-\frac{\alpha}{2})} \sqrt{MSE \left(\frac{1}{n} + \frac{\bar{x}^2}{s_{xx}} \right)} < \beta_0 < \widehat{\beta}_0 + t_{(n-2, 1-\frac{\alpha}{2})} \sqrt{MSE \left(\frac{1}{n} + \frac{\bar{x}^2}{s_{xx}} \right)} \right)$$

- فاصله اطمینان برای واریانس:

$$\frac{(n-2)MSE}{\sigma^2} \sim X^2_{(n-2)} \rightarrow \left(\frac{(n-2)MSE}{X^2_{(\frac{\alpha}{2}, n-2)}} < \sigma^2 < \frac{(n-2)MSE}{X^2_{(1-\frac{\alpha}{2}, n-2)}} \right)$$

فواصل اطمینان برای میانگین پاسخ به ازای x_0 و پیش‌بینی m مشاهده جدید y_0 .

• فاصله اطمینان برای میانگین پاسخ به ازای x_0 :

$$\widehat{y}_0 \sim N \left(E[y|x_0], \sigma^2 \left(\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{xx}} \right) \right)$$

$$\left(\widehat{y}_0 - t_{(n-2, 1-\frac{\alpha}{2})} \sqrt{MSE \left(\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{xx}} \right)} < E[y|x_0] < \widehat{y}_0 + t_{(n-2, 1-\frac{\alpha}{2})} \sqrt{MSE \left(\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{xx}} \right)} \right)$$

• فاصله اطمینان برای پیش‌بینی m مشاهده جدید y_0 :

$$\left(\widehat{y}_0 - t_{(n-2, 1-\frac{\alpha}{2})} \sqrt{MSE \left(\frac{1}{m} + \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{xx}} \right)} < y_0 < \widehat{y}_0 + t_{(n-2, 1-\frac{\alpha}{2})} \sqrt{MSE \left(\frac{1}{m} + \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{xx}} \right)} \right)$$

ضریب تعیین:

- نسبت تغییرات تبیین شده متغیر پاسخ به وسیله متغیر رگرسیونی X

$$R^2 = \frac{SSR}{S_{yy}} = 1 - \frac{SSE}{S_{yy}} \quad 0 \leq R^2 \leq 1$$

- نکته: اگر ρ ، ضریب همبستگی بین دومتغیر پاسخ و پیشگو ما باشد، آنگاه:

$$R^2 = \rho^2$$

معیارهای مناسبیت مدل:

- باقی مانده های ما پرت و دور افتاده نباشند.
- توزیع باقی مانده ها نرمال باشد.
- رسم نمودارهای باقی مانده استیودنت شده و استاندارد.
- رسم انواع نمودارهای باقی مانده ها در برابر متغیرهای پیشگو و پاسخ مدل.
- بررسی فرض ثبات واریانس خطاها.
- آزمون فقدان برازش خطا.

آزمون فرض بررسی فقدان برازش خط

• این آزمون تنها زمانی قابل اجراست که:

به ازای حداقل یک سطح X ، مشاهدات تکراری برای Y وجود داشته باشد.

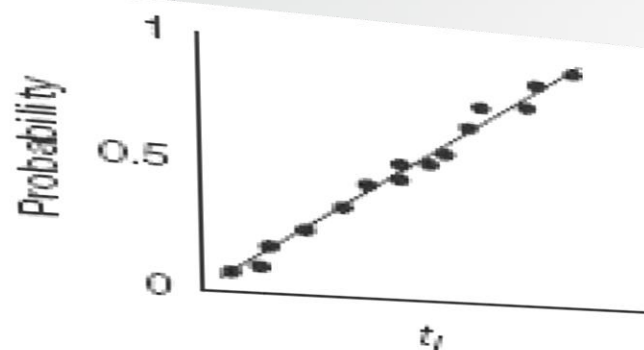
$$SSE = SS_{PE} + SS_{LOF}$$

$$(y_{ij} - \hat{y}_i) = (y_{ij} - \bar{y}_i) + (\bar{y}_i - \hat{y}_i)$$

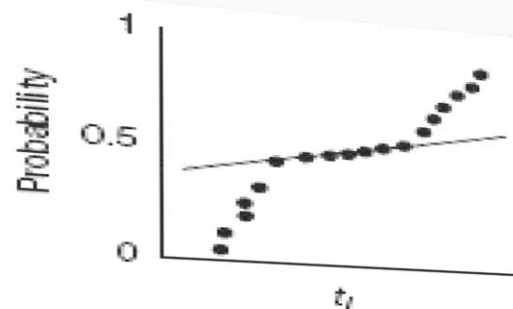
$$\sum_{i=1}^m \sum_{j=1}^{h_i} (y_{ij} - \hat{y}_i)^2 = \sum_{i=1}^m \sum_{j=1}^{n_i} (y_{ij} - \bar{y}_i)^2 + \sum_{i=1}^m n_i (\bar{y}_i - \hat{y}_i)^2$$

$$\begin{cases} H_0 : \beta_1 = 0 \\ H_1 : \beta_1 \neq 0 \end{cases}, F_0 = \frac{\frac{SS_{LOF}}{m-2}}{\frac{SS_{PE}}{n-m}} = \frac{MS_{LOF}}{MS_{PE}}, F_0 > F_{(\alpha, (m-2), (n-m))}$$

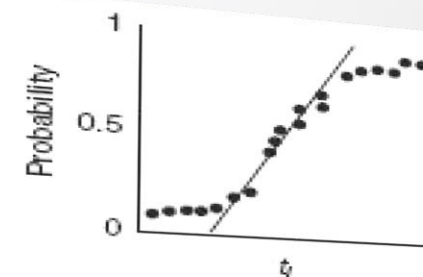
بررسی فرض نرمال بودن خطاها



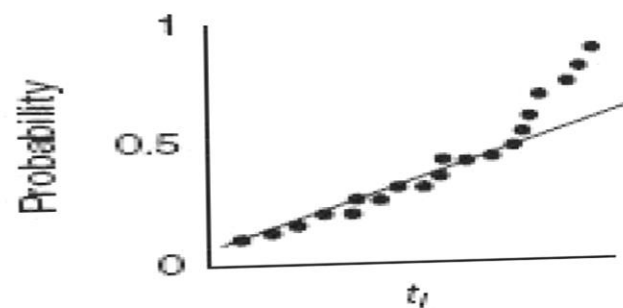
(a)



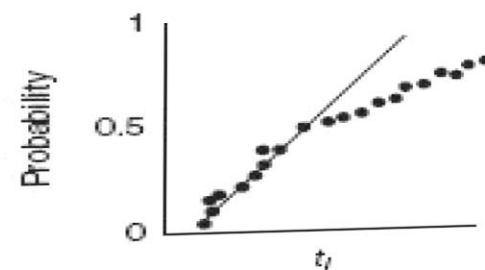
(b)



(c)

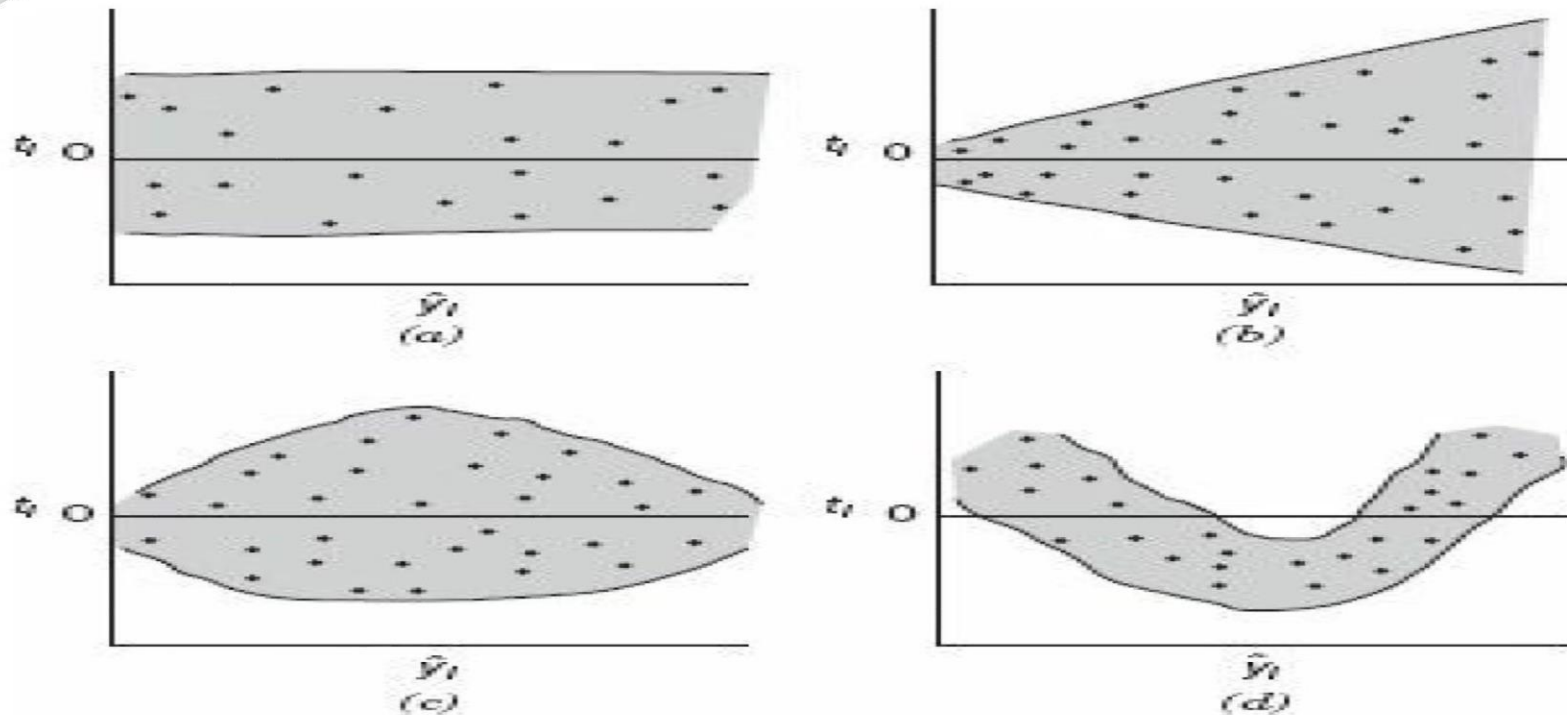


(d)



(e)

بررسی ثبات واریانس خطاها



انواع باقی مانده ها:

- باقی مانده ها:

$$e_i = y_i - \hat{y}_i, e_i \sim N(0, \sigma^2)$$

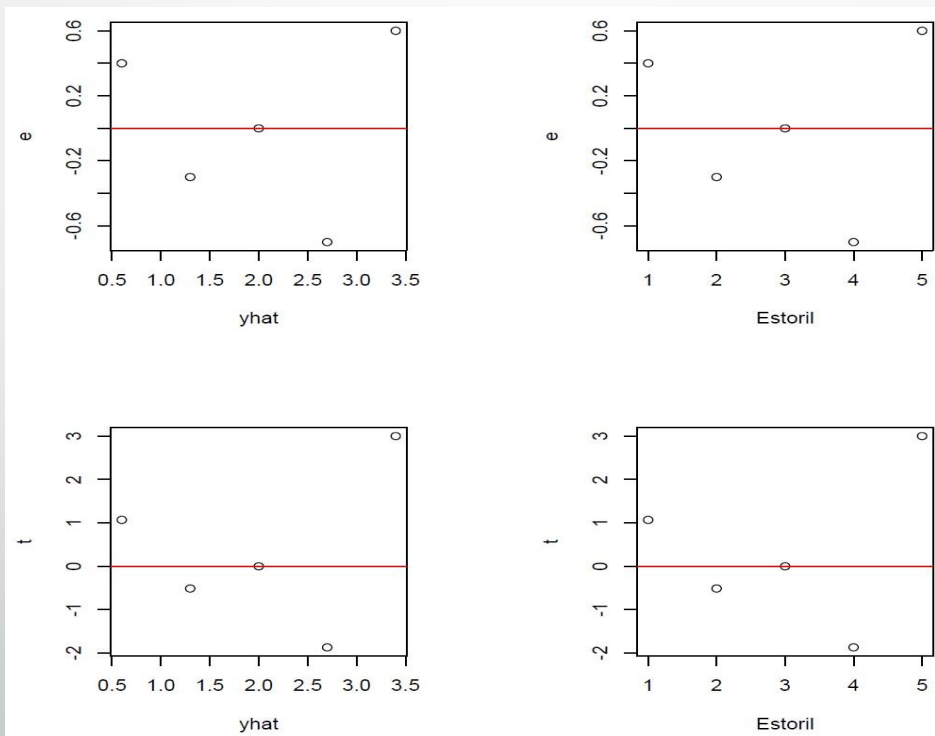
- باقی مانده های استاندارد شده:

$$d_i = \frac{e_i - 0}{\sqrt{(MSE)}}, d_i \sim N(0, 1)$$

- باقی مانده های استیودنت شده:

$$t_i = \frac{e_i}{\sqrt{MSE \left(1 - \left(\frac{1}{n} + \frac{(x_i - \bar{x})^2}{S_{xx}} \right) \right)}}$$

بررسی داده‌های پرت و مشاهدات دورافتاده



کدهای R برای مسایل مطرح شده :

- برای رسم نمودار پراکنش :

```
Plot(x, y)
```

- برای ساختن یک مدل رگرسیونی یگانه داریم :

```
Fit<-lm(y~x)
```

- برای رسم خط رگرسیونی داریم :

```
plot(x, y);abline(fit, col="red")
```

- برای بدست آوردن شیب خط و عرض از مبدا داریم :

```
Coefficients(fit) or coef(fit)
```

- برای بدست آوردن مقدار و پیشبینی به صورت فاصله اطمینان داریم :

```
New<-data.frame(x=c(...));predict(fit,newdata=New , interval="prediction", level=0.95)
```

- برای بدست آوردن متوسط \bar{y} درنقطه جدید یا قدیمی‌ها بصورت فاصله اطمینان داریم :

```
Predict(fit,newdata=New , interval="confidence" , level=0.95)
```

- برای رسم نمودار یک برآورد فاصله‌ای یا متوسط آنها و مقادیر پیشبینی شده بصورت یکجا داریم :

```
matplot(x, cbind(a, b), type="L")
```

- برای دیدن نتیجه آزمون وجود شیب خط رگرسیون داریم:

```
Summary(fit) and anova(fit)
```

در دستور `summary` می‌توانیم به ضریب تعیین، SS ها و درجه های آزادی و... پی ببریم .

- برای بدست آوردن فاصله اطمینان برای شیب و عرض از مبدا داریم:

```
Confit(fit, level=0.95)
```

- برای بدست آوردن ضریب همبستگی داریم:

```
Cor.test(x, y, method="pearson")
```

- برای مشاهده باقی مانده های اصلی، استاندارد و استیودنت شده داریم:

```
e<-residuals(fit);d<-rstandard(fit);t<-rstudent(fit)
```

- برای مشاهده مقادیر \hat{y} ، داریم:

```
Yhat<-fitted(fit)
```

- برای رسم نمودار احتمال نرمال داریم:

```
qqnorm(e)
```

- برای بررسی آزمون فقدان برازش خط داریم:

```
fit1<-lm(y~x);fit2<-aov(y~factor(x));anova(fit1,fit2)
```

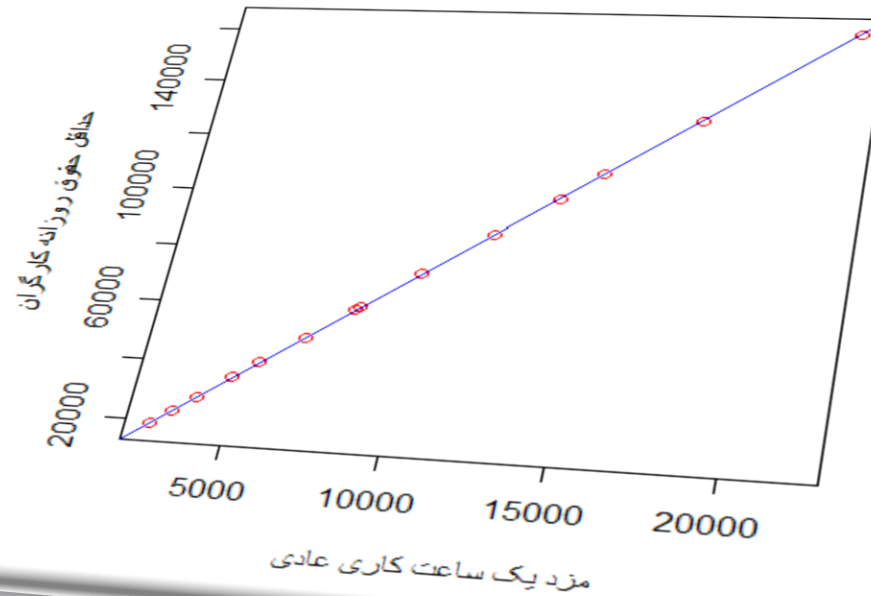
در ادامه به بررسی مثالی از مباحث مطرح شده می پردازیم:
در این قسمت می خواهیم به بررسی مزد یک ساعت کاری عادی کارگران
مشمول قانون کار، بروی حداقل حقوق آنها بپردازیم:

● ابتدا داده های خود را به شکل زیر فراخوانی می کنیم:

```
> data
      min_hoghogh  mozd_per_h
1      162370      22152
2      129900      17720
3      110100      15020
4      101000      13777
5       87840      11984
6       73200       9986
7       61000       8322
8       60000       8186
9       50000       6821
10      40864       5575
11      35534       4848
12      28446       3881
13      23282       3175
14      18930       2583
```

حال برای رسم نمودار پراکنش و خطرگرسیونی بین این دو متغیر مستقل و پیشگو داریم:

```
> fit1=lm(min_hoghogh~mozd_per_h)
> plot(mozd_per_h,min_hoghogh,xlab="مزد یک ساعت کاری عادی",ylab="حداقل حقوق روزانه کارگران",col="REd")
> abline(fit1,col="Blue")
```



برای بدست آوردن مقادیر شیب خط و عرض از مبدا آن داریم:

```
> coef(fit1)
(Intercept)  mozd_per_h
-0.4451111    7.3302412
```

برای بدست آوردن خلاصه‌ای از مدل یگانه خود داریم:

```
>summary(fit1)
```

```
Call:
lm(formula = min_hoghogh ~ mozd_per_h)
Residuals: Min 1Q Median 3Q Max
-9.0585 -3.3171 -1.7360 0.8163 11.7117
Coefficients:
Estimate Std. Error t value Pr(>|t|)
(Intercept) -0.4451111 3.2300262 -0.138 0.893
mozd_per_h 7.3302412 0.0002903 25246.590 <2e-16 ***
--- Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
Residual standard error: 6.156 on 12 degrees of freedom
Multiple R-squared: 1, Adjusted R-squared: 1
F-statistic: 6.374e+08 on 1 and 12 DF, p-value: < 2.2e-16
```

برای بدست آوردن یک فاصله اطمینان ۹۵٪ برای پیش‌بینی و متوسط مقدار متغیر وابسته خود داریم:

```
>predict.lm(fit1,interval="confidence")
```

	fit	lwr	upr
1	162379.06	162370.33	162387.79
2	129891.43	129885.15	129897.71
3	110099.78	110094.81	110104.75
4	100988.29	100983.83	100992.75
5	87845.17	87841.27	87849.06
6	73199.34	73195.75	73202.94
7	61001.82	60998.15	61005.49
8	60004.91	60001.22	60008.60
9	49999.13	49995.15	50003.12
10	40865.65	40861.26	40870.04
11	35536.56	35531.90	35541.23
12	28448.22	28443.14	28453.30
13	23273.07	23267.66	23278.48
14	18933.57	18927.88	18939.26

```
>predict.lm(fit1,interval = "prediction")
```

	fit	lwr	upr
1	162379.06	162363.06	162395.06
2	129891.43	129876.62	129906.24
3	110099.78	110085.47	110114.08
4	100988.29	100974.15	101002.42
5	87845.17	87831.20	87859.13
6	73199.34	73185.46	73213.23
7	61001.82	60987.92	61015.73
8	60004.91	59991.00	60018.82
9	49999.13	49985.14	50013.12
10	40865.65	40851.54	40879.76
11	35536.56	35522.36	35550.77
12	28448.22	28433.88	28462.56
13	23273.07	23258.61	23287.53
14	18933.57	18919.00	18948.14

برای بدست آوردن یک فاصله اطمینان ۹۵٪ برای پیش‌بینی و متوسط مقدار متغیر وابسته خود در نقاطی جدید داریم:

```
> New<-data.frame(c(22000,19000,18000,16000,18888,17980,13900,16578,15890,21000,17800,18000,15600,12780))
```

```
>predict.lm(fit1,interval="confidence",newdata = New)
```

	fit	lwr	upr
1	162379.06	162370.33	162387.79
2	129891.43	129885.15	129897.71
3	110099.78	110094.81	110104.75
4	100988.29	100983.83	100992.75
5	87845.17	87841.27	87849.06
6	73199.34	73195.75	73202.94
7	61001.82	60998.15	61005.49
8	60004.91	60001.22	60008.60
9	49999.13	49995.15	50003.12
10	40865.65	40861.26	40870.04
11	35536.56	35531.90	35541.23
12	28448.22	28443.14	28453.30
13	23273.07	23267.66	23278.48
14	18933.57	18927.88	18939.26

```
>predict.lm(fit1,interval = "prediction",newdata = New)
```

	fit	lwr	upr
1	162379.06	162363.06	162395.06
2	129891.43	129876.62	129906.24
3	110099.78	110085.47	110114.08
4	100988.29	100974.15	101002.42
5	87845.17	87831.20	87859.13
6	73199.34	73185.46	73213.23
7	61001.82	60987.92	61015.73
8	60004.91	59991.00	60018.82
9	49999.13	49985.14	50013.12
10	40865.65	40851.54	40879.76
11	35536.56	35522.36	35550.77
12	28448.22	28433.88	28462.56
13	23273.07	23258.61	23287.53
14	18933.57	18919.00	18948.14

برای بدست آوردن ضریب همبستگی متغیر مستقل و متغیر پاسخ خود داریم:

```
>cor.test(min_hoghogh, mozd_per_h, method = "pearson")
Pearson's product-moment correlation
data: min_hoghogh and mozd_per_h t = 25247, df = 12, p-value < 2.2e-16
alternative hypothesis: true correlation is not equal to 0
95 percent confidence interval: 1 1 sample estimates: cor 1
```

برای انجام آزمون آنوا یا آنالیز واریانس داریم:

```
>anova(fit1)

Analysis of Variance Table Response:
min_hoghogh Df Sum Sq Mean Sq F value Pr(>F)
mozd_per_h 1 2.4153e+10 2.4153e+10 637390311 < 2.2e-16 ***
Residuals 12 4.5500e+02 3.8000e+01
--- Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

برای انجام آزمون فقدان برازش داریم:

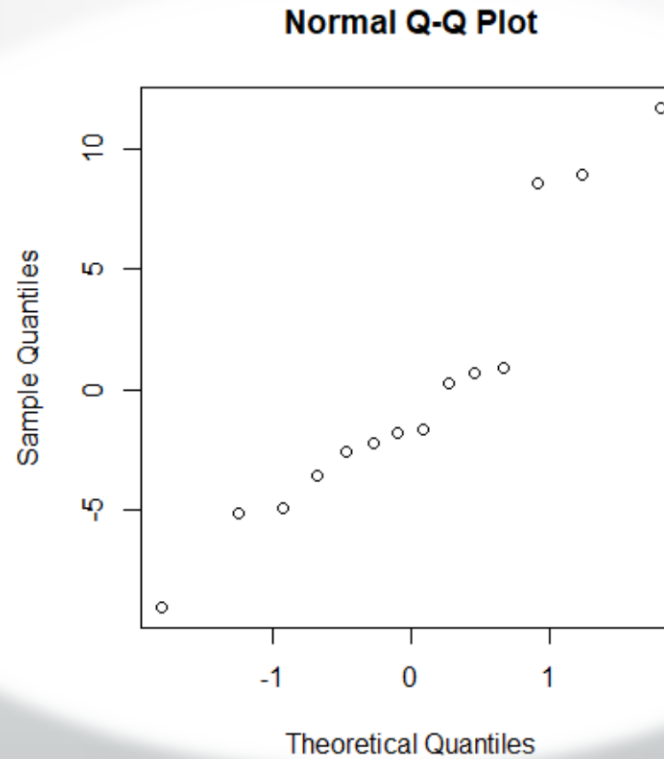
```
>model.reg<-lm(min_hoghogh~mozd_per_h)
>model.aov<-aov(min_hoghogh~factor(mozd_per_h))
```

```
anova(model.reg,model.aov)
Analysis of Variance Table Model
1: min_hoghogh ~ mozd_per_h Model
2: min_hoghogh ~ factor(mozd_per_h)
  Res.Df RSS      Df    Sum of Sq F Pr(>F)
1     12  454.71
2      0    0.00     12    454.71      0
```

```
> #the pure error value is:
> anova(model.reg,model.aov)[2,2]
[1] 0
> #the pure error df is :
> anova(model.reg,model.aov)[2,1]
[1] 0
> #the lack of fit value is:
> anova(model.reg,model.aov)[2,4]
[1] 454.7149
> #the lack of fit df is:
> anova(model.reg,model.aov)[2,3]
[1] 12
```

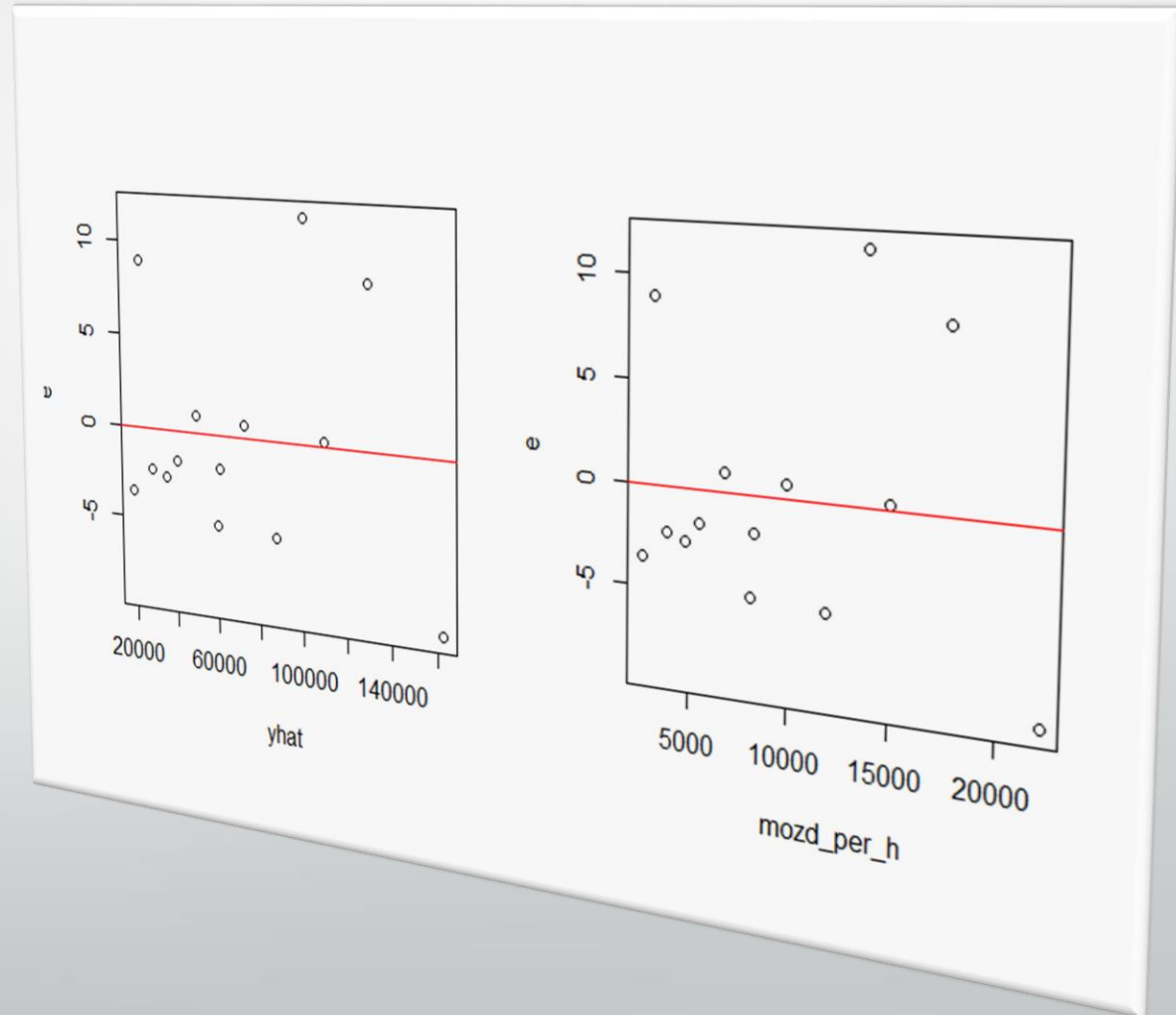
برای تعریف انواع باقی مانده ها و رسم نمودار نرمال احتمال باقی مانده ها داریم:

```
yhat<-fitted(fit1)  
> s<-rstandard(fit1)  
> t<-rstudent(fit1)  
> e<-residuals(fit1)  
> qqnorm(e)
```



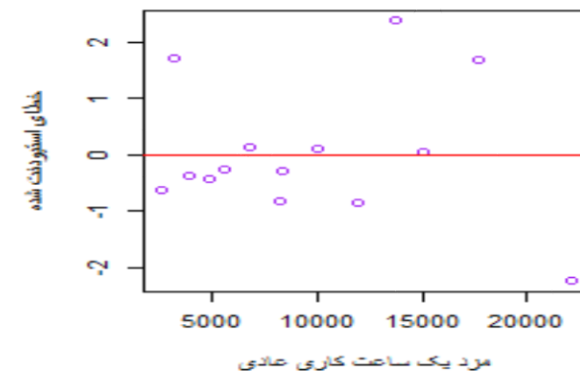
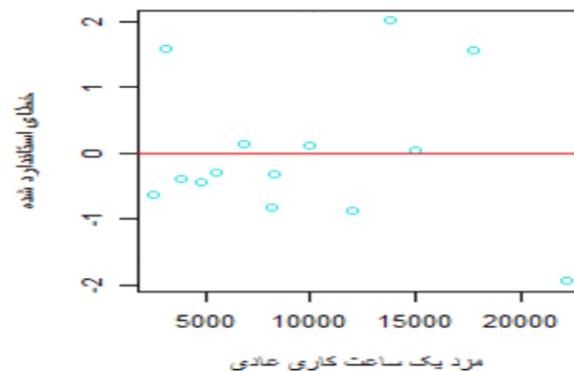
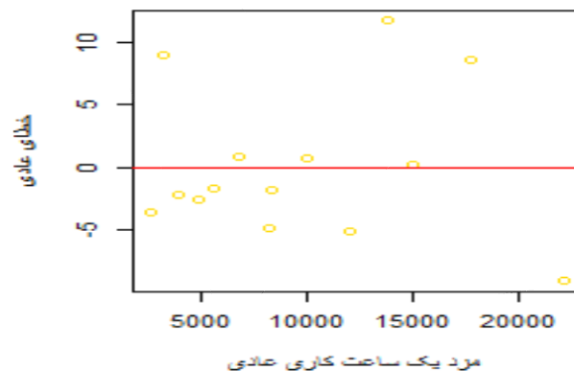
برای رسم نمودار، باقی مانده‌ها در برابر متغیر مستقب و متغیر پیشگو داریم:

```
> par(mfrow=c(1,2),pty="s")  
> plot(yhat,e)  
> abline(h=0,col="red")  
> plot(mozd_per_h,e)  
> abline(h=0,col="red")
```



برای رسم نمودارهای انواع باقی مانده ها داریم:

```
> par(mfrow=c(1, 3), pty="s")  
> plot(mozd_per_h, e, xlab="مزد یک ساعت کاری عادی", ylab="خطای عادی", col="Gold")  
> abline(h=0, col="red")  
> plot(mozd_per_h, s, xlab="مزد یک ساعت کاری عادی", ylab="خطای استاندارد شده", col=85)  
> abline(h=0, col="red")  
> plot(mozd_per_h, t, xlab="مزد یک ساعت کاری عادی", ylab="خطای استیودنت شده", col="purple")  
> abline(h=0, col="red")
```



با تشکر از همراهی شما عزیزان

