



دانشگاه صنعتی امیرکبیر  
(پلی تکنیک تهران)

دانشکده مهندسی کامپیوتر و فناوری اطلاعات

## تکلیف شماره ۱ درس یادگیری ماشین

استاد

دکتر ناظر فرد

امیرحسین کاشانی

## فهرست مطالب

Q1).....	4
a) .....	4
b) .....	5
c).....	5
d) .....	6
e) .....	6
f) .....	8
g) .....	9
h) .....	12
i) .....	12
Q2).....	13
a) .....	13
b) .....	14
c).....	15
D).....	16
e) .....	17
f) .....	18
g) .....	19
h) .....	20
Q3).....	21
a) .....	21
b) .....	22
c).....	23
d) .....	24
e) .....	27
F) .....	29
g) .....	30
H).....	31
Q4).....	33
a) .....	33
b) .....	33
c).....	35

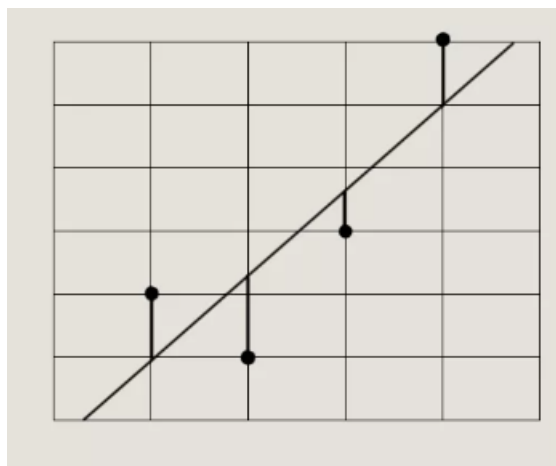
d) .....	35
e) .....	36
f) .....	37
g) .....	37
Q5).....	38
a) .....	38
b) .....	39
c).....	39
d) .....	41
e) .....	42

# Q1)

a)

Least square error:

در این حالت به دنبال یافتن یک خط با بیشترین تناسب نسبت به مجموعه داده ورودی هستیم به بیان دیگر به دنبال خطی می گردیم که مجموع توان دو فاصله نقاط از نقطه متناظر آن ها در آن خط کمترین مقدار ممکن باشد. این تسک دقیقاً مشابه عملی است که در رگرسیون صورت می گیرد و  $w$  های بدست آمده همان معادله خطی هستند که least square error به آن نیاز دارد.



برنامه نویسی خطی

یک مساله برنامه نویسی خطی به شکل ماکسیمم کردن یا مینیمم کردن یک معادله خطی با در نظر گرفتن یک قیود خطی بر روی متغیر های به کار رفته در معادله ی اولیه مطرح می شود.

$$z = 0.6x_1 + 0.35x_2$$

$$5x_1 + 7x_2 \geq 8$$

$$4x_1 + 2x_2 \geq 15$$

$$2x_1 + x_2 \geq 3$$

$$x_1 \geq 0, x_2 \geq 0.$$

مثال فوق یک نمونه از مقدار دهی به معادله  $ax_1 + bx_2$  می باشد که شرط های زیر را برقرار کرده است.

## بهینه سازی محدب

مساله ای است که همه قیود آن حالت محدب دارند و تابع هدفمان نیز بهینه نمودن یک مساله محدب است. از آنجا که ترکیب خطی یک مساله محدب است در نتیجه برنامه نویسی خطی نیز یک حالت خاص از بهینه سازی محدب محسوب می گردد.

b)

رگرسیون ساده از آنجایی که به شکل

$$y = \alpha + \beta x$$

تعریف می شود پارامتر complexity مدل برای آن مطرح نمی گردد اما برای تعداد iteration می توان از دوسیاست حداکثر iteration یا مینیمم خطای قابل قبول استفاده کرد و در صورت رسیدن به هر یک از این شروط اجرای الگوریتم را متوقف نمود.

\*\* بصورت کلی برای انتخاب میزان پیچیدگی میتوان از یک سیاست جست و جوی باینری بهره برد برای مثال در ابتدا یک مدل با درجه پیچیدگی یک و یک مدل با ماکسیمم درجه پیچیدگی مورد نظر ایجاد کرد (برای مثال M) سپس یک مدل با درجه  $M/2$  و بعد از آن در بخشی که نتیجه بهتری داشته را انتخاب کنیم (فرض کنیم بخش راست با پیچیدگی درجه M بهتر بوده) در گام بعدی درجه پیچیدگی ما روند شده  $0.75 M$  خواهد بود و این جستو جوی باینری را ادامه می دهیم تا به یک حالت بهینه نسبت به فضای موجود دست پیدا کنیم.

البته لازم به ذکر است که درحالت هایی که بازه گزینه ها برای انتخاب پیچیدگی کوچک است می توانیم ۳ تا ۳ (برای مثال ۳ فرض شده) پیچیدگی را زیاد کنیم و در نقطه ای که بهترین جواب وجود داشت تمام پیچیدگی های درجه ای نزدیک را بررسی کنیم.

c)

به این پارامتر، پارامتر بایاس گفته می شود. در صورت عدم وجود این پارامتر الگوریتم رگرسیون فقط توانایی محاسبه خطوطی که از مبدا می گذرند را خواهد داشت به بیان دیگر نشانگر عرض از مبدا خط پیشنهادی این روش می باشد و در صورت نبودن این پارامتر حاصل اجرای الگوریتم معادل شکل زیر می باشد.

$$\hat{y} = \mathbf{w}^T \mathbf{x}$$

اما وجود پارامتر  $\theta_0$  به ما این امکان را می دهد تا همه توابع خطی موجود را بتوانیم ایجاد کنیم.

$$\hat{y} = \mathbf{w}^T \mathbf{x} + b$$

d)

روش k-fold داده آماری را به k دسته تقسیم می کند و هر بار یکی را به عنوان تست بر می گزیند و بقیه را به عنوان داده آموزش مصرف می کند. (این کار را k بار تکرار می کند)

روش leave-one-out یک داده را به عنوان تست بر می دارد و کل مابقی دیتا ست را آموزش می بیند و این کار را k بار تکرار می کند. این کار باعث می شود که هر با آموزش دادن این مدل مدت زمان قابل توجهی به طول انجامد اما در زمانی که حجم داده کم است بهتر است از این روش استفاده شود.

I. در مقیاس کوچک leave one out و در مقیاس بزرگ k-fold به دلیل اینکه بار محاسباتی کمتری در هر یادگیری خواهد داشت بهتر است.

II. باتوجه به اینکه روش leave-one-out فقط یک داده را در هر دور ملاک قضاوت قرار می دهد نسبت به داده پرت مقاومت خوبی ندارد در نتیجه k-fold در مقابل نویز مقاوم تر است.

III. K-fold، در هر مرحله بخش کوچکتري از داده ها را به عنوان داده آموزش بر می گزیند

IV. به طور کلی نمی توان روش مشخصی برای تعیین k بیان نمود به عنوان یک روش دمه دستی k بین ۵ تا ۱۰ توصیه می شود که بیانگر این است که در هر تکرار بین ۱۰ تا ۲۰ درصد داده ها را به عنوان تست برداریم. لازم به ذکر است که در صورتی که حجم داده های ما کم است باید از k های بزرگتر استفاده کنیم تا میزان داده مناسب برای اجرای فرآیند آموزش داشته باشیم ولی در زمانی که حجم داده معقول است k با مقدار کمتر زمان کمتر در اجرای آموزش را برای ما به ارمغان می آورد و همین طور به فرآیند آموزش لطمه نمی زند.

e)

regularization یکی از راه های پر طرف دار برای جلوگیری از overfit در اجرای الگوریتم می باشد و به این صورت عمل می کند که در تابع هزینه وزن های انتخاب شده را جمع می کند. مانند شکل زیر

$$J(\theta) = \frac{1}{2m} \left[ \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^n \theta_j^2 \right]$$
$$\min_{\theta} J(\theta)$$

این کار باعث این می شود که از وزن های به شدت زیاد که نتیجه را به سمت خود نا متوازن می کنند جلوگیری شود. حال در صورتی که از توان دو برای  $\theta$  استفاده کنیم به آن L2 یا Ridge Regression گفته می شود و در صورتی که از قدر مطلق استفاده گردد به آن L1 یا Lasso Regression گفته می شود. به طور کلی با توجه به اینکه L2 توان ۲ ضرایب را استفاده می کند نتیجه وزن های نهایی کوچک خواهد بود و برای زمانی مناسب است که تعداد ویژگی ها زیاد هستند و feature engineering به شکل

مناسب انجام نشده است اما Lasso Regression سخت گیری کمتری در وزن ها استفاده می کند و وزن های کمتری را صفر می کند( این کار باعث حذف فیچر می شود) در نتیجه برای زمان هایی مناسب است که فیچرها رو با دقت انتخاب کرده ایم. علت استفاده از این نرمال سازی کمک به تشخیص فیچرهای مناسب برای بررسی مدل های مورد نظر و جلوگیری از بیش برآزش می باشد.

f)

از آن جهت که تابع در انتخاب نمونه های آزمایش و تست بصورت رندوم انتخاب می کنیم وجود نویز در دوسته تقسیم خواهد شده و نیازی تغییر در این نسبت نخواهیم داشت.

افزایش درجه برای یادگیری نویز کار چندان پسندیده ای نیست و به طور کلی اشتباه است (در صورتی که اصل موضوع آموزش یادگیری نویز نباشد) و با افزودن درجه، مدل را به سمتی سوق می دهیم که داده های آموزشی را با نویز هایشان حفظ کند ولی این کار باعث می شود که در بخش تست عملکرد ضعیف تر بدلیل **overfitting** از خود نمایش دهد. در صورتی که مدل قبل از افزودن نویز امکان یادگیری کامل داشته باید درجه قبلی را حفظ کند یا اینکه مقدار کمی به درجه آن اضافه گردد.



g)

$$g) \text{ I. } J(\theta) = \frac{1}{m} \sum_{i=1}^m (\underbrace{h(x^i)}_{w^T X_i} - y^i)^2$$

$$\nabla J(\theta) = \frac{2}{m} \underbrace{(h(x^i) - y^i)}_{\text{خطا}} \times \bar{X}$$

5

$$\nabla J(\theta) = \left( \frac{2}{m} \right) \times \text{Err} \times X_i$$

به طور کلی ضریب را نزخ یادگیری  
فرض می کنیم در معادله برای  
به روز رسانی

10

نزخ یادگیری

$$\eta = 1$$

$$w_{old} = \begin{bmatrix} 1 \\ 2 \end{bmatrix}$$

$$X = \begin{bmatrix} 1 & 0.19 \\ 1 & 0.15 \\ 1 & 0.12 \\ 1 & 0.17 \end{bmatrix}$$

$x_1$   $x_2$

$$Y = \begin{bmatrix} 1.1 \\ 1.25 \\ 1.12 \\ 1.40 \end{bmatrix}$$

داده های ورودی

15

$$w_{new} = \begin{bmatrix} 1 \\ 2 \end{bmatrix} + 1 \times \left( \frac{111 - 218}{-114} \right) \times \begin{bmatrix} 1 \\ 0.19 \end{bmatrix}$$

$$= \begin{bmatrix} -0.17 \\ 0.147 \end{bmatrix}$$

5

$$w_{new} = \begin{bmatrix} -0.17 \\ 0.147 \end{bmatrix} + 1 \times \left( \frac{+0.17 = -0.125 + 0.125}{0.185} \right) \times \begin{bmatrix} 1 \\ 0.12 \end{bmatrix}$$

$$= \begin{bmatrix} -0.115 \\ -0.1875 \end{bmatrix}$$

10

↓ این کار را به ترتیبی دیگر انجام می دهیم

~~$$w_c = \begin{bmatrix} -0.115 \\ -0.1875 \end{bmatrix}$$~~

15

نتیجه نهایی

$$w = \begin{bmatrix} 0.109 \\ 0.19287 \end{bmatrix}$$

20

9) II.

$$X = \begin{bmatrix} 1 & .19 \\ 1 & .10 \\ 1 & .15 \\ 1 & .17 \end{bmatrix} \quad Y = \begin{bmatrix} 1.1 \\ .150 \\ .15 \\ .140 \end{bmatrix}$$

$$\theta = (X^T X)^{-1} X^T Y$$

$$X^T X = \begin{bmatrix} 1 & 1 & 1 & 1 \\ .19 & .10 & .15 & .17 \end{bmatrix} \times \begin{bmatrix} 1 & .19 \\ 1 & .10 \\ 1 & .15 \\ 1 & .17 \end{bmatrix} = \begin{bmatrix} 4 & 2.15 \\ 2.15 & .1145 \end{bmatrix}$$

$$\det = .18$$

$$(X^T X)^{-1} = \frac{1}{.18} \begin{bmatrix} .1145 & -2.15 \\ -2.15 & 4 \end{bmatrix}$$

$$X^T Y = \begin{bmatrix} 1 & 1 & 1 & 1 \\ .19 & .10 & .15 & .17 \end{bmatrix} \times \begin{bmatrix} 1.1 \\ .150 \\ .15 \\ .140 \end{bmatrix} = \begin{bmatrix} 2.15 \\ .1145 \end{bmatrix}$$

10

$$\theta = \frac{1}{.18} \begin{bmatrix} .1145 & -2.15 \\ -2.15 & 4 \end{bmatrix} \begin{bmatrix} 2.15 \\ .1145 \end{bmatrix}$$

15

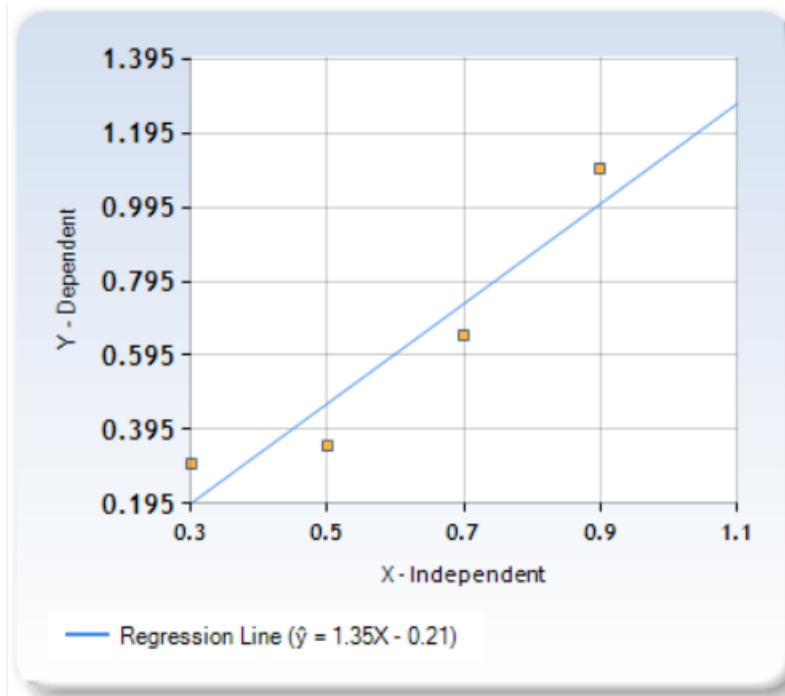
$$= \begin{bmatrix} 1.150 \\ -.121 \end{bmatrix} \quad y = 1.150x - .121$$

III)

normal question برای حالت

Sum of squares ( $SS_x$ ) = 0.2

Sum of products ( $SP$ ) = 0.27



h)

به طور کلی به روز رسانی وزن ها به ازای هر داده موجب ناپایداری در روند همگرایی می شود. احتمال اینکه این ناپایدار باعث عدم همگرایی در انتها شود کم است، اما روند همگرایی را کند تر می کند و نسبت به داده های پرت واکنش بدی نشان می دهد حال چه در اجرای *gradient descent* عادی به ازای هر داده این کار تکرار شود یا در *stochastic gradient descent* بر روی یک نمونه محدود از کل داده ها اجرا شود.

برای بهبود این مشکل می توان از *Batch gradient descent* استفاده نمود که به ازای هر داده به روز رسانی را انجام نمی دهد. بلکه به ازای هر دسته میانگین خطا هر دسته از ورودی ها اینکار را انجام می دهد. از مفروضات این الگوریتم این است که همه دیتاست حتما مورد بررسی قرار گرفت و این فرض دلیل اصلی کند بودن این روش در اجرا می باشد اما روند اجرا پایدار تر بوده و نوسانات شدید ندارد.

در آخر با ترکیب دو روش *Batch* , *stochastic* به روش *mini Batch* می رسیم که فرض استفاده کامل از دیتاست را حذف می کند ولی بصورت *Batch* وزن ها را به روز رسانی می کند.

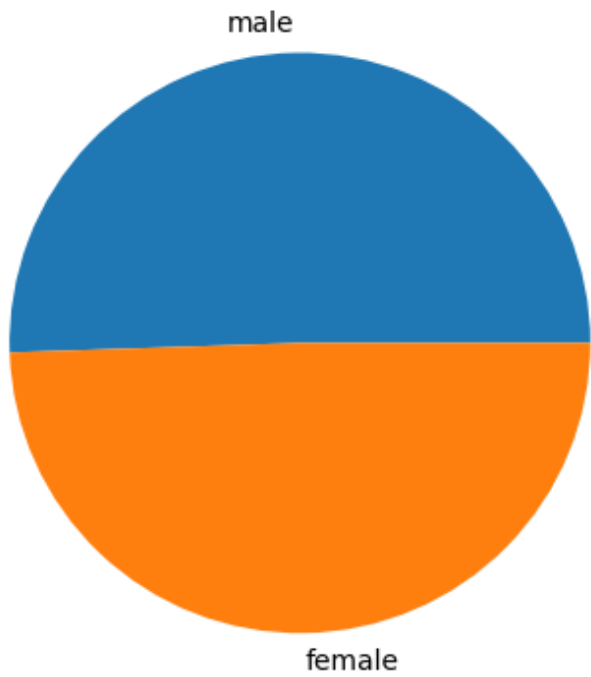
i)

در رابطه با شباهت این دو روش می توان گفت که هر دو آن ها در تلاش هستند تا یک مدل با کمترین خطا مجموع مربعات را برای داده های ورودی خود پیدا کنند تا بتوانند برخی ویژگی دیگر را استخراج نمایند اما تفاوت بنیادی در رفتار این دو این است که در *function estimation* می توان همه متغیرها را متغیرهای آزاد فرض کرد به بیان دیگر هدف یافتن وزن ها است اما در *regression* هدف پیش بینی مقدار متغیر وابسته (هدف) بر اساس سایر متغیرهای آزاد است .

## Q2)

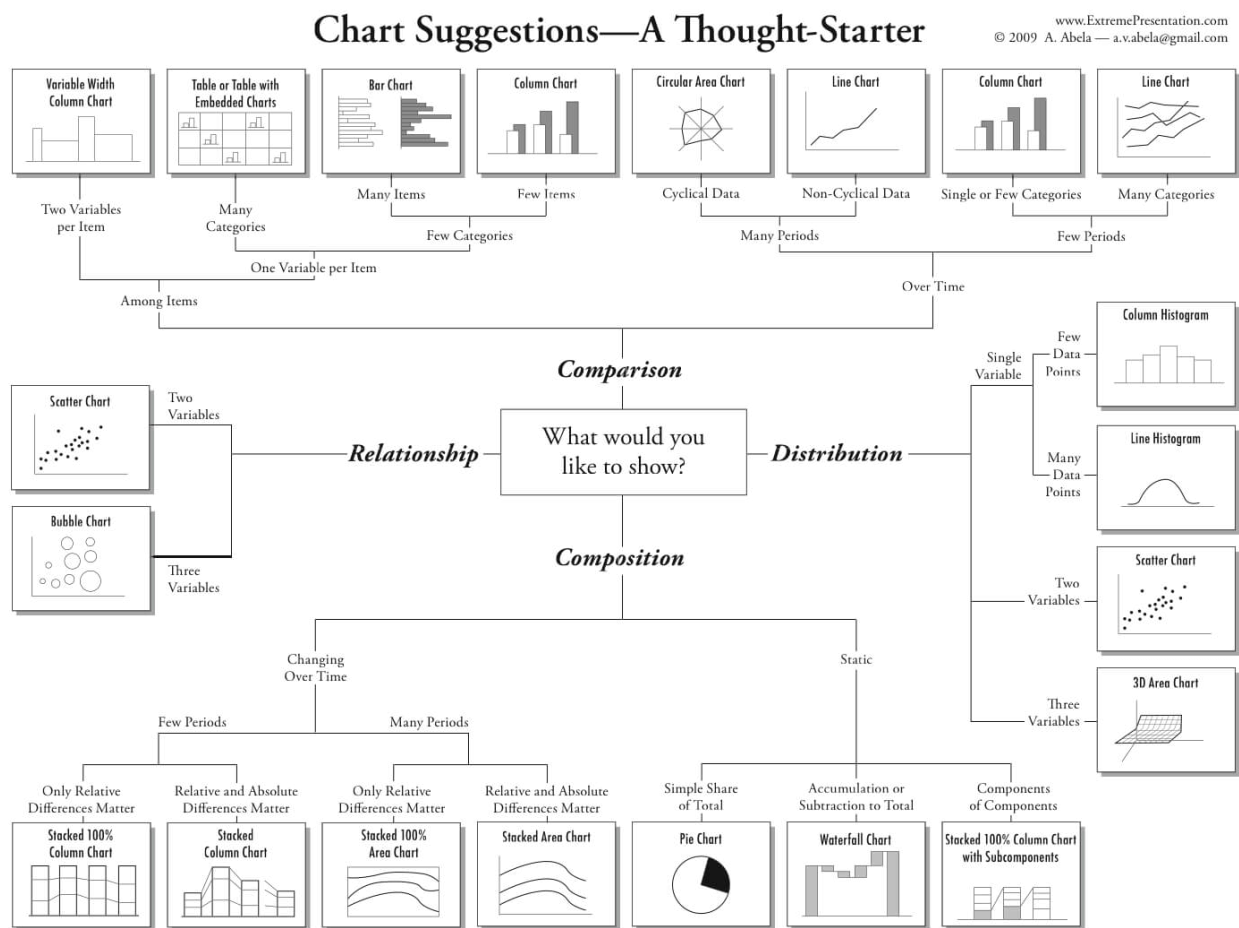
a)

```
C:\Anaconda\python.exe G:/university/master/term3/machineLearning/codes/Q2.py  
male      676  
female    662  
Name: sex, dtype: int64
```

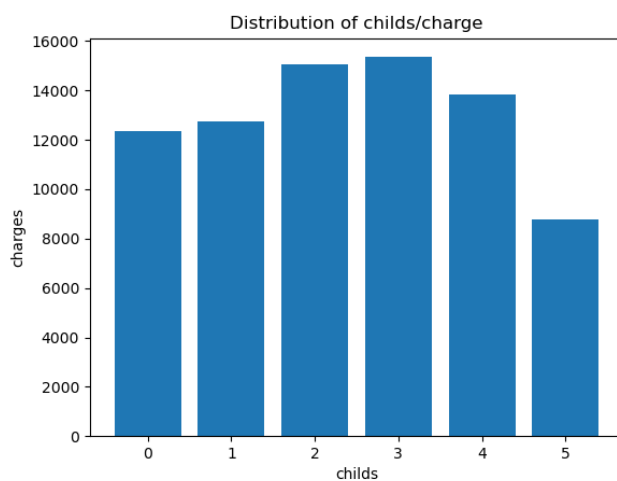


نتایج حاکی از این بود که مقدار اندکی تعداد آقایان بیشتر از خانم ها می باشد.

b)

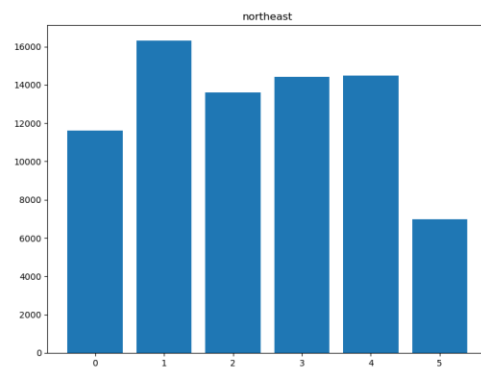
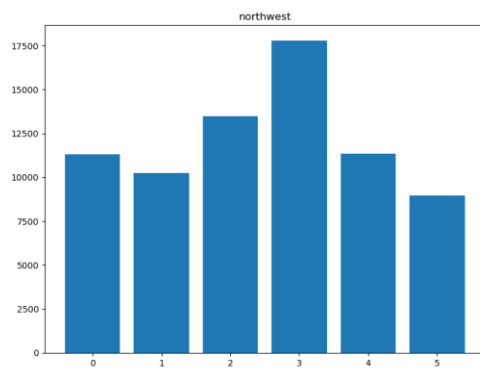
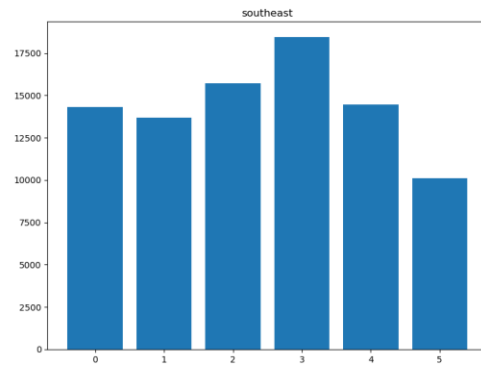
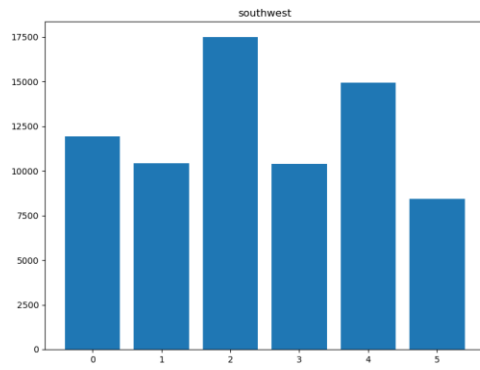


از بین نمودار های ممکن scatter plot نموداری است که برای بررسی ارتباط دو نمودار انتخاب گردیده است

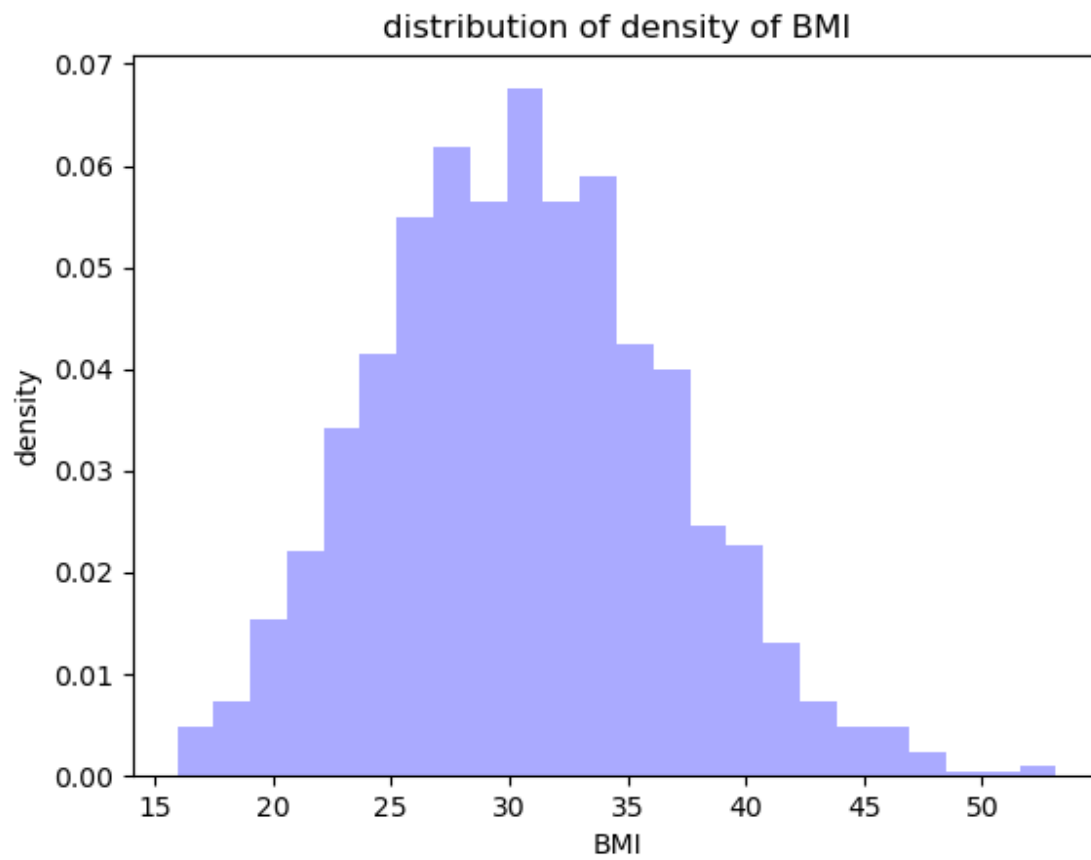


c)

Q 2 part c



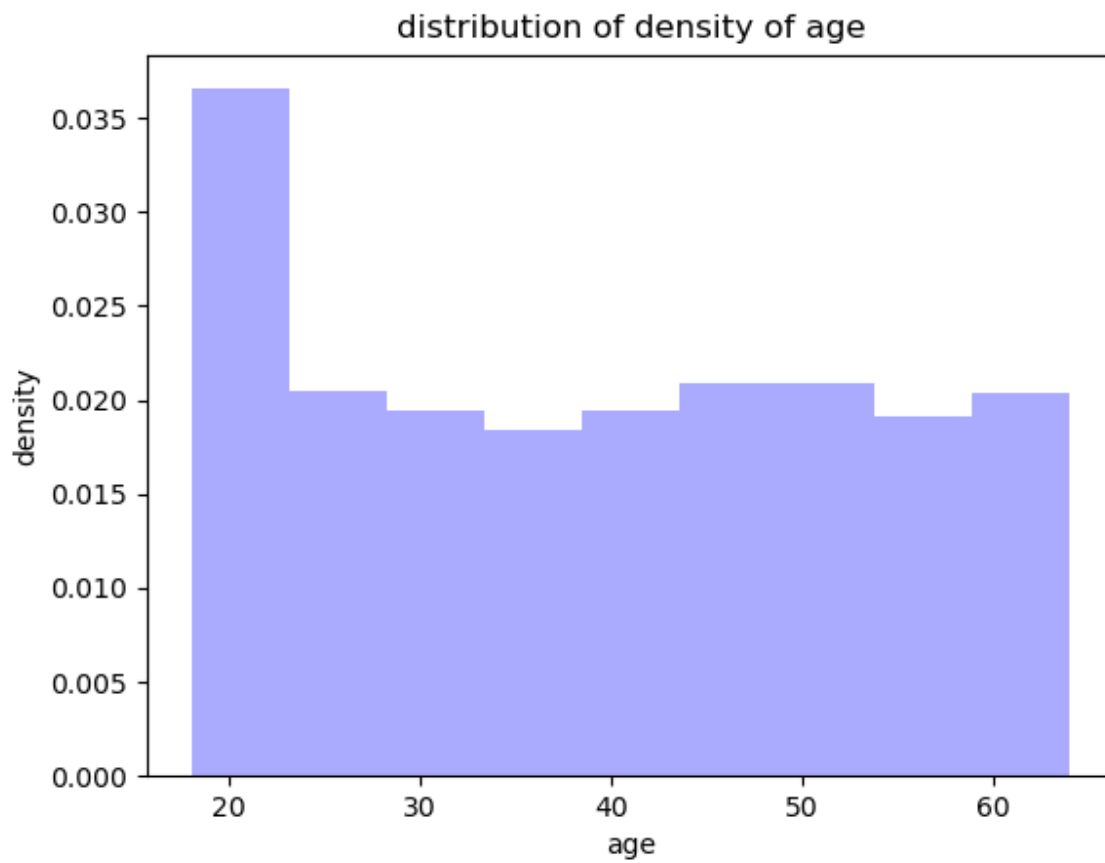
D)





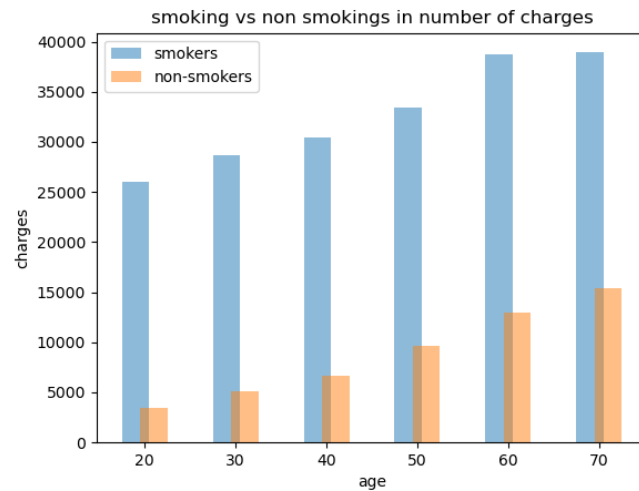
e)

توزیع داده های سن به گونه ای است که به جز بازه ی سنی حدود ۲۰ سال که ۳۵ درصد جامعه را به خود اختصاص داده در بقیه سنین توزیع افراد برابر می باشد.

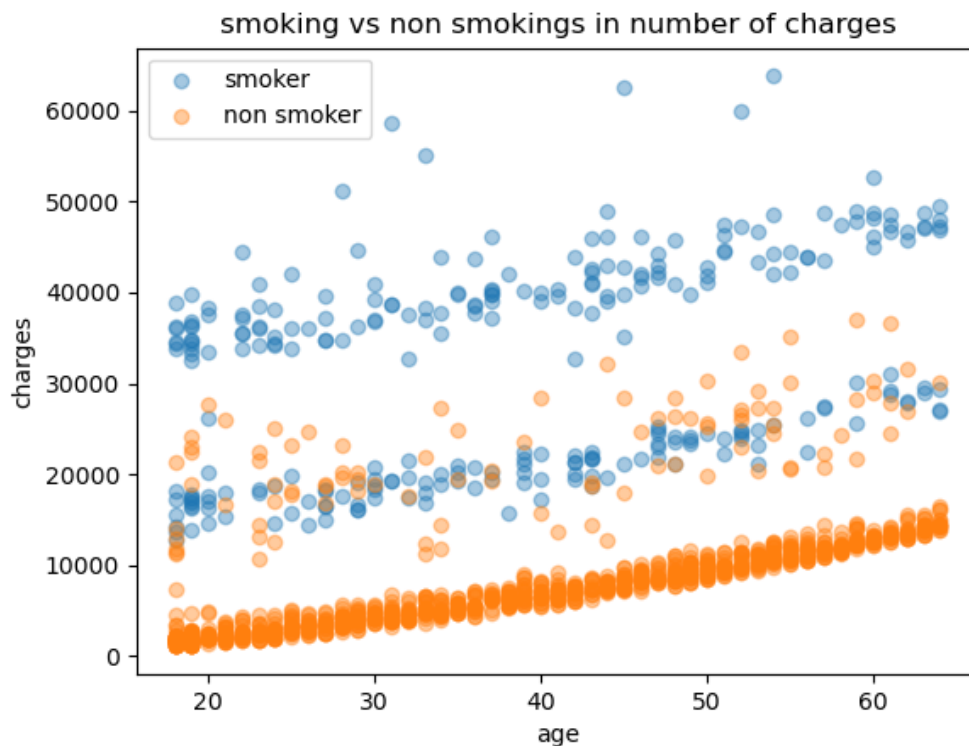


f)

نمودار اولیه بیانگر این است که هیچ وقت smokers از non-smokers میزان دریافتی (charges) کمتری نداشته اند اما برای بررسی دقیق تر می توانیم از نمودار دوم نیز استفاده کنیم.

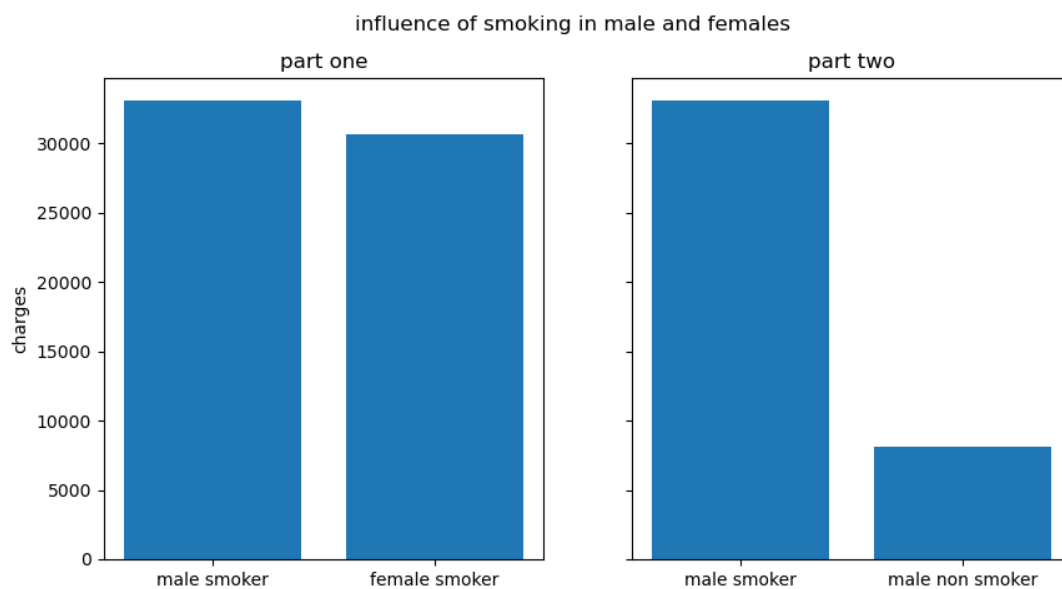


با توجه به توزیع داده ها می توان گفت در همه بازه های سنی کسانی که سیگار مصرف نمی کنند احتمال کمتری دارند اما در بازه ۲۰ تا ۳۰ تعداد افرادی که سیگار مصرف نمی کنند و charge بیشتری داشته اند قابل ملاحظه است همچنین در بازه ۵۰ تا ۶۰ سال نیز این پترن تکرار شده است اما شدت و تراکم آن کمتر است.



g)

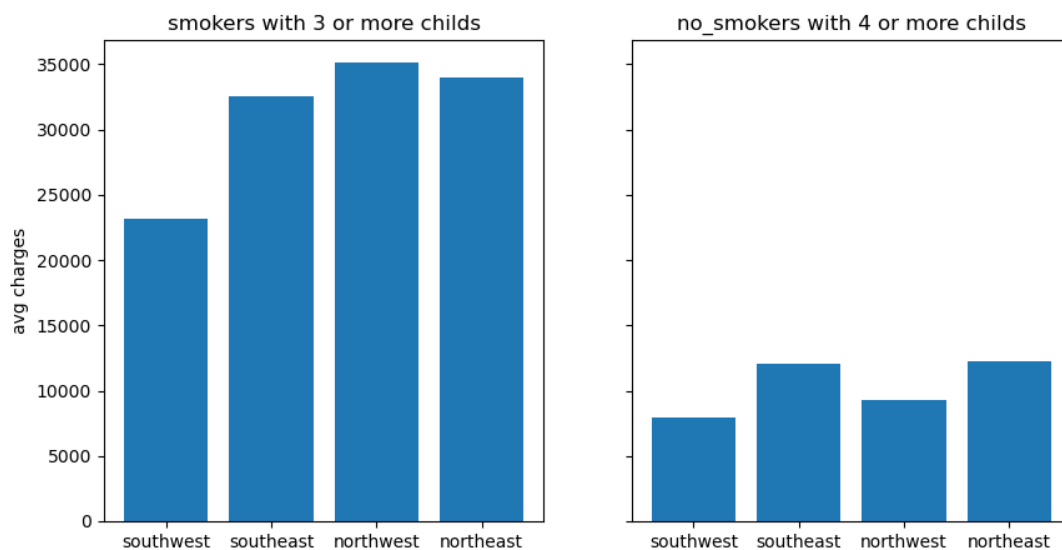
همان طور که در شکل ها مشاهده می کنید میانگین مردان و زنان سیگاری بسیار به یک دیگر نزدیک است (آقایان بیشتر هست البته ) اما مقدار میانگین در بین آقایان دسته ای که دخانیات مصرف نمی کنند آمار charges کمتری دارند در میانگین.



h)

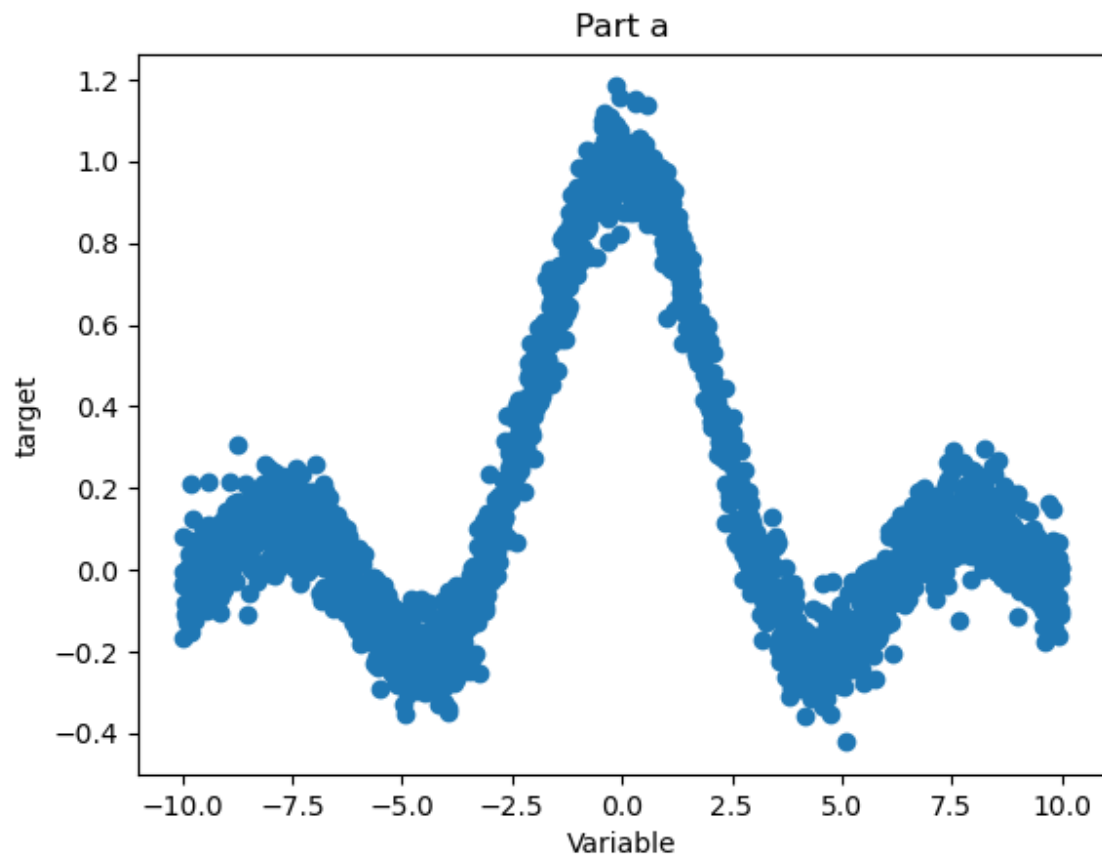
با توجه به شکل های بدست آمده

در همه بخش ها باتوجه به نتایج بدست آمده کسانی که دخانیات مصرف می کنند و بیشتر از ۳ فرزند دارند charges بیشتر خواهند داشت.



# Q3)

a)



b)

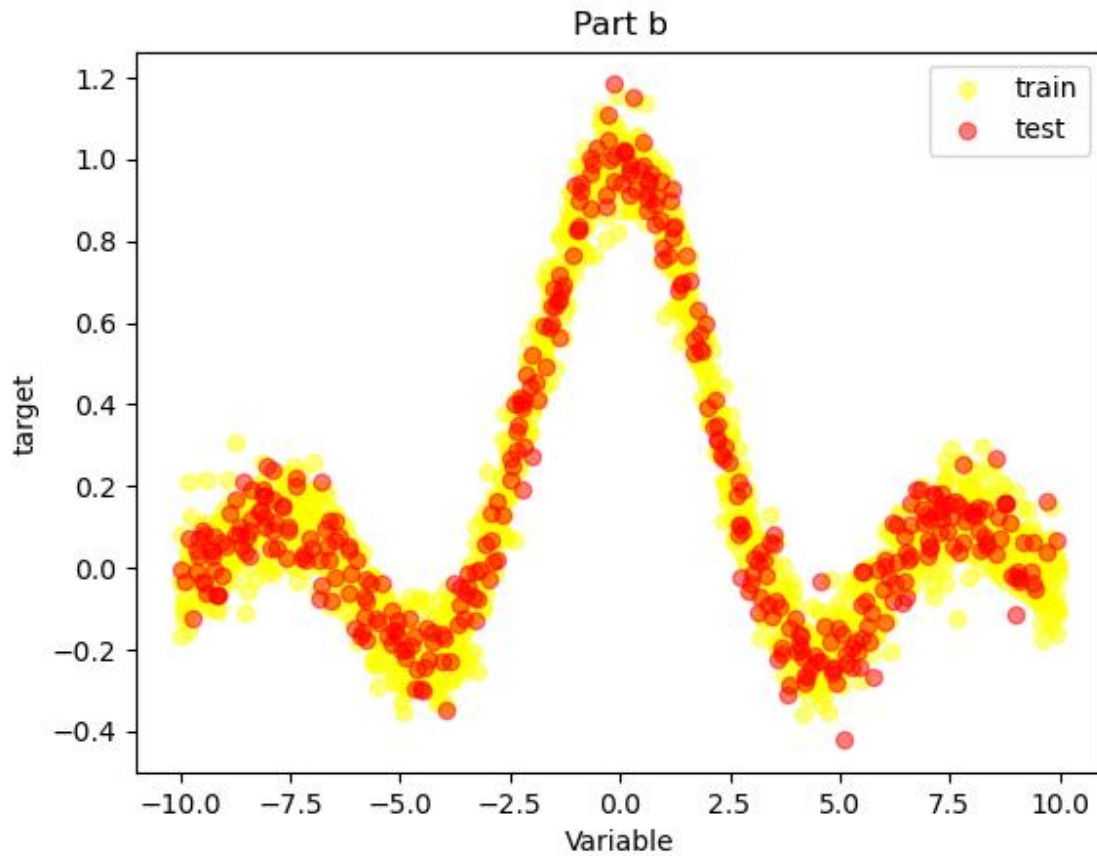
برای اینکه بتوانیم چک کنیم که آیا داده مورد نظر sort شده اند ابتدا دو المان اول را چک کرده که آیا صعودی هستند یا نزولی سپس بقیه المان های داخل سطر نیز باید از این قاعده تبعیت کنند، یعنی اگر صعودی هستند ، صعودی چیده شوند و اگر نزولی هستند نزولی چیده شوند . این فرآیند با یک بار چرخش بر روی داده ها در  $O(n)$  بدست می آید.

نتیجه کد

```
C:\Anaconda\python.exe G:/university/master/term3/machineLearning/HW1/codes/Q3.py  
its not sorted
```

c)

با توجه به شکل به دست آمده توزیع داده های موجود در آموزش و تست یکسان می باشد. لازم به ذکر است که برای اینکار ما از تابع sample استفاده کرده ایم که عملیات shuffling را به طور ضمنی با انتخاب رندوم از بین داده هال انجام می دهد.



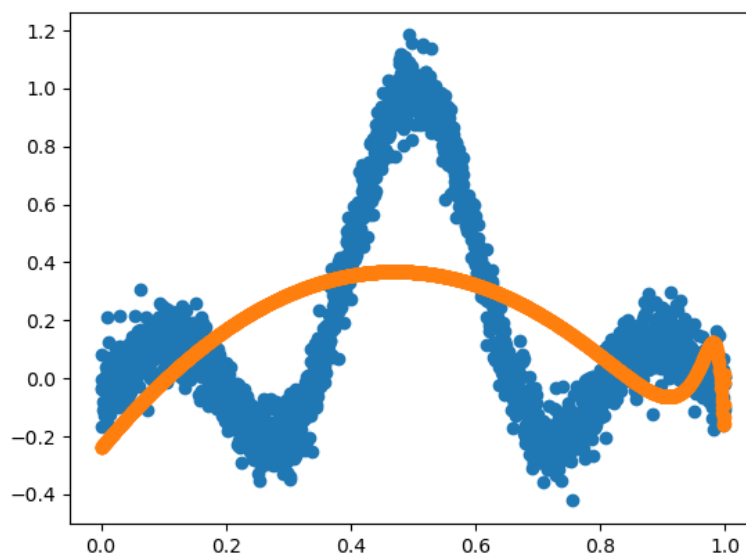
d)

repeat until convergence: {  
     $w_j = w_j - \alpha \frac{\partial J(\mathbf{w}, b)}{\partial w_j}$       for  $j = 0..n-1$   
     $b = b - \alpha \frac{\partial J(\mathbf{w}, b)}{\partial b}$   
}

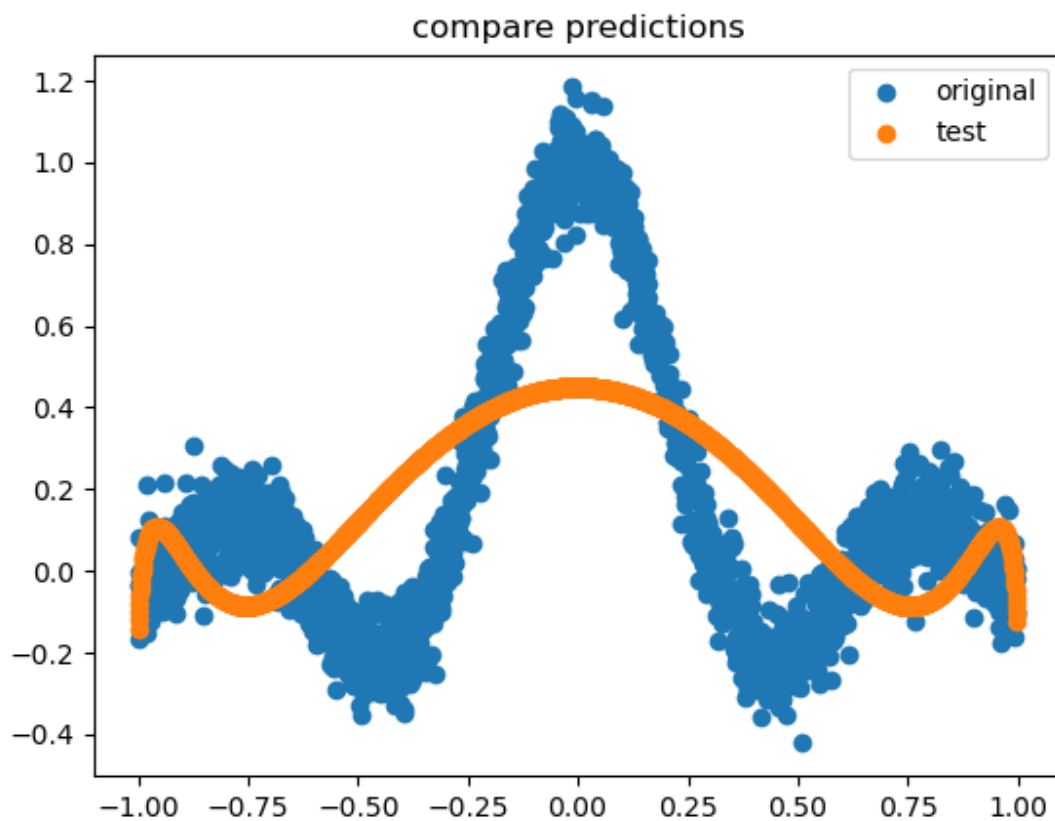
$$\frac{\partial J(\mathbf{w}, b)}{\partial w_j} = \frac{1}{m} \sum_{i=0}^{m-1} (f_{\mathbf{w}, b}(\mathbf{x}^{(i)}) - y^{(i)}) x_j^{(i)}$$

$$\frac{\partial J(\mathbf{w}, b)}{\partial b} = \frac{1}{m} \sum_{i=0}^{m-1} (f_{\mathbf{w}, b}(\mathbf{x}^{(i)}) - y^{(i)})$$

برای پیاده سازی این سوال ابتدا تابع `cost function` را پیاده سازی کردیم و سپس از روی آن `compute_gradient` را اجرا می‌کنیم و شیب تغییرات را به دست می‌آوریم. با استفاده از اعداد بدست آمده در تابع `gradient_descent` مساله را به شکل مورد نظر حل می‌کند. لازم به ذکر است تابع نهایی ما نیاز دارد تا ستون‌های مورد نظر قبل از ورود به تابع ایجاد شده باشند به این معنی که اگر تابع از درجه ۵ می‌باشد باید ستون‌های بعدی به آن اضافه شده باشند.

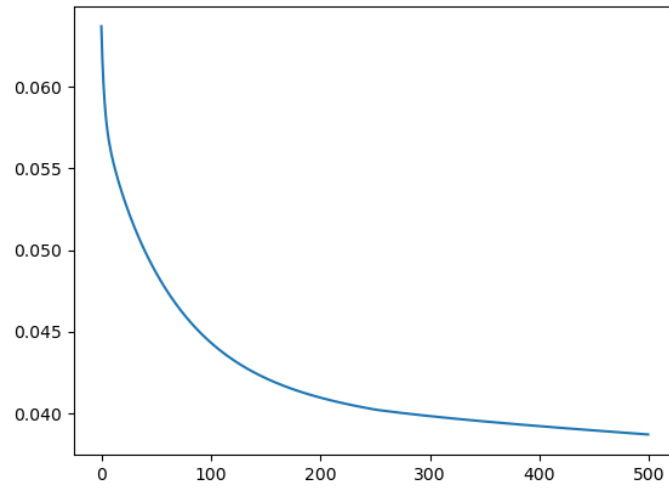




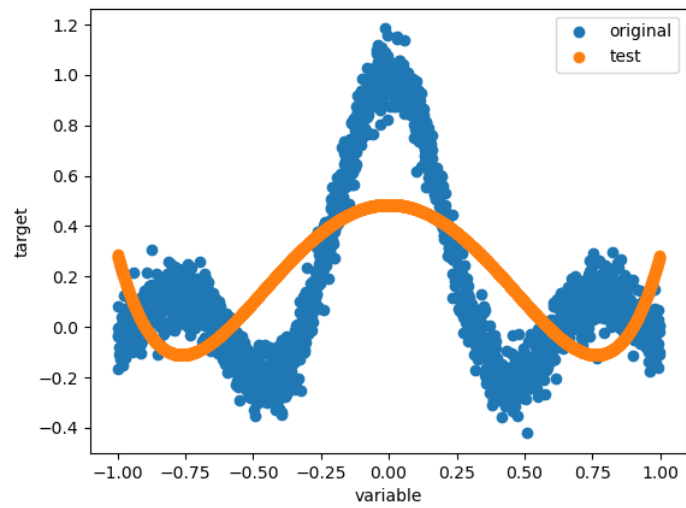


شکل های فوق چند نمونه از اجرای الگوریتم می باشد که شکل آخر حاصل اجرای الگوریتم با ۵ بعد (به علاوه مقدار ثابت که با فرض آن باید گفت شش بعد) در ۵۰۰ گام اجرا شده است. و همچنین به مرور زمان پارامتر آموزش تغییر پیدا می کند. زمانی که ۵۰ درصد iteration ها را انجام داده در ۰.۷۵ ضرب می شود و بعد از اینکه ۷۵ درصد iteration ها را انجام داد به ۰.۴ مقدار اولیه تبدیل می شود و در ۱۰ درصد انتهایی مقدار ۰.۱ مقدار اولیه را دارد. (لازم به ذکر است باتوجه به اینکه توان های داده های از ماکسیمم مقدار float قابل محاسبه بیشتر می شد مقادیر ورودی در ابتدای کار بر ۱۰ تقسیم می شوند، حالت های مختلف نورمال سازی بررسی شد اما در نهایت همین scale ساده بهترین نتیجه را داد، اولین نمودار این بخش مربوط به زمانی است که هر کامپلکسیتی را مستقل از بقیه نورمال کنیم این کار بعضا باعث از بین رفتن تقارن در شکل نهایی بدست آمده می شد)

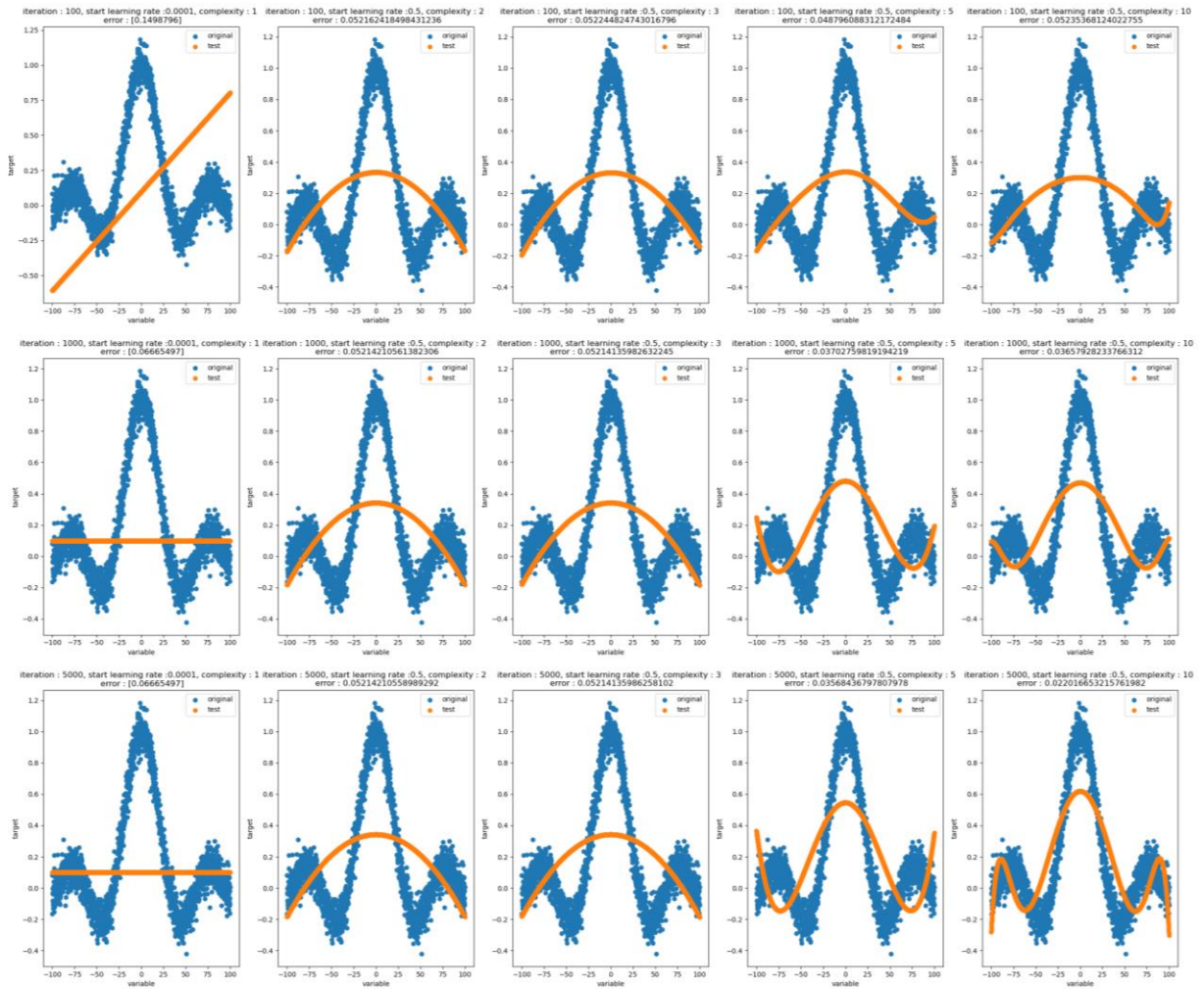
iteration : 500, start learning rate :0.99, complexity : 7  
error : 0.038727208404963076

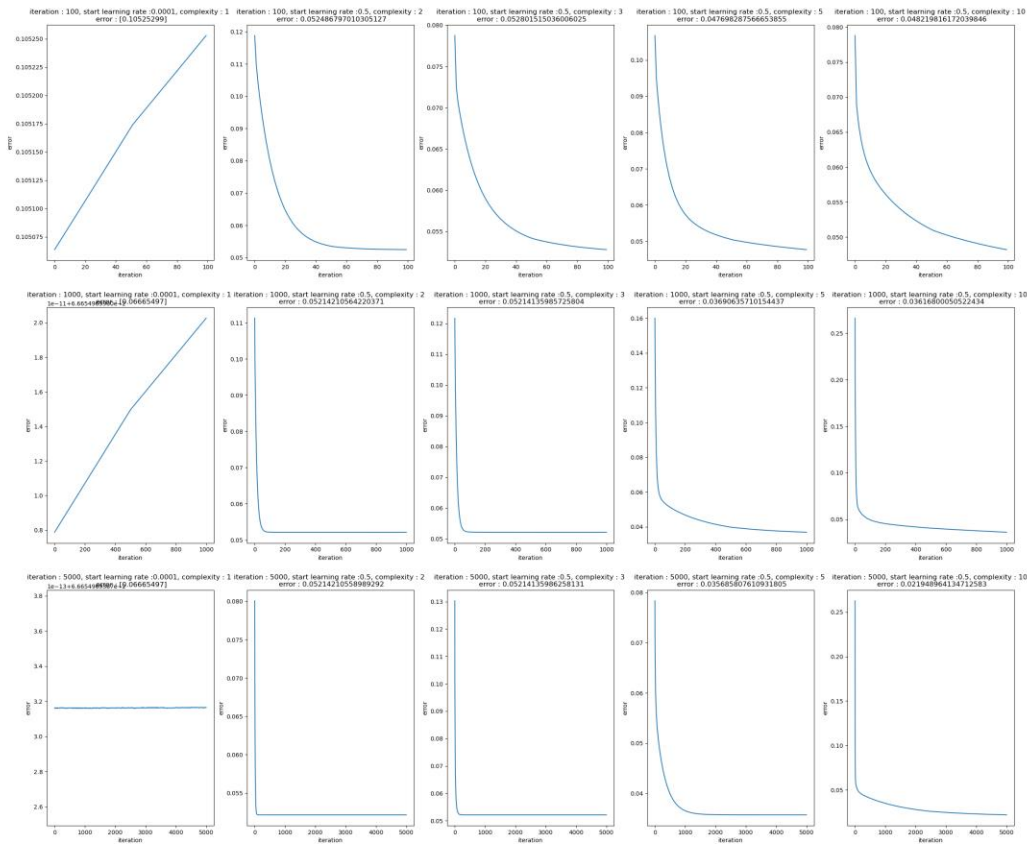


iteration : 500, start learning rate :0.99, complexity : 7  
error : 0.038727208404963076



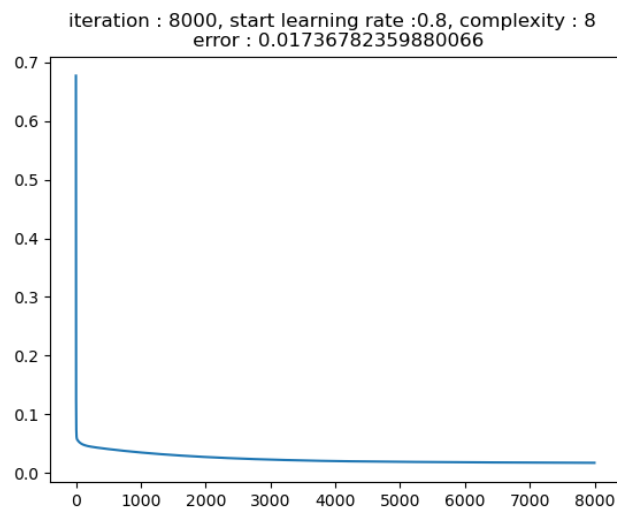
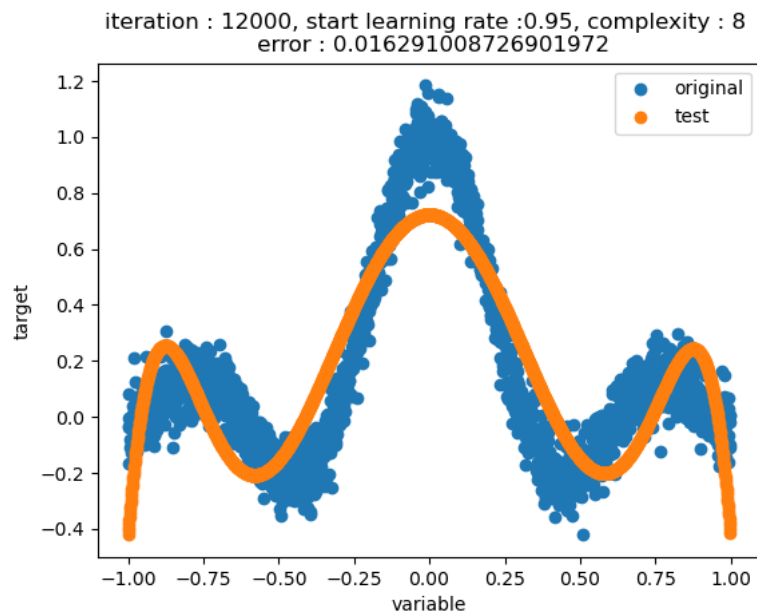
e)





F)

با توجه به مشاهدات و مقارن بودن تابع نسبت به محور  $y$  قطعاً توابع با درجه توانی زوج کاندید های مناسب تری هستند نظر شخصی بنده استفاده از درجه ۸ بود با iteration بیشتر که نتیجه نهایی را می توانید در زیر مشاهده کنید

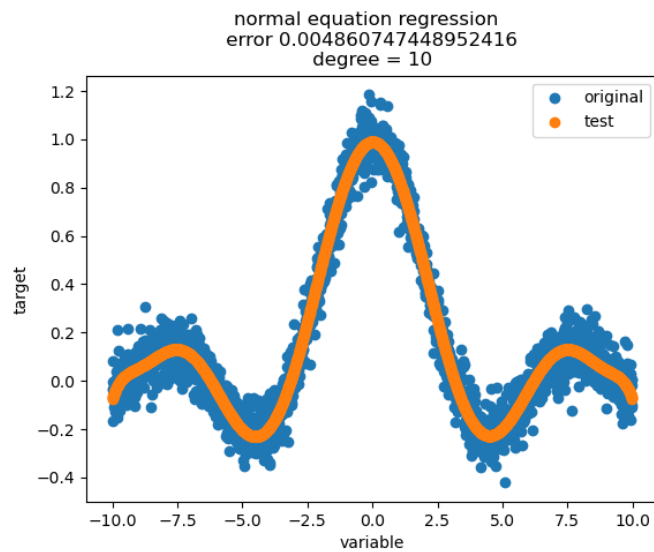
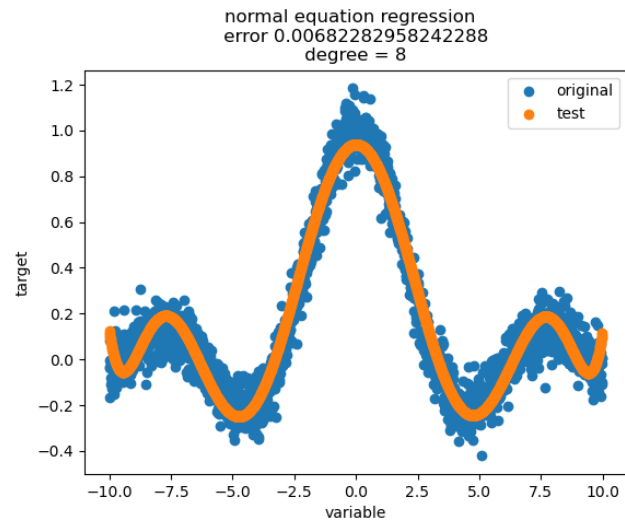


g)

خطای ذکر شده MSE می باشد

همانطور که در اسلاید های درس نیز گفته شد برای استفاده از normal equation به شکل زیر عمل شده است.

$$\theta = (X^T X)^{-1} X^T y$$

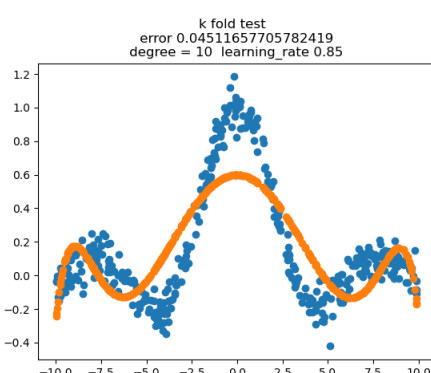
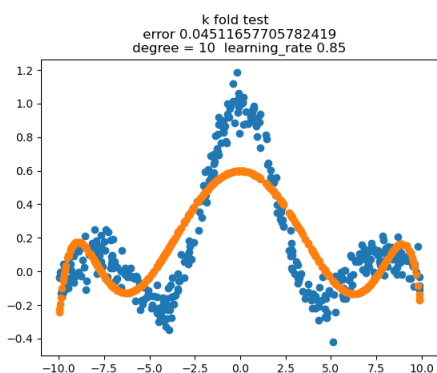
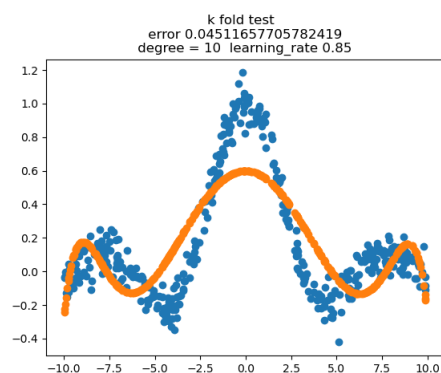
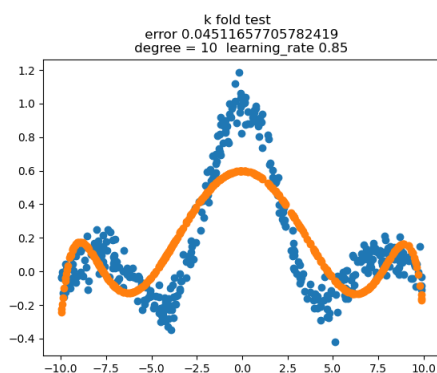


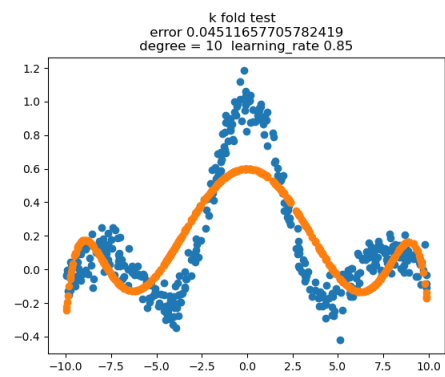
H)

برای حالات مختلف میانگین خطا را در جدول زیر مشاهده می کنید در همه حالات ۲۵۰۰ iteration مورد نظر قرار گرفته همچنین درجه مستقل از بایاس بیان شده است. دلیل استفاده نکردن از ۰.۹۹ در همه حالات این است که این کار در برخی مدل ها باعث ایجاد ناپایداری می شود.

Learning rate (initial)	complexity	err
0.99	5	۰/۰۷۱۷۷۵۱۶۴۹۸۳۴۲۹۸۳
0.99	6	۰/۰۵۹۹۰۱۵۳۳۰۳۴۳۳۴۹۴۴
0.99	8	۰/۰۴۶۳۰۲۳۶۰۴۵۷۶۵۸۱۴
0.75	10	۰/۰۴۷۱۹۰۸۸۷۱۹۴۰۸۲۱۵۴
0.85	10	۰/۰۴۵۶۸۵۸۱۰۱۸۳۲۵۶۱۴
0.75	12	۰/۰۴۸۰۲۶۷۴۷۵۱۴۶۵۲۱۹

۵ تصویر مرتبط با بهترین نتیجه می باشد که در آن ۰.۸۵ نرخ یادگیری و پیچیدگی ۱۰ می باشد.







## Q4)

a)

نتیجه اجرای الگوریتم با نرخ یادگیری ۰/۰۱ به شکل زیر بوده است

```
C:\Anaconda\python.exe G:/university/master/term3/machineLearning/HW1/codes/Q4.py
main is called
----- a -----
R2 erro :0.8288824795479934
R2 prediction 0.8110443153477049
***

Process finished with exit code 0
```

b)

در این جا یک مجموعه توانی (همه زیرمجموعه های ممکن از مجموعه ویژگی ها) را ایجاد می کنیم و مقدار خطاها را برای هر یک محاسبه می کنیم و در نهایت بهترین را نیز اعلام می کنیم. بدلیل کمبود جا بخشی از خروجی ها را در پایین آورده ایم

```
Q4 x
C:\Anaconda\python.exe G:/university/master/term3/machineLearning/HW1/codes/Q4.py
main is called
----- b -----
---
['Length1', 'Length2', 'Length3', 'Height', 'Width']
R2 : 0.8603779696520067
predicted R2 : 0.8473633736065198
---
---
['Length2', 'Length3', 'Height', 'Width']
R2 : 0.8637009216935334
predicted R2 : 0.8523416804123396
---
---
['Length1', 'Length3', 'Height', 'Width']
R2 : 0.8638204165379262
predicted R2 : 0.8523690140628792
---
---
['Length3', 'Height', 'Width']
R2 : 0.8648007739165349
predicted R2 : 0.8552745876256014
---
---
['Length1', 'Length2', 'Height', 'Width']
R2 : 0.8630970598981349
predicted R2 : 0.8525497102978596
---
```

```
Q4 x
---
['Length1', 'Length3']
R2 : 0.8490245172285553
predicted R2 : 0.8430527179768228
---
---
['Length3']
R2 : 0.8518319615261469
predicted R2 : 0.8486967521040832
---
---
['Length1', 'Length2']
R2 : 0.8405206260726797
predicted R2 : 0.8341513790768086
---
---
['Length2']
R2 : 0.8429201077124555
predicted R2 : 0.8395186777259532
---
---
['Length1']
R2 : 0.8374195298066425
predicted R2 : 0.833799947729293
---
****
best answer
['Length1', 'Height', 'Width']
0.872659063110441
```

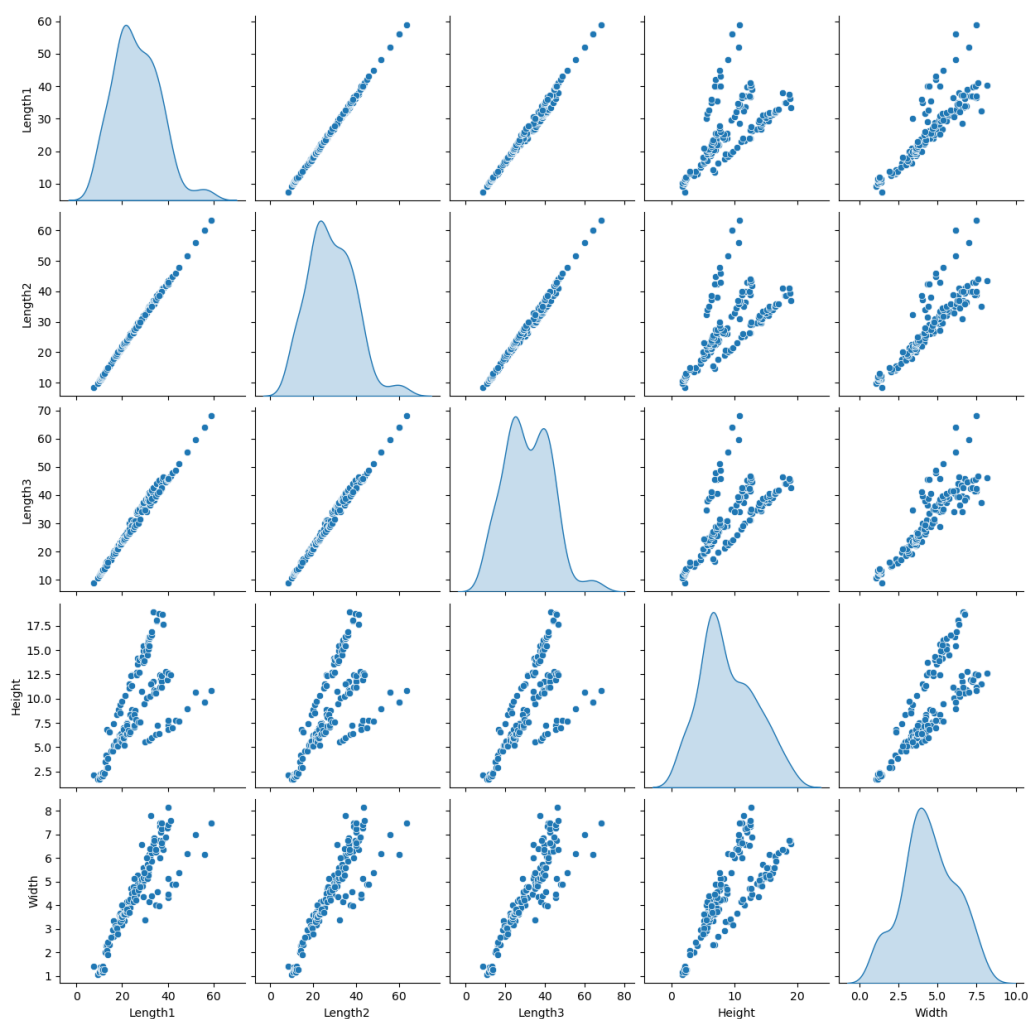
در نهایت بهترین جواب length1 , Height , Width

c)

```
Q4 ×
C:\Anaconda\python.exe G:/university/master/term3/machineLearning/HW1/codes/Q4.py
main is called
----- c -----
Weight ~ Length3 + Width + 1
0.8757636700691216
***
Process finished with exit code 0
```

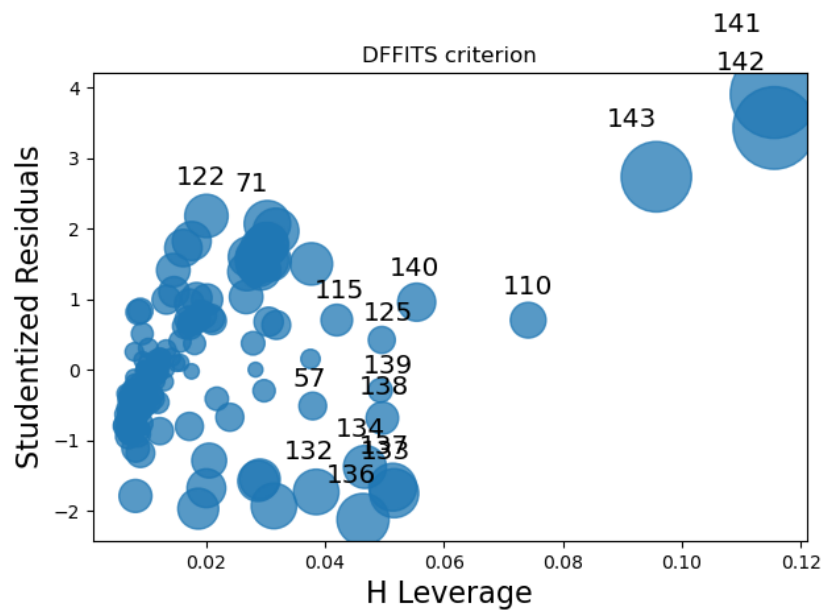
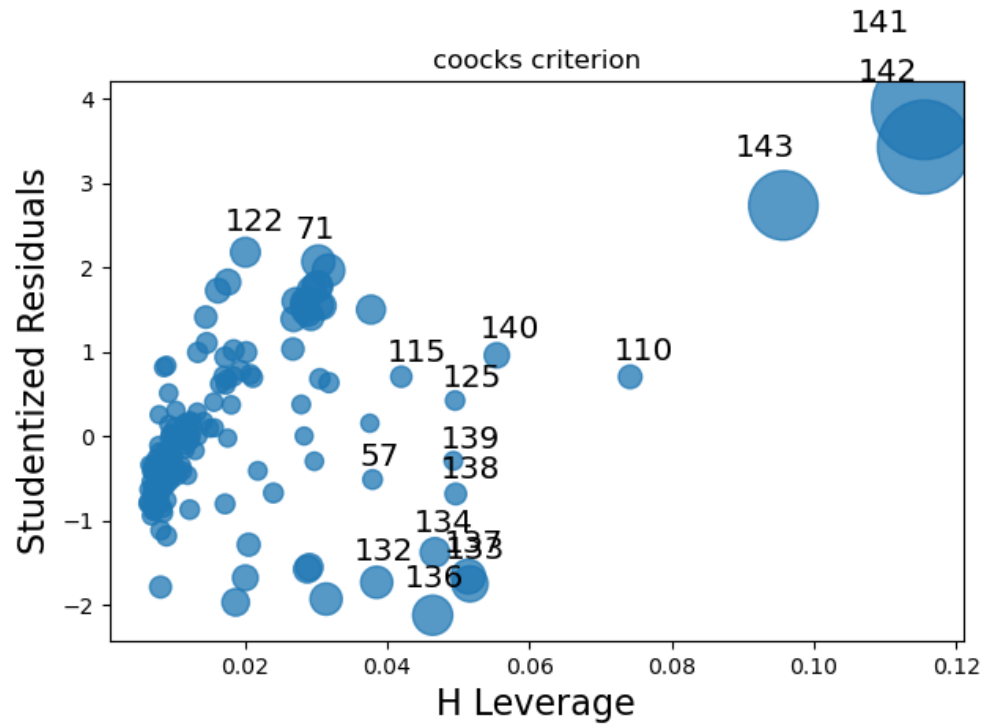
نتیجه این روش تقریباً با روش قبلی از نظر خطا برابر است اما با یک عنصر کمتر که این کار باعث بهبود الگوریتم می شود.

d)



نمودار فوق نسبت دو به دو پارامترهای ورودی ما را نمایش می دهد از آن جا که تعدادی از خانه های جدول حالت خطی به گرفته اند در نتیجه از روی یکی دیگری را می توان محاسبه نمود و در نتیجه multicollinearity در داده های ما وجود دارد.

e)



f)

هنوز ایده ای ندارم

g)

از دو فیچری که در بخش c بدست آمد استفاده کردیم و نورمال سازی نیز انجام دادیم و مساله را با درجه ۲ تحلیل نمودیم.  
حاصل مدل نهایی به شکل زیر بدست آمد.

OLS Regression Results			
=====			
Dep. Variable:	Weight	R-squared:	0.966
Model:	OLS	Adj. R-squared:	0.965
Method:	Least Squares	F-statistic:	862.1
Date:	Sat, 29 Oct 2022	Prob (F-statistic):	1.37e-109
Time:	13:44:54	Log-Likelihood:	-885.69
No. Observations:	158	AIC:	1783.
Df Residuals:	152	BIC:	1802.
Df Model:	5		
Covariance Type:	nonrobust		
=====			

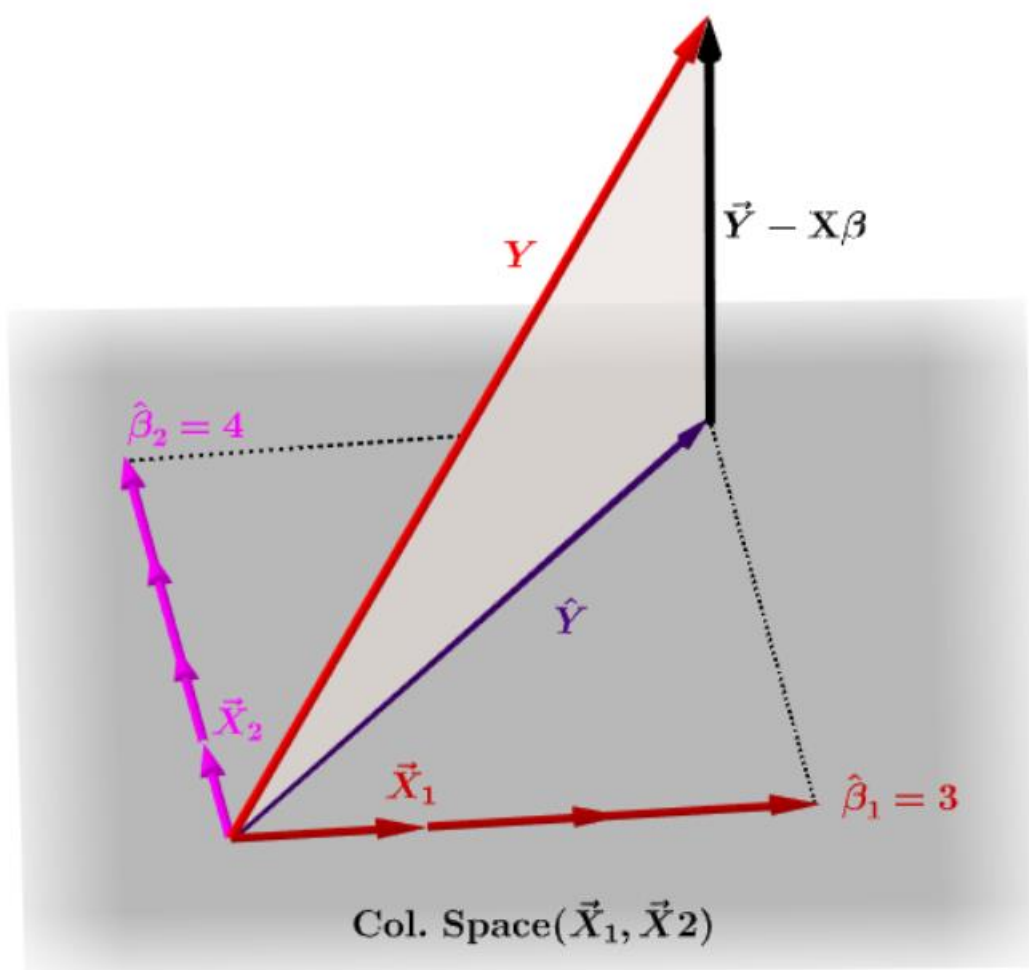
## Q5)

a)

ماتریس **Hat** ماتریسی است که با اعمال شدن به مقدار واقعی داده ورودی مقدار تخمین زده شده توسط رگرسیون را برای آن نشان می دهد. به بیان دیگر مقدار واقعی داده ورودی را بر روی شکل نهایی رگرسیون تصویر می کند. (orthogonal projection)

$$\hat{y} = Hy$$

به کمک این ماتریس می توانیم میزان تاثیر گذاری هر المان در مدل نهایی را با استفاده از مولفه عمود بر مدل نهایی بدست آوریم.



$$\mathbf{H} = \mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T.$$

همچنین برای بدست آوردن هر یک از  $\beta_i$  ها از حاصل ضرب زیر کمک می گیریم که همان رابطه ای است که در بالا آورده شده است.

$$(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{Y}.$$

b)

step wise regression یک مدل iterative در ساخت رگرسیون می باشد که در آن در هر مرحله به انتخاب متغیرهای مستقل (وزن فیچرها) یا حذف آن ها می پردازد. در این روش در هر مرحله سعی در انتخاب بهترین متغیر و یا حذف بدترین متغیر در جهت ساخت مدل رگرسیون ما انجام می شود.

به طور کلی ۳ روش برای استفاده از step wise بیان می شود: (منظور از میزان تاثیر گذاری در حذف یا افزودن متغیرها در اینجا F-test و T-test می باشد).

۱. Forward selection : در این حالت مدل از حالتی که هیچ متغیری ندارد آغاز می کند و در هر مرحله بهترین وزن یا متغیری که بیشترین تاثیر آماری در مدل را دارد به مساله اضافه می کند و این کار را تا زمانی انجام می دهد که به جواب optimal دست یابد. (به بیان دیگر تا زمانی ادامه می دهد که وزن های بعدی تاثیر چندانی در نتیجه نهایی نداشته باشند)
۲. Deleting one at time: بر خلاف روش قبلی فرض می کند همه فیچر ها در مدل رگرسیون حضور دارند سپس هر کدام را به تنهایی تست می کنیم. المانی که کمترین تاثیر را بر روی مدل داشت حذف می کنیم. (لازم به ذکر است که در این حالت یک threshold مد نظر می گیریم که حتما یال حذف شده باید از آن threshold تاثیر کمتری داشته باشد. همچنین مساله زمانی به پایان می رسد که هیچ متغیری کمتر از threshold تاثیر گذاری نداشته باشد).
۳. Bidirectional elimination: این حالت حالت ترکیبی از یک و دو می باشد که در برخی از گام ها ممکن است که یک متغیر انتخاب شده را حذف کند و متغیری جدید را به مجموعه در همان گام اضافه کند.

c)

### Multicollinearity

این مشکل زمانی رخ می دهد که دو یا چند تا از متغیرهای مستقل (ویژگی ها) همبستگی زیادی نسبت به هم داشته باشند و به بیان دیگر به توان از روی یکی دیگری را محاسبه نمود. این مشکل باعث ایجاد خطا چندانی در پیش بینی ما نخواهد شد و همچنین به قابل اطمینان بودن مدل ما نیز خدشه ای وارد نمی کند. اما ممکن است وزن هایی که به متغیرهایی که همبستگی زیاد دارند به شکل غیر قابل تفسیر اختصاص داده شود اما با وجود این شرایط نتیجه مجموع وزن ها به درستی کار می کند. (برای مثال فرض کنید ما یک مفهوم را دو بار به عنوان ویژگی به مدل می دهیم ممکن است که به یکی از آن ها وزن ۰.۳ و به دیگری ۰.۷ داده شود در صورتی که منطقاً با توجه به اینکه یک مفهوم بوده اند باید وزنی معادل به آن ها داده می شد)

راه شناسایی

۱. تغییرات زیاد در وزن های رگرسیون با حذف یا افزودن یک متغیر (VIF (Variable Inflation Factors)

۲. وجود ضرایب نامتناسب نسبت به بقیه ضرایب استفاده شده در مدل رگرسیون

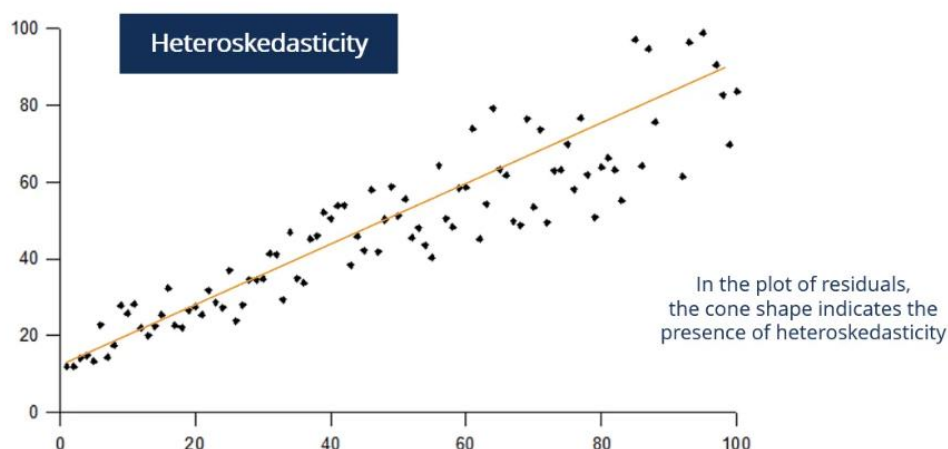
۳. بعضاً علامت ضرایب با آن تئوری منطبق نیست و پارامتر های تاثیر گذار در جهت مثبت، تاثیر منفی به خود بگیرند.

راه حل

همان طور که واضح هست باید ستون هایی که شناسایی شده اند که با بقیه همبستگی دارند را حذف کنیم، در مواردی که همبستگی نه به قدری نیست که بتوانیم حذف و نه به قدری است بتوانیم نگه داریم می توانیم از تعریف متغیر جدید که ترکیبی از این دو می باشد استفاده کنیم و برای آن یک ضریب در مدل رگرسیونمان اختصاص دهیم.

## Heteroscedasticity

این پدیده زمانی اتفاق می افتد که خطا (residual) متغیر پیش بینی شده در بازه های متفاوت یکسان نباشد برای مثال داده هایی که از مقدار آلفا بزرگ تر باشند احتمال خطای بیشتری داشته باشند نسبت به آنهایی که کوچکتر از آلفا پیش بینی شده اند. یکی از شایع ترین حالات این مشکل این است که نمودار خطا نسبت به مقدار پیش بینی شده به شکل مخروطی در می آید.



این پدیده برای حالتی که از روش ordinary least square (ols) regression برای بدست آوردن ضرایب استفاده می کنیم ایجاد خطا می کند زیرا این روش برای محاسبه وزن ها فرض بر ثابت بودن واریانس متغیرهای آزاد جدید در همه بازه های زمانی می کند و در صورتی که پدیده Heteroscedasticity اتفاق بیفتد پاسخ این روش دیگر مناسب نخواهد بود.

### شناسایی

راه اول برای شناسایی آن نمایش آن و یافتن رابطه مخروطی شکل در داده های موجود است.

راه دوم : heteroscedasticity به دو دسته pure , impure تقسیم می شود در حالت خالص روش ترسیم توصیه می شود اما در حالت impure علت بوجود آمدن مشکل وجود تعداد زیاد فیچرها یا تعداد خیلی کم فیچ ها است در این حالت می توان با اصلاح ویژگی ها این مشکل را برطرف نمود

راه سوم : استفاده از p\_value که در صورتی که برای ویژگی های انتخاب شده این مقدار پایین باشد می تواند از علائم این مشکل باشد.

### راه حل های حل مساله

- استفاده از مقادیر لگاریتمی به عنوان ویژگی که جلوی ایجاد واریانس های زیاد را در حالت نهایی می گیرد. ( این کار برای داده هایی که بازه مقادیری گسترده دارند بیشتر توصیه می شود )
- استفاده از تبدیل های غیر خطی برای ستون هایی که این مشکل را دارند (این حالت، در اصل حالت کلی تر پیشنهاد اول می باشد)
- استفاده از روش هایی به جز ols مانند MINQUE , Heteroscedasticity-consistent standard errors , weighted least squares



d)

#### DFFITS

معیار DFFITS اندازه گیری تغییر در مقدار پیش بینی شده (میزان تاثیرگذاری) برای مشاهده نام است و با حذف مشاهده نام محاسبه می شود. یک مقدار بزرگ نشان می دهد که داده ی ورودی جدید در این فضا برای مدل رگرسیون تاثیر بسزایی داشته است.

DFFIT حاصل مقداری پیشنهادی با حضور داده منهای مقدار پیشنهادی بدون داده مورد نظر است.

S خطای استاندارد در حالت بدون در نظر گرفتن داده

h برابر leverage می باشد که همان ماتریس hat می باشد که از قطر آن در فرمول DFFITS استفاده شده است.

$$DFFIT = \widehat{y_i} - \widehat{y_{i(i)}}$$

$$DFFITS = \frac{DFFIT}{s_{(i)} \sqrt{h_{ii}}}$$

$$h_{ii} = \frac{\partial \widehat{y}_i}{\partial y_i}$$

#### Cooks distance

یک معیار آماری که با آن در رگرسیون خطی به تشخیص outlier ها می پردازند. در این معیار بر اساس دو عامل اندازه گیری می شود. میزان فاصله مقدار تخمین زده شده با مقدار اصلی (residual) و میزان نفوذ leverage که در بالا تعریف شده است.

$$D_i = \frac{\sum_{j=1}^n (\hat{Y}_j - \hat{Y}_{j(i)})^2}{(p+1) \hat{\sigma}^2}$$

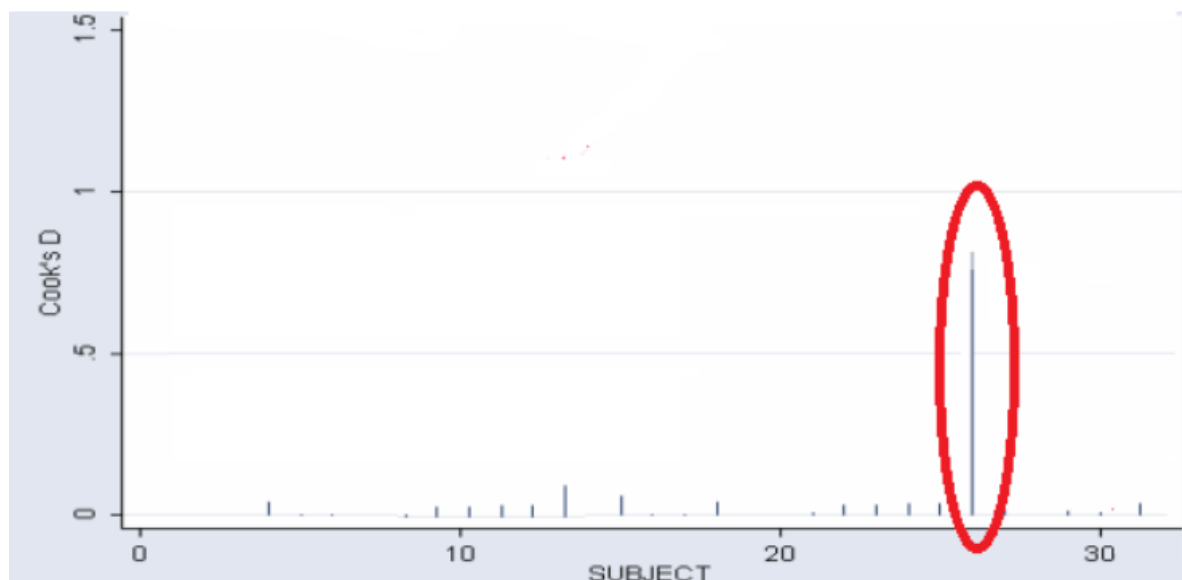
در صورت فرمول مربع خطای بین مقدار تخمین زده شده و مقدار حقیقی با حذف داده مورد نظر از دیتاست جمع می شود.

P تعداد ضرایب موجود در مدل می باشد.

$\sigma$  نیز انحراف معیار خطا به توان ۲ می باشد.

$$\hat{\sigma}^2 = \frac{1}{n-m} \sum_{j=1}^n \hat{\varepsilon}_j^2.$$

البته لازم به ذکر است که حالات مختلفی از این فرمول وجود دارد که فرق هایی جزئی با یکدیگر دارند.



نمونه کشف یک outlier

e)

بخاطر اینکه در این روش از نمایش یک خط در فضای حالات برای متغیر وابسته استفاده می کنیم، زمانی که واریانس خطا از تابع نرمال پیروی نمی کند به این معنا است که استفاده از رگرسیون خطی اشتباه بوده و جواب مساله در بازه های مختلف ضرایب خطاهای متفاوتی خواهد داشت (میزان خطای پیش بینی شده به شدت متفاوت می باشد) از این رو می بایست از روش های غیر خطی یا روش هایی مانند weighted least squares , Heteroscedasticity-consistent standard errors , MINQUE که این مشکل را ندارند استفاده نمود.