

به نام خدا

تمرین سوم یادگیری ماشین

غزل زمانی نژاد

.. ۱

1)

$$1.1) y = \beta_0 + \epsilon_i$$

برای این مسئله regression تابع ضرر را SSE تعریف میکنیم:

$$L = \sum_{i=1}^n (y_i - \beta_0)^2$$

$$\frac{\partial L}{\partial \beta_0} = -2 \sum_{i=1}^n (y_i - \beta_0) = 0 \rightarrow \sum_{i=1}^n y_i = n \beta_0 \rightarrow \beta_0 = \frac{\sum_{i=1}^n y_i}{n}$$

$$\rightarrow \beta_0 = \frac{1}{10} \sum_{i=1}^{10} y_i = \frac{1}{10} (561) = 56.1$$

$$1.2) y_i = \beta_1 x_i + \epsilon_i$$

$$L = \sum_{i=1}^n (y_i - \beta_1 x_i)^2 \rightarrow \frac{\partial L}{\partial \beta_1} = -2 \sum_{i=1}^n (y_i - \beta_1 x_i) x_i = 0$$

$$\rightarrow \sum_{i=1}^n y_i x_i - \beta_1 \sum_{i=1}^n x_i^2 = 0 \rightarrow \sum_{i=1}^n y_i x_i = \beta_1 \sum_{i=1}^n x_i^2 \rightarrow \beta_1 = \frac{\sum_{i=1}^n y_i x_i}{\sum_{i=1}^n x_i^2}$$

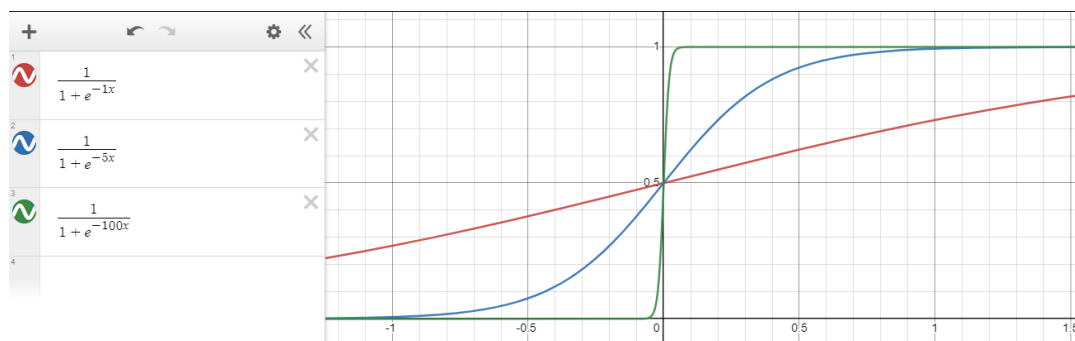
$$\rightarrow \beta_1 = \frac{\sum_{i=1}^{10} y_i x_i}{\sum_{i=1}^{10} x_i^2} = \frac{12521}{4173} \approx 3.0004$$

$$1.3) \hat{y} = 25 - 0.5(6) = 22$$

خطای رگرسیون مقداری غیر صفر است یعنی بعضی نمونه‌ها را نمی‌تواند به صورت دقیق پیش‌بینی کند و ممکن است مقدار این نمونه ختم کم‌تر یا 22 متفاوت باشد

$$1.4) \sigma^2 = \frac{\sum_{i=1}^n (y_i - \bar{y})^2}{n - k} = \frac{7}{16 - 2} = 0.5 \quad (k \text{ تعداد پارامترهای رگرشن (شیب و عرض از مبدا)})$$

۲. الف)



همانطور که در شکل مشاهده می‌شود، هرچه وزن بیشتر باشد تابع به تغییر مقدار ورودی حساس‌تر می‌شود. اگر این تغییرات به دلیل نویز در داده‌های آموزشی باشد، ممکن است مدل نویز را یاد بگیرد که منجر به کمبود تعمیم‌دهی می‌شود.

همچنین وزن‌های بزرگ اهمیت بیشتری به ویژگی‌های مرتبط می‌دهند. یعنی مدل ممکن است بیش از حد به مجموعه‌ای از ویژگی‌ها اعتماد کند، حتی اگر برای تعمیم‌دهی ضروری نباشند. وزن‌های بزرگ به پیچیدگی بیشتری در مدل منجر می‌شوند. یک مدل بیش از حد پیچیده ممکن است به جای یادگیری الگوهای موجود، داده‌های آموزشی را حفظ کند که باعث **overfit** می‌شود.

$$2.2) \hat{w}_{MAP} = \operatorname{argmax}_w f(w|D) = \operatorname{argmax}_w \sum_i \ln \frac{f(x_i, w)}{p(y_i|x)} + \ln f(w)$$

$$\begin{aligned} \text{Term I: } \log P(y_1, y_2, \dots, y_n | x_1, x_2, \dots, x_n) &= \log \prod_{i=1}^n P(y_i | x_i) = \log \prod_{i=1}^n g(x_i) \\ &= \sum_{i=1}^n \log P(y_i | x_i) = \sum \log [g(x_i)^{y_i} (1-g(x_i))^{1-y_i}] = \\ &= \sum y_i \log g(x_i) + (1-y_i) \log (1-g(x_i)) = \\ &= \sum y_i \log \sigma(w^T x_i) + (1-y_i) \log (1-\sigma(w^T x_i)) \end{aligned}$$

$$\text{Term II: } \log \frac{1}{\sqrt{2\pi} \sigma} \exp\left(-\frac{1}{2} \left(\frac{w-\mu}{\sigma}\right)^2\right) = -\frac{1}{2} \log(2\pi\sigma^2) - \frac{1}{2} \left(\frac{w-\mu}{\sigma}\right)^2$$

$$w \sim N(0, I) \rightarrow -\frac{1}{2} \log(2\pi I) - \frac{1}{2} \left(\frac{w}{I}\right)^2$$

$$\frac{\partial \text{Term I}}{\partial w} = \sum_{i=1}^n y_i \frac{x_i \sigma(w^T x_i) (1-\sigma(w^T x_i))}{\sigma(w^T x_i)} + (1-y_i) \frac{-\sigma(w^T x_i) (1-\sigma(w^T x_i)) x_i}{1-\sigma(w^T x_i)} =$$

$$= \sum_{i=1}^n x_i y_i (1-\sigma(w^T x_i)) - (1-y_i) \sigma(w^T x_i) x_i = \left[y_i - \sigma(w^T x_i) \right] x_i$$

$$\frac{\partial \text{Term II}}{\partial w} = -\frac{w}{I} \cdot \frac{1}{I} = -\frac{w}{I}$$

$$\Rightarrow \nabla_w L(w) = \sum_{i=1}^n (y_i - \sigma(w^T x_i)) x_i - \left(\frac{w}{I}\right) \quad \text{w/o regularization term it's}$$

$$\rightarrow w_{t+1} = w_t - \eta \left[\sum_{i=1}^n (y_i - \sigma(w_t^T x_i)) x_i - \frac{w_t}{I} \right]$$

$$\begin{aligned}
 g(\lambda) &= \frac{1}{2} \sum_{i=1}^3 \sum_{j=1}^3 \lambda_i \lambda_j y_i y_j x_i^T x_j + \sum_{i=1}^3 \lambda_i \\
 &= \frac{1}{2} \left[\lambda_1 \lambda_1 y_1 y_1 x_1^T x_1 + \lambda_1 \lambda_2 y_1 y_2 x_1^T x_2 + \lambda_1 \lambda_3 y_1 y_3 x_1^T x_3 + \right. \\
 &\quad \lambda_2 \lambda_1 y_2 y_1 x_2^T x_1 + \lambda_2 \lambda_2 y_2 y_2 x_2^T x_2 + \lambda_2 \lambda_3 y_2 y_3 x_2^T x_3 + \\
 &\quad \left. \lambda_3 \lambda_1 y_3 y_1 x_3^T x_1 + \lambda_3 \lambda_2 y_3 y_2 x_3^T x_2 + \lambda_3 \lambda_3 y_3 y_3 x_3^T x_3 \right] \\
 &\quad + \lambda_1 + \lambda_2 + \lambda_3 = \\
 &= \frac{1}{2} \left[\lambda_1^2 \cdot 1 \cdot \overbrace{[1,2] \begin{bmatrix} 1 \\ 2 \end{bmatrix}}^5 + \lambda_1 \lambda_2 \cdot 1 \cdot \overbrace{[1,2] \begin{bmatrix} 1 \\ 2 \end{bmatrix}}^3 + \lambda_1 \lambda_3 (-1) \overbrace{[1,2] \begin{bmatrix} -1 \\ -2 \end{bmatrix}}^{-5} + \right. \\
 &\quad \lambda_2 \lambda_1 \cdot 1 \cdot \overbrace{[-1,2] \begin{bmatrix} 1 \\ 2 \end{bmatrix}}^3 + \lambda_2^2 \cdot 1 \cdot \overbrace{[-1,2] \begin{bmatrix} 1 \\ 2 \end{bmatrix}}^5 + \lambda_2 \lambda_3 (-1) \overbrace{[-1,2] \begin{bmatrix} -1 \\ -2 \end{bmatrix}}^{-3} + \\
 &\quad \left. \lambda_3 \lambda_1 (-1) \overbrace{[-1,-2] \begin{bmatrix} 1 \\ 2 \end{bmatrix}}^{-5} + \lambda_3 \lambda_2 (-1) \overbrace{[-1,-2] \begin{bmatrix} 1 \\ 2 \end{bmatrix}}^{-3} + \lambda_3^2 \cdot 1 \cdot \overbrace{[-1,-2] \begin{bmatrix} -1 \\ -2 \end{bmatrix}}^5 \right] \\
 &\quad + \lambda_1 + \lambda_2 + \lambda_3
 \end{aligned}$$

$$g(\lambda) = \frac{1}{2} \left[5\lambda_1^2 + 6\lambda_1\lambda_2 + 10\lambda_1\lambda_3 + 5\lambda_2^2 + 6\lambda_2\lambda_3 + 5\lambda_3^2 \right] + \lambda_1 + \lambda_2 + \lambda_3$$

$$\begin{aligned}
 \max_{\lambda_1, \lambda_2, \lambda_3} g(\lambda) \\
 \text{s.t. } \begin{cases} \lambda_i \geq 0 & i=1,2,3 \\ \sum_{i=1}^3 \lambda_i y_i = -\lambda_1 - \lambda_2 + \lambda_3 = 0 \end{cases}
 \end{aligned}$$

$$\frac{\partial g}{\partial \lambda_1} = \frac{1}{2} [10\lambda_1 + 6\lambda_2 + 10\lambda_3] + 1 = -5\lambda_1 - 3\lambda_2 - 5\lambda_3 + 1 = 0$$

$$\rightarrow 5\lambda_1 + 3\lambda_2 + 5\lambda_3 = 1$$

$$\frac{\partial g}{\partial \lambda_2} = \frac{1}{2} [6\lambda_1 + 10\lambda_2 + 6\lambda_3] + 1 = -3\lambda_1 - 5\lambda_2 - 3\lambda_3 + 1 = 0$$

$$\rightarrow 3\lambda_1 + 5\lambda_2 + 3\lambda_3 = 1$$

$$\frac{\partial g}{\partial \lambda_3} = \frac{1}{2} [10\lambda_1 + 6\lambda_2 + 10\lambda_3] + 1 = -5\lambda_1 - 3\lambda_2 - 5\lambda_3 + 1 = 0$$

$$\rightarrow 5\lambda_1 + 3\lambda_2 + 5\lambda_3 = 1$$

$$\Rightarrow \begin{bmatrix} 5 & 3 & 5 \\ 3 & 5 & 3 \\ 5 & 3 & 5 \end{bmatrix} \begin{bmatrix} \lambda_1 \\ \lambda_2 \\ \lambda_3 \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} \rightarrow \begin{aligned} \lambda_1 &= \frac{1}{8} - \lambda_3 \\ \lambda_2 &= \frac{1}{8} \end{aligned}$$

$$w = \sum_{i=1}^3 \lambda_i y_i x_i = (\frac{1}{8} - \lambda_3)(-1)(1, 2) + \frac{1}{8}(-1)(-1, 2) + \lambda_3 \cdot 1(-1, -2)$$

$$= (\lambda_3 - \frac{1}{8}, 2\lambda_3 - \frac{1}{4}) + (\frac{1}{8}, \frac{1}{4}) + (-\lambda_3, -2\lambda_3)$$

$$= (0, -\frac{1}{2})$$

$$b \rightarrow x_3 \text{ است از } y_3 = +1 \xrightarrow{\text{است از } y_3 = +1} w^T x_3 + b = 1 \rightarrow (0, -\frac{1}{2}) \begin{pmatrix} -1 \\ 2 \end{pmatrix} + b = 1 \rightarrow b = 0$$

$$\xRightarrow{\text{معادله خط}} w^T x + b = 0 \rightarrow \begin{pmatrix} 0 \\ -\frac{1}{2} \end{pmatrix} x + 0 = 0 \rightarrow -\frac{1}{2} x_2 = 0$$

۴. الف) Kernelها در SVM برای تبدیل فضای ویژگی‌های ورودی به یک فضای با ابعاد بالاتر مورد استفاده قرار می‌گیرند.

این تبدیلات به SVM امکان را می‌دهد که در فضای با ابعاد بالاتر، مسائل غیرخطی را به شکل خطی حل کند. مفهوم اصلی این است که اگر داده‌ها در فضای ویژگی به شکل خطی جداپذیر نباشند، می‌توان با افزودن بعد جدید به فضا، داده‌ها را به شکل خطی جداپذیر کرد. کرنل تابعی است که ورودی را به یک فضای با ابعاد بالاتر تبدیل می‌کند، بدون این که نیاز به محاسبه و ذخیره فضای جدید باشد و این مورد از مزایای کرنل است. کافیسیت حاصل ضرب نقطه‌ای دو ویژگی را در فضای دیگر بدانیم. حل مسئله به صورت خطی هم مزیت دیگر آن است.

ب)

4.2)

از آنجایی که K_1 و K_2 دو کرنل معتبر هستند:

$$K_1(x, x') = \sum_i \phi^1(x_i) \phi^1(x'_i), \quad K_2(x, x') = \sum_i \phi^2(x_i) \phi^2(x'_i)$$

$$\begin{aligned} K(x, x') &= \left(\sum_i \phi^1(x_i) \phi^1(x'_i) \right) \times \left(\sum_j \phi^2(x_j) \phi^2(x'_j) \right) = \\ &= \sum_i \sum_j \underbrace{\phi^1(x_i) \phi^2(x_j)}_{\phi_k(x)} \underbrace{\phi^1(x'_i) \phi^2(x'_j)}_{\phi_k(x')} \end{aligned}$$

$$\xrightarrow{\phi_k(x) = \phi^1(x_i) \phi^2(x_j)} \sum_k \phi(x_k) \phi(x'_k) = K(x, x')$$

۵. .

(۵,۱) ابتدا دیتاست را لود می‌کنیم. این داده شامل ۸ ستون و ۵۰۰ سطر است. مقادیر آماری ستون‌های عددی در جدول

مشاهده می‌شود:

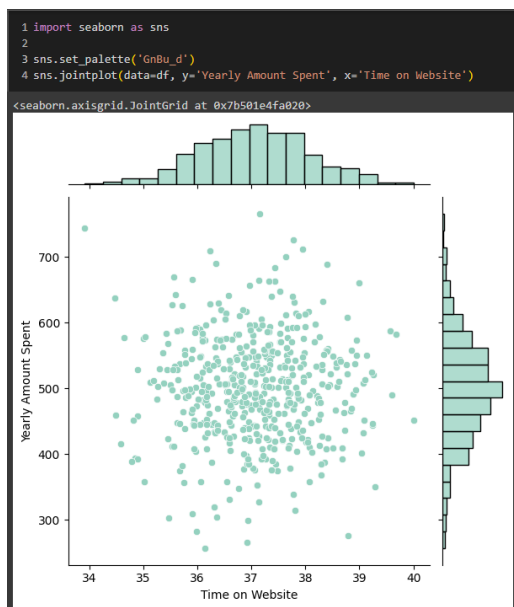
```
1 df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 500 entries, 0 to 499
Data columns (total 8 columns):
#   Column              Non-Null Count  Dtype
---  -
0   Email                500 non-null    object
1   Address              500 non-null    object
2   Avatar               500 non-null    object
3   Avg. Session Length  500 non-null    float64
4   Time on App          500 non-null    float64
5   Time on Website      500 non-null    float64
6   Length of Membership  500 non-null    float64
7   Yearly Amount Spent  500 non-null    float64
dtypes: float64(5), object(3)
memory usage: 31.4+ KB
```

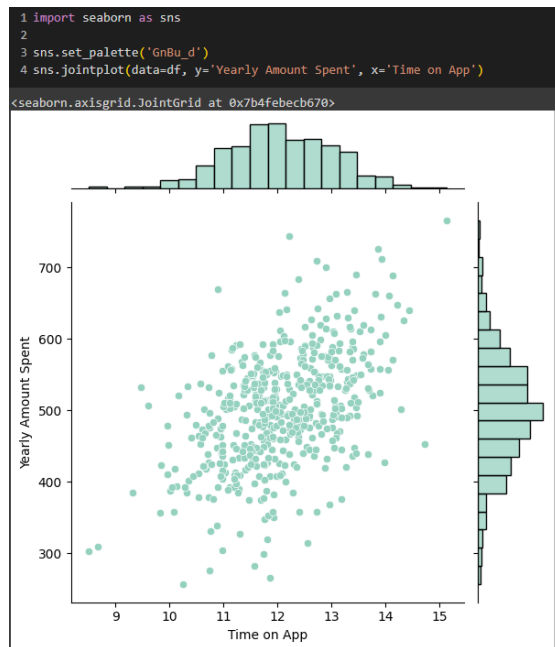
```
1 df.describe()
```

	Avg. Session Length	Time on App	Time on Website	Length of Membership	Yearly Amount Spent
count	500.000000	500.000000	500.000000	500.000000	500.000000
mean	33.053194	12.052488	37.060445	3.533462	499.314038
std	0.992563	0.994216	1.010489	0.999278	79.314782
min	29.532429	8.508152	33.913847	0.269901	256.670582
25%	32.341822	11.388153	36.349257	2.930450	445.038277
50%	33.082008	11.983231	37.069367	3.533975	498.887875
75%	33.711985	12.753850	37.716432	4.126502	549.313828
max	36.139662	15.126994	40.005182	6.922689	765.518462

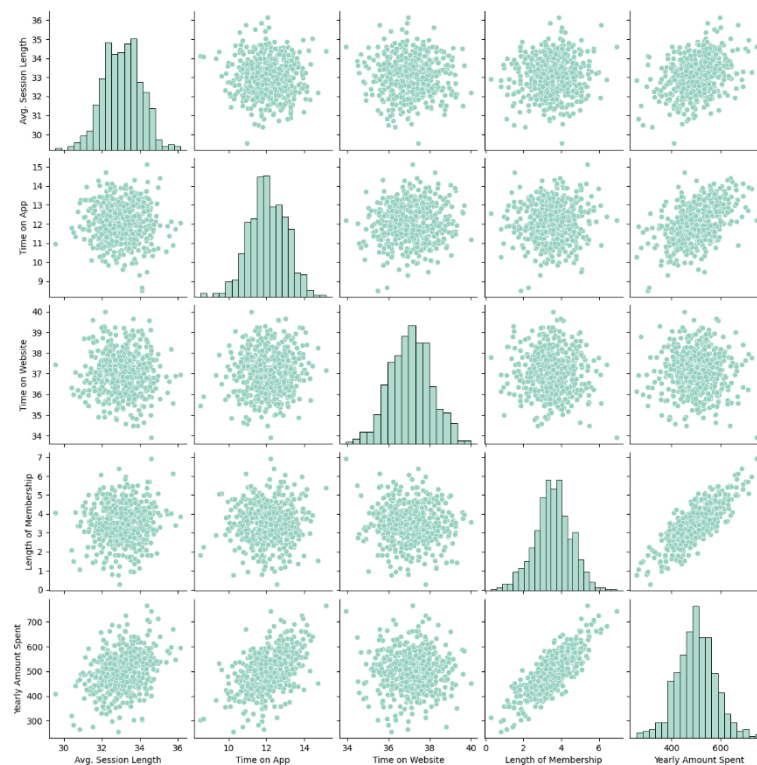
(۵,۲) پلات joint دو ویژگی خواسته شده:



(۵,۳)



(۵,۴) پلات رسم شده:



محور عمودی در نمودارهای ردیف آخر Yearly Amount Spent است. به نظر می‌رسد فیچر Length of Membership بیشترین تاثیر را روی میزان خرید دارد و با کم و زیاد شدن طول عضویت، مقدار خرج سالانه تغییر بیشتری داشته.

(۵,۵) ستون‌های عددی را به عنوان ویژگی در نظر می‌گیریم. البته ممکن است آدرس هم در میزان خرید سالیانه تاثیر داشته باشد ولی دسته‌بندی آن با توجه به تعداد کم داده دشوار است پس از آن صرف نظر می‌کنیم. سپس مقادیر ویژگی‌ها را نرمالیزه می‌کنیم زیرا مقادیر ممکن برای ویژگی‌های مختلف در بازه‌های مختلف قرار دارند و با نرمالیزه کردن همه مقادیر را به یک مقیاس می‌بریم.

```
1 print(df.columns)
2 # drop Email, Address and Avatar. Wouldn't help as feature
3 X = df.iloc[:, 3:-1]
4 y = df['Yearly Amount Spent']
5
6 X
```

Index(['Email', 'Address', 'Avatar', 'Avg. Session Length', 'Time on App', 'Time on Website', 'Length of Membership', 'Yearly Amount Spent'], dtype='object')

	Avg. Session Length	Time on App	Time on Website	Length of Membership
0	34.497268	12.655651	39.577668	4.082621
1	31.926272	11.109461	37.268959	2.664034
2	33.000915	11.330278	37.110597	4.104543
3	34.305557	13.717514	36.721283	3.120179
4	33.330673	12.795189	37.536653	4.446308
...
495	33.237660	13.566160	36.417985	3.746573
496	34.702529	11.695736	37.190268	3.576526
497	32.646777	11.499409	38.332576	4.958264
498	33.322501	12.391423	36.840086	2.336485
499	33.715981	12.418808	35.771016	2.735160

500 rows x 4 columns

```
1 from sklearn.preprocessing import StandardScaler
2
3 # normalize each column
4 scaler = StandardScaler()
5 X = scaler.fit_transform(X)
```

(۵,۶)

```
1 from sklearn.model_selection import train_test_split
2
3 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=101)
```

(۵,۷) مدل رگرسیون را آموزش می‌دهیم. دقت آن روی داده آموزشی حدود ۹۸ درصد است.

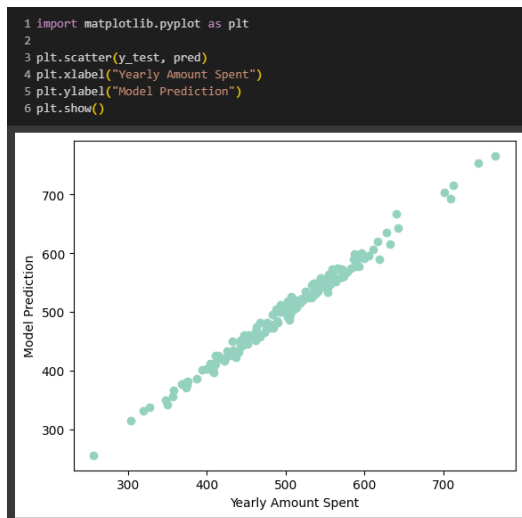
```
1 from sklearn.linear_model import LinearRegression
2
3 reg = LinearRegression().fit(X_train, y_train)
4 reg.score(X_train, y_train)
```

0.9817562058732432

(۵,۸)

```
1 pred = reg.predict(X_test)
```


(۵,۹



(۵,۱۰

```
1 import numpy as np
2
3 mae = np.mean(np.abs(y_test - pred))
4 mse = np.mean((y_test - pred)**2)
5 root_mse = np.sqrt(mse)
6
7 print('MAE:', mae, '\nMSE:', mse, '\nRoot MSE:', root_mse)
```

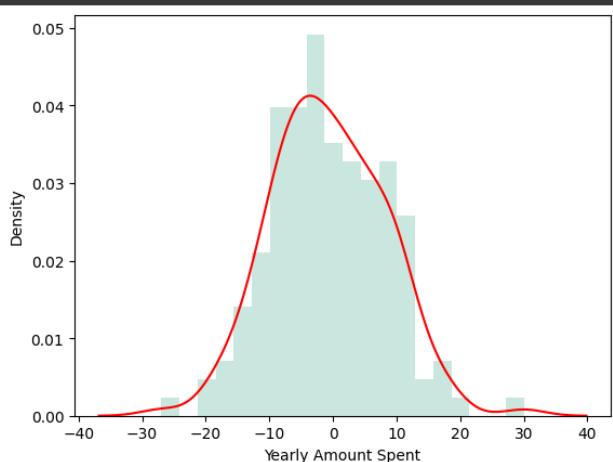
MAE: 7.228148653430832
MSE: 79.81305165097451
Root MSE: 8.933815066978637

(۵,۱۱) هیستوگرام و نمودار kde را در یک پلات رسم می‌کنیم. اگر مدل به خوبی فیت شده باشد، خطاهای پیش‌بینی باید به طور مساوی در اطراف مقدار واقعی توزیع شوند. یعنی خطاهای مثبت و منفی باید به طور مساوی توزیع شوند. علاوه بر این، شدت خطاها باید به طور مساوی توزیع شوند. میانگین خطا هم باید صفر شود که اینها توصیفی از توزیع نرمال است و در شکل هم نمودار تفاضل مقادیر به نرمال نزدیک است.

```

1 difference = y_test - pred
2
3 plt.hist(difference, bins=20, alpha=0.5, density=True)
4 sns.kdeplot(difference, color='red')
5 plt.show()

```



۵،۱۲ مدل برای هر یک از ویژگی‌ها وزنی در نظر گرفته است.

```

1 reg.coef_
array([25.76252659, 38.32855202,  0.19220992, 61.17355707])

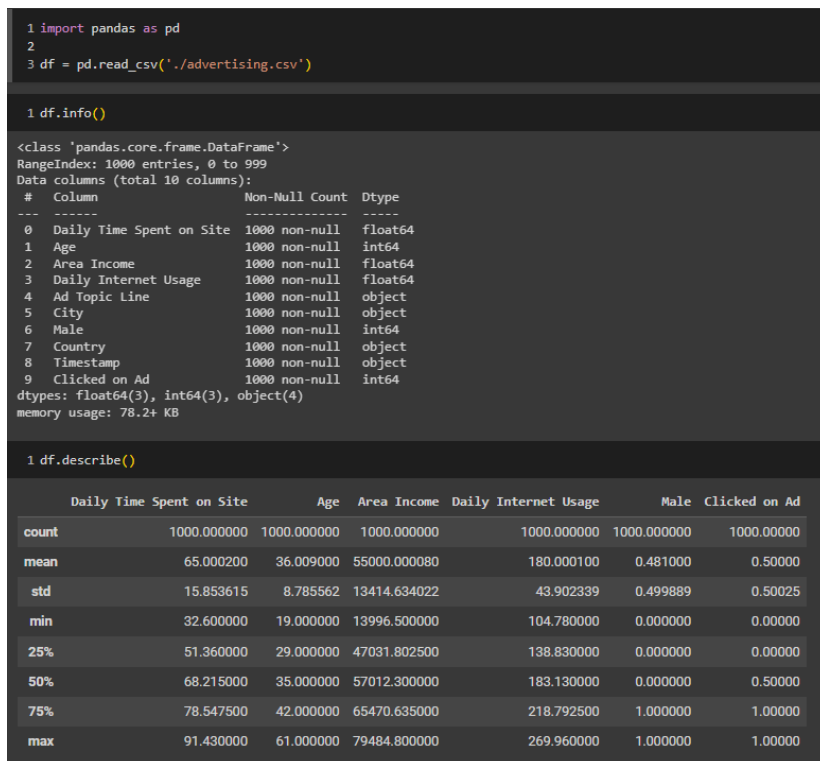
```

در رگرسیون خطی، **coef**، اهمیت هر فیچر را نشان می‌دهد. هر چه مقدار **coef** یک فیچر بیشتر باشد، آن فیچر در پیش‌بینی خروجی مدل تاثیر بیشتری دارد. اگر داده‌ها نرمالیزه نباشند، ممکن است ضرایب به طور غیرواقعی بزرگ یا کوچک شوند و منجر به پیش‌بینی اشتباه شود. این دلیل است که فیچرهایی که دامنه بزرگتری دارند، تاثیر بیشتری بر مدل خواهند داشت.

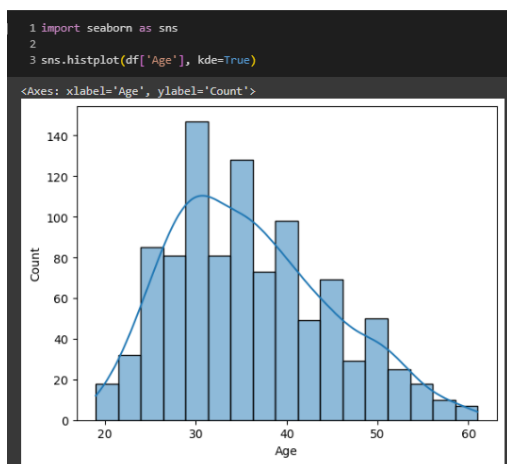
۵،۱۳ ضریب ویژگی آخر یعنی **Length of Membership** از همه بزرگتر است پس یعنی این ویژگی از اهمیت بیشتری برخوردار است. بهتر از سرمایه‌گذاری روی این طول عضویت مشتریان انجام شود.

۶.

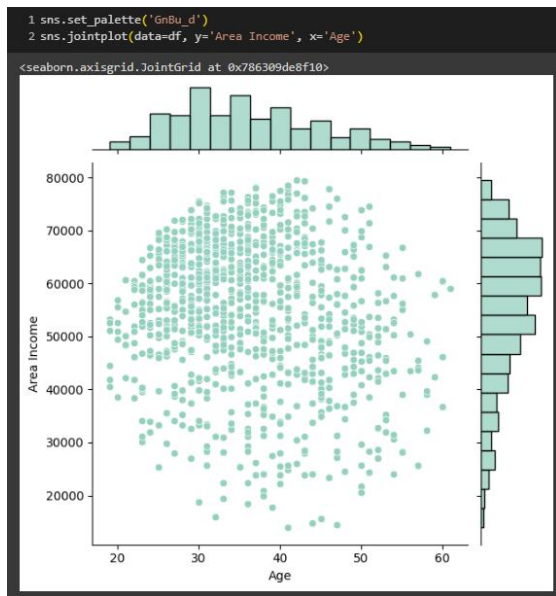
(۶,۱) این دیتاست از ۱۰ ستون و ۱۰۰۰ سطر تشکیل شده است.



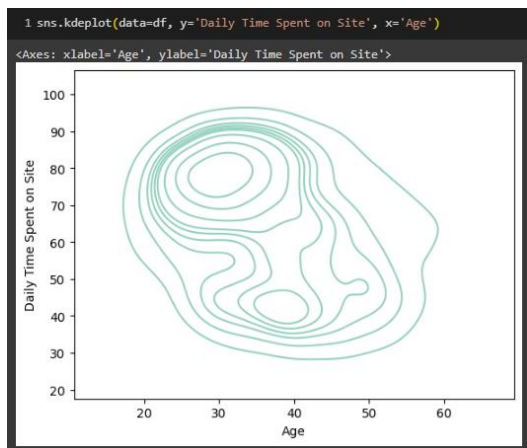
(۶,۲)



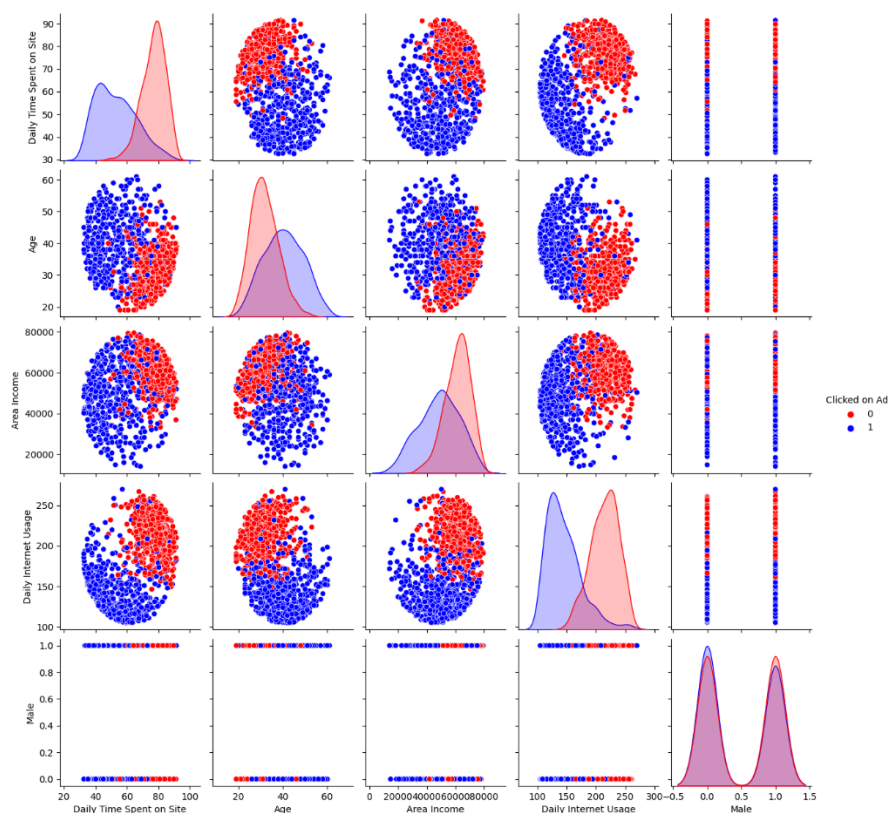
(۶,۳)



(۶,۴)



(۶,۵)



(۶,۶) افرادی که حوالی ۳۰ سال دارند بیشترین کلیک را روی تبلیغات کرده‌اند. میان سن و درآمد افراد تا حدی correlation وجود دارد. بیشتر جمعیت این دیتاست، بین ۳۰ تا ۴۰ سال سن دارند. همچنین از کانتورهای موجود در پلات kde می‌توان یافت که بیشتر افراد بین ۳۰ تا ۴۰ سال بیشترین زمان را در سایت صرف کرده‌اند و این میزان بیشتر بین ۷۰ تا ۸۰ دقیقه بوده‌است. به طور کلی با توجه به pair plot، با کمتر شدن زمان صرف شده در سایت، بالاتر رفتن سن، کمتر بودن درآمد، و مصرف اینترنت در طول روز، تعداد کلیک بر روی تبلیغ بیشتر شده‌است.

(۶,۷) ۵ ستون: میزان زمان صرف شده در سایت به طور روزانه، سن، درآمد، میزان اینترنت مصرفی به طور روزانه و مرد بودن را به عنوان ویژگی داده برای مسئله طبقه‌بندی در نظر می‌گیریم. سپس برای اینکه مقادیر عددی همگی در یک مقیاس قرار گیرند، داده‌ها را نرمالیزه می‌کنیم. ۳۰ درصد از داده را به عنوان تست جدا می‌کنیم.

```

1 from sklearn.preprocessing import StandardScaler
2 from sklearn.model_selection import train_test_split
3
4 X = df[['Daily Time Spent on Site', 'Age', 'Area Income', 'Daily Internet Usage', 'Male']]
5 y = df['Clicked on Ad']
6
7 scaler = StandardScaler()
8 X = scaler.fit_transform(X)
9
10 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=42)
11 print(X_train.shape)

```

(700, 5)

۶،۸) یک کلاس logistic regression تعریف می‌کنیم. برای حذف بایاس و اضافه کردن آن به بردار وزن، در ابتدا یک بعد به ورودی اضافه می‌کنیم. سپس بردار وزن را با ابعاد $(d+1) \times 1$ با 0 مقداردهی اولیه می‌کنیم. توابع مورد نیاز را تعریف می‌کنیم:

```

3 class LogisticReg:
4     def __init__(self, X, y, epochs=1000, alpha=0.01):
5         # add one more dim to input (bias)
6         self.X = np.hstack((np.ones((X.shape[0], 1)), X))
7         self.y = y
8
9         # initialize weights
10        self.W = np.zeros((self.X.shape[1], 1))
11
12        self.epochs = epochs
13        self.alpha = alpha
14
15        def sigmoid(self, z):
16            return 1 / (1 + np.exp(-z))
17
18        def forward(self, X):
19            z = np.dot(X, self.W)
20            return self.sigmoid(z)
21
22        def loss(self, y_hat, y):
23            l = -1 * np.average(y * np.log(y_hat) + (1-y) * np.log(1-y_hat))
24            return l
25
26        def grads(self, X, y_hat, y):
27            # round L/round w = x (y_hat - y)
28            return np.dot(X.T, np.expand_dims(y_hat-y, axis=1)) / len(y)
29

```

در تابع train به تعداد ایپاک باید مدل را آموزش دهیم. برای این کار ابتدا یک مرحله forward propagation انجام می‌دهیم و بعد میزان خطای cross entropy را حساب می‌کنیم. سپس گرادیان را محاسبه و بردار وزن را آپدیت می‌کنیم.

از تابع predict برای پیش‌بینی داده تست استفاده می‌کنیم. یک مرحله forward را انجام می‌دهیم و به داده با توجه به خروجی و ترشهولد (در اینجا ۰،۵) برچسب می‌زنیم.

```

30 def train(self):
31     for e in range(self.epochs):
32         # forward propagation
33         y_hat = np.squeeze(self.forward(self.X))
34         # compute loss
35         l = self.loss(y_hat, self.y)
36         # compute gradients
37         dw = self.grads(self.X, y_hat, self.y)
38         # update parameters
39         self.W -= self.alpha * dw
40
41     if e % 100 == 0:
42         print(f'Epoch {e}, Loss:', "{:.3f}".format(l))
43
44 def predict(self, X):
45     # add one more dimension
46     X = np.hstack((np.ones((X.shape[0], 1)), X))
47     probs = self.forward(X)
48     preds = (probs >= 0.5).astype(int)
49     return preds
50
51

```

مدل را در ۱۰۰۰ اپیاک و با learning rate ۰,۰۱ آموزش می‌دهیم. دقت مدل روی داده آموزش حدود ۹۷٪ است.

```

1 model = LogisticReg(X_train, y_train)
2 model.train()

Epoch 0, Loss: 0.693
Epoch 100, Loss: 0.424
Epoch 200, Loss: 0.311
Epoch 300, Loss: 0.253
Epoch 400, Loss: 0.219
Epoch 500, Loss: 0.196
Epoch 600, Loss: 0.180
Epoch 700, Loss: 0.168
Epoch 800, Loss: 0.159
Epoch 900, Loss: 0.151

1 from sklearn.metrics import accuracy_score, classification_report,
2
3 pred_train = model.predict(X_train)
4 acc_train = accuracy_score(y_train, pred_train)
5 print("Train Accuracy", "{:.3f}".format(acc_train))

Train Accuracy 0.969

```

(۶,۹)

```

1 pred_test = model.predict(X_test)
2 pred_test.shape

(300, 1)

```

(۶,۱۰) میزان دقت مدل روی داده آزمون ۹۵٪ است. مطابق confusion matrix، ۲ داده برچسب واقعی ۰ داشتند و به اشتباه ۱ پیش‌بینی شدند. در مقابل ۱۳ داده برچسب واقعی ۱ داشتند و به اشتباه ۰ پیش‌بینی شدند.

```

1 pred_test = model.predict(X_test)
2 acc_test = accuracy_score(y_test, pred_test)
3 print("Test Accuracy", "{:.3f}".format(acc_test))

Test Accuracy 0.950

1 print(confusion_matrix(y_test, pred_test))

[[144  2]
 [ 13 141]]

1 print(classification_report(y_test, pred_test))

              precision    recall  f1-score   support

     0       0.92       0.99       0.95       146
     1       0.99       0.92       0.95       154

 accuracy          0.95
 macro avg         0.95
 weighted avg      0.95

```

۷. ابتدا داده را لود می‌کنیم و توزیع آن را بررسی می‌کنیم. ۸۶٪ از داده ham است.

```
1 import pandas as pd
2
3 df = pd.read_csv('./spamSMS.csv', encoding='ISO-8859-1')

1 data = df['v2']
2 labels = df['v1']
3
4 print("data count:", len(data))

data count: 5572

1 classes = pd.unique(labels)
2
3 c0 = len(labels[labels == classes[0]])
4 c1 = len(labels[labels == classes[1]])
5 print(f'{classes[0]} count: {c0} ({c0/ len(labels)}%)')
6 print(f'{classes[1]} count: {c1} ({c1/ len(labels)}%)')

ham count: 4825 (0.8659368269921034%)
spam count: 747 (0.13406317300789664%)
```

از count vectorizer برای استخراج ویژگی استفاده می‌کنیم. این تابع تعداد تکرار هر کلمه موجود در vocabulary را می‌شمارد و در یک دیکشنری ذخیره می‌کند.

داده را به دو بخش آموزش و آزمون تقسیم می‌کنیم.

```
1 from sklearn.model_selection import train_test_split
2
3 X_train, X_test, y_train, y_test = train_test_split(features, labels, test_size=0.3, random_state=42)

1 print("train data shape:", X_train.shape)
2 print("test data shape:", X_test.shape)

train data shape: (3900, 8672)
test data shape: (1672, 8672)
```

Grid Search: در این روش فضای ابرپارامترهای ممکن مشخص شده و تمام ترکیب‌های ممکن آنها بررسی می‌شوند. با استفاده از grid search کل فضای حالت امتحان می‌شود و می‌توانیم حالت بهینه را پیدا کنیم (در صورتی که در محدود جستجو باشد). اما بدلیل بررسی کل فضای جستجو، این روش زمان‌بر است و برای فضای حالت بزرگ (حالات زیاد برای ابرپارامترهای زیاد) مناسب نیست.

```
1 from sklearn.svm import SVC
2 from sklearn.model_selection import GridSearchCV, RandomizedSearchCV
3
4 parameters = {'C': [0.1, 1, 10], 'kernel': ['linear', 'rbf'], 'gamma': [0.1, 1, 10]}
5 svc = SVC()
6 grid_search_svm = GridSearchCV(svc, parameters)
7 grid_search_svm.fit(X_train, y_train)
8
9 print(sorted(grid_search_svm.cv_results_.keys()))
```

پس از آموزش مدل با ابرپارامترهای گوناگون به روش grid به نتایج زیر رسیدیم:



در این روش svm با کرنل خطی و گاما ۰٫۱ و $C=1$ بهترین نتیجه را بدست آورد و به دقت ۹۷٪ روی داده آزمون رسید (نتایج سایر مدل‌ها نیز در نوت‌بوک موجود است).

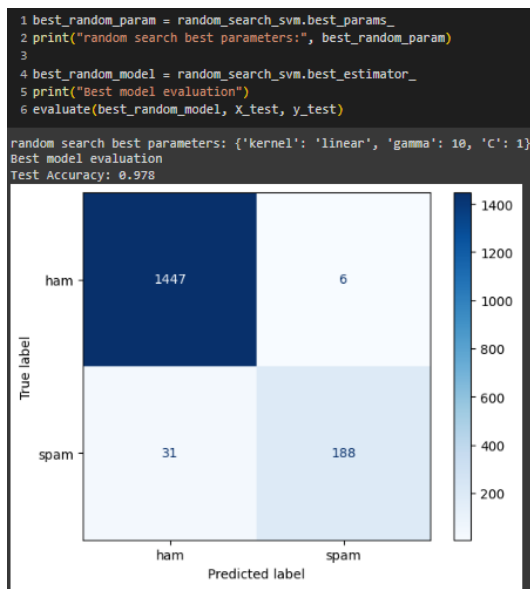
Random Search: در این روش تعدادی نمونه تصادفی از فضای ابرپارامترها انتخاب می‌شوند و مدل‌ها با آنها آموزش می‌بینند. این روش می‌تواند در فضای ابرپارامترهای پیچیده و تعداد زیاد ابعاد مفید باشد اما بدلیل اینکه تمامی فضای جستجو مورد بررسی قرار نمی‌گیرد، ممکن است به حالت بهینه نرسیم. در این سوال ابرپارامترها را برای ۱۰ حالت تصادفی انتخاب می‌کنیم و مدل آموزش می‌دهیم.

```

1 svm_model = SVC()
2
3 parameters = {'C': [0.1, 1, 10], 'kernel': ['linear', 'rbf'], 'gamma': [0.1, 1, 10]}
4 random_search_svm = RandomizedSearchCV(svm_model, param_distributions=parameters, n_iter=10)
5 random_search_svm.fit(X_train, y_train)

```

پس از آموزش مدل با ابرپارامترهای گوناگون به روش random به نتایج زیر رسیدیم:



در این روش svm با کرنل خطی و گاما ۱۰ و C=۱ بهترین نتیجه را بدست آورد و به دقت ۹۷٪ روی داده آزمون رسید (نتایج سایر مدل ها نیز در نوت بوک موجود است).

مقایسه دو روش: grid از نظر زمانی بسیار بیشتر از random طول می کشد. از لحاظ کارایی grid می تواند به حالت بهینه برسد در صورتی که ممکن است random به بهینه نرسد.

۸. ابتدا داده را لود می کنیم:

```
1 import pandas as pd
2
3 df = pd.read_csv('./housePricing.csv')
4 df
```

	Id	MSSubClass	MSZoning	LotFrontage	LotArea	Street	Alley	LotShape	LandContour	Utilities	...	PoolArea	PoolQC	Fence	MiscFeature	MiscVal	MoSold	YrSold	SaleType	SaleCondition	SalePrice
0	1	60	RL	65.0	8450	Pave	NaN	Reg	Lvl	AllPub	...	0	NaN	NaN	NaN	0	2	2008	WD	Normal	208500
1	2	20	RL	80.0	9600	Pave	NaN	Reg	Lvl	AllPub	...	0	NaN	NaN	NaN	0	5	2007	WD	Normal	181500
2	3	60	RL	68.0	11250	Pave	NaN	IR1	Lvl	AllPub	...	0	NaN	NaN	NaN	0	9	2008	WD	Normal	223500
3	4	70	RL	60.0	9550	Pave	NaN	IR1	Lvl	AllPub	...	0	NaN	NaN	NaN	0	2	2006	WD	Abnorml	140000
4	5	60	RL	84.0	14260	Pave	NaN	IR1	Lvl	AllPub	...	0	NaN	NaN	NaN	0	12	2008	WD	Normal	250000
...
1455	1456	60	RL	62.0	7917	Pave	NaN	Reg	Lvl	AllPub	...	0	NaN	NaN	NaN	0	8	2007	WD	Normal	175000
1456	1457	20	RL	85.0	13175	Pave	NaN	Reg	Lvl	AllPub	...	0	NaN	MnPrv	NaN	0	2	2010	WD	Normal	210000
1457	1458	70	RL	66.0	9042	Pave	NaN	Reg	Lvl	AllPub	...	0	NaN	GdPrv	Shed	2500	5	2010	WD	Normal	266500
1458	1459	20	RL	68.0	9717	Pave	NaN	Reg	Lvl	AllPub	...	0	NaN	NaN	NaN	0	4	2010	WD	Normal	142125
1459	1460	20	RL	75.0	9937	Pave	NaN	Reg	Lvl	AllPub	...	0	NaN	NaN	NaN	0	6	2008	WD	Normal	147500

1460 rows x 81 columns

سپس بر روی آن پیش پردازش های زیر را انجام می دهیم: ستون هایی که بیشتر از ۴۰٪ داده شان مقدار ندارند و nan هستند را حذف می کنیم. برای سایر ستون ها، اگر متنی باشند مقدار nan را به unknown و اگر عددی باشند با مقدار میانگین آن ستون جا به جا می کنیم.

```
1 # count the number of missing values of each column
2 nan = list(df.isna().sum())
3
4 # find the columns' names which need to be replaced or dropped
5 replaced_col_text = []
6 replaced_col_number = []
7 dropped_col = []
8
9 for idx, name in enumerate(df.columns):
10     if nan[idx] == 0:
11         continue
12     if nan[idx] > len(df) * 0.4:
13         dropped_col.append(name)
14     # find text columns which have nan value inside
15     elif df[name].dtype == object:
16         replaced_col_text.append(name)
17     # find digit columns which have nan value
18     else: replaced_col_number.append(name)
19
20 print("text columns with nan:", replaced_col_text)
21 print("number columns with nan:", replaced_col_number)
22 print("dropping columns:", dropped_col)
23
24 # replace missing value of a text column with 'unknown'
25 df[replaced_col_text] = df[replaced_col_text].fillna(value='unknown')
26 # replace missing value of a digit column with its column average
27 df[replaced_col_number] = df[replaced_col_number].fillna(value=df[replaced_col_number].mean())
28 # drop columns with more than 40% missing value
29 df_dropped = df.drop(columns=dropped_col, axis=1)
30
31 # check if any column has missing value
32 for col in df_dropped.columns:
33     if df_dropped[col].isna().any():
34         print(col, "has na value")
35
36 text columns with nan: ['MasVnrType', 'BsmtQual', 'BsmtCond', 'BsmtExposure', 'BsmtFinType1', 'BsmtFinType2', 'Electrical', 'GarageType', 'GarageFinish', 'GarageQual', 'GarageCond']
37 number columns with nan: ['LotFrontage', 'MasVnrArea', 'GarageYrBlt']
38 dropping columns: ['Alley', 'FireplaceQu', 'PoolQC', 'Fence', 'MiscFeature']
```

در خروجی نام ستون هایی که مورد تغییر یا حذف قرار گرفته اند مشاهده می شود.

سپس ستون‌هایی که مقدار غیر عددی دارند، باید مقادیر آنها را با یک مقدار عددی تعویض کنیم. برای این کار از one hot encoder کتابخانه sklearn استفاده می‌کنیم. برای این کار ستون مربوطه حذف می‌شود و چند ستون (تعداد بستگی به مقادیر منحصر به فرد آن ستون دارد) اضافه می‌شود و با روش one hot مقادیر را ذخیره می‌کنیم.

```
1 from sklearn.preprocessing import OneHotEncoder
2
3 text_columns = df_dropped.select_dtypes(include=['object']).columns
4 # apply one-hot encoding on text columns
5 one_hot = pd.get_dummies(df_dropped[text_columns])
6 # Drop the original text columns and add numeric ones
7 new_df = pd.concat([df_dropped.drop(text_columns, axis=1), one_hot], axis=1)
8 new_df
```

	Id	MSSubClass	LotFrontage	LotArea	OverallQual	OverallCond	YearBuilt	YearRemodAdd	MasVnrArea	BsmtFinSF1	...	SaleType_ConLw	SaleType_New	SaleType_Oth	SaleType_WD	SaleCondition_Abnormal	S
0	1	60	65.0	8450	7	5	2003	2003	196.0	706	...	0	0	0	1	0	
1	2	20	80.0	9600	6	8	1976	1976	0.0	978	...	0	0	0	1	0	
2	3	60	68.0	11250	7	5	2001	2002	162.0	486	...	0	0	0	1	0	
3	4	70	60.0	9550	7	5	1915	1970	0.0	216	...	0	0	0	1	1	
4	5	60	84.0	14260	8	5	2000	2000	350.0	655	...	0	0	0	1	0	
...
1455	1456	60	62.0	7917	6	5	1999	2000	0.0	0	...	0	0	0	1	0	
1456	1457	20	85.0	13175	6	6	1978	1988	119.0	790	...	0	0	0	1	0	
1457	1458	70	66.0	9042	7	9	1941	2006	0.0	275	...	0	0	0	1	0	
1458	1459	20	68.0	9717	5	6	1950	1996	0.0	49	...	0	0	0	1	0	
1459	1460	20	75.0	9937	5	6	1965	1965	0.0	830	...	0	0	0	1	0	

1460 rows x 283 columns

بدین ترتیب بعد از این کار، تعداد ستون‌ها به ۲۸۳ افزایش می‌یابد.

بعد ویژگی‌ها و همچنان قیمت خانه را نرمالیزه می‌کنیم.

```
1 from sklearn.preprocessing import StandardScaler
2
3 X = new_df.drop('SalePrice', axis=1) # features
4 y = new_df['SalePrice'] # labels
5
6 numeric_columns = X.select_dtypes(include=['float64', 'int64']).columns
7 scaler = StandardScaler()
8 X[numeric_columns] = scaler.fit_transform(X[numeric_columns])
9 scaler = StandardScaler()
10 y = scaler.fit_transform(y.values.reshape(-1, 1)).ravel()
```

۲۰ ویژگی برتر را انتخاب می‌کنیم و داده را به سه بخش آموزش، اعتبارسنجی و آزمون تقسیم می‌کنیم.

```
1 from sklearn.feature_selection import SelectKBest, f_regression
2
3 k_best = 20
4 selector = SelectKBest(score_func=f_regression, k=k_best)
5 X_selected = selector.fit_transform(X, y)
```

```
1 from sklearn.model_selection import train_test_split
2
3 X_train, X_temp, y_train, y_temp = train_test_split(X_selected, y, test_size=0.3, random_state=4)
4 X_val, X_test, y_val, y_test = train_test_split(X_temp, y_temp, test_size=0.5, random_state=4)
5
6 print("train features shape:", X_train.shape)
7 print("train labels shape:", y_train.shape)
```

train features shape: (1022, 20)
train labels shape: (1022,)

پس از تمیز کردن داده، مدل رگرسیون را آموزش می‌دهیم و نتیجه را بر روی داده‌ها گزارش می‌کنیم.

```
1 from sklearn.svm import SVR
2 from sklearn.metrics import mean_squared_error
3
4 # SVM Regression model
5 model = SVR()
6 model.fit(X_train, y_train)
7
8 # Make predictions on training, validation, and test sets
9 train_pred = model.predict(X_train)
10 val_pred = model.predict(X_val)
11 test_pred = model.predict(X_test)
12
13 # calculate loss on each split of data
14 loss_train = mean_squared_error(y_train, train_pred)
15 loss_val = mean_squared_error(y_val, val_pred)
16 loss_test = mean_squared_error(y_test, test_pred)
17
18 print('Loss on train data:', "{:.3f}".format(loss_train))
19 print('Loss on val data:', "{:.3f}".format(loss_val))
20 print('Loss on test data:', "{:.3f}".format(loss_test))
```

```
Loss on train data: 0.129
Loss on val data: 0.090
Loss on test data: 0.136
```

مطابق انتظار، میزان خطا روی داده آزمون از همه بیشتر است اما تفاوت چندانی با خطای آموزش ندارد و مدل دچار overfit نشده است.