

Лабораторна робота №1. Лінійна регресія

Виконав студент групи КМ-91мп

Галета М.С.

Завдання на лабораторну роботу

1. Обрати відповідний файл з даними.
2. Збудувати модель множинної лінійної регресії, використовуючи метод найменших квадратів, обравши в якості залежної змінної останній стовпчик, а всі інші – в якості незалежних змінних. Для побудови моделі використовувати перші 200 записів у файлі з даними.
3. Оцінити статистичну значимість коефіцієнтів отриманої моделі.
4. Оцінити адекватність збудованої моделі за допомогою коефіцієнту множинної детермінації.
5. Збудувати модель множинної лінійної регресії (з тими самими аргументами, що у п. 2), використовуючи метод градієнтного спуску.
6. Вивести обчислені похибки збудованих моделей регресії для 50 наступних записів у файлі з даними.

In [1]:

```
1 import numpy as np
2 import pandas as pd
3 import matplotlib.pyplot as plt
4 from IPython.display import display
5
6 %matplotlib inline
```

Зчитування датасету

In [2]:

```
1 df = pd.read_csv('MP-04-Galeta.csv', delimiter=';', names=['x1', 'x2', 'x3', 'x4', 'x5']
2
3 X_train = df.iloc[:200, :-1].values
4 X_test = df.iloc[200:250, :-1].values
5 y_train = df['y'].iloc[:200].values.reshape((200, 1))
6 y_test = df['y'].iloc[200:250].values.reshape((50, 1))
```

Центрування і стандартизація даних

$$X_{new} = \frac{X - \mu}{\sigma}$$

In [3]:

```
1 class Scaler:
2     def __init__(self):
3         self.mean = None
4         self.std = None
5
6     def fit(self, X):
7         self.mean = np.mean(X, axis=0, keepdims=True)
8         self.std = np.std(X, axis=0, keepdims=True)
9
10    def transform(self, X):
11        X_new = (X - self.mean)/self.std
12        return X_new
13
14    def fit_transform(self, X):
15        self.mean = np.mean(X, axis=0, keepdims=True)
16        self.std = np.std(X, axis=0, keepdims=True)
17        X_new = (X - self.mean)/self.std
18        return X_new
```

In [4]:

```
1 sc = Scaler()
2 X_train = sc.fit_transform(X_train)
3 X_test = sc.transform(X_test)
```

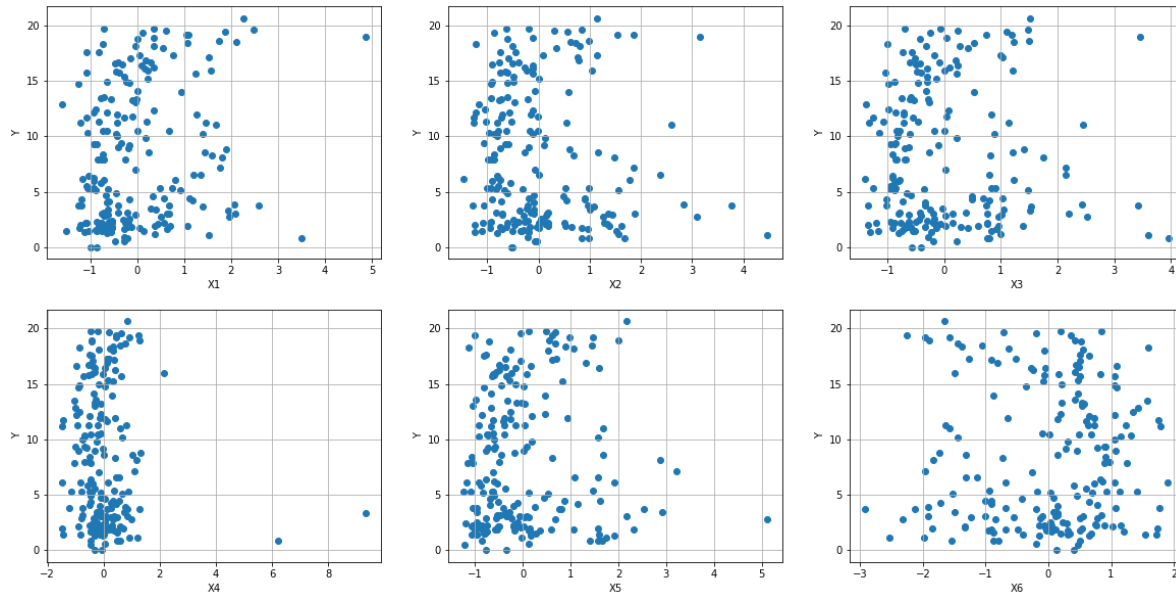
Графіки залежності окремо кожної з ознак від Y

In [5]:

```

1 fig, ax = plt.subplots(2, 3, figsize=(20,10))
2
3 xlabel = 1
4 for i in range(2):
5     for j in range(3):
6         ax[i][j].set_xlabel('X'+str(xlabel))
7         ax[i][j].set_ylabel('Y')
8         ax[i][j].grid(True)
9         ax[i][j].scatter(X_train[:, xlabel-1], y_train)
10        xlabel += 1
11
12 plt.show()

```



Додавання фіктивного стовпчика з одиниць для реалізації таким чином вільного коефіцієнта регресії

In [6]:

```

1 X_train = np.concatenate((np.ones((200,1)), X_train), axis=1)
2 X_test = np.concatenate((np.ones((50,1)), X_test), axis=1)

```

Модель лінійної регресії методом найменших квадратів

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

In [7]:

```

1 coefs = (np.linalg.pinv(X_train.T @ X_train) @ X_train.T) @ y_train
2 df_coefs = pd.DataFrame(columns=['b', 'w1', 'w2', 'w3', 'w4', 'w5', 'w6'], data=coefs.)
3 print("Коефіцієнти регресії")
4 display(df_coefs)

```

Коефіцієнти регресії

	b	w1	w2	w3	w4	w5	w6
0	8.189645	8.412788	-0.157077	-6.709854	-1.172603	0.535967	0.489362

Перевірка статистичної значимості коефіцієнтів за Стьюдентом

$$S = \sqrt{\frac{\sum_{i=0}^n (y_i - y_{pred_i})^2}{n - m - 1}}$$

$$S_\beta = S \sqrt{\text{diag}((X^T X)^{-1})}$$

$$t_\beta = \frac{coef_i}{S_\beta}$$

Якщо $|t_\beta| > t_{\text{табл.}}$ отже відповідний коефіцієнт є значимим

In [8]:

```

1 def statistical_significance(X, y, coefs, n, m):
2     y_pred = X @ coefs
3
4     s = np.sqrt(np.sum((y-y_pred)**2)/(n-m-1))
5     s_beta = s*np.sqrt(np.diag(np.linalg.pinv(X.T @ X)))
6
7     t_beta = np.hstack(coefs)/s_beta
8
9     return t_beta, y_pred

```

In [9]:

```

1 t_table = 1.6602
2 t_beta, y_pred = statistical_significance(X_train, y_train, coefs, 200, 6)

```

In [10]:

```
1 sign = pd.DataFrame(columns=['b', 'w1', 'w2', 'w3', 'w4', 'w5', 'w6'], data=np.array([
2 print("Статистична значимість (1 якщо значимий, 0 якщо незначимий)")
3 display(sign)
```

Статистична значимість (1 якщо значимий, 0 якщо незначимий)

	b	w1	w2	w3	w4	w5	w6
0	1	1	0	1	1	0	0

Перевірка адекватності моделі за допомогою коефіцієнта множинної детермінації

Коефіцієнт детермінації визначається наступним чином:

$$R^2 = 1 - \frac{V(y|x)}{V(y)} = 1 - \frac{\sigma^2}{\sigma_y^2},$$

де $V(y|x) = \sigma^2$ — умовна дисперсія залежної змінної.

Для розрахунку вибіркового коефіцієнта детермінації, використовують вибіркові оцінки значень відповідних дисперсій:

$$R^2 = 1 - \frac{\hat{\sigma}^2}{\hat{\sigma}_y^2} = 1 - \frac{ESS/n}{TSS/n} = 1 - \frac{ESS}{TSS},$$

де $ESS = \sum_{t=1}^n e_t^2 = \sum_{t=1}^n (y_t - \hat{y}_t)^2$ — сума квадратів залишків регресії, y_t, \hat{y}_t — фактичні та оціночні значення пояснювальної змінної.

$TSS = \sum_{t=1}^n (y_t - \bar{y})^2 = n\hat{\sigma}_y^2$ — загальна сума квадратів.

$$\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$$

У випадку класичної лінійної множинної регресії (регресії з константою):

$$TSS = RSS + ESS, \text{ де } RSS = \sum_{t=1}^n (\hat{y}_t - \bar{y})^2$$

І як наслідок:

$$R^2 = \frac{RSS}{TSS}$$

Коефіцієнт детермінації — це відношення поясненої дисперсії до загальної

In [11]:

```
1 rss = np.mean((y_pred - np.mean(y_train))**2)
2 ess = np.mean((y_train - y_pred)**2)
3 tss = np.var(y_train)
4 R2 = rss/tss
5 print("R2 = {}".format(R2))
6 print('Звідси випливає, що лише 33% дисперсії буде пояснено при використанні даної моделі')
```

R2 = 0.33120750473071003

Звідси випливає, що лише 33% дисперсії буде пояснено при використанні даної моделі лінійної регресії.

Модель лінійної регресії методом градієнтного спуску

In [12]:

```
1 class LinearRegression:
2     def __init__(self):
3         self.w = None
4         self.b = None
5         self.cost = None
6
7     def initialize_with_zeros(self, dim):
8         w = np.zeros((dim, 1))
9         b = 0
10        return w, b
11
12    def compute_cost(self, H, y):
13        cost = np.mean((H - y) ** 2)/2
14        return cost
15
16    def forward_propagation(self, w, b, X):
17        H = X @ w + b
18        return H
19
20    def backward_propagation(self, X, y, H):
21        dw = X.T @ (H - y) / X.shape[0]
22        db = np.sum(H - y) / X.shape[0]
23
24        grads = {'dw': dw, 'db': db}
25        return grads
26
27    def update_parametres(self, w, b, grads, learning_rate):
28        dw = grads['dw']
29        db = grads['db']
30
31        w = w - learning_rate*dw
32        b = b - learning_rate*db
33        return w, b
34
35    def fit(self, X, y, verbose=False, epochs=1000, l_rate=0.1):
36        w, b = self.initialize_with_zeros(X.shape[1])
37        for i in range(epochs):
38            H = self.forward_propagation(w, b, X)
39            self.cost = self.compute_cost(H, y)
40            if verbose:
41                print("Loss: "+str(self.cost))
42            grads = self.backward_propagation(X, y, H)
43            w, b = self.update_parametres(w, b, grads, l_rate)
44
45        self.w = w
46        self.b = b
47
48    def predict(self, X):
49        y_pred = self.forward_propagation(self.w, self.b, X)
50        return y_pred
```

In [13]:

```
1 lr = LinearRegression()
2 lr.fit(np.delete(X_train, 0, axis=1), y_train, False, 10000)
```

In [14]:

```
1 print("Loss: {}".format(lr.cost))
```

Loss: 12.459655405997378

Порівняння коефіцієнтів, отриманих двома методами

In [15]:

```
1 coefs_2 = np.concatenate((np.array([lr.b]), np.hstack(lr.w)))
2 df_coefs_2 = pd.DataFrame(columns=['b', 'w1', 'w2', 'w3', 'w4', 'w5', 'w6'], data=coefs_2)
3 print("Метод найменших квадратів")
4 display(df_coefs)
5 print("Гرادієнтний спуск")
6 display(df_coefs_2)
```

Метод найменших квадратів

	b	w1	w2	w3	w4	w5	w6
0	8.189645	8.412788	-0.157077	-6.709854	-1.172603	0.535967	0.489362

Градiєнтний спуск

	b	w1	w2	w3	w4	w5	w6
0	8.189645	8.412788	-0.157077	-6.709854	-1.172603	0.535967	0.489362

Похибки моделей на тестовій вибірці

In [16]:

```
1 y_pred_test_1 = X_test @ coefs
2 y_pred_test_2 = lr.predict(np.delete(X_test, 0, axis=1))
3
4 print("Середньоквадратичні похибки (MSE) для:")
5 err_1 = 2 * lr.compute_cost(y_pred_test_1, y_test)
6 err_2 = 2 * lr.compute_cost(y_pred_test_1, y_test)
7 print(" 1) метод найменших квадратів: {}".format(err_1))
8 print(" 2) градiєнтний спуск: {}".format(err_2))
```

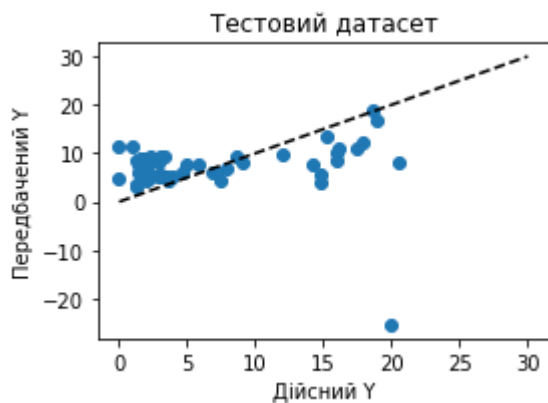
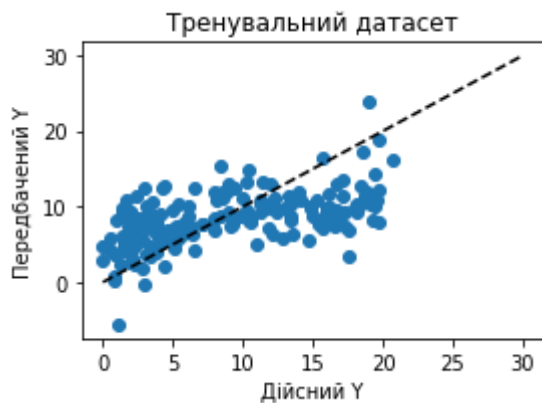
Середньоквадратичні похибки (MSE) для:

- 1) метод найменших квадратів: 70.00665695418553
- 2) градiєнтний спуск: 70.00665695418553

Візуалізація

In [17]:

```
1 plt.figure(figsize=(4, 3))
2 plt.title("Тренувальний датасет")
3 plt.scatter(y_train, y_pred)
4 plt.plot([0, 30], [0, 30], "--k")
5 plt.axis("tight")
6 plt.xlabel("Дійсний Y")
7 plt.ylabel("Передбачений Y")
8 plt.tight_layout()
9
10 plt.figure(figsize=(4, 3))
11 plt.title("Тестовий датасет")
12 plt.scatter(y_test, y_pred_test_1)
13 plt.plot([0, 30], [0, 30], "--k")
14 plt.axis("tight")
15 plt.xlabel("Дійсний Y")
16 plt.ylabel("Передбачений Y")
17 plt.tight_layout()
```



In [18]:

```

1 test_errors_1 = []
2 test_errors_2 = []
3 for i in range(50):
4     test_errors_1.append((np.linalg.norm(y_pred_test_1[i] - y_test[i]))**2)
5     test_errors_2.append((np.linalg.norm(y_pred_test_2[i] - y_test[i]))**2)

```

In [19]:

```

1 plt.figure(figsize=(17, 8))
2 ax = plt.subplot()
3
4 ax.bar(np.arange(1, y_test.shape[0]+1)-0.15, test_errors_1, width=0.25,
5         label="Похибки методу найменших квадратів для тестового набору")
6 ax.bar(np.arange(1, y_test.shape[0]+1)+0.15, test_errors_2, width=0.25,
7         label="Похибки методу градієнтного спуску для тестового набору")
8
9 plt.xticks(range(1, y_test.shape[0]+1), range(1, y_test.shape[0]+1))
10 plt.legend(loc='upper left')
11 plt.grid()
12 plt.show()

```

