

Міністерство освіти і науки України  
Національний технічний університет України «Київський політехнічний  
інститут імені Ігоря Сікорського»  
Факультет інформатики та обчислювальної техніки  
Кафедра інформатики та програмної інженерії

Звіт

з лабораторної роботи № 3 з дисципліни  
«Прикладні задачі машинного навчання»  
«Класифікація, регресія і кластеризація з  
використанням бібліотеки scikit-learn»

Варіант 7

Виконав студент:           ІІ-12 Васильєв Єгор Костянтинович

Перевірів:                 Нестерук Андрій Олександрович

Київ 2023

## Лабораторна робота №3

**Тема:** Класифікація, регресія і кластеризація з використанням бібліотеки scikit-learn

### Постановка завдання:

- 1) Первинна обробка завантажених даних середньо січневих температур.
- 2) Пошук регресійної прямої та тестування моделі.
- 3) Прогнозування температур.
- 4) Візуалізація набору даних та регресійної прямої.
- 5) Порівняння результату з минулою лабораторною роботою.
- 6) Класифікація згенерованого набору даних класифікатором SVC.
- 7) Порівняння класифікаторів KNeighborsClassifier, GaussianNB та SVC.

### Хід роботи:

- Первинна обробка даних

Contiguous U.S.				Date			
Average Temperature		January		Temperature		Anomaly	
0	Units: Degrees Fahrenheit	nan	nan	0	1895	-2.95000	-3.43
1	Base Period: 1901-2000	nan	nan	1	1896	-0.28889	1.36
2	Missing: -99	nan	nan	2	1897	-2.12778	-1.95
3	Date	Value	Anomaly	3	1898	-0.73889	0.55
4	189501	26.69	-3.43	4	1899	-1.28889	-0.44
5	189601	31.48	1.36	5	1900	1.17222	3.99
6	189701	28.17	-1.95	6	1901	-0.21111	1.50
7	189801	30.67	0.55	7	1902	-1.07778	-0.06
8	189901	29.68	-0.44	8	1903	-0.62222	0.76
9	190001	34.11	3.99	9	1904	-2.47222	-2.57
10	190101	31.62	1.50	10	1905	-2.72222	-3.02
11	190201	30.06	-0.06	11	1906	0.48889	2.76
12	190301	30.88	0.76				
13	190401	27.55	-2.57				

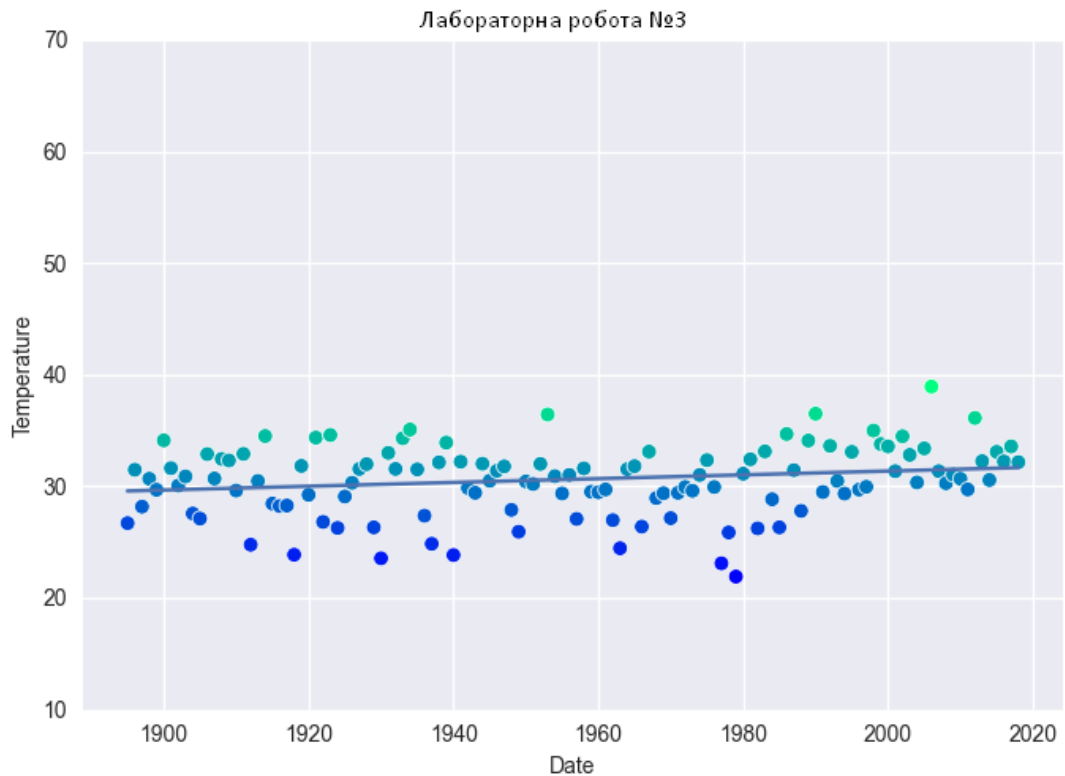
- Порівняння прогнозованих температур та фактичних

```
predicted: 30.82, expected: 28.94, difference: -6.09%
predicted: 31.55, expected: 29.71, difference: -5.82%
predicted: 30.07, expected: 26.26, difference: -12.66%
predicted: 30.29, expected: 24.84, difference: -17.99%
predicted: 30.99, expected: 25.86, difference: -16.54%
predicted: 30.63, expected: 27.07, difference: -11.62%
predicted: 30.02, expected: 34.36, difference: 14.47%
```

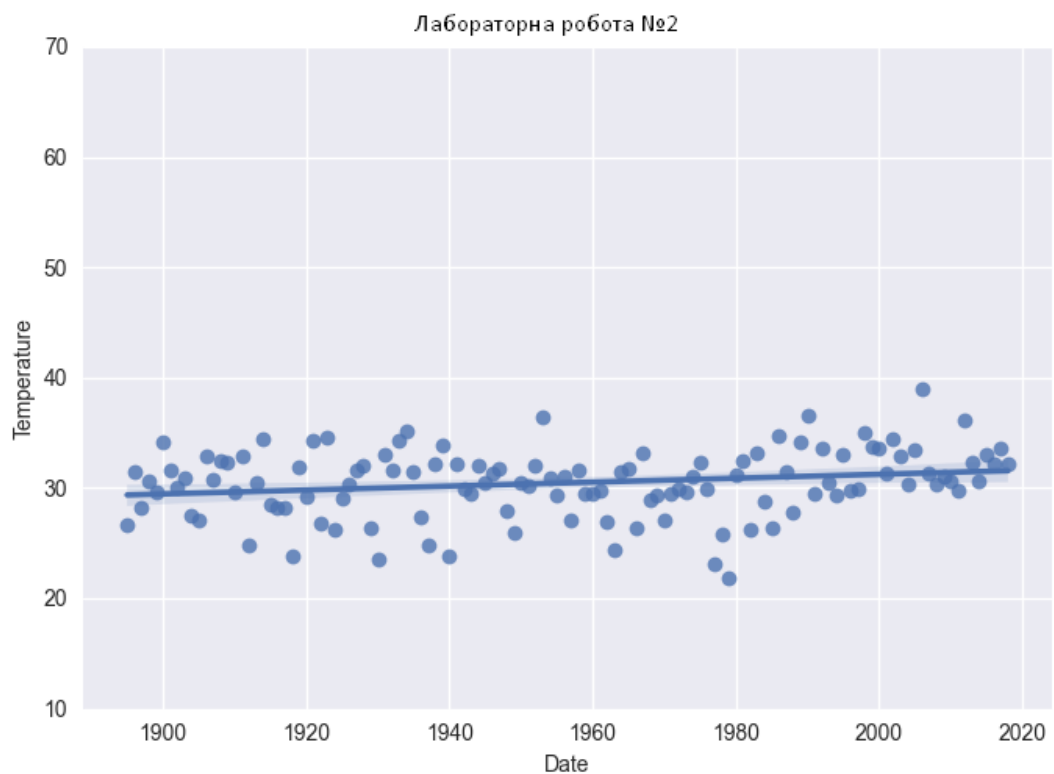
- Передбачення температур на 2019 та 1890 роки

```
Predicted temperature in 2019: [31.68369771]  
Predicted temperature in 1890: [29.48934931]
```

- Візуалізація даних

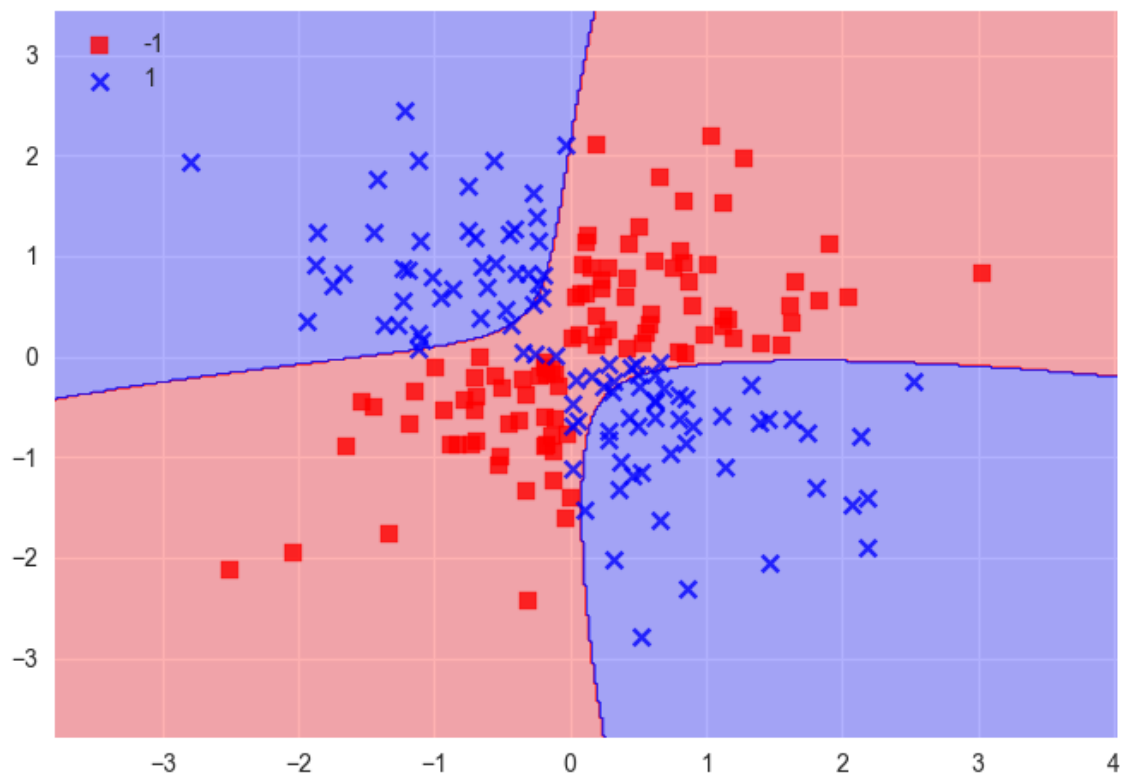


- Порівняння з регресійною прямою минулої лабораторної роботи

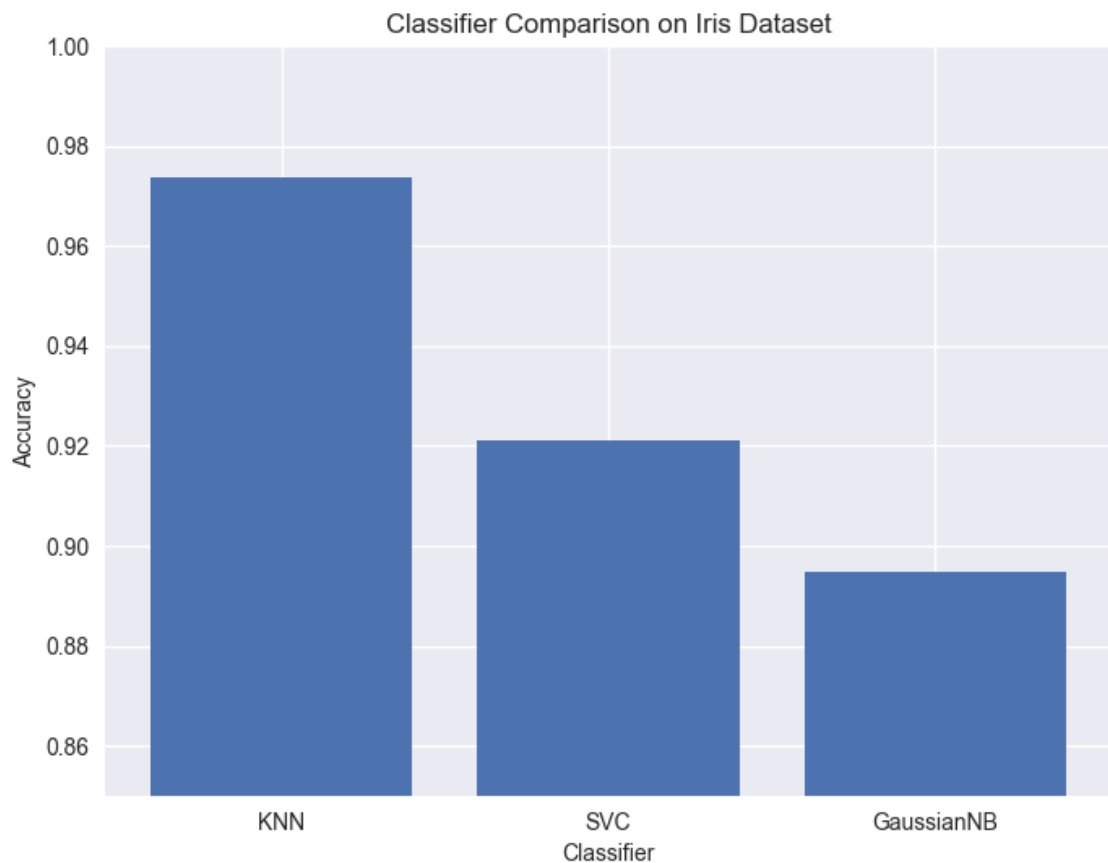


Як бачимо, лінійні регресії ідентичні, оскільки в обох бібліотеках використовується той самий метод пошуку прямої.

- Візуалізація класифікації опорних векторів згенерованих даних



- Порівняння класифікаторів на вбудованому наборі даних



Використаний алгоритм: KNeighborsClassifier використовує алгоритм k-найближчих сусідів для класифікації точок даних, SVC використовує машинний алгоритм опорних векторів, а GaussianNB використовує наївний Баєсів алгоритм. Ці алгоритми мають різні сильні та слабкі сторони та можуть по-різному працювати з різними типами наборів даних.

Швидкість навчання: KNeighborsClassifier і GaussianNB відносно швидкі для навчання, оскільки вони не потребують обчислення великої кількості параметрів. SVC, з іншого боку, може бути повільним для навчання, особливо на великих наборах даних.

Точність прогнозування: точність прогнозування кожного класифікатора може відрізнятись залежно від характеру набору даних, що аналізується. Загалом SVC добре працює на широкому діапазоні наборів даних і часто вважається хорошим універсальним вибором. KNeighborsClassifier також може добре працювати, особливо на менших наборах даних, але може бути більш чутливим до шуму в даних. GaussianNB часто використовується для класифікації тексту та інших програм, що включають розріджені дані.

Параметри налаштування: кожен класифікатор має інший набір параметрів налаштування, які можна налаштувати для оптимізації його продуктивності на певному наборі даних. Наприклад, KNeighborsClassifier має параметр `n_neighbors`, який керує кількістю сусідів, які використовуються в класифікації, тоді як SVC має параметр `C`, який керує штрафом за неправильну класифікацію. Налаштування цих параметрів може бути важливим кроком в оптимізації продуктивності кожного класифікатора.

## **Висновок**

Було поглиблено знання бібліотек `pandas`, `matplotlib`, `seaborn`, `numpy` та досліджено різні типи класифікаторів `Scikit-Learn`; було побудовано лінійну регресію та візуалізовано зміну середньої температури січня континентальної Америки; було досліджено класифікатор на базі методу опорних векторів для згенерованого набору даних; було порівняно декілька класифікаційних оцінювачів: KNeighborsClassifier, SVC та GaussianNB для вбудованого в `scikit-learn` Iris датасету; таким чином було вивчено тему «Класифікація, регресія і кластеризація з використанням бібліотеки `scikit-learn`»

```

1  import ...
13
14  pd.set_option('display.precision', 2)
15  plt.style.use('seaborn-v0_8')
16
17  initial_df = pd.read_csv(r'F:\Egor\Уроки\Машинне навчання\Ла62\1895-2018.csv')
18  nyc = initial_df.copy()
19  nyc.columns = ['Date', 'Temperature', 'Anomaly']
20  nyc.drop(index=range(4), inplace=True)
21  nyc.reset_index(drop=True, inplace=True)
22  nyc = nyc.astype({'Date': 'int64', 'Temperature': 'float'})
23  nyc.loc[:, 'Date'] = nyc.loc[:, 'Date'].floordiv(100)
24  print(nyc.head())
25  dates = nyc.loc[:, 'Date'].values
26  temps = nyc.loc[:, 'Temperature'].values
27
28  X_train, X_test, Y_train, Y_test = train_test_split(nyc.loc[:, 'Date'].values.reshape(-1, 1),
29                                                    nyc.loc[:, 'Temperature'].values, random_state=11)
30
31  linear_regression = LinearRegression()
32  linear_regression.fit(X=X_train, y=Y_train)
33  k = linear_regression.coef_
34  b = linear_regression.intercept_
35  predicted = linear_regression.predict(X_test)
36  expected = Y_test
37  for p, e in zip(predicted[:5], expected[:5]):
38      print(f'predicted: {p:.2f}, expected: {e:.2f}, difference: {(e-p)/p*100:.2f}%')
39  predict = (lambda x: linear_regression.coef_ * x + linear_regression.intercept_)
40  print('Predicted temperature in 2019:', predict(2019))
41  print('Predicted temperature in 1890:', predict(1890))
42  axes = sns.scatterplot(data=nyc, x='Date', y='Temperature', hue='Temperature', palette='winter', legend=False)
43  axes.set_ylim(10, 70)
44  x = np.array([min(nyc.loc[:, 'Date'].values), max(nyc.loc[:, 'Date'].values)])
45  y = predict(x)
46  line = plt.plot(x, y)
47  plt.show()
48  # part 2
49  np.random.seed(1)
50  X_xor = np.random.randn(200, 2)
51  Y_xor = np.logical_xor(X_xor[:, 0] > 0, X_xor[:, 1] > 0)
52  Y_xor = np.where(Y_xor, 1, -1)
53  plt.scatter(X_xor[Y_xor == 1, 0], X_xor[Y_xor == 1, 1], c='b', marker='x', label='1')
54  plt.scatter(X_xor[Y_xor == -1, 0], X_xor[Y_xor == -1, 1], c='r', marker='s', label='-1')
55  plt.xlim(-3, 3)
56  plt.ylim(-3, 3)
57  plt.legend(loc='best')
58  plt.show()

```

**Вихідний код**

```

60 def plot_decision_regions(X, Y, classifier, test_idx=None, resolution=0.02):
61     markers = ('s', 'x', 'o', '^', 'v')
62     colors = ('red', 'blue', 'lightgreen', 'gray', 'cyan')
63     cmap = ListedColormap(colors[:len(np.unique(Y))])
64     x1_min, x1_max = X[:, 0].min() - 1, X[:, 0].max() + 1
65     x2_min, x2_max = X[:, 1].min() - 1, X[:, 1].max() + 1
66     xx1, xx2 = np.meshgrid(np.arange(x1_min, x1_max, resolution), np.arange(x2_min, x2_max, resolution))
67     Z = classifier.predict(np.array([xx1.ravel(), xx2.ravel()]).T)
68     Z = Z.reshape(xx1.shape)
69     plt.contourf(xx1, xx2, Z, alpha=0.3, cmap=cmap)
70     plt.xlim(xx1.min(), xx1.max())
71     plt.ylim(xx2.min(), xx2.max())
72     for idx, c1 in enumerate(np.unique(Y)):
73         plt.scatter(x=X[Y == c1, 0], y=X[Y == c1, 1], alpha=0.8, c=colors[idx], marker=markers[idx], label=c1)
74
75     if test_idx:
76         X_test, Y_test = X[test_idx, :], Y[test_idx]
77         plt.scatter(X_test[:, 0], X_test[:, 1], c='', edgecolor='black', alpha=1.0, linewidth=1, marker='o', s=100)
78
79
80 svm = SVC(kernel='rbf', random_state=1, gamma=0.10, C=10.0)
81 svm.fit(X_xor, Y_xor)
82 plot_decision_regions(X_xor, Y_xor, classifier=svm)
83 plt.legend(loc='upper left')
84 plt.show()
85
86 # part 3
87
88 iris = load_iris()
89
90 # Split dataset into training and testing data
91 X_train, X_test, y_train, y_test = train_test_split(iris.data, iris.target, random_state=11)
92
93 # Initialize classifiers
94 knn = KNeighborsClassifier()
95 svc = SVC()
96 gnb = GaussianNB()
97
98 # Train classifiers
99 knn.fit(X_train, y_train)
100 svc.fit(X_train, y_train)
101 gnb.fit(X_train, y_train)
102
103 # Make predictions on testing data
104 knn_prediction = knn.predict(X_test)
105 svc_prediction = svc.predict(X_test)
106 gnb_prediction = gnb.predict(X_test)
107
108 # Calculate accuracy scores
109 knn_accuracy = accuracy_score(y_test, knn_prediction)
110 svc_accuracy = accuracy_score(y_test, svc_prediction)
111 gnb_accuracy = accuracy_score(y_test, gnb_prediction)
112

```

```

113 # Print accuracy scores
114 print(f'KNN Accuracy: {knn_accuracy:.3f}')
115 print(f'SVC Accuracy: {svc_accuracy:.3f}')
116 print(f'GaussianNB Accuracy: {gnb_accuracy:.3f}')
117
118 # Plot results
119 labels = ['KNN', 'SVC', 'GaussianNB']
120 accuracies = [knn_accuracy, svc_accuracy, gnb_accuracy]
121 plt.figure(figsize=(8, 6))
122 plt.bar(labels, accuracies)
123 plt.title('Classifier Comparison on Iris Dataset')
124 plt.xlabel('Classifier')
125 plt.ylabel('Accuracy')
126 plt.ylim(0.85, 1)
127 plt.show()
128

```