



Міністерство освіти і науки України
Національний технічний університет України
“Київський політехнічний інститут імені Ігоря Сікорського”
Факультет інформатики та обчислювальної техніки
Кафедра інформаційних систем і технологій

Лабораторна робота №3
Прикладні задачі машинного навчання
“Класифікація, регресія і кластеризація з використанням
бібліотеки scikit-learn”

Виконав
студент групи ІК – 33:
Вересоцький А. Ю.

Перевірив:
асистент кафедри ІСТ
Нестерук А.О.

Завдання:

1. Потрібно завантажити метеорологічні дані в 1895-2024 роках з CSV-файлу в DataFrame. Після цього дані треба буде відформатувати для використання.

Завантажую середні січніві температури в Лексінгтон штат Кентуккі з 1895 по 2024 рік через часові ряди NOAA «Climate at a Glance»:

```
1 lex = pd.read_csv('data.csv')
✓ [67] < 10 ms

1 lex.head()
✓ [68] < 10 ms
```

#	Lexington	Kentucky January Average Temperature
0	# Units: Degrees Fahrenheit	NaN
1	# Missing: -99	NaN
2	Date	Value
3	189501	28.8
4	189601	35.6

Форматування:

Видаляємо непотрібні рядки, оновлюємо індекси та змінюємо назви колонок:

```
1 lex.drop(index=[0, 1, 2], inplace=True)
✓ [69] < 10 ms

1 lex.reset_index(drop=True, inplace=True)
✓ [70] < 10 ms

1 lex.columns = ['Date', 'Temperature']
✓ [71] < 10 ms
```

```
1 lex.head()
✓ [72] < 10 ms
```

	Date	Temperature
0	189501	28.8
1	189601	35.6
2	189701	30.3
3	189801	39
4	189901	35.1

Перевіряємо типи даних та змінюємо їх за необхідності:

```
1 lex.dtypes
✓ [73] < 10 ms

2 rows | Length: 2, dtype: object
┌───┬───┐
│   │   │
├───┴───┘
Date    object
Temperature object

1 lex.Date = lex.Date.astype(int)
2 lex.Temperature = lex.Temperature.astype(float)
✓ [74] < 10 ms

1 lex.Date = lex.Date.floordiv(100)
✓ [75] < 10 ms
```

```
1 lex.head()
✓ [76] < 10 ms

5 rows | 5 rows x 2 cols
┌───┬───┬───┐
│   │ Date │ Temperature │
├───┴───┴───┘
0    1895    28.8
1    1896    35.6
2    1897    30.3
3    1898    39.0
4    1899    35.1
```

2. Розбиття даних для навчання і тестування.

Дані розбиваємо на навчальний і тестовий набори. Ключовий аргумент `random_state` використовується для забезпечення відтворюваності результатів:

```
1 X_train, X_test, y_train, y_test = train_test_split(lex.Date.values.reshape(-1, 1),
lex.Temperature.values, random_state=11)
✓ [349] < 10 ms
```

Для перевірки пропорції навчальних тестових даних (75% до 25%) задамо розміри `X_train` і `X_test`:

```
1 X_train.shape
✓ [350] < 10 ms
(97, 1)

1 X_test.shape
✓ [351] < 10 ms
(33, 1)
```

3. Навчання моделі.

Скористаємося оцінювачем `LinearRegression`:

```
1 linear_regression = LinearRegression()
2 linear_regression.fit(X_train, y_train)
✓ [352] < 10 ms
```

▼ LinearRegression ⓘ ?

LinearRegression()

Кут нахилу зберігається в атрибуті `coeff_` оцінювача, а точка перетину - в атрибуті `intercept_`:

```
1 linear_regression.coef_  
✓ [353] < 10 ms  
array([-0.00189958])  
  
1 linear_regression.intercept_  
✓ [354] < 10 ms  
np.float64(38.1286603204756)
```

Пізніше ці значення будуть використані для виведення регресійної прямої і прогнозування для конкретних дат.

4. Тестування моделі.

Протестуємо модель за даними з `X_test` і перевіримо прогнози по набору даних, виводячи прогнозовані і очікувані значення для кожного п'ятого елементу:

```
1 predicted = linear_regression.predict(X_test)  
✓ [355] < 10 ms  
  
1 expected = y_test  
✓ [356] < 10 ms  
  
1 for p, e in zip(predicted[::5], expected[::5]):  
2     print(f'predicted: {p:.2f}, expected: {e:.2f}')
```

✓ [357] < 10 ms

predicted:	34.53,	expected:	28.80
predicted:	34.41,	expected:	37.40
predicted:	34.38,	expected:	36.50
predicted:	34.32,	expected:	43.90
predicted:	34.30,	expected:	27.80
predicted:	34.33,	expected:	33.60
predicted:	34.49,	expected:	36.00

5. Прогнозування майбутніх температур і оцінка минулих температур.

Скористаємося отриманими значеннями кута нахилу і точки перетину для прогнозування середньої температури в січні 2019 року, а також оцінки середньої температури в січні 1890 року.

```

1 predict = (lambda x: linear_regression.coef_ * x + linear_regression.intercept_)
✓ [358] < 10 ms

1 predict(2019)
✓ [359] < 10 ms

array([34.29341246])

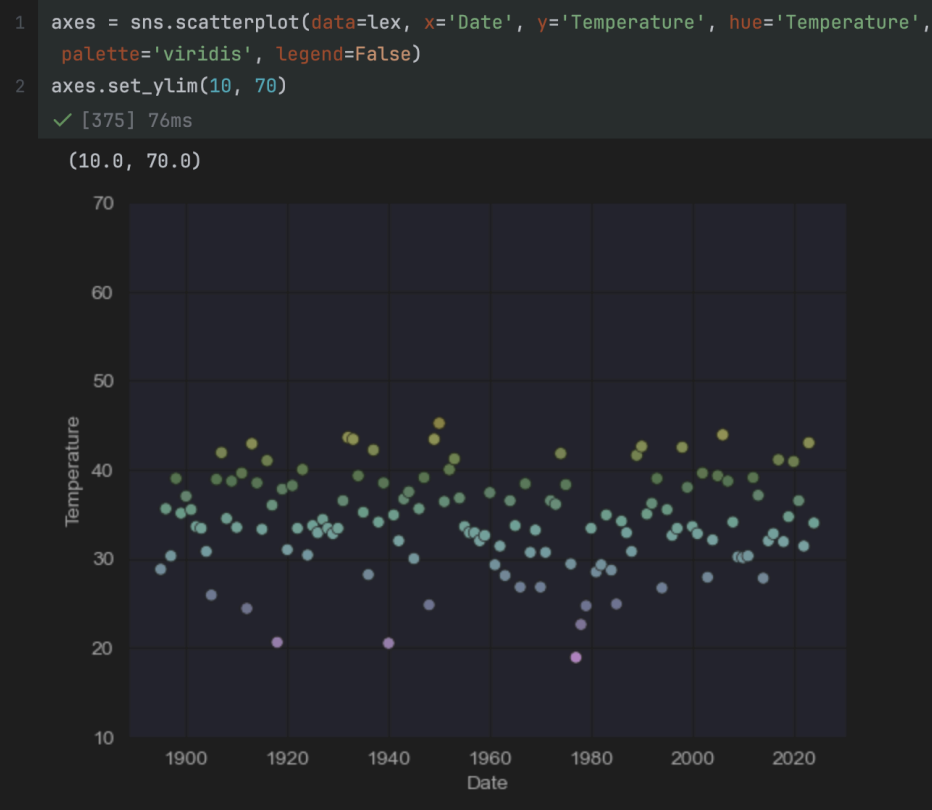
1 predict(1890)
✓ [360] < 10 ms

array([34.53845801])

```

6. Візуалізація набору даних з регресійній прямий.

Тепер побудуємо діаграму розкиду даних за допомогою функції `scatterplot` бібліотеки `Seaborn` і функції `plot` бібліотеки `Matplotlib`. Для виведення точок даних скористаємося методом `scatterplot` з колекцією `DataFrame` з ім'ям `lex` та змінимо масштаб осі `y`, щоб при виведенні регресійної прямої лінійність відносини була більш очевидною:



Перейдемо до висновку регресійної прямої. Почнемо зі створення масиву, що містить мінімальні і максимальні значення дати з `lex.Date`. Вони стануть координатами `x` початкової і кінцевої точок регресійної прямої.

В результаті передачі `predict` масиву `x` буде отримано масив відповідних прогнозованих значень, які будуть використовуватися в якості координат `y`.

Нарешті, функція `plot` бібліотеки `Matplotlib` малює лінію по масивам `x` та `y`, що зображує координати `x` та `y` точок відповідно:

Отримана діаграма розкиду даних з регресійній прямий зображена на наступній діаграмі.

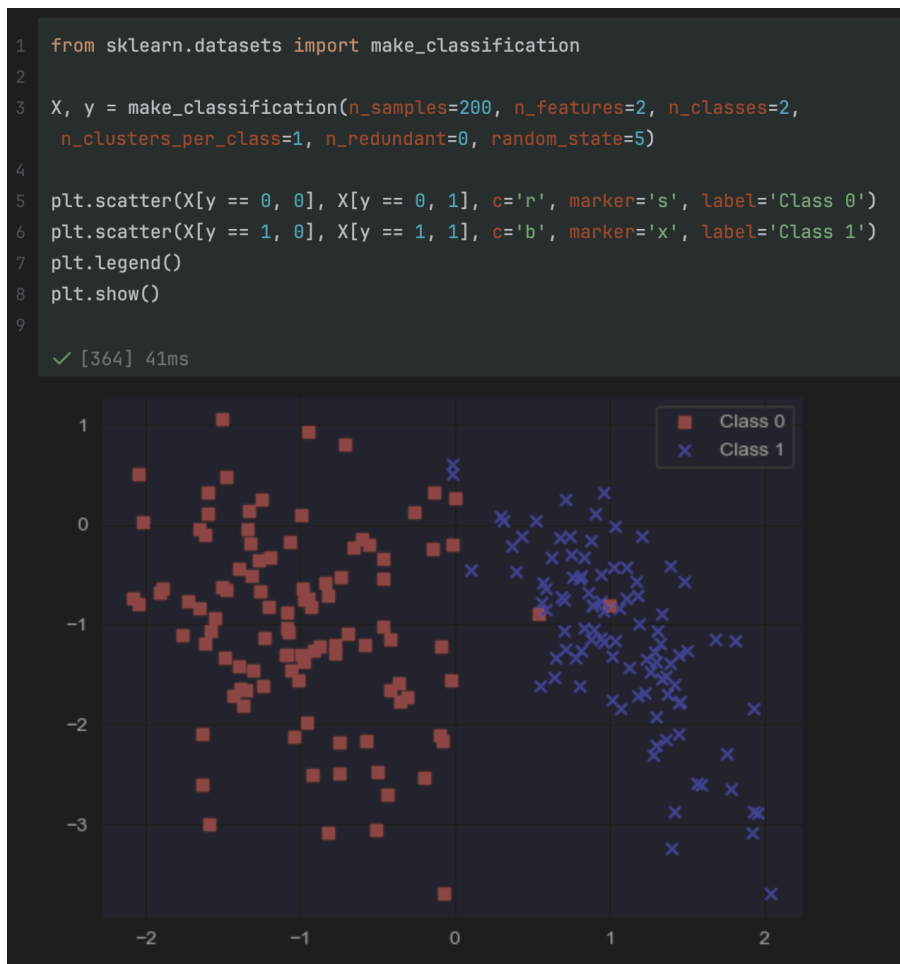


7. SVC класифікація.

генеруйте набір даних

та класифікуйте його використавши класифікатор SVC

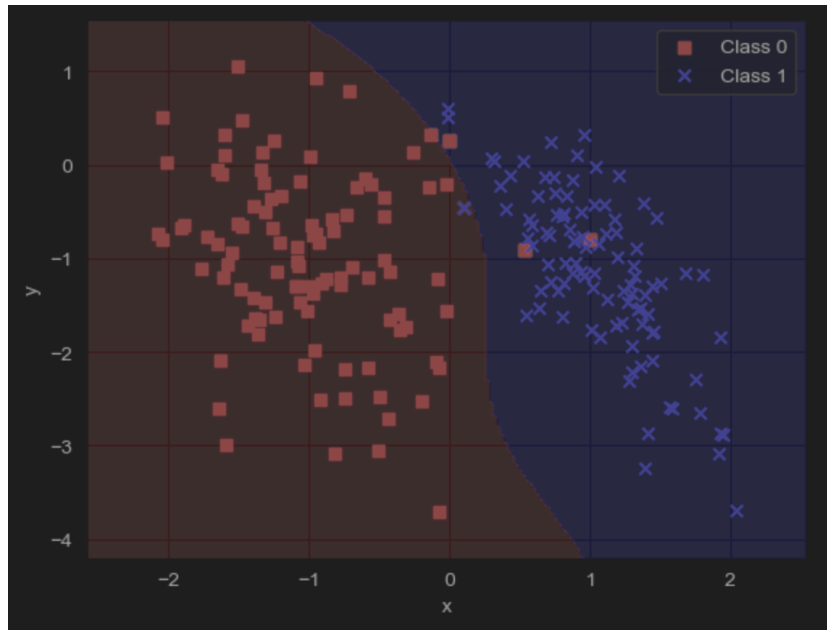
Генеруємо набір даних за допомогою вбудованої у бібліотеку `scikit-learn` функцію **`make_classification`** та візуалізуємо ці дані:



Використовуємо Support Vector Machine (SVM) для класифікації наших двовимірних даних та візуалізуємо рішення моделі:

```
1 svm = SVC(C=10).fit(X, y)
2
3 x_min, x_max = X[:, 0].min() - .5, X[:, 0].max() + .5
4 y_min, y_max = X[:, 1].min() - .5, X[:, 1].max() + .5
5
6 xx, yy = np.meshgrid(np.linspace(x_min, x_max, 200), np.linspace(y_min, y_max, 200))
7
8 Z = svm.predict(np.c_[xx.ravel(), yy.ravel()]).reshape(xx.shape)
9
10 plt.contourf(xx, yy, Z, alpha=0.2, cmap='bwr_r')
11 plt.scatter(X[y == 0, 0], X[y == 0, 1], c='r', marker='s', label='Class 0')
12 plt.scatter(X[y == 1, 0], X[y == 1, 1], c='b', marker='x', label='Class 1')
13
14 plt.xlabel('x')
15 plt.ylabel('y')
16 plt.xlim([x_min, x_max])
17 plt.ylim([y_min, y_max])
18 plt.legend()
19 plt.show()
```

✓ [1003] 81ms



8. Детальне порівняння класифікаційних оцінювачів KNeighborsClassifier, SVC та GaussianNB для вбудованого в scikit-learn набору даних breast cancer.

Спочатку завантажуюємо набір даних breast cancer та розбиваємо дані на навчальний та тестовий набори:

```
4 from sklearn.datasets import load_breast_cancer
5
6 cancer = load_breast_cancer()
7 X_train, X_test, y_train, y_test = train_test_split(cancer.data, cancer.target,
8     stratify=cancer.target, random_state=42)
```

Проводимо масштабування даних та навчання моделей:

```
8 scaler = StandardScaler()
9 X_train = scaler.fit_transform(X_train)
10 X_test = scaler.transform(X_test)
11
12 models = {
13     "KNN (k=5)": KNeighborsClassifier(n_neighbors=5),
14     "SVC (RBF Kernel)": SVC(kernel='rbf', C=1.0, gamma='scale', random_state=42),
15     "GaussianNB": GaussianNB()
16 }
17
18 results = {}
19 for name, model in models.items():
20     model.fit(X_train, y_train)
21     y_pred = model.predict(X_test)
22
23     acc = accuracy_score(y_test, y_pred)
24     report = classification_report(y_test, y_pred, target_names=cancer.target_names,
25     output_dict=True)
26
27     results[name] = {
28         "accuracy": acc,
29         "precision": report["macro avg"]["precision"],
30         "recall": report["macro avg"]["recall"],
31         "f1-score": report["macro avg"]["f1-score"],
32         "confusion_matrix": confusion_matrix(y_test, y_pred)
33     }
```


Виводимо результати та візуалізуємо матрицю помилок:

```
34 for name, metrics in results.items():
35     print(f"\n==== {name} ====")
36     print(f"Accuracy: {metrics['accuracy']:.4f}")
37     print(f"Precision: {metrics['precision']:.4f}")
38     print(f"Recall: {metrics['recall']:.4f}")
39     print(f"F1-score: {metrics['f1-score']:.4f}")
40
41 fig, axes = plt.subplots(1, 3, figsize=(15, 4))
42
43 for ax, (name, metrics) in zip(axes, results.items()):
44     sns.heatmap(metrics["confusion_matrix"], annot=True, fmt="d", cmap="Blues", ax=ax,
45                 xticklabels=cancer.target_names, yticklabels=cancer.target_names)
46     ax.set_title(name)
47     ax.set_xlabel("Predicted Label")
48     ax.set_ylabel("True Label")
49
50 plt.tight_layout()
51 plt.show()
```

✓ [403] 214ms

==== KNN (k=5) ====

Accuracy: 0.9790

Precision: 0.9839

Recall: 0.9717

F1-score: 0.9772

==== SVC (RBF Kernel) ====

Accuracy: 0.9790

Precision: 0.9759

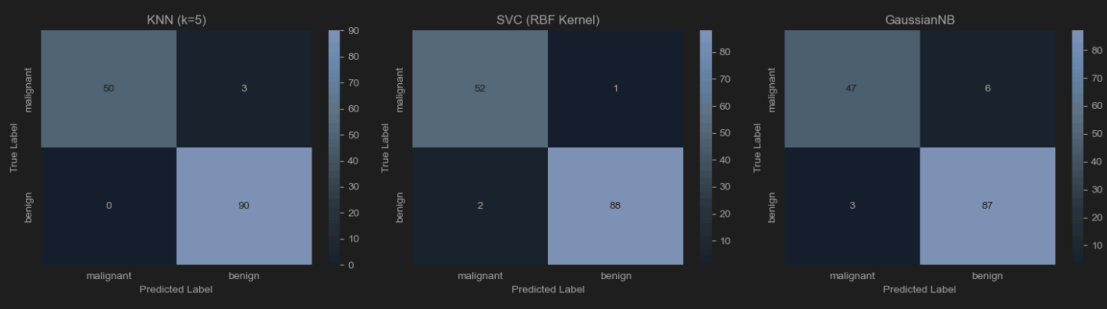
Recall: 0.9795

F1-score: 0.9776

==== GaussianNB ====

Accuracy: 0.9371

Precision: 0.9377



Висновки

KNN та SVC демонструють дуже схожі результати, майже ідеальні.

- KNN виграє за Precision, тому він робить менше False Positives (не буде марно турбувати пацієнтів).
- SVC виграє за Recall, тому він рідше пропускає реальні випадки хвороби (краще для медичних діагнозів).
- F1-score у них практично однаковий, тому можна вибрати будь-який.

GaussianNB значно поступається обом іншим алгоритмам, особливо за Recall (він пропускає більше хворих пацієнтів). Це може бути через те, що наївний байєсівський класифікатор припускає, що всі ознаки незалежні, що не завжди відповідає реальним даним.