

# Лабораторна робота №3. Байєсівський класифікатор

Виконав студент групи КМ-91мп

Галета М.С.

## Завдання на лабораторну роботу

1. Розділити файл з даними на навчальну та тестову вибірки.
2. Збудувати наївний байєсівський класифікатор для визначення значення цільової характеристики (останній стовпчик) на основі вхідних аргументів (початкові стовпчики).
3. Перевірити якість його роботи, використовуючи файл із даними тестової вибірки.
4. Надати графічне представлення результатів.

In [1]:

```
1 import numpy as np
2 import pandas as pd
3 import matplotlib.pyplot as plt
4 from sklearn.model_selection import train_test_split
```

## Зчитування датасету

In [2]:

```

1 dataset = pd.read_csv('MP-04-Galeta.csv', sep=';',
2                       names=['x1', 'x2', 'x3', 'x4', 'x5', 'x6', 'x7', 'x8', 'x9', 'x10', 'x11', 'x12', 'y'],
3                       dtype=str)
4 dataset.head(5)

```

Out[2]:

	x1	x2	x3	x4	x5	x6	x7	x8	x9	x10	x11	x12	y
0	x	y	g	t	n	c	b	n	t	s	n	d	edible
1	x	f	n	t	n	c	b	n	t	s	n	d	edible
2	f	s	w	f	n	w	b	k	t	f	n	g	edible
3	f	y	w	t	p	c	n	p	e	s	n	g	poisonous
4	f	y	n	t	p	c	n	k	e	s	n	g	poisonous

## 1) Розділення датасету на тренувальну та тестову вибірки

In [3]:

```

1 X_train, X_test, y_train, y_test = train_test_split(dataset.iloc[:, :-1].values,
2                                                     dataset.iloc[:, -1].values,
3                                                     test_size=0.4,
4                                                     random_state=42)

```

In [4]:

```

1 def categories_words(x_train, y_train):
2     """ Функція, яка повертає списки зі слів,
3         що відносяться до відповідних класів """
4
5     all_words_list = []
6     edible_words_list = []
7     poisonous_words_list = []
8
9     for i, y in enumerate(y_train):
10         if y == 'edible':
11             edible_words_list += list(x_train[i])
12         if y == 'poisonous':
13             poisonous_words_list += list(x_train[i])
14     all_words_list = set(edible_words_list + poisonous_words_list)
15
16     return all_words_list, edible_words_list, poisonous_words_list

```

## 2) Побудова наївного байєсівського класифікатора

In [5]:

```

1 class NaiveBayesClassifier:
2     def __init__(self, alpha):
3         self.alpha = alpha
4
5         self.train_set_x = None
6         self.train_set_y = None
7
8         self.all_words_list = []
9         self.ham_words_list = []
10        self.spam_words_list = []
11
12    def fit(self, train_set_x, train_set_y):
13
14        self.all_words_list, self.edible_words_list, self.poisonous_words_list = categorize(train_set_x)
15
16        self.edible_words_list = np.array(self.edible_words_list)
17        self.poisonous_words_list = np.array(self.poisonous_words_list)
18        self.theta_edible = {}
19        self.theta_poisonous = {}
20        k = len(self.all_words_list)
21        n_p, n_e = len(self.poisonous_words_list), len(self.edible_words_list)
22        unique_p, counts_p = np.unique(self.poisonous_words_list, return_counts=True)
23        x_p = dict(zip(unique_p, counts_p))
24        unique_e, counts_e = np.unique(self.edible_words_list, return_counts=True)
25        x_e = dict(zip(unique_e, counts_e))
26
27        for w in self.all_words_list:
28            self.theta_edible[w] = (x_e.get(w, 0) + self.alpha) / (n_e + self.alpha * k)
29            self.theta_poisonous[w] = (x_p.get(w, 0) + self.alpha) / (n_p + self.alpha * k)
30
31        self.p_p = (train_set_y == 'poisonous').sum() / len(train_set_y)
32        self.p_e = (train_set_y == 'edible').sum() / len(train_set_y)
33
34    def predict(self, test_set_x):
35
36        prediction = []
37        for sent in test_set_x:
38            pr_p = np.log(self.p_p)
39            pr_e = np.log(self.p_e)
40            for w in sent:
41                if w in self.theta_poisonous:
42                    pr_p += np.log(self.theta_poisonous.get(w, 1))
43                if w in self.theta_edible:
44                    pr_e += np.log(self.theta_edible.get(w, 1))
45            prediction += ['poisonous' if pr_p >= pr_e else 'edible']
46
47        return prediction

```

In [6]:

```

1 alpha = 1 # Коефіцієнт згладжування
2 nbc = NaiveBayesClassifier(alpha)
3 nbc.fit(X_train, y_train)

```

In [7]:

```

1 y_pred = np.array(nbc.predict(X_test))

```

### 3) Оцінка результату на тестовій вибірці

In [8]:

```
1 actual = list(y_test)
2 accuracy = (y_pred == y_test).mean()
3 print("Доля правильно вгаданих відповідей на тестовій вибірці:", accuracy, '%')
```

Доля правильно вгаданих відповідей на тестовій вибірці: 0.98 %

### 4) Графічне представлення результатів

In [9]:

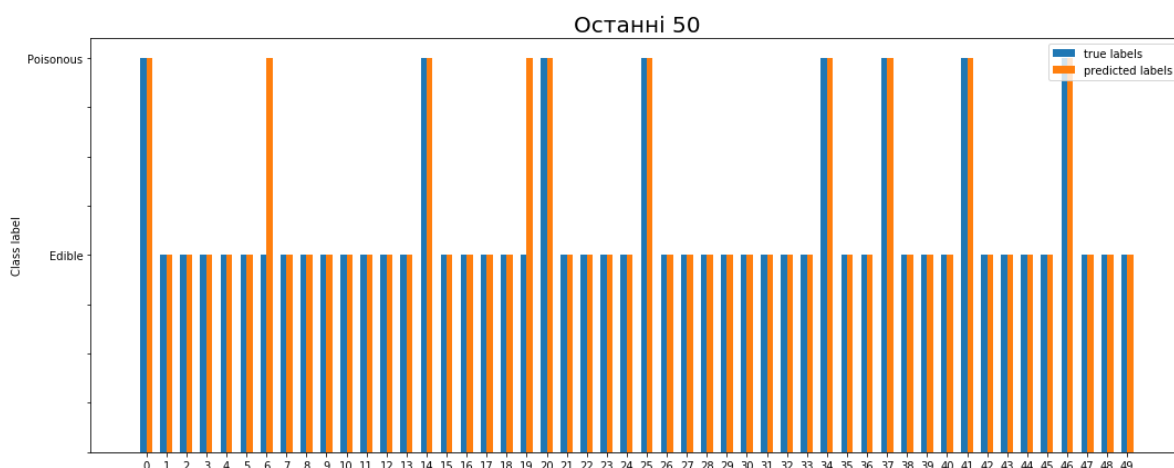
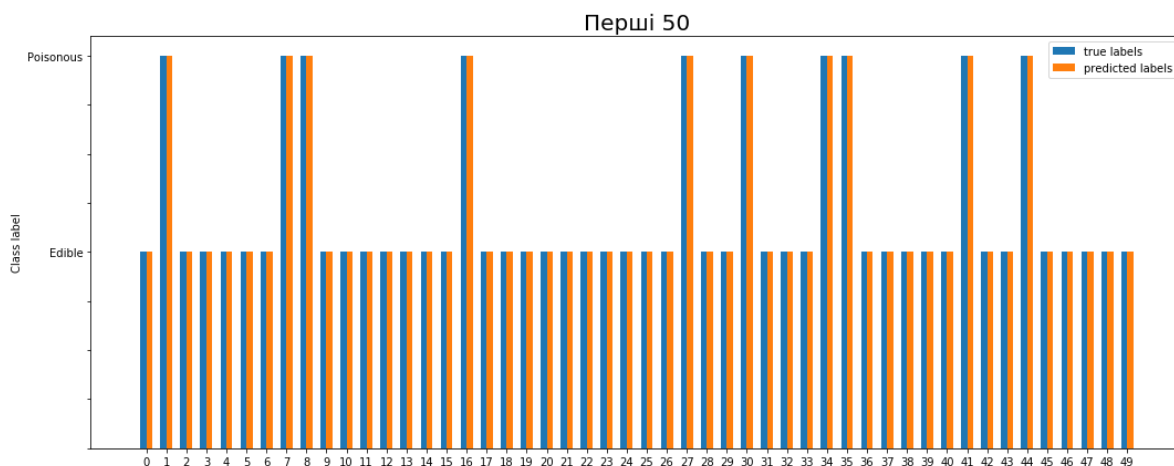
```
1 y_test_plot = np.where(y_test == 'edible', 1, y_test)
2 y_test_plot = np.where(y_test_plot == 'poisonous', 2, y_test_plot)
3 y_pred_plot = np.where(y_pred == 'edible', 1, y_pred)
4 y_pred_plot = np.where(y_pred_plot == 'poisonous', 2, y_pred_plot).astype(int)
```

In [10]:

```

1 x = np.arange(len(y_test[:50]))
2 width = 0.3
3 fig, ax = plt.subplots(figsize=(15,6))
4 plt.title("Перші 50", fontsize=20)
5 ax.bar(x - width/2, y_test_plot[:50], width, label='true labels')
6 ax.bar(x + width/2, y_pred_plot[:50], width, label='predicted labels')
7 ax.set_ylabel('Class label')
8 ax.set_xticks(x)
9 ax.set_yticklabels(('','','','Edible','Edible','','','Poisonous'))
10 ax.legend()
11 plt.tight_layout()
12
13 x = np.arange(len(y_test[50:]))
14 width = 0.3
15 fig, ax = plt.subplots(figsize=(15,6))
16 plt.title("Останні 50", fontsize=20)
17 ax.bar(x - width/2, y_test_plot[50:], width, label='true labels')
18 ax.bar(x + width/2, y_pred_plot[50:], width, label='predicted labels')
19 ax.set_ylabel('Class label')
20 ax.set_xticks(x)
21 ax.set_yticklabels(('','','','Edible','Edible','','','Poisonous'))
22 ax.legend()
23 plt.tight_layout()

```



## Висновки

1) Було збудовано наївний байєсівський класифікатор

**2) Класифікатор показує високу долю правильно вгаданих відповідей, помилився лише в 2-х випадках**