

**Московский Авиационный Институт
(Национальный исследовательский университет)**

Кафедра 806

**Отчет
по лабораторной работе №1
по курсу «Искусственный интеллект»**

Студент: Ильиных В. М.

Группа: М8О-301Б-19

Москва, 2022

Задача

1. Реализовать алгоритмы Linear/Logistic Regression, KNN, SVM, Naive Bayes в отдельных классах
2. Классы должны наследоваться от BaseEstimator и ClassifierMixin
3. Организовать процесс предобработки, обучения и тестирования с помощью Pipeline
4. Настроить гиперпараметры с помощью кросс-валидации (GridSearchCV, RandomSearchCV), вывести и сохранить эти гиперпараметры в файл вместе с обученными моделями
5. Для каждой модели получить оценки метрик: Confusion Matrix, Accuracy, Recall, Precision, ROC_AUC curve
6. Проанализировать полученные результаты и сделать выводы о применимости моделей

Ход работы

Логистическая регрессия имеет 4 параметра, значит используем поиск по сетке, но сначала нормализуем.

```
1. class Logistic_Regression(BaseEstimator, ClassifierMixin):
2.     def __init__(self, lr=0.1, batch=10, epochs=1, alpha=0.0001):
3.         self.lr = lr
4.         self.batch = batch
5.         self.epochs = epochs
6.         self.alpha = alpha
7.
8.     def fit(self, data, labels):
9.         self.w = np.random.normal(0, 1, (data.shape[1]+1,))
10.        data = np.concatenate((data, np.ones((data.shape[0],1))),
11.                               axis=1)
12.        for _ in range(self.epochs):
13.            for i in range(self.batch, len(data), self.batch):
14.                data_batch = data[i-self.batch:i]
15.                labels_batch = labels[i-self.batch:i]
16.                pred = self.sigmoid(np.dot(self.w, data_batch.T))
17.                grad = 2 * self.alpha * self.w + np.dot(pred -
18.                    labels_batch, data_batch)
19.                self.w -= self.lr * grad
```

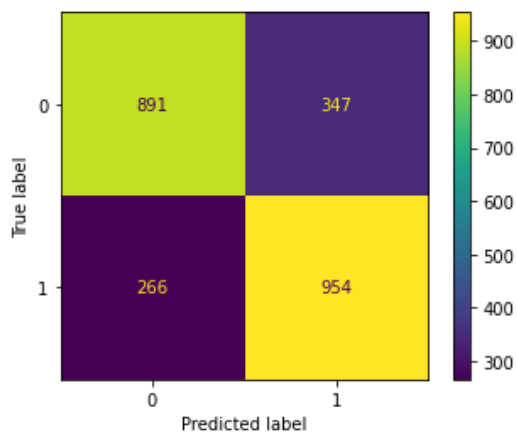
```

20.         return self
21.
22.     def sigmoid(self, x):
23.         warnings.filterwarnings('ignore')
24.         return 1 / (1 + np.exp(-x))
25.
26.     def predict(self, data):
27.         return (self.sigmoid(np.concatenate((data,
np.ones((data.shape[0],1))), axis=1).dot(self.w)) >
0.5).astype('int64')

```

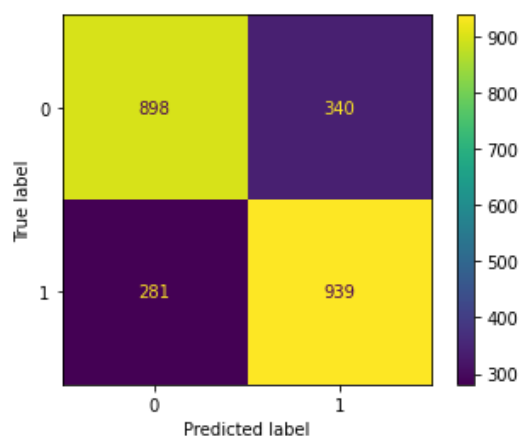
Результаты:

Accuracy: 0.7506102522375916
Precision: 0.7332820906994619
Recall: 0.7819672131147541



Результаты логистической регрессии из sklearn

Accuracy: 0.7473555736371034
Precision: 0.7341673182173573
Recall: 0.769672131147541

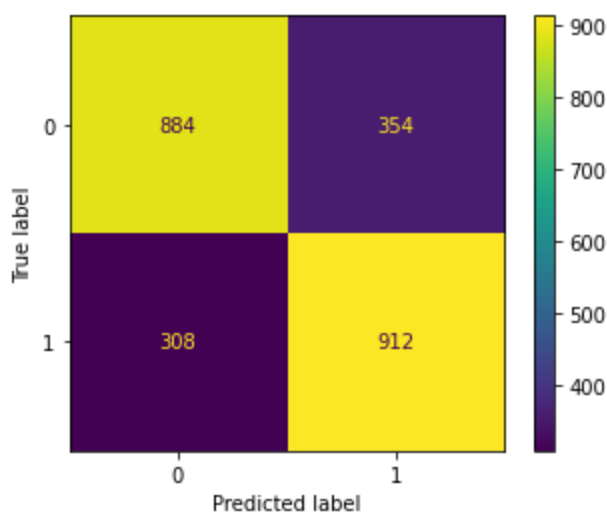


Метод ближайших соседей – KNN – имеет единственный параметр – собственно количество соседей. Для вычисления расстояния используется евклидова метрика.

```
1. class KNN(BaseEstimator, ClassifierMixin):
2.     def __init__(self, k=5):
3.         self.k = k
4.
5.     def fit(self, data, labels):
6.         self.data = data
7.         self.labels = labels
8.         return self
9.
10.    def predict(self, data):
11.        res = np.ndarray((data.shape[0],))
12.        for i, x in enumerate(data):
13.            neighbors = np.argpartition(((self.data - data[i]) **
14.                                         2).sum(axis=1), self.k - 1)[:self.k]
15.            values, counts = np.unique(self.labels[neighbors],
16.                                       return_counts=True)
17.            res[i] = values[counts.argmax()]
18.        return res
```

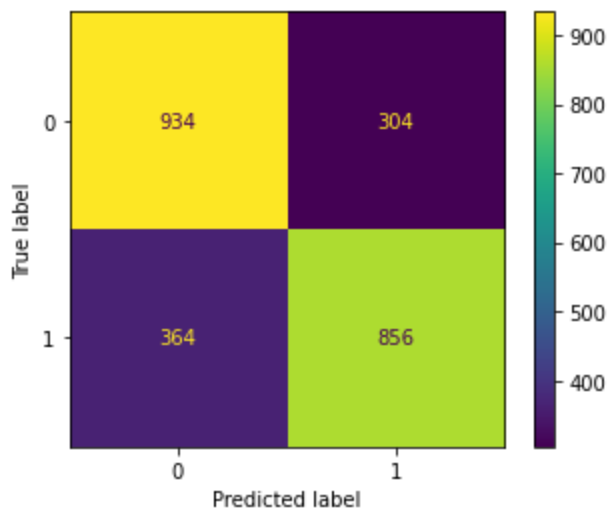
Результаты:

Accuracy: 0.7306753458096013
Precision: 0.7203791469194313
Recall: 0.7475409836065574



Результаты KNeighborsClassifier из sklearn

Accuracy: 0.7282343368592351
Precision: 0.7379310344827587
Recall: 0.7016393442622951



Метод Naive Bayes:

```
1. class NaiveBayes(BaseEstimator, ClassifierMixin):
2.     def __init__(self, bins):
3.         self.bins = bins
4.         pass
5.
6.     def fit(self, data, labels):
7.         self.data = data
8.         self.labels = labels
9.         self.classes = []
10.        for j in np.unique(labels):
11.
12.            self.classes.append([])
13.            for i in range(data.shape[1]):
14.                self.classes[j].append(*np.histogram(data[labels
15.                    == j, i], bins = self.bins))
16.
17.                self.classes[j][-1][0] = self.classes[j][-
18.                    1][0].astype('float64') / len(data[labels == j, i])
19.
20.            self.prclasses = np.unique(labels, return_counts =
21.                True)[1] / len(labels)
22.
23.    def predict(self, maindata):
```

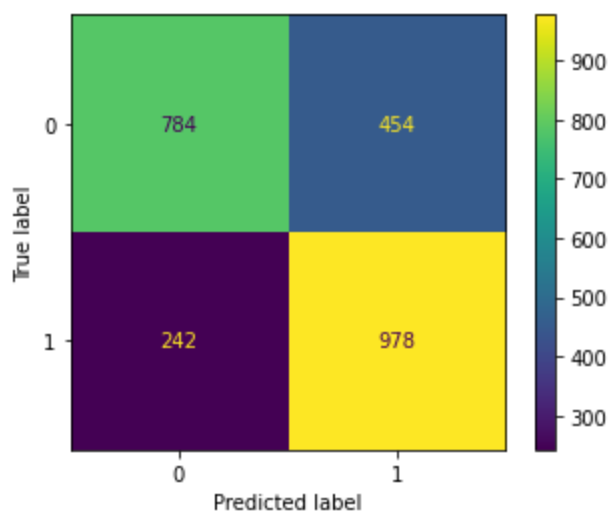
```

20.         res = np.ndarray((maindata.shape[0],))
21.         for j, data in enumerate(maindata):
22.             maximum = 0
23.             ans = 0
24.             for i in range(len(self.classes)):
25.                 p = self.prclasses[i]
26.                 for k in range(len(self.classes[i])):
27.                     ind = np.digitize(data[k],
28.                                     self.classes[i][k][1])
29.                     if ind >= len(self.classes[i][k][1]) or ind <=
30.                        0:
31.                         p = 0
32.                     else:
33.                         p *= self.classes[i][k][0][ind - 1]
34.                     if p > maximum:
35.                         maximum = p
36.                         ans = i
37.             res[j] = ans
38.         return res

```

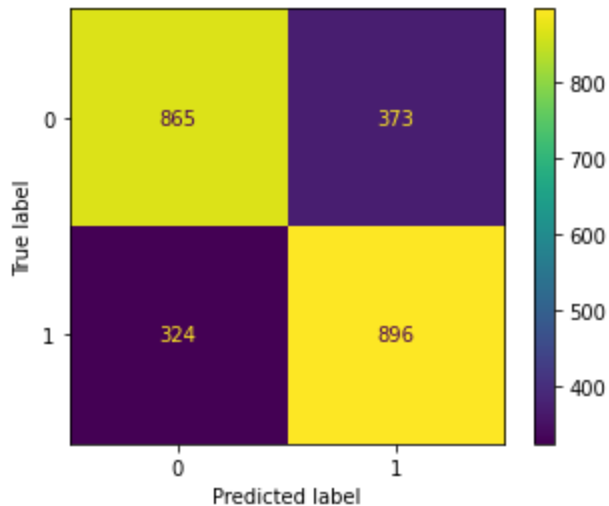
Результаты:

Accuracy: 0.7168429617575265
Precision: 0.6829608938547486
Recall: 0.8016393442622951



Результаты для реализации из sklearn:

Accuracy: 0.7164361269324654
Precision: 0.7060677698975572
Recall: 0.7344262295081967



Последний метод – SVM имеет 4 параметра, используем поиск по сетке:

```
1. class SVM(BaseEstimator, ClassifierMixin):
2.     def __init__(self, lr=0.1, batch=10, epochs=1, alpha=0.0001):
3.         self.lr = lr
4.         self.batch = batch
5.         self.epochs = epochs
6.         self.alpha = alpha
7.
8.     def fit(self, data, labels):
9.         self.w = np.random.normal(0, 1, (data.shape[1]+1,))
10.        data = np.concatenate((data, np.ones((data.shape[0],1))),
11.                               axis=1)
12.        labels = labels * 2 - 1
13.        for _ in range(self.epochs):
14.            for i in range(self.batch, len(data), self.batch):
15.                data_batch = data[i-self.batch:i]
16.                labels_batch = labels[i-self.batch:i]
17.
18.                grad = 2 * self.alpha * self.w
19.                for i, x in enumerate(data_batch):
20.                    if 1 - x.dot(self.w) * labels_batch[i] > 0:
21.                        grad -= x * labels_batch[i]
```

```

21.
22.             self.w -= self.lr * grad
23.         return self
24.
25.     def predict(self, data):
26.         return (np.sign(np.concatenate((data,
np.ones((data.shape[0],1))), axis=1).dot(self.w)) + 1) / 2

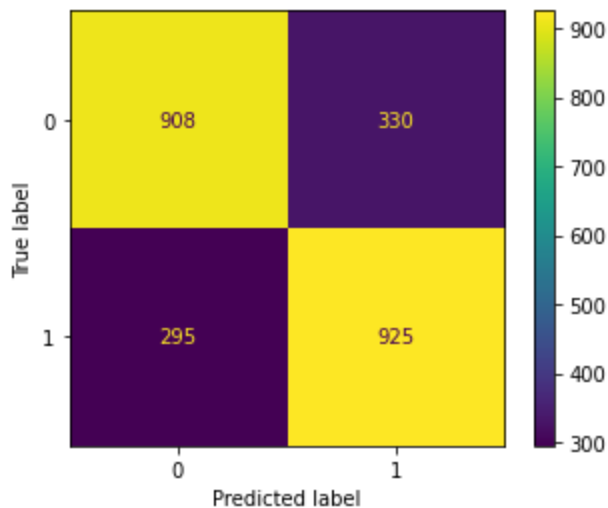
```

Результаты:

```

Accuracy: 0.7457282343368592
Precision: 0.7370517928286853
Recall: 0.7581967213114754

```



Выводы

Как видно из результатов, итоги не самые плохие – от 70 до 80%, однако и не самые хорошие. Лучше всего себя показал метод SVM, однако метрика Recall самая высокая – в районе 80% - у метода Naive Bayes. Реализации из sklearn оказались местами хуже, нежели написанные с нуля.

Можно сделать вывод, что данная задача не очень хорошо решается машинным обучением.