

**Московский авиационный институт
(национальный исследовательский университет)**

**Факультет информационных технологий и прикладной
математики**

Кафедра вычислительной математики и программирования

**Курсовой проект по курсу «Искусственный интеллект и глубокое
обучение»**

Студент: К. М. Воронов
Преподаватель: Б.В. Вишняков
Группа: М8О-407Б-19
Дата:
Оценка:
Подпись:

Москва, 2023

Задача

Выбрать задачу (классификации или регрессии), датасет и метрики качества. Выбранные данные необходимо проанализировать и визуализировать. Реализовать алгоритм (несколько алгоритмов), сравнить результаты работы с аналогом из sklearn.

1 Описание данных

1. Pregnancies - Number of times pregnant (Количество беременностей)
2. Glucose - Plasma glucose concentration a 2 hours in an oral glucose tolerance test (Концентрация глюкозы в плазме крови через 2 часа при пероральном тесте на толерантность к глюкозе)
3. BloodPressure - Diastolic blood pressure (mm Hg) (Диастолическое артериальное давление)
4. SkinThickness - Triceps skin fold thickness (mm) (Толщина кожной складки трицепса)
5. Insulin - 2-Hour serum insulin (mu U/ml) (2-Часовой сывороточный инсулин)
6. BMI - Body mass index (weight in kg/(height in m)²) (Индекс массы тела)
7. DiabetesPedigree - Diabetes pedigree function (Функция родословной диабета)
8. Age (Возраст)
9. Outcome - Class variable (Наличие диабета)

2 Описание алгоритмов

В данной работе реализованы следующие алгоритмы обучения:

- 1) k-Nearest Neighbors (KNN)

Идея заключается в определении класса объекта по классам k ближайших (каких больше - такой и класс). В качестве расстояния используется евклидова метрика.

- 2) Naïve Bayes Пусть у нас есть объект B , который описывается признаками b_i и классы A_k . Алгоритм основан на формуле Байеса:

$$P(A|B) = \frac{P(B|A) * P(A)}{P(B)},$$

Сделав предположение, что все признаки независимы между собой, мы можем разложить условную вероятность объекта при условии класса на произведение вероятностей наличия признаков при условии класса:

$$P(B|A_k) = \prod_{i=1}^n P(b_i|A_k)$$

$P(b_i|A)$ распределение каждого признака в зависимости от класса. Есть два варианта - предполагаем, что признаки по классу имеют нормальное распределение или строим гистограмму.

Так как все объекты разные, то $P(B)$ является одинаковой константой для всех B , поэтому ей можно пренебречь.

В результате берется класс, для которого $P(A|B)$ больше всего.

3 Предобработка данных

Для начала я проверил данные на присутствие в них дубликатов и пустых ячеек. Далее я проанализировал распределение данных. Оказалось, что в столбцах Glucose, BloodPressure, SkinThickness и BMI присутствуют выбросы в виде нулей, что невозможно для этих признаков. Из-за того, что данных не так много, я заполнил их средними значениями для каждого из исходов (есть диабет или нет). Также, так как здоровых было больше, чем больных, я сделал оверсепмлиг, путём копирования данных.

4 Реализация алгоритмов

Реализация KNN

```

1 from sklearn.base import BaseEstimator, ClassifierMixin
2
3 class KNN(BaseEstimator, ClassifierMixin):
4     def __init__(self, k):
5         self.k = k
6
7     def fit(self, data, labels):
8         self.data = data
9         self.labels = labels
10
11     def euclidean_distance(self, row1, row2):
12         distance = 0
13         for i in range(len(row1)):
14             distance += (row1[i] - row2[i]) ** 2
15         return math.sqrt(distance)
16
17     def predict(self, maindata):
18         res = np.ndarray((maindata.shape[0],))
19         for j, data in enumerate(maindata):
20             distances = []

```

```

21         for i, row in enumerate(self.data):
22             distances.append((self.euclidean_distance(data, row), self.labels[i]))
23         distances.sort(key = lambda tup: tup[0])
24         dictionary = collections.defaultdict(int)
25         for i in range(self.k):
26             dictionary[distances[i][1]] += 1
27         res[j] = max(dictionary.items(), key = lambda tup: tup[1])[0]
28     return res

```

Naive Bayes реализован двумя способами:

1) По гистограммам

2) Делая предположение, что данные имеют нормальное распределение

Первая реализация выглядит следующим образом:

```

1  from sklearn.base import BaseEstimator, ClassifierMixin
2
3  class NaiveBayes(BaseEstimator, ClassifierMixin):
4      def __init__(self, bins):
5          self.bins = bins
6          pass
7
8      def fit(self, data, labels):
9          self.data = data
10         self.labels = labels
11         self.classes = []
12         for j in np.unique(labels):
13
14             self.classes.append([])
15             for i in range(data.shape[1]):
16                 self.classes[j].append(*np.histogram(data[labels == j, i], bins = self.
17                     bins))
18                 self.classes[j][-1][0] = self.classes[j][-1][0].astype('float64') / len(
19                     data[labels == j, i])
20
21         self.prclasses = np.unique(labels, return_counts = True)[1] / len(labels)
22
23     def predict(self, maindata):
24         res = np.ndarray((maindata.shape[0],))
25         for j, data in enumerate(maindata):
26             maximum = 0
27             ans = 0
28             for i in range(len(self.classes)):
29                 p = self.prclasses[i]
30                 for k in range(len(self.classes[i])):
31                     ind = np.digitize(data[k], self.classes[i][k][1])
32
33                     if ind >= len(self.classes[i][k][1]) or ind <= 0:
34                         p = 0
35                     else:

```

```

34         p *= self.classes[i][k][0][ind - 1]
35
36         if p > maximum:
37             maximum = p
38             ans = i
39         res[j] = ans
40     return res

```

Вторая:

```

1  from sklearn.base import BaseEstimator, ClassifierMixin
2
3  class GaussianNaiveBayes(BaseEstimator, ClassifierMixin):
4      def __init__(self):
5          pass
6
7      def fit(self, data, labels):
8          self.data = data
9          self.labels = labels
10         self.mathexp = []
11         self.variance = []
12         self.classes = []
13         self.prob = []
14
15         for j in zip(*np.unique(labels, return_counts=True)):
16             self.prob.append(j[1] / len(labels))
17             self.classes.append(j[0])
18             self.mathexp.append(data[labels == j[0],].mean(axis = 0))
19             self.variance.append(data[labels == j[0],].var(axis = 0))
20
21     def predict(self, maindata):
22         res = np.ndarray((maindata.shape[0],))
23         for j, data in enumerate(maindata):
24             maximum = 0
25             ans = 0
26             for i in range(len(self.classes)):
27                 t = np.exp((-1/2) * ((data - self.mathexp[i]) ** 2) / (2 * self.variance
28                                     [i])) / np.sqrt(2 * np.pi * self.variance[i])
29                 t = np.cumprod(t)
30                 t[-1] *= self.prob[i]
31                 if t[-1] > maximum:
32                     maximum = t[-1]
33                     ans = self.classes[i]
34             res[j] = ans
35         return res

```

5 Метрики качества

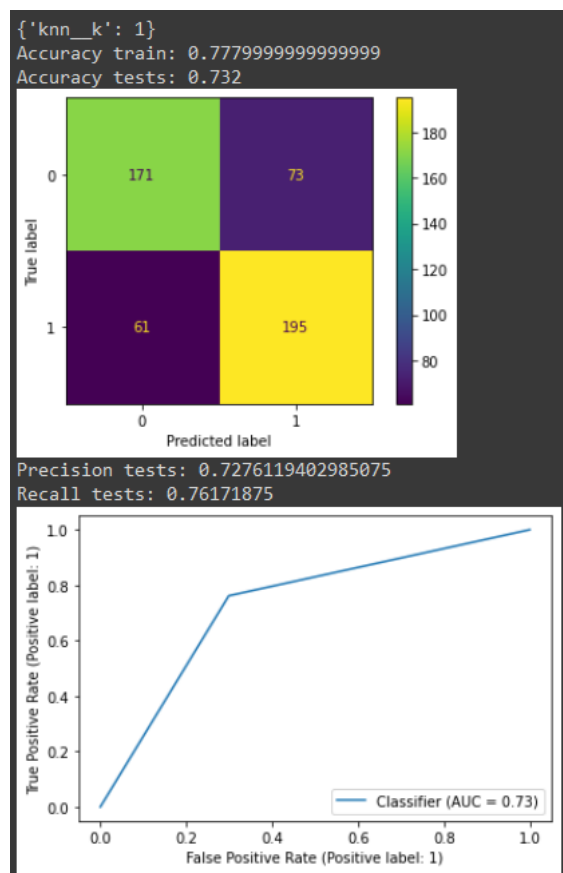
В качестве метрик качества я использовал:

1. accuracy - доля правильно предсказанных объектов
2. precision - доля правильно предсказанных положительных объектов среди всех объектов, предсказанных положительным классом
3. recall - доля правильно найденных положительных объектов среди всех объектов положительного класса.

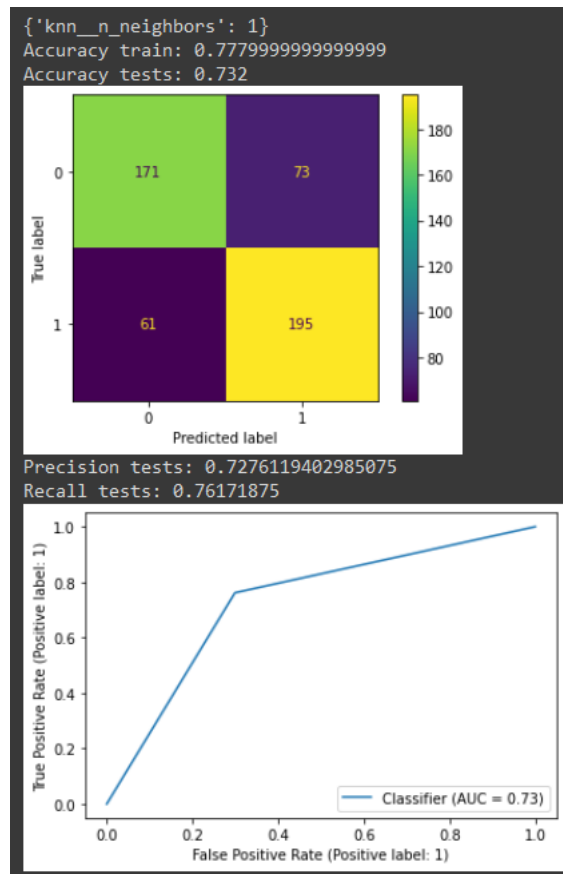
Использование всех этих метрик показывает более полную картину работы алгоритма.

6 Полученные результаты

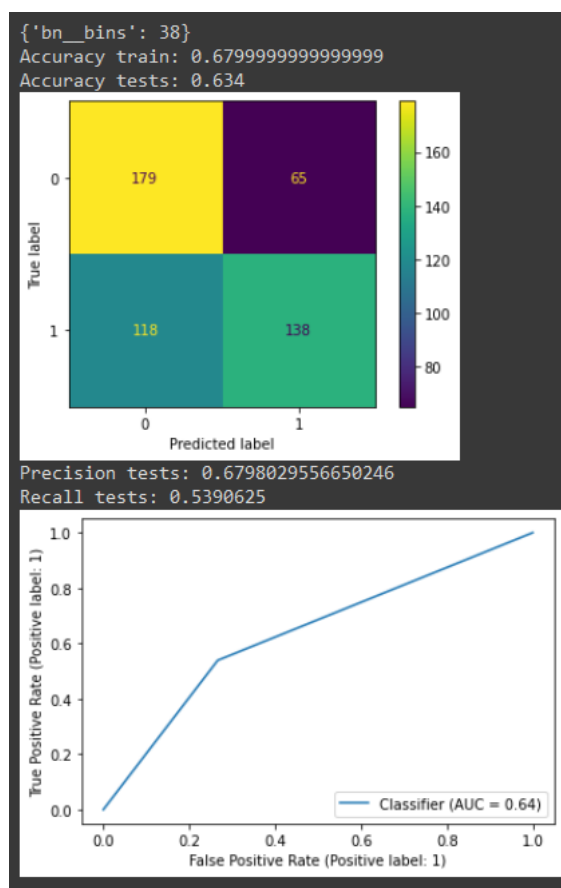
KNN показал себя неплохо, дав точность в районе 77 процентов. Моя реализация



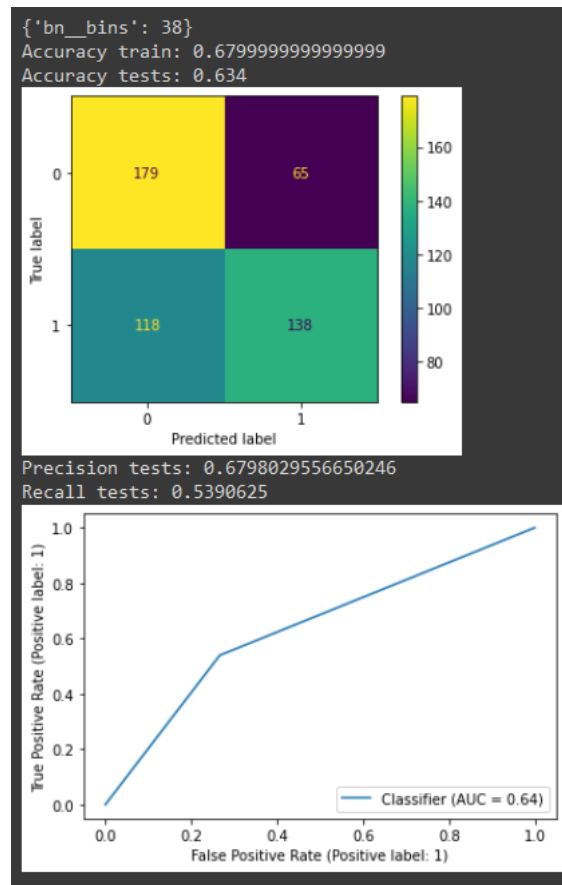
Реализация из sklearn дала точно такие же результаты



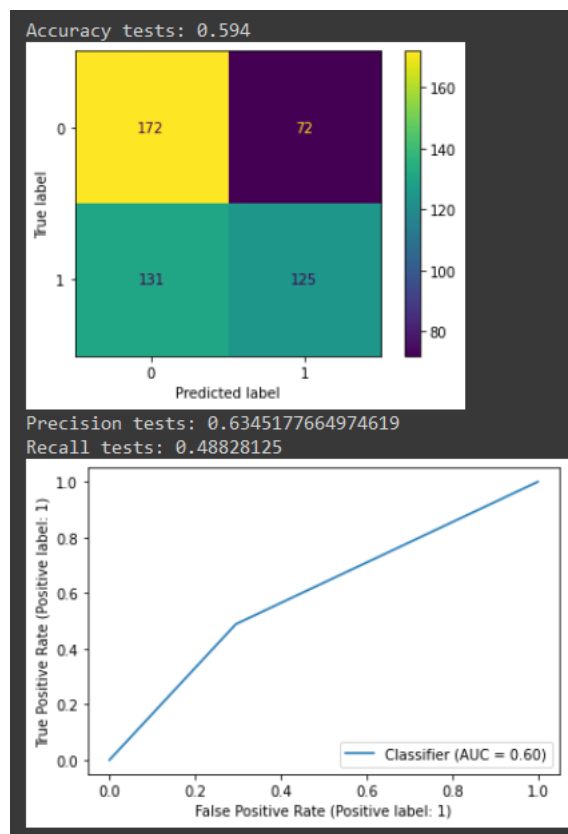
Байесовский классификатор показал себя хуже. Это связано с тем, что распределения признаков по классам почти совпадают. Однако байесовский классификатор с использованием гистограмм показал себя лучше, дав точность в районе 68 процентов.



Байесовский классификатор с нормальным распределением, моя реализация



Байесовский классификатор с нормальным распределением, sklearn



7 Выводы

Сделав данный курсовой проект, я еще раз потренировался обрабатывать данные, реализовал несколько алгоритмов, а именно KNN и Naive Bayes. К сожалению, байесовский классификатор показал себя не очень хорошо на этих данных, однако он может себя проявить в другой раз. Я думаю, что на этих данных могут хорошо работать дерево решений и случайный лес. Линейная модель здесь также не даст хороших результатов. Решение таких задач, как предсказывание болезней, может спасти чью-то жизнь, потому что очень часто на ранних этапах болезнь можно вылечить.