

# Министерство науки и высшего образования Российской Федерации Федеральное государственное бюджетное образовательное учреждение высшего образования

## «Московский государственный технический университет имени Н.Э. Баумана

(национальный исследовательский университет)» (МГТУ им. Н.Э. Баумана)

ФАКУЛЬТЕТ «Информатика и системы управления» (ИУ)

КАФЕДРА «Информационная безопасность» (ИУ8)

## ОТЧЕТ ПО ЛАБОРАТОРНОЙ РАБОТЕ №3 (3.1) МАШИННОЕ ОБУЧЕНИЕ

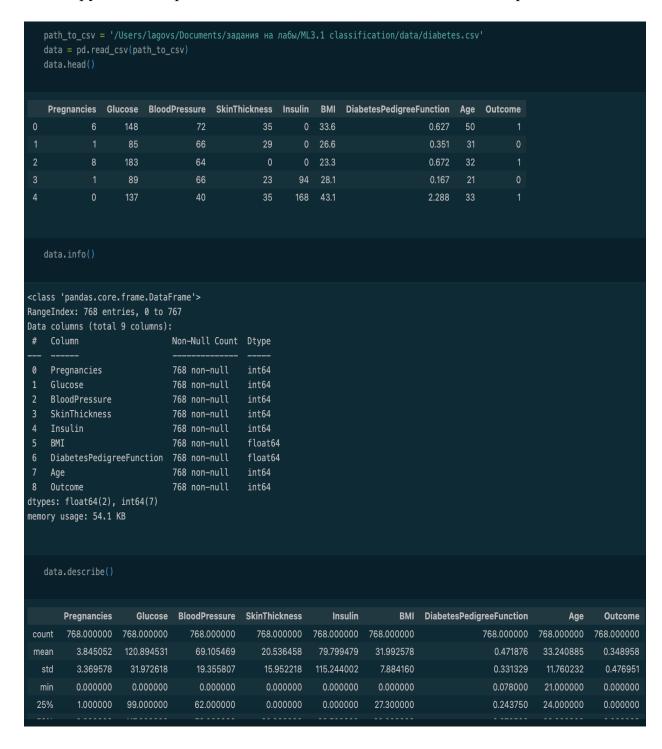
Студент	ИУ8-92	Лагов С. П.	
	(Группа)	(И.О.Фамилия)	
Преподаватель:		Коннова Н.С.	
		(И.О. Фамилия)	

### Цель работы:

Познакомиться с основными приемами работы с моделями классификации в scikit-learn.

Ход работы:

1. Загрузите встроенный датасет о диагностике сахарного диабета.



2. Постройте модель классификации для предсказания наличия заболевания. Оцените качество построенной модели с помощью отчета о классификации и матрицы классификации.

Доп. задание - Напишите функцию, которая автоматически обучает все перечисленные модели и для каждой выдает оценку точности.

```
def predict(model, x_train, y_train, x_test, y_test, method, show_coeffs=False):
    model.fit(x_train, y_train)
    if show_coeffs:
        print("Coefficients: \n", model.coef_)
        _ = [print(k, v) for k, v in zip(x_train.columns, model.coef_[0])]
        print("Intercept: \n", model.intercept_)

        y_pred = model.predict(x_test)

        score = metrics.accuracy_score(y_test, y_pred)
        print(f"Score: {score}")

        print(f'Confusion Matrix:\n{confusion_matrix(y_test, y_pred)}')
        print(f'Classification Report:\n{classification_report(y_test, y_pred)}')

        global accuracies
        accuracies[method] = score
```

```
table = PrettyTable()
table.field_names = ["Method", "Precision"] # метод / оценка точности
for method in accuracies.keys():
    table.add_row([method, round(accuracies[method],8)])
print(table)
```

```
y = data.Outcome
   X = data.drop(["Outcome"], axis=1)
   y.shape, X.shape
((768,), (768, 8))
   x_train, x_test, y_train, y_test = train_test_split(X, y, random_state=104, test_size=0.25, shuffle=True)
   predict(LogisticRegression(), x_train, y_train, x_test, y_test, "Logistic Regression", show_coeffs=False)
Score: 0.78645833333333333
Confusion Matrix:
[[112 16]
[ 25 39]]
Classification Report:
             precision recall f1-score support
                  0.82
                           0.88
                                     0.85
                                                128
                  0.71
                           0.61
                                     0.66
                                                64
                                     0.79
   accuracy
                                                192
  macro avg
                  0.76
                           0.74
                                     0.75
                                               192
weighted avg
                  0.78
                           0.79
                                     0.78
                                                192
```

4. Постройте альтернативную полиномиальную модель, сравните ее с предыдущей.

```
poly = PolynomialFeatures(2)
  x_poly_train = poly.fit_transform(x_train)
  x_poly_test = poly.fit_transform(x_test)
  predict(LogisticRegression(), x_poly_train, y_train, x_poly_test, y_test, "Polynomial Regression", show_coeffs=False)
Score: 0.69270833333333334
Confusion Matrix:
[[100 28]
[ 31 33]]
Classification Report:
            precision recall f1-score support
                0.76 0.78 0.77
         0
                                              128
                0.54 0.52
                                   0.53
                                   0.69
                                              192
   accuracy
                0.65
                          0.65
                                   0.65
                                              192
  macro avg
                 0.69
                          0.69
                                   0.69
                                              192
 eighted avg
```

Вывод: полиномиальная модель более точно предсказывает значения, чем линейная, это видно по значением параметра score.

5. Попробуйте применить к той же задаче другие модели регрессии. Для каждой из них выведите матрицу классификации и оценку точности.

Использовались следующие модели:

- 1. SVM (Метод опорных векторов) (RBF Radial Basis Function)
- 2. SVM (Метод опорных векторов) (linear)
- 3. SVM (Метод опорных векторов) (poly)
- 4. Метод ближайших соседей
- 5. Дерево решений
- 6. Случайный лес
- 7. Многослойный перцерптрон

```
Score: 0.80208333333333334
Confusion Matrix:
[ 23 41]]
Classification Report:
                         recall f1-score support
             precision
                  0.83
                            0.88
                                      0.86
                                                128
                  0.73
                            0.64
                                      0.68
                                      0.80
   accuracy
                  0.78
                            0.76
                                      0.77
                                                192
   macro avg
weighted avg
                  0.80
                            0.80
                                      0.80
```

```
Score: 0.7760416666666666
Confusion Matrix:
[[117 11]
[ 32 32]]
Classification Report:
             precision recall f1-score support
                  0.79
                            0.91
                                      0.84
                                                 128
                                      0.60
                  0.74
                            0.50
                                                  64
   accuracy
                                      0.78
                            0.71
  macro avg
                  0.76
                                      0.72
weighted avg
                  0.77
                            0.78
                                      0.76
                                                 192
```

```
predict(\textbf{KN} eighbors \texttt{Classifier()}, \cdot \textbf{x\_train,} \cdot \textbf{y\_train,} \cdot \textbf{x\_test,} \cdot \textbf{y\_test,} \cdot \texttt{"KN} eighbors \texttt{Classifier",} \cdot \textbf{show\_coeffs=False})
Score: 0.72395833333333334
Confusion Matrix:
[[103 25]
 [ 28 36]]
Classification Report:
                   precision
                                     recall f1-score support
                          0.79
                                        0.80
                                                       0.80
                                                                       128
               0
                          0.59
                                        0.56
                                                       0.58
                                                                       64
                                                       0.72
     accuracy
```

192

0.68

0.72

0.69

0.72

0.69

0.72

macro avg

weighted avg

```
Confusion Matrix:
[[95 33]
[31 33]]
Classification Report:
            precision
                       recall f1-score
                                         support
                0.75
                         0.74
         0
                                  0.75
                                            128
                         0.52
                0.50
                                  0.51
                                             64
                                  0.67
                                            192
   accuracy
                0.63
                         0.63
                                  0.63
                                            192
  macro avg
weighted avg
                0.67
                          0.67
                                  0.67
                                            192
```

```
predict(MLPClassifier(), x_train, y_train, x_test, y_test, "MLP (Многослойный перцерптрон)", show_coeffs=False)
Score: 0.72395833333333334
Confusion Matrix:
[[109 19]
Classification Report:
             precision
                         recall f1-score support
                  0.76
                            0.85
                                      0.80
                  0.61
                            0.47
                                      0.53
                                                 64
                                      0.72
   accuracy
                  0.69
                            0.66
                                      0.67
  macro avg
weighted avg
                  0.71
                            0.72
                                      0.71
```

Таблица со значениями коэффициента детерминации:

+   Method	++   Precision
Logistic Regression Polynomial Regression	0.78645833     0.69270833
SVM (Метод опорных векторов) (RBF — Radial Basis Function	n)   0.77604167
SVM (Метод опорных векторов) (linear)   SVM (Метод опорных векторов) (poly)	0.80208333     0.77604167
KNeighborsClassifier   Decision Tree	0.72395833     0.66666667
Random Forest   MLP (Многослойный перцерптрон)	0.74479167     0.72395833
+	+

Из всех моделей лучше всего себя показала модель по методу опорных векторов (linear)

#### Ответы на контрольные вопросы

1. <u>Чем отличается применение разных моделей классификации в</u> бибилиотеке sklearn ?

Ответ: библиотеке представлены В различные модели классификации, такие как логистическая регрессия, деревья решений и ансамблевые методы, каждая из которых имеет свои особенности и подходит для разных типов задач и данных. Выбор модели зависит от характера данных, требуемой интерпретируемости И производительности, поэтому рекомендуется тестировать несколько моделей для нахождения наилучшего решения

2. Что показывает метрика точности регрессии? Ответ: метрика точности регрессии показывает, насколько предсказанные значения модели близки к реальным значениям целевой переменной. Она оценивает способность модели объяснять вариацию данных и минимизировать ошибку предсказания. В зависимости от используемой метрики (MSE, MAE, R<sup>2</sup>) точность может выражаться в виде среднего отклонения, суммы квадратов ошибок или доли объяснённой дисперсии

3. Какое значение имеют коэффициенты логистической регрессии?

Ответ: коэффициенты логистической регрессии отражают влияние каждого признака на вероятность отнесения объекта к определённому классу. Они показывают, насколько изменение значения признака на одну единицу изменяет логарифм шансов принадлежности к классу, при условии, что остальные признаки остаются неизменными. Положительный коэффициент увеличивает шансы, отрицательный — уменьшает

4. Что показывает матрица классификации?

Ответ: матрица классификации показывает, как модель предсказывает классы для каждого наблюдения по сравнению с их истинными значениями. Она содержит информацию о количествах верных и неверных предсказаний для каждого класса, отображая, сколько объектов было правильно классифицировано и сколько ошибочно. Это позволяет оценить качество модели, выявить смещения и анализировать, какие классы модель путает между собой

5. Какие параметры имеет конструктор объекта логистической регрессии?

Ответ: основные penalty, который определяет регуляризации (L1, L2, ElasticNet или None), и C, который настраивает силу регуляризации. Также присутствует параметр solver. определяющий метод оптимизации (например, 'liblinear', 'lbfgs', 'newton-cg'). Параметр max iter задаёт максимальное количество итераций для сходимости алгоритма, a random state обеспечивает воспроизводимость результатов. Дополнительно можно настроить параметры fit intercept, чтобы определять, учитывать ли константу в модели, и class weight, чтобы учесть дисбаланс классов

6. <u>Какие атрибуты имеет объект логистической регрессии?</u>
Ответ: coef\_ содержит коэффициенты модели для каждого признака после обучения. intercept представляет свободный член

(смещение) модели. Атрибут classes\_ содержит список классов, которые модель может предсказывать. n\_iter\_ показывает количество итераций, выполненных алгоритмом для достижения сходимости. Также может присутствовать score

7. <u>Какие параметры и атрибуты имеют объекты других моделей машинного обучения библиотеки sklearn\_?</u>

Ответ: основные параметры включают гиперпараметры обучения, такие как регуляризация (C, penalty), метод оптимизации (solver), количество итераций (max\_iter) и вес классов (class\_weight). Атрибуты включают обученные коэффициенты (coef\_), смещение (intercept\_), список классов (classes ) и количество итераций (n iter)