

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

Кафедра ИУ5

Рубежный контроль №2

«Технологии машинного обучения»

Вариант: №3

Выполнил:Баркалова И.В.

Группа: ИУ5-63Б

Преподаватель: Гапанюк Ю. Е.

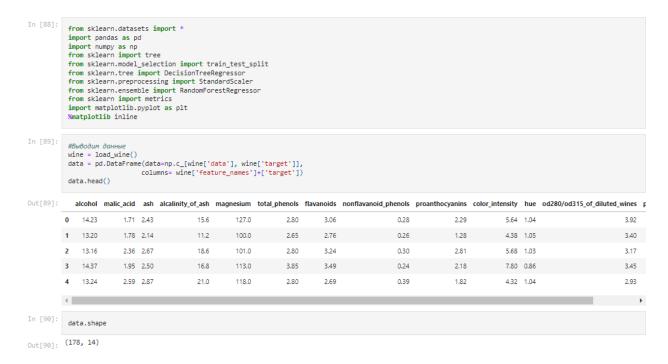
## Задание:

Для заданного набора данных (по Вашему варианту) постройте модели классификации или регрессии (в зависимости от конкретной задачи, рассматриваемой в наборе данных). Для построения моделей используйте методы 1 и 2 (по варианту для Вашей группы). Оцените качество моделей на основе подходящих метрик качества (не менее двух метрик). Какие метрики качества Вы использовали и почему? Какие выводы Вы можете сделать о качестве построенных моделей? Для построения моделей необходимо выполнить требуемую предобработку данных: заполнение пропусков, кодирование категориальных признаков, и т.д.

Группа	Метод №1	Метод №2
ИУ5-63б	Дерево решений	Случайный лес

Набор данных по варианту: <a href="https://scikit-learn.org/stable/modules/generated/sklearn.datasets.load\_wine.html#sklearn.datasets.load\_wine">https://scikit-learn.org/stable/modules/generated/sklearn.datasets.load\_wine.html#sklearn.datasets.load\_wine</a>

## Текст программы и результат ее выполнения:



```
In [91]: data.describe()
                              alcohol malic acid
                                                                   ash alcalinity of ash magnesium total phenols flavanoids nonflavanoid phenols proanthocyanins color intensity
                                                                                                                                                                                                                                         hue od280/od3
                count 178,000000 178,000000 178,000000
                                                                                 178,000000 178,000000
                                                                                                                     178,000000 178,000000
                                                                                                                                                                    178,000000
                                                                                                                                                                                           178,000000
                                                                                                                                                                                                                178,000000 178,000000
                 mean 13.000618 2.336348 2.366517
                                                                               19.494944 99.741573
                                                                                                                       2.295112 2.029270
                                                                                                                                                                      0.361854
                                                                                                                                                                                             1.590899
                                                                                                                                                                                                              5.058090 0.957449
                   std
                           0.811827 1.117146 0.274344
                                                                                   3.339564 14.282484
                                                                                                                        0.625851
                                                                                                                                       0.998859
                                                                                                                                                                      0.124453
                                                                                                                                                                                              0.572359
                                                                                                                                                                                                                   2.318286
                                                                                                                                                                                                                                  0.228572
                   min 11.030000 0.740000 1.360000
                                                                                  10.600000 70.000000
                                                                                                                       0.980000 0.340000
                                                                                                                                                                      0.130000
                                                                                                                                                                                              0,410000
                                                                                                                                                                                                                   1.280000 0.480000
                  25% 12.362500 1.602500 2.210000
                                                                                  17.200000 88.000000
                                                                                                                        1.742500
                                                                                                                                       1.205000
                                                                                                                                                                      0.270000
                                                                                                                                                                                               1.250000
                                                                                                                                                                                                                   3,220000
                                                                                                                                                                                                                                  0.782500
                  50% 13.050000 1.865000 2.360000 19.500000 98.000000
                                                                                                                       2.355000 2.135000
                                                                                                                                                                      0.340000
                                                                                                                                                                                               1.555000
                                                                                                                                                                                                                   4.690000 0.965000
                  75% 13.677500 3.082500 2.557500
                                                                                  21.500000 107.000000
                                                                                                                        2.800000 2.875000
                                                                                                                                                                       0.437500
                                                                                                                                                                                               1.950000
                                                                                                                                                                                                                   6.200000
                                                                                                                                                                                                                                  1.120000
                  max 14.830000 5.800000 3.230000 30.000000 162.000000 3.880000 5.080000
                                                                                                                                                                       0.660000
                                                                                                                                                                                               3.580000
                                                                                                                                                                                                                 13.000000 1.710000
                4
 In [92]: #Разделение признако
                  x = data.drop('target', axis=1)
                 v = data['target']
                 model=tree.DecisionTreeClassifier(criterion="entropy")
 Out[93]: DecisionTreeClassifier(criterion='entropy')
               In a Jupyter environment, please rerun this cell to show the HTML representation or trust the notebook.
               On GitHub, the HTML representation is unable to render, please try loading this page with nbviewer.org.
 In [94]:
                  #Оценка модели(Оценка равна 1 в случае очень маленького набора данных)
                 model.score(x,y)
Out[94]: 1.0
In [95]:
                 tree.plot_tree(model)
Out[95]: [Text(0.4230769230769231, 0.9, 'X[6] <= 1.575\nentropy = 1.567\nsamples = 178\nvalue = [59, 71, 48]'),
                [Text(0.4230769230769231, 0.9, 'X[6] <= 1.575\nentropy = 1.567\nsamples = 178\nvalue = [59, 71, 48]'),
Text(0.15384615384615385, 0.7, 'X[9] <= 3.825\nentropy = 0.771\nsamples = 62\nvalue = [0, 14, 48]'),
Text(0.0592307692307693, 0.5, 'entropy = 0.0\nsamples = 13\nvalue = [0, 13, 0]'),
Text(0.23076923076923078, 0.5, 'X[3] <= 17.15\nentropy = 0.144\nsamples = 49\nvalue = [0, 1, 48]'),
Text(0.15384615384615385, 0.3, 'entropy = 0.0\nsamples = 1\nvalue = [0, 1, 0]'),
Text(0.059623076923076923, 0.3, 'entropy = 0.0\nsamples = 48\nvalue = [0, 1, 0]'),
Text(0.6923076923076923, 0.7, 'X[12] <= 724.5\nentropy = 1.0\nsamples = 116\nvalue = [59, 57, 0]'),
Text(0.5384615384615384, 0.5, 'X[0] <= 13.175\nentropy = 0.133\nsamples = 54\nvalue = [1, 53, 0]'),
Text(0.6153846153846154, 0.3, 'entropy = 0.0\nsamples = 50\nvalue = [0, 50, 0]'),
Text(0.6153846153846154, 0.3, 'entropy = 0.0\nsamples = 3\nvalue = [0, 3, 0]'),
Text(0.5384615384615384, 0.1, 'entropy = 0.0\nsamples = 3\nvalue = [1, 0, 0]'),
Text(0.6303076923076923, 0.1, 'entropy = 0.0\nsamples = 1\nvalue = [1, 0, 0]'),
Text(0.8461538461538461, 0.5, 'X[9] <= 3.46\nentropy = 0.345\nsamples = 62\nvalue = [58, 4, 0]'),
Text(0.7692307692307693, 0.3, 'entropy = 0.0\nsamples = 58\nvalue = [0, 4, 0]'),
Text(0.923076923076931, 0.3, 'entropy = 0.0\nsamples = 58\nvalue = [58, 0, 0]'))
In [96]: data.isnull().sum() #В нашем наборе нет пропусков
Out[96]: alcohol
               malic_acid
               ash
               alcalinity_of_ash
               magnesium
total_phenols
               flavanoids
                nonflavanoid_phenols
               proanthocyanins
               color_intensity
               od280/od315_of_diluted_wines
               proline
                target
               dtype: int64
                 x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=0.2, random_state=0)
In [98]:
                regressor = DecisionTreeRegressor()
regressor.fit(x_train, y_train)
```

```
DucisionTreeRegressor()
In a Jupyter environment, please rerun this cell to show the HTML representation or trust the notebook.

On GitHub, the HTML representation is unable to render, please try loading this page with nbviewer.org.

In [99]:

y_pred = regressor.predict(x_test)
print('Mepea Naboulute Errors', metrics.mean_absolute_error(y_test, y_pred))
print('Mean Squared Errors', metrics.mean_squared_error(y_test, y_pred)))

Appens opewenHuh
Men Absolute Errors: 0.0833333333333333

Mean Squared Errors: 0.083333333333333333

Root Mean Squared Errors: 0.08333333333333333

Root Mean Squared Errors: 0.28867513459481287

In [100—

x_train, x_test, y_train, y_test = train_test_split(data.drop(['target'], axis=1), data['target'], test_size=0.5, random_state=17)

In [101—

# Feature Scaling
sc = StandardScaler()
x_train = sc.fit_transform(x_train)
x_test = sc.transform(x_test)

In [102—
regressor = RandomForestRegressor(n_estimators=20, random_state=0)
regressor.fit(x_train, y_train)
y_pred = regressor.predict(x_test)
```

Общая точность модели случайного леса больше, чем точность дерева решний.

print('Mean Absolute Error:', metrics.mean\_absolute\_error(y\_test, y\_pred))
print('Mean Squared Error:', metrics.mean\_squared\_error(y\_test, y\_pred))
print('Root Mean Squared Error:', np.sqrt(metrics.mean\_squared\_error(y\_test, y\_pred)))

in [103... print('Случайный лес')

Mean Absolute Error: 0.11797752808988764 Mean Squared Error: 0.04741573033707866 Root Mean Squared Error: 0.2177515334896144