## Рубежный контроль №2¶

Работа Карягин А.Д., группа ИУ5-62Б, вариант 8 (задача №2)

Задача 1

In [55]:

import numpy as np

Кластеризуйте данные с помощью двух алгоритмов кластеризации. Алгоритмы для студентов группы ИУ5-62Б: MeanShift и иерархическая кластеризация. Сравните качество кластеризации с помощью следующих метрик качества кластеризации (если это возможно для Вашего набора данных):

- 1. Adjusted Rand index 2. Adjusted Mutual Information
- 3. Homogeneity, completeness, V-measure
- 4. Коэффициент силуэта

Сделате выводы о том, какой алгоритм осуществляет более качественную кластеризацию на Вашем наборе данных.

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

```
import pandas as pd
from typing import Dict, Tuple
from scipy import stats
from IPython.display import Image
from sklearn import cluster, datasets, mixture
from sklearn.neighbors import kneighbors_graph
from sklearn.preprocessing import StandardScaler
from sklearn.metrics import adjusted_rand_score
from sklearn.metrics import adjusted_mutual_info_score
from sklearn.metrics import homogeneity_completeness_v_measure
from sklearn.metrics import silhouette_score
from sklearn.cluster import MeanShift, AgglomerativeClustering
from itertools import cycle, islice
import seaborn as sns
```

import matplotlib.pyplot as plt %matplotlib inline

sns.set(style="ticks") from sklearn.datasets import load\_breast\_cancer

In [37]:

dataset = load\_breast\_cancer() for x in dataset: print(x)

data target target\_names DESCR

feature\_names filename In [38]:

# Признаки print(dataset.feature\_names) ['mean radius' 'mean texture' 'mean perimeter' 'mean area' 'mean smoothness' 'mean compactness' 'mean concavity'

'mean concave points' 'mean symmetry' 'mean fractal dimension' 'radius error' 'texture error' 'perimeter error' 'area error' 'smoothness error' 'compactness error' 'concavity error' 'concave points error' 'symmetry error' 'fractal dimension error' 'worst radius' 'worst texture' 'worst perimeter' 'worst area' 'worst smoothness' 'worst compactness' 'worst concavity' 'worst concave points' 'worst symmetry' 'worst fractal dimension'] In [39]:

1 1 1 1 1 1 0 1 0 1 1 1 1 1 1 1 1 0 0 1 0 1 0 1 1 1 1 1 1 0 0 1 0 0 1 0 1 0 1 1 1 1 1 1 0 1 1 0 1 0 1 0 1

1 1 1 1 1 1 1 0 0 0 0 0 0 1]

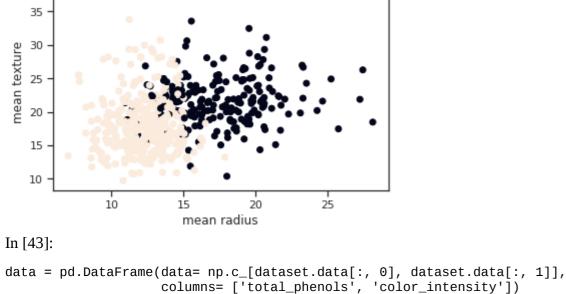
print(dataset.target)

In [40]: # Имена меток print(dataset.target\_names) ['malignant' 'benign']

In [41]: # Разделение набора данных x\_axis = dataset.data[:, 0] y\_axis = dataset.data[:, 1]

In [42]: # Построение plt.xlabel(dataset.feature\_names[0]) plt.ylabel(dataset.feature\_names[1])

plt.scatter(x\_axis, y\_axis, c=dataset.target) plt.show() 35



data.head()

Out[44]: total\_phenols color\_intensity **0** 17.99 10.38

1	20.57	17.77
2	19.69	21.25
3	11.42	20.38
4	20.29	14.34
In [45]:		
data.shape		

**MeanShift**¶

Out[45]:

In [44]:

In [46]:

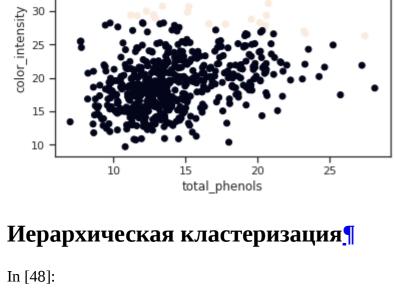
plt.show()

plt.show()

In [50]:

In [47]: plt.xlabel('total\_phenols')
plt.ylabel('color\_intensity')

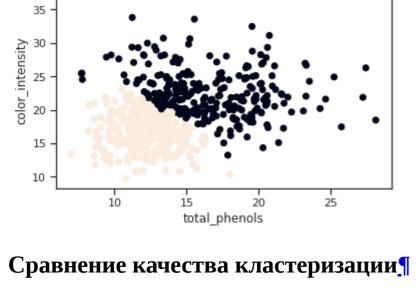
result\_MeanShift = MeanShift().fit\_predict(data)



## result\_AgglomerativeClustering = AgglomerativeClustering(n\_clusters=2).fit\_predict(data)

In [49]: plt.xlabel('total\_phenols') plt.ylabel('color\_intensity') plt.scatter(data['total\_phenols'], data['color\_intensity'], c=result\_AgglomerativeClustering)

plt.scatter(data['total\_phenols'], data['color\_intensity'], c=result\_MeanShift)



## def do\_clustering(cluster\_dataset, method): Выполнение кластеризации для данных примера

temp\_cluster = method.fit\_predict(cluster\_dataset) return temp\_cluster In [51]:

Вычисление метрик кластеризации

import warnings warnings.simplefilter(action='ignore', category=FutureWarning) def claster\_metrics(method, data, true\_y):

result\_Method = do\_clustering(data, method) list = []

list.append(adjusted\_rand\_score(true\_y, result\_Method))

list.append(adjusted\_mutual\_info\_score(true\_y, result\_Method)) h, c, v = homogeneity\_completeness\_v\_measure(true\_y, result\_Method) list.append(h) list.append(c) list.append(v) list.append(silhouette\_score(data, result\_Method)) names = ['ARI', 'AMI', 'Homogeneity', 'Completeness', 'V-measure', 'Silhouette'] for i in range(0,6): print('{}: {};'.format(names[i], list[i]))

In [52]: claster\_metrics(MeanShift(), data, dataset.target) ARI: 0.0259749909392261;

AMI: 0.01701444592127596; Homogeneity: 0.011959534897521055; Completeness: 0.04824827952634769;

V-measure: 0.019167843518729348; Silhouette: 0.46514827521556396; In [53]: claster\_metrics(AgglomerativeClustering(n\_clusters=2), data, dataset.target)

V-measure: 0.39442522838916433;

ARI: 0.48070586006925514; AMI: 0.3936323965676893; Homogeneity: 0.40262214017760883; Completeness: 0.3865554161475689;

Silhouette: 0.4127921205305984; Иерархическая кластеризация оказалась более качественной по сравнению с MeanShift.

In [ ]: