

## Домашнее задание 8

### 1. Можно ли отобрать наиболее значимые признаки с помощью PCA?

После того, как отработает метод главных компонент, матрица признаков переходит к новому базису. Причем процесс этот не обратим и обратно исходную матрицу не восстановить. Поэтому те компоненты-признаки, которые образуются после применения PCA, не имеют ничего общего с исходными признаками. Следовательно наиболее значимые признаки с помощью PCA подобрать нельзя.

### 2. (\*) Написать свою реализацию метода главных компонент с помощью сингулярного разложения с использованием функции [numpy.linalg.svd\(\)](https://docs.scipy.org/doc/numpy/reference/generated/numpy.linalg.svd.html) (<https://docs.scipy.org/doc/numpy/reference/generated/numpy.linalg.svd.html>)

In [9]:

```
1 import numpy as np
2 from sklearn import datasets
3 import matplotlib.pyplot as plt
```

In [10]:

```
1 # Загрузим ирисушечный датасет из sklearn
2 iris = datasets.load_iris()
3 X = iris.data
```

In [11]:

```

1 U, s, W = np.linalg.svd(X)
2 # Транспонируем матрицу W
3 V = W.T
4
5 # s - список диагональных элементов, его нужно привести к виду диагональной матрицы для
6 Sigma = np.zeros_like(X, dtype=float)
7 Sigma[np.diag_indices(min(X.shape))] = s
8 print(f'Матрица Sigma размер:{Sigma.shape} :\n{Sigma}')
9 print(f'Матрица U размер:{U.shape} :\n{U}')
10 print(f'Матрица V размер:{V.shape} :\n{V}')

```

Матрица Sigma размер:(150, 4) :

[illegible]

In [12]:

```
1 eig_sum = sum(s)
2 var_exp = [(i / eig_sum) * 100 for i in sorted(s, reverse=True)]
3 cum_var_exp = np.cumsum(var_exp)
4 print(f'Доля дисперсии, описываемая каждой из компонент \n{var_exp}')
5
6 # а теперь оценим кумулятивную (то есть накапливаемую) дисперсию при учетывании каждой
7 print(f'Кумулятивная доля дисперсии по компонентам \n{cum var exp}')
```

Доля дисперсии, описываемая каждой из компонент

[80.59340691495326, 14.916876798004958, 2.906715976729492, 1.583000310312297  
2]

Кумулятивная доля дисперсии по компонентам

[ 80.59340691 95.51028371 98.41699969 100. ]

Видим, что на первые два компонента приходится более 95% дисперсии. Занулим 3 и 4.

In [13]:

```
1 n_elements = 2
2 Sigma = Sigma[:, :n_elements]
3 Sigma
```

Out[13]:

[illegible]

In [14]:

```
1 V = V[:, :n_elements]
2 V
```

Out[14]:

```
array([[ -0.75110816,   0.2841749 ],
       [ -0.38008617,   0.5467445 ],
       [ -0.51300886,  -0.70866455 ],
       [ -0.16790754,  -0.34367081 ]])
```

In [15]:

```
1 # Создадим новую матрицу признаков с 4 признаками, но рангом 2
2 B = U.dot(Sigma.dot(V.T))
3 B
```

Out[15]:

```
array([[5.0952927 , 3.50597743, 1.40192232, 0.20165319],
       [4.74588049, 3.19610853, 1.46136967, 0.25800276],
       [4.68667405, 3.21586325, 1.30954904, 0.19452725],
       [4.61488457, 3.08894388, 1.46347879, 0.27002699],
       [5.07488651, 3.50623125, 1.36428119, 0.1863997 ],
       [5.52598407, 3.7330351 , 1.67566825, 0.28872322],
       [4.731593 , 3.2288014 , 1.36216771, 0.21446447],
       [5.00510918, 3.39830515, 1.47931372, 0.24418439],
       [4.37933538, 2.93134058, 1.38864652, 0.25618379],
       [4.80551481, 3.23360903, 1.48569239, 0.26393296],
       [5.39533378, 3.70766642, 1.49514864, 0.2183418 ],
       [4.89451945, 3.29088668, 1.51906398, 0.27146211],
       [4.67854319, 3.16443092, 1.41000708, 0.24052709],
       [4.30090163, 3.00174374, 1.08842179, 0.12739443],
       [5.73037625, 4.07476895, 1.2813685 , 0.09359359],
       [5.90310008, 4.12548386, 1.48153824, 0.16928199],
       [5.45240789, 3.81002345, 1.36951483, 0.15684908],
       [5.09813811, 3.49356935, 1.43489086, 0.21628344],
```

In [16]:

```
1 # Проведем сингулярное разложение матрицы B
2 U, s, W = np.linalg.svd(B)
3 # Транспонируем матрицу W
4 V = W.T
5 print("Элементы диагональной матрицы", s)
```

Элементы диагональной матрицы [9.59599139e+01 1.77610337e+01 2.07617549e-14  
2.38003498e-15]

Видим, что они практически равны 0. Создадим матрицу признаков C в новом двухмерном базисе через матрицу перехода

In [17]:

```
1 C = W.dot(B.T).T
2 C
```

Out[17]:

```
array([[ -5.91274714e+00,  2.30203322e+00,  1.79717352e-15,
        -1.11022302e-16],
       [ -5.57248242e+00,  1.97182599e+00,  1.86309301e-15,
        -5.55111512e-17],
       [ -5.44697714e+00,  2.09520636e+00,  1.15012166e-15,
        -1.11022302e-16],
       [ -5.43645948e+00,  1.87038151e+00,  2.02268757e-15,
        -1.11022302e-16],
       [ -5.87564494e+00,  2.32829018e+00,  2.05564732e-15,
         0.00000000e+00],
       [ -6.47759822e+00,  2.32464996e+00,  2.03656536e-15,
         1.66533454e-16],
       [ -5.51597520e+00,  2.07090423e+00,  1.42941214e-15,
        -2.77555756e-17],
       [ -5.85092859e+00,  2.14807482e+00,  1.60635394e-15,
        -2.77555756e-17],
       [ -5.15891972e+00,  1.77506408e+00,  1.48839274e-15,
         0.00000000e+00].
```

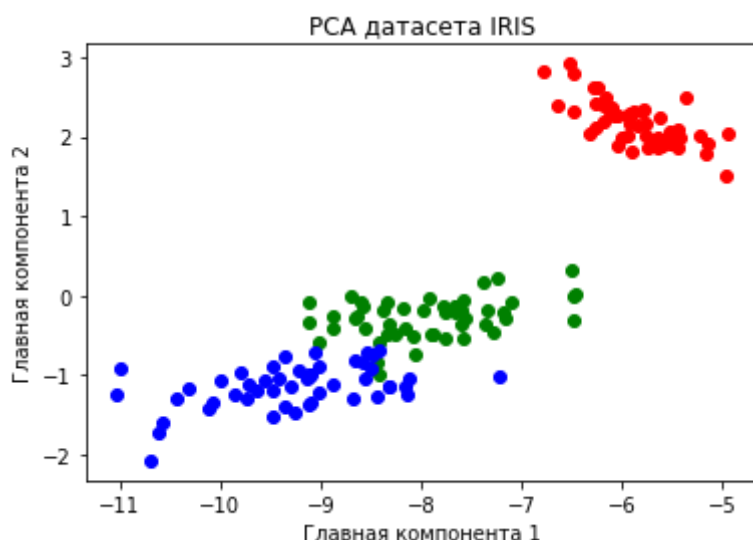
Видим, что последние два столбца признаков стали равны 0. Их можно удалить.

In [18]:

```
1 # Сформируем новую матрицу "объекты-признаки"
2 D = C[:, :2]
```

In [19]:

```
1 plt.figure()
2 y = iris.target
3 for c, i in zip("rgb", [0, 1, 2]):
4     plt.scatter(D[y==i, 0], D[y==i, 1], c=c)
5 plt.xlabel('Главная компонента 1')
6 plt.ylabel('Главная компонента 2')
7 plt.title('PCA датасета IRIS')
8 plt.show()
```





In [25]:

```
1 y_predict_train_PCA = model.predict(X_train_PCA)
2 y_predict_test_PCA = model.predict(X_test_PCA)
3 print("Точность на обучающей выборке", accuracy_score(y_train_PCA, y_predict_train_PCA))
4 print("Точность на тестовой выборке", accuracy_score(y_test_PCA, y_predict_test_PCA))
```

Точность на обучающей выборке 0.9714285714285714

Точность на тестовой выборке 1.0

**Вывод: Результаты работы модели после снижения размерности стали лучше**