

In [1]:

```
1 import numpy as np
2 import scipy as sp
3 from sklearn.decomposition import PCA
4
5 import seaborn as sns
6 import matplotlib.pyplot as plt
7 from mpl_toolkits.mplot3d import Axes3D
8 %matplotlib inline
9
10 sns.set(font_scale=1.3)
11
12 red = '#FF3300'
13 blue = '#0099CC'
14 green = '#00CC66'
```

## Метод главных компонент

### (Principal component analysis, PCA)

[It had to be U - the SVD song\\_\(https://www.youtube.com/watch?v=JEYLfIVvR9I\)](https://www.youtube.com/watch?v=JEYLfIVvR9I), --- веселая песенка :)

**Explained Visually (a setosa project):**

[Principal Component Analysis \(http://setosa.io/ev/principal-component-analysis/\)](http://setosa.io/ev/principal-component-analysis/).

[Eigenvectors and Eigenvalues \(http://setosa.io/ev/eigenvectors-and-eigenvalues/\)](http://setosa.io/ev/eigenvectors-and-eigenvalues/).

## SVD-разложение

Генерируем датасет

In [2]:

```
1 X = sp.stats.uniform.rvs(size=(100, 10))
2 print(X)
```

```
[2.24084070e-02 7.00103319e-01 2.63010135e-01 2.03580359e-01
 5.49954852e-01 7.02545097e-01 3.68183381e-01 2.60767053e-01
 4.95555721e-01 2.98442510e-01]
[8.80017678e-01 7.34578065e-01 7.48594445e-01 3.88825026e-01
 7.37351383e-01 4.21592240e-01 3.10308328e-01 8.05600492e-01
 2.46826164e-01 5.65853033e-01]
[4.57373636e-01 9.52827803e-01 1.28006858e-01 5.29167968e-01
 7.81841934e-01 4.54895379e-01 1.75524858e-01 7.23797660e-01
 1.52313322e-01 2.38744438e-01]
[8.79424355e-01 8.92249798e-01 7.41546486e-01 4.68253173e-01
 9.87763106e-01 3.71276482e-01 5.67146955e-02 4.01979515e-01
 5.93226129e-01 3.07978983e-01]
[6.69162116e-01 2.01043136e-01 1.40365214e-02 8.96574696e-01
 8.21204071e-01 1.43875152e-01 4.47468243e-01 4.74077162e-01
 8.29257112e-01 7.61010525e-01]
[5.69292042e-01 4.15699417e-01 2.34955400e-02 9.16514345e-01
 6.28465854e-01 6.65166514e-02 5.22224813e-01 7.34252138e-02
 6.70150425e-01 8.81686960e-02]
[1.41605564e-01 3.25534715e-01 8.17183697e-01 7.15028178e-01
 3.13000760e-01 6.46066040e-01 7.30007650e-01 8.50410005e-01
 3.13000760e-01 6.46066040e-01 7.30007650e-01 8.50410005e-01]
```

Вот так можно вычислять сингулярное разложение

In [3]:

```
1 U, D, V = sp.linalg.svd(X, full_matrices=False)
2 print(U.shape, D.shape, V.shape)
```

```
(100, 10) (10,) (10, 10)
```

## PCA

`sklearn.decomposition.PCA` (<https://scikit-learn.org/stable/modules/generated/sklearn.decomposition.PCA>) (`n_components=None`)

Методы:

- `fit(X)` --- обучиться на данных `X`;
- `fit_transform(X)` --- обучиться на данных `X` и вернуть сжатое представление `X`;
- `transform(X_new)` --- вернуть сжатое представление `X_new` для обученной ранее модели;
- `inverse_transform(Y)` --- восстановить сжатые данные `Y` в исходное пространство.

Атрибуты:

- `components_` --- главные компоненты в порядке убывания собственных чисел, размер (`n_components`, `n_features`);
- `explained_variance_` --- дисперсия вдоль главных компонент, равны собственным числам, размер (`n_components`);
- `explained_variance_ratio_` --- доля дисперсии, объясняемая каждой компонентой, размер (`n_components`);
- `mean_` --- среднее по данным, размер (`n_components`);
- `noise_variance_` --- оценка дисперсии шума для метода Probabilistic PCA.

Другие модификации, реализованные в sklearn:

- [KernelPCA](https://scikit-learn.org/stable/modules/generated/sklearn.decomposition.KernelPCA) [\\_](https://scikit-learn.org/stable/modules/generated/sklearn.decomposition.KernelPCA)(<https://scikit-learn.org/stable/modules/generated/sklearn.decomposition.KernelPCA>);
- [SparsePCA](https://scikit-learn.org/stable/modules/generated/sklearn.decomposition.SparsePCA) [\\_](https://scikit-learn.org/stable/modules/generated/sklearn.decomposition.SparsePCA)(<https://scikit-learn.org/stable/modules/generated/sklearn.decomposition.SparsePCA>);
- [IncrementalPCA](https://scikit-learn.org/stable/modules/generated/sklearn.decomposition.IncrementalPCA) [\\_](https://scikit-learn.org/stable/modules/generated/sklearn.decomposition.IncrementalPCA)(<https://scikit-learn.org/stable/modules/generated/sklearn.decomposition.IncrementalPCA>).

Генерируем двумерный датасет

In [4]:

```
1 X = sp.stats.multivariate_normal.rvs(size=150, mean=[0, 3],
2                                     cov=[[3, 1], [1, 1]])
```

Применяем PCA с одной главной компонентой

In [5]:

```
1 pca = PCA(n_components=1)
2 Y = pca.fit_transform(X)
3 X_hat = pca.inverse_transform(Y)
```

Его главные компоненты (точнее, одна компонента) -- двумерные векторы

In [6]:

```
1 pca.components_
```

Out[6]:

```
array([[ -0.91866203,  -0.3950444 ]])
```

Вектор средних

In [7]:

```
1 pca.mean_
```

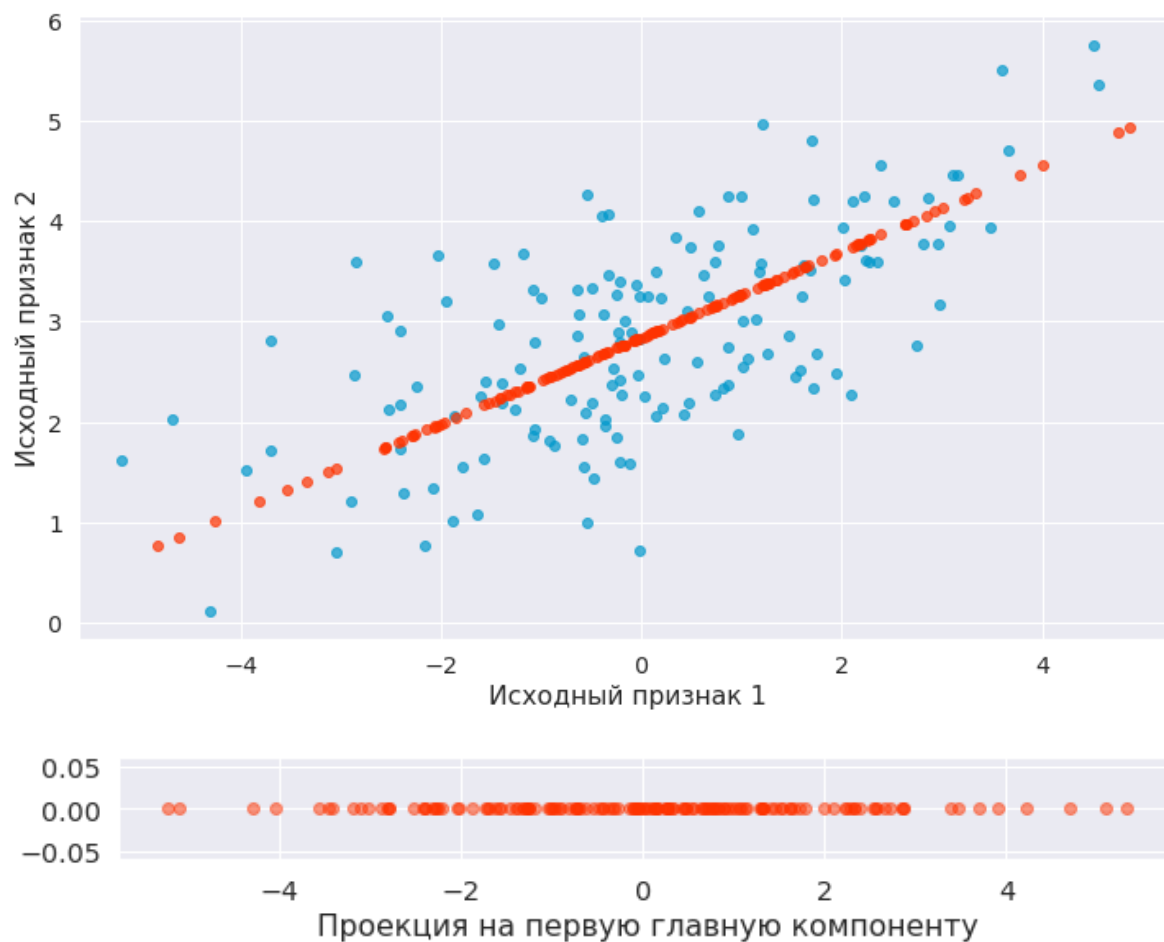
Out[7]:

```
array([0.0587315 , 2.86784748])
```

На первом графике синим отмечены исходные точки, красным - они же после проецирования и обратного преобразования. На втором графике точки в одномерном пространстве.

In [8]:

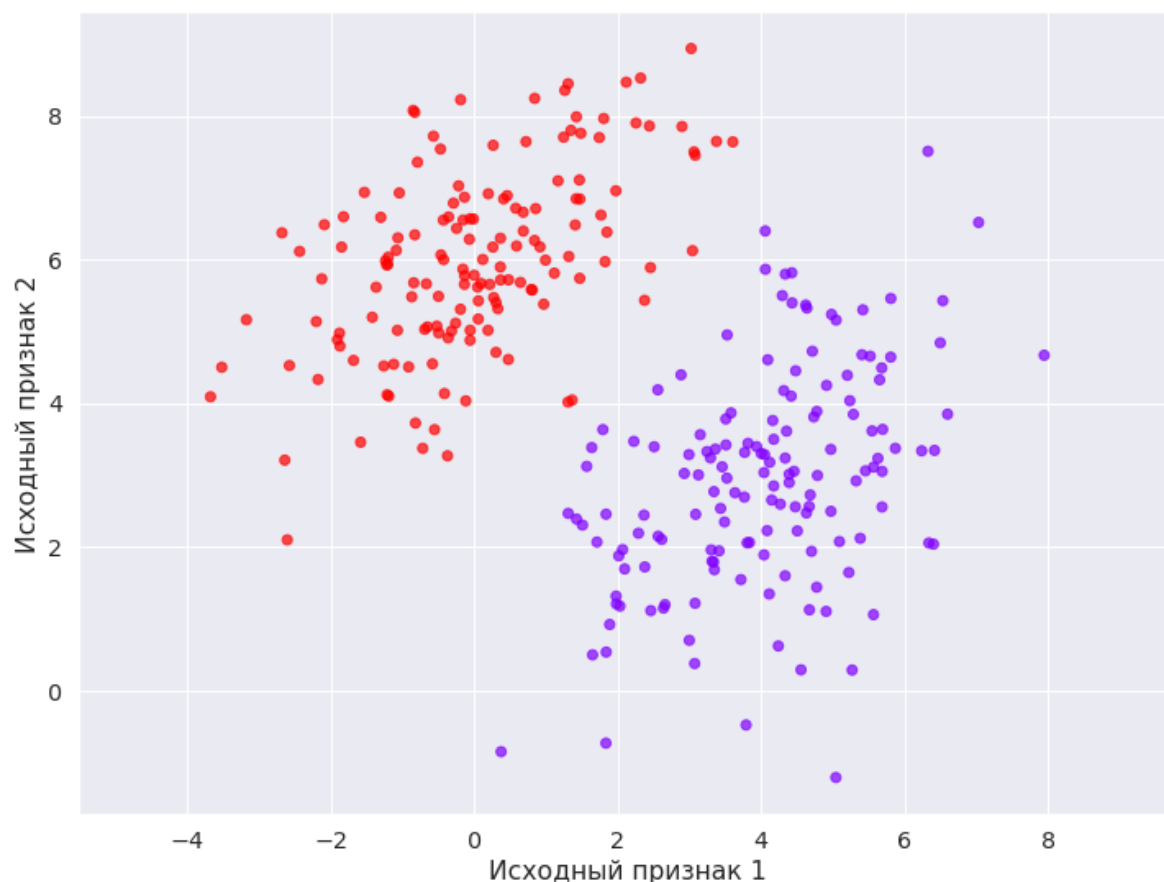
```
1 plt.figure(figsize=(12, 7))
2 plt.scatter(X[:, 0], X[:, 1], alpha=0.7, color=blue)
3 plt.scatter(X_hat[:, 0], X_hat[:, 1], color=red, alpha=0.7)
4 plt.xlabel('Исходный признак 1')
5 plt.ylabel('Исходный признак 2')
6 plt.axis('equal')
7 plt.show()
8
9 plt.figure(figsize=(10, 1))
10 plt.scatter(Y, np.zeros(len(Y)), alpha=0.5, color=red)
11 plt.xlabel('Проекция на первую главную компоненту')
12 plt.show()
```



Пусть есть два хорошо разделимых класса

In [9]:

```
1 a = sp.stats.multivariate_normal.rvs(size=150, mean=[4, 3],
2                                     cov=[[2, 1], [1, 2]])
3 b = sp.stats.multivariate_normal.rvs(size=150, mean=[0, 6],
4                                     cov=[[2, 1], [1, 2]])
5 X = np.vstack([a, b])
6 c = np.hstack([np.zeros(len(a)), np.ones(len(b))])
7
8 plt.figure(figsize=(12, 9))
9 plt.scatter(X[:, 0], X[:, 1], c=c, alpha=0.7, cmap='rainbow')
10 plt.xlabel('Исходный признак 1')
11 plt.ylabel('Исходный признак 2')
12 plt.xlim((-5, 9)), plt.ylim((-1, 10))
13 plt.axis('equal')
14 plt.show()
```



Проецируем на одномерное подпространство, тут хорошо заметна пространственная структура классов.

In [10]:

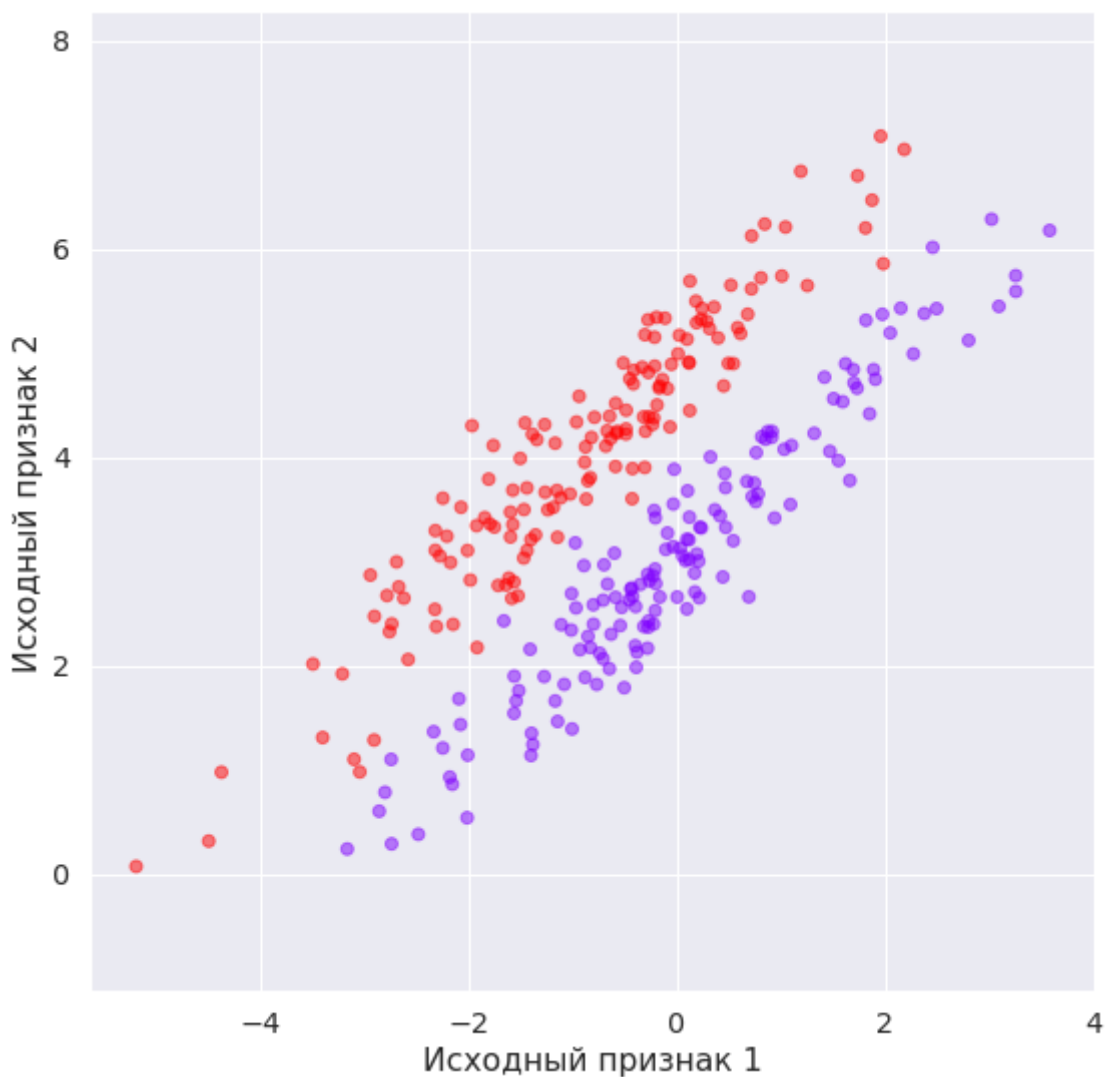
```
1  pca = PCA(n_components=1)
2  Y = pca.fit_transform(X)
3  X_hat = pca.inverse_transform(Y)
4
5  plt.figure(figsize=(12, 10))
6  plt.scatter(X[:, 0], X[:, 1], c=c, alpha=0.5, cmap='rainbow')
7  plt.scatter(X_hat[:, 0], X_hat[:, 1], c=green, alpha=0.7)
8  plt.xlabel('Исходный признак 1')
9  plt.ylabel('Исходный признак 2')
10 plt.xlim((-5, 9)), plt.ylim((-1, 10))
11 plt.axis('equal')
12 plt.show()
13
14 plt.figure(figsize=(10, 1))
15 plt.scatter(Y, np.zeros(len(Y)), c=c, alpha=0.3, cmap='rainbow')
16 plt.xlabel('Проекция на первую главную компоненту')
17 plt.show()
```



А что, если два вытянутых класса, расположенных близко, как на рисунке ниже?

In [11]:

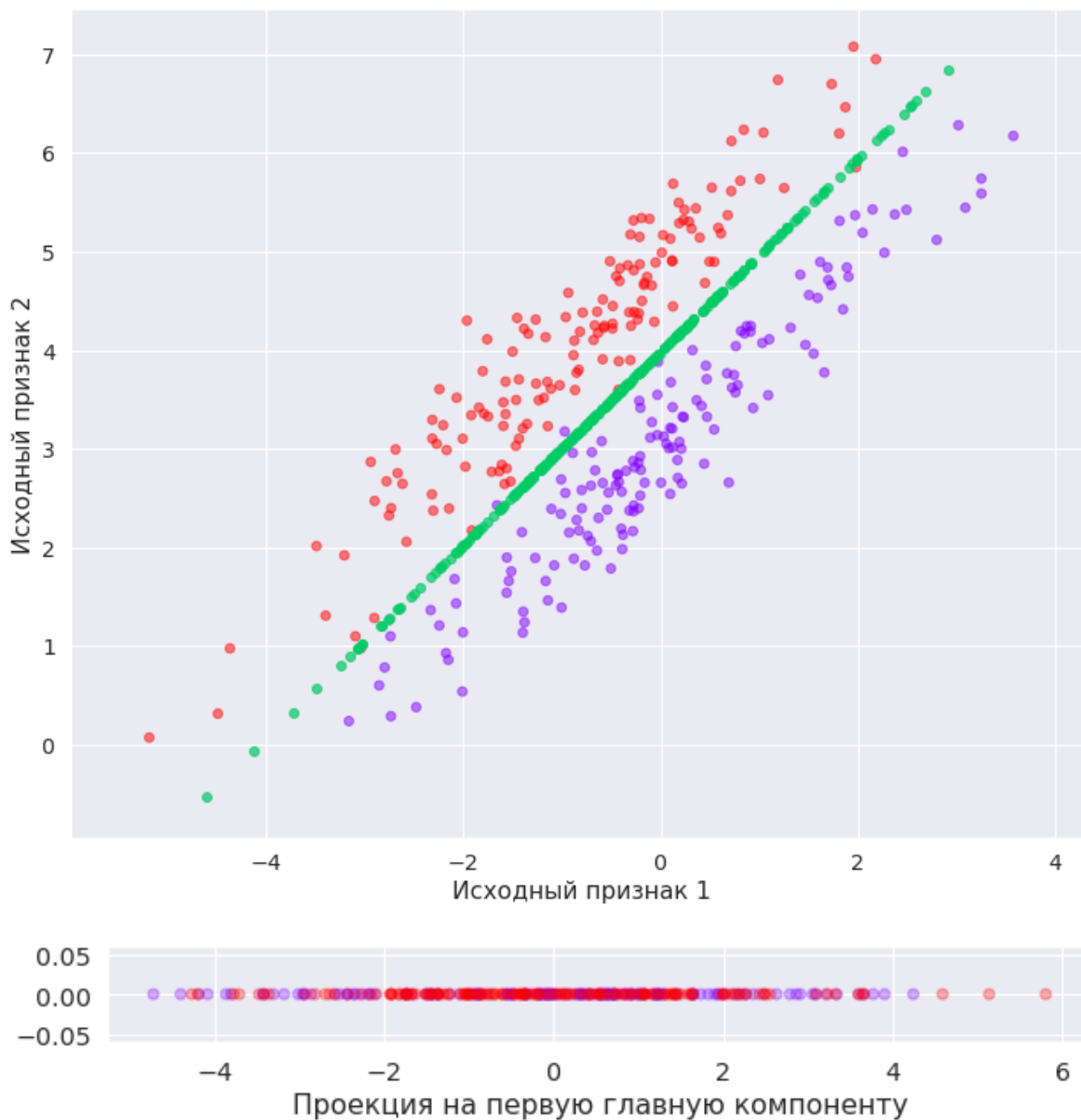
```
1 a = sp.stats.multivariate_normal.rvs(size=150, mean=[0, 3],
2                                     cov=[[2, 1.9], [1.9, 2]])
3 b = sp.stats.multivariate_normal.rvs(size=150, mean=[-1, 4],
4                                     cov=[[2, 1.9], [1.9, 2]])
5 X = np.vstack([a, b])
6 c = np.hstack([np.zeros(len(a)), np.ones(len(b))])
7
8 plt.figure(figsize=(9, 9))
9 plt.scatter(X[:, 0], X[:, 1], c=c, alpha=0.5, cmap='rainbow')
10 plt.xlabel('Исходный признак 1')
11 plt.ylabel('Исходный признак 2')
12 plt.xlim((-5, 5)), plt.ylim((-1, 9))
13 plt.axis('equal')
14 plt.show()
```



Главная компонента направлена вдоль этих классов, и при проецировании точки классов смешаются.

In [12]:

```
1  pca = PCA(n_components=1)
2  Y = pca.fit_transform(X)
3  X_hat = pca.inverse_transform(Y)
4
5  plt.figure(figsize=(12, 10))
6  plt.scatter(X[:, 0], X[:, 1], c=c, alpha=0.5, cmap='rainbow')
7  plt.scatter(X_hat[:, 0], X_hat[:, 1], c=green, alpha=0.7)
8  plt.xlabel('Исходный признак 1')
9  plt.ylabel('Исходный признак 2')
10 plt.xlim((-5, 5)), plt.ylim((-1, 9))
11 plt.axis('equal')
12 plt.show()
13
14 plt.figure(figsize=(10, 1))
15 plt.scatter(Y, np.zeros(len(Y)), c=c, alpha=0.3, cmap='rainbow')
16 plt.xlabel('Проекция на первую главную компоненту')
17 plt.show()
```



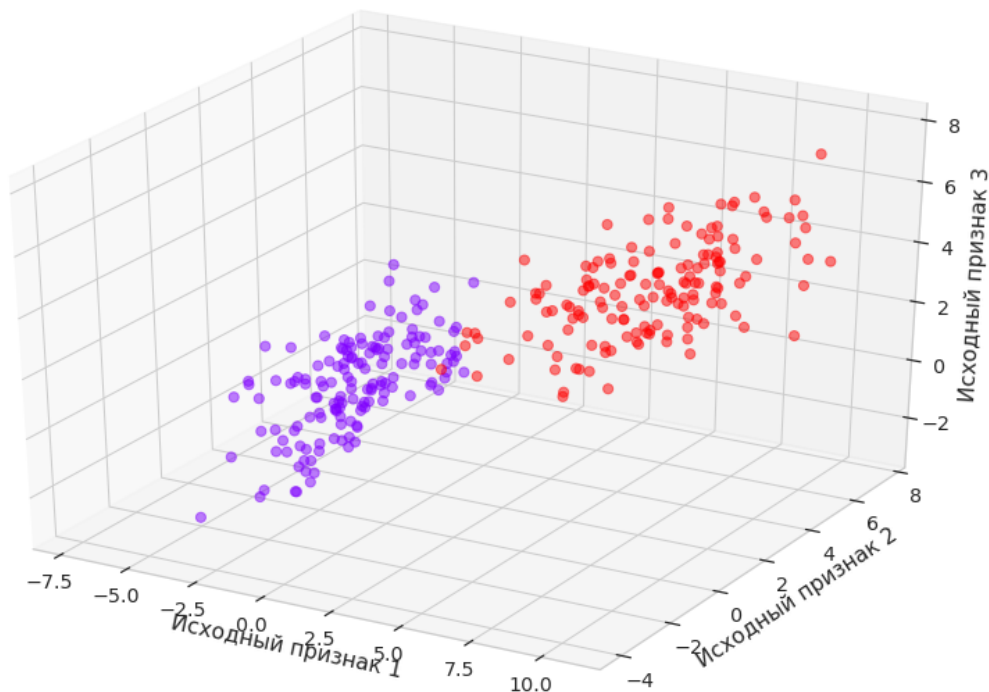


In [13]:

```
1 a = sp.stats.multivariate_normal.rvs(size=150, mean=[-2, 1, 0],
2                                     cov=[[3, -1, 1],
3                                         [-1, 4, 0.5],
4                                         [1, 0.5, 1.5]])
5 b = sp.stats.multivariate_normal.rvs(size=150, mean=[7, 3, 4],
6                                     cov=[[3, 2, 1],
7                                         [2, 4, 0.5],
8                                         [1, 0.5, 1.5]])
9 X = np.vstack([a, b])
10 c = np.hstack([np.zeros(len(a)), np.ones(len(b))])
```

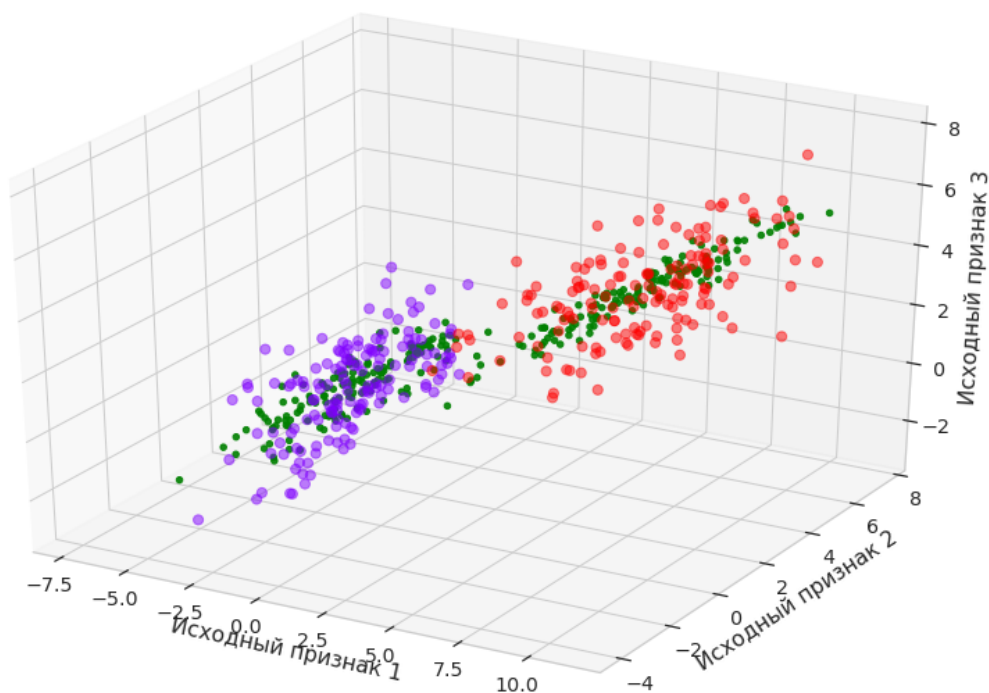
In [14]:

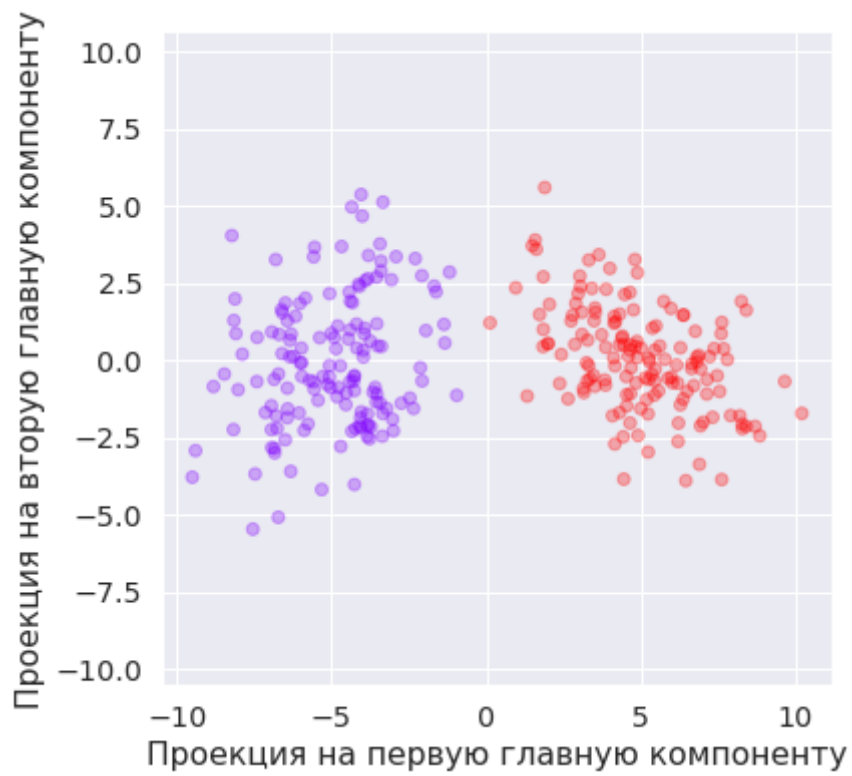
```
1 sns.set(font_scale=1.3, style='whitegrid')
2
3 fig = plt.figure(figsize=(15, 10))
4 ax = fig.add_subplot(111, projection='3d')
5 ax.scatter(X[:, 0], X[:, 1], X[:, 2],
6           s=50, c=c, alpha=0.5, cmap='rainbow')
7 ax.set_xlabel('Исходный признак 1')
8 ax.set_ylabel('Исходный признак 2')
9 ax.set_zlabel('Исходный признак 3')
10 # ax.axis('equal')
11 plt.show()
```



In [15]:

```
1  pca = PCA(n_components=2)
2  Y = pca.fit_transform(X)
3  X_hat = pca.inverse_transform(Y)
4
5  fig = plt.figure(figsize=(15, 10))
6  ax = fig.add_subplot(111, projection='3d')
7  ax.scatter(X[:, 0], X[:, 1], X[:, 2],
8             s=50, c=c, alpha=0.5, cmap='rainbow')
9  ax.scatter(X_hat[:, 0], X_hat[:, 1], X_hat[:, 2],
10            s=20, c='green', alpha=0.9)
11 ax.set_xlabel('Исходный признак 1')
12 ax.set_ylabel('Исходный признак 2')
13 ax.set_zlabel('Исходный признак 3')
14 # ax.axis('equal')
15 plt.show()
16
17 sns.set(font_scale=1.3)
18 plt.figure(figsize=(6, 6))
19 plt.scatter(Y[:, 0], Y[:, 1], c=c, alpha=0.3, cmap='rainbow')
20 plt.xlabel('Проекция на первую главную компоненту')
21 plt.ylabel('Проекция на вторую главную компоненту')
22 plt.axis('equal')
23 plt.show()
```





In [16]:

```
1 import plotly
2 import plotly.graph_objs as go
3 plotly.offline.init_notebook_mode()
```

In [17]:

```
1 fig = go.Figure()
2 fig.update_layout(
3     autosize=False,
4     width=1000,
5     height=1000,
6 )
7
8 fig.add_trace(go.Scatter3d(
9     x=X[c==0, 0], y=X[c==0, 1], z=X[c==0, 2], mode='markers',
10    marker={'size': 7, 'opacity': 0.65, 'color': red}
11 ))
12
13 fig.add_trace(go.Scatter3d(
14     x=X[c==1, 0], y=X[c==1, 1], z=X[c==1, 2], mode='markers',
15    marker={'size': 7, 'opacity': 0.65, 'color': blue}
16 ))
17
18 fig.add_trace(go.Scatter3d(
19     x=X_hat[:, 0], y=X_hat[:, 1], z=X_hat[:, 2], mode='markers',
20    marker={'size': 4, 'opacity': 0.9, 'color': green}
21 ))
22
23 fig.update_layout(
24     margin={'l': 0, 'r': 0, 'b': 0, 't': 0},
25     scene = dict(
26         xaxis_title='Признак 1',
27         yaxis_title='Признак 2',
28         zaxis_title='Признак 3',
29     )
30 )
31
32 # Отображение.
33 plotly.offline.iplot(fig)
```



## Другие методы снижения размерности

<http://scikit-learn.org/stable/modules/manifold.html#manifold> (<http://scikit-learn.org/stable/modules/manifold.html#manifold>).

Примеры:

[http://scikit-learn.org/stable/auto\\_examples/manifold/plot\\_compare\\_methods.html#sphx-glr-auto-examples-manifold-plot-compare-methods-py](http://scikit-learn.org/stable/auto_examples/manifold/plot_compare_methods.html#sphx-glr-auto-examples-manifold-plot-compare-methods-py) ([http://scikit-learn.org/stable/auto\\_examples/manifold/plot\\_compare\\_methods.html#sphx-glr-auto-examples-manifold-plot-compare-methods-py](http://scikit-learn.org/stable/auto_examples/manifold/plot_compare_methods.html#sphx-glr-auto-examples-manifold-plot-compare-methods-py)).

[http://scikit-learn.org/stable/auto\\_examples/manifold/plot\\_manifold\\_sphere.html#sphx-glr-auto-examples-manifold-plot-manifold-sphere-py](http://scikit-learn.org/stable/auto_examples/manifold/plot_manifold_sphere.html#sphx-glr-auto-examples-manifold-plot-manifold-sphere-py) ([http://scikit-learn.org/stable/auto\\_examples/manifold/plot\\_manifold\\_sphere.html#sphx-glr-auto-examples-manifold-plot-manifold-sphere-py](http://scikit-learn.org/stable/auto_examples/manifold/plot_manifold_sphere.html#sphx-glr-auto-examples-manifold-plot-manifold-sphere-py)).