

In [301]:

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 import scipy.stats as sps
4 from tqdm import tqdm_notebook as tqdm
5 from sklearn.metrics import roc_auc_score
6
7 %matplotlib inline
```

In [2]:

```
1 %set_env OCTAVE_EXECUTABLE=C:\Octave\octave-5.1.0-w64-64\mingw64\bin\octave.exe
```

env: OCTAVE\_EXECUTABLE=C:\Octave\octave-5.1.0-w64-64\mingw64\bin\octave.exe

In [3]:

```
1 from oct2py import Oct2Py
2 oc = Oct2Py()
```

In [193]:

```
1 n_obj = 10
2 dim = 3
3 n_models = 5
4 k = 10
5
6 X = np.random.rand(n_obj, dim)
7 w = np.random.rand(dim, n_models)
8 pi = np.zeros(n_models)
9 pi[0] = pi[1] = 0.5
10
11 a = oc.generate_mixture_logistic(X, w, pi.reshape(-1, 1))
12 b = oc.generate_single_logistic(X, w[:, 0].reshape(-1, 1))
```

In [219]:

```
1 X, y, idx = oc.generate_syntethic_cluster(100, 2, 0, 1., 0, 0, nout = 3)
2 y = y.reshape(y.size)
```

```
In [195]: 1 if X.shape[1] == 2:
2         X = np.hstack([X, np.ones((X.shape[0], 1))])
3         X[:,10]
```

```
Out[195]: array([[ -2.61266672, -2.23986673,  1.          ],
[ -4.25576278, -2.96898152,  1.          ],
[ -2.85738672, -2.18585945,  1.          ],
[ -1.00902584, -1.63937964,  1.          ],
[ -3.10686002, -1.93344944,  1.          ],
[ -3.74156942, -4.19801759,  1.          ],
[ -3.84263114, -2.86860189,  1.          ],
[ -3.08322392, -3.2292212 ,  1.          ],
[ -2.90667065, -2.77943204,  1.          ],
[ -0.81491968, -3.35739661,  1.          ]])
```

```
In [196]: 1 A_single, w, hessian_single = oc.maximize_evidence_single_logistic_laplace(X, y.reshape(-1, 1), nout = 3)
2         w_single = w.reshape(3)
```

```
In [197]: 1 A, w, hessian = oc.maximize_evidence_multilevel_logistic_laplace(X, y.reshape(-1, 1), idx, nout = 3)
```

```
In [198]: 1 w
```

```
Out[198]: Cell([array([[ 1.49701475],
[ -1.61026852],
[  0.52547267]]),
array([[ 1.30066182e+00],
[ -1.23768325e+00],
[ -5.40015117e-06]])])
```

In [199]:

1 A

```
Out[199]: Cell([array([[0.41133528, 0.          , 0.          ],
                        [0.          , 0.35494992, 0.          ],
                        [0.          , 0.          , 1.17733025]]),
               array([[5.48416951e-01, 0.00000000e+00, 0.00000000e+00],
                       [0.00000000e+00, 6.05139884e-01, 0.00000000e+00],
                       [0.00000000e+00, 0.00000000e+00, 6.23489183e+04]]))
```

In [200]:

1 hessian

```
Out[200]: Cell([array([[ 63.44196901,  56.29933652, -19.8839977 ],
                        [ 56.29933652,  54.84080098, -18.54343351],
                        [-19.8839977 , -18.54343351,   8.14188954]]),
               array([[7.51722318e+01, 7.44696527e+01, 2.37880220e+01],
                       [7.44696527e+01, 8.20620524e+01, 2.46806036e+01],
                       [2.37880220e+01, 2.46806036e+01, 6.23430103e+04]]))
```

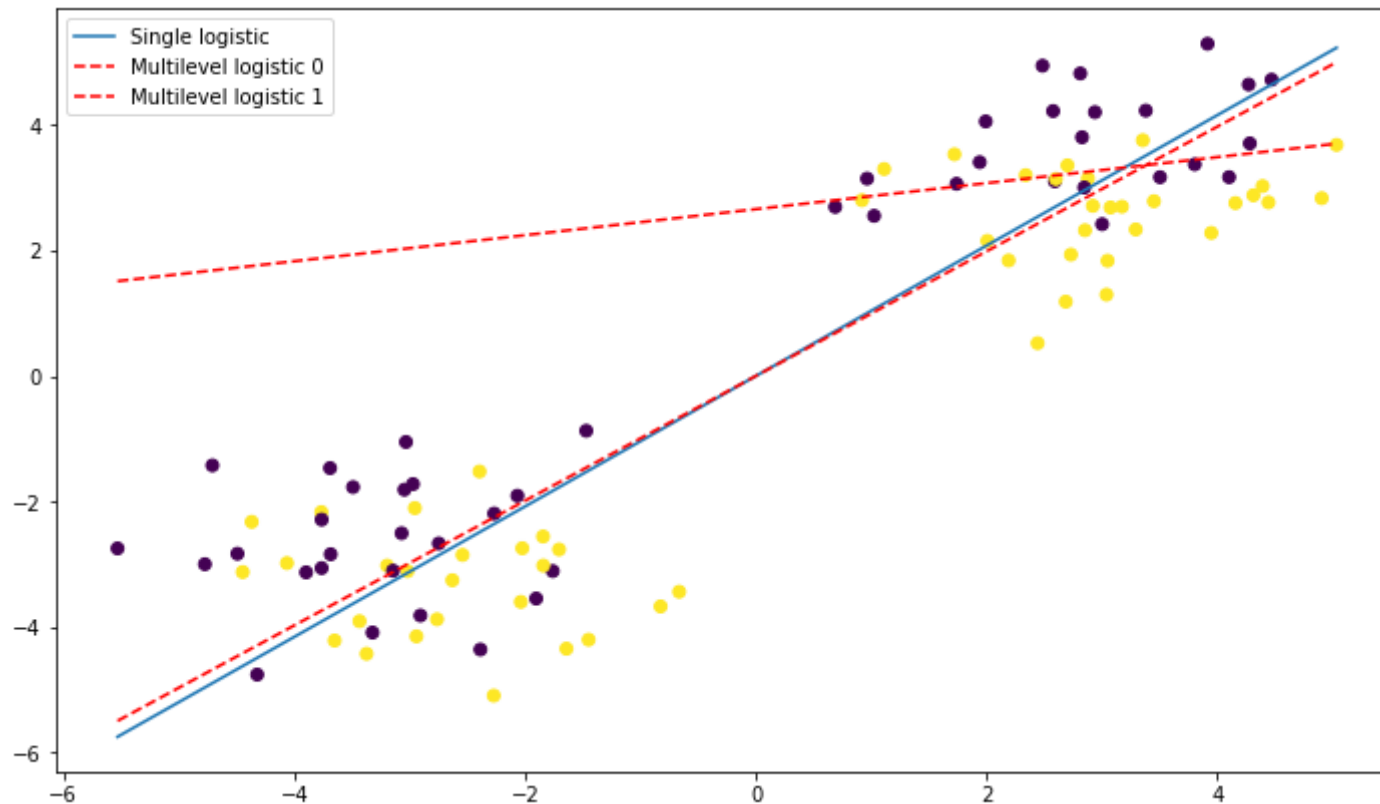
In [98]:

```
1 def get_significance_level(w0, hess0, w1, hess1):
2     score, n = oc.get_significance_level_no_intersect(w0, hess0, w1, hess1, nout = 2)
3     return 1 - sps.chi2(df = n).cdf(score)
4
5 def get_logistic_bound(grid, w):
6     return - w[0] / w[1] * grid - w[2] / w[1]
```

In [76]:

```
1 plt.figure(figsize = (12, 7))
2 plt.scatter(X[:, 0], X[:, 1], c = (y + 1) // 2)
3 grid = np.linspace(X[:, 0].min(), X[:, 0].max(), 500)
4 plt.plot(grid, get_logistic_bound(grid, w_single), label = 'Single logistic')
5 for i in range(2):
6     w_part = w[i].reshape(3)
7     print(w_part)
8     plt.plot(grid, get_logistic_bound(grid, w_part),
9             label = 'Multilevel logistic ' + str(i), color = 'red',
10            ls = '--')
11 plt.legend()
12 plt.show()
```

```
[ 4.95510684e-01 -4.98719087e-01 -1.53697521e-06]
[ 0.37691551 -1.81725113  4.83463227]
```



Видим, что на двух кластерах имеется различие. Сравним две модели на кластерах, а также общую при помощи `s_score`

```
In [99]: 1 get_significance_level(w[0], hessian[0], w[1], hessian[1])
```

```
Out[99]: 0.03726747732555202
```

Критерий отвергает гипотезу неразличимости моделей!

Обучим логистическую регрессию в нетривиальном случае

```
In [57]: 1 y_fake = np.array([int(X[i][1] > 3) for i in range(X.shape[0])])
2 y_fake[y_fake == 0] = -1
3 y_fake
```

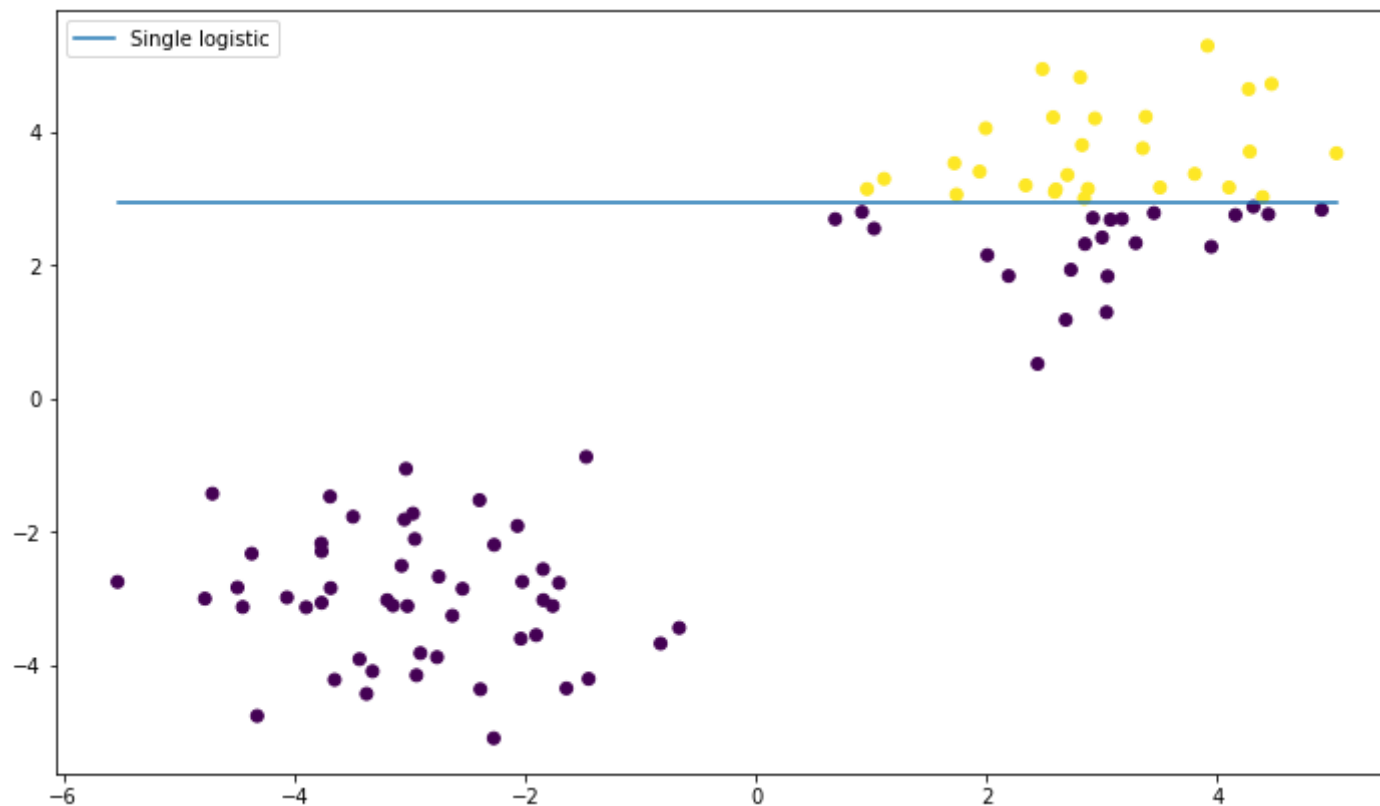
```
Out[57]: array([-1, -1, -1, -1, -1, -1, -1, -1, -1, -1, -1, -1, -1, -1, -1, -1, -1,
-1, -1, -1, -1, -1, -1, -1, -1, -1, -1, -1, -1, -1, -1, -1, -1, -1,
-1, -1, -1, -1, -1, -1, -1, -1, -1, -1, -1, -1, -1, -1, -1, -1, -1,
 1, -1, -1, -1, -1,  1, -1,  1,  1, -1,  1,  1,  1,  1, -1, -1, -1,
 1, -1,  1, -1,  1,  1, -1,  1,  1,  1, -1,  1, -1,  1, -1,  1,
-1, -1,  1,  1,  1,  1, -1,  1, -1, -1, -1,  1, -1,  1,  1])
```

```
In [58]: 1 func_args = {}
2 _, w = oc.maximize_evidence_single_logistic_laplace(X, y_fake.reshape(-1, 1), nout = 2)
```

```
In [59]: 1 w
```

```
Out[59]: array([[ -4.48667571e-06],
 [ 2.41739099e+01],
 [-7.10542972e+01]])
```

```
In [62]: 1 plt.figure(figsize = (12, 7))
2 plt.scatter(X[:, 0], X[:, 1], c = (y_fake + 1) // 2)
3 grid = np.linspace(X[:, 0].min(), X[:, 0].max(), 500)
4 plt.plot(grid, - w[0] / w[1] * grid - w[2] / w[1], label = 'Single logistic')
5 plt.legend()
6 plt.show()
```



Получается, что за месяц я наконец-таки научился нормально обучать линейную регрессию и считать s-score. Это уже почти победа

## Исследование ошибки первого рода :)

Проведём следующее исследование: повторим эксперимент, проведённый выше, но с различным размером выборки в кластерах. Поймем,

в какой момент начнет отвергаться гипотеза. Полученное число будем использовать в качестве изначальной оценки для размера блока в предложенном алгоритме.

Поскольку гипотеза неверна, то есть имеем справедливость альтернативной гипотезы  $H_1$ , интересующая нас величина - мощность критерия, которую можно проверить при помощи бутстрепа.

Постараемся добиться мощности критерия, равной 0.8

```
In [123]: 1 def get_significance_clustered(n_points = 10):
2     X, y, idx = oc.generate_syntethic_cluster(n_points * 2, 2, 0, 1., 0, 0, nout = 3)
3     X = np.hstack([X, np.ones((X.shape[0], 1))])
4     A, w, hessian = oc.maximize_evidence_multilevel_logistic_laplace(X, y.reshape(-1, 1), idx, nout = 3)
5     return get_significance_level(w[0], hessian[0], w[1], hessian[1])
6
7 def get_criterion_power(n_points, alpha = 0.05):
8     n_rejected = 0
9     n_samples = 100 # Вычисляем с точностью до 10-1
10    for i in range(n_samples):
11        pval = get_significance_clustered(n_points)
12        n_rejected += (pval < alpha)
13    return n_rejected / n_samples
```

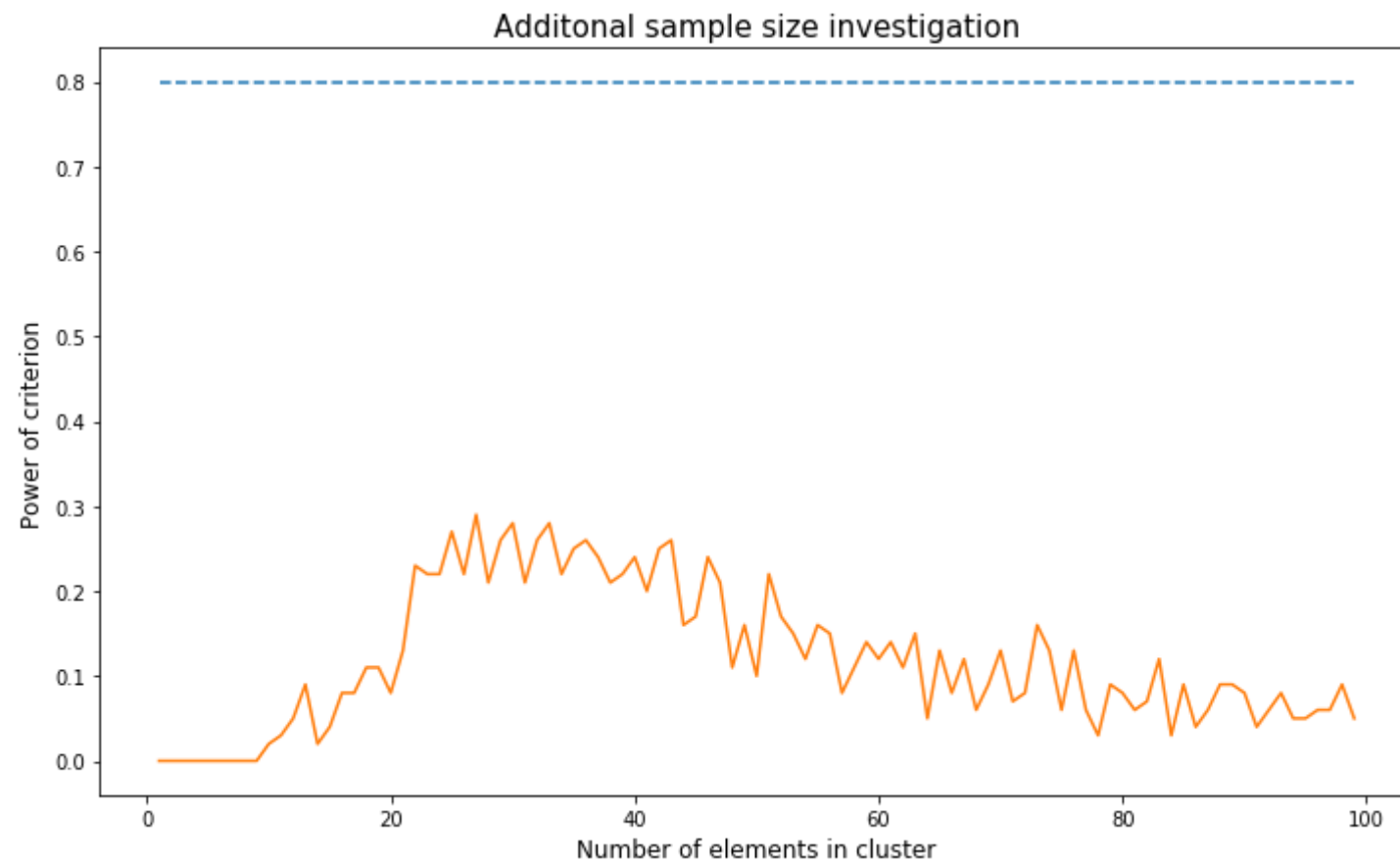
Переберем размер от 1 до 100

In [125]:

```
1 grid = np.arange(1, 100)
2 power = []
3
4 for n in tqdm(grid):
5     power.append(get_criterion_power(n))
6
7 plt.figure(figsize = (12, 7))
8 plt.title('Additonal sample size investigation', fontsize = 15)
9 plt.xlabel('Number of elements in cluster', fontsize = 12)
10 plt.ylabel('Power of criterion', fontsize = 12)
11 plt.plot(grid, np.ones_like(grid) * 0.8, ls = '--')
12 plt.plot(grid, power)
13 plt.show()
```

100% 99/99 [22:26<00:00, 13.76s/it]



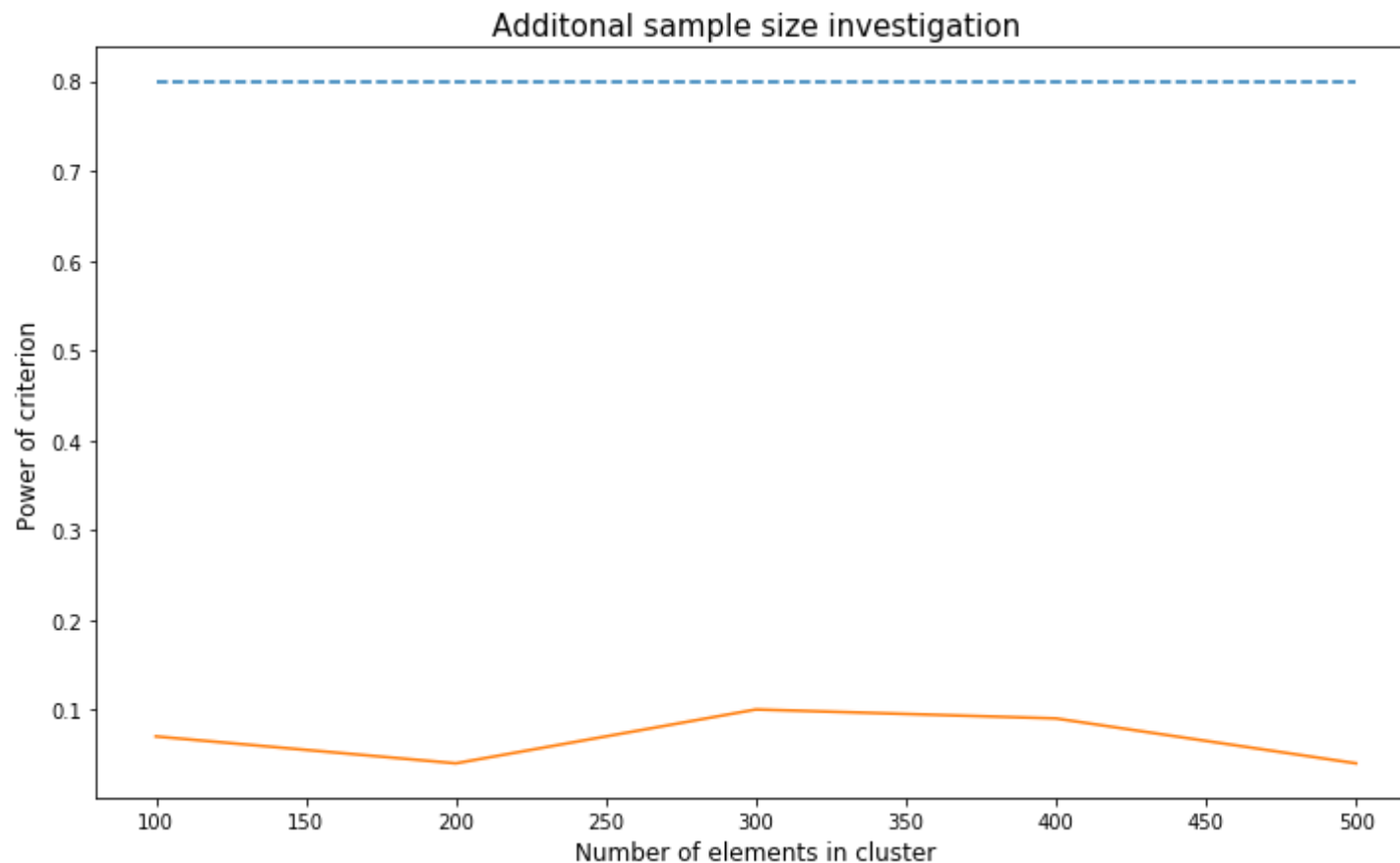


Теперь посмотрим, что будет, если добавлять по 100 объектов

In [126]:

```
1 grid = np.arange(1, 6) * 100
2 power = []
3
4 for n in tqdm(grid):
5     power.append(get_criterion_power(n))
6
7 plt.figure(figsize = (12, 7))
8 plt.title('Additonal sample size investigation', fontsize = 15)
9 plt.xlabel('Number of elements in cluster', fontsize = 12)
10 plt.ylabel('Power of criterion', fontsize = 12)
11 plt.plot(grid, np.ones_like(grid) * 0.8, ls = '--')
12 plt.plot(grid, power)
13 plt.show()
```

100% 5/5 [00:57<00:00, 11.61s/it]



Получаем, что оптимальный размер блока равен где-то 30. Именно его и возьмем для исследования ошибки алгоритма

## Генерация одного кластера данных

Чтобы убедиться в том, что предложенный алгоритм может дать результат, сгенерируем набор синтетических данных аналогично на двух кластерах, но после каждого блока разделяющие прямые будут поворачиваться вокруг заданных точек. Таким образом, это проэмулирует ситуацию, в которой одному признаковому описанию соответствуют разные ответы.

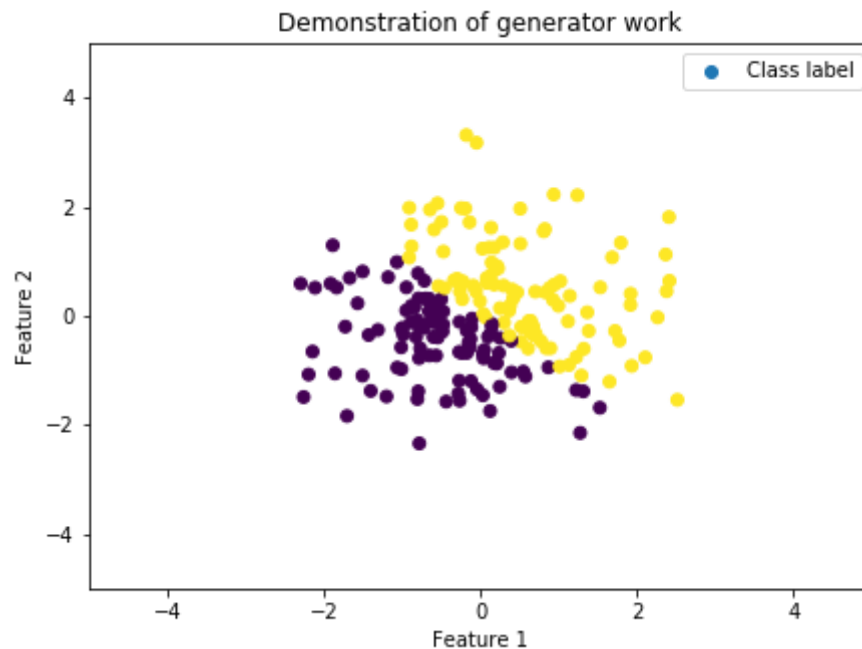
Будем случайно генерировать матрицу ковариаций и генерировать выборку нужного размера.

Определим сначала функцию генерации одного кластера вокруг заданной точки и генерации классов с учетом нормали к разделяющей прямой

```
In [156]: 1 def get_one_cluster(sample_size, point, normal):  
2     alphas = [1, 1]  
3     sigma = np.diag(alphas)  
4     X_sample = sps.multivariate_normal(mean = point, cov = sigma).rvs(sample_size)  
5     y = np.dot(X_sample - point, normal) > 0  
6     return X_sample, y * 2 - 1
```

In [157]:

```
1 X, y = get_one_cluster(200, [0, 0], [1, 1])
2
3 plt.figure(figsize = (7, 5))
4 plt.title('Demonstration of generator work')
5 plt.xlabel('Feature 1')
6 plt.ylabel('Feature 2')
7 plt.xlim(-5, 5)
8 plt.ylim(-5, 5)
9 plt.scatter(X[:, 0], X[:, 1], c = (y + 1) // 2, label = 'Class label')
10 plt.legend()
11 plt.show()
```



Возможно, следует добавить в распределение некоторый шум, но на первом этапе рассмотрим идеальное разделение. Обучим на кластере логистическую регрессию.

```
In [158]: 1 A, w = oc.maximize_evidence_single_logistic_laplace(X, y.reshape(-1, 1), nout = 2)
          2 w
```

```
Out[158]: array([[68.11527655],
                 [67.23790067]])
```

Видим, что вектор нормали (1, 1) почти угадан, качество разделения получится хорошим

## Исследование мощности критерия

Для предлагаемого алгоритма требуется отыскать оптимальный размер блока, на котором строится новая модель. Этот размер - гиперпараметр, однако в конкретном случае можно подобрать его, исходя из мощности критерия.

В зависимости от того, насколько различаются разделяющие прямые (по косинусному расстоянию), мощности могут получаться различными. Поэтому необходимо зафиксировать угол поворота нормали и посмотреть, на каком размере выборки мощность критерия станет достаточно высокой.

Для этого рассмотрим одиночную модель логистической регрессии, для различения будем использовать s-score

```
In [161]: 1 def rotate_origin_only(xy, radians):
          2     """Only rotate a point around the origin (0, 0)."""
          3     x, y = xy
          4     x_new = x * np.cos(radians) + y * np.sin(radians)
          5     y_new = -x * np.sin(radians) + y * np.cos(radians)
          6     return [x_new, y_new]
```

```
In [162]: 1 rotate_origin_only([1, 1], np.pi / 4.)
```

```
Out[162]: (1.4142135623730951, 0.0)
```

Напишем функцию, которая сгенерирует 2 кластера вокруг одной и той же точки с поворотом нормали на заданный угол

```
In [170]: 1 def generate_rotation_samples(sample_size, radians, point = [0, 0], normal1 = None):
2         if normal1 is None:
3             alpha_start = np.pi * np.random.rand()
4             normal1 = [np.cos(alpha_start), np.sin(alpha_start)]
5             normal2 = rotate_origin_only(normal1, radians)
6             X1, y1 = get_one_cluster(sample_size, point, normal1)
7             X2, y2 = get_one_cluster(sample_size, point, normal2)
8         return X1, y1, X2, y2
```

Теперь напишем вычисление мощности критерия по размеру выборки. Для этого будем генерировать 100 примеров пар выборок и смотреть, в каком проценте случаев отвергается гипотеза совпадения моделей. Мощность таким образом будет получена с точностью до 0.1

Возьмём угол поворота  $\alpha = \frac{\pi}{12}$

```
In [182]: 1 def get_criterion_power(sample_size, alpha = np.pi / 12):
2         n_samples = 100
3         n_rejections = 0
4         for i in range(n_samples):
5             X1, y1, X2, y2 = generate_rotation_samples(sample_size, alpha)
6             _, w1, hess1 = oc.maximize_evidence_single_logistic_laplace(X1, y1.reshape(-1, 1), nout = 3)
7             _, w2, hess2 = oc.maximize_evidence_single_logistic_laplace(X2, y2.reshape(-1, 1), nout = 3)
8             pval = get_significance_level(w1, hess1, w2, hess2)
9             n_rejections += (pval < 0.05)
10        return n_rejections / n_samples
```

Теперь переберем размер выборки от 1 до 100.

In [185]:

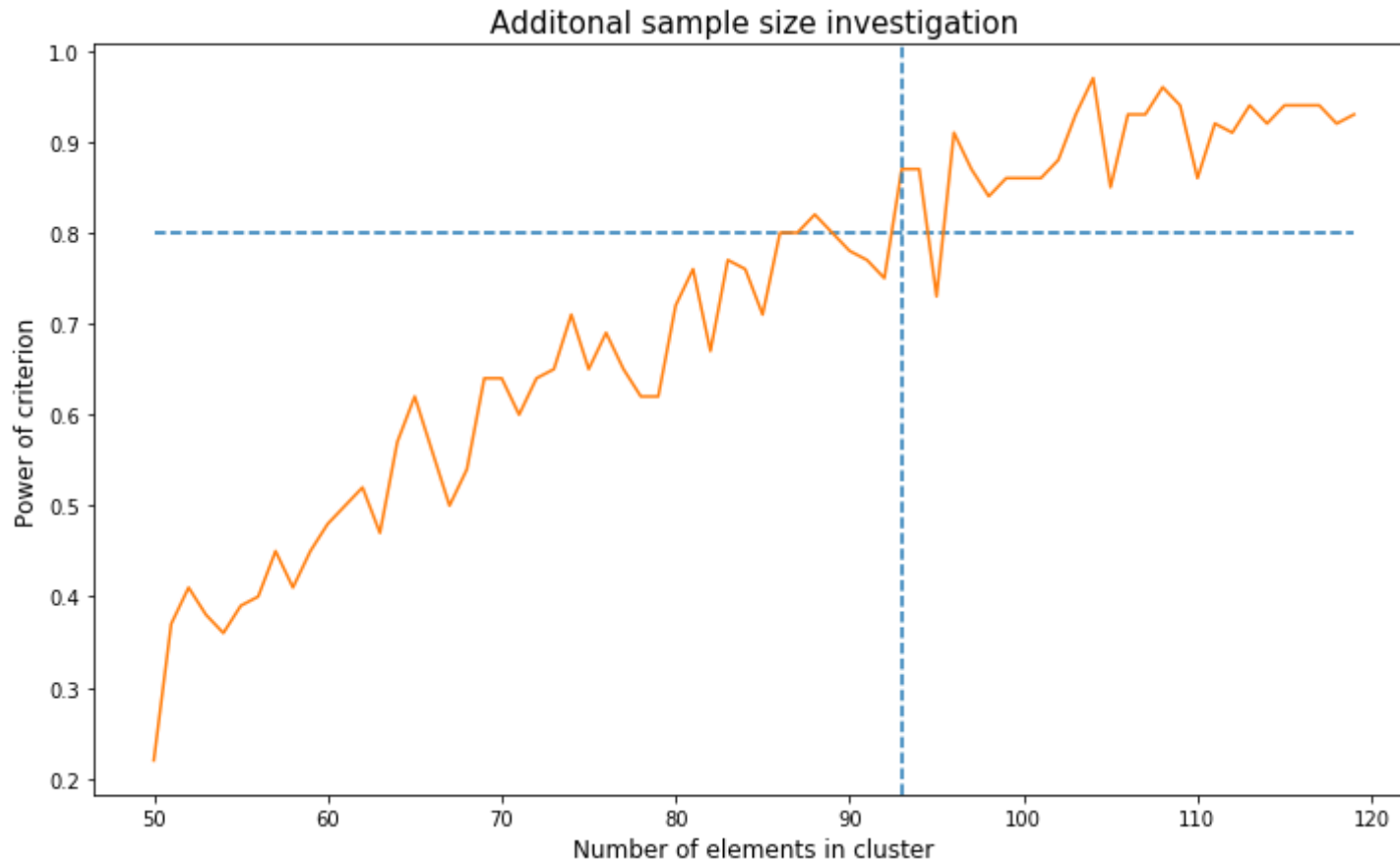
```
1 grid = np.arange(50, 120)
2 power = []
3
4 for n in tqdm(grid):
5     p_new = get_criterion_power(n)
6     print(p_new)
7     power.append(p_new)
```

```
0.6
0.64
0.65
0.71
0.65
0.69
0.65
0.62
0.62
0.72
0.76
0.67
0.77
0.76
0.71
0.8
0.8
0.82
0.8
0.78
```



In [189]:

```
1 plt.figure(figsize = (12, 7))
2 plt.title('Additional sample size investigation', fontsize = 15)
3 plt.xlabel('Number of elements in cluster', fontsize = 12)
4 plt.ylabel('Power of criterion', fontsize = 12)
5 plt.plot(grid, np.ones_like(grid) * 0.8, ls = '--')
6 plt.plot(grid, power)
7 plt.axvline(93, ls = '--')
8 plt.show()
```



Видим, что на размере выборки около 90 достигается значение мощности 0.8, а значит можно использовать данный критерий. При этом, как видно из анализа ошибки первого рода, такое значение допускает также мало ложных отвержений. Таким образом будем использовать такое значение, а именно размер блока  $C = 93$

Впоследствии необходимо запустить этот же анализ с точностью до 0.01

## **Реализация алгоритма**

Напишем функцию, которая будет делить пришедшую выборку на блоки заданного размера и давать предсказания вероятностей с помощью нового алгоритма. Затем сравним имеющиеся результаты по AUC-ROC с алгоритмом, обученным по всей выборке, а также обучаемый лишь на префиксах перед тем, как предсказывать следующий элемент (эмуляция процесса во времени)

In [325]:

```
1 def predict_logistic_proba(X, w):
2     scores = X.dot(w)
3     probas = np.exp(scores) / (1 + np.exp(scores))
4     return probas
5
6 def predict_proba_multilevel(X, w, idx):
7     X1 = X[idx == 1]
8     X2 = X[idx == 2]
9     pred1 = predict_logistic_proba(X1, w[0])
10    pred2 = predict_logistic_proba(X2, w[1])
11    pred = np.zeros(X.shape[0])
12    pred[idx == 1] = pred1
13    pred[idx == 2] = pred2
14    return pred
15
16 def get_new_algo_predictions(X, y, block_size, idx):
17     curr_start_ind = 0
18     curr_learning_start = 0
19     curr_w = np.zeros((2, 3))
20     curr_hess = np.zeros((2, 3, 3))
21     pred = []
22     while curr_start_ind < X.shape[0]:
23         curr_end_ind = min(X.shape[0], curr_start_ind + block_size)
24         X_new = X[curr_start_ind: curr_end_ind]
25         y_new = y[curr_start_ind: curr_end_ind]
26         idx_new = idx[curr_start_ind: curr_end_ind]
27         pred_new = predict_proba_multilevel(X_new, curr_w, idx_new)
28         pred += pred_new.tolist()
29
30         _, w_new, hess_new = oc.maximize_evidence_multilevel_logistic_laplace(X_new, y_new.reshape(-1, 1),
31                                                                                   idx_new.reshape(-1, 1),
32                                                                                   nout = 3)
33
34         w_new = np.array([w.reshape(3) for w in w_new])
35         pval1 = get_significance_level(curr_w[0].reshape(-1, 1), curr_hess[0],
36                                       w_new[0].reshape(-1, 1), hess_new[0])
37         pval2 = get_significance_level(curr_w[1].reshape(-1, 1), curr_hess[1],
38                                       w_new[1].reshape(-1, 1), hess_new[1])
39
40         if pval1 < 0.05 or pval2 < 0.05:
41             curr_w = w_new
42             curr_hess = hess_new
```

```

42         curr_learning_start = curr_start_ind
43     else:
44         X_new = X[curr_learning_start: curr_end_ind]
45         y_new = y[curr_learning_start: curr_end_ind]
46         idx_new = idx[curr_learning_start: curr_end_ind]
47         _, w_new, hess_new = oc.maximize_evidence_multilevel_logistic_laplace(X_new, y_new.reshape(-1, 1),
48                                                                                   idx_new.reshape(-1, 1),
49                                                                                   nout = 3)
50         w_new = np.array([w.reshape(3) for w in w_new])
51         curr_w = w_new
52         curr_hess = hess_new
53
54     curr_start_ind = curr_end_ind
55     return pred
56
57 def get_old_algo_predictions(X, y, block_size, idx):
58     curr_start_ind = 0
59     curr_w = np.zeros((2, 3))
60     curr_hess = np.zeros((2, 3, 3))
61     pred = []
62     while curr_start_ind < X.shape[0]:
63         curr_end_ind = min(X.shape[0], curr_start_ind + block_size)
64         X_new = X[curr_start_ind: curr_end_ind]
65         y_new = y[curr_start_ind: curr_end_ind]
66         idx_new = idx[curr_start_ind: curr_end_ind]
67         pred_new = predict_proba_multilevel(X_new, curr_w, idx_new)
68         pred += pred_new.tolist()
69
70         X_new = X[: curr_end_ind]
71         y_new = y[: curr_end_ind]
72         idx_new = idx[: curr_end_ind]
73         _, w_new, hess_new = oc.maximize_evidence_multilevel_logistic_laplace(X_new, y_new.reshape(-1, 1),
74                                                                                   idx_new.reshape(-1, 1),
75                                                                                   nout = 3)
76         w_new = np.array([w.reshape(3) for w in w_new])
77         curr_w = w_new
78         curr_hess = hess_new
79
80     curr_start_ind = curr_end_ind
81     return pred

```

## Генерация итогового датасета

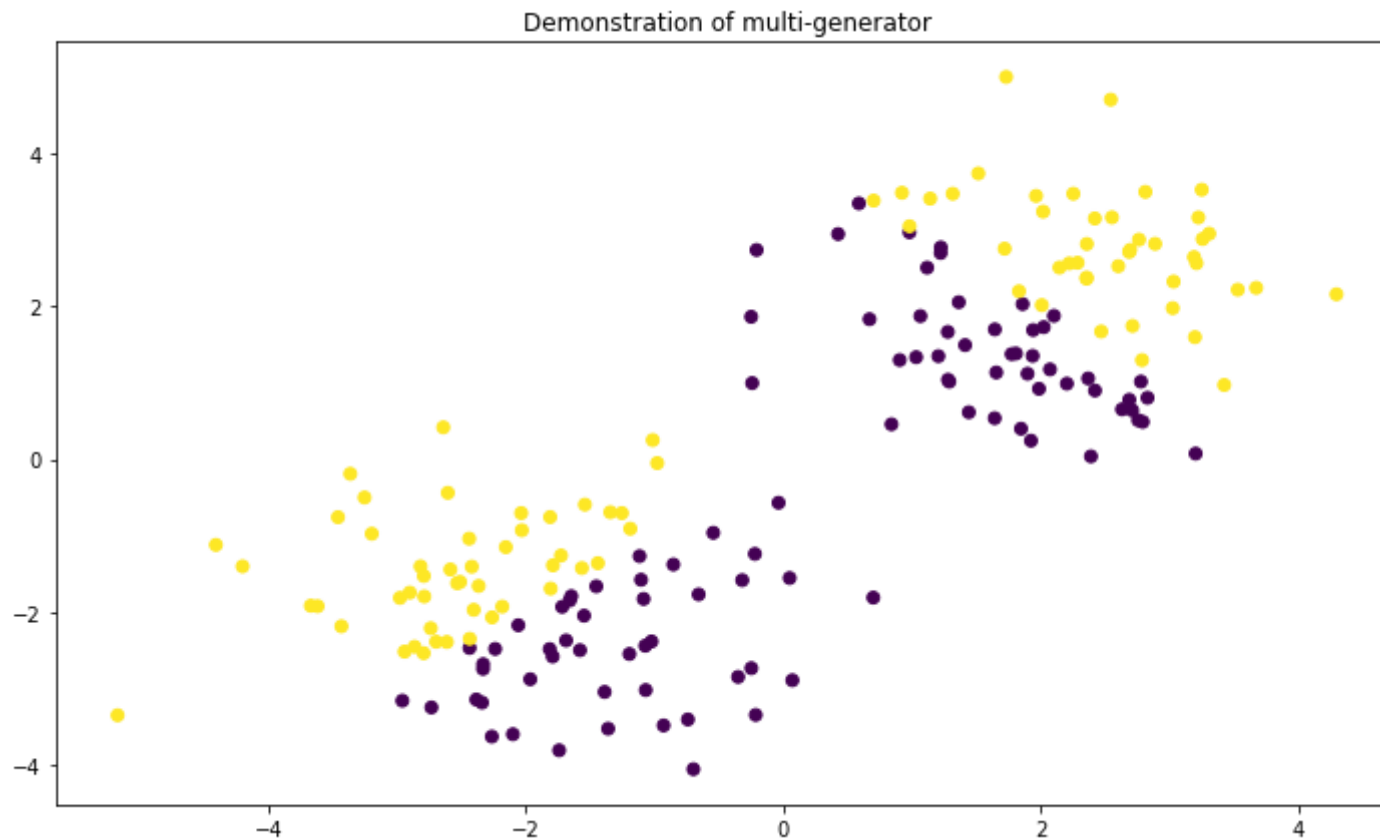
Теперь будем генерировать данные для двух кластеров сразу. Сделаем так, чтобы обе разделяющие прямые сделали полный оборот (для этого надо сделать 24 генерации).

В каждом положении будем генерировать 2 блока, чтобы модель успевала учитывать изменения во времени. Итого должно получиться 48 блоков по 93 примера в каждом, и всё это на 2 кластерах. То есть получится около 10000 примеров.

Для начала сгенерируем один пример в качестве демонстрации

In [213]:

```
1 normal_left = [1, 1]
2 normal_right = [-1, 1]
3
4 X1, y1 = get_one_cluster(93, [2, 2], normal_left)
5 X2, y2 = get_one_cluster(93, [-2, -2], normal_right)
6 X = np.vstack([X1, X2])
7 y = np.concatenate([y1, y2])
8
9 plt.figure(figsize = (12, 7))
10 plt.title('Demonstration of multi-generator')
11 plt.scatter(X[:, 0], X[:, 1], c = (y + 1) // 2)
12 plt.show()
```



Видим, что в обоих кластерах есть разделяющие прямые, при этом они перпендикулярны друг другу

```

In [365]: 1 X_all = np.empty(shape = (0, 2))
          2 y_all = np.empty(shape = 0)
          3 idx_all = np.empty(shape = 0)
          4 sample_size = 93
          5 alpha = np.pi / 12
          6
          7 normal_left = [1, 1]
          8 normal_right = [-1, 1]
          9
         10 for alpha_ind in range(24):
         11     for i in range(2):
         12         X1_l, y1_l, X2_l, y2_l = generate_rotation_samples(sample_size, alpha, [-3, -3], normal_left)
         13         X1_r, y1_r, X2_r, y2_r = generate_rotation_samples(sample_size, alpha, [3, 3], normal_right)
         14         idx_new = np.concatenate([np.zeros(sample_size), np.ones(sample_size)]) + 1
         15         X_all = np.vstack([X_all, X2_l, X2_r])
         16         y_all = np.concatenate([y_all, y2_l, y2_r])
         17         idx_all = np.concatenate([idx_all, idx_new])
         18
         19     normal_left = rotate_origin_only(normal_left, alpha)
         20     normal_right = rotate_origin_only(normal_right, alpha)

```

```

In [366]: 1 X_all = np.hstack([X_all, np.ones((X_all.shape[0], 1))])
          2 X_all.shape, y_all.shape, idx_all.shape

```

Out[366]: ((8928, 3), (8928,), (8928,))

## Тестирование алгоритма

Остаётся только получить предсказания и вычислить AUC-ROC для двух алгоритмов

```

In [372]: 1 pred_new = get_new_algo_predictions(X_all, y_all, 93 * 2, idx_all)

```

```

In [373]: 1 roc_auc_score((y_all + 1) // 2, pred_new)

```

Out[373]: 0.9834559960247349

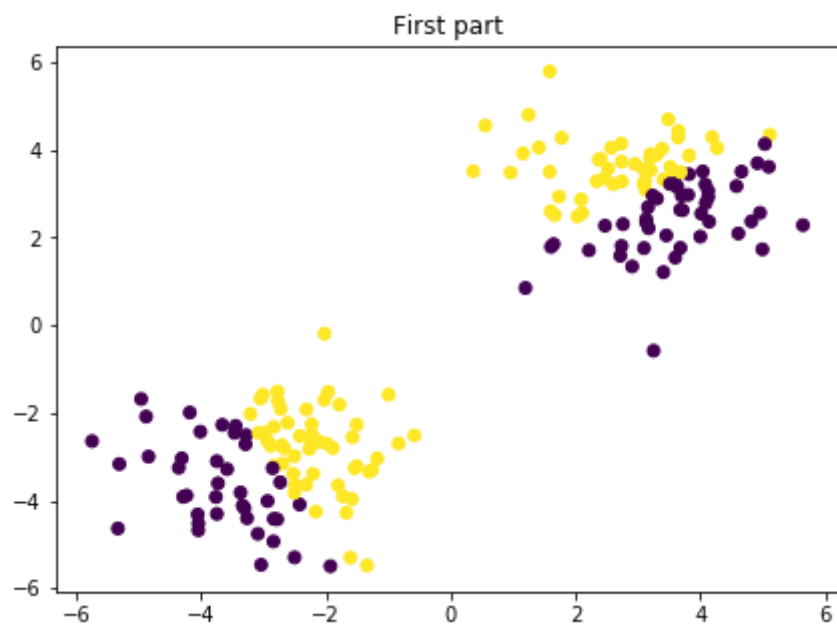
```
In [374]: 1 pred_old = get_old_algo_predictions(X_all, y_all, 93 * 2, idx_all)
```

```
In [375]: 1 roc_auc_score((y_all + 1) // 2, pred_old)
```

```
Out[375]: 0.611747159090909
```

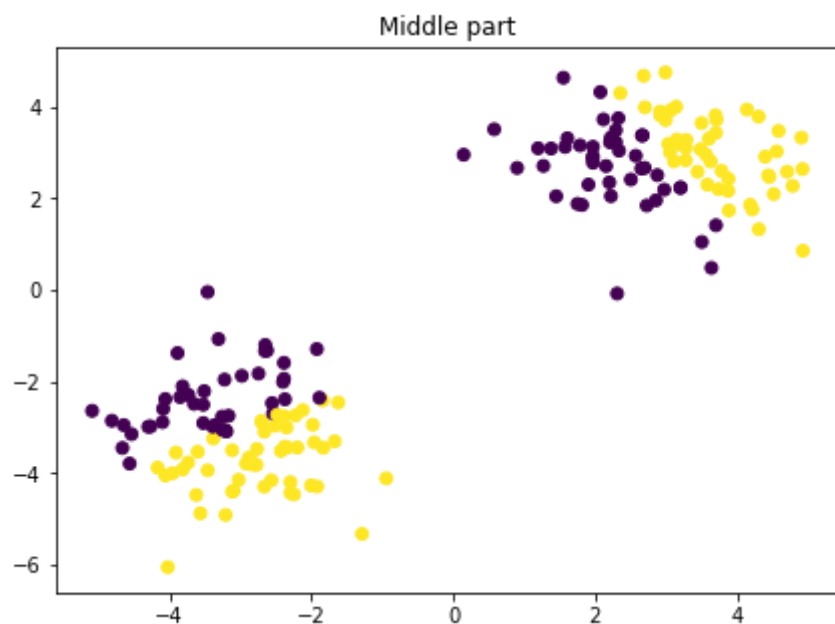
## Визуализация датасета

```
In [379]: 1 plt.figure(figsize = (7, 5))  
2 plt.title('First part')  
3 plt.scatter(X_all[:93 * 2, 0], X_all[:93 * 2, 1], c = (y_all[:93 * 2] + 1) // 2)  
4 plt.show()
```





```
In [381]: 1 plt.figure(figsize = (7, 5))
          2 plt.title('Middle part')
          3 plt.scatter(X_all[93 * 24:93 * 26, 0], X_all[93 * 24:93 * 26, 1], c = (y_all[93 * 24:93 * 26] + 1) // 2)
          4 plt.show()
```



```
In [382]: 1 plt.figure(figsize = (7, 5))
          2 plt.title('Everything')
          3 plt.scatter(X_all[:, 0], X_all[:, 1], c = (y_all + 1) // 2)
          4 plt.show()
```

