

In [1]:

```
1 import os
2
3 import PIL
4 import numpy as np
5 import scipy.stats as sps
6 import pandas as pd
7
8 import seaborn as sns
9 import matplotlib.pyplot as plt
10 import matplotlib.gridspec as gridspec
11
12 from sklearn.pipeline import Pipeline
13 from sklearn.svm import SVC, SVR, LinearSVC
14 from sklearn.metrics import accuracy_score, mean_squared_error
15 from sklearn.preprocessing import PolynomialFeatures, StandardScaler
16 from sklearn.model_selection import train_test_split, RandomizedSearchCV
17 from sklearn.datasets import make_moons, make_blobs, make_circles
18
19 sns.set(style='dark', font_scale=1.5)
```

started 08:39:23 2020-03-27, finished in 1.33s

SVM в scikit-learn

- Эффективная реализация SVM с **линейным ядром**: [классификация \(https://scikit-learn.org/stable/modules/generated/sklearn.svm.LinearSVC.html#sklearn.svm.LinearSVC\)](https://scikit-learn.org/stable/modules/generated/sklearn.svm.LinearSVC.html#sklearn.svm.LinearSVC), [регрессия \(https://scikit-learn.org/stable/modules/generated/sklearn.svm.LinearSVR.html#sklearn.svm.LinearSVR\)](https://scikit-learn.org/stable/modules/generated/sklearn.svm.LinearSVR.html#sklearn.svm.LinearSVR)
- Обычный SVM с произвольными ядрами: [классификация \(https://scikit-learn.org/stable/modules/generated/sklearn.svm.SVC.html#sklearn.svm.SVC\)](https://scikit-learn.org/stable/modules/generated/sklearn.svm.SVC.html#sklearn.svm.SVC), [регрессия \(https://scikit-learn.org/stable/modules/generated/sklearn.svm.SVR.html#sklearn.svm.SVR\)](https://scikit-learn.org/stable/modules/generated/sklearn.svm.SVR.html#sklearn.svm.SVR)

Описание наиболее важных параметров:

- `C` : Параметр регуляризации. Сила регуляризации обратно пропорциональна `C` . Должна быть строго положительной.
- `kernel` : Тип ядра, который будет использоваться в алгоритме. Это должен быть `linear` , `poly` , `rbf` , `sigmoid` , `precomputed` или `callable`.
- `degree` : Степень функции ядра полинома (`poly`). Игнорируется другими ядрами.
- `gamma` : Коэффициент ядра для `rbf` , `poly` и `sigmoid` .

Описание наиболее важных атрибутов:

- `support_` : Индексы опорных векторов.
- `support_vectors_` : Опорные вектора.
- `n_support_` : Количество опорных векторов для каждого класса.
- `coef_` : Оценка коэффициентов (коэффициенты в основной задаче). Это доступно только в случае линейного ядра.
- `intercept_` : Константы в функции принятия решения.
- `decision_function()` : Вычисляет функцию принятия решения для объекта.

Советы по практическому использованию

- **Размер кеша ядра:** для SVC, SVR размер кеша ядра сильно влияет на время выполнения для более крупных выборок. Если у вас достаточно RAM, рекомендуется установить для параметра `cache_size` более высокое значение, чем значение по умолчанию, равное 200(МБ), например, 500(МБ) или 1000(МБ).
- **Значение C :** Величина C по умолчанию равна 1, и это разумный выбор по умолчанию. Если у вас много шумных наблюдений, вы должны уменьшить его. Это соответствует более регуляризованной оценке. `LinearSVC` и `LinearSVR` менее чувствительны к C , когда он становится большим, и результаты прогнозирования перестают улучшаться после определенного порога. Между тем, большие значения C будут занимать больше времени для обучения, иногда до 10 раз дольше.
- Модели SVM не инвариантны к масштабированию, поэтому **настоятельно рекомендуется масштабировать ваши данные.**

Ядерные функции

В sklearn-реализации ядром может быть любая функция из следующих:

- **linear** : $K(x, z) = \langle x, z \rangle$
- **polynomial** : $K(x, z) = (\gamma \langle x, z \rangle + r)^k$, где k определяется ключевым словом `degree`, r ключевым словом `coef0`.
- **rbf** : $K(x, z) = \exp(-\gamma \|x - z\|^2)$, где $\gamma > 0$ определяется ключевым словом `gamma`
- **sigmoid** : $K(x, z) = \tanh(\gamma \langle x, z \rangle + r)$ где r определяется ключевым словом `coef0`.

Сложность вычислений

Support Vector Machines являются мощными инструментами, но их требования к вычислительным ресурсам быстро растут с ростом обучающей выборки. SVM решает задачу квадратичного программирования (QP), отделяя опорные векторы от остальной части обучающих данных. QP solver, используемый этой реализацией на основе `libsvm`, масштабируется между $O(dn^2)$ и $O(dn^3)$ (где n -- размер выборки, d --- число признаков) в зависимости от того, насколько эффективно на практике используется кеш `libsvm` (зависит от данных).

Также обратите внимание, что для линейного случая алгоритм, используемый в [LinearSVC \(https://scikit-learn.org/stable/modules/generated/sklearn.svm.LinearSVC.html\)](https://scikit-learn.org/stable/modules/generated/sklearn.svm.LinearSVC.html) в реализации `liblinear`, гораздо эффективнее, чем его аналог SVC на основе `libsvm`, и может масштабироваться почти линейно до миллиона выборок и/или признаков.

Простые примеры классификации

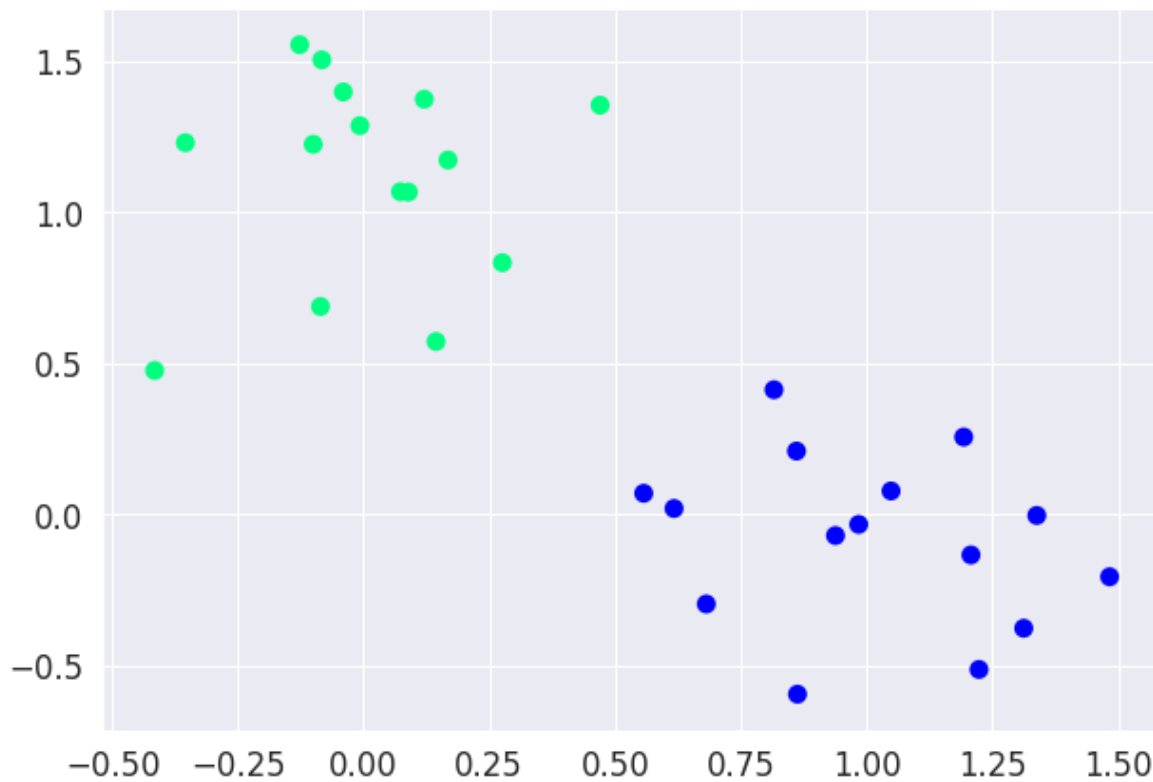
Рассмотрим SVM для решения задачи классификации:

Для начала сгенерируем данные

In [2]:

```
1 X, Y = make_blobs(n_samples=30, n_features=2, centers=[(1, 0), (0, 1)],
2                   cluster_std=0.3, random_state=21)
3
4 plt.figure(figsize=(10, 7))
5 plt.scatter(X[:, 0], X[:, 1], c=Y, linewidths=0, s=100, cmap='winter')
6 plt.grid()
7 plt.show()
```

started 08:39:26 2020-03-27, finished in 330ms



Вспомогательная функция для отрисовки

In [3]:

```
1  ▼ def draw_graphics(models, X, Y, support_size=800,
2      point_size=50, margin=False):
3      ...
4      Визуализирует решающие правила для каждой модели.
5      models --- все обученные SVM-модели, которые нужно визуализировать.
6      X --- объекты для визуализации (предполагается обучающая выборка)
7      Y --- ответы для визуализации (предполагается обучающая выборка)
8      X_test --- объекты, на которых необходимо посчитать качество
9      Y_test --- соответствующие им ответы.
10     point_size --- размер точки на графике
11     ncol --- количество колонок у таблицы графиков
12  ▼ margin --- если True, то визуализируется решающая функция,
13      иначе решающее правило
14     params --- список параметров SVM, которые нужно напечатать на графике
15     ...
16
17     # определение количества строк таблицы графиков в зависимости
18     # от количества колонок и количества моделей
19     n_rows = (len(models) + 1) // 2
20
21     plt.figure(figsize=(16, 7*n_rows))
22  ▼ for i_model, model in enumerate(models):
23     plt.subplot(n_rows, 2, i_model+1)
24
25     # Значения гиперпараметров
26     C = model.get_params()['C']
27
28     # Визуализация опорных векторов model.support_vectors_
29  ▼ plt.scatter(
30         model.support_vectors[:, 0], model.support_vectors[:, 1],
31         edgecolor='black', s=np.abs(model._dual_coef_) * support_size,
32         alpha=0.5, zorder=10, facecolor='none', linewidths=3
33     )
34
35     # Визуализация остальных точек
36  ▼ plt.scatter(
37         X[:, 0], X[:, 1], c=Y, zorder=10, s=point_size, alpha=0.7,
38         cmap=plt.cm.Set3, edgecolor='black', linewidths=1.5
39     )
40
41     # печатаем коэффициенты для опорных векторов
42  ▼ for number, support_vector in enumerate(model.support_vectors_):
43     coef = model._dual_coef_[0, number]
44  ▼ plt.text(
45         support_vector[0] + 0.01, support_vector[1] + 0.01,
46         '{:.2}'.format(coef), fontsize=18,
47  ▼     bbox=dict(
48         facecolor='FF3300' if np.abs(coef) == C else 'w',
49         alpha=0.7
50     ))
51
52     # Определение границ графика
53     x_min = X[:, 0].min() - 0.5
54     x_max = X[:, 0].max() + 0.5
55     y_min = X[:, 1].min() - 0.5
56     y_max = X[:, 1].max() + 0.5
57
58     # Сетка точек в пространстве
59     XX, YY = np.mgrid[x_min:x_max:500j, y_min:y_max:500j]
```

```

60 # Значения решающей функции для этой сетки.
61 # Чтобы их получить, нужно передать матрицу размера (N, 2)
62 Z = model.decision_function(np.c_[XX.ravel(), YY.ravel()])
63
64 # Количество классов
65 K = 2 if len(Z.shape) == 1 else Z.shape[-1]
66
67 if K == 2:
68     # Ответы -- вектор. Переводим их к размеру сетки
69     Z = Z.reshape(XX.shape)
70
71     # Отрисовка решающей функции
72     if margin:
73         plt.pcolormesh(XX, YY, Z, cmap=plt.cm.RdBu)
74     else:
75         plt.pcolormesh(XX, YY, Z > 0, cmap=plt.cm.Set3)
76
77     # Отрисовка разделяющей прямой и разделяющей полосы
78     plt.contour(
79         XX, YY, Z, colors=['k', 'k', 'k'],
80         linestyles=['--', '-', '--'],
81         levels=[-1, 0, 1]
82     )
83 else:
84     # Отрисовка решающей функции
85     plt.pcolormesh(
86         XX, YY, np.argmax(Z, axis=1).reshape(XX.shape),
87         cmap=plt.cm.Set3
88     )
89
90 plt.xlim(x_min, x_max)
91 plt.ylim(y_min, y_max)
92 plt.xticks(())
93 plt.yticks(())
94
95 plt.title('$C$ = {}, kernel = {}'.format(
96     C, model.get_params()['kernel']
97 ))
98
99 plt.tight_layout()
100 plt.show()

```

started 08:39:27 2020-03-27, finished in 31ms

Визуализируем результаты для SVM с линейным ядром.

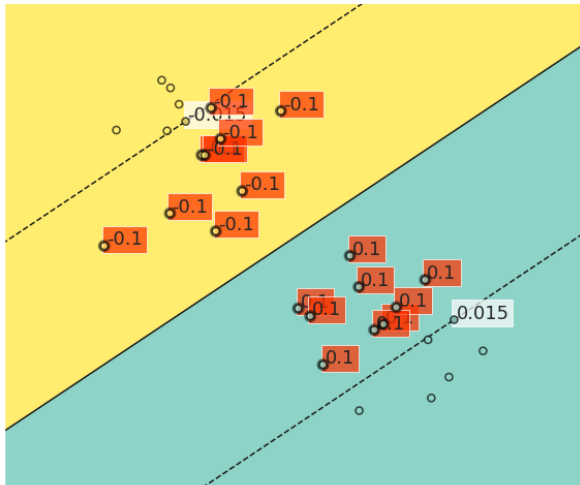
На графике желтым и синим цветами показана оценка классов по построенному классификатору. Черная сплошная линия отвечает за разделяющую поверхность, пунктирные линии -- границы разделяющей полосы, определяемой моделью SVM. Для каждого опорного объекта указан коэффициент λ_i , с которым он входит в решающее правило. Красный цвет метки имеют объекты-нарушители, которые выходят за пунктирную линию. Для таких объектов $\lambda_i = C$ -- установленный гиперпараметр. Пограничные объекты имеют белый цвет метки, причем их значение λ_i лежит в пределах от 0 до C , что полностью соответствует теории.

In [4]:

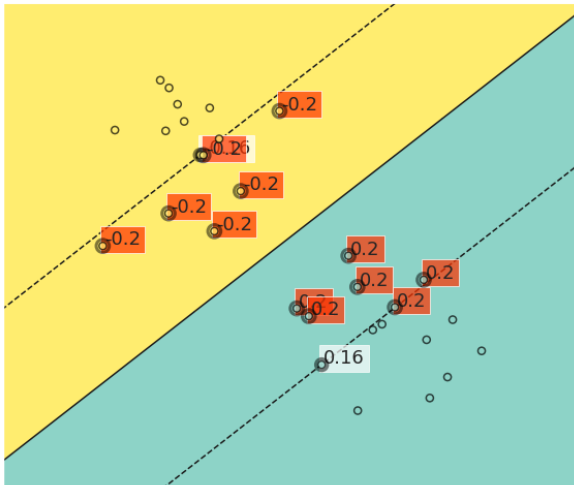
```
1 ▼ draw_graphics(  
2 ▼     [SVC(kernel='linear', C=C).fit(X, Y)  
3     for C in [0.1, 0.2, 0.5, 0.75, 1, 2, 5, 10]],  
4     X, Y  
5 )
```

started 08:39:28 2020-03-27, finished in 3.93s

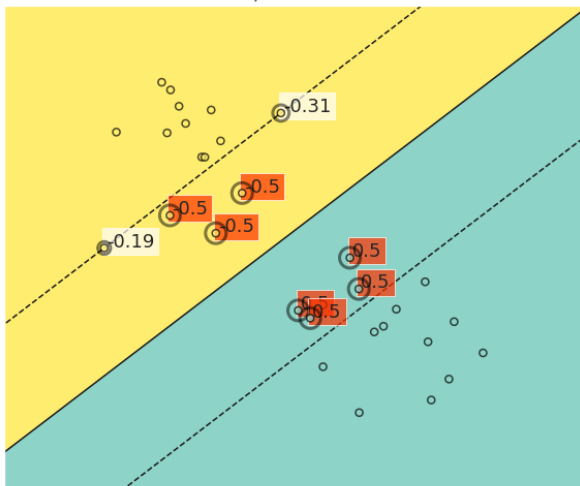
$C = 0.1$, kernel = linear



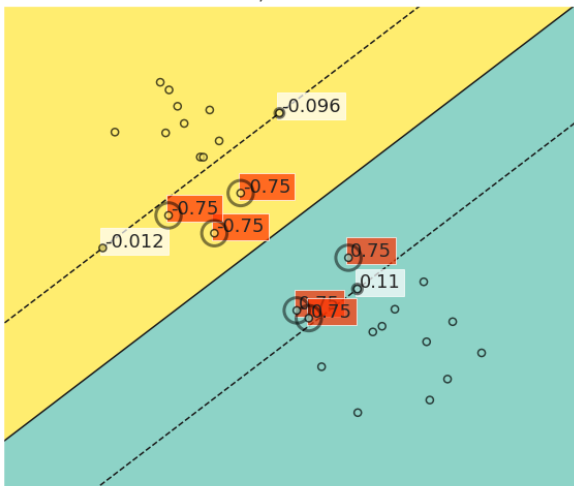
$C = 0.2$, kernel = linear



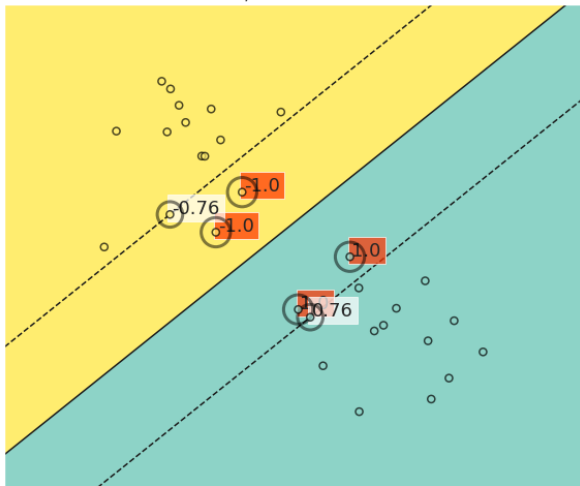
$C = 0.5$, kernel = linear



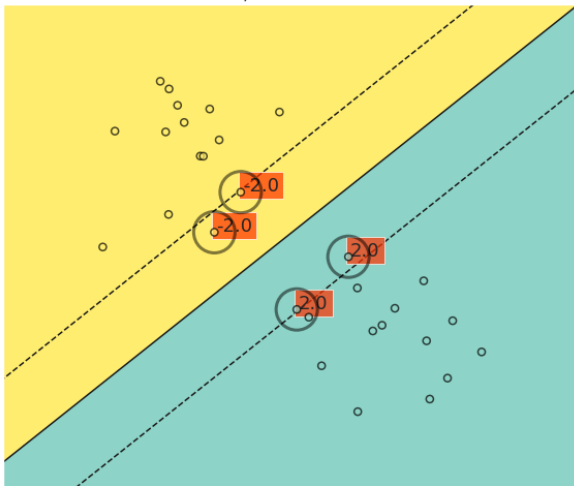
$C = 0.75$, kernel = linear



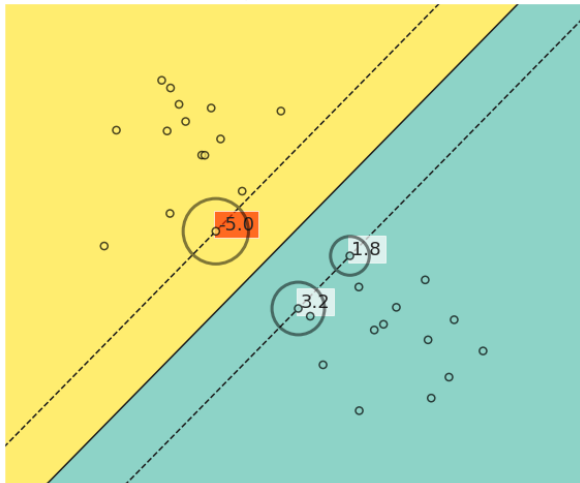
$C = 1$, kernel = linear



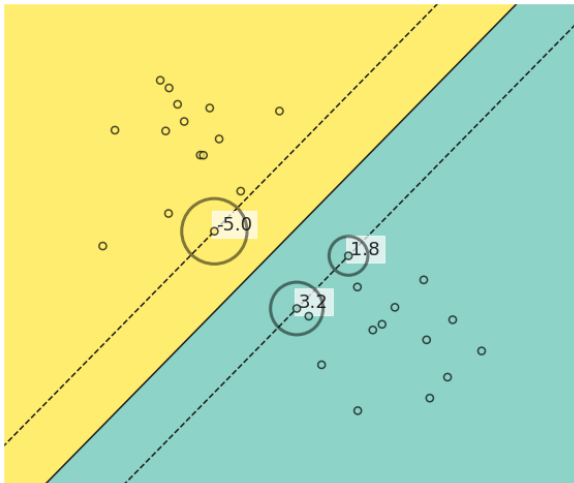
$C = 2$, kernel = linear



$C = 5$, kernel = linear



$C = 10$, kernel = linear



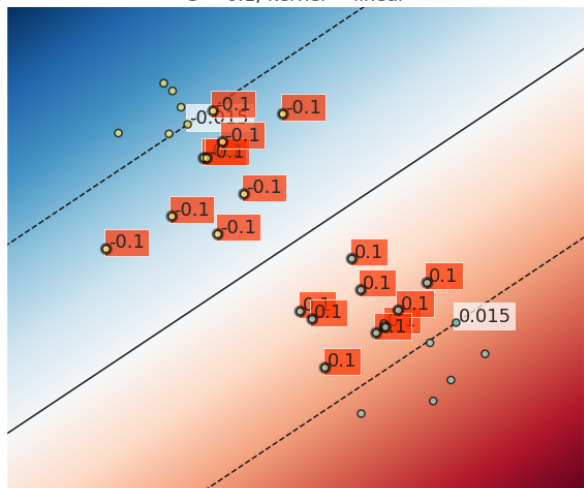
На следующих графиках цветом визуализирована решающая функция $f(x) = \langle \hat{\theta}, x \rangle + \hat{\theta}_0$

In [5]:

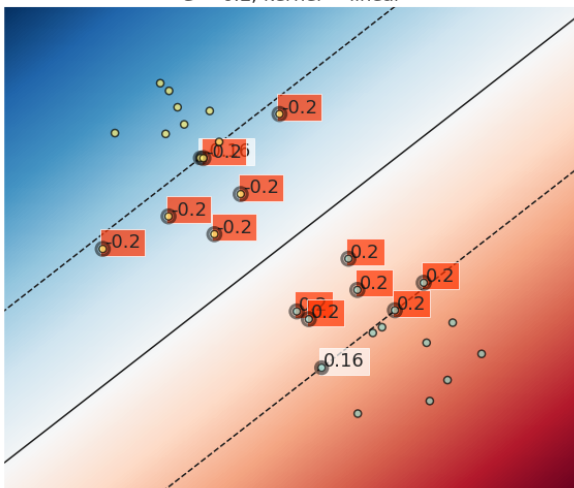
```
1 ▾ draw_graphics(  
2 ▾     [SVC(kernel='linear', C=C).fit(X, Y)  
3     for C in [0.1, 0.2, 0.5, 0.75, 1, 2, 5, 10]],  
4     X, Y, margin=True  
5 )
```

started 08:39:32 2020-03-27, finished in 4.03s

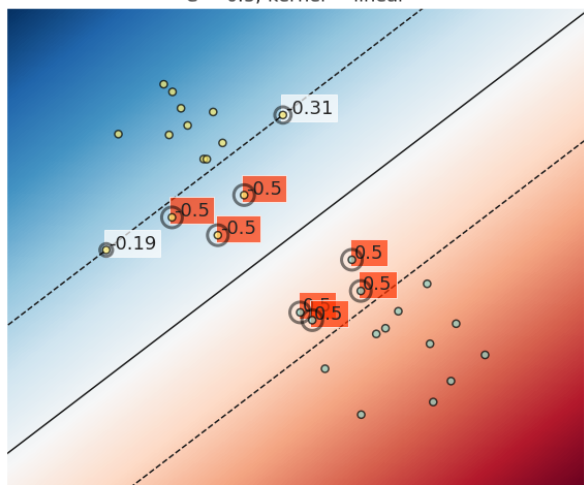
$C = 0.1$, kernel = linear



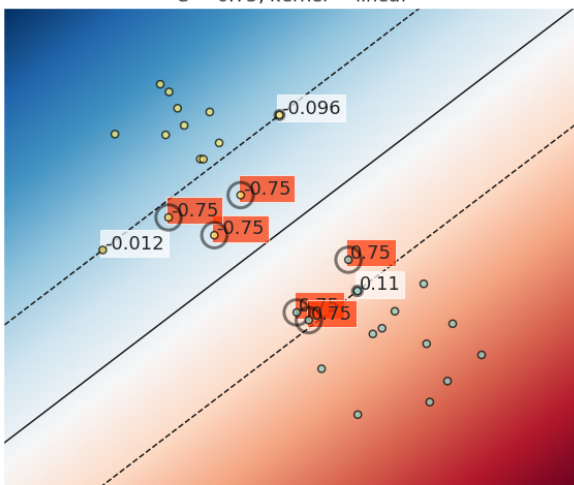
$C = 0.2$, kernel = linear



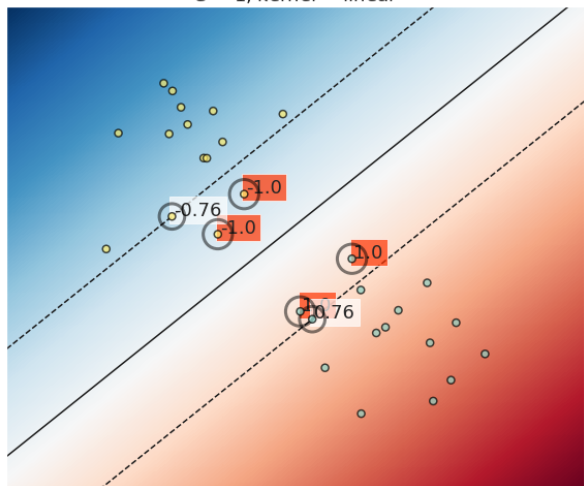
$C = 0.5$, kernel = linear



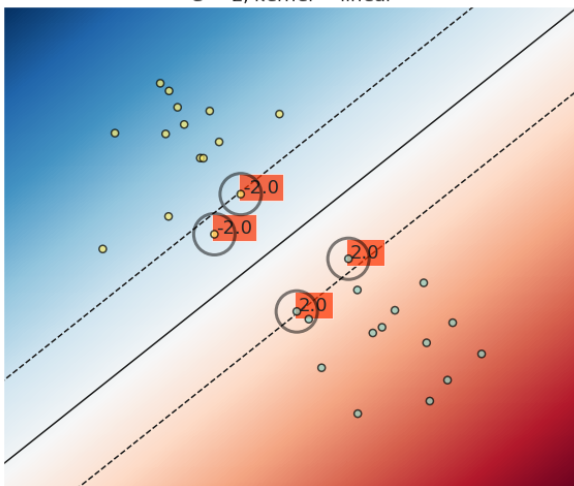
$C = 0.75$, kernel = linear



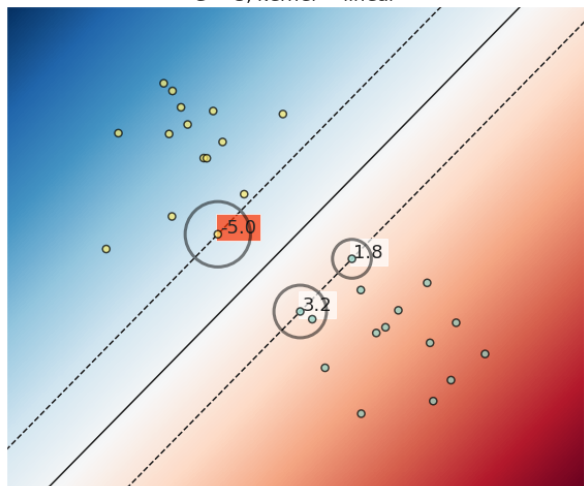
$C = 1$, kernel = linear



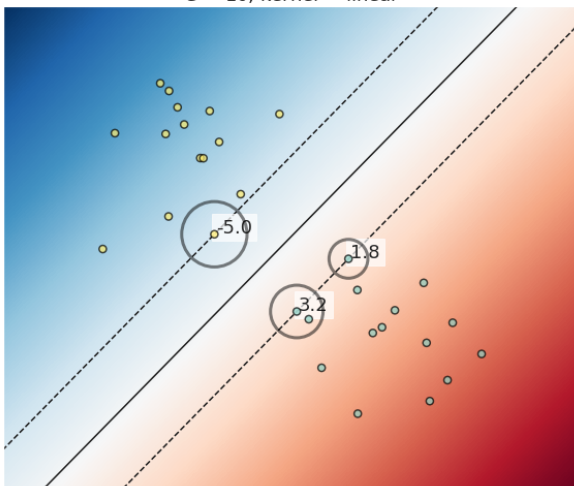
$C = 2$, kernel = linear



$C = 5$, kernel = linear



C = 10, kernel = linear



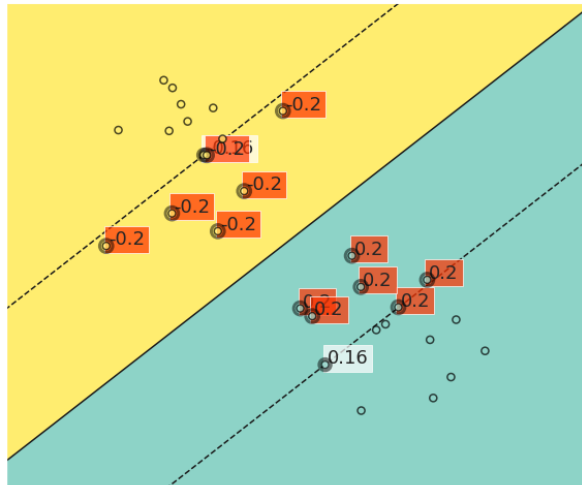
Аналогично визуализируем результаты для SVM с линейным, полиномиальным и радиально базисным ядром:

In [6]:

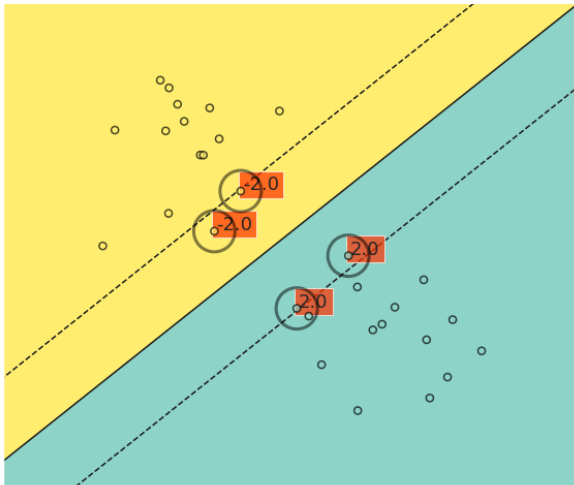
```
1 ▾ draw_graphics(  
2 ▾     [SVC(kernel=kernel, gamma=2, C=C).fit(X, Y)  
3     for kernel in ('linear', 'poly', 'rbf', 'sigmoid') for C in [0.2, 2]],  
4     X, Y  
5 )
```

started 08:39:36 2020-03-27, finished in 3.80s

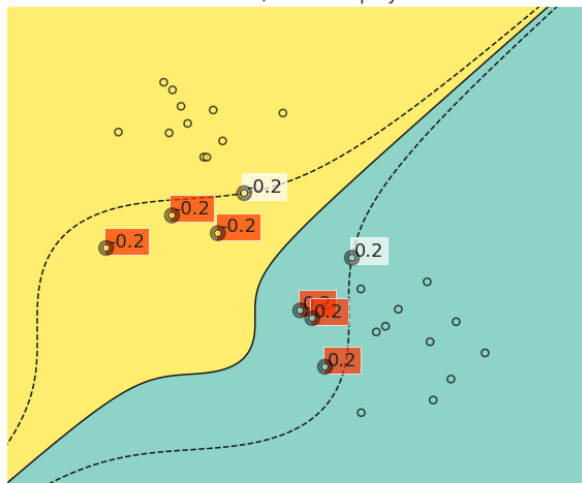
$C = 0.2$, kernel = linear



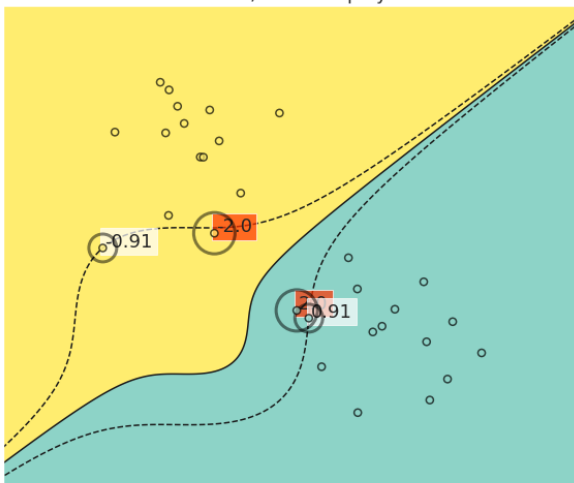
$C = 2$, kernel = linear



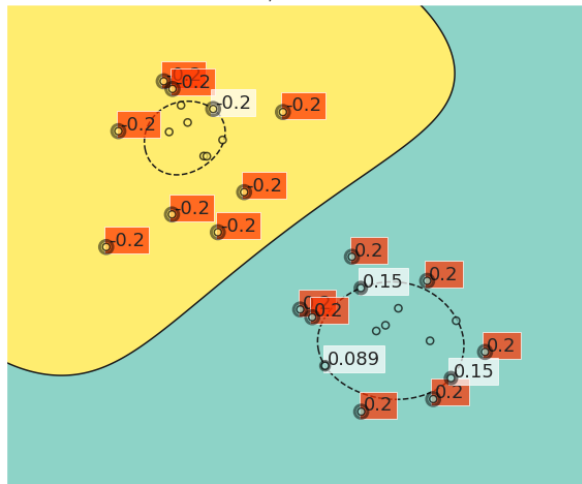
$C = 0.2$, kernel = poly



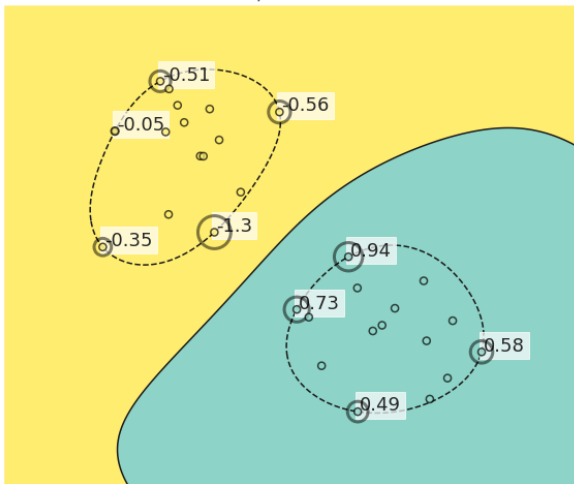
$C = 2$, kernel = poly



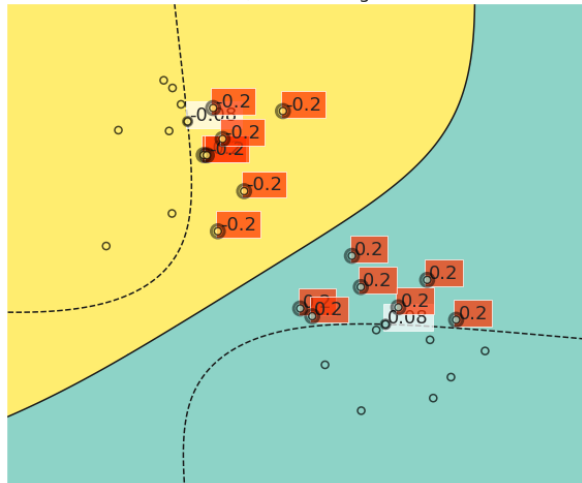
$C = 0.2$, kernel = rbf



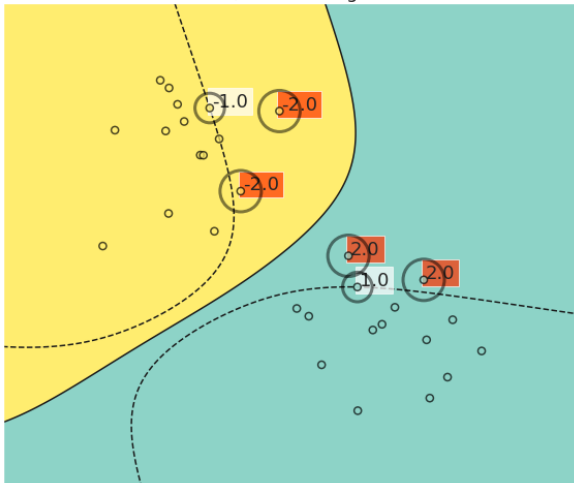
$C = 2$, kernel = rbf



$C = 0.2$, kernel = sigmoid



$C = 2$, kernel = sigmoid

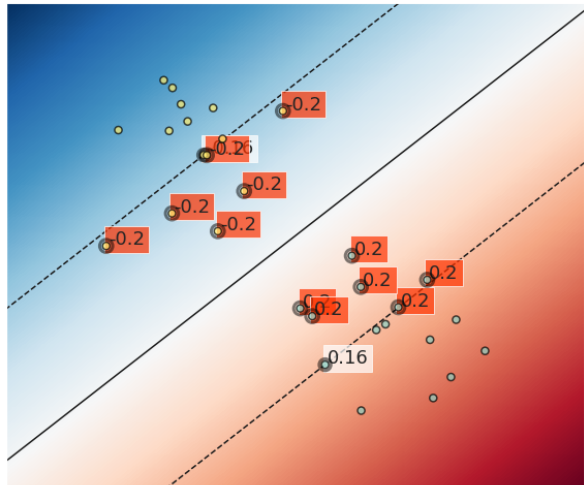


In [7]:

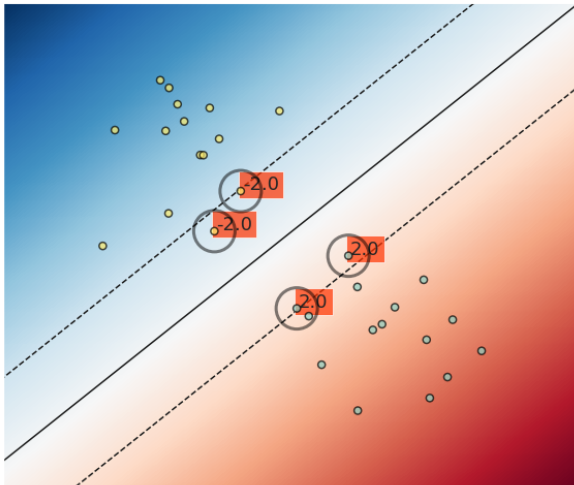
```
1 ▾ draw_graphics(  
2 ▾     [SVC(kernel=kernel, gamma=2, C=C).fit(X, Y)  
3     for kernel in ('linear', 'poly', 'rbf', 'sigmoid') for C in [0.2, 2]],  
4     X, Y, margin=True  
5 )
```

started 08:39:40 2020-03-27, finished in 4.23s

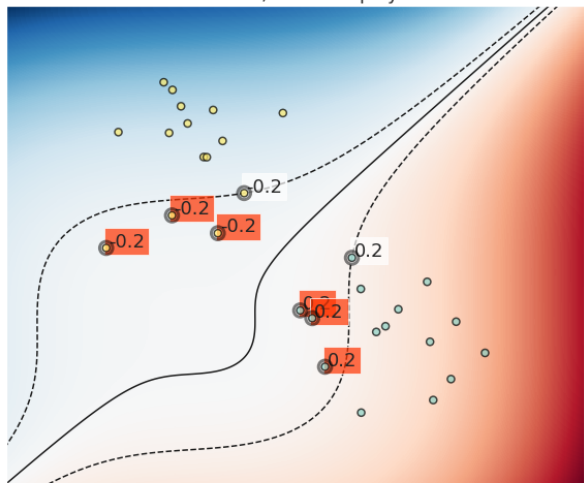
$C = 0.2$, kernel = linear



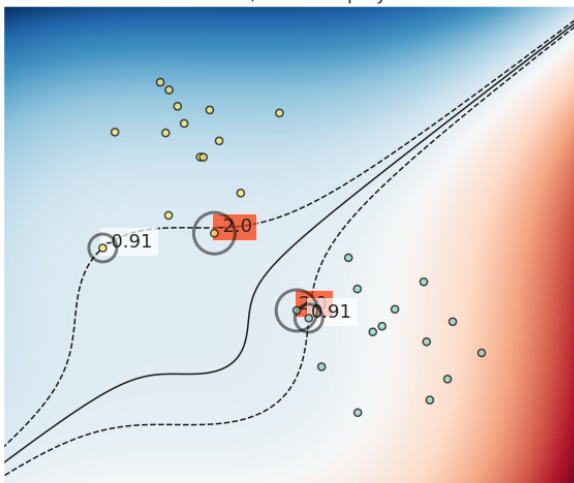
$C = 2$, kernel = linear



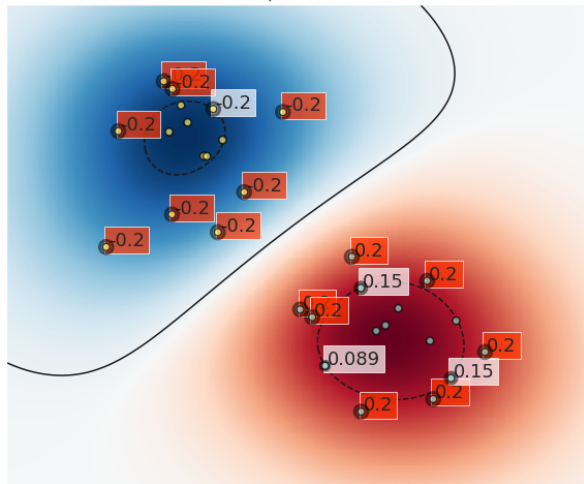
$C = 0.2$, kernel = poly



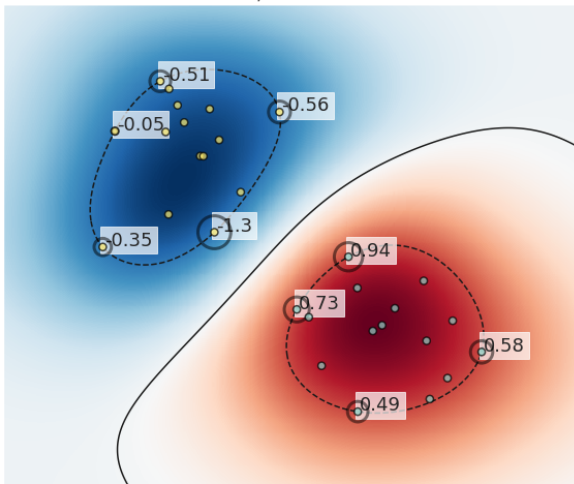
$C = 2$, kernel = poly



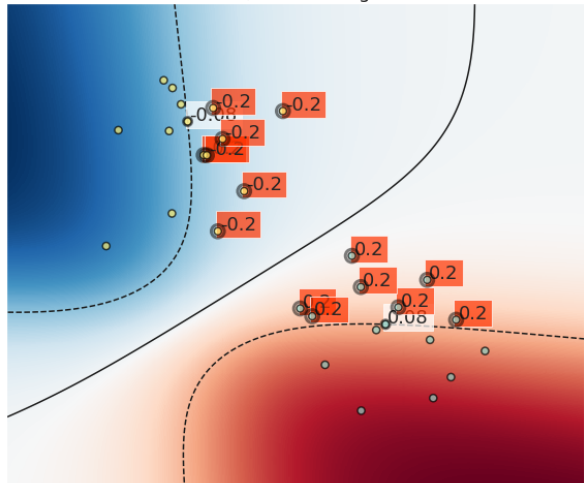
$C = 0.2$, kernel = rbf



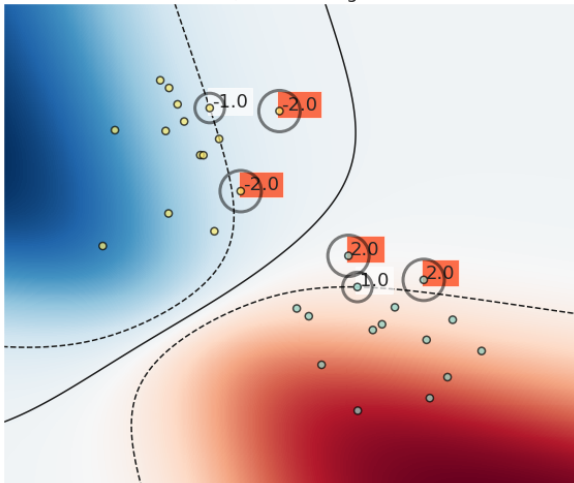
$C = 2$, kernel = rbf



$C = 0.2$, kernel = sigmoid



$C = 2$, kernel = sigmoid

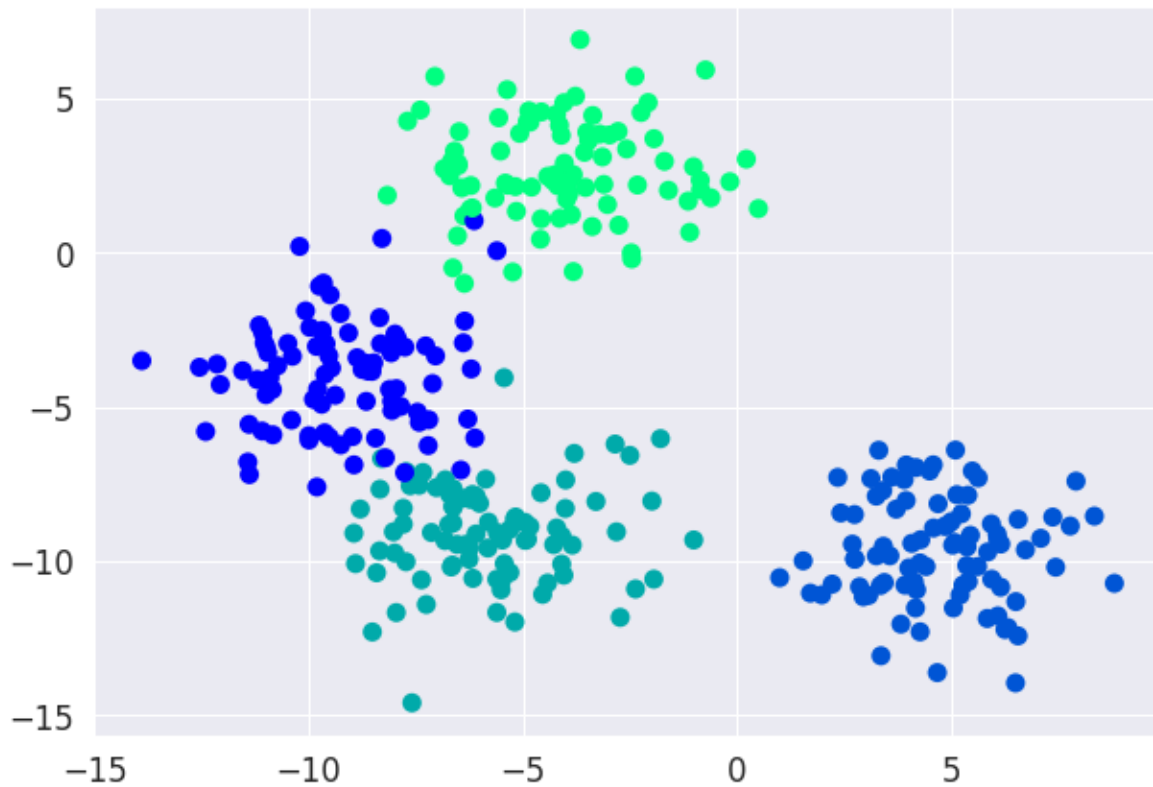


Рассмотрим случай нескольких классов. Сгенерируем выборку

In [8]:

```
1 X, Y = make_blobs(n_samples=350, n_features=2, centers=4,  
2                   cluster_std=1.7, random_state=21)  
3  
4 plt.figure(figsize=(10, 7))  
5 plt.scatter(X[:, 0], X[:, 1], c=Y, linewidths=0, s=100, cmap='winter')  
6 plt.grid()  
7 plt.show()
```

started 08:39:44 2020-03-27, finished in 359ms

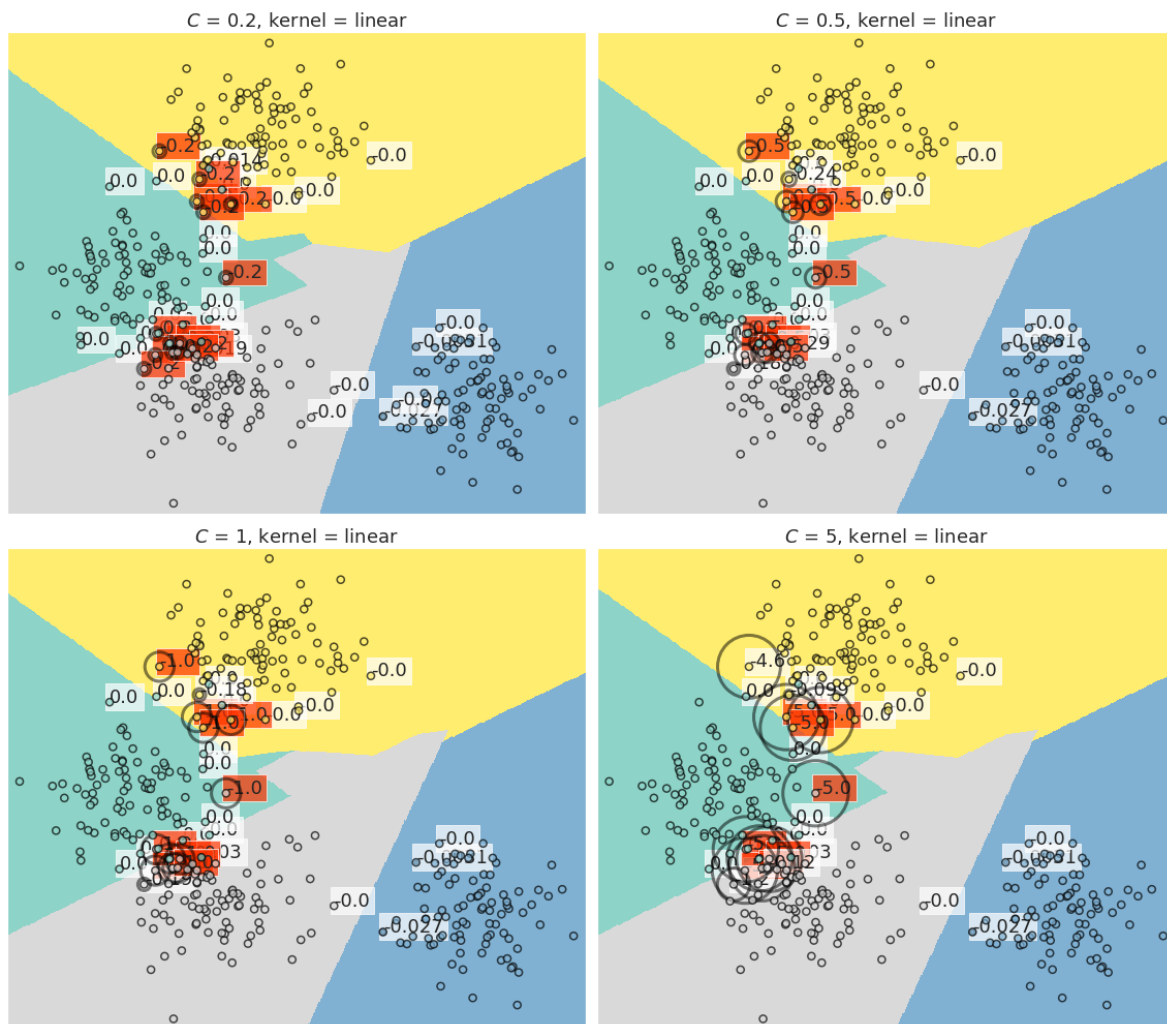


Посмотрим на то, как SVM классифицирует такую выборку

In [9]:

```
1 draw_graphics(  
2     [SVC(kernel='linear', C=C).fit(X, Y)  
3     for C in [0.2, 0.5, 1, 5]],  
4     X, Y  
5 )
```

started 08:39:44 2020-03-27, finished in 2.79s



Простые примеры регрессии

Сгенерируем выборку:

In [13]:

```
1 X = np.sort(5 * sps.uniform.rvs(size=(40, 1)), axis=0)  
2 Y = np.sin(X).ravel()  
3  
4 Y[::5] += 3 * (0.5 - sps.uniform.rvs(size=8))
```

started 08:41:38 2020-03-27, finished in 16ms

Инициализируем модели:

In [14]:

```
1  svr_lin = SVR(kernel='linear', C=100, gamma='auto').fit(X, Y)
2  svr_poly = SVR(kernel='poly', C=100, gamma='auto',
3                degree=3, epsilon=.1, coef0=1).fit(X, Y)
4  svr_rbf = SVR(kernel='rbf', C=100, gamma=0.1, epsilon=.1).fit(X, Y)
```

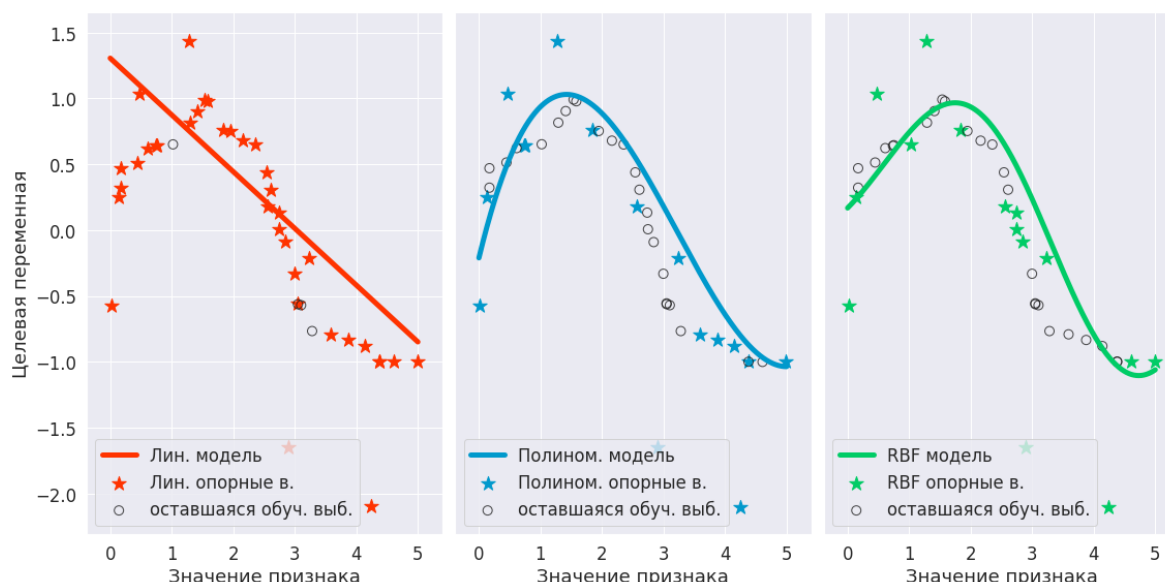
started 08:41:44 2020-03-27, finished in 4.37s

Визуализируем результаты:

In [15]:

```
1 svrs = [svr_lin, svr_poly, svr_rbf]
2 kernel_label = ['Лин.', 'Полином.', 'RBF']
3 model_color = ['#FF3300', '#0099CC', '#00CC66']
4
5 plt.figure(figsize=(16, 8))
6 for ix, svr in enumerate(svrs):
7
8     plt.subplot(1, 3, ix + 1)
9
10    grid = np.linspace(0, 5, 1001).reshape((-1, 1))
11    plt.plot(grid, svrs[ix].predict(grid), color=model_color[ix], lw=5,
12             label='{} модель'.format(kernel_label[ix]))
13    plt.scatter(X[svr.support_], y[svr.support_], marker='*',
14               color=model_color[ix], s=200,
15               label='{} опорные в.'.format(kernel_label[ix]))
16    plt.scatter(X[np.setdiff1d(np.arange(len(X)), svr.support_)],
17               y[np.setdiff1d(np.arange(len(X)), svr.support_)],
18               facecolor="none", edgecolor="k", s=80, alpha=0.8,
19               label='оставшаяся обуч. выб.')
```

started 08:41:50 2020-03-27, finished in 891ms



Классификация автодорожных знаков

Рассмотрим задачу классификации изображений (многоклассовая классификация на 42 класса), а именно дорожных знаков.

Загрузим данные с диска: <https://drive.google.com/file/d/115JLjH3nRg-FczEUM7mCx13olb6wTfiA/view>
(<https://drive.google.com/file/d/115JLjH3nRg-FczEUM7mCx13olb6wTfiA/view>)

In [0]:

```
1 from google.colab import drive
2
3 drive.mount('/content/drive')
```

In [0]:

```
1 ! unzip 'drive/My Drive/road_signs.zip'
2 ! rm '00_test_img_input/train/gt.csv'
```

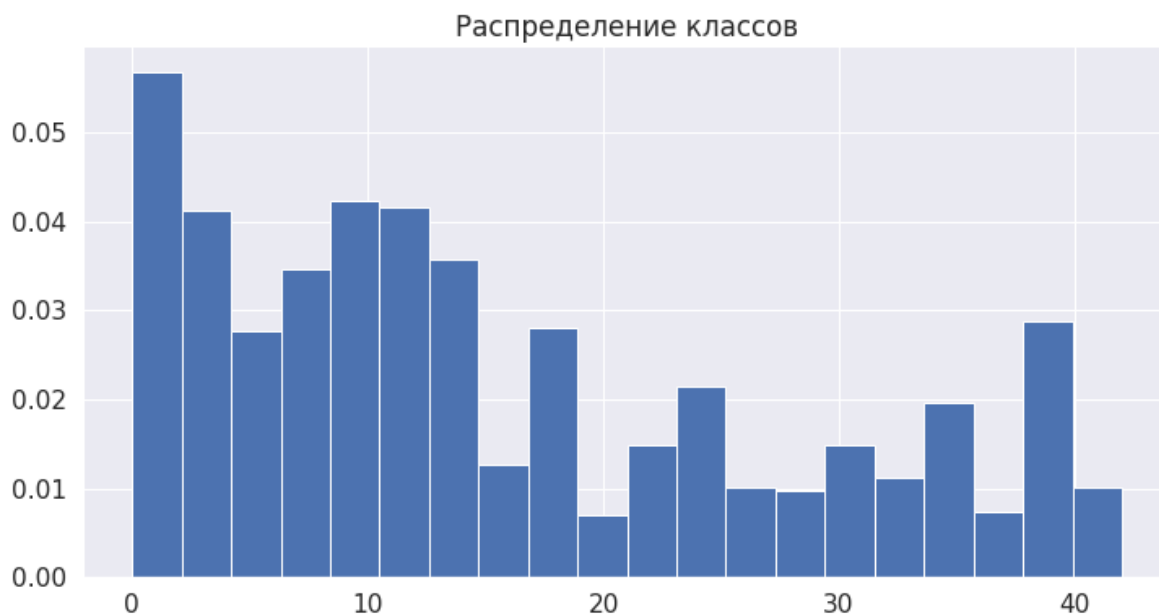
In [0]:

```
1 gt = pd.read_csv('00_test_img_gt/gt.csv', index_col='filename')
2 path = '00_test_img_input/train'
```

Посмотрим на распределение классов в выборке:

In [0]:

```
1 plt.figure(figsize=(12, 6))
2 gt.class_id.hist(density=True, bins=20)
3 plt.title('Распределение классов')
4 plt.show()
```



Видим, что явного дисбаланса классов нет.

Посмотрим на данные:

In [0]:

```
1 plt.figure(figsize=(8, 8))
2 gs = gridspec.GridSpec(4, 4)
3 gs.update(wspace=0.05, hspace=0.05)
4
5 for idx, img_path in enumerate(os.listdir(path)):
6
7     if idx == 16:
8         break
9
10    img = PIL.Image.open(os.path.join(path, img_path))
11
12    plt.subplot(gs[idx])
13    plt.imshow(img)
14    plt.xticks([])
15    plt.yticks([])
16
17 plt.show()
```



Теперь из каждой картинке получим [HOG дескрипторы](#)

(https://en.wikipedia.org/wiki/Histogram_of_oriented_gradients), которые являются по сути признаками, на которых мы и будем учить классификатор.

Схема вычисления HOG:

1. Конвертация цветной картинке в черно-белую, путем взятия яркости изображения, как линейной комбинации каналов. (Точнее компоненты Y в формате [YCbCr](#) (https://en.wikipedia.org/wiki/YCbCr#ITU-R_BT.601_conversion))
2. Вычисляются производные изображения I_x и I_y путем свертки с обычными разностными ядрами или ядрами Собеля:

$$D_x = (-1 \ 0 \ 1), D_y = (-1 \ 0 \ 1)^T,$$

$$S_x = \begin{pmatrix} -1 & 0 & 1 \\ -2 & 0 & 2 \\ -1 & 0 & 1 \end{pmatrix}, S_y = \begin{pmatrix} 1 & 2 & 1 \\ 0 & 0 & 0 \\ -1 & -2 & -1 \end{pmatrix}$$

3. Вычисляется модуль градиента и направление градиента по формуле:

$$|G| = \sqrt{I_x^2 + I_y^2}, \Theta = \text{atan2}(I_y, I_x)$$

где atan2 — знаковый арктангенс, принимающий значения от $-\pi$ до π .

4. Изображение разбивается на ячейки размером `cell_rows` × `cell_cols` пикселей и для каждой ячейки строится гистограмма направлений с `bin_count` корзин. Пиксель ячейки входит в одну из корзин гистограммы с весом, равным модулю градиента в данном пикселе. В простейшем случае ячейки не пересекаются.
5. Ячейки объединяются в блоки размером `block_row_cells` × `block_col_cells`, блоки могут пересекаться. Гистограммы различных ячеек в блоке конкатенируются в вектор v и нормируются:

$$v = \frac{v}{\sqrt{|v|^2 + \epsilon}}$$

6. Конкатенация векторов v из всех блоков является дескриптором изображения.

In [0]:

```
1  def extract_hog(img, bins=9):
2      """ Извлечение hog-дескрипторов из изображения. """
3
4      # приводим картинки к одному и тому же размеру с помощью бикубик интерпол.
5      img = img.resize((64, 64), resample=PIL.Image.BICUBIC)
6      img = np.array(img, dtype=np.float32)
7
8      # из трехканальной картинки получаем одноканальную, считая ее яркость
9      brightness = 0.299 * img[:, :, 0] + 0.587 * img[:, :, 1] \
10                  + 0.114 * img[:, :, 2]
11
12      # считаем норму и угол градиента
13      dx, dy = np.gradient(brightness)
14      G = np.sqrt(dx ** 2 + dy ** 2)
15      Theta = np.arctan2(dy, dx)
16
17      # считаем гистограммы ориентированных градиентов
18      cells = np.zeros((64, 9))
19      for i in range(8):
20          for j in range(8):
21              hist, _ = np.histogram(
22                  a=Theta[i * 8:(i + 1) * 8, j * 8:(j + 1) * 8],
23                  bins=bins,
24                  range=(-np.pi, np.pi),
25                  weights=G[i * 8:(i + 1) * 8, j * 8:(j + 1) * 8]
26              )
27
28              cells[8 * i + j] = hist
29
30      cells = cells.reshape(-1)
31
32      # нормировка
33      for i in range(0, 64 * 9, 9):
34          cells[i:i + 9] /= (np.linalg.norm(cells[i:i + 9], ord=2) + 1e-5)
35
36      return cells
```

Теперь для каждой картинки посчитаем ее гистограмму ориентированных градиентов (в моем случае это вектор размерности 576) и сохраним все в память:

In [0]:

```
1  X = np.zeros((39209, 576))
2  y = np.zeros(39209)
3
4  for idx, img_path in enumerate(os.listdir(path)):
5      img = PIL.Image.open(os.path.join(path, img_path)).convert('RGB')
6      X[idx] = extract_hog(img)
7      y[idx] = gt.loc[img_path].values[0]
```

Разделим выборку на обучающую и тестовую:

In [0]:

```
1 X_train, X_test, y_train, y_test = train_test_split(  
2     X, y, test_size=0.2, random_state=42  
3 )
```

В данном случае обычный SVC будет очень долго работать, поэтому воспользуемся LinearSVC :

In [0]:

```
1 clf = LinearSVC(max_iter=5000)  
2 clf.fit(X_train, y_train)  
3  
4 accuracy_score(y_test, clf.predict(X_test))
```

Out[17]:

0.979979597041571

Без нейросетей и даже без подбора гиперпараметров мы получили очень высокое качество! Именно так и классифицировали изображения, когда не было нейросетей. Напоминаю, что было 42 класса, а мы получили качество 98%.