

In [1]:

```
1 import numpy as np
2 from matplotlib import pyplot as plt
3 from matplotlib import animation as animation
4 import seaborn as sns
5
6 sns.set(font_scale=1.6, palette='RdBu')
```

## Методы оптимизации

Методы оптимизации широко используются в анализе данных и машинном обучении. Например, они применяются при обучении линейной регрессии с **lasso** или **elastic net** регуляризацией. Кроме того, на оптимизации функции потерь основано обучение любых нейронных сетей. Поэтому каждый, кто планирует заниматься анализом данных, должен знать основные методы оптимизации.

### Базовые элементы оптимизации.

1. Переменные - параметры, по которым требуется оптимизировать функцию.
2. Ограничения - границы, в которых могут варьироваться переменные.
3. Функция потерь - функция, которую минимизирует метод.

Постановка задачи оптимизации - определение функции потерь и переменных, по которым будет минимизироваться функция потерь и ограничений на эти переменные.

### Основные методы оптимизации.

Пусть задача оптимизации имеет вид  $f(\theta) \rightarrow \min_{\theta}$ , и  $\nabla_{\theta} f(\theta)$  - градиент функции  $f(\theta)$ .

Все методы, рассматриваемые здесь, являются итеративными. Они последовательно приближают текущее значение параметра  $\theta$  к оптимальному  $\theta^*$ .

#### 1. Градиентный спуск (GD - Gradient Descent).

$$\theta_{t+1} = \theta_t - \eta \nabla_{\theta} f(\theta_t)$$

Здесь изменение параметра  $\theta$  происходит по направлению антиградиента. Метод градиентного спуска основан на том факте, что антиградиент - направление наибольшего локального убывания функции. Поскольку, это свойство антиградиента локально, на каждом шаге антиградиент вычитается с заданным коэффициентом  $\eta$ , как правило, меньшим 1. Подбор  $\eta$  осуществляется пользователем.

#### 2. Покоординатный градиентный спуск.

Часто приходится оптимизировать функции от большого числа переменных. В таком случае можно не ждать вычисления градиента сразу по всем переменным, а обновлять их значения после вычисления частной производной по компонентам. Иначе говоря, шаг не по направлению градиента, а в направлении отдельных компонент градиента.

Таким образом, шаг покоординатного спуска можно записать так:

$$\theta_{t+1}^{(i)} = \theta_t^{(i)} - \eta \frac{\partial f}{\partial \theta_i} f(\theta_t).$$

### 3. Метод тяжелого шара (Momentum).

SGD считается относительно неплохим методом оптимизации и применяется на практике. Но у него есть ряд недостатков. К таким недостаткам относятся - застревание в локальных минимумах или седловых точках (при слишком маленьком learning rate), а также "пролетание" узких глобальных минимумов (при слишком большом learning rate).

Для частичного исправления этих недостатков используют метод momentum, в котором направление шага метода постепенно накапливается.

$$v_t = \gamma v_{t-1} + \eta \nabla_{\theta} f(\theta_t),$$

$$\theta_{t+1} = \theta_t - v_t.$$

Чем больше значение  $\gamma$ , тем больше метод ориентируется на свои предыдущие действия. Как правило, значение  $\gamma$  берут не менее 0.9.

### 4. Метод Ньютона-Рафсона (в англоязычной литературе - Newton's Method).

Этот метод, в отличие от всех рассмотренных выше, является методом оптимизации второго порядка. Для того, чтобы лучше понять его суть, рассмотрим его вывод. Для оптимизируемой функции  $f_{\theta}$  рассмотрим аппроксимацию 2 порядка:

$$f(\theta + h) \approx g(\theta + h) := f(\theta) + h^T f'(\theta) + \frac{1}{2} h^T f''(\theta) h.$$

Пусть мы хотим минимизировать  $g(\theta)$ . Тогда из необходимого условия локального минимума:

$$f'(\theta) + f''(\theta)h = 0,$$

$$h = -f''(\theta)^{-1} f'(\theta).$$

Шаг метода Ньютона-Рафсона выглядит так:

$$\theta_{t+1} = \theta_t - \eta \cdot f''(\theta)^{-1} f'(\theta).$$

## Эксперименты.

Нет универсального метода оптимизации, который всегда работает лучше, чем остальные. Поэтому для выбора наилучшего метода оптимизации и оптимальных гиперпараметров для него проводят ряд экспериментов. Ниже приведена визуализация нескольких экспериментов и сравнение скорости сходимости различных методов оптимизации, запущенных из одной точки.

Реализуем методы оптимизации.

In [2]:

```
1 def gradient_descent(theta0, func_grad, eta, iter_count=15_0):
2     '''
3     Градиентный спуск.
4
5     Параметры.
6     1) theta0 - начальное приближение theta,
7     2) func_grad - функция, задающая градиент оптимизируемой функции,
8     3) eta - скорость обучения,
9     4) iter_count - количество итераций метода.
10    '''
11
12    theta = theta0
13    history = [theta0]
14
15    for iter_id in range(iter_count):
16        theta = theta - eta * func_grad(theta)
17        history.append(theta)
18    return history
19
20
21 def coord_descent(theta0, func_grad, eta, iter_count=150):
22     '''
23     Покоординатный градиентный спуск.
24
25     Параметры.
26     1) theta0 - начальное приближение theta,
27     2) func_grad - функция, задающая градиент оптимизируемой функции,
28     3) eta - скорость обучения,
29     4) iter_count - количество итераций метода.
30    '''
31
32    d = len(theta0)
33
34    theta = theta0
35    history = [theta0.copy()]
36
37    for iter_id in range(iter_count):
38        for coord_id in range(d):
39            theta[coord_id] = theta[coord_id] \
40                - eta * func_grad(theta)[coord_id]
41            history.append(theta.copy())
42    return history
43
44
45 def momentum(theta0, func_grad, eta, gamma, iter_count=150):
46     '''
47     Метод тяжелого Шарика.
48
49     Параметры.
50     1) theta0 - начальное приближение theta,
51     2) func_grad - функция, задающая градиент оптимизируемой функции,
52     3) eta - скорость обучения,
53     4) gamma - параметр инерции,
54     5) iter_count - количество итераций метода.
55    '''
56
57    theta = theta0
58    history = [theta0]
59    v = theta0
```

```

60
61     for iter_id in range(iter_count):
62         v = gamma * v + eta * func_grad(theta)
63         theta = theta - v
64         history.append(theta)
65     return history
66
67
68 def newton(theta0, func_grad, func_hessian, eta, iter_count=150):
69     '''
70     Метод Ньютона-Равсона.
71
72     Параметры.
73     1) theta0 - начальное приближение theta,
74     2) func_grad - функция, задающая градиент оптимизируемой функции,
75     3) func_hessian - функция, задающая гессиан оптимизируемой функции,
76     3) eta - скорость обучения,
77     4) iter_count - количество итераций метода.
78     '''
79
80     theta = theta0
81     history = [theta0]
82
83     for iter_id in range(iter_count):
84         theta = theta \
85             - eta * (np.linalg.inv(func_hessian(theta)) @ func_grad(theta))
86         history.append(theta)
87     return history

```

Реализуем функции, которые будем оптимизировать.

In [3]:

```

1 def square_sum(x):
2     return 5 * x[0]**2 + x[1]**2
3
4 def square_sum_grad(x):
5     return np.array([10, 2]) * x
6
7 def square_sum_hessian(x):
8     return np.diag([10, 2])

```

Создадим директорию, в которой будем хранить визуализацию экспериментов.

In [4]:

```

1 !rm -rf saved_gifs
2 !mkdir saved_gifs

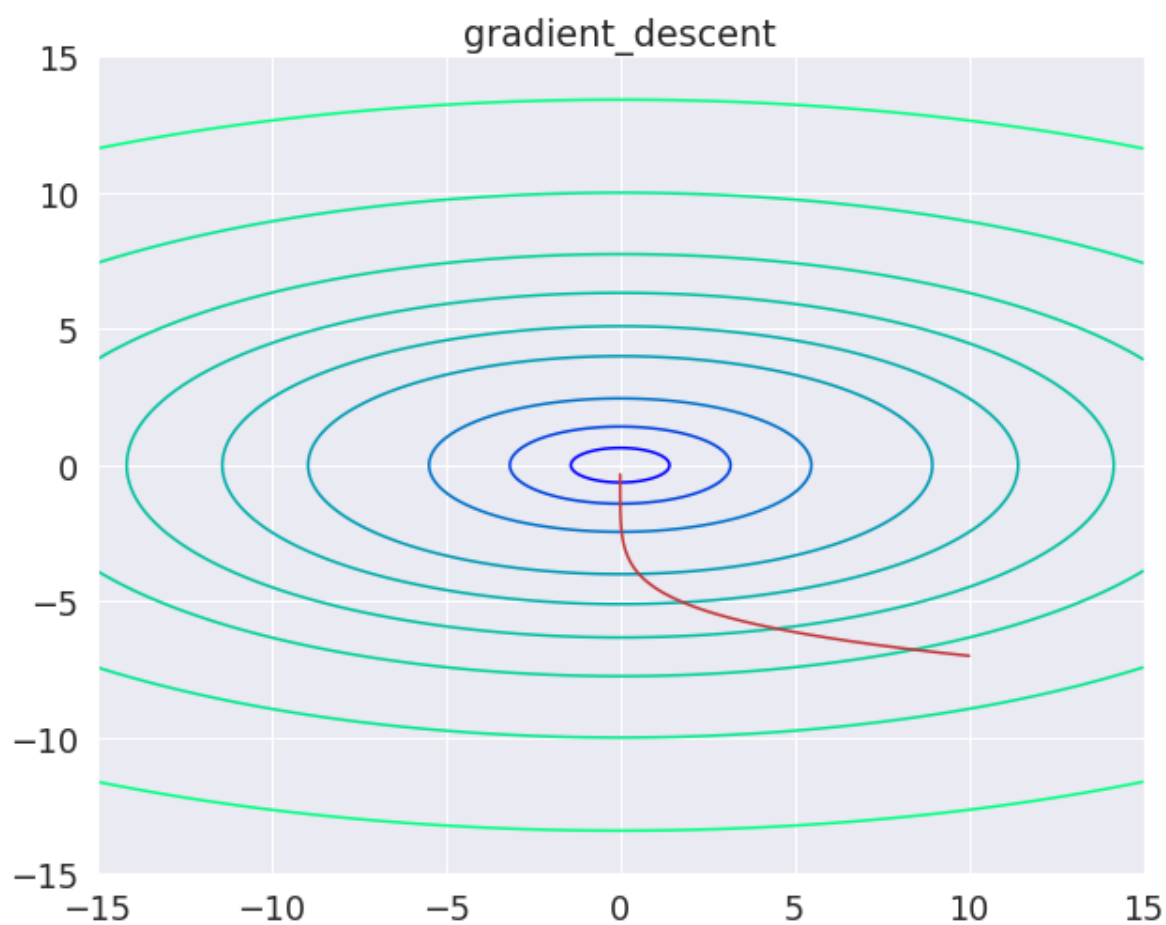
```

In [5]:

```
1 def make_experiment(func, trajectory, graph_title):
2     '''
3     Функция, которая для заданной функции рисует её линии уровня, а также
4     траекторию сходимости метода оптимизации.
5
6     Параметры.
7     1) func - оптимизируемая функция,
8     2) trajectory - траектория метода оптимизации,
9     3) graph_name - заголовок графика.
10    '''
11
12    fig, ax = plt.subplots(figsize=(10, 8))
13    xdata, ydata = [], []
14    ln, = plt.plot([], [])
15
16    mesh = np.linspace(-15.0, 15.0, 300)
17    X, Y = np.meshgrid(mesh, mesh)
18    Z = np.zeros((len(mesh), len(mesh)))
19
20    for coord_x in range(len(mesh)):
21        for coord_y in range(len(mesh)):
22            Z[coord_x][coord_y] = func(np.array((mesh[coord_x],
23                                                    mesh[coord_y])))
24
25    def init():
26        ax.contour(X, Y, np.log(Z),
27                   np.log([2, 10, 30, 80, 130, 200, 300, 500, 900]),
28                   cmap='winter')
29        ax.set_title(graph_title)
30        return ln,
31
32    def update(frame):
33        xdata.append(trajectory[frame][0])
34        ydata.append(trajectory[frame][1])
35        ln.set_data(xdata, ydata)
36        return ln,
37
38    ani = animation.FuncAnimation(
39        fig, update, frames=range(len(trajectory)),
40        init_func=init, repeat=True
41    )
42    ani.save(f'saved_gifs/{graph_title}.gif',
43            writer='imagemagick', fps=5)
```

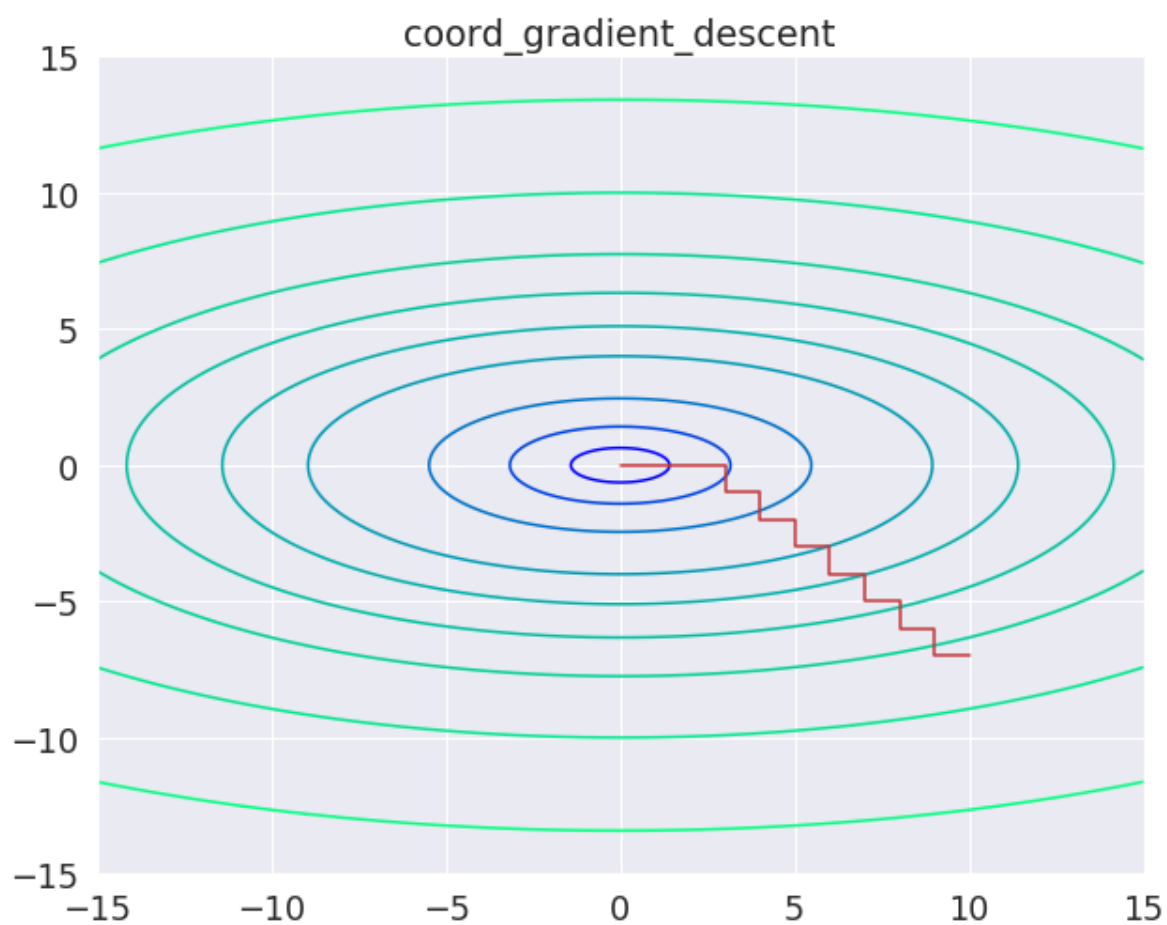
In [6]:

```
1 gd_trajectory = gradient_descent(np.array((10, -7)), square_sum_grad, 0.01)  
2 make_experiment(square_sum, gd_trajectory, 'gradient_descent')
```



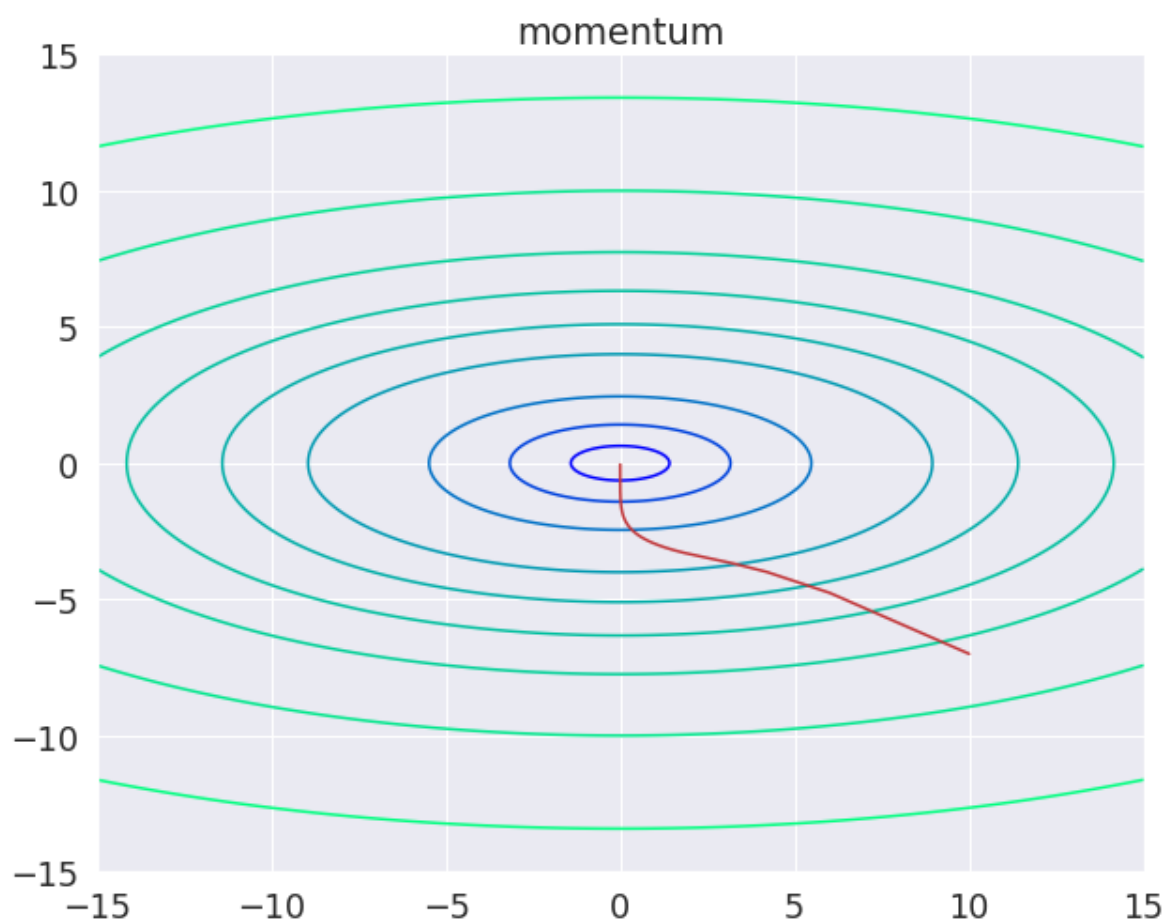
In [7]:

```
1 coord_gd_trajectory = coord_descent(np.array((10, -7)), square_sum_grad, 0.01,  
2 make_experiment(square_sum, coord_gd_trajectory, 'coord_gradient_descent'))
```



In [8]:

```
1 momentum_trajectory = momentum(np.array((10, -7)), square_sum_grad, 0.01, 0.3)
2 make_experiment(square_sum, momentum_trajectory, 'momentum')
```





In [9]:

```
1 newton_trajectory = newton(np.array((10, -7)), square_sum_grad, square_sum_hess  
2 make_experiment(square_sum, newton_trajectory, 'newton')
```

