

In []:

```
1 # ! apt-get install imagemagick
```

In []:

```
1 import numpy as np
2 import seaborn as sns
3
4 from IPython.display import Image, clear_output
5 from matplotlib import pyplot as plt
6 from matplotlib.animation import FuncAnimation, ImageMagickFileWriter
7
8 sns.set(font_scale=1.6, palette='RdBu')
```

Продвинутые методы оптимизации

Вы уже познакомились с основными методами оптимизации, которые широко используются в классическом машинном обучении. С развитием нейронных сетей и активным внедрением нейросетевого подхода, методы оптимизации стали ещё более актуальными. Но стандартные методы оптимизации, SGD и метод тяжёлого шара, имеют ряд недостатков, из-за чего их редко применяют в чистом виде. Для обучения современных нейросетей используют более продвинутые методы.

Ключевая особенность всех рассматриваемых ниже методов в том, что они являются адаптивными. Т.е. для различных параметров оптимизируемой функции обновление происходит с различной скоростью.

Пусть задача оптимизации имеет вид $f(x) \longrightarrow \min_x$, и $\nabla_x f(x)$ — градиент функции $f(x)$.

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.

1. SGD

Обычный и стохастический градиентный спуск.

$$x_{t+1} = x_t - \eta v_t,$$

где $v_t = \nabla f(x_t)$ — аналогия со скоростью.

In []:

```
1 def sgd(init_parameters, func_grad, lr, n_iter):
2     '''
3     Метод оптимизации SGD.
4
5     Параметры:
6     - parameters - начальное приближение параметров,
7     - func_grad - функция, задающая градиент оптимизируемой функции,
8     - lr - скорость обучения,
9     - n_iter - количество итераций метода.
10
11     Возвращает историю обновлений параметров.
12     '''
13
14     parameters = init_parameters.copy()
15     history = [parameters.copy()]
16
17     for i in range(n_iter):
18         diff = -lr * func_grad(parameters)
19         parameters += diff
20         history.append(parameters.copy())
21
22     return history
```

2. SGD + Momentum

Сгладим градиент, используя информацию о том, как градиент изменялся раньше. Физическая аналогия — добавляем инерцию.

$$x_{t+1} = x_t + v_t,$$

где $v_t = \mu v_{t-1} - \eta \nabla f(x_t)$ — сглаживаем градиенты.

In []:

```
1 def sgd_momentum(init_parameters, func_grad, lr, mu, n_iter):
2     '''
3     Метод оптимизации SGD Momentum.
4
5     Параметры:
6     - parameters - начальное приближение параметров,
7     - func_grad - функция, задающая градиент оптимизируемой функции,
8     - lr - скорость обучения,
9     - mu - коэффициент сглаживания,
10    - n_iter - количество итераций метода.
11
12    Возвращает историю обновлений параметров.
13    '''
14    parameters = init_parameters.copy()
15    history = [parameters.copy()]
16    diff = np.zeros_like(parameters)
17
18    for i in range(n_iter):
19        diff = mu * diff - lr * func_grad(parameters)
20        parameters += diff
21        history.append(parameters.copy())
22
23    return history
```

3. Adagrad

Adagrad — один из самых первых адаптивных методов оптимизации.

Во всех изученных ранее методах есть необходимость подбирать шаг метода (коэффициент η). На каждой итерации все компоненты градиента оптимизируемой функции домножаются на одно и то же число η . Но использовать одно значение η для всех параметров не оптимально, так как они имеют различные распределения и оптимизируемая функция изменяется с совершенно разной скоростью при небольших изменениях разных параметров.

Поэтому гораздо логичнее **изменять значение каждого параметра с индивидуальной скоростью**. При этом, *чем с большей степени от изменения параметра меняется значение оптимизируемой функции, тем с меньшей скоростью стоит обновлять этот параметр*. Иначе высок шанс расходимости метода. Получить такой результат удастся, если разделить градиент на сумму квадратов скорости изменений параметров.

Пусть $x^{(i)}$ — i -я компонента вектора x .

$$x_{t+1,i} = x_{t,i} - \frac{\eta}{\sqrt{g_{t,i} + \epsilon}} \cdot \nabla f_i(x_t)$$
$$g_t = g_{t-1} + \nabla f(x_t) \odot \nabla f(x_t)$$

В матрично-векторном виде шаг алгоритма можно переписать так:

$$x_{t+1} = x_t - \frac{\eta}{\sqrt{g_t + \epsilon}} \odot \nabla f(x_t).$$

Здесь \odot обозначает произведение Адамара, т.е. поэлементное перемножение векторов.

In []:

```
1 def adagrad(init_parameters, func_grad, lr, eps, n_iter):
2     '''
3     Метод оптимизации Adagrad.
4
5     Параметры:
6     - parameters - начальное приближение параметров,
7     - func_grad - функция, задающая градиент оптимизируемой функции,
8     - lr - скорость обучения,
9     - eps - минимальное значение нормирующего члена,
10    - n_iter - количество итераций метода.
11
12    Возвращает историю обновлений параметров.
13    '''
14    parameters = init_parameters.copy()
15    history = [parameters.copy()]
16    norm = np.zeros_like(parameters)
17
18    for iter_id in range(n_iter):
19        cur_grad = func_grad(parameters)
20        norm += cur_grad ** 2
21        parameters -= lr * cur_grad / np.sqrt(norm + eps)
22        history.append(parameters.copy())
23
24    return history
```

4. RMSProp

Алгоритм RMSProp основан на той же идее, что и алгоритм Adagrad — адаптировать learning rate отдельно для каждого параметра $\theta^{(i)}$. Однако Adagrad имеет серьёзный недостаток. Он с одинаковым весом учитывает квадраты градиентов как с самых первых итераций, так и с самых последних. Хотя, на самом деле, наибольшую значимость имеют модули градиентов на последних нескольких итерациях.

Для этого предлагается использовать **экспоненциальное сглаживание**.

$$x_{t+1} = x_t - \frac{\eta}{\sqrt{g_t + \varepsilon}} \odot \nabla f(x_t).$$
$$g_t = \mu g_{t-1} + (1 - \mu) \nabla f(x_t) \odot \nabla f(x_t)$$

Стандартные значения гиперпараметров: $\mu = 0.9, \eta = 0.001$.

In []:

```
1 def rmsprop(init_parameters, func_grad, lr, mu, eps, n_iter):
2     '''
3     Метод оптимизации RMSProp.
4
5     Параметры:
6     - parameters - начальное приближение параметров,
7     - func_grad - функция, задающая градиент оптимизируемой функции,
8     - lr - скорость обучения,
9     - mu - коэффициент сглаживания,
10    - eps - минимальное значение нормирующего члена,
11    - n_iter - количество итераций метода.
12
13    Возвращает историю обновлений параметров.
14    '''
15
16    parameters = init_parameters.copy()
17    history = [parameters.copy()]
18    norm = np.zeros_like(parameters)
19
20    for iter_id in range(n_iter):
21        cur_grad = func_grad(parameters)
22        norm = mu * norm + (1 - mu) * cur_grad ** 2
23        parameters -= lr * cur_grad / np.sqrt(norm + eps)
24        history.append(parameters.copy())
25
26    return history
```

5. Adam

Этот алгоритм совмещает в себе 2 идеи:

- идею алгоритма Momentum о *накапливании градиента*,
- идею методов Adadelta и RMSProp об *экспоненциальном сглаживании* информации о предыдущих значениях квадратов градиентов.

Благодаря использованию этих двух идей, метод имеет 2 преимущества над большей частью методов первого порядка, описанных выше:

1. Он обновляет все параметры θ не с одинаковым `learning rate`, а выбирает для каждого θ_i индивидуальный `learning rate`, что *позволяет учитывать разреженные признаки с большим весом*.
2. Adam за счёт применения экспоненциального сглаживания к градиенту *работает более устойчиво в окрестности оптимального значения параметра θ^** , чем методы, использующие градиент в точке x_t , не накапливая значения градиента с прошлых шагов.

Формулы шага алгоритма выглядят так:

$$\begin{aligned}v_{t+1} &= \beta v_t + (1 - \beta) \nabla x(x_t) \\ g_t &= \mu g_{t-1} + (1 - \mu) \nabla x(x_t) \odot \nabla x(x_t)\end{aligned}$$

Чтобы эти оценки не были смещёнными, нужно их отнормировать:

$$\begin{aligned}\hat{v}_{t+1} &= \frac{v_{t+1}}{1 - \beta^t}, \\ \hat{g}_t &= \frac{g_t}{1 - \mu^t}.\end{aligned}$$

Тогда получим итоговую формулу шага:

$$x_{t+1} = x_t - \frac{\eta}{\sqrt{\hat{g}_t + \varepsilon}} \odot \hat{v}_{t+1}.$$

In []:

```
1 def adam(init_parameters, func_grad, eps, lr, beta, mu, n_iter):
2     '''
3     Adam.
4
5     Параметры.
6     1) theta0 - начальное приближение theta,
7     2) func_grad - функция, задающая градиент оптимизируемой функции,
8     3) eps - мин. возможное значение знаменателя,
9     4) lr - скорость обучения,
10    5) beta - параметр экспоненциального сглаживания,
11    6) mu - параметр экспоненциального сглаживания,
12    7) n_iter - количество итераций метода.
13    '''
14
15    parameters = init_parameters.copy()
16    history = [parameters.copy()]
17    cumulative_m = np.zeros_like(parameters)
18    cumulative_v = np.zeros_like(parameters)
19    pow_beta, pow_mu = beta, mu
20
21    for iter_id in range(n_iter):
22        current_grad = func_grad(parameters)
23        cumulative_m = beta * cumulative_m + (1 - beta) * current_grad
24        cumulative_v = mu * cumulative_v + (1 - mu) * current_grad ** 2
25
26        scaled_m = cumulative_m / (1 - pow_beta)
27        scaled_v = cumulative_v / (1 - pow_mu)
28        parameters = parameters - lr * cumulative_m / (np.sqrt(cumulative_v) +
29        history.append(parameters)
30
31        pow_beta *= beta
32        pow_mu *= mu
33
34    return history
```

Статьи.

Для того, чтобы подробнее познакомиться с представленными выше методами, мотивацией их авторов и теоретическими оценками сходимости, можно прочитать оригинальные статьи.

- Adagrad: <http://www.jmlr.org/papers/volume12/duchi11a/duchi11a.pdf>
(<http://www.jmlr.org/papers/volume12/duchi11a/duchi11a.pdf>),
- Adam: <https://arxiv.org/abs/1412.6980> (<https://arxiv.org/abs/1412.6980>).

Как можно заметить, для нейросетей мы рассматриваем *только методы оптимизации первого порядка*. Это связано с тем, что эффективные архитектуры нейронных сетей имеют большое количество параметров, из-за чего методы второго порядка, требующие на одну итерацию $O(d^2)$ памяти и выполняющие $O(d^3)$ операций, работают слишком долго и их преимущество в количестве итераций до сходимости утрачивает смысл.

Реализуем функции, которые будем оптимизировать.

In []:

```
1 def square_sum(x):
2     ''' f(x, y) = x^2 + y^2 '''
3
4     return 5 * x[0]**2 + x[1]**2
5
6 def square_sum_grad(x):
7     ''' grad f(x, y) = (10x, 2y) '''
8
9     return np.array([10, 2]) * x
10
11
12 def complex_sum(x):
13     ''' f(x, y) = (x-3)^2 + 8(y-5)^4 + sqrt(x) + sin(xy)'''
14
15     return (x[0]-3)**2 + 8 * (x[1]-5)**4 + x[0]**0.5 + np.sin(x[0]*x[1])
16
17 def complex_sum_grad(x):
18     ''' grad f(x, y) = (2(x-3) + 1/(2 sqrt(x)) + ycos(xy), 32(y-5)^3 + xcos(xy)
19
20     partial_x = 2*(x[0]-3) + 0.5*x[0]**(-0.5) + x[1]*np.cos(x[0]*x[1])
21     partial_y = 32*(x[1]-5)**3 + x[0]*np.cos(x[0]*x[1])
22
23     return np.array([partial_x, partial_y])
```

Создадим директорию, в которой будем хранить визуализацию экспериментов.

In []:

```
1 !rm -rf saved_gifs
2 !mkdir saved_gifs
```

In []:

```
1 def make_experiment(func, trajectory, graph_title,
2                     min_y=-7, max_y=7, min_x=-7, max_x=7):
3     ...
4     Функция, которая для заданной функции рисует её линии уровня,
5     а также траекторию сходимости метода оптимизации.
6
7     Параметры.
8     1) func - оптимизируемая функция,
9     2) trajectory - траектория метода оптимизации,
10    3) graph_name - заголовок графика.
11    ...
12
13    fig, ax = plt.subplots(figsize=(10, 8))
14    xdata, ydata = [], []
15    ln, = plt.plot([], [])
16
17    mesh_x = np.linspace(min_x, max_x, 300)
18    mesh_y = np.linspace(min_y, max_y, 300)
19    X, Y = np.meshgrid(mesh_x, mesh_y)
20    Z = np.zeros((len(mesh_x), len(mesh_y)))
21
22    for coord_x in range(len(mesh_x)):
23        for coord_y in range(len(mesh_y)):
24            Z[coord_y, coord_x] = func(
25                np.array((mesh_x[coord_x],
26                        mesh_y[coord_y])))
27
28
29    def init():
30        ax.contour(
31            X, Y, np.log(Z),
32            np.log([0.5, 10, 30, 80, 130, 200, 300, 500, 900]),
33            cmap='winter'
34        )
35        ax.set_title(graph_title)
36        return ln,
37
38    def update(frame):
39        xdata.append(trajectory[frame][0])
40        ydata.append(trajectory[frame][1])
41        ln.set_data(xdata, ydata)
42        return ln,
43
44    ani = FuncAnimation(
45        fig, update, frames=range(len(trajectory)),
46        init_func=init, repeat=True
47    )
48    writer = ImageMagickFileWriter(fps=10)
49    ani.save(f'saved_gifs/{graph_title}.gif',
50            writer=writer)
51    plt.show()
```

Простая функция $f(x, y) = x^2 + y^2$

In []:

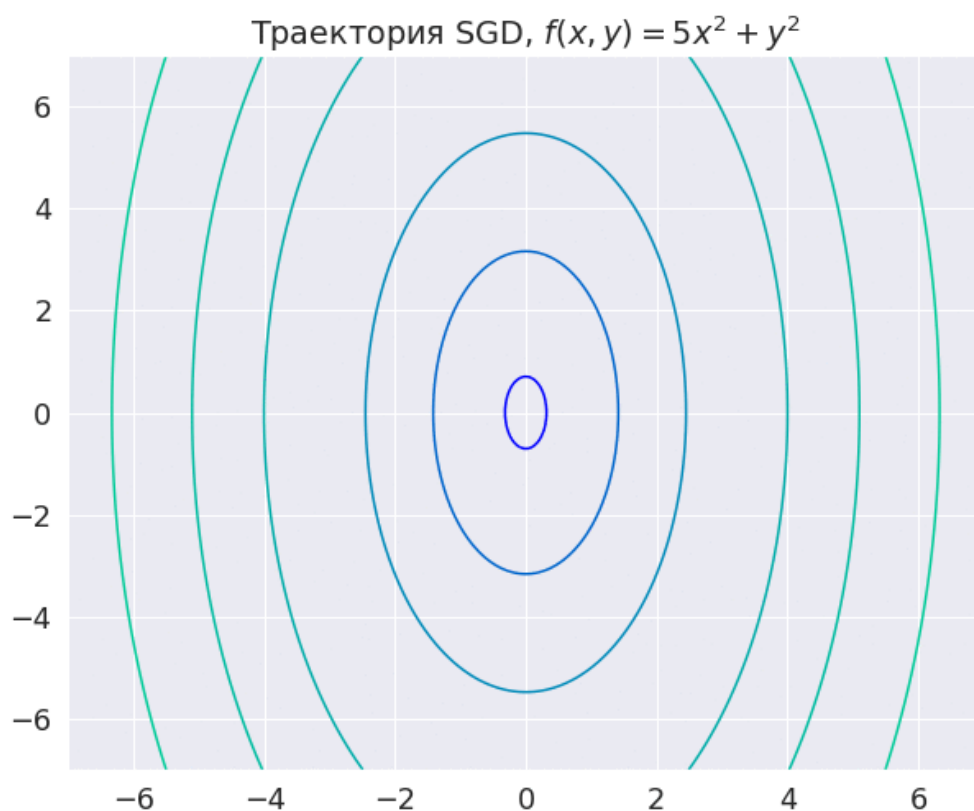
```
1 parameters = np.array((5, 5), dtype=float)
2 func_name = '$f(x, y) = 5x^2 + y^2$'
3 func_grad = square_sum_grad
4 func = square_sum
5 n_iter = 100
```

SGD

In []:

```
1 sgd_trajectory = sgd(
2     init_parameters=parameters,
3     func_grad=func_grad,
4     lr=0.01,
5     n_iter=n_iter
6 )
7 graph_title = 'Траектория SGD, ' + func_name
8 make_experiment(
9     func,
10    sgd_trajectory,
11    graph_title,
12 )
13 clear_output()
14 Image(open(f'saved_gifs/{graph_title}.gif', 'rb').read())
```

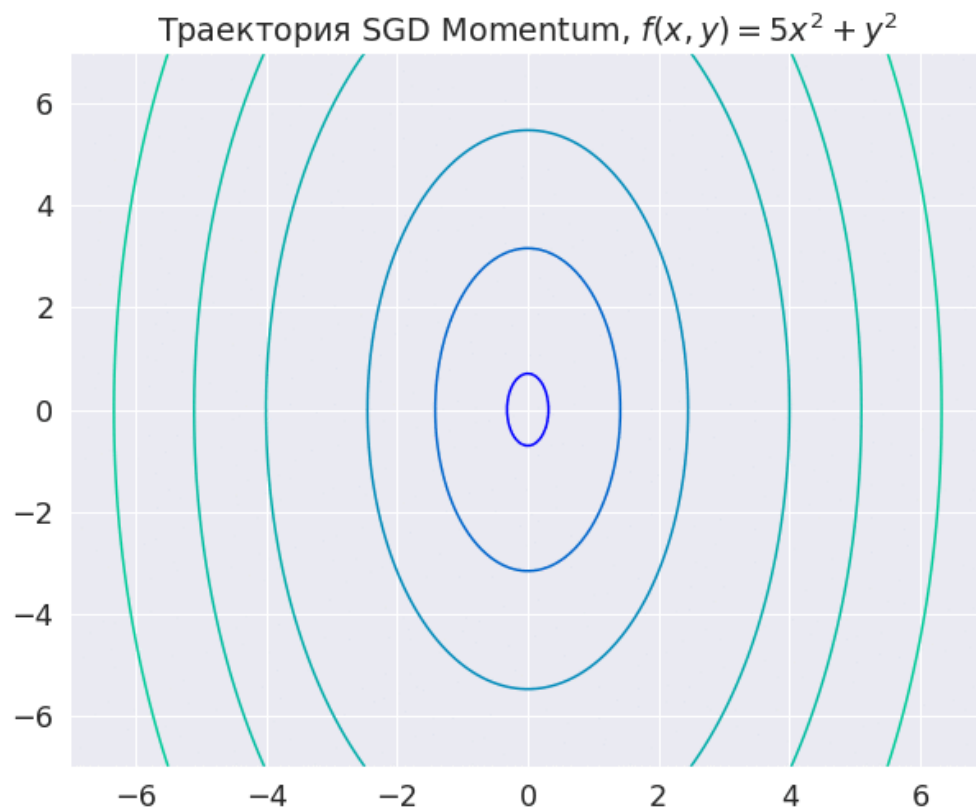
Out[26]:



In []:

```
1 sgd_momentum_trajectory = sgd_momentum(  
2     init_parameters=parameters.copy(),  
3     func_grad=func_grad,  
4     lr=0.01,  
5     n_iter=n_iter,  
6     mu=0.9  
7 )  
8 graph_title = 'Траектория SGD Momentum, ' + func_name  
9 make_experiment(  
10     func,  
11     sgd_momentum_trajectory,  
12     graph_title,  
13 )  
14 clear_output()  
15 Image(open(f'saved_gifs/{graph_title}.gif', 'rb').read())
```

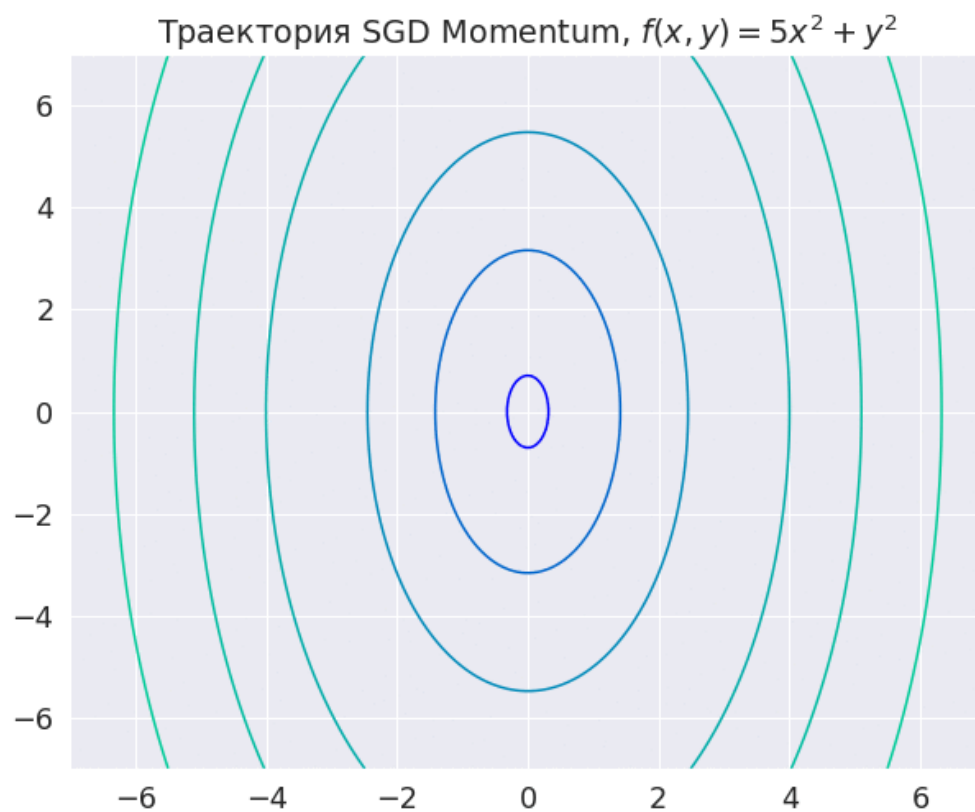
Out[27]:



In []:

```
1 sgd_momentum_trajectory = sgd_momentum(  
2     init_parameters=parameters.copy(),  
3     func_grad=func_grad,  
4     lr=0.01,  
5     n_iter=n_iter,  
6     mu=0.7  
7 )  
8 graph_title = 'Траектория SGD Momentum, ' + func_name  
9 make_experiment(  
10     func,  
11     sgd_momentum_trajectory,  
12     graph_title,  
13 )  
14 clear_output()  
15 Image(open(f'saved_gifs/{graph_title}.gif', 'rb').read())
```

Out[28]:

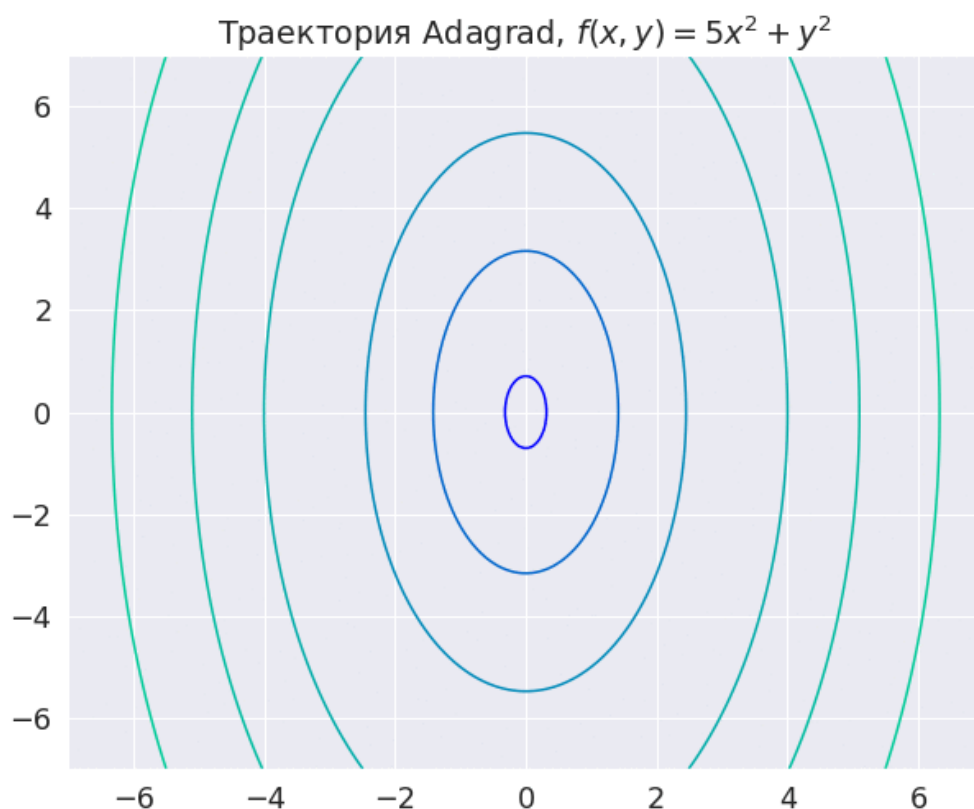


Adagrad

In []:

```
1 adagrad_trajectory = adagrad(  
2     init_parameters=parameters.copy(),  
3     func_grad=func_grad,  
4     lr=0.1,  
5     n_iter=n_iter,  
6     eps=1e-6,  
7 )  
8 graph_title = 'Траектория Adagrad, ' + func_name  
9 make_experiment(  
10     func,  
11     adagrad_trajectory,  
12     graph_title,  
13 )  
14 clear_output()  
15 Image(open(f'saved_gifs/{graph_title}.gif', 'rb').read())
```

Out[29]:



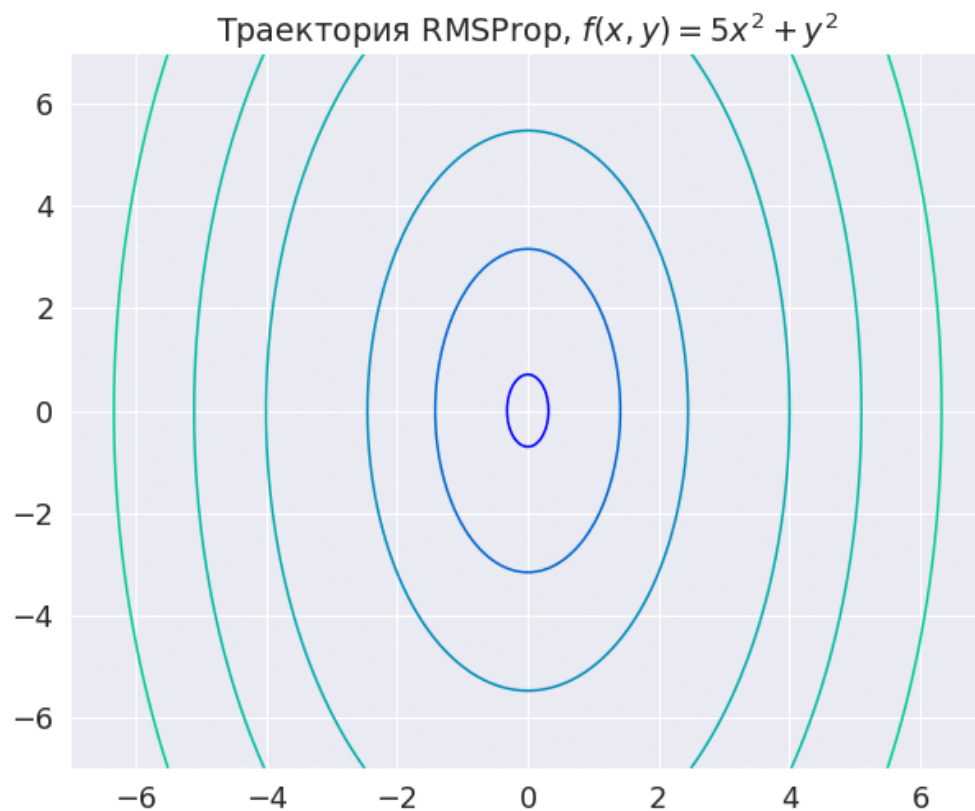
По графику траектории можно заметить, что метод успел сделать очень небольшой шаг. Это связано с очень быстрой скоростью убывания адаптивного шага (learning rate). Поэтому для получения адекватных результатов с Adagrad стоит брать значение η значительно больше чем при использовании SGD.

RMSPProp

In []:

```
1 rmsprop_trajectory = rmsprop(  
2     init_parameters=parameters.copy(),  
3     func_grad=func_grad,  
4     lr=0.1,  
5     n_iter=n_iter,  
6     eps=1e-6,  
7     mu=0.9  
8 )  
9 graph_title = 'Траектория RMSProp, ' + func_name  
10 make_experiment(  
11     func,  
12     rmsprop_trajectory,  
13     graph_title,  
14 )  
15 clear_output()  
16 Image(open(f'saved_gifs/{graph_title}.gif', 'rb').read())
```

Out[30]:

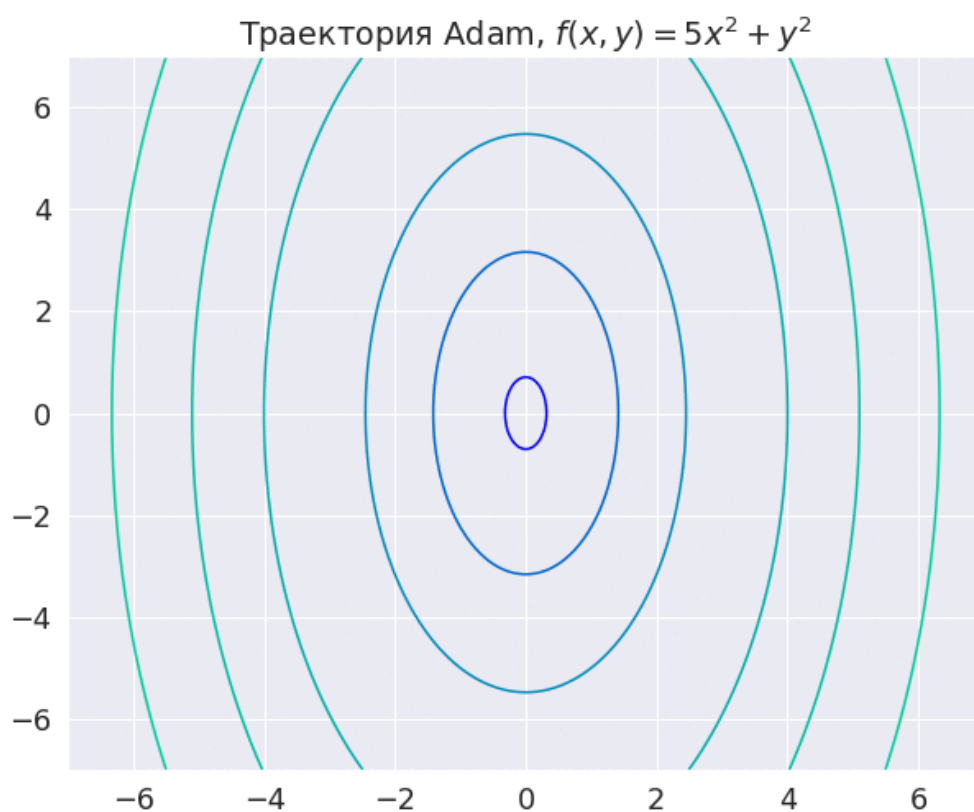


Adam

In []:

```
1 adam_trajectory = adam(  
2     init_parameters=parameters.copy(),  
3     func_grad=func_grad,  
4     lr=0.1,  
5     n_iter=n_iter,  
6     eps=1e-6,  
7     mu=0.9,  
8     beta=0.9  
9 )  
10 graph_title = 'Траектория Adam, ' + func_name  
11 make_experiment(  
12     func,  
13     adam_trajectory,  
14     graph_title,  
15 )  
16 clear_output()  
17 Image(open(f'saved_gifs/{graph_title}.gif', 'rb').read())
```

Out[31]:



Сложная функция

Typesetting math: 100%

$$f(x, y) = (x - 3)^2 + 8(y - 5)^4 + \sqrt{x} + \sin(xy)$$

In []:

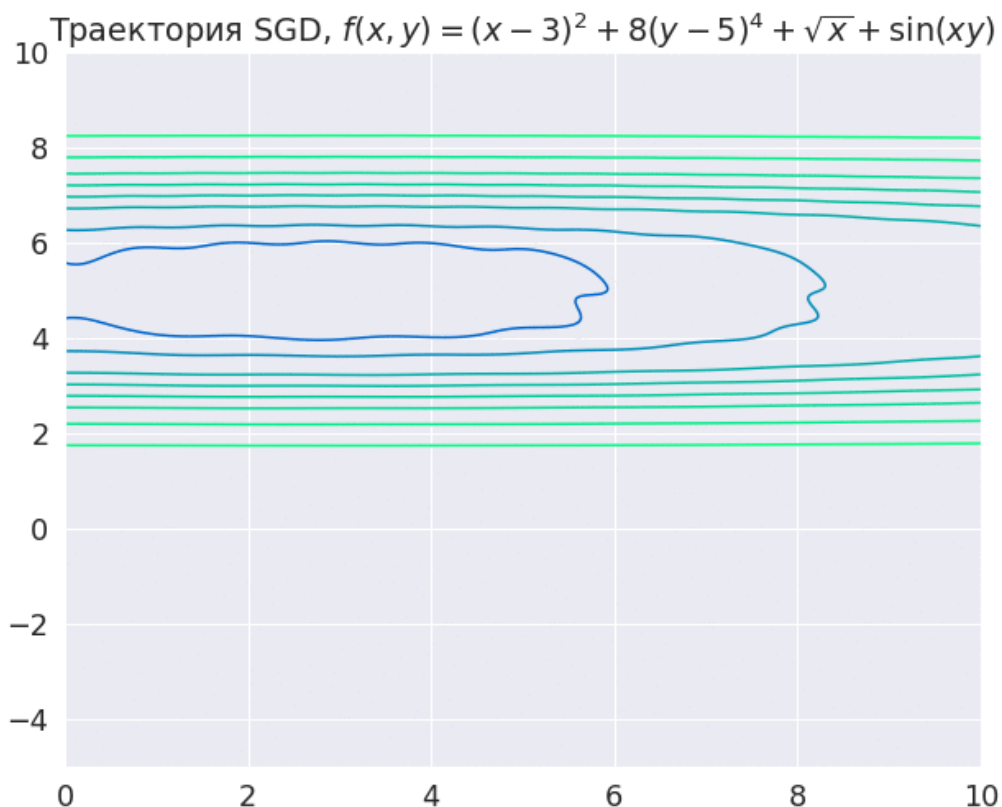
```
1 parameters = np.array((5, -2), dtype=float)
2 func_name = '$f(x, y) = (x-3)^2 + 8(y-5)^4 + \sqrt{x} + \sin(xy)$'
3 func_grad = complex_sum_grad
4 func = complex_sum
5 n_iter = 100
```

SGD

In []:

```
1 sgd_trajectory = sgd(
2     init_parameters=parameters,
3     func_grad=func_grad,
4     lr=0.0002,
5     n_iter=n_iter
6 )
7 graph_title = 'Траектория SGD, ' + func_name
8 make_experiment(
9     func,
10    sgd_trajectory,
11    graph_title,
12    -5, 10, 0, 10
13 )
14 clear_output()
15 Image(open(f'saved_gifs/{graph_title}.gif', 'rb').read())
```

Out[17]:



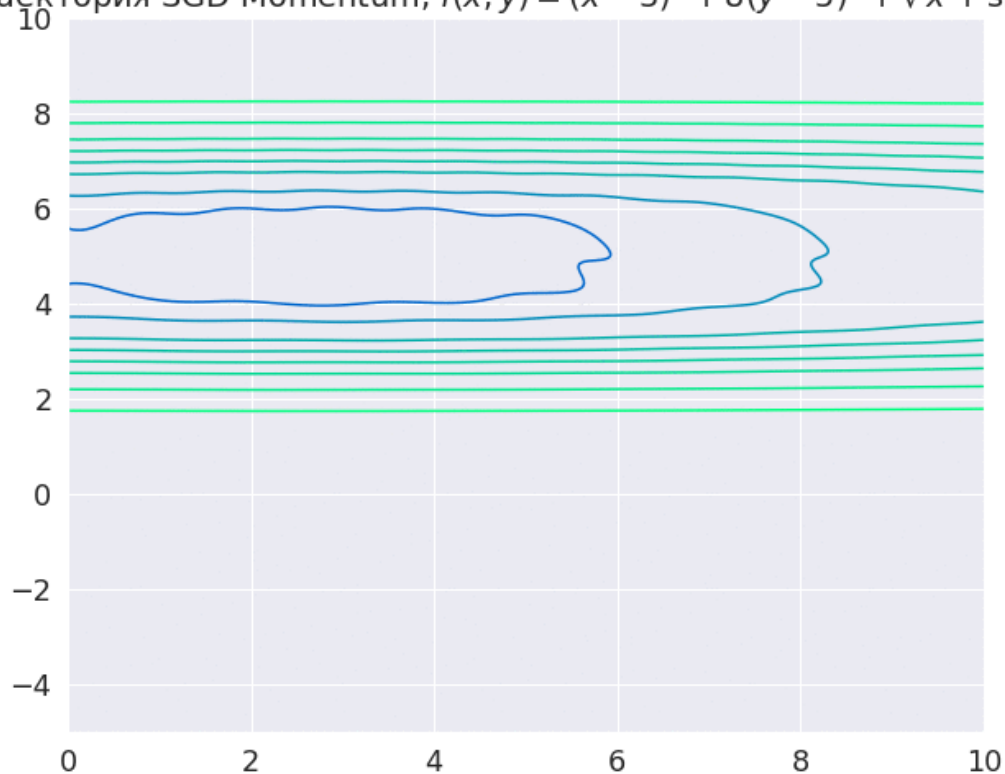
SGD Momentum

In []:

```
1 sgd_momentum_trajectory = sgd_momentum(  
2     init_parameters=parameters,  
3     func_grad=func_grad,  
4     lr=0.0002,  
5     n_iter=n_iter,  
6     mu=0.9  
7 )  
8 graph_title = 'Траектория SGD Momentum, ' + func_name  
9 make_experiment(  
10     func,  
11     sgd_momentum_trajectory,  
12     graph_title,  
13     -5, 10, 0, 10  
14 )  
15 clear_output()  
16 Image(open(f'saved_gifs/{graph_title}.gif', 'rb').read())
```

Out[18]:

Траектория SGD Momentum, $f(x, y) = (x - 3)^2 + 8(y - 5)^4 + \sqrt{x} + \sin(xy)$

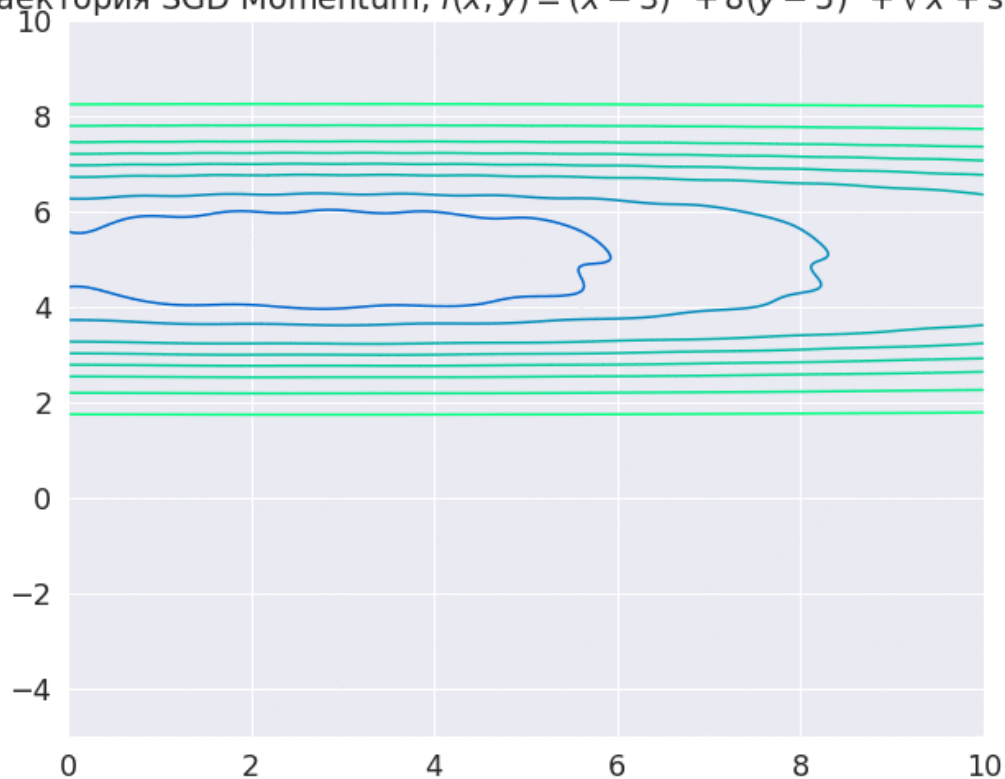


In []:

```
1 sgd_momentum_trajectory = sgd_momentum(  
2     init_parameters=parameters.copy(),  
3     func_grad=func_grad,  
4     lr=0.0002,  
5     n_iter=n_iter,  
6     mu=0.7  
7 )  
8 graph_title = 'Траектория SGD Momentum, ' + func_name  
9 make_experiment(  
10     func,  
11     sgd_momentum_trajectory,  
12     graph_title,  
13     -5, 10, 0, 10  
14 )  
15 clear_output()  
16 Image(open(f'saved_gifs/{graph_title}.gif', 'rb').read())
```

Out[19]:

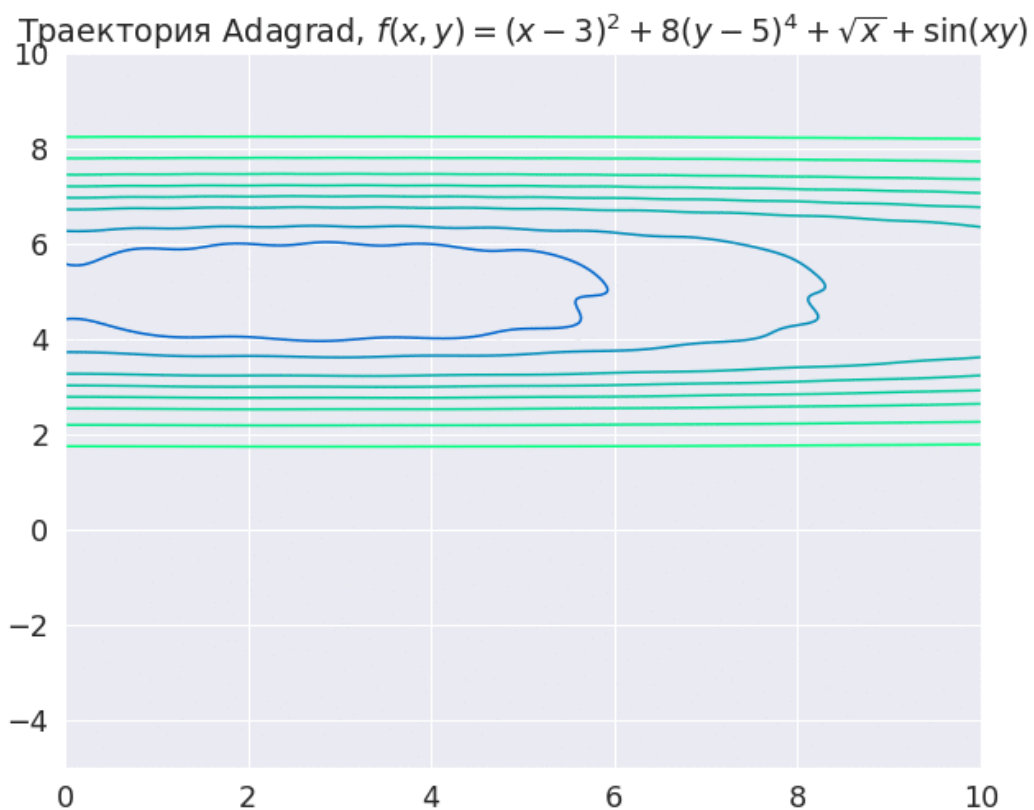
Траектория SGD Momentum, $f(x, y) = (x - 3)^2 + 8(y - 5)^4 + \sqrt{x} + \sin(xy)$



In []:

```
1 adagrad_trajectory = adagrad(  
2     init_parameters=parameters.copy(),  
3     func_grad=func_grad,  
4     lr=0.1,  
5     n_iter=n_iter,  
6     eps=1e-6,  
7 )  
8 graph_title = 'Траектория Adagrad, ' + func_name  
9 make_experiment(  
10     func,  
11     adagrad_trajectory,  
12     graph_title,  
13     -5, 10, 0, 10  
14 )  
15 clear_output()  
16 Image(open(f'saved_gifs/{graph_title}.gif', 'rb').read())
```

Out[20]:

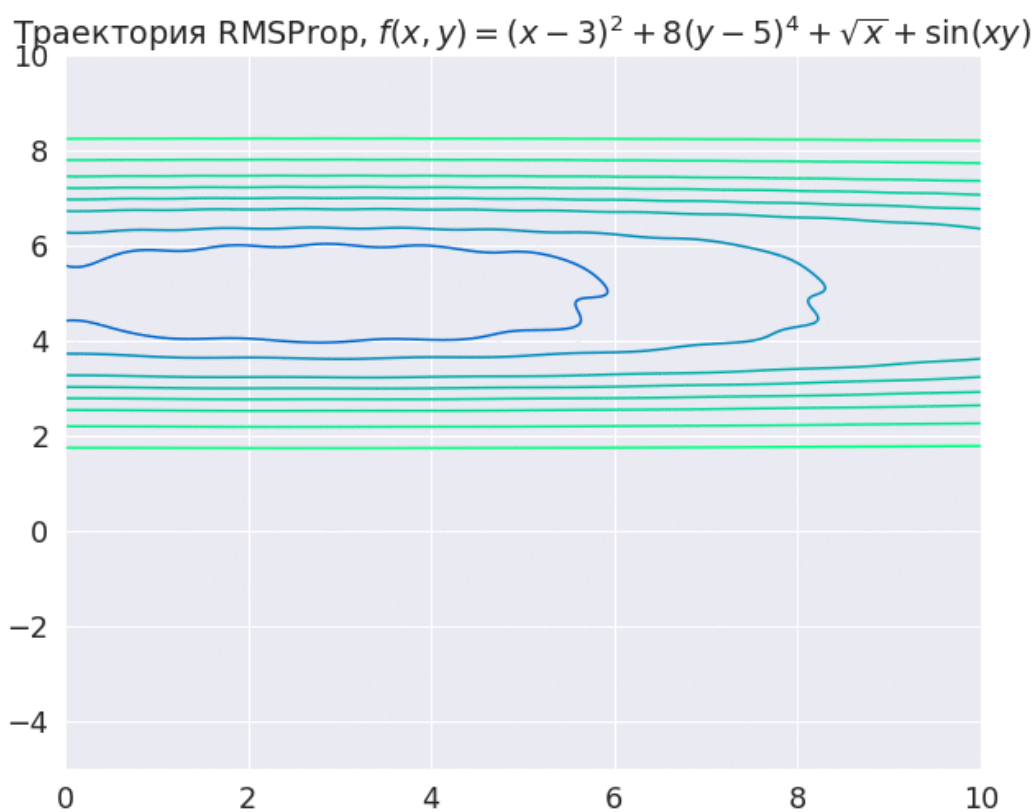


RMSPProp

In []:

```
1 rmsprop_trajectory = rmsprop(  
2     init_parameters=parameters.copy(),  
3     func_grad=func_grad,  
4     lr=0.1,  
5     n_iter=n_iter,  
6     eps=1e-6,  
7     mu=0.9  
8 )  
9 graph_title = 'Траектория RMSProp, ' + func_name  
10 make_experiment(  
11     func,  
12     rmsprop_trajectory,  
13     graph_title,  
14     -5, 10, 0, 10  
15 )  
16 clear_output()  
17 Image(open(f'saved_gifs/{graph_title}.gif', 'rb').read())
```

Out[21]:



In []:

```
1 adam_trajectory = adam(  
2     init_parameters=parameters.copy(),  
3     func_grad=func_grad,  
4     lr=0.1,  
5     n_iter=n_iter,  
6     eps=1e-6,  
7     mu=0.9,  
8     beta=0.9  
9 )  
10 graph_title = 'Траектория Adam, ' + func_name  
11 make_experiment(  
12     func,  
13     adam_trajectory,  
14     graph_title,  
15     -5, 10, 0, 10  
16 )  
17 clear_output()  
18 Image(open(f'saved_gifs/{graph_title}.gif', 'rb').read())
```

Out[22]:

