

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')
```

started 00:26:20 2020-05-03, finished in 858ms

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 $Q(x) \rightarrow \min_x$, и $\nabla_x Q(x)$ - градиент функции $Q(x)$.

1. 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 в большей степени от изменения параметра меняется значение оптимизируемой функции, тем с меньшей скоростью стоит обновлять этот параметр. Иначе высок шанс расходимости метода.

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

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

В матрично-векторном виде шаг алгоритма можно переписать так:

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

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

2. 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 значение квадраты градиентов как с самых первых итераций, так и с самых последних. Хотя, на самом деле, наибольшую значимость имеет модули градиентов на последних нескольких итерациях. Для этого предлагается использовать экспоненциальное сглаживание.

$$g_t = \mu g_{t-1} + (1 - \mu) \nabla Q(x_t) \odot \nabla Q(x_t)$$

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

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

3. Adadelta

Этот метод по формуле шага и по смыслу очень похож на RMSProp. Авторы метода заметили, что в различных методах 1 порядка для оптимальной сходимости нужно брать совершенно разные значения learning rate (η), а иногда - подбирать значение η в зависимости от решаемой задачи. Чтобы избавиться от необходимости находить наилучшее значение η . Для этого корень среднеквадратичной ошибки обновления параметра (RMS) считается теперь и для $\Delta\theta$.

$$\begin{aligned} d_t &= \mu d_{t-1} + (1 - \mu) \Delta x_t \odot \Delta x_t \\ \Delta x_t &= - \frac{\sqrt{d_{t-1} + \epsilon}}{\sqrt{g_t + \epsilon}} \odot \nabla Q(x_t) \\ x_{t+1} &= x_t + \Delta x_t \end{aligned}$$

Преимущество данного метода по сравнению с RMSProp - отсутствие необходимости подбирать значения параметра η . Экспериментальным путём выяснено, что для Adadelta наилучшее значение $\gamma \sim 0.9$.

4. Adam

Этот алгоритм совмещает в себе 2 идеи: идею алгоритма Momentum о накоплении 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:

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

Чтобы эти оценки не были смещёнными, нужно их отнормировать:

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

Тогда получим итоговую формулу шага:

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

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>),

Adadelta -- <https://arxiv.org/abs/1212.5701> (<https://arxiv.org/abs/1212.5701>),

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.

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 adagrad(theta0, func_grad, eps=1e-6, eta=0.01, iter_count=150):
2      '''
3      Adagrad.
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) eta - скорость обучения,
10     5) iter_count - количество итераций метода.
11     '''
12
13     theta = theta0
14     history = [theta0]
15     cumulative_grad = np.zeros(theta0.shape)
16
17  ▾     for iter_id in range(iter_count):
18         current_grad = func_grad(theta)
19         cumulative_grad += current_grad**2
20         theta = theta - eta * current_grad / np.sqrt(cumulative_grad + eps)
21         history.append(theta)
22
23     return history
24
25
26  ▾ def rmsprop(theta0, func_grad, eps=1e-6, eta=0.01, mu=0.9, iter_count=150):
27     '''
28     RMSProp.
29
30     Параметры.
31     1) theta0 - начальное приближение theta,
32     2) func_grad - функция, задающая градиент оптимизируемой функции,
33     3) eps - мин. возможное значение знаменателя,
34     4) eta - скорость обучения,
35     5) mu - параметр экспоненциального сглаживания,
36     6) iter_count - количество итераций метода.
37     '''
38
39     theta = theta0
40     history = [theta0]
41     cumulative_grad = np.zeros(theta0.shape)
42
43  ▾     for iter_id in range(iter_count):
44         current_grad = func_grad(theta)
45         cumulative_grad = mu * cumulative_grad + (1-mu) * current_grad**2
46         theta = theta - eta * current_grad / np.sqrt(cumulative_grad + eps)
47         history.append(theta)
48
49     return history
50
51  ▾ def adadelta(theta0, func_grad, eps=1e-6, mu=0.9, iter_count=150):
52     '''
53     Adadelta.
54
55     Параметры.
56     1) theta0 - начальное приближение theta,
57     2) func_grad - функция, задающая градиент оптимизируемой функции,
58     3) eps - мин. возможное значение знаменателя,
59     4) mu - параметр экспоненциального сглаживания,
```

```

60     5) iter_count - количество итераций метода.
61     ...
62
63     theta = theta0.astype(float)
64     history = [theta0]
65     cumulative_grad = np.zeros(theta0.shape)
66     cumulative_delta = np.zeros(theta0.shape)
67
68     for iter_id in range(iter_count):
69         current_grad = func_grad(theta)
70         cumulative_grad = mu * cumulative_grad + (1-mu) * current_grad**2
71         delta_theta = current_grad / np.sqrt(cumulative_grad + eps)
72         delta_theta *= -np.sqrt(cumulative_delta + eps)
73         theta += delta_theta
74         cumulative_delta = mu * cumulative_delta + (1-mu) * delta_theta**2
75         history.append(theta)
76
77     return history
78
79     def adam(theta0, func_grad, eps=1e-6, eta=0.01, beta=0.9, mu=0.95, iter_count
80         ...
81         Adam.
82
83         Параметры.
84         1) theta0 - начальное приближение theta,
85         2) func_grad - функция, задающая градиент оптимизируемой функции,
86         3) eps - мин. возможное значение знаменателя,
87         4) eta - скорость обучения,
88         5) beta - параметр экспоненциального сглаживания,
89         6) mu - параметр экспоненциального сглаживания,
90         7) iter_count - количество итераций метода.
91         ...
92
93         theta = theta0
94         history = [theta0]
95         cumulative_m = np.zeros(theta0.shape)
96         cumulative_v = np.zeros(theta0.shape)
97         pow_beta, pow_mu = beta, mu
98
99         for iter_id in range(iter_count):
100             current_grad = func_grad(theta)
101             cumulative_m = beta * cumulative_m + (1 - beta) * current_grad
102             cumulative_v = mu * cumulative_v + (1 - mu) * current_grad**2
103
104             scaled_m = cumulative_m / (1 - pow_beta)
105             scaled_v = cumulative_v / (1 - pow_mu)
106             theta = theta - eta * scaled_m / (np.sqrt(cumulative_v) + eps)
107             history.append(theta)
108
109             pow_beta *= beta
110             pow_mu *= mu
111
112         return history

```

started 00:26:23 2020-05-03, finished in 18ms

Реализуем функции, которые будем оптимизировать.

In [3]:

```
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(x
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])
```

started 00:26:24 2020-05-03, finished in 13ms

Создадим директорию, в которой будем хранить визуализацию экспериментов.

In [4]:

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

started 00:26:25 2020-05-03, finished in 245ms

In [5]:

```
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 = animation.FuncAnimation(
45        fig, update, frames=range(len(trajectory)),
46        init_func=init, repeat=True
47    )
48 ▾    ani.save(f'saved_gifs/{graph_title}.gif',
49            writer='imagemagick', fps=5)
```

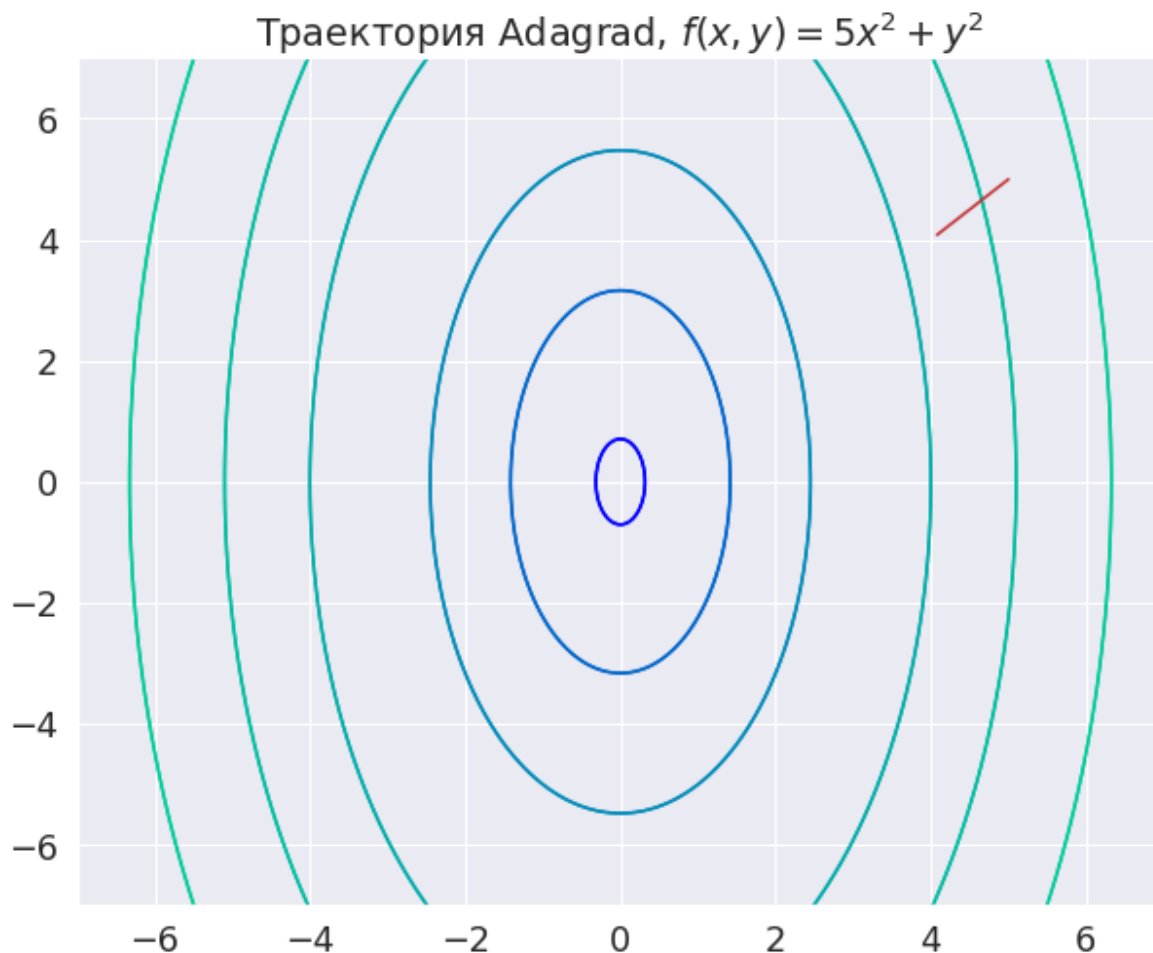
started 00:26:28 2020-05-03, finished in 12ms

Adagrad

In [6]:

```
1 ▾ adagrad_trajectory = adagrad(  
2     np.array((5, 5)), square_sum_grad, eta=0.1, iter_count=30  
3 )  
4 ▾ make_experiment(  
5     square_sum, adagrad_trajectory,  
6     'Траектория Adagrad,  $f(x, y) = 5x^2 + y^2$ '  
7 )
```

started 00:26:28 2020-05-03, finished in 6.31s

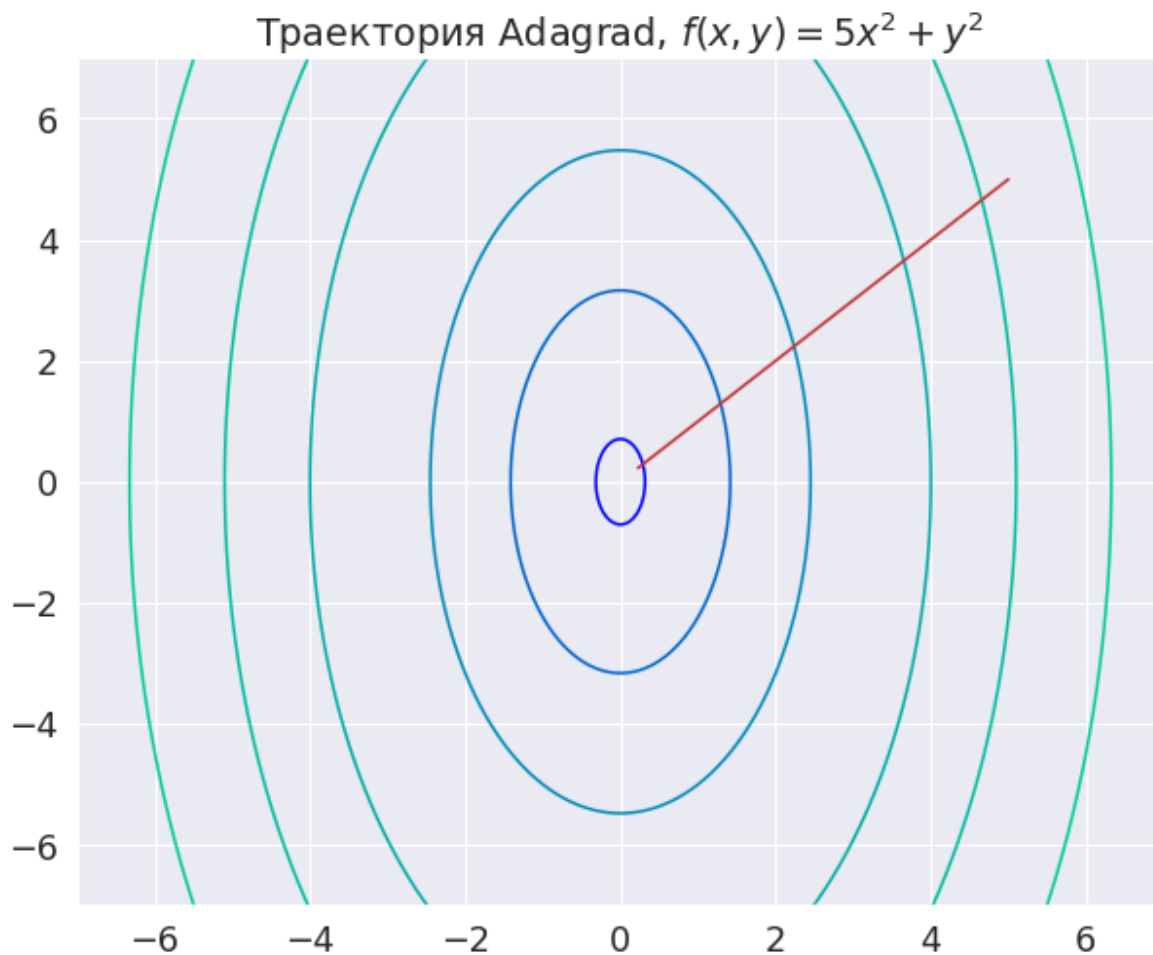


По графику траектории можно заметить, что метод успел сделать очень небольшой шаг. Это связано с очень быстрой скоростью убывания адаптивного шага (learning rate). Поэтому для получения адекватных результатов с Adagrad стоит брать значение η значительно больше чем при использовании SGD.

In [7]:

```
1 ▾ adagrad_trajectory = adagrad(  
2     np.array((5, 5)), square_sum_grad, eta=0.9, iter_count=30  
3 )  
4 ▾ make_experiment(  
5     square_sum, adagrad_trajectory,  
6     'Траектория Adagrad,  $f(x, y) = 5x^2 + y^2$ '  
7 )
```

started 00:26:35 2020-05-03, finished in 6.58s



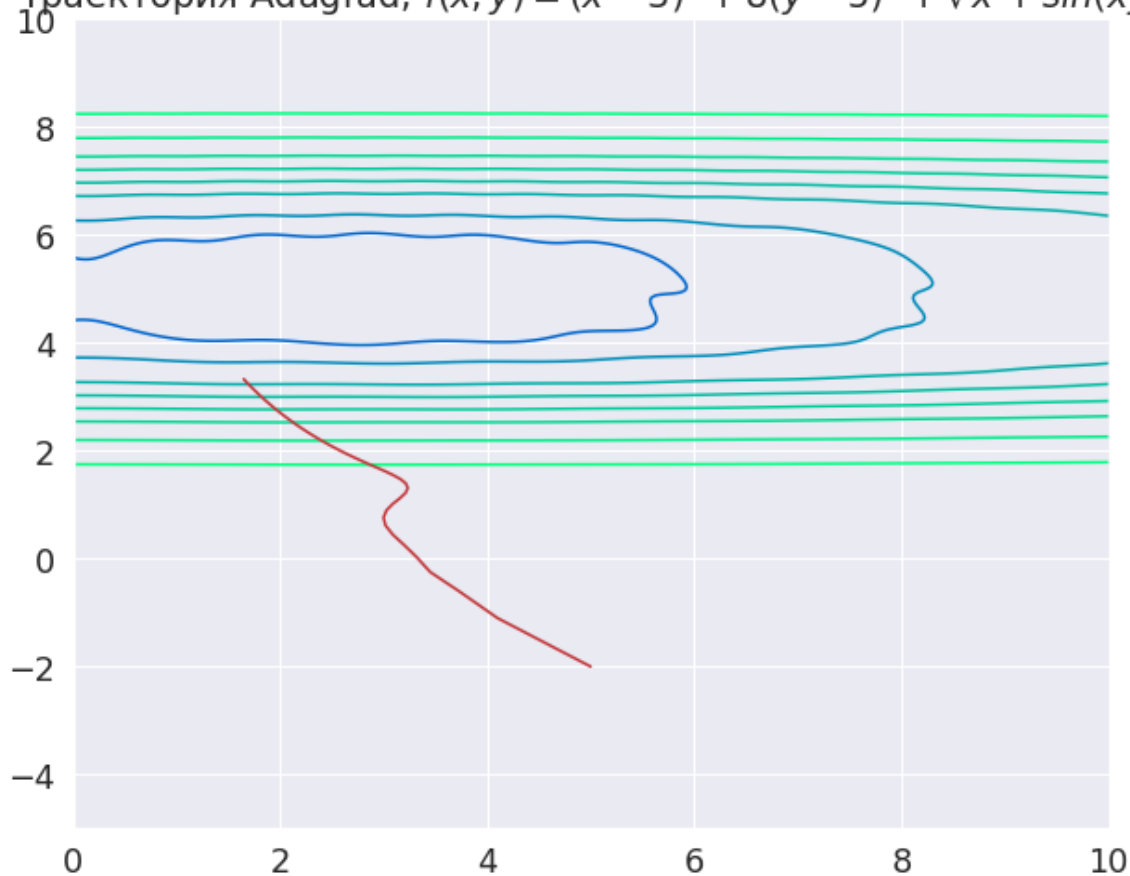
Опробуем Adagrad на другой функции.

In [8]:

```
1 ▾ adagrad_trajectory = adagrad(  
2     np.array((5, -2)), complex_sum_grad, eta=0.9, iter_count=100  
3 )  
4 ▾ make_experiment(  
5     complex_sum, adagrad_trajectory,  
6     'Траектория Adagrad, $f(x, y) = (x-3)^2 + 8(y-5)^4 + \sqrt{x} + \sin(xy)$  
7     -5, 10, 0, 10  
8 )
```

started 00:26:41 2020-05-03, finished in 55.8s

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

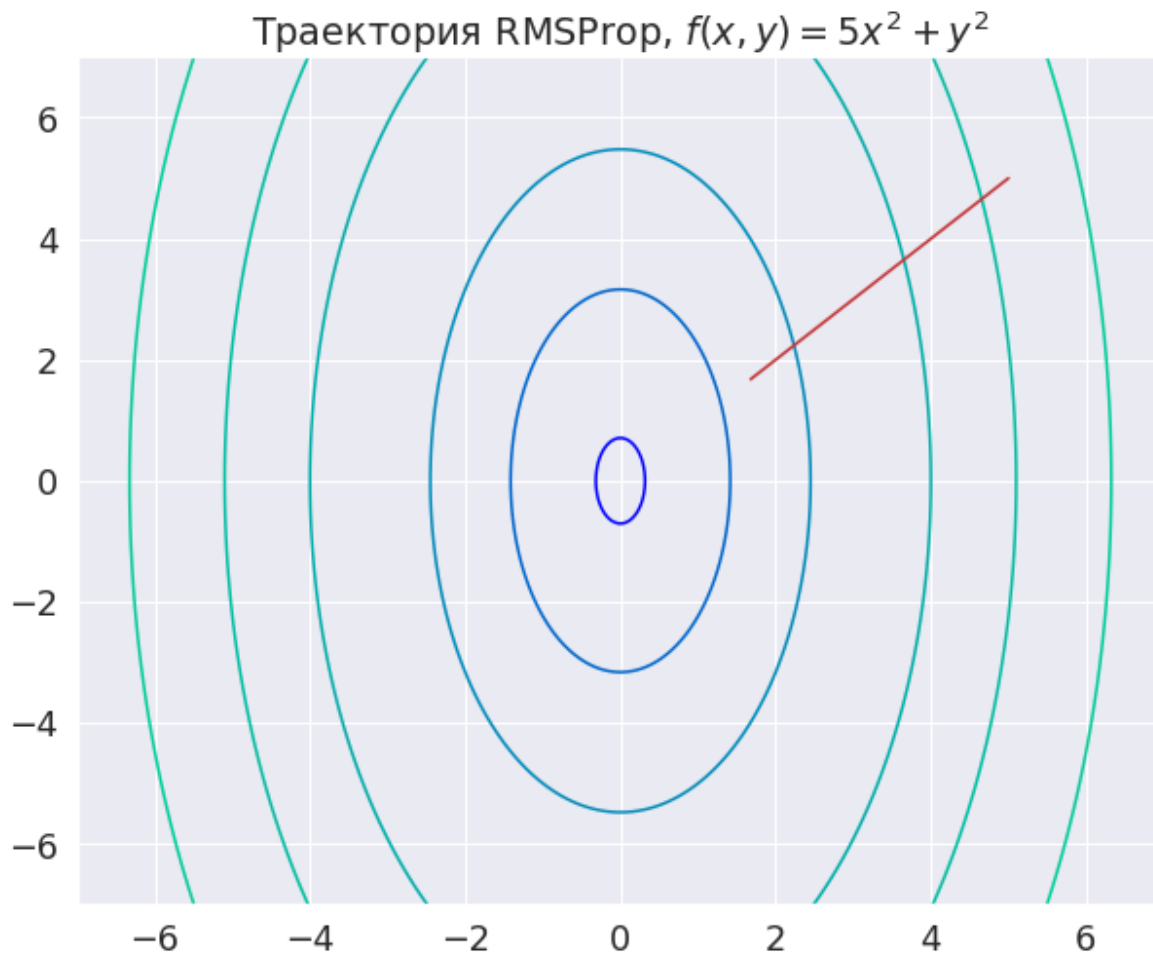


RMSProp

In [9]:

```
1 ▾ rmsprop_trajectory = rmsprop(  
2     np.array((5, 5)), square_sum_grad, eta=0.1, iter_count=30  
3 )  
4 ▾ make_experiment(  
5     square_sum, rmsprop_trajectory,  
6     'Траектория RMSProp,  $f(x, y) = 5x^2 + y^2$ '  
7 )
```

started 00:27:37 2020-05-03, finished in 6.93s

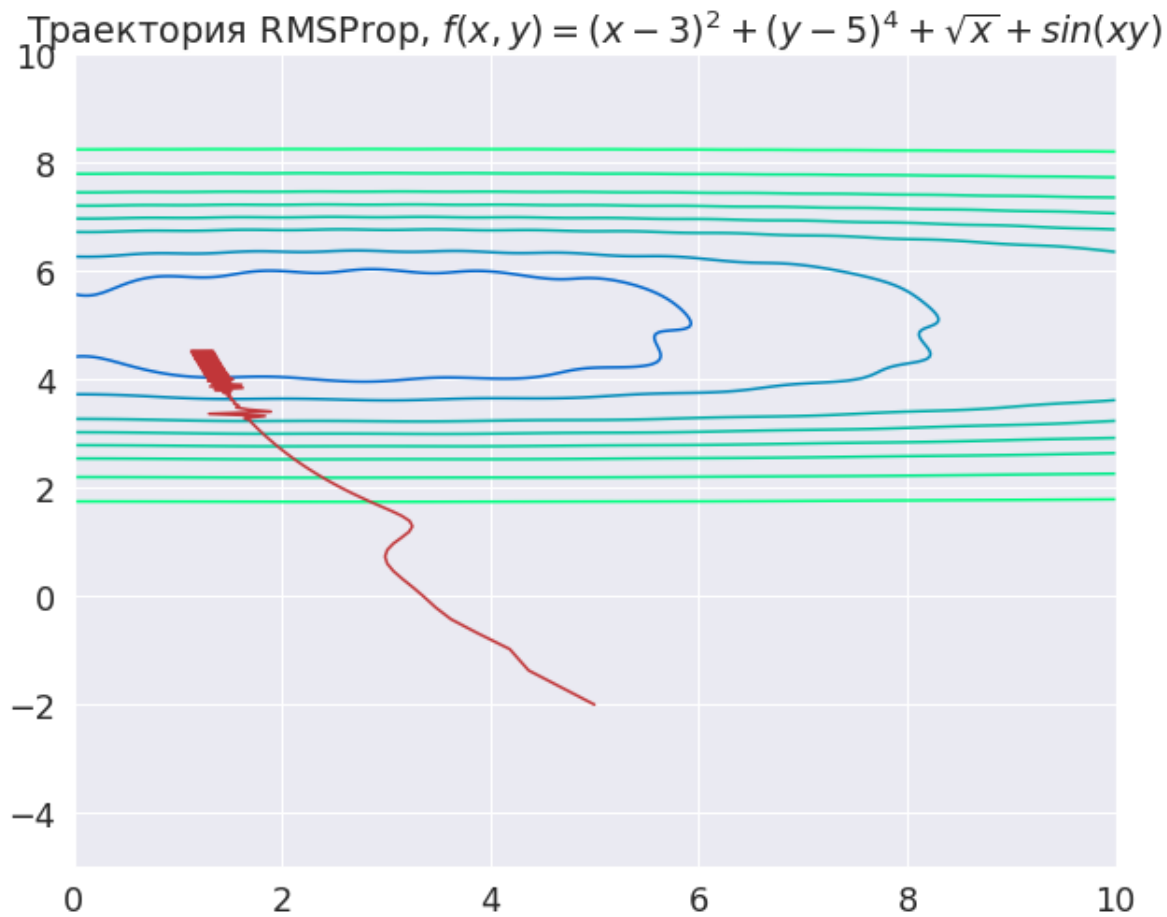


В RMSProp адаптивный шаг убывает медленнее, что делает этот метод более устойчивым.

In [10]:

```
1 rmsprop_trajectory = rmsprop(np.array((5, -2)), complex_sum_grad, eta=0.2, ito
2 ▼ make_experiment(
3     complex_sum, rmsprop_trajectory,
4     'Траектория RMSProp, $f(x, y) = (x-3)^2 + (y-5)^4 + \sqrt{x} + \sin(xy)$',
5     -5, 10, 0, 10
6 )
```

started 00:27:44 2020-05-03, finished in 54.5s

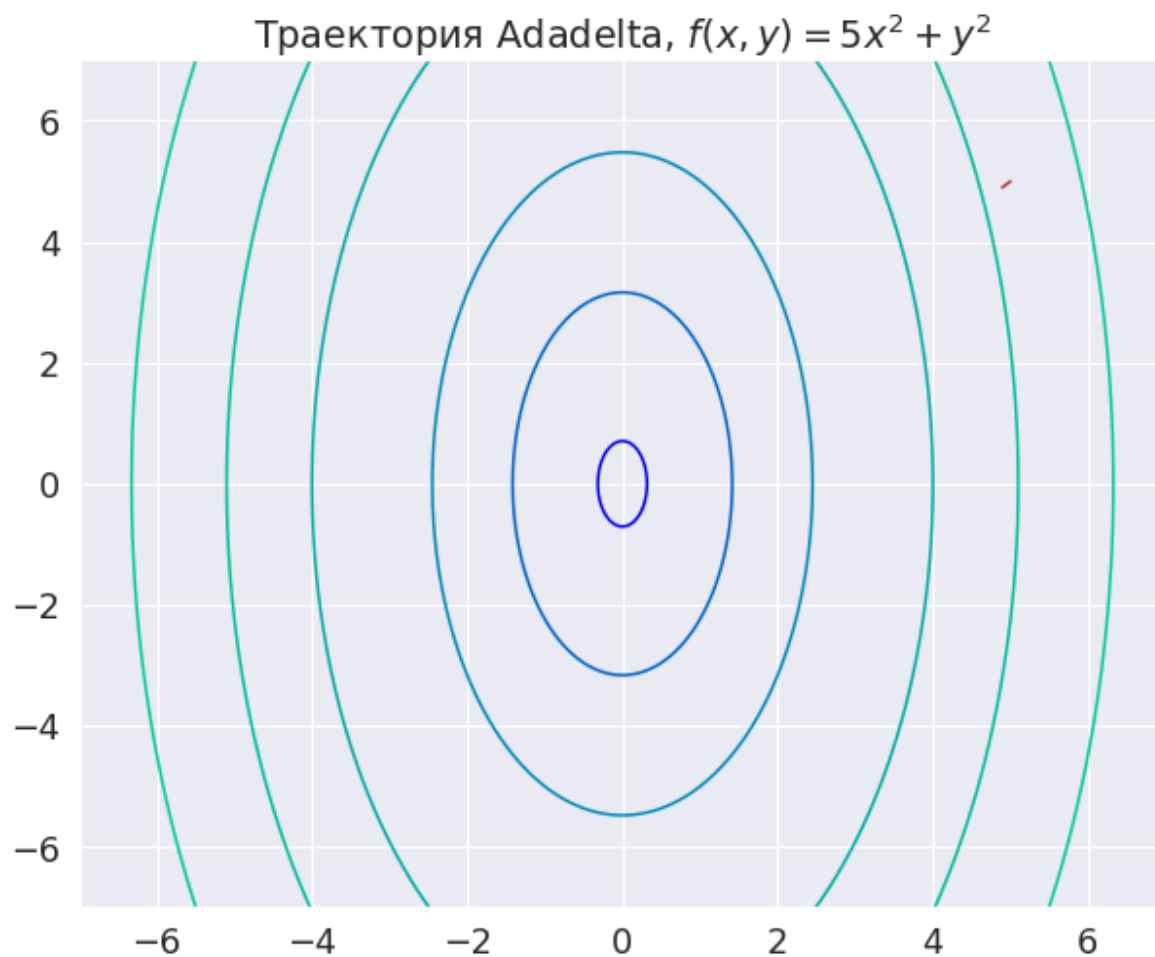


Adadelta

In [11]:

```
1 ▾ adadelta_trajectory = adadelta(  
2     np.array((5, 5)), square_sum_grad, iter_count=30  
3 )  
4 ▾ make_experiment(  
5     square_sum, adadelta_trajectory,  
6     'Траектория Adadelta,  $f(x, y) = 5x^2 + y^2$ '  
7 )
```

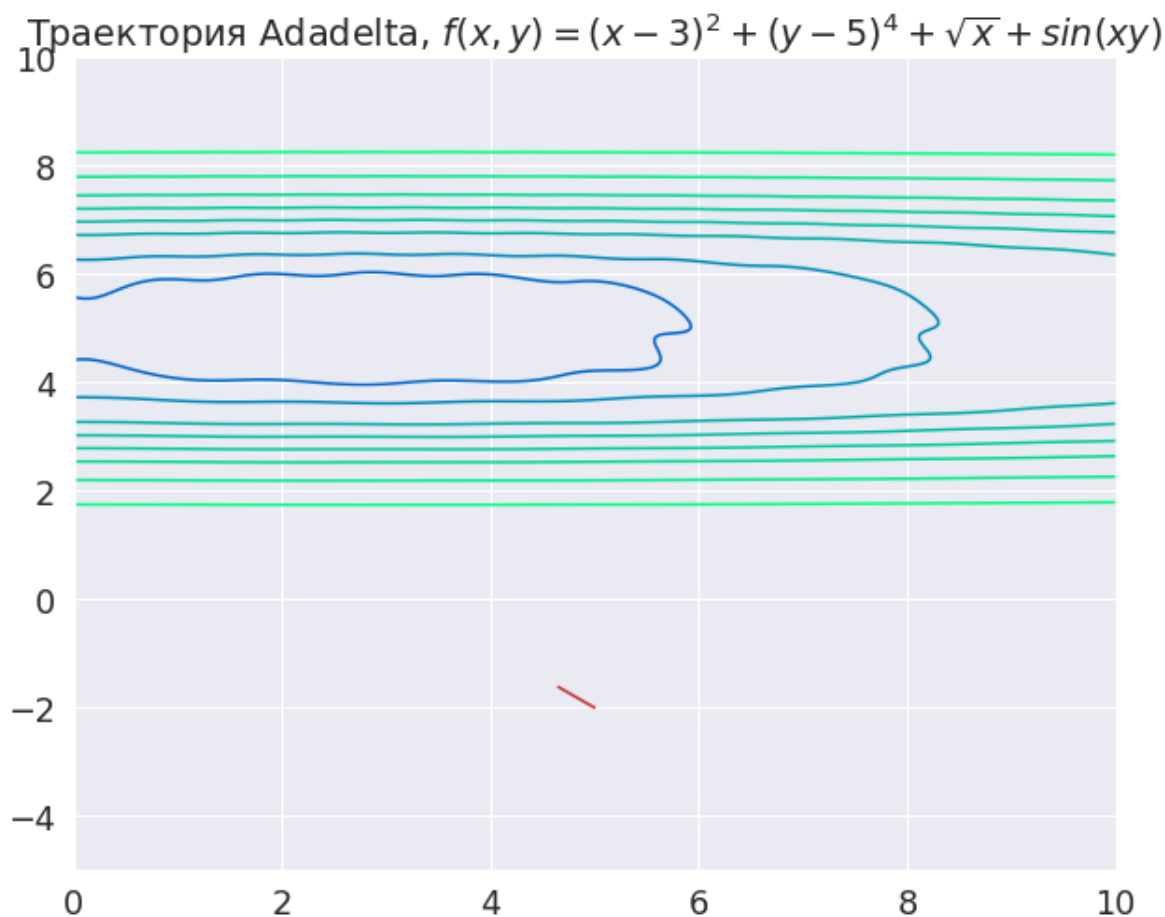
started 00:28:38 2020-05-03, finished in 6.20s



In [12]:

```
1 ▾ adadelata_trajectory = adadelata(  
2     np.array((5, -2)), complex_sum_grad, iter_count=100  
3 )  
4 ▾ make_experiment(  
5     complex_sum, adadelata_trajectory,  
6     'Траектория Adadelata, $f(x, y) = (x-3)^2 + (y-5)^4 + \sqrt{x} + \sin(xy)$',  
7     -5, 10, 0, 10)
```

started 00:28:45 2020-05-03, finished in 53.8s



На наших оптимизируемых функциях метод Adadelata показал очень низкую скорость сходимости. Тем не менее, на многих архитектурах нейросетей он работает гораздо лучше.

Adam

In [13]:

```
1 ▾ adam_trajectory = adam(  
2     np.array((5, -2)), complex_sum_grad, eta=0.2, iter_count=100  
3 )  
4 ▾ make_experiment(  
5     complex_sum, adam_trajectory,  
6     'Траектория Adam, $f(x, y) = (x-3)^2 + (y-5)^4 + \sqrt{x} + \sin(xy)$',  
7     -5, 10, 0, 10  
8 )
```

started 00:29:38 2020-05-03, finished in 57.4s

