

Optimisation Algorithms for Data Analysis - Assignment 2

CS7DS2

Maurice Kiely - 23350601

03/03/2024

Question a)

I began by defining 2 classes to represent my 2 functions.

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 import sympy
```



```
1 class func1:
2     def __init__(self):
3         self.x_sym, self.y_sym = sympy.symbols('x y', real=True)
4     def f(self, x, y):
5         return 5 * (x - 3)**4 + 7 * (y - 9)**2
6     def dx(self, x, y):
7         df_dx = sympy.diff(self.f(self.x_sym, self.y_sym),
8                             self.x_sym)
9         return df_dx.subs({self.x_sym: x, self.y_sym: y})
10    def dy(self, x, y):
11        df_dy = sympy.diff(self.f(self.x_sym, self.y_sym),
12                            self.y_sym)
13        return df_dy.subs({self.x_sym: x, self.y_sym: y})
```



```
12
13 class func2:
14     def __init__(self):
15         self.x_sym, self.y_sym = sympy.symbols('x y', real=True)
16     def f(self, x, y):
17         return sympy.Max(x - 3, 0) + 7 * sympy.Abs(y - 9)
18     def dx(self, x, y):
19         df_dx = sympy.diff(self.f(self.x_sym, self.y_sym),
20                             self.x_sym)
21         return df_dx.subs({self.x_sym: x, self.y_sym: y})
22     def dy(self, x, y):
23         df_dy = sympy.diff(self.f(self.x_sym, self.y_sym),
24                             self.y_sym)
25         return df_dy.subs({self.x_sym: x, self.y_sym: y})
```

These classes are used to conveniently find the gradients at each (x, y) point. This obviously saves a lot of time when implementing the different methods. My functions were:

$$1. \ 5(x - 3)^4 + 7(y - 9)^2$$

$$2. \ (x - 3)^+ + 7|y - 9|$$

i) Polyak

```
1 def polyak_GD(func, x0, y0, f_star=0, epsilon=0.0001,
2     max_iters=1000, tol=100):
3     # Initialize x,y and path and f(x,y) lists
4     x, y = x0, y0
5     path = [(x, y)]
6     f_vals = []
7
8     # Implement formula to find alpha and update at each step.
9     for i in range(max_iters):
10         df_dx = func.dx(x, y)
11         df_dy = func.dy(x, y)
12         alpha = (func.f(x, y) - f_star) / (df_dx**2 + df_dy**2 +
13             epsilon)
14         x -= alpha * df_dx
15         y -= alpha * df_dy
16         # Append both (x, y) point and f value to lists
17         path.append((x, y))
18         f_vals.append(func.f(x, y))
19
20         # Check stopping conditions for alpha too high or too low
21         if alpha > tol:
22             print(f"Stopping after {i} Iterations as alpha too
23                 large")
24             return path, f_vals
25
26         if sympy.sqrt(func.dx(x, y)**2 + func.dy(x, y)**2) <
27             epsilon:
28             return path, f_vals
29
30
31     print(f"No sufficient convergence by iteration
32         {max_iters}\nFinished at (x, y): ({x:0.4f}, {y:0.4f})")
33
34     return path, f_vals
```

The polyak function begins by initializing the variables for the first pass of calculating alpha, as well as those to store the path taken by $f(x,y)$. The algorithm works by using the formula $\alpha = \frac{f(x)-f^*}{\nabla f(x)^T \nabla f(x) + \epsilon}$. This function calculates the step size by transforming $f(x - \alpha \nabla f(x)) \approx f^*$. We update x iteratively as such.

As stopping criterion, we use $\|\nabla f(x)\| < \epsilon$. Here, if the gradient isn't changing, we stop. This is because, this leads to an extremely large alpha otherwise which can introduce 'noise' around the solution.

ii) RMSProp

```
1 def rmsprop(func, x0, y0, alpha0, beta, epsilon=0.0001,
2             max_iters=1000):
3     # Initialize x,y, sums and path,f(x,y) lists
4     x, y = x0, y0
5     sum_dx = 0
6     sum_dy = 0
7     path = [(x, y)]
8     f_vals = []
9
10    for t in range(1, max_iters + 1):
11        # Take gradients for current x, y
12        grad_x = func.dx(x, y)
13        grad_y = func.dy(x, y)
14
15        # Update denominator for x and y using the following
16        sum_dx = beta * sum_dx + (1 - beta) * grad_x**2
17        sum_dy = beta * sum_dy + (1 - beta) * grad_y**2
18
19        # Compute the step sizes for x and y
20        step_size_x = alpha0 / sympy.sqrt(sum_dx) + epsilon
21        step_size_y = alpha0 / sympy.sqrt(sum_dy) + epsilon
22
23        # Update x and y
24        x -= step_size_x * grad_x
25        y -= step_size_y * grad_y
26
27        path.append((x, y))
28        f_vals.append(func.f(x, y))
29
30    return path, f_vals
```

The RMSProp method uses the gradient from each variable to calculate a vector of step sizes. This idea is taken from the Adagard, however, the RMSProp method adds a decaying memory property to put more emphasis on step sizes based on more recent gradients. This algorithm begins by initializing all necessary variables including the sums of gradients for each variable. At each iteration, the gradient is calculated and the sum variable is updated. This is done by taking a β weighted average of current and past gradients. β can be adjusted accordingly to emphasise more recent gradients or not. α is a learning rate that dictates how impactful more recent gradients are. We use this calculated sum to get the step size and update x accordingly.

iii) Heavy Ball

```
1 def heavy_ball(func, x0, y0, alpha, beta, num_iters=1000,
2     tol=1e-5, stopping=True):
3     # Initialize x,y, momentum variables and path,f(x,y) lists
4     x, y = x0, y0
5     ux, uy = 0, 0
6     path = [(x, y)]
7     f_vals = []
8
9     for i in range(num_iters):
10         # Compute the gradients
11         grad_x = func.dx(x, y)
12         grad_y = func.dy(x, y)
13         # Update momentum variable with current and past momentum
14         ux = beta * ux + alpha * grad_x
15         uy = beta * uy + alpha * grad_y
16         # Update x and y
17         x -= ux
18         y -= uy
19         path.append((x, y)) # Store the new values
20         f_vals.append(func.f(x, y))
21
22         if func.f(x, y) > 1e20:
23             break
24         # Check if the norm of the gradient is below the
25         # tolerance
26         if sympy.sqrt(grad_x**2 + grad_y**2) < tol and stopping
27             == True:
28             break
29     return path, f_vals
```

Heavy Ball/Polyak Momentum is an algorithm that uses 1 step size for all variables at each iteration. However, unlike standard polyak, its step size is influenced by values of past gradients. At step t , the step size is a β weighted average of past gradients and learning rate $\alpha \times$ gradient. We implement this algorithm by initializing the variables including the momentum variables. We then begin iteration. Each iteration begins by calculating the gradient at each (x,y) . These are then added to the momentum variables with weight $0 \leq \beta \leq 1$ which is taken as an argument. These momentum-influenced step sizes are then taken from x and y .

Similar to the Polyak algorithm, we implement a stopping criterion of $\|\nabla f(x)\| < \epsilon$.

This avoids noise around convergence although a good choice of β will minimise this.

iii) Adam

```
1 def adam(func, x0, y0, alpha, beta1, beta2, epsilon=1e-6,
2     num_iters=1000):
3     # Initialize x,y, momentum variables and path,f(x,y) lists
4     x, y = x0, y0
5     m_x, m_y = 0, 0
6     v_x, v_y = 0, 0
7     path = [(x, y)]
8     f_vals = []
9
10    for t in range(1, num_iters + 1):
11        # Take gradients for current x, y
12        grad_x = func.dx(x, y)
13        grad_y = func.dy(x, y)
14        # Update momentum variable for gradient
15        m_x = beta1 * m_x + (1 - beta1) * grad_x
16        m_y = beta1 * m_y + (1 - beta1) * grad_y
17        # Update momentum variable for squared gradient
18        v_x = beta2 * v_x + (1 - beta2) * grad_x**2
19        v_y = beta2 * v_y + (1 - beta2) * grad_y**2
20        # Compute bias-corrected estimates
21        m_x_hat = m_x / (1 - beta1**t)
22        m_y_hat = m_y / (1 - beta1**t)
23
24        v_x_hat = v_x / (1 - beta2**t)
25        v_y_hat = v_y / (1 - beta2**t)
26
27        # Update variables
28        x -= alpha * m_x_hat / (sympy.sqrt(v_x_hat) + epsilon)
29        y -= alpha * m_y_hat / (sympy.sqrt(v_y_hat) + epsilon)
30        if func.f(x, y) > 1e20:
31            break
32
33        path.append((x, y))
34        f_vals.append(func.f(x, y))
35
36    return path, f_vals
```

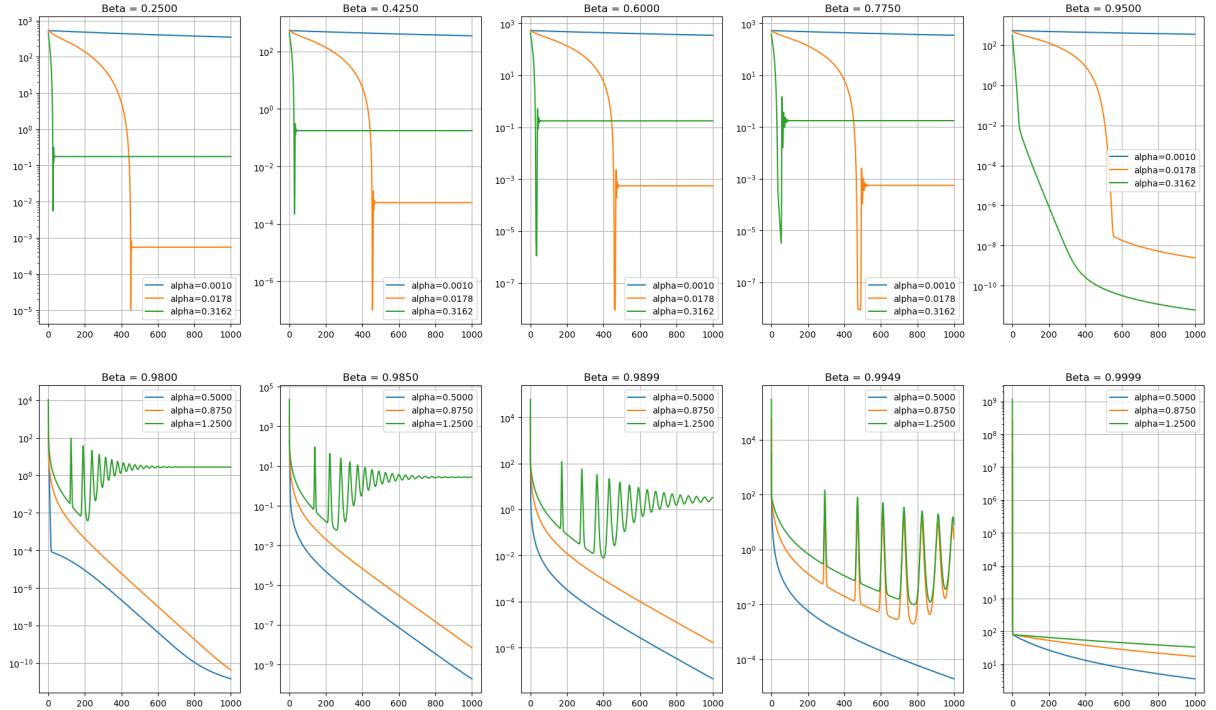
The Adam algorithm is a combination of RMSProp and HeavyBall. It takes the momentum aspect of Heavy ball and applies it to a vector of step sizes corresponding to each variable. This algorithm is implemented by initializing all variables including the momentum variables. Once iteration begins, the first step is to calculate the gradients of x and y . These are used to update the momentum variables based on the weight of β_1 . This is also done with respect to the squared gradients using β_2 . These updated momentum variables are then bias-corrected and used along with the learning rate α to update x and y .

Question b)

i) RMSProp

To find a rough idea of what values of α and β to use, I used an iterative grid search on a linear scale for β and a log scale for α . This was done twice with the second iteration narrowing around the most accurate values.

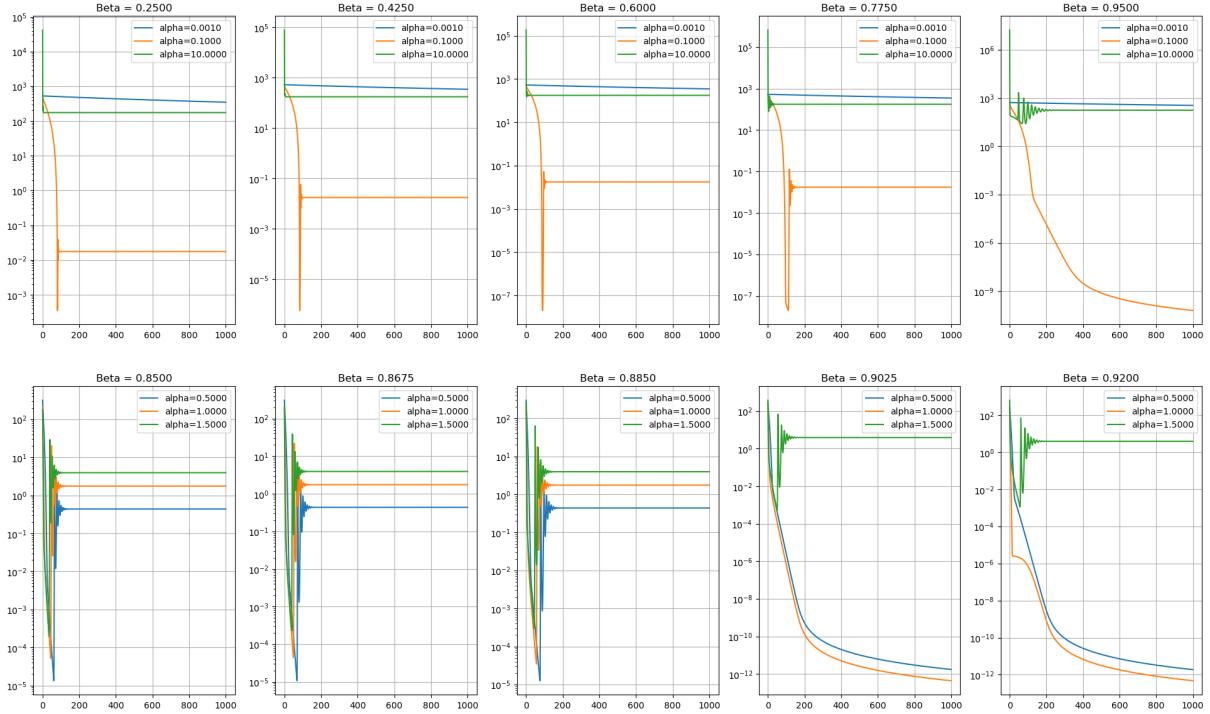
Function 1: alpha vs beta



The top 5 images all show a large spread of β . This gives Adagard-like results. There is a clear difference at later iterations of the chatter being far more prominent in RMSProp than Adagard. This is because the step size isn't reduced enough and too much influence is taken from past gradients. We can see, however, that a higher β does cause a quicker

conversion overall. This is to be expected as the gradient slowing doesn't impact the algorithms with higher β as strongly as those with lower β .

Function 2: alpha vs beta



We can also see that lower values of α give a more stable result. We can see that very low value lead to a slow convergence when β is small, however, for higher β , the stability is far more valuable. The same results remain true for the piecewise function.

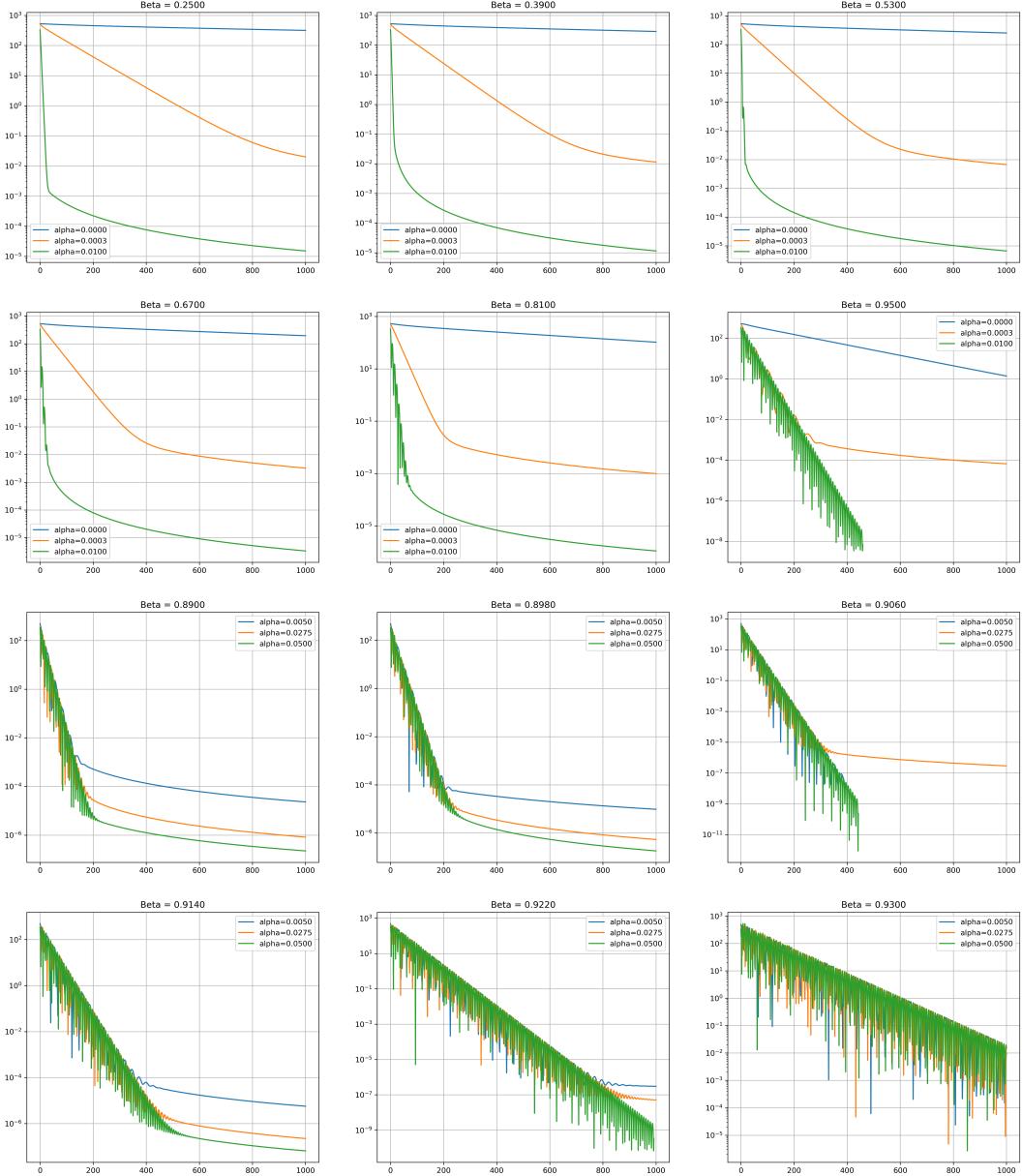
Also, a contour plot here would show a more direct path to the local minimum than a constant or Polyak.

ii) Heavy Ball

Like the RMSProp implementation I used an iterative grid search to find values of α and β that performed well.

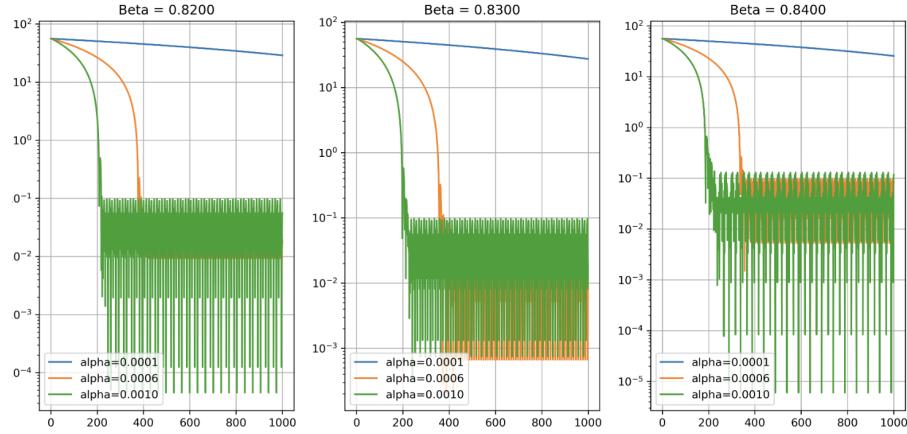
We can see in the initial 6 plots that a lower value of β gives a more stable result. This is because there are going to be less extreme overshooting oscillations of the x and y values. For higher values of β , it takes longer to change the direction of the momentum due to the large impact that a variable will have in the other direction once it hits its minimum for a given variable. We can see that this introduces a lot of oscillation which can be seen by the 'bumps' in the graph. α here is the learning rate which controls rate of convergence to a minimum. However, a larger value of α also introduces a higher likelihood of instability and noise.

Function 1: alpha vs beta



The plot of the second function shows results essentially consistent with those above. In this plot, the bumps/oscillations are more extreme due to the kinks in the function and oscillations don't decay to 0. This is due to the kink around the minimum causing the same gradients to be called repeatedly as (x,y) cross it. These results are all unstable and heavy ball is not a good method for piecewise functions.

A contour plot would show oscillations over and back along the path of the gradient with intensity relative to β , however, I had no room to include these.

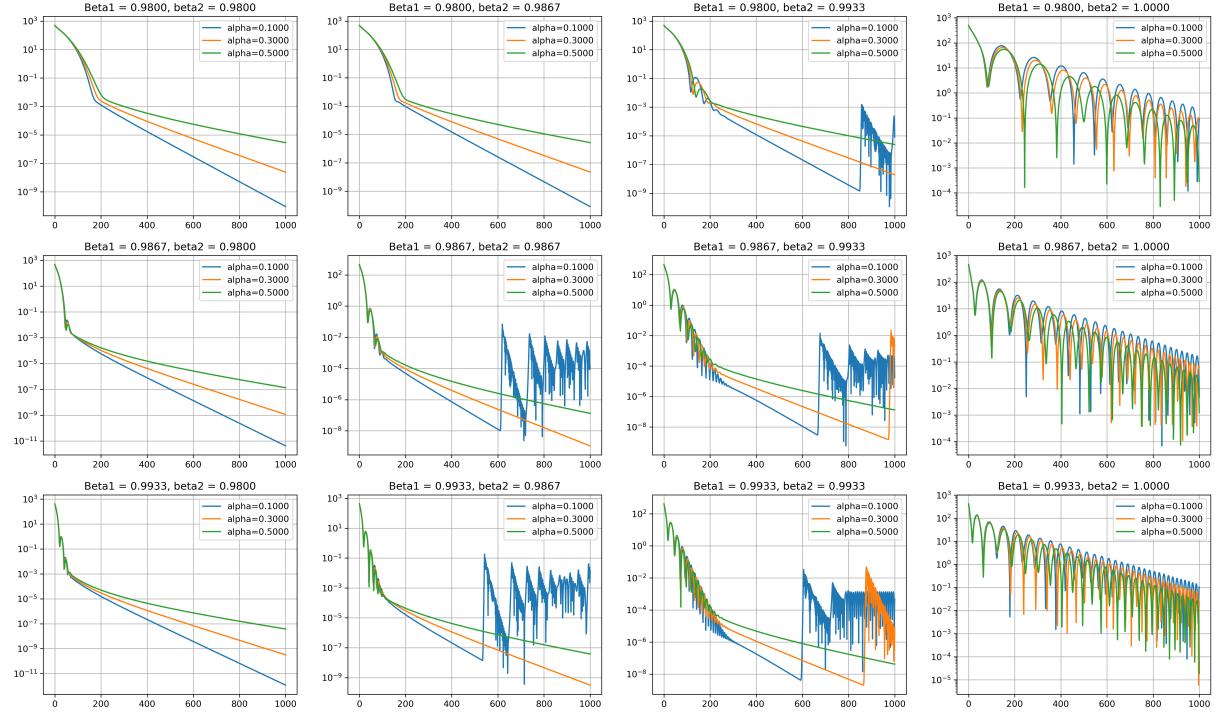


iii) Adam

Like the previous 2 algorithms, we can see that a lower learning rate, α , leads to slower convergence. Although it can't be seen in this example, if α is too big it will cause the solution to diverge.

It can be seen that the general 'bumpy' behaviour only occurs for $\beta_2 \approx 1$. This happens because the increased weight given to the past squared gradients amplifies the oscillations. When $\beta_2 \approx 1$, these values essentially don't decay, or at least decay extremely slowly. The β_1 values also cause more oscillations as they increase but it is far less impactful than β_2 .

Function 1: alpha vs beta1 vs beta2



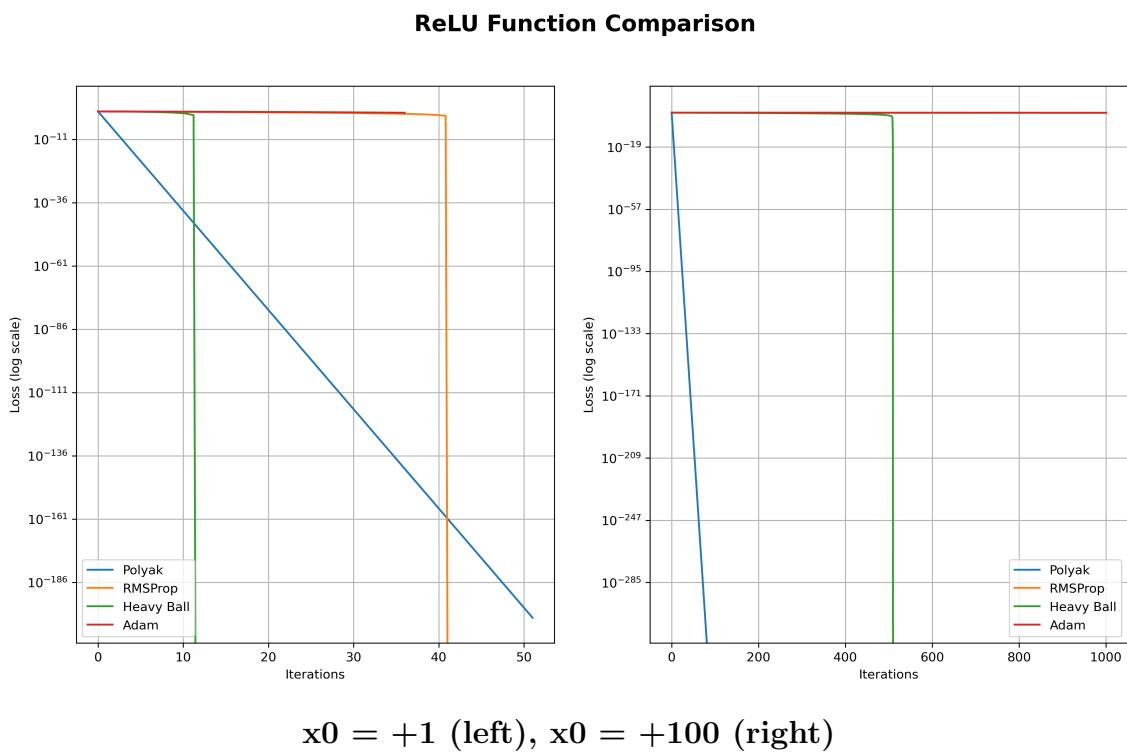
These oscillations would also be very apparent in a contour plot.

Question c)

i) ReLU (-1)

Each of the four algorithms give values of $(x, f(x)) = (-1, 0)$ for all iterations. This is because every one of the functions require the $\frac{df(x)}{dx} \neq 0$ for the x value to change. Each change to the step size is some multiplicative function of the gradient and hence each time the step size changes, it changes by 0.

ii,iii) ReLU (+1, +100)



On the left image, we would expect polyak and heavy ball to fall to 0 essentially after the first iteration. $1 - (\frac{1-0}{1^2})(1) = 0$. However, the reason it takes longer is due to the ϵ that makes it logarithmic. This could be omitted. The adam and RMSProp functions would require a higher value of α to increase the stepsize and decrease to 0 quicker.

Code

```
1 import numpy as np
2 import sympy
3 import matplotlib.pyplot as plt
4
```

```

5
6 class func1:
7     def __init__(self):
8         self.x_sym, self.y_sym = sympy.symbols('x y', real=True)
9
10    def f(self, x, y):
11        return 5 * (x - 3)**4 + 7 * (y - 9)**2
12
13    def dx(self, x, y):
14        df_dx = sympy.diff(self.f(self.x_sym, self.y_sym),
15                            self.x_sym)
16        return df_dx.subs({self.x_sym: x, self.y_sym: y})
17
18    def dy(self, x, y):
19        df_dy = sympy.diff(self.f(self.x_sym, self.y_sym),
20                            self.y_sym)
21        return df_dy.subs({self.x_sym: x, self.y_sym: y})
22
23 class func2:
24     def __init__(self):
25         self.x_sym, self.y_sym = sympy.symbols('x y', real=True)
26
27     def f(self, x, y):
28         return sympy.Max(x - 3, 0) + 7 * sympy.Abs(y - 9)
29
30     def dx(self, x, y):
31         df_dx = sympy.diff(self.f(self.x_sym, self.y_sym),
32                             self.x_sym)
33         return df_dx.subs({self.x_sym: x, self.y_sym: y})
34
35     def dy(self, x, y):
36         df_dy = sympy.diff(self.f(self.x_sym, self.y_sym),
37                             self.y_sym)
38         return df_dy.subs({self.x_sym: x, self.y_sym: y})

```

```

1 def polyak_GD(func, x0, y0, f_star=0, epsilon=0.0001,
2               max_iters=1000, tol=100):
3     # Initialize x,y and path and f(x,y) lists
4     x, y = x0, y0
5     path = [(x, y)]

```

```

6
7     # Implement formula to find alpha and update at each step.
8     for i in range(max_iters):
9         df_dx = func.dx(x, y)
10        df_dy = func.dy(x, y)
11        alpha = (func.f(x, y) - f_star) / (df_dx**2 + df_dy**2 +
12            epsilon)
13        # Update x and y based on polyak step - optimal alpha
14        x -= alpha * df_dx
15        y -= alpha * df_dy
16        # Append both (x, y) point and f value to lists
17        path.append((x, y))
18        f_vals.append(func.f(x, y))

19        # Check stopping conditions for alpha too high or too low
20        if alpha > tol:
21            print(f"Stopping after {i} Iterations as alpha too
22                  large")
23            return path, f_vals
24        if sympy.sqrt(func.dx(x, y)**2 + func.dy(x, y)**2) <
25            epsilon:
26            return path, f_vals

27        print(f"No sufficient convergence by iteration
28              {max_iters}\nFinished at (x, y): ({x:0.4f}, {y:0.4f})")
29
30
31    def rmsprop(func, x0, y0, alpha0, beta, epsilon=0.0001,
32               max_iters=1000):
33        # Initialize x,y, sums and path,f(x,y) lists
34        x, y = x0, y0
35        sum_dx = 0
36        sum_dy = 0
37        path = [(x, y)]
38        f_vals = []

39        for t in range(1, max_iters + 1):
40            # Take gradients for current x, y
41            grad_x = func.dx(x, y)

```

```

42     grad_y = func.dy(x, y)

43

44     # Update denominator for x and y using the following
45     sum_dx = beta * sum_dx + (1 - beta) * grad_x**2
46     sum_dy = beta * sum_dy + (1 - beta) * grad_y**2

47

48     # Compute the step sizes for x and y
49     step_size_x = alpha0 / sympy.sqrt(sum_dx + epsilon)
50     step_size_y = alpha0 / sympy.sqrt(sum_dy + epsilon)

51

52     # Update x and y
53     x -= step_size_x * grad_x
54     y -= step_size_y * grad_y

55

56     path.append((x, y))
57     f_vals.append(func.f(x, y))

58

59     return path, f_vals

60

61

62 def heavy_ball(func, x0, y0, alpha, beta, num_iters=1000,
63               tol=1e-5, stopping=True):
64     # Initialize x,y, momentum variables and path,f(x,y) lists
65     x, y = x0, y0
66     ux, uy = 0, 0
67     path = [(x, y)]
68     f_vals = []

69     for i in range(num_iters):
70         # Take gradients for current x, y
71         grad_x = func.dx(x, y)
72         grad_y = func.dy(x, y)

73

74         # Update momentum variable with current and past momentum
75         ux = beta * ux + alpha * grad_x
76         uy = beta * uy + alpha * grad_y

77

78         # Update x and y
79         x -= ux
80         y -= uy

```

```

82     path.append((x, y)) # Store the new values
83     f_vals.append(func.f(x, y))
84
85     if func.f(x, y) > 1e20:
86         break
87
88     # Check if the norm of the gradient is below the
89     # tolerance
90     if sympy.sqrt(grad_x**2 + grad_y**2) < tol and stopping
91     == True:
92         break
93
94
95
96 def adam(func, x0, y0, alpha, beta1, beta2, epsilon=1e-6,
97 num_iters=1000):
98     # Initialize x,y, momentum variables and path,f(x,y) lists
99     x, y = x0, y0
100    m_x, m_y = 0, 0
101    v_x, v_y = 0, 0
102    path = [(x, y)]
103    f_vals = []
104
105    for t in range(1, num_iters + 1):
106        # Take gradients for current x, y
107        grad_x = func.dx(x, y)
108        grad_y = func.dy(x, y)
109        # Update momentum variable for gradient
110        m_x = beta1 * m_x + (1 - beta1) * grad_x
111        m_y = beta1 * m_y + (1 - beta1) * grad_y
112        # Update momentum variable for squared gradient
113        v_x = beta2 * v_x + (1 - beta2) * grad_x**2
114        v_y = beta2 * v_y + (1 - beta2) * grad_y**2
115        # Compute bias-corrected estimates
116        m_x_hat = m_x / (1 - beta1**t)
117        m_y_hat = m_y / (1 - beta1**t)
118
119        v_x_hat = v_x / (1 - beta2**t)
120        v_y_hat = v_y / (1 - beta2**t)

```

```

120
121     # Update variables
122     x -= alpha * m_x_hat / (sympy.sqrt(v_x_hat) + epsilon)
123     y -= alpha * m_y_hat / (sympy.sqrt(v_y_hat) + epsilon)
124     if func.f(x, y) > 1e20:
125         break
126
127     path.append((x, y))
128     f_vals.append(func.f(x, y))
129
130 return path, f_vals

```

```

1 def RMS_grid_search(func, x0, y0, alphas, betas,
2                     epsilon=1e-4, num_iters=1000):
3     best_alpha = None
4     best_beta = None
5     best_f_value = 1e9
6
7     for alpha in alphas:
8         for beta in betas:
9             _, f_vals = rmsprop(func, x0, y0, alpha, beta,
10                      epsilon, num_iters)
11             f_value = f_vals[-1]
12             if f_value < best_f_value:
13                 best_f_value = f_value
14                 best_alpha = alpha
15                 best_beta = beta
16
17     return best_alpha, best_beta, best_f_value
18
19
20 alphas = np.logspace(-3, 2, 5)
21 betas = np.linspace(0.25, 0.99, 5)
22
23 # Perform the grid search
24 best_alpha, best_beta, best_f_value = RMS_grid_search(function1,
25             x0, y0, alphas, betas)
26
27 print(f"Best alpha: {best_alpha:.4f}")
28 print(f"Best beta: {best_beta:.4f}")
29
30 alphas = np.linspace(0.5, 2, 5)
31 betas = np.linspace(0.98, 0.999, 5)

```

```

28
29 # Perform the grid search
30 best_alpha, best_beta, best_f_value = RMS_grid_search(function1,
31             x0, y0, alphas, betas)
32
33 print(f"Best alpha: {best_alpha:.4f}")
34 print(f"Best beta: {best_beta:.4f}")
35
36 alphas = np.logspace(-3, 2, 5)
37 betas = np.linspace(0.25, 0.95, 5)
38 paths_1 = [[rmsprop(function1, x0, y0, alpha0=alpha,
39                     beta=beta)[0] for beta in betas] for alpha in alphas]
40 f_vals_1 = [[rmsprop(function1, x0, y0, alpha0=alpha,
41                     beta=beta)[1] for beta in betas] for alpha in alphas]
42
43 alphas = np.linspace(0.5, 2, 5)
44 betas = np.linspace(0.98, 0.9999, 5)
45 paths_2 = [[rmsprop(function1, x0, y0, alpha0=alpha,
46                     beta=beta)[0] for beta in betas] for alpha in alphas]
47 f_vals_2 = [[rmsprop(function1, x0, y0, alpha0=alpha,
48                     beta=beta)[1] for beta in betas] for alpha in alphas]
49
50 x = np.linspace(0, 1000, 1000)
51 alphas = np.logspace(-3, 2, 5)
52 betas = np.linspace(0.25, 0.95, 5)
53
54 plt.figure(figsize=(24, 14))
55 plt.suptitle("Function 1: alpha vs beta", fontsize=24,
56               fontweight='bold')
57
58 for i in range(5):
59     plt.subplot(2, 5, i+1)
60     plt.semilogy(x, f_vals_1[0][i],
61                   label=f'alpha={alphas[0]:0.4f}')
62     plt.semilogy(x, f_vals_1[1][i],
63                   label=f'alpha={alphas[1]:0.4f}')
64     plt.semilogy(x, f_vals_1[2][i],
65                   label=f'alpha={alphas[2]:0.4f}')
66     plt.title(f'Beta = {betas[i]:0.4f}')
67     plt.legend()
68     plt.grid(visible=True)

```

```

60
61
62 alphas = np.linspace(0.5, 2, 5)
63 betas = np.linspace(0.98, 0.9999, 5)
64 for i in range(5):
65     plt.subplot(2, 5, i+6)
66     plt.semilogy(x, f_vals_2[0][i],
67                   label=f'alpha={alphas[0]:0.4f}')
68     plt.semilogy(x, f_vals_2[1][i],
69                   label=f'alpha={alphas[1]:0.4f}')
70     plt.semilogy(x, f_vals_2[2][i],
71                   label=f'alpha={alphas[2]:0.4f}')
72     plt.title(f'Beta = {betas[i]:0.4f}')
73     plt.legend()
74     plt.grid(visible=True)

75
76 plt.show()

```

```

1 alphas = np.logspace(-3, 1, 5)
2 betas = np.linspace(0.25, 0.95, 5)
3
4 # Perform the grid search
5 best_alpha, best_beta, best_f_value = RMS_grid_search(function2,
6             x0, y0, alphas, betas)
7
8 print(f"Best alpha: {best_alpha:0.4f}")
9 print(f"Best beta: {best_beta:0.4f}")
10
11 alphas = np.linspace(0.5, 1.5, 5)
12 betas = np.linspace(0.85, 0.92, 5)
13
14 # Perform the grid search
15 best_alpha, best_beta, best_f_value = RMS_grid_search(function2,
16             x0, y0, alphas, betas)
17
18 print(f"Best alpha: {best_alpha:0.4f}")
19 print(f"Best beta: {best_beta:0.4f}")
20
21 alphas = np.logspace(-3, 1, 3)
22 betas = np.linspace(0.25, 0.95, 5)
23 paths_1 = [[rmsprop(function1, x0, y0, alpha0=alpha,
24                     beta=beta)[0] for beta in betas] for alpha in alphas]

```

```

22 f_vals_1 = [[rmsprop(function1, x0, y0, alpha0=alpha,
23   beta=beta)[1] for beta in betas] for alpha in alphas]
24
25 alphas = np.linspace(0.5, 1.5, 3)
26 betas = np.linspace(0.85, 0.92, 5)
27 paths_2 = [[rmsprop(function1, x0, y0, alpha0=alpha,
28   beta=beta)[0] for beta in betas] for alpha in alphas]
29 f_vals_2 = [[rmsprop(function1, x0, y0, alpha0=alpha,
30   beta=beta)[1] for beta in betas] for alpha in alphas]
31
32 x = np.linspace(0, 1000, 1000)
33 alphas = np.logspace(-3, 1, 3)
34 betas = np.linspace(0.25, 0.95, 5)
35 plt.figure(figsize=(24, 14))
36 plt.suptitle("Function 2: alpha vs beta", fontsize=24,
37   fontweight='bold')
38
39 for i in range(5):
40   plt.subplot(2, 5, i+1)
41   plt.semilogy(x, f_vals_1[0][i],
42     label=f'alpha={alphas[0]:0.4f}')
43   plt.semilogy(x, f_vals_1[1][i],
44     label=f'alpha={alphas[1]:0.4f}')
45   plt.semilogy(x, f_vals_1[2][i],
46     label=f'alpha={alphas[2]:0.4f}')
47   plt.title(f'Beta = {betas[i]:0.4f}')
48   plt.legend()
49   plt.grid(visible=True)
50
51
52 alphas = np.linspace(0.5, 1.5, 3)
53 betas = np.linspace(0.85, 0.92, 5)
54 for i in range(5):
55   plt.subplot(2, 5, i+6)
56   plt.semilogy(x, f_vals_2[0][i],
57     label=f'alpha={alphas[0]:0.4f}')
58   plt.semilogy(x, f_vals_2[1][i],
59     label=f'alpha={alphas[1]:0.4f}')
60   plt.semilogy(x, f_vals_2[2][i],
61     label=f'alpha={alphas[2]:0.4f}')
62   plt.title(f'Beta = {betas[i]:0.4f}')

```

```

53     plt.legend()
54     plt.grid(visible=True)
55
56 plt.show()

1 def HB_grid_search(func, x0, y0, alphas, betas, tol=1e-4,
2     num_iters=1000):
3     best_alpha = None
4     best_beta = None
5     best_f_value = 1e9
6
7     for alpha in alphas:
8         for beta in betas:
9             _, f_vals = heavy_ball(func, x0, y0, alpha, beta,
10                 num_iters=num_iters, tol=tol, stopping=False)
11             f_value = f_vals[-1]
12             if f_value < best_f_value:
13                 best_f_value = f_value
14                 best_alpha = alpha
15                 best_beta = beta
16
17     return best_alpha, best_beta, best_f_value
18
19
20 # Perform the grid search
21 best_alpha, best_beta, best_f_value = HB_grid_search(function1,
22             x0, y0, alphas, betas)
23
24 print(f"Best alpha: {best_alpha:.4f}")
25 print(f"Best beta: {best_beta:.4f}")
26
27 alphas = np.linspace(0.001, 0.02, 6)
28 betas = np.linspace(0.91, 0.99, 6)
29
30 # Perform the grid search
31 best_alpha, best_beta, best_f_value = HB_grid_search(function1,
32             x0, y0, alphas, betas)
33
34 print(f"Best alpha: {best_alpha:.4f}")
35 print(f"Best beta: {best_beta:.4f}")

```

```

34
35 alphas = np.logspace(-5, -2, 3)
36 betas = np.linspace(0.25, 0.95, 6)
37 paths_1 = [[heavy_ball(function1, x0, y0, alpha=alpha,
38     beta=beta, stopping=True)[0] for beta in betas] for alpha in
39     alphas]
40 f_vals_1 = [[heavy_ball(function1, x0, y0, alpha=alpha,
41     beta=beta, stopping=True)[1] for beta in betas] for alpha in
42     alphas]
43
44 alphas = np.logspace(-5, -2, 3)
45 betas = np.linspace(0.25, 0.95, 6)
46
47 plt.figure(figsize=(24, 28), dpi=300)
48 plt.suptitle("Function 1: alpha vs beta", fontsize=24,
49     fontweight='bold')
50
51 for i in range(6):
52     plt.subplot(4, 3, i+1)
53     plt.semilogy(np.linspace(0, len(f_vals_1[0][i]),
54         len(f_vals_1[0][i])), f_vals_1[0][i],
55         label=f'alpha={alphas[0]:0.4f}')
56     plt.semilogy(np.linspace(0, len(f_vals_1[1][i]),
57         len(f_vals_1[1][i])), f_vals_1[1][i],
58         label=f'alpha={alphas[1]:0.4f}')
59     plt.semilogy(np.linspace(0, len(f_vals_1[2][i]),
60         len(f_vals_1[2][i])), f_vals_1[2][i],
61         label=f'alpha={alphas[2]:0.4f}')
62     plt.title(f'Beta = {betas[i]:0.4f}')
63     plt.legend()
64     plt.grid(visible=True)

```

```

60
61
62 alphas = np.linspace(0.005, 0.05, 3)
63 betas = np.linspace(0.89, 0.93, 6)
64 for i in range(6):
65     plt.subplot(4, 3, i+7)
66     plt.semilogy(np.linspace(0, len(f_vals_2[0][i]),
67                   len(f_vals_2[0][i])), f_vals_2[0][i],
68                   label=f'alpha={alphas[0]:0.4f}')
69     plt.semilogy(np.linspace(0, len(f_vals_2[1][i]),
70                   len(f_vals_2[1][i])), f_vals_2[1][i],
71                   label=f'alpha={alphas[1]:0.4f}')
72     plt.semilogy(np.linspace(0, len(f_vals_2[2][i]),
73                   len(f_vals_2[2][i])), f_vals_2[2][i],
74                   label=f'alpha={alphas[2]:0.4f}')
75     plt.title(f'Beta = {betas[i]:0.4f}')
76     plt.legend()
77     plt.grid(visible=True)
78
79 plt.show()

```

```

1      alphas = np.logspace(-5, -2, 6)
2      betas = np.linspace(0.25, 0.95, 6)
3
4 # Perform the grid search
5 best_alpha, best_beta, best_f_value = HB_grid_search(function2,
6           x0, y0, alphas, betas)
7
8 print(f"Best alpha: {best_alpha:0.4f}")
9 print(f"Best beta: {best_beta:0.4f}")
10
11 alphas = np.linspace(0.0001, 0.001, 6)
12 betas = np.linspace(0.82, 0.86, 6)
13
14 # Perform the grid search
15 best_alpha, best_beta, best_f_value = HB_grid_search(function2,
16           x0, y0, alphas, betas)
17
18 print(f"Best alpha: {best_alpha:0.4f}")
19 print(f"Best beta: {best_beta:0.4f}")
20
21 alphas = np.logspace(-5, -2, 3)

```

```

20 betas = np.linspace(0.25, 0.95, 5)
21 paths_1 = [[heavy_ball(function2, x0, y0, alpha=alpha,
22     beta=beta, stopping=False)[0] for beta in betas] for alpha in
23     alphas]
24 f_vals_1 = [[heavy_ball(function2, x0, y0, alpha=alpha,
25     beta=beta, stopping=False)[1] for beta in betas] for alpha in
26     alphas]
27
28 alphas = np.linspace(0.0001, 0.001, 3)
29 betas = np.linspace(0.82, 0.86, 5)
30 paths_2 = [[heavy_ball(function2, x0, y0, alpha=alpha,
31     beta=beta, stopping=False)[0] for beta in betas] for alpha in
32     alphas]
33 f_vals_2 = [[heavy_ball(function2, x0, y0, alpha=alpha,
34     beta=beta, stopping=False)[1] for beta in betas] for alpha in
35     alphas]
36
37 alphas = np.logspace(-5, -2, 3)
38 betas = np.linspace(0.25, 0.95, 5)
39
40 plt.figure(figsize=(24, 14), dpi=300)
41 plt.suptitle("Function 2: alpha vs beta", fontsize=24,
42     fontweight='bold')
43
44 for i in range(5):
45     plt.subplot(2, 5, i+1)
46     plt.semilogy(np.linspace(0, len(f_vals_1[0][i]),
47         len(f_vals_1[0][i])), f_vals_1[0][i],
48         label=f'alpha={alphas[0]:0.4f}')
49     plt.semilogy(np.linspace(0, len(f_vals_1[1][i]),
50         len(f_vals_1[1][i])), f_vals_1[1][i],
51         label=f'alpha={alphas[1]:0.4f}')
52     plt.semilogy(np.linspace(0, len(f_vals_1[2][i]),
53         len(f_vals_1[2][i])), f_vals_1[2][i],
54         label=f'alpha={alphas[2]:0.4f}')
55     plt.title(f'Beta = {betas[i]:0.4f}')
56     plt.legend()
57     plt.grid(visible=True)
58
59
60 alphas = np.linspace(0.0001, 0.001, 3)

```

```

46 betas = np.linspace(0.82, 0.86, 5)
47 for i in range(5):
48     plt.subplot(2, 5, i+6)
49     plt.semilogy(np.linspace(0, len(f_vals_2[0][i]),
50                   len(f_vals_2[0][i])), f_vals_2[0][i],
51                   label=f'alpha={alphas[0]:0.4f}')
52     plt.semilogy(np.linspace(0, len(f_vals_2[1][i]),
53                   len(f_vals_2[1][i])), f_vals_2[1][i],
54                   label=f'alpha={alphas[1]:0.4f}')
55     plt.semilogy(np.linspace(0, len(f_vals_2[2][i]),
56                   len(f_vals_2[2][i])), f_vals_2[2][i],
57                   label=f'alpha={alphas[2]:0.4f}')
58 plt.title(f'Beta = {betas[i]:0.4f}')
59 plt.legend()
60 plt.grid(visible=True)
61
62 plt.show()

```

```

1 def adam_grid_search(func, x0, y0, alphas, betas, betas2,
2                     tol=1e-4, num_iters=1000):
3     best_alpha = None
4     best_beta1 = None
5     best_beta2 = None
6     best_f_value = 1e9
7
8     for alpha in alphas:
9         for beta1 in betas:
10            for beta2 in betas2:
11                _, f_vals = adam(func, x0, y0, alpha, beta1,
12                                  beta2, epsilon=1e-6, num_iters=1000)
13                f_value = f_vals[-1]
14                if f_value < best_f_value:
15                    best_f_value = f_value
16                    best_alpha = alpha
17                    best_beta1 = beta1
18                    best_beta2 = beta2
19
20
21     return best_alpha, best_beta1, best_beta2, best_f_value
22
23 alphas = np.logspace(-5, -2, 4)
24 betas = np.linspace(0.25, 0.9999, 4)

```

```

23 beta2s = np.linspace(0.25, 0.9999, 4)
24
25 # Perform the grid search
26 best_alpha, best_beta1, best_beta2, best_f_value =
27     adam_grid_search(function1, x0, y0, alphas, betas=beta1s,
28                       beta2s=beta2s)
29
30 print(f"Best alpha: {best_alpha:0.4f}")
31 print(f"Best beta1: {best_beta1:0.4f}")
32 print(f"Best beta2: {best_beta2:0.4f}")
33
34 alphas = np.logspace(-5, -2, 4)
35 beta1s = np.linspace(0.8, 92, 4)
36 beta2s = np.linspace(0.58, 0.68, 4)
37
38 # Perform the grid search
39 best_alpha, best_beta1, best_beta2, best_f_value =
40     adam_grid_search(function1, x0, y0, alphas, betas=beta1s,
41                       beta2s=beta2s)
42
43 print(f"Best alpha: {best_alpha:0.4f}")
44 print(f"Best beta1: {best_beta1:0.4f}")
45 print(f"Best beta2: {best_beta2:0.4f}")
46
47 alphas = np.linspace(0.1, 0.5, 3)
48 beta1s = np.linspace(0.8, 0.99, 4)
49 beta2s = np.linspace(0.98, 0.99999, 4)
50 paths_1 = [[[adam(function1, x0, y0, alpha, beta1, beta2,
51                   epsilon=1e-6, num_iters=1000)[0] for beta1 in beta1s] for
52                   beta2 in beta2s] for alpha in alphas]
53 f_vals_1 = [[[adam(function1, x0, y0, alpha, beta1, beta2,
54                   epsilon=1e-6, num_iters=1000)[1] for beta1 in beta1s] for
55                   beta2 in beta2s] for alpha in alphas]
56
57
58 plt.figure(figsize=(24, 14), dpi=300)
59 plt.suptitle("Function 1: alpha vs beta1 vs beta2", fontsize=24,
60               fontweight='bold')
61
62 for i in range(4):
63     plt.subplot(3, 4, i+1)

```

```

55     plt.semilogy(np.linspace(0, len(f_vals_1[0][0][i]),
56                   len(f_vals_1[0][0][i])), f_vals_1[0][0][i],
57                   label=f'alpha={alphas[0]:0.4f}')
58     plt.semilogy(np.linspace(0, len(f_vals_1[0][1][i]),
59                   len(f_vals_1[0][1][i])), f_vals_1[0][1][i],
60                   label=f'alpha={alphas[1]:0.4f}')
61     plt.semilogy(np.linspace(0, len(f_vals_1[0][2][i]),
62                   len(f_vals_1[0][2][i])), f_vals_1[0][2][i],
63                   label=f'alpha={alphas[2]:0.4f}')
64     plt.title(f'Beta1 = {beta2s[0]:0.4f}, beta2 =
65               {beta2s[i]:0.4f}')
66     plt.legend()
67     plt.grid(visible=True)
68
69 for i in range(4):
70     plt.subplot(3, 4, i+5)
71     plt.semilogy(np.linspace(0, len(f_vals_1[1][0][i]),
72                   len(f_vals_1[1][0][i])), f_vals_1[1][0][i],
73                   label=f'alpha={alphas[0]:0.4f}')
74     plt.semilogy(np.linspace(0, len(f_vals_1[1][1][i]),
75                   len(f_vals_1[1][1][i])), f_vals_1[1][1][i],
76                   label=f'alpha={alphas[1]:0.4f}')
77     plt.semilogy(np.linspace(0, len(f_vals_1[1][2][i]),
78                   len(f_vals_1[1][2][i])), f_vals_1[1][2][i],
79                   label=f'alpha={alphas[2]:0.4f}')
80     plt.title(f'Beta1 = {beta2s[1]:0.4f}, beta2 =
81               {beta2s[i]:0.4f}')
82     plt.legend()
83     plt.grid(visible=True)
84
85 for i in range(4):
86     plt.subplot(3, 4, i+9)
87     plt.semilogy(np.linspace(0, len(f_vals_1[2][0][i]),
88                   len(f_vals_1[2][0][i])), f_vals_1[2][0][i],
89                   label=f'alpha={alphas[0]:0.4f}')
90     plt.semilogy(np.linspace(0, len(f_vals_1[2][1][i]),
91                   len(f_vals_1[2][1][i])), f_vals_1[2][1][i],
92                   label=f'alpha={alphas[1]:0.4f}')
93     plt.semilogy(np.linspace(0, len(f_vals_1[2][2][i]),
94                   len(f_vals_1[2][2][i])), f_vals_1[2][2][i],
95                   label=f'alpha={alphas[2]:0.4f}')

```

```

76     plt.title(f'Beta1 = {beta2s[2]:0.4f}, beta2 =
77                 {beta2s[i]:0.4f}')
78     plt.legend()
79     plt.grid(visible=True)
80

```

```

1      class ReLU:
2          def f(self, x):
3              return max(0, x)
4          def dx(self, x):
5              return 1 if x > 0 else 0
6
7          def polyak_GD_relu(func, x0, f_star=0, epsilon=0.0001,
8                  max_iters=50, tol=100):
9              # Initialize x,y and path and f(x,y) lists
10             x = x0
11             path = [x]
12             f_vals = []
13
14             # Implement formula to find alpha and update at each step.
15             for i in range(max_iters):
16                 df_dx = func.dx(x)
17                 alpha = (func.f(x) - f_star) / (df_dx**2 + epsilon)
18                 # Update x and y based on polyak step - optimal alpha
19                 x -= alpha * df_dx
20                 # Append both (x, y) point and f value to lists
21                 path.append(x)
22                 f_vals.append(func.f(x))
23
24             print(f"No sufficient convergence by iteration
25                   {max_iters}\nFinished at x: ({x:0.4f}, f(x) =
26                   {f_vals[-1]:0.4f})")
27             return path, f_vals
28
29         def rmsprop_relu(func, x0, alpha0, beta, epsilon=0.0001,
30                         max_iters=50):
31             # Initialize x,y, sums and path,f(x,y) lists
32             x = x0
33             sum_dx = 0
34             path = [x]
35             f_vals = []

```

```

32
33     for t in range(1, max_iters + 1):
34         # Take gradients for current x, y
35         grad_x = func.dx(x)
36
37         # Update denominator for x and y using the following
38         sum_dx = beta * sum_dx + (1 - beta) * grad_x**2
39
40         # Compute the step sizes for x and y
41         step_size_x = alpha0 / sympy.sqrt(sum_dx + epsilon)
42
43         # Update x and y
44         x -= step_size_x * grad_x
45
46         path.append(x)
47         f_vals.append(func.f(x))
48
49     return path, f_vals
50
51 def heavy_ball_relu(func, x0, alpha, beta, num_iters=50,
52                     tol=1e-5, stopping=True):
53     # Initialize x,y, momentum variables and path,f(x,y) lists
54     x = x0
55     ux = 0
56     path = [x]
57     f_vals = []
58
59     for i in range(num_iters):
60         # Take gradients for current x, y
61         grad_x = func.dx(x)
62
63         # Update momentum variable with current and past momentum
64         ux = beta * ux + alpha * grad_x
65
66         # Update x and y
67         x -= ux
68
69         path.append(x) # Store the new values
70         f_vals.append(func.f(x))
71
72     print(f"Finished at x: ({x:.4f}), f(x) = {f_vals[-1]:.4f}")

```

```

72     return path, f_vals
73
74 def adam_relu(func, x0, alpha, beta1, beta2, epsilon=1e-6,
75   num_iters=35):
76     # Initialize x,y, momentum variables and path,f(x,y) lists
77     x = x0
78     m_x = 0
79     v_x = 0
80     path = [x]
81     f_vals = []
82
83     for t in range(1, num_iters + 1):
84         grad_x = func.dx(x)
85
86         m_x = beta1 * m_x + (1 - beta1) * grad_x
87         v_x = beta2 * v_x + (1 - beta2) * grad_x**2
88
89         m_x_hat = m_x / (1 - beta1**t)
90         v_x_hat = v_x / (1 - beta2**t)
91
92         x -= alpha * m_x_hat / (sympy.sqrt(v_x_hat) + epsilon)
93
94         path.append(x)
95         f_vals.append(func.f(x))
96
97     return path, f_vals
98
99 relu_func = ReLU()
100 x0 = 1
101 path, f_vals = polyak_GD_relu(relu_func, x0, max_iters=30)
102 f_vals
103 path, f_vals = rmsprop_relu(relu_func, x0, 0.02, 0.9,
104   max_iters=30)
105 f_vals
106 path, f_vals = heavy_ball_relu(relu_func, x0, 0.02, 0.9,
107   num_iters=30)
108 f_vals
109 path

```

```

109 plt.figure(figsize=(15, 8), dpi=300)
110 plt.suptitle("ReLU Function Comparison", fontsize=18,
111     fontweight='bold')
112
113 x0 = 1
114 plt.subplot(1,2,1)
115 plt.semilogy(np.linspace(0, len(polyak_GD_relu(relu_func,
116     x0)[0]), len(polyak_GD_relu(relu_func, x0)[0])),
117     polyak_GD_relu(relu_func, x0)[0], label='Polyak')
118 plt.semilogy(np.linspace(0, len(rmsprop_relu(relu_func, x0,
119     0.02, 0.9)[0]), len(rmsprop_relu(relu_func, x0, 0.02,
120     0.9)[0])), rmsprop_relu(relu_func, x0, 0.02, 0.9)[0],
121     label='RMSProp')
122
123 x0 = 100
124
125 plt.subplot(1,2, 2)
126 plt.semilogy(np.linspace(0, len(polyak_GD_relu(relu_func, x0,
127     max_iters=1000)[0]), len(polyak_GD_relu(relu_func, x0,
128     max_iters=1000)[0])), polyak_GD_relu(relu_func, x0,
129     max_iters=1000)[0], label='Polyak')
130 plt.semilogy(np.linspace(0, len(rmsprop_relu(relu_func, x0,
131     0.02, 0.9,max_iters=1000)[0]), len(rmsprop_relu(relu_func,
132     x0, 0.02, 0.9,max_iters=1000)[0])), rmsprop_relu(relu_func,
133     x0, 0.02, 0.9,max_iters=1000)[0], label='RMSProp')
134 plt.semilogy(np.linspace(0, len(heavy_ball_relu(relu_func, x0,
135     0.02, 0.9, num_iters=1000)[0]),
136     len(heavy_ball_relu(relu_func, x0, 0.02, 0.9,
137     num_iters=1000)[0])), heavy_ball_relu(relu_func, x0, 0.02,
138     num_iters=1000)[0], label='Heavy Ball')
```

```
    0.9, num_iters=1000)[0], label='Heavy Ball')
129 plt.semilogy(np.linspace(0, len(adam_relu(relu_func, x0, 0.02,
    0.9, 0.999, num_iters=1000)[0]), len(adam_relu(relu_func, x0,
    0.02, 0.9, 0.999, num_iters=1000)[0])), adam_relu(relu_func,
    x0, 0.02, 0.9, 0.999, num_iters=1000)[0], label='Adam')
130 plt.legend()
131 plt.grid(visible=True)
132 plt.xlabel('Iterations')
133 plt.ylabel('Loss (log scale)')
134 plt.show()
```