

Department of Electrical Engineering

Indian Institute of Technology, Kharagpur

Algorithms, AI and ML Laboratory (EE22202)

Spring, 2025-26

Report 1: Gradient-based Algorithms for Optimization

Name: Arijit Dey

Roll No: 24IE10001

Gradient-based Algorithms for Optimization

Consider the following unconstrained optimization problem:

$$\min_{x \in R} f(x) = 0.5 \left(\sum_{i=1}^{n-1} x_i^2 \right) + 0.5 \kappa x_n^2$$

where κ is a scalar parameter.

1 Initial Assumptions & Function Definition

Unless explicitly stated by the question, we use the following values for the unknown variables:

- $n = 5$
- $\kappa = 5$

Initial setup

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3
4 n = 5
5 K = 5
6
7 def f(x):
8     return 0.5 * np.sum(x[:-1] ** 2) + 0.5 * K * x[-1] ** 2
```

2 Determining the Optimal Solution

The optimal solution for $f(x)$ is obtained when $\nabla_x f(x) = 0$.

$$\nabla_x f(x) = [x_1 \ x_2 \ x_3 \ \cdots \ \kappa x_n]$$

This is satisfied for $x_i = 0$ for $i \in \{1 \dots n\}$. Hence the optimal solution is $\begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$.

3 Gradient and Hessian of $f(x)$.

Mathematically analyzing $f(x)$, we obtain:

$$\nabla_x f(x) = [x_1 \ x_2 \ x_3 \ \cdots \ \kappa x_n]$$

$$H = \begin{pmatrix} 1 & 0 & 0 & \cdots & 0 \\ 0 & 1 & 0 & \cdots & 0 \\ 0 & 0 & 1 & \cdots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \cdots & \kappa \end{pmatrix}$$

Gradient & Hessian Functions

```
1 x = np.random.uniform(0, 10, n)
2
3 def grad_f(x):
4     g = np.copy(x)
5     g[-1] = K * x[-1]
6     return g
7
8 print(x)
9 print(grad_f(x))
10
11 def hessian():
12     hessian = np.eye(n)
13     hessian[n - 1, n - 1] = K
14     return hessian
15
16 print(hessian())
```

Result

```
1 [2.82754714 0.74166732 9.03752034 3.66116654 8.57032453]
2 [ 2.82754714  0.74166732  9.03752034  3.66116654 42.85162266]
3 [[1. 0. 0. 0. 0.]
4  [0. 1. 0. 0. 0.]
5  [0. 0. 1. 0. 0.]
6  [0. 0. 0. 1. 0.]
7  [0. 0. 0. 0. 5.]]
```

4 Convexity, Smoothness and Strong Convexity of $f(x)$

The eigenvalues of the Hessian of $f(x)$ are 1 and κ . Since both are non-negative, the Hessian is positive semi-definite and hence the function is convex for all x in the domain.

Also the eigenvalues are bounded by $[\alpha, B]$ where $\alpha = 1$ and $B = \kappa$, $f(x)$ is guaranteed to be smooth and strongly-convex.

Since the eigenvalues are finite, we can say the eigenvalues are bounded between $[\alpha, \beta]$ where $\alpha = 1$ and $\beta = \kappa$. This concludes the fact that $f(x)$ is smooth and strongly convex.

$$\text{Condition Number} = \frac{B}{\alpha} = \frac{\kappa}{1} = \kappa$$

5 Generating the initial candidate solution

For $n = 100$ and $\kappa = 5$, we generate an initial candidate solution x_0 by sampling a Gaussian random vector, with zero mean and identity matrix as the covariance matrix, and multiplying 100 to the sampled vector.

Generation of x_0

```
1 n = 100
2 K = 5
3
4 x0 = 100*np.random.multivariate_normal(np.zeros(n), np.eye(n))
5 print(x0[:5])
```

Result

```
1 [-151.47236778 -154.43656717 175.74317823 -164.83626122 103.17357891]
```

6 Implement gradient descent, accelerated gradient descent and momentum method

For all the learning methods, we keep learning rate to $\frac{1}{\beta}$. In this case, this will equal to $\frac{1}{\kappa} = 0.2$

Implementation of various learning algorithms

```
1  T = 200
2  mus = [0.4, 0.6, 0.8]
3  L = float(max(1, K))
4  eta = 1.0 / L
5
6  def run_gd():
7      x = x0.copy()
8      history = []
9      for t in range(T):
10         history.append(f(x))
11         x = x - eta * grad_f(x)
12     return history
13
14 def run_agd():
15     x = x0.copy()
16     y = x0.copy()
17     history = []
18     m = min(1, K)
19     beta = (np.sqrt(L) - np.sqrt(m)) / (np.sqrt(L) + np.sqrt(m))
20
21     for t in range(T):
22         history.append(f(x))
23         x_next = y - eta * grad_f(y)
24         y = x_next + beta * (x_next - x)
25         x = x_next
26     return history
27
28 def run_momentum_method():
29     history = {}
30     for mu in mus:
31         x_m = x0.copy()
32         v = np.zeros(n)
33         history_m = []
34         for t in range(T):
35             history_m.append(f(x_m))
36             v = mu * v - eta * grad_f(x_m)
37             x_m = x_m + v
38         history[mu] = history_m
39     return history
40
41 history_gd = run_gd()
```

```

42 history_agd = run_agd()
43 history_momentum = run_momentum_method()

```

7 Plotting

We plot the values of objective function obtained while running the above optimization algorithms in both linear and semilog scales. Along with that, we also plot the theoretical upper bounds for for gradient descent and accelerated gradient descent algorithms.

$$\text{Upper bound in gradient descent} = \left(1 + \frac{\alpha}{\beta - \alpha}\right)^{-t} f(x_0)$$

$$\text{Upper bound in accelerated gradient descent} = \left(1 + \frac{1}{\sqrt{\frac{\beta}{\alpha}} - 1}\right)^{-t} \frac{\alpha + \beta}{2} \|x_0 - x^*\|^2$$

```

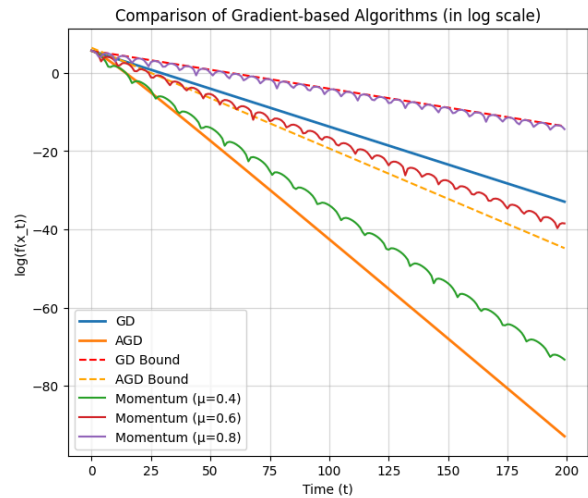
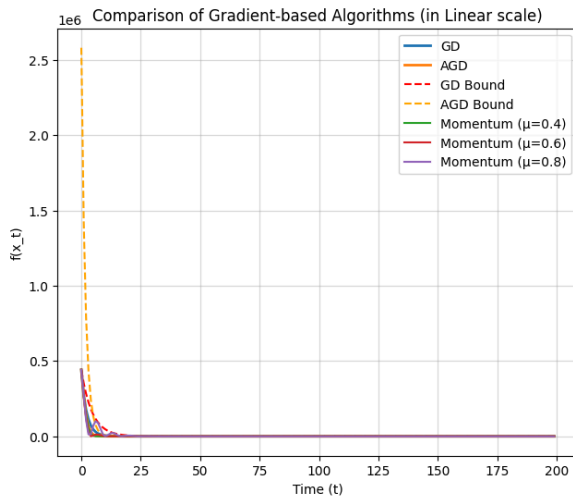
Plotting using matplotlib
1  def plot_histories(history_gd, history_agd, history_momentum):
2      t_range = np.arange(T)
3
4      beta = float(max(1, K))
5      alpha = float(min(1, K))
6      cond_num = beta / alpha
7
8      # Bound for Gradient Descent
9      gamma_gd = alpha / (beta - alpha)
10     bound_gd = (1 + gamma_gd)**(-t_range) * f(x0)
11
12     # Bound for AGD
13     x_opt = np.zeros(n)
14     gamma_agd = 1.0 / (np.sqrt(cond_num) - 1)
15     prefactor = (alpha + beta) / 2.0 * np.sum((x0 - x_opt)**2)
16     bound_agd = (1 + gamma_agd)**(-t_range) * prefactor
17
18     def power_formatter(value, _):
19         return f"10^{value}"
20
21     plt.style.use('dark_background')
22     fig, (ax1, ax2) = plt.subplots(1, 2, figsize=(16, 6))
23
24     for ax, is_log in zip([ax1, ax2], [False, True]):
25         if is_log:
26             history_gd = np.log10(np.array(history_gd))
27             history_agd = np.log10(np.array(history_agd))
28             bound_gd = np.log10(np.array(bound_gd))
29             bound_agd = np.log10(np.array(bound_agd))
30
31             history_momentum2 = {}
32             for mu, history_m in history_momentum.items():
33                 history_momentum2[mu] = np.log10(np.array(history_m))
34             history_momentum = history_momentum2

```

```

35         ax.set_ylabel('log(f(x_t))')
36         ax.set_title('Comparison of Gradient-based Algorithms (in log scale)')
37     else:
38         ax.set_ylabel('f(x_t)')
39         ax.set_title('Comparison of Gradient-based Algorithms (in Linear scale)')
40
41     ax.plot(t_range, history_gd, label='GD', linewidth=2)
42     ax.plot(t_range, history_agd, label='AGD', linewidth=2)
43     ax.plot(t_range, bound_gd, label='GD Bound', color='red', linestyle='--')
44     ax.plot(t_range, bound_agd, label='AGD Bound', color='orange', linestyle='--')
45
46     # Plot Momentum results for each mu
47     for mu, history_m in history_momentum.items():
48         ax.plot(t_range, history_m, label=f'Momentum ( $\mu$ = {mu})')
49
50     ax.set_xlabel('Time (t)')
51     ax.legend()
52     ax.grid(True, which="both", ls="-", alpha=0.5)
53
54     plt.show()
55
56 plot_histories(history_gd, history_agd, history_momentum)

```



8 Results

From the above observations, we conclude the following facts:

1. All the algorithms are able to perfectly converges to the optimal solution x^* within the defined iteration limit.
2. The accelerated gradient descent algorithm converges fastest to the optimal solution followed by momentum method with low values of μ followed by gradient descent and closing with the momentum method with higher values of μ .

Accelerated Gradient Descent > Momentum ($\mu = 0.4$ > Momentum ($\mu = 0.6$ > Gradient Descent > Momentum ($\mu = 0.8$)

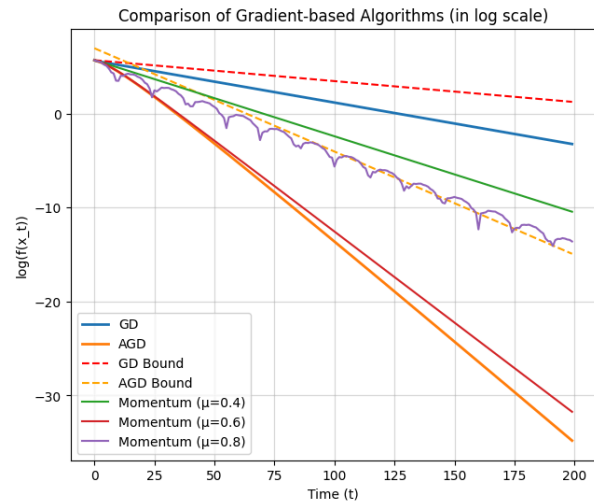
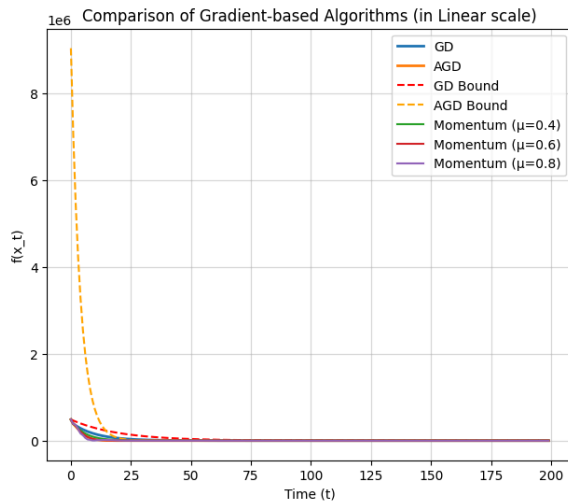
- For $\kappa = 5$, both gradient descent and accelerated gradient descent stay well within the theoretical upper bound

9 For $\kappa = 20$

```

1 K = 20
2 L = float(max(1, K))
3 eta = 1.0 / L
4
5 history_gd = run_gd()
6 history_agd = run_agd()
7 history_momentum = run_momentum_method()
8
9 plot_histories(history_gd, history_agd, history_momentum)

```



9.1 Results

From the above plots for $\kappa = 20$, we come to the following conclusions:

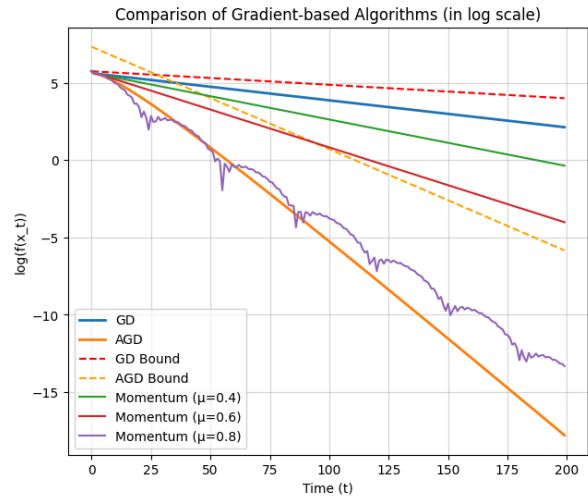
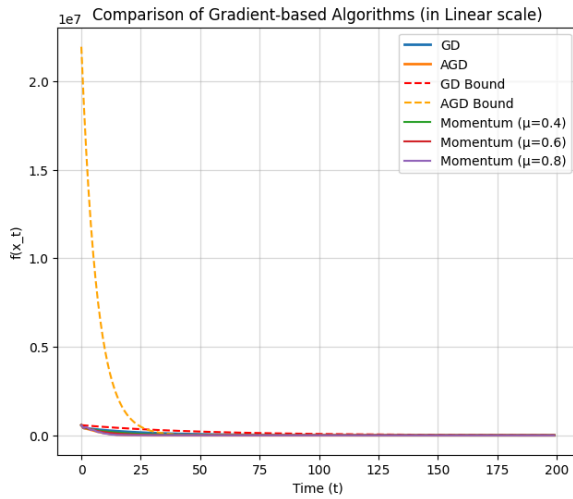
- All the the algorithms are still able to perfectly converges to the optimal solution x^* within the defined iteration limit.
- However the order of convergence of the algorithms change in this case

Accelerated Gradient Descent $>$ Momentum ($\mu = 0.6 >$ Momentum ($\mu = 0.8 >$ Momentum ($\mu = 0.4 >$ Gradient Descent

- Both gradient descent and accelerated gradient descent still stay well within the theoretical upper bound.

10 For $\kappa = 50$

```
1 K = 50
2 L = float(max(1, K))
3 eta = 1.0 / L
4
5 history_gd = run_gd()
6 history_agd = run_agd()
7 history_momentum = run_momentum_method()
8
9 plot_histories(history_gd, history_agd, history_momentum)
```



10.1 Results

From the above plots for $\kappa = 50$, we come to the following conclusions:

1. All the the algorithms are still able to perfectly converges to the optimal solution x^* within the defined iteration limit.
2. The order of convergence of the algorithms again change in this case

Accelerated Gradient Descent $>$ Momentum ($\mu = 0.8 >$ Momentum ($\mu = 0.6 >$ Momentum ($\mu = 0.4 >$ Gradient Descent

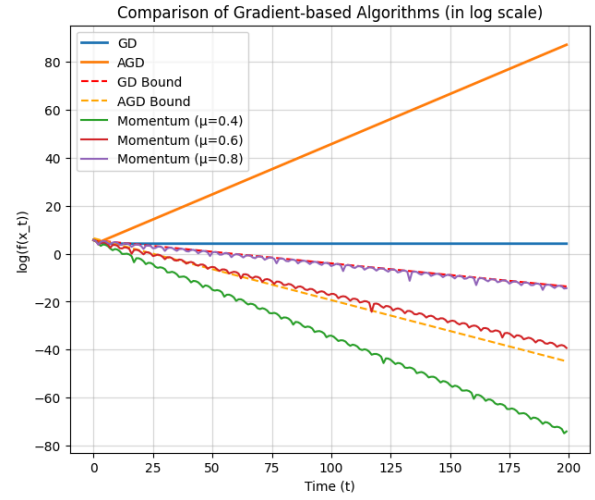
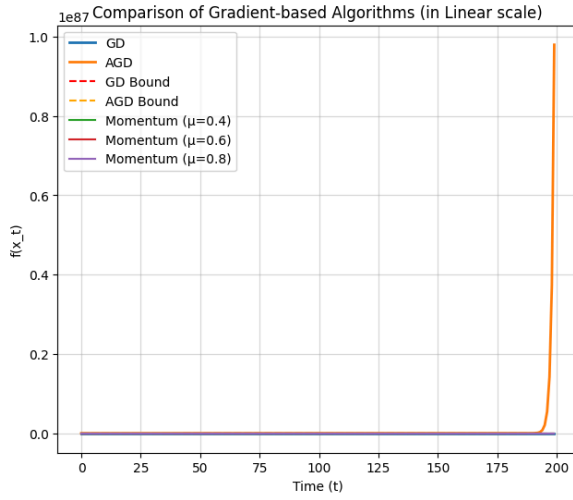
3. Both gradient descent and accelerated gradient descent still stay well within the theoretical upper bound.

11 For $\eta = \frac{2}{\beta}$

```

1  K = 5
2  L = float(max(1, K))
3  eta = 2 / L
4
5
6  history_gd = run_gd()
7  history_agd = run_agd()
8  history_momentum = run_momentum_method()
9
10 plot_histories(history_gd, history_agd, history_momentum)

```



11.1 Results

From the above plots for $\eta = \frac{2}{\beta}$, we come to the following conclusions:

1. All the the algorithms except **accelerated gradient descent** are able to perfectly converges to the optimal solution x^* within the defined iteration limit.
2. However, both gradient descent and accelerated gradient descent exceed their theoretical upper bounds for all x in the domain.
3. The accelerated gradient descent overshoot to ∞ towards the very end of the iterations

12 Conclusion

1. The given function $f(x)$ has a optimal solution at $x = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$.
2. The function is convex, smooth and strongly-convex with smoothness and strongly-convex parameters $\beta = \kappa$ and $\alpha = 1$.
3. For various values of κ , the function can be converged to its optimal solution by various algorithms like gradient descent, accelerated gradient descent and momentum method.
4. Gradient Descent (GD) shows smooth and steady decrease in $f(x)$, which matches its expected exponential convergence. Accelerated Gradient Descent (AGD) converges significantly faster than GD and reaches low function values in fewer iterations. The Momentum method converges faster than GD for moderate μ values. Among the tested values, $\mu = 0.6$ gives the best balance between speed and stability.

5. As k increases, the condition number $\kappa = \frac{\beta}{\alpha}$ increases, making the optimization problem more ill-conditioned. GD converges more slowly, which is consistent with its dependence on the condition number. AGD still remains the fastest method even for larger κ , although its convergence also slows down compared to smaller k values. The Momentum method shows stronger oscillations for larger μ when k is large, indicating sensitivity to parameter choice.
6. When the step size is doubled, GD becomes unstable and exhibits oscillatory or divergent behavior. AGD also fails to converge with the doubled step size, despite its faster convergence under the correct step size. Momentum diverges the fastest because the momentum term amplifies oscillations when the step size is too large. This experiment confirms that $\eta = 1/\beta$ is the maximum stable step size for these methods, as predicted by the theory.