

②

①

Gradient descent is a foundational optimization technique used in machine learning to minimize a cost function iteratively.

The three gradient descent which is explored in the assignment are as:

Vanilla Gradient Descent; It is the simplest form of gradient descent

i) It computes the gradient of ~~each~~ cost function w.r.t ^{all} ~~each~~ data points in each iteration

ii) The update rule for parameter (weight) θ is given by,

$$\theta_{\text{new}} = \theta_{\text{old}} - \text{learning rate} \times \nabla J(\theta_{\text{old}})$$

where $\nabla J(\theta_{\text{old}})$ is the gradient of cost function with respect to θ_{old} .

Mini-Batch Gradient descent.

i) This is a compromise between Vanilla and Stochastic Gradient descent.

ii) It divides dataset into small batch (mini-batches) & computes gradient w.r.t one mini batch at a time.

iii) Update rule is similar to

(2)

Vanilla Gradient descent, but it's sum for the gradients of the mini-batch.

$$\theta_{\text{new}} = \theta_{\text{old}} - \text{learning rate} \times (1/\text{batch size}) \times \sum (\nabla J(\theta_{\text{old}}) \text{ for data in mini-batch})$$

Stochastic Gradient Descent (SGD) :

- In stochastic gradient descent the gradient is computed and parameters are updated for each individual data point in the dataset.
- The update rule is similar to Vanilla Gradient descent but it uses a single data point at a time :

$$\theta_{\text{new}} = \theta_{\text{old}} - \text{learning rate} \times \nabla J(\theta_{\text{old}}) \text{ for a randomly selected data point.}$$

Observations from Graphs:

- (1) a) Vanilla Gradient descent converges more smoothly.
- b) Vanilla gradient descent takes computationally more intensive as it uses whole dataset, so we can see cost is very high at the start.



② Stochastic Gradient descent - exhibits faster convergence.

- a) SGD's path looks more erratic as it introduces significant noise into parameter updates.
- b) It's less computationally intensive & ~~often~~ as it processes one data point at a time.

③ Mini-Batch Gradient descent looks like a mid-path as far as convergence is concerned.

- a) In case of mini batch path is less erratic as it introduces controlled noise.
- b) It's versatile choice as it balances the trade off between ~~faster~~ convergence ~~and~~ smoothness and ~~faster~~ computational efficiency.

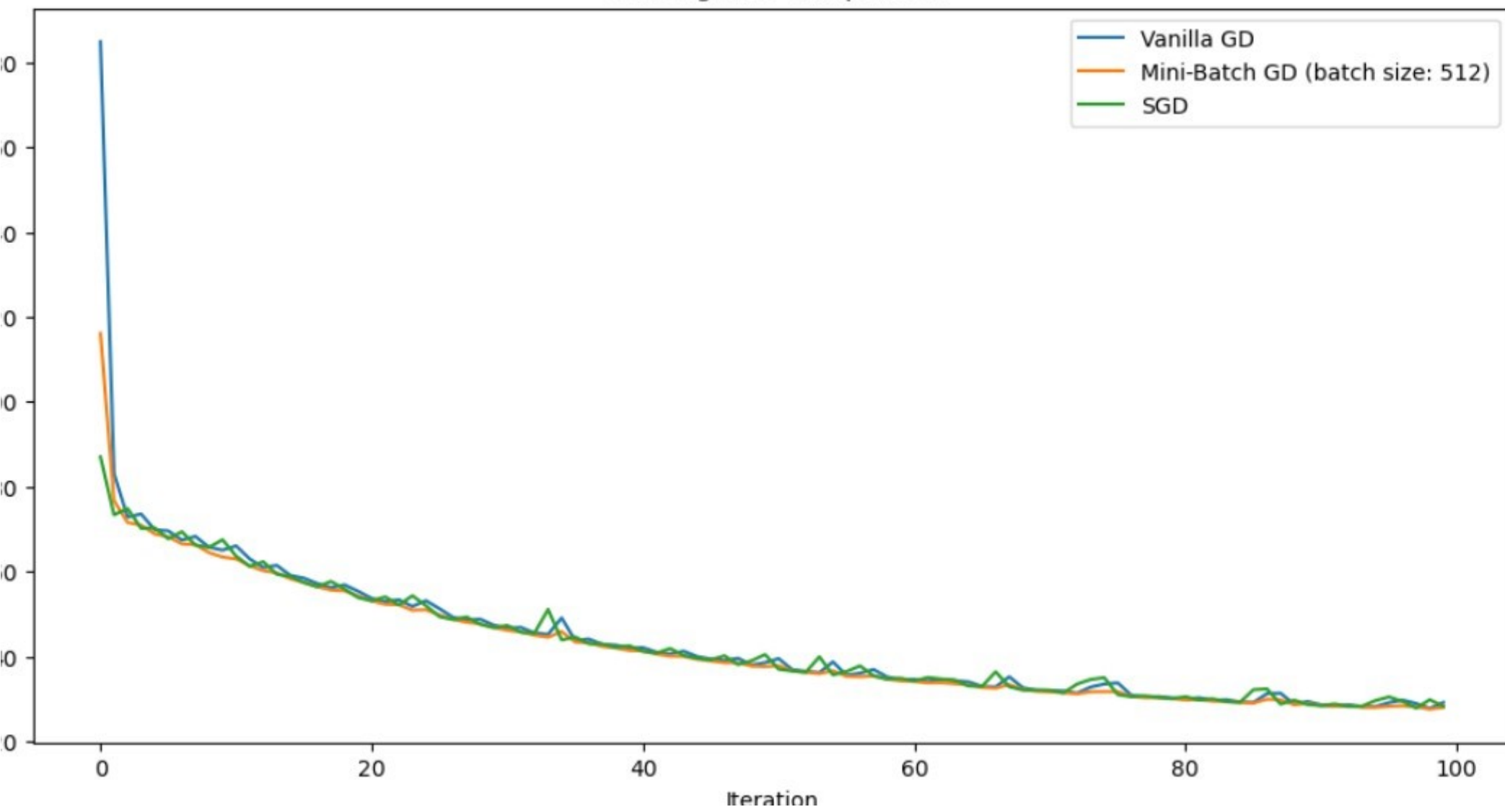
Overall observation:

SGD is useful when computational resources are limited and faster convergence is needed ~~with~~ even with noise.

Vanilla works best when computational resources are abundant & smooth convergence path is needed.

Mini-Batch Gradient descent is a versatile choice that balances the trade off b/w SGD & Vanilla.

Convergence Comparison



Experiment with mini-batch Gradient descent

- Load the dataset with at least 1500 rows and 15 columns. (Transportation dataset in exp)
- Carefully pre-process by removing nulls, missing values & added scaling feature.
- Split data into training set & validation set to assess performance.

Initialization: Initialize model parameters (0) and set hyperparameters including learning rate, number of iterations and convergence criteria. In our case, we have kept:
 no. of iterations = 1000, learning rate = 0.01.

Batch size selection: Define batch size to experiment with. In our case, we have used batch size of 16, 32, 64, 128 and 256.

- The optimal size depends on a number of factors such as dataset size, hardware resources and convergence speed etc. Smaller batch size is generally used when fast convergence is a priority while large batch size is needed when

(6)
computational efficiency is critical or when dealing with noisy data.

Training loops:

- Implement the loop for each variant.
- for each batch size; perform following steps
 - a) Divide training data into mini-batches of selected data batch size.
 - b) Iterate thru mini-batches, update parameters using appropriate gradient descent algo.
 - c) Record & track (plot on graph) the cost (or loss) on the validation set.

Analysis:

→ A graph is plotted by plotting the cost (or loss) on validation set for each batch size over iterations. We see best batch size is 64 for given experiment with minimum loss.

Conclusion:

Batch size there is always a trade off between faster convergence & loss.

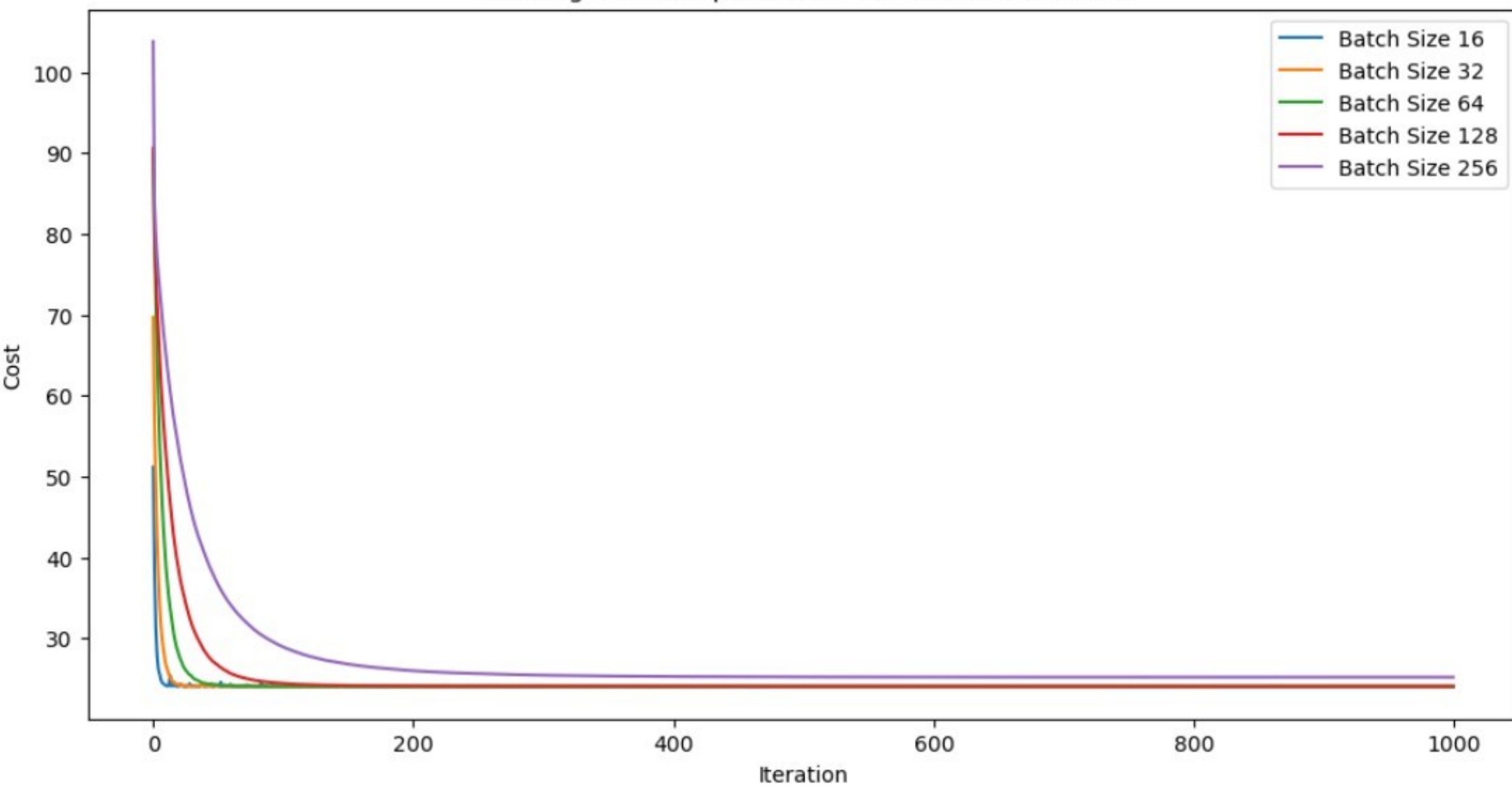
Large batch size → Smoother & Stable convergence
Smaller batch size → Faster convergence but can lead to unstable convergence

→ A small batch size can introduce more noise but may lead to faster convergence while a large batch size can reduce noise but might converge more slowly.

→ Mini-batch is suitable for a wide range of dataset size, can leverage parallelism for faster training on multi-core processors or GPUs.

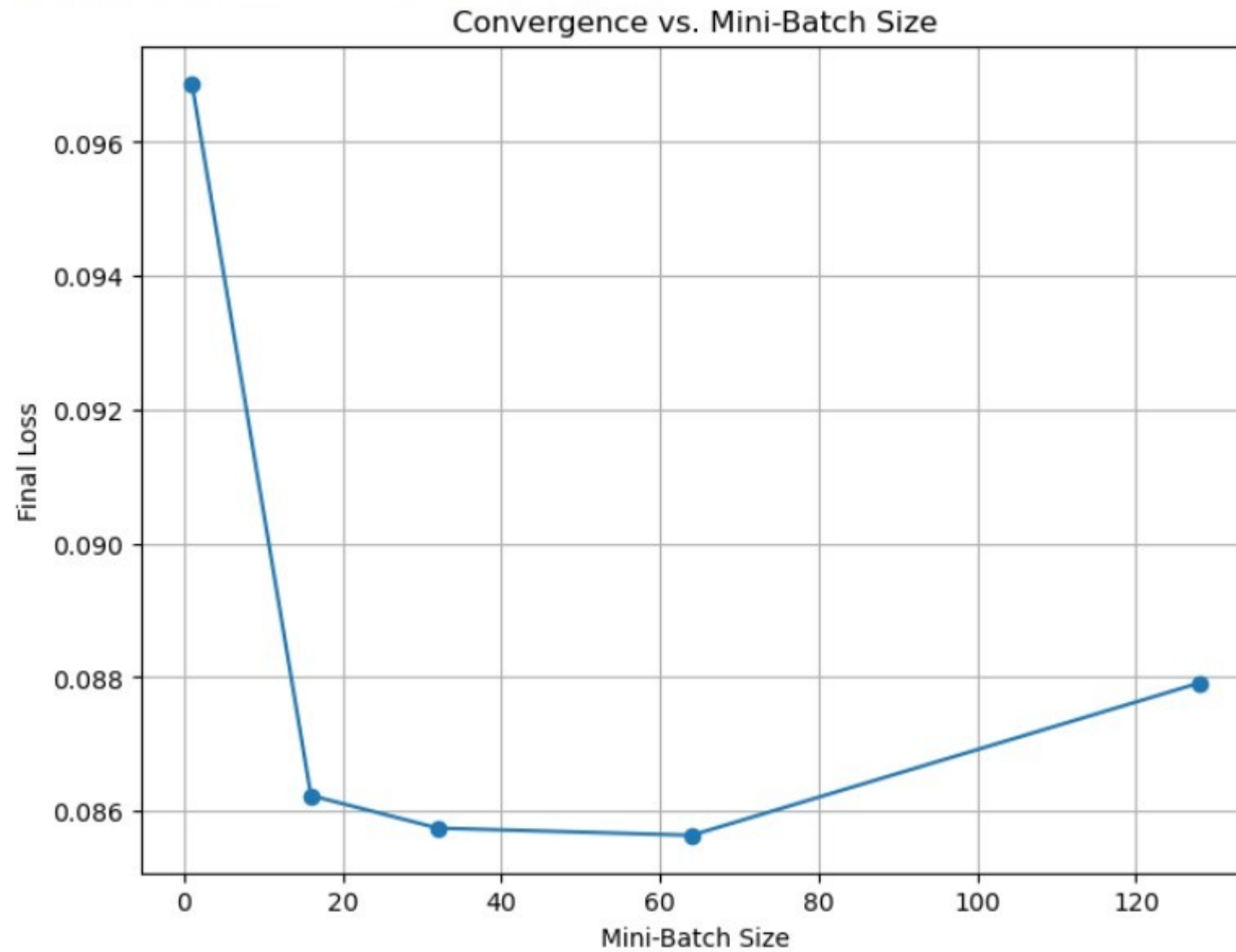
(All the graphs are in next pages)

Convergence Comparison for Different Batch Sizes



Best Batch Size: 64

Final Loss with Best Batch Size: 0.08563883618900707



② In optimization line search strategy is one of the two basic iterative approaches to find local minimum of a function. The method starts at a point x_0 and then move in a direction d by a step α . The step length α is chosen so that to minimize $f(x_0 + \alpha d)$.

Bisection Search: This is a simple and robust method that works out by repeatedly bisecting intervals of possible values until a step size is found ~~given~~ the condition which satisfies the condition. Here - $C = (a+b)/2$.

Golden Search: This is more efficient variant of bisection search that uses golden ratios instead of mid-point bisection to find the suitable step size.

Here mid-point, is found such that.

$$\frac{b}{a} = \psi$$

$$\text{where } \psi = \frac{1+\sqrt{5}}{2} \approx 1.6180$$

You start with interval $[a, b]$ and calculate two points

$$x_1 = a + (1-\psi)(b-a)$$

$$x_2 = a + \psi(b-a)$$

Evaluate function at x_1 & x_2 & compare function values.

Depending on function value, update the interval & continue until its sufficiently small.

Armijo rule: This is more sophisticated method that uses the gradient of the objective function to ~~find the~~ guide the search for a suitable α . Its based on the idea of making sure that the function decreases sufficiently during each iteration.

→ Rule involves a parameter typically a small number between (0 & 1) and a decrease factor.

→ At each iteration, check if new function value is less than previous one or equal to current value. minimize a certain prop. of gradient times the step.

Step

⇒ If condition is met, step size is accepted. Otherwise it's reduced and the check is repeated until the condition is satisfied.

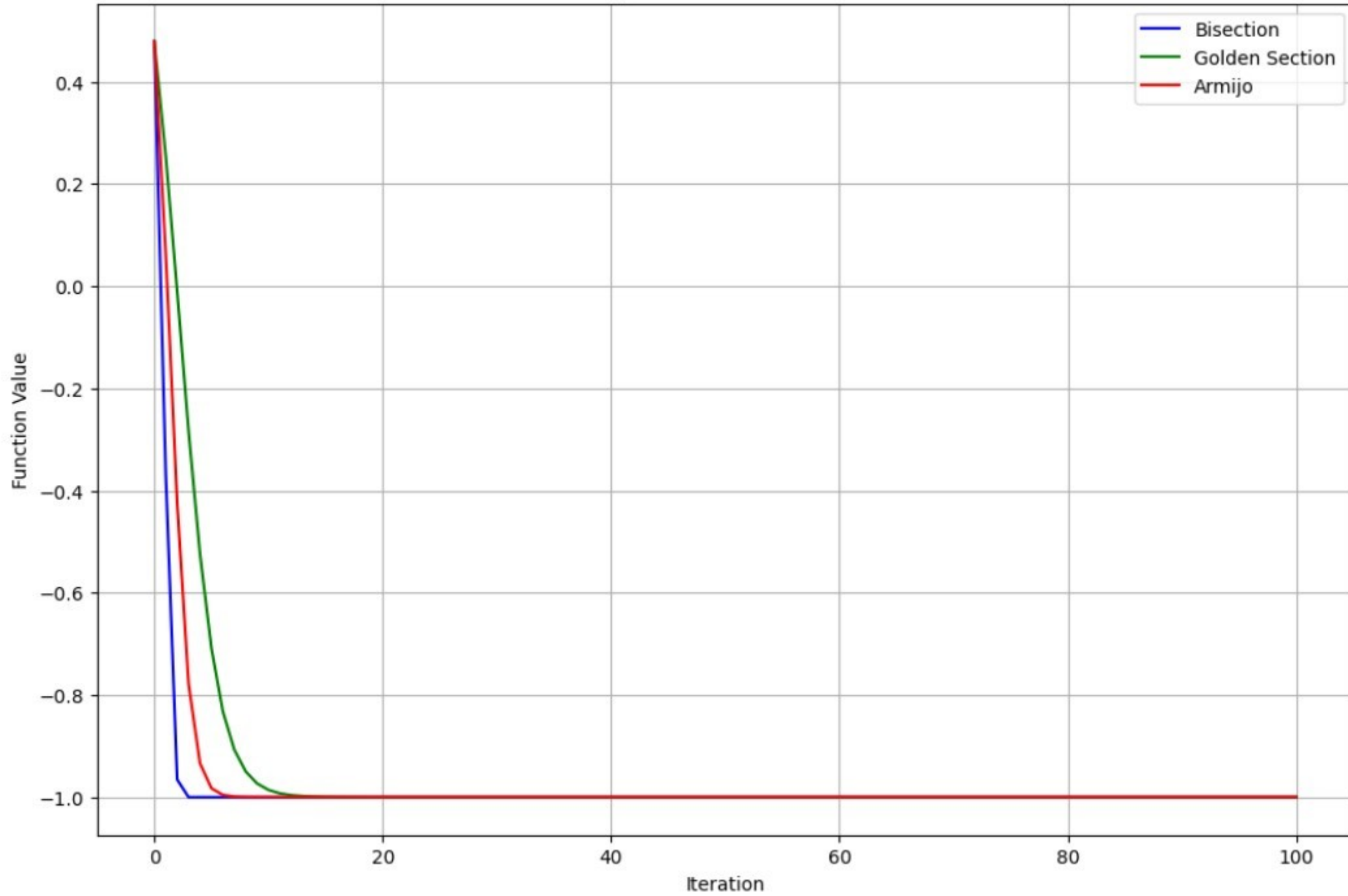
These line searches are essential tools in optimization algorithms for finding optimal solⁿ. Choice of ~~at~~ search depends on problem statement.

Bisection & Golden Search are suitable for one dimensional & Armijo is more suitable for multi-dimensional optimized solⁿ.

In this case, Bisection works the best, then Armijo & last is Golden Search.

Method	Adv.	Disadv.
Bisection Search.	Simple, Robust	Slow for large prob.
Golden.	More eff. than Bisection (generally)	Less robust than Bisection
Armijo	Most sophisticated & works best for complicated problem.	May not converge & expensive computationally

Convergence of Gradient Descent



(3)

(1)

Function $f(x) = (x^2 \cos(x) + \sin(x) - x)$

Gradient descent is a first-order iterative optimization algorithm for finding a local minimum of a differentiable function.

Its the learning rate which determines how much we adjust our current position which greatly affects the convergence of algorithm.

$$f(x) = x^2 \cos(x) + \sin(x) - x$$

$$f'(x) = 2x \cos(x) - x^2 \sin(x) + \cos(x) - 1$$

Constant Step Size: In this case, we simply use a fixed step throughout the optimization process. This is simple but it leads to slow convergence or overshoot the minima.

Mathematically $x_{old} =$

$$x_{new} = x_{old} - \gamma f'(x_{old})$$

Choice of γ is crucial here. If its too large, it will overshoot the minimum and diverge. If its too large, small, it will converge too slowly.

Decaying step size: In this case, we start with an initial step size α_0 and decrease it over time using a schedule like $\alpha_k = \alpha_0/k$ where k is positive constant.

A common strategy of decaying is

$$\alpha_t = \frac{\alpha_0}{1+kt}$$

where $\alpha_t \rightarrow$ step size at iteration t

$\alpha_0 \rightarrow$ initial step size

$k \rightarrow$ decay constant

Choice of α_0 and k will determine how quickly step size decays. A smaller k will result in smaller decay while large k will result in faster decay.

Bold driver Algorithm: The bold driver

algorithm dynamically adjusts the size of step if progress is good and decreases if it's not. Exact rule of step size can vary. One of the way is to multiply step size by a factor greater than 1 after a successful iteration & a factor less than 1 after unsuccessful one.

method

Conclusion: The best choice of step size will vary depend on specific characteristics of function and initial condition. It's often good to use diff. methods and step size to see what works best.

Overall:

Constant Step Size: Its simplest approach but might not be most efficient. Here, convergence depends on gamma (γ) & it doesn't adapt to function.

Decaying step size: It can be more efficient as it allows for larger step initially and then smaller steps later on. However choice of α & k is crucial.

Bold driver: Its adaptive approach & here step size varies based on function's behaviour. It's far more sensitive to initial step size & choice of increase/decrease factor.

Clearly in this example, bold driver's algorithm worked best.

Comparison of Step Size Strategies

