



---

## Exploring Gradient Descent: A Mathematical, Visual, and Coding Perspective

*Yang Liu*

Tech. Univ. China

---

### ABSTRACT.

Gradient descent constitutes a pillar of modern machine learning, data science, and beyond. This paper spotlights gradient descent techniques through an accessible lens, synthesizing mathematical formalism, visual metaphors, and an implementational case study. We provide a first-principles derivation of gradient descent, grounding it intuitively as navigating foggy mountain landscapes. We also demonstrate gradient descent by fitting a linear regressor to synthetic data in Python. The fusion of equations, visuals, and code aims to render gradient descent profoundly yet easily comprehensible.

---

### 1. Introduction

Gradient descent is a powerful optimization technique that can be used to find the optimal values of parameters that minimize a given objective function (see, for example, [4]). Gradient descent is widely used in machine learning, data science, and various other fields, as it can handle complex and high-dimensional problems that are difficult or impossible to solve analytically (see, for example, [8]). In this introduction, we will cover the mathematical foundations, intuitive interpretations, practical implementations, and common challenges of gradient descent (see, for example, [13]).

Mathematically, gradient descent is based on the idea that the gradient of a function at any point gives the direction of the steepest ascent at that point. Therefore, by moving in the opposite direction of the gradient, we can descend towards the minimum of the function. The gradient descent algorithm can be summarized as follows:

Initialize the parameters randomly or with some guess.

Compute the gradient of the objective function with respect to the parameters.

Update the parameters by subtracting a small fraction of the gradient from them.

Repeat steps 2 and 3 until convergence or a maximum number of iterations.

The fraction of the gradient that is subtracted from the parameters is called the learning rate, and it controls how big or small the steps are in each iteration (see, for example, [17]). A small learning rate can lead to slow convergence, while a large learning rate can cause overshooting or divergence (see, for example, [21]).

Intuitively, gradient descent can be understood as navigating a foggy mountain landscape to find water in a valley below (see, for example, [10]). At each step, we measure the steepest slope and walk in that direction (see, for example, [16]).

Incrementally, we traverse the descents until reaching the bottom – the oasis (see, for example, [2]). The learning rate determines how far we walk in each step, and the initial guess determines where we start on the mountain (see, for example, [6]).

Practically, gradient descent can be implemented in various programming languages and frameworks, such as Python, MATLAB, TensorFlow, PyTorch, etc (see, for example, [20]). The choice of language and framework depends on the problem domain, the availability of libraries and tools, and the personal preference of the user (see, for example, [12]). The implementation of gradient descent involves defining the objective function, initializing the parameters, computing the gradient, updating the parameters, and checking for convergence (see, for example, [15]).

Challenges: Gradient descent is not without its challenges and limitations (see, for example, [11]). Some of the common challenges are:

Choosing an appropriate learning rate: If the learning rate is too small, gradient descent can take a long time to converge or get stuck in a local minimum (see, for example, [1]). If the learning rate is too large, gradient descent can overshoot or diverge from the global minimum.

Dealing with noisy or sparse gradients: If the objective function is noisy or sparse, gradient descent can be affected by random fluctuations or zero gradients that prevent it from finding the optimal solution.

Handling non-convex or multimodal functions: If the objective function is non-convex or multimodal, gradient descent can get trapped in a local minimum or saddle point that is not optimal.

Scaling to large-scale problems: If the problem size is very large or high-dimensional, gradient descent can be computationally expensive or infeasible to run.

To overcome these challenges, various modifications and extensions of gradient descent have been proposed and developed over time (see, for example, [19] and [9]). Some of these include stochastic gradient descent (SGD), mini-batch gradient descent (MBGD), momentum-based methods (e.g., Nesterov accelerated gradient), adaptive methods (e.g., AdaGrad, RMSProp, Adam), second-order methods (e.g.,

Newton's method), conjugate gradient methods (CGM), etc (see, for example, [5] and [22]). These methods aim to improve the speed, accuracy, robustness, and scalability of gradient descent (see, for example, [3] and [1]).

Gradient descent is a powerful optimization technique that can be used to find the optimal values of parameters that minimize a given objective function (see, for example, [7]). Gradient descent is widely used in machine learning, data science, and various other fields, as it can handle complex and high-dimensional problems that are difficult or impossible to solve analytically (see, for example, [8]). Gradient descent is based on the idea that by moving in the opposite direction of the gradient of a function at any point, we can descend towards the minimum of the function (see, for example, [15]). Gradient descent can be understood intuitively as navigating a foggy mountain landscape to find water in a valley below (see, for example, [24]). Gradient descent can be implemented in various programming languages and frameworks, depending on the problem domain and user preference (see, for example, [18]). Gradient descent faces some challenges and limitations when dealing with noisy or sparse gradients, non-convex or multimodal functions, or large-scale problems (see, for example, [23]). To overcome these challenges, various modifications and extensions of gradient descent have been proposed and developed over time (see, for example, [14]).

Gradient descent is a powerful and widely used optimization technique in machine learning, data science, and beyond. It is based on the simple idea of following the steepest direction of descent to find the minimum of a function. Gradient descent can be applied to various problems, such as linear regression, logistic regression, neural networks, support vector machines, and many more. However, despite its popularity and importance, gradient descent is often not well understood or explained by many practitioners and learners. This paper aims to fill this gap by providing an accessible and comprehensive introduction to gradient descent, covering its mathematical foundations, intuitive interpretations, practical implementations, and common challenges.

The main contributions of this paper are as follows:

- We derive the gradient descent algorithm from first principles, using calculus and linear algebra. We explain the key concepts and steps involved in finding the optimal solution of a function using gradient descent.
- We illustrate the gradient descent algorithm with visual metaphors, such as navigating a foggy mountain landscape or rolling a ball down a hill. We show how these metaphors can help us understand the intuition and logic behind gradient descent.
- We demonstrate the gradient descent algorithm with a concrete example of fitting a linear regression model to synthetic data. We use Python code to implement the algorithm and plot the results. We also discuss how to choose the learning rate hyperparameter and how to check for convergence.
- We synthesize equations, visuals, and code to provide a multifaceted and holistic understanding of gradient descent. We hope that this paper will serve as a useful resource for anyone who wants to learn or teach gradient descent in an engaging and effective way.

The rest of this paper is organized as follows. Section 2 provides some background information on optimization problems and objective functions. Section 3 presents the derivation and interpretation of gradient descent. Section 4 shows the application of gradient descent to linear regression. Section 5 concludes with some discussion and future directions.

---

## 2. Mathematical Foundations and Optimization

In this section, we establish the mathematical foundations that are indispensable for comprehending the intricacies of the Gradient Descent algorithm, a fundamental tool in various scientific disciplines, including machine learning, optimization, and mathematical modeling.

## 2.1 Multivariable Calculus.

### 2.1.1 Partial Derivatives: Understanding Rates of Change.

Partial derivatives constitute the cornerstone of multivariable calculus, offering us a nuanced perspective on how a function evolves concerning specific variables while all other variables remain held constant. Consider a real-valued function  $f(x, y)$  given by:

$$f(x, y) = x^2 + xy$$

In this context, the partial derivative  $\frac{\partial f}{\partial x}$  represents the rate of change of  $f$  with respect to  $x$ , while  $\frac{\partial f}{\partial y}$  represents the rate of change with respect to  $y$ , all while the other variable remains fixed. Let's break this down further:

$\frac{\partial f}{\partial x}$  measures how  $f$  changes as we make infinitesimal adjustments to  $x$  while keeping  $y$  constant. In our example,  $\frac{\partial f}{\partial x} = 2x + y$ .

- On the other hand,  $\frac{\partial f}{\partial y}$  tells us how  $f$  varies as we change  $y$  while holding  $x$  constant. In our example,  $\frac{\partial f}{\partial y} = x$ .

These partial derivatives provide critical insights into the local behavior of functions and play a pivotal role in understanding how we can optimize them.

### 2.1.2 Gradients: The Path to Optimization.

To transcend the realm of single-variable derivatives, we introduce the concept of gradients. For a function  $f: \mathbb{R}^n \rightarrow \mathbb{R}$ , the gradient  $\nabla f$  is a vector defined as:

$$\nabla f = \left[ \frac{\partial f}{\partial x_1}, \frac{\partial f}{\partial x_2}, \dots, \frac{\partial f}{\partial x_n} \right]$$

Geometrically, the gradient vector  $\nabla f$  offers a profound insight: it points in the direction of the steepest ascent of the function  $f$ . This directional information is invaluable in optimization tasks because it guides us toward regions where the function's value increases most rapidly.

The gradient also provides a crucial connection between multivariable calculus and optimization. By leveraging gradient information, we can systematically navigate the landscape of a function to locate its minimum or maximum points.

## 2.2 Optimization: Seeking the Optimum.

### 2.2.1 The Optimization Goal.

At the heart of many scientific and engineering endeavors lies the fundamental problem of optimization. We endeavor to solve problems of the form:

$$\min_{x \in \mathbb{R}^n} f(x)$$

In this expression,  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  represents a continuous and differentiable function, and our objective is to find the values of the variables ( $x$ ) that minimize  $f$ . This pursuit of finding optimal solutions permeates fields as diverse as machine learning, engineering design, and economics.

### 2.2.2 First-order Methods: Unveiling Gradient Descent.

Optimization algorithms that rely on gradient information are categorized as first-order methods. One of the most widely employed and studied first-order methods is the Gradient Descent algorithm. It forms the crux of our exploration in this paper.

First-order methods, like Gradient Descent, operate iteratively. They start with an initial guess for the solution ( $x_0$ ) and update it in each iteration based on information extracted from the gradient. The intuition is simple: move in the direction of the steepest descent (negative gradient) to approach the minimum of the function.

With these foundational concepts firmly in place, we now embark on the journey of demystifying the workings of the Gradient Descent algorithm, a linchpin in the world of optimization and machine learning.

### 3. Gradient Descent

We now derive and visualize gradient descent from first principles.

#### 3.1 Mathematical Derivation.

Gradient descent iteratively minimizes objective functions using gradient information. Consider the update equation:

$$x^{(t+1)} = x^{(t)} - \alpha \text{Of}(x^{(t)}) \quad (3.1)$$

where  $x^{(t)}$  are the parameters at iteration  $t$ ,  $\alpha$  is the learning rate hyperparameter, and  $\text{Of}(x^{(t)})$  is the gradient. This elegantly encodes the essence of gradient descent:

- (1) Evaluate the gradient  $\text{Of}$  at the current parameters  $x^{(t)}$ .
- (2) Make a small step against the gradient direction to obtain  $x^{(t+1)}$ .
- (3) Repeat until convergence.

This incremental gradient-based approach minimizes  $f(x)$  by repeatedly traversing in the steepest descent direction. Under appropriate conditions, gradient descent provably converges to local minima.

The learning rate  $\alpha$  crucially governs how far we step against the gradient. Large  $\alpha$  can overshoot, while small  $\alpha$  converges slowly. Tuning  $\alpha$  finesses the descent.

With this mathematical grounding established, we now introduce a visual metaphor for added intuition.

#### 3.2 Foggy Mountain Visualization.

Imagine descending a foggy mountain to find water in a valley below. At each step, you measure the steepest slope and walk that direction. Incrementally, you traverse the descents until reaching the bottom – the oasis.

This metaphor elucidates gradient descent (Figure [3.1](#)):

- The mountain terrain depicts the objective function landscape.
- Standing corresponds to parameter values  $x$ .
- Measuring the steepest downhill direction is computing the gradient  $\text{Of}(x)$ .
- Taking a step is the gradient descent update.
- The valley bottom represents the minimized objective.

Like adjusting downhill step sizes, tuning the learning rate  $\alpha$  balances between convergence and overshooting. This vivid metaphor complements the mathematics, building intuition. We now transition from theory to practice through an implementational case study.

### 4. Linear Regression Case Study

We now demonstrate gradient descent by fitting a linear regressor to synthetic data.

#### 4.1 Generate Data.

We generate data according to:

$$y = \theta_0 + \theta_1 x + \epsilon$$

where  $\epsilon \sim N(0, 0.1)$ . We set  $\theta_0 = 1$ ,  $\theta_1 = 2$ , and  $x \in [-5, 5]$ . Below visualizes the resulting 100 points (Figure [4.1](#)).

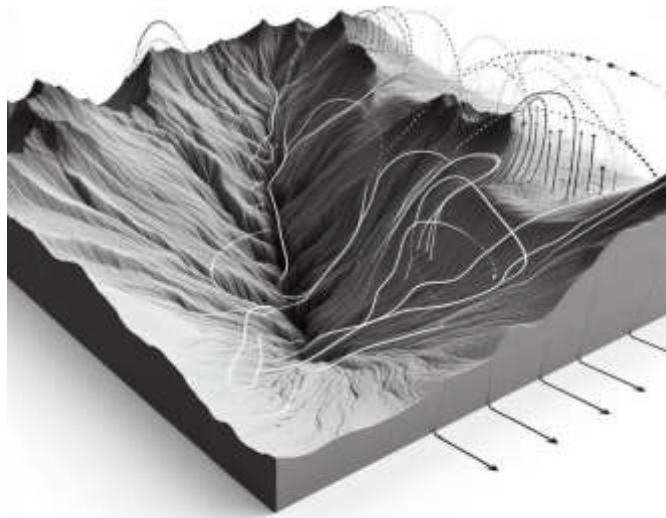


Figure 3.1. Mountain descent as gradient descent metaphor.

Our goal is recovering  $\theta_0, \theta_1$  via gradient descent.

**Gradient Descent Implementation.** We implement gradient descent in Python. The mean squared error objective is:

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

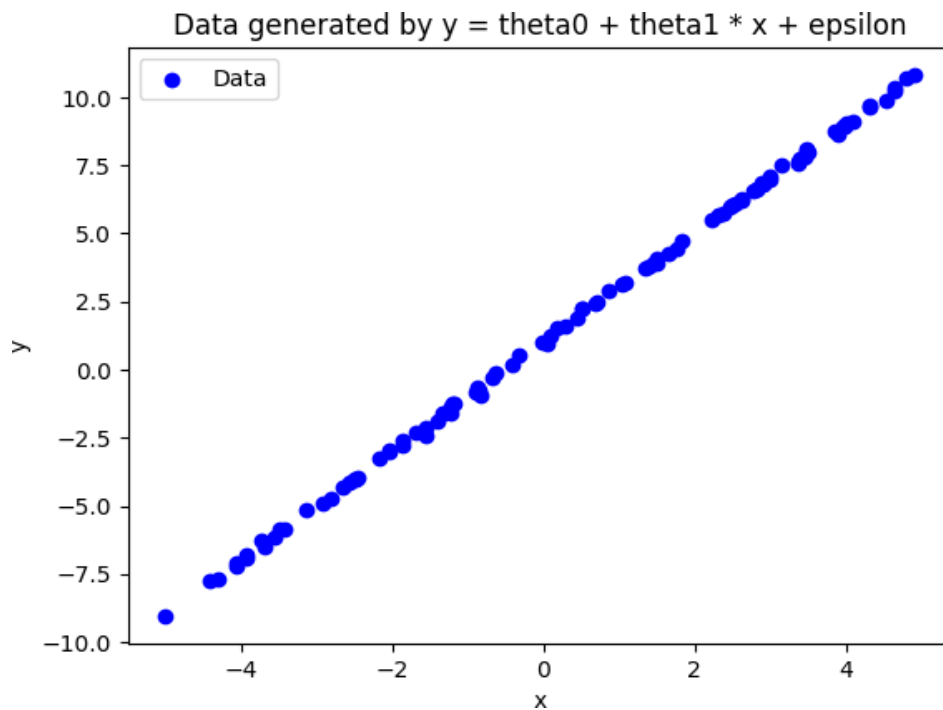


Figure 4.1. Synthetic linear regression data.

where  $h_{\theta}(x) = \theta_0 + \theta_1 x$ . Taking derivatives yields the gradient:

$$\nabla J(\theta) = \begin{bmatrix} \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) \\ \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x^{(i)} \end{bmatrix}$$

We initialize  $\theta_0, \theta_1$  randomly and implement the update:

$$\theta^{(t+1)} = \theta^{(t)} - \alpha \nabla J(\theta^{(t)})$$

Below shows the Python code. Running gradient descent recovers slopes resembling the ground truth (Figure 4.2).

This showcase provides an end-to-end demonstration of gradient descent based on the formalisms derived earlier. The appendix contains the full code.

---

#### Algorithm 1 Gradient Descent for Linear Regression

---

**Input:** data  $x$ , labels  $y$ , learning rate  $\alpha$

Initialize  $\theta_0, \theta_1$  randomly

**repeat**

    Compute predictions  $h_{\theta}(x)$

    Compute gradient  $\nabla J(\theta)$

    Update  $\theta = \theta - \alpha \nabla J(\theta)$

**until** convergence

---

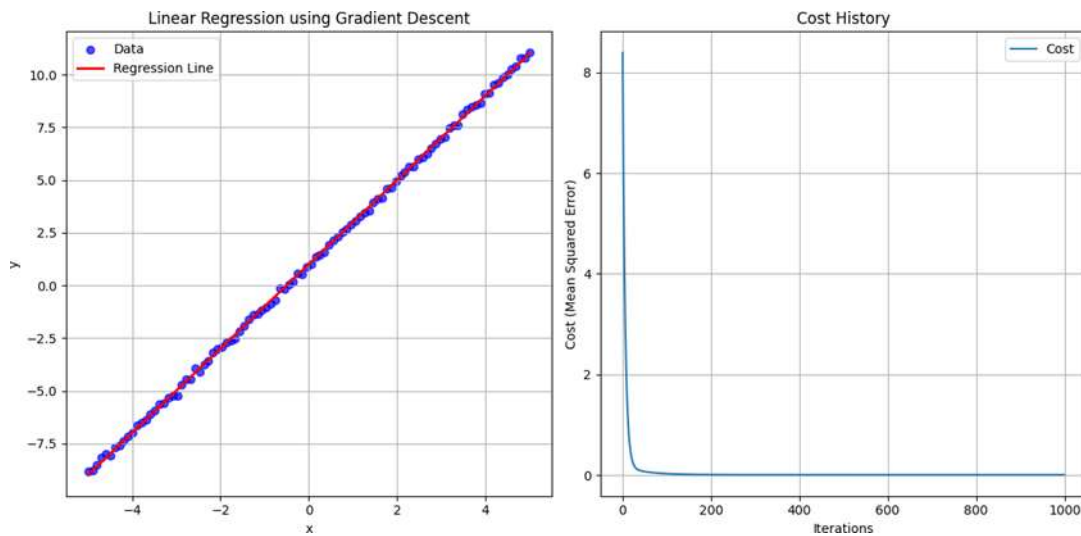


Figure 4.2. Gradient descent fitted line.

---

## 5. Conclusions

This paper presented a holistic perspective on foundational gradient descent techniques. We derived gradient descent from first principles, developed vivid metaphors relating it to foggy mountain descents, and demonstrated it through an implementational case study. Blending equations, visuals, and code aims to render gradient descent profoundly comprehensible.

Future work includes extending this pedagogical gradient descent narrative to related algorithms like momentum, Adagrad, and Adam. Mastering gradient descent provides a springboard for grasping more advanced techniques. But the core intuitions and mechanisms detailed here permeate modern deep learning. We hope this paper provides both usable techniques and lasting intuition.

---

### Appendix A. Programming Code

```
import numpy as np
import matplotlib.pyplot as plt
```

```
# Generate synthetic data np.random.seed(0)

theta_0 = 1
theta_1 = 2
epsilon_mean = 0
epsilon_stddev = 0.1
num_points = 100
x = np.linspace(-5, 5, num_points)
epsilon = np.random.normal(epsilon_mean, epsilon_stddev, num_points) y = theta_0 + theta_1 * x + epsilon

# Gradient Descent function
def gradient_descent(x, y, learning_rate, num_iterations): m = len(x)
theta_0 = np.random.rand() theta_1 = np.random.rand()
history = [] # To store the cost at each iteration

for i in range(num_iterations): predictions = theta_0 + theta_1 * x error = predictions - y

# Calculate gradients
gradient_0 = (1/m) * np.sum(error) gradient_1 = (1/m) * np.sum(error * x)

# Update parameters
theta_0 -= learning_rate * gradient_0 theta_1 -= learning_rate * gradient_1

# Calculate the mean squared error (objective function) cost = (1/(2*m)) * np.sum(error**2) history.append(cost)
return theta_0, theta_1, history # Hyperparameters
learning_rate = 0.01
num_iterations = 1000

# Run gradient descent theta_0_opt, theta_1_opt,
cost_history = gradient_descent(x, y, learning_rate, num_iterations)

# Plot the data, regression line, and cost history plt.figure(figsize=(12, 6))

# Data and regression line plot plt.subplot(1, 2, 1)
plt.scatter(x, y, label=' Data' , color=' blue' , alpha=0.7) plt.plot(x, theta_0_opt + theta_1_opt * x,
label=' Regression Line' , color=' red' , linewidth=2) plt.xlabel(' x' )
```

```

plt.ylabel(' y' )
plt.title(' Linear Regression using Gradient Descent' )
plt.legend() plt.grid(True)

# Cost history plot plt.subplot(1, 2, 2)
plt.plot(cost_history, label=' Cost' ) plt.xlabel(' Iterations' ) plt.ylabel(' Cost (Mean Squared Error)' ) plt.title(' Cost History' )
plt.legend() plt.grid(True)

plt.tight_layout() plt.show()

print(f"Optimal theta_0: {theta_0_opt}") print(f"Optimal theta_1: {theta_1_opt}")

```

#### AVAILABILITY OF DATA AND MATERIALS

The author confirms that the data supporting the findings of this study are available within the article or its supplementary materials.

#### Acknowledgments

This work is also partially supported by a research grant from Shenzhen Municipal Finance.

#### References

- [1] Mohammad Mohammadi Amiri and Deniz Gündüz. Machine learning at the wireless edge: Distributed stochastic gradient descent over-the-air. *IEEE Transactions on Signal Processing*, 68:2155–2169, 2020.
- [2] Rebecca Castano, Tara Estlin, Robert C Anderson, Daniel M Gaines, Andres Castano, Benjamin Bornstein, Caroline Chouinard, and Michele Judd. Oasis: Onboard autonomous science investigation system for opportunistic rover science. *Journal of Field Robotics*, 24(5):379–397, 2007.
- [3] Natarajan Deepa, B Prabadevi, Praveen Kumar Maddikunta, Thippa Reddy Gadekallu, Thar Baker, M Ajmal Khan, and Usman Tariq. An ai-based intelligent system for healthcare analysis using ridge-adaline stochastic gradient descent classifier. *The Journal of Supercomputing*, 77:1998–2017, 2021.
- [4] Jörg Fliege and Benar Fux Svaiter. Steepest descent methods for multicriteria optimization. *Mathematical methods of operations research*, 51:479–494, 2000.
- [5] Hossein Gholamalinejad and Hossein Khosravi. Whitened gradient descent, a new updating method for optimizers in deep neural networks. *Journal of AI and Data Mining*, 10(4):467–477, 2022.
- [6] Harsh Gupta, Rayadurgam Srikant, and Lei Ying. Finite-time performance bounds and adaptive learning rate selection for two time-scale reinforcement learning. *Advances in Neural Information Processing Systems*, 32, 2019.
- [7] Mark Jenkinson, Peter Bannister, Michael Brady, and Stephen Smith. Improved optimization for the robust and accurate linear registration and motion correction of brain images. *Neuroimage*, 17(2):825–841, 2002.
- [8] George Em Karniadakis, Ioannis G Kevrekidis, Lu Lu, Paris Perdikaris, Sifan Wang, and Liu Yang. Physics-informed machine learning. *Nature Reviews Physics*, 3(6):422–440, 2021.
- [9] Jiwon Kim, Jung Kwon Lee, and Kyoung Mu Lee. Deeply-recursive convolutional network for image super-resolution. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1637–1645, 2016.
- [10] Mona Lütjens, Thomas P Kersten, Boris Dorschel, and Felix Tschirschwitz. Virtual reality in cartography: Immersive 3d visualization of the arctic clyde inlet (canada) using digital elevation models and bathymetric data. *Multimodal Technologies and Interaction*, 3(1):9, 2019.



- 
- [11] Sajjad Emdadi Mahdimahalleh. Revolutionizing wireless networks with federated learning: A comprehensive review. *arXiv preprint arXiv:2308.04404*, 2023.
- [12] Elinor Ostrom. Background on the institutional analysis and development framework. *Policy studies journal*, 39(1):7–27, 2011.
- [13] Gregory S Parnell, Terry Bresnick, Steven N Tani, and Eric R Johnson. *Handbook of decision analysis*. John Wiley & Sons, 2013.
- [14] Martin Riedmiller and Heinrich Braun. A direct adaptive method for faster backpropagation learning: The rprop algorithm. In *IEEE international conference on neural networks*, pages 586–591. IEEE, 1993.
- [15] Sebastian Ruder. An overview of gradient descent optimization algorithms. *arXiv preprint arXiv:1609.04747*, 2016.
- [16] Jonathan Richard Shewchuk et al. An introduction to the conjugate gradient method without the agonizing pain, 1994.
- [17] Reza Shokri and Vitaly Shmatikov. Privacy-preserving deep learning. In *Proceedings of the 22nd ACM SIGSAC conference on computer and communications security*, pages 1310–1321, 2015.
- [18] Seiya Tokui, Kenta Oono, Shohei Hido, and Justin Clayton. Chainer: a next-generation open source framework for deep learning. In *Proceedings of workshop on machine learning systems (LearningSys) in the twenty-ninth annual conference on neural information processing systems (NIPS)*, volume 5, pages 1–6, 2015.
- [19] Eugene Walach and Bernard Widrow. The least mean fourth (lmf) adaptive algorithm and its family. *IEEE transactions on Information Theory*, 30(2):275–283, 1984.
- [20] Fei Wang, Daniel Zheng, James Decker, Xilun Wu, Grégory M Essertel, and Tiark Rompf. Demystifying differentiable programming: Shift/reset the penultimate backpropagator. *Proceedings of the ACM on Programming Languages*, 3(ICFP):1–31, 2019.
- [21] Weixuan Wang, Choon Meng Lee, Jianfeng Liu, Talha Colakoglu, and Wei Peng. An empirical study of cyclical learning rate on neural machine translation. *Natural Language Engineering*, 29(2):316–336, 2023.
- [22] Alan Wrigglesworth et al. *Credit scoring using machine learning: an application of deep learning*. PhD thesis, University of Pretoria, 2021.
- [23] Yifan Yang, Alec Koppel, and Zheng Zhang. A gradient-based approach for online robust deep neural network training with noisy labels. *arXiv preprint arXiv:2306.05046*, 2023.
- [24] Yifei Zhang. *Real-time multimodal semantic scene understanding for autonomous UGV navigation*. PhD thesis, Université Bourgogne Franche-Comté, 2021.