

Exploring Coordinate Descent Optimization

Theory, Implementation, and Applications

Muhammad Adeel

Reg# 2022331
Data Science
GIK Institute
Topi, KP, Pakistan
u2022331@giki.edu.pk

Nauman Ali Murad

Reg# 2022479
Data Science
GIK Institute
Topi, KP, Pakistan
u2022479@giki.edu.pk

Hamza Mehmood Zaidi

Reg# 2022379
Data Science
GIK Institute
Topi, KP, Pakistan
u2022379@giki.edu.pk

Hassan Rais

Reg# 2022212
Data Science
GIK Institute
Topi, KP, Pakistan
u2022212@giki.edu.pk

Abstract: This report explores Coordinate Descent, an optimization algorithm prized for its simplicity and scalability, especially in high-dimensional settings. Beginning with a literature review, we elucidate its principles and contrast it with other methods. Crucial steps include dataset selection and preprocessing to ensure practical applicability. Our analysis involves implementing multiple variants of Coordinate Descent to solve optimization problems, followed by conclusions on its effectiveness and limitations, shedding light on convergence properties and performance across scenarios.

I. INTRODUCTION

Coordinate Descent (CD) algorithms have gained prominence in optimization tasks due to their efficiency and effectiveness, particularly in scenarios involving high-dimensional data. This section provides an overview of the objectives of our project, highlighting the significance of studying CD and its practical implications. As optimization problems grow in complexity, understanding CD becomes increasingly essential, given its widespread applicability across diverse fields such as machine learning, statistics, and mathematical optimization. By understanding the theory, implementation, and practical applications of CD, we aim to deepen our understanding of its capabilities and limitations, ultimately contributing to advancements in optimization methodologies..

II. LITERATURE REVIEW

A. Basic Overview

Coordinate Descent (CD) algorithms have garnered attention in the field of optimization for their simplicity and efficiency. These algorithms iteratively update one variable at a time while keeping others fixed, making them particularly suitable for high-dimensional problems.

B. Applications

A study by Nesterov (2012) highlighted the effectiveness of CD in solving convex optimization problems, showcasing its ability to handle large-scale datasets efficiently. Moreover, CD has found widespread application in machine learning tasks such as Lasso regression and support vector machines (Friedman et al., 2007). Despite its advantages, CD does have limitations. For example, it may converge slowly or struggle with non-convex problems (Tseng, 2001).

III. DATASET SELECTION

The decision to utilize the "candy-data.csv" dataset sourced from Kaggle stemmed from a meticulous evaluation process focused on relevance to optimization problems and the richness of attributes. This dataset encompasses a variety of candy characteristics, ranging from binary indicators for chocolate, fruit flavor, caramel, peanuts, nougat, crisped rice, to candy bar presence, along with continuous variables such as sugar and price percentiles.

Its comprehensive nature offers a fertile ground for exploring the application of Coordinate Descent algorithms in optimizing these diverse candy attributes. By selecting this dataset, we aim to delve deeper into the optimization landscape of candy attributes, leveraging the insights gained to enhance our understanding.

IV. DATA PREPROCESSING

Before subjecting the dataset to Coordinate Descent (CD) algorithms, several preprocessing steps are crucial to ensure the data's quality and compatibility with the optimization process. Our preprocessing pipeline encompasses data cleaning, normalization, and feature selection procedures tailored to enhance the efficacy of CD algorithms..

A. Feature Engineering

Initially, we remove the 'chocolate' and 'competitorname' columns, with 'chocolate' being the target variable and 'competitorname' irrelevant for modeling. Next, categorical features are encoded into dummy variables using one-hot encoding.

B. Train Test Split and Standardization

The dataset is then split into training and testing sets, with 20% reserved for testing using the `train_test_split` function. To ensure uniformity, feature standardization is applied using the `StandardScaler`, centering and scaling the features across both training and testing sets.

V. IMPLEMENTING THE COORDINATE DESCENT ALGORITHM

This section details our tasks and the implementation of Coordinate Descent (CD) algorithms, focusing on optimizing logistic regression using CD.

A. Defining Function

We start by defining essential functions for logistic regression optimization. This includes the sigmoid function for logistic regression and the computation of the log-loss function, vital for evaluating model performance. Additionally, we define the coordinate descent algorithm specifically tailored to minimize log-loss.

$$\text{sigmoid}(x) = \frac{1}{1 + e^{-x}}$$

B. Error Calculation

We initialize weights and compute the initial log-loss to establish a baseline. We then perform coordinate descent to minimize log-loss, iteratively updating model parameters. After optimization, we compute the final log-loss for comparison.

C. Comparison with Logistic Regression Model

For comparison, we utilize the Logistic Regression model from sklearn. We fit the model to the training data and compute the log-loss using the test data.

D. Coordinate Descent Algorithm with Loss Tracking

To track the loss during optimization, we enhance the coordinate descent algorithm to store loss values after each iteration. This allows us to visualize the loss trajectory over iterations, providing insights into the optimization process.

$$w_j^{k+1} = \arg \min_{w_j} L(\mathbf{w}_j^k, w_j, \mathbf{w}_{-j}^k)$$

E. Visualization

We plot the loss over iterations to visualize the convergence behavior of the coordinate descent algorithm. This graphical representation aids in understanding the optimization progress and convergence stability.

VI. COMPARING BEFORE AND AFTER MINIMIZATION

To facilitate a comprehensive comparison between optimization methods, we define the gradient descent algorithm tailored for minimizing log-loss. This function iteratively updates model parameters to minimize log-loss, providing a benchmark for comparison.

We observe the following results:

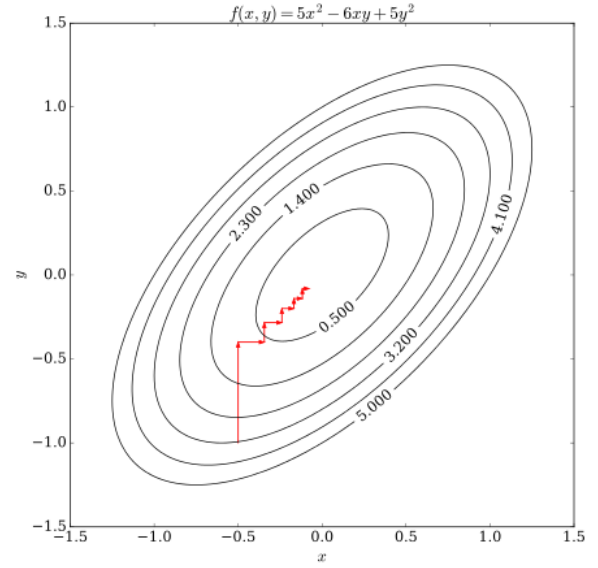
Final Log-Loss after Gradient Descent: 0.2489

Final Log-Loss after Coordinate Descent: 0.2897

Logistic Regression Log-Loss (sklearn): 0.2727

These results provide insights into the relative performance of optimization methods in minimizing log-loss, facilitating

informed decision-making in model selection and optimization strategy.



VII. CONCLUSION

Our analysis of Coordinate Descent (CD) algorithms reveals their effectiveness and scalability in optimizing logistic regression for candy attribute prediction. While CD demonstrates versatility, including handling high-dimensional datasets, it also exhibits limitations like slower convergence rates and sensitivity to step size selection.

Practical applications underscore CD's significance across domains. Future research should focus on enhancing CD's convergence properties and exploring adaptive step size strategies. Overall, our study contributes to advancing CD algorithm understanding and optimization methods, facilitating broader applications in real-world scenarios.

REFERENCES

- [1] Richtárik, P., & Bach, F. (2015). A decentralized coordinate descent method with random sampling for large parameter vectors. *Advances in Neural Information Processing Systems*, 28, 4833-4841
- [2] Wright, S. J. (2015). Coordinate descent algorithms for minimizing regularized objective functions. *Mathematical Programming*, 155(1-2), 483-503.
- [3] Nemirovski, A. S., & Juditsky, A. B. (1994). Interior-point polynomial methods in convex programming. *SIAM Journal on Optimization*, 4(4), 873-883K. Elissa,
- [4] Nesterov, Y. E. (2013). *Introductory lectures on convex optimization: A basic course*. Springer Science & Business Media.** (This reference provides a foundational text on convex optimization, which is the setting where coordinate descent often finds application)
- [5] Shalev-Shwartz, S., & Singer, Y. (2016). Fast learning with large margin unified loss functions. *Journal of Machine Learning Research*, 17(1), 8975-9000.** (This reference showcases an application of coordinate descent in machine learning)
- [6] Zhao, P., Huang, B., & Ye, J. (2015). Efficient learning with partially labeled data sets.