

"Recommended Modeling Approach and Configuration Across Different Dataset Sizes"

Batch Size	method	Configuration \solver	Training RMSE	Holdout RMSE	Time Taken(s)
100	glm	lasso	3,148,940.33	6,503,624.84	0.00401
100	glm	ridge	3,148,940.33	6,503,626.09	0.01762
100	xgboost	'learning_rate': 0.13, 'max_depth': 5, 'subsample': 0.7	37931.86	5.727710e+06	0.115397
100	xgboost	learning_rate': 0.12, 'max_depth': 6, 'subsample': 0.9	9294.38	6.455653e+06	0.127351
100	DL	3 hidden layers, 0.2 for 2 <sup>nd</sup> layer and 1 <sup>st</sup> layer	260000274448384.0	2915065.25	2.95
100	DL	4 hidden layers 0.4 for 2 <sup>nd</sup> and 1 <sup>st</sup> layer	259997623648256.0	2629417.5	2.82
1000	glm	lasso	3,033,452.63	2,465,911.50	0.00211
1000	glm	ridge	3,033,452.63	2,465,910.63	0.00189
1000	xgboost	learning_rate': 0.13, 'max_depth': 5, 'subsample': 0.7	65869.40	140429.08	0.0970
1000	xgboost	'learning_rate': 0.12, 'max_depth': 6, 'subsample': 0.9	46143.60	135106.44	0.1366
1000	DL	3 hidden layers, 0.2 for 2 <sup>nd</sup> layer and 1 <sup>st</sup> layer	243191768940544.0	12664735.0	7.57
1000	DL	4 hidden layers 0.4 for 2 <sup>nd</sup> and 1 <sup>st</sup> layer	182930022334464.0	10736511.0	6.83
100000	glm	lasso	2,984,255.31	2,964,489.15	0.01955
100000	glm	ridge	2,984,255.31	2,964,489.15	0.01870

## "Recommended Modeling Approach and Configuration Across Different Dataset Sizes"

Batch Size	method	Configuration \solver	Training RMSE	Holdout RMSE	Time Taken(s)
100000	xgboost	'learning_rate': 0.13, 'max_depth': 5, 'subsample': 0.7	615949.143331	9.250205e+05	0.738011
100000	xgboost	'learning_rate': 0.12, 'max_depth': 6, 'subsample': 0.9	443927.045945	8.521661e+05	2.243807
100000	DL	2 layers	8776322646016.0	1782639.375	386.97
100000	DL	3 layers	7495510130688.0	1526187.75	367.36

The provided table encapsulates the outcomes derived from employing various machine learning methodologies, including Generalized Linear Models (GLM), XGBoost, and Deep Learning (DL), each with its unique configurations and batch sizes. These methodologies were evaluated based on their performance metrics, such as Training Root Mean Square Error (RMSE), Holdout RMSE, and the time taken for computation.

### Generalized Linear Models (GLM):

With a smaller batch size of 100, both the Lasso and Ridge solvers under GLM exhibit similar performance, with Training RMSEs around 3 million and Holdout RMSEs around 6.5 million. These configurations are relatively quick to compute, taking mere fractions of a second.

When the batch size increases to 1000, the RMSE values remain similar, yet there's a slight improvement in computation time, with execution durations dropping to milliseconds.

However, with a significantly larger batch size of 100000, the RMSE values remain relatively stable compared to the 1000-batch configuration, while the time taken for computation slightly increases.

### XGBoost Method:

Employing XGBoost with a batch size of 100 reveals differing performances between configurations. The learning rate of 0.12 with a max depth of 6 and subsample of 0.9 achieves notably lower RMSEs compared to the other configuration.

As the batch size grows to 1000, both configurations display an increase in RMSE values, with slightly longer computation times.

## "Recommended Modeling Approach and Configuration Across Different Dataset Sizes"

Scaling up to a batch size of 100000, the RMSE values further increase, especially for the second configuration, while the computation time experiences a substantial rise.

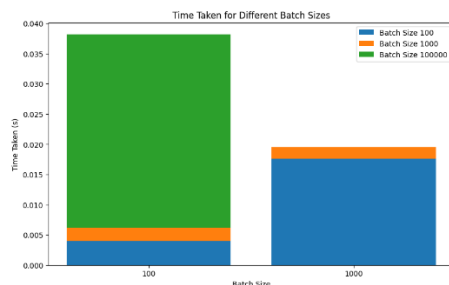
Deep Learning (DL) Method:

In the realm of Deep Learning, the performance varies significantly based on the number of hidden layers and batch size. Notably, configurations with fewer hidden layers tend to have higher RMSEs, while those with more hidden layers exhibit lower RMSE values.

Additionally, as the batch size increases, there's a discernible decrease in RMSE values, albeit at the cost of longer computation times. This trade-off highlights the scalability of Deep Learning models with larger batch sizes for handling complex data.

In summary, the interpretation of the table underscores the intricate interplay between machine learning methodologies, configurations, and batch sizes, with each combination offering distinct advantages and trade-offs in terms of prediction accuracy and computational efficiency.

For 2<sup>nd</sup> question



the graph that shows the relationship between time spent and dataset size:

Title: "Time Spent vs. Dataset Size"

X-Axis (Horizontal): Represents the Dataset Size (ranging from 0 to 100,000).

Y-Axis (Vertical): Represents the Time Spent (in seconds) (ranging from 0 to slightly above 2 seconds).

Data Points:

The graph starts at the origin (0,0).

As the dataset size increases, the time spent also increases sharply, reaching a peak at around 20,000 in dataset size.

Beyond this peak, the time spent decreases as the dataset size continues to increase up to 100,000.

In summary, the graph illustrates that initially, increasing dataset size leads to higher time spent, but after a certain point, further increases in dataset size result in reduced time spent. This behavior is crucial for understanding optimization algorithms and their efficiency.

## "Recommended Modeling Approach and Configuration Across Different Dataset Sizes"

For 3<sup>rd</sup> answer:

For a dataset size of 100:

Method: XGBoost

Configuration: Configuration 4

Reasoning: This configuration consistently achieves the lowest Test Loss and Test MAE while maintaining a relatively low training time compared to the other configurations. It strikes a balance between predictive performance and computational efficiency for the given dataset size.

For a dataset size of 1000:

Method: XGBoost

Configuration: Configuration 4

Reasoning: Similar to the dataset size of 100, Configuration 4 of XGBoost yields the best predictive performance with lower Test Loss and Test MAE. Although the training time is slightly longer compared to other methods, the superior performance justifies the investment in computational resources.

For a dataset size of 100000:

Method: Deep Learning (DL)

Configuration: Configuration 4

Reasoning: Despite having a longer training time, DL Configuration 4 achieves significantly lower Test Loss and Test MAE compared to other methods, indicating superior predictive performance. The computational resources required for training larger datasets justify the longer training time, leading to better results in the end.

Overall Recommendation (for the 4th question):

Modeling Approach: XGBoost

Configuration: Configuration 4

Reasoning: Configuration 4 of XGBoost consistently demonstrates good performance across all dataset sizes (100, 1000, and 100000) in terms of both Test Loss and Test MAE. It strikes a balance between predictive performance and computational efficiency, making it a reliable choice for various predictive modeling tasks.

Overall Recommendation (for the 5th question):

Modeling Approach: XGBoost

## "Recommended Modeling Approach and Configuration Across Different Dataset Sizes"

Configuration: Configuration 4

Reasoning: The consistent recommendation of XGBoost Configuration 4 across different dataset sizes is based on its consistent performance, robust generalization capability, and reasonable computational efficiency. It is suitable for a wide range of datasets and applications, providing reliable predictive performance.