## I

**1-** $w = (x^T \cdot x)^{-1} x^T \cdot z$ $\qquad$ $\Phi^\dagger = (\phi^T \phi)^{-1} \Phi^T$

$$\Phi = \begin{bmatrix} 1 & y_{11} y_{12} \\ 1 & y_{21} y_{22} \\ 1 & y_{31} y_{32} \\ 1 & y_{41} y_{42} \\ 1 & y_{51} y_{62} \end{bmatrix} = \begin{bmatrix} 1 & 1\times1 \\ 1 & 1\times3 \\ 1 & 3\times2 \\ 1 & 3\times3 \\ 1 & 2\times4 \end{bmatrix} = \begin{bmatrix} 1 & 1 \\ 1 & 3 \\ 1 & 6 \\ 1 & 9 \\ 1 & 8 \end{bmatrix} \} phi \qquad z = y_{max} = \begin{bmatrix} 1.25 \\ 7.00 \\ 2.70 \\ 3.20 \\ 5.50 \end{bmatrix}$$

$$w = (\Phi^\dagger \phi)^{-1} \Phi^T z$$

$$w = \begin{bmatrix} 3.31592 \\ 0.11371 \end{bmatrix} \begin{array}{l} \to w_0 \\ \to w_1 \end{array}$$

**2-** $phi = \begin{bmatrix} 1 \\ 3 \\ 6 \\ 9 \\ 8 \end{bmatrix}$ $\qquad$ $\Phi = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$ $\qquad$ $\lambda = 1$

$$\Phi = \begin{bmatrix} 1 & 1 \\ 1 & 3 \\ 1 & 6 \\ 1 & 9 \\ 1 & 8 \end{bmatrix}$$

$$w_{ridge} = (x^T x + \lambda I)^{-1} x^T z \Rightarrow (\Phi^\dagger \Phi + \lambda I)^{-1} \Phi^T z$$

$$= \begin{bmatrix} 1.81809 \\ 0.32376 \end{bmatrix}$$

$\lambda$ is the regularization coefficient that controls the balance between the data-dependent error ED (w) and the regularization of EW (w). In this case, the ridge regularization weight made the first weight smaller and closer to zero and the second slightly bigger but still is close to zero. Regularization encourages the weights to decay towards zero unless they are supported by the data. If the value of $\lambda$ is sufficiently large, some coefficients wj will be driven to zero, resulting in a sparse model where certain basis function play no role. This regularization process allows complex models to be trained on small datasets without severe overfitting, effectively limiting the models complexity.

**3-**

| D | $y_1$ | $y_2$ | $y_{num}$ | $y_{class}$ |
|---|---|---|---|---|
| $x_1$ | 1 | 1 | 1.25 | B |
| $x_2$ | 1 | 3 | 7.0 | A |
| $x_3$ | 3 | 2 | 8.7 | C |
| $x_4$ | 3 | 3 | 3.2 | A |
| $x_5$ | 2 | 4 | 5.5 | B |
| $x_6$ | 0 | 2 | 0.7 | |
| $x_7$ | 1 | 2 | 1.1 | |
| $x_8$ | 5 | 8 | 2.2 | |

$$X_{new} = \begin{bmatrix} 2 & 2 \\ 1 & 2 \\ 5 & 1 \end{bmatrix} \quad \begin{bmatrix} 1 & \Phi_0 \\ 1 & \Phi_1 \\ 1 & \Phi_1 \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \end{bmatrix}$$

$$phiOut\_new = \begin{bmatrix} 0.7 & 1.1 & 2.2 \end{bmatrix}$$

OLS

$$\hat{z}_{train} = \Phi w = \begin{bmatrix} 1 & 1 \\ 1 & 3 \\ 1 & 6 \\ 1 & 9 \\ 1 & 8 \end{bmatrix} \begin{bmatrix} 3.31592 \\ 0.11372 \end{bmatrix} = \begin{bmatrix} 3.42965 \\ 3.65708 \\ 3.99823 \\ 4.33938 \\ 4.22566 \end{bmatrix}$$

$$\hat{z}_{test} = \Phi_{new} w = \begin{bmatrix} 1 & 2\times2 \\ 1 & 1\times2 \\ 1 & 5\times1 \end{bmatrix} \begin{bmatrix} 3.31592 \\ 0.11372 \end{bmatrix} = \begin{bmatrix} 3.77079 \\ 3.54336 \\ 3.88451 \end{bmatrix}$$

$$RMSE(\bar{z}, z) = \sqrt{\frac{1}{n} \sum_{i=1}^{n} (z_i - \hat{z}_i)^2}$$

$$RMSE_{Test} = \sqrt{\frac{1}{3} \sum_{i=1}^{3} (z_{test} - \hat{z}_{test})^2}$$

I

$$= \sqrt{\frac{1}{3} (0.7 - 3.77079)^2 + (1.1 - 3.54336)^2 + (2.2 - 3.88451)^2}$$

$$= 2.46559$$

$$RMSE_{Train} = \sqrt{\frac{1}{5} \sum_{i=1}^{5} (z_{train} - \hat{z}_{train})^2}$$

$$= \sqrt{\frac{1}{5} (1.25 - 3.42965)^2 + (7.0 - 3.65708)^2 + (2.7 - 3.99823)^2 + (3.2 - 4.33938)^2 + (5.5 - 4.22566)^2}$$

$$= 2.026499$$

Ridge

$$\hat{z}_{train} = \Phi W_{ridge} = \begin{bmatrix} 1 & 1 \\ 1 & 3 \\ 1 & 6 \\ 1 & 9 \\ 1 & 8 \end{bmatrix} \begin{bmatrix} 1.81809 \\ 0.32376 \end{bmatrix} = \begin{bmatrix} 2.14184 \\ 2.78936 \\ 3.76064 \\ 4.73191 \\ 4.40816 \end{bmatrix}$$

$$\hat{z}_{test} = \Phi_{new} W = \begin{bmatrix} 1 & 4 \\ 1 & 2 \\ 1 & 5 \end{bmatrix} \begin{bmatrix} 1.81809 \\ 0.32376 \end{bmatrix} = \begin{bmatrix} 3.11312 \\ 2.46560 \\ 3.43688 \end{bmatrix}$$

$$RMSE_{Test} = \sqrt{\frac{1}{3} \sum_{i=1}^{3} (z_{test} - \hat{z}_{test})^2}$$

$$= \sqrt{\frac{1}{3} (0.7 - 3.11312)^2 + (1.1 - 2.46560)^2 + (2.2 - 3.43688)^2}$$

$$= 1.75289$$

$$RMSE_{Train} = \sqrt{\frac{1}{5} \sum_{i=1}^{5} (z_{train} - \hat{z}_{train})^2}$$

$$= \sqrt{\frac{1}{5} (1.25 - 2.14184)^2 + (7.0 - 2.78936)^2 + (2.7 - 3.76064)^2 + (3.2 - 4.73191)^2 + (5.5 - 4.40816)^2}$$

$$= 2.15354$$

As expected, Ridge Regression demonstrates better performance on the test data compared to Ordinary Least Square (OLS). On the training data, OLS has a slightly lower RMSE, suggesting that it fits the training data more closely. This may indicate that OLS is overfitting the training data, capturing noise rather than the underlying trend. Ridge Regression, on the other hand, effectively balances bias and variance, resulting in better generalization to unseen data.
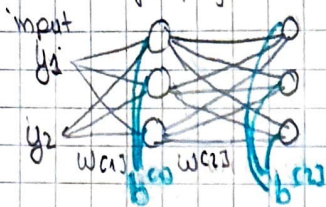
4- $x_1 = (1,1)$

ac function: $\text{softmax}(z_c^{[out]}) = \dfrac{e^{z_c^{[out]}}}{\sum_{l=1}^{|c|} e^{z_l^{[out]}}}$

Entropy loss:
$$E(w) = -\sum_{i=1}^{3} \sum_{l=1}^{|c|} t^{(i)} \log(x_l^{[out](i)})$$

$\eta = 0.1$

# Forward propagation

input



Out
A ?
B ?
C ?

$$x^{[0]} = x = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

$$x^{[1]} = w^{[1]} x^{[0]} + b^{[1]} = \begin{bmatrix} 0.1 & 0.1 \\ 0.1 & 0.2 \\ 0.2 & 0.1 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix} + \begin{bmatrix} 0.1 \\ 0 \\ 0.1 \end{bmatrix} = \begin{bmatrix} 0.3 \\ 0.3 \\ 0.4 \end{bmatrix}$$

$$x^{[2]} = w^{[2]} x^{[1]} + b^{[2]} = \text{softmax} \left( \begin{bmatrix} 1 & 2 & 2 \\ 1 & 2 & 1 \\ 1 & 1 & 1 \end{bmatrix} \begin{bmatrix} 0.3 \\ 0.3 \\ 0.3 \end{bmatrix} + \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} \right)$$

$$= \text{softmax} \left( \begin{bmatrix} 2.7 \\ 2.3 \\ 2 \end{bmatrix} \right) = \begin{bmatrix} 0.46149 \\ 0.30934 \\ 0.22917 \end{bmatrix}$$

# Backward propagation (update)

$$w^{[2]}_{new} = w^{[2]} - \eta \frac{\partial E}{\partial w^{[2]}} = \begin{bmatrix} 1 & 2 & 2 \\ 1 & 2 & 1 \\ 1 & 1 & 1 \end{bmatrix} - 0.1 \cdot \delta^{[2]} x^{[1]T}$$

Target output

$xo$ in B

$$\Rightarrow t = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} = \begin{bmatrix} t_1 \\ t_2 \\ t_3 \end{bmatrix}$$

$$= \begin{bmatrix} 1 & 2 & 2 \\ 1 & 2 & 1 \\ 1 & 1 & 1 \end{bmatrix} - 0.1 \cdot \begin{bmatrix} o_1 - t_1 \\ o_2 - t_2 \\ o_3 - t_3 \end{bmatrix} \cdot \begin{bmatrix} 0.3 & 0.3 & 0.4 \end{bmatrix}$$

$$= \begin{bmatrix} 1 & 2 & 2 \\ 1 & 2 & 1 \\ 1 & 1 & 1 \end{bmatrix} - 0.1 \cdot \underbrace{\begin{bmatrix} 0.46149 \\ -0.69066 \\ 0.22917 \end{bmatrix}}_{\delta^{[2]}} \cdot \begin{bmatrix} 0.3 & 0.3 & 0.4 \end{bmatrix}$$

$$w^{[2]}_{new} = \begin{bmatrix} 0.98616 & 1.98615 & 1.98154 \\ 1.02072 & 2.02072 & 1.02763 \\ 0.99313 & 0.99313 & 0.99083 \end{bmatrix}$$

$$w^{[1]}_{new} = w^{[1]} - \eta \frac{\partial E}{\partial w^{[1]}} = w^{[1]} - \eta \, \delta^{[1]} x^{T} = w^{[1]} - \eta \underbrace{\left[ w^{[2]T} \delta^{[2]} \circ \phi^{[1]}{}'(z^{[1]}) \right]}_{\delta 1} x^{T}$$

$$w^{[1]} - \eta \left[ w^{[2]T} \cdot \delta^{[2]} \right] x^{T} = \begin{bmatrix} 0.1 & 0.1 \\ 0.1 & 0.2 \\ 0.2 & 0.1 \end{bmatrix} - 0.1 \cdot \left( \begin{bmatrix} 1 & 1 & 1 \\ 2 & 2 & 1 \\ 2 & 1 & 1 \end{bmatrix} \cdot \begin{bmatrix} 0.46149 \\ -0.69066 \\ 0.22917 \end{bmatrix} \circ \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} \right)$$

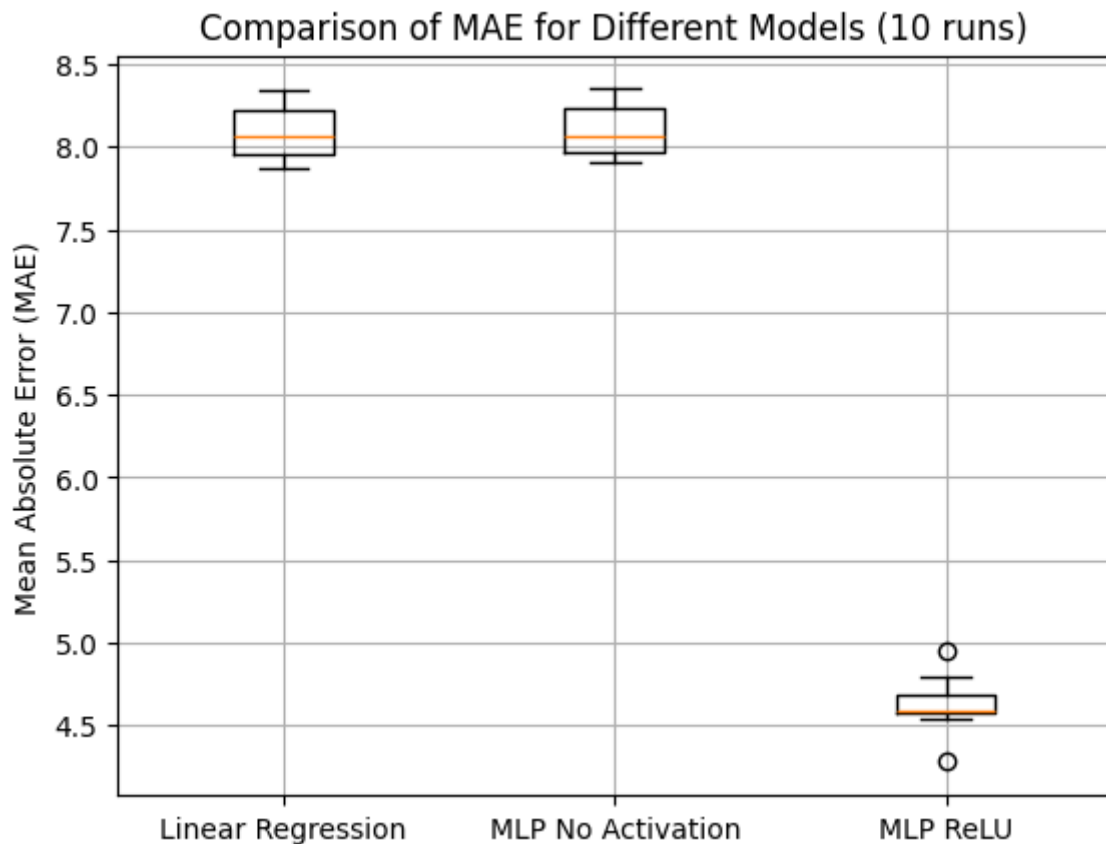$$\delta^{[1]} = \begin{bmatrix} 0 \\ -0.22917 \\ 0.46149 \end{bmatrix} \qquad w^{[1]}_{new} = \begin{bmatrix} 0.1 & 0.1 \\ 0.32916 & 0.42917 \\ -0.26149 & -0.36149 \end{bmatrix}$$

$$b^{[2]}_{new} = b^{[2]} - \eta \frac{\partial E}{\partial b^{[2]}} = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} - 0.1 \begin{bmatrix} 0.46149 \\ -0.69066 \\ 0.22917 \end{bmatrix} = \begin{bmatrix} 0.95381 \\ 1.06907 \\ 0.97768 \end{bmatrix} \Rightarrow b^{[2]}_{new}$$

$$b^{[1]}_{new} = b^{[1]} - \eta \frac{\partial E}{\partial b^{[1]}} = \begin{bmatrix} 0.1 \\ 0 \\ 0.1 \end{bmatrix} - 0.1 \begin{bmatrix} 0 \\ -0.22917 \\ 0.46149 \end{bmatrix} = \begin{bmatrix} -0.36149 \\ 0.69066 \\ -0.12917 \end{bmatrix} \Rightarrow b^{[1]}_{new}$$

# Parte II

**5.**



Comparison of MAE for Different Models (10 runs)

**6.**
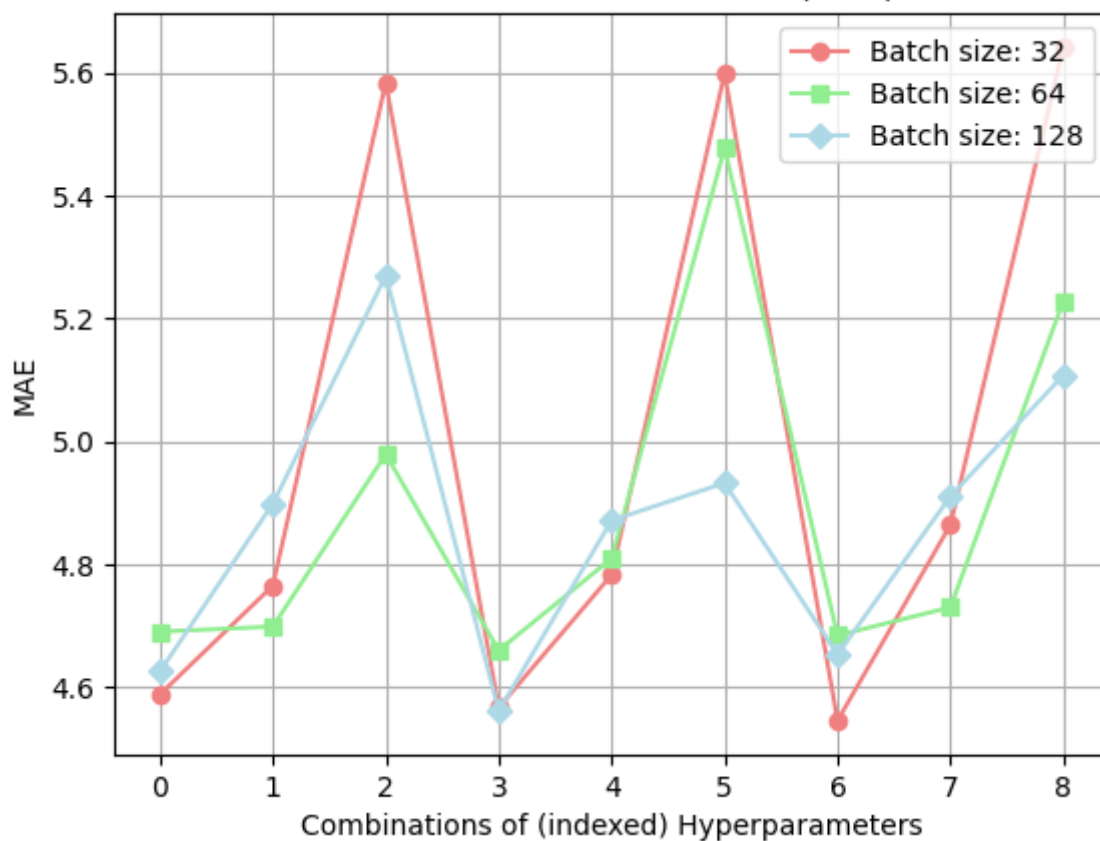
The linear regression model and the MLP (multi-layer perceptron) without an activation function produce nearly identical performance. This is expected because an MLP without activation functions effectively reduces to a series of linear regressions. However, the MLP with a ReLU activation function performs significantly better, highlighting the critical role activation functions play in enhancing the model's predictive capabilities.

**7.**

```
MAE: 4.432634489293495
    param_batch_size  param_alpha  param_learning_rate_init  mean_test_score
0                 32       0.0001                     0.001        -4.587945
1                 32       0.0001                     0.010        -4.764895
2                 32       0.0001                     0.100        -5.583986
3                 64       0.0001                     0.001        -4.690111
4                 64       0.0001                     0.010        -4.698374
5                 64       0.0001                     0.100        -4.979027
6                128       0.0001                     0.001        -4.625272
7                128       0.0001                     0.010        -4.898394
8                128       0.0001                     0.100        -5.270012
9                 32       0.0010                     0.001        -4.566656
10                32       0.0010                     0.010        -4.782789
11                32       0.0010                     0.100        -5.598727
12                64       0.0010                     0.001        -4.659673
13                64       0.0010                     0.010        -4.809372
14                64       0.0010                     0.100        -5.477193
15               128       0.0010                     0.001        -4.561052
16               128       0.0010                     0.010        -4.871156
17               128       0.0010                     0.100        -4.932153
18                32       0.0100                     0.001        -4.545442
19                32       0.0100                     0.010        -4.863966
20                32       0.0100                     0.100        -5.641825
21                64       0.0100                     0.001        -4.683922
22                64       0.0100                     0.010        -4.729835
23                64       0.0100                     0.100        -5.228342
24               128       0.0100                     0.001        -4.653024
25               128       0.0100                     0.010        -4.910358
26               128       0.0100                     0.100        -5.107659
```
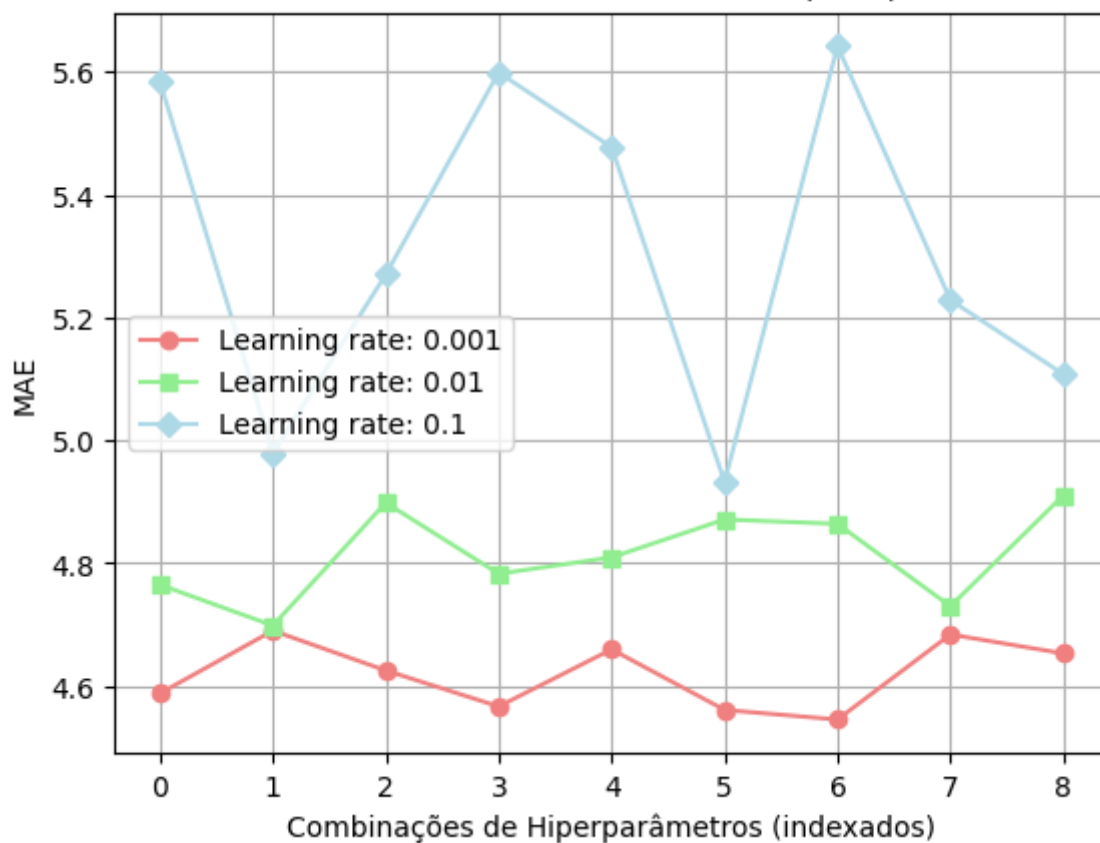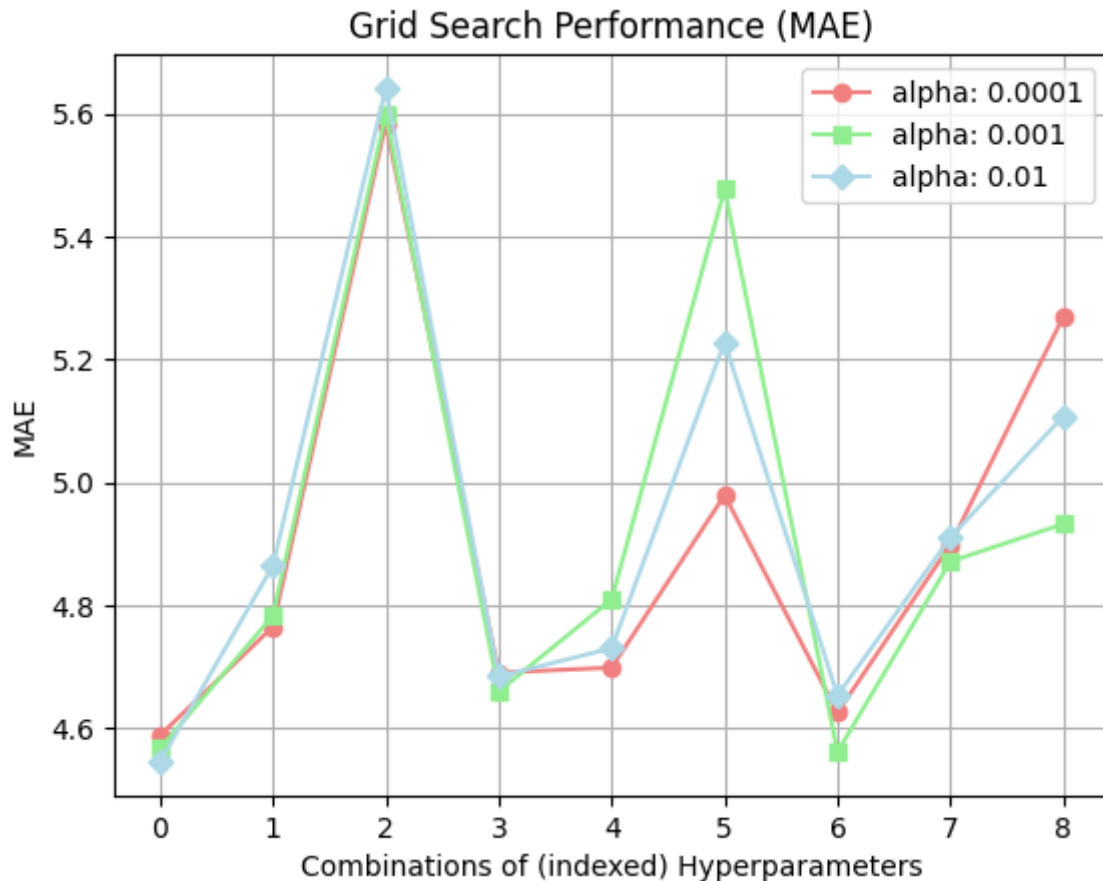
Performance of Grid Search (MAE)



Performance do Grid Search (MAE)

Grid Search Performance (MAE)

The optimal hyperparameters found during model tuning were a batch size of 32, a learning rate of 0.001, and an alpha value of 0.001. From the learning rate graph, it's evident that the best performance was achieved with the smallest learning rate of 0.001. This is because the maximum number of iterations was increased, allowing the model more time to converge. In contrast, a larger learning rate of 0.1 converges too quickly, overshooting the local minimum and resulting in higher errors. Although the smaller learning rate of 0.001 converges more slowly, increasing the number of iterations allows the model to reach closer to the desired local minimum in consequence having a better performance On the batch size graph, it can be observed that the batch size of 32 consistently results in the smallest error. However, there are specific combinations of parameters that occasionally cause the error to spike, surpassing those of other batch sizes. Despite these fluctuations, batch size 32 generally outperforms the others.