

I. Pen and Paper

1)

Pan - and - Paper

$$\Phi(y_1, y_2) = \begin{bmatrix} 1 & 1 & 1 \\ 1 & 3 & 6 \\ 1 & 6 & 9 \\ 1 & 9 & 12 \\ 1 & 12 & 15 \end{bmatrix} = \begin{bmatrix} 1 \\ 2 \\ 6 \\ 9 \\ 15 \end{bmatrix}$$

$$y = \begin{bmatrix} 1, 1, 1 \\ 1, 3, 6 \\ 1, 6, 9 \\ 1, 9, 12 \\ 1, 12, 15 \end{bmatrix}$$

$$\beta = (\mathbf{x}^T \mathbf{x})^{-1} \mathbf{x}^T y$$

$$= \left(\begin{bmatrix} 1 & 1 & 1 & 1 & 1 \\ 1 & 3 & 6 & 9 & 15 \end{bmatrix} \right)^{-1} \begin{bmatrix} 1 \\ 2 \\ 6 \\ 9 \\ 15 \end{bmatrix} = \begin{bmatrix} 1 & 1 & 1 & 1 & 1 \\ 1 & 3 & 6 & 9 & 15 \end{bmatrix} \begin{bmatrix} 1 \\ 2 \\ 6 \\ 9 \\ 15 \end{bmatrix}$$

$$= \begin{bmatrix} 5 & 27 \\ 27 & 191 \end{bmatrix}^{-1} \begin{bmatrix} 19,65 \\ 111,25 \end{bmatrix} = \frac{1}{726} \begin{bmatrix} 191 & 27 \\ 27 & 5 \end{bmatrix} \begin{bmatrix} 19,65 \\ 111,25 \end{bmatrix}$$

$$= \begin{bmatrix} 3,31593 \\ 0,11372 \end{bmatrix} \quad \beta_0 = 3,31593, \quad \beta_1 = 0,11372$$

$\hat{y}_{num} = 3,31593 + 0,11372 \cdot \Phi(y_1, y_2)$

 Resposta: $\hat{y}_{num} = 3,31593 + 0,11372 * \Phi(y_1, y_2)$

2)

$$\begin{aligned}
 & \text{Ridge: } \hat{w} = (x^T x + \lambda I)^{-1} x^T y \\
 & x^T x + \lambda I = \begin{bmatrix} 5 & 27 \\ 27 & 192 \end{bmatrix} + \lambda \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} = \begin{bmatrix} 6 & 27 \\ 27 & 192 \end{bmatrix} \\
 & \hat{w} = \begin{bmatrix} 6 & 27 \\ 27 & 192 \end{bmatrix}^{-1} \begin{bmatrix} 15,65 \\ 11,25 \end{bmatrix} = \frac{1}{423} \begin{bmatrix} 192 & 27 \\ -27 & 6 \end{bmatrix} \begin{bmatrix} 15,65 \\ 11,25 \end{bmatrix} \\
 & = \begin{bmatrix} 1,81808 \\ 0,32376 \end{bmatrix} \quad (\Rightarrow \beta_0 = 1,81808, \beta_1 = 0,32376) \\
 \therefore \hat{y}_{\text{Ridge}}(\text{Ridge}) &= 1,81808 + 0,32376 \cdot \phi(y_1, y_2)
 \end{aligned}$$

By using Ridge's formula, the value for the intercept decreased and the value for the slope increased. This is the result of the decrease of large coefficients, a characteristic of this regularization method. Since the data set is pretty small and therefore noisy, applying Ridge's regularization helps reduce variance in the data, preventing overfitting at the cost of introducing bias. It can have a stronger effect by increasing the λ parameter.

3)

$$\hat{y}_{\text{num}} = 2,31593 + 0,11372 \phi(y_1, y_2)$$

Grām without ridge:

$$\hat{y}_1 = 3,31593 + 0,11372 \times 1 = 3,42965$$

$$\hat{y}_2 = 3,31593 + 0,11372 \times 3 = 3,65709$$

$$\hat{y}_3 = 3,31593 + 0,11372 \times 6 = 3,99825$$

$$\hat{y}_4 = 3,31593 + 0,11372 \times 9 = 4,33941$$

$$\hat{y}_5 = 3,31593 + 0,11372 \times 8 = 4,22569$$

$$RMSE = \sqrt{\frac{(1,25 - 3,42965)^2 + (1,0 - 3,65709)^2 + (2,2 - 3,99825)^2 + (3,2 - 4,33941)^2 + (5,5 - 4,22569)^2}{5}}$$

$$= 2,02650$$

$$\hat{y}_{\text{num}} (\text{ridge}) = 1,81808 + 0,32376 + \phi(y_1, y_2)$$

Grām with Ridge:

$$\hat{y}_1 = 1,81808 + 0,32376 \times 1 = 2,14184$$

$$\hat{y}_2 = 1,81808 + 0,32376 \times 3 = 2,78936$$

$$\hat{y}_3 = 1,81808 + 0,32376 \times 6 = 3,76064$$

$$\hat{y}_4 = 1,81808 + 0,32376 \times 9 = 4,73192$$

$$\hat{y}_5 = 1,81808 + 0,32376 \times 8 = 4,40816$$

$$RMSE = \sqrt{\frac{(1,25 - 2,14184)^2 + (1,0 - 2,78936)^2 + (2,2 - 3,76064)^2 + (3,2 - 4,73192)^2 + (5,5 - 4,40816)^2}{5}}$$

$$= 2,15354$$

3.

$$\hat{y}_0 = 3,31593 + 0,11372 \times 4 = 3,44081$$

$$\hat{y}_1 = 3,31593 + 0,11372 \times 2 = 3,35433$$

$$\hat{y}_2 = 3,31593 + 0,11372 \times 5 = 3,88453$$

$$\text{RMSE} = \sqrt{\frac{(0,4 - 3,77081)^2 + (0,1 - 3,55333)^2 + (2,2 - 3,88453)^2}{3}}$$

$$= 2,46560$$

$$\hat{y}_0 (\text{Ridge}) = 1,81808 + 0,22376 \times 4 = 3,11212$$

$$\hat{y}_1 (\text{Ridge}) = 1,81808 + 0,22376 \times 2 = 2,46560$$

$$\hat{y}_2 (\text{Ridge}) = 1,81808 + 0,22376 \times 5 = 3,43688$$

$$\text{RMSE} = \sqrt{\frac{(0,4 - 3,11212)^2 + (0,1 - 2,46560)^2 + (2,2 - 3,43688)^2}{3}}$$

$$= 1,45285$$

For the model without Ridge we can see that the train RMSE score is lower than the test RMSE score. This indicates slight levels of overfitting.

The train RMSE for the model after applying Ridge's regularization is higher when compared with the train RMSE without Ridge, as predicted. This happens because the model is more generalized and less fit to that training data. However, the RMSE score for the test with Ridge applied is significantly lower because of the effects of this regularization.

In conclusion, the results are expected. By using Ridge's formula, we were able to reduce test variance in the model.

4)

~~Diagram of a neural network with three layers: input layer (3 nodes), hidden layer (2 nodes), and output layer (1 node). The input layer has values $x^{(1)} = [1, 1, 1]$. The hidden layer has weights $w^{(1)} = \begin{bmatrix} 0.1 & 0.1 \\ 0.1 & 0.2 \\ 0.2 & 0.1 \end{bmatrix}$ and bias $b^{(1)} = [0, 0, 1]$. The output layer has weights $w^{(2)} = \begin{bmatrix} 0.3 & 0.3 \\ 0.3 & 2.3 \end{bmatrix}$ and bias $b^{(2)} = [2.0]$. The output of the hidden layer is $z^{(1)} = [1, 1]$. The output of the output layer is $x^{(2)} = \text{softmax}(z^{(2)}) = \begin{bmatrix} 0.16149 \\ 0.30934 \\ 0.22917 \end{bmatrix}$.~~

$$\begin{aligned}
 & x^{(1)} = w^{(1)} x^{(0)} + b^{(1)} \\
 & \begin{bmatrix} 0.1 & 0.1 \\ 0.1 & 0.2 \\ 0.2 & 0.1 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix} + \begin{bmatrix} 0.1 \\ 0 \\ 0.1 \end{bmatrix} = \begin{bmatrix} 0.3 \\ 0.3 \\ 0.4 \end{bmatrix} \\
 & z^{(1)} = w^{(2)} x^{(1)} + b^{(2)} = \\
 & \begin{bmatrix} 0.3 & 0.3 \\ 0.3 & 2.3 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix} + \begin{bmatrix} 2.0 \end{bmatrix} = \begin{bmatrix} 1.7 \\ 1.3 \\ 2.0 \end{bmatrix} \\
 & x^{(2)} = \text{softmax}(z^{(2)}) = \begin{bmatrix} \frac{e^{2.0}}{e^{2.0} + e^{1.7} + e^{1.3}} \\ \frac{e^{1.7}}{e^{2.0} + e^{1.7} + e^{1.3}} \\ \frac{e^{1.3}}{e^{2.0} + e^{1.7} + e^{1.3}} \end{bmatrix} = \begin{bmatrix} 0.16149 \\ 0.30934 \\ 0.22917 \end{bmatrix}
 \end{aligned}$$

Back propagation:

$$w^{(2)} = w^{(2)} - \mu \frac{\partial E}{\partial w}$$

Last layer:

$$\frac{\partial E}{\partial w} = s^{(2)} \times (x^{(2)})^T$$

$$s^{(2)} = \frac{\partial E}{\partial x^{(2)}} = \frac{\partial E}{\partial z^{(2)}} \cdot \frac{\partial z^{(2)}}{\partial x^{(2)}} = \frac{\partial E}{\partial z^{(2)}} = x^{(2)} \cdot t$$

$$\begin{bmatrix} 0.16149 \\ 0.30934 \\ 0.22917 \end{bmatrix} - \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} = \begin{bmatrix} 0.16149 \\ -0.30934 \\ 0.22917 \end{bmatrix}$$

$$\frac{\partial C}{\partial W} = \delta^{(2)} \times (y^{(2)})^T = \begin{bmatrix} -0,98616 \\ -0,63066 \\ 0,22917 \end{bmatrix} \begin{bmatrix} 0,2 & 0,2 & 0,1 \end{bmatrix}$$

$$= \begin{bmatrix} 0,13848 & 0,13848 & 0,13848 \\ -0,20920 & -0,20920 & -0,20920 \\ 0,06875 & 0,06875 & 0,05167 \end{bmatrix}$$

$$W^{(3)} = W^{(2)} - 0,1 \frac{\partial E}{\partial W} = \begin{bmatrix} 0,98616 & 1,98616 & 1,98156 \\ 1,02072 & 2,02072 & 1,02162 \\ 0,99312 & 0,99312 & 0,99083 \end{bmatrix}$$

$$\delta^{(3)} = b^{(2)} - 0,1 \frac{\partial C}{\partial b^{(2)}} = b^{(2)} - 0,1 \delta^{(2)}$$

$$= \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} - 0,1 \begin{bmatrix} 0,46149 \\ -0,63066 \\ 0,22917 \end{bmatrix} = \begin{bmatrix} 0,53285 \\ 1,06907 \\ 0,57703 \end{bmatrix}$$

- / / -

Layer intermediária.

$$W^{(3)} = W^{(2)} - \mu \frac{\partial E}{\partial W^{(2)}} \quad \frac{\partial E}{\partial W^{(2)}} = \delta^{(2)} \times (x^{(0)})^T$$

$$\delta^{(2)} = \left(\frac{\partial z^{(2)}}{\partial x^{(0)}} \right)^T \times \delta^{(3)} \circ \frac{\partial x^{(2)}}{\partial z^{(2)}}$$

$$= (W^{(2)})^T \times \delta^{(3)}$$

$$= \begin{bmatrix} -0,02227 \\ -0,25144 \\ 0,43179 \end{bmatrix}$$

$$\frac{\partial C}{\partial W^{(1)}} = \delta^{(2)} \times (x^{(0)})^T = \begin{bmatrix} -0,02227 & -0,22227 \\ -0,25144 & -0,25144 \\ 0,43179 & 0,431795 \end{bmatrix}$$

$$W^{(1)} = W^{(1)} - 0,1 \frac{\partial E}{\partial W^{(1)}} = \begin{bmatrix} 0,10223 & 0,10220 \\ 0,12514 & 0,22514 \\ 0,15682 & 0,05682 \end{bmatrix}$$

$$b^{(1)} = b^{(1)} - 0,1 \delta^{(1)} = \begin{bmatrix} 0,1 \\ 0 \\ 0 \end{bmatrix} - 0,1 \times \begin{bmatrix} -0,02227 \\ -0,25144 \\ 0,43179 \end{bmatrix}$$

$$= \begin{bmatrix} 0,10223 \\ 0,07514 \\ 0,05682 \end{bmatrix}$$

II. Programming and critical analysis

Here we have an overview of all the imports used and how we imported the dataset, which are both common to all questions. Therefore, we will be omitting this code from the beginning of each task to avoid redundancy.

```

import pandas as pd, matplotlib.pyplot as plt, numpy as np
import seaborn as sns
import warnings

from sklearn.model_selection import train_test_split, GridSearchCV
from sklearn.neural_network import MLPRegressor
from sklearn.metrics import mean_absolute_error
from sklearn.linear_model import LinearRegression
from sklearn.model_selection import train_test_split

warnings.filterwarnings("ignore")

# Load the data
df = pd.read_csv('parkinsons.csv')
X = df.drop('target', axis=1)
y = df['target']

```

5) Performing 10 runs for each model and calculating the MAE for each one

```

mae_linear = []
mae_mlp_no_activation = []
mae_mlp = []

for i in range(1, 11):
    X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=i)

    # Linear Regression Model
    linear = LinearRegression()
    linear.fit(X_train, y_train)
    y_pred_linear = linear.predict(X_test)
    mae_linear.append(mean_absolute_error(y_test, y_pred_linear))

    # MLP Model no activation
    mlp_no_activation = MLPRegressor(hidden_layer_sizes=(10,10), activation='identity', random_state=0)
    mlp_no_activation.fit(X_train, y_train)
    y_pred_mlp_no_activation = mlp_no_activation.predict(X_test)
    mae_mlp_no_activation.append(mean_absolute_error(y_test, y_pred_mlp_no_activation))

    # MLP Model ReLU
    mlp = MLPRegressor(hidden_layer_sizes=(10,10), activation='relu', random_state=0)
    mlp.fit(X_train, y_train)
    y_pred_mlp = mlp.predict(X_test)
    mae_mlp.append(mean_absolute_error(y_test, y_pred_mlp))

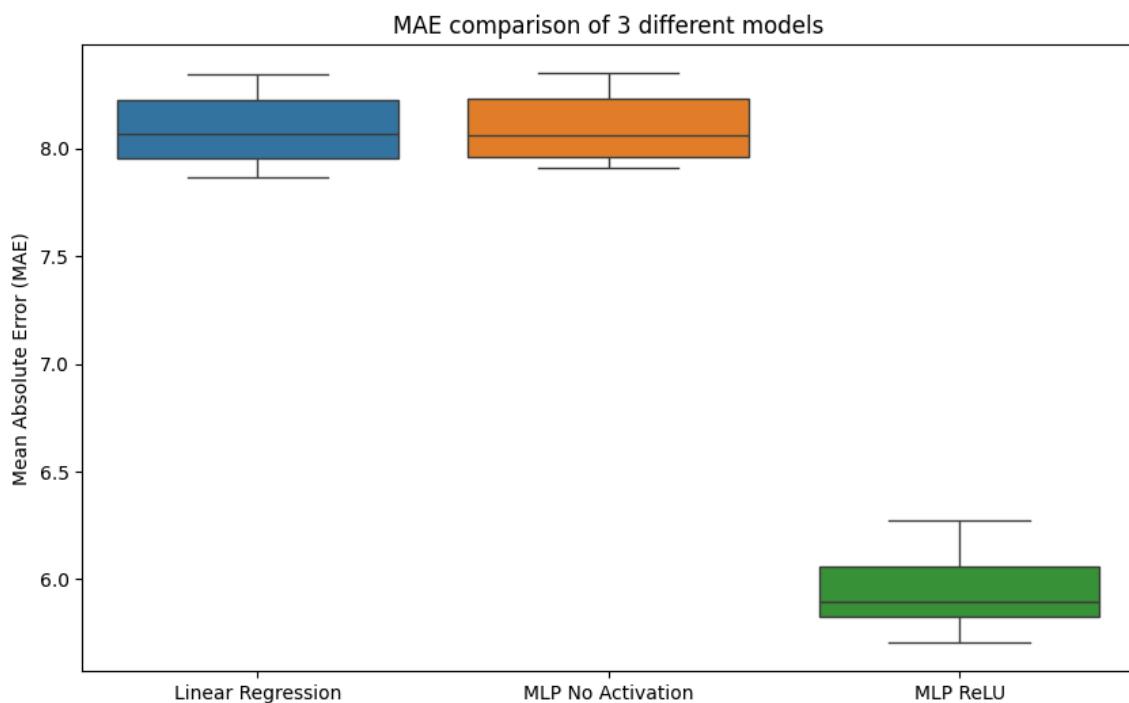
```

Plotting the data

```
boxplot_data = [mae_linear, mae_mlp_no_activation, mae_mlp]

plt.figure(figsize=(10, 6))
sns.boxplot(data=boxplot_data)
plt.xticks([0, 1, 2], ['Linear Regression', 'MLP No Activation', 'MLP ReLU'])
plt.ylabel('Mean Absolute Error (MAE)')
plt.title('MAE comparison of 3 different models')
plt.show()
```

Output:



- 6) If we analyze both the Linear Regression and the MLP with no activation boxplots, we notice them to be very similar to each other.

This similarity makes sense if we consider the nature and calculus behind both models. The Linear Regression model is a simple model that captures linear relationships between the features and the target variable. On the other hand, the MLP model with no activation essentially is a more complex linear regression, despite having two hidden layers. Since there are no non-linear activation functions in the process, the model will not be able to catch other patterns or relationships in the data other than the linear ones, and so, the result will end up very similar, except with a lot of extra steps. Because of this, both models create two very similar MAE boxplots, as indicated in the first paragraph.

We can finally understand the activation functions' importance in introducing non-linearity to the model so that it can capture non-linear and more complex relationships between the features and, therefore, achieve a higher generalization capacity and better test results. Without them, MLP's would just be linear transformations of the input.

7)

Splitting the data, creating the MLP regressor, setting up the Grid Search and assessing the best combination of parameters

```
# Split the data
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=0)

# Grid Search parameters
param_grid = {
    'alpha': [0.0001, 0.001, 0.01], #L2 penalty
    'learning_rate_init': [0.001, 0.01, 0.1], #Learning rate
    'batch_size': [32, 64, 128] #Batch size
}

# Create the MLP Regressor
mlp_regressor = MLPRegressor(hidden_layer_sizes=(10,10), random_state=0)

# Grid Search
grid_search = GridSearchCV(mlp_regressor, param_grid=param_grid,
                           scoring='neg_mean_absolute_error', cv = 5)
grid_search.fit(X_train, y_train)

# Evaluate the best model
results = pd.DataFrame(grid_search.cv_results_)
best_model = grid_search.best_estimator_
y_pred = best_model.predict(X_test)
test_mae = mean_absolute_error(y_test, y_pred)

print("Best hyperparameters found:", "\nL2_penalty:", grid_search.best_params_['alpha'], "\nLearning rate:", grid_search.best_params_['learning_rate_init'],
      "\nBatch size: ", grid_search.best_params_['batch_size'])
print("Test MAE of the best model: ", test_mae)
```

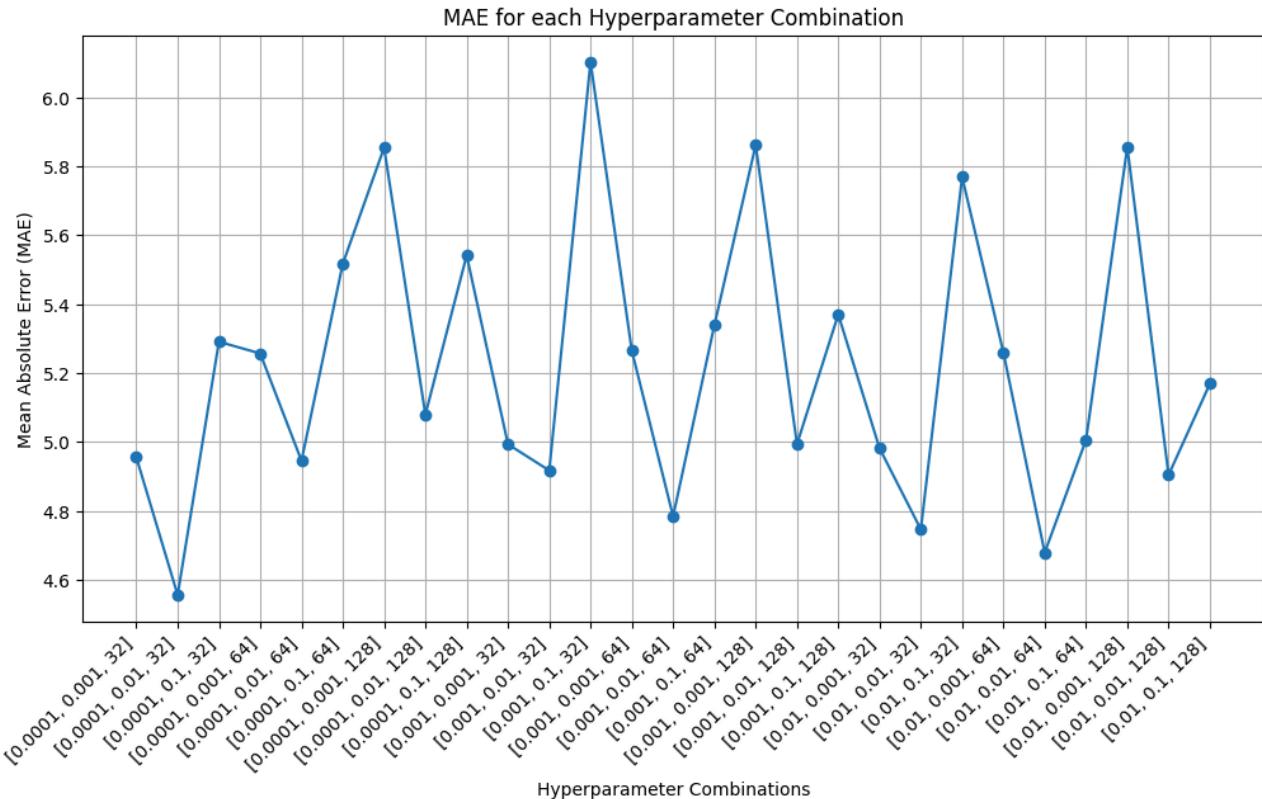
Plotting the MAE for every parameter combination

```
# Extract relevant columns
params = results['params']
mean_mae = -results['mean_test_score'] # MAE(since neg_mean_absolute_error is used)

# Labels
param_labels = [
    f"[{param['alpha']}], {param['learning_rate_init']}, {param['batch_size']}]"
    for param in params
]

# Plot the results
plt.figure(figsize=(12, 6))
plt.plot(range(len(mean_mae)), mean_mae, marker='o')
plt.xticks(range(len(mean_mae)), param_labels, rotation=45, ha = 'right')
plt.ylabel('Mean Absolute Error (MAE)')
plt.xlabel('Hyperparameter Combinations')
plt.title('MAE for each Hyperparameter Combination')
plt.grid()
plt.show()
```

Output:



The best combination of hyperparameters observed after this experience is:

$$\{\text{L2 penalty} = 0.0001; \text{Learning Rate} = 0.01; \text{Batch} = 32\}$$

As for the trade-offs of each hyperparameter:

- L2 penalty (alpha):

This parameter corresponds to the level of penalization to large coefficients we desire by applying Ridge's regularization in our model. Higher values of alpha correspond to more regularization.

Therefore, we can associate overly lower values of this parameter with overfitting and low MAE scores in the test compared to the training, as low regularization is being applied, and excessively high levels with underfitting. The balance was found at alpha = 0.0001 in this model.

- Learning Rate

This parameter represents how big are the updates the model does to its weights and biases on each iteration.

Values too low might cause the model to take too much time to converge to the optimal solution while values too high might cause overshooting and divergence in the values.

This was the hyperparameter that had the most impact on the MAE scores with Learning Rate = 0.01 being the optimal value for this model.

- Batch

This parameter is the number of instances the model has in account before making the next update to the weights and biases.

Generally lower values will be always better for the model's performance, since it means it is updating values more frequently. However they cause the training to take more time for the same reason, so the trade-off here is in balancing training time and effectiveness of the MLP.

As expected, the model had better MAE scores for the lower value of batch = 32.