

Part I: Pen and paper

1. Consider the problem of learning a regression model from 4 bivariate observations $\{(0.7, -0.3), (0.4, 0.5), (0.2, 0.8), (-0.4, 0.3)\}$ and their targets $\{0.8, 0.6, 0.3, 0.3\}$.

- (a) Given the radial basis function, $\phi_j(x) = e^{-\frac{\|x-c_j\|^2}{2}}$, that transforms the original space onto a new space characterized by the similarity of the original observations to the following data points:

$$c_1 = (0, 0)$$

$$c_2 = (1, -1)$$

$$c_3 = (-1, 1)$$

Learn the Ridge regression (l_2 regularization) using the closed solution with $\lambda = 0.1$.

We begin by applying the radial basis function and data points to our observations, to get our transformed matrix:

First we calculate our Φ matrix after the transformation

$$\Phi = \begin{bmatrix} 1 & 0,748 & 0,748 & 0,101 \\ 1 & 0,815 & 0,271 & 0,331 \\ 1 & 0,712 & 0,963 & 0,712 \\ 1 & 0,882 & 0,161 & 0,654 \end{bmatrix}$$

Loss function with regularization:

using Ridge regression

$$E(w) = \frac{1}{2} \sum_{i=1}^n (z_i - w^T \cdot x_i)^2 + \frac{\lambda}{2} \|w\|^2$$

To obtain the best weights, we minimize the loss function by taking its derivative and finding where it equals zero. Since the loss function has positive coefficients, this zero point is a minimum.

To minimize w , we need to find the $E(w)$ minimum

$$\nabla E(w) = \nabla \left(\frac{1}{2} (z - x \cdot w)^T (z - x \cdot w) + \frac{\lambda}{2} w^T w \right) = 0$$

$$\Rightarrow -x^T z + x^T \cdot x \cdot w + \lambda \cdot w = 0$$

$$\Rightarrow x^T z = (x^T \cdot x + \lambda I) \cdot w$$

$$\Rightarrow w = (x^T \cdot x + \lambda I)^{-1} \cdot x^T \cdot z$$

After having developed the above expression we get to the closed form solution for Ridge Regression.

By inputting our transformed matrix into the above expression we learn the weights:

$$w = (\Phi^T \cdot \Phi + \lambda I)^{-1} \cdot \Phi^T \cdot z = \begin{bmatrix} 0,33914 \\ 0,19945 \\ 0,40096 \\ -0,29600 \end{bmatrix}$$

(b) Compute the training RMSE for the learnt regression.

After calculating the model's prediction for each observation we can compare them with the original output using the RMSE's formula.

6-)

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

$$= \sqrt{\frac{1}{4} (0,04^2 + 0,09^2 + 0,07^2 + 0,09^2)}$$

$$= 0.06508$$

$$\hat{z}_i = w^T x_i$$

$$\hat{z}_0 = 0,75844$$

$$\hat{z}_1 = 0,51232$$

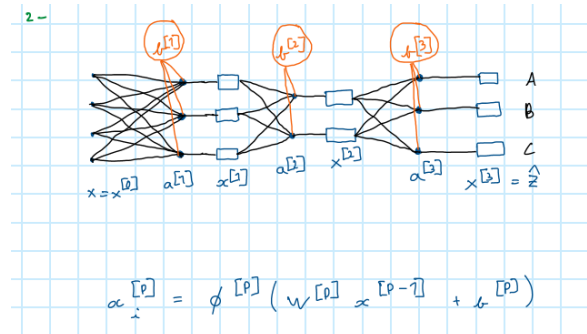
$$\hat{z}_2 = 0,30905$$

$$\hat{z}_3 = 0,38629$$

$$\hat{z}_i - z_i \approx (-0,04 \quad -0,09 \quad -0,07 \quad 0,09)$$

2. Consider a MLP classifier of three outcomes –A, B and C. Perform one batch gradient descent update (with learning rate $\eta = 0.1$) for training observations x_1 and x_2 .

Before beginning, here is how the MLP looks like and the notation we chose to tackle this problem.



For the upcoming calculations, we are going to calculate many derivatives using these rules:

$$a = wX$$

$$C = AB$$

$$\frac{\partial E}{\partial A} = \frac{\partial E}{\partial C} \cdot B^T$$

$$\frac{\partial E}{\partial B} = A^T \cdot \frac{\partial E}{\partial C}$$

$$a = wX + b$$

$$C = A + B$$

$$\frac{\partial E}{\partial A} = \frac{\partial E}{\partial B} = \frac{\partial E}{\partial C}$$

$$X = f(a)$$

$$C = f(A)$$

$$\frac{\partial E}{\partial A} = \frac{\partial E}{\partial C} \odot f'(A)$$

1. Forward Propagation

For starters, let's grab our input observations (x_1 and x_2) and propagate them across our MLP.

x_1 :

$$x_1^{[1]} = f(w^{[1]} x_1^{[0]} + b^{[1]})$$

$$= f\left(\begin{bmatrix} 5 \\ 6 \\ 5 \end{bmatrix}\right) = \begin{bmatrix} 0.46211 \\ 0.76759 \\ 0.46211 \end{bmatrix}$$

$a_1^{[1]}$

$$x_1^{[2]} = f(w^{[2]} x_1^{[1]} + b^{[2]})$$

$$= f\left(\begin{bmatrix} 4.97058 \\ 2.68581 \end{bmatrix}\right) = \begin{bmatrix} 0.45047 \\ -0.57642 \end{bmatrix}$$

$a_1^{[2]}$

$$\hat{x}_1 = x_1^{[3]} = f(w^{[3]} x_1^{[2]} + b^{[3]})$$

$$= f\left(\begin{bmatrix} 0.87405 \\ 1.77499 \\ 0.87405 \end{bmatrix}\right) = \begin{bmatrix} -0.91590 \\ -0.80494 \\ -0.91590 \end{bmatrix}$$

$a_1^{[3]}$

x_2 :

$$x_2^{[1]} = f(w^{[1]} x_2^{[0]} + b^{[1]})$$

$$= f\left(\begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}\right) = \begin{bmatrix} -0.90515 \\ -0.90515 \\ -0.90515 \end{bmatrix}$$

$a_2^{[1]}$

$$x_2^{[2]} = f(w^{[2]} x_2^{[1]} + b^{[2]})$$

$$= f\left(\begin{bmatrix} -4.43089 \\ -1.71544 \end{bmatrix}\right) = \begin{bmatrix} -0.99956 \\ -0.99343 \end{bmatrix}$$

$a_2^{[2]}$

$$\hat{x}_2 = x_2^{[3]} = f(w^{[3]} x_2^{[2]} + b^{[3]})$$

$$= f\left(\begin{bmatrix} -0.99300 \\ -2.99212 \\ -0.99300 \end{bmatrix}\right) = \begin{bmatrix} -0.98652 \\ -0.99816 \\ -0.98652 \end{bmatrix}$$

$a_2^{[3]}$

2. Back Propagation

Before we start this step, here are some auxiliary calculations/derivatives that will help us along the way.

$$\begin{aligned}\phi(x) &= \tanh(0.5x - 2) \\ \phi'(x) &= \frac{1}{2} \operatorname{sech}^2(0.5x - 2) = \frac{1}{2} (1 - \tanh^2(0.5x - 2)) \\ &\Downarrow \\ \phi'(a^{[3]}) &= \frac{1}{2} (1 - x^{[3]^2}) \\ E(w) &= \frac{1}{2} (z - \hat{z})^2 \\ \frac{dE}{dz} &= \frac{1}{2} \times \frac{d}{dz} (z - \hat{z})^2 \\ &= \frac{1}{2} \times 2 (z - \hat{z}) \times (-1) \\ &= \hat{z} - z\end{aligned}$$

With that done, let's begin by finding the expressions required to calculate the gradient for the $W^{[3]}$ and the $b^{[3]}$ matrix.

Back Propagation

Layer 3

$$\begin{aligned}\frac{dE}{dw_i^{[3]}} &= \frac{dE}{da_i^{[3]}} \cdot x_i^{[2]T} \\ \delta_i^{[3]} &= \frac{dE}{dx_i^{[3]}} \odot \phi'(a_i^{[3]}) \\ \frac{dE}{dx_i^{[3]}} &= \frac{dE}{d\hat{z}_i} = \hat{z}_i - z_i \\ \frac{dE}{dw_i^{[3]}} &= (\hat{z}_i - z_i) \odot \phi'(a_i^{[3]}) \cdot x_i^{[2]T} \\ \frac{dE}{da_i^{[3]}} &= \frac{dE}{d\hat{z}_i} = \delta_i^{[3]}\end{aligned}$$

$$\begin{aligned}\phi(x) &= \tanh(0.5x - 2) \\ \phi'(x) &= \frac{1}{2} \operatorname{sech}^2(0.5x - 2) = \frac{1}{2} (1 - \tanh^2(0.5x - 2)) \\ &\Downarrow \\ \phi'(a^{[3]}) &= \frac{1}{2} (1 - x^{[3]^2}) \\ E(w) &= \frac{1}{2} (z - \hat{z})^2 \\ \frac{dE}{dz} &= \frac{1}{2} \times \frac{d}{dz} (z - \hat{z})^2 \\ &= \frac{1}{2} \times 2 (z - \hat{z}) \times (-1) \\ &= \hat{z} - z\end{aligned}$$

Since we are training with a single-batch gradient descent update, we will compute the gradients for

all our observations and then sum these gradients to determine the resulting adjustment. (Although not shown, we have also stored the delta values for future uses).

$$\frac{dE}{dw_1^{[3]}} = \begin{bmatrix} 0,00305 & -0,00391 \\ -0,14313 & 0,15315 \\ 0,00305 & -0,00391 \end{bmatrix} \quad \frac{dE}{dw_2^{[3]}} = \begin{bmatrix} 0,02653 & 0,02621 \\ 0,000003 & 0,000003 \\ 0,00018 & 0,00018 \end{bmatrix}$$

$$\text{therefore } \frac{dE}{dw^{[3]}} = \begin{bmatrix} 0,02964 & 0,02252 \\ -0,14314 & 0,18315 \\ 0,00287 & -0,0408 \end{bmatrix}$$

$$\frac{dE}{dt^{[3]}} = \delta_1^{[3]} + \delta_2^{[3]} = \begin{bmatrix} -0,01932 \\ -0,31773 \\ 0,00696 \end{bmatrix}$$

Afterwards, we proceed with layer 2, the logic being the same. Notice, how we know make use of the delta already calculated in the previous layer.

Layer 2

$$\frac{dE}{dw_1^{[2]}} = \frac{dE}{da_1^{[2]}} \cdot x_1^{[2]T}$$

$$\delta_1^{[2]} = \frac{dE}{da_1^{[2]}} \odot \phi'(a_1^{[2]})$$

$$\frac{dE}{dx_1^{[2]}} = w^{[2]T} \cdot \frac{dE}{da_1^{[2]}} = \delta_1^{[2]}$$

$$\bullet \frac{dE}{dw_2^{[2]}} = w^{[2]T} \cdot \delta_1^{[2]} \odot \phi'(a_2^{[2]}) \cdot x_2^{[2]T}$$

$$\bullet \frac{dE}{dt_1^{[2]}} = \frac{dE}{da_1^{[2]}} = \delta_1^{[2]}$$

therefore

$$\frac{dE}{dw^{[2]}} = \frac{dE}{dw_1^{[2]}} + \frac{dE}{dw_2^{[2]}} = \begin{bmatrix} -0,17304 & -0,27519 & -0,17304 \\ -0,04678 & -0,07719 & -0,04678 \end{bmatrix}$$

$$\frac{dE}{dt^{[2]}} = \frac{dE}{dt_1^{[2]}} + \frac{dE}{dt_2^{[2]}} = \begin{bmatrix} -0,37449 \\ -0,10773 \end{bmatrix}$$

And finally, the adjustments for layer 1.

Layer 1

$$\frac{dE}{dw_2[1]} = \left(\frac{dE}{da_2[1]} \right) \cdot x_1^T$$

$$\delta_2[1] = \left(\frac{dE}{da_2[1]} \right) \odot \phi'(a_2[1])$$

$$\frac{dE}{dx_1[1]} = w_2[1]^T \cdot \left(\frac{dE}{da_2[1]} \right) = \delta_2[1]$$

- $\frac{dE}{dw_1[1]} = w_1[1]^T \cdot \delta_1[1] \odot \phi'(a_1[1]) \cdot x_1^T$
- $\frac{dE}{da_2[1]} = \frac{dE}{da_2[1]} = \delta_2[1]$

therefore

$$\frac{dE}{dw[1]} = \frac{dE}{dw_1[1]} + \frac{dE}{dw_2[1]} = \begin{bmatrix} -0,18721 & -0,18721 & -0,18721 & -0,18721 \\ -0,33589 & -0,33589 & -0,33589 & -0,33589 \\ -0,18721 & -0,18721 & -0,18721 & -0,18721 \end{bmatrix}$$

$$\frac{dE}{db[1]} = \frac{dE}{db_1[1]} + \frac{dE}{db_2[1]} = \begin{bmatrix} -0,18721 \\ -0,33589 \\ -0,18721 \end{bmatrix}$$

3. Updating the Weights/Biases

With all the previous steps done, we can obtain the new weights/biases by multiplying the computed adjustments with the learning rate and adding the result to the old value of the weights/biases.

updates

$$w[3] = w[3] - \eta \frac{dE}{dw[3]} = \begin{bmatrix} 0,99704 & 0,99775 \\ 3,07431 & 0,98169 \\ 0,99971 & 1,00041 \end{bmatrix}$$

$$w[2] = w[2] - \eta \frac{dE}{dw[2]} = \begin{bmatrix} 1,01730 & 4,02852 & 1,01730 \\ 1,00468 & 1,00772 & 1,00468 \end{bmatrix}$$

$$w[1] = w[1] - \eta \frac{dE}{dw[1]} = \begin{bmatrix} 1,01872 & 1,01872 & 1,01872 & 1,01872 \\ 1,03359 & 1,03359 & 1,03359 & 1,03359 \\ 1,01872 & 1,01872 & 1,01872 & 1,01872 \end{bmatrix}$$

$$b[3] = b[3] - \eta \frac{dE}{db[3]} = \begin{bmatrix} 1,00193 \\ 1,03177 \\ 0,99930 \end{bmatrix}$$

$$b[2] = b[2] - \eta \frac{dE}{db[2]} = \begin{bmatrix} 1,03745 \\ 1,01077 \end{bmatrix}$$

$$b[1] = b[1] - \eta \frac{dE}{db[1]} = \begin{bmatrix} 1,01872 \\ 1,03359 \\ 1,01872 \end{bmatrix}$$

Part II: Programming

In the next image, we have provided an overview of all the imports utilized and our approach to importing the dataset. Specifically, 'X' stores all of the input variables values, and 'y' stores the Output Class values (Wine quality). To streamline the code, we've omitted this information from the beginning of each section to avoid redundancy.

```
import matplotlib.pyplot as plt
import pandas as pd
import numpy as np
import warnings

from sklearn.model_selection import train_test_split
from sklearn.neural_network import MLPRegressor
from sklearn.metrics import mean_absolute_error
from sklearn.metrics import mean_squared_error
from math import sqrt

warnings.filterwarnings("ignore")

data = pd.read_csv('../data/winequality-red.csv', delimiter=';')

# Split the data into features (X) and target (y)
X = data.drop(columns=['quality'])
y = data['quality']

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

1. Consider the *winequality – red.csv* dataset where the goal is to estimate the quality of a wine based on physicochemical inputs. Using a 80-20 training-test split with a fixed seed (`random_state = 0`), you are asked to learn MLP regressors to answer the following questions.
 - (a) Learn a MLP regressor with 2 hidden layers of size 10, reLU activation and early stopping with 20% of training data set aside for validation. Plot the distribution of the residues (in absolute value) using a histogram.

Calculating the residuals:

```
# Initialize lists to store residuals from each run
residuals = []

for run in range(1, 11):

    # Create and train the MLP regressor
    mlp = MLPRegressor(hidden_layer_sizes=(10, 10), activation='relu', validation_fraction=0.2, random_state=run, early_stopping=True)
    mlp.fit(X_train, y_train)

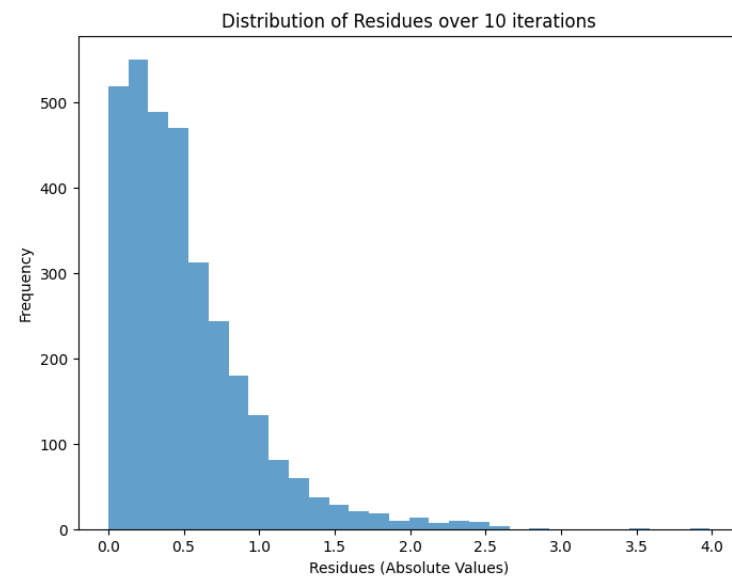
    # Make predictions on the test set
    y_pred = mlp.predict(X_test)

    # Calculate the absolute residuals
    residuals.append(np.abs(y_test - y_pred))
```

Plotting them:

```
# Plot a histogram of the residuals
plt.figure(figsize=(8, 6))
plt.hist(np.concatenate(residuals), bins=30, alpha=0.7)
plt.title('Distribution of Residues over 10 iterations')
plt.xlabel('Residues (Absolute Values)')
plt.ylabel('Frequency')
plt.show()
```

Output:



- (b) Assess the impact of rounding and bounding on the MAE of the MLP learnt in the previous question.

Calculating the Mean Absolute Error after rounding and bounding the predictions:

```
original_mae = []
rounded_bounded_mae = []

for run in range(1, 11):

    # Create and train the MLP regressor
    mlp = MLPRegressor(hidden_layer_sizes=(10, 10), activation="relu", validation_fraction=0.2, random_state=run, early_stopping=True)
    mlp.fit(X_train, y_train)

    # Make predictions on the test set
    y_pred = mlp.predict(X_test)

    # Round and bound the predictions
    y_pred_rounded_bounded = np.clip(np.round(y_pred), 1, 10)

    mae_original = mean_absolute_error(y_test, y_pred)
    mae_rounded_bounded = mean_absolute_error(y_test, y_pred_rounded_bounded)

    original_mae.append(mae_original)
    rounded_bounded_mae.append(mae_rounded_bounded)
```

Plotting:

```
# Calculate the average MAE for both original and rounded/bounded predictions
average_original_mae = np.mean(original_mae)
average_rounded_bounded_mae = np.mean(rounded_bounded_mae)

print("Average Original MAE:", average_original_mae)
print("Average Rounded/Bounded MAE:", average_rounded_bounded_mae)
```

Results:

MLP Model	Average MAE
Original MAE	0.50972
Rounded/Bounded MAE	0.43875

Assessment: The model using rounding/bounding is more accurate as we expected. This difference in performance can be attributed to the fact that bounding restricts predictions to values within the valid domain, while rounding enhances the model's accuracy by aligning predictions with the discrete nature of the target variable.

- (c) Assess the impact on RMSE from replacing early stopping by a well-defined number of iterations in $\{20, 50, 100, 200\}$.

Calculating the averaged RMSEs for each max iteration.

```
rmse_values = []

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=0)
iteration_values = [20, 50, 100, 200]

for iteration in iteration_values:
    rmse_scores = []

    for run in range(1, 11):
        mlp = MLPRegressor(hidden_layer_sizes=(10, 10), activation='relu', max_iter=iteration, random_state=run)

        mlp.fit(X_train, y_train)

        prediction = mlp.predict(X_test)

        rmse = sqrt(mean_squared_error(y_test, prediction))
        rmse_scores.append(rmse)

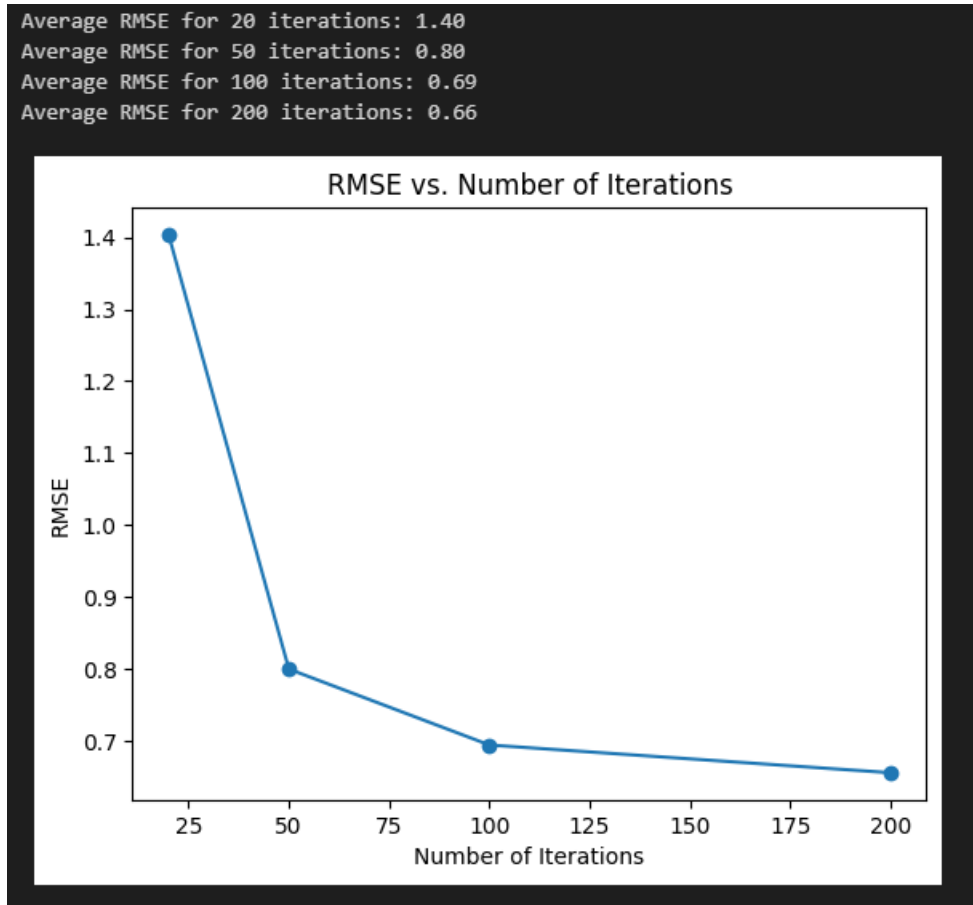
    average_rmse = np.mean(rmse_scores)
    rmse_values.append(average_rmse)
```

Printing and plotting the averaged RMSEs for each iteration.

```
for i, num_iterations in enumerate(iteration_values):
    print(f"Average RMSE for {num_iterations} iterations: {rmse_values[i]:.2f}")

plt.plot(iteration_values, rmse_values, marker='o')
plt.title('RMSE vs. Number of Iterations')
plt.xlabel('Number of Iterations')
plt.ylabel('RMSE')
plt.show()
```

Output:



Comparison:

MLP Model	Averaged RMSE
Early Stopping	0.67
200 Iterations	0.66
100 Iterations	0.69
50 Iterations	0.80
20 Iterations	1.40

- (d) Critically comment the results obtained in the previous question, hypothesizing at least one reason why early stopping favors and/or worsens performance.

Early stopping uses a small validation set to evaluate the model's accuracy. The moment the model's performance starts to deteriorate, early stopping terminates training. This way, the model generalizes better by not memorizing the entire dataset, helping reducing overfitting.

In question 3, we arrived at the same conclusion. When replacing early stopping with a fixed number of iterations, it's evident that the model's performance improves with the increase in the number of iterations. With too few iterations, the model tends to be underfitted, while with more iterations, it becomes more precise (until it starts overfitting).

We also calculated the average Root Mean Square Error (RMSE) for early stopping and compared it to the model with a fixed number of iterations. Interestingly, the early stopping approach resulted in a slightly larger average RMSE (0.67 vs. 0.66) when compared to the model with 200 iterations. This suggests that while early stopping generally enhances performance, there might be cases where the model achieves greater efficiency with a fixed number of iterations.