

Homework 3

13 de outubro de 2023 10:52

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

- a. [4v] Given the radial basis function, $\phi_j(x) = \exp\left(-\frac{\|x - c_j\|^2}{2}\right)$, that transforms the original space onto a new space characterized by the similarity of the original observations to the following data points, $\{c_1 = \begin{pmatrix} 0 \\ 0 \end{pmatrix}, c_2 = \begin{pmatrix} 1 \\ -1 \end{pmatrix}, c_3 = \begin{pmatrix} -1 \\ 1 \end{pmatrix}\}$. Learn the Ridge regression (L_2 regularization) using the closed solution with $\lambda = 0.1$.

$$\begin{array}{|c|c|c|c|} \hline & y_1 & y_2 & z \\ \hline x_1 & 0.7 & -0.3 & 0.8 \\ \hline x_2 & 0.4 & 0.5 & 0.6 \\ \hline x_3 & -0.2 & 0.8 & 0.3 \\ \hline x_4 & -0.4 & 0.3 & 0.3 \\ \hline \end{array} \quad \xrightarrow{\phi} \quad \begin{array}{|c|c|c|c|} \hline & y_1 & y_2 & y_3 & z \\ \hline x_1' & e^{-0.29} & e^{-0.29} & e^{-2.29} & 0.8 \\ \hline x_2' & e^{0.205} & e^{1.305} & e^{1.105} & 0.6 \\ \hline x_3' & e^{-0.34} & e^{-2.34} & e^{-0.34} & 0.3 \\ \hline x_4' & e^{-0.125} & e^{-1.825} & e^{-0.425} & 0.3 \\ \hline \end{array}$$

Regressão de Ridge: $w = (X^T X + \lambda I)^{-1} X^T z$, $\lambda = 0.1$

$$w = \left(\begin{bmatrix} 1 & 1 & 1 & 1 \\ e^{-0.29} & e^{0.205} & e^{0.34} & e^{-0.125} \\ e^{-0.29} & e^{1.305} & e^{2.34} & e^{1.825} \\ e^{-2.29} & e^{1.105} & e^{-0.34} & e^{-0.425} \end{bmatrix} \begin{bmatrix} 1 & e^{-0.29} & e^{0.29} & e^{-2.29} \\ 1 & e^{0.205} & e^{1.305} & e^{1.105} \\ 1 & e^{-0.34} & e^{-2.34} & e^{-0.34} \\ 1 & e^{-0.125} & e^{-1.825} & e^{-0.425} \end{bmatrix} + \begin{bmatrix} 0.1 & 0 & 0 & 0 \\ 0 & 0.1 & 0 & 0 \\ 0 & 0 & 0.1 & 0 \\ 0 & 0 & 0 & 0.1 \end{bmatrix} \right)^{-1} X^T z$$

$$= \begin{bmatrix} 4.10000 & 3.157178 & 1.276981 & 1.798017 \\ 3.157178 & 2.68966 & 0.991646 & 1.429161 \\ 1.276981 & 0.991646 & 0.768703 & 0.339552 \\ 1.798017 & 1.429161 & 0.339552 & 1.153987 \end{bmatrix} \begin{bmatrix} 1 & e^{-0.29} & e^{0.205} & e^{0.34} & e^{-0.125} \\ e^{-0.29} & e^{1.305} & e^{2.34} & e^{1.825} & e^{-0.425} \\ e^{-2.29} & e^{1.105} & e^{-0.34} & e^{-0.425} & e^{-0.29} \\ e^{0.205} & e^{1.305} & e^{2.34} & e^{1.825} & e^{0.29} \end{bmatrix} \begin{bmatrix} 0.8 \\ 0.6 \\ 0.3 \\ 0.3 \end{bmatrix}$$

$$= \begin{bmatrix} 0.33914 \\ 0.19945 \\ 0.46096 \\ -0.29600 \end{bmatrix} w_0 \\ w_1 \\ w_2 \\ w_3 \quad \underline{f(x)} = 0.33914 + 0.19945 x_1 + 0.46096 x_2 - 0.296 x_3$$

- b. [2v] Compute the training RMSE for the learnt regression.

$$\hat{z} = x^T w = \begin{bmatrix} 1 & e^{-0.29} & e^{0.205} & e^{0.34} & e^{-0.125} \\ 1 & e^{0.205} & e^{1.305} & e^{1.105} & e^{-0.34} \\ 1 & e^{0.34} & e^{2.34} & e^{0.34} & e^{-0.425} \\ 1 & e^{-0.125} & e^{-1.825} & e^{-0.425} & e^{0.29} \end{bmatrix} \begin{bmatrix} 0.33914 \\ 0.19945 \\ 0.46096 \\ -0.29600 \end{bmatrix} = \begin{bmatrix} 0.75844 \\ 0.51232 \\ 0.30905 \\ 0.38629 \end{bmatrix}$$

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

$$= \sqrt{\frac{1}{4} ((0.8 - 0.75844)^2 + (0.6 - 0.51232)^2 + (0.3 - 0.30905)^2 + (0.3 - 0.38629)^2)}$$

$$\approx 0.06508 //$$

2. [6v] Consider a MLP classifier of three outcomes – A, B and C – characterized by the weights,

$$W^{[1]} = \begin{pmatrix} 1 & 1 & 1 & 1 \\ 1 & 1 & 2 & 1 \\ 1 & 1 & 1 & 1 \end{pmatrix}, b^{[1]} = \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix}, W^{[2]} = \begin{pmatrix} 1 & 4 & 1 \\ 1 & 1 & 1 \end{pmatrix}, b^{[2]} = \begin{pmatrix} 1 \\ 1 \end{pmatrix}, W^{[3]} = \begin{pmatrix} 1 & 1 \\ 3 & 1 \\ 1 & 1 \end{pmatrix}, b^{[3]} = \begin{pmatrix} 1 \\ 1 \end{pmatrix}$$

the activation $f(x) = \frac{e^{0.5x-2}-e^{-0.5x+2}}{e^{0.5x-2}+e^{-0.5x+2}} = \tanh(0.5x-2)$ for every unit, and squared error loss

$\frac{1}{2} \|z - \hat{z}\|_2^2$. Perform one batch gradient descent update (with learning rate $\eta = 0.1$) for training observations $x_1 = \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix}$ and $x_2 = \begin{pmatrix} 1 \\ 0 \\ -1 \end{pmatrix}$ with targets B and A, respectively.

$$A = \begin{pmatrix} 1 \\ -1 \\ 1 \end{pmatrix} \quad B = \begin{pmatrix} -1 \\ 1 \\ 1 \end{pmatrix} \quad E(w) = \frac{1}{2} \sum_{i=1}^n (t_i - x_i)^2, \quad \kappa = \phi(z)$$

C.A:

$$\frac{\partial E}{\partial w^{[3]}} = \frac{\partial E}{\partial x^{[3]}} \frac{\partial x^{[3]}}{\partial z^{[3]}} \frac{\partial z^{[3]}}{\partial w^{[3]}} = \delta^{[3]} \frac{\partial z^{[3]}}{\partial w^{[3]}} ; \frac{\partial E}{\partial b^{[3]}} = \delta^{[3]}$$

$$\delta^{[p]} = \left(\frac{\partial z^{[p+1]}}{\partial x^{[p]}} \right)^T \delta^{[p+1]} \frac{\partial x^{[p+1]}}{\partial z^{[p]}}$$

$$\frac{\partial E}{\partial x^{[3]}} = -\sum_{i=1}^3 (t_i - z_i)$$

$$\frac{\partial z^{[p]}}{\partial x^{[p]}} = \frac{\partial f(z^{[p]})}{\partial z^{[p]}} = 0.5 \cdot \text{sech}^2(0.5 z^{[p]} - 2)$$

$$\frac{\partial z^{[p]}}{\partial w^{[p]}} = w^{[p]}$$

$$\frac{\partial z^{[p]}}{\partial b^{[p]}} = 1$$

$$\frac{\partial z^{[p]}}{\partial x^{[p-1]}} = w^{[p-1]}$$

Forward propagation

$$z^{[1][1]} = W^{[1]} x^{[0][0]} + b^{[1]} = \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & 1 & 2 & 1 \\ 1 & 1 & 1 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix} + \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} = \begin{bmatrix} \xi \\ \eta \\ \zeta \end{bmatrix}$$

$$x^{[1][0]} = \tanh(0.5 \begin{bmatrix} \xi \\ \eta \\ \zeta \end{bmatrix} - 2) = \begin{bmatrix} 0.462117 \\ 0.462117 \\ 0.462117 \end{bmatrix}$$

$$z^{[2][1]} = W^{[2]} x^{[1][0]} + b^{[2]} = \begin{bmatrix} 1 & 4 & 1 \\ 1 & 1 & 1 \end{bmatrix} \begin{bmatrix} 0.462117 \\ 0.462117 \\ 0.462117 \end{bmatrix} + \begin{bmatrix} 1 \\ 1 \end{bmatrix} = \begin{bmatrix} 4.970611 \\ 2.685828 \end{bmatrix}$$

$$x^{[2][1]} = \tanh(0.5 \begin{bmatrix} 4.970611 \\ 2.685828 \end{bmatrix} - 2) = \begin{bmatrix} 0.450483 \\ 0.450483 \end{bmatrix}$$

$$z^{[3][1]} = W^{[3]} x^{[2][1]} + b^{[3]} = \begin{bmatrix} 1 & 1 \\ 3 & 1 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} 0.450483 \\ 0.450483 \end{bmatrix} + \begin{bmatrix} 1 \\ 1 \end{bmatrix} = \begin{bmatrix} 0.824062 \\ 1.375027 \\ 0.824062 \end{bmatrix}$$

$$x^{[3][1]} = \tanh(0.5 \begin{bmatrix} 0.824062 \\ 1.375027 \\ 0.824062 \end{bmatrix} - 2) = \begin{bmatrix} -0.915900 \\ -0.804940 \\ -0.915900 \end{bmatrix}$$

$$z^{[3][2]} = W^{[3]} x^{[2][1]} + b^{[3]} = \begin{bmatrix} 1 & 1 \\ 3 & 1 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} 0.450483 \\ 0.450483 \end{bmatrix} + \begin{bmatrix} 1 \\ 1 \end{bmatrix} = \begin{bmatrix} -0.999564 \\ -0.999564 \\ -0.999564 \end{bmatrix}$$

$$x^{[3][2]} = \tanh(0.5 \begin{bmatrix} -0.999564 \\ -0.999564 \\ -0.999564 \end{bmatrix} - 2) = \begin{bmatrix} -0.999564 \\ -0.999564 \\ -0.999564 \end{bmatrix}$$

$$z^{[3][3]} = W^{[3]} x^{[2][1]} + b^{[3]} = \begin{bmatrix} 1 & 1 \\ 3 & 1 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} -0.999564 \\ -0.999564 \\ -0.999564 \end{bmatrix} + \begin{bmatrix} 1 \\ 1 \end{bmatrix} = \begin{bmatrix} -0.999564 \\ -0.999564 \\ -0.999564 \end{bmatrix}$$

$$x^{[3][3]} = \tanh(0.5 \begin{bmatrix} -0.999564 \\ -0.999564 \\ -0.999564 \end{bmatrix} - 2) = \begin{bmatrix} -0.999564 \\ -0.999564 \\ -0.999564 \end{bmatrix}$$

backwards propagation

$$\delta^{[3][1]} = \frac{\partial E}{\partial x^{[3]}} \frac{\partial x^{[3]}}{\partial z^{[3]}} = (x^{[3][1]} - t^{[1]}) \circ 0.5 \cdot \text{sech}^2(0.5 z^{[3]} - 2)$$

$$= \left(\begin{bmatrix} -0.915900 \\ -0.804940 \\ -0.915900 \end{bmatrix} - \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} \right) \circ \begin{bmatrix} 0.080563 \\ 0.176036 \\ 0.080563 \end{bmatrix} = \begin{bmatrix} 0.006775 \\ 0.317735 \\ 0.006775 \end{bmatrix}$$

$$\delta^{[2][1]} = \left(\frac{\partial z^{[2]}}{\partial x^{[1]}} \right)^T \delta^{[3][1]} \frac{\partial x^{[3]}}{\partial z^{[2]}} = W^{[3]T} \delta^{[3][1]} \frac{\partial f(z^{[2]})}{\partial z^{[2]}}$$

$$= \begin{bmatrix} 1 & 3 & 1 \\ 1 & 1 & 1 \end{bmatrix} \cdot \begin{bmatrix} 0.006775 \\ 0.317735 \\ 0.006775 \end{bmatrix} \cdot 0.5 \cdot \text{sech}^2(0.5 z^{[2]} - 2)$$

$$= \begin{bmatrix} -0.374482 \\ -0.905148 \\ -0.905148 \end{bmatrix}$$

$$\delta^{[2][2]} = \left(\frac{\partial z^{[2]}}{\partial x^{[1]}} \right)^T \delta^{[3][1]} \frac{\partial x^{[3]}}{\partial z^{[2]}} = W^{[3]T} \delta^{[3][1]} \frac{\partial f(z^{[2]})}{\partial z^{[2]}}$$

$$= \begin{bmatrix} 1 & 4 & 1 \\ 1 & 1 & 1 \end{bmatrix} \cdot \begin{bmatrix} -0.374482 \\ -0.905148 \\ -0.905148 \end{bmatrix} \cdot 0.5 \cdot \text{sech}^2(0.5 z^{[2]} - 2)$$

$$= \begin{bmatrix} -0.999564 \\ -0.999564 \\ -0.999564 \end{bmatrix}$$

$$\delta^{[1][1]} = \left(\frac{\partial z^{[1]}}{\partial x^{[0]}} \right)^T \delta^{[2][1]} \frac{\partial x^{[2]}}{\partial z^{[1]}} = W^{[2]T} \delta^{[2][1]} \frac{\partial f(z^{[1]})}{\partial z^{[1]}}$$

$$= \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix} \cdot \begin{bmatrix} -0.999564 \\ -0.999564 \\ -0.999564 \end{bmatrix} \cdot 0.5 \cdot \text{sech}^2(0.5 z^{[1]} - 2)$$

$$= \begin{bmatrix} -0.999564 \\ -0.999564 \\ -0.999564 \end{bmatrix}$$

$$\delta^{[1][2]} = \left(\frac{\partial z^{[1]}}{\partial x^{[0]}} \right)^T \delta^{[2][2]} \frac{\partial x^{[2]}}{\partial z^{[1]}} = W^{[2]T} \delta^{[2][2]} \frac{\partial f(z^{[1]})}{\partial z^{[1]}}$$

$$= \begin{bmatrix} 1 & 1 & 1 \\ 3 & 1 & 1 \end{bmatrix} \cdot \begin{bmatrix} -0.999564 \\ -0.999564 \\ -0.999564 \end{bmatrix} \cdot 0.5 \cdot \text{sech}^2(0.5 z^{[1]} - 2)$$

$$= \begin{bmatrix} -0.999564 \\ -0.999564 \\ -0.999564 \end{bmatrix}$$

$$\delta^{[1][3]} = \left(\frac{\partial z^{[1]}}{\partial x^{[0]}} \right)^T \delta^{[2][3]} \frac{\partial x^{[2]}}{\partial z^{[1]}} = W^{[2]T} \delta^{[2][3]} \frac{\partial f(z^{[1]})}{\partial z^{[1]}}$$

$$= \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix} \cdot \begin{bmatrix} -0.999564 \\ -0.$$

I. Pen-and-paper

I.

Cálculo da Regressão de Ridge

```
In [ ]: import numpy as np
import math

e = math.e

X = np.array([[1, e ** -0.29, e ** -0.29, e ** -2.29],
              [1, e ** -0.205, e ** -1.305, e ** -1.105],
              [1, e ** -0.34, e ** -2.34, e ** -0.34],
              [1, e ** -0.125, e ** -1.825, e ** -0.425]])

Z = np.array([0.8, 0.6, 0.3, 0.3]).T

W = np.linalg.pinv(X.T @ X + 0.1 * np.identity(4)) @ X.T @ Z

Z_pred = X @ W
```

II.

Cálculo do Delta 3

```
In [ ]: X = np.array([-0.98652085, -0.9981635, -0.98652085])
T = np.array([1, -1, -1])
Z = np.array([-0.99299631, -2.99212439, -0.99299631])

delta3 = (X - T) * (0.5 * (1 - np.tanh((0.5 * Z) - 2) ** 2))
```

Cálculo dos Deltas 1 e 2

```
In [ ]: W_T = np.array([[1, 1],
                      [4, 1],
                      [1, 1]])
delta_prev = np.array([[ -1.15092543e-05],
                      [-1.72899256e-04]])
Z = np.array([[1],
              [1],
              [1]])

delta = np.multiply(np.dot(W_T, delta_prev), 0.5 * (1 - np.tanh(0.5 * Z - 2) ** 2))
```

Cálculo dos dE / dW

```
In [ ]: delta_X1 = np.array([[0.006775],
                           [-0.317735],
                           [0.006775]])
X1 = np.array([[0.450483, -0.576421]])
delta_X2 = np.array([[ -2.659614e-02],
```

```

[3.369628e-06],
[1.804629e-04]])
X2 = np.array([[ -0.999564, -0.993432]])

dEdW = delta_X1 @ X1 + delta_X2 @ X2

```

Cálculo dos pesos atualizados

```

In [ ]: W = np.array([[1, 1, 1, 1],
                     [1, 1, 2, 1],
                     [1, 1, 1, 1]])
dEdW = np.array([[ -0.18720683, -0.18719017, -0.18719017, -0.18717351],
                 [-0.33589135, -0.33587157, -0.33587157, -0.33585179],
                 [-0.18720683, -0.18719017, -0.18719017, -0.18717351]])

W_new = W - 0.1 * dEdW

```

Cálculo dos biases atualizados

```

In [ ]: B = np.array([[1],
                     [1],
                     [1]])
delta_X1 = np.array([[ -0.187190],
                     [-0.335872],
                     [-0.187190]])
delta_X2 = np.array([[ -1.6661921e-05],
                     [-1.9781619e-05],
                     [-1.6661921e-05]])

B_new = B - 0.1 * (delta_X1 + delta_X2)

```

II. Programming and critical analysis

Consider the winequality-red.csv dataset (available at the webpage) where the goal is to estimate the quality (sensory appreciation) 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. Given their stochastic behavior, average the performance of each MLP from 10 runs (for reproducibility consider seeding the MLPs with random_state $\in \{1..10\}$).

- 1. [3.5v] Learn a MLP regressor with 2 hidden layers of size 10, rectifier linear unit activation on all nodes, and early stopping with 20% of training data set aside for validation. All remaining parameters (e.g., loss, batch size, regularization term, solver) should be set as default. Plot the distribution of the residues (in absolute value) using a histogram.**

```
In [ ]: import pandas as pd
from sklearn.neural_network import MLPRegressor
from sklearn.model_selection import train_test_split
from sklearn.metrics import mean_absolute_error
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns

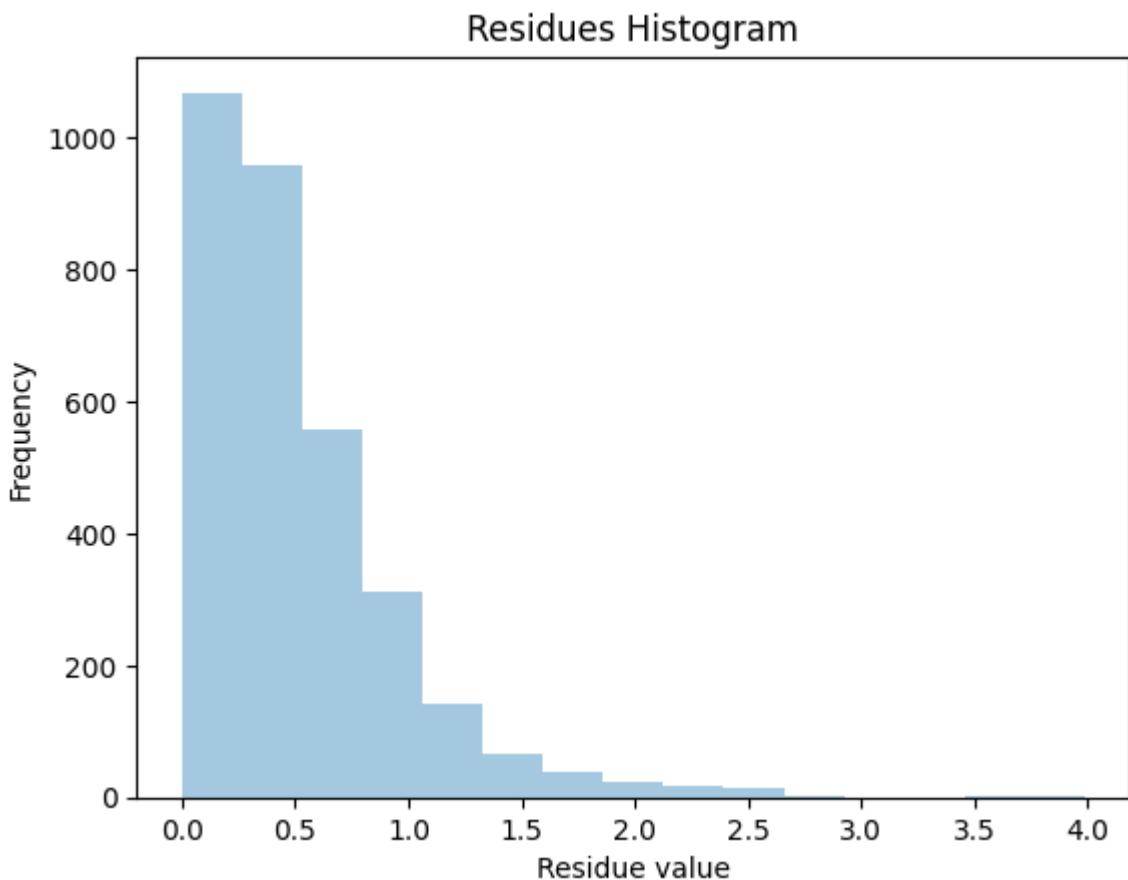
df = pd.read_csv('winequality-red.csv', delimiter=';')
X = df.drop('quality', axis=1)
y = df['quality']

X_train, X_test, y_train, y_test = train_test_split(X, y,
train_size=0.8, random_state=0)
residues, errors_original, errors_rounded = [], [], []

for i in range(1, 11):
    mlp = MLPRegressor(hidden_layer_sizes=(10, 10), activation='relu',
early_stopping=True, random_state=i, validation_fraction=0.2)
    mlp.fit(X_train, y_train)
    y_pred = mlp.predict(X_test)
    residues.append(np.abs(np.array(y_pred) - np.array(y_test)))

    # for next question
    rounded_predictions = np.round(y_pred)
    rounded_predictions = np.clip(rounded_predictions, 1, 10)
    errors_original.append(mean_absolute_error(y_test, y_pred))
    errors_rounded.append(mean_absolute_error(y_test, rounded_predictions))

sns.distplot(residues, kde=False, bins=15)
plt.xlabel('Residue value')
plt.ylabel('Frequency')
plt.title('Residues Histogram')
plt.show()
```



2. [1.5v] Since we are in the presence of a integer regression task, a recommended trick is to round and bound estimates. Assess the impact of these operations on the MAE of the MLP learnt in previous question.

```
In [ ]: mae_original = np.mean(errors_original)
mae_rounded = np.mean(errors_rounded)

print(f"Original MAE: {round(mae_original, 5)}")
print(f"MAE after rounding and bounding: {round(mae_rounded, 5)}")
```

Original MAE: 0.50972
MAE after rounding and bounding: 0.43875

3. [1.5v] Similarly assess the impact on RMSE from replacing early stopping by a well-defined number of iterations in {20,50,100,200} (where one iteration corresponds to a batch).

```
In [ ]: from sklearn.metrics import mean_squared_error

for iter_num in (20, 50, 100, 200):
    errors = []

    for i in range(1, 11):
        mlp = MLPRegressor(hidden_layer_sizes=(10, 10), activation='relu',
max_iter=iter_num, random_state=i)
        mlp.fit(X_train, y_train)
        y_pred = mlp.predict(X_test)

        errors.append(mean_squared_error(y_test, y_pred, squared=False))
```

```
rmse = np.mean(errors)
print(f"RMSE for MLP with max {iter_num} iterations: {round(rmse, 5)}")
```

RMSE for MLP with max 20 iterations: 1.40398
RMSE for MLP with max 50 iterations: 0.79961
RMSE for MLP with max 100 iterations: 0.69404
RMSE for MLP with max 200 iterations: 0.65545

4. [1.5v] Critically comment the results obtained in previous question, hypothesizing at least one reason why early stopping favors and/or worsens performance.

R: Se não utilizarmos *early stopping*, o fim do treino é baseado nos dados de treino, o que pode levar tanto a *underfitting* (caso pare num máximo local) como a *overfitting* (caso se ajuste demasiado aos dados de treino). O *early stopping* consiste em terminar o treino quando o *validation score* não é melhorado durante um certo número consecutivo de épocas, o que previne ambos os problemas mencionados anteriormente.

Analizando os resultados obtidos na pergunta anterior, podemos observar que com o aumento do número de iterações o erro diminui. Isto indica-nos que o *threshold* do *early stopping* deveria ser de pelo menos 200 iterações, visto que, até este ponto o erro diminui, pelo que não temos evidências de *overfitting* (neste caso, o erro aumentaria com o número de iterações).