

I. Pen-and-paper

1)

Homework 3

Using $\Phi(y_1, y_2) = y_1 \times y_2$:

1) $\hat{z} = XW$ ↗ pesos

↘ observações
↘ valor previsto

D	$\Phi(y_1, y_2)$
u_1	1
u_2	3
u_3	6
u_4	9
u_5	8

Using OLS closed form:

$$W = (X^T X)^{-1} X^T Z$$

↗ valores reais do output

$$X = \begin{bmatrix} 1 & 1 \\ 1 & 3 \\ 1 & 6 \\ 1 & 5 \\ 1 & 8 \end{bmatrix} \quad X^T = \begin{bmatrix} 1 & 1 & 1 & 1 & 1 \\ 1 & 3 & 6 & 5 & 8 \end{bmatrix}$$

$$X^T X = \begin{bmatrix} 1 & 1 & 1 & 1 & 1 \\ 1 & 3 & 6 & 5 & 8 \end{bmatrix} \begin{bmatrix} 1 & 1 \\ 1 & 3 \\ 1 & 6 \\ 1 & 5 \\ 1 & 8 \end{bmatrix} = \begin{bmatrix} 5 & 27 \\ 27 & 101 \end{bmatrix}$$

$$(X^T X)^{-1} = \begin{bmatrix} 0,84513 & -0,11047 \\ -0,11047 & 0,02211 \end{bmatrix}$$

$$(X^T X)^{-1} X^T Z = \begin{bmatrix} 3,31593 \\ 0,11372 \end{bmatrix}$$

$$\hat{z} = 3,31593 + 0,11372 \Phi(y_1, y_2)$$

2)

$\hat{z} = XW$ Using Ridge regression with penalty factor $\lambda = 1$:

$$W = (X^T X + \lambda I)^{-1} X^T z$$

$$X^T X = \begin{bmatrix} 5 & 27 \\ 27 & 191 \end{bmatrix} \quad (X^T X + \lambda I) = \begin{bmatrix} 6 & 27 \\ 27 & 192 \end{bmatrix}$$

$$(X^T X + \lambda I)^{-1} = \begin{bmatrix} 0,45390 & -0,06383 \\ -0,06383 & 0,014184 \end{bmatrix}$$

$$(X^T X + \lambda I)^{-1} X^T z = \begin{bmatrix} 1,81805 \\ 0,32376 \end{bmatrix}$$

$$\hat{z} = 1,81805 + 0,32376 \phi(y_1, y_2)$$

A análise das matrizes de pesos mostra que a regularização da Ridge Regression teve um impacto significativo na redistribuição dos valores dos pesos. No modelo de Regressão Linear, o primeiro peso era bastante elevado (3.3159), enquanto o segundo era muito pequeno (0.1137). Com a aplicação da Ridge Regression, o primeiro peso diminuiu substancialmente (1.8181), e o segundo aumentou ligeiramente (0.3238). Isto acontece porque a regularização de Ridge penaliza pesos grandes para evitar o overfitting, reduzindo a dependência excessiva de variáveis específicas.

Ao mesmo tempo, a Ridge Regression redistribui a importância entre as variáveis, permitindo que o segundo peso, que anteriormente era subestimado, se torne mais relevante. Desta forma, a esta equilibra melhor a influência das variáveis no modelo, melhorando a sua capacidade de generalização.

3)

3)

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

Para Regressão linear utilizando a forma fechada OLS:

$$\hat{z} = Xw = \begin{bmatrix} 1 & 1 \\ 1 & 3 \\ 1 & 6 \\ 1 & 9 \\ 1 & 8 \\ 1 & 4 \\ 1 & 2 \\ 1 & 5 \end{bmatrix} \begin{bmatrix} 3,31593 \\ 0,11372 \end{bmatrix} = \begin{bmatrix} 3,42965 \\ 3,65708 \\ 3,99823 \\ 4,33938 \\ 4,22567 \\ 3,77080 \\ 3,54336 \\ 3,88451 \end{bmatrix} \rightarrow \begin{matrix} \hat{z}_1 \\ \hat{z}_2 \\ \hat{z}_3 \\ \hat{z}_4 \\ \hat{z}_5 \\ \hat{z}_6 \\ \hat{z}_7 \\ \hat{z}_8 \end{matrix}$$

u_6, u_7, u_8

$$RMSE_{com} = \sqrt{\frac{\sum_{i=1}^5 (z_i - \hat{z}_i)^2}{5}} =$$

DataSet de treino = 2,026499

$$RMSE_{com} = \sqrt{\frac{\sum_{i=6}^8 (z_i - \hat{z}_i)^2}{3}} =$$

DataSet novo = 2,465590

Para a Ridge Regression com fator de penalidade $\lambda=1$:

$$\hat{z} = Xw = \begin{bmatrix} 1 & 1 \\ 1 & 3 \\ 1 & 6 \\ 1 & 9 \\ 1 & 8 \\ 1 & 4 \\ 1 & 2 \\ 1 & 5 \end{bmatrix} \begin{bmatrix} 1,81805 \\ 0,32376 \end{bmatrix} = \begin{bmatrix} 2,14184 \\ 2,78936 \\ 3,76064 \\ 4,73191 \\ 4,40816 \\ 3,11312 \\ 2,46560 \\ 3,14368 \end{bmatrix} \rightarrow \begin{matrix} \hat{z}_1 \\ \hat{z}_2 \\ \hat{z}_3 \\ \hat{z}_4 \\ \hat{z}_5 \\ \hat{z}_6 \\ \hat{z}_7 \\ \hat{z}_8 \end{matrix}$$

u_6, u_7, u_8

$$RMSE_{com} = \sqrt{\frac{\sum_{i=1}^5 (z_i - \hat{z}_i)^2}{5}} =$$

DataSet de treino = 2,15354

$$RMSE_{com} = \sqrt{\frac{\sum_{i=6}^8 (z_i - \hat{z}_i)^2}{3}} =$$

DataSet novo = 1,75289

Como o RMSE nos dados de teste para o modelo regularizado (Ridge) é menor do que o RMSE do modelo não regularizado, podemos concluir que a regularização melhorou a capacidade de generalização do modelo, reduzindo o erro e tornando as previsões mais precisas, o que era esperado devido à penalização de coeficientes excessivamente altos.

Por outro lado, o RMSE nos dados de treino é maior no modelo regularizado do que no modelo não regularizado. Isto ocorre porque a regularização impõe restrições nos coeficientes, evitando que o modelo se ajuste demasiado aos dados de treino (overfitting), o que leva a um maior erro nesse conjunto.

No entanto, esta ligeira perda de desempenho para os dados de treino é compensada por uma melhor performance nos dados de teste, indicando que o modelo regularizado está menos suscetível a overfitting e, portanto, tem maior capacidade de generalização.

4)

$$W^{[1]} = \begin{bmatrix} 0,1 & 0,1 \\ 0,1 & 0,2 \\ 0,2 & 0,1 \end{bmatrix} \quad b^{[1]} = \begin{bmatrix} 0,1 \\ 0 \\ 0,1 \end{bmatrix} \quad W^{[2]} = \begin{bmatrix} 1 & 2 & 2 \\ 1 & 2 & 1 \\ 1 & 1 & 1 \end{bmatrix} \quad b^{[2]} = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}$$

Output activation function:

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

Hidden layers activation function: None $\Rightarrow z^{[1]} = x^{[1]}$

Error function \rightarrow Cross-Entropy loss:

$$\text{Error} = - \sum_{i=1}^N \sum_{l=1}^{|c|} t_l^{(i)} \log(z_l^{[out]}(i))$$

Back propagation:

$$W^{[i]} = W^{[i]} - \eta \frac{\partial \mathcal{E}}{\partial W^{[i]}} \quad \rightarrow = 0,1$$

$$b^{[i]} = b^{[i]} - \eta \frac{\partial \mathcal{E}}{\partial b^{[i]}}$$

$$\begin{aligned} \frac{\partial \mathcal{E}}{\partial W^{[i]}} &= \frac{\partial \mathcal{E}}{\partial X^{[i]}} \circ \frac{\partial X^{[i]}}{\partial z^{[i]}} \cdot \left(\frac{\partial z^{[i]}}{\partial W^{[i]}} \right)^T \\ &= \delta^{[i]} \cdot (X^{[i-1]})^T \end{aligned}$$

$$\frac{\partial \mathcal{E}}{\partial b^{[i]}} = \frac{\partial \mathcal{E}}{\partial X^{[i]}} \circ \frac{\partial X^{[i]}}{\partial z^{[i]}} \cdot \frac{\partial z^{[i]}}{\partial b^{[i]}} = \delta^{[i]}$$

Forward Propagation

$$z^{[i]} = W^{[i]} \cdot x^{[i-1]} + b^{[i]}$$

$$x^{[i]} = \text{ativação}(z^{[i]})$$

$$z^{[1]} = W^{[1]} x^{[0]} + b^{[1]}$$

$$= \begin{bmatrix} 0,1 & 0,1 \\ 0,1 & 0,2 \\ 0,2 & 0,1 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix} + \begin{bmatrix} 0,1 \\ 0 \\ 0,1 \end{bmatrix} = \begin{bmatrix} 0,3 \\ 0,3 \\ 0,4 \end{bmatrix}$$

$$h_1 = \begin{bmatrix} y_1 & y_2 \\ 1 & 1 \end{bmatrix}$$

$$\Phi(y_1, y_2) = 1 \times 1 = 1$$

$$x_0 = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

$$z^{[2]} = W^{[2]} x^{[1]} + b^{[2]}$$

$$= \begin{bmatrix} 1 & 2 & 2 \\ 1 & 2 & 1 \\ 1 & 1 & 1 \end{bmatrix} \begin{bmatrix} 0,3 \\ 0,3 \\ 0,4 \end{bmatrix} + \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} = \begin{bmatrix} 2,7 \\ 2,3 \\ 2 \end{bmatrix}$$

$$x^{[2]} = \text{softmax}(z_c^{[2]}) = \frac{e^{z_c^{[out]}}}{\sum_{l=1}^3 e^{z_l^{[out]}}}$$

$$\sum_{l=1}^3 e^{z_l^{[out]}} = e^{2,7} + e^{2,3} + e^2 = 32,24297$$

$$x^{[2]} = \begin{bmatrix} \frac{e^{2,7}}{32,24297} \\ \frac{e^{2,3}}{32,24297} \\ \frac{e^2}{32,24297} \end{bmatrix} = \begin{bmatrix} 0,4615 \\ 0,3053 \\ 0,2292 \end{bmatrix}$$

$$W^{[2]} = W^{[2]} - 0,1 \frac{\partial E}{\partial W^{[2]}}$$

$$E = - \sum_{i=1}^N \sum_{l=1}^L t_l^{(i)} \log(z_l^{[out] (i)}), N=1, \log 0$$

$$= - \sum_{l=1}^L t_l \log(z_l^{[out]})$$

$$\frac{\partial E}{\partial W^{[2]}} = \sum^{[2]} X^{[1]} \text{ isto porque classificação: softmax + cross entropy loss}$$

$$\sum^{[2]} = X^{[out]} - t \text{ onde } t \text{ para } x_1 = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}$$

$$= \begin{bmatrix} 0,4615 \\ 0,3093 \\ 0,2292 \end{bmatrix} - 1 = \begin{bmatrix} 0,4615 \\ -0,6907 \\ 0,2292 \end{bmatrix}$$

$$\frac{\partial E}{\partial W^{[2]}} = \begin{bmatrix} 0,4615 \\ -0,6907 \\ 0,2292 \end{bmatrix} \begin{bmatrix} 0,3 & 0,3 & 0,4 \end{bmatrix} = \begin{bmatrix} 0,13845 & 0,13845 & 0,18468 \\ -0,20720 & -0,20720 & -0,27626 \\ 0,06875 & 0,06875 & 0,091667 \end{bmatrix}$$

$$W^{[2]} = \begin{bmatrix} 1 & 2 & 2 \\ 1 & 2 & 1 \\ 1 & 1 & 1 \end{bmatrix} - 0,1 \frac{\partial E}{\partial W^{[2]}} = \begin{bmatrix} 0,98616 & 1,98616 & 1,98154 \\ 1,02072 & 2,02072 & 1,02763 \\ 0,99312 & 0,99312 & 0,99083 \end{bmatrix}$$

$$b^{[2]} = b^{[2]} - 0,1 \frac{\partial E}{\partial b^{[2]}} =$$

$$= \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} - 0,1 \sum^{[2]} = \begin{bmatrix} 0,95385 \\ 1,06907 \\ 0,97708 \end{bmatrix}$$

$$W^{[1]} = W^{[1]} - 0,1 \frac{\partial \mathcal{E}}{\partial W^{[1]}}$$

$$\frac{\partial \mathcal{E}}{\partial W^{[1]}} = \delta^{[1]} X^{[0]}$$

Não é função de ativação

$$\delta^{[1]} = \left(W^{[2]} \right)^T \cdot \delta^{[2]} \odot \frac{\partial X^{[1]}}{\partial z^{[1]}}$$

$$\delta^{[1]} = \begin{bmatrix} 5,55111 \times 10^{-17} \\ -0,22912 \\ 0,46149 \end{bmatrix}$$

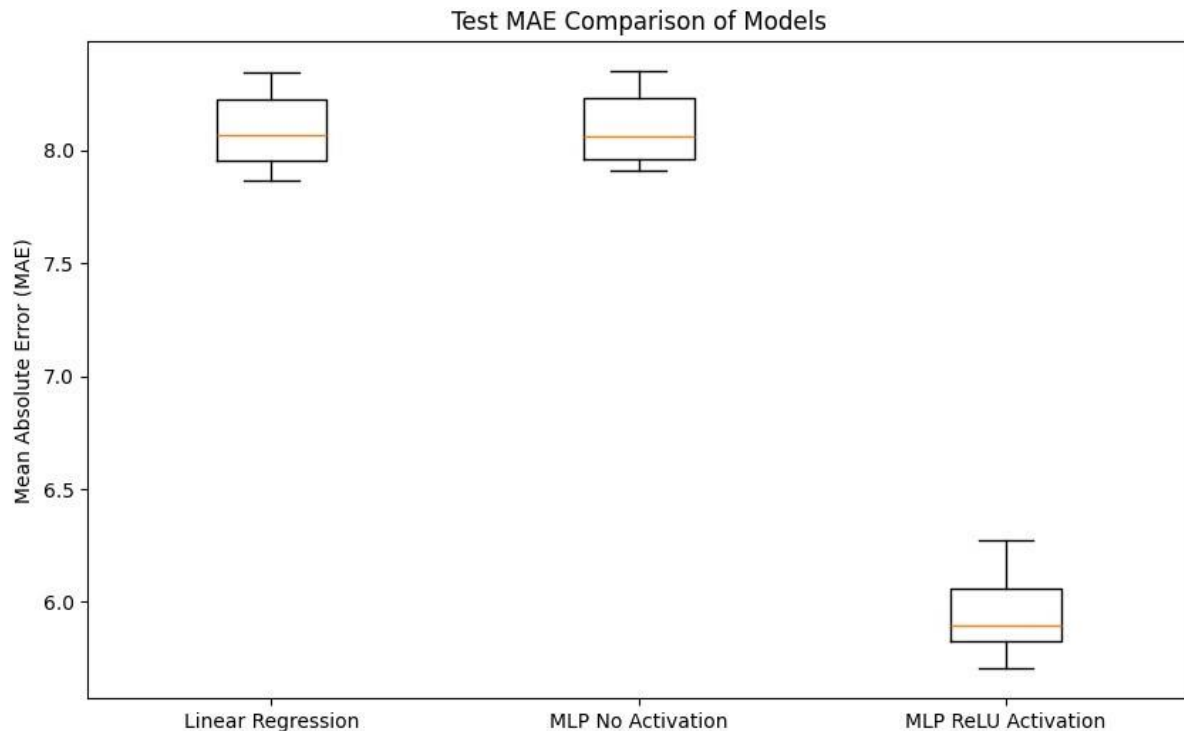
$$W^{[1]} = \begin{bmatrix} 0,1 & 0,1 \\ 0,1 & 0,2 \\ 0,2 & 0,1 \end{bmatrix} - 0,1 \left(\delta^{[1]} X^{[0]} \right)$$

$$= \begin{bmatrix} 0,1 & 0,1 \\ 0,12292 & 0,2229 \\ 0,15385 & 0,0539 \end{bmatrix}$$

$$b^{[1]} = \begin{bmatrix} 0,1 \\ 0 \\ 0,1 \end{bmatrix} - 0,1 \delta^{[1]} = \begin{bmatrix} 0,1615 \\ 0,02292 \\ 0,05385 \end{bmatrix}$$

II. Programming and critical analysis

5)



6)

Uma regressão linear é um modelo simples capaz de estabelecer relações entre as features e o target. Este modelo assume que os dados podem ser modelados através de uma reta (ou plano, etc. para dimensões maiores). Está, portanto, limitada para capturar relações não lineares complexas.

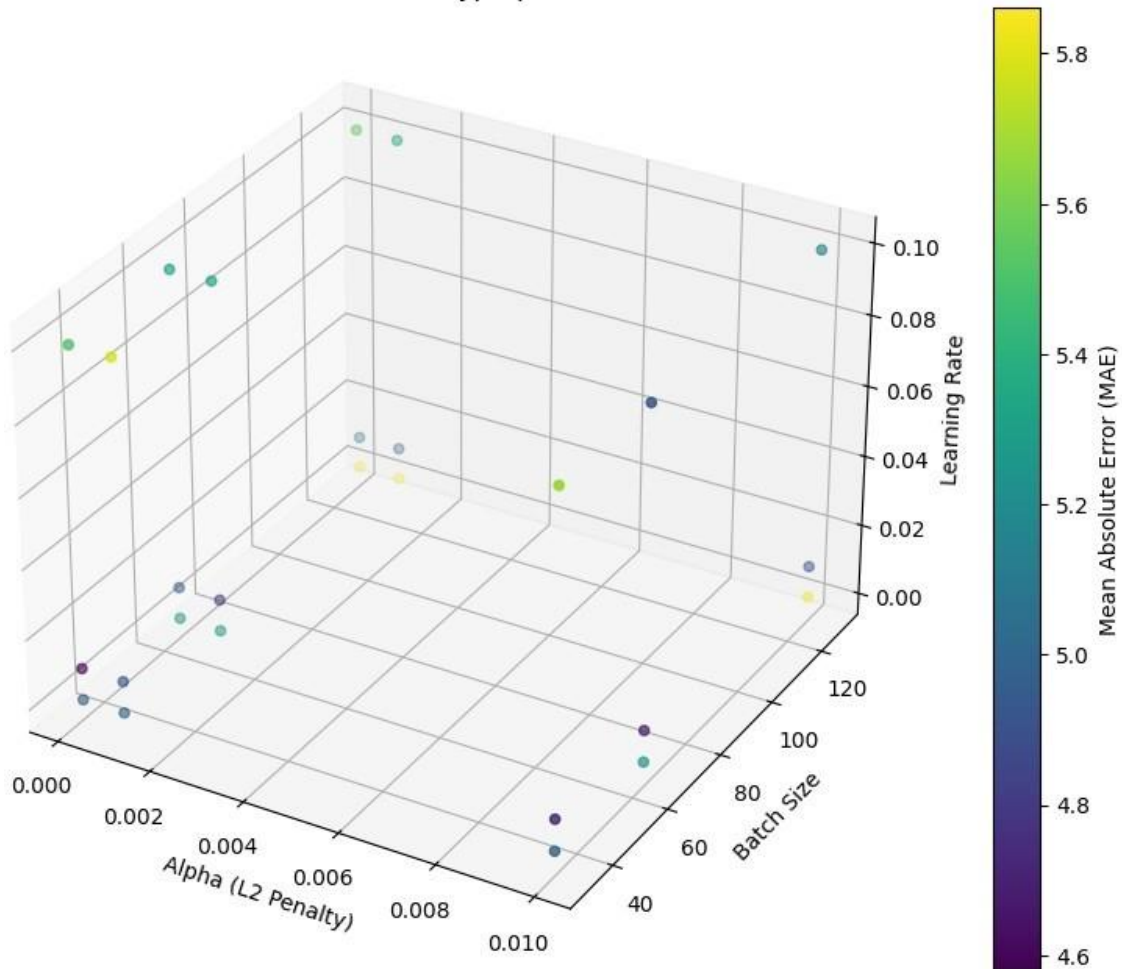
Um modelo MLP sem funções de ativação é equivalente a uma regressão linear. Se não tiver funções de ativação não lineares, o output da rede neuronal será sempre uma transformação linear do input (independentemente do número de camadas).

Assim, e tal como podemos observar no gráfico do exercício anterior, uma rede neuronal e um modelo MLP sem funções de ativação vão ter resultados idênticos.

Entende-se, portanto, a importância de funções de ativação não lineares nos modelos MLP. Ao incluir funções de ativação, o modelo torna-se capaz de capturar padrões mais complexos, incluindo relações não lineares (interações entre as variáveis que seriam impossíveis de identificar com uma simples transformação linear). Isto permite que o modelo MLP aprenda representações mais completas e flexíveis dos dados, melhorando a sua capacidade de generalizar para novos inputs. Embora esta flexibilidade possa aumentar o risco de overfitting, também possibilita ao modelo extrair padrões complexos que a regressão linear não consegue capturar.

7)

MAE Visualization for Different Hyperparameter Combinations



Nesta representação em 3D, é possível observar como o desempenho do modelo, medido pelo Erro Absoluto Médio (MAE), varia em função de diferentes combinações de hiperparâmetros: alpha, learning rate e batch size. A barra de cores representa os valores de MAE, com tons mais escuros indicando erros menores - melhor desempenho do modelo.

Uma learning rate mais alta permite que o modelo convirja rapidamente, mas pode resultar em "saltos" sobre a solução ótima, prejudicando a generalização. Em contrapartida, uma learning rate mais baixa torna o treino mais demorado, mas frequentemente resulta num desempenho superior ao explorar soluções de forma mais detalhada.

Valores elevados de alpha proporcionam uma forte regularização, ajudando a reduzir o overfitting ao penalizar pesos excessivos. No entanto, se o valor de alpha for demasiado alto, o modelo pode deixar de identificar padrões relevantes nos dados. Por outro lado, valores baixos de alpha diminuem a regularização, permitindo um melhor ajuste aos dados de treinamento, mas aumentando o risco de overfitting.

Quanto ao batch size, um valor reduzido atualiza os pesos com mais frequência, gerando uma aprendizagem mais "ruidosa" que pode ajudar o modelo a evitar mínimos locais, embora demore mais para convergir. Em

contraste, valores grandes de batch size processam mais dados antes de realizar atualizações, resultando em mudanças mais precisas e estáveis, mas que podem tornar a aprendizagem menos variável, possivelmente fazendo com que o modelo fique preso em mínimos locais.

No gráfico, observamos claramente a vantagem de utilizar valores mais baixos de batch size e learning rate, como esperado. No caso de alpha, no entanto, abordagens mais extremas parecem funcionar melhor. Tanto um alpha baixo (que permite ao modelo aprender relações mais complexas entre os dados) como um alpha alto (que simplifica o modelo) produzem melhores resultados do que valores intermédios.

END