



POLYTECHNIQUE MONTRÉAL

INF8245AE – MACHINE LEARNING

## Assignment 1 – Linear Regression

*Mattéo Colavita - 2142009*

# 1 Question 1 : Linear and Weighted Ridge Regression

$$\begin{aligned}
 L(w) &= \|Xw - y\|_2^2 + w^T \Lambda w \\
 1. \quad L'(w) &= [Xw - y](Xw - y)^T + w^T \Lambda w \\
 \begin{matrix} Xw: n \times 1 \\ y: n \times 1 \\ X^T X: n \times n \\ X^T y: n \times 1 \end{matrix} & \quad L'(w) = [w^T X^T X w - y^T X w - y^T X w + y^T y + w^T \Lambda w] \\
 \begin{matrix} (Xw)^T (Xw) \\ w^T X^T X w \end{matrix} & \quad L'(w) = [w^T X^T X w - 2y^T X w + y^T y + w^T \Lambda w] \\
 \begin{matrix} \text{Shape of } y \text{ is } n \times 1 \\ \text{so} \end{matrix} & \quad L'(w) = (\partial w)^T X^T X w + w^T X^T X (\partial w) - 2y^T X (\partial w) + (\partial w)^T \Lambda w + w^T \Lambda (\partial w) \\
 \begin{matrix} Xw: n \times 1 \\ y: n \times 1 \\ y^T X w: 1 \times 1 \end{matrix} & \quad L'(w) = 2(X^T X w) - 2y^T X (\partial w) + (y^T y) + (w^T \Lambda w) \\
 & \quad = 2(X^T X w) - 2y^T X (\partial w) + 0 + (\partial w)^T \Lambda w + w^T \Lambda (\partial w) \\
 & \quad = 2(X^T X w) - 2y^T X (\partial w) + (\partial w)^T \Lambda w + (\partial w)^T \Lambda w \\
 & \quad = 2(X^T X w) - (\partial w)^T 2X^T y + 2(\Lambda w) \\
 \frac{\partial L(w)}{\partial w} &= 2(X^T X w + \Lambda w) - 2X^T y \\
 2. \text{ Minimize with } \frac{\partial L(w)}{\partial w} &= 0 \\
 2(X^T X w + \Lambda w) &= 2X^T y \\
 X^T X w + \Lambda w &= X^T y \\
 w^* &= (X^T X w + \Lambda)^{-1} X^T y
 \end{aligned}$$

FIGURE 1 – Derivative of the loss function for ridge regression with respect to the weights.

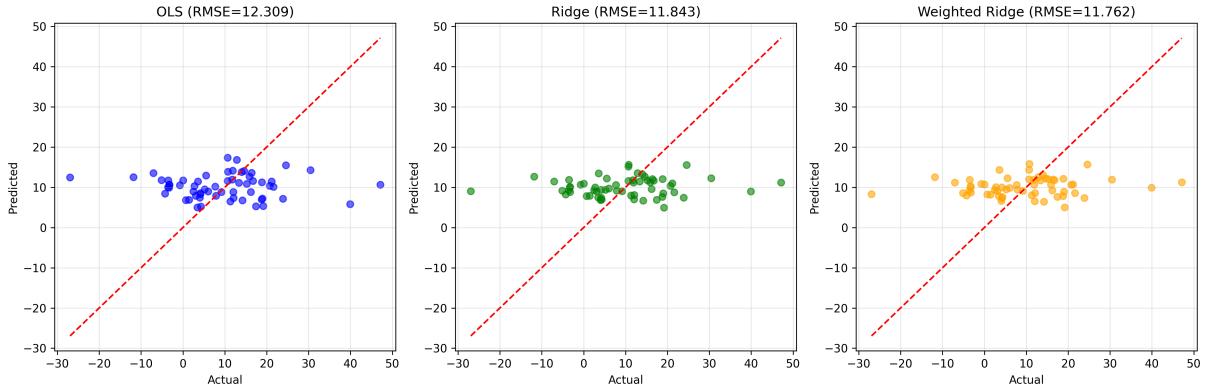


FIGURE 2 – Comparison of predictions from linear, ridge and weighted ridge regression on the test set.

## 2 Question 2 : Cross-Validation

Metric	Best $\lambda$	$\lambda=0.01$	$\lambda=0.1$	$\lambda=1$	$\lambda=10$	$\lambda=100$
MAE	10	7.381	7.316	7.140	7.110	7.817
MaxError	100	27.758	27.681	27.559	27.476	27.095
RMSE	10	9.855	9.772	9.577	9.532	10.101

TABLE 1 – Mean MAE, MaxError, and RMSE scores obtained via 5-fold cross-validation for different values of  $\lambda$  in ridge regression.

### 3 Question 3 : Gradient Descent for Ridge Regression

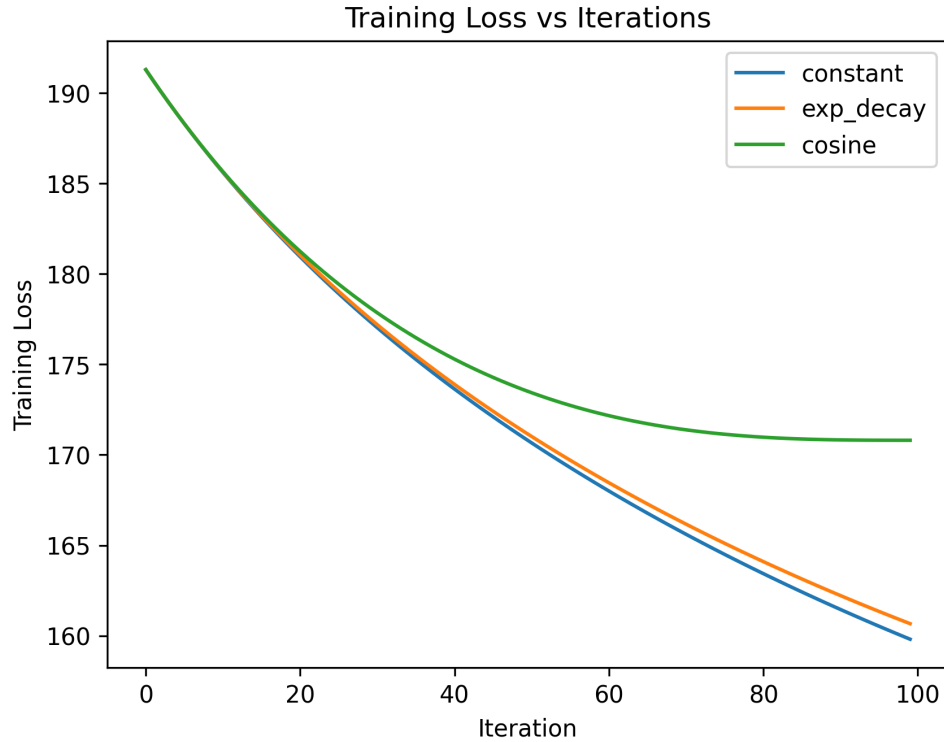


FIGURE 3 – Training loss vs iterations for different learning rate schedules in gradient descent for ridge regression.

Learning Rate Schedule	RMSE
Constant	13.8910
Exponential Decay	13.9283
Cosine Annealing	14.3473

TABLE 2 – RMSE results for different learning rate schedules in gradient descent for ridge regression.

Among the tested learning rate schedules, the constant learning rate shows the fastest decrease in training loss and achieves the lowest test RMSE of 13.89, suggesting better generalization. Similarly, exponential decay showed a fast decrease in training loss, with a slightly higher RMSE of 13.93, while cosine annealing plateaued early, indicating slower convergence and less effective learning within the given number of iterations. None of the schedules fully converged within 100 iterations, as the training loss was still decreasing, but constant and exponential decay showed a clear downward trend that could improve further with more iterations.