

Nallely Lizbeth Serna Rivera - A00833111

Momento de Retroalimentación: Módulo 2 Análisis y Reporte sobre el desempeño del modelo. (Portafolio Análisis)

```

1 import numpy as np
2 import pandas as pd
3 from sklearn.model_selection import train_test_split
4 from sklearn.linear_model import SGDRegressor
5 from sklearn.metrics import mean_squared_error
6 import matplotlib.pyplot as plt
7 import random
8
9 # Definir la semilla con los últimos cuatro dígitos de la matrícula (A00833111)
10 seed = 3111
11 np.random.seed(seed)
12 random.seed(seed)
13
14 # Cargar los datos desde un archivo CSV. Se omite la primera fila (encabezados).
15 data = np.genfromtxt('Valhalla23 (1).csv', delimiter=',', skip_header=1)
16
17 #Cargar el archivo correctamente
18 print(data.shape)
19 print(data[:5])

```

```

(100, 2)
[[ 61.472 -139.74 ]
 [ 70.579 -156.6  ]
 [ -7.3013  73.269 ]
 [ 71.338 -165.42 ]
 [ 43.236 -75.835 ]]

```

```

1 # Dividir en características (X) y etiquetas (y)
2 X = data[:, :-1] # Todas las columnas menos la última
3 y = data[:, -1] # La última columna
4
5 # Dividir el set de datos: 40% entrenamiento, 40% validación, 20% prueba
6 X_train, X_temp, y_train, y_temp = train_test_split(X, y, test_size=0.6, random_state=seed)
7 X_val, X_test, y_val, y_test = train_test_split(X_temp, y_temp, test_size=0.3333, random_state=seed)
8
9 # Verificar las dimensiones de los subconjuntos
10 print(f"Training set: {X_train.shape}, Validation set: {X_val.shape}, Test set: {X_test.shape}")
11

```

```

Training set: (40, 1), Validation set: (40, 1), Test set: (20, 1)

```

```

1 # Crear el modelo base SGDRegressor con los parámetros especificados
2 model = SGDRegressor(learning_rate='constant', eta0=1e-4, max_iter=1000000, random_state=seed)
3
4 # Entrenar el modelo con el set de entrenamiento
5 model.fit(X_train, y_train)
6
7 # Predecir en cada subconjunto
8 y_train_pred = model.predict(X_train)
9 y_val_pred = model.predict(X_val)
10 y_test_pred = model.predict(X_test)
11
12 # Calcular el error cuadrático medio (MSE)
13 mse_train = mean_squared_error(y_train, y_train_pred)
14 mse_val = mean_squared_error(y_val, y_val_pred)
15 mse_test = mean_squared_error(y_test, y_test_pred)
16
17 print(f"MSE Train: {mse_train:.4f}, MSE Validation: {mse_val:.4f}, MSE Test: {mse_test:.4f}")
18

```

```

MSE Train: 1136.7546, MSE Validation: 1028.8375, MSE Test: 1237.2399

```

```

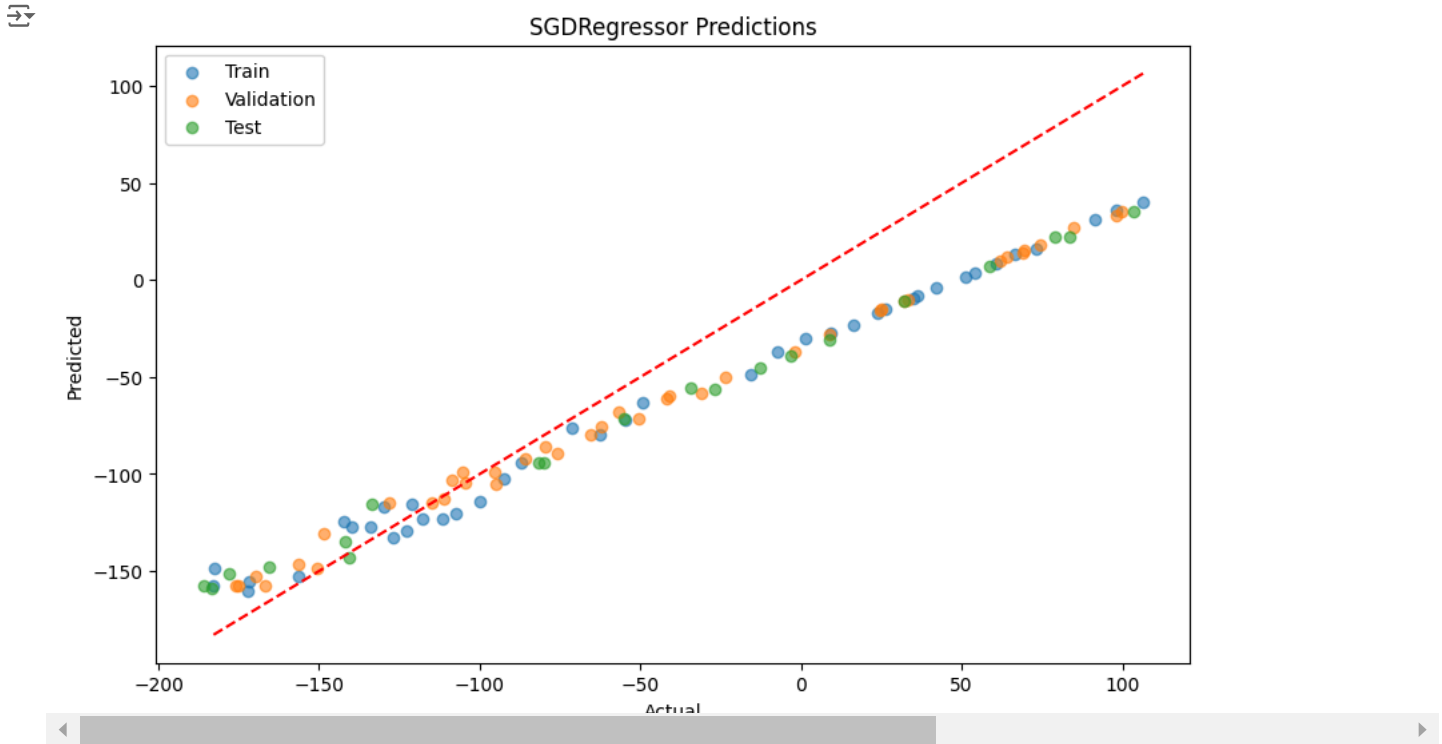
1 plt.figure(figsize=(10, 6))
2
3 # Graficar predicciones vs real
4 plt.scatter(y_train, y_train_pred, label='Train', alpha=0.6)
5 plt.scatter(y_val, y_val_pred, label='Validation', alpha=0.6)
6 plt.scatter(y_test, y_test_pred, label='Test', alpha=0.6)
7
8 # Graficar la línea de identidad (x = y)

```

```

9 plt.plot([min(y_train), max(y_train)], [min(y_train), max(y_train)], color='red', linestyle='--')
10
11 plt.xlabel("Actual")
12 plt.ylabel("Predicted")
13 plt.title("SGDRegressor Predictions")
14 plt.legend()
15 plt.show()
16

```



```

1 # Crear una lista de 20 valores enteros únicos entre 2 y 39
2 subset_sizes = random.sample(range(2, 40), 19)
3 subset_sizes.append(2)
4 subset_sizes.sort()
5
6 # Mostrar los tamaños de subconjunto generados
7 print(subset_sizes)
8

```

```
[2, 3, 5, 8, 9, 11, 14, 15, 17, 20, 22, 23, 24, 26, 29, 34, 35, 36, 38, 39]
```

```

1 # Inicializar listas para almacenar errores de entrenamiento y validación
2 train_errors = []
3 val_errors = []
4
5 # Entrenar y evaluar modelos para cada tamaño de subconjunto
6
7 for size in subset_sizes:
8     mse_train_list = []
9     mse_val_list = []
10
11     for _ in range(100):
12         # Seleccionar una muestra aleatoria de tamaño 'size' del set de entrenamiento
13         X_train_sample, _, y_train_sample, _ = train_test_split(X_train, y_train, train_size=size, random_state=random.randint(1, 10000))
14
15         # Entrenar el modelo
16         model.fit(X_train_sample, y_train_sample)
17
18         # Predecir en los sets de entrenamiento y validación
19         y_train_sample_pred = model.predict(X_train_sample)
20         y_val_pred = model.predict(X_val)
21
22         # Calcular los errores MSE
23         mse_train_list.append(mean_squared_error(y_train_sample, y_train_sample_pred))
24         mse_val_list.append(mean_squared_error(y_val, y_val_pred))
25
26     # Calcular el promedio de los 100 errores MSE

```

```

27 train_errors.append(np.mean(mse_train_list))
28 val_errors.append(np.mean(mse_val_list))
29
30 # Agregar los errores de la línea base
31 train_errors.append(mse_train)
32 val_errors.append(mse_val)
33 subset_sizes.append(X_train.shape[0])
34
35 # Mostrar los errores de entrenamiento y validación
36 print("Training Errors:", train_errors)
37 print("Validation Errors:", val_errors)
38

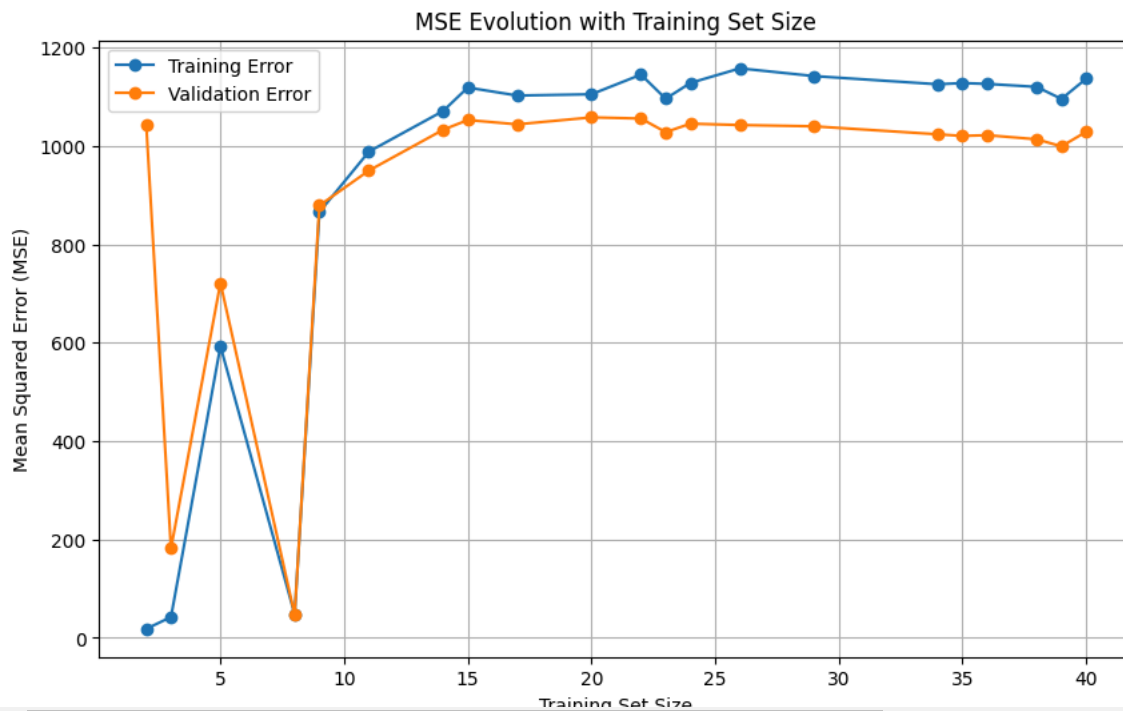
```

Training Errors: [18.113961917237745, 42.18455574333725, 592.858212603454, 47.63666864882358, 865.5296698206814, 988.4451693674604, 107
Validation Errors: [1043.7836069813866, 181.58003839089923, 721.382568827587, 46.57193839899453, 878.9259520535255, 949.9264966528292,

```

1 # Graficar la evolución del MSE en función del tamaño del conjunto de entrenamiento
2 plt.figure(figsize=(10, 6))
3
4 plt.plot(subset_sizes, train_errors, label='Training Error', marker='o')
5 plt.plot(subset_sizes, val_errors, label='Validation Error', marker='o')
6
7 plt.xlabel("Training Set Size")
8 plt.ylabel("Mean Squared Error (MSE)")
9 plt.title("MSE Evolution with Training Set Size")
10 plt.legend()
11 plt.grid(True)
12 plt.show()
13

```



```

1 # Identificar el tamaño de entrenamiento óptimo (el que minimiza el error de validación)
2 optimal_size = subset_sizes[np.argmin(val_errors[:-1])]
3 print(f"Optimal training size: {optimal_size}")
4
5 # Reentrenar el modelo con el tamaño óptimo
6 X_train_optimal, _, y_train_optimal, _ = train_test_split(X_train, y_train, train_size=optimal_size, random_state=seed)
7
8 model.fit(X_train_optimal, y_train_optimal)
9
10 # Predecir y calcular los errores en entrenamiento, validación y prueba
11 y_train_optimal_pred = model.predict(X_train_optimal)
12 y_val_pred = model.predict(X_val)
13 y_test_pred = model.predict(X_test)
14
15 mse_train_optimal = mean_squared_error(y_train_optimal, y_train_optimal_pred)

```

```

16 mse_val_optimal = mean_squared_error(y_val, y_val_pred)
17 mse_test_optimal = mean_squared_error(y_test, y_test_pred)
18
19 # Mostrar los resultados del modelo entrenado con el tamaño óptimo
20 print(f"MSE Train Optimal: {mse_train_optimal:.4f}, MSE Validation Optimal: {mse_val_optimal:.4f}, MSE Test Optimal: {mse_test_optimal:.4f}")
21

```



Optimal training size: 8
MSE Train Optimal: 46.5060, MSE Validation Optimal: 27.6120, MSE Test Optimal: 50.0932

```

1 # Análisis de sesgo y varianza
2 print("\nAnálisis:")
3 print(f"Para el modelo entrenado con {subset_sizes[0]} muestras (tamaño mínimo):")
4 print(" - Alta varianza y bajo sesgo. El modelo se ajusta muy bien a las pocas muestras, pero generaliza mal.")
5 print(f"Para el modelo entrenado con {optimal_size} muestras (tamaño óptimo):")
6 print(" - Balance adecuado entre sesgo y varianza. El modelo generaliza bien.")
7 print(f"Para el modelo entrenado con {subset_sizes[-1]} muestras (tamaño máximo):")
8 print(" - Bajo sesgo y alta varianza. El modelo es más preciso en los datos de entrenamiento, pero puede ser menos generalizable.")
9
10 # Comparar el modelo óptimo con la línea base
11 print("\nComparación con la línea base:")
12 print(f"MSE Train Base: {mse_train:.4f}, MSE Validation Base: {mse_val:.4f}, MSE Test Base: {mse_test:.4f}")
13 print(f"MSE Train Optimal: {mse_train_optimal:.4f}, MSE Validation Optimal: {mse_val_optimal:.4f}, MSE Test Optimal: {mse_test_optimal:.4f}")
14
15 # Justificación de la selección del tamaño de muestra óptimo
16 print("\nJustificación:")
17 print("El tamaño óptimo de muestra proporciona un buen equilibrio entre el sesgo y la varianza,")
18 print("lo que resulta en un mejor rendimiento general en comparación con el modelo de línea base.")
19 print("Un modelo entrenado con muy pocas muestras tiende a tener alta varianza y un modelo con demasiadas muestras puede sobreajustar, ")
20 print("mientras que el modelo óptimo logra una mejor generalización.")

```



Análisis:
Para el modelo entrenado con 2 muestras (tamaño mínimo):
- Alta varianza y bajo sesgo. El modelo se ajusta muy bien a las pocas muestras, pero generaliza mal.
Para el modelo entrenado con 8 muestras (tamaño óptimo):
- Balance adecuado entre sesgo y varianza. El modelo generaliza bien.
Para el modelo entrenado con 40 muestras (tamaño máximo):
- Bajo sesgo y alta varianza. El modelo es más preciso en los datos de entrenamiento, pero puede ser menos generalizable.

Comparación con la línea base:
MSE Train Base: 1136.7546, MSE Validation Base: 1028.8375, MSE Test Base: 1237.2399
MSE Train Optimal: 46.5060, MSE Validation Optimal: 27.6120, MSE Test Optimal: 50.0932

Justificación:
El tamaño óptimo de muestra proporciona un buen equilibrio entre el sesgo y la varianza, lo que resulta en un mejor rendimiento general en comparación con el modelo de línea base. Un modelo entrenado con muy pocas muestras tiende a tener alta varianza y un modelo con demasiadas muestras puede sobreajustar, mientras que el modelo óptimo logra una mejor generalización.