

Entrega 6: diseño de redes neuronales

Aprendizaje Automatico y Big Data- Alejandro Barrachina Argudo

Introducción

En este documento se explicará el código del entregable 6 y el proceso de diseño de redes neuronales.

Para esta práctica se usarán los siguientes *imports* vistos en la figura 0.1. Parte del código se reutiliza de la práctica anterior.

```
1 from typing import Union
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import sklearn.linear_model as lm
5 import sklearn.preprocessing as sp
6 import sklearn.model_selection as ms
7 import cmdline # Custom command line parser
8 import os
9 import sys
```

Figura 0.1: Código de las bibliotecas usadas

También usaremos una serie de constantes para todo el programa (figura 0.2).

```
1 # Constants
2 # Path to save the plots
3 plot_folder = "./memoria/images"
4 # Path to save the csv files
5 csv_folder = "./memoria/csv"
6 # Random state for reproducibility
7 RANDOM_STATE = 1
```

Figura 0.2: Constantes del programa

El *dataset* para esta práctica lo generamos aleatoriamente con la función *gen_data* (figura 0.4). El dataset se compone de una linea de datos “ideales” y datos con ruido para comprobar la eficacia de la red neuronal.

Para dibujar estos datos usaremos la función *plot_dataset* (figura 0.5).

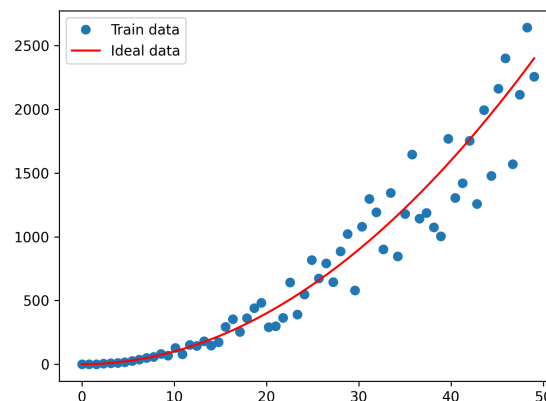


Figura 0.3: Ejemplo de los dígitos del *dataset*

```

1 def gen_data(m: int, seed: int = 1, scale: float = 0.7) -> tuple[np.ndarray, np.ndarray, np.
  ndarray, np.ndarray]:
2     """Generates a dataset with noise
3
4     Args:
5         m (int): number of samples
6         seed (int, optional): random seed. Defaults to 1.
7         scale (float, optional): scale of the noise. Defaults to 0.7.
8
9     Returns:
10         tuple[np.ndarray, np.ndarray, np.ndarray, np.ndarray]: x_train, y_train, x_ideal,
11         y_ideal
12     """
13     c: int = 0
14     x_train: np.ndarray = np.linspace(0, 49, m)
15     np.random.seed(seed)
16     y_ideal: np.ndarray = x_train**2 + c
17     y_train: np.ndarray = y_ideal + scale * \
18         y_ideal * (np.random.sample((m,)) - 0.5)
19     x_ideal: np.ndarray = x_train
20     return x_train, y_train, x_ideal, y_ideal

```

Figura 0.4: Función *gen_data*

```

1 def plot_dataset(x: np.ndarray, y: np.ndarray, x_ideal: np.ndarray, y_ideal: np.ndarray, name:
  str) -> None:
2     """Plots the dataset and the ideal data
3
4     Args:
5         x (np.ndarray): x values of the dataset with noise
6         y (np.ndarray): y values of the dataset with noise
7         x_ideal (np.ndarray): x ideal values of the dataset
8         y_ideal (np.ndarray): y ideal values of the dataset
9         name (str): name of the file
10    """
11    plt.plot(x, y, 'o', label='Train data')
12    plt.plot(x_ideal, y_ideal, label='Ideal data', c='red')
13    plt.legend()
14    plt.savefig(f'{plot_folder}/{name}.png', dpi=300)
15    plt.clf()

```

Figura 0.5: Función *plot_dataset*

1. Sobreajuste a los ejemplos de entrenamiento

En este apartado vamos a analizar lo que ocurre si tienes pocos datos de test, haciendo así que el problema se amolde solo a los datos de entrenamiento y no consiga generalizar.

Para ello usaremos la función *overfitting* (figura 1.2). Dentro de esta función primero hacemos una separación de los datos para dejar un porcentaje del 67% a los datos de entrenamiento y un 33% a los datos de test. Seguido esto, utilizamos las funciones *train*(1.3) para entrenar el modelo lineal y *test*(1.4) para sacar los costes (función *cost*, figura 1.5) de ambos conjuntos de datos. La figura 1.1 nos muestra en gráfica el sobreajuste producido por esta función. Este gráfico se genera con la función *plot_linear_data* (figura 1.6).

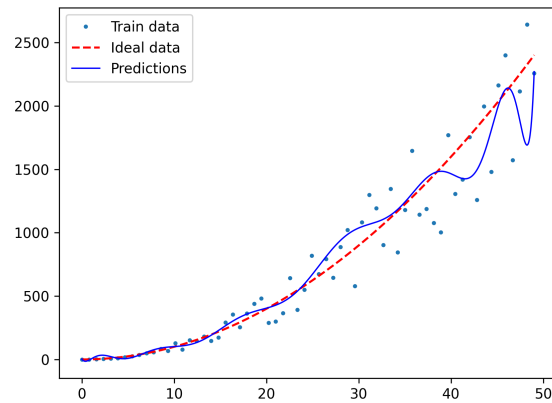


Figura 1.1: Gráfica del sobreajuste

```

1 def overfitting(x: np.ndarray, y: np.ndarray, x_i: np.ndarray, y_i: np.ndarray) -> None:
2     """Tests the overfitting of the model
3
4     Args:
5         x (np.ndarray): x values of the dataset with noise
6         y (np.ndarray): y values of the dataset with noise
7         x_i (np.ndarray): x ideal values of the dataset
8         y_i (np.ndarray): y ideal values of the dataset
9     """
10    print("Overfitting")
11    x_train, x_test, y_train, y_test = ms.train_test_split(
12        x, y, test_size=0.33, random_state=RANDOM_STATE)
13    pol, scal, model, x_train_aux = train(x_train, y_train, 15)
14    range_x: np.ndarray = np.linspace(np.min(x), np.max(x), 1000)
15    range_x = range_x[:, None]
16    range_x_p: np.ndarray = pol.transform(range_x)
17    range_x_p = scal.transform(range_x_p)
18    y_pred: np.ndarray = model.predict(range_x_p)
19    plot_linear_data(x, y, x_i, y_i, range_x, y_pred, 'overfitting')
20    test_cost, train_cost = test(
21        x_test, y_test, x_train_aux, y_train, pol, scal, model)
22
23    print(f"Train cost: {train_cost}")
24    print(f"Test cost: {test_cost}")

```

Figura 1.2: Función *overfitting*

```

1 def train(x_train: np.ndarray, y_train: np.ndarray, grado: int) -> tuple[sp.PolynomialFeatures
  , sp.StandardScaler, lm.LinearRegression, np.ndarray]:
2     """ Trains a model given the training data with polynomial features
3
4     Args:
5         x_train (np.ndarray): x values of the training data
6         y_train (np.ndarray): y values of the training data
7         grado (int): degree of the polynomial
8
9     Returns:
10        tuple[sp.PolynomialFeatures, sp.StandardScaler, lm.LinearRegression, np.ndarray]:
11        _description_
12
13    poly: sp.PolynomialFeatures = sp.PolynomialFeatures(
14        degree=grado, include_bias=False)
15    x_train = poly.fit_transform(x_train[:, None])
16    scal: sp.StandardScaler = sp.StandardScaler()
17    x_train = scal.fit_transform(x_train)
18    model: lm.LinearRegression = lm.LinearRegression()
19    model.fit(x_train, y_train)
20    return poly, scal, model, x_train

```

Figura 1.3: Función *train*

```

1 def test(x_test: np.ndarray, y_test: np.ndarray, x_train_aux: np.ndarray, y_train: np.ndarray,
  poly: sp.PolynomialFeatures, scal: sp.StandardScaler, model: Union[lm.LinearRegression,
  lm.Ridge]) -> tuple[float, float]:
2     """Tests the model with the test data
3
4     Args:
5         x_test (np.ndarray): x values of the test data
6         y_test (np.ndarray): y values of the test data
7         x_train_aux (np.ndarray): x values of the training data
8         y_train (np.ndarray): y values of the training data
9         poly (sp.PolynomialFeatures): polynomial features
10        scal (sp.StandardScaler): standard scaler
11        model (Union[lm.LinearRegression, lm.Ridge]): model to test
12
13    Returns:
14        tuple[float, float]: test cost, train cost
15
16    x_test = poly.transform(x_test[:, None])
17    x_test = scal.transform(x_test)
18
19    y_pred_test: np.ndarray = model.predict(x_test)
20    test_cost: float = cost(y_test, y_pred_test)
21
22    y_pred_train: np.ndarray = model.predict(x_train_aux)
23    train_cost: float = cost(y_train, y_pred_train)
24
25    return test_cost, train_cost

```

Figura 1.4: Función *test*

```

1 def cost(y: np.ndarray, y_hat: np.ndarray) -> float:
2     """Calculates the cost of the model
3
4     Args:
5         y (np.ndarray): real values
6         y_hat (np.ndarray): predicted values
7
8     Returns:
9         float: cost of the model
10
11    return np.mean((y_hat - y)**2) / 2

```

Figura 1.5: Función *cost*

```

1 def plot_linear_data(x: np.ndarray, y: np.ndarray, x_ideal: np.ndarray, y_ideal: np.ndarray,
2   model_range: np.ndarray, model: np.ndarray, name: str) -> None:
3     """Plots the dataset, the ideal data and the model
4
5     Args:
6         x (np.ndarray): x values of the dataset with noise
7         y (np.ndarray): y values of the dataset with noise
8         x_ideal (np.ndarray): x ideal values of the dataset
9         y_ideal (np.ndarray): y ideal values of the dataset
10        model_range (np.ndarray): x values of the model
11        model (np.ndarray): y values of the model
12        name (str): file name
13    """
14    plt.plot(x, y, 'o', label='Train data', markersize=2)
15    plt.plot(x_ideal, y_ideal, label='Ideal data',
16             c='red', linestyle='dashed', linewidth=1.5)
17    plt.plot(model_range, model, label='Predictions', c='blue', linewidth=1)
18    plt.legend()
19    plt.savefig(f'{plot_folder}/{name}.png', dpi=300)
20    plt.clf()

```

Figura 1.6: Función *plot_linear_data*

2. Elección del grado del polinomio usando un conjunto de validación

En este apartado vamos a escoger el grado del polinomio basándonos en el menor coste de validación de grados entre 1 y 10. Primero dividimos los datos en 60 % de entrenamiento, 20 % de validación y 20 % de test. Entrenamos el modelo con las funciones del apartado anterior (1.3, 1.4) y escogemos el grado que menos coste de validación nos de.

La función que lleva todo este proceso es *seleccion_grado* (figura 2.2). La figura 2.1 nos muestra en la comparativa entre el modelo del grado escogido y la recta ideal. El grado escogido es 2.

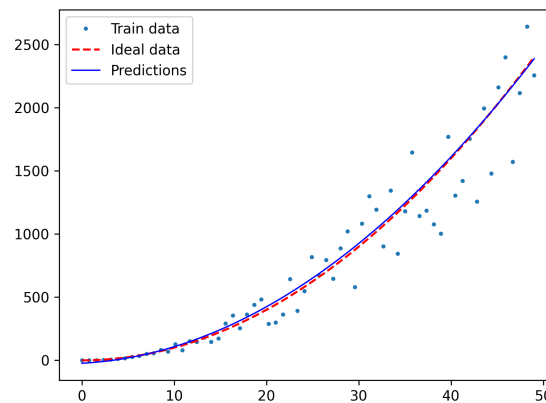


Figura 2.1: Gráfica de la selección del grado

```

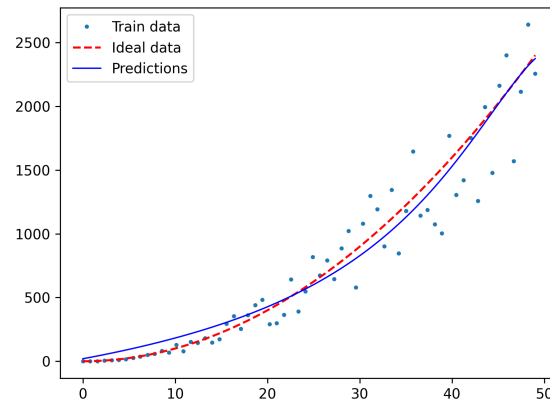
1 def seleccion_grado(x: np.ndarray, y: np.ndarray, x_i: np.ndarray, y_i: np.ndarray) -> None:
2     """Selects the best degree for the model
3
4     Args:
5         x (np.ndarray): x values of the dataset with noise
6         y (np.ndarray): y values of the dataset with noise
7         x_i (np.ndarray): x ideal values of the dataset
8         y_i (np.ndarray): y ideal values of the dataset
9     """
10    print("Selección de grado")
11    x_train, x_test, y_train, y_test = ms.train_test_split(
12        x, y, test_size=0.4, random_state=RANDOM_STATE)
13    x_test, x_cv, y_test, y_cv = ms.train_test_split(
14        x_test, y_test, test_size=0.5, random_state=RANDOM_STATE)
15    min_cost: float = 0
16    min_grado: float = 0
17    models: np.ndarray = np.empty(10, dtype=object)
18    for grado in range(10):
19        pol, scal, model, x_train_aux = train(x_train, y_train, grado + 1)
20        cv_cost, train_cost = test(
21            x_cv, y_cv, x_train_aux, y_train, pol, scal, model)
22        models[grado] = (pol, scal, model, x_train_aux)
23        if min_cost == 0 or cv_cost < min_cost:
24            min_cost = cv_cost
25            min_grado = grado + 1
26    print(f"Grado seleccionado: {min_grado}")
27
28    x_range: np.ndarray = np.linspace(np.min(x), np.max(x), 1000)
29    x_range: np.ndarray = x_range[:, None]
30    x_range_p: np.ndarray = models[min_grado - 1][0].transform(x_range)
31    x_range_p: np.ndarray = models[min_grado - 1][1].transform(x_range_p)
32    y_pred: np.ndarray = models[min_grado - 1][2].predict(x_range_p)
33
34    plot_linear_data(x, y, x_i, y_i, x_range, y_pred, 'grado')
35
36    test_cost, train_cost = test(
37        x_test, y_test, models[min_grado - 1][3], y_train, models[min_grado - 1][0], models[
38        min_grado - 1][1], models[min_grado - 1][2])
39    print(f"Train cost: {train_cost}")
40    print(f"CV cost: {min_cost}")
41    print(f"Test cost: {test_cost}")

```

Figura 2.2: Función *seleccion_grado*

3. Elección del parámetro λ

Como en este apartado usamos regularización, incluimos la función *train_reg* (figura 3.2) para entrenar el modelo con regularización, la función de test se mantiene igual. La función *seleccion_lambda* (figura 3.3) nos ayudará a escoger el mejor λ entre $[1e-6, 1e-5, 1e-4, 1e-3, 1e-2, 1e-1, 1, 10, 100, 300, 600, 900]$. La figura 3.1 nos muestra la comparativa entre el modelo con el λ escogido y la recta ideal. El λ escogido es 10.

Figura 3.1: Gráfica de la selección de λ

```

1 def train_reg(x_train: np.ndarray, y_train: np.ndarray, grado: int, l: float) -> tuple[sp.
  PolynomialFeatures, sp.StandardScaler, lm.Ridge, np.ndarray]:
2     """ Trains a model given the training data with polynomial features and regularization
3
4     Args:
5         x_train (np.ndarray): x values of the training data
6         y_train (np.ndarray): y values of the training data
7         grado (int): degree of the polynomial
8         l (float): lambda value for the regularization
9
10    Returns:
11        tuple[sp.PolynomialFeatures, sp.StandardScaler, lm.Ridge, np.ndarray]: _description_
12    """
13    poly: sp.PolynomialFeatures = sp.PolynomialFeatures(
14        degree=grado, include_bias=False)
15    x_train = poly.fit_transform(x_train[:, None])
16    scal: sp.StandardScaler = sp.StandardScaler()
17    x_train = scal.fit_transform(x_train)
18    model: lm.Ridge = lm.Ridge(alpha=l)
19    model.fit(x_train, y_train)
20    return poly, scal, model, x_train

```

Figura 3.2: Función *train_reg*

```

1 def seleccion_lambda(x: np.ndarray, y: np.ndarray, x_i: np.ndarray, y_i: np.ndarray) -> None:
2     """Selects the best lambda for the model
3
4     Args:
5         x (np.ndarray): x values of the dataset with noise
6         y (np.ndarray): y values of the dataset with noise
7         x_i (np.ndarray): x ideal values of the dataset
8         y_i (np.ndarray): y ideal values of the dataset
9     """
10    print("Selección de lambda")
11    lambdas: list[float] = [1e-6, 1e-5, 1e-4, 1e-3,
12                           1e-2, 1e-1, 1, 10, 100, 300, 600, 900]
13    alpha: float = 0
14    min_cost: float = -1
15    x_train, x_test, y_train, y_test = ms.train_test_split(
16        x, y, test_size=0.4, random_state=RANDOM_STATE)
17    x_test, x_cv, y_test, y_cv = ms.train_test_split(
18        x_test, y_test, test_size=0.5, random_state=RANDOM_STATE)
19    models: np.ndarray = np.empty(len(lambdas), dtype=object)
20
21    for l in lambdas:
22        pol, scal, model, x_train_aux = train_reg(x_train, y_train, 15, l)
23        test_cost, train_cost = test(
24            x_cv, y_cv, x_train_aux, y_train, pol, scal, model)
25        models[lambdas.index(l)] = (pol, scal, model, x_train_aux)
26        if min_cost == -1 or test_cost < min_cost:
27            min_cost = test_cost
28            alpha = l
29        print(f"Lambda: {l}-> Cost: {test_cost}")
30    print(f"Lambda seleccionado: {alpha}")
31
32    x_range: np.ndarray = np.linspace(np.min(x), np.max(x), 1000)
33    x_range = x_range[:, None]
34    x_range_p: np.ndarray = models[lambdas.index(alpha)][0].transform(x_range)
35    x_range_p = models[lambdas.index(alpha)][1].transform(x_range_p)
36    y_pred: np.ndarray = models[lambdas.index(alpha)][2].predict(x_range_p)
37
38    plot_linear_data(x, y, x_i, y_i, x_range, y_pred, 'lambda')
39
40    test_cost, train_cost = test(
41        x_test, y_test, models[lambdas.index(alpha)][3], y_train, models[lambdas.index(alpha)]
42        [0], models[lambdas.index(alpha)][1], models[lambdas.index(alpha)][2])
43    print(f"Train cost: {train_cost}")
44    print(f"CV cost: {min_cost}")
45    print(f"Test cost: {test_cost}")

```

Figura 3.3: Función *seleccion_lambda*

4. Elección de hiperparámetros

En este apartado vamos a escoger grado y regularización mirando el mínimo coste de validación. Para ello usaremos los posibles parámetros de los dos apartados anteriores. La función *seleccion_hiperparametros* (figura 4.2) nos ayudará a escoger el mejor grado y λ . La figura 4.1 nos muestra la comparativa entre el modelo con los hiperparámetros escogidos y la recta ideal. Los hiperparámetros escogidos son grado 12 y λ 1e-6.

Podemos ver los resultados por λ (vertical) y grado (horizontal) en las tablas 4.1 4.2 y 4.3.

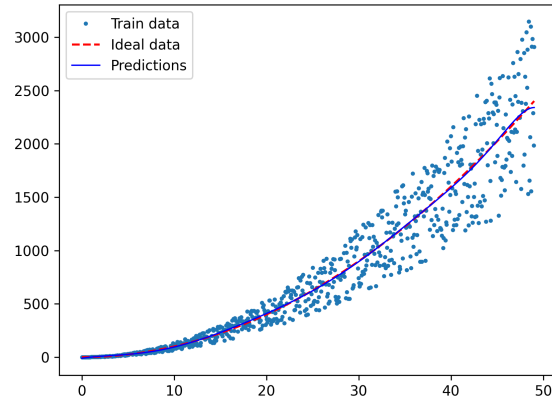


Figura 4.1: Gráfica de la selección de hiperparámetros

0.00000e+00	1.00000e+00	2.00000e+00	3.00000e+00	4.00000e+00
1.00000e-06	6.94186e-310	3.73427e+04	2.00062e+04	2.00069e+04
1.00000e-05	6.94186e-310	3.73427e+04	2.00062e+04	2.00069e+04
1.00000e-04	5.43817e-310	3.73427e+04	2.00062e+04	2.00069e+04
1.00000e-03	5.43817e-310	3.73426e+04	2.00062e+04	2.00069e+04
1.00000e-02	5.88785e+00	3.73419e+04	2.00060e+04	2.00073e+04
1.00000e-01	4.35701e+01	3.73343e+04	2.00046e+04	2.00140e+04
1.00000e+00	1.71402e+01	3.72597e+04	2.00104e+04	2.00877e+04
1.00000e+01	5.49533e+00	3.66353e+04	2.06875e+04	2.01468e+04
1.00000e+02	3.82710e+01	3.87107e+04	2.41588e+04	2.08719e+04
3.00000e+02	5.56075e+00	6.21938e+04	3.48534e+04	2.68364e+04
6.00000e+02	3.70935e+01	9.73603e+04	5.62252e+04	4.08061e+04
9.00000e+02	4.34393e+01	1.23253e+05	7.65306e+04	5.58958e+04

Cuadro 4.1: Resultados de los hiperparámetros

0.00000e+00	5.00000e+00	6.00000e+00	7.00000e+00	8.00000e+00
1.00000e-06	2.00187e+04	2.00036e+04	2.00037e+04	2.00015e+04
1.00000e-05	2.00187e+04	2.00038e+04	2.00029e+04	2.00026e+04
1.00000e-04	2.00185e+04	2.00053e+04	2.00024e+04	2.00020e+04
1.00000e-03	2.00171e+04	2.00116e+04	2.00048e+04	2.00020e+04
1.00000e-02	2.00093e+04	2.00118e+04	2.00104e+04	2.00063e+04
1.00000e-01	2.00036e+04	1.99979e+04	1.99980e+04	2.00001e+04
1.00000e+00	2.00724e+04	2.00344e+04	2.00135e+04	2.00049e+04
1.00000e+01	2.03006e+04	2.03256e+04	2.02652e+04	2.01927e+04
1.00000e+02	2.09946e+04	2.16251e+04	2.20570e+04	2.22388e+04
3.00000e+02	2.51627e+04	2.56445e+04	2.66362e+04	2.75402e+04
6.00000e+02	3.51190e+04	3.36068e+04	3.38446e+04	3.47103e+04
9.00000e+02	4.66041e+04	4.27799e+04	4.16433e+04	4.17849e+04

Cuadro 4.2: Resultados de los hiperparámetros

0.00000e+00	9.00000e+00	1.00000e+01	1.10000e+01	1.20000e+01
1.00000e-06	1.99926e+04	1.99706e+04	1.99346e+04	1.99032e+04
1.00000e-05	1.99986e+04	1.99872e+04	1.99680e+04	1.99431e+04
1.00000e-04	2.00011e+04	1.99973e+04	1.99879e+04	1.99720e+04
1.00000e-03	2.00008e+04	1.99993e+04	1.99964e+04	1.99909e+04
1.00000e-02	2.00023e+04	1.99994e+04	1.99968e+04	1.99938e+04
1.00000e-01	2.00016e+04	2.00013e+04	1.99994e+04	1.99962e+04
1.00000e+00	2.00019e+04	2.00014e+04	2.00021e+04	2.00035e+04
1.00000e+01	2.01386e+04	2.01060e+04	2.00893e+04	2.00820e+04
1.00000e+02	2.22515e+04	2.21748e+04	2.20634e+04	2.19493e+04
3.00000e+02	2.81917e+04	2.85895e+04	2.87844e+04	2.88343e+04
6.00000e+02	3.56910e+04	3.65691e+04	3.72692e+04	3.77821e+04
9.00000e+02	4.24743e+04	4.33353e+04	4.41806e+04	4.49245e+04

Cuadro 4.3: Resultados de los hiperparámetros

```

1 def seleccion_hiperparametros(x: np.ndarray, y: np.ndarray, x_i: np.ndarray, y_i: np.ndarray)
  -> None:
2     """Selects the best hyperparameters for the model
3
4     Args:
5         x (np.ndarray): x values of the dataset with noise
6         y (np.ndarray): y values of the dataset with noise
7         x_i (np.ndarray): x ideal values of the dataset
8         y_i (np.ndarray): y ideal values of the dataset
9     """
10    print("Selección de hiperparametros")
11    x_train, x_test, y_train, y_test = ms.train_test_split(
12        x, y, test_size=0.4, random_state=RANDOM_STATE)
13    x_test, x_cv, y_test, y_cv = ms.train_test_split(
14        x_test, y_test, test_size=0.5, random_state=RANDOM_STATE)
15
16    lambdas: list[float] = [1e-6, 1e-5, 1e-4, 1e-3,
17                            1e-2, 1e-1, 1, 10, 100, 300, 600, 900]
18
19    models: np.ndarray = np.empty((16, len(lambdas)), dtype=object)
20
21    min_cost: float = -1
22    elec_lambda: float = 0
23    eled_grado: int = 0
24    costs = np.empty((16, len(lambdas)))
25
26    for i in range(1, 16):
27        for l in lambdas:
28            pol, scal, model, x_train_aux = train_reg(x_train, y_train, i, l)
29            models[i][lambdas.index(l)] = (pol, scal, model, x_train_aux)
30            cv_cost, train_cost = test(
31                x_cv, y_cv, x_train_aux, y_train, pol, scal, model)
32            costs[i][lambdas.index(l)] = cv_cost
33            if min_cost == -1 or cv_cost < min_cost:
34                min_cost = cv_cost
35                elec_lambda = l
36                eled_grado = i
37            print(f"Grado: {i} Lambda: {l}-> Cost: {cv_cost}")
38    print(f"Grado seleccionado: {eled_grado}")
39    print(f"Lambda seleccionado: {elec_lambda}")
40
41    costs = np.append(np.array(lambdas)[None, :], costs, axis=0)
42    costs = np.append(np.array(range(0, 17))[None, :], costs.T, axis=0)
43    write_csv(costs, 'hiperparametros')
44
45    x_range: np.ndarray = np.linspace(np.min(x), np.max(x), 10000)
46    x_range = x_range[:, None]
47    x_range_p: np.ndarray = models[eled_grado][lambdas.index(
48        elec_lambda)][0].transform(x_range)
49    x_range_p = models[eled_grado][lambdas.index(
50        elec_lambda)][1].transform(x_range_p)
51    y_pred: np.ndarray = models[eled_grado][lambdas.index(
52        elec_lambda)][2].predict(x_range_p)
53
54    plot_linear_data(x, y, x_i, y_i, x_range, y_pred,
55                    'hiperparametros')
56
57    test_cost, train_cost = test(
58        x_test, y_test, models[eled_grado][lambdas.index(elec_lambda)][3], y_train, models[
59        eled_grado][lambdas.index(elec_lambda)][0], models[eled_grado][lambdas.index(elec_lambda)
60        ][1], models[eled_grado][lambdas.index(elec_lambda)][2])
61    print(f"Train cost: {train_cost}")
62    print(f"CV cost: {min_cost}")
63    print(f"Test cost: {test_cost}")

```

Figura 4.2: Función *seleccion_hiperparametros*

5. Curvas de aprendizaje

En este apartado vamos a hacer una comparativa de coste de validación y de entrenamiento variando el tamaño de muestra de entrenamiento. Para ello usaremos la función *learning_curve* (figura 5.2). La figura 5.1

nos muestra la comparativa entre el coste de validación y de entrenamiento. Para pintar la gráfica se ha usado la función `draw_learning_curve` (figura 5.3).

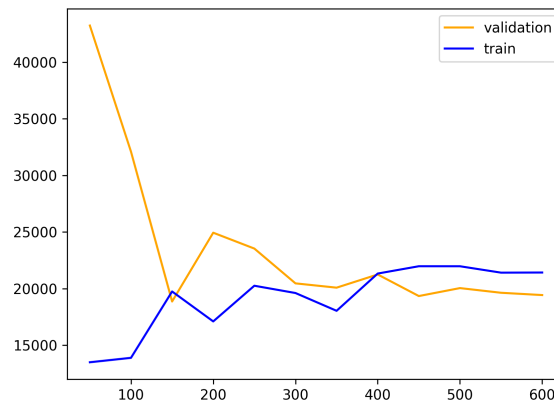


Figura 5.1: Gráfica de las curvas de aprendizaje

```

1 def learning_curve():
2     """Plots the learning curve of the model based on the number of training sample sizes
3     """
4     print("Learning curve")
5     J_cost_train: list[float] = []
6     J_const_val: list[float] = []
7     params: list[int] = range(50, 601, 50)
8     x, y, x_i, y_i = gen_data(1000)
9     x_train_t, x_test, y_train_t, y_test = ms.train_test_split(
10         x, y, test_size=0.4, random_state=RANDOM_STATE)
11     x_test, x_cv, y_test, y_cv = ms.train_test_split(
12         x_test, y_test, test_size=0.5, random_state=RANDOM_STATE)
13
14     for i in range(len(params)):
15         indexes = np.linspace(
16             0, len(x_train_t) - 1, params[i], dtype=int)
17         x_train = x_train_t[indexes]
18         y_train = y_train_t[indexes]
19         pol, scal, model, x_train_aux = train(x_train, y_train, 16)
20         cv_cost, train_cost = test(
21             x_cv, y_cv, x_train_aux, y_train, pol, scal, model)
22         J_cost_train.append(train_cost)
23         J_const_val.append(cv_cost)
24
25     draw_learning_curve(params, J_cost_train, params, J_const_val)
26
27     print(f'Mejor valor de entrenamiento: {np.min(J_cost_train)}')
28     print(f'Mejor valor de validacion: {np.min(J_const_val)}')

```

Figura 5.2: Función `learning_curve`

```

1 def draw_learning_curve(x: np.ndarray, y: np.ndarray, x_v: np.ndarray, y_v: np.ndarray):
2     """Plots the learning curve of the model based on the number of training sample sizes
3
4     Args:
5         x (np.ndarray): number of samples
6         y (np.ndarray): cost of the training set
7         x_v (np.ndarray): number of samples
8         y_v (np.ndarray): cost of the validation set
9     """
10    plt.figure()
11    plt.plot(x_v, y_v, c='orange', label='validation')
12    plt.plot(x, y, c='blue', label='train')
13    plt.legend()
14    plt.savefig(f'{plot_folder}/learning_curve.png', dpi=300)

```

Figura 5.3: Función *draw_learning_curve*

6. Utilidades

Algunas utilidades que se han usado en el código son las siguientes:

- *test_overfitting* (figura 6.1): Función que testea el sobreajuste.
- *test_seleccion_grado* (figura 6.2): Función que testea la selección del grado.
- *test_regularización* (figura 6.3): Función que testea la selección de λ .
- *test_hiperparametros* (figura 6.4): Función que testea la selección de hiperparámetros.
- *test_learning_curve* (figura 6.5): Función que testea las curvas de aprendizaje.
- *write_csv* (figura 6.6): Función que escribe los datos en un archivo *csv*.
- *main* (figura 6.7): Función que ejecuta todas las pruebas.
- *CommandLine* (figura ??): Clase que nos permite añadirle parámetros a la ejecución del programa.

```

1 def test_overfitting(x: np.ndarray, y: np.ndarray, x_i: np.ndarray, y_i: np.ndarray,
2     commandLine: CommandLine) -> None:
3     """Tests the overfitting of the model
4
5     Args:
6         x (np.ndarray): x values of the dataset with noise
7         y (np.ndarray): y values of the dataset with noise
8         x_i (np.ndarray): x ideal values of the dataset
9         y_i (np.ndarray): y ideal values of the dataset
10        commandLine (CommandLine): command line arguments
11    """
12    if (not os.path.exists(f'{plot_folder}/overfitting.png') or commandLine.overfitting or
13        commandLine.all):
14        overfitting(x, y, x_i, y_i)
15    else:
16        if (commandLine.interactive):
17            answer = ''
18            while (answer != 'y' and answer != 'n'):
19                print("Recreate overfitting test? [y/n]")
20                answer = input()
21            if (answer == 'y'):
22                overfitting(x, y, x_i, y_i)

```

Figura 6.1: Función *test_overfitting*

```

1 def test_seleccion_grado(x: np.ndarray, y: np.ndarray, x_i: np.ndarray, y_i: np.ndarray,
2   commandLine: CommandLine) -> None:
3     """Selects the best degree for the model
4     """
5     if (not os.path.exists(f'{plot_folder}/grado.png') or commandLine.grado or commandLine.all
6     ):
7         seleccion_grado(x, y, x_i, y_i)
8     else:
9         if (commandLine.interactive):
10            answer = ''
11            while (answer != 'y' and answer != 'n'):
12                print("Recreate grado test? [y/n]")
13                answer = input()
14            if (answer == 'y'):
15                seleccion_grado(x, y, x_i, y_i)

```

Figura 6.2: Función *test_seleccion_grado*

```

1 def test_regularizacion(x: np.ndarray, y: np.ndarray, x_i: np.ndarray, y_i: np.ndarray,
2   commandLine: CommandLine) -> None:
3     """Selects the best lambda for the model
4     """
5     if (not os.path.exists(f'{plot_folder}/lambda.png') or commandLine.reg or commandLine.all)
6     :
7         seleccion_lambda(x, y, x_i, y_i)
8     else:
9         if (commandLine.interactive):
10            answer = ''
11            while (answer != 'y' and answer != 'n'):
12                print("Recreate regularizacion test? [y/n]")
13                answer = input()
14            if (answer == 'y'):
15                seleccion_lambda(x, y, x_i, y_i)

```

Figura 6.3: Función *test_regularizacion*

```

1 def test_hiperparametros(x: np.ndarray, y: np.ndarray, x_i: np.ndarray, y_i: np.ndarray,
2   commandLine: CommandLine) -> None:
3     """Selects the best hyperparameters for the model
4     """
5     if (not os.path.exists(f'{plot_folder}/hiperparametros.png') or commandLine.hyperparam or
6     commandLine.all):
7         seleccion_hiperparametros(x, y, x_i, y_i)
8     else:
9         if (commandLine.interactive):
10            answer = ''
11            while (answer != 'y' and answer != 'n'):
12                print("Recreate hiperparametros test? [y/n]")
13                answer = input()
14            if (answer == 'y'):
15                seleccion_hiperparametros(x, y, x_i, y_i)

```

Figura 6.4: Función *test_hiperparametros*

```

1 def test_learning_curve(x: np.ndarray, i: np.ndarray, x_i: np.ndarray, y_i: np.ndarray,
2   commandLine: CommandLine) -> None:
3     """Plots the learning curve of the model based on the number of training sample sizes
4     """
5     if (not os.path.exists(f'{plot_folder}/learning_curve.png') or commandLine.learning_curve
6     or commandLine.all):
7         learning_curve()
8     else:
9         if (commandLine.interactive):
10            answer = ''
11            while (answer != 'y' and answer != 'n'):
12                print("Recreate learning curve test? [y/n]")
13                answer = input()
14            if (answer == 'y'):
15                learning_curve()

```

Figura 6.5: Función *test_learning_curve*

```

1 def write_csv(data: np.ndarray, name: str, split: int = 4) -> None:
2     """Writes a numpy array to a csv file
3
4     Args:
5         data (np.ndarray): data to write
6         name (str): name of the file
7         split (int, optional): number of columns to split the data. Defaults to 4.
8     """
9     np.printoptions(suppress=False)
10    # Format string to display numbers in scientific notation with 2 decimal places
11    fmt_float = '%.5e'
12    for i in range(1, len(data), split):
13        data_acc = np.concatenate((data[:, :1], data[:, i:i+split]), axis=1)
14
15        filename = f'{csv_folder}/{name}{i//split}.csv'
16        np.savetxt(filename, data_acc, delimiter=',', fmt=fmt_float)

```

Figura 6.6: Función *write_csv*

```

1 def main() -> None:
2     """Main function
3     """
4     commandLine = CommandLine()
5     commandLine.parse(sys.argv[1:])
6     x, y, x_i, y_i = gen_data(64)
7     funcs = [test_overfitting, test_seleccion_grado,
8             test_regularizacion, test_hiperparametros, test_learning_curve]
9
10    prepare_folder()
11    if not os.path.exists(f'{plot_folder}/dataset.png'):
12        plot_dataset(x, y, x_i, y_i, 'dataset')
13
14    for func in funcs:
15        if func.__name__ == 'test_hiperparametros':
16            x, y, x_i, y_i = gen_data(750)
17        func(x, y, x_i, y_i, commandLine)

```

Figura 6.7: Función *main*

```

1 def prepare_folder() -> None:
2     """Creates the folder to save the plots
3     """
4     if not os.path.exists(plot_folder):
5         os.makedirs(plot_folder)

```

Figura 6.8: Función *prepare_folder*

```

1 import argparse
2
3
4 class CommandLine:
5     interactive: bool = False
6     overfitting: bool = False
7     grado: bool = False
8     reg: bool = False
9     hyperparam: bool = False
10    learning_curve: bool = False
11    all: bool = False
12
13    def __init__(self):
14        self.parser = argparse.ArgumentParser(
15            description='Practica 6 - Aprendizaje Automatico')
16        self.parser.add_argument('-I', "--Interactive",
17                                help='Interactive mode', required=False, default="", action='
18store_true')
19        self.parser.add_argument('-O', "--Overfitting", help='runs the Overfitting test',
20                                required=False, default="", action='store_true')
21        self.parser.add_argument(
22            '-G', "--Grado", help='Grado', required=False, default="", action='store_true')
23        self.parser.add_argument('-R', "--Regularization", help='runs the Regularization test'
24                                ,
25                                required=False, default="", action='store_true')
26        self.parser.add_argument('-HP', "--Hyperparam", help='runs the Hyperparam test',
27                                required=False, default="", action='store_true')
28        self.parser.add_argument('-L', "--LearningCurve", help='runs the LearningCurve test',
29                                required=False, default="", action='store_true')
30        self.parser.add_argument('-A', "--All", help='runs all tests',
31                                required=False, default="", action='store_true')
32
33    def parse(self, sysargs):
34        args = self.parser.parse_args(sysargs)
35        if args.Interactive:
36            self.interactive = True
37        if args.Overfitting:
38            self.overfitting = True
39        if args.Grado:
40            self.grado = True
41        if args.Regularization:
42            self.reg = True
43        if args.Hyperparam:
44            self.hyperparam = True
45        if args.LearningCurve:
46            self.learning_curve = True
47        if args.All:
48            self.all = True

```

Figura 6.9: Clase *CommandLine*