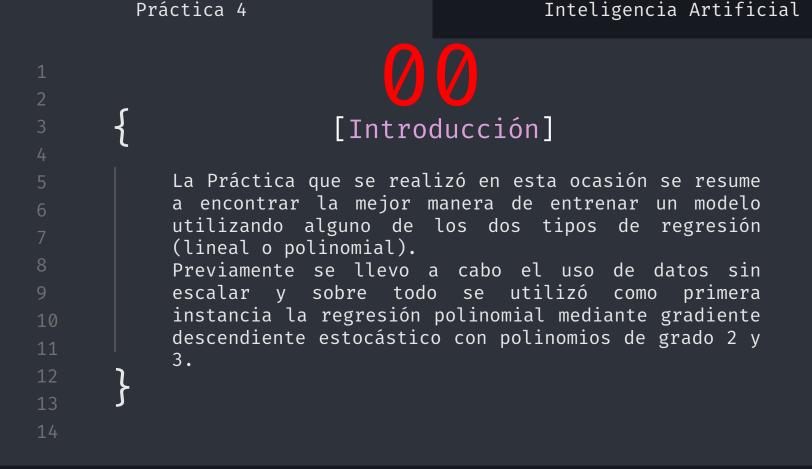
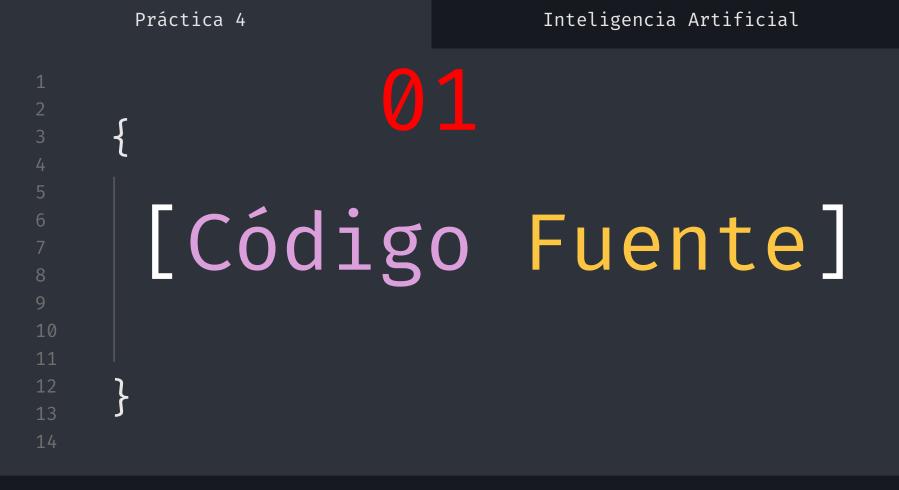
Inteligencia Artificial

Práctica 4



Inteligencia Artificial

Práctica 4



11

12

13

14

Código Fuente

Importación de las librerías que se ocuparan a lo largo de la programación de la regresión polinomial

```
In [1]:

# Importación de librerias
import pandas as pd
import numpy as np
import operator
import os
from sklearn import preprocessing
from sklearn.linear_model import LinearRegression
from sklearn.linear_model import SGDRegressor
from sklearn.model_selection import KFold
from sklearn.model_selection import train_test_split
from sklearn.metrics import mean_squared_error
from sklearn.metrics import r2_score
from sklearn.preprocessing import PolynomialFeatures
import matplotlib.pyplot as plt
```

```
Creación del dataset "cal_housing.csv" en un dataframe para poder trabajar de mejor manera con él
```

```
In [2]: 1 # Creación del archivo a trabajar en un dataframe
2 dataframe = pd.read_csv('./cal_housing.csv',sep = ',', engine = 'python')
```


Código Fuente

Mediante el dataframe ya creado hacemos la debida identificación del corpus y las etiquetas

Creación del conjunto de entrenamiento y el conjunto de prueba

```
In [4]: 1 # Generación Data Test y Data Train
2
3 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.2, random_state = 0)
```

Creamos una clase validación

```
1 class validation set:
     def __init__(self, X_train, y_train, X_test, y_test):
        self.X_train = X_train
        self.y_train = y_train
        self.X test = X test
        self.y test = y test
10 validation sets = []
11 kf = KFold(n splits = 10)
12 for train index, test index in kf.split(X train):
     X_train_v, X_test_v = X_train[train_index], X_train[test_index]
     y_train_v, y_test_v = y_train[train_index], y_train[test_index]
15
     validation sets.append(validation set(X train v, y train v, X test v, y test v))
18
19 regression_type_list = list()
20 scale_type_list = list()
21 | learning rate list = list()
22 eta0 list = list()
23 iterations_list = list()
24 mse list = list()
25 r2 list = list()
```

Creamos la "clase" o más bien definimos una variable para la regresión lineal

```
In [6]:
         1 # Definición de la regresión Lineal
            def linear regression(X train, X test, y train, y test, eta0 = 0.001, learning rate = 'constant', iterations = 10000, standa
                if standar scale:
                    X train = preprocessing.StandardScaler().fit transform(X train)
                    X test = preprocessing.StandardScaler().fit transform(X test)
                if robust scale:
                    X train = preprocessing.RobustScaler().fit transform(X train)
                    X test = preprocessing.RobustScaler().fit transform(X test)
                regr = SGDRegressor(learning_rate = learning_rate, eta0 = eta0, max iter = iterations)
                regr.fit(X train, y train)
         11
         12
                y test pred = regr.predict(X test)
                mse = mean squared error(y test, y test pred)
         13
                r2 = r2 score(y test, y test pred)
                return [mse, r2]
```

Creamos la "clase" o más bien definimos una variable para la regresión polinomial

```
1 # Definición de la regresión Polinomial
In [7]:
          3 def polynomial regression(X train, X test, y train, y test, degree = 2, eta0 = 0.001, learning rate = 'constant', iterations
                 polynomial features= PolynomialFeatures(degree = degree)
                 X train = polynomial features.fit transform(X train)
                 X test = polynomial features.fit transform(X test)
                 if standar scale:
                    X train = preprocessing.StandardScaler().fit transform(X train)
          9
                    X test = preprocessing.StandardScaler().fit transform(X test)
                if robust scale:
         10
                     X train = preprocessing.RobustScaler().fit transform(X train)
         11
         12
                     X test = preprocessing.RobustScaler().fit transform(X test)
                 regr = SGDRegressor(learning rate = learning rate, eta0 = eta0, max iter = iterations)
         13
         14
                 regr.fit(X train, y train)
                y_test_pred = regr.predict(X_test)
                 mse = mean squared error(y test, y test pred)
         17
                 r2 = r2_score(y_test, y_test_pred)
                 return [mse, r2]
```

```
In [8]: 1 # Definición para obtener los mejores parametros
         3 def get_best_params(validation_sets, eta0, learnin_rate, iterations, type_regression, scale_type):
               c = 1
               mses = list()
               r2s = list()
               n = len(validation_sets)
               for validation set in validation sets:
                  if scale type == 'none':
        10
                      if type regression == 1:
                          mse, r2 = linear regression(validation set.X train, validation set.X test, validation set.y train, validatio
        12
                          mse, r2 = polynomial regression(validation set.X train, validation set.X test, validation set.y train, valid
        14
                   elif scale type == 'std':
                      if type_regression == 1:
        16
                          mse, r2 = linear regression(validation set.X train, validation set.X test, validation set.y train, validatio
        18
                          mse, r2 = polynomial regression(validation set.X train, validation set.X test, validation set.y train, valid
        19
                   elif scale_type == 'robust':
                       if type regression == 1:
                          mse, r2 = linear regression(validation set.X train, validation set.X test, validation set.y train, validatio
        22
                          mse, r2 = polynomial regression(validation set.X train, validation set.X test, validation set.y train, valid
        24
                   print('kfold:', c)
        25
                   print('\tmse:', mse)
        26
                   print('\tr2:', r2)
                   c = c + 1
        28
                   mses.append(mse)
        29
                  r2s.append(r2)
               mses mean = sum(mses) / n
        31
               r2s mean = sum(r2s) / n
        32
               return [mses_mean, r2s_mean]
        35
        36 | scale_types = ['none', 'std', 'robust']
        37 type resgressions = [1, 2, 3]
        38 # learning_rates = ['constant', 'optimal', 'invscaling', 'adaptive']
        39 | learning_rates = ['constant']
        40 eta0s = [0.0001, 0.00001, 0.000001]
        41 iterationss = [200000, 300000]
```

Definimos
una variable
para obtener
los mejores
parámetros,
los tipos de
escalamiento
y el
escalamiento

```
44
45 for scale type in scale types:
      for type resgression in type resgressions:
47
          for learning rate in learning rates:
48
             for eta0 in eta0s:
                 for iterations in iterationss:
49
50
                    print(scale type)
                    print('\t', type resgression)
51
                    print('\t\t', learning rate)
52
53
                    print('\t\t\t', eta0)
                    print('\t\t\t', iterations)
54
                    mses_mean, r2s_mean = get_best_params(validation_sets, eta0, learning_rate, iterations, type_resgression
55
                    scale type list.append(scale type)
56
57
                    regression type list.append(type resgression)
58
                    learning rate list.append(learning rate)
59
                    eta0 list.append(eta0)
60
                    iterations list.append(iterations)
61
                    mse_list.append(mses_mean)
                    r2_list.append(r2s_mean)
62
63
                    print("")
64
                    print("")
```

Práctica 4

Ahora
definimos
nuestra
tabla para
poder
obtener los
datos con un
orden

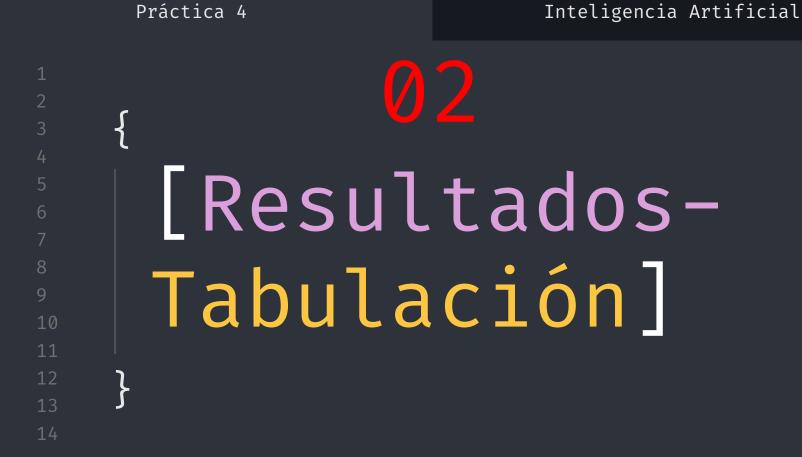
Ahora bien,
dentro de
esta tabla,
no vemos un
orden
lógico,
ahora
actuamos a
ordenarlo
por los
valores de
mse

```
Código Fuente

In [11]: 1 mse, r2 = polynomial_regression(X_train, X_t
```

```
In [11]: 1 mse, r2 = polynomial_regression(X_train, X_test, y_train, y_test, degree = 3, eta0 = 0.00001, learning_rate = 'constant', it

In [12]: 1 print('mse con X_train y X_test:', mse)
2 print('r2 con X_train y X_test:', r2)
```



```
none
                 constant
                         0.0001
                                 200000
kfold: 1
        mse: 7.192700199140325e+29
        r2: -5.18121508822569e+19
kfold: 2
        mse: 5.56001709344023e+28
        r2: -4.138950193827837e+18
kfold: 3
        mse: 1.6493170640592684e+30
        r2: -1.2242896921727535e+20
kfold: 4
        mse: 1.0782327142660636e+31
        r2: -7.971556722470891e+20
kfold: 5
        mse: 1.309845956394241e+29
```

```
r2: -2.1361167949617705e+19
kfold: 10
        mse: 2.604410766671443e+31
        r2: -1.9599077750244953e+21
none
                 constant
                         0.0001
                                  300000
kfold: 1
        mse: 4.1685799386679566e+29
        r2: -3.002809609287313e+19
kfold: 2
        mse: 8.256949510862957e+30
        r2: -6.146582322333697e+20
kfold: 3
        mse: 9.884170936495558e+29
        r2: -7.337029887656681e+19
```

```
r2: -3.742898020088328e+19
none
                 constant
                         1e-05
                                 200000
kfold: 1
        mse: 7.750716008430426e+27
        r2: -5.58317816892074e+17
kfold: 2
        mse: 1.6486369597457387e+25
        r2: -1227267137747102.5
kfold: 3
        mse: 2.3522516947697042e+26
        r2: -1.746078765604132e+16
kfold: 4
        mse: 5.258357963208061e+26
```

```
none
                 constant
                         1e-05
                                 300000
kfold: 1
        mse: 9.613248908153987e+26
       r2: -6.924841702112051e+16
kfold: 2
        mse: 1.5117898080636857e+27
        r2: -1.125396309751364e+17
kfold: 3
        mse: 4.6675715271335034e+26
        r2: -3.4647429731214104e+16
kfold: 4
        mse: 5.6349983911565445e+26
        r2: -4.1660495653494216e+16
kfold: 5
```

```
mse: 1.3067685604027836e+28
        r2: -1.0329650602820649e+18
kfold: 10
        mse: 1.603908345287333e+28
        r2: -1.206995638546246e+18
none
                 constant
                         1e-06
                                  200000
kfold: 1
        mse: 8.272678465022025e+26
        r2: -5.959170450081433e+16
kfold: 2
        mse: 2.5811463419087993e+26
        r2: -1.9214394439084796e+16
kfold: 3
        mca: 7 33913309061193a±25
```

```
kfold: 10
        mse: 6.500778935856916e+24
        r2: -489205748307702.56
none
         1
                 constant
                         1e-06
kfold: 1
        mse: 6.292768366591612e+24
        r2: -453295501063572.56
kfold: 2
        mse: 4.149308682473071e+26
        r2: -3.088800211753892e+16
kfold: 3
        mse: 1.846659452049196e+24
        r2: -137077714246940.9
kfold: 4
```

```
r2: -121142192019051.22
kfold: 8
        mse: 5.613946991027348e+24
        r2: -418664666909265.4
kfold: 9
        mse: 2.0876333292182605e+26
        r2: -1.6502174548017188e+16
kfold: 10
        mse: 3.0735165197082735e+23
        r2: -23129258259598.26
robust
                 constant
                         1e-05
                                  200000
kfold: 1
        mse: 9.117600006513666e+22
        r2 - -6567804220144 507
```

RESULTADOS - TABULACIÓN SIN ORDENAR

	Tipos de escalamiento	Tipo de regresión	Learning	eta0	Iteraciones	mse	r2
0	none	1	constant	0.000100	200000	4.620386e+30	-3.451906e+20
1	none	1	constant	0.000100	300000	1.565061e+30	-1.175406e+20
2	none	1	constant	0.000010	200000	6.243362e+27	-4.611774e+17
3	none	1	constant	0.000010	300000	8.296358e+27	-6.384893e+17
4	none	1	constant	0.000001	200000	1.264673e+26	-9.216813e+15
5	none	1	constant	0.000001	300000	1.761400e+26	-1.310978e+16
6	none	2	constant	0.000100	200000	1.705548e+47	-1.262719e+37
7	none	2	constant	0.000100	300000	2.775312e+47	-2.073150e+37
8	none	2	constant	0.000010	200000	3.883418e+45	-2.855971e+35
9	none	2	constant	0.000010	300000	1.320227e+46	-9.511375e+35
10	none	2	constant	0.000001	200000	1.524542e+43	-1.151256e+33
11	none	2	constant	0.000001	300000	1.276942e+43	-9.510445e+32
12	none	3	constant	0.000100	200000	1.122011e+66	-8.250235e+55
13	none	3	constant	0.000100	300000	1.321542e+66	-1.013416e+56
14	none	3	constant	0.000010	200000	2.516924e+63	-1.909706e+53
15	none	3	constant	0.000010	300000	3.556790e+64	-2.799108e+54
16	none	3	constant	0.000001	200000	4.245958e+61	-3.316343e+51
17	none	3	constant	0.000001	300000	2.289084e+61	-1.708496e+51
18	std	1	constant	0.000100	200000	4.823021e+09	6.393440e-01
19	std	1	constant	0.000100	300000	4.824937e+09	6.391974e-01
20	std	1	constant	0.000010	200000	4.822685e+09	6.393663e-01
21	std	1	constant	0.000010	300000	4.822069e+09	6.394121e-01
22	std	1	constant	0.000001	200000	4.821078e+09	6.394871e-01
23	std	1	constant	0.000001	300000	4.820945e+09	6.394966e-01
24	std	2	constant	0.000100	200000	4.430263e+09	6.687677e-01
25	std	2	constant	0.000100	300000	4.411492e+09	6.701361e-01
26	std	2	constant	0.000010	200000	4.433751e+09	6.684951e-01
27	std	2	constant	0.000010	300000	4.427348e+09	6.689693e-01

28	std	2	constant	0.000001	200000	4.413337e+09	6.700165e-01
29	std	2	constant	0.000001	300000	4.412737e+09	6.700638e-01
30	std	3	constant	0.000100	200000	2.256206e+18	-1.664487e+08
31	std	3	constant	0.000100	300000	3.878283e+19	-2.861396e+09
32	std	3	constant	0.000010	200000	4.264831e+09	6.809793e-01
33	std	3	constant	0.000010	300000	4.263254e+09	6.811216e-01
34	std	3	constant	0.000001	200000	4.254667e+09	6.817725e-01
35	std	3	constant	0.000001	300000	4.252322e+09	6.819538e-01
36	robust	1	constant	0.000100	200000	4.900985e+09	6.334022e-01
37	robust	1	constant	0.000100	300000	4.892553e+09	6.340690e-01
38	robust	1	constant	0.000010	200000	4.900094e+09	6.334710e-01
39	robust	1	constant	0.000010	300000	4.898692e+09	6.335806e-01
40	robust	1	constant	0.000001	200000	4.890141e+09	6.342141e-01
41	robust	1	constant	0.000001	300000	4.890130e+09	6.342147e-01
42	robust	2	constant	0.000100	200000	2.300894e+20	-1.712873e+10
43	robust	2	constant	0.000100	300000	2.449452e+19	-1.884110e+09
44	robust	2	constant	0.000010	200000	5.363482e+09	5.980839e-01
45	robust	2	constant	0.000010	300000	5.221100e+09	6.089986e-01
46	robust	2	constant	0.000001	200000	5.062148e+09	6.202077e-01
47	robust	2	constant	0.000001	300000	5.095500e+09	6.176499e-01
48	robust	3	constant	0.000100	200000	8.206808e+25	-6.286588e+15
49	robust	3	constant	0.000100	300000	6.506099e+26	-4.820852e+16
50	robust	3	constant	0.000010	200000	1.323675e+25	-9.869357e+14
51	robust	3	constant	0.000010	300000	8.708027e+24	-6.459236e+14
52	robust	3	constant	0.000001	200000	4.357090e+22	-3.229512e+12
53	robust	3	constant	0.000001	300000	6.341891e+22	-4.858064e+12

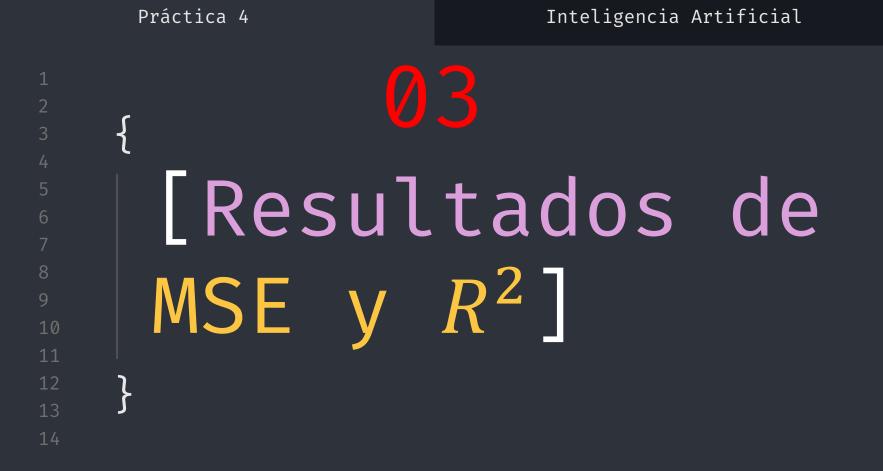
10

RESULTADOS - TABULACIÓN ORDENADA

):							
	Tipos de escalamiento	Tipo de regresión	Learning	eta0	Iteraciones	mse	r2
35	std	3	constant	0.000001	300000	4.252322e+09	6.819538e-01
34	std	3	constant	0.000001	200000	4.254667e+09	6.817725e-01
33	std	3	constant	0.000010	300000	4.263254e+09	6.811216e-01
32	std	3	constant	0.000010	200000	4.264831e+09	6.809793e-01
25	std	2	constant	0.000100	300000	4.411492e+09	6.701361e-01
29	std	2	constant	0.000001	300000	4.412737e+09	6.700638e-01
28	std	2	constant	0.000001	200000	4.413337e+09	6.700165e-01
27	std	2	constant	0.000010	300000	4.427348e+09	6.689693e-01
24	std	2	constant	0.000100	200000	4.430263e+09	6.687677e-01
26	std	2	constant	0.000010	200000	4.433751e+09	6.684951e-01
23	std	1	constant	0.000001	300000	4.820945e+09	6.394966e-01
22	std	1	constant	0.000001	200000	4.821078e+09	6.394871e-01
21	std	1	constant	0.000010	300000	4.822069e+09	6.394121e-01
20	std	1	constant	0.000010	200000	4.822685e+09	6.393663e-01
18	std	1	constant	0.000100	200000	4.823021e+09	6.393440e-01
19	std	1	constant	0.000100	300000	4.824937e+09	6.391974e-01
41	robust	1	constant	0.000001	300000	4.890130e+09	6.342147e-01
40	robust	1	constant	0.000001	200000	4.890141e+09	6.342141e-01
37	robust	1	constant	0.000100	300000	4.892553e+09	6.340690e-01
39	robust	1	constant	0.000010	300000	4.898692e+09	6.335806e-01
38	robust	1	constant	0.000010	200000	4.900094e+09	6.334710e-01
36	robust	1	constant	0.000100	200000	4.900985e+09	6.334022e-01
46	robust	2	constant	0.000001	200000	5.062148e+09	6.202077e-01
47	robust	2	constant	0.000001	300000	5.095500e+09	6.176499e-01
45	robust	2	constant	0.000010	300000	5.221100e+09	6.089986e-01
44	robust	2	constant	0.000010	200000	5.363482e+09	5.980839e-01
30	std	3	constant	0.000100	200000	2.256206e+18	-1.664487e+08

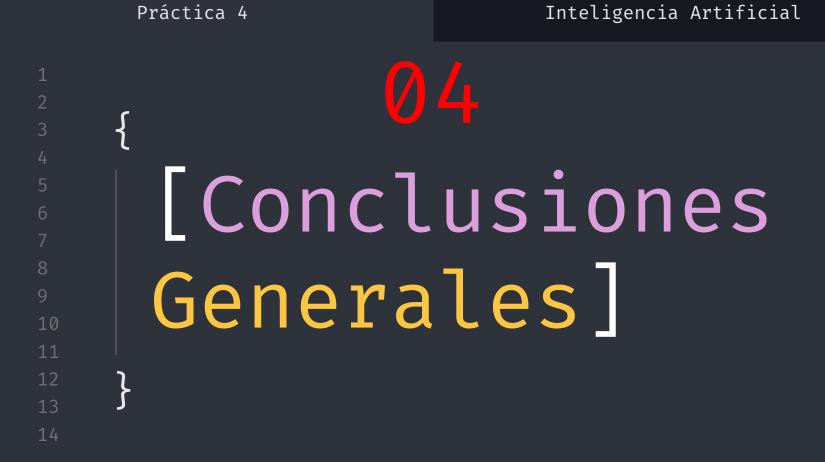
43	robust	2	constant	0.000100	300000	2.449452e+19	-1.884110e+09
31	std	3	constant	0.000100	300000	3.878283e+19	-2.861396e+09
42	robust	2	constant	0.000100	200000	2.300894e+20	-1.712873e+10
52	robust	3	constant	0.000001	200000	4.357090e+22	-3.229512e+12
53	robust	3	constant	0.000001	300000	6.341891e+22	-4.858064e+12
51	robust	3	constant	0.000010	300000	8.708027e+24	-6.459236e+14
50	robust	3	constant	0.000010	200000	1.323675e+25	-9.869357e+14
48	robust	3	constant	0.000100	200000	8.206808e+25	-6.286588e+15
4	none	1	constant	0.000001	200000	1.264673e+26	-9.216813e+15
5	none	1	constant	0.000001	300000	1.761400e+26	-1.310978e+16
49	robust	3	constant	0.000100	300000	6.506099e+26	-4.820852e+16
2	none	1	constant	0.000010	200000	6.243362e+27	-4.611774e+17
3	none	1	constant	0.000010	300000	8.296358e+27	-6.384893e+17
1	none	1	constant	0.000100	300000	1.565061e+30	-1.175406e+20
0	none	1	constant	0.000100	200000	4.620386e+30	-3.451906e+20
11	none	2	constant	0.000001	300000	1.276942e+43	-9.510445e+32
10	none	2	constant	0.000001	200000	1.524542e+43	-1.151256e+33
8	none	2	constant	0.000010	200000	3.883418e+45	-2.855971e+35
9	none	2	constant	0.000010	300000	1.320227e+46	-9.511375e+35
6	none	2	constant	0.000100	200000	1.705548e+47	-1.262719e+37
7	none	2	constant	0.000100	300000	2.775312e+47	-2.073150e+37
17	none	3	constant	0.000001	300000	2.289084e+61	-1.708496e+51
16	none	3	constant	0.000001	200000	4.245958e+61	-3.316343e+51
14	none	3	constant	0.000010	200000	2.516924e+63	-1.909706e+53
15	none	3	constant	0.000010	300000	3.556790e+64	-2.799108e+54
12	none	3	constant	0.000100	200000	1.122011e+66	-8.250235e+55
13	none	3	constant	0.000100	300000	1.321542e+66	-1.013416e+56

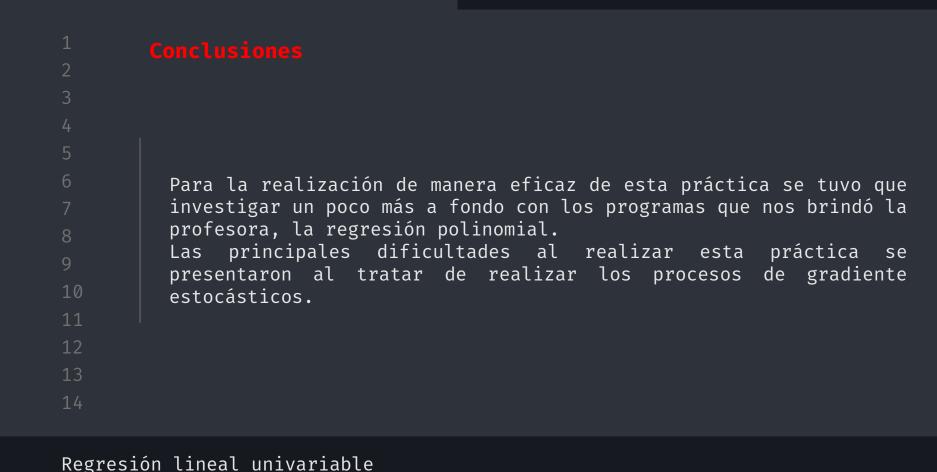
10



RESULTADOS

```
1 mse, r2 = polynomial_regression(X_train, X_test, y_train, y_test, degree = 3, eta0 = 0.00001, learning_rate = 'constant', it
In [12]:
          print('mse con X_train y X_test:', mse)
          2 print('r2 con X_train y X_test:', r2)
         mse con X_train y X_test: 4281796831.20584
         r2 con X_train y X_test: 0.6716299909141945
```





Inteligencia Artificial

Práctica 3