In [438	<pre>import pandas as pd import matplotlib.pyplot as plt</pre>
In [439	<pre>import numpy as np from sklearn.model_selection import train_test_split  #Traemos el dataset de Boston y armaremos nuestro modelo con todas las variables boston = pd.read_csv('boston.csv') boston.head()</pre>
Out[439	CRIM         ZN         INDUS         CHAS         NOX         RM         AGE         DIS         RAD         TAX         PTRATIO         B         LSTAT         MEDV           0         0.00632         18.0         2.31         0         0.538         6.575         65.2         4.0900         1         296         15.3         396.90         4.98         24.0           1         0.02731         0.0         7.07         0         0.469         6.421         78.9         4.9671         2         242         17.8         396.90         9.14         21.6           2         0.02729         0.0         7.07         0         0.469         7.185         61.1         4.9671         2         242         17.8         392.83         4.03         34.7           3         0.03237         0.0         2.18         0         0.458         6.998         45.8         6.0622         3         222         18.7         394.63         2.94         33.4           4         0.06905         0.0         2.18         0         0.458         7.147         54.2         6.0622         3         222         18.7         396.90         5.33         36.2
Out[440	X
In []:	504 6.48 505 7.88 Name: LSTAT, Length: 506, dtype: float64
Out[441	test = df.drop(train.index)\nX_train = train['LSTAT']\ny_train = train['MEDV']\nX_test = test['LSTAT']\ny_test
In [442 In [443	Test:',X_test.size)  Size of dataset: 506   Size of Train: 404   Size of Test: 102  #print('Size of dataset:',len(boston),'  Size of Train:',len(train),'  Size of Test:',len(test))
In [444	<pre>print("Train: ", X_train,"\nTest: ", X_test)  Train: 477</pre>
In [445	Ajustaremos nuestro modelo hasta que el error del modelo entrenado con train tienda a 0 y analizaremos el error que retornan las predicciones  Vamos a usar un ejemplo: https://towardsdatascience.com/bias-variance-trade-off-with-python-example-6519d2084be4 Ya que es nuestras librerias hacen malas predicciones  import numpy as np import matplotlib.pyplot as plt from random import randint from math import pi as PI
In [446	<pre>from sklearn.metrics import mean_squared_error from sklearn.model_selection import train_test_split  plt.style.use('ggplot')</pre>
	<pre># predictors x = np.linspace(0, 2, num = NUM_OBS) # noise eps = np.random.normal(0, 1, NUM_OBS) # outcome y = np.sin(PI*x) + eps  fig = plt.figure(figsize=(7,7)) ax = plt.axes() ax.set_title('y = sin(\pi x) + \text{E'}) ax.set_xlabel("x") ax.set_ylabel("y")</pre>
Out [446	<pre>ax.scatter(x, y, c = 'k')  <pre></pre></pre>
	1 -
In [447 In [448	<pre># armamos los dataset de training y test en proporcion 80-20 x_train, x_test, y_train, y_test = train_test_split(x, y, test_size = 0.2, random_state = 1)  # grafiquemos la dispersion de puntos para cada set fig = plt.figure(figsize = (7,7)) ax = plt.axes()</pre>
Out[448	<pre>ax.set_title('y = sin(\pix) + E') ax.set_xlabel("x") ax.set_ylabel("y")  ax.scatter(x_train, y_train, label = "Train set") ax.scatter(x_test, y_test, label = "Test set") ax.legend()  <matplotlib.legend.legend 0x7fec0184e730="" at=""></matplotlib.legend.legend></pre>
	$y = \sin(\pi x) + \varepsilon$ $2 - \frac{1}{2} - 1$
In [449	# number of polynomial models to investigate  NUM_MODELS = 25 # seteamos los colores para cada grado de polinomio  colors = iter(plt.cm.rainbow(np.linspace(0, 1, NUM_MODELS))) # creamos los graficos  fig, axs = plt.subplots(2, 2, figsize = (20,20)) # titulos
	<pre>axs[0, 0].set_title('Train Set Predictions') axs[1, 0].set_title('Test Set Predictions') axs[0, 1].set_title('Train Set MSE') axs[1, 1].set_title('Test Set MSE') # etiquetas de los ejes axs[0, 0].set_xlabel("x train") axs[0, 0].set_ylabel("y") axs[0, 0].set_ylim([-3, 3]) axs[1, 0].set_xlabel("x test")</pre>
	<pre>axs[1, 0].set_ylabel("y") axs[1, 0].set_ylim([-3, 3])  axs[0, 1].set_xlabel("polynomial degree") axs[0, 1].set_ylabel("train MSE")  axs[1, 1].set_xlabel("polynomial degree") axs[1, 1].set_ylabel("test MSE") # dispersion de puntos axs[0,0].scatter(x_train, y_train, c = 'k', label = "y train") axs[1,0].scatter(x_test, y_test, c = 'k', label = "y test")</pre>
	<pre>#para almacenar los mse que salgan de cada modelo train_mse = [] test_mse = []  for k in range(0, NUM_MODELS):</pre>
	<pre># polyfit nos da los coeficientes para cada grado # k-th degree polynomial coefficients fit_coeff = np.polyfit(x_train, y_train, deg = k+1)  # nos da los valores y para cada coeficiente # train and test k-th degree polynomial fit y_train_pred = np.polyval(fit_coeff, x_train) y_test_pred = np.polyval(fit_coeff, x_test)  # plot train and test k-th degree polynomial fit axs[0,0].scatter(x_train, y_train_pred, color = c, label = "deg: {}".format(k+1))</pre>
	<pre>axs[1,0].scatter(x_test, y_test_pred, color = c, label = "deg: {}".format(k+1))  # obtengo los errores  # train and test MSE of k-th degree polynomial fit  iter_train_mse = mean_squared_error(y_train_pred, y_train)  iter_test_mse = mean_squared_error(y_test_pred, y_test)  train_mse.append(iter_train_mse)  test_mse.append(iter_test_mse)</pre>
	<pre># plot train and test MSE of k-th degree polynomial fit axs[0,1].plot(k+1,</pre>
	<pre>color = c, label = "deg: {}".format(k+1), marker = 'D', markersize = 12, markeredgecolor = 'black', markeredgewidth = 3)  # plot dashed line to interpolate MSE measures axs[0,1].plot(range(1,NUM_MODELS+1), train_mse, 'k') axs[1,1].plot(range(1,NUM_MODELS+1), test_mse, 'k')</pre>
	<pre># draw legends axs[0,1].legend(loc = "upper right",</pre>
	/Users/tomasnozica/opt/anaconda3/lib/python3.8/site-packages/IPython/core/interactiveshell.py:3437: RankWarnin g: Polyfit may be poorly conditioned exec(code_obj, self.user_global_ns, self.user_ns) /Users/tomasnozica/opt/anaconda3/lib/python3.8/site-packages/IPython/core/interactiveshell.py:3437: RankWarnin g: Polyfit may be poorly conditioned exec(code_obj, self.user_global_ns, self.user_ns) /Users/tomasnozica/opt/anaconda3/lib/python3.8/site-packages/IPython/core/interactiveshell.py:3437: RankWarnin g: Polyfit may be poorly conditioned exec(code_obj, self.user_global_ns, self.user_ns) /Users/tomasnozica/opt/anaconda3/lib/python3.8/site-packages/IPython/core/interactiveshell.py:3437: RankWarnin g: Polyfit may be poorly conditioned exec(code_obj, self.user_global_ns, self.user_ns) /Users/tomasnozica/opt/anaconda3/lib/python3.8/site-packages/IPython/core/interactiveshell.py:3437: RankWarnin g: Polyfit may be poorly conditioned exec(code_obj, self.user_global_ns, self.user_ns) /Users/tomasnozica/opt/anaconda3/lib/python3.8/site-packages/IPython/core/interactiveshell.py:3437: RankWarnin g: Polyfit may be poorly conditioned exec(code_obj, self.user_global_ns, self.user_ns) /Users/tomasnozica/opt/anaconda3/lib/python3.8/site-packages/IPython/core/interactiveshell.py:3437: RankWarnin g: Polyfit may be poorly conditioned
Out[449	exec(code_obj, self.user_global_ns, self.user_ns) /Users/tomasnozica/opt/anaconda3/lib/python3.8/site-packages/IPython/core/interactiveshell.py:3437: RankWarnin g: Polyfit may be poorly conditioned exec(code_obj, self.user_global_ns, self.user_ns) <matplotlib.legend.legend 0x7fec015d5a00="" at="">  Train Set Predictions  Train Set MSE  deg: 1 deg: 2 deg: 3 deg: 4 deg: 5 deg: 6 deg: 7 deg: 8 deg: 9 deg: 10</matplotlib.legend.legend>
	deg: 11 deg: 12 deg: 13 deg: 14 deg: 15 deg: 16 deg: 17 deg: 18 deg: 17 deg: 18 deg: 20 deg: 21 deg: 22 deg: 23 deg: 21 deg: 22 deg: 23 deg: 25 deg: 25 deg: 25 deg: 25 deg: 25 deg: 26 deg: 26 deg: 26 deg: 27 deg: 26 deg: 26 deg: 27 deg: 26 deg: 27 deg: 26 deg: 26 deg: 27 deg: 27 deg: 28 deg: 2
	Test Set Predictions  Test Set MSE   deg: 1 deg: 2 deg: 3 deg: 4 deg: 5 deg: 6 deg: 7 deg: 8 deg: 9 deg: 10 deg: 11 deg: 12 deg: 13 deg: 14
In [450	105 deg: 16 deg: 17 deg: 18 deg: 20 deg: 21 deg: 22 deg: 22 deg: 23 deg: 24 deg: 25 x test
	Best fit polynomial degree: 2  Usamos este espacio para agradecerle a Nicolo Cosimo Albanese por haber hecho este artículo ya que estaba dificil armar esta seccion práctica  Cross Validation  Modelo 1: Y = B1 * X1  Modelo 2: Y = B2 * X2  Modelo 3: Y = B1 * X1 + B2 * X2
In [451 In [452	<pre>MODELO 1  import sklearn boston = pd.read_csv('boston.csv')  X = boston[['LSTAT']]</pre>
In [454 In [455 Out[455	x_crain, x_cesc, y_crain, y_cesc = crain_cesc_spire(x, y, cesc_srze=0.2, random_scace=0)
Out[455  In [456  In [457	¿ Por que neg MSE ? https://stackoverflow.com/questions/48244219/is-sklearn-metrics-mean-squared-error-the-larger-the-better-negated  mse['model1'] = abs(scores.mean())  MODELO 2  X = boston[['RM']] y = boston[['MEDV']]
In [459 Out[459 In [460	<pre>model_2 = linear_model.LinearRegression().fit(X_train, y_train) scores = cross_val_score(model_2, X, y, cv=5, scoring='neg_mean_squared_error') abs(scores) array([ 10.11581024,</pre>
In [461 In [462	<pre>MODELO 3  X = boston[['RM', 'LSTAT']] y = boston[['MEDV']]  X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=0)</pre>
In [463 Out[463 In [464	<pre>model_3 = Tilleal_model.EllealReglession().Tit(x_train, y_train) scores = cross_val_score(model_3, X, y, cv=5, scoring='neg_mean_squared_error') abs(scores)  array([11.78879381, 28.97215949, 47.84395882, 71.77053822, 36.6099291 ])  mse['model3'] = abs(scores.mean())</pre>
In [469 Out[469	mise