

# MODELOS Y SIMULACION

En la clase de hoy aplicaremos las técnicas de selección discutidas la última clase: Forward Stepwise Selection, Backward Stepwise Selection y Backward Stepwise Selection with p-values.

Para esta semana, siguen abiertas las Office Hours con el mismo link de siempre.

Ya se encuentra disponible la solución a la practica 1 en Google Drive.

## Inteligencia Artificial y espionaje

El tema para discutir esta semana es cómo esta cambiando el espinoaje a nivel mundial y como el aprovechamiento de la Inteligencia Artificial es una oportunidad y un desafío para las distintos países u organizaciones. Esta conversacion surge a partir de un entrevista publicada por el HAI (Human-centered Artificial Intelligence) de la Universidad de Stanford. El articulo se encuentra siguiendo el siguiente link: <https://hai.stanford.edu/news/re-imagining-espionage-era-artificial-intelligence>

Discutimos en clase sobre el mismo.

```
In [104.]: #Importamos las librerias que vamos a usar hoy
import pandas as pd
import numpy as np
from sklearn.neighbors import LinearModel
from sklearn.metrics import r2_score
from statsmodels.api import sm
import operator
import matplotlib.pyplot as plt

In [105.]: #Traemos el dataset de Boston y armaremos nuestro modelo con todas las variables
boston = pd.read_csv('boston.csv')
var_explcativas = boston.drop('MEDV', )
var_objetivo = boston['MEDV']

In [106.]: #Cargamos el modelo y lo entrenamos
regr = LinearModel.LinearRegression()
regr.fit(var_explcativas, var_objetivo)

Out[106.]:

In [107.]: #Vemos los coeficientes
regr.coef_

Out[107.]: array([ 0.0011326e+01,  0.0028974e+01,  0.0000000e+00,
        2.68673382e+00, -1.77666112e+01,  3.80989521e+00,
        6.92224640e+00, -1.47556685e+00,  3.06049479e+01,
        1.23446393e+00, -9.52747232e+01,  8.31168370e+01,
        -5.24738379e+01])

In [108.]: #Vemos el termino independiente
regr.intercept_

Out[108.]: -1.0774311471147114

In [109.]: #Corremos el modelo para obtener las predicciones
prediccion_precios = regr.predict(var_explcativas)

In [110.]: #Obtenemos el R2
r2 = r2_score(var_objetivo, prediccion_precios)
print('R2: ', round(r2, 1))

R2: 0.741
```

## Calculemos R<sup>2</sup> Ajustado

Segun lo visto en clase, definimos a R<sup>2</sup> como:

$$R^2 = 1 - \frac{RSS}{TSS}$$

donde RSS es Residual Sum of Squares y se calcula como:

$$RSS = \sum (y_i - \hat{y}_i)^2$$

y TSS:

$$TSS = \sum (y_i - y_{med-i})^2$$

Para R<sup>2</sup><sub>ajus</sub> queremos penalizar la incorporacion de nuevas variables:

$$R^2_{ajus} = 1 - \frac{RSS \cdot (n-d-1)}{TSS \cdot (n-1)}$$

reemplazando terminos:

$$R^2_{ajus} = 1 - (1 - R^2) \left( \frac{n-1}{n-d-1} \right)$$

```
In [111.]: #En base a lo definido anteriormente calculamos R2 ajustado
#N es el numero de ocurrencias
n = var_objetivo.size
#d es el numero de variables explicativas
d = var_explcativas.columns.size
print('n: ', n, 'd: ', d)

n: 506
d: 13

In [112.]: #Planteamos nuestra formula y vemos el resultado
r2_ajust = 1 - (1 - r2) * (n - 1) / (n - d - 1)
print('R2 Ajustado: ', round(r2_ajust, 1))
print('R2: ', round(r2, 1))

R2 Ajustado: 0.74
R2: 0.741

Una nueva libreria: StatsModels
```

Esta libreria permite hacer estimaciones para diferentes modelos estadísticos, además de permitir hacer pruebas estadísticas y exploración de data <https://www.statsmodels.org/stable/index.html>.

Respecto a SKLearn, para Regresión Lineal las diferencias no son significativas. Podemos leer este articulo donde se comparan ambos módulos: <https://becominghuman.ai/stats-models-vs-sklearn-for-linear-regression-f10d95ad90b>

¿Por que lo utilizaremos? porque nos brinda muchos detalles estadísticos que SKLearn no.

```
In [113.]: #En StatsModels podemos agregar nosotros mismo la constante para el armado de nuestro modelo.
var_explcativas = sm.add_constant(var_explcativas)
var_explcativas.head()

Out[113.]:
   const  CRIM  ZN  INDUS  CHAS  NOX  RM  AGE  DIS  RAD  TAX  PTRATIO  B  LSTAT
0  1.0  0.00632  18.0  2.31  0  0.538  6.575  65.2  4.0900  1  296  15.3  396.90  4.98
1  1.0  0.02731  0.0  7.07  0  0.469  6.421  78.9  4.9671  2  242  17.8  396.90  9.14
2  1.0  0.02729  0.0  7.07  0  0.469  6.185  61.1  4.9671  2  242  17.8  392.83  4.03
3  1.0  0.03237  0.0  2.18  0  0.458  6.998  45.8  6.0622  3  222  18.7  394.63  2.94
4  1.0  0.06905  0.0  2.18  0  0.458  7.147  54.2  6.0622  3  222  18.7  396.90  5.33

In [114.]: #Definimos nuestro modelo, OLS significa Ordinary Least Squares. Esto indica la forma en la que
#se calculan los parametros
model = sm.OLS(var_objetivo, var_explcativas)

In [115.]: #Entrenamos el modelo (notemos las diferencias con sklearn)
regr = model.fit()

In [116.]: #Vemos todos los datos que nos entrega StatsModels
print(regr.summary())

OLS Regression Results
=====
Dep. Variable:  MEDV    R-squared:  0.741
Model:  OLS    Adj. R-squared:  0.734
Method:  Least Squares    F-statistic:  408.1
Date:  Tue, 31 Aug 2021    Prob(F-statistic):  6.72e-125
Time:  16:16:37    Log-Likelihood:  -1498.9
No. Observations:  506    AIC:  3026.
No. Residuals:  494    BIC:  3085.
Model:  OLS
Covariance Type:  nonrobust

=====
coef    std err    t    P>|t|    [0.025    0.975]
-----
const    36.4595    5.103    7.144    0.000    26.432    46.487
CRIM    -0.1080    0.033    -3.287    0.001    -0.173    -0.043
ZN     0.0464    0.014    3.382    0.001    0.019    0.073
INDUS    0.0206    0.061    0.334    0.738    -0.100    0.141
CHAS     2.6867    0.862    3.118    0.002    0.994    4.380
NOX    -17.7746    3.820   -4.653    0.000   -25.372   -10.165
RM     3.8099    0.418    9.116    0.000    2.989    4.631
AGE    -0.0607    0.013    -0.052    0.958    -0.025    -0.027
DIS    -1.4756    0.139   -10.598    0.000   -1.747   -1.204
RAD     3.0605    0.066    4.613    0.000    2.926    3.200
PTRATIO -0.0123    0.004   -0.320    0.901   -0.020   -0.005
B     0.2587    0.131    1.979    0.000   -0.010   -0.696
LSTAT   0.0093    0.003    3.467    0.001    0.004    0.015
STATAT  -0.5248    0.051   -10.347    0.000   -0.624   -0.425
=====
Durbin-Watson:  1.78041
Prob(Cointibus):  0.000    Jarque-Bera (JB):  787.125
Skew:  1.521    Prob(JB):  8.84e-171
Autocorr:  8.281    Cond. No.:  1.51e+04

=====
Notes:
[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.
[2] The condition number is large, 1.51e+04. This might indicate that there are
strong multicollinearity or other numerical problems.
```

```
In [117.]: #Para acceder al R2 y R2 Ajustado
r2 = regr.rsquared
r2_adj = regr.rsquared_adj
print('R2: ', round(r2, 1))
print('R2 Ajustado: ', round(r2_adj, 1))

R2: 0.741
R2 Ajustado: 0.74
```

```
In [118.]: #Voy a ver mas interesantes: los p-values
p_values = regr.pvalues
print(regr.pvalues)

const      7.143191e-13
CRIM       7.888100e-03
INDUS      7.382801e-01
CHAS       1.925030e-03
NOX        4.245644e-05
AGE        1.879441e-18
RM         6.013491e-13
DIS        5.070529e-06
RAD        1.116137e-03
PTRATIO    1.308835e-12
B          1.979441e-18
LSTAT      7.709152e-03
dtype: float64
```

```
In [119.]: #Cuales son mayores a 0.05
print("Variables which p value is > 0.05:")
for index, value in p_values.items():
    if value > 0.05:
        print(index)

Variables which p value is > 0.05:
INDUS
DIS
```

```
In [120.]: #Ordenemoslos de mayor a menor
regr.pvalues.sort_values(ascending=False)
```

```
Out[120.]:
const      7.143191e-13
INDUS      7.382801e-01
CHAS       1.925030e-03
CRIM       7.888100e-03
ZN         1.116137e-03
RM         6.013491e-13
DIS        5.070529e-06
RAD        1.116137e-03
PTRATIO    1.308835e-12
B          1.979441e-18
LSTAT      7.709152e-03
dtype: float64
```

## Forward Stepwise Selection

Partimos de un modelo sin variables predictoras. Luego comenzamos a agregar de a 1 variable, eligiendo en cada ronda aquel que tenga mayor R2 A medida que vamos conformando nuestros modelos, almacenamos el R2 ajustado Terminando la iteracion, graficamos los R2 ajustados y seleccionamos el mejor modelo.

```
In [121.]: #Traemos el dataset de Boston y armaremos nuestro modelo con todas las variables
boston = pd.read_csv('boston.csv')
var_explcativas = boston.drop('MEDV', )
var_explcativas = sm.add_constant(var_explcativas)
var_objetivo = boston['MEDV']
```

```
In [122.]: # Seteamos las condiciones iniciales necesarias
variables = ['const']
iterate_columns = var_explcativas.columns.drop('const', )
r2_adj = []
vars_size = iterate_columns.size
var_model = []

# Iteramos k veces, siendo k la cantidad de variables
for k in range(1, vars_size):
    # Acumulamos los valores de R2 para cada variable en los modelos que armaremos
    r2 = []

    # Iteramos sobre todas las variables disponibles para la ronda
    for var in iterate_columns:
        # Fijamos que variables seran las que definiran nuestro modelo
        var_explcativa = var_explcativas.columns + [var]
        model = sm.OLS(var_objetivo, var_explcativa)
        regr = model.fit()

        # Almacenamos el R2 para cada set de variables que se prueba
        r2_var = regr.rsquared_adj

    # Seleccionamos aquella con mayor R2
    var_max_r2 = max(r2.items(), key=operator.itemgetter(1))[1]

    # Almacenamos las variables que describen a ese modelo
    var_model[k] = var_max_r2

    # Almacenamos el valor de R2
    r2_adj.append(r2_var_max_r2)

    # Eliminamos la variable para que deje de estar en consideracion en el ciclo siguiente
    iterate_columns = iterate_columns.drop(var_max_r2, )

    # Agregamos la variable seleccionada en esta vuelta para que la cuente en el modelo siguiente
    variables.append(var_max_r2)

#Seteamos el valor de la ronda k donde R2 es maximo
r2_max_index = r2_adj.index(max(r2_adj))

r2_variation vars_size r2_adj title x_label y_label:
***
Plot the scatter and the curve that shows how R2 value vary over
every iteration. A title and labels for the graph must be included.
Also, highlight the point where there is a maximum R2 value. Add information
about that point. The amount of independent variables is given.
***
# Graficamos la solucion
# Los valores para señalar cada iteracion
x = np.arange(vars_size)
# El punto donde R2 es maximo
r2_max_x = r2_adj.index(max(r2_adj))
r2_max_y = max(r2_adj)

# Titulo de nuestro grafico y nombre de los ejes
plt.title title, loc = 'left'
plt.xlabel x_label
plt.ylabel y_label

# Graficamos la recta con la variacion del R2 en cada vuelta
plt.plot(x, r2_adj)

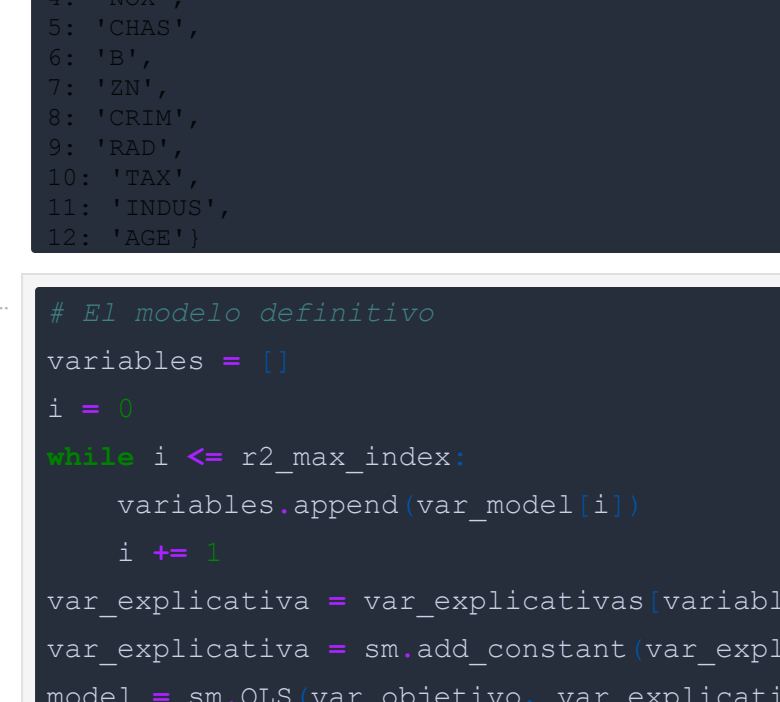
# Graficamos dispersion de puntos de X vs R2 ajustado
plt.scatter(x, r2_adj)

# Graficamos el punto donde R2 es maximo y añadamos su valor en el grafico
plt.scatter(r2_max_x, r2_max_y, marker='*', s=100, color='red')
plt.text(r2_max_x * (1 + 0.05), r2_max_y * (1 + 0.05) , round(r2_max_y, 1), fontsize=12)
plt.text(r2_max_x * (1 + 0.05), r2_max_y * (1 - 0.05) , 'k = ' + str(r2_max_x), fontsize=12)
plt.show()
```

```
In [124.]: r2_variation vars_size r2_adj 'Forward Stepwise Selection' 'R2' 'R2'

Forward Stepwise Selection
=====
coef    std err    t    P>|t|    [0.025    0.975]
-----
const    36.4595    5.103    7.144    0.000    26.432    46.487
CRIM    -0.1080    0.033    -3.287    0.001    -0.173    -0.043
ZN     0.0464    0.014    3.382    0.001    0.019    0.073
INDUS    0.0206    0.061    0.334    0.738    -0.100    0.141
CHAS     2.6867    0.862    3.118    0.002    0.994    4.380
NOX    -17.7746    3.820   -4.653    0.000   -25.372   -10.165
RM     3.8099    0.418    9.116    0.000    2.989    4.631
AGE    -0.0607    0.013    -0.052    0.958    -0.025    -0.027
DIS    -1.4756    0.139   -10.598    0.000   -1.747   -1.204
RAD     3.0605    0.066    4.613    0.000    2.926    3.200
PTRATIO -0.0123    0.004   -0.320    0.901   -0.020   -0.005
B     0.2587    0.131    1.979    0.000   -0.010   -0.696
LSTAT   0.0093    0.003    3.475    0.001    0.004    0.015
STATAT  -0.5248    0.051   -10.347    0.000   -0.624   -0.425
=====
Durbin-Watson:  1.78040
Prob(Cointibus):  0.000    Jarque-Bera (JB):  787.125
Skew:  1.523    Prob(JB):  8.60e-172
Autocorr:  8.305    Cond. No.:  1.47e+04

=====
Notes:
[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.
[2] The condition number is large, 1.47e+04. This might indicate that there are
strong multicollinearity or other numerical problems.
```



```
In [125.]: var_model

Out[125.]:
[0]: 'const',
[1]: 'PTRATIO',
[2]: 'CRIM',
[3]: 'INDUS',
[4]: 'CHAS',
[5]: 'RM',
[6]: 'ZN',
[7]: 'CRIM',
[8]: 'RM',
[9]: 'PTRATIO',
[10]: 'INDUS',
[11]: 'INDUS'
```

```
In [126.]: # El modelo definitivo
variables = []
i = 0
while i <= r2_max_index:
    variables.append(var_model[i])
    i += 1
var_explcativa = var_explcativas[variables]
var_explcativa = sm.add_constant(var_explcativa)
model = sm.OLS(var_objetivo, var_explcativa)
regr = model.fit()
print(regr.summary())

OLS Regression Results
=====
Dep. Variable:  MEDV    R-squared:  0.741
Model:  OLS    Adj. R-squared:  0.735
Method:  Least Squares    F-statistic:  408.1
Date:  Tue, 31 Aug 2021    Prob(F-statistic):  5.54e-137
Time:  16:16:37    Log-Likelihood:  -1498.9
No. Observations:  506    AIC:  3026.
No. Residuals:  494    BIC:  3072.
Model:  OLS
Covariance Type:  nonrobust

=====
coef    std err    t    P>|t|    [0.025    0.975]
-----
const    36.4595    5.103    7.144    0.000    26.432    46.487
CRIM    -0.1080    0.033    -3.287    0.001    -0.173    -0.044
ZN     0.0464    0.014    3.382    0.001    0.019    0.073
INDUS    0.0206    0.061    0.334    0.738    -0.100    0.141
CHAS     2.6867    0.862    3.118    0.002    0.994    4.380
NOX    -17.7746    3.820   -4.653    0.000   -25.372   -10.165
RM     3.8099    0.418    9.116    0.000    2.989    4.631
AGE    -0.0607    0.013    -0.052    0.958    -0.025    -0.027
DIS    -1.4756    0.139   -10.598    0.000   -1.747   -1.204
RAD     3.0605    0.066    4.613    0.000    2.926    3.200
PTRATIO -0.0123    0.004   -0.320    0.901   -0.020   -0.005
B     0.2587    0.131    1.979    0.000   -0.010   -0.696
LSTAT   0.0093    0.003    3.475    0.001    0.004    0.015
STATAT  -0.5248    0.051   -10.347    0.000   -0.624   -0.425
=====
Durbin-Watson:  1.78040
Prob(Cointibus):  0.000    Jarque-Bera (JB):  787.125
Skew:  1.523    Prob(JB):  8.60e-172
Autocorr:  8.305    Cond. No.:  1.47e+04

=====
Notes:
[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.
[2] The condition number is large, 1.47e+04. This might indicate that there are
strong multicollinearity or other numerical problems.
```

## Backward Stepwise Selection

Partimos de un modelo definido por todas variables predictoras. Luego comenzamos a eliminar de a 1 variable, eligiendo en cada ronda aquel que tenga mayor R2 A medida que vamos conformando nuestros modelos, almacenamos el R2 ajustado Terminando la iteracion, graficamos los R2 ajustados y seleccionamos el mejor modelo.

```
In [127.]: boston = pd.read_csv('boston.csv')
var_explcativas = boston.drop('MEDV', )
var_explcativas = sm.add_constant(var_explcativas)
var_objetivo = boston['MEDV']
```

```
In [128.]: # Seteamos las condiciones iniciales necesarias
variables = var_explcativas.columns
iterate_columns = variables.drop('const', )
vars_size = iterate_columns.size
r2_adj = []
var_model = []

# Iteramos k veces, siendo k la cantidad de variables
for k in range(1, vars_size):
    # Acumulamos los valores de R2 para cada variable en los modelos que armaremos
    r2 = []

    # Indice para insertar las variables temporalmente removidas en su lugar
    i = 0

    # Iteramos sobre todas las variables disponibles para la ronda
    for var in iterate_columns:
        # Eliminamos la variable para probar
        variables = variables.drop(var, )

        # Fijamos que variables seran las que definiran nuestro modelo
        var_explcativa = var_explcativas[variables]
        var_explcativa = sm.add_constant(var_explcativa)

        # Insertamos nuevamente la variable que sacamos en el lugar que estaba para la prox
        # ronda
        variables = variables.insert(i, var)

        # Entrenamos el modelo
        model = sm.OLS(var_objetivo, var_explcativa)
        regr = model.fit()

        # Almacenamos el R2 para cada set de variables que se prueba
        r2_var = regr.rsquared_adj

    # Actualizador de indice
    i += 1

    # Seleccionamos aquella con mayor R2
    var_max_r2 = max(r2.items(), key=operator.itemgetter(1))[1]

    # Almacenamos el valor de R2
    r2_adj.append(r2_var_max_r2)

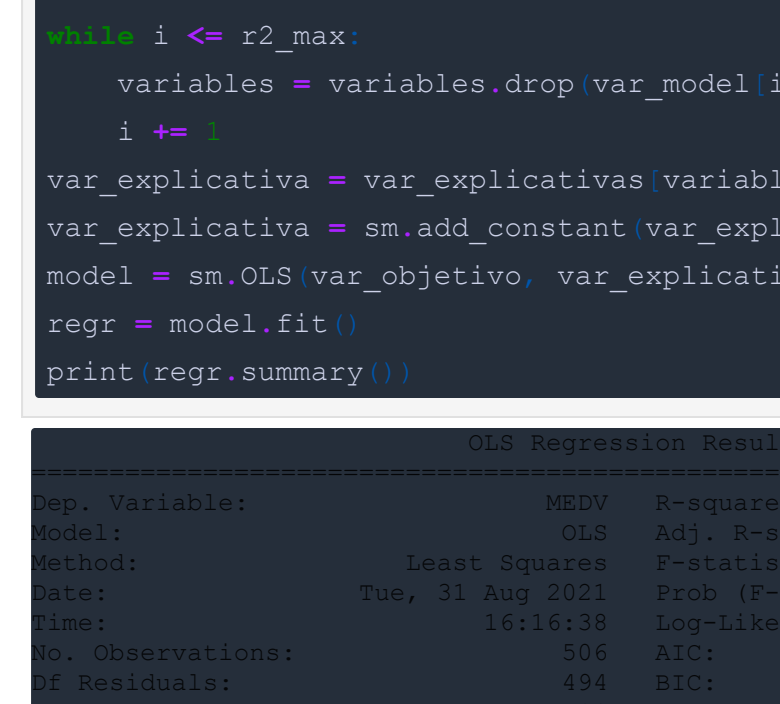
    # Almacenamos las variables que describen a ese modelo
    var_model[k] = var_max_r2

    # No itera mas sobre esa variable tampoco
    iterate_columns = iterate_columns.drop(var_max_r2, )

    # Dejamos de considerar como una variable posible para definir al modelo
    variables = variables.drop(var_max_r2, )

#Seteamos el valor de la ronda k donde R2 es maximo
r2_max_index = r2_adj.index(max(r2_adj))
```

```
In [129.]: r2_variation vars_size r2_adj 'Backward Stepwise Selection with p-values' 'R2' 'R2'
```



```
In [130.]: # El modelo definitivo
variables = var_explcativas.columns
i = 0
while i <= r2_max_index:
    variables = variables.drop(var_model[i], )
    i += 1
var_explcativa = var_explcativas[variables]
var_explcativa = sm.add_constant(var_explcativa)
model = sm.OLS(var_objetivo, var_explcativa)
regr = model.fit()
print(regr.summary())

OLS Regression Results
=====
Dep. Variable:  MEDV    R-squared:  0.741
Model:  OLS    Adj. R-squared:  0.735
Method:  Least Squares    F-statistic:  408.1
Date:  Tue, 31 Aug 2021    Prob(F-statistic):  5.54e-137
Time:  16:16:38    Log-Likelihood:  -1498.9
No. Observations:  506    AIC:  3026.
No. Residuals:  494    BIC:  3072.
Model:  OLS
Covariance Type:  nonrobust

=====
coef    std err    t    P>|t|    [0.025    0.975]
-----
const    36.4595    5.103    7.144    0.000    26.432    46.487
CRIM    -0.1080    0.033    -3.287    0.001    -0.173    -0.044
ZN     0.0464    0.014    3.382    0.001    0.019    0.073
INDUS    0.0206    0.061    0.334    0.738    -0.100    0.141
CHAS     2.6867    0.862    3.118    0.002    0.994    4.380
NOX    -17.7746    3.820   -4.653    0.000   -25.372   -10.165
RM     3.8099    0.418    9.116    0.000    2.989    4.631
AGE    -0.0607    0.013    -0.052    0.958    -0.025    -0.027
DIS    -1.4756    0.139   -10.598    0.000   -1.747   -1.204
RAD     3.0605    0.066    4.613    0.000    2.926    3.200
PTRATIO -0.0123    0.004   -0.320    0.901   -0.020   -0.005
B     0.2587    0.131    1.979    0.000   -0.010   -0.696
LSTAT   0.0093    0.003    3.475    0.001    0.004    0.015
STATAT  -0.5248    0.051   -10.347    0.000   -0.624   -0.425
=====
Durbin-Watson:  1.78040
Prob(Cointibus):  0.000    Jarque-Bera (JB):  787.125
Skew:  1.523    Prob(JB):  8.60e-172
Autocorr:  8.305    Cond. No.:  1.47e+04

=====
Notes:
[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.
[2] The condition number is large, 1.47e+04. This might indicate that there are
strong multicollinearity or other numerical problems.
```

## Backward Stepwise Selection with P-Values

Partimos de un modelo definido por todas variables predictoras. Luego eliminamos, en cada vuelta, a la variable con mayor P-Value A medida que vamos conformando nuestros modelos, almacenamos el R2 ajustado Terminando la iteracion, graficamos los R2 ajustados y seleccionamos el mejor modelo.

```
In [131.]: boston = pd.read_csv('boston.csv')
var_explcativas = boston.drop('MEDV', )
var_explcativas = sm.add_constant(var_explcativas)
var_objetivo = boston['MEDV']
```

```
In [132.]: variables = var_explcativas.columns
r2_adj = []
vars_out = []
for k in range(1, vars_size):
    var_explcativa = var_explcativas[variables]
    # Entrenamos el modelo
    model = sm.OLS(var_objetivo, var_explcativa)
    regr = model.fit()

    # Almacenamos el valor de R2
    r2_adj.append(regr.rsquared_adj)

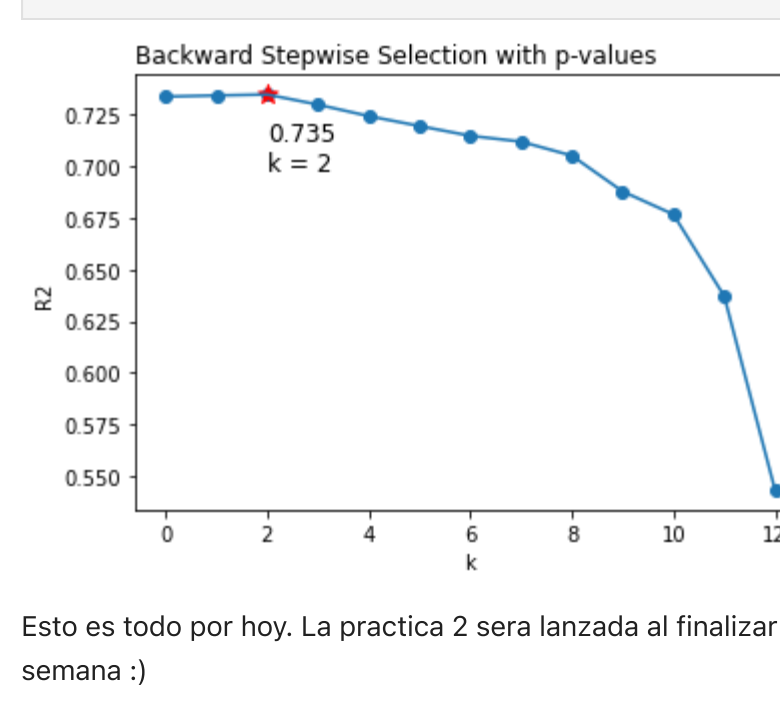
    # Obtenemos los p-values sin considerar la constante
    p_values = regr.pvalues.drop('const', )
    pvalue_index = p_values.argmax()
    pvalue_var = p_values.keys()[pvalue_index]

    # Eliminamos la variable cuyo p-value es el mas grande
    variables = variables.drop(pvalue_var)

    # Almacenamos la variable que queda fuera en cada ronda
    vars_out.append(pvalue_var)
```

```
Out[132.]:
['CRIM',
 'INDUS',
 'CHAS',
 'RM',
 'PTRATIO',
 'CRIM',
 'PTRATIO',
 'INDUS',
 'RM']

In [134.]: r2_variation vars_size r2_adj 'Backward Stepwise Selection with p-values' 'R2' 'R2'
```



Esto es todo por hoy. La practica 2 sera lanzada al finalizar esta clase. Lo visto hoy quizas les presente dificultad. Es normal. Buena semana :)