

ML_REGRESION_03

Modelo de Regresión Lineal Múltiple

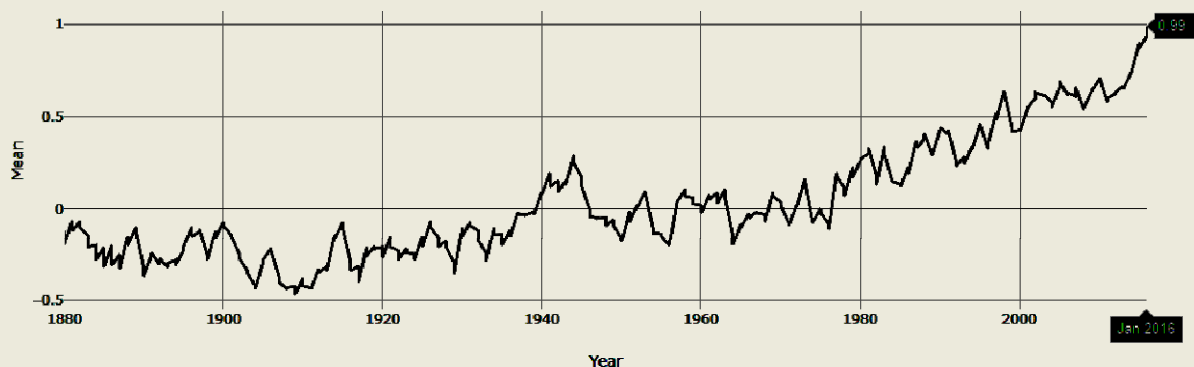
ML

La temperatura media del planeta, la concentración de dióxido de carbono CO₂ de la atmósfera, el nivel medio del mar y la masa glaciaria media están evolucionando durante las últimas décadas.

El problema a realizar con esta práctica consiste en realizar un modelo de aprendizaje supervisado basado en regresión, verificar si la relación del CO₂ con la temperatura, nivel medio del mar y masa glaciaria se ajusta bien a un modelo de regresión lineal múltiple o no, y utilizarlo para predecir valores futuros.

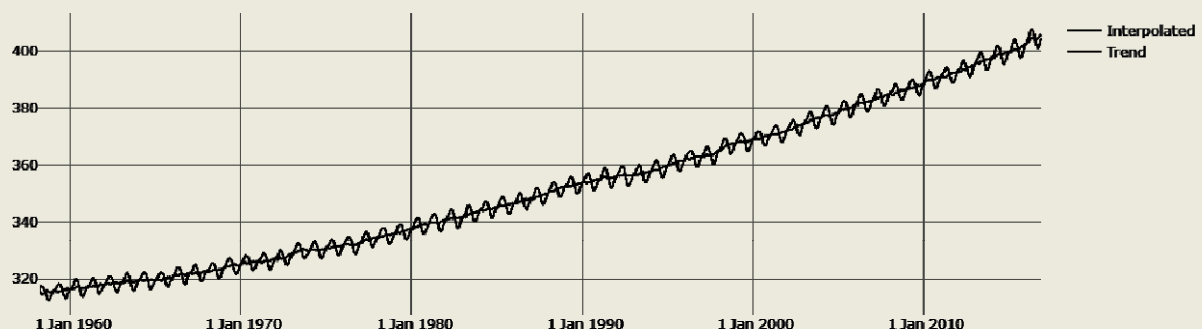
$$\hat{y} = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3$$

$\hat{y} \rightarrow \text{CO}_2$ / $x_1 \rightarrow \text{Temp}$ / $x_2 \rightarrow \text{NMM}$ / $x_3 \rightarrow \text{Glaciaria}$



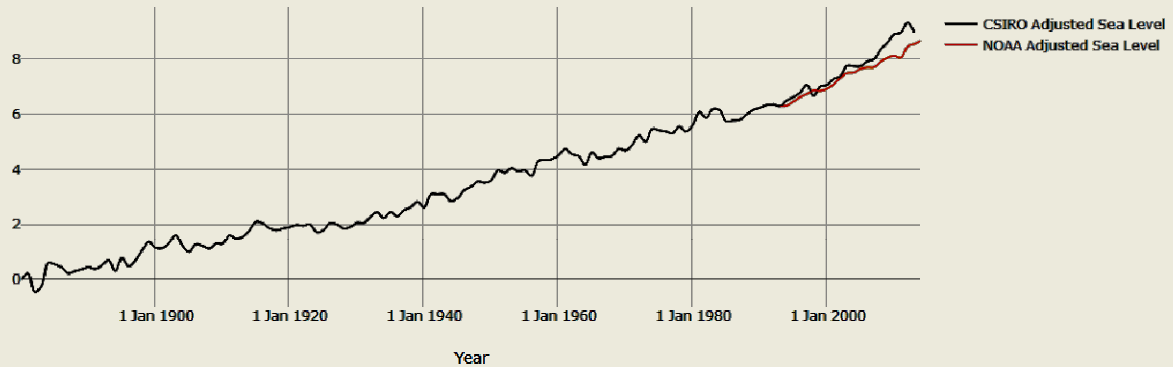
Anomalía de la temperatura media mensual en grados Celsius relativos a un periodo base

Fuente: <https://datahub.io/core/global-temp#resource-annual>



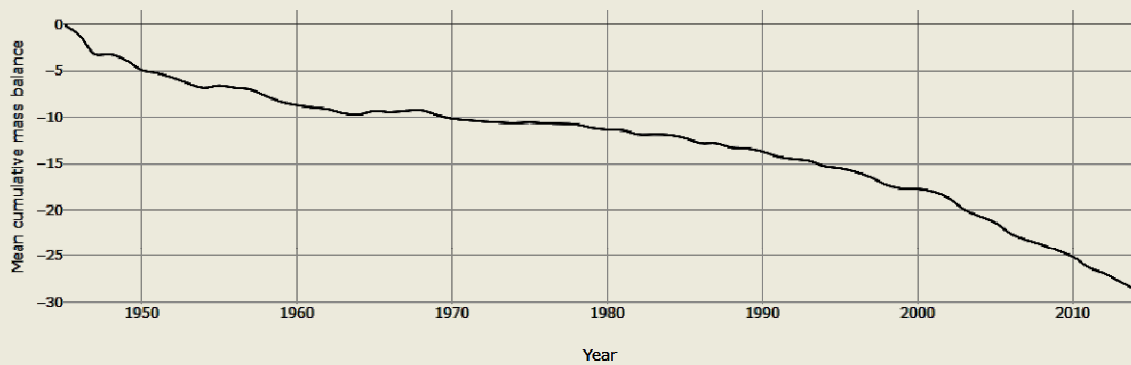
Tendencia del dióxido de carbono (CO₂) de la atmósfera

Fuente: https://pkgtstore.datahub.io/core/co2-ppm/co2-annmean-mlo_json/data/31185d494d1a6f6431aee8b8004b6164/co2-annmean-mlo_json.json



Nivel medio del mar (Average Sea Level)

Fuente: <https://datahub.io/core/sea-level-rise>



Masa Glaciar Media (Mean glacier mass)

Fuente: <https://datahub.io/core/glacier-mass-balance>

SOLUCIÓN

Importar las librerías necesarias para realizar la práctica.

```
# Importar librerías
import urllib.request, json
import matplotlib.pyplot as plt
import pandas as pd
import numpy as np
```

Descargar y leer la información de partida (Datasets) y leer las variables

```
# Información de partida (Datasets)
# Fuente:

# Temperatura: https://datahub.io/core/global-temp#data
# Temperatura anual media (formato json)
# https://pkgstore.datahub.io/core/global-temp/annual_json/data/529e69dbd597709e36ce11a5d0bb7243/annual_json.json
temp_url="https://pkgstore.datahub.io/core/global-temp/annual_json/data/529e69dbd597709e36ce11a5d0bb7243/annual_json.json"

# CO2: https://datahub.io/core/co2-ppm#data
# Concentración anual media (formato json)
# https://pkgstore.datahub.io/core/co2-ppm/co2-annmean-mlo_json/data/31185d494d1a6f6431aee8b8004b6164/co2-annmean-mlo_json.json
co2_url="https://pkgstore.datahub.io/core/co2-ppm/co2-annmean-mlo_json/data/31185d494d1a6f6431aee8b8004b6164/co2-annmean-mlo_json.json"

# Nivel medio del mar: https://datahub.io/core/sea-level-rise#data
# Nivel medio del mar anual (formato json)
# https://pkgstore.datahub.io/core/sea-level-rise/epa-sea-level_json/data/ac016d75688136c47a04ac70298e42ec/epa-sea-level_json.json
sea_url="https://pkgstore.datahub.io/core/sea-level-rise/epa-sea-level_json/data/ac016d75688136c47a04ac70298e42ec/epa-sea-level_json.json"

# Masa glaciar media: https://datahub.io/core/glacier-mass-balance
# Masa glaciar media anual (formato json)
# https://pkgstore.datahub.io/core/glacier-mass-balance/glaciers_json/data/6270342ca6134dadf8f94221be683bc6/glaciers_json.json
glaciar_url="https://pkgstore.datahub.io/core/glacier-mass-balance/glaciers_json/data/6270342ca6134dadf8f94221be683bc6/glaciers_json.json"

# lectura de la información de los ficheros JSON
with urllib.request.urlopen(temp_url) as url:
    temp_data = json.loads(url.read().decode())

with urllib.request.urlopen(co2_url) as url:
    co2_data = json.loads(url.read().decode())

with urllib.request.urlopen(sea_url) as url:
    sea_data = json.loads(url.read().decode())

with urllib.request.urlopen(glaciar_url) as url:
    glaciar_data = json.loads(url.read().decode())
```

Registrar los valores en listas.

```
# Registro de las variables (listas) de temperatura y co2
temp=[]
co2=[]
sea=[]
glaciar=[]
year=[]

ntemp=len(temp_data)
nco2=len(co2_data)
nsea=len(sea_data)
nglaciar=len(glaciar_data)

# Registro de temperaturas desde el año 1880
for i in range(ntemp):
    if temp_data[i]["Source"]=="GISTEMP":
        # Se utiliza la temperatura media en superficie (NASA)
        # GISTEMP: https://data.giss.nasa.gov/gistemp/
        temp.append(temp_data[i]["Mean"])
        year.append(temp_data[i]["Year"])

# Las listas de temperatura y años están en orden decreciente (de 2016 a 1880)
# Las ordenamos en orden creciente
temp.reverse()
year.reverse()
# y nos quedamos con la serie desde 1959 hasta 2013
# En total son 58 registros
temp=temp[1959-1880:2013-1880+1]
year=year[1959-1880:2013-1880+1]

# Registro de CO2 desde el año 1959
for i in range(nco2):
    co2.append(co2_data[i]["Mean"])
# nos quedamos con la serie desde 1959 hasta 2013
co2=co2[:2013-1959+1]

# Registro del nivel medio del mar desde el año 1880 hasta 2013
for i in range(nsea):
    sea.append(sea_data[i]["CSIRO Adjusted Sea Level"])

# nos quedamos con la serie desde 1959 hasta 2013
# En total son 55 registros
sea=sea[1959-1880:2013-1880+1]

# Registro de la masa media de los glaciares desde el año 1945 hasta 2014
for i in range(nglaciar):
    glaciar.append(glaciar_data[i]["Mean cumulative mass balance"])
```

```
# nos quedamos con la serie desde 1959 hasta 2013
# En total son 55 registros
glaciar=glaciar[1959-1945:2013-1945+1]
```

Visualizar la temperatura y CO2 por años

```
# Visualizamos la temperatura y CO2
fig, axs = plt.subplots(4,1)
```

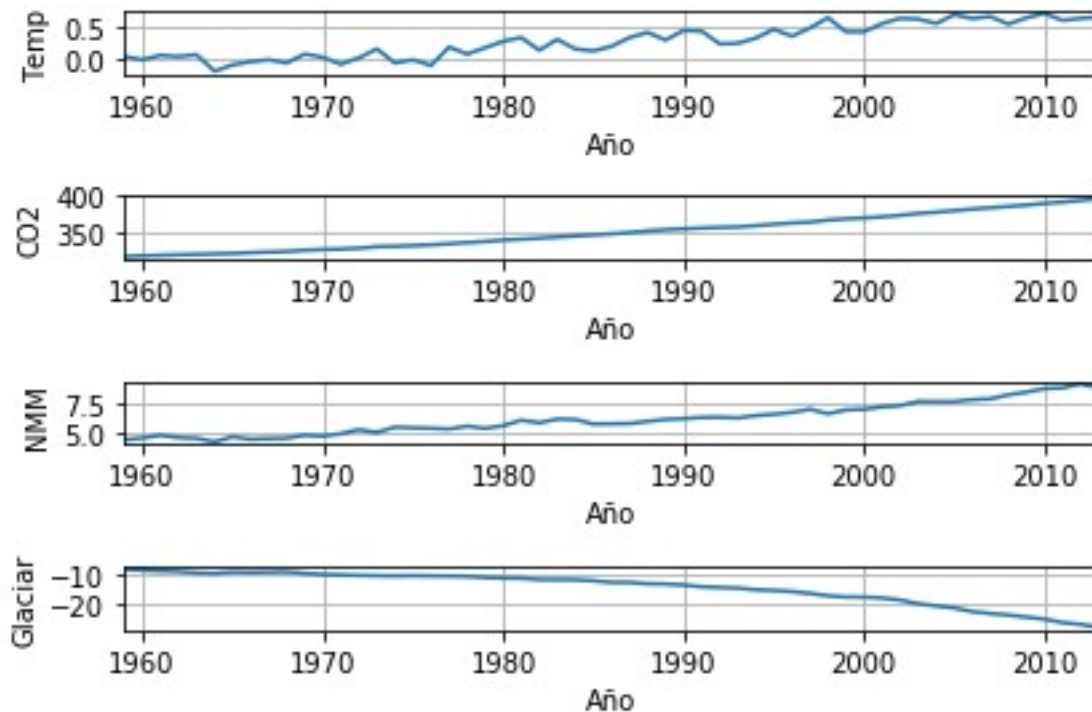
```
# Temperatura
axs[0].plot(year,temp)
axs[0].set_xlim(1959,2013)
axs[0].set_xlabel("Año")
axs[0].set_ylabel("Temp")
axs[0].grid(True)
```

```
# CO2
axs[1].plot(year,co2)
axs[1].set_xlim(1959,2013)
axs[1].set_xlabel("Año")
axs[1].set_ylabel("CO2")
axs[1].grid(True)
```

```
# Nivel medio del mar (NMM)
axs[2].plot(year,sea)
axs[2].set_xlim(1959,2013)
axs[2].set_xlabel("Año")
axs[2].set_ylabel("NMM")
axs[2].grid(True)
```

```
# Masa media de los glaciares
axs[3].plot(year,glaciar)
axs[3].set_xlim(1959,2013)
axs[3].set_xlabel("Año")
axs[3].set_ylabel("Glaciar")
axs[3].grid(True)
```

```
fig.tight_layout()
plt.show()
```



Realizar el modelo de regresión lineal. Previamente convertimos las listas de temperatura, CO2, nivel medio del mar y masa glaciar media en un DataFrame utilizando la librería Pandas.

```
# REGRESIÓN LINEAL (APRENDIZAJE SUPERVISADO)

# Utilizamos la serie temporal de temperaturas, CO2, nivel medio mar y
# glaciar desde 1959 hasta el año 2013
# para construir un modelo de regresión lineal múltiple

# Creamos un Dataframe con los datos utilizando la librería Pandas
datos={'temp':temp,'co2':co2,'sea':sea,'glaciar':glaciar}
df=pd.DataFrame(datos,columns=['temp','co2','sea','glaciar'])

# Asignamos las variables X (atributos) e y (etiquetas)
X=df[['temp','sea','glaciar']]
y=df[['co2']]

Importar la librería 'sklearn' para realizar el modelo de regresión.

# importamos las librerías para realizar regresión lineal
# Utilizamos sklearn (http://scikit-learn.org/stable/)
# Aprendizaje Supervisado: http://scikit-learn.org/stable/supervised\_learning.html#supervised-learning
# Ejemplo de Regresión Lineal: http://scikit-learn.org/stable/supervised\_learning.html#supervised-learning
```

learn.org/stable/auto_examples/linear_model/plot_ols.html#sphx-glr-auto-examples-linear-model-plot-ols-py

```
from sklearn import linear_model
from sklearn.metrics import mean_squared_error, r2_score
```

Dividimos el conjunto de datos para entrenamiento y test

```
# Dividimos el conjunto de datos para entrenamiento y test
# Elegimos a priori el 70 % (40 registros) para entrenamiento
# y el resto 30 % (18 registros) para test
```

```
X_train = np.array(X[:40])
y_train = np.array(y[:40])
```

```
X_test = np.array(X[40:])
y_test = np.array(y[40:])
```

Realizar el ajuste del modelo de regresión lineal, y la predicción para los datos de entrenamiento.

```
# Creamos el objeto de Regresión Lineal
regr=linear_model.LinearRegression()
```

```
# Entrenamos el modelo
regr.fit(X_train,y_train)
```

```
# Realizamos predicciones sobre los atributos de entrenamiento
y_pred = regr.predict(X_train)
```

Obtener los parámetros del modelo de regresión.

```
# Recta de Regresión Lineal (y=t0+t1*X)
# Pendiente de la recta
t1=regr.coef_
print('Pendiente: \n', t1)
# Corte con el eje Y (en X=0)
t0=regr.intercept_
print('Término independiente: \n', t0)

# Ecuación de la recta
print('El modelo de regresión es: y = %f + %f * X1 + %f * X2 + %f * X3'%(t0,t1[0][0],t1[0][1],t1[0][2]))
```

Pendiente:

```
[[ 4.39312079  7.21539814 -3.83224388]]
```

Término independiente:

```
[252.83908309]
```

*El modelo de regresión es: y = 252.839083 + 4.393121 * X1 + 7.215398 * X2 + -3.832244 * X3*

Cálculo del error (pérdida) del ajuste del modelo de regresión

```
# Error (pérdida)
print("Error o pérdida del modelo de regresión lineal para valores de
entrenamiento")
# Error Cuadrático Medio (Mean Square Error)
print("ECM : %.2f" % mean_squared_error(y_train, y_pred))
# Puntaje de Varianza. El mejor puntaje es un 1.0
print('Coeficiente Correlación: %.2f' % r2_score(y_train, y_pred))

Error o pérdida del modelo de regresión lineal para valores de entrenamiento
ECM : 4.92
Coeficiente Correlación: 0.98
```

Comprobar con los valores de test y prueba si el ajuste es bueno o no.

```
# Con el modelo de regresión ajustado con los valores de entrenamiento
# se aplica a los valores para test y validación
y_pred_test = regr.predict(X_test)
# Comprobar el error del modelo con los valores para test
# Error (pérdida)
print("Error o pérdida del modelo de regresión lineal para valores de test")
# Error Cuadrático Medio (Mean Square Error)
print("ECM : %.2f" % mean_squared_error(y_test, y_pred_test))
# Puntaje de Varianza. El mejor puntaje es un 1.0
print('Coeficiente Correlación: %.2f' % r2_score(y_test, y_pred_test))

Error o pérdida del modelo de regresión lineal para valores de test
ECM : 390.19
Coeficiente Correlación: -4.09
```

Utilizar el modelo para predecir valores de la concentración de CO2

```
# Predecir la concentración de CO2 para una anomalía de 0.8, nivel medio del
mar de 0.3 y masa media de los glaciares del mundo de -6.8
y_pred2 = regr.predict([[0.8,0.3,-6.8]])
print('La predicción de CO2 para X1=0.8, X2=0.3 y X3=-6.8 es: ',y_pred2)

La predicción de CO2 para X1=0.8, X2=0.3 y X3=-6.8 es:  [[284.57745751]]
```