

# Regresion Lineal

## Nombre: David Egas

Desarrollo de un ejemplo simple de regresión lineal en python.

### Importaciones necesarias

In [4]:

```
# Imports necesarios
import numpy as np
import pandas as pd
import seaborn as sb
import matplotlib.pyplot as plt
%matplotlib inline
from mpl_toolkits.mplot3d import Axes3D
from matplotlib import cm
plt.rcParams['figure.figsize'] = (16, 9)
plt.style.use('ggplot')
from sklearn import linear_model
from sklearn.metrics import mean_squared_error, r2_score
```

### Carga del dataset

In [6]:

```
#cargamos los datos de entrada
data = pd.read_csv("./articulos_ml.csv")
#veamos cuantas dimensiones y registros contiene
data.shape
```

Out[6]:

(161, 8)

In [7]:

data.head()

Out[7]:

	Title	url	Word count	# of Links	# of comments	# Images video	Elapsed days
0	What is Machine Learning and how do we use it ...	https://blog.signals.network/what-is-machine-l...	1888	1	2.0	2	34
1	10 Companies Using Machine Learning in Cool Ways	NaN	1742	9	NaN	9	5
2	How Artificial Intelligence Is Revolutionizing...	NaN	962	6	0.0	1	10
3	Dbrain and the Blockchain of Artificial Intell...	NaN	1221	3	NaN	2	68
4	Nasa finds entire solar system filled with eig...	NaN	2039	1	104.0	4	131

## Una descripción de la información que contiene el dataset

In [8]:

data.describe()

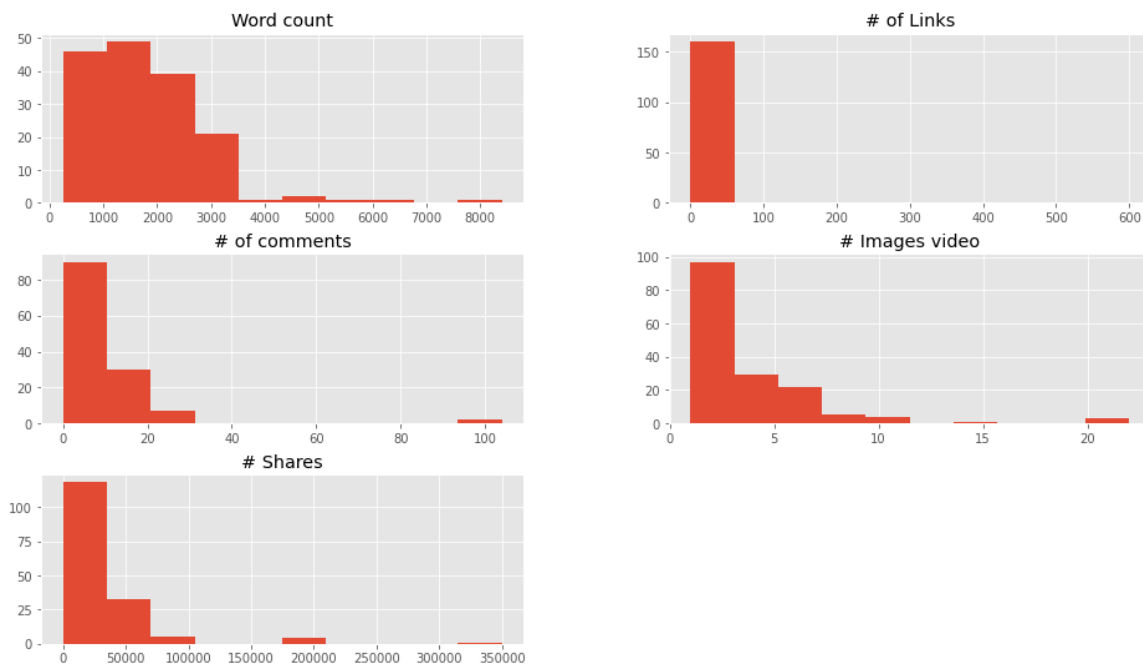
Out[8]:

	Word count	# of Links	# of comments	# Images video	Elapsed days	# Shares
<b>count</b>	161.000000	161.000000	129.000000	161.000000	161.000000	161.000000
<b>mean</b>	1808.260870	9.739130	8.782946	3.670807	98.124224	27948.347826
<b>std</b>	1141.919385	47.271625	13.142822	3.418290	114.337535	43408.006839
<b>min</b>	250.000000	0.000000	0.000000	1.000000	1.000000	0.000000
<b>25%</b>	990.000000	3.000000	2.000000	1.000000	31.000000	2800.000000
<b>50%</b>	1674.000000	5.000000	6.000000	3.000000	62.000000	16458.000000
<b>75%</b>	2369.000000	7.000000	12.000000	5.000000	124.000000	35691.000000
<b>max</b>	8401.000000	600.000000	104.000000	22.000000	1002.000000	350000.000000

## Características de entrada

In [9]:

```
# Visualizamos rápidamente las características de entrada
data.drop(['Title', 'url', 'Elapsed days'],1).hist()
plt.show()
```



In [10]:

```

# Vamos a RECORTAR Los datos en La zona donde se concentran más Los puntos
# esto es en el eje X: entre 0 y 3.500
# y en el eje Y: entre 0 y 80.000
filtered_data = data[(data['Word count'] <= 3500) & (data['# Shares'] <= 80000)]

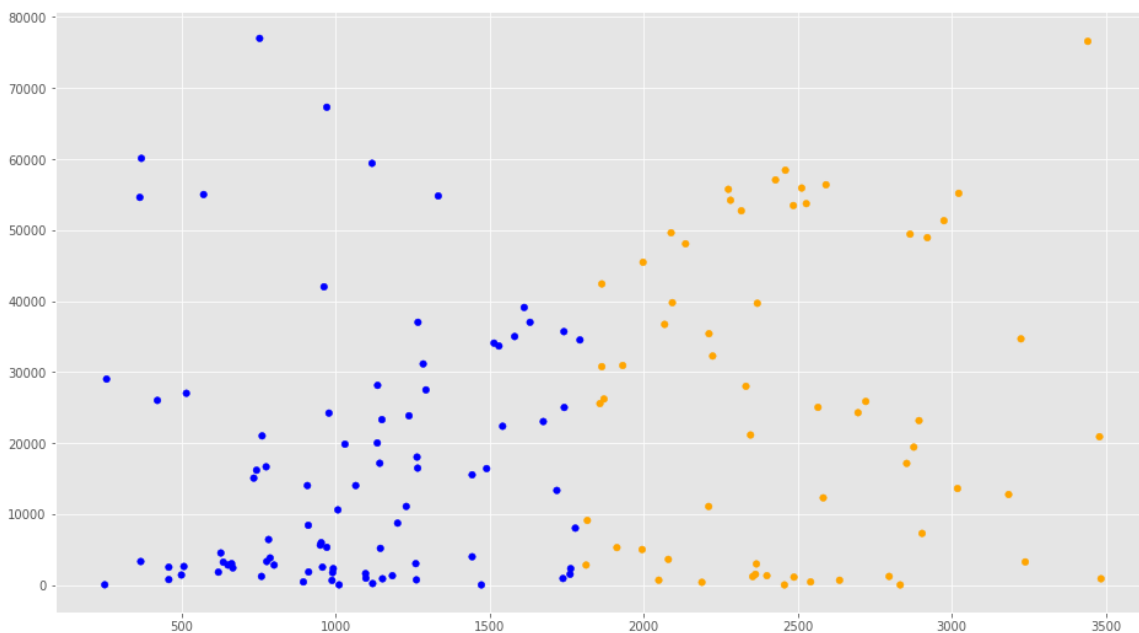
colores=['orange','blue']
tamanios=[30,60]

f1 = filtered_data['Word count'].values
f2 = filtered_data['# Shares'].values

# Vamos a pintar en colores los puntos por debajo y por encima de la media de Cantidad de Palabras
asignar=[]
for index, row in filtered_data.iterrows():
    if(row['Word count']>1808):
        asignar.append(colores[0])
    else:
        asignar.append(colores[1])

plt.scatter(f1, f2, c=asignar, s=tamanios[0])
plt.show()

```



## Entrenamiento del Dataset

In [12]:

```
# Asignamos nuestra variable de entrada X para entrenamiento y Las etiquetas Y.
dataX =filtered_data[["Word count"]]
X_train = np.array(dataX)
y_train = filtered_data['# Shares'].values

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

# Entrenamos nuestro modelo
regr.fit(X_train, y_train)

# Hacemos Las predicciones que en definitiva una línea (en este caso, al ser 2D)
y_pred = regr.predict(X_train)

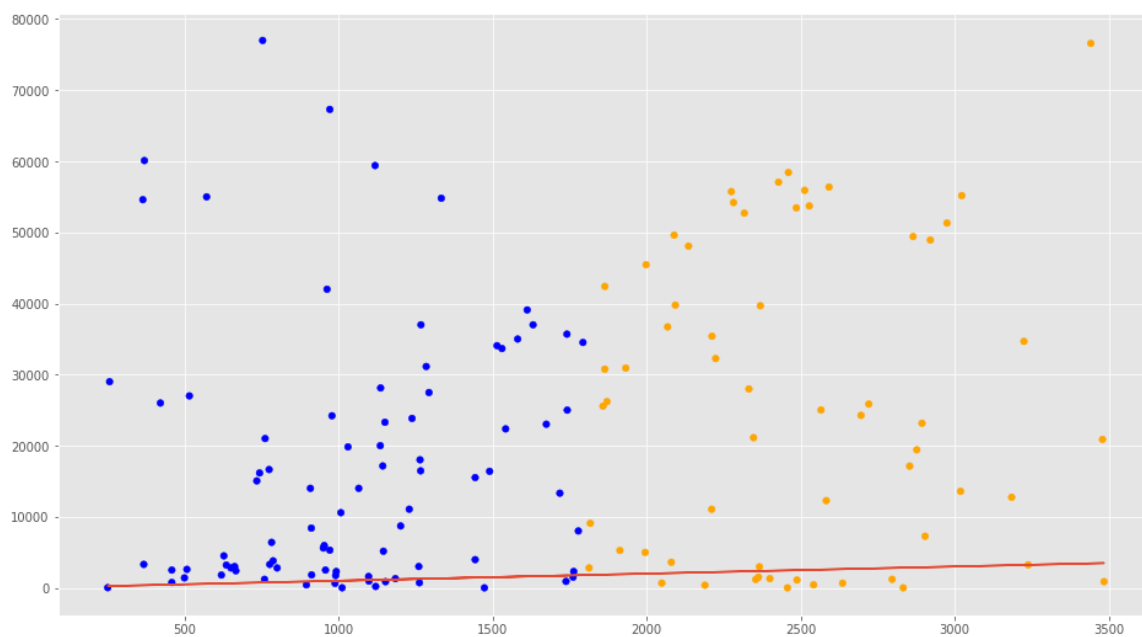
# Veamos Los coeficienetes obtenidos, En nuestro caso, serán La Tangente
print('Coefficients: \n', regr.coef_)
# Este es el valor donde corta el eje Y (en X=0)
print('Independent term: \n', regr.intercept_)
# Error Cuadrado Medio
print("Mean squared error: %.2f" % mean_squared_error(y_train, y_pred))
# Puntaje de Varianza. El mejor puntaje es un 1.0
print('Variance score: %.2f' % r2_score(y_train, y_pred))
```

```
Coefficients:
 [5.69765366]
Independent term:
 11200.30322307416
Mean squared error: 372888728.34
Variance score: 0.06
```

## Línea de separación

In [29]:

```
plt.scatter(f1, f2, c=asignar, s=tamamos[0])  
plt.plot(dataX,f1)  
plt.show()
```



In [ ]: