

Practicum

Gráficos con Python

Tips database

Sobre la base

- Utilizaremos un archivo de valores separados por comas ~> `tips.csv`.
- Las variables (columnas) empleadas en esta base son:
 - `total_bill`: total de la cuenta.
 - `tip`: propina para el camarero/a.
 - `sex`: sexo del camarero/a.
 - `smoker`: la mesa atendida fue en la sección de fumadores?
 - `day`: día de la semana del servicio.
 - `size`: cantidad de comensales atendidos.
- Aplicaremos algunas técnicas descriptivas para explorar la distribución de las propinas.

Cargando módulos

First things first, importemos módulos estandar en el análisis de datos con Python.

```
1 # Esta línea permite renderizar los gráficos en el mismo notebook.
2 %matplotlib inline
3 # Librería para la manipulación y limpieza de bases de datos.
4 import pandas as pd
5 # Librería para el procesamiento numérico
6 import numpy as np
7 # Librería para graficar
8 import matplotlib.pyplot as plt
9 # Modificar el estilo default de gráficos de Matplotlib.
10 plt.style.use('seaborn')
```

El primer punto a atacar, es inspeccionar la base. Saber con qué estamos trabajando.

```
1 # cargamos nuestro archivo .csv que se encuentra en la misma carpeta del
  notebook.
2 df = pd.read_csv('./tips.csv')
3 # examinemos las primeras 5 filas del archivo.
4 print(df.head())
5 # describamos las columnas. esto sólo funciona para las columnas numéricas (int
  o float).
6 print(df.describe())
7
8 # Describamos las frecuencias de cada una de las columnas strings. Agregamos
  len(df) para obtener porcentajes
9 # .value_counts nos permite obtener frecuencias discretas para cada nivel de
  nuestra columna.
10 print(df.sex.value_counts()/len(df))
11 print(df.smoker.value_counts()/len(df))
```

Practicum

Gráficos con Python

```
12 print(df.day.value_counts()/len(df))
13 print(df.time.value_counts()/len(df))
```

```
1   total_bill  tip    sex smoker  day    time  size
2  0      16.99  1.01  Female    No  Sun  Dinner    2
3  1      10.34  1.66   Male    No  Sun  Dinner    3
4  2      21.01  3.50   Male    No  Sun  Dinner    3
5  3      23.68  3.31   Male    No  Sun  Dinner    2
6  4      24.59  3.61  Female    No  Sun  Dinner    4
7      total_bill      tip      size
8 count  244.000000  244.000000  244.000000
9 mean    19.785943    2.998279    2.569672
10 std     8.902412    1.383638    0.951100
11 min     3.070000    1.000000    1.000000
12 25%    13.347500    2.000000    2.000000
13 50%    17.795000    2.900000    2.000000
14 75%    24.127500    3.562500    3.000000
15 max    50.810000   10.000000    6.000000
16 Male         0.643443
17 Female        0.356557
18 Name: sex, dtype: float64
19 No         0.618852
20 Yes        0.381148
21 Name: smoker, dtype: float64
22 Sat         0.356557
23 Sun         0.311475
24 Thur        0.254098
25 Fri         0.077869
26 Name: day, dtype: float64
27 Dinner       0.721311
28 Lunch        0.278689
29 Name: time, dtype: float64
```

- Python por sí solo, procesa todo en formato numérico \Rightarrow A diferencia de otros lenguajes o software diseñados para el análisis que presentan una serie de declaraciones para transformar strings a factores (`recode` en Stata y `as.factor(x)` en R).
- Debemos convertir las cadenas a integrales. `pandas` es muy útil para ello y presenta una serie de métodos para ello. Vamos a ocupar `replace`, que reemplaza valores antiguos de una serie (columna) por otros.

```
1 # Para traspasar nuestros valores de string a float, utilizamos .replace
2 df['sex'].replace(['Female', 'Male'], [1, 0], inplace=True)
3 df['smoker'].replace(['Yes', 'No'], [1, 0], inplace=True)
```

- Veamos como quedaron nuestras nuevas variables.

```
1 print(df.head())
```

Practicum

Gráficos con Python

		total_bill	tip	sex	smoker	day	time	size
2	0	16.99	1.01	1	0	Sun	Dinner	2
3	1	10.34	1.66	0	0	Sun	Dinner	3
4	2	21.01	3.50	0	0	Sun	Dinner	3
5	3	23.68	3.31	0	0	Sun	Dinner	2
6	4	24.59	3.61	1	0	Sun	Dinner	4

- Genial, ahora necesitamos agrupar categorías de la variable `day` en un indicador binario `weekend`. Para ello utilizaremos lo siguiente.

```
1 # Para recodificar múltiples valores en un nuevo valor, utilizamos .assign
2 df = df.assign(weekend = df.day.map({'Sat': 1, 'Sun': 1, 'Thur': 0, 'Fri': 0}))
3 print(df['weekend'].value_counts())
```

```
1 1    163
2 0     81
3 Name: weekend, dtype: int64
```

```
1 df.total_bill.plot(kind='hist')
```

```
1 <matplotlib.axes._subplots.AxesSubplot at 0x1c20c6cc50>
```

Practicum

Gráficos con Python

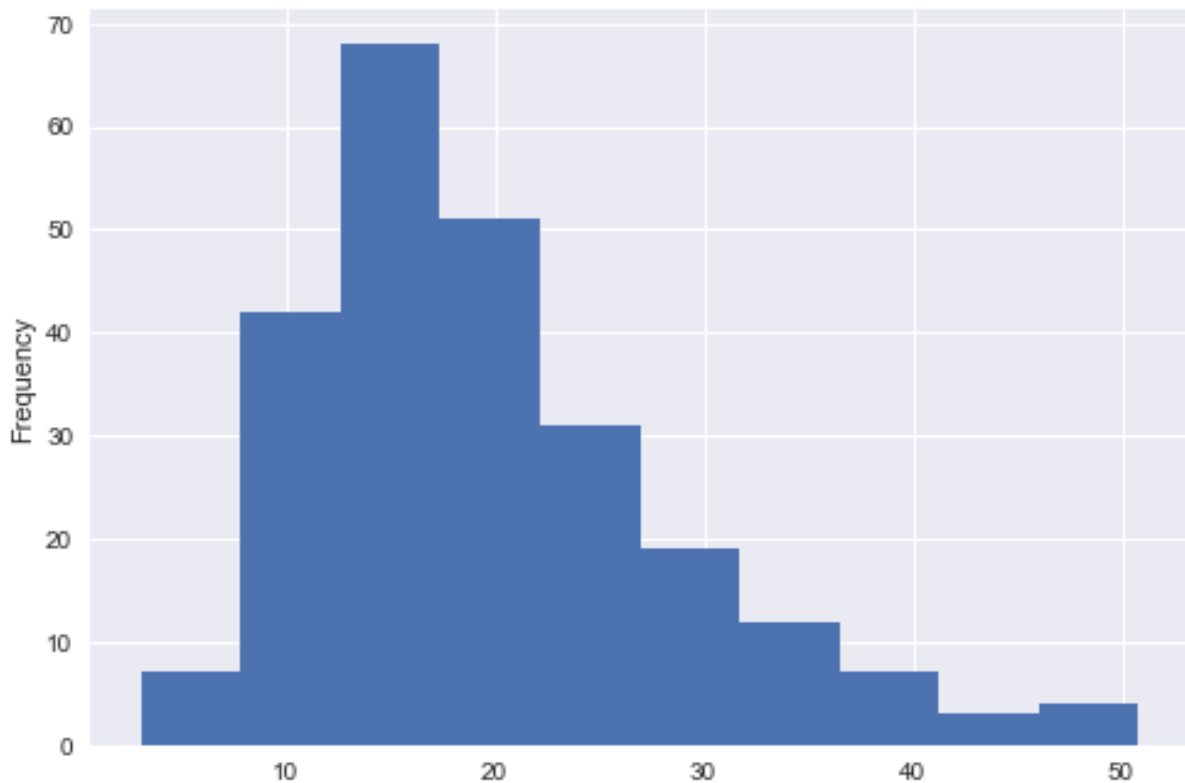


Figure 1: png

```
1
2 # Para obtener la media de una columna, utilizamos .mean()
3 mu = df['total_bill'].mean()
4 # Para obtener la desviación estandar de una columna, utilizamos .std()
5 sigma = df['total_bill'].std()
6
7 # Con mu y sigma estimados, estamos listos para poder generar una distribución
  normalizada del total de la cuenta.
8 df = df.assign(bill_std = (df['total_bill'] - mu) / sigma)
9
10 # Grafiquemos esta nueva variable
11 df.bill_std.hist()
```

```
1 <matplotlib.axes._subplots.AxesSubplot at 0x1c20c66eb8>
```

Practicum

Gráficos con Python

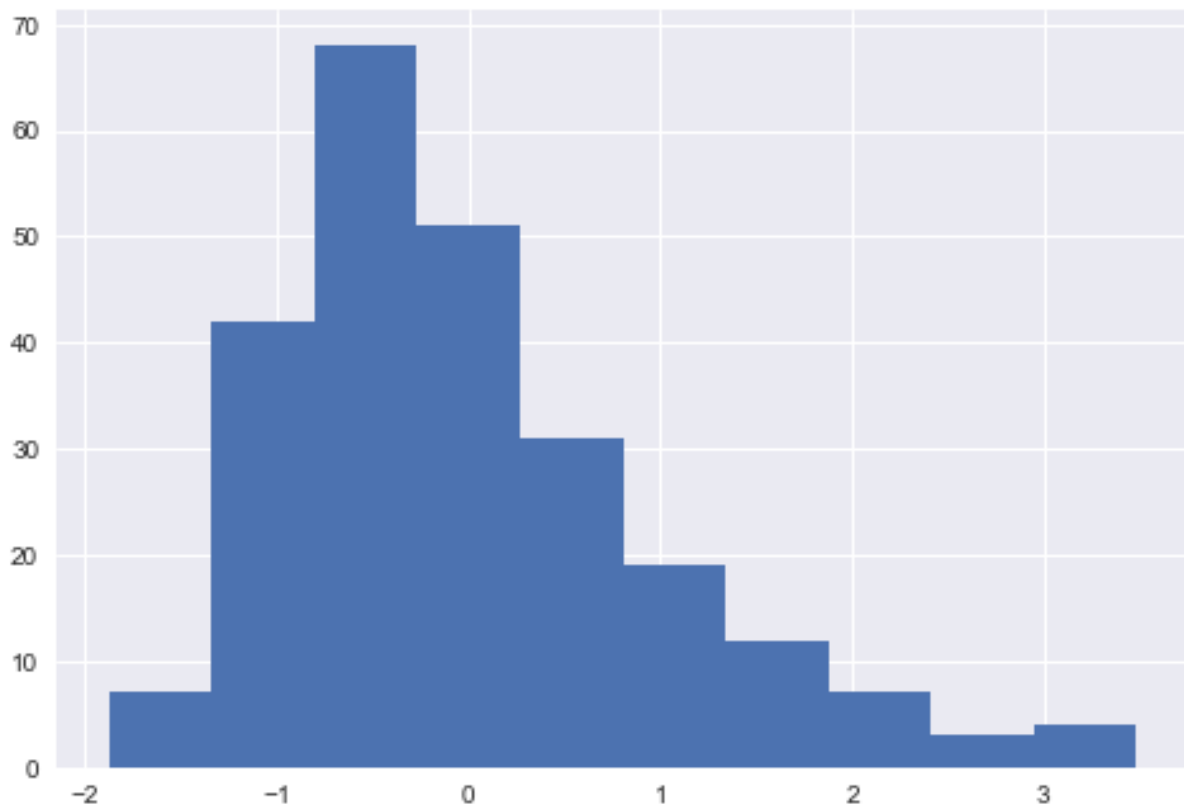


Figure 2: png

- La versión estandarizada aclara el hecho que los consumos en el restaurant tienen a concentrarse bajo la media.
- Ahora que estandarizamos la variable, tenemos una mejor apreciación por la tendencia a consumir bajo la media.
- Para estandarizar las demás variables, refactorizaremos el código en una función.

```
1 def standarize(x):
2     mu = x.mean()
3     sigma = x.std()
4     tmp = (x - mu) / sigma
5     return tmp
6
7 df = df.assign(tips_std = standarize(df['tip']))
8 print(df.head())
```

	total_bill	tip	sex	smoker	day	time	size	weekend	bill_std \
0	16.99	1.01	1	0	Sun	Dinner	2	1	-0.314066
1	10.34	1.66	0	0	Sun	Dinner	3	1	-1.061054
2	21.01	3.50	0	0	Sun	Dinner	3	1	0.137497
3	23.68	3.31	0	0	Sun	Dinner	2	1	0.437416

Practicum

Gráficos con Python

```
6 4      24.59  3.61    1      0  Sun  Dinner    4      1  0.539635
7
8      tips_std
9 0 -1.436993
10 1 -0.967217
11 2  0.362610
12 3  0.225291
13 4  0.442111
```

Gráficos

```
1 cross = pd.crosstab(df['sex'],df['size'] )
2 cross.plot(kind='bar')
```

```
1 <matplotlib.axes._subplots.AxesSubplot at 0x1c20cea3c8>
```

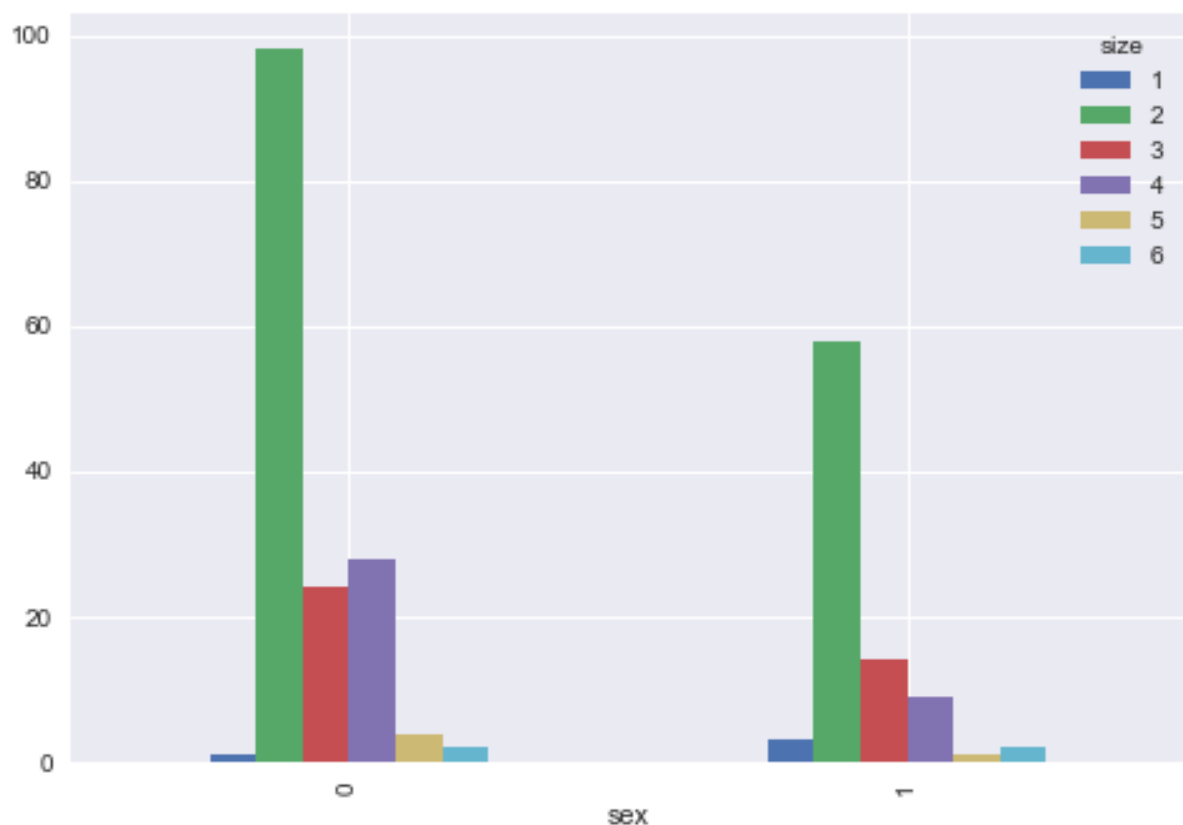


Figure 3: png

```
1 cross_2 = pd.crosstab(df['smoker'], df['size'])
```

Practicum

Gráficos con Python

```
2 cross_2.plot(kind="bar")
```

```
1 <matplotlib.axes._subplots.AxesSubplot at 0x1c20e1bef0>
```

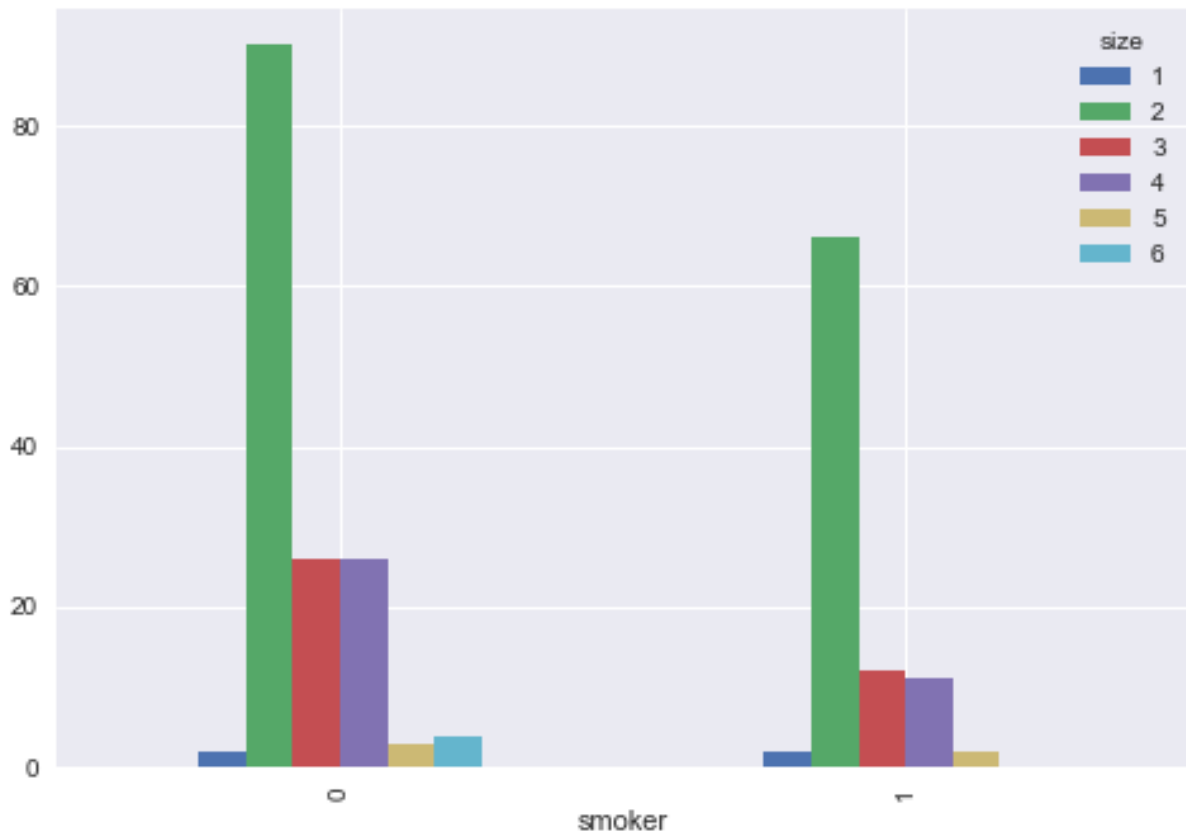


Figure 4: png

```
1 cross_3 = pd.crosstab(df['tip'], df['sex'])
2 cross_3.plot(kind='box')
3 pd.crosstab(df['tip'], df['smoker']).plot(kind = 'box')
```

```
1 <matplotlib.axes._subplots.AxesSubplot at 0x1c210a5400>
```

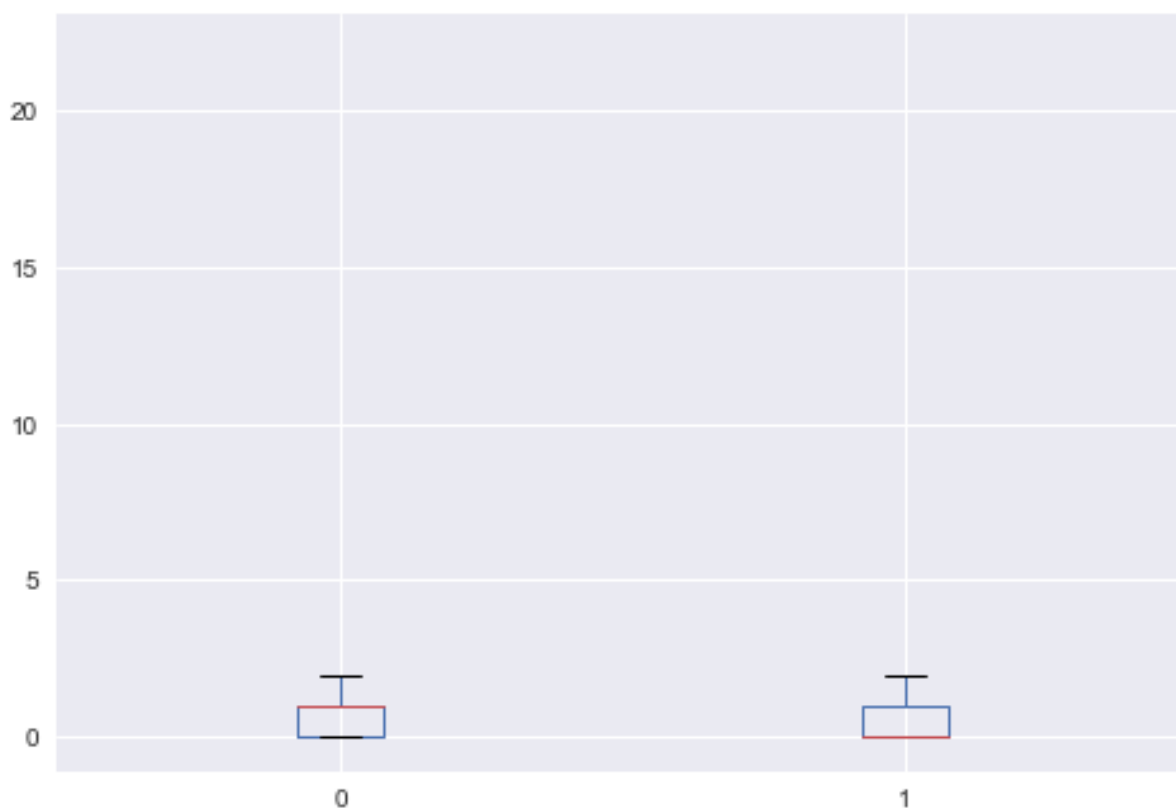


Figure 5: png

Practicum

Gráficos con Python

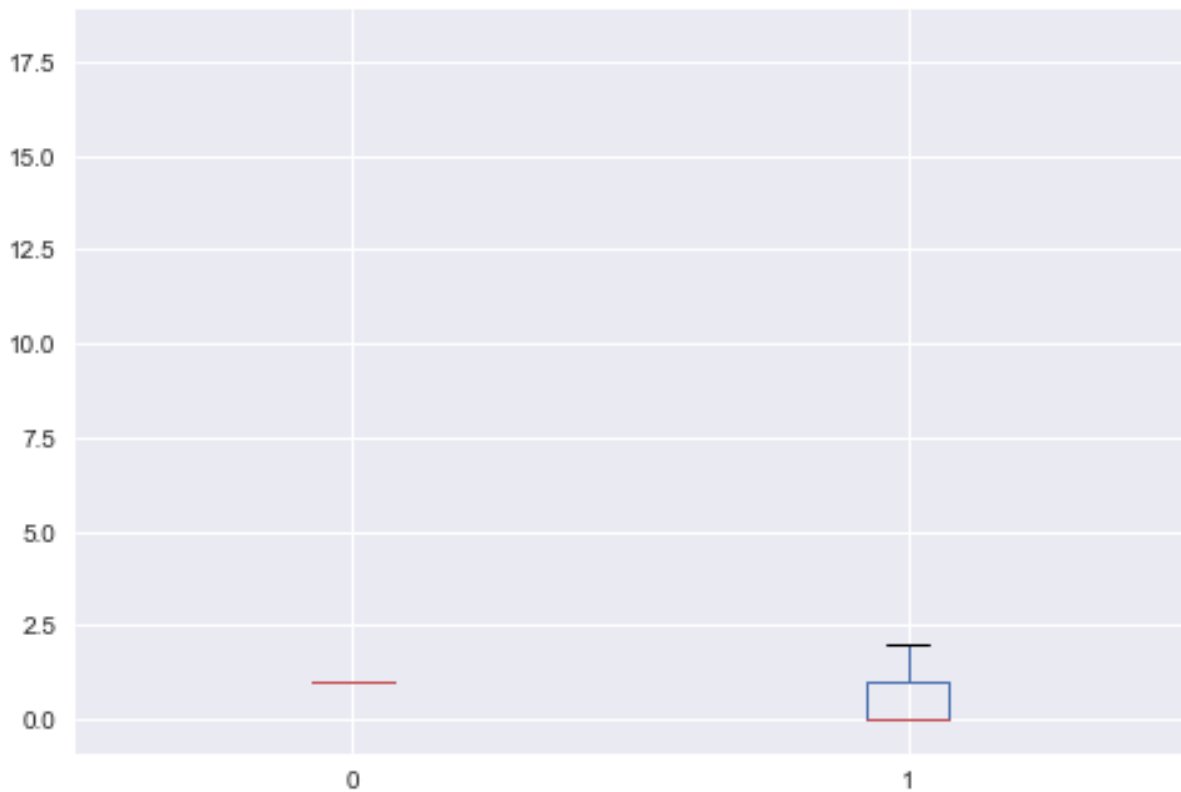


Figure 6: png

```
1 demo = df[['total_bill', 'tip']]
2 print(demo.corr())
```

```
1          total_bill      tip
2 total_bill    1.000000  0.675734
3 tip           0.675734  1.000000
```

```
1 aaa = df.corr()
2
3
4 # Existen módulos mas avanzados que vienen con una serie de diseños muy útiles
5 import seaborn as sns
6
7 sns.heatmap(aaa, annot=True)
8
9 sns.pairplot(df)
```

```
1 <seaborn.axisgrid.PairGrid at 0x1195e39e8>
```

Practicum

Gráficos con Python

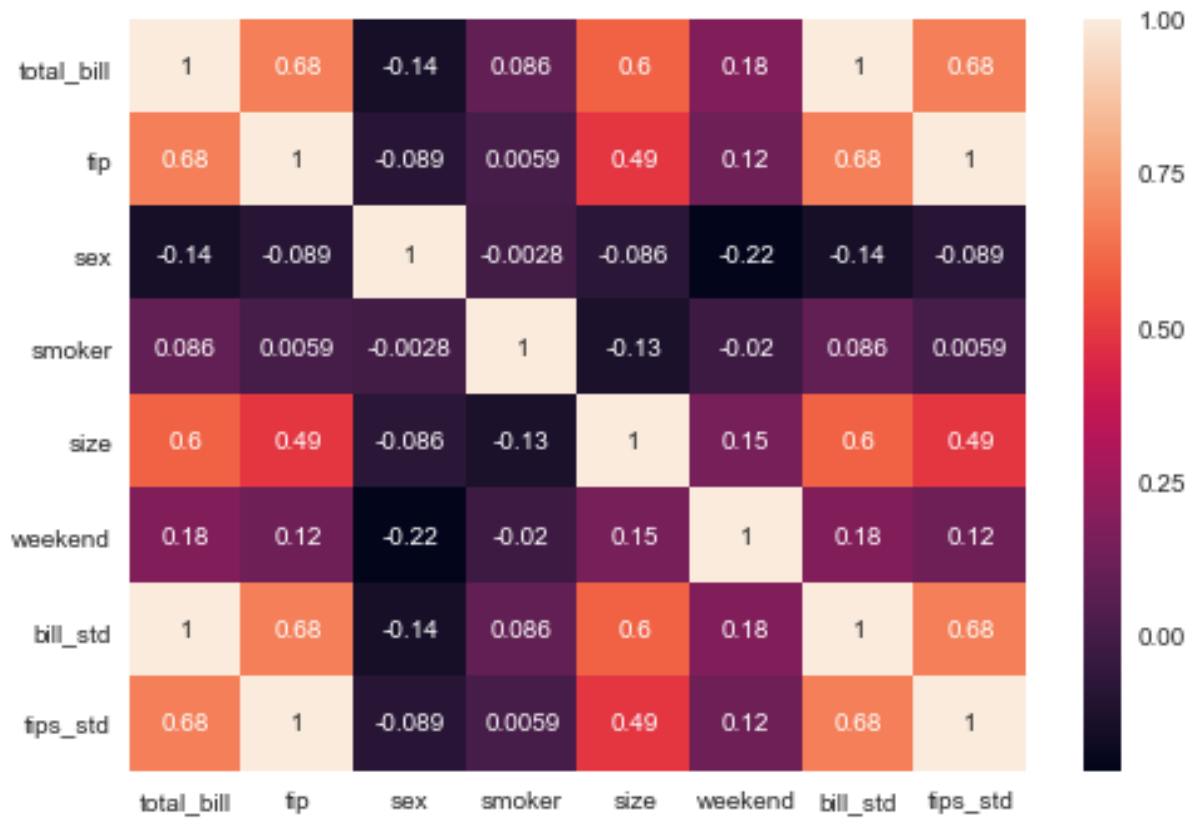


Figure 7: png

Practicum

Gráficos con Python

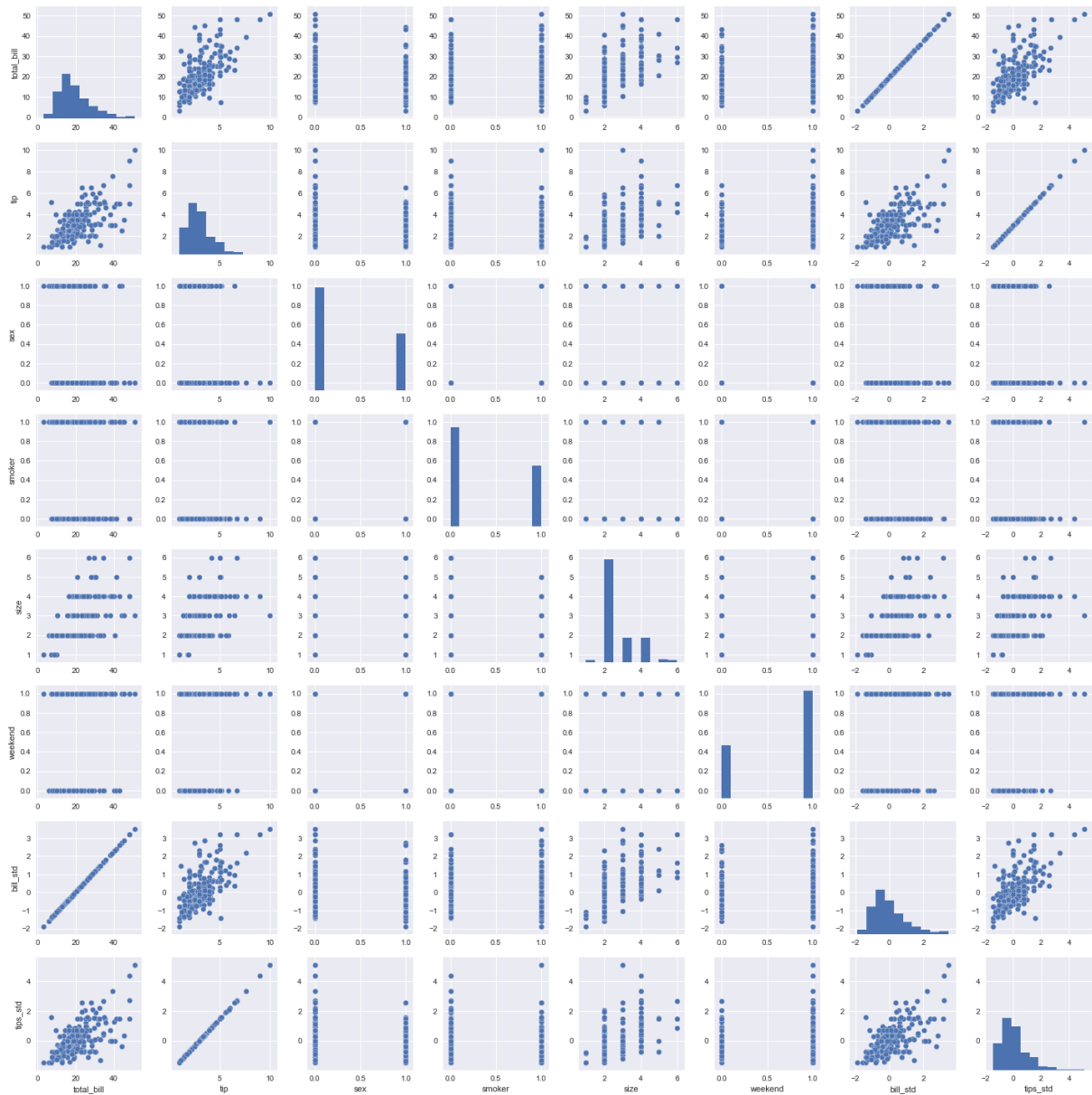


Figure 8: png

```
1 import plotly.plotly as py
2 import plotly.graph_objs as go
3
4 data = [ go.Scatter(x=df.total_bill, y=df.tip, mode='markers',)]
5 py.iplot(data)
```

```
1 data = [ go.Scatter(x = df['size'], y= df.total_bill, mode='markers', opacity
    =.5, )]
2 py.iplot(data)
```

Practicum

Gráficos con Python

Modelación

```
1 import statsmodels.api as statm
2 import statsmodels.formula.api as frm
```

```
1 model_1 = frm.ols('tip ~ size', data=df).fit()
2 print(model_1.summary2())
```

```
1                      Results: Ordinary least squares
2  =====
3  Model:                OLS                Adj. R-squared:    0.236
4  Dependent Variable: tip                AIC:                787.1273
5  Date:                2018-03-02 16:13 BIC:                794.1216
6  No. Observations:    244                Log-Likelihood:    -391.56
7  Df Model:            1                  F-statistic:       76.18
8  Df Residuals:        242                Prob (F-statistic): 4.30e-16
9  R-squared:           0.239              Scale:             1.4621
10 -----
11                Coef.    Std.Err.    t      P>|t|    [0.025    0.975]
12 -----
13 Intercept        1.1691     0.2234    5.2330   0.0000    0.7290    1.6092
14 size             0.7118     0.0816    8.7279   0.0000    0.5512    0.8725
15 -----
16 Omnibus:          81.369          Durbin-Watson:       1.820
17 Prob(Omnibus):    0.000          Jarque-Bera (JB):    273.339
18 Skew:            1.393          Prob(JB):           0.000
19 Kurtosis:        7.373          Condition No.:      9
20 =====
```

```
1 model_2 = frm.ols('tip ~ total_bill', data=df).fit()
2 print(model_2.summary2())
```

```
1                      Results: Ordinary least squares
2  =====
3  Model:                OLS                Adj. R-squared:    0.454
4  Dependent Variable: tip                AIC:                705.0762
5  Date:                2018-03-02 16:13 BIC:                712.0705
6  No. Observations:    244                Log-Likelihood:    -350.54
7  Df Model:            1                  F-statistic:       203.4
8  Df Residuals:        242                Prob (F-statistic): 6.69e-34
9  R-squared:           0.457              Scale:             1.0446
10 -----
```

Practicum

Gráficos con Python

```
11          Coef.    Std.Err.    t    P>|t|    [0.025    0.975]
12 -----
13 Intercept    0.9203    0.1597    5.7612    0.0000    0.6056    1.2349
14 total_bill    0.1050    0.0074   14.2604    0.0000    0.0905    0.1195
15 -----
16 Omnibus:            20.185    Durbin-Watson:            2.151
17 Prob(Omnibus):        0.000    Jarque-Bera (JB):        37.750
18 Skew:                0.443    Prob(JB):                0.000
19 Kurtosis:            4.711    Condition No.:            53
20 =====
```

```
1 model_3 = frm.ols('tip ~ total_bill + sex', data=df).fit()
2 print(model_3.summary2())
```

```
1          Results: Ordinary least squares
2 =====
3 Model:            OLS            Adj. R-squared:    0.452
4 Dependent Variable: tip            AIC:            707.0387
5 Date:            2018-03-02 16:13 BIC:            717.5302
6 No. Observations: 244            Log-Likelihood:    -350.52
7 Df Model:            2            F-statistic:        101.3
8 Df Residuals:        241            Prob (F-statistic): 1.18e-32
9 R-squared:            0.457            Scale:            1.0488
10 -----
11          Coef.    Std.Err.    t    P>|t|    [0.025    0.975]
12 -----
13 Intercept    0.9067    0.1750    5.1817    0.0000    0.5620    1.2513
14 total_bill    0.1052    0.0075   14.1097    0.0000    0.0905    0.1199
15 sex          0.0266    0.1383    0.1924    0.8476   -0.2459    0.2991
16 -----
17 Omnibus:            20.499    Durbin-Watson:            2.149
18 Prob(Omnibus):        0.000    Jarque-Bera (JB):        38.652
19 Skew:                0.447    Prob(JB):                0.000
20 Kurtosis:            4.733    Condition No.:            63
21 =====
```

```
1 plt.scatter(model_2.predict(), df['tip'])
```

```
1 <matplotlib.collections.PathCollection at 0x1c23acd4e0>
```

Practicum

Gráficos con Python

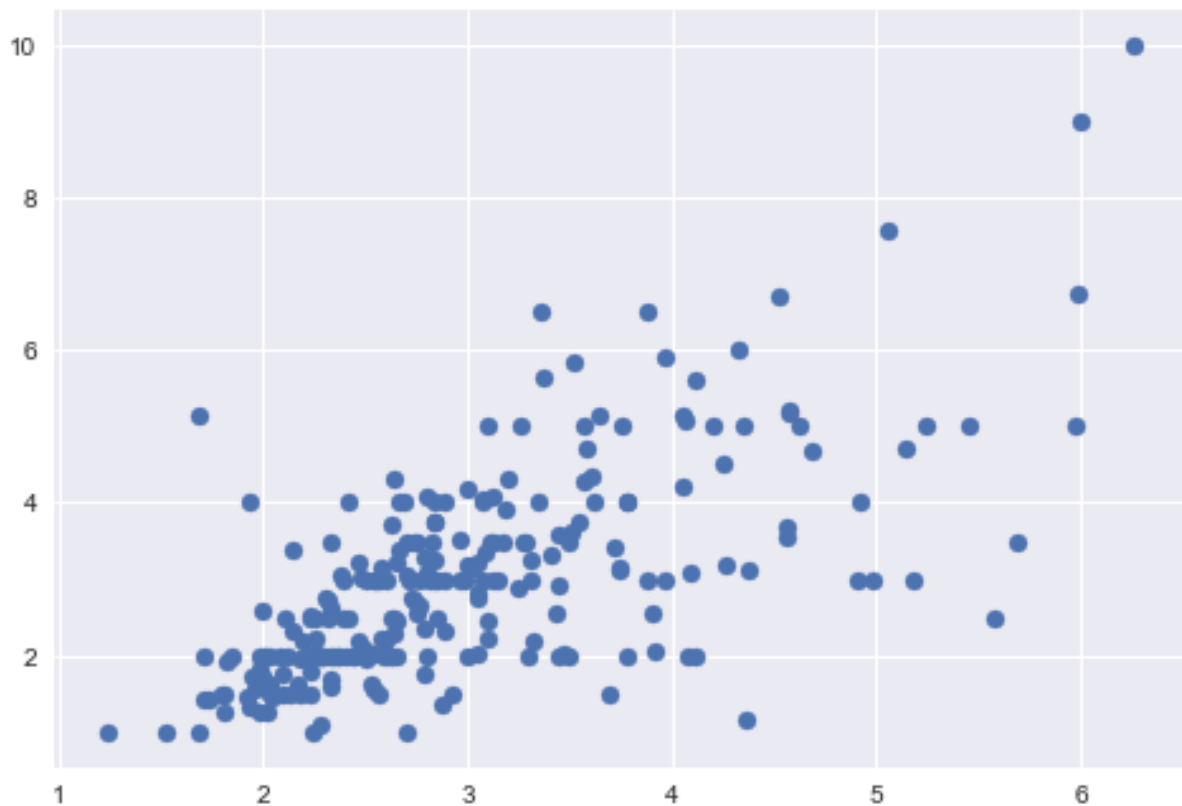


Figure 9: png

```
1 print(frm.ols('tip ~ smoker', data=df).fit().summary())
```

1	OLS Regression Results					
2	=====					
3	Dep. Variable:	tip	R-squared:	0.000		
4	Model:	OLS	Adj. R-squared:	-0.004		
5	Method:	Least Squares	F-statistic:	0.008506		
6	Date:	Fri, 02 Mar 2018	Prob (F-statistic):	0.927		
7	Time:	16:14:00	Log-Likelihood:	-424.95		
8	No. Observations:	244	AIC:	853.9		
9	Df Residuals:	242	BIC:	860.9		
10	Df Model:	1				
11	Covariance Type:	nonrobust				
12	=====					
13		coef	std err	t	P> t	[0.025 0.975]
14	-----					
15	Intercept	2.9919	0.113	26.517	0.000	2.770 3.214
16	smoker	0.0169	0.183	0.092	0.927	-0.343 0.377
17	=====					
18	Omnibus:	79.337	Durbin-Watson:	1.932		

Practicum

Gráficos con Python

```
19 Prob(Omnibus):          0.000   Jarque-Bera (JB):          213.973
20 Skew:                   1.456   Prob(JB):              3.44e-47
21 Kurtosis:               6.545   Cond. No.              2.43
22 =====
23
24 Warnings:
25 [1] Standard Errors assume that the covariance matrix of the errors is
    correctly specified.
```