



Correlation

Hypothesis testing+ Graphic Methods

Session 9

Programación Estadística con Python

Alberto Sanz, Ph.D

albertosanz@bigwaveanalytics.es

MASTER EN DATA ANALYTICS PARA LA EMPRESA

- Hypothesis testing over the relationship of two quantitative variables by the means of regression.
 - Graphic approach (Scatterplot)
 - Numeric approach (r coefficient & p.values)

- Beyond the numbers and plots:
 - Reflections on correlation and non linearity
 - The False Discovery Rate experiment

- Always **DESCRIBE** the two variables involved in the correlation.
 - ▣ Check and validate the integrity of the data prior to any analysis.

- **EXPLORE** of bivariate relation:
 - ▣ Graphically: **Scatterplot**
 - ▣ Numerically: Pearson's **r** & **p.value**

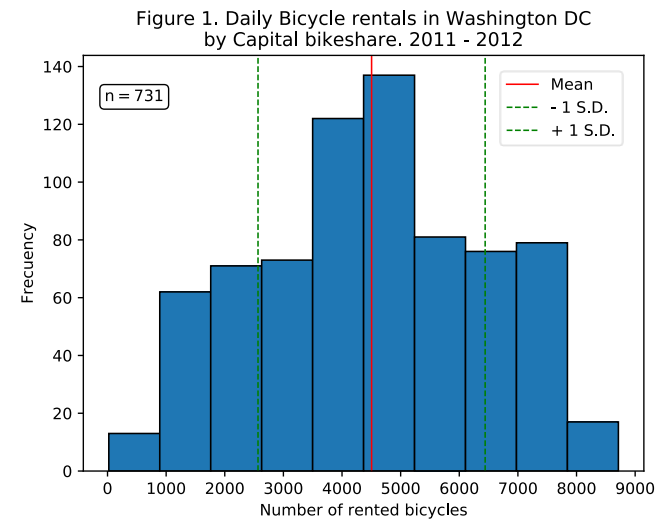
Research Question

4

Why some days are rent *more* bikes?

■ Temperature ?

- H0.: There is no linear association ($r=0$) between the *number of rentals* and the *temperature*.
- H1.: There is a linear association ($r \neq 0$) between the *number of rentals* and the *temperature*.



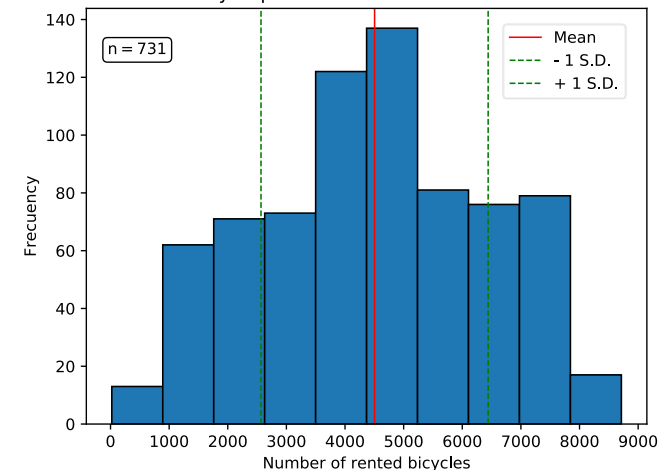
Describing quantitative variables

5

```
x=wbr['cnt']
plt.hist(x, bins=10,
edgecolor='black')
plt.xticks(np.arange(0, 10000,
step=1000))
plt.title('Figure 4. Daily Bicycle
rentals in Washington DC'
'\n'
'by Capital bikeshare.
2011 - 2012')
plt.ylabel('Frecuency')
plt.xlabel('Number of rented
bicycles')
```

```
props = dict(boxstyle='round',
facecolor='white', lw=0.5)
textstr = '$\mathrm{n}=%.0f$'%(n)
plt.text (-50,128, textstr ,
bbox=props)
```

Figure 1. Daily Bicycle rentals in Washington DC by Capital bikeshare. 2011 - 2012



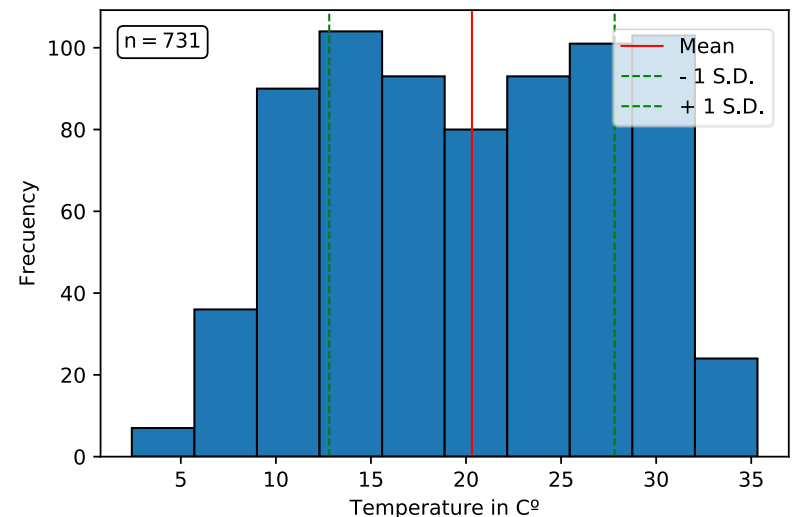
Describing quantitative variables

6

```
##histogram ver4
x=wbr['temp_celsius']
plt.hist(x, bins=10,
edgecolor='black')
#plt.xticks(np.arange(0, 10000,
step=1000))
plt.title('Figure 5. Temperature in
Celsius'

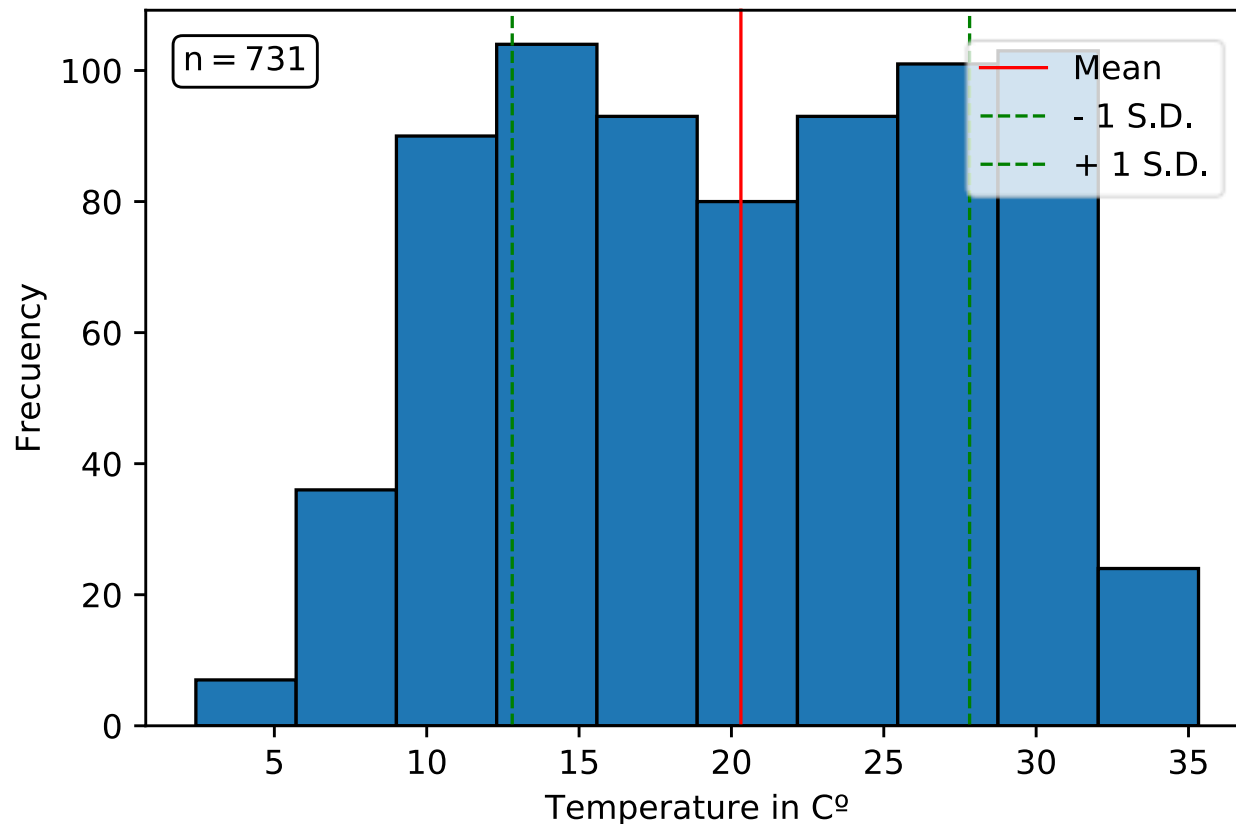
        '\n')
plt.ylabel('Frecuency')
plt.xlabel('Temperature in C°')
props = dict(boxstyle='round',
facecolor='white', lw=0.5)
textstr = '$\mathrm{n}=%.0f$'%(n)
plt.text (2,100, textstr ,
bbox=props)
```

Figure 5. Temperature in Celsius



Describing quantitative variables

Figure 5. Temperature in Celsius



1. Describe the two variables involved in hypothesis

Temperature

Rentals

Figure 5. Temperature in Celsius

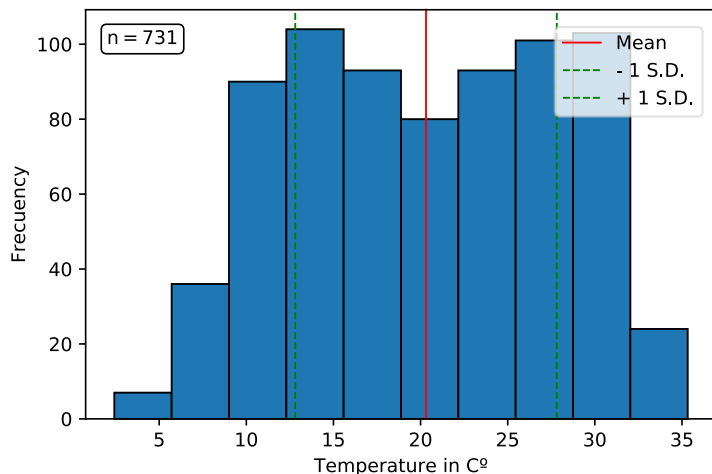
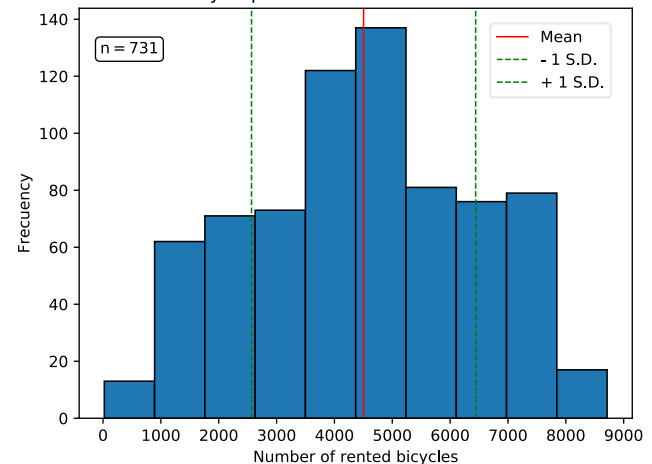


Figure 1. Daily Bicycle rentals in Washington DC by Capital bikeshare, 2011 - 2012



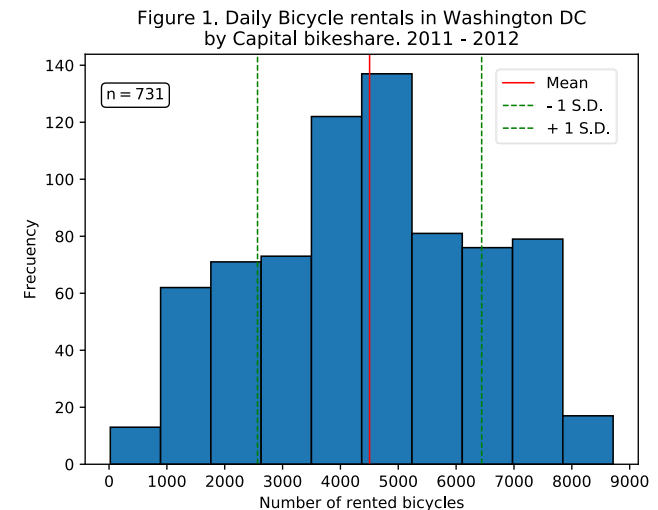
Correlation

9

1. Describe the two variables involved in hypothesis

Windspeed

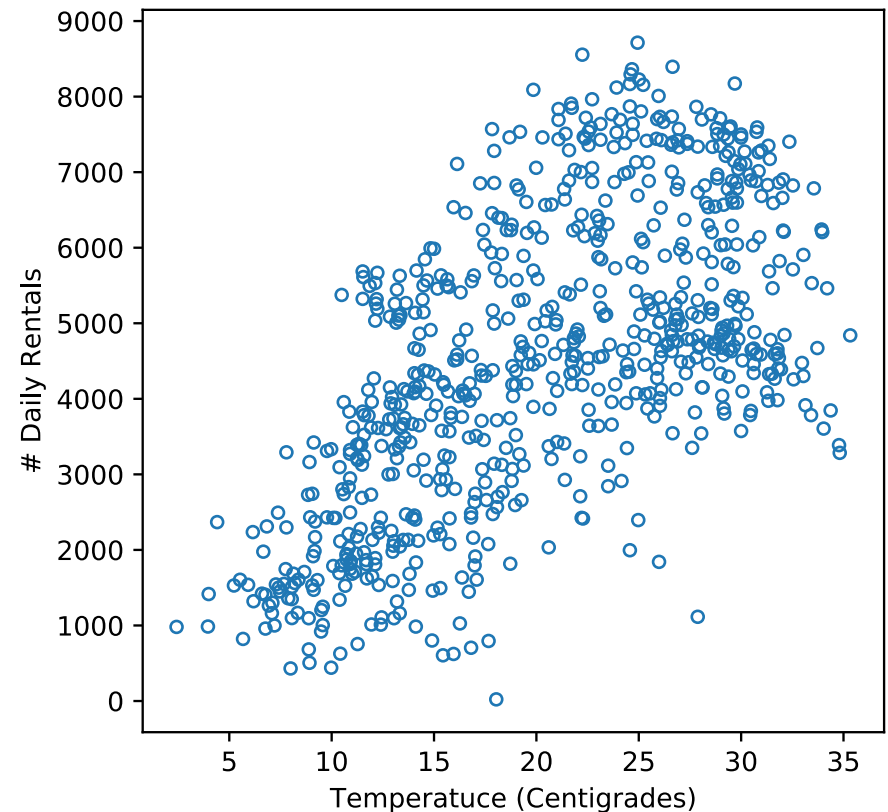
Rentals



2. Scatterplot

```
x=wbr.temp_Celsius  
y=wbr.cnt  
plt.scatter (x,y)
```

Figure 9. Daily bicycle rentals, by temperature.



Regression

3. Pearson's r

```
from scipy.stats.stats import pearsonr  
res = pearsonr(x, y)  
print (res)
```

```
[1] (0.62749400903349195, 2.8106223975901415e-81)
```



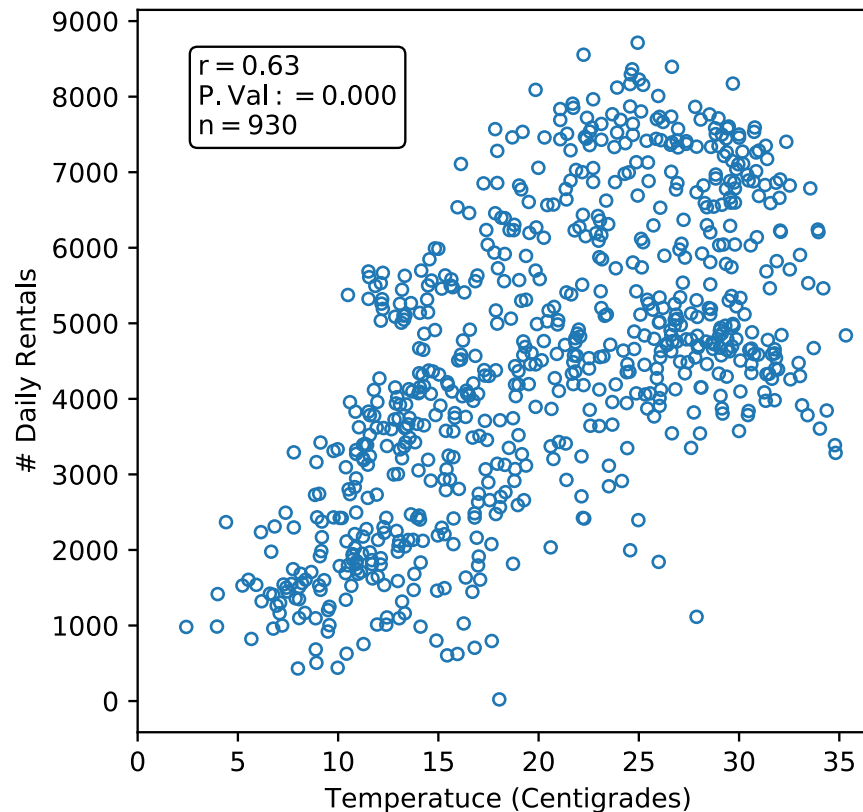
This is
Pearson's r



This is
The P.Value

Scatterplot + Pearson's r + test

Figure 9. Daily bicycle rentals, by temperature.



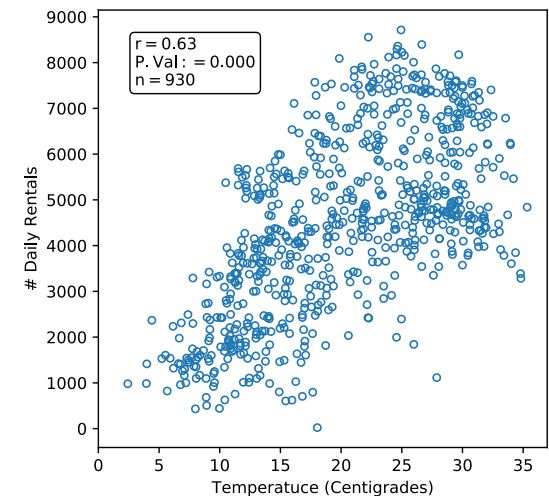
Conclusion

Conclusion:

As P. Value < 0.000

We can reject H_0 with a confidence higher than 99.9

Figure 9. Daily bicycle rentals, by temperature.



✗ H_0 .: There is no linear association between the *number of rentals* and the *temperature*.

✓ H_1 .: There is a linear association between the *number of rentals* and the *temperature*.

Questions?

Thank you !

Alberto Sanz