



Unidad 05: Análisis exploratorio de un conjunto de datos

Aprendizaje Automático

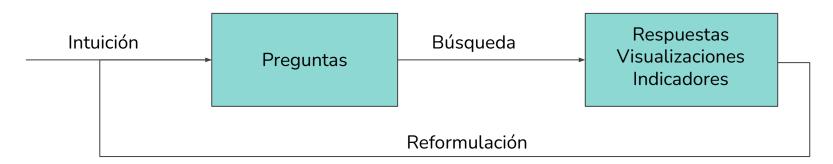
Docentes:
Diego P. Durante
ddurante@frba.utn.edu.ar

Ramiro Verrastro gramiro verrastro gramiro verrastro gramiro verrastro gramiro verrastro gramino verrastro verrastro gramino verrastro v

Tema 1 - Introducción

Proceso de exploración

- Generar preguntas utilizando la intuición
- Buscar las respuestas mediante visualización u otros indicadores
- Reformular y generar nuevas preguntas utilizando lo aprendido



- Tema 1 Introducción
- Tema 2 Búsqueda de insights

EDA



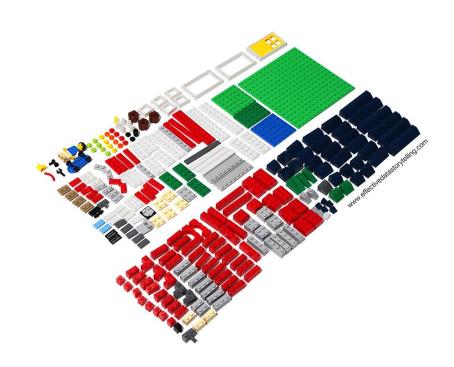
1 - Colección de datos



2 - Preparación de datos



3 - Visualización de datos



4 - Análisis de datos



5 - Conclusiones desde datos



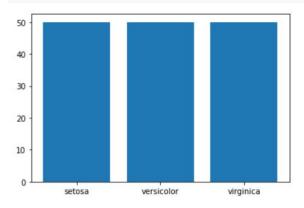
- Tema 1 Introducción
- Tema 2 Búsqueda de insights
- Tema 3 Análisis de distribuciones





Distribución de clases

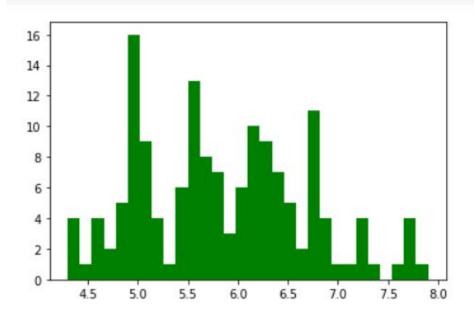
```
cantidad_de_ejemplos = df.groupby(['target_name']).size()
plt.bar(cantidad_de_ejemplos.index,cantidad_de_ejemplos.values)
plt.show()
```







Distribución de un feature



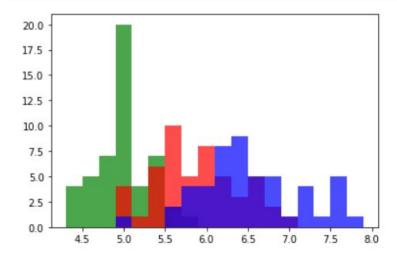


Distribución



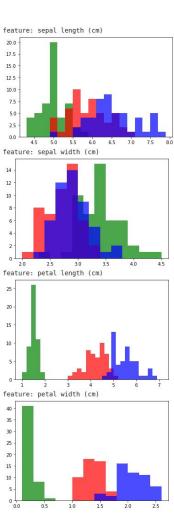
Distribución de features, separadas por clase

```
for clase,color in zip(['setosa', 'versicolor', 'virginica'],['g','r','b']):
   data = df[df['target_name']==clase]['sepal length (cm)']
   binwidth = 0.2
   plt.hist(data, bins=np.arange(min(data), max(data) + binwidth, binwidth), facecolor=color, alpha=0.7)
```



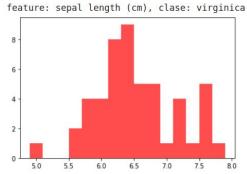


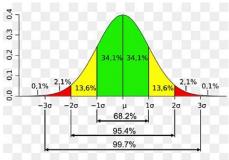
```
for feature in ['sepal length (cm)', 'sepal width (cm)', 'petal length (cm)', 'petal width (cm)']:
  data = df[feature]
  var feat = data.var()
  print("Varianza del feature {}: {:.3f}".format(feature,var feat))
  for clase in ['setosa', 'versicolor', 'virginica']:
    data = df[df['target name']==clase][feature]
    var feat = data.var()
    print("Varianza del feature {}, para la clase {}: {:.3f}".format(feature,clase,var feat))
Varianza del feature sepal length (cm): 0.686
Varianza del feature sepal length (cm), para la clase setosa: 0.124
Varianza del feature sepal length (cm), para la clase versicolor: 0.266
Varianza del feature sepal length (cm), para la clase virginica: 0.404
Varianza del feature sepal width (cm): 0.190
Varianza del feature sepal width (cm), para la clase setosa: 0.144
Varianza del feature sepal width (cm), para la clase versicolor: 0.098
Varianza del feature sepal width (cm), para la clase virginica: 0.104
Varianza del feature petal length (cm): 3.116
Varianza del feature petal length (cm), para la clase setosa: 0.030
Varianza del feature petal length (cm), para la clase versicolor: 0.221
Varianza del feature petal length (cm), para la clase virginica: 0.305
Varianza del feature petal width (cm): 0.581
Varianza del feature petal width (cm), para la clase setosa: 0.011
Varianza del feature petal width (cm), para la clase versicolor: 0.039
Varianza del feature petal width (cm), para la clase virginica: 0.075
```

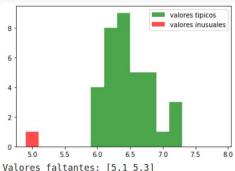


Valores Típicos, inusuales, faltantes

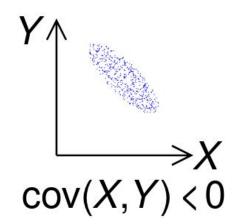
```
feature='sepal length (cm)'
clase = 'virginica'
print('feature: {}, clase: {}'.format(feature, clase))
data = df[df['target name']==clase][feature]
data typ = data[(data > (data.mean()-data.std()))) & (data < (data.mean()+data.std()))]</pre>
data out = data[(data < (data.mean()-2*data.std()))) & (data < (data.mean()+2*data.std()))]</pre>
binwidth = 0.2
n, bins, patches = plt.hist(data, bins=np.arange(min(data), max(data) + binwidth, binwidth), facecolor='r', alpha=0.7)
plt.show()
plt.hist(data typ, bins=np.arange(min(data), max(data) + binwidth, binwidth), facecolor='g', alpha=0.7, label='valores tipicos')
plt.hist(data out, bins=np.arange(min(data), max(data) + binwidth, binwidth), facecolor='r', alpha=0.7, label='valores inusuales')
plt.show()
print("Valores faltantes: {}".format(bins[np.where(n == 0)]))
```

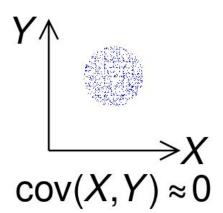


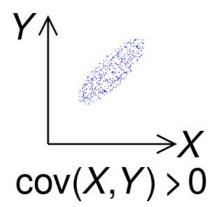




Covarianza







Matriz de covarianza

$$cov_{x,y} = rac{\sum (x_i - ar{x})(y_i - ar{y})}{N-1}$$

 $cov_{x,y}$ = covariance between variable a and y

 x_i = data value of x

i = data value of y

 $ar{x}$ = mean of x

 $ar{y}$ = mean of y

N = number of data values

df.drop(columns=['target']).cov()

	sepal length (cm)	sepal width (cm)	petal length (cm)	petal width (cm)
sepal length (cm)	0.685694	-0.042434	1.274315	0.516271
sepal width (cm)	-0.042434	0.189979	-0.329656	-0.121639
petal length (cm)	1.274315	-0.329656	3.116278	1.295609
petal width (cm)	0.516271	-0.121639	1.295609	0.581006

Matriz de correlación

$$rac{\displaystyle \sum_{i=1}^{N} (x_i - ar{x}) \cdot (y_i - ar{y})}{\sqrt{\displaystyle \sum_{i=1}^{N} (x_i - ar{x})^2} \cdot \sqrt{\displaystyle \sum_{i=1}^{N} (y_i - ar{y})^2}}$$

df.drop(columns=['target']).corr()

	sepal length (cm)	sepal width (cm)	petal length (cm)	petal width (cm)
sepal length (cm)	1.000000	-0.117570	0.871754	0.817941
sepal width (cm)	-0.117570	1.000000	-0.428440	-0.366126
petal length (cm)	0.871754	-0.428440	1.000000	0.962865
petal width (cm)	0.817941	-0.366126	0.962865	1.000000

- Tema 1 Introducción
- Tema 2 Búsqueda de insights
- Tema 3 Análisis de distribuciones
- Tema 4 Visualización

Visualización de datos

Gráfico de líneas

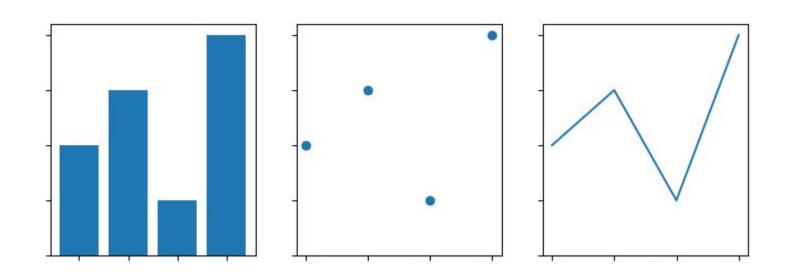


Gráfico de líneas

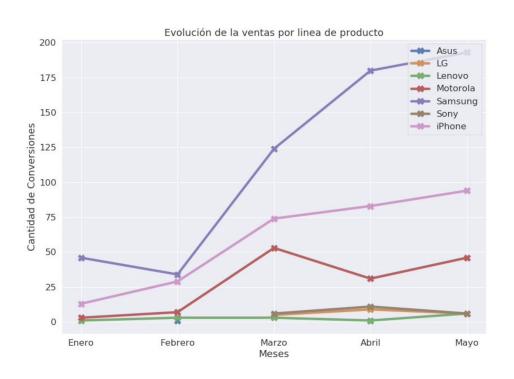


Gráfico de barras

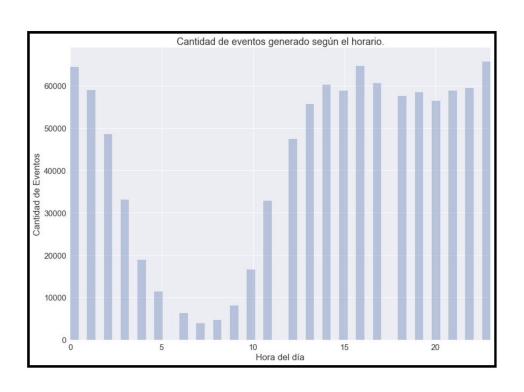


Gráfico de torta

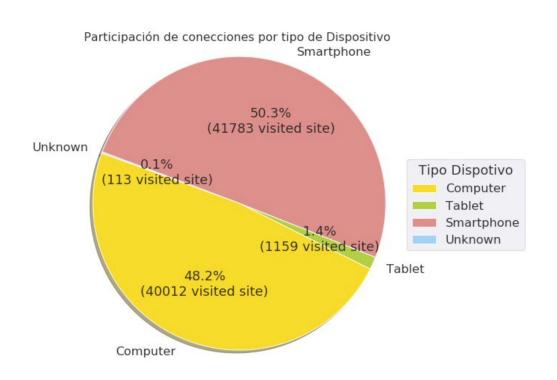


Gráfico de torta

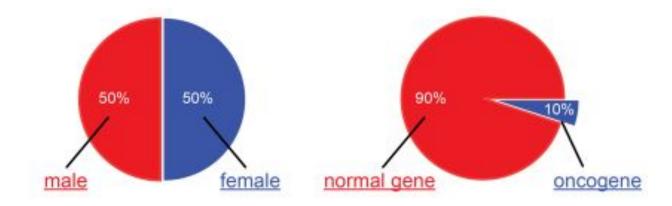
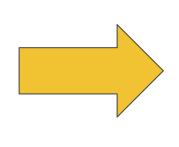


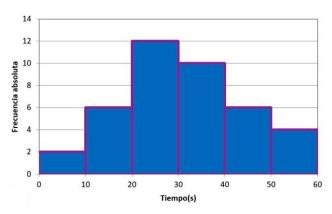
Gráfico en mosaico



Histograma

[0 - 10)	5
[10 - 20)	15
[20 - 30)	25
[30 - 40)	35
[40 - 50)	45
[50 - 60]	55
Total	





Histograma

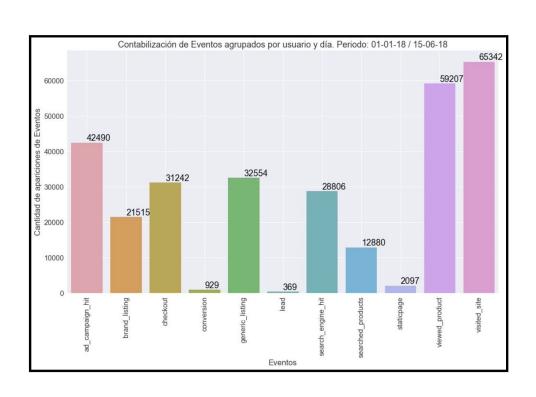


Diagrama de Pareto

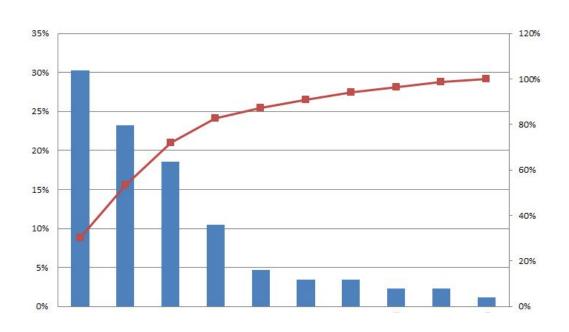
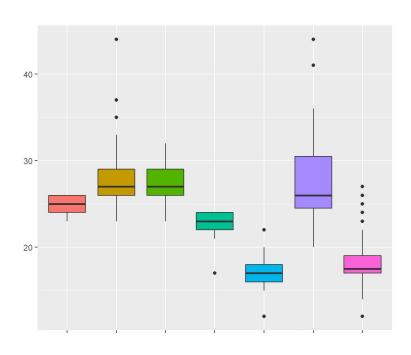
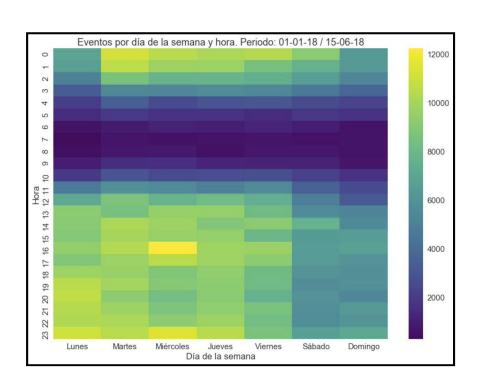


Gráfico de caja



Mapa de calor



Mapa de calor





Mapa de calor

Mapa con los Quiosues



Mapa de calor con Zoom sobre la zona de San Pablo. Eventos generados.

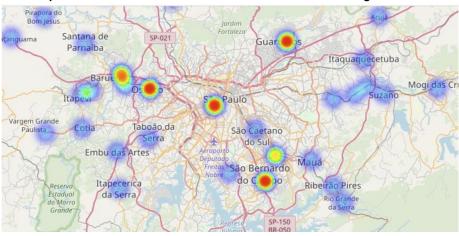


Diagrama de árbol

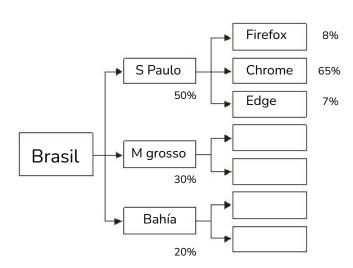
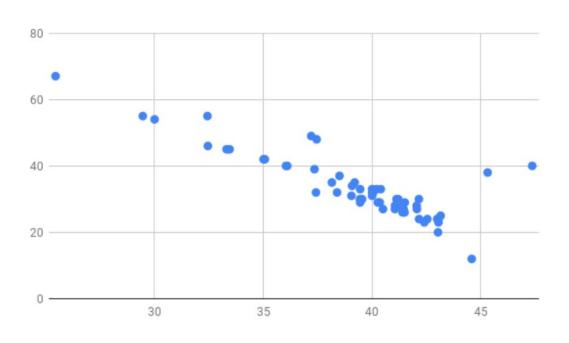
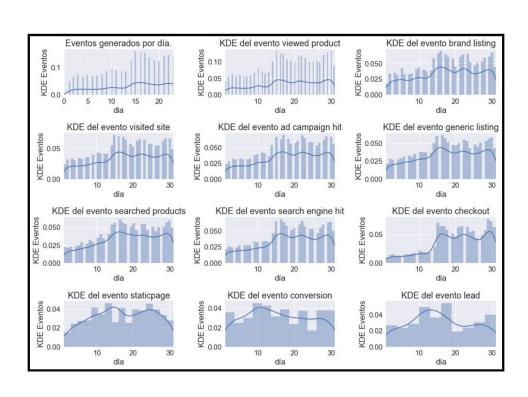


Gráfico de dispersión



Subplots





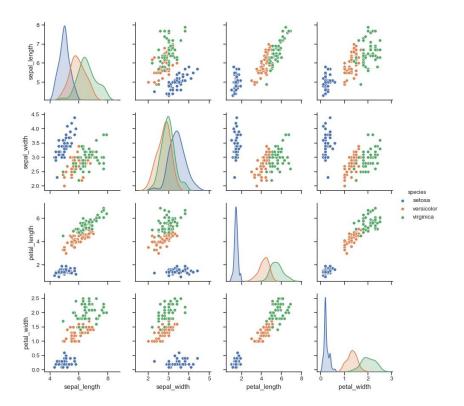


```
Trosy cheery fortuitous convenient entertaining apposite promising suitable suitable
```

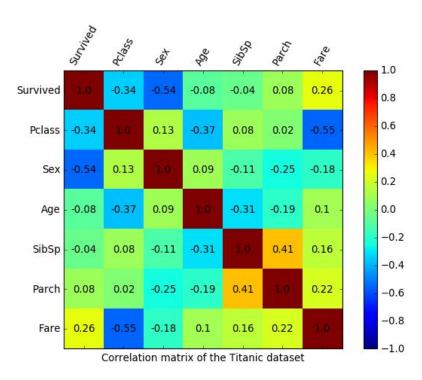
shutterstock.com · 53816176

Matriz de gráficos de dispersión

- setosa
- versicolor
- virginica

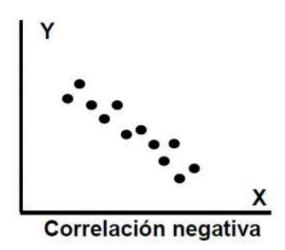


Matriz de correlación



Matriz de correlación





Preguntas?



