

“DATOS y PRE-PROCESADO”

Reconocimiento de Patrones

MAESTRÍA EN ELECTRÓNICA

Profesor: MSc. Felipe Meza

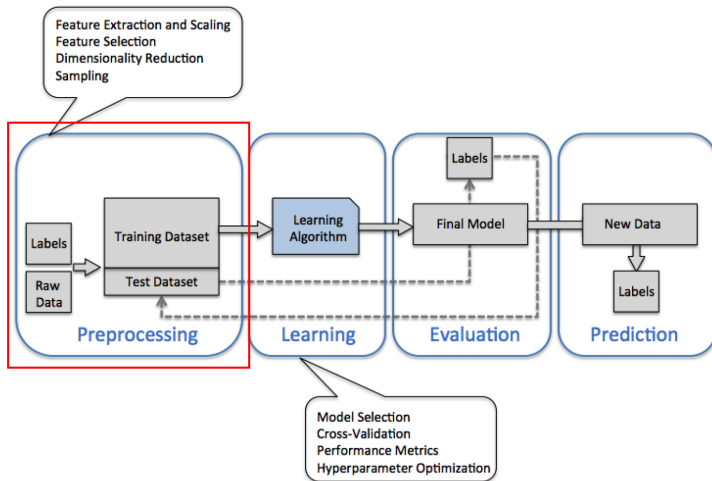


June 3, 2021

Agenda

- 1 Metodología de Diseño.
- 2 Pre-procesado.
- 3 Preparación de los datos.
- 4 Algunas Tareas del Pre-procesado.
- 5 Análisis exploratorio de los datos (EDA).
- 6 Valores faltantes.
- 7 Outliers.
- 8 Datos no-balanceados.
- 9 Transformación de datos.
- 10 Reducción de dimensiones.

Metodología de Diseño



Pre-procesado

- Consiste en identificar partes o componentes del conjunto de datos que sean **incompletas**, imprecisas, incorrectas o irrelevantes, de manera tal que puedan ser **reemplazadas**, modificadas o removidas.
- Puede implicar también la **transformación** de los datos “puros” a otros formatos que faciliten su manejo por parte de los algoritmos de minería de datos.
- Incluye también la **reducción** de los datos a menores dimensiones para agilizar su procesamiento.
- En inglés varios **términos** se refieren a tales tareas: data preparation, cleaning, pre-processing, cleansing, wrangling.

Preparación de los datos

- En las metodologías de diseño generalmente en las **primeras etapas**, corresponde llevar a cabo las tareas de selección de datos, pre-procesado o transformación.
- En `python`, se recurre al uso de **librerías** como `pandas` que resultan ser muy buenas para las tareas asociadas a la preparación de los datos.
- Las labores de preparación no son un componente integral de los algoritmos de aprendizaje, sin embargo puede tomar un **tiempo** considerable dependiendo del conjunto de datos a analizar (80%-90% del proceso), por lo que se debe prestar especial atención.

Algunas Tareas del Pre-procesado

- Estandarizar, Normalizar.
- Análisis exploratorio de los datos.
- Valores faltantes.
- Outliers.
- Datos no-balanceados.
- Transformación de datos.

Normalizar, Estandarizar

- **Normalizar** (1) es llevar los datos a una nueva escala en un rango entre 0 y 1. Recomendado en casos donde los datos tengan múltiples escalas y donde los algoritmos sean sensibles a la escala.
- **Estandarizar** (2) consiste en llevar la distribución de los datos a una media de 0 y una desviación estándar de 1. Recomendado en casos donde el algoritmo es sensible a una distribución normal.

$$z = \frac{x - \min(x)}{\max(x) - \min(x)} \quad (1)$$

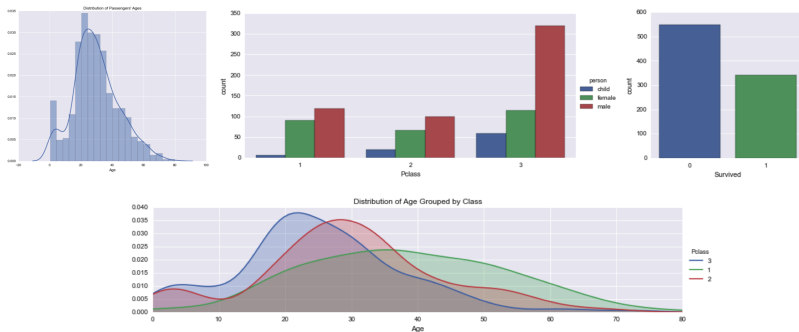
$$z = \frac{x - \mu}{\sigma} \quad (2)$$

Análisis exploratorio de los datos (EDA)

- Es la práctica del uso de métodos cuantitativos y **visuales** para comprender mejor un conjunto de datos sin tener que asumir hechos.
- Arrojar el conjunto de datos a un algoritmo y esperar los mejores resultados, **NO** es la mejor estrategia.
- Usualmente se lleva cabo una o varias de las siguientes actividades:
 - **Visualización** de un resumen estadístico del conjunto de datos.
 - **Exploración** visual de cualquier relación que pueda tener cada atributo con la clase que nos interesa predecir.
 - Mediante diagramas de dispersión **observar** cualquier tipo de agrupamiento que se pueda presentar en los datos.

Análisis exploratorio de los datos (EDA)

Ejemplo de EDA con conjunto de datos TITANIC



Valores faltantes

- No existe un **método universal** para instancias con valores faltantes.
- Algunas técnicas **comúnes** son:
 - Eliminar instancias.
 - Eliminar atributos.
 - Calcular “media” del atributo faltante.
 - Calcular “mediana” del atributo faltante.
 - Calcular “moda” del atributo faltante.
 - Usar regresión para estimar el valor del atributo faltante.



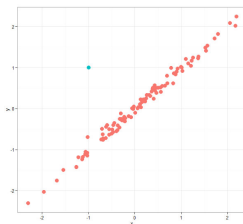
Valores faltantes

pandas

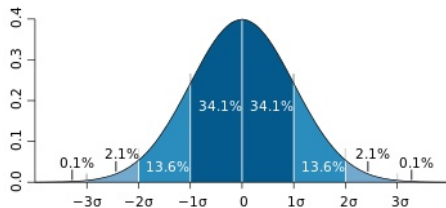
```
1 # Drop the col where all elements are missing values:
2 df.dropna(axis=1, how='all')
3
4 # Drop the col where any of the elements're missing values
5 df.dropna(axis=1, how='any')
6
7 # Keep only the rows which contain 2 missing values max
8 df.dropna(thresh=2)
9
10 # Fill all missing values with the mean of the column
11 df.fillna(df.mean())
12
13 # Fill any missing value in col 'A' with the col median
14 df['A'].fillna(df['A'].median())
```

Outliers

- Así como hay situaciones donde es importante mantenerlos en otros casos es necesario **eliminarlos**.
- Se recomienda hacer un análisis previo basado en la **naturaleza** de los datos.
- Algunas técnicas comunes son:
 - Removerlos usando desviación estándar (PYTHON).
 - Removerlos usando percentiles (pandas).



Outliers - Desviación Estándar



Regla en distribución normal: Remover datos en $(mean + 2\sigma)$ y $(mean - 2\sigma)$

Regla 68-95-99 ó 3σ

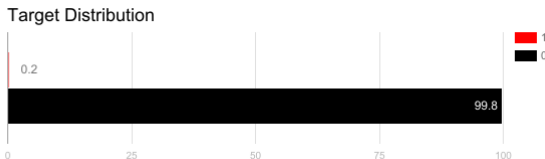
Outliers - Desviación Estándar

PYTHON

```
1 import numpy
2
3 arr = [10, 386, 479, 491, 501 ... 411, 399, 363, 19, 543]
4
5 elements = numpy.array(arr)
6
7 mean = numpy.mean(elements, axis=0)
8 sd = numpy.std(elements, axis=0)
9
10 final_list = [x for x in arr if (x > mean - 2 * sd)]
11 final_list = [x for x in final_list if (x < mean + 2 * sd)]
12 print(final_list)
```

Datos no-balanceados

- Ocurre cuando una clase de datos en el conjunto, posee una mayoría importante de la cantidad de datos e.g un conjunto de datos de 2 clases, donde: $CLASE1 = 98\%$ y $CLASE2 = 2\%$.



Datos no-balanceados

- Algunas técnicas comunes son:
 - Usar otras métricas diferentes al porcentaje de exactitud, por ejemplo:
 - Precision/Specificity: cuantas instancias seleccionadas son relevantes.
 - Recall/Sensitivity: cuantas instancias relevantes son seleccionadas.
 - F1 score: media armónica de "precision" y "recall".
 - Muestreo de datos:
 - sub-muestreo: Eliminar instancias abundantes (sólo si hay suficientes datos).
 - sobre-muestreo: Generar instancias faltantes (mediante métodos de repetición o generación, solo en caso de que sea posible)
 - Descomponer el conjunto de datos en subconjuntos.
 - Hacer clustering de grupo abundante.

Transformación de datos

- Ocurre cuando **transformamos** un valor z_i en y_i mediante una función $f()$ de forma tal que $y_i = f(z_i)$.
- Se hace con el fin de alinear los datos con alguna suposición estadística, mejorar la interpretación de los datos o bien obtener gráficos de mejor apariencia.
- Técnica muy común: **One Hot Encode**, permite convertir datos categóricos en numéricos (vectores binarios).

Sample	Category	Numerical
1	Human	1
2	Human	1
3	Penguin	2
4	Octopus	3
5	Alien	4
6	Octopus	3
7	Alien	4

Datos Categóricos

Sample	Human	Penguin	Octopus	Alien
1	1	0	0	0
2	1	0	0	0
3	0	1	0	0
4	0	0	1	0
5	0	0	0	1
6	0	0	1	0
7	0	0	0	1

Vectores Binarios

$$\begin{bmatrix} 0 & 0 & 0 & 1 & 0 \end{bmatrix} \times \begin{bmatrix} 17 & 24 & 1 \\ 23 & 5 & 7 \\ 4 & 6 & 13 \\ 10 & 12 & 19 \\ 11 & 18 & 25 \end{bmatrix} = \begin{bmatrix} 10 & 12 & 19 \end{bmatrix}$$

One Hot Encode

Algunas Fuentes de DATASETS

- Kaggle Datasets
 - www.kaggle.com
- UCI Machine Learning Repository
 - <https://archive.ics.uci.edu/ml/index.php>
- Google's Datasets
 - <https://toolbox.google.com/datasetsearch>
- Microsoft Datasets
 - <https://msropendata.com/>
- Awesome Public Datasets
 - <https://github.com/awesomedata/awesome-public-datasets>

Questions?



Felipe Meza - fmeza@itcr.ac.cr