

Juan Pablo Padilla Martin
Examen Minería de Datos

① **Categorías**

Nominal: símbolos o nombre de las cosas

Ordinal: Orden o posición en el conjunto de datos, evalúa objetivamente cualidades del objeto

Binario: representa solo dos categorías, simétricas y asimétricas (falso-verdadero, macho-hembra, etc)

Númericos

Intervalo: medida en escala de igual tamaño, tienen orden, pueden ser positivos o negativos, no poseen cero verdadero

Razón: Posee cero verdadero, es un valor múltiplo de otros, los valores están ordenados y se pueden calcular diferencias entre ellos.

②

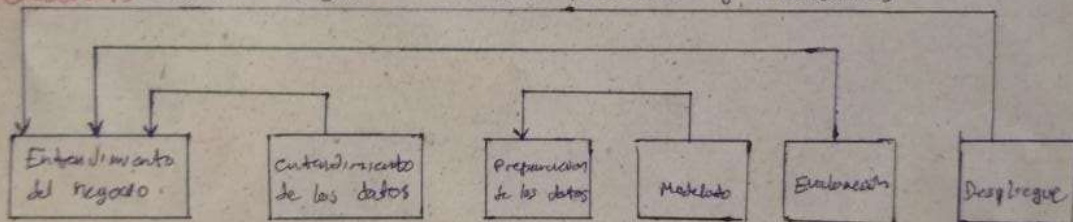
Clasificación: clasifica y etiqueta datos en clases predefinidas

Regresión: Partiendo de los datos genera una función predictiva

Agrupamiento: creación de grupos y conjuntos

Asociación: encuentra conjuntos de datos frecuentes y relacionados

③



CRISP-DM

⑤

1- **Si** hay correlación

2- Escolaridad e ingreso tienen poca correlación
Es pertinente eliminar los de 37 años por que es el más lejano al conjunto de datos

$$X < 50K$$

$$X > 100K$$

$$50K < X < 100K$$

| prepa | uni | sec | maes | Σ |
|-------|-----|-----|------|----------|
| 1 | 0 | 1 | 1 | 3 |
| 1 | 2 | 0 | 0 | 3 |
| 0 | 1 | 1 | 0 | 2 |

$$n = 8$$

$$\text{media A} = 55.64$$

$$\text{media B} = 252.11$$

$$\text{D. standar A} = 11.42$$

$$\text{D. standar B} = 30.53$$

$$\sum_{i=1}^n \frac{(x_i - 55.64)(y_i - 252.11)}{(D_{stand A})(D_{stand B})} = \frac{1625.56}{2780.22} = 0.5828$$

$$C = 0.5828$$

$$\text{Correlation} = 0.5828$$