# EXPLORATORY DATA ANALYSIS

## Exploracion inicial

`.head()` y `.info()`

Para variables categoricos, ver categorias y cuenta → `.value_counts()`

`.describe()`

Ver variables numéricos con histogramas

## Data Validation

Para ver rapidamente los tipos de datos → `books.dtypes`

Convertimos al tipo deseado

Para comparar → `books['col'].isin(['valor_A', 'valor_B'])`

for each filo retorna true si el valor esta en la liste.

Para ver solo las columnas numéricas → `books.select_dtypes('number')`

## Metiendo datos en los nulos

Para ver la cantidad de valores nulos en cada columna → `df.isna().sum()`

Para encontrar el numero de valores máximo por columna nulo →

```
threshold = len(df) * 0.05
cols_to_drop = df.columns[df.isna().sum() <= threshold]
df.dropna(subset = cols_to_drop, inplace=True)
```

## Convirtiendo y analizando datos categóricos

Para comprobar que una fila contiene esto ó lo otro → `df['col'].str.contains("esto | lo otro")`

Para encontrar multiples frases →

```
categorias = ['A', 'B', 'C'..... ]
conditions = [ (df['col'].str.contains(frase_A)),
             _ _ _ _ _ _ _ ]
df['New_col'] = np.select(conditions, categorias, default='Other')
```

Deben estar en el mismo orden

## Trabajando los outliers

Preguntarnos → ¿Porque estan?
            ↳ ¿Estan bien los datos?

# EXPLORATORY DATA ANALYSIS

## Patrones temporales

Para extraer datos de un datetime →

```
df['col'].dt.month
 "    . dt.year
 "    . dt.day
```

## Correlacion

Usar el coeficiente de Pearson

```
df.corr()
```

## Consideraciones para datos categoricos

Para ver el porcentaje de cada categoria → `df['col'].value_counts(normalize=True)`

Para agregar valores usado crosstab → `pd.crosstab(pbres['Salida'], pbres['Destino'], values=pbres['Precio'], aggfunc='median')`