

# CLEANING DATA

Para ver fechas fuera del rango posible

Ver fecha de hoy  $\rightarrow$  `datetime.date.today()`  
 $\rightarrow$  Filtramos fuera de esto

## Valores duplicados

Para encontrar duplicados  $\rightarrow$

`dup = df.duplicated()`  $\rightarrow$  Selecciona booleanos.  
 $\rightarrow$  Lista de columnas donde haya duplicados  
`subset = [...]`  
`keep = ...`  $\rightarrow$  que valor se elija 'first', 'last' ...

Para borrar duplicados  $\rightarrow$

`df.drop_duplicates(inplace=True)`  $\rightarrow$  Así solo tenemos duplicados completos  
 $\rightarrow$  Tiene también `subset` y `keep`

## Datos categóricos

Para borrar datos con categorías que no existen

`valos = set(df['colcat']).difference(categorías posibles)`

Consistencia  $\rightarrow$  Mayúsculas  
 $\rightarrow$  Espacios

Para agrupar columnas numéricas en categorías

`df['colcat'] = pd.cut(df['col'], q=3, labels=['a', 'b', 'c'])`  $\rightarrow$  Cuantos grupos  
Si tenemos muchos los bloques  $\rightarrow$   
`bloques [0, 100, 250, ...]`  
`... = pd.cut(df['col'], bins = bloques, labels = ...)`

## Uniformidad en los datos

Tienen que estar en las mismas unidades y formato (en caso de fechas)

Para hacer en una columna de datetime algo en concreto  $\rightarrow$

`df['col'].dt.strftime('%Y')`

## Cross-field Validation

Consiste en usar varias columnas para comprobar la integridad de los datos  $\rightarrow$  por ej

`col.1 + col.2 mod.3 = col. total`

Por ello podemos sumarle a axis=1

# CLEANING DATA

## Data NAs

Por ver su distribución ->

```
import missing as msno
msno.matrix(df)
plt.show()
```

## Comparando strings

Para calcular la distancia entre strings ->

```
from difflib import process
matches = process.extract('...', df[col], limit=df.shape[0])
```

## Creando pares

Para generarlos por una columna ->

```
import recordlinkage
indexer = recordlinkage.Indexer()
indexer.block('col')
pairs = indexer.index(df1, df2)

comp = recordlinkage.Compare()
comp.exact('col.df1', 'col.df2',)
comp.string('...', '...', threshold=0.85)
potential_matches = comp.compare(pairs, df1, df2)
```

Para mi, los matches ->

```
matches = get_matches(matches)
dup = matches.index.get_loc(values[i])
df1 = df2[df2.index.isin(dup)]
df = df1.append(df2)
```