# Data Cleaning → how to make the dataset ready to be Consumed.
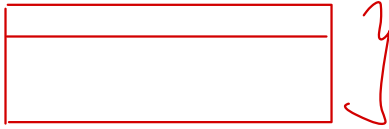
Analysis

Format

txt → Csv

Cout, groupby

right

# Challenges

1. Missing Values
2. Date format
3. Comma Separated Values (Nested Values)
4. Duration

| Title | Director | Cost | listed-in | Country |
|-------|----------|------|-----------|---------|
| ABC | $D_1, D_2$ | $C_1, C_2, C_3$ | $G_1, G_2$ | Cont1, Cont2 |
| XYZ | $D_2, D_3$ | $C_2, C_5$ | $G_{11}, G_2$ | Cont1, Cont2 |

Nested

[ Sort on the basis of movies directed by each director.
→ Director with most number of movies

$df_{final} \cdot groupby \, ("Director") \, ["Title"] \cdot nunique() \cdot sort\_values$
(ascending = false)

Step 1    split (c_, '

| Title | Cost0 | Cost1 | Cost2 |
|-------|-------|-------|-------|
| ABC | $C_1$ | $C_2$ | $C_3$ |

Step 2  ) Stack ( col into rows

| ABC | $C_1$ |
|-----|-------|
| ABC | $C_2$ |
| ABC | $C_3$ |

| Title | Director | Cast |
|-------|----------|------|
| ABC | D1 | C1 |
| ABC | D1 | C2 |
| ABC | D2 | C1 |
| ABC | D2 | C2 |
| ABC | D1 | C3 |
| ABC | D2 | C3 |

**Step 3**

a) | Title | Cast | $df_c$

b) | Title | Director | $df_D$

c) | Title | Country | $df_{cnt}$

D) | Title | Listed-in | $Df_{G}$

$df_{CD}$

$df_{CDcnt}$

$Df_{CDCGcnt}$

**Step 4**  Merge → left, inner

## Step5

Merge it back with remaining cols of the original $\underline{DF}$ using "Title"

$$df_{final}$$

→ Most popular actor & director pair.

→ Genre with most number of movies.

⊛ { Date_added } → datetime → Extract month, year, day / day of week

↳ pd.to_datetime( _____ ) →

what insights you can find if you fix this ?

$\rightarrow$ month when most movies were added

year

$\rightarrow$ day of week $\rightarrow$

Duration

Split ('_,') [0]

| Title | Type | Duration | Duration [0] |
|-------|------|----------|--------------|
| ABC | MOVIE | 90 min | 90 |
| XYZ | TV | 2 seasons | 2 |
| MNO | MOVIE | 110 minutes | 110 |

Average runtime of movies & tv shows
Median runtime                 min    seasons

df. groupby (["Type"] (Duration[0]]) . median ( )

so we can say what is the most popular format in each country or genre..do people watch longer movies in some countries?

# Missing values

(*) Mode → most occuring Cat

① Unknown
② Delete
③ Mode

DFA

| Actor | Director | Mode (Director) |
|-------|----------|-----------------|
| $A_1$ | $D_1$ | $D_1$ |
| $A_1$ | $D_1$ | $D_1$ |
| $A_1$ | $D_1$ | $D_1$ |
| $A_1$ | | $D_1$ |
| $A_1$ | $D_3$ | $D_1$ |
| $A_2$ | $D_5$ | $D_1$ |
| $A_2$ | $D_5$ | $D_5$ |
| $A_2$ | $D_5$ | $D_5$ |
| $A_1$ | $D_1$ | $D_5$ |
| $A_1$ | $D_1$ | $D_1$ |
| $A_2$ | $D_4$ | $D_1$ |
| $A_2$ | | $D_5$ |
| | | $D_5$ |

Mode (Director) → $D_1$

By Method → you are right ?

"Most Appropriate" imputation

How to do this

Reference column I will see → accordingly
I will create mode &
then inputs

$\Big\{$ df.groupby ("Actor") (Director) . mode ( )

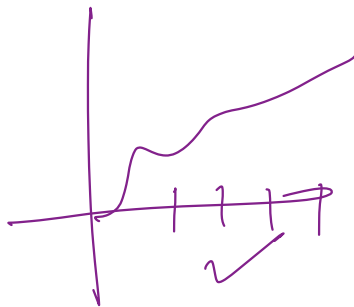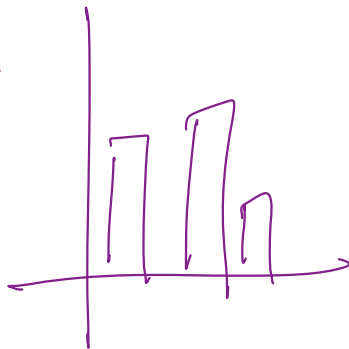| Actor | Mode(Director) |
|-------|----------------|
| A1 | D1 |
| A2 | D5 |

DFB

Soly    Merge DFA & DFB on the basis of Actor

Analyse →



Insights