(Python) , (SQL) , (Pandas , Numpy) , Seabon , Matplotlib.

Riturik.

How to solve Netflix biz case ⟶ Numpy , pandas

Solution oriented

trend for prices of an
ecomerce company

tableau , SQL ⟵

Netflix Biz

→ Explore
→ some insights, import observations
→ Recommendation

Concerns
① Missing value
② Duration
③ Date format
④ Nested data

Pre-processing
↳ Cleaning
↳ Sanitizing

Exploratory data
Analysis
(EDA)

80%

# Nested data

→ processing issues

| Title | Cast | Director |
|-------|------|----------|
| ABC | M, N | $D_1, D_2$ |
| 3 idiots | Amir, Madhwan, Sharma Joshi | $D_7$ |
| XYZ | W, M | $D_1$ |

Q who's the most popular actor on Netflix platform

→ genre, country, Actor . . . .

M

Q who's/that is the most popular Actor-director pair.

df.groupby(["Cast", "Director"])[Title].unique()

.sort value( )

① 

| Title | Cast_0 | Cast_1 | Cast_2 |
|-------|--------|--------|--------|
| ABC | M | N | Null |
| Bidi- | Anil | Madh | Sharma |
| XY2 | W | M | |

Split will not solve my problem

50%

X will not be able to solve

② Stack (unpivot) (transpose)

Convert colums into rows

df

| Title | cost | value | Directes | genve |corny |
|-------|------|-------|----------|-------|------|
| ABC | Cast_0 | M → D1 | | |
| Abs C | Cast_1 | N → D2 | | |
| Bidiots | Cast_0 | Aniv → D1 | | |
| Jidiots | Last_1 | Madh → D2 | | |
| Bidiots | Cast_2 | Sharman | | |
| XY2 | | | | |

delete    Rename to cast

df . groupby ( "Cast" ) [ Title ] . count ( )

n unique( )

| Cast | Direct |
|------|--------|
| A B C | M | D1 |
| ABC | N | D1 |
| ABC | M | D2 |
| AB C | N | D2 |

| | | | M | 1 |
|---|---|---|---|---|
| M | ~~2~~ | | N | 1 |
| N | ~~2~~ | | | |

Unnest¹ ①   Split & Stack  →  Cast  ⟹ df1

Unest² ②        → Director ⟹ df2

Unest³ ③  _____ → Country ⟹ df3

④  _____ → genre ⟹ df4

**DF₁**

| Title | Cast |
|-------|------|
| ABC | M |
| A>C | N |

**DF₂**

| Title | Di |
|-------|-----|
| A>C | D₁ |
| ABC | D₂ |

merge → **DFₐ**

**DF3**

| | Country |
|---|---------|
| | |

**DF4**

| | Genre |
|---|-------|
| | |

merge → **DF_b**

merge → **DF final**

Remaing col.

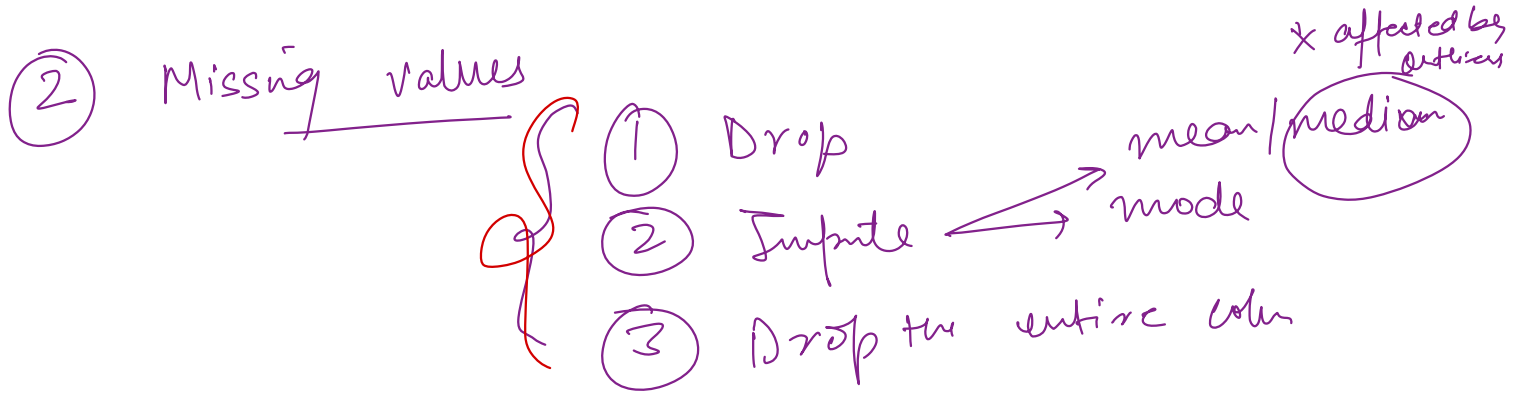| Title | Cast | Direct | Countries | genre | ✓ | ✓✓ |
|-------|------|--------|-----------|-------|---|----|
| | | | | | | |

merge → Original data (to bring the remaing cols)
Title

Summary:
step 1 : split into dataframes with titles into centre and other features as the specific dataframes
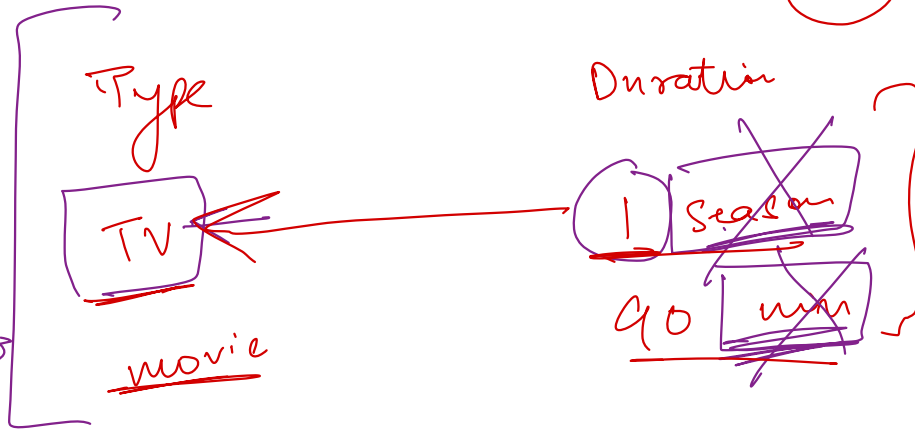step 2: join these different dataframes using title

② Missing values

{
  ① Drop
  ② Impute $\longrightarrow$ mean/median    ✗ affected by outliers
                          mode
  ③ Drop the entire colm
}

Medical

IQR

1.5 ✗ IQR

3 ✗ IQR

⑧ Duration

Split [0]

Type

Duration

Duration_new

TV

1 | Season

1

movie

90 | min

90

Tell me the average emting of the movies in which SRK is present

1 get the average emting of all the actor & sort it in descending

df. groupby ("Cast") [Duration] . average
                  v                        ↑
                                        mean ( ) . Sort van
            "Type"

| | | |
|---|---|---|
| Anpu Khur | movie | 100 |
| Anpu Khur | TV | 2 |

④ Date format

Pd. to_datetime (Date_added)

↳ Extract year
         month
         day of the week  }

① : which is the most busy month for
              netflix

↳ Increase server capacity
         before they comes

⌐————————————————————————————————————⌐
Seabon
matplot  ↳ Charts to visualize

Pick up
the Insight & give recommendations.
imp

$+$ ipynb $\longrightarrow$ pdf $\longrightarrow$ upload to
platform

Pls do come for review session

| Cast | |
|------|--|
| SRK ✓ | ✓ |
| salma | ✓ |
| SRK | ✓ |

Numerical → Mean/Mean
Categorical → mode
SRK

$x = df[cast].mode()$

$df[cast].fillna(x)$ 〕 Simple imputation

Bonus

→ Multiple imputation

India

Module → ajay.shenoy@scale.com
Lead        └ Batch name