

Netflix

Programming → Python

Numpy
Pandas

Seaborn

Practice

medium to tough

Follow the steps :

Insights & recommendation is true.

Assumption working good for netflix

Challenges

- ① Missing (NaN ...)
- ② Duration col.
- ③ ↪, ↪, ↪, ↪

Model inputation
median / mean
9 seasons → Avg runtime for actors
challenge & why?

A → Rajpal
Yadav

| Actors | movies |
|---------|--------|
| A, B, C | M1 |
| A, B | M2 |
| B | M3 |
| X, Y, A | M4 |
| D | M5 |

Get me the no. of movies ever an actor has worked in.

Group by

A

None

Nested

Unnested data

(M) Data type issues

Id. • Date time (1)

- dt • years
- dt • month
- dt • week

Nested data

[Cast , Director , Country , listed-in]

↳ if we do not fin it.

↙ ↗ ↗
 ↗ ↗ ↗
 ↗ ↗ ↗

Most popular actor in our DA.
 Best Actor - Director Combo.
 Famous actor in a country.

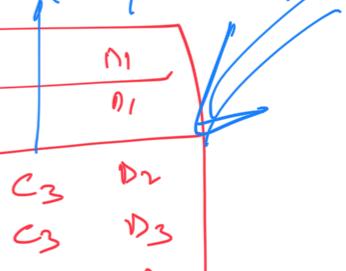
| Title | Cast | Director | Country |
|-------|---------------------------------|---------------------------------|---------|
| ABC | C ₁ , C ₂ | D ₁ | |
| X Y Z | C ₃ , C ₄ | D ₂ , D ₃ | |
| M N O | C ₅ , G ₁ | D ₄ | |

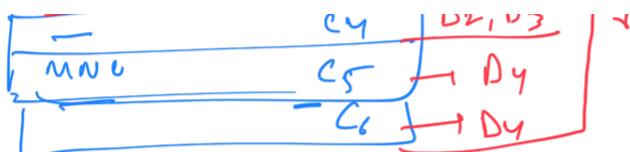
Vineet → How many rows are total.

| | | |
|-------|----------------|---------------------------------|
| ABC | C ₁ | D ₁ |
| ABC | C ₂ | D ₁ |
| X Y Z | C ₃ | D ₂ , D ₃ |

| | | | |
|---|-------|----------------|----------------|
| ① | ABC | | D ₁ |
| ② | ABC | | D ₁ |
| ③ | X Y Z | C ₃ | D ₂ |
| ④ | X Y Z | C ₃ | D ₃ |

Cast
↓
Direct and good

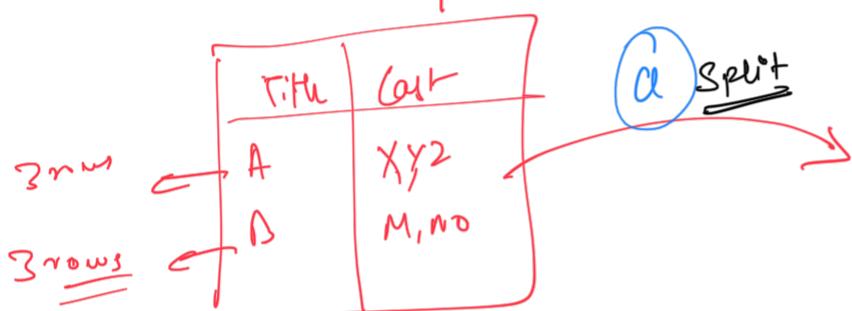




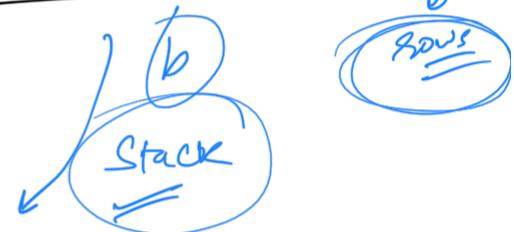
| | | | |
|---|-------|-------|-----|
| ⑤ | X Y 1 | C 4 | U 2 |
| ⑥ | X Y 2 | C 4 . | D 3 |
| ⑦ | M N O | C 6 | D 4 |
| ⑧ | M N I | C 6 | D 4 |

Steps to solve uncertainty

① Treat every column differently



| Title | Col-0 | Col-1 | Col-2 |
|-------|-------|-------|-------|
| A | X | Y | Z |
| B | M | N | O |



T1

| Title | col | val | P1 | P2 | P3 |
|-------|-------|-----|----|----|----|
| A | Col-0 | X | | | |
| | Col-1 | Y | | | |
| | Col-2 | Z | | | |
| B | Col-0 | M | | | |
| | Col-1 | N | | | |
| | Col-2 | O | | | |
| C | Col-0 | -0 | | | |
| | Col-1 | -1 | | | |
| | Col-2 | -2 | | | |

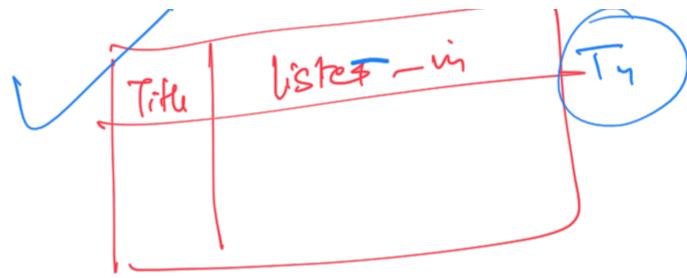
② Repeat this step for all the col

| Title | Direct |
|-------|--------|
| A | D1 |
| A | D2 |

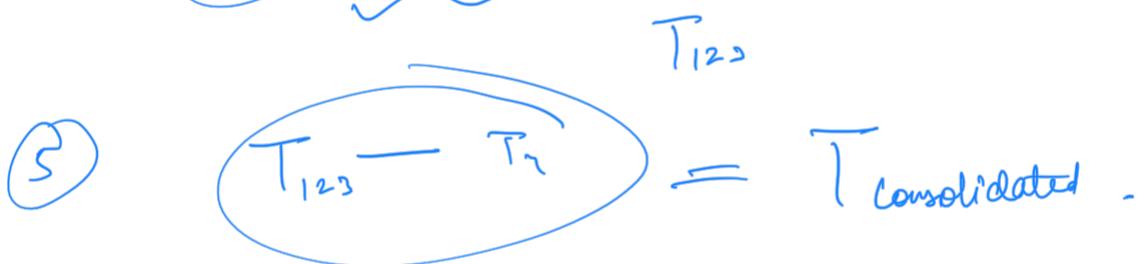
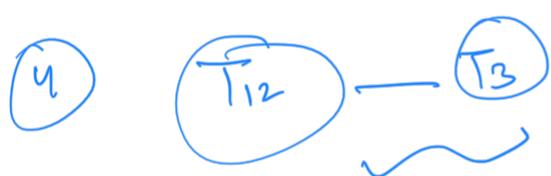
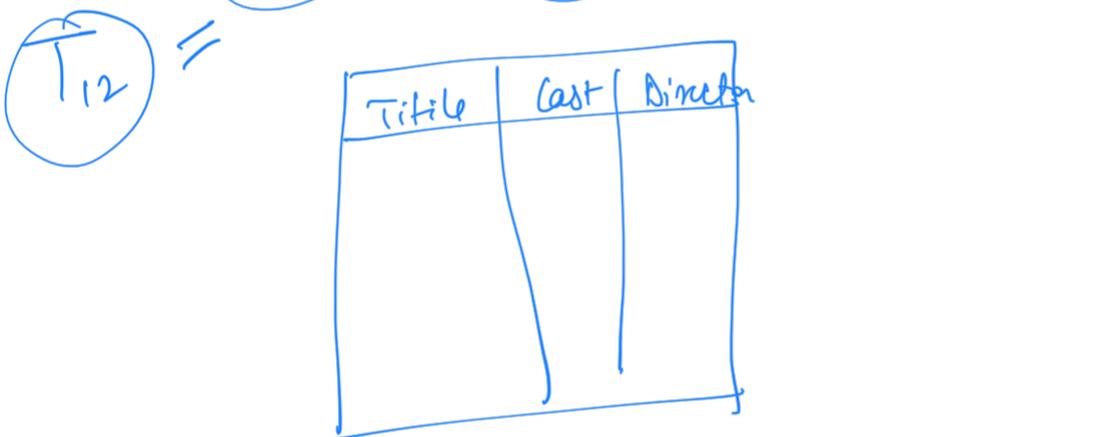
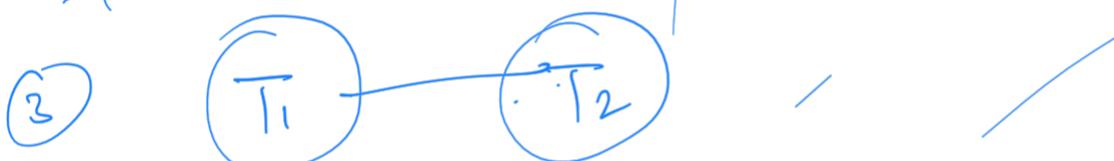
| Title | Country |
|-------|---------|
| A | USA |
| A | India |

T2 ✓

2



Separate table for each col



Tconsolidated

{ Is T_c unique at each row level.

↳ ... i.e. the most popular actor in

↳ ~~get ...~~

~~df.groupby([["Country", "Cast"], ["Title"]]).count()~~

~~↳ unique~~

~~↳ listed in multiple Directors multiple~~

~~[df.groupby([["Country", "Cast"], ["listed in", "Directors"]]).size()]~~

~~↳ Cast~~

- ④ Join with original dataset on 'Title'
& fetch the remaining column.

~~Easy~~

Questions

Missing value imputation

- ① Mean / median / mode .

Right .

Romus industry. (Optional)

Advanced way

of input

| Country | Title | Genre | Actor | Mode |
|---------|-------|-------|-------|------|
| US | A | Comd. | BST | X |
| US | B | | | X |
| US | C | | | X |
| US | D | | | X |

Tom Cruise

usually

*multiple
input*

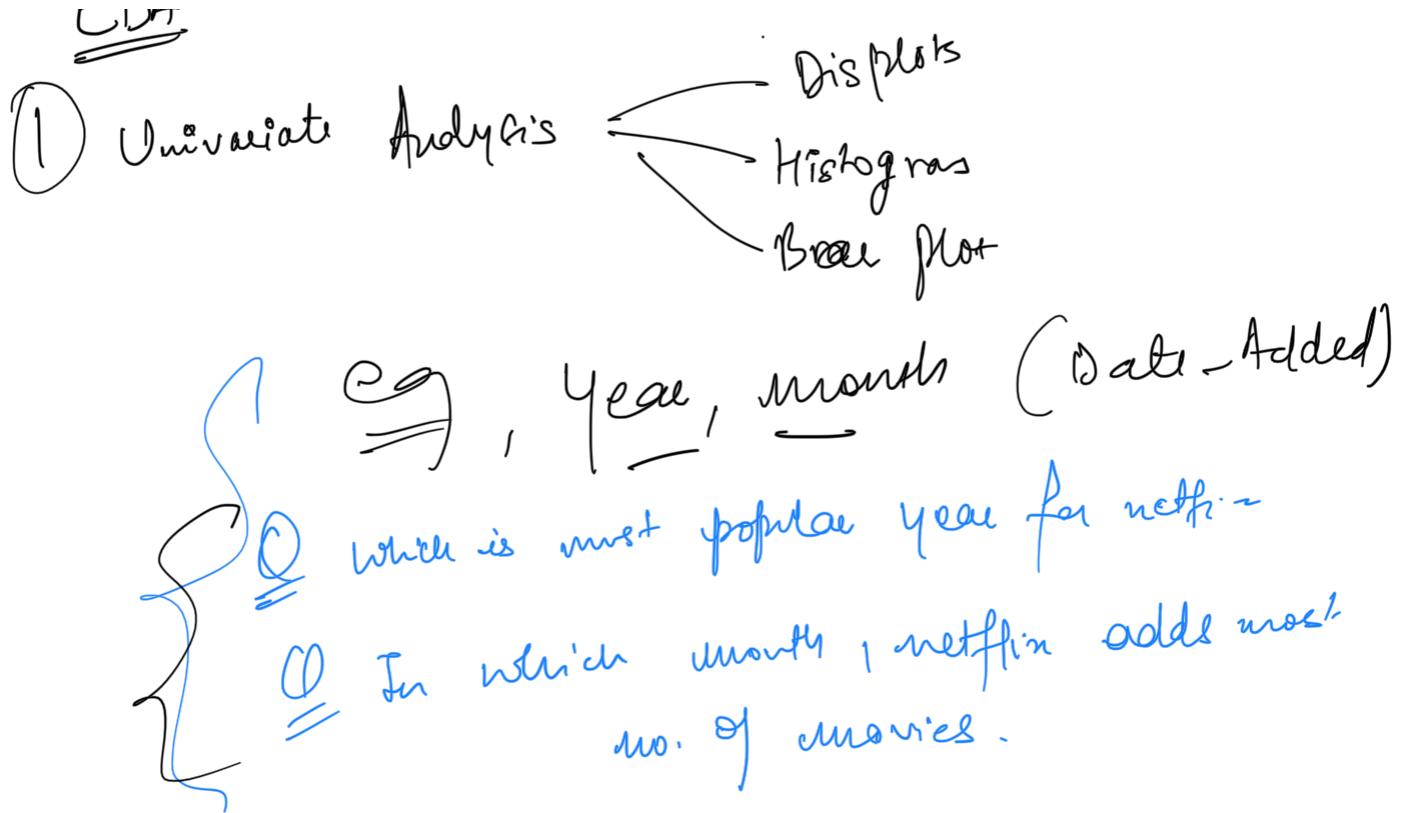
(1)

Director

| Country | Model(Actor) |
|---------|--------------|
| India | SRK |
| USA | Tom |
| China | Bruce Lee |

Processing of the data is complete

END



② Bivariate analysis → X

→ ←

Insights & Recommendations