

How to perform EDA process

Why we perform EDA? → Solve a business problem

Clear i) Business understanding: Correlate the insights from the data with your business understanding

Credit card : Data
↓
[$\leftarrow [age \geq 18]$]

| age |
|-----|
| |
| |
| |

$\frac{20,000}{\downarrow}$
[$age < 18$]

| age | usage |
|------|--------------------------|
| < 18 | <input type="checkbox"/> |

→ client

ii) Data understanding :

| x | y | z |
|---|---|---|
| ✓ | | |

+ Data dictionary

✓ x →
 ✓ y →
 ✓ z →

very imp

oil & gas

| boilerid | operatorname | age | min temp | max temp | ... |
|----------|--------------|-----|----------|----------|-----|
| — | | | | | |

X age of the operator
 age of the boilers →

- Data cleaning → Check for duplicated rows

→ Data type correction

→ Category check (check the list of unique categories)

Ecomm

city

- New Delhi
- Delhi
- Deli
- South Delhi
- NCR

→ New Delhi

Python → Pandas

- Integer → int
- decimal → float
- String/text/alphanumeric → object
- Datetime → Datetime

complete duplicated row : keep one & drop other

Duplicates in a certain column

Investigate

There will always be a reason behind it

according take next step

| c ₁ | c ₂ | c ₃ | c ₄ | c ₅ |
|----------------|----------------|----------------|----------------|----------------|
| a ₁ | b ₁ | c ₁ | d ₁ | e ₁ |
| a ₁ | b ₁ | c ₁ | d ₁ | e ₁ |

duplicated rows

| c.id | o.id | p.id | v.id |
|----------------|------------------|------------------|------------------|
| c ₁ | c ₂ | c ₃ | c ₄ |
| a ₁ | b ₁ ✓ | c ₁ ✓ | d ₁ ✓ |
| a ₁ | b ₂ ✓ | c ₂ ✓ | d ₂ ✓ |

Here the data is duplicated because there are two orders placed by one customer

| customer amt | order date |
|-----------------|------------|
| 1.0 | 01/01/2011 |
| 2.0 | 01/02/2011 |
| 3.0 | |
| 7.0 | |
| 9.0 | |

→ object → Pandas Datetime format

→ Since this is amt (wrong data type)

$$\frac{50 + 29 + 28 + 21}{4}$$

$$\frac{29 + 25 + 28 + 21}{4} =$$

| Age | Gender | Sal |
|-----|--------|-----|
| 29 | M | 50k |
| 25 | F | NA |
| 28 | M | 29k |
| 28 | M | 28k |
| 21 | F | 21k |

$\frac{50 + 29 + 28 + 21}{4}$
 $\frac{29 + 25 + 28 + 21}{4} =$
 $\text{avg}(\text{Sal}) = 32\text{k}$
 $\text{avg}(\text{age}) = 25.75$

25.75 → NA
 32 → NA
 Garbage value
 NOISE

Analysis ✓
 actual insights?
 Wrong
 Biased
 Not True

• objective → EDA → we don't perform missing value imputation

• objective → Machine Learning modelling → various ML-Algo (They can't operate on a data with missing values)

↙
 Necessity to perform
 missing value imputation

df.duplicated() →

| | c ₁ | c ₂ | c ₃ | c ₄ |
|---|----------------|----------------|----------------|----------------|
| 0 | a ₁ | b ₁ | c ₁ | d ₁ |
| 1 | a ₁ | b ₁ | c ₁ | d ₁ |
| 2 | a ₂ | b ₂ | c ₂ | d ₂ |
| 3 | a ₁ | b ₁ | c ₁ | d ₁ |

True → 1

false → 0

0 false →
 ✓ 1 True →
 2 false
 ✓ 3 True →

0
 1
 0
 1

→ ②
2
 Count of
 True

df.drop_duplicates()

Subset :

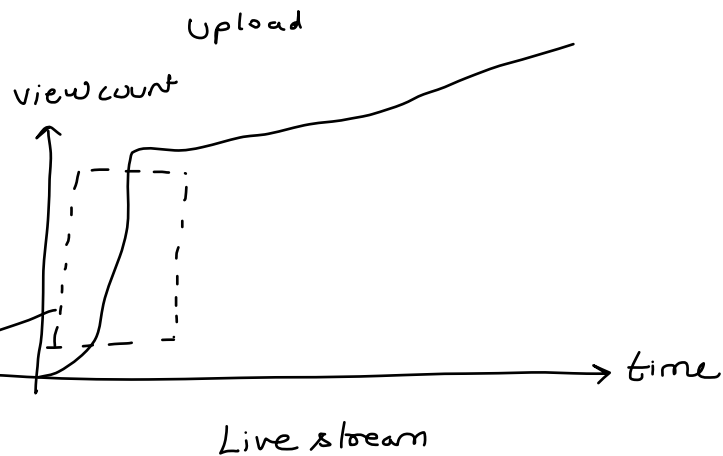
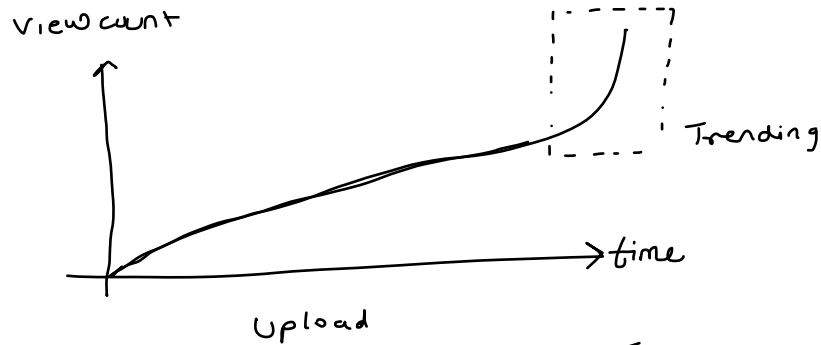
| | | | | | |
|-------|----------|-------|---------|---------|---------|
| | | ↓ | ↓ | ↓ | ↓ |
| | | c_1 | c_2 | c_3 | c_4 |
| first | ← 2478 ✓ | a_1 | ✓ b_1 | ✓ c_1 | ✓ d_1 |
| last | ← 2499 ✓ | a_1 | ✓ b_1 | ✓ c_1 | ✓ d_1 |
| | 2530 | a_2 | b_2 | c_2 | d_2 |

| | | | |
|-------|-------|-------|-------|
| | ↓ | | |
| c_1 | c_2 | c_3 | c_4 |
| a_1 | b_1 | c_1 | d_1 |
| a_1 | b_1 | c_2 | d_1 |

complete duplicated rows

We are looking at the data by considering all the columns?

Subset → column name on which you want to check for duplicated entries
by default it will check all the columns
(complete duplicated rows)



This significant increase in the no. of viewership makes you vide become trending