


DAN-II

| Start at 9:03pm

⇒ Agenda →

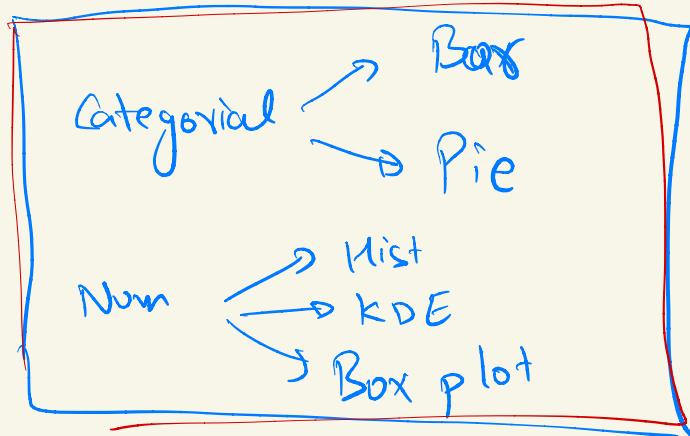
- Revision
- Bivariate Analysis
- Subplots

⇒ Revision →

- Univariate → 1 variable
- Bivariate →
- Multivariate →

⇒ Univariate Data viz

↓
1 variable



- ⇒ Bivariate → 2 variable
- Num - Num → Cat - Cat → Cat - Num
- Num - Num / Continuous - Continuous

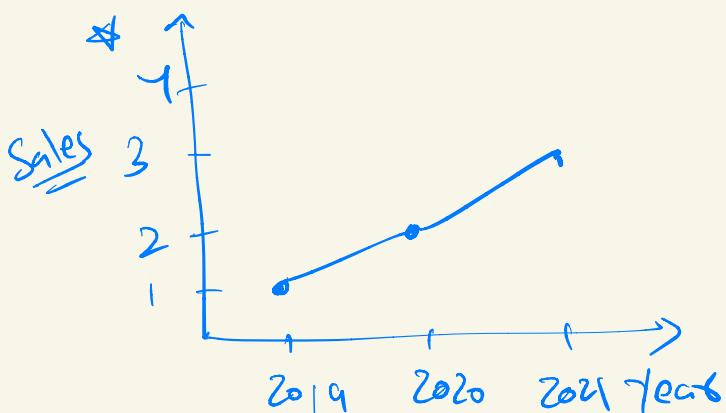
⊕ What question can we ask?

- Sales, year
- ↓ ↓
- i) How does sales vary over years
 - ii) How are features associated
↳ Correlation →
→ Year ↑ ~ ↑ ↓
 - iii)

(i) → Sales vs Year →

⊕ what chart should we use?

↳ line



① Dataviz →

Date / Time ⇒ go to
x axis

Num-Num →
 → line
 → Scatter

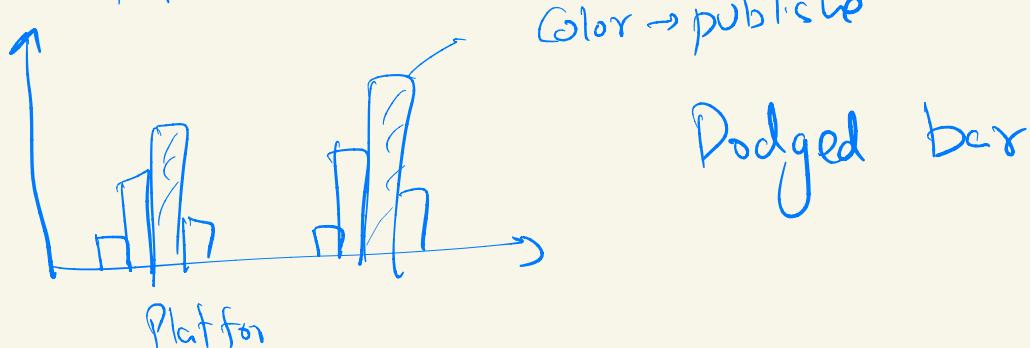
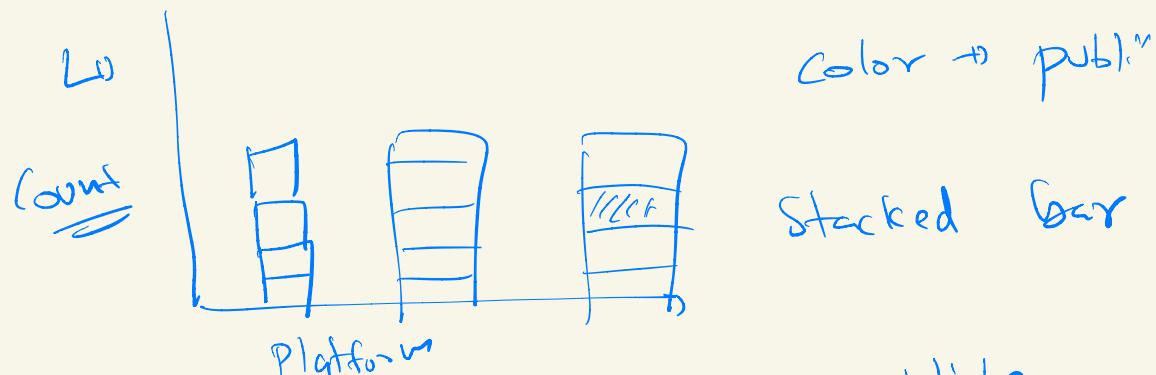
⇒ Cont - Cont →
→ Publisher, Platform

Q What types of questions can be asked?

→ preferred platform for publisher.

→ distribution of publisher for top 3 platforms

⇒ Distribution (Count) of one wrt other category.



⇒ Sunburst →



⇒ Stacked bar +

Stacked

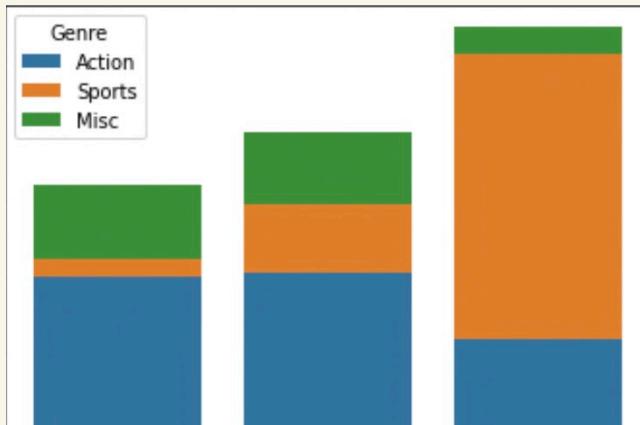
→ it is not readable

↳ If you don't label
for

Date

dodge

→ for comparison



→ Helps only with total

⇒ Total is more important

2 Cat - Cat

→ Dodge Bar chart

→ Stacked Bar

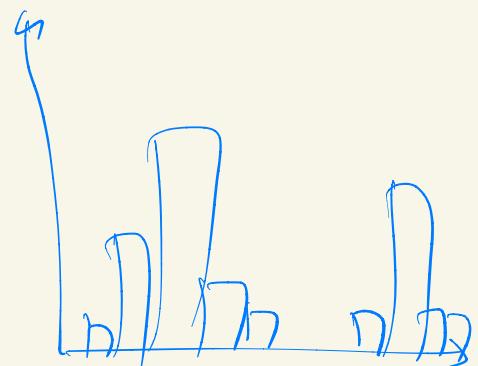
\Rightarrow Cat - Num \rightarrow

* Publisher vs Sales \rightarrow

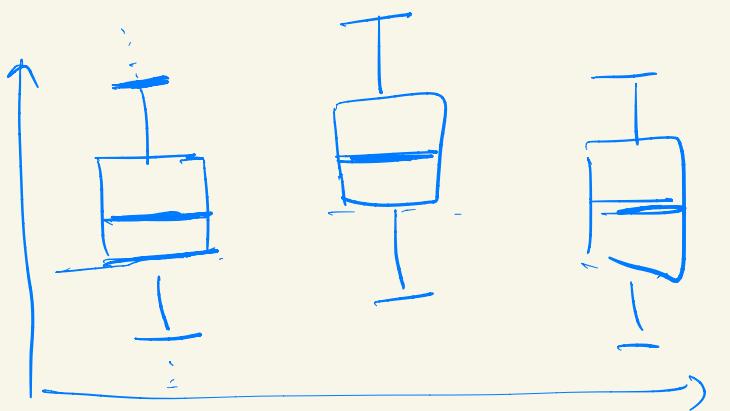
$\oplus \rightarrow$ what question can be asked??

- 1) What is avg/sum sales for every publisher
- 2) Sales distribution for top 3 publisher.

\Rightarrow Distribution of a numerical value.



\Rightarrow Box whiskers-



\Rightarrow Can we use bar plot as well for
Num-Cat variables.

\Rightarrow Sum of Sales | avg sales } median Sales



X ✓ -

Doubts

→ Salary of DS →

→ 30, 35, 25, 28, 40, 42, 38, 39, 300

→ avg → $\frac{31.5}{9} \rightarrow 72$

median → $\frac{35+38}{2} = \underline{\underline{36.5}} \rightarrow 38$

25, 28, 30, 35, 38, 39, 40, 42, 300

percentile → x value lower or equal to the given x

$$25 \rightarrow \frac{1}{9} \times 100$$

$$28 \rightarrow \frac{2}{9} \times 100 \quad \dots \dots \quad 300 \rightarrow 100\%$$

$$30 \rightarrow \frac{3}{9} \times 100$$

=D • Distribution) viz →

• lower value

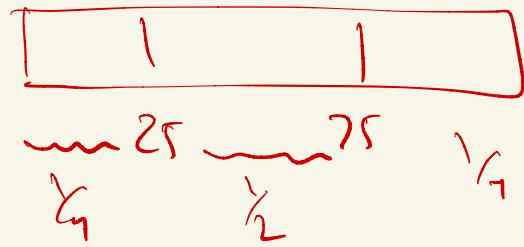
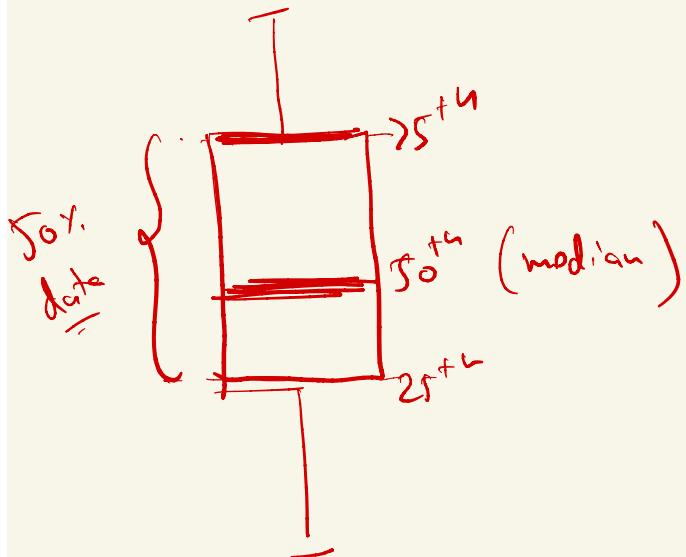
• median

• outlier

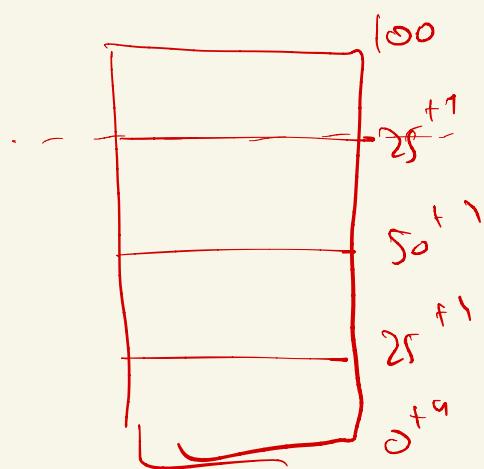
• upper value

• 25th percentile

• 25 percentil



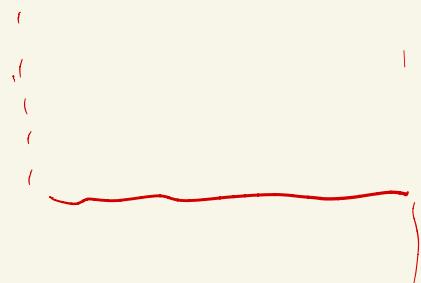
outliers $25^{th} + 1.5 \text{ IQR}$



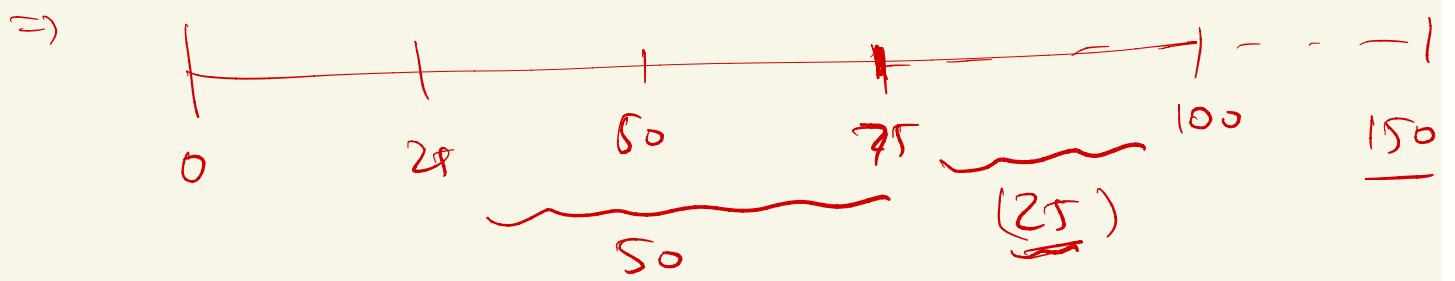
\rightarrow

$$\Rightarrow \boxed{25^{th} + 1.5 (75^{th} - 25^{th})}$$

anything above this
is an outlier



outlier



$$75 + \underbrace{(50 \times 1.5)}_{75}$$