

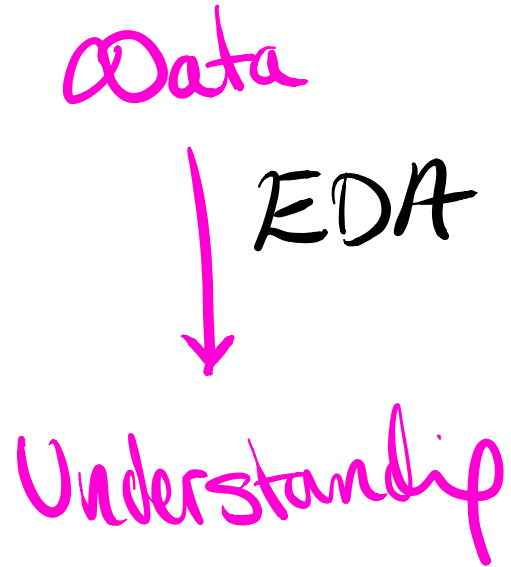
# Exploratory Data Analysis

---

# What is EDA?

---

- ⊗ Investigation
- ⊗ Description
- ⊗ Explanation



# Challenges to Successful EDA?

- ⊗ Data Quality
- ⊗ Process / Method  
Z.B. Eg.  $t$ -test
- ⊗ Organizational challenges
- ⊗ Business / Strategic Challenges.

# Data Quality

---

- ⊗ Missp → NULL
- ⊗ Unknown Values
- ⊗ Representative / Non R. sample.
- ⊗ Inconsistent Values (Data entry)

# "Explanation" Issues

⊗ Cause → Effect Analysis

→  $E \rightarrow C; C \rightarrow E$

$H \rightarrow E; H \rightarrow C$

$E \text{ assoc. } C$

Distinguishing?



⊗

⊗

→ ⊗ Hypothesis testing (or Bad Data?)

# Organisational Challenges

- ① Bureaucracy → Access, Speed  
→ Change Process  
Ex. Data Collection
- ② Ownership → Gate Keeping Hoarding
- ③ "Politics" → status, Hierarchy, Duplication, unclear Responsibilities
- ④ Governance → help  
→ Internal Ethics

NB. ① HIPPO - "Senior" Opinions [Refute]

# Business Challenges

---

↳ ⊗ define question

⊗ ROI

⊗ Time Available

⊗ Staff & "Human" Resource

⊗ Strategic Value?

# EDA Workflow

⇒ assuming most challenges solved  
universal

⊗ Data Structure

⊗ Distributions } All Cols

⊗ Correlations

⊗ Filters + Distributions  
+ Correlation } i.e. "important"  
Ranges

⊗ Groups

⇒ i.e. Factor Analysis



# Cont.

---

## Domain-Specific

⊗ Formula (Eg.  $BMI = \frac{w^2}{h}$ )

⊗ Predictive Models

$$\text{Eg. } y_{\text{Exper.}} = x_{\text{BMI}} + \dots$$

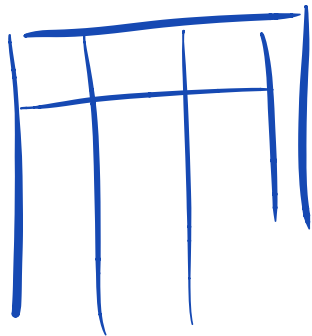
⊗ Specialized Visuals  
Eg. Maps.

Aim: Understand the Data,

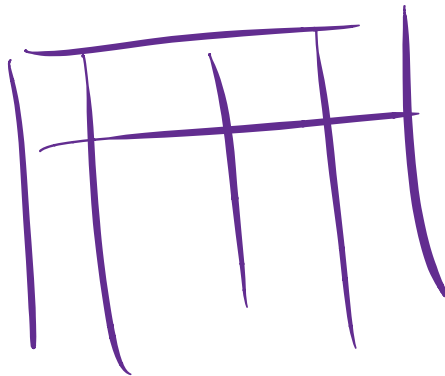
---

# Example: Education

---



School 7  
Year



... 2 ...

5

---

Data  $\rightarrow$  ① Structure

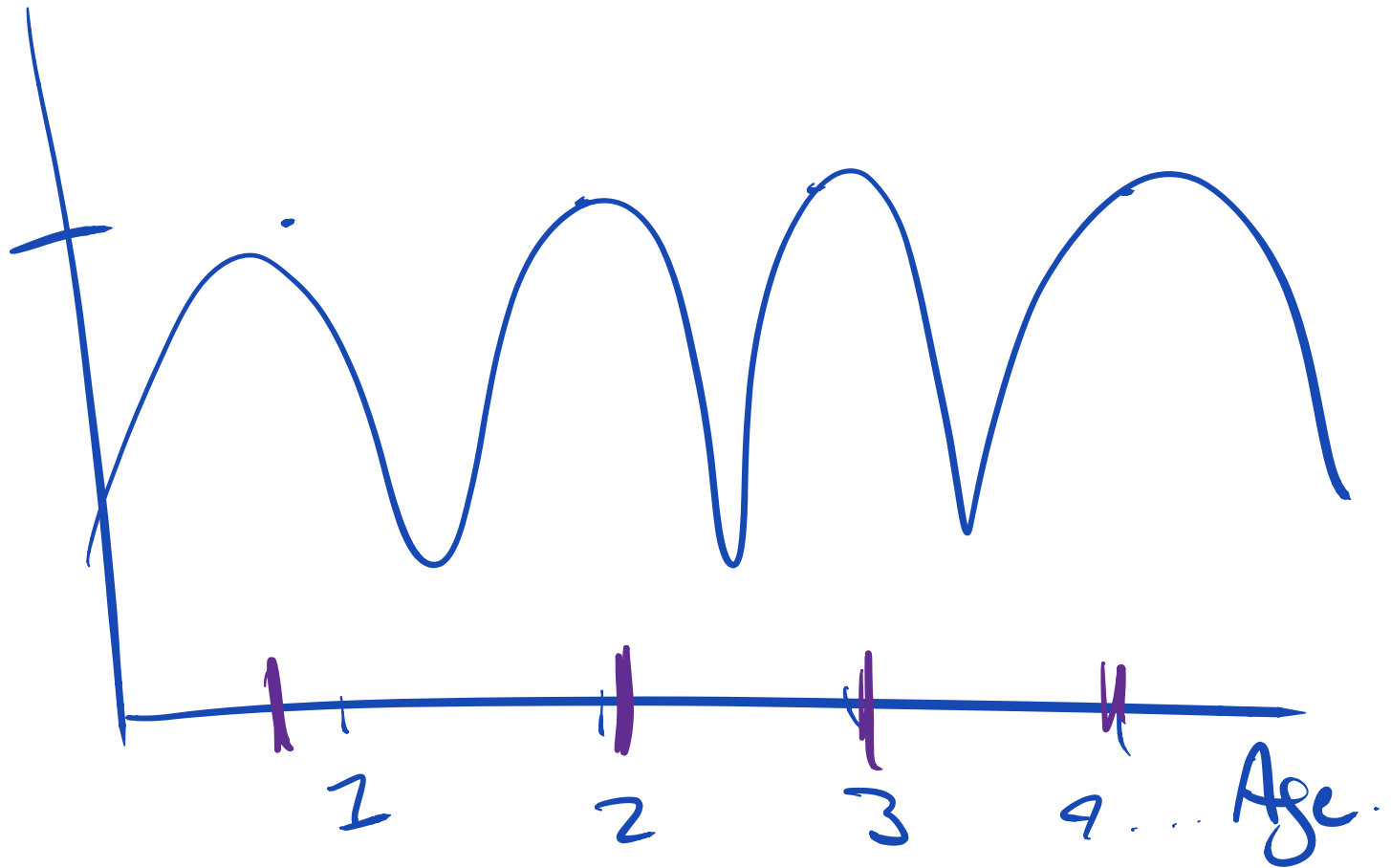
Student ID

Year

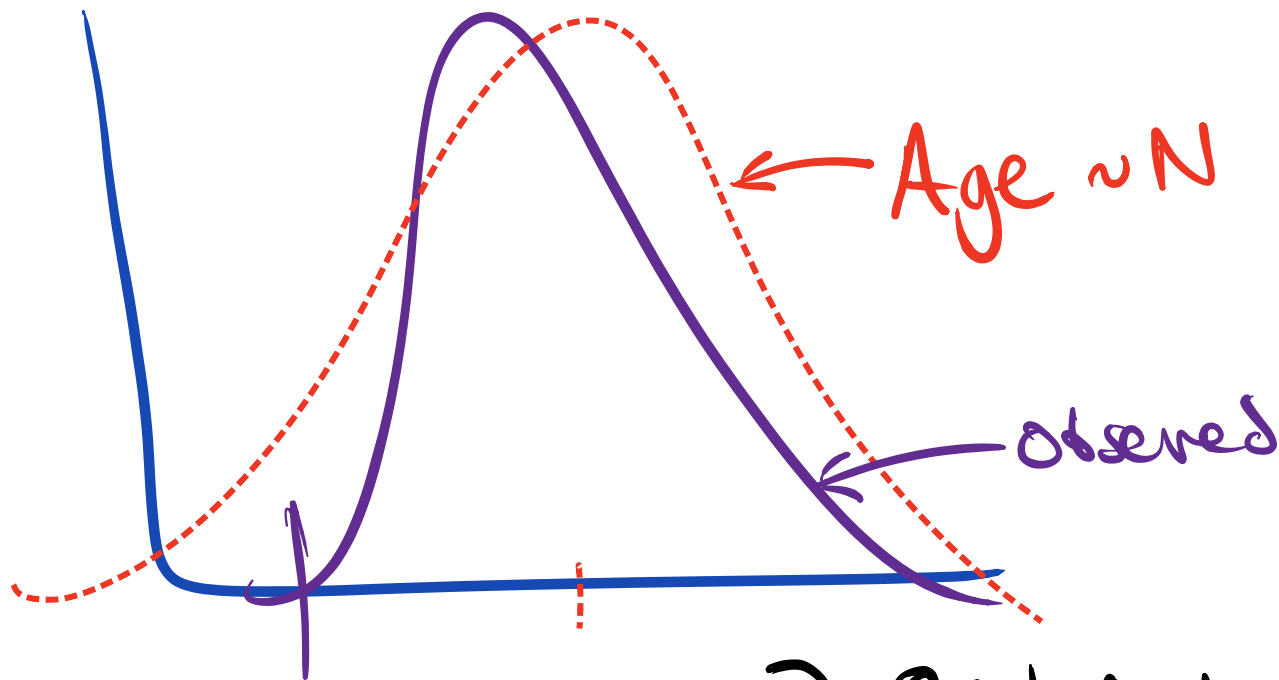
Mean SAT

## ② Distribution

---

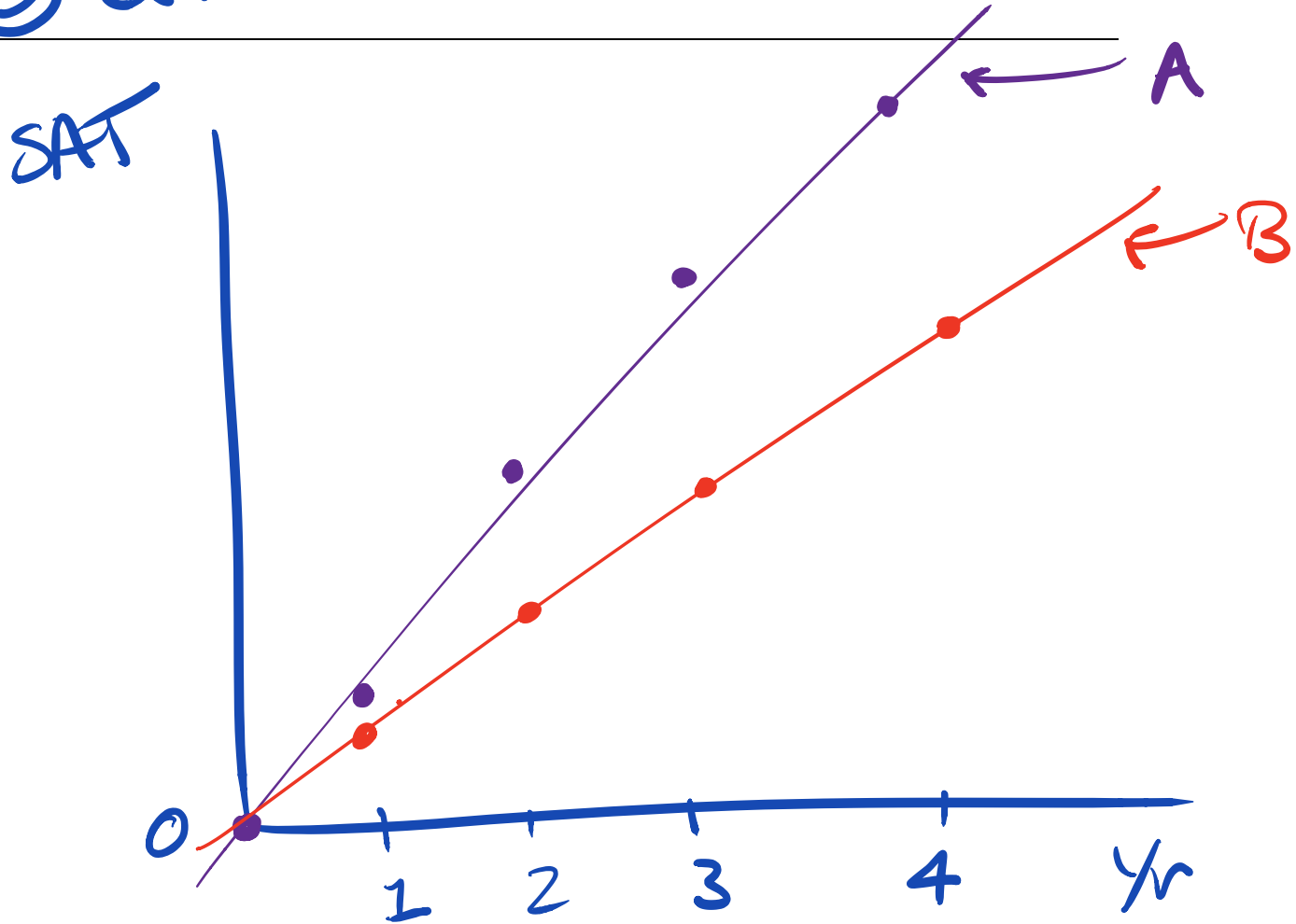


## ② Dist.



... } Problems happen  
if you skip

### ③ Correlahan



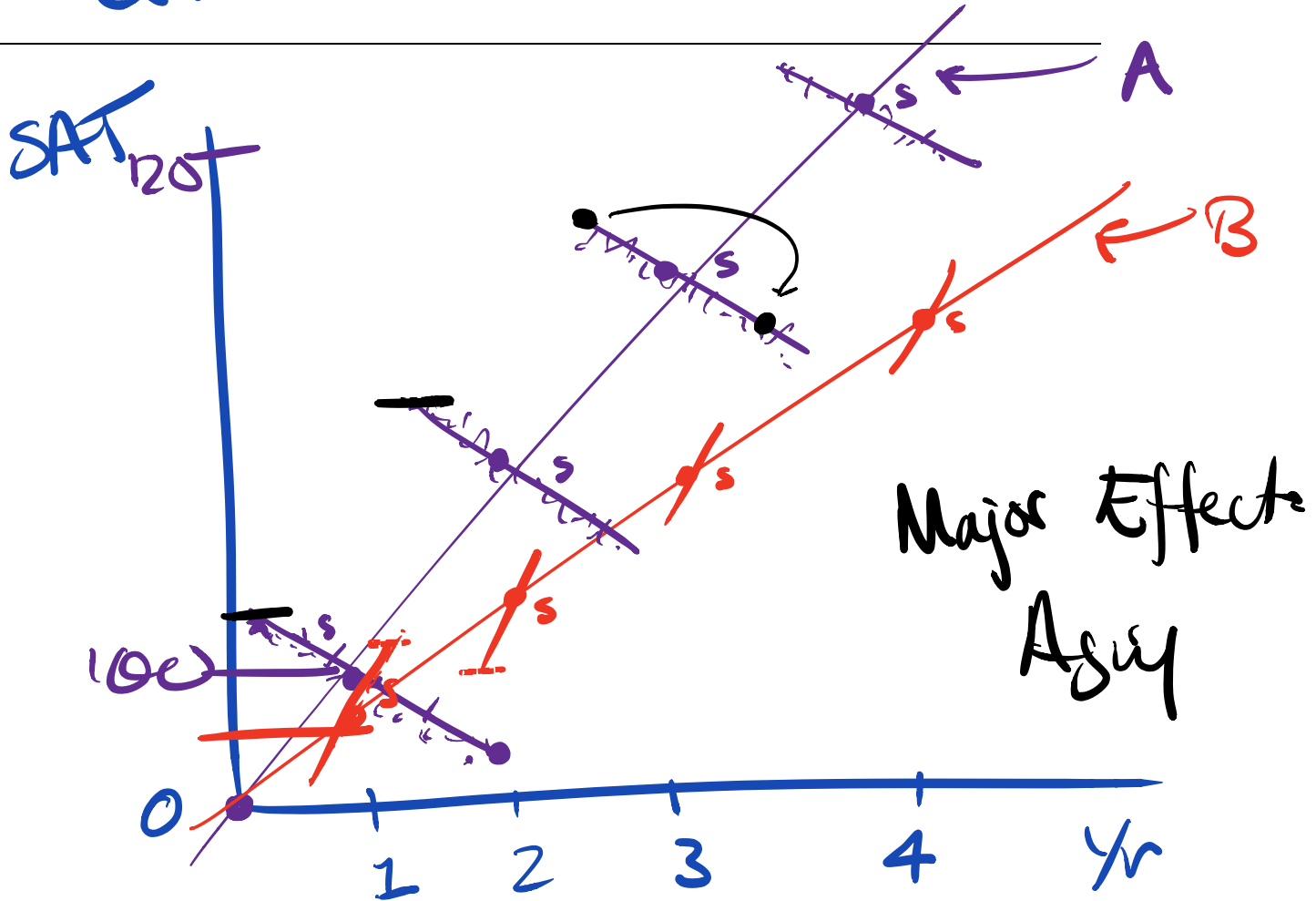
... Explain?

---

Is School A >  
School B ?



# Correlation



# Problem Causal Analysis

---

Age  $\xrightarrow{a}$  Score.

School  
Score  $\xrightarrow{s}$

Penalty  $\xrightarrow{p}$

;

$a \gg s$

$(i + r + p) \gg s$

# Simpsons Paradox

~~infer~~

though

$A > B$

every subgroup  $A$

$<$   
every subgroup  $B$ .

---

Q2A













