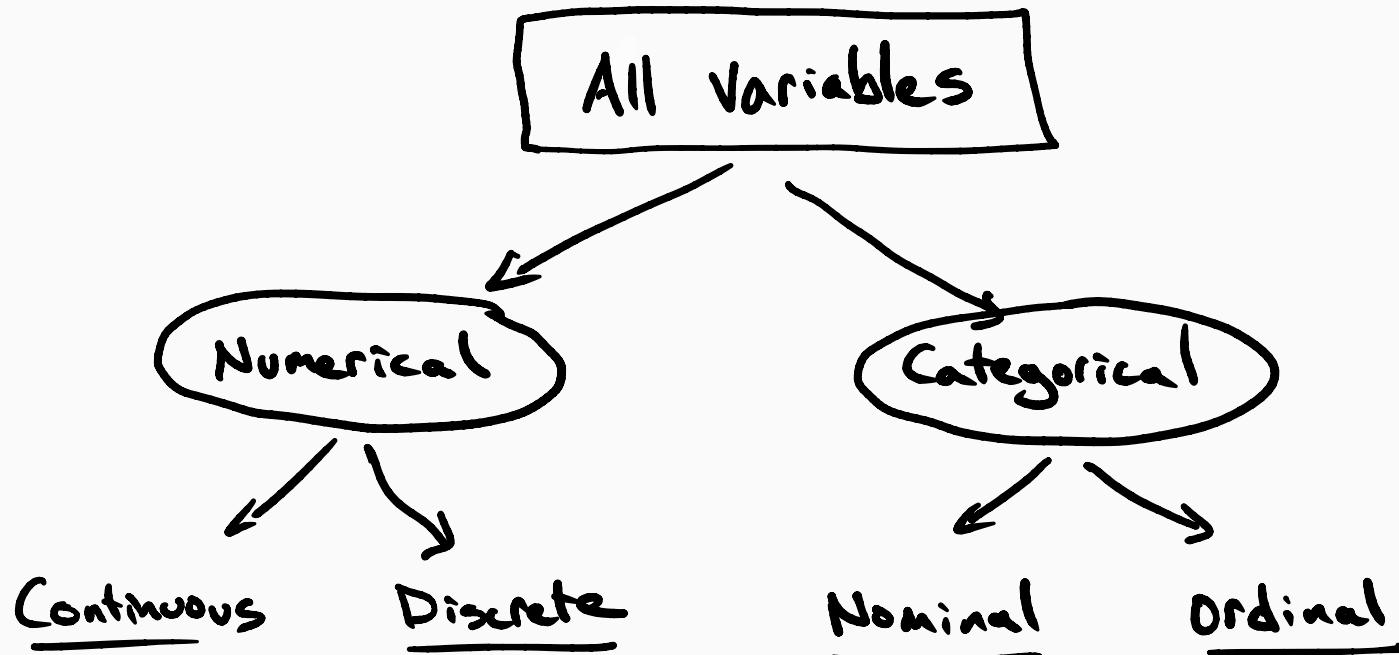


Taxonomy of Inference

Recall: Taxonomy of Data



Taxonomy of Inference

The method of analysis should be driven by:

1. Research Question
2. Data

Let

X: 2-level categorical variable

Y: >2-level categorical variable

W: Numerical (continuous) variable

Categorical Data

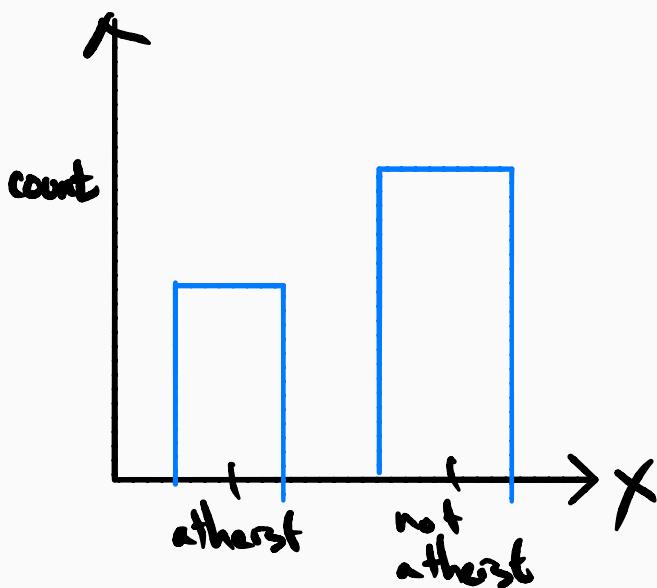
visualize: bar plots

parameters
of
interest : P

breakfast
"cold cereal"
"cold cereal"
"pancakes"
.

X : One 2-level categorical variable

E.g. Ask a person if they are "atheist" or "not atheist".
Visualize: simple bar plot



Parameter
of
interest) P : proportion of
"successes"

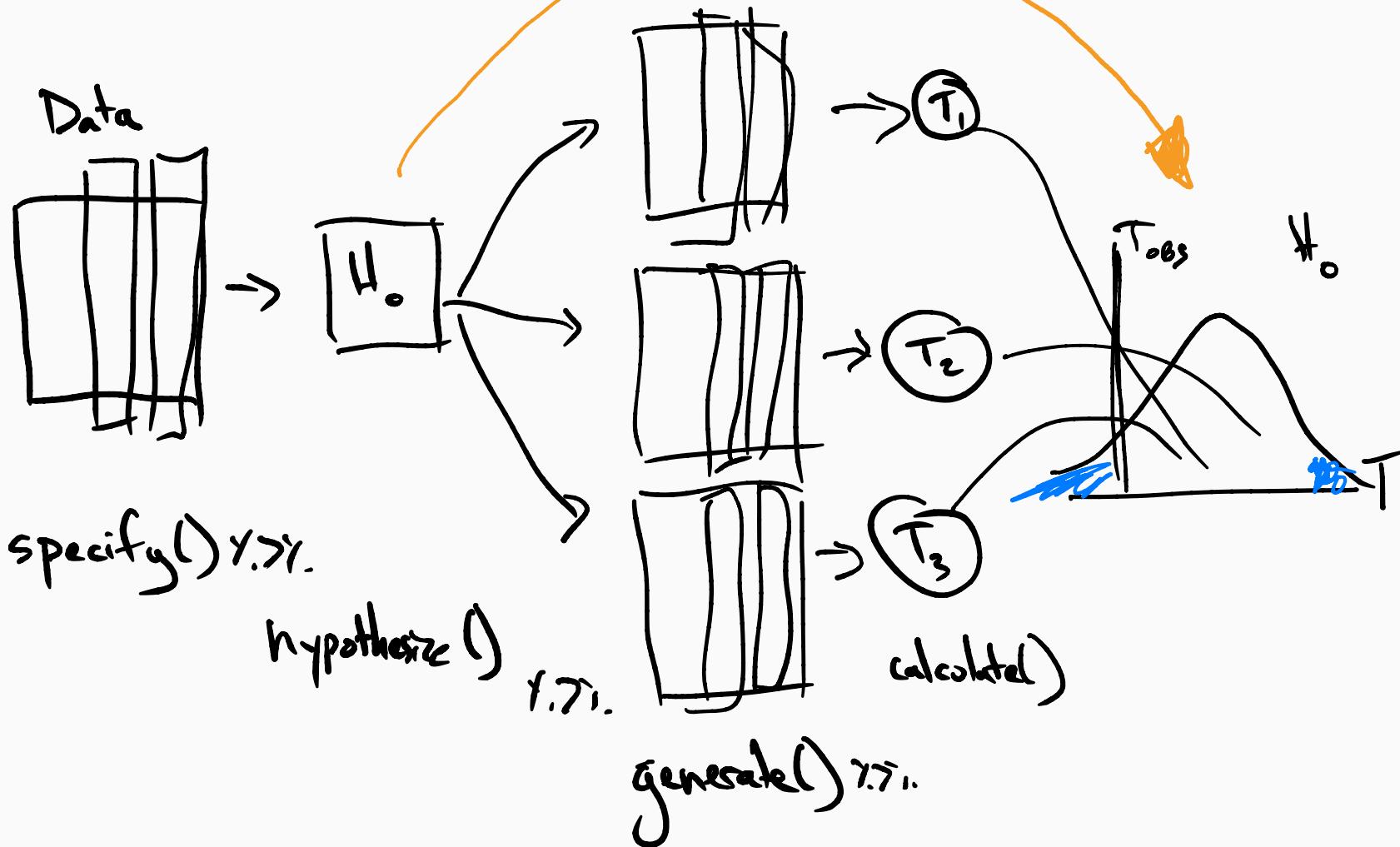
Sampling
Dist. of \hat{P}) • Exact: Binomial
• Simulation

$$\frac{\hat{P} - P_0}{\text{SE}} \sim \text{Normal}(0,1)$$

if $np \geq 10$
 $n(1-p) \geq 10$

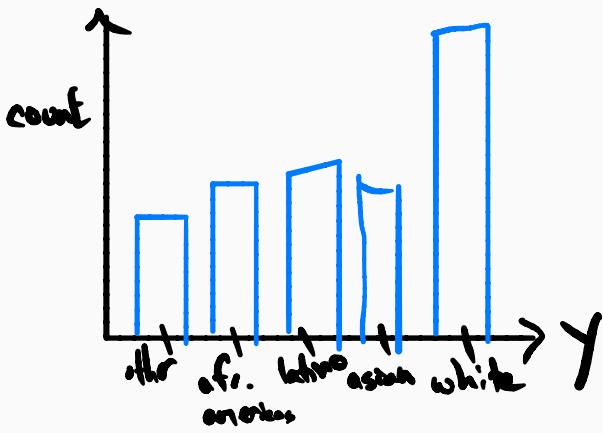
Recall: There is only one test

Approximation
shortcut



Y : One > 2-level categorical variable

E.g. Ethnicity of Reed's First year class.
visualize: bar plot



parameters of interest : P_1, P_2, \dots, P_k

Inferential Technique

Goodness of Fit test.

→ calculate χ^2 stat

$$H_0: P_1 = P_1^0, P_2 = P_2^0, \dots, P_k = P_k^0$$

Get null dist. of χ^2 using

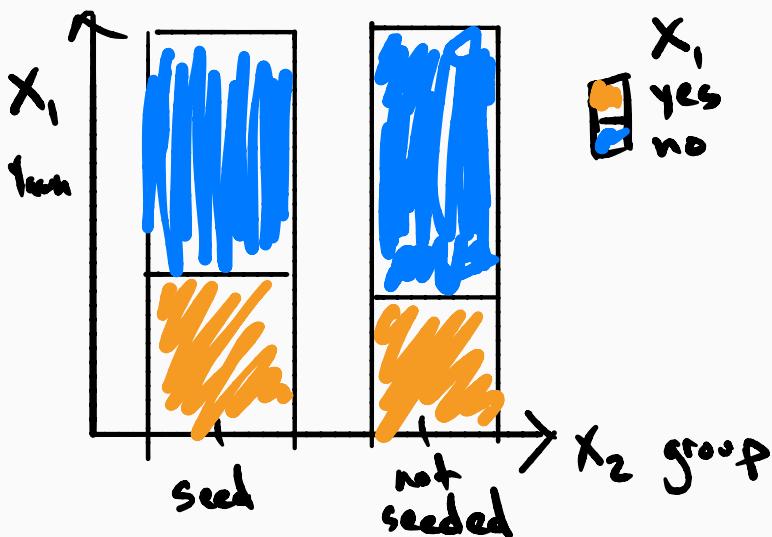
- Simulation

- χ^2 dist. $df = k - 1$

$X_1 \sim X_2$: Two 2-level categorical variables

Eg: Yawning and being in the seeded groups.

Visualize: barplot



X_1
yes
no

parameters
of
interest

$P_1 - P_2$ ← diff in props
of people in
2 group that
yawned.

Null dist
 $\hat{P}_1 - \hat{P}_2$

$$\frac{\text{• permute} \quad \text{• } \hat{P}_1 - \hat{P}_2 - (P_1^o - P_2^o)}{\text{SE}} \sim N(0, 1)$$

if - $np \geq 10$ - ind.
- $n(1-p) \geq 10$ obs.

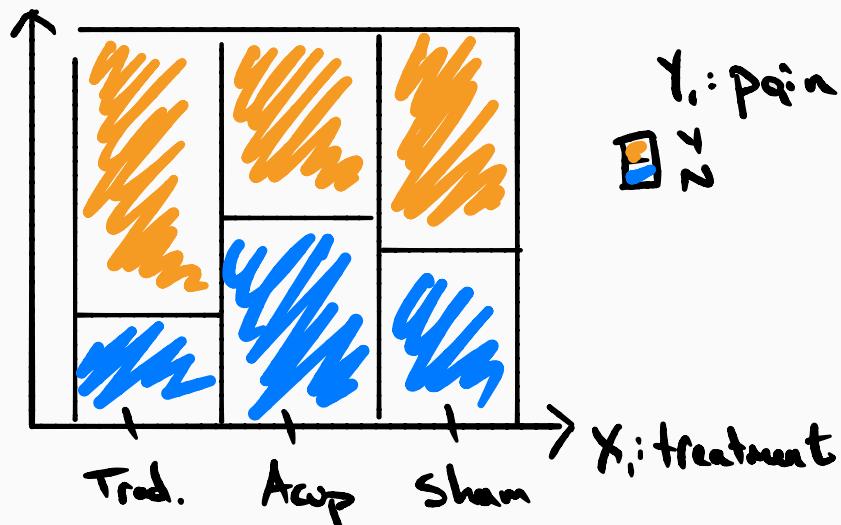
for both groups.

$Y_1 \sim Y_2$: Two > 2-level categorical variables

$Y_i \sim X_i$

Eg: Does Acupuncture relieve pain?

Visualize: bar plot



Inferential Test of Independence

Question:

- Use χ^2 statistic

$$\chi^2 = \sum_{i=1}^k \frac{(O_i - E_i)^2}{E}$$

$$P(A, B) = P(A)P(B)$$

- Find null:

- permutation

- χ^2 dist $df = (R-1)(C-1)$

Numerical Data

- visualize:
 - boxplot
 - density
 - histogram

- parameters
 - μ
 - σ
 - q_i^*
 - max
 - .
 - .
 - .

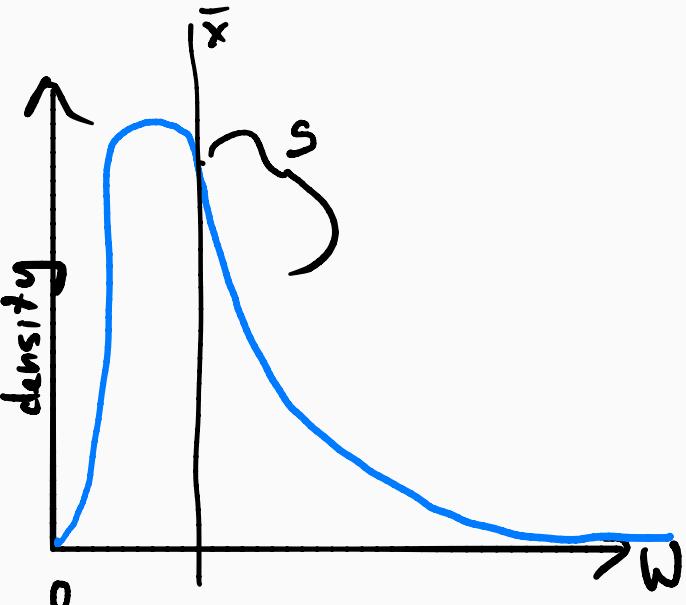
weight in kg

1.035
.917
.984
.
.
.

W : One numerical variable

E.g. Individual incomes in U.S.

Visualize: density, histogram, boxplot.



parameter
of
interest
 μ

sampling
dist
of
 \bar{x}

$$\frac{\bar{x} - \mu_0}{s/\sqrt{n}} \sim t_{df=n-1}$$

if - ind obs.

- W being nearly normal

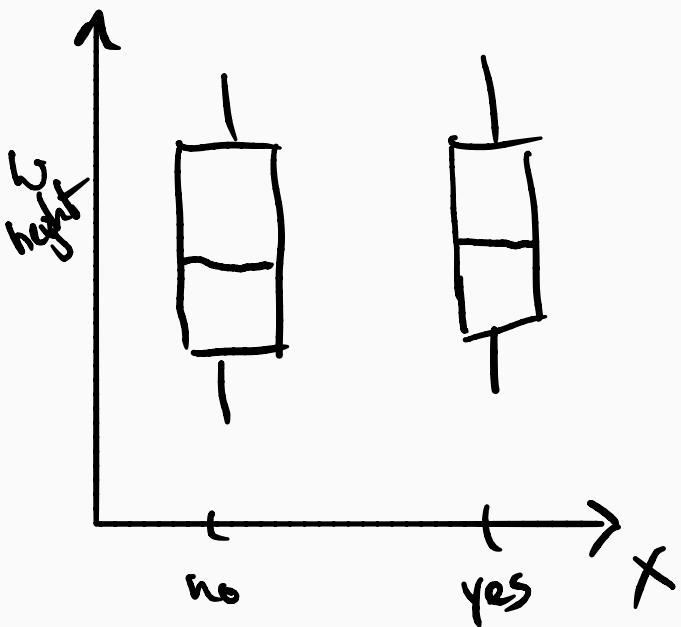
CLT

- as n increases the "nearly normal" assumption to can be relaxed.

$W \sim X$: One numerical, one 2-level categorical variable

E.g. Does average height differ between physically active and not " teens?

visualize: side by side boxplot



parameter
of
interest
 $\mu_1 - \mu_2$

sampling dist
$$\frac{(\bar{x}_1 - \bar{x}_2) - (\mu_1 - \mu_2)}{\sqrt{\frac{s_e^2}{n_1} + \frac{s_e^2}{n_2}}}$$
 n.t df =

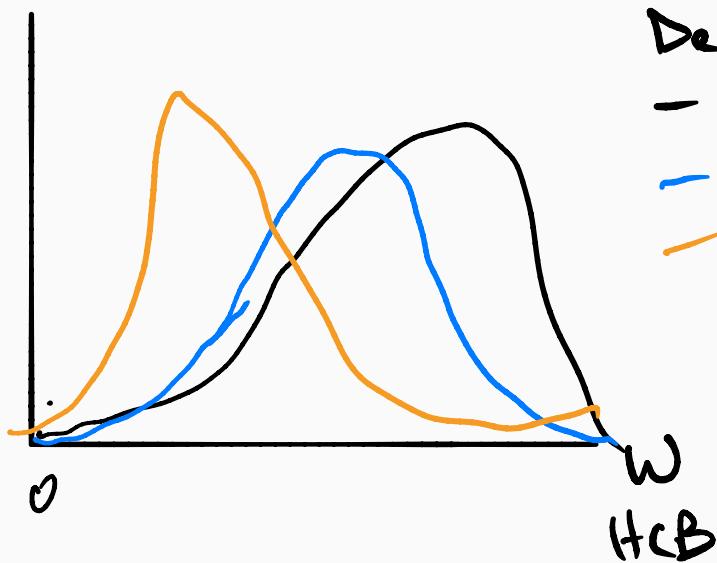
$$\min(n_1-1, n_2-1)$$

if - ind obs.
~ W nearly normal

$W \sim Y$: One numerical, one > 2-level categorical variable

E.g.: HCB and Wolf River Depth.

Visualize: Overlaid density plots.



Inferential
Question

Test
Statistic

Null
Dist. f F

$$H_0: \mu_1 = \mu_2 = \mu_3$$

$$F = \frac{MSB}{MSE}$$

1. Permute
2. F-dist

- ind obs.
- normal in each group
- constant variance