

POPULATION

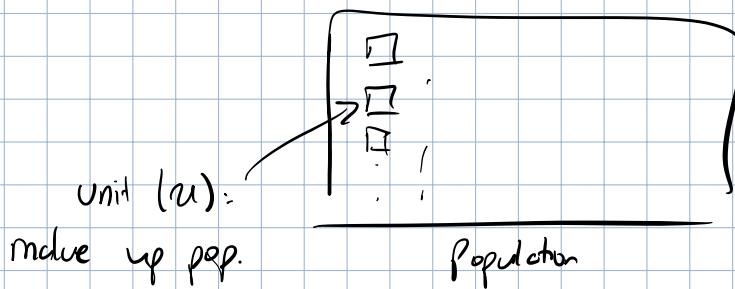
Defn: finite set of elements, like entire universe

Notation: P

Ex://

1. Population of world

2. All posts on Insta



Variable: attribute of unit

↪ Notation: $x_u \rightarrow$ attr. x of unit u

Population attribute: property of P , summarizes

↪ Notation: $a(P)$

Ex://

1. avg. income for all humans on Earth

2. total # of likes on Insta posts.

Population attributes

Pop. attributes are fn:

$$a(P) = f(y_1, \dots, y_N)$$

Ex://

1. Pop. total: $a(P) = \sum_{i=1}^N y_i$

2. Count: $a(P) = \sum_{u \in P} I_A(y_u)$, $I_A = \begin{cases} 1, & y \in A \\ 0, & y \notin A \end{cases}$

Numerical / graphical

Types of attr.

Location

Defn: describes centre of all variates

- Avg.

- Median

Skewness

Asymmetry

- o Pearson's moment of coeff.

Spread

Defn: variability of variate.

- o Population variance:

$$\sigma^2(P) = \frac{1}{N} \sum (y_i - \bar{y})^2$$

- o S. D.

- o Coefficient of variation

$$c(P) = \frac{\text{SD}_P(b)}{b}$$

Order of attributes can yield insights

↳ Notation: $y_{(1)} \leq \dots \leq y_{(N)}$ $\Rightarrow y_{(i)}$ is i^{th} smallest variate in pop.

↳ Location Attributes: min, max, median, quartiles ..

↳ Variability "": range, IQR, Median abs. deviation

$$d = \text{median}_{\text{all } p} (|y_{(p)} - \text{median}_{\text{all } p} y_i|)$$

↑ median of diff. b/w median & data points

Attribute Properties

Q: What happens if variates change? Will attribute also change?

2 important terms:

(1) Invariance: no change

(2) Equivariance: commensurate change

LOCATION

↳ Invariance:

$$a(y_1 + b, \dots, y_N + b) = a(y_1, \dots, y_N)$$

↳ Equivariance:

$$a(y_1 + b, \dots, y_N + b) = a(y_1, \dots, y_N) + b$$

Ex:// Show pop. avg. is location equivalent.

① Define original $a(p)$

$$a(p) = \frac{1}{N} \sum_{i=1}^N y_i$$

② Change the units

$$a(p^*) = \frac{1}{N} \sum_{i=1}^N (y_i + b)$$

$$\{y_1 + b, \dots, y_N + b\}$$

③ Put RHS in form of original attr.

$$a(p^*) = \frac{1}{N} \sum_{i=1}^N y_i + \frac{1}{N} \cdot N b$$

$$= a(p) + b \quad \blacksquare$$

SCALE:

↳ Invariance:

$$a(m y_1, \dots, m y_N) = a(y_1, \dots, y_N)$$

↳ Equivalent

$$a(m y_1, \dots, m y_N) = m \cdot a(y_1, \dots, y_N)$$

LOCATION - SCALE

↳ Invariance

$$a(m y_1 + b, \dots, m y_N + b) = a(y_1, \dots, y_N)$$

↳ Equivalent:

$$a(m y_1 + b, \dots, m y_N + b) = m \cdot a(y_1, \dots, y_N) + b$$

Ex:// Show pop. avg. is location-scale invariant.

$$a(m \cdot y_1 + b, \dots, m \cdot y_N + b) = \frac{1}{N} \sum_{i=1}^N (m y_i + b)$$

$$= \frac{m}{N} \sum_{i=1}^N y_i + b$$

$$= m \cdot a(p) + b \quad \blacksquare$$

REPLICATION

Duplicate pop variants k times \rightarrow new pop.

$$P^k = \{ \underbrace{y_1, \dots, y_1}_{k}, \dots, \underbrace{y_N, \dots, y_N}_{k \text{ copies}} \}$$

\hookrightarrow Invariance:

$$a(P^k) = a(P)$$

\hookrightarrow Equivalence:

$$a(P^k) = k \cdot a(P)$$

Ex:// Show pop. avg is replication invariant

$$\begin{aligned} a(P^k) &= \frac{1}{Nk} \sum_{i=1}^{Nk} y_i \\ &= \frac{1}{Nk} \sum_{i=1}^N k \cdot y_i \quad \left. \begin{array}{l} \text{diff. way of} \\ \text{writing sum} \end{array} \right. \\ &= \frac{1}{N} \sum_{i=1}^N y_i \\ &= a(P) \end{aligned}$$

Influence, Sensitivity Curves, Breakdown Points

Influence: $\Delta(a, u) = a(1, \dots, y_{u-1}, y_u, y_{u+1}, \dots, y_n)$ w/ unit u
 $\quad \quad \quad - a(1, \dots, y_{u-1}, y_{u+1}, \dots, y_n)$ w/out unit u

• If some unit u has high infl. compared to other units:

- ① Error point
- ② Interesting point

Ex:// $a(P) = \text{mean.}$

w/ u : \bar{y}

$$\text{w/out } u: \frac{1}{N-1} \sum_{\substack{l \in P \\ l \neq u}} y_l = \frac{N\bar{y} - y_u}{N-1}$$

$$\Delta(a, u) = \bar{y} - \frac{N\bar{y} - y_u}{N-1} = \boxed{\frac{y_u - \bar{y}}{N-1}}$$

In R: simulate diff. u and find Δ for all \rightarrow plot

Sensitivity Curves:

Q: how sensitive is attribute to data addition?

Defn:

$$SC(a, y) = N \left[a(y_1, \dots, y_{n-1}, y) - a(y_1, \dots, y_{n-1}) \right]$$

↑
choose y in range ($f(y)$)

Ex:// Sensitivity curve for avg.

① Define pop:

$$P = \{y_1, \dots, y_{n-1}\} \quad P^* = \{y_1, \dots, y_{n-1}, y\}$$

② SC:

$$SC(a, y) = N [a(P^*) - a(P)]$$

③ Solve for $a(P^*)$, $a(P)$

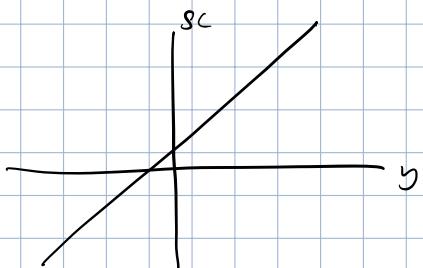
$$a(P^*) = \frac{1}{N} \left(\sum_{i=1}^{N-1} y_i + y \right) = \frac{(n-1) \bar{y}_{n-1} + y}{N}$$

$$a(P) = \frac{1}{N-1} \sum_{i=1}^{N-1} y_i = \bar{y}_{n-1}$$

④ Plug:

$$\begin{aligned} SC(a, y) &= N \cdot \left(\frac{(n-1) \bar{y}_{n-1} + y}{N} - \bar{y}_{n-1} \right) \\ &= y - \bar{y}_{n-1} \end{aligned}$$

⑤ Plot:



⑥ Interpret:

Avg. is very sensitive to extreme values

Other examples:

- | | | | |
|-------------------|---|--------------|--------------|
| ① S.C. for median | } | Review notes | |
| ② II | | | 2nd quartile |
| ③ II | | | maximum |

Breakdown points:

Q: how large of pop in data must be bad for attr. to break?

$-\infty, \infty$

error (new - orig) = Inf.

For mean: 1 datapoint $\Rightarrow 1/N$

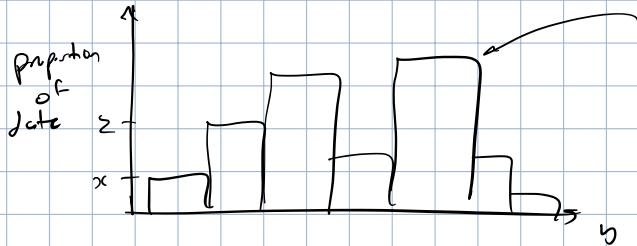
For median: $N/2$

If prop. high \Rightarrow attr. is robust to bad data.

Graphical Attributes

Visual repr. of pop is an attribute \Leftarrow summarizes pop.

Histogram:



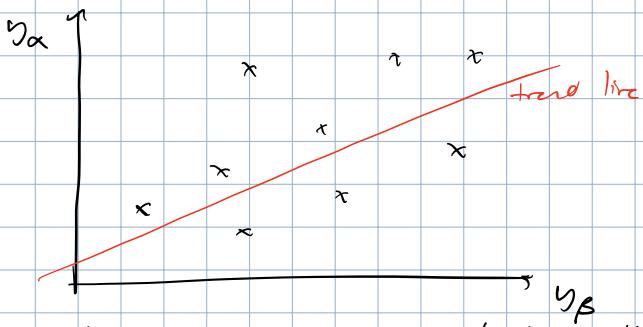
Bucket defn:

A: equal sized bins

B: same # of datapoints

L: bins width changes (lob of datapoints \rightarrow width \downarrow)

Scatter plot:



Might need to change opacity / add jitter if multiple datapoints w/ same coord.

Power transformations

Defn:

$$T_\alpha(y) = \begin{cases} y^\alpha, & \alpha > 0 \\ \log(y), & \alpha = 0 \\ -|y^\alpha|, & \alpha < 0 \end{cases}$$

Monotonic transformation:

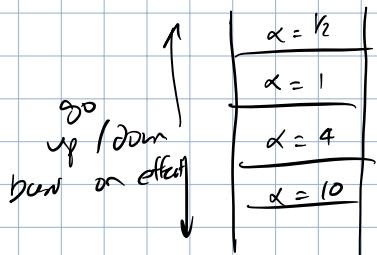
$$y_a < y_b \Rightarrow T_\alpha(y_a) < T_\alpha(y_b)$$

Why transform?

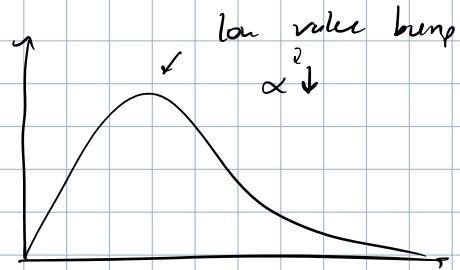
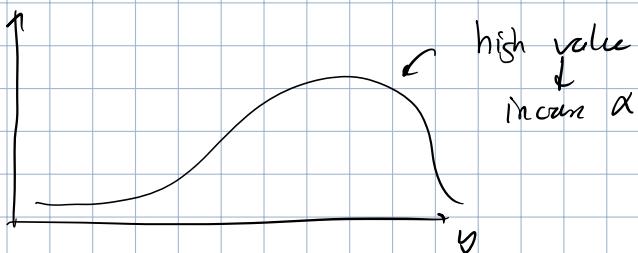
A: More symmetric histograms

B: Linearize scatterplots

Tukey ladder:

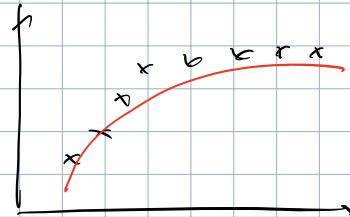


Bump rule #1: more symmetric histograms.

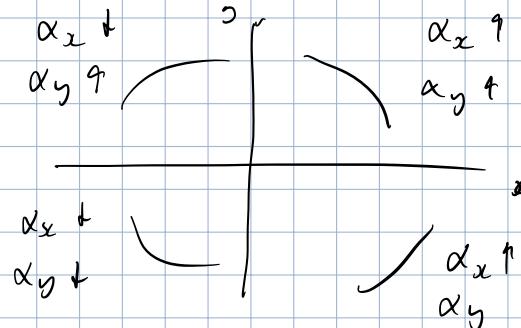


Bump rule #2: linear scatterplot

1. Find the "bump" in your graph



2. Use Tukey's ladder



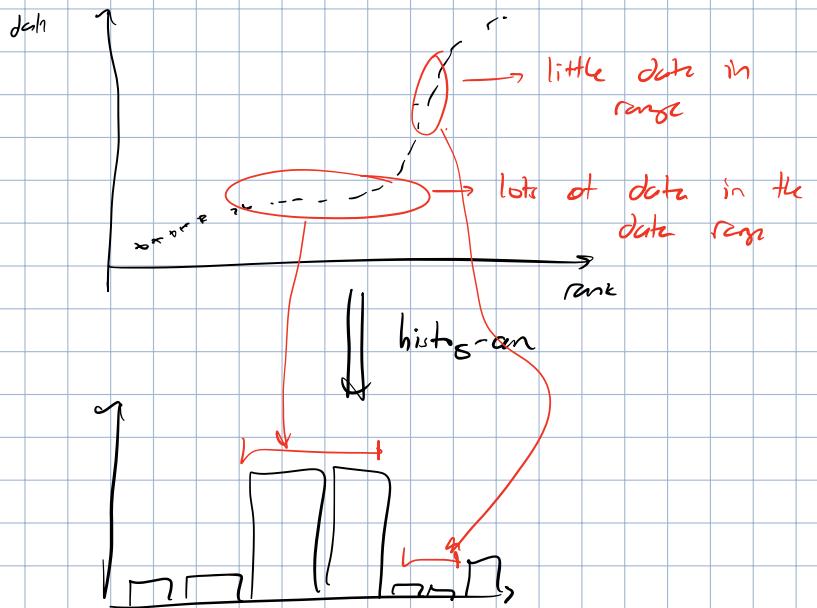
Order, Rank, Quantiles

Rank:

If y_i is $y_{(k)}$ (es. k^{th} smallest), then $r_i = k$

Ex:// $y_i : 100 \quad -50 \quad 1024 \quad 6$
 $r_i : 3 \quad 1 \quad 4 \quad 2$

Scatterplot of rank vs. data



Quantiles:

$$Q_y(p) \Rightarrow \text{datapoint w rank } p \cdot N$$

Ex: $Q_y(1/2) \Rightarrow \text{datapoint w rank } N/2 \Rightarrow \text{median}$

$Q_y(0.3) \Rightarrow 30^{\text{th}} \text{ percentile } (30\% \text{ of total} < \text{datapoint})$

Measure of location

Spread quantiles:

$$\text{IQR: } Q_y(0.75) - Q_y(0.25)$$

$$\text{Range: } Q_y(1) - Q_y(1/N)$$

Slope of 2 quantiles \propto spread of data



Implicit Attributes

Defn: attr. that summarizes pop., but usually an optimization soln.

Minimum of function: $\hat{\theta} = \underset{\theta \in \Theta}{\operatorname{argmin}} p(\theta; P) \leftarrow$ optimization func.

\uparrow
1st of all possible θ

$$= \underset{\theta \in \Theta}{\operatorname{argmin}} \sum_{u \in P} p(\theta; u)$$

argmin can also be used for max:

$$\hat{\theta}_{\max} = \underset{\theta \in \Theta}{\operatorname{argmin}} -p(\theta; P)$$

Ex://

Least squares: $\hat{\theta} = \underset{\theta \in \Theta}{\operatorname{argmin}} \sum_{u \in P} (y_u - \theta)^2$

⋮
⋮
= \bar{y}

Least absolute deviation: $\hat{\theta} = \underset{\theta \in \Theta}{\operatorname{argmin}} \sum_{u \in P} |y_u - \theta|$

⋮
⋮
= $Q_y(0.5)$

θ can be vector valued!

Linear regression:

$$\hat{\theta} = (\hat{\alpha}, \hat{\beta}) = \underset{(\alpha, \beta) \in \mathbb{R}^2}{\operatorname{argmin}} \sum_{u \in P} [y_u - (\alpha + \beta x_u)]^2$$

↳ Influence of datapoint still applicable in simplified attr.

$$\Delta(\hat{\theta}, u) = \|\hat{\theta} - \tilde{\theta}_{[u]}\|_2 \Rightarrow \text{Euclidean norm}$$

Robust Regression

Methods of dealing w/ influential datapoints?

① Don't do anything

② Remove datapoints

↳ Drastic change to line of best fit

↳ Line doesn't summarize all datapoints

③ Weight layer inflows smaller (weighted linear regression)

↳ Manual tuning

④ Robust regression: weight based off residuals (high residuals \rightarrow low weight)

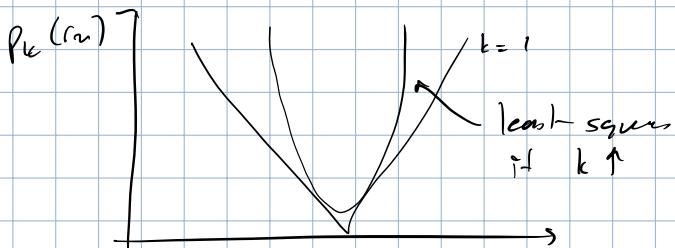
$$\hat{\theta} = \underset{\theta \in \mathbb{R}^n}{\operatorname{argmin}} \sum_{u \in P} \rho(r_u)$$

Linear reg.: $\rho(r_u) = r_u^2$

Weighted u: $\rho(r_u) = w_u \cdot r_u^2$

Huber loss:

$$\rho_k(r_u) = \begin{cases} \frac{r_u^2}{2}, & |r_u| \leq k \\ k|r_u| - \frac{k^2}{2}, & |r_u| > k \end{cases} \Rightarrow \text{lower weight for high residuals}$$



How to choose k?

Theoretical: $k = 1.3$

Common choice: $k = 1, 2, \dots$

Con: no closed form soln. for θ

Gradient Descent

Objective: $\hat{\theta} = \underset{\theta \in \mathbb{R}^n}{\operatorname{argmin}} \sum_{u \in P} \rho(\theta)$

Also:

1. Initialize $i \leftarrow 0, \theta_0$.

2. Loop:

i) Gradient:

$$s_i = \nabla \rho(\theta; p) \Big|_{\theta_i} = \begin{bmatrix} \frac{\partial \rho}{\partial \theta_1} \\ \vdots \\ \frac{\partial \rho}{\partial \theta_n} \end{bmatrix} \Big|_{\theta=\theta_i}$$

ii) Normalization:

$$d_i = \frac{s_i}{\|s_i\|}$$

iii) Line search:

$$\hat{\alpha}_i = \underset{\alpha > 0}{\operatorname{argmin}} \rho(\alpha_i - \alpha d_i)$$

iv) Make step:

$$\theta_{i+1} = \theta_i - \alpha d_i \quad \text{step of descent}$$

v) Check converge:

If θ is convex \Rightarrow return θ_{i+1}

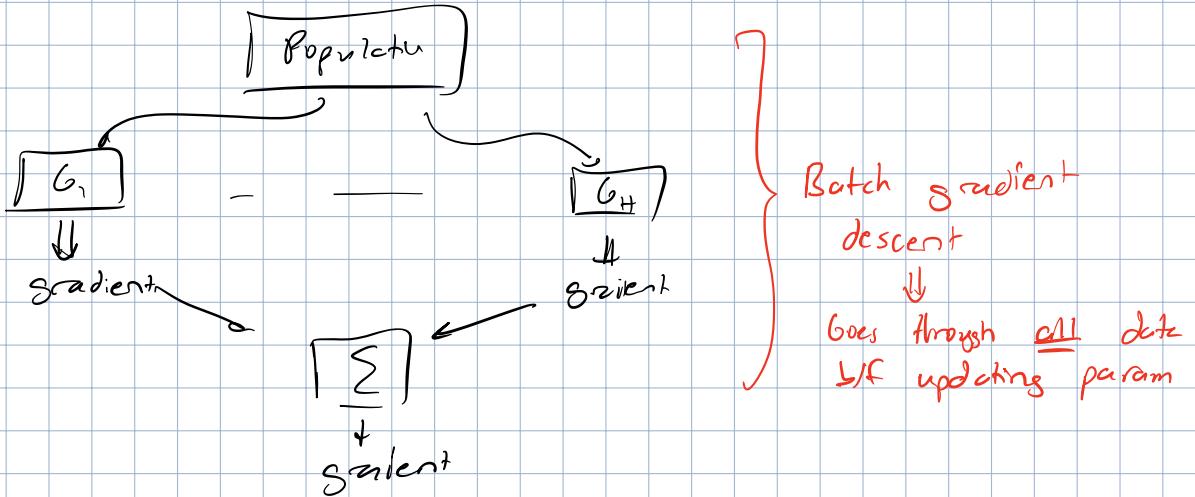
Else: $i \leftarrow i + 1$, loop

\checkmark l. nom of
 θ_{i+1} & θ_i
is $< \epsilon$

Doesn't always get global minima

Batch gradient descent

Batch gradient calculations b/c not dependent on each other



Another optimization? calculate gradient off sample of pop. & use fixed step size.

Batch sequential:

1. Divide pop. into batches
 2. Run G.D. on batch
 3. Update θ
 4. Move to next batch
- until θ converge.

Batch stochastic:

1. Randomly sample pop.
2. Run G.D. on sample

3. Update θ

4. Repeat 1 \rightarrow 4 until converge.

Batch algos take more steps, but less time/step

Systems of Equations

All defn of implicit attribute: θ that solves

$$\psi(\theta; \rho) = \sum_{u \in P} \psi_u(\theta; u) = \vec{0}$$

sys. of equations

Ex:// 1.

$$\psi(\theta; \rho) = \sum_{u \in P} \psi_u(\theta; u) = \sum_{u \in P} w_u(y_u - \theta)$$

θ that best solves \Rightarrow weighted avg.

2. Weighted least squares problem:

$$\psi(\theta; \rho) = \sum_{u \in P} \psi_u(\theta; u) = \sum_{u \in P} \rho_u' (y_u - \alpha - \beta(x_u - c)) \begin{pmatrix} 1 \\ x_u - c \end{pmatrix} = \vec{0}$$

Q: How do we find θ that solves sys. of equations?

A: Root finding methods!

① Newton's Method

1. $i \leftarrow 0$

2. Loop:

a. Update: $\hat{\theta}_{i+1} = \hat{\theta}_i - \frac{\psi(\hat{\theta}_i; \rho)}{\psi'(\hat{\theta}_i; \rho)}$

b. Check convergence: $\hat{\theta}_i \approx \hat{\theta}_{i+1}$?

 ↳ Yes: return $\hat{\theta}_{i+1}$

 ↳ No: Loop

3. Return $\hat{\theta}_i$

Iteratively Reweighted Least Squares

Obj: Find $\hat{\theta} = (\alpha, \beta)$ that solves:

$$\hat{\theta} = \underset{\theta \in \mathbb{R}^2}{\operatorname{argmin}} \sum_{u \in P} \rho \left(\underbrace{y_u - \alpha - \beta(x_u - c)}_{\text{residual } r_u} \right)$$

Solving derivative:

$$g = \nabla p(\theta; \rho) = \sum_{u \in P} \rho'(r_u) \begin{bmatrix} 1 \\ z_u \end{bmatrix} = 0$$

Recasting into weighted least square:

$$\begin{aligned} g &= \sum_{u \in P} \left(\frac{\rho'(r_u)}{r_u} \right) r_u z_u \\ &= \sum_{u \in P} w_u r_u z_u = 0 \end{aligned}$$

- Why? WLS has closed form solution

$$\hat{\alpha} = \bar{y} - \hat{\beta}(\bar{x} - c), \quad \hat{\beta} = \frac{\sum w_u (x_u - \bar{x})(y_u - \bar{y})}{\sum w_u (x_u - \bar{x})^2}$$

Q: How do we find $\hat{\alpha}, \hat{\beta}$?

A: Choose initial values \rightarrow update residuals \rightarrow keep going!

Also:

1. Init: $i \leftarrow 0, \hat{\theta}_i = (\hat{\alpha}_i, \hat{\beta}_i)$

2. Loop:

a) Get residuals & weights $\forall u \in P$

$$r_u = y_u - z_u' \hat{\theta}_i \rightarrow w_u = \frac{\rho'(r_u)}{r_u}$$

b) Solve WLS problem via closed form expr. for $\hat{\alpha}$ & $\hat{\beta}$

$$\sum_{u \in P} w_u r_u z_u = 0$$

c) Update param.

$$\hat{\theta}_{i+1} = \hat{\theta}_i \text{ from b)}$$

d) Check convergence $\forall u \hat{\theta}_{i+1}$ & $\hat{\theta}_i$

Works for any θ vector in linear response model: $y_u = z_u' \theta + r_u$

Newton-Raphson

Goal: Find $\theta \in \mathbb{R}^n$ st. $\Psi(\theta; \rho) = 0 \Rightarrow \Psi(\theta; \rho) = \begin{bmatrix} \Psi_1 \\ \vdots \\ \Psi_n \end{bmatrix} = k \text{ equations}$

Q: What's derivative of Ψ ?

$$\Psi'(\theta; \rho) = \frac{\partial \Psi(\theta; \rho)}{\partial \theta} = \begin{bmatrix} \frac{\partial \Psi_1}{\partial \theta_1} & \cdots & \frac{\partial \Psi_1}{\partial \theta_k} \\ \vdots & & \vdots \\ \frac{\partial \Psi_n}{\partial \theta_1} & \cdots & \frac{\partial \Psi_n}{\partial \theta_k} \end{bmatrix} \rightarrow \text{Jacobian of } \Psi$$

If we have a guess $\hat{\theta}_i$, then near $\hat{\theta}_i$, we can say

$$\Psi(\theta; \rho) \approx \Psi(\hat{\theta}_i; \rho) + \Psi'(\hat{\theta}_i; \rho) \times (\theta - \hat{\theta}_i)$$

Setting this to 0

$$\Psi(\hat{\theta}_i; \rho) + \Psi'(\hat{\theta}_i; \rho) \times (\theta - \hat{\theta}_i) = 0$$

$$\therefore \hat{\theta}_{i+1} = \hat{\theta}_i - [\Psi'(\hat{\theta}_i; \rho)]^{-1} \Psi(\hat{\theta}_i; \rho)$$

Algorithm:

1. Initialize: $i \leftarrow 0, \hat{\theta}_0$

2. Loop:

a) Update iterate:

$$\hat{\theta}_{i+1} = \hat{\theta}_i - [\Psi'(\hat{\theta}_i; \rho)]^{-1} \Psi(\hat{\theta}_i; \rho)$$

b) Check convergence or $i \leftarrow i+1$

3. Return $\hat{\theta}_i$

Note: If $\Psi = \nabla p$ (in grad descent), we are practically doing the same thing but descent controlled by Ψ'

SAMPLES

Idea: we cannot expect to have all population data to calculate attributes

↳ use samples! Attributes of sample estimate pop. attributes

Samples introduce:

① Sample error:

$$\text{Sample error} = a(p) - a(s)$$

Depends on attribute & sample

We want to reduce this

⑦ Consistency (Fischer)

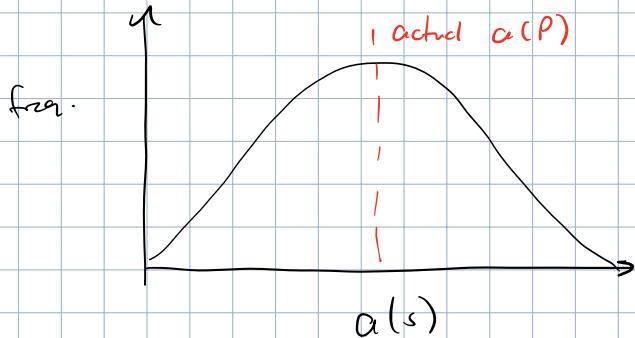
If sample was pop \Rightarrow should yield same result
 ↳ i.e. sample error $\rightarrow 0$ as $n \uparrow$

All possible samples

Idea: error depends on sample, so look @ all possible samples

of samples: $\binom{N}{n}$ total pop. size
 $\binom{n}{n}$ desired sample size

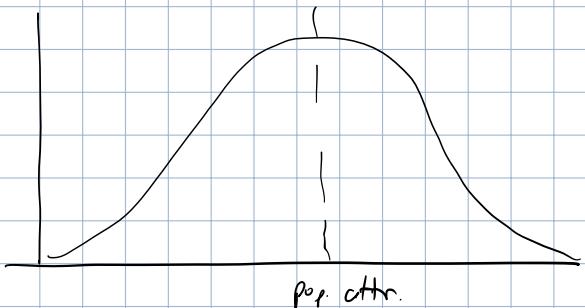
Calculate attribute on all possible samples



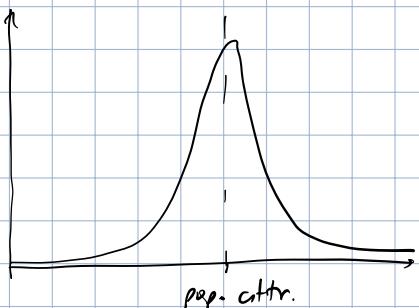
Taking sample error \forall samples \Rightarrow avg. sample error ~ 0

Consistency & Effect of Sample Size

Sample error is dependent on sample size



$\xrightarrow{\text{sample size } \uparrow}$
 \downarrow dispersion



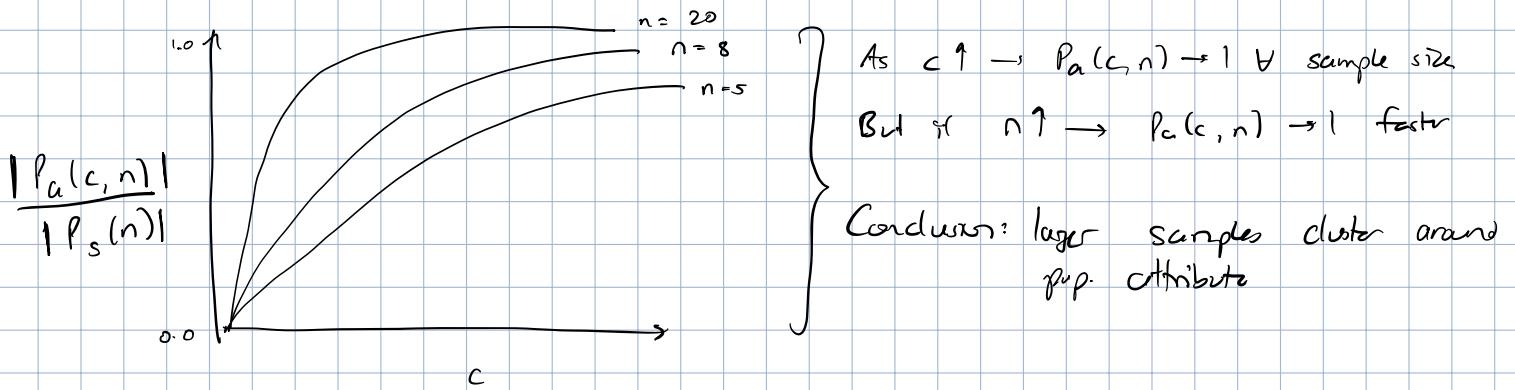
Q: How do we quantify this concentration?

Cont # of samples s s.t. $\|a(s) - a(p)\|_1 < c$

\therefore Let $P_s(n) = \{s : S \subset P \wedge |S| = n\}$

$P_a(c, n) = \{s : S \subset P_s(n) \wedge \|a(s) - a(p)\|_1 < c\}$

Then interesting graph:



Not necessary that all attributes behave like this \rightarrow could have skew, variability or consistent bias

However, almost all approach pop. attr. as sample size \uparrow

Comparison across attributes

We use relative absolute sample error to compare diff. attributes

Defn: For any $c > 0$, let

$$P_{a^*}(c, n) = \{s: s \in P_s(n) \wedge \frac{\|a(s) - a(p)\|_1}{\|a(p)\|_1} < c\}$$

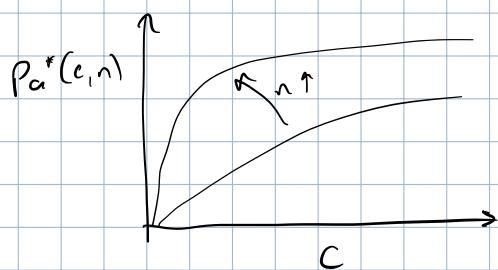
Define proportion $\forall c > 0, n \leq N$

w.r.t. some pop. attribute

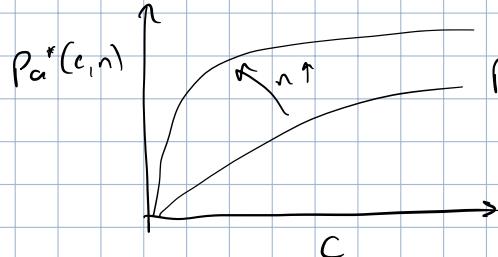
$$\left\{ \begin{array}{l} \text{consistency} \\ P_a^*(c, n) = \frac{|P_{a^*}(c, n)|}{|P_s(n)|} \end{array} \right. \begin{array}{l} \rightarrow \# \text{ of samples of size } n \text{ w/} \\ \text{relat. absolute sample error} < c \\ \hookrightarrow \# \text{ of samples of size } n \end{array}$$

For cross-attribute comparisons, we evaluate how samples track pop. value

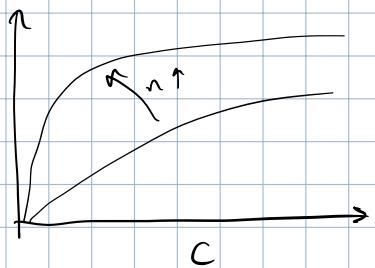
Attr. 1



Attr. 2



Attr. 3



The more left & above attr. & curves are compared to P , a outperforms R in consistency

These results may change based on population

Selecting Samples

Objective: select samples s.t. sample error is minimized

To help us, we use a sampling distribution. Properties are

- ① Exact \Leftrightarrow all possible samples available
- ② Approximate \Leftrightarrow subset of all possible samples available
- ③ In expectation \Leftrightarrow randomly choose 1 sample } -> most realistic
- ④ is often really good!
 - ↳ Frequency histogram is scaled down but proportional density is identical!

Quantifying Sample Error

Recall that avg. sample error is:

$$\# \text{ of samples in } P_S \rightarrow \frac{1}{M} \sum_{S \in P_S} [a(S) - a(P)] \quad \left\{ \begin{array}{l} \\ P_S: \text{population of all possible samples} \end{array} \right.$$

We can quantify conc. of sample errors if we select a sample S via:

① Sampling bias

Expected sample error from repeated sampling of S from P_S

$$\begin{aligned} \text{Sampling bias} &= E[a(S) - a(P)] = E[a(S)] - a(P) \\ &= \sum_{S \in P_S} [a(S) - a(P)] p(S) \end{aligned}$$

Result: If $p(S) = 1/M$, bias = avg. sample error of $a(P)$

Result: If bias = 0 $\Rightarrow a(S)$ is unbiased estimator of $a(P)$

② Sampling variance

Dispersion of sample errors

$$\begin{aligned} \text{Var} \{a(S)\} &= E[(a(S) - E[a(S)])^2] \\ &= \sum_{S \in P_S} [a(S) - E[a(S)]]^2 p(S) \end{aligned}$$

Sampling standard deviation: $SD \{a(S)\} = \sqrt{\text{Var} \{a(S)\}}$

⑦ Sampling mean squared error

To quantify distance b/w $a(s)$ & $c(p)$:

$$\text{MSE}[a(s)] = \text{Var}[a(s)] + \text{Bias}[a(s)]^2$$

Should be as small as possible!

We can write all of this in terms of random variable where:

$$P(A = a) = \sum_{s \in P(a)} p(s) I[a(s) = a]$$

Sampling Mechanisms

We want to create a mechanism to sample based on units of a pop., not selecting a sample from a population of samples

Sampling mechanism defined by:

① $P(u)$: selecting unit u

② $P(u | k, s_{k-1})$: selecting u on k^{th} draw given we're observing sequence $s_{k-1} = (u_1, \dots, u_{k-1})$, u is k^{th} one in sample

To determine $p(s)$ (prob. of selecting sample s): sum $P(s_n)$ over all pairs of s_n .

Mechanism 1: Simple Random Sampling w/out Replacement

$$\textcircled{1} \quad P(u) = 1/N$$

$$\textcircled{2} \quad P(u | k, s_{k-1}) = \frac{1}{N-k+1} = \frac{1}{\underbrace{N-(k-1)}_{\substack{\text{Choose } k-1 \\ \text{units already, so } N-(k-1) \text{ units left}}}}$$

\therefore Probability of sequence s_n :

$$P(s_n) = \frac{1}{N} \cdot \frac{1}{N-1} \cdot \frac{1}{N-2} \cdots \cdot \frac{1}{N-n+1}$$

Probability of selecting sample of size n

$$\binom{N}{n} = \frac{1}{M}$$

Mechanism 2: Simple Random Sampling w/ Replacement

$$\textcircled{1} \quad P(u) = 1/N$$

$$\textcircled{2} \quad P(u | k, s_{k-1}) = \frac{1}{N}$$

\therefore Probability of S_n : $\Pr(S_n) = \left(\frac{1}{N}\right)^n$

\therefore Probability of sample s : $\frac{\binom{n!}{n_1! n_2! \dots n_N!}}{N^n} \Leftarrow n_k$ is # of duplicates of unit k

Mechanism 3: Simple Random Sampling - Weird Hybrid

Idea: do sampling w/ replacement but remove duplicates

\hookrightarrow Sample size b/w 1 & n

Unit Inclusion Probabilities

Inclusion probability of unit u :

$$\pi_u = P(u \in S) = \sum_{s \in P_S} p(s) I[u \in s]$$

$$D_u = \begin{cases} 1 & \text{if } u \in s \\ 0 & \text{o.w.} \end{cases}$$

$$\text{Thus: } P(D_u = 1) = \pi_u, \quad P(D_u = 0) = 1 - \pi_u$$

$$\hookrightarrow E[D_u] = 1 \cdot P(D_u = 1) + 0 \cdot P(D_u = 0) \\ = \pi_u$$

$$\hookrightarrow \text{Var}[D_u] = E[D_u^2] - E[D_u]^2$$

$$= [1^2 \cdot P(D_u = 1) + 0^2 \cdot P(D_u = 0)] - \pi_u^2 \\ = \pi_u - \pi_u^2$$

Joint inclusion probability: both u & v are in sample

$$\pi_{uv} = \Pr(u \in S \wedge v \in S) \\ = \sum_{s \in P_S} p(s) I[u \in s, v \in s]$$

$$\therefore \text{Cov}[D_u, D_v] = E[D_u D_v] - E[D_u] E[D_v] \\ = \pi_{uv} - \pi_u \pi_v$$

Simple random sampling w/out replacement:

$$\pi_u = \frac{1 \cdot \binom{N-1}{n-1}}{\binom{N}{n}}, \quad \pi_{uv} = \frac{1 \cdot 1 \cdot \binom{N-2}{n-2}}{\binom{N}{n}} = \frac{n(n-1)}{N(N-1)}$$

Simple random sampling w replacement & hybrid

$$\pi_u = 1 - P(u \notin S) \\ = 1 - \left(\frac{N-1}{N} \right)^n$$

$$\pi_{uv} = 1 - P(u \notin S \vee v \notin S \vee u, v \in S) \\ = 1 - 2 \left(\frac{N-1}{N} \right)^n + \left(\frac{N-2}{N} \right)^n$$

Estimating Totals

Many attributes are either totals / function of totals

We can estimate these population total attr ($a(P) = \sum_{u \in P} y_u$) using Horvitz-Thompson estimator:

$$\hat{a}_{HT}(P) = \sum_{u \in S} \frac{y_u}{\pi_u} = \sum_{u \in P} \frac{y_u}{\pi_u} \cdot D_u$$

Bias of estimator:

$$\begin{aligned} E[\hat{a}_{HT}(S) - a(P)] &= E[\hat{a}_{HT}(S)] - a(P) \\ &= E\left[\sum_{u \in P} \frac{y_u}{\pi_u} D_u\right] - a(P) \\ &= \sum_{u \in P} \frac{y_u}{\pi_u} E[D_u] - a(P) \\ &= \sum_{u \in P} \frac{y_u}{\pi_u} \cdot \pi_u - a(P) \\ &= a(P) - a(P) \\ &= 0 \end{aligned}$$

Variance:

$$\text{Recall: } \text{Var}[\sum a_i x_i] = \sum_j \sum a_i a_j \text{Cov}(x_i, x_j)$$

$$\begin{aligned} \therefore \text{Var}[\hat{a}_{HT}(S)] &= \text{Var}\left[\sum_{u \in P} \frac{y_u}{\pi_u} D_u\right] \\ &= \sum_{u \in P} \sum_{v \in P} \frac{y_u}{\pi_u} \cdot \frac{y_v}{\pi_v} \cdot \text{Cov}(D_u, D_v) \\ &= \sum_{u \in P} \sum_{v \in P} \frac{y_u}{\pi_u} \cdot \frac{y_v}{\pi_v} \underbrace{(\pi_{uv} - \pi_u \pi_v)}_{\Delta_{uv}} \end{aligned}$$

Yates-Grundy Formulation:

$$\text{Var} [\tilde{\alpha}_{HT}(s)] = \frac{1}{2} \sum_{u \in P} \sum_{v \in P} \Delta_{uv} \left(\frac{y_u}{\pi_u} - \frac{y_v}{\pi_v} \right)^2$$

* Note: since $\text{bias} = 0$, $MSE = \text{Variance}$ for this estimator

We can actually estimate the variance of HT by using an HT estimate:

$$\text{Var} [\tilde{\alpha}_{HT}(s)] = \sum_{u \in P} \sum_{v \in P} \Delta_{uv} \frac{y_u}{\pi_u} \frac{y_v}{\pi_v} = \sum_{(u,v) \in P(u,v)} q_{u,v} \quad \Delta_{uv} \frac{y_u}{\pi_u} \cdot \frac{y_v}{\pi_v}$$

Thus:

$$\text{Var} [\tilde{\alpha}_{HT}(s)] = \sum_{(u,v) \in S_{u,v}} \frac{q_{u,v}}{\pi_{uv}} = \sum_{u \in S} \sum_{v \in S} \left(\frac{\pi_{uv} - \pi_u \pi_v}{\pi_{uv}} \right) \frac{y_u}{\pi_u} \frac{y_v}{\pi_v}$$

↳ sample of all possible pairs

Note:

$$SE [\tilde{\alpha}_{HT}(s)] = \widehat{SD} [\tilde{\alpha}_{HT}(s)] = \sqrt{\widehat{\text{Var}} [\tilde{\alpha}_{HT}(s)]}$$

Why do we care? Via HT, we get:

- ① Estimate of $f(\text{pop. to tot})$
- ② Estimate of variance of $f(\text{pop. to tot})$ } unbiased!

Sampling Design

Defn:

$$(P_s, p(s))$$

Pop. of samples

Probability of sample selected

We can choose samples & $p(s)$ s.t. MSE is low

From HT:

$$MSE [\tilde{\alpha}_{HT}(s)] = \text{Var} [\tilde{\alpha}_{HT}(s)] = -\frac{1}{2} \sum_{u \in P} \sum_{v \in P} \Delta_{uv} \left(\frac{y_u}{\pi_u} - \frac{y_v}{\pi_v} \right)^2$$

Choose sampling design to make this close to 0

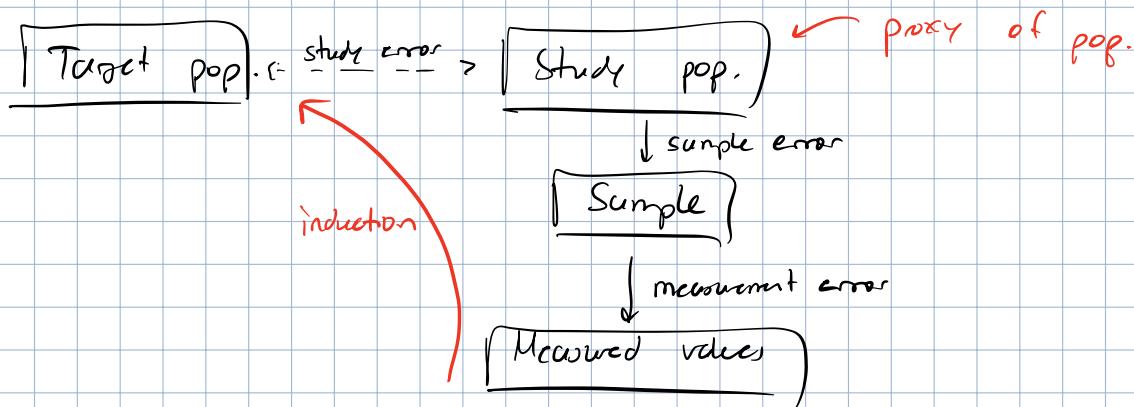
Ways to make MSE $\rightarrow 0$:

- ① $\pi_u \propto y_u$
- ② $\pi_u \propto x_u$ if x_u correlates w/ y
- ③ $y_u \approx y_v \Rightarrow \pi_u \approx \pi_v$

INDUCTIVE INFERENCE

If we have probabilistic sampling, we know relative freq. of certain attr. being selected → like on insurance policy

Sources of Error



Study error: $\alpha(P_{\text{study}}) - \alpha(P_{\text{target}})$ ← cannot be controlled by prob sampling

Total error from sample S:

$$\alpha(S) - \alpha(P_{\text{target}}) = [\alpha(S) - \alpha(P_{\text{study}})] + [\alpha(P_{\text{study}}) - \alpha(P_{\text{target}})]$$

Sources of measurement error:

- ① Measuring device bias/variability
- ② Operator error
- ③ Method error (eg. measuring length horizontal vs vertically)

Comparing Sub-populations

Goal: compare attributes of 2 populations, $\alpha(P_1)$ and $\alpha(P_2)$

First, let's define a measure of diff. b/w 2 attr.:

A: Measure of location attr.: $\alpha(P_1) - \alpha(P_2)$

B: Measure of spread IL: $\alpha(P_1) / \alpha(P_2)$

Q: If we want to figure out the diff., how do we know if it's bigger than the diff. observed if $P_1 \approx P_2$?

↳ Q: What does it mean for $P_1 \approx P_2$?

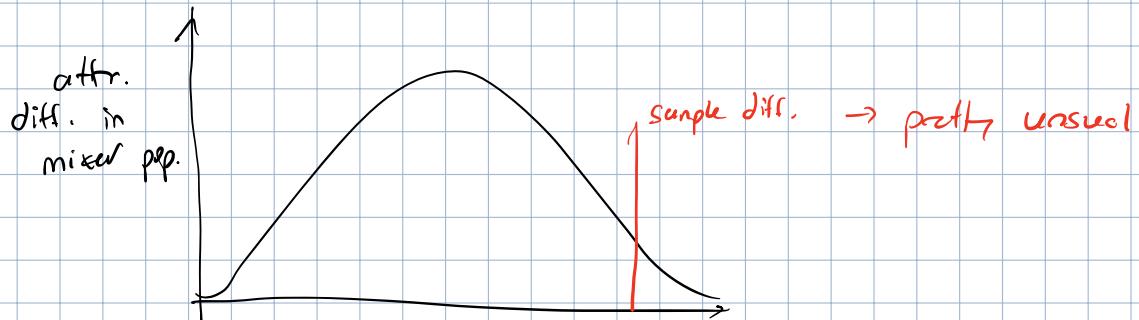
$$P_1 = \{x_1, \dots, x_m\} \quad \xrightarrow{\text{random mix}} \quad P_1^* = \{y_1, x_2, \dots, y_n\}$$

$$P_2 = \{y_1, \dots, y_n\} \quad \xrightarrow{\text{into new pop. \& draw 2 new pop.}} \quad P_2^* = \{x_1, y_2, \dots, x_m\}$$

If attr. don't change, then $P_1 \approx P_2$

A: To know if diff. is unusual \rightarrow hypothesis test for quantification

Intuition: draw many samples from mixed pop. & put attr. diff. in histogram. If sample diff. is unusual, then we won't see very many instances of sample diff. in histogram



Anatomy of a Significance Test

Goal: quantify diff. in $\alpha(P_1)$ vs. $\alpha(P_2)$ in a randomly mixed pop.

Steps:

① Null hypothesis (H_0):

Hypothesize P_1 & P_2 from same pop. $\rightarrow \alpha(P_1) \sim \alpha(P_2)$

H_A : alternative hypothesis (\overline{H}_0)

* Note: we do not state $H_0: \alpha(P_1) = \alpha(P_2)$, we say it to is 2 pop. drawn randomly from same pop. *

② Discrepancy measure

Quantifies inconsistency b/w observed data & H_0

↳ large values of measure \rightarrow evidence against H_0

Notation: $D(P_1, P_2)$

Ex:// 1. Average:

$$D(P_1, P_2) = |\bar{Y}_1 - \bar{Y}_2|$$

2. SD:

$$D(P_1, P_2) = \left| \frac{SD(P_1)}{SD(P_2)} - 1 \right|$$

3. Pop. avg of P_1 is greater than pop. avg. of P_2

$$D(P_1, P_2) = \bar{Y}_2 - \bar{Y}_1$$

③ Calculate observed discrepancy

D_{obs}

* Note: this is only looks at 1 attr of diff. *

④ Compare observed discrepancy against sampling distr. of discrepancy from randomly mixed pop.

⑤ Calculate p-val

$$p\text{-val} = \Pr(D \geq D_{\text{obs}} \mid H_0 \text{ is true})$$

If p-val small, 2 options:

A: H_0 is true & we observe very unusual data

} Choose to believe B

B: H_0 is false

To calculate p-val exactly: take all permutations of pop & derive sample distr.

L: Approximation: just take M random sub-pops,

$$\hat{p\text{-val}} = \frac{1}{M} \sum_{i=1}^M I[D(p_{1,i}^*, p_{2,i}^*) \geq D_{\text{obs}}]$$

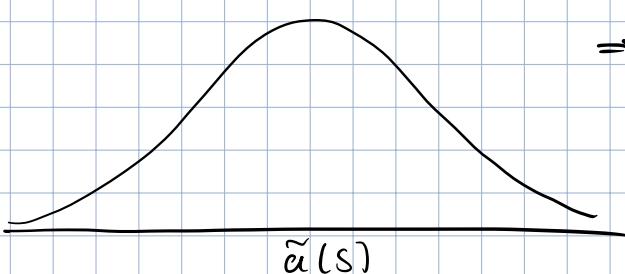
Count # of random mixed sub-pops
have discrepancy $\geq D_{\text{obs}}$

Errors in sig. tests:

① Type I: reject H_0 but actually H_0 true

② Type II: accept H_0 but actually H_0 false

Revisiting Sampling Distributions



⇒ ① Construct all possible samples of particular size

② Took $\alpha(S)$ per sample

③ Distr.

To find 95% of all possible $\hat{\alpha}(S)$ ⇒ take 2.5th & 97.5th quantiles of sampling distr.
Too much!

Q: How can we recalc sampling distr without finding all possible samples

A: 3 methods:

① Take logz # of samples & approx.

② Approx to normal distr. (why? CLT!)

↳ Useful if attr. is a function of total. O.U. not general.

③ Approx via resampling method

Random vs. Observed Intervals

If pop. attr. of interest is $\alpha(P) = \mu = \frac{1}{N} \sum_{u \in P} y_u$, then estimator from sample is:

$$\bar{Y} \sim N(\mu, \frac{\sigma^2}{n}) \quad \sigma^2 = \frac{1}{N} \sum_{u \in P} (y_u - \mu)^2$$

CLT

Standardized:

$$Z = \frac{\bar{Y} - \mu}{\frac{\sigma}{\sqrt{n}}} \sim N(0, 1)$$

* Important: if pop. is finite, include finite pop. corr.

$$\text{Var} \{ \bar{Y} \} = \frac{\sigma^2}{n} \left(\frac{N-n}{N-1} \right) \quad (\text{assuming } \sigma^2 \text{ is known from pop.})$$

↳ If $N \gg n$, $\frac{N-n}{N-1} \rightarrow 1$

↳ For completeness:

$$\bar{Y} \sim N\left(\mu, \frac{\sigma^2}{n} \left(\frac{N-n}{N-1} \right)\right), \quad Z = \frac{\bar{Y} - \mu}{\frac{\sigma}{\sqrt{n}} \sqrt{\frac{N-n}{N-1}}} \sim N(0, 1)$$

Random interval:

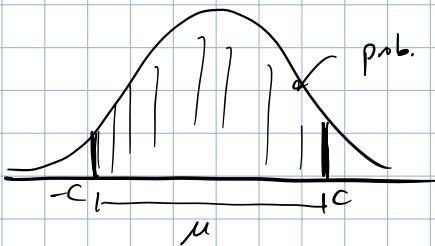
Want to find interval that contains μ w/ prob. $\underline{1-p}$

$$\left[\bar{Y} - c \times \frac{\sigma}{\sqrt{n}} \sqrt{\frac{N-n}{N-1}}, \bar{Y} + c \times \frac{\sigma}{\sqrt{n}} \sqrt{\frac{N-n}{N-1}} \right]$$

Random blc defined by R.V. Can substitute \bar{Y} w/ \bar{y}

coverage
prob.

Derivation:



$$1 - \rho = \Pr \left(-c \leq \frac{\bar{Y} - \mu}{\frac{\sigma}{\sqrt{n}} \sqrt{\frac{n-1}{n}}} \leq c \right)$$
$$= \Pr \left(\bar{Y} - c \dots \leq \mu \leq \bar{Y} + c \dots \right)$$

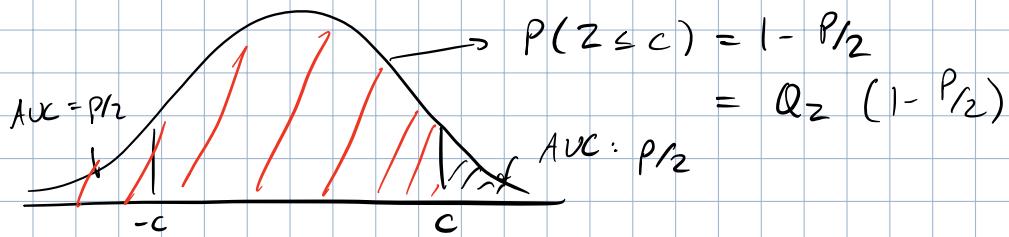
Meaning? μ is covered by interval $100(1-\rho)\%$ of time

How to determine c ?

$$c = Q_z \left(1 - \frac{\rho}{2} \right)$$

$\uparrow z \text{ distr.}$

Why?



Observed interval

We usually only have 1 sample → switch \bar{Y} to \bar{y}

$$\left[\bar{y} - c \times \frac{\sigma}{\sqrt{n}} \sqrt{\frac{n-1}{n}}, \bar{y} + c \times \frac{\sigma}{\sqrt{n}} \sqrt{\frac{n-1}{n}} \right]$$

Known as confidence interval. Not random ⇒ either has or doesn't have μ

If normality assumption T ⇒ 100(1- ρ)% of intervals will contain μ

Student t-based intervals

Assumption: knew $\text{SD}[\bar{Y}]$ → unrealistic

We can estimate!

L, In HT, we can $\text{SE}[\hat{\alpha}(S)] = \tilde{\text{SD}}[\hat{\alpha}(S)]$

$$\therefore \text{Pivotal quantity: } \frac{\alpha(s) - \alpha(p)}{SE[\hat{\alpha}(s)]}$$

If data is normally distr., then:

$$T = \frac{\bar{Y} - \mu}{SE[\bar{Y}]} = \frac{\bar{Y} - \mu}{\frac{\hat{\sigma}}{\sqrt{n}} \sqrt{\frac{n-n}{n-1}}} \sim t_{(n-1)}$$

↑
sample size

Interval via t-distr.:

$$[\hat{\alpha}(s) - c \times \hat{SD}[\hat{\alpha}(s)], \hat{\alpha}(s) + c \times \hat{SD}[\hat{\alpha}(s)]] \Rightarrow \begin{matrix} \text{Random Ctr} \\ \downarrow \text{random} \\ \text{random} \end{matrix} \quad \begin{matrix} \text{Random Ctr} \\ \downarrow \text{random} \\ \text{length!} \end{matrix}$$

How do we calculate c ?

↳ Same idea, but use $Q_T(1 - \alpha/2)$ instead of Q_Z

For pop. vars.:

Random interval:

$$\text{R.I.} = \bar{Y} \pm c \times \frac{\hat{\sigma}}{\sqrt{n}} \sqrt{\frac{n-n}{n-1}} \quad (\hat{\sigma} \text{ is R.V. from sampling distr.})$$

Observed:

$$\text{C.I.} = \bar{y} \pm c \times \frac{\hat{\sigma}}{\sqrt{n}} \sqrt{\frac{n-n}{n-1}} \Rightarrow (\hat{\sigma} \text{ is observed: } \hat{\sigma} = \sqrt{\frac{\sum (y_{in} - \bar{y})^2}{n}})$$

Pivotal quantities are important!

↳ For scale attr.: $\frac{\bar{s}(s)}{s(p)} \sim \chi^2_{\text{from}}$

Resampling Intuition

Problem: assuming we know sampling distr. of discrepancy / pivot

Q: How can we figure out sampling distr. without taking all possible samples?

A: ① Take M samples & take attr. → find distr.

↳ Problem: requires repeated sampling. We usually only have 1 sample.

② Take B samples from original sample \leq ratio than pop. \leq

Using S as approx. of study pop.!

Need to sample w/ replacement since S is small

Each sample from S is bootstrap sample

Take attr. from each sample & find distr.

Bootstrap Method

Bootstrap estimate of distr. = bootstrap sample \rightarrow take attr. \rightarrow find distr. of attr.

$$P_{\text{true}} \leftrightarrow P_{\text{study}} \leftrightarrow P^* \approx S$$
$$\downarrow \quad \quad \quad \downarrow$$
$$S^*$$

V. powerful method of estimating dist. b/c only 1 sample

Depends if S is good approx of P \rightarrow can only be as good as $a(S)$

Bootstrap sample error: $a(S^*) - a(S)$

\uparrow
bootstrap sample

Bootstrap standard deviation:

$$SE = \hat{SD}_* \left[\bar{\alpha}(S) \right] = \sqrt{\frac{\sum_{b=1}^B [\alpha(S_b^*) - \bar{\alpha}^*]^2}{B-1}} \Rightarrow \bar{\alpha}^* = \frac{1}{B} \sum_{b=1}^B \alpha(S_b^*)$$

↳ For pop. avg:

$$\hat{SD}_* \left[\bar{Y} \right] = \sqrt{\frac{\sum_{b=1}^B (\bar{y}_b^* - \bar{y}^*)^2}{B-1}} \Rightarrow \bar{y}_b^* = \text{avg of bootstrap sample } b$$
$$\bar{y}^* = \frac{1}{B} \sum_{b=1}^B \bar{y}_b^*$$

Compare to the original method of estimating $SD[\bar{Y}]$, bootstrap slightly overestimates

Pros of bootstrap distribution:

① Estimating sampling bias

$$\text{Sampling bias} = \text{avg. bootstrap sample error} = \bar{\alpha}^* - a(S)$$

If we know estimator is biased, then we can correct by defining new estn.

$$a^*(s) = a(s) - \text{bias}$$

Using estimate of bias:

$$\begin{aligned} a^*(s) &= a(s) - \bar{a}^* + a(s) \\ &= 2a(s) - \bar{a}^* \end{aligned}$$

② C.I. & hypothesis testing

Bootstrap Confidence Intervals

We will use bootstrap distr. as proxy for sampling \Rightarrow approx. C.I.

Naive normal theory intervals:

$$\text{C.I. now defined as } \bar{y} \pm c * \underbrace{\hat{\text{SD}}[\bar{Y}]}_{\substack{\text{still from} \\ \text{normal dist.}}} \quad \downarrow \quad \text{bootstrap estimate}$$

Percentile method

- ① Generate B bootstrap samples from sample S
- ② For each sample, calculate $a_b = a(S_l^*) \rightarrow$ estn.
- ③ Take quantiles using a_1, \dots, a_B
- ④ Create C.I.

This is equivalent to any transformation of attributes

Pro: simple & transform equivariant

PREDICTION

Accuracy of Prediction

Obj: determine if model $\hat{y} = \mu(\hat{x})$ ($y = \mu(x) + \text{error}$) is accurate

Average prediction squared error (APSE):

$$\frac{1}{N} \sum_{u \in S} (y_u - \hat{\mu}(x_u))^2 \quad \left. \right\} \begin{array}{l} \text{quantifies distance b/w} \\ \text{actual & predicted val.} \end{array}$$

↳ Proportional to residual mean squared error $\left(\sum_{i=1}^n r_i^2 \right)$

↳ Lower APSE \rightarrow better prediction

We often notice that making model more complex leads to lower APSE but becomes over-tuned

We should use APSE on data not used to train the model:

$$\text{APSE}(P, \hat{\mu}_s) = \frac{1}{N} \sum_{u \in P} (y_u - \hat{\mu}_s(x_u))^2$$

↑
data to calculate APSE
↑
data for model calc.

$$= \underbrace{\left(\frac{n}{N}\right) \text{APSE}(S, \hat{\mu}_s)}_{\text{APSE from sample}} + \underbrace{\left(\frac{N-n}{N}\right) \text{APSE}(T, \hat{\mu}_s)}_{\text{APSE from data not used in training}}$$

Train model on sample S

Usually we use:

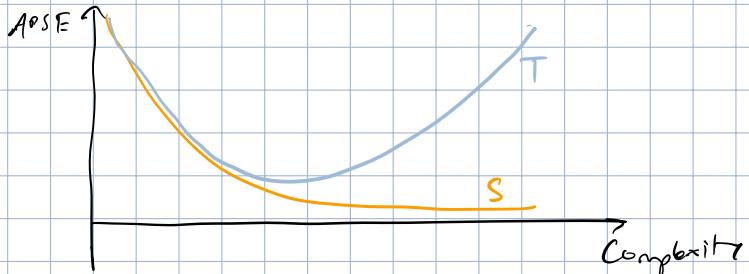
$$\text{APSE}(T, \hat{\mu}_s) = \frac{1}{N-n} \sum_{u \in T} (y_u - \hat{\mu}_s(x_u))^2$$

$$\downarrow \quad \text{If } n \ll N \rightarrow \text{APSE}(T, \hat{\mu}_s) = \text{APSE}(P, \hat{\mu}_s)$$

T: data not used in training

Usually $\text{APSE}(T, \hat{\mu}_s) > \text{APSE}(S, \hat{\mu}_s)$

Overfitting: make model too complex \rightarrow not great @ predicting outside of sample



APSE highly depends on quality of sample we choose to use for training

Predictions over Multiple Samples

Problem w/ naive APSE: error depends on initial sample picked for fitting model

Sdn:

① Take multiple samples & fit model $\hat{\mu}_{s_j}(\vec{x})$

② Take $\text{APSE}(P, \hat{\mu}_{s_j}) \forall j \rightarrow$ sampling distr. of APSE

③ Take avg. APSE

$$\text{APSE}(P, \tilde{\mu}) = \frac{1}{M} \sum_{j=1}^M \text{APSE}(P, \hat{\mu}_{s_j})$$

We can decompose APSE($\rho, \bar{\mu}$) into 3 different parts:

↳ Preamble: $\tau(\bar{x})$ is the relationship b/w y & x observed in pop.

Since x can have duplicate vals. w diff. y -vals, we partition pop. into:

$$\rho = \bigcup_{k=1}^K A_k \quad \text{where } A_k \text{ is clusters of units w/ same } x \text{ value}$$

For each cluster:

$$\tau(x_k) = \frac{1}{n_k} \sum_{u \in A_k} y_u$$

of unib
in cluster

↳ Decomposition:

Note: $\bar{\mu}(x) = \frac{1}{M} \sum_{j=1}^M \hat{\mu}_{sj}(x) \xrightarrow{\text{estimated}} \text{Avg. predictor func. across } M \text{ samples}$

$$\begin{aligned} \text{APSE}(\rho, \bar{\mu}) &= \frac{1}{M} \sum_{j=1}^M \frac{1}{N} \sum_{u \in \rho} (y_u - \hat{\mu}_{sj}(x_u))^2 \\ &= \frac{1}{M} \sum_{j=1}^M \frac{1}{N} \sum_{u \in \rho} (y_u - \tau(x_u))^2 + \overbrace{\frac{1}{M} \sum_{j=1}^M \frac{1}{N} \sum_{u \in \rho} (\hat{\mu}_{sj}(x_u) - \tau(x_u))^2}^{\text{MSE-like!}} \\ &= \text{Ave}_x (\text{Var}[Y|x]) + \overbrace{\frac{1}{M} \sum_{j=1}^M \frac{1}{N} \sum_{u \in \rho} (\hat{\mu}_{sj}(x_u) - \bar{\mu}(x_u))^2}^{\text{Var}[\bar{\mu}]} + \overbrace{\frac{1}{N} \sum_{u \in \rho} (\bar{\mu}(x_u) - \tau(x_u))^2}^{\text{Bias}^2[\bar{\mu}]} \end{aligned}$$

\downarrow
avg. of conditional variance of
 y given x
Integ. of $\mu(x)$

Note that:

$$\text{APSE}(\rho, \bar{\mu}) = \left(\frac{n}{N} \right) (\text{APSE}(\rho, \bar{\mu}) \text{ on samples used by } \mu)$$

$$+ \underbrace{\left(\frac{N-n}{N} \right)}_{\text{Dominates if } n \ll N, \text{ also}} (\text{APSE}(\rho, \bar{\mu}) \text{ on samples not used by } \mu)$$

more fair since out-of-sample perf.