

QUANTITATIVE METHODS 1A

* → Eric Scott + DESCRIBE MY QUANTITATIVE BACKGROUND
 TA: Cameron Sces

Wed. lectures when Monday is a holiday

PROJECTS - Weeks 7-10 (DATASET PRESENTATIONS)

GRADES: HW 50% ; Final 50%

RESOURCES: IDRE, UCLLA EDD / SPSS

↳ Courses on R, SPSS, PYTHON, ETC.

INDIVIDUALS - UNIT OF ANALYSIS

VARIABLES (FACTOR) - CHARACTERISTICS OF INDIVIDUALS

4 TYPES OF VARIABLES (STEVENS 1946)

• NOMINAL - Gender, Ethnicity, Political Party (UNORDERED CATEGORICAL)

• ORDINAL - Distance b/w values is not constant (RANKING)

• INTERVAL - No meaningful zero → Temperatures, Time, Democracy Score

• RATIO - Meaningful zero - AGE

DUMMIES VARIABLE - Binary, Categorical (0,1)

GRAPHS - Purpose: EXAMINE THE CHARACTERISTICS OF THE DISTRIBUTION OF A SINGLE VARIABLE.

↳ SEE TUFTE, THE VISUAL DISPLAY OF QUANTITATIVE INFORMATION
 ILLUMINATE UNDERLYING PATTERNS OF INTEREST, HIGH DATA-TO-INK RATIO, NO CHART SUNK, NO CAVESOME.

"What's wrong?" ?

George Napoleon's Army Grade Distribution

Get Tufte - The Visual Display of Quantitative Information.

Pie Charts - Drawbacks: DIFFICULT TO SKILL ZERO
HARD TO DISTINGUISH B/W. CATEGORIES

Bar Charts - X-Axis is Nominal, Then Gross Area & Meaningful.

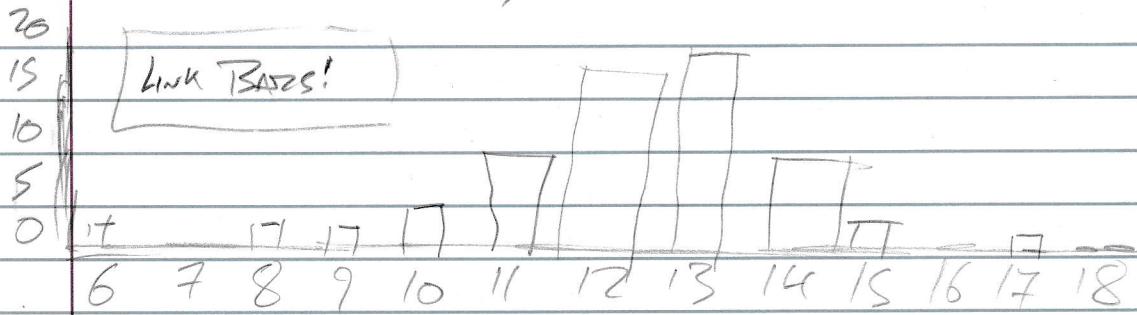
* Stem + Leaf Plots - CHOOSE THE STEM (TENS, ONES, ETC)

| | |
|--------|------------------------|
| Review | 1 3 3 4 : 13, 13, 14 |
| | 2 6 2 : 26, 22 |
| | 3 5 9 : 35, 39 |
| | 4 9 : 49 |

TENS ONES

DISTINCTION BETWEEN GRAPHS THAT ARE USEFUL AND THOSE THAT ARE MERELY TECHNICALLY CORRECT.

Histogram: Pick Bins, Accurate Data to Bins



Describing A Graph: CENTER, SHAPE, SPREAD, OUTLIERS

DESCRIPTIVE STATISTICS

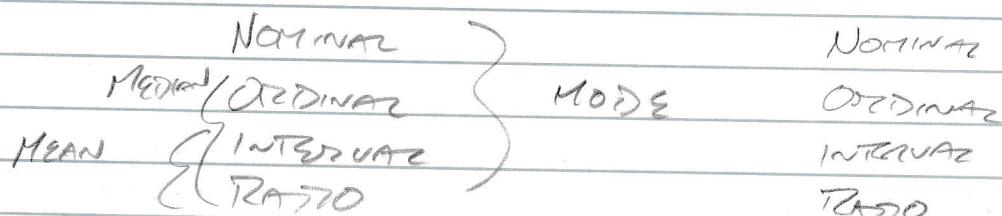
MEAN: $\bar{x} = \frac{\sum_{i=1}^n x_i}{n}$

x_i = individual obs.
 n = # of obs.
 \rightarrow uses all data points
 \bar{x} = average of x

MEDIAN: MIDDLE VALUE (OR AVG. OF 2 MIDDLE VALUES)
 \rightarrow less vulnerable to outliers "INSUFFICIENT STATISTIC"

MODE: MOST COMMON ENTRY (REPEATED VALUES)
 \rightarrow DEPENDS ON

ALL OF THE ABOVE ARE MEASURES OF CENTRIC TENDENCY



VARIANCE: $\sigma^2 = \sum_{i=1}^n (x_i - \bar{x})^2$

STANDARD DEVIATION: $\sigma = \sqrt{\sigma^2}$ ("measures the average distance from the mean")

RANGE: $x_m - x_1$

IQR: $x_{75\%} - x_{25\%}$ (INTER-QUARTILE RANGE)

VARIANCE CAN OFTEN BE MORE SIGNIFICANT THAN AVERAGES
 (Ex - Polarization)

* THE SAMPLE MEAN IS THE VALUE THAT MINIMIZES THE VARIANCE.

QUANTILES : THE i TH QUANTILE OF A VARIABLE IS
THE VALUE SUCH THAT $i\%$ ARE SMALLER THAN THAT
VALUE AND $100 - i\%$ ARE GREATER THAN THAT VALUE.

SKEWNESS :

KURTOSIS :

BOX PLOTS :

CROSS TABS : IDENTIFY CAUSE + EFFECT

QUANTITATIVE METHODS WEEK 2

PROBABILITY - THE LONG-RUN FREQUENCY OF EVENTS

SIMULATION - USING SIMULATIONS TO EXPLORE PROBABILITY

↳ Many important results are reached through simulation

PERSONALISTIC POINT: "THE PROBABILITY OF SEEING A RESULT THIS FAR (OR FURTHER) FROM THE NULL HYPOTHESIS, IF IN FACT THE NULL IS TRUE."

FOCUS: WE WANT TO KNOW THE PROBABILITY THAT OUR RESULT HAPPENED RANDOMLY - BUT WE NEVER KNOW FOR SURE.
 (BUT WE CAN REPEAT).

TYPES OF PROBABILITY: "CLASSICAL" \Rightarrow EQUALLY LIKELY OUTCOMES
 MATHEMATICAL
 EMPIRICAL (FREQUENTIST)
 SUBJECTIVE (BAYESIAN)

"COUNT + DIVIDE" \Rightarrow DIVIDE # EVENTS BY # POSSIBLE EVENTS

MATHEMATICAL: 3 AXIOMS (KOLMOGOROV)
 GIVEN EVENTS A_i , EVENT SPACE S
 (SEE SIDES)

FREQUENTIST : The Probability of An Event
is its long-run Relative Frequency.
 $P(E)$ is the Probability of Event E.

SUBJECTIVE : "Beliefs" of Probability (Bayesian)

S - SAMPLE SPACE : $S = S_1, \dots, S_n$

TRAIL : Random Trial

Outcome : Result

$$S = 1$$



INTERSECTION : $P(A \cap B) = 1 = \text{Circles}$

UNION : $P(A \cup B) = 1 = \text{Circles}$

CONDITIONAL : $P(A|B)$

$$P(\emptyset) = 0$$

$$P(A') = P(A^c) = 1 - P(A) \Rightarrow \text{COMPLEMENT}$$

$P(A \cap B) = 0 = \text{MUTUALLY EXCLUSIVE}$

* IF TWO EVENTS ARE MUTUALLY EXCLUSIVE, THE PROB.

THAT ONE OR THE OTHER WILL OCCUR IS THE SUM OF THE PROBABILITIES THAT EACH OCCURS,

$$P(A \cup B) = P(A) + P(B)$$

* IF EVENTS ARE MUTUALLY EXCLUSIVE & COLLECTIVELY

EXHAUSTIVE, THE SUM OF THEIR PROBABILITIES MUST

ADD TO ONE.

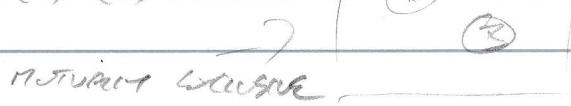


INDEPENDENCE : IF TWO EVENTS ARE INDEPENDENT,

THE PROBABILITY OF BOTH A & B OCCURRING

AS THE PRODUCT OF THEIR INDIVIDUAL PROBABILITIES

$$P(A \cap B) = P(A)P(B)$$



Get a 6-sided die w/
3.3 (expected value) on
all 6 sides.

IF EVENTS ARE NOT INDEPENDENT, THEN:

$$P(A \cap B) = \frac{P(A \cap B)}{P(B)}$$

Concurrent Sample Probabilities

Finite Sample Space: The Probability of
Any s_i is $\frac{1}{n}$. $P(s_i) = \frac{1}{n}$

Given some event A (subset of S) s_1, \dots, s_m (max)

$$P(A) = \frac{m}{n}$$

Probability Distributions: Discrete + Continuous

(See slide)

EXPECTED VALUE: "Averages" after Many Trials

$$\mathbb{E}(x) = \sum (x_i P(x_i))$$

Variance:

$$\sum (x_i - \mathbb{E}(x))^2 \times P(x_i)$$

Bernoulli Distribution: (coin toss)

$$\mathbb{E}(\text{Toss}) = 1(p) + 0(1-p) \quad (\text{Heads} = 1)$$

$$V(\text{Toss}) = p(1-p)$$

Random Variable K - THE SUM OF INDEPENDENT BERNoulli TRIALS

$$\boxed{\frac{(N)}{K} p^k (1-p)^{n-k}}$$

$$\frac{(N)}{K} = \frac{n!}{K!(n-K)!} = \frac{n(n-1)(n-2)\dots}{K(K-1)\dots \times (n-K-1)\dots}$$

CUMULATIVE DISTRIBUTION FUNCTIONS

INSTEAD OF $P(X = x_i)$ WE GETTING $P(X \leq x_i)$

PDF \Rightarrow PROB. OF SEEING EXACTLY X

CDF \Rightarrow PROB. OF SEEING LESS THAN X

DISJOINT EVENTS CAN BE ADDED INTO A CDF

Central Limit Theorem - As we increase N,

Normal Approximation:

$$P(K \geq 40) = 1 - \Phi(z) \quad \left. \begin{array}{l} \text{So:} \\ P(K \geq 40) = \\ 1 - \Phi(2.56) \end{array} \right\}$$
$$z = \frac{40 - 20}{\sqrt{\frac{40 \cdot 60}{40}}} = 2.5612...$$

GEOMETRIC DISTRIBUTION

1) USEFUL FOR PART FAILURE (TIME TO FAILURE)

HYPERGEOMETRIC DISTRIBUTION

$$P(K = k) = \frac{\binom{R}{k} \binom{N-R}{n-k}}{\binom{N}{n}}$$

K = THE NUMBER OF SUCCESSES IN n DRAWS w/o
REPLACEMENT FROM A POPULATION OF SIZE N,
CONTAINING R SUCCESSES AND N-R FAILURES

Poisson Distribution : Number of Events in A Discrete Time Period. (Often Infrequent).

$$P(K=k) = \frac{\lambda^k e^{-\lambda}}{k!}$$

Continuous Distributions:

Since the Probability of Any Given (Continuous) Point, Is Zero, we have to talk about Density Functions instead of Probability Functions.

Uniform Distribution

We Integrate the PDF to Get the CDF

PDF - Prob. Density Function = $f(x)$

CDF - Cumulative Density Function = $P(X < x) = F(x)$
 $= \int f(x) dx$

$$E(x) = \int x \cdot (f(x) dx) \quad (\text{Uniform Distribution})$$

$$V(x) = \int (x - E(x))^2 \cdot (f(x) dx)$$

Exponential Distribution:

$$f(x) = \lambda e^{-\lambda x}, x \geq 0$$

$$F(x) = 1 - e^{-\lambda x}, x \geq 0$$

$$F(x) = 0, x < 0$$

$$E(x) = \lambda^{-1} \quad V(x) = \lambda^{-2}$$

Beta Distribution

$$f(x) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} x^{\alpha-1} (1-x)^{\beta-1}$$

If $\alpha = \beta$, then
Dist. is centered
on 0.5, if not
then skewed

$$E(x) = \frac{\alpha}{\alpha + \beta} \quad V(x) = \frac{\alpha\beta}{(\alpha + \beta)^2(\alpha + \beta + 1)}$$

The Central Limit Theorem

THE MEAN OF A SAMPLE OF INDEPENDENT DATA FROM THE SAME DISTRIBUTION WITH EXPECTED VALUE μ AND FINITE VARIANCE σ^2 WILL BE NORMALLY DISTRIBUTED WITH AN EXPECTED VALUE OF μ AND A VARIANCE OF $\frac{\sigma^2}{n}$, OR STANDARD DEVIATION OF $\frac{\sigma}{\sqrt{n}}$

This means that as long as our variables "look like something like a mean", we can use the Normal Dist.

NORMAL DISTRIBUTION:

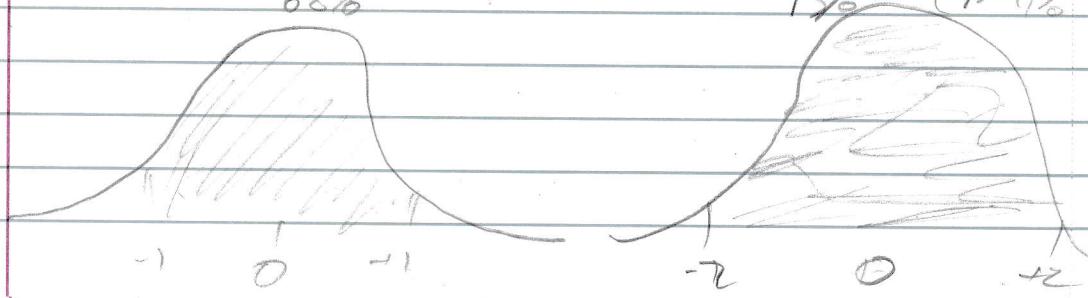
$$f(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2} \frac{(x-\mu)^2}{\sigma^2}} \quad (\text{Maximized when } x=\mu)$$

$$\mathbb{E}(x) = \mu \quad V(x) = \sigma^2 \quad D(x)$$

$$\begin{aligned}\mathbb{E}(x+y) &= \mathbb{E}(x) + \mathbb{E}(y) \\ \mathbb{E}(cx) &= c\mathbb{E}(x)\end{aligned}$$

$$\begin{aligned}\text{IF INDEPENDENT: } V(x+y) &= V(x) + V(y) \\ V(cx) &= c^2 V(x)\end{aligned}$$

$$X \sim N(\mu, \sigma^2); Z = \frac{x-\mu}{\sigma}; Z \sim N(0, 1)$$



Chi - Square Distribution

Sum of Many Normal Distributions (Squared)

Always > 0

Approaches Normal as we add more Normal Ds.

F Distribution

RATIO OF CHI-SQUARE DISTRIBUTIONS

$$K_1 \sim \chi^2_{d_1}$$

$$K_2 \sim \chi^2_{d_2} \quad G = \frac{K_1/d_1}{K_2/d_2}$$

$$G \sim F_{d_1, d_2}$$

Quantitative Methods Week 3

2 REVOLUTIONS IN SOCIAL SCIENCE

- 1) EXPERIMENTS
 - 2) CAUSAL EFFECTS
- } DRASTICALLY CHANGED RESEARCH

NULL HYPOTHESIS - INNOCENT UNTIL PROVEN GUILTY (H_0)

ALTERNATIVE HYPOTHESIS - GUILTY BEYOND A REASONABLE DOUBT
(H_A)

WHAT IS THE PROBABILITY OF SEEING A SAMPLE LIKE THE ONE WE HAVE, IF THE NULL HYPOTHESIS WERE TRUE?

WE CAN: REJECT NULL

OR FAIL TO REJECT NULL

HYPOTHESIS TESTING:

P-VALUE: THE PROBABILITY OF SEEING DATA AS EXTREME OR MORE EXTREME THAN OUR SAMPLE.

IF THE P-VALUE IS REASONABLY SMALL, WE PRESUME THAT THE RESULT ISN'T JUST RANDOM ERROR: REASON TO REJECT THE NULL.

ALPHA (α): THE PROBABILITY (P-VALUE) BELOW

WHICH WE REJECT THE NULL. TYPICAL: .05

CAN BE MANIPULATED, AVOID HA PRE-REGISTRATION

$$Z\text{-SCORE}: Z = \frac{X - \mu}{\sigma}$$

STEPS:

1) SET α

2) STATE NULL, ALTERNATIVE HYPOTHESES

3) CHECK P-VALUE AGAINST CRITICAL VALUE

(CAN REQUIRE CALCULATE Z-SCORE)

4) CONCLUDE

IF Z-SCORE $\geq z$, WE'LL REJECT THE NULL
AT A 0.05 CONFIDENCE INTERVAL

Errors : TYPE I (False Positive)
TYPE II (False Negative)

Non-Parametric Tests - Make Fewer Assumptions

Power : THE POWER OF A TEST IS THE PROBABILITY
OF CORRECTLY REJECTING THE NULL WHEN IT
IS FALSE.

204B - Week 4

Simple Bivariate Regression (OLS)

a) Minimize the sum of vertical squared deviation from the line.

b) Produce unbiased + efficient estimates

$$Y = \alpha + \beta X + \epsilon$$

L) Hypotheses invoke β (Does X affect Y ? Is Beta 0?)

Regressions give you (approximately) the Maximum Likelihood + Bayesian + Least-Squares result.

X: (AUSG (Independent)) Y: Effect (Dependent)

Relationship: Type, Direction, Strength (Accuracy of Prediction)
L) Positive/Negative

Dangerous to Project Beyond the Range of the Data

Outliers need not deviate from the trend

Outliers exert a relatively stronger effect than other data points

Correlation Coefficient: A single number that measures the strength + direction of a straight-line relationship

Line Relationship: $-1 \leq r \leq 1$

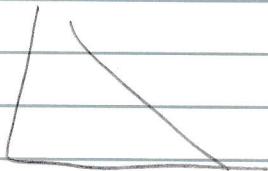
(Only looks at straight lines, vulnerable to outliers)

$$r = \frac{\sum z_{xi} z_{yi}}{n-1} \quad \text{where } z_{xi} = \frac{(x_i - \bar{x})}{s_x}$$

$$z_{yi} = \frac{(y_i - \bar{y})}{s_y}$$

$$r = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{s_x s_y (n-1)}$$

$$r = -1$$



$$r = +1$$



Start by Re-scaling X + Y to Z-scores.

That covers Four Quadrants.

| | X | Y | $\frac{14}{12}$ | $\frac{10}{10}$ | $\frac{4}{2}$ | $\frac{6}{4}$ | $\frac{6}{6}$ | $(X-\bar{x})^2$ |
|-----------------|---|----|-----------------|-----------------|---------------|---------------|---------------|-----------------|
| $n=5$ | 2 | 10 | | | | | | 4 |
| | 2 | 12 | | | | | | 4 |
| | 4 | 12 | | | | | | 0 |
| | 6 | 12 | | | | | | 4 |
| | 6 | 14 | | | | | | 4 |
| Mean: \bar{x} | 4 | 12 | | | | | | 16 |

$$s_x = \sqrt{\frac{\sum (x_i - \bar{x})^2}{n-1}} = \sqrt{\frac{16}{4}} = 2$$

$$s_y = \sqrt{\frac{\sum (y_i - \bar{y})^2}{n-1}} = \sqrt{\frac{8}{4}} = \sqrt{2}$$

Then calculate Z-scores

$$\frac{4}{8}$$

$$\underline{z_x} \quad \underline{z_y} \quad \underline{z_x * z_y}$$

1/28/19

Linear Regression

$$Y = \alpha + \beta x + \epsilon = y = \beta_0 + \beta_1 x + \epsilon$$

α : intercept ($y=0, x=A$)

Minimize the sum of squared vertical deviations from the best-fit line.

e = error term (deviations from the line)

$$y_i = \alpha + \beta x_i + e_i$$

$$\sum_{i=0}^n e_i^2 \quad (\text{minimize this quantity})$$

$$e_i = y_i - \alpha - \beta x_i$$

$$\sum_{i=0}^n (e_i)^2 = \sum_{i=0}^n (y_i - \alpha - \beta x_i)^2$$

(This says the sum of the error terms is the same as the sum of the equations)

$$b = r * \frac{\sum y}{\sum x}$$

Least-Squares Regression

$\alpha = \beta_0$ = Intercept

$\beta = \beta_1$ = Slope

| PARAMETER | ESTIMATES |
|-----------|--------------------------|
| β_0 | $\hat{\beta}_0$ or b_0 |
| β_1 | $\hat{\beta}_1$ or b_1 |
| σ | |

Vertical Distance to the Regression Line:

RESIDUALS - THE PART OF Y WE DIDN'T GET

ϵ = True Score (Unknown)

e = Estimate (Residual)

$$\text{MODEL: } Y_i = \beta_0 + \beta_1 X_i + \epsilon_i$$

$$\text{ESTIMATE: } \hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 X_i + e_i$$

The sum of squared vertical deviations is always too small (this is why we choose the line).

$$a = \bar{Y} - b \bar{X}$$

$$\hat{Y} = a + b x$$

$$Y_i - \hat{Y}_i \quad (\text{ACTUAL MINUS PREDICTED})$$

$$b = r \frac{s_y}{s_x}$$

$$S_e = \sqrt{\frac{\sum \epsilon^2}{n-2}}$$

$$\left(S_e = \sqrt{\frac{\sum \epsilon^2}{n-2}} \right)$$

Measuring the Fit of the Model:

Null Hypothesis = X doesn't affect Y ; $\beta = 0$

To measure the fit, compare the residuals from both models.

$$\sum e_{\text{REGRESS}}^2 = \sum (y_i - \hat{y}_i)^2$$

$$\sum e_{\text{NULL}}^2 = \sum (y_i - \bar{y})^2$$

OR: Partition the variance in Y into
the part we have explained by X (REG, TSS)

With the part that we haven't explained

$$TSS = \sum (y_i - \bar{y})^2$$

$$REGSS = \sum (\hat{y}_i - \bar{y})^2$$

$$TSS = \sum (y_i - \bar{y})^2$$

$$\frac{\text{Explained Variance}}{\text{Total Variance}} = \frac{REGSS}{TSS} = \frac{\sum (\hat{y}_i - \bar{y})^2}{\sum (y_i - \bar{y})^2} = R^2$$

$$\frac{REGSS}{TSS} = \frac{\sum (\hat{y}_i - \bar{y})^2}{\sum (y_i - \bar{y})^2} = R^2$$

R^2 is the percentage of variance in Y explained by X .

Regression is unbiased, efficient, consistent

1/29/19

Regression by Hand

Formulae: $S_x = \sqrt{\frac{\sum (x - \bar{x})^2}{n-1}}$ $S_y = \sqrt{\frac{\sum (y - \bar{y})^2}{n-1}}$

$$r_{xy} = \frac{\sum (x - \bar{x})(y - \bar{y})}{\sqrt{S_x S_y (n-1)}} \quad b = r_{xy} S_y \quad \alpha = \bar{y} - b\bar{x}$$

$$\alpha = \bar{y} - b\bar{x} \quad \hat{y} = \alpha + bx \quad \hat{\sigma} = S_E = \sqrt{\frac{\sum (y - \hat{y})^2}{n-2}}$$

Steps: 1) Calculate Standard Errors

- (a) Find \bar{x}
- (b) Subtract \bar{x} from x
- (c) Square resulting Product
- (d) Sum the column of Squares

Step 4

Steps

- (e) Find \bar{y}
- (f) Subtract \bar{y} from y
- (g) Square Resulting Product
- (h) Sum the column of Squares
- (i) Multiply $(x - \bar{x})(y - \bar{y})$
- (j) Calculate Standard