

Distributions and Samples

Clicker Question

The major difference between an observational study and an experiment is that

- A. An experiment manipulates features of the situation
- B. An experiment does not employ observation
- C. An observational study records what happens
- D. An observational study employs a coding system

Review

- Observational research involves careful recording and analysis of what is observed
 - Without an attempt to manipulate what happens
- In **naturalistic observation** the observer seeks to remain unobtrusive whereas
- In **participant observation** the observer becomes part of the situation
- Risks that must be minimized:
 - Observer bias
 - Reactivity
 - Anthropomorphizing

Review

- Recording observations
 - Must extract that which is to be analyzed: coding systems, etc.
- Distinguish **continuous observation** from
 - **Time sampling**
 - **Event sampling**
 - **Situation sampling**
- Analyze observations in terms of **variables**—a characteristic or feature that varies and takes on different values

Clicker Question

To determine how many vehicles travel a given road, a researcher installs a camera that takes a picture of traffic every 15 minutes. This researcher is using

- Continuous observation
- Time sampling
- Event sampling
- Situation sampling

Variables

- A variable is a characteristic or feature of an event that varies—takes on different values.
- Variables of a thrown ball:
 - velocity, momentum, direction, spin, . . .
- Variables of a World Series:
 - winner, number of games, fights, strikeouts, . . .
- Variables of human hair:
 - color, length, texture, . . .
- Variables of human cognition:
 - memory span, speed of reasoning, emotional state, . . .

Types of variables

- Variables differ in the type of measurement of the values of the variable that is possible. Sometimes one refers to types of scales rather than types of variable.
- **Categorical or nominal variables:** items can be assigned to a category (whose members can then be counted, or compared on another variable).
 - Gender: male/female
 - Major: psychology, political science, economics, . . .
 - Stellar spectra: O, B, A, F, G, K, and M
 - Organisms: Plant, Animal, Bacteria, Virus, . . .

7

Types of variables - 2

- **Ordinal or rank variables:** There is a rank-order to the values the variable may take.
 - Numbers might be assigned to the items, but since there is no metric
 - one cannot compare how much higher or lower one item on the scale is than another
- Movies: *, **, ***, ****
Class rank: top 10, next 10, etc.
Patient condition: resting and comfortable, stable, guarded, and critical
Socio-economic class: low, middle, high

8

Types of variables - 3

- **Interval variables:** equal differences between numbers assigned to items reflect equal differences between the values being measured.
 - Allows additive comparison— x is three more than y
 - But lacking a natural 0, does not permit multiplicative comparison— x is three times y
- Intelligence: IQ score
- Temperature: in degrees Celsius or Fahrenheit
- Personality: degree of extroversion

9

Types of variables - 4

- **Ratio variables:** items are rated on a scale with equal intervals and a natural 0-point.

- Allows for both additive and multiplicative comparison

Age: in year, months, days, . . .

Temperature: in degrees Kelvin

Time: in milliseconds, seconds, years, . . .

Velocity, acceleration, etc.

- Interval and ratio data often treated similarly and counted as score data

10

Clicker Question

The variable MINUTES OF COMERCIALS PER HOUR is

- A categorical or nominal variable
- An ordinal or rank variable
- An interval variable
- A ratio variable

Clicker Question

The variable PARTY AFFILIATION (Libertarian, Green, Republican, Democrat, other) is

- A categorical or nominal variable
- An ordinal or rank variable
- An interval variable
- A ratio variable

Types of Variables

- Categorical or nominal variables:** major
- Ordinal or rank variables:** patient condition
- Interval variables:** temperature in degrees Fahrenheit
- Ratio variables:** age

Score variables

Distributions of values

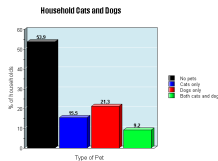
Since the values of a variable vary, they will be *distributed*

A major part of understanding a domain of objects is to describe *how* they are distributed on a given variable

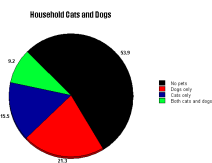
One of the best ways to present a distribution is to graph it

Nominal variables and bar graphs

Example: Profile of pet ownership in San Diego County



Value of graphs: provide an intuitive appreciation of the data

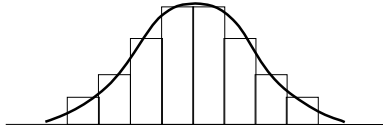


Bar graphs and pie charts work well with nominal and ordinal variables

Score variables and histograms

Since score variables are continuous, *histograms* rather than bar graphs are used

This is done by creating *bins* and tabulating the number of items in each bin



The size of bins can create radically different pictures of the distribution!

Normal and non-normal distributions

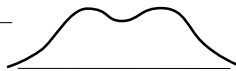
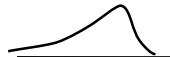
Normal distributions

Have a *single* peak
Scores *equally* distributed around the peak
Fewer scores *further* from the peak



Non-normal distributions:

- *Skewed*
- *Bimodal*



Clicker Question

The distribution below is

< 70	70-74	75-79	80-84	85-89	90-94	95+
21	12	9	23	3	21	18

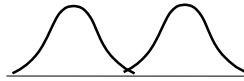
- Normal since it has one peak
- Normal since scores are equally distributed around the peak
- Not normal since the scores are not equally distributed around the peak
- Not normal since there are not fewer scores further from the peak

Describing distributions

Two principal measures:

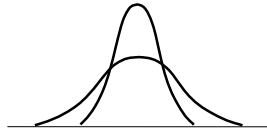
Central tendency

Two comparable distributions differing in central tendency



Variability

Two distributions with same central tendency but differing in variability



Three measures of central tendency

- Consider this distribution of values
2, 6, 9, 7, 9, 9, 10, 8, 6, 7
- Mean: the arithmetic average
 $73 / 10 = 7.3$
- Median: the score of which half are higher and half are lower = 7.5
- Mode: the most frequent score = 9

20

Which measure to use?

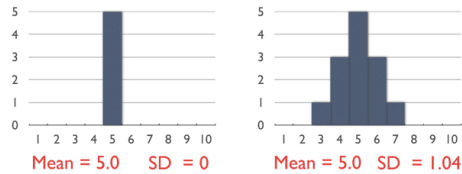
- If the distribution is normal, all three measures of central tendency give the same result
 - The mean is the easiest to calculate and the most frequently reported
- If there are extreme outliers in one direction, the mean may be distorted
 - Exam scores: 21, 72, 76, 79, 82, 84, 87, 88, 90, 91, 95
 - Mean: 78.6
 - Median: 84
- In such a case, the median gives a better picture of the central tendency of the class

Measures of variability

- How much do the scores vary?
 - Range: the lowest value to the highest value
 - Variance: $\frac{\sum (X-\text{mean})^2}{N}$
 - Standard Deviation (SD): $\sqrt{\text{Variance}}$
- Intuitive interpretation (with normal distributions):
 - One standard deviation: the part of the range in which 68% of the scores fall
 - Two standard deviations: the part of the range in which 95% of the scores fall
 - Three standard deviations: the part of the range in which 99% of the scores fall.

22

Same Mean, Different SD



23

Variance and Standard Deviation

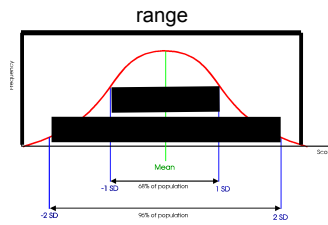
- Consider a distribution

4	5	5	6	6	6	7	7	8	Mean = 6
-2	-1	-1	0	0	0	1	1	2	X - Mean
4	1	1	0	0	0	1	1	4	(X-mean) ²
$\sum (X-\text{mean})^2 = 12$									Variance
N									9
$\sqrt{1.33} = 1.15$									SD

Range of 1 SD = $6 \pm 1.15 = 4.85$ to 7.15

Range of 2 SD = $6 \pm 2.30 = 3.70$ to 8.30

Range and Standard Deviation



Clicker Question

On the exam on which scores were distributed normally and the mean was 86 and the SD was 3.5,

- A. 68% of the scores were between 82.5 and 89.5
- B. 95% of the scores were between 82.5 and 89.5
- C. 99% of the scores were between 79 and 93
- D. 68% of the scores were between 79 and 93

Populations

- The group about which we seek to draw conclusions in a study are known as the population.
- Sometimes one can study each member of the population of interest
- But if the population is large
 - It may be impossible to study the whole population
 - There may be no need to study the whole population

Samples

- A sample is a subset of the population chosen for study.
- From studying the distribution of a variable in a sample one makes an estimate of the distribution in the actual population
- Sometimes the estimate from a sample may be more accurate than trying to study the population itself
 - U.S. Census

28

Does the sample reflect the population?

- Does the mean of the sample reflect the mean of the actual population?
 - Very unlikely that the mean of the same will exactly equal the mean of the population
 - Given the mean of a sample, what is the range within which the mean of the actual population lies?
 - Bottom line—with larger samples this range becomes smaller and smaller
 - **And this effect depends only on the size of the sample, not the size of the population sampled!**

29

Is the sample biased?

- If information about the sample is to be informative about the actual population, the sample must be representative
 - Randomization: attempt to insure that the sample is representative by avoiding bias in selecting the sample
- Risk: inadvertently developing a misrepresentative sample
 - E.g., using telephone numbers in the phonebook to sample electorate

30

Distribution on nominal variables

- Take the special case of a variable with two values (exhaustive and exclusive)
 - Heads/Tails
 - True/False
 - Born in January/not-born in January
 - Male/Femalewhere the value for each item is independent of that for other items
- Consider the likely distributions

31

Discussion Question

Consider these to be orders of births of babies in a hospital. Which is more likely?

M F M F F M F M F F

M M M M M M M M M M

F F F F F M M M M M

Each pattern is equally likely

A very different question

Consider these to be totals of births in a hospital on a given day. Which of these outcomes is more likely?

5 males / 5 females

7 males / 3 females

10 males / 0 females

From populations to samples

Start from the situation in which we know the distribution in the actual population: $p(M) = .5$

We draw a sample of a given size, say 10.

Is it possible that we could get a sample of all males?

Yes, the probability is about .001

What is the probability that we could get a sample of 7 males and 3 females?

It is about .117

What is the probability that we could get a sample of 5 males and 5 females?

It is about .246

What happens as sample size gets larger?

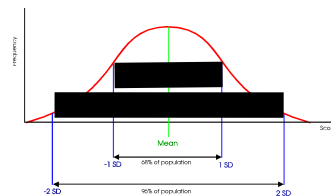
With larger sample sizes, the probability of a distribution in the sample closely approximating the distribution in the actual population increases

The important question is how much the mean of the samples will vary from the mean of the actual population

To determine this, we need to know the standard deviation (SD) of the sample.

Standard deviation and mean

In $\approx 68\%$ of samples, **the mean of the sample** will fall within 1 standard deviation of the **mean of the population**



In $\approx 95\%$ of samples, **the mean of the sample** will fall within 2 standard deviations from the **mean of the population**

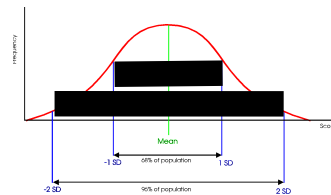
Inferring Mean of Population from Mean of the Sample

- Just as we can determine from the mean and standard deviation of the actual population where the mean of the sample is likely to be
 - We can infer from the mean and standard deviation of the sample how far from the mean of the sample the mean of the actual population will likely be
 - 68% of the time it will be within one SD
 - 95% of the time it will be within two SD
 - 99% of the time it will be within three SD
- These percentages express our confidence that the mean will be in the range specified

37

Standard deviation and mean

≈ 68% of the time, **the mean of the population** will fall within 1 standard deviation of the **mean of the sample**



≈ 95% of time, the **mean of the population** will fall within 2 standard deviations from the **mean of the sample**

SD and larger sample size

As sample size grows, the SD of the sample shrinks.

So with larger samples, the variability around the mean shrinks

Assume mean in the sample is 50%

Sample size	+/- 2 SD (95%)	+/- 3 SD (99%)
10	34.5-65.5	29.5-70.5
20	39-61	35.6-64.4
50	43-57	40.9-59.1
100	45-55	43.5-56.5
500	47.8-52.2	47.1-52.9
1000	48.4-51.6	48-52

Clicker Question

Why do most election polls study approx. 500 people even if the population is many million?

- A. It gets hard to analyze data when too much is collected
- B. It costs too much to survey more than about 500 people
- C. With 500 people the SD is already small enough to make a good estimate of the actual population
- D. With 500 people the SD is already large enough to make a good estimate of the actual population

Generalize to Score Variables

Score variables: Interval and ratio variables

With score variables, it is the scores that are distributed (not the items in a given category)

Example: age of person eating at the Food Court

Draw a sample to make inference of average age of person eating at the Food Court

<17	17	18	19	20	21	22	23	24	25	>25
6	18	23	34	32	18	26	29	14	10	10
	2	1	3	1	2		1			

Estimating real distribution

<17	17	18	19	20	21	22	23	24	25	>25
6	18	23	34	32	18	26	29	14	10	10
	2	1	3	1	2		1			
	1	2	4	6	3	2	2			

Mean of the actual population: 20.63

Mean of the sample: 19.4 20.1 Want to predict more accurately?

SD of the sample: 1.9 1.6

Range of 1 SD = 17.5-22.3 18.5-21.7 Use a larger sample size

Range of 2 SD = 15.9-24.2 16.9-23.3
