STAT/SOC/CSSS 221 Statistical Concepts and Methods for the Social Sciences

Course Introduction

Christopher Adolph

Department of Political Science and Center for Statistics and the Social Sciences University of Washington, Seattle

Our first case: John Snow's celebrated cholera map

Course details

Some examples to ponder

A typical encounter with "statistics" takes places when we read a news item on a survey:

In 1995 USA Weekend magazine asked readers to return a survey with a variety of questions about sex and violence on television. Of the 65,142 readers who responded, 98% were "very or somewhat concerned about violence on TV". Based on this survey, can we conclude that about 98% of U.S. citizens are concerned about violence on TV? Why or why not?

Some examples to ponder

A typical encounter with "statistics" takes places when we read a news item on a survey:

- In 1995 USA Weekend magazine asked readers to return a survey with a variety of questions about sex and violence on television. Of the 65,142 readers who responded, 98% were "very or somewhat concerned about violence on TV". Based on this survey, can we conclude that about 98% of U.S. citizens are concerned about violence on TV? Why or why not?
- In the November 2004 presidential election, many media outlets reported that exit polls of bellwether districts showed George Bush winning 44% of the Hispanic vote. But telephone polls days before the election showed Bush winning only 32% of the Hispanic vote. Why the discrepancy? Which number would you trust?

Popular and traditional meaning: "Statistics are numbers measured for some purpose"

Popular and traditional meaning: "Statistics are numbers measured for some purpose"

Really, numbers measured for some purpose are not statistics, but data

Actual scientific meaning:

Statistics is the collection of procedures and principles for gaining and processing information in order to make decisions when faced with uncertainty

Popular and traditional meaning: "Statistics are numbers measured for some purpose"

Really, numbers measured for some purpose are not statistics, but data

Actual scientific meaning:

Statistics is the collection of procedures and principles for gaining and processing information in order to make decisions when faced with uncertainty

Hence statistics is concerned with:

The process of data collection

Popular and traditional meaning: "Statistics are numbers measured for some purpose"

Really, numbers measured for some purpose are not statistics, but data

Actual scientific meaning:

Statistics is the collection of procedures and principles for gaining and processing information in order to make decisions when faced with uncertainty

Hence statistics is concerned with:

- The process of data collection
- Summarizing the information within data

Popular and traditional meaning: "Statistics are numbers measured for some purpose"

Really, numbers measured for some purpose are not statistics, but data

Actual scientific meaning:

Statistics is the collection of procedures and principles for gaining and processing information in order to make decisions when faced with uncertainty

Hence statistics is concerned with:

- The process of data collection
- Summarizing the information within data
- Proper interpretation of data to answer a scientific research question

Why Statistics?

Why should you learn statistics?

Statistics is the science of uncertainty and almost everything we learn is uncertain

Statististics helps us navigate/summarize/infer from oceans of data and the computer age has produced vast datasets like never before

Statistics is ubiquitous in scientific fields and in most grad programs (business, policy, medicine)

Helps us understand the big picture (general laws) *and* how each individual varies from that big picture

Key elements of a statistical study

- The individuals or objects studied (unit of analysis) and how they were selected
- The variables measured about each object
- The setting or context in which the measurement were made
- Unmeasured variables on which the subjects vary
- The magnitude of any claimed effects of differences in measured variables
- The uncertainty of these claimed effects

Statistics as part of the Scientific Method

- Observe the world / Read past studies
- Porm a research question
- Build a theory, preferably causal, to answer the question
- Choose an area to test theory
- Operationalize theory: Measure variables, generate hypotheses
- Explore and analyze the data
- Report results: is the hypothesis confirmed, or rejected?
- Replicate & repeat...

Statistics as part of the Scientific Method

- Observe the world / Read past studies
- Porm a research question
- Build a theory, preferably causal, to answer the question
- Choose an area to test theory [SELECTION]
- Operationalize theory: Measure variables, generate hypotheses [MEASUREMENT]
- Explore and analyze the data [ANALYSIS]
- Report results: is the hypothesis confirmed, or rejected? [INTERPRETATION]
- Replicate & repeat...

Populations vs. Sample

Population: Complete set of units of interest in a study

e.g., all American voters; all students at UW; all friendships of the students in this class

VS.

Populations vs. Sample

Population: Complete set of units of interest in a study

e.g., all American voters; all students at UW; all friendships of the students in this class

vs.

Sample: The subset of the population actually studied, which may be random or non-random; representative or non-representative.

e.g., 1000 voters dialed at random; 500 UW students choosen by random ID number; the first friendship mentioned by each student in this class.

When the sample includes the entire population, we call it a census

Observation: A study of the relationship between several variables based on their natural variation

e.g., a longitudinal survey tracks 1000 children over several years, noting how much violent TV each watches and whether they committed violent crimes

VS.

Observation: A study of the relationship between several variables based on their natural variation

e.g., a longitudinal survey tracks 1000 children over several years, noting how much violent TV each watches and whether they committed violent crimes

vs.

Experiment: A study of the relationship between two or more variables, one of which is controlled by the experimenters

e.g., scientists randomly assign 500 children to a group which watches violent TV, and 500 to a group which does not, then records their rates of criminal activity.

Internal validity: Whether a study is conducted well enough to make valid inferences about the relationship of variables in the population from the sample

Consider the TV violence & crime example: a study with high internal validity is one that correctly estimates the effect of TV violence on the criminal activity of those studied

Internal validity: Whether a study is conducted well enough to make valid inferences about the relationship of variables in the population from the sample

Key threats to internal validity:

Measurement error: e.g., mistaking children with missing criminal records for non-criminals

Internal validity: Whether a study is conducted well enough to make valid inferences about the relationship of variables in the population from the sample

Key threats to internal validity:

- Measurement error: e.g., mistaking children with missing criminal records for non-criminals
- Selection bias: e.g., if parents of generally well-behaved children are more likely to be enrolled in the study

Internal validity: Whether a study is conducted well enough to make valid inferences about the relationship of variables in the population from the sample

Key threats to internal validity:

- Measurement error: e.g., mistaking children with missing criminal records for non-criminals
- Selection bias: e.g., if parents of generally well-behaved children are more likely to be enrolled in the study
- Confounders: e.g., omitted variables like parental income which also affect crime

Internal validity: Whether a study is conducted well enough to make valid inferences about the relationship of variables in the population from the sample

Key threats to internal validity:

- Measurement error: e.g., mistaking children with missing criminal records for non-criminals
- Selection bias: e.g., if parents of generally well-behaved children are more likely to be enrolled in the study
- Confounders: e.g., omitted variables like parental income which also affect crime
- Reverse causation: e.g., a taste for violence leads to watching violent TV

Well-run experiments tend to have high internal validity (randomization)

But even well-run observational studies are vulnerable to above

External validity: Whether a study's findings apply to other similar situations in the real world (not just the lab)

External validity: Whether a study's findings apply to other similar situations in the real world (not just the lab)

Possibilities for failure here are endless, especially for experiments:

 Artificiality of treatment: assigned TV may have less effect than self-selected TV

External validity: Whether a study's findings apply to other similar situations in the real world (not just the lab)

Possibilities for failure here are endless, especially for experiments:

- Artificiality of treatment: assigned TV may have less effect than self-selected TV
- Selection bias: what if the participants were recruited by TV ads?

External validity: Whether a study's findings apply to other similar situations in the real world (not just the lab)

Possibilities for failure here are endless, especially for experiments:

- Artificiality of treatment: assigned TV may have less effect than self-selected TV
- Selection bias: what if the participants were recruited by TV ads?
- Ouration of treatment: what if it takes 1000s of hours of TV violence?
- Many more...

Observational studies and "natural" experiments have higher potential external validity

Cholera outbreaks were common in 19th century London; 10,000s of deaths

Contemporary theories:

- Cholera caused by "miasma" in the air coming from swamps
- Or a "poison" slowly losing strength as it passes from victim to victim?
- London doctor John Snow thought contaminated water the cause

Outbreak in 1854: 500 deaths in 10 days in Soho

Snow has Broad Street pump handle removed

Did he stop the epidemic? Prove disease can be spread by germs?

How might the newspaper "analyze" John Snows's intervention?



(plot from Tufte, Visual Explanations)

- Overwhelming tendency to view time series data this way Doesn't help us make inferences about the data
- The data aren't being compared to any other variables: time series plots don't help us devise a model of the data

How might the newspaper "analyze" John Snows's intervention?



(plot from Tufte, Visual Explanations)

- Can we do better? Specify a research question?
- Translate it into variables? Formulate some hypotheses?
- Think about internal and external validity?

In 1954, London water was provided by competing private firms Residents would walk to the nearest street pump for water Snow recorded the location of each death in real time He placed these spatial data on a map, along with the *water pumps* Was one pump, from a particular company, contaminated with cholera?

Snow's spatial analysis: Tufte redrawing



Chris Adolph (UW)

Course Introduction

Snow's Cholera Map of London



What kind of sample did Snow collect?

Snow's Cholera Map of London



What kind of sample did Snow collect?

A census of cholera victims—but what about non-victims?

Snow's Cholera Map of London



What kind of sample did Snow collect?

A census of cholera victims—but what about non-victims?

Is this an observational study or experiment?

Snow's Cholera Map of London



What kind of sample did Snow collect?

A census of cholera victims—but what about non-victims?

Is this an observational study or experiment?

Combines features of both: a natural experiment

Snow's Cholera Map of London



How do we assess the relationship between deaths (red dots) and pumps (blue triangles)?

Snow's Cholera Map of London



How do we assess the relationship between deaths (red dots) and pumps (blue triangles)?

Are we convinced that a relationship exists?

Oxford St #2 Gt Marlborou Crown Dean St So Soho Briddle St Warwick Vigo St Coventry St

Snow's Cholera Map of London

How do we assess the relationship between deaths (red dots) and pumps (blue triangles)?

Are we convinced that a relationship exists?

What additional variables should we measure?



Snow's Cholera Map of London

Fact: For any spot on the map, there is a closest pump or pumps



Snow's Cholera Map of London

Fact: For any spot on the map, there is a closest pump or pumps

Modeling Assumptions:

Some (not all) pumps are contaminated

People use the closest pump



Snow's Cholera Map of London

Fact: For any spot on the map, there is a closest pump or pumps

Modeling Assumptions:

Some (not all) pumps are contaminated

People use the closest pump

Model prediction: Pattern of deaths should correspond to nearest-pump boundaries (in blue)

20 100 m ford St #2 Oxford 15 Gt Marlborough Crown Chat 9 Dean St So Soho Briddle St Warwick ŝ -Vigo St Coventry St 5 10 15 20

Snow's Cholera Map of London

Problems?



Snow's Cholera Map of London

Problems?

Distance in a city is *complicated*—the built environment lengthens some paths.



Snow's Cholera Map of London

Problems?

Distance in a city is *complicated*—the built environment lengthens some paths.

What about outliers? Can our theory be right if some cases lie outside "nearest pump zone" of Broad St. Pump?



Snow's Cholera Map of London

Problems?

Distance in a city is *complicated*—the built environment lengthens some paths.

What about outliers? Can our theory be right if some cases lie outside "nearest pump zone" of Broad St. Pump?

Outliers could point to missing variables *or* simple randomness



Snow's Cholera Map of London

Problems?

Distance in a city is *complicated*—the built environment lengthens some paths.

What about outliers? Can our theory be right if some cases lie outside "nearest pump zone" of Broad St. Pump?

Outliers could point to missing variables *or* simple randomness

Is our model deterministic, or probabilistic?

What explains outliers in this map?



Three cases:

- A prison (work house) with its own well.
- A brewery with its own water source. Saved by the beer.
- Some distant deaths attributed to preference for Broad St. water.

Snow used his data and map to convince officials to remove the handle from the Broad Street pump.

Credited with stopping the outbreak and providing first experimental evidence for germs

Some questions to consider later:

- Did the Broad Street Pump really cause the cholera outbreak?
- 2 Did removing the handle stop it?
- Oan we measure our uncertainty about our answers to 1 and 2?