

# Choosing Cases

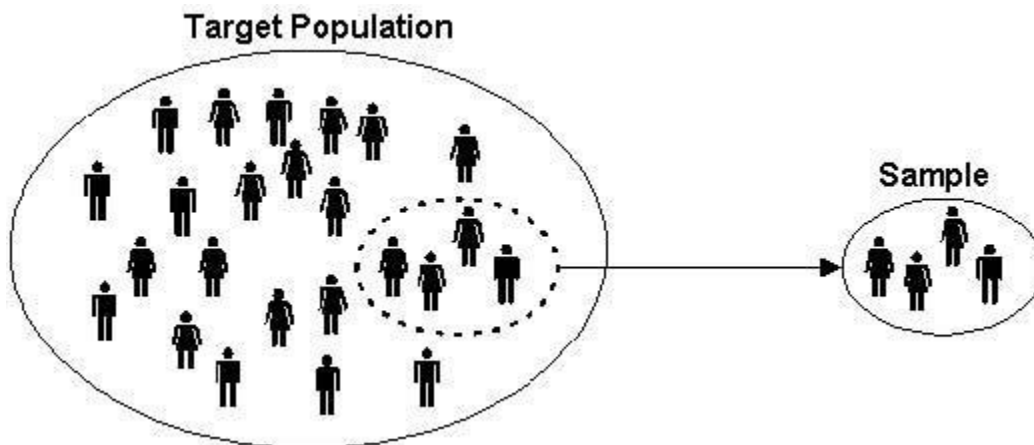
Once we have a hypothesis to explore or test, and once we have settled on a general research design, then we need to choose specific cases to analyze. Case selection is important whether we are conducting a case study of a single civil war, an experiment involving a few dozen college students, or a statistical comparison of hundreds of elections. All social scientists should learn how to select cases with care.

Done well, case selection can enhance the external validity of our research, making us more confident that our results would hold true beyond our particular study. Case selection can also help the internal validity of our research, making us more confident that our conclusions hold true within the confines of our study. Done poorly, case selection can compromise our research or even render it useless.

The purpose of this module is to help you think systematically and intelligently about case selection.

## Population vs. sample

The first question to ask is whether you plan to study the entire population of cases (often referred to as  $N$ ) or a smaller sample ( $n$ ) taken from that population.



Studying the entire population is appealing because it essentially guarantees the external validity of our research. We don't need to make inferences about what happened or why; we have analyzed every relevant case. Over time, teams of scholars have developed datasets with information about every recorded vote in the history of

Congress, and about every single interstate war over the last two centuries. Some research projects analyze entire populations like these.

Nevertheless, we usually lack the time, money, or skills needed to analyze the entire population of relevant cases. We might not even have a good way to identify the entire population; there is no master list, for example, of Kurdish rebels or newspaper stories about Senate elections. As a result of these various constraints, we typically pick a smaller sample. This is why polling firms interview 1500 people instead of 250 million. This is why undergraduates write their research papers about, say, democratization in India during the 20th century, and not about democratization in every country of the world over the last three centuries.

However, the choice between population and sample also depends on how we define the larger population. If we believe that “modern world wars” are conceptually distinct from “interstate wars” or “militarized conflicts,” then the population of modern world wars might consist of just two cases – World War I and World War II. Someone who wanted to study the origins of these world wars might actually have the resources needed to examine the entire population of cases. One can imagine other examples, such as “Communist nations in the 21st century” or “U.S. presidential elections decided by the Supreme Court” where the total population (N) is pretty small. In each of these examples, the author would have to justify the boundaries of their population. Are Communist nations in the 21st century really all that different from those that existed in the 20th century?

Suppose that someone wanted to study the relationship between motorcycle helmet laws and motorcycle fatalities in the American states. One could gather data for the most recent year available in all 50 states, which certainly sounds like the entire population of cases. If our aim is to generalize across a wider time period, though, then we would be dealing with a one-year sample. And if we are trying to generalize to some larger population of traffic laws, such as speed limits and seatbelts, and to a larger set of traffic fatalities, then our motorcycle helmet cases would also qualify as a sample.

Frankly, it seems unlikely that all you would want to accomplish in this example is to figure out what happened with one specific kind of law in one year. Doing so would really limit the larger significance of your work. Thus, the choice of sample versus population connects back to the larger aims of the study. When choosing cases we always need to ask ourselves, “What puzzle am I hoping to solve? To what scholarly literature or policy debate am I trying to contribute? What, then, is the population of relevant cases?”

## Sampling: random vs. deliberate

The vast majority of the time, for practical or conceptual reasons, we are dealing with samples. At the most general level, we need to decide whether to choose a sample of cases randomly or deliberately. One might think that random selection would always be preferred because the sample would more likely resemble the entire population, thus giving our study added external validity. This intuition is correct – as long as the number of cases is pretty large. If the number is small, then one might randomly select an atypical sample, which would actually hurt external validity.

You can take a real course about probability and statistics to understand why, or you can accept the following example as a rough proof. Let's imagine that a polling firm wanted to know what American adults think about a controversial issue like immigration. If the firm randomly selected just two people – let's call them Border Wall Bob and No Amnesty Nancy – it might conclude that all Americans have strongly negative views toward immigrants. And those conclusions would be wrong. If that same firm chose 1000 or 1500 Americans at random, it would be much more likely to identify the full range of attitudes, as well as the correct distribution. (The sample would rarely look exactly like the population, but it would probably be close if the firm sampled correctly.) With so many cases, a few extreme values in any direction will not distort the entire sample.

Thus, if the research design is based on a statistical comparison of many cases, scholars will probably choose their cases randomly. (The large number of cases will have the added benefit of helping us to establish the internal validity of our research: we can become more confident in concluding whether our measures are correlated, and whether any apparent relationships could be spurious.) If the research design is a detailed case study, however, the cases will almost always be chosen deliberately.

With experimental designs, the cases could be chosen deliberately or randomly. A lab experiment will probably not rely on a random sample of individuals; researchers will usually have to take whoever is willing or required to participate in the experiment. A survey or field experiment, on the other hand, might select at random a large number of individuals, voting precincts, villages, development projects, or some other unit of analysis. Such random selection of cases will help the external validity of the study, while experimental controls and random assignment of cases will generate internal validity.

Whether we choose cases randomly or deliberately, we are concerned about generating a biased sample. Some types of bias originate with the researcher. Suppose

you wanted to sample opinions from the entire college campus, but you only distributed surveys to three freshmen dorms. That sample would not reflect the full range of students on campus, and could bias the results if freshmen held different opinions from upperclassmen. Other types of sample bias are beyond the researchers' control — sometimes just bad luck. We might distribute surveys to a variety of dorms on campus, yet the main people who filled them out and returned them might be freshmen. Therefore, after taking a sample, it often makes sense to compare it to whatever is known about the larger population.

## Generating a random sample

To learn different ways of choosing cases randomly, you can consult standard research methods textbooks, which often do a good job of teaching this skill. See, for example, chapter 7 in Johnson and Reynolds, *Political Science Research Methods* 7th edition, or chapter 6 in Kellstedt and Whitten, *The Fundamentals of Political Science Research* 2nd edition. There you will encounter simple, systematic, stratified, and cluster random samples. You can also find helpful videos on-line, such as these two:

### Sampling Methods

<https://www.youtube.com/watch?v=FtZavr0eaw>

### Sampling: Simple Random, Convenience, systematic, cluster, stratified - Statistics Help

<https://www.youtube.com/watch?v=be9e-Q-jC-0>

[Note: both of these videos discuss “convenience sampling,” which they don’t exactly endorse. Convenience samples and snowball samples are both nonprobability samples in which each element or group within the population does not have an equal chance of being selected. The external validity of such samples is thus highly suspect. Nonprobability samples are used occasionally in social scientific research, but not often.]

It is certainly possible to combine strategies as well. A survey research firm conducting an exit poll on Election Day could start with a simple random sample of congressional districts, then a systematic random sample of voting precincts within those districts, and finish with a stratified random sample of individuals who showed up to vote at those precincts. Someone analyzing trends in media coverage of terrorism might

analyze only those years ending in 0, 2, 5, and 8, and then collect a simple random sample of stories for each year.

## Picking cases deliberately

Standard methods textbooks are pretty useless if you intend to choose cases deliberately. That's too bad, for it means that students planning to conduct case studies receive practically no guidance about a crucial step in the research process. One reason for this gap, I suspect, is that many social scientists view deliberate case selection with suspicion. A crafty researcher could pick one or two cases to prove, well, just about anything. What is supposed to be reputable social science could easily degenerate into intellectual sleight-of-hand or trickery. For a playful analogy, watch how master magician Ricky Jay manages to reveal just the right cards from a full deck:

### **Ricky Jay – Card Control**

<https://www.youtube.com/watch?v=y8lKh8YB9uQ>

### **Ricky Jay: 4 Queens 3 Ways**

<https://www.youtube.com/watch?v=JNUepjt6Qml>

While we might be delighted to watch someone manipulate cards so effortlessly, we could be outraged to discover a political scientist doing something similar when he or she picked cases to study. We would seriously doubt the study's internal validity, external validity, or both. Carefully selecting a few vivid examples to “prove” a general point is common among policy advocates and strong partisans, but it is not good practice for social scientists.

To learn more about deliberate case selection, especially for case studies, I would recommend reading chapter 6 in Lipson, *How to Write a BA Thesis*; pages 77-88 of Van Evera, *Guide to Methods for Students of Political Science*; and, if you're feeling ambitious, chapters 3 and 5 in Gerring, *Case Study Research*.

Because snappy YouTube videos about deliberate case selection are so rare, I will highlight some of their advice in the table below. One general strategy is to look for ways of maximizing the number of observations within each case; a related strategy is to find cases with analytically-useful variation (i.e., variation linked to the hypotheses we wish to explore or test). In both instances, we are trying to approximate the analytic leverage that comes with large-n statistical comparisons. But another strategy is simply

to emphasize a distinctive strength of the case study design – identifying causal links and mechanisms through careful process tracing – and to choose a single case or a few cases that will enable the researcher to study a piece of the political world in real depth.