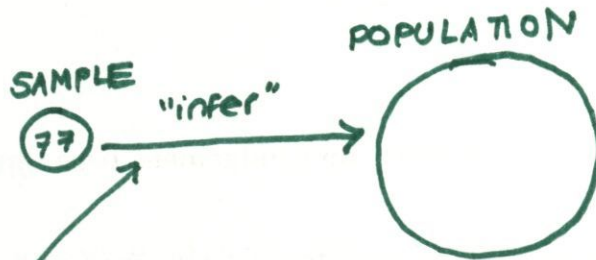


2 Data collection



Data collection is one of the most important parts of engineering statistics. If collected properly, data can make formal inferences easy to complete and easy to understand. On the other hand, if data is collected poorly, it can become nearly impossible to salvage a badly designed study and gain insights.

This chapter covers the general principles of data collection, ideas for effective experimentation, and examples of common experimental setups.

2.1 Sampling

Q: The most common question engineers ask about data collection is
How many observations do I need?

A: The answer depends on the variation in response that one expects.

Often we want to answer a question (conduct a study) about an identifiable, concrete population of items, but we want to use a **sample** to represent this (typically) much larger population.

Why? **save time, save money, maybe measurements destroy the sample, impossible (population too large/unattainable)**

Example 2.1. Measuring some characteristics of a sample of 20 electrical components (note: this is one sample with 20 units; the sample size is 20) from an incoming lot of 200.

If a sample is to be used to stand for a population, how that sample is chosen becomes very important.

A sample should

be representative of the population

SAMPLE

$n=20$

POPULATION

$n=200$

2.1.1 Systematic and judgement based methods

Definition 2.1. In systematic sampling, create a list of every member of the population. From the list, randomly select the first sample element from the first k elements on the population list. Thereafter, we select every k^{th} element on the list.

Disadvantage: **It can fail when cyclical patterns are present.**

Definition 2.2. In judgement-based sampling, select based on the opinion of an expert.

Disadvantage: **subject to unconscious/subconscious bias and preconceptions**

2.1.2 Simple random sampling

(SRS)

Definition 2.3. A simple random sample of size n from a population is a sample selected in such a manner that every collection of n items in the population is a priori equally likely to compose the sample.

equal chance of being selected
all possible combinations of sample units have an

Example 2.2. A statistics instructor wanted to know how many hours per week her students spend watching cat videos on YouTube. Rather than asking each one of them, she puts all of their names in a hat and draws out 10. This is a simple random sample of size 10.

Steps to randomly sample mechanically:

1. Let M be the number of digits in the number N , where N is the population size.
2. Give each member of the population an M -digit label.
3. Move through the table of random digits (Table B.1) from left to right, top to bottom, selecting population members for the sample when you encounter their indices (ignoring indices that have already been chosen) until you have selected n units for the sample.

Table B.1
Random Digits

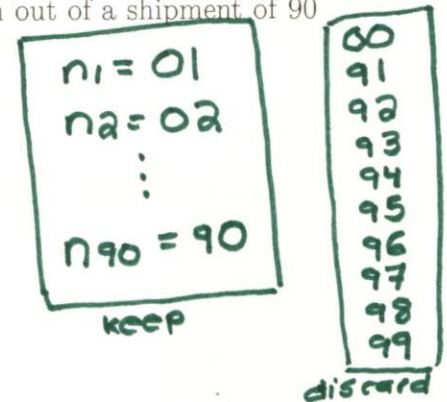
17159	66144	05091	13446	45633	13684	66024	91410	51351	22772
30156	90519	95785	47544	66735	35754	11088	67310	19720	08379
59069	01722	53338	41942	65118	71236	01932	70343	25812	62275
54107	58081	82470	59407	13475	95872	16268	78436	39251	64247
99681	81295	06315	28212	45029	57701	96327	85436	33614	29070

$$N = 90$$

$$M = 2$$

Example 2.3. Take a simple random sample of 12 units of pig iron out of a shipment of 90 units.

12, 15, 61, 44, 05, 09, 11, 34, 46,
45, 65, 31



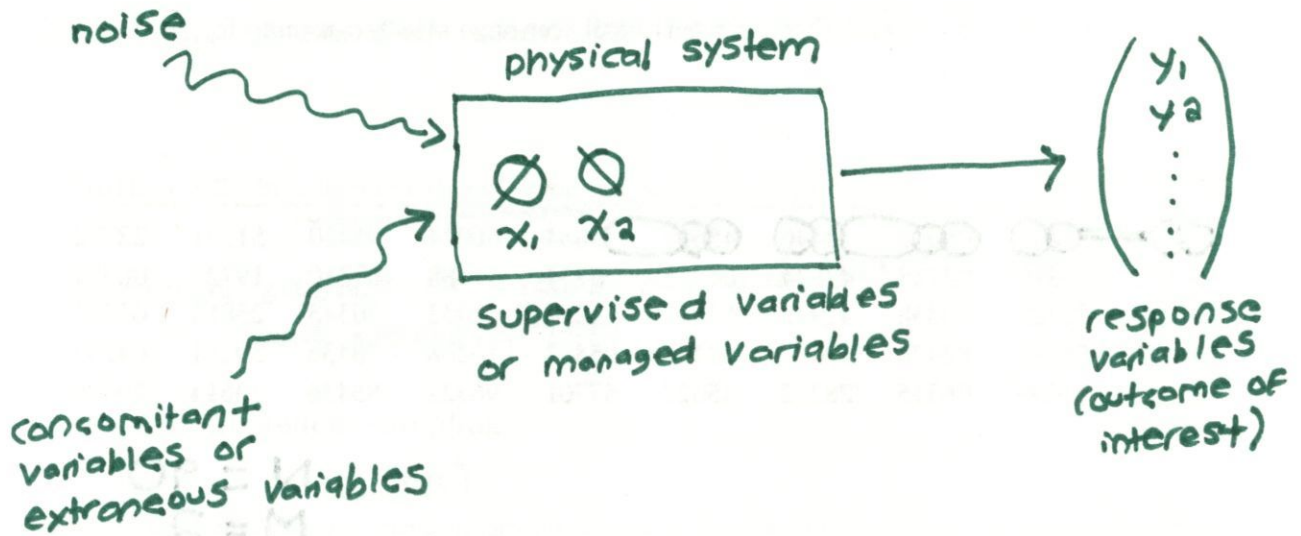
Alternatively: Use a computer.

R → free software, open source
 JMP → free for this class
 Excel → others

2.2 Effective experimentation

Purposefully changing a system and observing what happens as a result is a principled way of learning how a system works.

A typical experimental situation:



Example 2.4 (Chemical purity). Suppose you want to know about the effect of two different reactants (A and B) on the purity of a chemical for a given mixing speed and batch size. Reactant A has 2 levels (a_1 and a_2) and reactant B also has 2 levels (b_1 and b_2).

2.2.1 Taxonomy of variables

Planning an experiment is complicated. There are typically many different characteristics of the system an engineer is interested in improving and many variables that might influence them. Some terminology is needed.

Definition 2.4. A response variable in an experiment is one that is monitored as characterizing system performance/behavior.

Definition 2.5. A supervised (or managed) variable in an experiment is one over which an investigator exercises power, choosing a setting or settings for use in the study. When a supervised variable is held constant (has only one setting), it is called a control variable. When a supervised variable is given several settings in a study, it is called an experimental variable.

Definition 2.6. A concomitant (or accompanying) variable in an experiment is one that is observed but is neither a primary response variable nor a managed variable. Such a variable can change in relation to either experimental or unobserved causes and may or may not itself have an impact on a response variable.

Example 2.5 (Chemical purity, cont'd). What are the response variables, controlled variables, experimental variables, and concomitant variables?