

Workbook: Statistical Thinking for Forensic Practitioners

Sam Tyner, Hal Stern, Alicia Carriquiry

Contents

1	Introduction	5
2	Statistical Preliminaries	7
2.1	Probability	8
2.2	Probability to Statistical Inference	19
2.3	Statistical Inference - Estimation	24
2.4	Statistical Inference - Hypothesis Testing	27
3	Statistics for Forensic Science	29
3.1	Brief Review of Probability and Statistics	29
3.2	The Forensic Examination	29
3.3	Common Approaches to Assessing Forensic Evidence	29
	References	31

Chapter 1

Introduction

This workbook is intended to accompany the Statistical Thinking for Forensic Practitioners workshop taught by members of the Center for Statistics and Applications in Forensic Evidence (CSAFE). The slides for this workshop were constructed initially by Dr. Hal Stern of UC-Irvine and Dr. Alicia Carriquiry of Iowa State University.

When taking the workshop, please follow along with the slides handout (if given) and this workbook. The workbook contains the same material as the slides, with room for you to take notes and to fill in the missing material.

Chapter 2

Statistical Preliminaries

Briefly, this section contains a broad review of probability concepts and of statistical inference concepts, with examples from the forensic science context. We will cover probability, data collection, statistical distributions, estimation, and hypothesis testing.

2.0.1 Definitions

- **population:** _____
- **sample:** _____
- **probability:** Using knowledge about the _____ to make statements describing the _____. Probability can loosely be thought of as a type of deductive reasoning, where we are applying general knowledge about the population of interest to make conclusions about a small part of that population.
- **statistics:** Using knowledge about the _____ to make statements describing the _____. Statistics can loosely be thought of as a type of inductive reasoning, where we are applying knowledge about a sample to state that something *may* be true about the population generally.

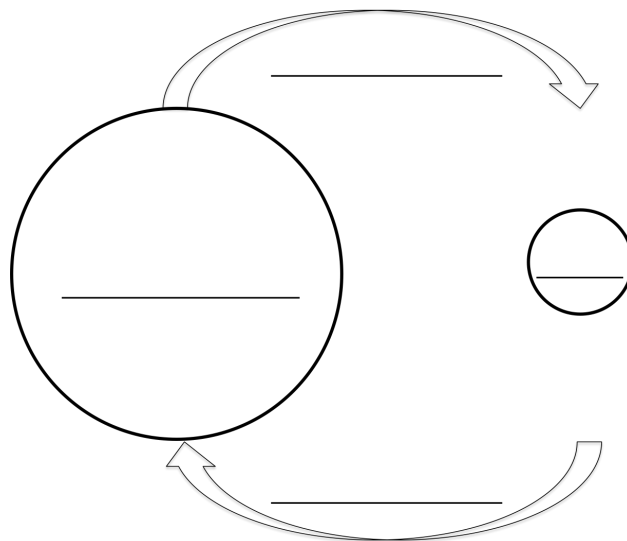


Figure 2.1: "The Big Picture"

2.0.2 Forensic Science Examples

- Suppose 100 1-pound bags of cocaine are seized on the US-Mexico border, and the FBI want to know the chemical composition of the confiscated drugs to store in their database.
 - Population: _____
 - Sample: _____
- A window was broken in a robbery, and the suspect who was apprehended nearby had glass fragments lodged in the soles of their shoes. Do the fragments from the suspect's shoes have the same or similar chemical composition as the broken window?
 - Population 1: _____
 - Sample 1: _____
 - Population 2: _____
 - Sample 2: _____
- A city government employee is suspected of embezzling funds from the city's coffers. Forensic accountants examine a subset of the city's transactions to determine whether embezzling occurred and how much money was lost.
 - Population: _____
 - Sample: _____

How do you think this pertains to pattern evidence? List some possible relevant populations and samples below.

- Population 1: _____
- Sample 1: _____
- Population 2: _____
- Sample 2: _____
- Population 3: _____
- Sample 3: _____

2.1 Probability

Probability concerns the *uncertainty* of outcomes. The set of all possible outcomes is called the _____ space, and a particular outcome or set of outcomes of interest is referred to as an _____.

2.1.1 Examples

1. Footwear

- Sample Space = All shoe sizes e.g. $\{6, 6.5, 7, 7.5, 8, 8.5, \dots\}$
- Event = Shoe of size 9

2. Footwear

- Sample Space = Brand of shoe e.g. { Nike, Vans, Converse, ... }
- Event = Nike sneaker

3. Firearms

- Sample Space = CMS (consecutive matching striae) for a pair of bullets e.g. {0, 1, 2, 3, 4, ... }
- Event = CMS of 10 or more

2.1.2 Interpretation

The probability of observing an event in a sample space is a number less than or equal to 1 and greater than or equal to 0 that describes the _____ that the event will occur.

There are two primary interpretations of probability:

1. The long run _____ of occurrence of an event.
2. The _____ belief of likelihood of an event occurring.

2.1.3 Basic Notation and Laws of Probability

Let an event of interest be denoted by _____. The probability of this event occurring is then denoted _____. Recall that the probability of an event is always between 0 and 1. When $P(Y) = 0$, the event Y will never happen. When $P(Y) = 1$, the event Y will always happen. The sum of the probabilities of all possible outcomes in the sample space always equal to _____.

The event of interest, Y , also has a complement event, \bar{Y} , which is read as “not Y ”. The complement, \bar{Y} , of an event, Y , is itself an event containing all outcomes in the sample space other than that initial event of interest, Y .

$$P(Y) + P(\bar{Y}) = \underline{\hspace{1cm}}$$

The above equation also gives us the following rules:

$$\begin{aligned} P(Y) &= 1 - P(\bar{Y}) \\ P(\bar{Y}) &= 1 - P(Y) \end{aligned} \tag{2.1}$$

2.1.4 Probability and Odds

The probability of an event defines the odds of the event. The odds *in favor* of an event Y are defined as the probability of Y divided by the probability of everything except Y (“not Y ”):

$$O(Y) = \frac{P(Y)}{P(\bar{Y})} = \frac{P(Y)}{1 - \underline{\hspace{1cm}}}.$$

Conversely, the odds *against* a event Y are defined as the probability of everything except Y (“not Y ”) divided by the probability of Y :

$$O(\bar{Y}) = \frac{P(\bar{Y})}{P(Y)} = \frac{1 - \underline{\hspace{1cm}}}{P(Y)}.$$

When we typically talk about odds, like in horse racing, the odds reported are the odds *against* the outcome of interest. Let's construct a horse race scenario using our probability notation to find the probability of a horse winning a race from the reported odds:

- Suppose you want to place a bet on a horse name Cleopatra winning the race. Odds for Cleopatra are reported as 4:1.
- Y = Cleopatra wins the race
- \bar{Y} = Any horse in the race *other than* Cleopatra wins the race.
- $O(\bar{Y}) = \frac{P(\bar{Y})}{P(Y)} = \frac{4}{1} = 4$
- We know that $P(Y) + P(\bar{Y}) = 1$. With this information, we can determine $P(Y)$, which is the probability that Cleopatra wins the race:

$$\begin{aligned}
 O(\bar{Y}) &= \frac{P(\bar{Y})}{P(Y)} = 4 \\
 \Rightarrow \frac{P(\bar{Y})}{P(Y)} &= 4 \\
 \Rightarrow \frac{1 - P(Y)}{P(Y)} &= 4 \quad (\text{See Equation 2.1}) \\
 \Rightarrow \frac{1}{P(Y)} - 1 &= 4 \\
 \Rightarrow \frac{1}{P(Y)} &= 5 \\
 \Rightarrow P(Y) &= \frac{1}{5} = 0.2 \\
 \Rightarrow P(\bar{Y}) &= 0.8
 \end{aligned}$$

- So, the odds for Cleopatra (4:1) mean that Cleopatra has a probability of 0.2 of winning the race. Because this outcome is not very likely (it will only happen in 1 race out of 5), you win money if Cleopatra wins simply because that is not a likely outcome.
- **Betting:** Suppose you bet \$1 on Cleopatra to win the race with 4:1 odds. You will win \$4 if Cleopatra wins, otherwise you've lost \$1.
- The amount you win (\$4) is determined so that you break even in the long run.
- Suppose 5 identical races are run. In 1 of those races, Cleopatra wins, and in the other 4, Cleopatra loses. If you bet \$1 on Cleopatra in each race, you will lose that \$1 4 of 5 times. So, in order for you to break even, the designated amount you'll win when Cleopatra wins is \$4.
- This is a statistical concept known as *expected value*. Your expected value when placing the bet is \$0. We compute expected value by multiplying each possible outcome value by its probability and adding them all together:

$$\begin{aligned}
 \$4 \cdot P(Y) + (-\$1) \cdot P(\bar{Y}) &= 0 \\
 \$4 \cdot 0.2 + (-\$1) \cdot 0.8 &= 0 \\
 \$0.8 - \$0.8 &= 0
 \end{aligned}$$

2.1.5 Probability Math

Up until now, we have only considered one event, Y . Now, suppose we have another event that we are interested in, Z .

Let's consider the possibility of *either* of these two events, Y and Z , occurring. We'd write this as $Y \cup Z$, which is mathematical notation for " Y or Z occurs". There are two scenarios that arise:

1. Y and Z cannot occur together: they are _____
2. Y and Z can occur together.

In scenario #1, computing the probability of either Y or Z happening is easy: we just add their respective probabilities together:

$$Y, Z \text{ mutually exclusive} \Rightarrow P(Y \cup Z) = P(Y) + P(Z)$$

In scenario #2, computing the probability of either Y or Z happening is more complicated because we know there is a chance that Y and Z can happen together. We'd write this as $Y \cap Z$, which is mathematical notation for " Y and Z occurs". In scenario #1, this event never occurred, so $P(Y \cap Z) = 0$ there. To compute the probability of Y or Z occurring in scenario #2, we have to consider the probability of Y , the probability of Z , and the probability of $Y \cap Z$. If we just add $P(Y) + P(Z)$ as in scenario #1, we include the event $Y \cap Z$ twice, so we have to subtract one instance of it:

$$Y, Z \text{ not mutually exclusive} \Rightarrow P(Y \cup Z) = P(Y) + P(Z) - P(Y \cap Z).$$

This probability is much easier to think about when illustrated. In Figure 2.2, we consider human blood types. There are four groups: A, B, O, and AB, and there are two RH types: + and -. We first consider the blood types A and B, represented by the two non-overlapping circles. Define:

- Event Y = a person has blood type A
- Event Z = a person has blood type B
- Event $Y \cup Z$ = a person has blood type A or blood type B

These two events are *mutually exclusive* because one person cannot have both blood type A and blood type B. (The circles don't overlap in the venn diagram) So, the probability that a randomly selected person has blood type A or B is:

$$P(Y \cup Z) = ____ + ____$$

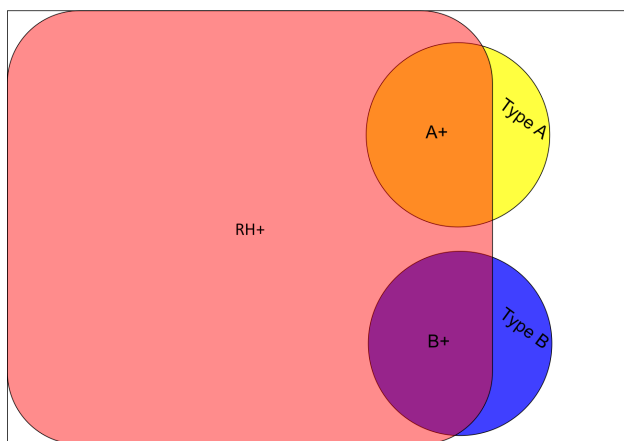


Figure 2.2: Probabilities of blood types in humans

Return to Figure 2.2 and consider two other events: a person having blood type A or having the Rh factor (RH+). We see in Figure 2.2 that someone can have both type A blood and the Rh factor (blood type A+). Define:

- Event Y = a person has blood type A
- Event Z = a person has the Rh factor

- Event $Y \cup Z$ = a person has blood type A or the Rh factor
- Event $Y \cap Z$ = a person has blood type A and the Rh factor (they have A+ blood)

So, the probability that someone has either type A blood or has the Rh factor is the sum of probability of having type A blood (represented by the yellow circle) and the probability of having the Rh factor (represented by the red rectangle) minus the probability of having A+ blood (represented by the orange area of overlap that is counted twice) in Figure 2.2. So, the probability that a randomly selected person has blood type A or the Rh factor is:

$$P(Y \cup Z) = \text{_____} + \text{_____} - \text{_____}$$

2.1.6 Conditional Probability

Let's consider an event of interest Y which has probability $P(Y)$. Then, suppose we learn of another event of interest Z that has occurred. Knowing that Z has occurred already may change our opinion about the likelihood of _____ occurring. The key idea here is that the probability of an event often depends on other information, leading us to the definition of *conditional probability*:

$$P(Y|Z),$$

which is the conditional _____ that Y occurs given that we know Z has occurred. Return to Figure 2.2. Suppose we want to know the probability of a person having type A blood, represented by the yellow circle. But, if we already know that a person has the Rh factor, we are only interested in the part of the type A circle that overlaps with the Rh+ rectangle. Thus the probability of having type A blood is different with different knowledge. The formula for calculating conditional probability is:

$$P(Y|Z) = \frac{P(Y \cap Z)}{P(Z)} \quad (2.2)$$

Returning to the venn diagram, the value $P(Y \cap Z)$ is represented by the overlap of the type A circle and the Rh+ rectangle, and the value $P(Z)$ is represented by the Rh+ rectangle. Then, the value $P(Y|Z)$ is the ratio of the overlap (A+) to the Rh+ rectangle.

Equation 2.2 also gives us a multiplication rule for computing probabilities:

$$P(Y \cap Z) = P(Y|Z) \cdot P(Z) \quad (2.3)$$

Philosophically speaking, it can be helpful to think of *all* probabilities as conditional. It is just a question of what information is assumed to be _____.

2.1.6.1 Examples

Death Penalty Convictions

A study of sentencing of 362 black people convicted of murder in Georgia in the 1980s found that 59 were sentenced to death (Baldus, Pulaski, and Woodworth (1983)). They also examined the race of the murder victim, either black or white, and found some disparities. In Table 2.1, DP means the defendant received the death penalty, NDP means the defendant did not receive the death penalty. The race of the victim (RV) is either black (B) or white (W).

Returning to Section 2.0.1, let's define the problem:

- **Population:** All black people convicted of murder in Georgia in the 1980s
- **Sample:** N/A (the whole population was studied)

RV	DP	NDP	Total
W	45	85	130
B	14	218	232
Total	59	303	362

Table 2.1: The results of the Baldus et al study for black defendants convicted of murder.

Using the numbers from Table 2.1, compute the following probabilities:

- $P(DP) = \text{---} = 0.\text{---}$
- $P(DP|RV = W) = \text{---} = 0.\text{---}$
- $P(DP|RV = B) = \text{---} = 0.\text{---}$

Note: These numbers are selected from the study, and should not be considered a comprehensive summary of its results. There are a number of things not discussed here. The entire publication can be found online¹

Consecutive Matching Striae

In firearms and toolmark analysis, the number of consecutive matching striae (CMS) between a crime scene sample and a lab sample is often used to help determine a match. Generally speaking, the higher the maximum number of CMS found in a pair, the more likely the two samples came from the same source. Several known match (KM) pairs and known non-match (KNM) pairs of bullets were examined, and the results are shown in Figure 2.3 (Hare, Hofmann, and Carriquiry (2017)). What is the probability of seeing two known matches (or two known non-matches) given the maximum number of CMS? Here, we condition on _____. Again, we briefly return to Section 2.0.1, let's define the problem:

- **Population:** All pairs of fired bullets from unknown sources
- **Sample:** A sample of pairs of known matches and known non-matches

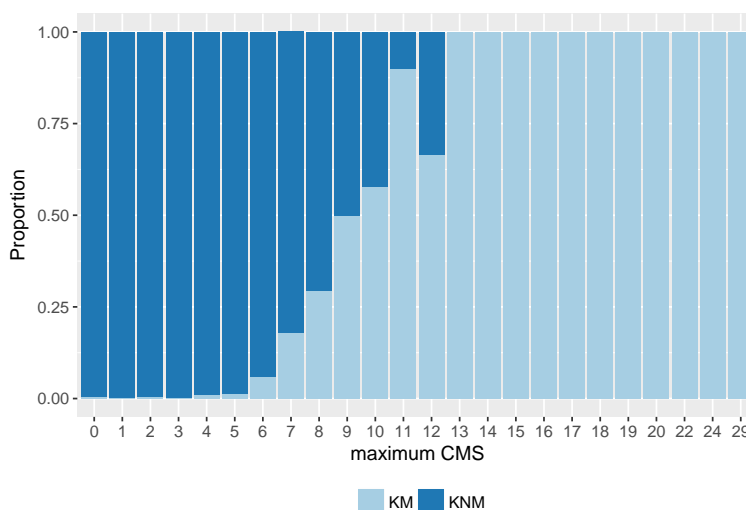


Figure 2.3: This bar chart represents the conditional probabilities of two bullets matching given the maximum number of CMS. The light blue represents known matches, while the dark blue represents known non-matches.

Generally, as seen in Figure 2.3, the probability of finding a match tends to increase with then number of maximum CMS. For _____ maximum CMS values is it much more likely that we have a _____ pair.

¹<http://scholarlycommons.law.northwestern.edu/cgi/viewcontent.cgi?article=6378&context=jelc>.

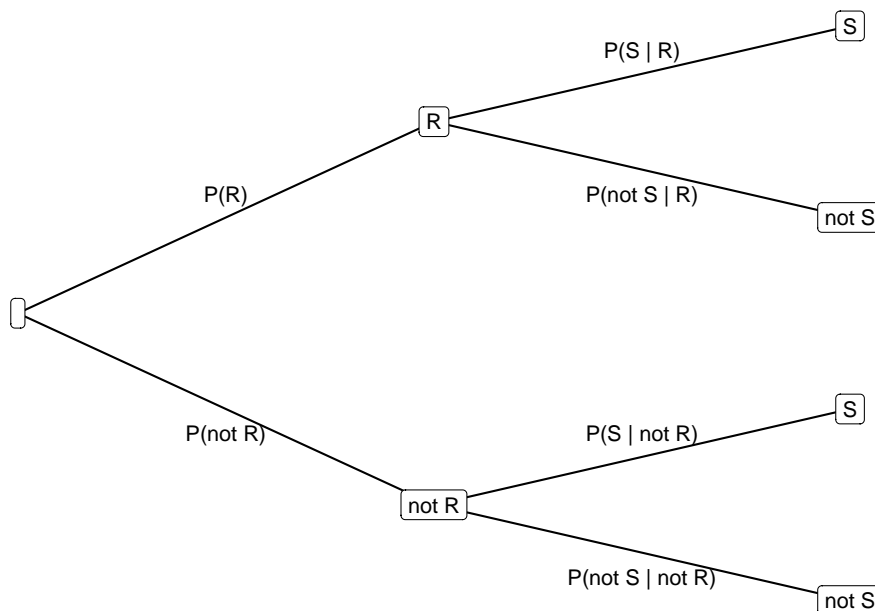


Figure 2.4: A probability tree showing the direction of flow when updating probabilities. Move from left to right on the tree. Events are in boxes.

probability for the events S and \bar{S} . Then, a witness says they saw a tall Caucasian male running from the scene, and the defendant is a tall Caucasian male. After hearing the witness' testimony, the jurors _____ their probabilities. Next, an expert witness testifies that fragments from a window broken during the crime and fragments found on the defendant's clothing match. Again, the jurors update their _____. This process continues throughout the trial. There are some key questions to consider:

- How should jurors update their probabilities?
- Do jurors *actually* think this way?

2.1.9 Bayes' Rule

Bayes' Rule provides an _____ formula for probabilities. Like in the trial scenario above, suppose we have an initial estimate for the probability of event S , $P(S)$. Then, we learn that an event R has occurred and we want to update or probability of event S . To do this, we need to know about the _____ of R and S . To update the probability of S , we can use Bayes' Rule, also called Bayes' _____:

$$\begin{aligned}
 P(S|R) &= \frac{P(R \cap S)}{P(R)} = \frac{P(R|S)P(S)}{P(R)} \\
 &= \frac{P(R|S)P(S)}{P(R|S)P(S) + P(R|\bar{S})P(\bar{S})}
 \end{aligned}
 \tag{2.4}$$

2.1.9.1 Examples

Consider performing diagnostic tests for gunshot residue.

- Let G denote the presence of gunshot residue

- Let \bar{G} denote the _____ of gunshot residue
- Let T denote a _____ diagnostic test
- Let \bar{T} denote a negative diagnostic test

Truth	T	\bar{T}
G	True Positive	False Negative
\bar{G}	False Positive	True Negative

Table 2.2: All potential outcomes of a diagnostic test for gunshot residue.

The values in the table can also be thought of as conditional probabilities:

- The value $P(T|G)$ is the _____ rate, also called *sensitivity* of the test
- The value $P(\bar{T}|\bar{G})$ is the _____ rate, also called the *specificity* of the test
- The value $P(T|\bar{G})$ is the _____ rate, the Type I error rate
- The value $P(\bar{T}|G)$ is the _____ rate, the Type II error rate

Studies of the diagnostics test usually tell us $P(T|G)$, _____, and $P(\bar{T}|\bar{G})$, _____.
 . Examiners may begin with some idea of $P(G)$, or the _____ of gunshot residue in a similar situation. What is most relevant for the case is the *positive predictive value*, or in probability notation, _____. We can use _____ to obtain this value:

$$P(G|T) = \frac{P(T|G)P(G)}{P(T|G)P(G) + P(T|\bar{G})P(\bar{G})}$$

Generally speaking, the most important thing to remember is that, in general, $P(T|G)$ _____ $P(G|T)$.

The careful application of Bayes' Rule can sometimes lead to surprising, non-intuitive results. Continuing with the gunshot residue test example, assume

- sensitivity is 98% ($P(T|G) = 0.98$)
- specificity is 96% ($P(\bar{T}|\bar{G}) = 0.96$)
- prevalence is 90% ($P(G) = 0.90$)
- Plug values into the Bayes' Rule formula to find $P(G|T)$:

$$\begin{aligned}
 P(G|T) &= \frac{P(T|G)P(G)}{P(T|G)P(G) + P(T|\bar{G})P(\bar{G})} \\
 &= \frac{0.98 \cdot 0.9}{0.98 \cdot 0.9 + (1 - 0.96) \cdot (1 - 0.9)} \\
 &= \frac{0.882}{0.882 + 0.004} \\
 &= 0.995
 \end{aligned} \tag{2.5}$$

- Now assume prevalence is 10% ($P(G) = 0.10$) and plug in the values again

$$\begin{aligned}
 P(G|T) &= \frac{P(T|G)P(G)}{P(T|G)P(G) + P(T|\bar{G})P(\bar{G})} \\
 &= \frac{0.98 \cdot 0.1}{0.98 \cdot 0.1 + (1 - 0.96) \cdot (1 - 0.1)} \\
 &= \frac{0.098}{0.098 + 0.036} \\
 &= \frac{0.098}{0.134} \\
 &= 0.731
 \end{aligned} \tag{2.6}$$

- So, even if there is a positive test, we are not really sure about whether gunshot residue is *actually* present.
- Why does this happen?? See Figure 2.5.

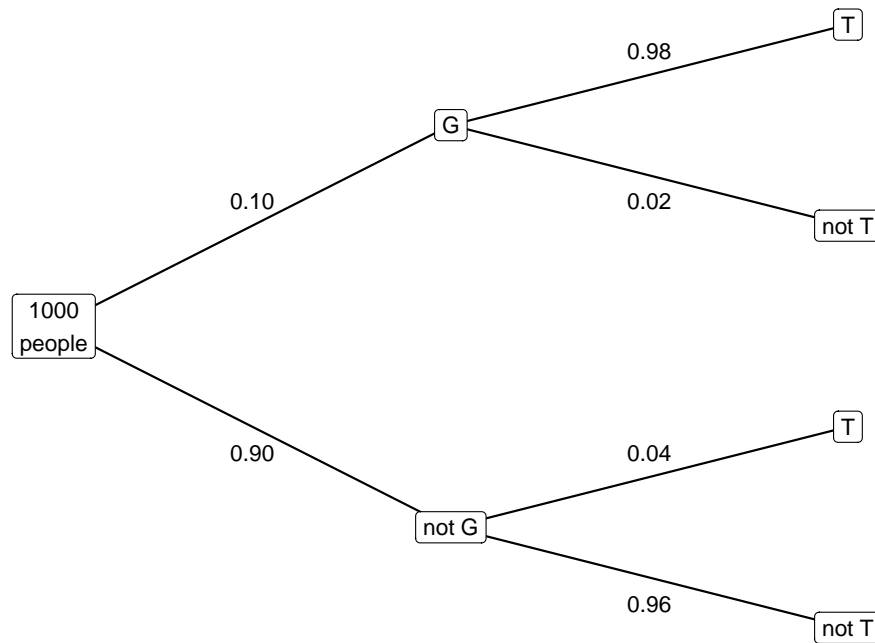


Figure 2.5: A probability tree showing the direction of flow when updating probabilities for the presence of gunshot residue. Suppose there are 1,000 people in the population you’re considering. Write the number of people in the groups throughout the tree according to the probabilities indicated on the branches of the tree

2.1.10 Bayes’ Rule to the Likelihood Ratio

In the general forensic setting, let S denote the event that the evidence from the scene and comparison sample are from the same source. Let E denote the evidence found at the scene. The formulation of Bayes’ Rule for this situation is:

$$P(S|E) = \frac{P(E|S)P(S)}{P(E|S)P(S) + P(E|\bar{S})P(\bar{S})}$$

We can rewrite Bayes’ Rule in terms of odds:

$$\frac{P(S|E)}{P(\bar{S}|E)} = \frac{P(E|S)P(S)}{P(E|\bar{S})P(\bar{S})} \quad (2.7)$$

Derivation of Equation 2.7 is shown in Equation 2.8. For now, just consider Equation 2.7:

- On the left, $\frac{P(S|E)}{P(\bar{S}|E)}$ are the odds in favor of S given the evidence E .
- The last term on the right, $\frac{P(S)}{P(\bar{S})}$ are the odds in favor of S before seeing the evidence E (the “prior odds”)
- The first term on the right $\frac{P(E|S)}{P(E|\bar{S})}$, is known as the _____ ratio
- The likelihood ratio (LR) is the factor by which we _____ prior odds of two samples being from the same source to get _____ odds (after seeing evidence) of the same source.

$$\begin{aligned}
 P(S|E) &= \frac{P(E|S)P(S)}{P(E|S)P(S) + P(E|\bar{S})P(\bar{S})} \\
 \Rightarrow \frac{1}{P(S|E)} &= \frac{P(E|S)P(S) + P(E|\bar{S})P(\bar{S})}{P(E|S)P(S)} \\
 &= 1 + \frac{P(E|\bar{S})P(\bar{S})}{P(E|S)P(S)} \\
 \Rightarrow \frac{1}{P(S|E)} - 1 &= \frac{P(E|\bar{S})P(\bar{S})}{P(E|S)P(S)} \\
 \frac{1}{P(S|E)} - \frac{P(S|E)}{P(S|E)} &= \\
 \frac{1 - P(S|E)}{P(S|E)} &= \\
 \frac{P(\bar{S}|E)}{P(S|E)} &= \frac{P(E|\bar{S})P(\bar{S})}{P(E|S)P(S)} \\
 \Rightarrow \frac{P(S|E)}{P(\bar{S}|E)} &= \frac{P(E|S)P(S)}{P(E|\bar{S})P(\bar{S})}
 \end{aligned} \quad (2.8)$$

2.1.10.1 Examples

Return to the gunshot residue (GSR) test example. Define:

- E = evidence = a positive test for (GSR)
- S = suspect has GSR on them

$$LR = \frac{P(E|S)}{P(E|\bar{S})} = \frac{0.98}{0.04} = 24.5$$

In a high prevalence case ($P(G) = 0.9$), the prior odds are $\frac{0.9}{0.1} = 9$. The posterior odds are $LR \times \text{prior odds} = 24.5 \times 9 = 220.5 : 1$.

In a low prevalence case ($P(G) = 0.1$), the prior odds are $\frac{0.1}{0.9} = \frac{1}{9}$. The posterior odds are $LR \times \text{prior odds} = 24.5 \times \frac{1}{9} = 24.5 : 9 = 2.72 : 1$.

We can also compute the likelihood ratio if the evidence were a negative test. This value turns out to be $\frac{1}{48}$, which is **not** the reciprocal of the LR for the positive test.

2.1.11 Recap

- Probability is the _____ language of _____
- Provides a common scale, from _____ to _____, for describing the chance that an event will occur
- **Conditional** probability is a key concept! The probability of an event depends on what _____ is available
- Independent events can be powerful! They allow us to _____ probabilities of events *directly*, as is common in _____.
- _____ is a mathematical result showing how we should _____ our probabilities when available information changes.
 - This will later lead us to the likelihood ratio as a numerical _____ of the evidence.
 - Bayes' Rule does not necessarily describe how people operate in practice.

2.1.12 Probability and the Courts

Sally Clark was the only person in the house when her first child died unexpectedly at 3 months old. The cause of death was determined to be SIDS, sudden infant death syndrome. One year later, Sally and her husband had a second child, who died at 2 months old under similar circumstances. Sally was convicted of murder.

During her trial, a pediatrician testified that the probability of a single SIDS death for a family like the Clarks (similar income, etc.) was $\frac{1}{8500} \approx 0.0001$, and thus the probability of two SIDS death in the family was $\frac{1}{8500^2} = \frac{1}{73 \times 10^6} \approx 1.37 \times 10^{-8}$. There are several problems with this approach to evidence. What do you think? Jot down a few ideas below:

Issues with the evidence presented by the pediatrician:

1. Is the probability of a child dying of SIDS given, $\frac{1}{8500}$, correct for “families like the Clarks”?
2. The use of direct multiplication of probabilities assumes independence of the two deaths in the family. (Independence within the family is not a reasonable assumption.)
3. Alternative hypotheses (causes of death of the infants) were not considered. Did something else with perhaps a higher likelihood cause the children's deaths?

2.2 Probability to Statistical Inference

Probability is important, but it is only one tool in our toolbox. Another, more powerful tool is statistical inference.

2.2.1 Collecting Data

First, we consider data collection. Where do data come from? One data source is an *experiment*. An investigator designs a study and collects information on and maybe applies treatments to a *sample*, a subset of

the population of interest. _____ can tell us a great deal about how to design an _____ or choose a _____.

The area of statistics concerned with creating studies is called *experimental design*. The experimental design literature is extensive (see for example Morris (2011)). Here are a few crucial points:

- The goal of an experiment is to compare _____
- Those _____ must be _____ assigned to units
- The _____ in the experiment must be large enough to be able to make informed conclusions
- Blinding plays an important role in avoiding _____. e.g. “double-blind” studies in medicine, where neither the patient nor the doctor administering the treatment know which treatment the patient is receiving

How is experimental design relevant to forensic science?

- Experiments are used to evaluate process improvements
- Blinding is used in “black box” studies, where examiners do not know ground truth

Experiments almost always involve *sampling* from the population of interest. Why?

- We sample because it is too _____ or _____ to study the *entire* population
- A _____ sample allows us to use the laws of _____ to describe how certain we are that our _____ answer reflects the _____.
- There are many famous failures (cautionary tales) with _____ sampling. (See Figure 2.6.)

How is sampling relevant to forensic science?

- Sampling techniques used to determine which and how many bags of suspect powder collected from a crime scene to test.

All data collected can be divided into one of two groups: qualitative or quantitative.

- **Qualitative** data describe qualities about the observations. For example, the race of a suspect, or their level of education. There are two subcategories of qualitative data:
 - _____: the data belong to one of a discrete number of groups or categories. For example: blood type (A, B, AB, or O)
 - _____: the data belong to one group in a set of ordered values. For example, the evaluation of a teacher (poor, average, excellent). The categories have an inherent ordering, unlike in categorical data.
- **Quantitative** data describe quantities that can be measured on the observations. These are numerical observations. There are also two subcategories of quantitative data:
 - _____: the values are distinct or separate. An easy-to-understand example is integer observations: $\{0, 1, 2, 3, 4, \dots\}$. A forensic science example is consecutive matching striations on bullets or toolmarks. (See Figure 2.3)
 - _____: the values can take on any value in a finite or infinite interval. Continuous values fall anywhere on the number line. A forensic science example is the refractive index of a glass fragment.

2.2.2 Probability Distributions

Suppose we are to collect data on some characteristic for a sample of individuals or objects (e.g. weight, trace element concentration). A probability _____ is used to describe these possible values and how



Figure 2.6: This picture from the US presidential election of 1948 shows President Harry Truman, who won the election, holding a newspaper that went to print with the headline "Dewey Defeats Truman!" The headline was based on biased sampling that favored typically Republican demographics. Image Source: <https://blogs.loc.gov/loc/2012/11/stop-the-presses/>

_____ each value is to occur. There are many, many possible probability distributions, but some of the most common are the Binomial, Poisson, Normal, and Lognormal distributions. The probabilities associated with each of these distributions and their possible outcomes are plotted in Figures

- Discrete distributions

- _____: counts the number of _____ in a fixed number (n) of _____. Possible values are $\{0, 1, 2, \dots, n\}$.
- _____: counts the number of _____ occurring. Possible values are $\{0, 1, 2, 3, 4, 5, \dots\}$

- Continuous distributions

- _____: the famous, symmetric “bell-shaped” _____. Possible values are all real numbers, $(-\infty, \infty)$
- _____: the (natural) logarithm of observations from this distribution follow a _____ distribution. Possible values are all positive real numbers, $(0, \infty)$

2.2.2.1 Normal

You may already be familiar with this distribution with the bell-shaped curve. Measurement error is one example of something often assumed to follow a normal distribution. The normal distribution is described by two parameters: the _____, denoted by μ , and the _____, denoted by σ . If we have a variable (say, something observed in our data like weight), we give the variable a capital letter, typically X . If this variable X is normally distributed with mean μ and standard deviation σ , we write this as:

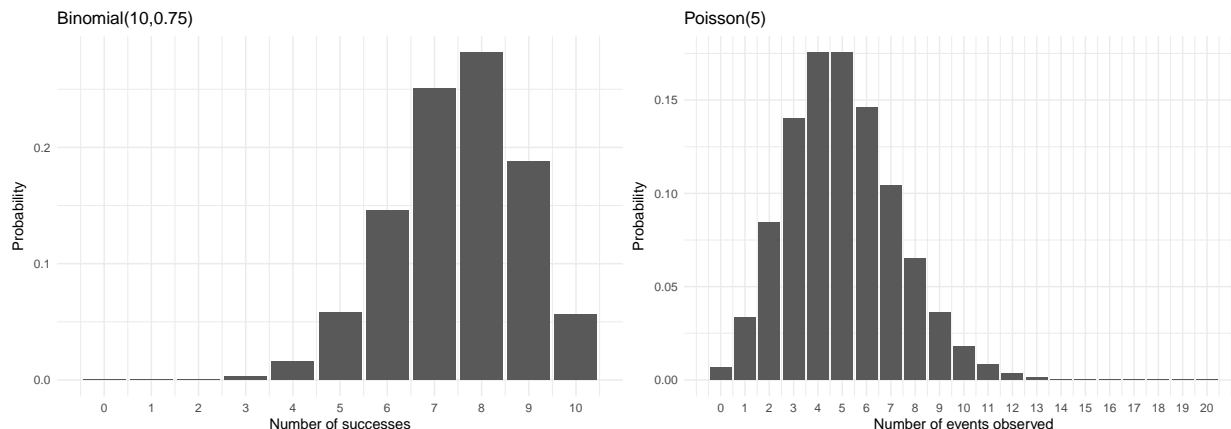


Figure 2.7: On the left, the probability of each possible outcome for variable with binomial distribution with 10 trials and probability of success 0.75. On the right, the probability of each possible outcome for variable with Poisson distribution with mean value 5.

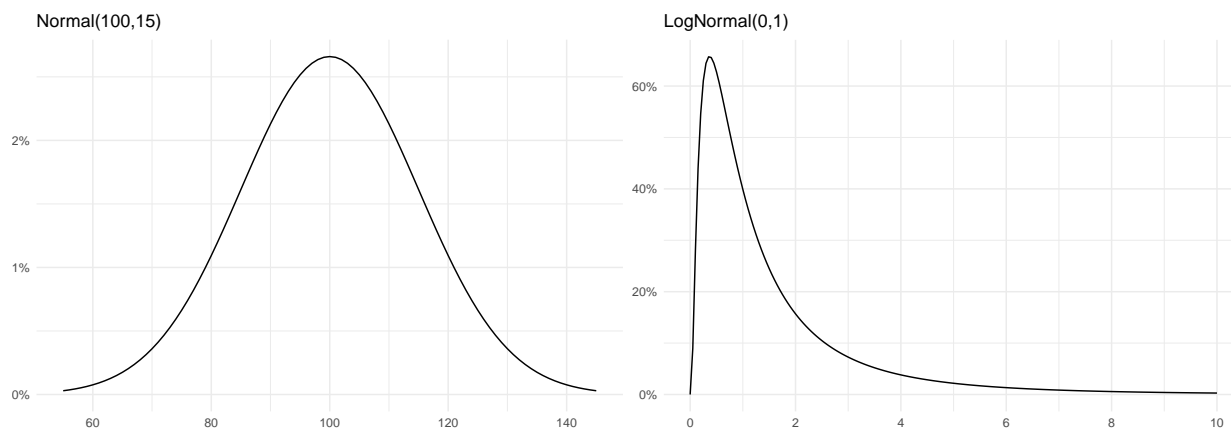


Figure 2.8: On the left, the probability distribution curve of possible outcomes for a variable with Normal distribution with mean value 100 and standard deviation 15. On the right, the probability distribution curve of possible outcomes for a variable with Lognormal distribution with mean value 0 and standard deviation 1.

$$X \sim N(\text{____}, \text{____})$$

In measurement error, for example, we typically assume that the mean is 0. So, if X represents measurement error, we'd write $X \sim N(0, \sigma)$.

There are many nice properties of the normal distribution. For instance, we know that _____% of observable values lie within _____ standard deviations of the mean ($\mu \pm 2\sigma$), and also that _____% of observable values lie within _____ standard deviations of the mean ($\mu \pm 3\sigma$). When working with the normal distribution, we use software (such as Excel, Matlab, R, SAS, etc.), tables², or websites like onlinestatbook.com or stattrek.com to compute probabilities of events.

2.2.2.2 Lognormal

We often act as if everything is normally distributed, but of course this is not true. For instance, a quantity that is certain to be _____ (greater than or equal to zero) cannot possibly be normally distributed.

²See for example <http://www.stat.ufl.edu/~athienit/Tables/Ztable.pdf>

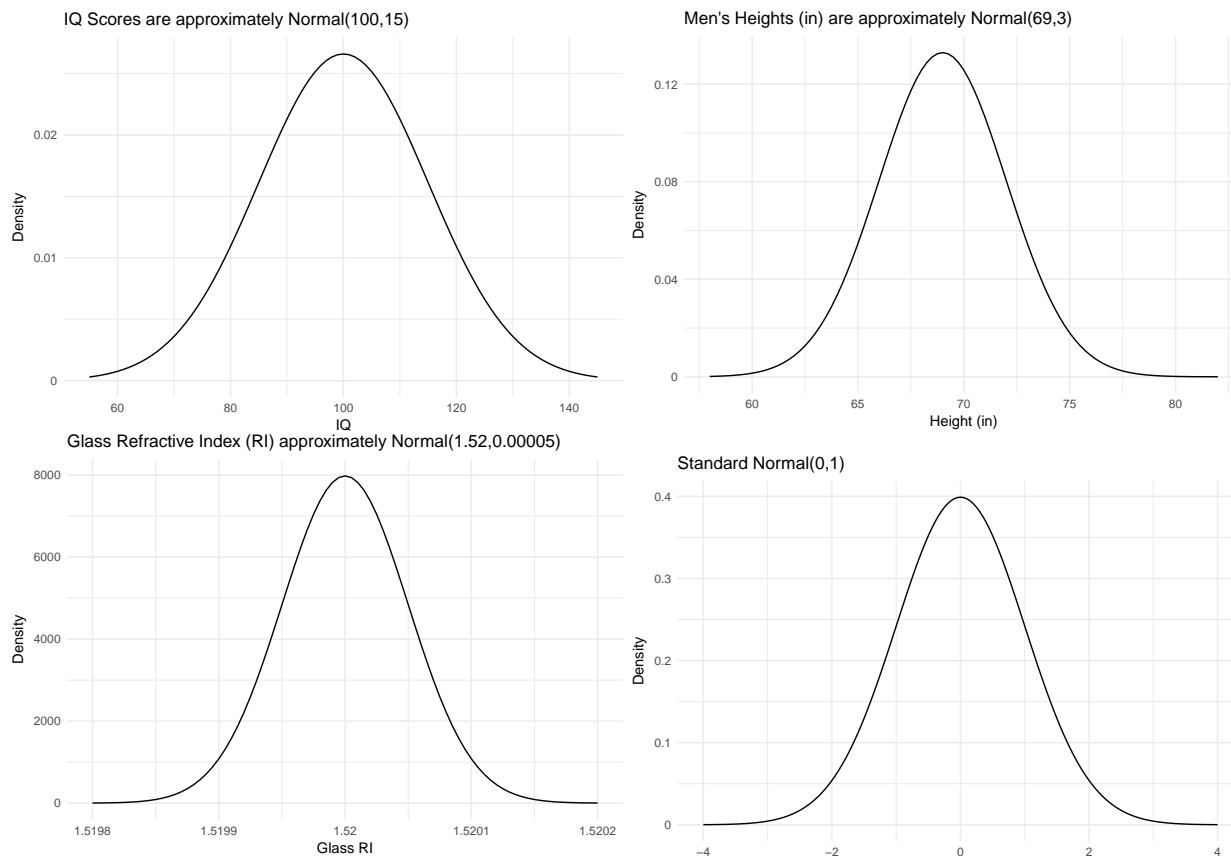


Figure 2.9: Four examples of the probability distribution functions for normally distributed variables

Consider trace element concentration: either none is detected, or there is some amount greater than 0 detected.

In cases where nonnegative values are not possible, we may believe that the (natural) _____ of the quantity is normal, which gives us a _____ distribution for the quantity itself. The lognormal distribution, like the normal, has two parameters: mean (on the log scale), denoted _____, and standard deviation (on the log scale), denoted _____.

2.2.2.3 Discrete

Some quantities take on very few possible values. These are *discrete* data.

Recall the two common discrete distributions from section 2.2.2:

- Binomial:
 - Data are _____ (two categories: “success” or “failure”)
 - Data are a result of n independent _____
 - $P(\text{success}) = p$ on each trial. (Same _____ of success each time)
 - Expected number of successes you expect to see out of n trials: _____ \times _____
 - Example: Suspect a student of cheating on an exam, response is the number of correct answers.
- Poisson:
 - Data are counts: number of events occurring in a _____ time

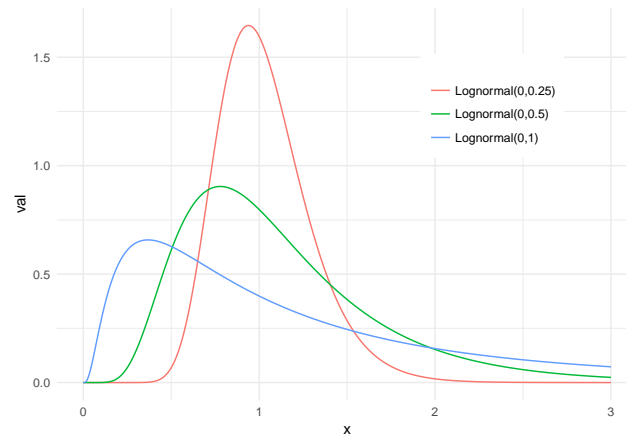


Figure 2.10: Three lognormal distributions with the same mean (on the log scale) and different standard deviations (on the log scale)

- The mean and the _____ of this distribution are the same, so the variability in responses increases as the _____ increases.
- Example: number of calls to 911 between 10:00 and midnight on Friday nights. See Figure 2.11 for a forensics example.

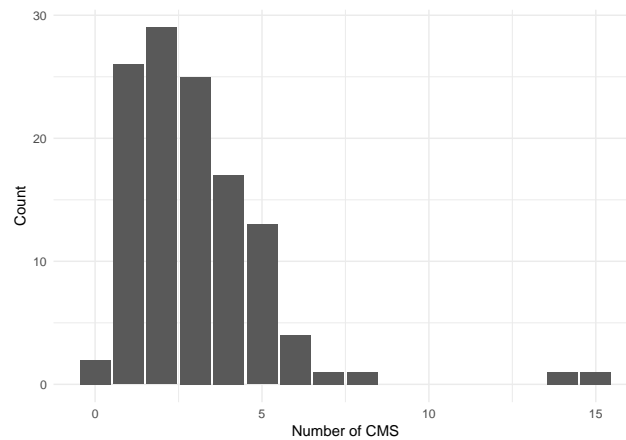


Figure 2.11: Distribution of the maximum number of CMS for a randomly selected bullet compared to 118 known lands approximately follows a Poisson distribution.

2.3 Statistical Inference - Estimation

Recall from Section 2.0.1:

- The _____ is the universe of objects of interest.
- The _____ is comprised of the objects available for study.
- _____ is deductive: use knowledge about the population to make statements describing the sample
- _____ is inductive: use knowledge about the sample to make statements describing the population

- Probability and statistics are used together!
 1. Build or assume a _____ for a population
 2. Assess the _____ using the model
 3. Refine the model, return to step 2.

2.3.1 Background

A _____ is a numerical characteristic of the population, e.g. the population mean. Statistical methods are usually concerned with learning about population parameters from _____.

Note: The mean of a *sample* and the mean of a *population* are different concepts. The mean of a sample can be calculated exactly, while the mean of a population is (usually) unknown, because there are too many objects in the population to record and calculate the mean.

The idea underlying statistical inference is that we can apply laws of probability to draw _____ about a population from a sample. This process is briefly summarized below:

- Observe _____ mean
- If we have a “good” sample this sample mean should be close to the _____ mean.
- The laws of _____ tell us how close we can expect them to be.

For example, suppose we are interested in the average height of the adult population in the U.S.

- Population: _____
- Sample: _____
- We can take the average height of everyone here and use this sample _____ to make _____ about the _____ mean of all U.S. adults.
- Note: This approach will work if our sample is a _____ sample from the population. This assumption may be questionable, so it should be verified.

The *goal* of statistical inference is _____ about a _____. Different possible parameters are:

- Mean
- Variance
- Proportion

We can also make different types of inferential statements, depending on what question we are trying to answer and how we are going to report our results. We will talk about:

- _____ estimate: an estimate of a parameter value
- _____ estimate: a range of plausible values for a parameter
- Hypothesis _____: examine a specific hypothesis about the true value of a parameter

When you want to do statistical inference, it is always important to look at your sample data before proceeding directly to inference. We do this because we want to

1. See general _____ in the data
2. Get an idea of the _____ of the distribution of the data
3. Identify _____ values and/or errors.

How we look at our data to check for these three things? If our data are _____, we look at a table of frequencies or a bar chart of the different outcomes. If our data are _____, we look at histograms of the values, or numerical summaries such as mean, median, standard deviation, or percentiles.

A quick example shows why it can be important to examine your data before a formal statistical analysis:

- Suppose the data are (19,20,21,22,23). Then, the mean is $\frac{19+20+21+22+23}{5} = 21$, the median is 21, and the standard deviation is 1.58.
- But what if the data you receive are (19,20,21,22,93)? Then, the mean is $\frac{19+20+21+22+93}{5} = 35$, the median is 21, and the standard deviation is 32.4.
- There could have been a typo, or someone interpreted some handwriting wrong, etc. The moral of the story is ALWAYS look at your data first!

2.3.2 Point Estimation

An _____ is a rule for estimating a population _____ from a sample. We evaluate the quality of the estimator by considering two key properties:

- Bias: how close _____ an estimator is to the true population mean
- Variability: how _____ is the estimate?

For the population mean, we might use sample mean as an estimator because it has _____ bias and _____ variability is the sample is _____. There are other possible estimators for the mean such as:

- The median (good for skewed data or data with outliers)
- The midrange ($\frac{\max + \min}{2}$)
- 47 (obviously this is just guessing and is not advised)

Let θ denote an unknown population parameter that we wish to estimate. The letter θ represents the true value of the parameter. In Figure 2.12, we see what would happen in many repeated attempts to estimate θ using estimators with different properties.

2.3.3 Standard Errors

One limitation of just providing a point estimate is that it doesn't give us any indication of _____. As we saw in Figure 2.12, a point estimate alone can be very different from the true mean. We can do better than this!

The _____ of an estimator measures the uncertainty in our estimate. When looking at a summary statistic, like mean, median, or percentiles, that statistic is also a _____ quantity. This means that if we had observed a different set of sample values, we would observe different values of the summary statistics. The idea of standard error is similar to the idea of standard deviation. Both are measures of spread, or variability. The difference is that standard deviation is a measure of variability of a sample or population, while standard error is a measure of variability of an _____.

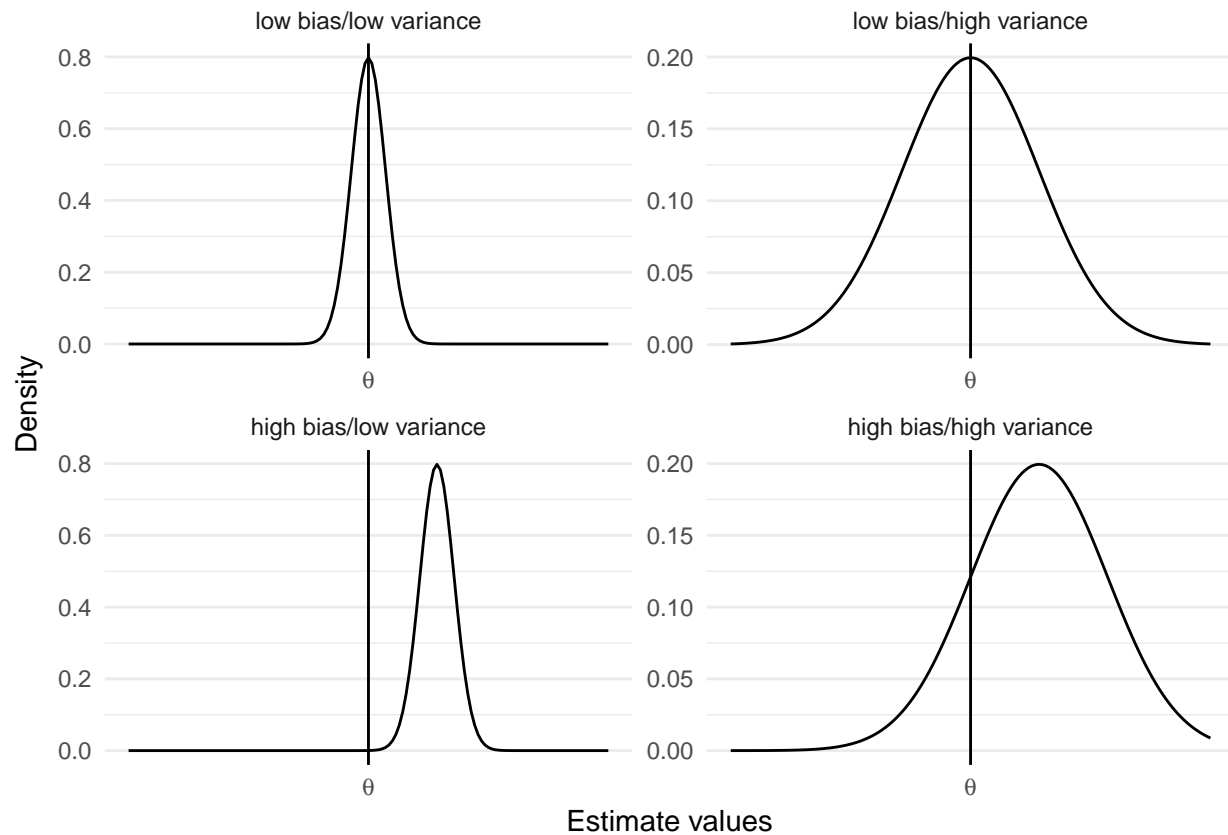


Figure 2.12: The curve in each plot shows the distribution of estimates we would see under each condition shown.

2.3.4 Sample Size

2.3.5 Interval Estimation

2.4 Statistical Inference - Hypothesis Testing

2.4.1 Background

2.4.2 Normal Data

2.4.3 Confidence Intervals

2.4.4 Comparing Two Means

2.4.5 Discussion

Chapter 3

Statistics for Forensic Science

3.1 Brief Review of Probability and Statistics

3.2 The Forensic Examination

3.3 Common Approaches to Assessing Forensic Evidence

3.3.1 Significance Testing / Coincidence Probability

3.3.2 Likelihood Ratio

References

- Baldus, David C., Charles Pulaski, and George Woodworth. 1983. “Comparative Review of Death Sentences: An Empirical Study of the Georgia Experience.” *Journal of Criminal Law and Criminology*.
- Hare, Eric, Heike Hofmann, and Alicia Carriquiry. 2017. “Automatic Matching of Bullet Land Impressions.” *The Annals of Applied Statistics* Upcoming.
- Morris, Max D. 2011. *Design of Experiments: An Introduction Based on Linear Models*. Chapman; Hall. 1968.