

Lecture Notes for Math Statistics

Vladislav Kargin

June 16, 2022

Contents

1	Basics of Statistical Estimation	3
1.1	Probability Reminder	3
1.2	An estimator, its bias and variance	6
1.3	Some common unbiased estimators	12
1.3.1	An estimator for the population mean μ	12
1.3.2	An estimator for the population proportion p	13
1.3.3	An estimator for the difference in population means $\mu_1 - \mu_2$	13
1.3.4	An estimator for the difference in population proportions $p_1 - p_2$	14
1.3.5	An estimator for the variance	15
1.4	The error of estimation and the 2-standard-error bound	16
1.5	Confidence intervals and pivotal quantities	20
1.6	Large-sample confidence intervals	25
1.7	How to determine the sample size	28
1.8	Small-sample confidence intervals	31
1.8.1	Small sample CIs for μ and $\mu_1 - \mu_2$	31
1.8.2	Small sample CIs for population variance σ^2	40
2	Properties of Point Estimators and Methods of Estimation	43
2.1	Relative efficiency.	43
2.2	Consistency	45
2.3	Sufficient statistics	53

2.4	Rao-Blackwell Theorem and Minimum-Variance Unbiased Estimator	58
2.5	Method of Moments Estimation	61
2.6	Maximum Likelihood Estimation (MLE).	68
2.7	Cramer-Rao Lower Bound and large sample properties of MLE	82
3	Hypothesis testing	88
3.1	Basic definitions	88
3.2	Calculating the Level and Power of a Test	95
3.2.1	Additional examples.	105
3.3	Determining the sample size	107
3.4	Relation with confidence intervals	108
3.5	p -values.	109
3.6	Small-sample hypothesis tests for population means . . .	115
3.7	Hypothesis testing for population variances	119
3.8	Neyman - Pearson Lemma and Uniformly Most Powerful Tests 124	
3.9	Likelihood ratio test.	129
3.9.1	An Additional Example	135
4	Bayesian Inference	138
4.1	Estimation	138
4.2	Hypothesis testing	141

Chapter 1

Basics of Statistical Estimation

1.1 Probability Reminder

Statistics consider a sample of data which depends on an unknown parameter. For example we can have n observations for the time before a smart-phone breaks.

$$X_1, X_2, \dots, X_n.$$

These are random variables. Data in a particular sample are usually denoted by lowercase letters:

$$x_1, x_2, \dots, x_n.$$

Since every datapoint is a random variable, it has a probability distribution.

We assume for simplicity that every datapoint has the same distribution as others and that they are independent of each other.

For example, we can model the lifetime of an iPhone as an exponential random variable with mean θ . That means that the density of X_1 is

$$f_{X_1}(x_1) = \frac{1}{\theta} e^{x_1/\theta},$$

the density of X_2 is

$$f_{X_2}(x_2) = \frac{1}{\theta} e^{x_2/\theta},$$

and so on. For simplicity, we usually say that the density of a datapoint is

$$f_X(x) = \frac{1}{\theta} e^{x/\theta}.$$

This is OK since we assumed that all data points have the same distribution.

(Note that the density is useful for continuous random variables. For discrete random variables, we use probability mass function.)

This allows us to compute the probabilities for individual random variables X_i , their expectations, variances and so on. [Examples]

If we want to calculate probabilities or expectation of functions that depends on several random variables, we need a joint distribution function of data - points. Since we assume that they are independent, this is easy. The joint density of independent datapoints is simply product of the individual densities for each datapoint. In our example,

$$\begin{aligned} f_{X_1, \dots, X_n}(x_1, \dots, x_n) &= \frac{1}{\theta} e^{x_1/\theta} \times \frac{1}{\theta} e^{x_2/\theta} \times \dots \times \frac{1}{\theta} e^{x_n/\theta} \\ &= \frac{1}{\theta^n} e^{(\sum_{i=1}^n x_i)/\theta} \end{aligned}$$

In statistics, if we think about this joint density as a function of the model parameter θ , we call it the *likelihood function* and denote it by letter L . So, in our example, we have

$$L(\theta|\vec{x}) = \frac{1}{\theta^n} e^{(\sum_{i=1}^n x_i)/\theta},$$

where we used notation \vec{x} to denote the vector of observed datapoints: $\vec{x} = (x_1, \dots, x_n)$.

Theoretically, if we know the joint density of random variables X_1, \dots, X_n , we can calculate the density of any function of these data. (A function of the data is called a *statistic*.) We have spent some time doing examples of these calculations in Math 447. However, often these calculations are difficult or

impossible. For example, even if we are interested in the simpler average of data-points:

$$\overline{X} := \frac{X_1 + \dots + X_n}{n} = \frac{1}{n} \sum_{i=1}^n X_i,$$

often it is difficult to obtain its exact distribution. For example, for exponentially distributed random variables, it is possible but non-trivial, and the result is a Gamma distribution.

One exception is the maximum or minimum function of data-points. For them, it is easy to calculate the cumulative distribution function. However, we will recall this calculation later.

The average \overline{X} is often very natural guess for the parameter value, and it is quite unfortunate that we sometimes are not able to find its exact distribution. Luckily, for many purposes, it is enough to know the expectation and the variance of the average. This is much easier because of two important properties. First, the expectation is linear:

$$\mathbb{E}(c_1 X_1 + \dots + c_n X_n) = c_1 \mathbb{E}X_1 + \dots + c_n \mathbb{E}X_n.$$

Second, if the random variables X_1, \dots, X_n are independent, then

$$\text{Var}(c_1 X_1 + \dots + c_n X_n) = c_1^2 \text{Var}(X_1) + \dots + c_n^2 \text{Var}(X_n).$$

Example 1.1.1. What are the expectation and variance of \overline{X} ?

Other great results from probability theory are the Chebyshev Inequality, Law of Large Numbers, and the Central Limit Theorem.

The Chebyshev inequality helps us to estimate the probabilities even if we only know the expectation and the variance of a random variable. It says that for every random variable Y with mean $\mu = \mathbb{E}Y$ and variance $\sigma^2 = \text{Var}Y$, we have

$$\mathbb{P}(|Y - \mu| > k\sigma) \leq \frac{1}{k^2},$$

or we can also write it as

$$\mathbb{P}(|Y - \mu| > \varepsilon) \leq \frac{\sigma^2}{\varepsilon^2}.$$

The Law of Large Numbers is an immediate consequences of the Chebyshev's inequality applied to random variable $Y = \bar{X}$. It says that for every $\varepsilon > 0$.

$$\mathbb{P}(|\bar{X} - \mu_X| > \varepsilon) \rightarrow 0,$$

where μ_X is the expectation of X . In other words, as the sample goes large, the probability that the average of X_i will be different from the expectation of X_i by more than ε converges to zero. The theorem follows by noting that the probability on the left-hand side can be estimated from above by $\text{Var}(\bar{X})/\varepsilon^2 = \text{Var}(X)/(n\varepsilon^2)$, which goes to zero as n grows to infinity.

In fact, as the sample becomes large, we can approximate the distribution of \bar{X} by the normal distribution with appropriate mean and variance. This result is known as the Central Limit Theorem. This approximation is usually much better than the rough estimates by the Chebyshev inequality. Namely, the Central Limit Theorem says that

$$\mathbb{P}\left[\frac{\bar{X} - \mu_X}{\sigma_X/\sqrt{n}} < t\right] = \mathbb{P}[Z < t] = \int_{-\infty}^t \frac{1}{\sqrt{2\pi}} e^{-x^2/2} dx.$$

where Z is a random variable with the standard normal distribution. Note that the probability on the left can also be written as

$$\mathbb{P}\left[\bar{X} < \mu_X + t\sigma_X/\sqrt{n}\right],$$

so the Central Limit Theorem allows us to approximate the distribution of \bar{X} without actually calculating exactly what is the density of this distribution. This is why it is an extremely useful theorem.

1.2 An estimator, its bias and variance

One of the main goals in statistics is to guess the value of an unknown parameter θ , given the realization of the data sample. Namely, we are given the realization of random variables X_1, \dots, X_n , and we want to have a rule how to guess θ . Mathematically, this means that we look for a function of the X_1, \dots, X_n , $f(X_1, \dots, X_n)$, which we call an *estimator* and which for

most of the sample realizations (x_1, \dots, x_n) would be a good guess for a parameter θ . (Any function of values in the data sample is called a *statistic*, so an estimator is a statistic which is meant to be a good guess for the true value of the parameter.)

Note on notation: If θ is a parameter to be estimated, then $\hat{\theta}$ denotes its estimator or a value of the estimator in a given sample. So $\hat{\theta}$ is a shortcut notation for a function of the data: $\hat{\theta} = \hat{\theta}(X_1, \dots, X_n)$, or for its value in a current sample $\hat{\theta} = \hat{\theta}(x_1, \dots, x_n)$.

Examples of estimators: $\hat{\theta} = \bar{X} := (X_1 + \dots + X_n)/n$ or $\hat{\theta} = X_{(n)} := \max(X_1, \dots, X_n)$. Even very unnatural functions such as $\sin(X_1 \times X_2 \times \dots \times X_n)$ can be thought as estimators. So how do we distinguish between good and bad estimators?

What do we mean by saying that $\hat{\theta}$ is a good guess for θ ?

Is $\hat{\theta} = \hat{\theta}(X_1, \dots, X_n)$ random? Is its distribution depends on the true value of the parameter θ ? What do we want from the distribution of the estimator? What do we want from the expected value and the variance of the estimator $\hat{\theta}$?

Bias of an estimator

Def: $Bias(\hat{\theta}) = \mathbb{E}\hat{\theta} - \theta$; (The bias of an estimator is its expected value minus the true value of the parameter).

Note that the bias can depend on the true value of the parameter. A good estimator should have $|Bias(\hat{\theta})|$ small, or even 0 for all values of the true parameter.

Def: when $Bias(\hat{\theta}) = 0$, we say that the estimator $\hat{\theta}$ is unbiased.

Example 1.2.1. What is the bias of \bar{X} , X_1 for our previous example about the lifetime of smartphones.

Why \bar{X} is better than X_1 as an estimator?

Variance of an estimator

Def: $\text{Var}(\hat{\theta}) = \mathbb{E}(\hat{\theta} - \mathbb{E}\hat{\theta})^2 = \mathbb{E}\hat{\theta}^2 - (\mathbb{E}\hat{\theta})^2$;

We want that $\text{Var}(\hat{\theta})$ be small for all values of the true parameter θ .

(Note: in Machine Learning the concepts of “bias” and “variance” are related to the complexity of the statistical model which is fit to the data. A model with small complexity has large bias but small variance in its

predictions and the model with larger complexity has small bias but large variance. We do not discuss it here.)

These are two natural requirements. Sometimes we value unbiasedness more than anything else. We want to make sure that an estimator is unbiased and only after this condition is ensured we look for estimators with low variance. However, sometimes we can tolerate that an estimator is a bit biased and look for estimators that have small bias and small variance. In this case, it is useful to define one measure of the quality of an estimator.

Mean Square Error (MSE) of an estimator

Define the Mean Squared Error as

$$MSE(\hat{\theta}) = \mathbb{E} \left\{ (\hat{\theta} - \theta)^2 \right\}.$$

It can be shown that

$$MSE(\hat{\theta}) = \text{Var}(\hat{\theta}) + [\text{Bias}(\hat{\theta})]^2$$

Proof of the decomposition.

•

$$\begin{aligned} \mathbb{E} \left\{ (\hat{\theta} - \theta)^2 \right\} &= \mathbb{E} \left\{ (\hat{\theta} - \mathbb{E}\hat{\theta} + \mathbb{E}\hat{\theta} - \theta)^2 \right\} \\ &= \text{Var}(\hat{\theta}) + [\text{Bias}(\hat{\theta})]^2 + 2\mathbb{E}\{(\hat{\theta} - \mathbb{E}\hat{\theta})(\mathbb{E}\hat{\theta} - \theta)\} \end{aligned}$$

• But

$$\mathbb{E}\{(\hat{\theta} - \mathbb{E}\hat{\theta})(\mathbb{E}\hat{\theta} - \theta)\} = (\mathbb{E}\hat{\theta} - \theta)\mathbb{E}\{(\hat{\theta} - \mathbb{E}\hat{\theta})\} = 0.$$

Example 1.2.2. Suppose that for a certain estimator $\hat{\theta}$ of θ we know that $\mathbb{E}\hat{\theta} = a\theta + b$ for some constant $a \neq 0$ and $b \neq 0$.

- What is $\text{Bias}(\hat{\theta})$, in terms of a , b and θ ?
- Find a function of $\hat{\theta}$ that is an unbiased estimator for θ .

Comment: if you find an biased estimator $\hat{\theta}$, you can sometimes easily correct the bias to get an unbiased estimator. **However, e.g. if $\mathbb{E}\hat{\theta} = \sqrt{\theta}$, then the estimator $\tilde{\theta} = \hat{\theta}^2$ is not unbiased for θ !**

$$\mathbb{E}\{\hat{\theta}^2\} = (\mathbb{E}\hat{\theta})^2 + \text{Var}(\hat{\theta}) = \theta + \text{Var}(\hat{\theta})$$

In fact, it is often quite difficult to find an unbiased estimator.

[Quiz]

Example 1.2.3. The reading on a voltage meter connected to a test circuit is uniformly distributed over the interval $(\theta, \theta + 1)$, where θ is the true but unknown voltage of the circuit. Suppose that Y_1, Y_2, \dots, Y_n denote a random sample of such readings.

- Calculate the bias of \bar{Y} as an estimator of θ .
- Find an unbiased estimator of θ (based on \bar{Y}).
- Find $MSE(\bar{Y})$.

Solution. It is straightforward to calculate the bias:

$$\begin{aligned} bias(\bar{Y}) &= \mathbb{E}(\bar{Y}) - \theta = \\ &= \frac{1}{n} \sum_{i=1}^n \mathbb{E}(Y_i) - \theta \\ &= \mathbb{E}(Y) - \theta = 1/2. \end{aligned}$$

Then $MSE(\bar{Y}) = bias(\bar{Y})^2 + \text{Var}(\bar{Y})$, and since Y_i independent,

$$\text{Var}(\bar{Y}) = \frac{1}{n^2} \sum_{i=1}^n \text{Var}(Y_i) = \frac{1}{n} \text{Var}(Y).$$

In order to calculate the variance of Y , which is uniform on $[\theta, \theta + 1]$, note that the variance of the shifted variable is the same, so we can calculate variance of X which is uniform on $[-1/2, 1/2]$,

$$\text{Var}(X) = \int_{-1/2}^{1/2} x^2 dx = \left| \frac{x^3}{3} \right|_{-1/2}^{1/2} = \frac{1}{12}.$$

Altogether $MSE(\bar{Y}) = \frac{1}{4} + \frac{1}{12n}$.

[Quiz] [Quiz]

Example 1.2.4 (Reminder about the distribution of the minimum and the minimum). Recall the notation $Y_{(1)} = \min\{Y_1 \dots Y_n\}$. Then for the CDF of

$Y_{(1)}$, we have:

$$\begin{aligned} F_{Y_{(1)}}(y) &\equiv \Pr(Y_{(1)} \leq y) = 1 - \Pr(Y_{(1)} > y) \\ &= 1 - \Pr(\text{all } Y_i \text{'s are } > y) \\ &= 1 - [1 - F(y)]^n \end{aligned}$$

And the PDF is

$$f_{Y_{(1)}}(y) = n[1 - F(y)]^{n-1}f(y).$$

Similarly for the maximum we have notation $Y_{(n)} = \max\{Y_1 \dots Y_n\}$. The CDF is

$$\begin{aligned} F_{Y_{(n)}}(y) &\equiv \Pr(Y_{(n)} \leq y) = \Pr(Y_i \leq y, \text{ for all } i) \\ &= [F(y)]^n \end{aligned}$$

and the PDF is

$$f_{Y_{(n)}}(y) = n[F(y)]^{n-1}f(y).$$

Example 1.2.5. Calculate the distribution of the minimum and the maximum for the sample X_1, \dots, X_n from the exponential distribution with parameter θ . Use one of them to obtain an unbiased estimate of the parameter θ . What is the variance of this estimator?

Solution First, we calculate the CDF of each observation as

$$F_{X_i}(x) = \int_0^x \frac{1}{\theta} e^{-t/\theta} dt = 1 - e^{-x/\theta}.$$

Then, by using formulas above we calculate the densities of the maximum and the minimum:

$$f_{X_{(n)}}(x) = n(1 - e^{-x/\theta})^{n-1} \times \frac{1}{\theta} e^{-x/\theta},$$

and

$$\begin{aligned} f_{X_{(1)}}(x) &= n(e^{-x/\theta})^{n-1} \times \frac{1}{\theta} e^{-x/\theta} \\ &= \frac{n}{\theta} e^{-nx/\theta}. \end{aligned}$$

The density of the maximum is not particularly nice but the density of the minimum shows that the minimum $X_{(1)} = \min\{X_1, \dots, X_n\}$ is distributed as the exponential random variable with parameter θ/n .

So the expectation of $\hat{\theta} = nX_{(1)}$ is θ and it gives an unbiased estimator of θ .

What is its variance?

$$\mathbb{V}\text{ar}(nX_{(1)}) = n^2 \mathbb{V}\text{ar}X_{(1)} = n^2(\theta/n)^2 = \theta^2,$$

so it is not a particularly good estimator of θ . Its variance does not decline as the sample size grows.

Example 1.2.6. The sample values Y_1, Y_2, \dots, Y_n are uniform on $(\theta, \theta + 1)$. Consider the estimator $\hat{\theta} = Y_{(1)} := \min\{Y_1, \dots, Y_n\}$. Calculate the bias of $\hat{\theta}$.

Can we correct the bias?

Solution. We want to calculate $\mathbb{E}Y_{(1)}$. It is convenient to define shifted variables $X_i = Y_i - \theta$, since then $\mathbb{E}Y_{(1)} = \mathbb{E}X_{(1)} + \theta$ and it is easier to calculate $\mathbb{E}X_{(1)}$ because X_i are simply uniform random variables on $[0, 1]$. The expectation can be calculated without this transformation but the formulas would be more cumbersome.

Then, since the density and cdf of X_i are $f_X(x) = 1$ and $F_X(x) = x$ supported on $[0, 1]$, then we can use the formulas from above and calculate the pdf of the minimum $X_{(1)}$, $f_{X_{(1)}}(x) = n(1-x)^{n-1}$. In other words, $X_{(1)}$ has Beta distribution with parameters $\alpha = 1$ and $\beta = n$. By the facts about the Beta distribution, it follows that the expectation is $\mathbb{E}X_{(1)} = \alpha/(\alpha + \beta) = 1/(n+1)$.

Alternatively, we can simply integrate using the density of $X_{(1)}$, and calculate

$$\mathbb{E}X_{(1)} = n \int_0^1 x(1-x)^{n-1} dx = n \frac{\Gamma(2)\Gamma(n)}{\Gamma(n+2)} = n \frac{1!(n-1)!}{(n+1)!} = \frac{1}{n+1}.$$

(The integral can be calculated by doing integration by parts or by using a very useful formula for Beta integrals:

$$\int_0^1 x^{\alpha-1}(1-x)^{\beta-1} dx = \frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha+\beta)},$$

where $\Gamma(x)$ is the Gamma function. For integer argument x , $\Gamma(x) = (x-1)!$.

Hence the bias of $Y_{(1)} = \mathbb{E}Y_{(1)} - \theta = \mathbb{E}X_{(1)} = \frac{1}{n+1}$. Note that in this example the bias $\rightarrow 0$ as the sample size increases. In addition, we can easily correct the bias by using $\hat{\theta} = Y_{(1)} - \frac{1}{n+1}$.

What is the MSE of $\hat{\theta} = Y_{(1)} - \frac{1}{n+1}$?

Since there is no bias, we only need to calculate $\mathbb{V}\text{ar}(\hat{\theta}) = \mathbb{V}\text{ar}(Y_{(1)}) = \mathbb{V}\text{ar}(X_{(1)})$.

From the facts about the Beta distribution we have

$$\mathbb{V}\text{ar}(X_{(1)}) = \frac{\alpha\beta}{(\alpha+\beta)^2(\alpha+\beta+1)} = \frac{n}{(n+1)^2(n+2)} \sim \frac{1}{n^2}$$

for large n . (We could also calculate it directly from the density.)

This is the MSE of $\hat{\theta} = Y_{(1)} - \frac{1}{n+1}$. For large n , it is much smaller than the MSE of the estimator $\bar{Y} - \frac{1}{2}$, which we calculated as $\frac{1}{12n}$.

Summary: In this section we introduced simple measures that help us to evaluate how good an estimator is, – its bias, variance and mean squared error.

1.3 Some common unbiased estimators

1.3.1 An estimator for the population mean μ

Let Y_1, Y_2, \dots, Y_n denote a random sample of n independent identically distributed observations from a population with mean μ (that is, in our statistical model one of the parameters is $\mathbb{E}Y_i = \mu$) and variance σ^2 (another parameter is $\mathbb{V}\text{ar}(Y_i) = \sigma^2$). Then the most popular **unbiased estimator** for μ is the sample mean:

$$\hat{\mu} = \bar{Y} = \frac{1}{n} \sum_{i=1}^n Y_i.$$

In the following, for convenience we use Y to denote a random variable that has the same distribution as each of the observations Y_i .

- Unbiased: $\mathbb{E}(\bar{Y}) = \mathbb{E}Y = \mu$; so $\text{bias}(\bar{Y}) = 0$.
- Variance: $\mathbb{V}\text{ar}(\bar{Y}) = \frac{1}{n} \text{Var}(Y) = \frac{\sigma^2}{n}$;

- $MSE(\bar{Y}) = bias^2 + \text{Var} = \frac{\sigma^2}{n}$

Definition 1.3.1. The **standard error** of an estimator is the square root of its variance $SE(\hat{\theta}) = \sqrt{\text{Var}(\hat{\theta})}$

Another notation for the variance and the standard error of an estimator $\hat{\theta}$ is $\sigma_{\hat{\theta}}^2$ and $\sigma_{\hat{\theta}}$, respectively.

So the variance of the sample mean \bar{Y} is $\sigma_{\bar{Y}}^2 = \sigma^2/n$ and the standard error is $\sigma_{\bar{Y}} = \frac{\sigma}{\sqrt{n}}$.

1.3.2 An estimator for the population proportion p

Let Y_1, Y_2, \dots, Y_n denote a random sample of size n from a population with a Bernoulli distribution

$$P(Y_i = 1) = p, \quad P(Y_i = 0) = 1 - p.$$

This is a special case of the situation in the previous section and we can use the same estimator, the sample mean. The special feature of this situation is that there is only one parameter, p , not two as in the previous case. In this case, the sample mean has a special name, the **sample proportion**. It is an unbiased estimator for the parameter p (the population proportion).

$$\hat{p} = \frac{1}{n} \sum_{i=1}^n Y_i,$$

that is, it is the proportion of $Y_i = 1$'s in the sample.

The estimator is unbiased since $\mathbb{E}(\hat{p}) = \mathbb{E}Y_1 = p$. Its variance and standard error are $\sigma_{\hat{p}}^2 = \frac{pq}{n}$ and $\sigma_{\hat{p}} = \sqrt{\frac{pq}{n}}$, respectively, with $q = 1 - p$.

1.3.3 An estimator for the difference in population means

$$\mu_1 - \mu_2$$

Suppose we have two samples:

- $\{Y_1^{(1)}, Y_2^{(1)}, \dots, Y_{n_1}^{(1)}\}$ of size n_1 from Population 1: with mean μ_1 and variance σ_1^2 ;

- $\{Y_1^{(2)}, Y_2^{(2)}, \dots, Y_{n_2}^{(2)}\}$ of size n_2 from Population 2: with mean μ_2 and variance σ_2^2 ;

An **unbiased estimator** for the **difference in population means** $\theta = \mu_1 - \mu_2$ (it is the parameter of interest) is the

difference in sample means: $\hat{\theta} = \bar{Y}_1 - \bar{Y}_2 = \frac{1}{n_1} \sum_{i=1}^{n_1} Y_i^{(1)} - \frac{1}{n_2} \sum_{i=1}^{n_2} Y_i^{(2)}.$

- Unbiased: $\mathbb{E}(\hat{\theta}) = \mathbb{E}\bar{Y}_1 - \mathbb{E}\bar{Y}_2 = \mathbb{E}Y_1^{(1)} - \mathbb{E}Y_1^{(2)} = \mu_1 - \mu_2$;
- Variance: $\sigma_{\hat{\theta}}^2 = \text{Var}(\hat{\theta}) = \frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}$;
 – Proof: Since the two samples are independent, $\text{Var}(\bar{Y}_1 - \bar{Y}_2) = \text{Var}(\bar{Y}_1) + \text{Var}(\bar{Y}_2)$.
- Standard error: $\sigma_{\hat{\theta}} = \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}$.

1.3.4 An estimator for the difference in population proportions $p_1 - p_2$

Suppose we have two samples:

- $\{Y_1^{(1)}, Y_2^{(1)}, \dots, Y_{n_1}^{(1)}\}$ of size n_1 from Population 1: $P(Y_i^{(1)} = 1) = p_1$ and $P(Y_i^{(1)} = 0) = 1 - p_1$;
- $\{Y_1^{(2)}, Y_2^{(2)}, \dots, Y_{n_1}^{(2)}\}$ of size n_2 from Population 2: $P(Y_i^{(2)} = 1) = p_2$ and $P(Y_i^{(2)} = 0) = 1 - p_2$;

An **unbiased point estimator** for the difference in population means $\theta = p_1 - p_2$ (the parameter of interest) is the

difference in sample proportions: $\hat{\theta} = \hat{p}_1 - \hat{p}_2 = \frac{1}{n_1} \sum_{i=1}^{n_1} Y_i^{(1)} - \frac{1}{n_2} \sum_{i=1}^{n_2} Y_i^{(2)}$

- Unbiased: $\mathbb{E}(\hat{\theta}) = \mathbb{E}(\hat{p}_1 - \hat{p}_2) = p_1 - p_2 = \theta$;
- Variance: $\sigma_{\hat{\theta}}^2 = \text{Var}(\hat{\theta}) = \frac{p_1 q_1}{n_1} + \frac{p_2 q_2}{n_2}$;

- Standard error: $\sigma_{\hat{\theta}} = \sqrt{\frac{p_1 q_1}{n_1} + \frac{p_2 q_2}{n_2}}$.

where $q_1 = 1 - p_1$ and $q_2 = 1 - p_2$.

[Quiz]

[Quiz]

1.3.5 An estimator for the variance

Let Y_1, Y_2, \dots, Y_n denote a sample of size n from a population with mean μ and variance σ^2 , then an unbiased estimator for σ^2 is the **sample variance** (note the **$n - 1$** in the denominator).

$$S^2 \equiv \frac{\sum_{i=1}^n (Y_i - \bar{Y})^2}{n - 1}$$

The square root of S^2 is called the **sample standard deviation** and denoted, as could be expected, S .

Let us check that S^2 is indeed unbiased.

$$\begin{aligned} \mathbb{E} \sum_{i=1}^n (Y_i - \bar{Y})^2 &= \mathbb{E} \left[\sum_{i=1}^n Y_i^2 - 2 \left(\sum_{i=1}^n Y_i \right) \bar{Y} + n \bar{Y}^2 \right] = \mathbb{E} \left[\sum_{i=1}^n Y_i^2 - n \bar{Y}^2 \right] \\ &= \sum_{i=1}^n \mathbb{E} Y_i^2 - n \mathbb{E} (\bar{Y}^2) \end{aligned}$$

$$\text{for the first term: } \sum_{i=1}^n \mathbb{E} Y_i^2 = \sum_{i=1}^n [\mu^2 + \sigma^2] = n\mu^2 + n\sigma^2$$

$$\begin{aligned} \text{for the second term: } n \mathbb{E} (\bar{Y}^2) &= n \mathbb{E} (\bar{Y}^2) = n [\text{Var} \bar{Y} + (\mathbb{E} \bar{Y})^2] \\ &= n(\sigma^2/n + \mu^2) \end{aligned}$$

$$\Rightarrow \mathbb{E} \sum_{i=1}^n (Y_i - \bar{Y})^2 = (n - 1) \sigma^2$$

Hence $\mathbb{E} S^2 = \mathbb{E} \frac{A}{n-1} = \sigma^2$ and therefore S^2 is **unbiased** for σ^2 .

The variance of this estimator is more complicated to compute and we will not do it here.

Note that **S is NOT an unbiased estimator for σ**

The identity used in the first line of the proof is worth remembering:

$$\frac{1}{n} \sum_{i=1}^n (Y_i - \bar{Y})^2 = \frac{1}{n} \sum_{i=1}^n \mathbb{E} Y_i^2 - \mathbb{E}(\bar{Y}^2).$$

1.4 The error of estimation and the 2-standard-error bound

Definition 1.4.1. The error of estimation ε is the distance between an estimator and its target parameter. That is, $\varepsilon = |\hat{\theta} - \theta|$.

The error of estimation is a random quantity that change from sample to sample. The Chebyshev inequality give a method to estimate the probability that the error of estimation is larger than a multiple of the standard error of the estimator.

Reminder: $\sigma_{\hat{\theta}}$ is the standard error, that is, the standard deviation of the estimator $\hat{\theta}$. By the Chebyshev inequality:

$$\mathbb{P}\{|\hat{\theta} - \theta| > k\sigma_{\hat{\theta}}\} \leq \frac{1}{k^2}$$

- For $b = 2 \cdot \sigma_{\hat{\theta}}$, the RHS of the Chebyshev inequality is = 25%. [This is a bound, the true probability that $|\varepsilon| \geq b$ is smaller, often as small as 5%]
- $2\sigma_{\hat{\theta}}$ is called the **2-standard-error bound** on the error of the estimator. The meaning is that the error of estimation is lower than the 2-standard-error bound with large probability.

The Central Limit Theorem for sums of independent random variables says that a sum of large number of these variables has a distribution, which is closed to the normal distribution.

Since the estimator for the mean, \bar{Y} is such a sum (only divided by n), it becomes approximately normal when n is large, so if sample size is large, then the estimation error, $|\bar{Y} - \mu|$, is less than 2-standard-error, $2\sigma_{\bar{Y}}$, with probability 95% (instead of 75%).

This observation holds also for the other standard estimators that we considered in the previous section.

Example 1.4.2 (Titanic survivors). In a random sample of 136 Titanic **first class** passengers that survived the Titanic ship accident, 91 were women. In a random sample of 119 **third class** survivors, 72 were women. Assume that these are small samples from two large populations of “survivors”: first-class survivors and third-class survivors.

What is an unbiased estimate for the difference in proportions of females in these populations? What is the two-standard error bound?

Solution. $\hat{p}_1 = 66.9\%$; $\hat{p}_3 = 60.5\%$; $\hat{p}_1 - \hat{p}_3 = 6.4\%$

Two standard error bound is:

$$2\sqrt{\frac{\hat{p}_1(1-\hat{p}_1)}{n_1} + \frac{\hat{p}_3(1-\hat{p}_3)}{n_3}} = 12.1\%$$

Does the data suggest that the total first class and third class survivor populations had approximately the same proportions of females?

Does the data suggest that women from the first and third classes had approximately the same chances to survive? [Quiz]

Example 1.4.3 (Titanic survivors II). In a random sample of 95 female passengers in the first class, 91 survived the Titanic ship accident. In a random sample of 145 women in the third class, 72 survived.

What is an unbiased estimate for the difference in proportions of survivors in the populations of the first and the third class female passengers? What is the two-standard error bound?

Solution. $\hat{p}_1 = 95.8\%$; $\hat{p}_3 = 49.7\%$; $\hat{p}_1 - \hat{p}_3 = 46.1\%$

Two standard error bound is:

$$2\sqrt{\frac{\hat{p}_1(1-\hat{p}_1)}{n_1} + \frac{\hat{p}_3(1-\hat{p}_3)}{n_3}} = 9.3\%$$

Does the data suggest that female passengers from the third class had lower chances to survive than female passengers from the first class?

Some additional info on this example: The chances of a man in the first class to survive: 36.9%. The chances of a man in the third class: 13.5%.

Example 1.4.4 (Elementary school IQ). The article “A Longitudinal Study of the Development of Elementary School Children’s Private Speech” (Merrill-Palmer Q., 1990: 443–463) reported on a study of children talking to themselves (private speech). [The idea of the study is that private speech could be related to IQ, because IQ is supposed to measure mental maturity, and it was known that private speech decreases as students progress through the primary grades.] The study included 33 students whose first-grade IQ scores are given here:

082 096 099 102 103 103 106 107 108 108 108 108 109 110 110 111 113
113 113 113 115 115 118 118 119 121 122 122 127 132 136 140 146

(a) Suppose we want an estimate of the average value of IQ for the first graders served by this school. What is an unbiased estimate for this parameter?

[Hint: Sum is 3753.]

Solution. $\hat{\mu} = \bar{X} = 3753/33 = 113.7273$

(b) Calculate and interpret a point estimate of the population standard deviation σ . [Hint: Sum of squared observations is 432,015]

Solution.

$$S^2 = \frac{1}{32} \left(432,015 - 33 \times (113.7273)^2 \right) = 162.3856$$

$$\hat{\sigma} = S = \sqrt{162.3856} = 12.7431$$

While S^2 is an unbiased estimator of σ^2 , the estimator S for σ is biased.

(c) What is the two-standard-error bound on the error of estimation?

Solution.

The two-standard-error bound is $2\sigma_{\hat{\mu}} = 2\sigma_{\bar{X}} = 2\sigma/\sqrt{n}$. We use the estimate of σ to calculate the bound.

$$2S/\sqrt{n} = 2 \times 12.7431/\sqrt{33} = 4.4366$$

Since the estimate of μ is 113.7273 and the two-error bound for the error of estimation is 4.4366, the data suggest that this is an above average class, because the nationwide IQ average is around 100.

(d) Calculate a point estimate of the proportion of all such students whose IQ exceeds 100. [Hint: Think of an observation as a “success” if it exceeds 100.]

Solution.

The number of students with IQ above 100 is 30. So the point estimate is $\hat{p} = 30/33 = 90.91\%$.

Example 1.4.5 (Elementary school IQ II). The data set mentioned in the previous example also includes these third grade verbal IQ observations for males:

117 103 121 112 120 132 113 117 132 149 125 131 136 107 108 113 136
114

(18 observations) and females:

114 102 113 131 124 117 120 90 114 109 102 114 127 127 103

(15 observations)

Let the male values be denoted X_1, \dots, X_m and the female values Y_1, \dots, Y_n .

(a) Calculate the point estimate for the difference between male and female verbal IQ.

Solution.

$$\bar{X} - \bar{Y} = \frac{2186}{18} - \frac{1707}{15} = 121.4444 - 113.8 = 7.6444$$

(b) What is the standard error of the estimator?

Solution. First we calculate the sample variances S_x^2 and S_y^2 for these two samples.

$$S_x^2 = \frac{1}{17} \left(268,046 - 18 \times 121.4444^2 \right) = 151.0964$$

$$S_y^2 = \frac{1}{14} \left(196,039 - 15 \times 113.8^2 \right) = 127.3143$$

Then, we calculate the estimate of the standard error:

$$\hat{\sigma}_{\hat{\theta}} = \sqrt{\frac{S_x^2}{m} + \frac{S_y^2}{n}} = \sqrt{\frac{151.0964}{18} + \frac{127.3143}{15}} = 4.1088$$

So we see that the estimate of the difference 7.6444. However, the two-standard-error bound is 8.2176 and the data does not give an evidence that the difference is positive.

[End of section Quiz]

1.5 Confidence intervals and pivotal quantities

A **point estimator** is a function of data sample that gives a single number that is our “best guess” for the parameter, for examples: \bar{Y} for μ and \hat{p} for p . It is difficult to gauge its error unless we add some additional info like standard error.

The **interval estimator**: an interval $(\hat{\theta}_L, \hat{\theta}_U)$ (a lower end and an upper end) which the parameter is **believed** to fall into. Both ends of the interval are functions of the data. In particular, they are random.

For example, for μ , we could take $(\bar{Y} - 2S, \bar{Y} + 2S)$. Alternatively, we could take $(0.8\bar{Y}, 1.2\bar{Y})$.

The quality of the interval estimator can be judged by two characteristics: its length and the probability that it covers the true value of the parameter.

Interval estimators are commonly called *confidence intervals*. The **confidence coefficient** $(1 - \alpha)$ is the probability that a confidence interval $(\hat{\theta}_L, \hat{\theta}_U)$ covers the true parameter, that is

$$P(\hat{\theta}_L \leq \theta \leq \hat{\theta}_U) = 1 - \alpha$$

(α should be small, like 5%, then the confidence coefficient is large: 95%.)

Question: for a given $(1 - \alpha)$, how can we find the desired confidence interval $(\hat{\theta}_L, \hat{\theta}_U)$??

The difficulty is that we do not know true value of the parameter, so for example we cannot use the Chebyshev inequality to build the interval estimate: The interval $(\theta - 2\sigma_{\hat{\theta}}, \theta + 2\sigma_{\hat{\theta}})$ is **not an interval estimator** because we do not know neither θ , no $\sigma_{\hat{\theta}}$.

Of course, we can use $\hat{\theta}$ instead of θ in this interval and estimate $\sigma_{\hat{\theta}}$, but

why is this OK? We need a more rigorous and more difficult argument: **the pivotal quantity method**.

The **pivotal quantity** or **pivot** is a quantity which is a function of both the sample data and the parameter θ , but whose **distribution does not depend** on the parameter θ !

Note that the distribution of the sample data also depends on the parameter θ . So the pivotal quantity depends on θ both directly and indirectly, through the data, and our aim is to make sure that the dependences of the pivot on the data and on the parameter somehow cancel each other on the distribution level.

Let X denote the data sample. (It is a vector of observations.) Find a function $T(X, \theta)$ (the pivot) so that its distribution does not depend on θ and so it is known. [This is the crucial point and often it is a tough problem.]

Use the distribution of T to find a pair of L and U such that

$$\Pr(L \leq T \leq U) = 1 - \alpha$$

Manipulate the inequalities $L \leq T(X, \theta) \leq U$ so that they become

$$L^*(X) \leq \theta \leq U^*(X),$$

so that **θ is in the middle!!!**

- Note that both the LHS and the RHS are random as both are function of the data X .
- Often not the whole data sample is used in this process but a summary of the data in the form of a statistic $f(X)$ and the pivot is searched in the form $T(f(X), \theta)$.

Example 1.5.1. Recollect our example about the smartphone lifetime, where we had a sample X_1, \dots, X_n of random variables distributed according to the exponential distribution with parameter θ . Recall that then the minimum of the sample $X_{(1)} = \min\{X_1, \dots, X_n\}$ is distributed according to the exponential distribution with parameter θ/n . We can use $X_{(1)}$ to construct an

unbiased point estimator $\hat{\theta} = nX_{(1)}$. although this is a rather bad estimator since its variance does not decrease when n grows. Here we illustrate the pivotal quantity method by using the statistic $X_{(1)}$ to develop an interval estimator (“confidence interval”) for θ with the confidence coefficient 0.90 (so $\alpha = 0.1$).

We are looking for a pivotal quantity, that is a function that depends only on $X_{(1)}$ and θ so that this function $T(X_{(1)}, \theta)$ has the distribution which is independent of θ .

We claim that $T = nX_{(1)}/\theta$ is a good function for this purpose and that its distribution is an exponential with parameter 1.

Proof. For conciseness of notation, let us write Y to denote $X_{(1)}$. The density of Y is

$$f_Y(y) = \frac{n}{\theta} e^{-ny/\theta}.$$

We are making transformation $T = nY/\theta$ that has the inverse transformation $Y = (\theta/n)T$. Let us use notation $y(t)$ for the function $y = (\theta/n)t$. By using the density transformation method, we calculate the density of T as follows :

$$\begin{aligned} f_T(t) dt &= f_Y(y(t)) dy(t) = f_Y(y(t)) \frac{dy(t)}{dt} dt \\ &= \frac{n}{\theta} e^{-n(\theta/n)t/\theta} \times \frac{\theta}{n} dt = e^{-t} dt \end{aligned}$$

So, T has the exponential density with parameter 1. □

Now we look for L and U , so that

$$0.90 = \Pr(L \leq T \leq U) = \int_L^U e^{-t} dt = e^{-L} - e^{-U}.$$

There are infinitely many combinations of L and U which satisfy this. One possibility is to let

$$\Pr(T > U) = 0.05 \text{ and } \Pr(T < L) = 0.05$$

Solutions are $L = 0.051$ and $U = 2.996$ Now we have

$$0.051 \leq T = n \frac{X_{(1)}}{\theta} \leq 2.996$$

We manipulate these two inequalities to put θ in the middle:

$$\frac{nX_{(1)}}{2.996} \leq \theta \leq \frac{nX_{(1)}}{0.051}$$

Remark: This example is only meant to illustrate the method. However, note that the resulting confidence interval is not very good. Since the variance of nX_1 does not go to zero as n grows we cannot expect that the length of this confidence interval goes to 0 as n grows.

Example 1.5.2 (A sample from a normal distribution). Let X_1, \dots, X_n be a sample from a normal distribution $\mathcal{N}(\mu, \sigma^2)$, where the parameter σ is known. Then \bar{X} is a normally distributed random variable with mean μ and variance σ^2/n .

Then the quantity

$$T = \frac{\bar{X} - \mu}{\sigma/\sqrt{n}}$$

has the normal distribution with mean 0 and variance 1. In particular, its distribution does not depend on μ . (Here T depends not only on the data and parameter μ but also on σ but we assume that σ is known.) By convention, in this important example, the quantity T is denoted Z .

The next task is to look for L and U so that $\Pr(L \leq Z \leq U) = 1 - \alpha$. There are 3 standard ways to do it. One of them is to choose L and U so that $\Pr(Z < L) = \alpha/2$ and $\Pr(Z > U) = \alpha/2$. By symmetry of the normal distribution $L = -U$ and by convention this U is denoted $z_{\alpha/2}$.

Then we have

$$-z_{\alpha/2} \leq Z = \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \leq z_{\alpha/2}.$$

This inequality can be converted to the desired confidence interval for parameter μ :

$$\bar{X} - z_{\alpha/2} \frac{\sigma}{\sqrt{n}} \leq \mu \leq \bar{X} + z_{\alpha/2} \frac{\sigma}{\sqrt{n}}$$

Alternatively, we can look for a “one-sided interval”. This means that we take $L = -\infty$ and look for U such that $\Pr\{Z > U\} = \alpha$. By definition this U is denoted z_α and it can be found from a table or by using software.

Then, the inequality is

$$Z = \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \leq z_\alpha$$

and it can be transformed to the *lower confidence bound* on the parameter μ

$$\bar{X} - z_\alpha \frac{\sigma}{\sqrt{n}} \leq \mu.$$

Note the difference from the previous inequality. Here we use the factor z_α before the standard error $\frac{\sigma}{\sqrt{n}}$ while in the previous inequality we used $z_{\alpha/2}$.

Similarly, by using $U = \infty$ and looking for L such that $\mathbb{P}\{Z < L\} = \alpha$, we can derive the *upper confidence bound* on μ :

$$\bar{\mu} < X - z_\alpha \frac{\sigma}{\sqrt{n}}.$$

Example 1.5.3 (Confidence interval for σ^2). Suppose that $X_1, \dots, X_n \sim N(\mu, \sigma^2)$ but now we know μ and we are interested in deriving a confidence interval for σ^2 .

The quantity $T = \sum_{i=1}^n ((X_i - \mu)/\sigma)^2$ is known to be distributed according to the χ^2 distribution with n degrees of freedom, $\chi^2(n)$. Therefore, it is a valid pivotal quantity and we can use it to derive a confidence interval for σ^2 .

Again, it can be done in three possible ways. One of them is to find the quantities L and U such that $\mathbb{P}(T < L) = \alpha/2$ and $\mathbb{P}(T > U) = \alpha/2$. In this case, the distribution is not symmetric and we need to find 2 really different quantities. The quantity U is $\chi_{\alpha/2}^2(n)$ and the quantity L is $\chi_{1-\alpha/2}^2(n)$. They can be found from a table or by using software.

Then, we get

$$\chi_{1-\alpha/2}^2(n) \leq \frac{1}{\sigma^2} \sum_{i=1}^n (X_i - \mu)^2 \leq \chi_{\alpha/2}^2(n),$$

and, by putting σ^2 in the middle,

$$\frac{1}{\chi_{\alpha/2}^2(n)} \sum_{i=1}^n (X_i - \mu)^2 \leq \sigma^2 \leq \frac{1}{\chi_{1-\alpha/2}^2(n)} \sum_{i=1}^n (X_i - \mu)^2$$

Note that for large n , $\chi_{\alpha/2}^2(n)$ is relatively close to n .

The upper and lower confidence bounds can be found similarly. For example, $(1 - \alpha)$ **upper** confidence bound for σ^2 is

$$\sigma^2 \leq \frac{1}{\chi_{1-\alpha}^2(n)} \sum_{i=1}^n (X_i - \mu)^2$$

Note that in this inequality we use α instead of $\alpha/2$.

1.6 Large-sample confidence intervals

Finding a good pivotal quantity for a parameter is NOT an easy job! The good news is that for large sample we can easily find an **approximate pivotal quantity**. The Central Limit Theorem (CLT) ensures that when **the sample size n is large enough** an appropriate estimator is in many cases is approximately normal random variable.

- for $\theta = \mu$, the estimator $\hat{\theta} = \bar{Y}$ is approximately $\sim N(\mu, \frac{\sigma^2}{n})$;
- for $\theta = p$, the estimator $\hat{\theta} = \hat{p}$ is approximately $\sim N(p, \frac{p(1-p)}{n})$;
- for $\theta = \mu_1 - \mu_2$, the estimator $\hat{\theta} = \bar{Y}_1 - \bar{Y}_2$ is approximately $\sim N(\mu_1 - \mu_2, \frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2})$;
- for $\theta = p_1 - p_2$, the estimator $\hat{\theta} = \hat{p}_1 - \hat{p}_2$ is approximately $\sim N(p_1 - p_2, \frac{p_1(1-p_1)}{n_1} + \frac{p_2(1-p_2)}{n_2})$;

In general, if we are willing to assume that the approximation is sufficiently good and an estimator $\hat{\theta}$ has a normal distribution,

$$\hat{\theta} \sim N(\theta, \sigma_{\hat{\theta}}^2),$$

where $\sigma_{\hat{\theta}}^2 = \text{Var}(\hat{\theta})$, then we can write a pivotal quantity:

$$Z \equiv \frac{\hat{\theta} - \theta}{\sigma_{\hat{\theta}}} \sim N(0, 1).$$

Here $\hat{\theta}$ is a function of the data and $\sigma_{\hat{\theta}}$ is either given as known, or is assumed to be estimated with good precision.

Since $\hat{\theta}$ is usually only approximately normal, then Z is an approximately pivotal quantity and the resulting confidence intervals will be only approximate. Their quality depends on how well the CLT works for the estimator $\hat{\theta}$.

Example 1.6.1 (Two-sided confidence interval). By using our results for the normally distributed variables X_1, \dots, X_n , we write the (approximate) two-sided interval for θ , based on the point estimator $\hat{\theta}$ which is assumed to be approximately distributed as $N(\theta, \sigma_{\hat{\theta}}^2)$.

The two-sided confidence interval for θ with confidence coefficient $1 - \alpha$ is

$$\left[\hat{\theta} - z_{\alpha/2} \sigma_{\hat{\theta}}, \quad \hat{\theta} + z_{\alpha/2} \sigma_{\hat{\theta}} \right]$$

Notations

- z_c : $P(Z > z_c) = c$ where $Z \sim N(0, 1)$
- $t_c(n)$: $P(T > t_c) = c$ where $T \sim t$ distribution with n degrees of freedom.
- $\chi_c^2(n)$: $P(X > \chi_c^2) = c$ where $X \sim \chi^2$ distribution with n degrees of freedom.

Example 1.6.2 (Upper and lower confidence bounds). The one-sided large sample confidence intervals are as follows:

- The upper bound confidence interval with confidence coefficient $1 - \alpha$ is

$$(-\infty, \quad \hat{\theta} + z_{\alpha} \sigma_{\hat{\theta}} \quad]$$

- The lower bound confidence interval with confidence coefficient $1 - \alpha$ is

$$[\hat{\theta} - z_{\alpha} \sigma_{\hat{\theta}} \quad , +\infty)$$

Example 1.6.3. A study was done on 41 first-year medical students to see if their anxiety levels changed during the first semester. One measure used was the level of serum cortisol, which is associated with stress. For each of the 41 students the level was compared during finals at the end of the

semester against the level in the first week of classes. The average difference was 2.08 with a standard deviation of 7.88. Find a 95% lower confidence bound for the population mean difference μ . Does the bound suggest that the mean population stress change is necessarily positive?

Example 1.6.4. A random sample of 539 households from a mid-western city was selected, and it was determined that 133 of these households owned at least one firearm (“The Social Determinants of Gun Ownership: Self-Protection in an Urban Environment,” *Criminology*, 1997: 629–640). Using a 95% confidence level, calculate a lower confidence bound for the proportion of all households in this city that own at least one firearm.

Example 1.6.5. Two brands of refrigerators, denoted A and B, are each guaranteed for 1 year. In a random sample of 50 refrigerators of brand A, 12 were observed to fail before the guarantee period ended. An independent random sample of 60 brand B refrigerators also revealed 12 failures during the guarantee period. Estimate the true difference $(p_1 - p_2)$ between proportions of failures during the guarantee period, with the confidence coefficient approximately .98.

Solution.

- $n_A = 50$ and $Y_A = 12$, hence $\hat{p}_1 = 12/50 = 0.24$
- $n_B = 60$ and $Y_B = 12$, hence $\hat{p}_2 = 12/60 = 0.2$
- Use $\hat{\theta} = \hat{p}_1 - \hat{p}_2 = 0.04$ as the (point) estimator.
- $\hat{\theta}$ is approximately normal and we have

$$\hat{\theta} \pm z_{\alpha/2} \sigma_{\hat{\theta}}$$

as the $100(1 - \alpha)\%$ confidence interval.

- Note that

$$\begin{aligned} \sigma_{\hat{\theta}} &= \sqrt{\text{Var}\hat{\theta}} \\ &= \sqrt{p_1(1 - p_1)/n_A + p_2(1 - p_2)/n_B}, \end{aligned}$$

where p_1 and p_2 are unknown but can be approximated by \hat{p}_1 and \hat{p}_2 .

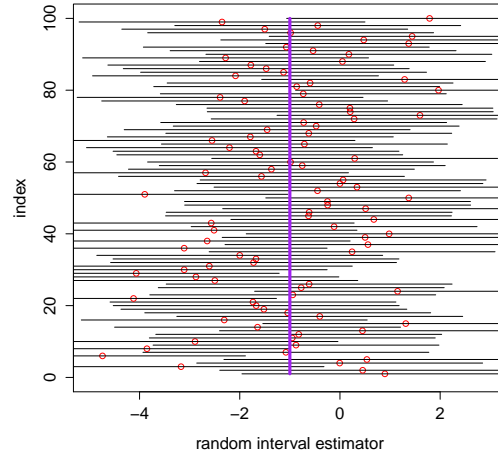


Figure 1.1: Interval estimator for 100 different samples. The confidence interval $(\hat{\theta} - 1.96\sigma_{\hat{\theta}}, \hat{\theta} + 1.96\sigma_{\hat{\theta}})$, – centered at the point estimators $\hat{\theta}$ shown by red circles, – in 95% cases covers the true parameter (shown by the purple line).

So,

$$\hat{\sigma}_{\hat{\theta}} = \sqrt{0.24(1 - 0.24)/50 + 0.2(1 - 0.2)/60} = 0.0795.$$

We also have $z_{\alpha/2} = z_{0.01} = 2.326348$, and therefore the confidence interval is

$$(0.04 - 2.326348 \times 0.0795, 0.04 + 2.326348 \times 0.0795) = (-0.1449, 0.2249)$$

[Quiz]

1.7 How to determine the sample size

The sample size dilemma

- We need to collect samples to make inference about the population parameter. Question: how large n should be? 30? 40? 50?
- On the one hand, the more data you have, the more accurate is your estimator $\hat{\theta}$ for θ

- On the other hand, collecting samples is NOT free: it costs money, time, personnel..

Conclusion: We want minimal sample size which allows us to achieve given precision with a given level of confidence.

The formulas for confidence intervals provide an easy way to find the required size of the sample.

Precision + Confidence level \rightarrow Required minimal sample size.

Rather than give a bunch of formulas we illustrate the method in examples.

Example 1.7.1. The reaction of an individual to a stimulus in a psychological experiment may take one of two forms, A or B. If an experimenter wishes to estimate the probability p that a person will react in manner A, how many people must be included in the experiment? Assume that the experimenter will be satisfied if the error of estimation is less than .04 with probability equal to .90. Assume also that he expects p to lie somewhere in the neighborhood of .6.

- This is an estimating p in a binomial distribution problem. n is to be found.
- Although n is yet to be found, let's assume that it is large enough, in which case $\hat{\theta} = \hat{p}$ is approximately normal, and then $Z \equiv \frac{\hat{p}-p}{\sqrt{p(1-p)/n}} \approx \frac{\hat{p}-p}{\sqrt{\tilde{p}(1-\tilde{p})/n}}$ is approximately standard normal and hence, $\mathbb{P}(-z_{0.05} \leq \frac{\hat{p}-p}{\sqrt{\tilde{p}(1-\tilde{p})/n}} \leq z_{0.05}) = 0.9$. This is to say that with probability 0.9,

$$|\hat{p} - p| < z_{0.05} \sqrt{p(1-p)/n}$$

- If we want this error to be smaller than 0.04, only need to have $z_{0.05} \sqrt{\tilde{p}(1-\tilde{p})/n} \leq 0.04$. Solve n and we have $n \geq 406$.
- The prior information $\tilde{p} = 0.6$ is used to approximate the standard error

Example 1.7.2. A state wildlife service wants to estimate the mean number of days that each licensed hunter actually hunts during a given season, with a bound on the error of estimation equal to 2 hunting days. If data collected in earlier surveys have shown σ to be approximately equal to 10, how many hunters must be included in the survey?

- The client is not sophisticated and does not formulate explicitly what is the level of confidence required. In this case, it is typical to set the confidence level at 95% and use the 2-standard-error bound. If we want the error of estimation to be less than 2, then

$$2 > 2\sigma_{\hat{\theta}} = 2\frac{\sigma}{\sqrt{n}} = 2\frac{10}{\sqrt{n}} \Rightarrow n > 100.$$

Example 1.7.3. Telephone pollsters often interview between 1000 and 1500 individuals regarding their opinions on various issues. A survey question asks if a person believes that the performance of their athletics teams has a positive impact on the perceived prestige of the institutions. The goal of the survey is to see if there is a difference between the opinions of men and women on this issue. Suppose that you design the survey and wish to estimate the difference in a pair of proportions, correct to within .02, with probability .9. How many interviewees should be included in each sample?

1. What is the standard error of $\hat{\theta}$? $\sigma_{\hat{\theta}} = \sqrt{\frac{p_1(1-p_1)}{n_1} + \frac{p_2(1-p_2)}{n_2}}$
2. $\hat{\theta} = \hat{p}_1 - \hat{p}_2$ is approximately normal with mean $\theta = p_1 - p_2$ and variance $\sigma_{\hat{\theta}}^2$;
3. Hence, with probability 0.9, $|\hat{\theta} - \theta| < z_{0.05}\sigma_{\hat{\theta}}$
4. Take $n_1 = n_2 = n$. For p_1 and p_2 , since we have no prior information, we just replace them by 0.5, as the most conservative guess.
5. If we want this error to be smaller than 0.02, only need to solve

$$z_{0.05}\sqrt{\frac{1/4}{n} + \frac{1/4}{n}} < 0.02$$

6.

$$n \geq \frac{1}{2} \left(\frac{z_{0.05}}{0.02} \right)^2 = \frac{1}{2} \left(\frac{1.645}{0.02} \right)^2 = 3382.5$$

So we should take $n = 3383$.

[Quizzes]

1.8 Small-sample confidence intervals

1.8.1 Small sample CIs for μ and $\mu_1 - \mu_2$

Suppose, the parameter of interest is the **population mean** μ and we have a sample Y_1, \dots, Y_n . When the sample size is large, the Central Limit Theorem ensures that \bar{Y} is approximately normal with distribution $N(\mu, \frac{\sigma^2}{n})$.

In addition, the parameter σ^2 can be reliably estimated by the sample variance. Thus, the quantity

$$Z = \frac{\bar{Y} - \mu}{S/\sqrt{n}}$$

is approximately pivotal with the distribution $N(0, 1)$. Based on Z , we can find

- 2-sided confidence interval for μ : $[\bar{Y} - z_{\alpha/2} \frac{S}{\sqrt{n}}, \bar{Y} + z_{\alpha/2} \frac{S}{\sqrt{n}}]$
- 1-sided lower confidence interval for μ : $[\bar{Y} - z_{\alpha} \frac{S}{\sqrt{n}}, \infty]$
- 1-sided upper confidence interval for μ : $[-\infty, \bar{Y} + z_{\alpha} \frac{S}{\sqrt{n}}]$

Now suppose that the sample size n is small, say less than 30. Then, if the data is not normally distributed, the CLT does not help us and even if the variance of the population σ^2 is known,

$$T = \frac{\bar{Y} - \mu}{\sigma/\sqrt{n}}$$

is not normally distributed. Moreover, it is not a pivotal quantity. Its distribution can still depend on μ .

If the data IS normally distributed and σ^2 is known, then

$$T = \frac{\bar{Y} - \mu}{\sigma/\sqrt{n}}$$

is pivotal and normal. However, in most cases this fact cannot be used to construct a confidence interval for μ since σ^2 is not known.

Finally, if the data is normal and we try to use the sample variance S^2 instead of σ^2 , then

$$T = \frac{\bar{Y} - \mu}{S/\sqrt{n}}$$

can be proved to be pivotal but its distribution is not normal. The reason is that S^2 is a random quantity and dividing the normal random variable $\bar{Y} - \mu$ by a random quantity instead of a deterministic coefficient breaks normality of the quotient.

Conclusion: **Using the normal distribution in the case of small samples leads to erroneous intervals!**

A reminder from Probability Theory course.

In order to explain how to find a pivotal quantity for a small sample that consists of normally distributed random variables, we need to recall some results from probability theory.

Definition 1.8.1. If $X_1, \dots, X_n \sim N(0, 1)$, then the random variable $Y \equiv X_1^2 + \dots + X_n^2$ has the χ^2 *distribution* with n degrees of freedom.

Theorem 1.8.2 (Joint distribution of sample mean and sample variance). *Let $X_1, \dots, X_n \sim N(\mu, \sigma^2)$. Then the sample variance S^2 is **independent** of the sample mean \bar{X} and*

$$\frac{(n-1)S^2}{\sigma^2} = \sum_{i=1}^n \left(\frac{X_i - \bar{X}}{\sigma} \right)^2$$

has the χ^2 distribution with $n-1$ degrees of freedom.

Note that the random variable has χ^2 distribution with one less degree of freedom than if we would add together n independent standard normal

random variables. The reason for this reduction is that the random variables $Y_i - \bar{Y}$ are not independent. Intuitively they can be expressed in terms of $n - 1$ independent normal random variables, hence the reduction in the degree of freedom. The most surprising in this theorem is the independence of \bar{Y} and S^2 .

Definition 1.8.3. Suppose random variables $Z \sim N(0, 1)$ and $S \sim \chi^2(n)$ are independent. Then the random variable

$$T \equiv \frac{Z}{\sqrt{\frac{S}{n}}}$$

is distributed according to the t -distribution with n degrees of freedom, denoted as

$$T \sim t(n)$$

This distribution is also often called [Student's](#) t -distribution. It was discovered by William Gosset who worked for Guinness brewery and wrote his papers under the pen name Student.

It can be shown that for large n , – say for $n > 30$, the $t(n)$ -distribution is very close to the standard normal distribution $N(0, 1)$.

In fact t distribution has many properties that are similar to the normal distribution. For example, its PDF is symmetric with respect to 0.

However, t distribution has heavier tails: When n is small, a r.v with the $t(n)$ -distribution takes large values with much larger probability than a standard normal random variable.

Now, we know that if the variables Y_i are i.i.d and have normal distribution $N(\mu, \sigma^2)$, then, first,

$$\frac{\bar{Y} - \mu}{\sigma/\sqrt{n}} \sim N(0, 1),$$

and, second,

$$\frac{(n-1)S^2}{\sigma^2} \sim \chi^2(n-1)$$

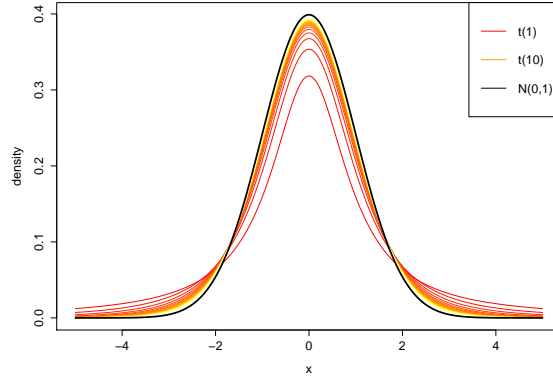


Figure 1.2: Densities of t distribution for several values of parameter df .

and these two random quantities are independent. Hence,

$$T = \frac{\frac{\bar{Y} - \mu}{\sigma/\sqrt{n}}}{\sqrt{\frac{\frac{n-1}{\sigma^2} S^2}{n-1}}} = \frac{\bar{Y} - \mu}{S/\sqrt{n}} \sim t(n-1),$$

that is T has a t -distribution with **degrees freedom** $df = n - 1$; Based on T ,

- 2-sided confidence interval for μ : $[\bar{Y} - t_{\alpha/2}\sigma_{\bar{Y}}, \bar{Y} + t_{\alpha/2}\sigma_{\bar{Y}}]$
- 1-sided lower confidence interval for μ : $[\bar{Y} - t_{\alpha}\sigma_{\bar{Y}}, \infty]$
- 1-sided upper confidence interval for μ : $[-\infty, \bar{Y} + t_{\alpha}\sigma_{\bar{Y}}]$;

t_{α} can be found using the t -table (Table 5, look for subscript α with $df = n - 1$).

Comparison of the t distribution and the standard normal distribution.

The densities of the t - distribution for several values of the degree of freedom are shown in Figure 1.2. We can note that

1. t distribution often has fatter tails than normal distribution. That means more chance to see more extreme values.

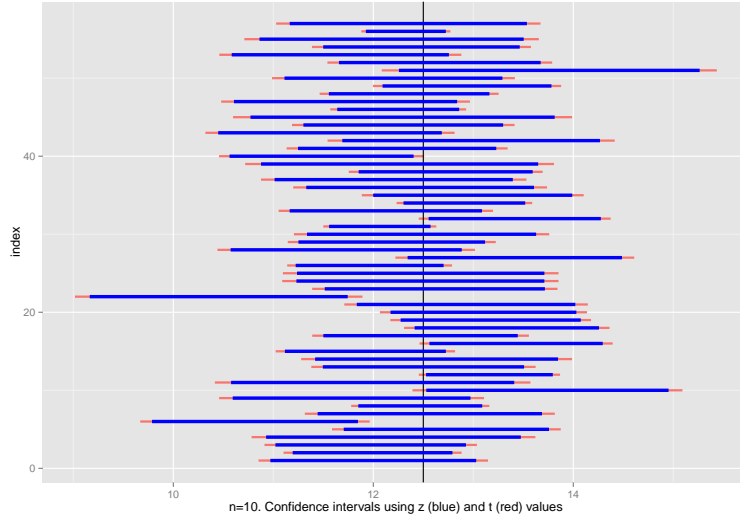


Figure 1.3: Comparison of z and t confidence intervals for $n = 10$.

2. When df increases, $t(df)$ distribution becomes closer and closer to normal.

Note that the statistic in both large sample and small sample cases is the same:

$$T \equiv \frac{\bar{Y} - \mu_Y}{S/\sqrt{n}}$$

The only difference is that in the case of a large sample, it is distributed as normal random variable, and in the case of a small sample it is distributed as a t -random variable, and only if the data is normal.

Figure 1.3 compares the z and t confidence intervals in a sample with n observations.

- Red: using the $t_{\alpha/2}$ value – the correct one
- Blue: using the $z_{\alpha/2}$ value – the incorrect one
 - Fail to take the uncertainty of S^2 into consideration
 - Fail to deliver the promised coverage probability

– Shorter than the red one (for small α , $z_{\alpha/2} < t_{\alpha/2}$)

Example 1.8.4. The reaction time (RT) to a stimulus is the interval of time commencing with stimulus presentation and ending with the first discernible movement of a certain type. The article “Relationship of Reaction Time and Movement Time in a Gross Motor Skill” (Percept. Motor Skills, 1973: 453–454) reports that the sample average RT for 16 experienced swimmers to a pistol start was .214 s and the sample standard deviation was .036 s.

Making any necessary assumptions, derive a 90% CI for true average RT for all experienced swimmers.

Now let us consider the two sample case. So we have samples X_1, \dots, X_{n_1} and Y_1, \dots, Y_{n_2} , with observations distributed according to $N(\mu_1, \sigma_1)$ and, $N(\mu_2, \sigma_2)$ respectively. We assume that n_1 and n_2 are small and want to find C.I. for $\mu_1 - \mu_2$.

In the large-sample case we used the fact that $\bar{X} - \bar{Y}$ is approximately normal with mean $\mu_1 - \mu_2$ and variance

$$\sigma_{\bar{X}-\bar{Y}} = \frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}$$

So we could write the approximate C.I. as

$$\left(\bar{X} - \bar{Y} - z_{\alpha/2} \sqrt{\frac{S_1^2}{n_1} + \frac{S_2^2}{n_2}}, \bar{X} - \bar{Y} + z_{\alpha/2} \sqrt{\frac{S_1^2}{n_1} + \frac{S_2^2}{n_2}} \right),$$

where S_1^2 and S_2^2 are sample variances for the two samples.

In principle, in the small sample case, this formula holds with $z_{\alpha/2}$ changed to $t_{\alpha/2}(\nu)$. However, the number of degrees of freedom ν is a rather complicated function of n_1 and n_2 .

Here we consider only the most simple case when it is assumed that $\sigma_1 = \sigma_2 = \sigma$. Then we can define the pooled-sample estimator for the common variance σ^2 ,

$$\begin{aligned} S_p^2 &\equiv \frac{\sum_{i=1}^{n_1} (X_i - \bar{X})^2 + \sum_{i=1}^{n_2} (Y_i - \bar{Y})^2}{n_1 + n_2 - 2} \\ &= \frac{(n_1 - 1)S_1^2 + (n_2 - 1)S_2^2}{n_1 + n_2 - 2} \end{aligned}$$

One can prove that S_p^2 is independent of $\bar{X} - \bar{Y}$ and that it is a multiple of a r.v. with the χ^2 distribution with $n_1 + n_2 - 2$ degrees of freedom.

- When the sample size n is small and σ replaced by S_p , the pivotal quantity for $\mu_1 - \mu_2$ becomes

$$T = \frac{\bar{Y}_1 - \bar{Y}_2 - (\mu_1 - \mu_2)}{\widehat{\sigma}_{\bar{Y}_1 - \bar{Y}_2}} = \frac{\bar{Y}_1 - \bar{Y}_2 - (\mu_1 - \mu_2)}{S_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}} \sim t(n_1 + n_2 - 2),$$

that is, T has the t -distribution with $df = n_1 + n_2 - 2$;

So in this simple case, we have the following confidence interval for $\mu_1 - \mu_2$

$$\left(\bar{X} - \bar{Y} - t_{\alpha/2}^{(n_1+n_2-2)} S_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}, \bar{X} - \bar{Y} + t_{\alpha/2}^{(n_1+n_2-2)} S_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}} \right),$$

Similarly, the lower bound confidence interval for $\mu_1 - \mu_2$ is

$$\left(\bar{X} - \bar{Y} - t_{\alpha} S_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}, \infty \right)$$

and the upper bound confidence interval for $\mu_1 - \mu_2$ is

$$\left(-\infty, \bar{X} - \bar{Y} + t_{\alpha} S_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}} \right)$$

Example 1.8.5. The carapace lengths of ten lobsters examined in a study of the infestation of the *Thenus orientalis* lobster by two types of barnacles, *Octolasmis tridens* and *O. lowei*, are given in the following table. Find a 95% confidence interval for the mean carapace length (in millimeters, mm) of *T.orientalis* lobsters caught in the seas in the vicinity of Singapore.

Lobster Field Number	A061	A062	A066	A070	A067	A069	A064	A068	A065	A063
Carapace Length (mm)	78	66	65	63	60	60	58	56	52	50

This is a small sample estimation problem for μ . Below are calculations done in R.

```
> x = c(78,66,65,63,60,60,58,56,52,50)
> x
```

```

[1] 78 66 65 63 60 60 58 56 52 50
> mean(x) # sample mean
[1] 60.8
> sum((x-60.8)^2)/(10-1) # sample variance
[1] 63.51111
> sqrt(sum((x-60.8)^2)/(10-1)) # sample standard deviation
[1] 7.969386
> qt(0.975,9) # 0.025 percentage point
               # for t distribution with df=9
[1] 2.262157
> qt(0.975,9)*sqrt(sum((x-60.8)^2)/(10-1))/sqrt(10)
[1] 5.700955

```

Answer: 60.8 ± 5.700955 . Note that we have divided by $\sqrt{9}$ when we estimated S and then again by $\sqrt{10}$ when we estimated $\sigma_{\bar{X}}$. Do not forget the second division.

Example 1.8.6. Chronic anterior compartment syndrome is a condition characterized by exercise-induced pain in the lower leg. Swelling and impaired nerve and muscle function also accompany the pain, which is relieved by rest. Susan Beckham and her colleagues conducted an experiment involving ten healthy runners and ten healthy cyclists to determine if pressure measurements within the anterior muscle compartment differ between runners and cyclists. The data—compartment pressure, in millimeters of mercury—are summarized in the following table:

Condition	Runners		Cyclists	
	Mean	s	Mean	s
Rest	14.5	3.92	11.1	3.98
80% maximal O ₂ consumption	12.2	3.49	11.5	4.95

1. Construct a 95% confidence interval for the difference in mean compartment pressures between runners and cyclists under the resting condition.

2. Construct a 90% confidence interval for the difference in mean compartment pressures between runners and cyclists who exercise at 80% of maximal oxygen (O₂) consumption.

This is the two sample problem for estimation of $\mu_1 - \mu_2$. We solve it using R.

```
> 14.5-11.1 # diff in sample means.
[1] 3.4
> (3.92^2*9+3.98^2*9)/18 # pooled sample variance
[1] 15.6034
> sqrt((3.92^2*9+3.98^2*9)/18) # pooled sample standard deviation
[1] 3.950114
> # 0.025 percentage point for t distribution with df=18
> qt(0.975,18)
[1] 2.100922
>
> qt(0.975,18) * sqrt((3.92^2*9+3.98^2*9)/18) * sqrt(1/10+1/10)
[1] 3.711373
```

Answer: 3.4 ± 3.711373 . The R command `qt(0.975,18)` gives 2.100922 which can also be obtained from $t_{0.025}$ for $df=18$ in Table 5

Now let us solve the second part of the exercise.

```
> 12.2-11.5 # diff in sample means.
[1] 0.7
> (3.49^2*9+4.95^2*9)/18 # pooled sample variance
[1] 18.3413
> sqrt((3.49^2*9+4.95^2*9)/18) # pooled sample standard deviation
[1] 4.282674
> # 0.05 percentage point for t distribution with df=18
> qt(0.95,18)
[1] 1.734064
> qt(0.95,18) * sqrt((3.49^2*9+4.95^2*9)/18) * sqrt(1/10+1/10)
[1] 3.3212
```


Answer: 0.7 ± 3.3212 .

[Quizzes]

1.8.2 Small sample CIs for population variance σ^2

Population variance σ^2 quantifies the amount of **variability** in the population. We have already shown that if we observe the data sample (X_1, \dots, X_n) , then σ^2 can be estimated by an **unbiased point** estimator

$$S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2.$$

How do we get a confidence interval for σ^2 ?

If the sample is large then S^2 is approximately normally distributed. It is not difficult to derive a formula for $\text{Var}(S^2)$ and develop an estimator for this variance. Then the standard method for large sample confidence intervals works. In practice, however, we are usually interested in confidence intervals for σ^2 when we have a small data sample.

So assume that the data sample is small. In this case we must restrict ourself to the situation when the data is normally distributed.

Assume that all sample data points $X_1, X_2, \dots, X_n \sim N(\mu, \sigma^2)$.

The pivotal quantity (see Theorem 1.8.2) is

$$T = \frac{\sum_{i=1}^n (X_i - \bar{X})^2}{\sigma^2} = \frac{(n-1)S^2}{\sigma^2} \sim \chi_{(n-1)}^2$$

We need to find L and U so that

$$\mathbb{P}(L \leq \frac{(n-1)S^2}{\sigma^2} \leq U) = 1 - \alpha$$

A usual choice is $L = \chi_{1-(\alpha/2)}^2$ and $U = \chi_{\alpha/2}^2$, both corresponding to $(n-1)$ d.f.

Hence,

$$\begin{aligned} \mathbb{P}\left(\chi_{1-(\alpha/2)}^2(n-1) \leq \frac{(n-1)S^2}{\sigma^2} \leq \chi_{\alpha/2}^2(n-1)\right) &= 1 - \alpha \\ \Leftrightarrow \mathbb{P}\left(\frac{(n-1)S^2}{\chi_{\alpha/2}^2(n-1)} \leq \sigma^2 \leq \frac{(n-1)S^2}{\chi_{1-\alpha/2}^2(n-1)}\right) &= 1 - \alpha \end{aligned}$$

- Suppose we want a one-sided bound for σ^2 , say **lower bound**. We will want to make use of the pivotal quantity:

$$\frac{(n-1)S^2}{\sigma^2} \sim \chi_{(n-1)}^2$$

- Note that if $b = \chi_{\alpha}^2$ for $(n-1)$ d.f., then

$$\mathbb{P}\left(\frac{(n-1)S^2}{\sigma^2} \leq b\right) = 1 - \alpha$$

- Then we have, with $100(1 - \alpha)\%$ probability,

$$\frac{(n-1)S^2}{\chi_{\alpha, n-1}^2} \leq \sigma^2.$$

Hence $\frac{(n-1)S^2}{\chi_{\alpha, n-1}^2}$ is a $(1 - \alpha)$ confidence lower bound for σ^2 .

Similarly, $\frac{(n-1)S^2}{\chi_{1-\alpha, n-1}^2}$ is a $(1 - \alpha)$ confidence **upper bound** for σ^2 .

What if we want to build a confidence interval for the standard deviation $\sigma = \sqrt{\sigma^2}$, instead of σ^2 ? This is simple:

- Since we know that

$$\mathbb{P}\left(\frac{(n-1)S^2}{\chi_{\alpha/2}^2} \leq \sigma^2 \leq \frac{(n-1)S^2}{\chi_{1-\alpha/2}^2}\right) = 1 - \alpha$$

- This is equivalent to

$$\mathbb{P}\left(\sqrt{\frac{(n-1)S^2}{\chi_{\alpha/2}^2}} \leq \sigma \leq \sqrt{\frac{(n-1)S^2}{\chi_{1-\alpha/2}^2}}\right) = 1 - \alpha$$

- Therefore we can easily obtain a C.I. for σ .

Example 1.8.7. Suppose that you wished to describe the variability of the carapace lengths of this population of lobsters. Find a 90% confidence interval for the population variance σ^2 .

Lobster Field Number	A061	A062	A066	A070	A067	A069	A064	A068	A065	A063
Carapace Length (mm)	78	66	65	63	60	60	58	56	52	50

```

> x = c(78,66,65,63,60,60,58,56,52,50)
> x
[1] 78 66 65 63 60 60 58 56 52 50
> sum(( x-mean(x) )^2) # the numerator of sample variance and the CI LB and UB
[1] 571.6
> # and then we calculate the denominators of the CI LB and UB
> qchisq(0.05,9)
[1] 3.325113
> qchisq(0.95,9)
[1] 16.91898

```

The answer is $(571.6/16.91898, 571.6/3.325113)$. Note that $\chi^2_{0.95,9} = 3.325113 = \text{qchisq}(0.05,9)$ and $\chi^2_{0.05,9} = 16.91898 = \text{qchisq}(0.95,9)$.

Both can also be obtained from Table 6.

Chapter 2

Properties of Point Estimators and Methods of Estimation

2.1 Relative efficiency

Suppose that two estimators are both unbiased. Which one is better? Usually, an estimator with smaller variance is preferable. The quantitative measure used to compare the estimators is relative efficiency.

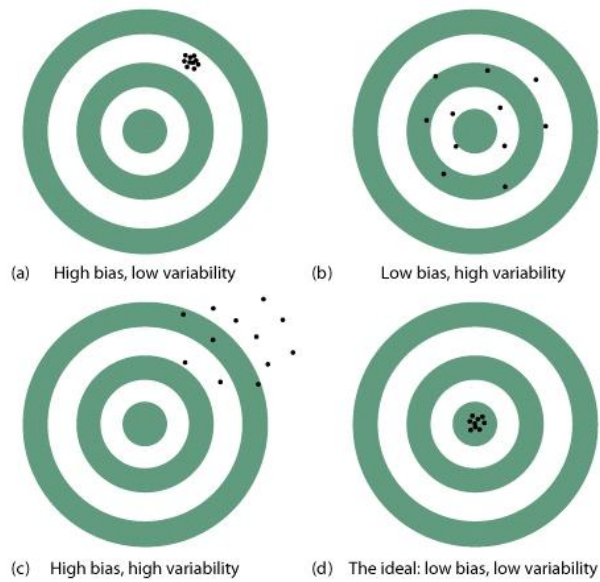
Definition 2.1.1. Given two unbiased estimators $\hat{\theta}_1$ and $\hat{\theta}_0$ of the same parameter θ , the relative efficiency of $\hat{\theta}_1$ relative to $\hat{\theta}_0$ is defined to be the ratio of their variances

$$eff(\hat{\theta}_1, \hat{\theta}_0) = \frac{Var(\hat{\theta}_0)}{Var(\hat{\theta}_1)}.$$

We can think about $\hat{\theta}_0$ as a reference estimator. The estimator $\hat{\theta}_1$ with relative efficiency which is greater than 1 is better than $\hat{\theta}_0$ since its variance is smaller and so $\hat{\theta}_1$ is a more accurate estimator of θ than $\hat{\theta}_0$! Note that we assumed from outset that both estimators are unbiased.

Quiz:

- Relative efficiency of (a) to (b):



- Relative efficiency of (c) to (d):
- Relative efficiency of (b) to (d):
- Relative efficiency of (d) to (b):
- Simple example: $Y_1, \dots, Y_9 \sim N(\mu_Y, 1)$. Want to estimate μ_Y .

1. $\hat{\theta}_1 = \frac{1}{9} \sum_{i=1}^9 Y_i$,
2. $\hat{\theta}_2 = \frac{1}{2}(Y_1 + Y_2)$.

- Both $\hat{\theta}_1$ and $\hat{\theta}_2$ are unbiased estimator for θ !

$$- \mathbb{E}\hat{\theta}_1 = \theta \text{ and } \mathbb{E}\hat{\theta}_2 = \theta ;$$

- $\text{Var}(\hat{\theta}_1) = \frac{1}{9}\text{Var}(Y_1) = \frac{1}{9}$;
- $\text{Var}(\hat{\theta}_2) = \frac{1}{2}\text{Var}(Y_1) = \frac{1}{2}$;
- $\text{Ref}(\hat{\theta}_1, \hat{\theta}_2) = \frac{9}{2} > 1$;
- $\hat{\theta}_1$ is better than $\hat{\theta}_2$; why?

- because $\hat{\theta}_2$ did not use all information available, and hence is less efficient.

Example 2.1.2. Let Y_1, Y_2, \dots, Y_n denote a random sample from the uniform distribution on the interval $[0, \theta]$. Two unbiased estimators for θ are

$$\hat{\theta}_1 = 2\bar{Y} \quad \text{and} \quad \hat{\theta}_2 = \frac{n+1}{n}Y_{(n)},$$

where $Y_{(n)} = \max\{Y_1, Y_2, \dots, Y_n\}$. Find the efficiency of $\hat{\theta}_1$ relative to $\hat{\theta}_2$.

2.2 Consistency

Suppose again that we have a sample (X_1, \dots, X_n) from a probability distribution that depends on parameter θ . Note that although we speak about *an* estimator $\hat{\theta} = \hat{\theta}(X_1, \dots, X_n)$, in fact the distribution of the estimator depends on n , so it would be more correctly speak about a sequence of random variables $\hat{\theta}_n$.

Usually, we expect that when the size of the sample becomes larger, that is, n grows, the distribution of the estimator $\hat{\theta}_n$ become concentrated more and more around the true value of the parameter θ . Technically this property called *consistency* and we are giving its mathematical definition below.

Plots show a simulation study. A sample X_1, X_2, \dots from the distribution $N(\theta, 1/4)$ was generated with $\theta = 10$ and we computed $\hat{\theta}_k = (X_1 + \dots + X_k)/k$. Figure 2.1 shows a path of $\hat{\theta}_k$. It suggests that if we get more and more

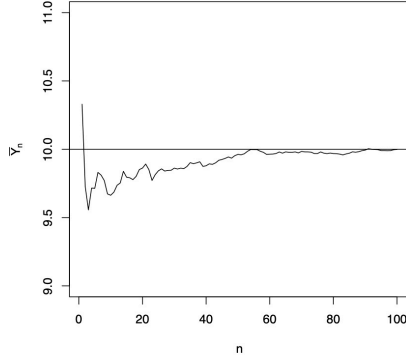


Figure 2.1

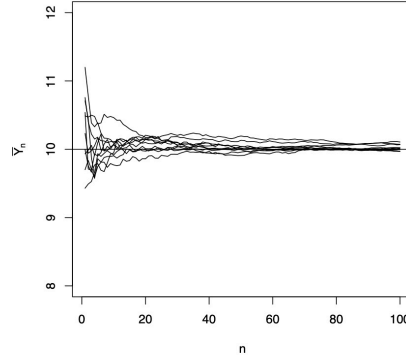


Figure 2.2

data, $\hat{\theta}_k$ converges to the true value of θ . In fact, this is a consequence of the strong law of large numbers, which says that this behavior is observed with probability 1.

What about several different samples? Figure 2.2 shows the situation when the sample X_1, X_2, \dots was generated 10 times and 10 paths of $\hat{\theta}_k$ were plotted. This picture suggest that when the sample size grows the distribution of $\hat{\theta}_k$ around the true value of the parameter θ . Mathematically this is a consequence of the weak law of large numbers.

Definition 2.2.1 (Convergence in probability). A sequence of random variables, $X_1, X_2, \dots, X_n, \dots$, is **convergent in probability** to a random variable X if, for any $\epsilon > 0$, as $n \rightarrow \infty$

$$\mathbb{P}(|X_n - X| < \epsilon) \rightarrow 1,$$

that is,

$$\lim_{n \rightarrow \infty} \mathbb{P}(|X_n - X| < \epsilon) = 1$$

This is denoted either as $X_n \xrightarrow{\mathbb{P}} X$ or as $\text{plim}_{n \rightarrow \infty} X_n = X$.

Note that $a_n \equiv \mathbb{P}(|X_n - X| < \epsilon)$ is simply a number (it is not random). Hence, $\{a_1, a_2, \dots, a_n, \dots\}$ form a sequence of numbers, and their limit is defined in the usual “calculus” sense.

Definition 2.2.2 (Consistency). An estimator $\hat{\theta}_n$ is a **consistent estimator** of θ , if $\hat{\theta}_n$ converges in probability to θ

$$\hat{\theta}_n \xrightarrow{\mathbb{P}} \theta.$$

By writing out the definition of the convergence in probability in detail, we see that this definition can be also written as saying that an estimator $\hat{\theta}_n$ is a **consistent estimator** of θ , if for any $\varepsilon > 0$, as $n \rightarrow \infty$,

$$\mathbb{P}(|\hat{\theta}_n - \theta| < \varepsilon) \rightarrow 1.$$

The consistency of the estimator means that as the sample size goes to infinity, we become more and more sure that the distance between $\hat{\theta}_n$ and θ is smaller **smaller than any positive ε** !

Consistency describes a property of the estimator in the $n \rightarrow \infty$ limit. It is NOT meant to describe the property of the estimator for a fixed n .

An unbiased estimator can be inconsistent, and a biased estimator can be consistent! Consistency is more important than unbiasedness because it ensures that if collect enough data we will eventually learn the true value of the parameter.

So, how can I tell if the estimator is consistent? One way is to see how MSE changes with n .

Theorem 2.2.3. *If $MSE(\hat{\theta}_n) \rightarrow 0$ as $n \rightarrow \infty$, then the estimator $\hat{\theta}_n$ is consistent.*

Proof. $MSE(\hat{\theta}_n) \rightarrow 0$ if and only if $bias(\hat{\theta}_n) \rightarrow 0$ and $\text{Var}(\hat{\theta}_n) \rightarrow 0$. Fix an $\varepsilon > 0$ and choose n_0 so that $|bias(\hat{\theta}_n)| < \varepsilon/2$ for all $n > n_0$. Then $\mathbb{P}(|\hat{\theta}_n - \theta| > \varepsilon) \leq \mathbb{P}(|\hat{\theta}_n - \mathbb{E}\hat{\theta}_n| > \varepsilon/2)$ for all $n > n_0$, – because $|\mathbb{E}\hat{\theta}_n - \theta| < \varepsilon/2$ for these n and therefore the event $|\hat{\theta}_n - \theta| > \varepsilon$ can occur only if $|\hat{\theta}_n - \mathbb{E}\hat{\theta}_n| > \varepsilon/2$ occurred. Now apply the Chebyshev inequality,

$$\mathbb{P}(|\hat{\theta}_n - \mathbb{E}\hat{\theta}_n| > \varepsilon/2) \leq \frac{\text{Var}(\hat{\theta}_n)}{(\varepsilon/2)^2}$$

By our assumption it can be made arbitrarily small for all sufficiently large n because $\text{Var}(\hat{\theta}_n) \rightarrow 0$. We showed that $\mathbb{P}(|\hat{\theta}_n - \theta| > \varepsilon) \rightarrow 0$ for any $\varepsilon > 0$. \square

- If $Bias(\hat{\theta}_n) \rightarrow 0$ as $n \rightarrow \infty$, then the estimator is called **asymptotically unbiased**.
- The theorem says that any estimator which is asymptotically unbiased and has its variance converging to 0 as $n \rightarrow \infty$ is a consistent estimator.

Example 2.2.4 (Sample mean is a consistent estimator of the population mean). Let Y_1, Y_2, \dots be a sample from a population with mean μ and variance σ^2 .

- Sample mean $\bar{Y}_n = \frac{1}{n} \sum_{i=1}^n Y_i$. Its expectation is μ and its variance is $\text{Var}(\bar{Y}_n) = \sigma^2/n \rightarrow 0$.
- So \bar{Y}_n is an **unbiased and consistent** estimator of μ .

Example 2.2.5 (Biased and consistent estimator of the mean). For parameter $\theta = \mu$, consider a modified sample mean $\hat{\theta}_n = \frac{n}{n-1} \bar{Y}$

- $Bias(\hat{\theta}) = \mathbb{E}\hat{\theta}_n - \theta = \frac{n}{n-1}\mu - \mu = \frac{1}{n-1}\mu \rightarrow 0$ as $n \rightarrow \infty$; (**$\hat{\theta}_n$ is a biased estimator of $\mu \neq 0$ for every n . It is, however, asymptotically unbiased**)
- $\text{Var}(\hat{\theta}) = \frac{n^2}{(n-1)^2} \sigma^2/n \rightarrow 0$.
- **Conclusion:** $\hat{\theta}_n$ is a biased but consistent estimator.

Example 2.2.6. • $Y_i \sim \text{Unif}[\theta, \theta + 1]$

- $\hat{\theta}_1 = \bar{Y} - 1/2$; $\hat{\theta}_2 = Y_{(n)} - n/(n+1)$
 - We show in Ex 9.3 that they are both unbiased. Want to show that they are consistent.
 - ONLY NEED TO SHOW THAT BOTH VARIANCES GOES TO 0 AS $n \rightarrow \infty$.
 - The variance of $\hat{\theta}_1$ is $\frac{\sigma_{Y_1}^2}{n}$ while $\sigma_{Y_1}^2 = 1/12$.
 - The variance of $\hat{\theta}_2$ is the same as the variance of $T \equiv Y_{(n)}$.
- What's the PDF or CDF of T ?

Unbiasedness, efficiency, consistency

- Unbiasedness: concerns expectation; for fixed n
 - Efficiency: concerns the variance; for fixed n ; **meaningful only for unbiased estimators.**
 - Consistency:
 - only care about $n \rightarrow \infty$;
 - concerns bias and variance (and whether they vanish for large n);
 - however, **does not necessarily imply unbiasedness for finite n .**
1. Can biased estimator be consistent? Yes!
 2. Can unbiased estimator be inconsistent? Yes!
 3. Can an unbiased **consistent** estimator be less efficient than an unbiased **inconsistent** estimator? Yes!
 - Up to a certain point.

[Quizzes]

In practice, it is often difficult to calculate the variance of the estimator. So it is useful that in some cases consistency can be established without actually calculating the variance.

Recall that the consistency is the convergence of the estimator to the true value of the parameter in the probability. So here are some properties of the convergence in probability.

If there are two estimator sequences, $\hat{\theta}_n \xrightarrow{p} \theta$ and $\hat{\theta}'_n \xrightarrow{p} \theta'$, then:

1. $\hat{\theta}_n + \hat{\theta}'_n \xrightarrow{p} \theta + \theta'$;
2. $\hat{\theta}_n \times \hat{\theta}'_n \xrightarrow{p} \theta \times \theta'$;
3. $\hat{\theta}_n / \hat{\theta}'_n \xrightarrow{p} \theta / \theta'$ provided that $\theta' \neq 0$;
4. For any continuous function $g(u)$, $g(\hat{\theta}_n) \xrightarrow{p} g(\theta)$;

5. For any continuous function $g(u, v)$, $g(\hat{\theta}_n, \hat{\theta}'_n) \xrightarrow{p} g(\theta, \theta')$;
6. For a sequence of numbers $\{a_n, n = 1, \dots\}$, $a_n \rightarrow a$ (in the calculus sense) implies that $a_n \xrightarrow{\mathbb{P}} a$ (a_n 's are viewed as special random variables).

Properties 4 and 5 are called the continuous mapping theorem for the convergence in probability.

We omit the proof.

We mention here also another very useful result which is called Slutsky's theorem. We will not need it for the proof of consistency but it is behind all the proofs that we had for the confidence intervals. If X_n converges *in distribution* to a variable X , and Y_n converges *in probability* to a constant c , then

- $X_n + Y_n \rightarrow X + c$ in distribution;
- $X_n Y_n \rightarrow cX$ in distribution;
- $X_n / Y_n \rightarrow X / c$ in distribution, provided that $c \neq 0$.

Example 2.2.7 (S^2 is a consistent estimator of σ^2). By definition

$$S_n^2 = \frac{1}{n-1} \left(\sum_{i=1}^n Y_i^2 - n\bar{Y}^2 \right) = \frac{n}{n-1} \left(\frac{1}{n} \sum_{i=1}^n Y_i^2 - \bar{Y}^2 \right)$$

We want to show that

$$S_n^2 \xrightarrow{p} \sigma^2$$

- By our previous result about the means of independent variables, we have (1):

$$\frac{1}{n} \sum_{i=1}^n Y_i^2 \xrightarrow{p} \mathbb{E}(Y_i^2) = \sigma^2 + \mu^2,$$

- and (2) $\bar{Y} \xrightarrow{p} \mathbb{E}(Y_1) = \mu$. (Note that \bar{Y} is actually a sequence of random variables that depends on the sample size n . This dependence is suppressed in the notation but we should remember about it to understand what it means that $\bar{Y} \xrightarrow{\mathbb{P}} \mathbb{E}(Y_i)$.)

- By an application of [continuous mapping theorem](#) with $g(u) = u^2$, we have (3):

$$(\bar{Y})^2 \xrightarrow{p} \mu^2$$

- Consider $g(u, v) = u - v$ and apply the continuous mapping theorem again,

$$\frac{1}{n} \sum_{i=1}^n Y_i^2 - (\bar{Y}_n)^2 \xrightarrow{p} \sigma^2 + \mu^2 - \mu^2 = \sigma^2$$

Recall that

$$S_n^2 = \frac{n}{n-1} \left(\frac{1}{n} \sum_{i=1}^n Y_i^2 - \bar{Y}_n^2 \right),$$

and we have just shown that

$$\frac{1}{n} \sum_{i=1}^n Y_i^2 - (\bar{Y}_n)^2 \xrightarrow{p} \sigma^2$$

It remains to notice that

$$\frac{n}{n-1} \rightarrow 1 \text{ implies } \frac{n}{n-1} \xrightarrow{p} 1,$$

and another application of the continuous mapping theorem (with function $g(u, v) = uv$, and variables $u = n/(n-1)$ and $v = \left(\frac{1}{n} \sum_{i=1}^n Y_i^2 - \bar{Y}_n^2 \right)$ shows that

$$S_n^2 = \frac{n}{n-1} \left(\frac{1}{n} \sum_{i=1}^n Y_i^2 - \bar{Y}_n^2 \right) \xrightarrow{p} 1 \times \sigma^2 = \sigma^2.$$

Example 2.2.8 (Another estimator of σ^2). Another estimator of σ can be defined as

$$S_0^2 \equiv \frac{1}{n} \left(\sum_{i=1}^n Y_i^2 - n\bar{Y}_n^2 \right)$$

The denominator of S^2 is $(n-1)$ while that of his brother is n . Being $(n-1)$ makes the estimator S^2 unbiased (as an estimator of σ^2). Is S_0^2 a consistent estimator of σ^2 as well? Yes! (By continuous mapping theorem.)

Note that S_0^2 is biased. But it is still a consistent estimator of σ^2 .

Example 2.2.9. • $Y_i \sim N(\mu, \sigma^2)$

- $\hat{\sigma}_2^2 = \frac{1}{2}(Y_1 - Y_2)^2$
- You probably should have called it

$$\hat{\sigma}_{2,n}^2$$

since we want to talk about consistency here.

- But something is wrong with such a notation – is n really necessary?
- $\hat{\sigma}_{2,n}^2$ is not consistent because however many observations we have in the sample, $\hat{\sigma}_{2,n}^2$ is always based on the first two observations only, and hence more observations do not add anything (will not change the variance, let alone making it smaller).

Example 2.2.10. • Two samples. Possibly different means. Same n . Suppose they share the same variance σ^2

- Show that the pooled sample variance is a consistent estimator of σ^2

Example 2.2.11. • $Y_1, \dots, Y_n \sim N(\mu, \sigma^2)$

- Assume that $n = 2k$ is even.
- Let $\hat{\sigma}^2 = \frac{1}{2k} \sum_{i=1}^k (Y_{2i} - Y_{2i-1})^2$

(a) Show that it is unbiased

- Note that $Y_{2i} - Y_{2i-1} \sim N(0, 2\sigma^2)$
- $\mathbb{E}[(Y_{2i} - Y_{2i-1})^2] = \text{Var}(Y_{2i} - Y_{2i-1}) + [\mathbb{E}(Y_{2i} - Y_{2i-1})]^2 = 2\sigma^2 + 0^2 = 2\sigma^2$
- $\mathbb{E}\hat{\sigma}^2 = \mathbb{E}\left(\frac{1}{2k} \sum_{i=1}^k (Y_{2i} - Y_{2i-1})^2\right) = \frac{1}{2k} \sum_{i=1}^k \mathbb{E}[(Y_{2i} - Y_{2i-1})^2] = \frac{1}{2k} \cdot k \cdot 2\sigma^2 = \sigma^2$

(b) Show that it is consistent.

- Let $X_i = (Y_{2i} - Y_{2i-1})^2$, then $\hat{\sigma}^2 = \frac{1}{2} \cdot \frac{1}{k} \sum_{i=1}^k X_i = \frac{1}{2} \bar{X}_k$
- $\bar{X}_k \xrightarrow{p} \mathbb{E}X_1 = 2\sigma^2$ (See 2nd bullet pt. above)
- $\frac{1}{2} \bar{X}_k \xrightarrow{p} \frac{1}{2} 2\sigma^2 = \sigma^2$

- What if Y_1, \dots, Y_n are not normal?

[Quizzes]

Is consistency really important?

Yes. If an estimator is not consistent, then it will not produce the correct estimation even if we are given the luxury of getting unlimited amount of data for free. It's a shame if one cannot get the correct answer in this situation. **An inconsistent estimator is a waste of time.**

Does consistency guarantee good performance?

Not necessarily. We still live in a finite sample world. Something that is **ultimately** good for very large sample, may not be good enough for a realistic sample size.

2.3 Sufficient statistics

There is a huge multitude of functions of the data that we can consider in the search for a good estimator. So it is worthwhile to check if we can reduce the data to one or just a few of summary statistics. This is the main idea behind the concept of sufficient statistics.

Definition 2.3.1. Let X_1, X_2, \dots, X_n denote a random sample from a distribution with unknown parameter θ . A statistic $T = T(X_1, X_2, \dots, X_n)$ is said to be **sufficient for θ** if the **conditional distribution** of (X_1, X_2, \dots, X_n) , given $T = t$, does not depend on θ . That is,

$$\Pr_{\theta}(X_1, \dots, X_n | T(X_1, \dots, X_n) = t)$$

depends only on t but does not depend on θ .

Intuitively, if we know that $T = t$, then revealing the complete information about X_1, X_2, \dots, X_n does not give us any additional information about θ .

In principle, a sufficient statistic can be a **vector**, that is, it can consist of several functions. For example, if you take a vector of order statistics, $T = (X_{(1)}, X_{(2)}, \dots, X_{(n)})$, then it is a sufficient statistic. However, we don't

gain very much by considering such statistics since they do not reduce the data.

If we take a function of a vector of sufficient statistics and reduce the dimension, then it can potentially break the sufficiency, however in some cases the resulting function is still sufficient. For example any invertible function of a sufficient statistic is sufficient, – it does not lose any information.

A sufficient statistic is **minimal** if it can be written as a function of any other of sufficient statistics. (A minimal sufficient statistic exists under mild conditions on the distribution of the data but there are some counterexamples.)

Why do we care about sufficient statistics?

In some cases, we can find a good estimator of a parameter θ by a 2-step procedure:

1. Find a sufficient statistic $T(X_1, X_2, \dots, X_n)$ for parameter θ . Informally, it contains all information about the parameter which is available in the data.
2. Find an unbiased estimator of θ , which is a function of T . We can hope that all relevant information in data was used and the estimator cannot be further improved.

How we can find a good sufficient statistics? We should try to factorize the **likelihood function**.

Recall that the **likelihood function** is the *same thing* as the **joint distribution density** or **joint distribution pmf** of data X_1, \dots, X_n . In this course we consider only the situation when X_1, \dots, X_n are independent and identically distributed, so

1. If X_1, \dots, X_n are discrete random variables, and $p(x) := \mathbb{P}(X = x|\theta)$ is the probability mass function of each of them, then

$$L(\theta \mid x_1, \dots, x_n) = p_{X_1, \dots, X_n}(x_1, \dots, x_n \mid \theta) = \prod_{i=1}^n p(x_i \mid \theta).$$

2. If X_1, \dots, X_n are continuous random variables and $f(x|\theta)$ is the density of each of them, then

$$L(\theta | x_1, \dots, x_n) = f_{X_1, \dots, X_n}(x_1, \dots, x_n | \theta) = \prod_{i=1}^n f(x_i | \theta).$$

To simplify the notation, we sometimes write $L(\theta)$ instead of $L(\theta | x_1, \dots, x_n)$.

Theorem 2.3.2 (Fisher's Factorization Criterion). *A statistic T is a sufficient statistic for parameter θ if and only if $L(\theta)$ can be factorized into two nonnegative functions:*

$$L(\theta | x_1, \dots, x_n) = g(\theta, t(x_1, \dots, x_n)) \times h(x_1, \dots, x_n)$$

Here $g(\theta, t)$ is a function only of t (the observed value of T) and θ and the function $h(x_1, \dots, x_n)$ does not depend on θ at all.

Example 2.3.3. • $Y_1, \dots, Y_n \sim B(p)$. Find a sufficient statistic for p .

- Likelihood:

$$\begin{aligned} L(p) &= \prod_{i=1}^n \{p^{y_i} (1-p)^{(1-y_i)}\} \\ &= p^{\sum_{i=1}^n y_i} (1-p)^{(n - \sum_{i=1}^n y_i)} \\ &= \left\{ \frac{p}{1-p} \right\}^{\sum_{i=1}^n y_i} (1-p)^n \times 1 \end{aligned}$$

- We can define a statistic $T = \sum_{i=1}^n Y_i$. Then we will have

- $g_p(t) = \left\{ \frac{p}{1-p} \right\}^t (1-p)^n$, and $h(y_1, \dots, y_n) = 1$
- The first term only depends on p and T (or t)
- The second term does not depend on p

Example 2.3.4. • $Y_i \sim_{iid} \text{Poisson}(\lambda)$, i.e. $\sim p(y) = e^{-\lambda} \frac{\lambda^y}{y!}$

- Likelihood (use independence):

$$L(\lambda) = \prod_{i=1}^n e^{-\lambda} \frac{\lambda^{y_i}}{y_i!}$$

•

$$L(\lambda) = e^{-n\lambda} \frac{\lambda^{\sum_{i=1}^n y_i}}{\prod_{i=1}^n (y_i!)} = e^{-n\lambda} \lambda^{\sum_{i=1}^n y_i} \times \frac{1}{\prod_{i=1}^n (y_i!)}$$

Example 2.3.5. Let X_1, X_2, \dots, X_n be a random sample in which X_i is exponentially distributed (remember life of smartphones?) with parameter θ . Find a sufficient statistic for θ .

Example 2.3.6. • $Y_i \sim_{iid} N(\mu, \sigma^2)$, i.e.

$$\sim f(y) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{(y-\mu)^2}{2\sigma^2}\right\}$$

- Likelihood:

$$\begin{aligned} L(\cdot) &= \prod_{i=1}^n \left[\frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{(y_i-\mu)^2}{2\sigma^2}\right\} \right] \\ &= (2\pi\sigma^2)^{-\frac{n}{2}} \exp\left\{-\frac{\sum_{i=1}^n (y_i-\mu)^2}{2\sigma^2}\right\} \end{aligned}$$

- Note that $\sum_{i=1}^n (y_i - \mu)^2 = \sum_{i=1}^n (y_i - \bar{y} + \bar{y} - \mu)^2 = \sum_{i=1}^n (y_i - \bar{y})^2 + n(\bar{y} - \mu)^2 = (n-1)s^2 + n(\bar{y} - \mu)^2$

- Thus we have

$$L(\cdot) = (2\pi\sigma^2)^{-\frac{n}{2}} \cdot \exp\left\{-\frac{(n-1)s^2}{2\sigma^2}\right\} \cdot \exp\left\{-\frac{n(\bar{y} - \mu)^2}{2\sigma^2}\right\}$$

$$L(\cdot) = (2\pi\sigma^2)^{-\frac{n}{2}} \cdot \exp\left\{-\frac{(n-1)s^2}{2\sigma^2}\right\} \cdot \exp\left\{-\frac{n(\bar{y} - \mu)^2}{2\sigma^2}\right\}$$

- The argument in $L(\cdot)$ is not specified because there are two situations.

1. μ is unknown and σ^2 is known $\Rightarrow L(\mu)$

$$L(\cdot) = \exp \left\{ -\frac{n(\bar{y} - \mu)^2}{2\sigma^2} \right\} \cdot (2\pi\sigma^2)^{-\frac{n}{2}} \cdot \exp \left\{ -\frac{(n-1)s^2}{2\sigma^2} \right\}$$

\bar{Y} is a sufficient statistic for μ

2. Both μ and σ^2 are unknown $\Rightarrow L(\mu, \sigma^2)$

$$L(\cdot) = (2\pi\sigma^2)^{-\frac{n}{2}} \cdot \exp \left\{ -\frac{(n-1)s^2}{2\sigma^2} \right\} \cdot \exp \left\{ -\frac{n(\bar{y} - \mu)^2}{2\sigma^2} \right\} \times 1$$

The pair (\bar{Y}, S^2) is a sufficient statistic for μ and σ^2

Quiz 2.3.7. Every function of a sufficient statistics is a sufficient statistic.

A. True

B. False

Quiz 2.3.8. Every strictly decreasing function of a sufficient statistics is a sufficient statistic.

A. True

B. False

Even a minimal sufficient statistic is not unique!

Example 2.3.9 (Uniform distribution). Let X_1, X_2, \dots, X_n be uniformly distributed in $(0, \theta)$. What is a sufficient statistic for θ ?

- Density of one random variable:

$$\begin{aligned} f_{X_i}(x_i) &= \begin{cases} 1/\theta, & 0 < x_i < \theta, \\ 0, & \text{otherwise} \end{cases} \\ &= \frac{1}{\theta} \mathbb{1}_{0 < x_i < \theta}, \end{aligned}$$

where $\mathbb{1}_{0 < x_i < \theta}$ is the indicator function of the event $\{0 < x_i < \theta\}$.

- Likelihood:

$$L(\theta) = \prod_{i=1}^n \left[\mathbb{1}_{0 \leq x_i \leq \theta \frac{1}{\theta}} \right]$$

- Note that $\prod_{i=1}^n \mathbb{1}_{\theta \geq x_i} = \mathbb{1}_{\theta \geq x_{(n)}}$. Not very obvious. Think!
- Thus the likelihood $L(\theta) = \mathbb{1}_{\theta \geq x_{(n)}} \frac{1}{\theta^n} = \mathbb{1}_{\theta \geq x_{(n)}} \frac{1}{\theta^n} \times 1$
- Hence $T = X_{(n)}$ is a sufficient statistic for θ .
- **Any 1-to-1 function of a sufficient statistic is also a sufficient statistic for the same parameter.** \Rightarrow the unbiased estimator $\hat{\theta}^{\frac{n+1}{n}} X_{(n)}$ is also a sufficient statistic for θ .
- Note that previously that this estimator is much more efficient than another unbiased estimator $2\bar{X}$. This is a reflection of a general fact which is called the Rao theorem.

Exercise 2.3.10. • $f(y) = \exp[-(y - \theta)]$ for $y > \theta$

- Sufficient statistic?

2.4 Rao-Blackwell Theorem and Minimum-Variance Unbiased Estimator

This section relates sufficiency with efficiency and unbiasedness. It will tell us why a sufficient statistic is very useful in statistical inference.

- We have learned two good qualities of an estimator $\hat{\theta}$ for a parameter θ :
 - Unbiasenss: $E\hat{\theta} = \theta$;
 - Low variance: $Var(\hat{\theta})$ is small;

- Relative efficiency of two estimators $\hat{\theta}_1$ and $\hat{\theta}_2$

$$\text{Reff}(\hat{\theta}_1, \hat{\theta}_2) = \frac{\text{Var}(\hat{\theta}_2)}{\text{Var}(\hat{\theta}_1)};$$

whichever has the smaller variance is more efficient (better)!

Definition 2.4.1. An *MVUE* estimator is an unbiased estimator that has the smallest variance for each value of the parameter θ . **MVUE = Minimal Variance Unbiased Estimator.**

How can we find a MVUE estimator?

The main idea of the following theorem is that there is an MVUE estimator which is a function of a sufficient statistic.

Theorem 2.4.2 (Rao-Blackwell Theorem). *Let $\hat{\theta}$ be an unbiased estimator for θ such that $\text{Var}(\hat{\theta}) < \infty$. If T is a sufficient statistic for θ , define $\hat{\theta}^* = \mathbb{E}(\hat{\theta}|T)$. Then, for all θ , $\mathbb{E}\hat{\theta}^* = \theta$ and $\text{Var}(\hat{\theta}^*) \leq \text{Var}(\hat{\theta})$.*

- Given an unbiased estimator $\hat{\theta}$ and a sufficient statistic T , we can find a modified estimator $\hat{\theta}^*$, which is improved in the sense that

- $\hat{\theta}^*$ is still unbiased;
- $\hat{\theta}^*$ has a smaller (or at least no larger) variance than $\hat{\theta}$;

- Remarks:

- $\hat{\theta}^*$ is a function of T
- $\hat{\theta}^*$ is random
- If $\hat{\theta}$ is already a function of T , then $\mathbb{E}(\hat{\theta} | T) = \hat{\theta}$, i.e., taking the conditional expectation does not change anything, in particular it does not improve the efficiency.

Proof. Because T is sufficient for θ , the conditional distribution of any statistic (including $\hat{\theta}$), given T , does not depend on θ . Thus, $\hat{\theta}^* = \mathbb{E}(\hat{\theta}|T)$ is not a function of θ and is therefore a statistic. The fact that $\hat{\theta}^*$ is almost obvious from the law of repeated expectation.

$$\mathbb{E}\hat{\theta}^* = \mathbb{E}\left[\mathbb{E}(\hat{\theta}|T)\right] = \mathbb{E}\hat{\theta} = \theta.$$

For the variance we use another theorem about conditional expectations:

$$\begin{aligned}\text{Var}(\hat{\theta}) &= \text{Var}\left[\mathbb{E}(\hat{\theta}|T)\right] + \mathbb{E}\left[\text{Var}(\hat{\theta}|T)\right] \\ &= \text{Var}(\hat{\theta}^*) + \mathbb{E}\left[\text{Var}(\hat{\theta}|T)\right].\end{aligned}$$

Since the second term is non-negative, we find that $\text{Var}(\hat{\theta}^*) < \text{Var}(\hat{\theta})$. \square

This theorem implies that if an unbiased estimator $\hat{\theta}$ is NOT a function of a sufficient statistics T then we can find another unbiased estimator $\hat{\theta}^* = \mathbb{E}(\hat{\theta}|T)$ which is a function of T at least as good as $\hat{\theta}$. However, it does NOT imply that this estimator is an MVUE. Perhaps we can find another sufficient statistic T_2 and improve $\hat{\theta}^*$ by taking a conditional expectation with respect to T_2 .

A natural conjecture is that if T is a *minimal* sufficient statistic and some function $\hat{\theta} = \hat{\theta}(T)$ is an unbiased estimator of θ , then $\hat{\theta}$ is an MVUE. This is again not quite correct. One has to impose a stronger requirement that T is a *complete minimal* sufficient statistic. In this case a function $\hat{\theta} = \hat{\theta}(T)$ which is unbiased for θ is indeed MVUE.

This raises questions about what is a complete sufficient statistic and how to check that a sufficient statistic is minimal and complete. We will not be concerned with this question here and simply promise that in all our examples the sufficient statistics obtained by factorization theorem will be minimal and sufficient.

Routine to find the MVUE

1. Factorize the likelihood function and find a (minimal) sufficient statistic T ;
2. find a function of T which is unbiased for the parameter of interest θ ;

Example 2.4.3 (Exponential). Suppose X_1, X_2, \dots, X_n all from the exponential distribution with the parameter θ . What is an MVUE for θ ?

We showed that $T = (X_1 + \dots + X_n)$ is a sufficient statistic. In fact it is a minimal and complete statistic, and since $\bar{X} = T/n$ is unbiased for θ hence it is an MVUE.

Example 2.4.4 (Bernoulli). Let X_i 's are iid Bernoulli with parameter p . What is an MVUE for p ?

Same argument works as in the previous example. We already proved that $T = \sum_{i=1}^n X_i$ is a sufficient statistic for p . Hence $\hat{p} = \bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$ is MVUE (for p).

Example 2.4.5 (Normal). X_1, \dots, X_n are i.i.d. from normal distribution $N(\mu, \sigma^2)$. What is the MVUE for μ and σ^2 .

We have shown that \bar{X} and S^2 are joint sufficient statistics for μ and σ^2 . In addition, we know that \bar{X} and S^2 are unbiased estimators of μ and σ . Therefore, \bar{X} and S^2 as the MVUEs for μ and σ^2 , respectively.

[Quizzes]

2.5 Method of Moments Estimation

Suppose that we are looking for a good estimator $\hat{\theta}(X_1, \dots, X_n)$ for a parameter θ . The method in the previous section asks us to find a minimal complete sufficient statistic $T = T(X_1, \dots, X_n)$ by using a factorization criterion and then find a function of T which would be an unbiased estimator of θ . This will lead us to MVUE. Unfortunately, even if T is known, it is often difficult to find $\hat{\theta}(T)$ which would be unbiased for θ . In fact in some cases no unbiased estimator for θ exists.

For this reason we are looking for other methods to construct an estimator, which would be if not optimal then at least easy to construct and which would have small MSE in large samples.

We will consider two such methods, Method of Moments Estimation (MME) and Maximum Likelihood Estimation (MLE).

So we consider the situation when data X_1, \dots, X_n is from a distribution with parameter θ . The parameter θ can be a vector $\theta = (\theta_1, \dots, \theta_s)$ that

consists from several components. We want to estimate θ on the basis of the data sample.

Recall that the k -th population moment is simply the theoretical expectation of the observation X_i . We denote it μ_k .

$$\mu_k(\theta) = \mathbb{E}(X_i^k) = \begin{cases} \int x^k f(x|\theta) dx & \text{if } X_i \text{ are continuous r.v.}, \\ \sum_x x^k p(x|\theta) & \text{if } X_i \text{ are discrete r.v..} \end{cases}$$

Note that the population moments are all functions of the parameter θ (and do not depend on the sample data).

In contrast, the sample moments are functions of the data sample X_i . (They are random quantities and their distribution depends on θ .) We denote the k -th sample moment m_k .

$$m_k = m_k(X_1, \dots, X_n) = \frac{1}{n} \sum_{i=1}^n X_i^k.$$

For example, the first sample moment equals to the sample mean $m_1 = \bar{X}$, the second sample moment can be expressed in terms of the sample variance and the sample mean: $m_2 = (n-1)S^2 + n(\bar{X})^2$.

It should be emphasized that the sample moments are all functions of the data, i.e., they are can be calculated. And they are all random.

The main idea behind Method of Moments Estimator is that by the Law of Large Numbers the sample moments converge to population moments:

$$m_k = \frac{1}{n} \sum_{i=1}^n X_i^k \xrightarrow{\mathbb{P}} \mathbb{E}X_i^k = \mu_k(\theta),$$

in probability as $n \rightarrow \infty$.

So for large n we have,

$$m_k(X_1, \dots, X_n) = \mu_k(\theta) + \varepsilon_k,$$

where ε_k is very small with probability close to 1. This means that empirical moments are consistent estimators of the population moments. In addition, we know the form of the functions $\mu_k(\theta)$, although we do not know the value

of θ . Hence we can invert these functions and get an estimator of θ from the estimator $\hat{\mu}_k = m_k(X_1, \dots, X_n)$ of $\mu_k(\theta)$.

So, the idea is to ignore ε_k and solve the equations

$$m_k(X_1, \dots, X_n) = \mu_k(\hat{\theta})$$

for $\hat{\theta}$. The solution is given by the inverse function $\hat{\theta} = \mu_k^{(-1)}(m_k)$. If the inverse function $\mu_k^{(-1)}$ is a continuous function then we can use the continuous mapping theorem and show that $\hat{\theta} \rightarrow \theta$ in probability as $n \rightarrow \infty$, in other words, that $\hat{\theta}$ is a consistent estimator of θ .

How many moments we need for estimation? Typically, if we need to estimate a vector that consists of s parameters $\theta_1, \dots, \theta_s$, we use the first s moments. However it might happen that one of the first theoretical moments $\mu_k(\theta)$ actually does not depend on the parameter of interest. For example, a distribution with symmetric density function always has the first population moment (population mean or expectation) equal to zero and this first moment is not going to help us in estimation of parameters.

Practical steps:

1. Calculate the first K **population moments** $\mu_k(\vec{\theta})$ as functions of the vector of unknown parameters $\vec{\theta} = (\theta_1, \dots, \theta_s)$. (Typically, if you don't know s parameters, you need to calculate the first s moments.)
2. Write the first s **sample moments** m_k as functions of data vector X_i .
3. **Match the moments**: Solve the system of K equations $\mu_k(\hat{\theta}) = m_k(\vec{X})$ for $\hat{\theta}$.
4. The solution is a (vector) function $\hat{\theta}$ of the sample moments m_k and hence of the data X_i , because the sample moments are functions of the data. Since we believe that the equations are approximately true, we also believe that the solutions $\hat{\theta}$ are close to the true values of the parameter θ . These solutions give us the desired Method of Moments Estimator ($\hat{\theta}_{MME}$).

MME and Consistency Under some mild regularity conditions on the distribution of data, Method of Moment estimators are consistent. This is roughly because

- Under the conditions, we have

$$m_k(X_1, \dots, X_n) \xrightarrow{\mathbb{P}} \mu_k(\theta_1, \dots, \theta_s), \text{ for } k = 1, \dots, s,$$

when $n \rightarrow \infty$.

- The solution to the sistem of equations “ $m_k = \mu_k(\hat{\theta}_1, \dots, \hat{\theta}_s)$, $k = 1, \dots, s$ ” is a continuous map

$$\hat{\theta}_k = \hat{\theta}_k(m_1, \dots, m_s),$$

where $k = 1, \dots, s$. It is the inverse of the moment transformation μ , that sends parameters $\theta_1, \theta, \theta_s$ to moments μ_1, \dots, μ_s , so we could write (in vector form) $\hat{\theta}(m_1, \dots, m_s) = \mu^{(-1)}(m_1, \dots, m_s)$.

- By the continuous mapping theorem, we may conclude (again for conciseness using vector notation for the parameter θ and the estimator $\hat{\theta}$).

$$\hat{\theta} = \hat{\theta}(m_1, \dots, m_s) \xrightarrow{\mathbb{P}} \mu^{-1}(\mu_1(\theta), \dots, \mu_s(\theta)) = \theta.$$

Example 2.5.1 (Normal). The data X_1, \dots, X_n are distributed according the normal distribution $N(\mu, \sigma^2)$

- There are 2 parameters to estimate. The vector parameter θ has two components, $\theta = (\mu, \sigma^2)$
- The first population moment is $\mu_1(\theta) := \mathbb{E}(X_i) = \mu$. (There is a clash of notation here: $\mu_1(\theta)$ is the first population moment and μ also denotes the first component of the parameter vector θ .)
- The second population moment is $\mu_2(\theta) := \mathbb{E}X_i^2 = \text{Var}(X_i) + (\mathbb{E}X_i)^2 = \sigma^2 + \mu^2$.
- The first sample moment is $m_1 := \frac{1}{n} \sum_{i=1}^n X_i = \bar{X}$.

- The second sample moment is $m_2 := \frac{1}{n} \sum_{i=1}^n X_i^2$.

By matching the population and sample moments we obtain equations:

$$\begin{aligned}\hat{\mu} &= m_1 = \bar{X} \\ \hat{\sigma}^2 + \hat{\mu}^2 &= m_2 = \frac{1}{n} \sum_{i=1}^n X_i^2\end{aligned}$$

After solving these equations we get:

$$\begin{aligned}\hat{\mu} &= \bar{X}; \\ \hat{\sigma}^2 &= \frac{1}{n} \sum_{i=1}^n X_i^2 - (\bar{X})^2\end{aligned}$$

These are the MME estimators of the parameters. Note that the MME estimator for the variance is different from the standard estimator, the sample variance, which is

$$S^2 = \frac{1}{n-1} \left(\sum_{i=1}^n X_i^2 - n(\bar{X})^2 \right)$$

It is easy to see that they are related by the formula:

$$\hat{\sigma}_{MME}^2 = \frac{n-1}{n} S^2.$$

In particular, the MME estimator is biased.

Example 2.5.2 (Poisson with unusual parameter). The data sample X_1, \dots, X_n is distributed according to the Poisson distribution with parameter λ . Find the estimator of the parameter $\theta = 1/\lambda$.

- Only 1 parameter to be estimated.
- The first population moment is $\mu(\theta) = \mathbb{E}X_i = \lambda = 1/\theta$
- The first sample moment is $m_1 = \bar{X}$
- Match the two quantities above and solve for $\hat{\theta}$:

$$1/\hat{\theta} = \bar{X}$$

- Hence,

$$\hat{\theta}^{MME} = 1/\bar{X}$$

If you attentively examine the example above then you will see that we could estimate the parameter λ (the population mean) by the sample mean \bar{X} . This would give an MME estimator of λ . Then, since parameters θ and λ are in one-to-one correspondence to each other, therefore solving MME equations for $\hat{\theta}$ can be done by solving MME equations for $\hat{\lambda}$ and then using the one-to-one relation between the parameters. This would give us the MME estimator $\hat{\theta} = 1/\hat{\lambda} = 1/\bar{X}$.

This is a manifestation of the general principle valid for MME estimators.

Plug-in or Invariance property of MME

- If the parameter of interest ψ is a function of another parameter θ whose MME is relative easy to find, i.e., if

$$\psi = h(\theta)$$

and $\hat{\theta}_{MME}$ is easy to obtain, then

$$\hat{\psi}_{MME} = h(\hat{\theta}_{MME}),$$

i.e. we can apply the function h to $\hat{\theta}_{MME}$ to obtain $\hat{\psi}_{MME}$

- This is often easier than the “standard” procedure to find MME of ψ , where you need to redo the whole process.

Note that the previous statement assumes that MME for ψ and θ use the same set of moments.

Example 2.5.3. Let X_1, \dots, X_n be i.i.d observation from the uniform distribution on the interval $[0, \theta]$. What is the Method of Moments estimator of θ ?

The first population moment is $\mu_1 = \theta/2$. The first sample moment is \bar{X} . Therefore, the MM estimator is $\hat{\theta} = 2\bar{X}$. It is an unbiased estimator since $\mathbb{E}\hat{\theta} = 2\mathbb{E}\bar{X} = 2\mathbb{E}X_i = \theta$.

Note, however, that \bar{X} is not a sufficient statistic for θ . Indeed, the minimal sufficient statistic in this example is $X_{(n)} = \max\{X_1, \dots, X_n\}$.

So, this is an example of an MME estimator which is not a function of a sufficient statistic. So it has no chance to be an MVUE. Even though it is unbiased, its variance could be reduced by conditioning on the sufficient statistic.

Reflection

- The four point estimators back in Chapter 8 were all MMEs.
- In practice, Method of Moments is a very intuitive way to find an estimator. It requires only the ability to calculate the moments in terms of the parameters and invert this relation.
- Sometimes, MME gives biased estimators but at least it is consistent (under very mild conditions).
- One of its deficiencies is that it is not always a function of a sufficient statistic.
- It might be not as efficient as MLE which we consider next.
- One of the biggest deficiencies of MME is that it needs a modification if the data is not i.i.d. A generalized method of moments (GMM) was developed in 1980s by an econometrist Lars Peter Hansen who received Nobel Prize in Economics in 2013 in part for this work.

Here is a couple of additional examples.

Example 2.5.4. X_i 's are Gamma(α, β) distribution.

- First moment $\mu_1(\alpha, \beta) = \mathbb{E}X_i = \alpha\beta$
- Second moment $= \text{Var}(X_i) + (\mathbb{E}X_i)^2 = \alpha\beta^2 + (\alpha\beta)^2$
- Set

$$\begin{aligned}\alpha\beta &= m_1 = \bar{X} \\ \alpha\beta^2 + (\alpha\beta)^2 &= m_2\end{aligned}$$

- Treat α and β as unknowns. Treat the right-hand-sides as known numbers (you can calculate their actual values from the data). Whether you can solve this equation system is a high school algebra problem.

Example 2.5.5.

$$Y_i \sim f(y) = \left(\frac{2}{\theta^2}\right) (\theta - y) \mathbb{1}_{0 \leq y \leq \theta}$$

- Use MME to find an estimator for θ

—

$$\begin{aligned} \mu_1(\theta) &= \mathbb{E}Y_i = \int_0^\theta y \cdot f(y) dy = \left(\frac{2}{\theta^2}\right) \int_0^\theta (\theta y - y^2) dy \\ &= \left(\frac{2}{\theta^2}\right) \left(\theta \frac{y^2}{2} - \frac{y^3}{3}\right) \Big|_0^\theta \\ &= \left(\frac{2}{\theta^2}\right) \left(\theta \frac{\theta^2}{2} - \frac{\theta^3}{3}\right) \\ &= \theta - \frac{2\theta}{3} = \frac{\theta}{3} \end{aligned}$$

— Set $\mathbb{E}Y_1 = \frac{\theta}{3} = m_1 = \bar{Y}$

— Solve and obtain $\hat{\theta} = 3\bar{Y}$.

- Note that in this example \bar{Y} is also not a sufficient statistic for θ .

2.6 Maximum Likelihood Estimation (MLE)

The most popular method to find an estimator is the method of maximum likelihood, MLE.

Definition 2.6.1. The maximum likelihood estimator $\hat{\theta}$ is the value of the parameter θ , at which the likelihood function $L(\theta|x_1, \dots, x_n)$ takes its maximum value.

Informally, if we know that the probability to observe data sample (x_1, \dots, x_n) is 90% if the value of the parameter were θ_1 versus 10% if the value of the

parameter were θ_2 then we might prefer θ_1 as an estimator of true unknown θ .

The maximum likelihood estimator is the best estimator from this point of view.

The steps in finding MLE are easy: write down the maximum likelihood function and find its maximum with respect to the parameter θ . It turns out that in many cases, it is technically easier to look for maximum of **log-likelihood function** $\ell(\theta|\vec{x}) := \log(L(\theta|\vec{x}))$ instead $L(\theta)$. (It is maximized at the same value of θ .)

The main technical question in ML estimation is how to find the global maximum.

The maximization can be done by a computer algorithm or analytically. If the maximization is done by computer, there are many specialized algorithms suitable for statistical applications. One of them is EM algorithm. If the maximization is done analytically, we aim to solve equations $\frac{d}{d\theta}\ell(\theta|\vec{y}) = 0$. (They are called the first order conditions “FOC” for the extremal points.) Note, however, that the solution of these equations are all local extrema: local maxima, local minima and saddle points, so one should be careful to choose the global maximum among all possible solutions. In addition, the maximum can occur on the boundary of the set of all possible values of θ . In this case, the FOC will not give you information about the global maximum.

Consistency: The important fact about MLE is that this estimator is consistent under the mild distributional conditions. So it shares this good property of the Method of Moments estimators.

Let us consider our usual example.

Example 2.6.2 (Exponential). Let X_1, \dots, X_n be i. i. d. observations distributed according to the exponential distribution with parameter θ . What is the ML estimator for θ ?

The density of an individual observation is

$$f(x_i) = \frac{1}{\theta} e^{-x_i/\theta}.$$

Since the observations are independent, the likelihood is just the product of

the density functions:

$$L(\theta|x_1, \dots, x_n) = \frac{1}{\theta^n} \prod_{i=1}^n e^{-x_i/\theta} = \frac{1}{\theta^n} e^{-(\sum_{i=1}^n x_i)/\theta}$$

Hence the log-likelihood is

$$\ell(\theta|x_1, \dots, x_n) = \log L(\theta|x_1, \dots, x_n) = -n \log \theta - \frac{1}{\theta} \sum_{i=1}^n x_i$$

The first-order condition equation is

$$\begin{aligned} \frac{d}{d\theta} \ell(\theta|x_1, \dots, x_n) &= 0, \\ -n \frac{1}{\theta} + \frac{1}{\theta^2} \sum_{i=1}^n x_i &= 0, \\ \sum_{i=1}^n x_i &= n\theta, \end{aligned}$$

and the ML estimator is $\hat{\theta} = \bar{X}$. Note that it coincides with the MM estimator.

Now, consider our second standard example.

Example 2.6.3 (Normal). Let X_1, \dots, X_n be i.i.d observations from a normal distribution, $N(\mu, \sigma^2)$. What are the ML estimators of μ and σ^2 ?

- Likelihood function:

$$\begin{aligned} L(\mu, \sigma^2) &= \prod_{i=1}^n \left[\frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{(x_i - \mu)^2}{2\sigma^2}\right\} \right] \\ &= (2\pi\sigma^2)^{-\frac{n}{2}} \exp\left\{-\frac{\sum_{i=1}^n (x_i - \mu)^2}{2\sigma^2}\right\} \end{aligned}$$

- Log-likelihood:

$$\begin{aligned} \ell(\mu, \sigma^2) &= \log(L(\mu, \sigma^2)) = -\frac{n}{2} \log(2\pi\sigma^2) - \frac{\sum_{i=1}^n (x_i - \mu)^2}{2\sigma^2} \\ &= -\frac{n}{2} \log(2\pi) - \frac{n}{2} \log(\sigma^2) - \frac{\sum_{i=1}^n (x_i - \mu)^2}{2\sigma^2} \end{aligned}$$

- There are two unknown variables in this function. We want to calculate the partial derivative separately and set each to be zero, and then solve the unknowns.

$$\frac{\partial\{\ell(\mu, \sigma^2)\}}{\partial\mu} = \frac{2 \sum_{i=1}^n (x_i - \mu)}{2\sigma^2} = \frac{\sum_{i=1}^n y_i - n\mu}{\sigma^2} = \frac{\bar{x} - \mu}{\sigma^2/n}$$

$$\frac{\partial\{\ell(\mu, \sigma^2)\}}{\partial\sigma^2} = -\frac{n}{2} \frac{1}{\sigma^2} + \frac{\sum_{i=1}^n (x_i - \mu)^2}{2} \frac{1}{\sigma^4}$$

- Set both to zero.

$$\frac{\bar{x} - \mu}{\sigma^2/n} = 0 \Rightarrow \mu = \bar{x}$$

$$-\frac{n}{2} \frac{1}{\sigma^2} + \frac{\sum_{i=1}^n (x_i - \mu)^2}{2} \frac{1}{\sigma^4} = 0 \Rightarrow \sigma^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$$

- Hence, $\hat{\mu}_{MLE} = \bar{X}$, $\hat{\sigma}_{MLE}^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2 = S_0^2$. These are exactly the sample mean, and the alternative version of sample variance (but not S^2 , the original sample variance). We know that \bar{X} is unbiased and both \bar{X} and S_0^2 are consistent. However, S_0^2 is a biased estimator of σ^2 .

In fact, as in the previous example, the ML estimator coincide with the MM estimator.

Here is an example, where we get an ML estimator which is different from the MM estimator. It also illustrates that it is sometimes not enough to solve the first order conditions.

Example 2.6.4 (Uniform on $(0, \theta)$). Let X_1, \dots, X_n be i.i.d observations from the uniform distribution on the interval $[0, \theta]$. What is the ML estimator for θ .

- The density of X_i is $f(x_i) = 1/\theta$ if $x_i \in [0, \theta]$, and 0 otherwise.
- The likelihood function is the product of the densities, so

$$L(\theta|x_1, \dots, x_n) = \frac{1}{\theta^n},$$

if $\min x_1, \dots, x_n \geq 0$, and $\max x_1, \dots, x_n \leq \theta$ and 0 otherwise.

- We can also write this

$$L(\theta) = \theta^{-n} \mathbb{1}_{\{x_{(n)} \leq \theta\}} \mathbb{1}_{\{x_{(1)} \geq 0\}}$$

- Or equivalently $L(\theta) = \theta^{-n}$ for $\theta \geq x_{(n)}$, and 0, otherwise. [← better here for MLE.](#)
- The log-likelihood is $\ell(\theta) = -n \log(\theta)$ for $\theta \geq x_{(n)}$, and $-\infty$, otherwise.
- $\ell'(\theta) = -\frac{n}{\theta}$ for $\theta \geq x_{(n)}$
- But $-\frac{n}{\theta} = 0$ does not have a solution! What can be wrong?
- The maximum of the function $\ell(\theta)$ is reached at the boundary point!

Indeed $\ell'(\theta) = -\frac{n}{\theta} < 0$, hence $\ell(\theta)$ is a decreasing function. The smallest possible value for θ in the domain $\theta \geq x_{(n)}$ is its left boundary point $x_{(n)}$, which is the maximum point for the likelihood.

So, $\hat{\theta}_{ML} = X_{(n)}$. Note that it is a function of a minimal sufficient statistic $X_{(n)}$. Is it biased? How does its variance compare with the MM estimator?

In order to calculate the expectation of $X_{(n)}$, let us repeat some information from Example 1.2.4. If X_i are from distribution with CDF $F(x)$ and density $f(x)$, then the CDF of $X_{(n)}$ is

$$\begin{aligned} F_{X_{(n)}}(y) &\equiv \Pr(X_{(n)} \leq y) = \Pr(X_i \leq y, \text{ for all } i) \\ &= [F(x)]^n \end{aligned}$$

and the PDF is

$$f_{X_{(n)}}(x) = n[F(x)]^{n-1}f(x).$$

In our case the density $f(x) = 1/\theta$ and cdf $F(x) = x/\theta$, so we have $f_{X_{(n)}}(x) = nx^{n-1}/\theta^n$. By integrating, we find $\mathbb{E}X_{(n)} = \frac{n}{n+1}\theta$ and $\mathbb{E}X_{(n)}^2 = \frac{n}{n+2}\theta^2$, so

$$\text{Var}(X_{(n)}) = \theta^2 \left(\frac{n}{n+2} - \frac{n^2}{(n+1)^2} \right) = \theta^2 \left(\frac{n}{(n+2)(n+1)^2} \right) < \theta^2 \frac{1}{(n+1)^2}.$$

So the ML estimator is biased with bias = $\theta/(n+1)$. However its variance is much smaller than the variance of the unbiased MM estimator. ($\text{Var}(2\bar{X}) = 4\text{Var}(X_i)/n = \theta^2/(3n)$.)

Here we see a clear difference between ML and MM estimators. The MM estimator $2\bar{X}$ is unbiased but it is not a function of a sufficient statistic. The ML estimator is biased but it is a function of a minimal sufficient statistic. In particular, the bias of the ML estimator can be corrected and this will lead to an MVUE estimator.

Note, by the way, that if the domain of the density function depends on the parameter θ , it is often a warning sign that the boundary point of θ may play a major role in finding the MLE.

MLE is always a function of a sufficient statistic!

Theorem 2.6.5. *Suppose $T(X_1, \dots, X_n)$ is a sufficient statistic for θ . Then $\hat{\theta}_{MLE}$ be written as a function of T , $\hat{\theta}_{MLE}(X_1, \dots, X_n) = \hat{\theta}(T)$.*

In particular, $\hat{\theta}(T)$ is a function of a complete minimal sufficient statistic as long as such a beast exists. This is rather appealing because this means that if are able to find a function of $\hat{\theta}_{MLE}$ which is unbiased, then this function is an MVUE.

Proof. If T is a sufficient statistic, then by the factorization criterion,

$$L(\theta) = g(t, \theta)h(x_1, \dots, x_n) \text{ and so } \ell(\theta) = \log(g(t, \theta)) + \log(h(y_1, \dots, y_n))$$

- Since $\log(h(x_1, \dots, x_n))$ has nothing to do with θ , as far as θ is concerned, $\log(h(x_1, \dots, x_n))$ is a constant; hence the maximizer of $\ell(\theta)$ over θ is the same as the maximizer of $\log(g(t, \theta))$.
- The maximizer of $\log(g(t, \theta))$, over all possible θ , has to depend only on t .

Thus the ML estimator is a function of T . □

Example 2.6.6 (Poisson with usual and unusual parameter). Suppose that X_1, \dots, X_n are i.i.d. from the Poisson distribution with parameter λ . What is the ML estimator for λ ?

The Poisson distribution is discrete so we work with probability mass functions (“pmf”s). The pmf of one observation X_i is

$$p_{X_i}(x_i|\lambda) := \mathbb{P}[X_i = x_i] = e^{-\lambda} \frac{\lambda^{x_i}}{x_i!},$$

where x_i can take values $0, 1, 2, \dots$. Since the observations are independent, the likelihood function is simply the product of pmfs of individual observations.

$$L(\lambda|x_1, \dots, x_n) = \lambda^{\sum_{i=1}^n x_i} e^{-n\lambda} / \prod_{i=1}^n (x_i!).$$

Hence, the log-likelihood is

$$\ell(\lambda) = \left(\sum_{i=1}^n x_i \right) \log \lambda - n\lambda - \log \left(\prod_{i=1}^n (x_i!) \right).$$

So we can write the first order condition as

$$\ell'(\lambda) = \left(\sum_{i=1}^n x_i \right) \frac{1}{\lambda} - n = 0$$

The solution gives us the ML estimator $\hat{\lambda}_{ML} = \bar{X}$. It coincides with the MM estimator.

Now suppose we use a different parameter in the model $\theta = 1/\lambda$ and want to find an ML estimator for θ . This simply means that now we write the distribution function for the observations in terms of θ not λ :

$$p_{X_i}(x_i|\theta) := e^{-1/\theta} \frac{(1/\theta)^{x_i}}{x_i!},$$

So the likelihood function will be

$$L(\theta|x_1, \dots, x_n) = (1/\theta)^{\sum_{i=1}^n x_i} e^{-n(1/\theta)} / \prod_{i=1}^n (x_i!).$$

Now it is rather obvious that if $L(\lambda|x_1, \dots, x_n)$ is maximized at $\lambda = \hat{\lambda}$, then $L(\theta|x_1, \dots, x_n)$ is maximized at the point that corresponds to this point, namely at $\theta = 1/\hat{\lambda}$. Therefore,

$$\hat{\theta}_{ML} = 1/\hat{\lambda}_{ML} = 1/\bar{X}.$$

The principle that the relation between parameters are transferred to their estimates is called the invariance, or plug-in, principle.

Theorem 2.6.7. Suppose that X_1, \dots, X_n are observations from the distribution that depends on parameter θ . If $\hat{\theta}_{ML} = \hat{\theta}_{ML}(X_1, \dots, X_n)$ is the maximum likelihood estimator for θ and $g(\cdot)$ is a one-to-one function, then $g(\hat{\theta}_{ML})$ is the maximum likelihood for parameter $\psi := g(\theta)$, i.e.,

$$\hat{\psi}_{ML} = g(\hat{\theta}_{ML})$$

Proof. We have the identity

$$L(\theta|y) = L(g^{-1}(\psi)|y),$$

and the expression on the right is the likelihood for ψ . If the MLE of ψ were $\psi^* \neq g(\hat{\theta}_{MLE})$, then it would follow that

$$L(g^{-1}(\psi^*)) > L(g^{-1}(g(\hat{\theta}_{MLE}))) = L(\hat{\theta}_{MLE})$$

But this would contradict the fact that $\hat{\theta}_{MLE}$ maximizes $L(\theta)$. \square

- The invariance property actually holds for any function ψ of a parameter θ (and not only for the one-to-one functions), once we define appropriately, what do we mean by the maximum likelihood estimator of $\psi(\theta)$ in this case.
- A discussion and a proof can be seen in Casella and Berger (2002) [Math 502]

Example 2.6.8. For example, suppose you want to estimate the probability that a random variable $X > 2$, you know that it is a Poisson r.v. and have a data sample $X_i, i = 1, \dots, n$.

Note that

$$\begin{aligned} \Pr\{X > 2\} &= 1 - \Pr\{X = 0\} - \Pr\{X = 1\} - \Pr\{X = 2\} \\ &= 1 - e^{-\lambda} - e^{-\lambda} \frac{\lambda}{1!} - e^{-\lambda} \frac{\lambda^2}{2!} \end{aligned}$$

The invariance principle tells us that the ML estimator for $\Pr\{X > 2\}$ is simply

$$1 - e^{-\hat{\lambda}} - e^{-\hat{\lambda}} \frac{\hat{\lambda}}{1!} - e^{-\hat{\lambda}} \frac{\hat{\lambda}^2}{2!},$$

where $\hat{\lambda} = \hat{\lambda}_{ML}$ is the maximum likelihood estimator for λ .

Since we know that $\hat{\lambda}_{ML} = \bar{X}$, therefore the ML estimator for $\Pr\{X > 2\}$ is

$$1 - e^{-\bar{X}} \left(1 + \frac{\bar{X}}{1!} + \frac{(\bar{X})^2}{2!} \right).$$

Quiz 2.6.9. Which of the following statements is/are correct?

- I. The method of moment estimator for θ is always unbiased for θ .
 - II. The method of moment estimator for θ is always a function of a sufficient statistic for θ .
 - III. The method of moment estimator for θ is always the minimum variance unbiased estimator for θ .
- A. I and II
- B. II and III
- C. I only
- D. II only
- E. None of above

Quiz 2.6.10. Which of the following statements is/are correct?

- I. There is only one sufficient statistic for θ .
 - II. A sufficient statistic for θ is always consistent for θ .
 - III. The maximum likelihood estimator for θ is always a function of a minimal sufficient statistic for θ .
 - IV In most cases, the maximum likelihood estimator is consistent.
- A. I and III
 - B. II and IV
 - C. III and IV
 - D. only IV
 - E. None of above

Here is a couple of other examples.

Example 2.6.11. The data Y_1, \dots, Y_n are from the Bernoulli distribution with parameter p . What is the MLE for p ?

- The likelihood function is

$$L(p|\vec{y}) = \prod_{i=1}^n p^{y_i} (1-p)^{1-y_i} = p^t (1-p)^{n-t},$$

where $t = \sum_{i=1}^n y_i$.

- The log-likelihood function is

$$\ell(p|\vec{y}) := \log(L(p|\vec{y})) = t \log(p) + (n-t) \log(1-p)$$

-

$$\frac{d}{dp} \ell(p|\vec{y}) = \frac{t}{p} - \frac{n-t}{1-p}$$

- Set

$$\ell'(p) = \frac{t}{p} - \frac{n-t}{1-p} = 0$$

We obtain

$$\frac{t}{p} = \frac{n-t}{1-p} \Rightarrow t - tp = np - tp \Rightarrow p = \frac{t}{n}$$

- Hence $\hat{p}_{MLE} = \frac{T}{n}$

- Recall that we proved that this is an unbiased estimator, and since it is a function of a minimal sufficient statistic, it is the MVUE. It is also consistent since its variance goes to 0.

Example 2.6.12. Let Y_1, \dots, Y_n be taken from distribution with the following density:

$$f_Y(y) = \frac{1}{\theta} r y^{r-1} e^{-y^r/\theta} \mathbb{1}_{y>0}, \quad \text{where } \theta > 0 \text{ and } r \text{ is known.}$$

Find a sufficient statistic for θ . Find the MLE of θ . Is it MVUE?

- $L(\theta) = \prod_{i=1}^n \left\{ \frac{1}{\theta} r y_i^{r-1} e^{-y_i^r/\theta} \right\} = \frac{1}{\theta^n} r^n (\prod_{i=1}^n y_i)^{r-1} e^{-\frac{1}{\theta} \sum_{i=1}^n y_i^r}$
- Clearly the sufficient statistic is $\sum_{i=1}^n Y_i^r$
- $L(\theta) = C \cdot \frac{1}{\theta^n} e^{-\frac{1}{\theta} \sum_{i=1}^n y_i^r}$ where C has nothing to do with θ .
- $\ell(\theta) = \log(C) - n \log(\theta) - \frac{1}{\theta} \sum_{i=1}^n y_i^r$
- $\ell'(\theta) = -\frac{n}{\theta} + \frac{1}{\theta^2} \sum_{i=1}^n y_i^r$. Note that $\log(C)$ disappears.
- Set $\ell'(\theta) = 0 \Rightarrow \frac{n}{\theta} = \frac{1}{\theta^2} \sum_{i=1}^n y_i^r \Rightarrow \theta^* = \frac{1}{n} \sum_{i=1}^n y_i^r$

- So, $\frac{1}{n} \sum_{i=1}^n Y_i^r$ is the MLE for θ .
- MVUE??? We know that the estimator is a sufficient statistic. Need to check unbiasedness.
- $\mathbb{E}(\frac{1}{n} \sum_{i=1}^n Y_i^r) = \frac{1}{n} \sum_{i=1}^n \mathbb{E}(Y_i^r)$
- Note that

$$\begin{aligned}\mathbb{E}(Y_i^r) &= \int_0^\infty \frac{1}{\theta} r y^{r-1} e^{-y^r/\theta} \cdot y^r dy \\ &= \int_0^\infty e^{-y^r/\theta} \cdot y^r d(y^r/\theta) \\ &\quad (\text{Let } u = y^r/\theta) \\ &= \theta \int_0^\infty e^{-u} \cdot u du\end{aligned}$$

- One can either calculate the integral explicitly or note that e^{-u} is the density of the exponential distribution with parameter 1 and the integral $\int_0^\infty e^{-u} \cdot u du$ calculate its expectation, which, as we already know, equals 1. Hence, we have

$$\mathbb{E}(Y_i^r) = \theta$$

- Therefore, $\mathbb{E}(\frac{1}{n} \sum_{i=1}^n Y_i^r) = \frac{1}{n} \sum_{i=1}^n \mathbb{E}(Y_i^r) = \frac{1}{n} n\theta = \theta$

It follows that the maximum likelihood estimator $\frac{1}{n} \sum_{i=1}^n Y_i^r$ is the MVUE for θ .

Exercise 2.6.13. Y_1, Y_2, \dots, Y_n is a sample of observations from $N(5, \theta)$ where the **variance θ is unknown** and is the parameter of interest:

$$f(y) = \frac{1}{\sqrt{2\pi\theta}} \exp\left[-\frac{(y-5)^2}{2\theta}\right].$$

- Find the sufficient statistic for θ .
- Find the Method of Moment Estimator (MME) $\hat{\theta}^{MM}$ for θ .

- (c). Find the Maximum Likelihood Estimator (MLE) $\hat{\theta}^{ML}$ for θ
- (d). Show directly (without using the general theorem that MM estimators are consistent) that $\hat{\theta}^{MM}$ is consistent for θ .
- (e). Show directly that $\hat{\theta}^{ML}$ is consistent for θ .
- (f). Prove that $\hat{\theta}^{ML}$ is the minimal variance unbiased estimator (MVUE) for θ

Exercise 2.6.14. Let Y_1, Y_2, Y_3 be three i.i.d. observations from the distribution with density:

$$f_Y(y) = \frac{ye^{-y/\theta}}{\theta^2} \mathbb{1}_{y>0}$$

In a data sample these random variables are observed to be 120, 130 and 128, respectively.

- Find the ML estimator of θ
- Is the ML estimator unbiased in this model? Explain.
- What is the ML estimator for the variance of Y_1 ?

More examples:

Example 2.6.15. Consider the situation when we have two samples. One of them is X_1, X_2, \dots, X_m from normal distribution $N(\mu_1, \sigma^2)$. The other is Y_1, Y_2, \dots, Y_n from normal distribution $N(\mu_2, \sigma^2)$. Here we assumed that the variance in both distributions is the same. What is the ML estimators for the parameters μ_1, μ_2 and σ^2 ?

- Given the observations $x_1, \dots, x_m, y_1, \dots, y_n$, the likelihood for μ_1, μ_2, σ^2

is the product of all the densities (including the X 's and the Y 's)

$$\begin{aligned} L(\mu_1, \mu_2, \sigma^2) &= \prod_{i=1}^m \left[\frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{(x_i - \mu_1)^2}{2\sigma^2}\right\} \right] \times \\ &\quad \prod_{i=1}^n \left[\frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{(y_i - \mu_2)^2}{2\sigma^2}\right\} \right] \\ &= (2\pi\sigma^2)^{-\frac{m}{2}} \exp\left\{-\frac{\sum_{i=1}^m (x_i - \mu_1)^2}{2\sigma^2}\right\} \times \\ &\quad (2\pi\sigma^2)^{-\frac{n}{2}} \exp\left\{-\frac{\sum_{i=1}^n (y_i - \mu_2)^2}{2\sigma^2}\right\} \end{aligned}$$

- So the log-likelihood is

$$\begin{aligned} \ell(\mu_1, \mu_2, \sigma^2) &= -\frac{m}{2} \log(2\pi) - \frac{m}{2} \log(\sigma^2) - \frac{\sum_{i=1}^m (x_i - \mu_1)^2}{2\sigma^2} \\ &\quad - \frac{n}{2} \log(2\pi) - \frac{n}{2} \log(\sigma^2) - \frac{\sum_{i=1}^n (y_i - \mu_2)^2}{2\sigma^2} \end{aligned}$$

Let us use the notation $\theta = \sigma^2$. The partial derivatives of the log-likelihood function with respect to parameters are:

$$\begin{aligned} \ell_{\mu_1} &= \frac{\bar{x} - \mu_1}{\theta/n} \\ \ell_{\mu_2} &= \frac{\bar{y} - \mu_2}{\theta/n} \\ \ell_{\theta} &= \left(-\frac{m}{2} \frac{1}{\theta} + \frac{\sum_{i=1}^m (x_i - \mu_1)^2}{2} \frac{1}{\theta^2} \right) + \left(-\frac{n}{2} \frac{1}{\theta} + \frac{\sum_{i=1}^n (y_i - \mu_2)^2}{2} \frac{1}{\theta^2} \right) \end{aligned}$$

- Set all three to zero and get the solution.

$$\begin{aligned} \hat{\mu}_1 &= \bar{X}, \hat{\mu}_2 = \bar{Y}, \\ \hat{\sigma}^2 &= \frac{\sum_{i=1}^m (x_i - \bar{X})^2 + \sum_{i=1}^n (y_i - \bar{Y})^2}{m + n} \end{aligned}$$

- Does $(\hat{\sigma}^2)_{ML}$ look a bit familiar?

Example 2.6.16. Let Y_1, \dots, Y_n be from the distribution with density

$$f(y) = (\theta + 1)y^\theta, \quad 0 < y < 1,$$

where $\theta > -1$. Find the MLE.

1. $L(\theta) = (\theta + 1)^n (\prod_{i=1}^n y_i)^\theta$ for $\theta > -1$
2. $\ell(\theta) = n \log(\theta + 1) + \theta \log(\prod_{i=1}^n y_i)$ for $\theta > -1$
3. $\ell'(\theta) = \frac{n}{\theta+1} + \log(\prod_{i=1}^n y_i)$ for $\theta > -1$
4. Let $\ell'(\theta) = 0 \Rightarrow \theta^* = -\frac{n}{\log(\prod_{i=1}^n y_i)} - 1$

$$\hat{\theta}_{MLE} = -\frac{n}{\log(\prod_{i=1}^n Y_i)} - 1$$

Note that we can show that $\prod_{i=1}^n Y_i$ is sufficient, which confirm that every ML estimator is a function of the sufficient statistic.

2.7 Cramer-Rao Lower Bound and large sample properties of MLE

In this section we learn about the Cramer-Rao lower bound on the variance of any *unbiased* estimator of θ . It is not possible to get smaller variance even if you use the MVUE. We also learn that in the limit, for $n \rightarrow \infty$, the maximum likelihood estimator achieves this bound. In this sense, the ML estimator is an *asymptotically* Minimal Variance Unbiased Estimator.

The idea behind the Cramer-Rao bound is that if the likelihood function is flat and does not depend on the parameter θ then it will be difficult to estimate the parameter from the data. The measure of the likelihood function flatness that the Cramer-Rao bound uses is the Fisher information. It is the average squared sensitivity of the log-density function to the parameter.

For the formal definition, let us define the score function $s(x, \theta)$ of a random variable X as $\log f(x, \theta)$ if X is continuous with probability density $f(x, \theta)$ and as $\log p(x, \theta)$ if it is discrete with probability mass function $p(x, \theta)$.

Definition 2.7.1. The **Fisher information** of a random variable X is defined as

$$I_X(\theta) = \mathbb{E} \left[\left(\frac{d}{d\theta} s(X, \theta) \right)^2 \right]$$

Example 2.7.2. Let us calculate the Fisher information for exponential random variable with mean θ . The density is $f(x, \theta) = \frac{1}{\theta}e^{-x/\theta}$, so the score is

$$s(x, \theta) = -\log \theta - x/\theta.$$

By definition,

$$\begin{aligned} I(\theta) &= \mathbb{E} \left[\left(\frac{d}{d\theta} (-\log \theta - X/\theta) \right)^2 \right] \\ &= \mathbb{E} \left[\left(-1/\theta + X/\theta^2 \right)^2 \right] \\ &= \mathbb{E} \left[1/\theta^2 - 2X/\theta^3 + X^2/\theta^4 \right] \end{aligned}$$

Recollect that the exponential distribution with parameter θ has mean θ and variance θ^2 . Hence $\mathbb{E}X = \theta$ and $\mathbb{E}X^2 = \theta^2 + \theta^2 = 2\theta^2$. So after substitution we get

$$I(\theta) = 1/\theta^2 - 2\theta/\theta^3 + 2\theta^2/\theta^4 = 1/\theta^2.$$

In many cases, it is easier to calculate the Fisher information by using a different formula. Namely, under some regularity conditions, one has:

$$I_X(\theta) = -\mathbb{E} \left[\frac{d^2}{d\theta^2} s(X, \theta) \right]. \quad (2.1)$$

For example if then the following result holds.

Lemma 2.7.3. *Suppose that X is a continuous random variable and the range of X (the set where the density is positive) does not depend on θ . Suppose also that the density is continuously differentiable in θ on the range. Then the equality (2.1) holds*

This lemma can be proved by rather challenging integral manipulations. Identity 2.1 often gives a shorter path to calculate Fisher's information. In the previous example, this identity gives

$$\begin{aligned} I(\theta) &= -\mathbb{E} \left[\frac{d^2}{d\theta^2} (-\log \theta - X/\theta) \right] \\ &= -\mathbb{E} \left[1/\theta^2 - 2X/\theta^3 \right] = 1/\theta^2. \end{aligned}$$

We formulate the following result as a theorem although we do not specify the regularity conditions precisely. Check for them in graduate level textbooks.

Theorem 2.7.4 (Cramer-Rao bound). *Let X_1, \dots, X_n be a sample of independent identically distributed observations from the distribution that depends on parameter θ . Under certain regularity conditions on the distribution, for every unbiased estimator $\hat{\theta}$,*

$$\text{Var}(\hat{\theta}) \geq \frac{1}{nI_X(\theta)}$$

For example, the regularity conditions are satisfied if X_1, \dots, X_n are i.i.d. continuous random variables with density $f(x, \theta)$ provided that the support of the distribution (that is, the set of x where $f(x, \theta)$ is positive) does not depend on θ and that the density is continuously differentiable in θ .

Note 1: If we are able to find an unbiased estimator that the Cramer-Rao bound is reached by its variance and the regularity conditions hold, then this estimator is MVUE (minimal variance unbiased estimator).

Note 2: Sometimes the Cramer - Rao bound is not sharp. That is, sometimes the variance of the MVUE will be larger than the bound given by the Cramer - Rao inequality.

Note 3: The regularity conditions are often violated if the parameter involves the domain of the density. (Like estimation of the θ for the random variable uniformly distributed on $[0, \theta]$.) In this case, the Cramer-Rao bound is invalid: there can be an estimator with smaller variance than the bound predicts.

Example 2.7.5 (Exponential). Suppose that X_1, \dots, X_n are i.i.d. observations from the exponential distribution with parameter θ . We have calculated that the Fisher Information of this distribution is $I_X(\theta) = 1/\theta^2$. Hence, the Cramer-Rao inequality says that every unbiased estimator $\hat{\theta}$ has variance $\geq \theta^2/n$. On the other hand the estimator $\hat{\theta} = \bar{X}$ is unbiased and has variance $\text{Var}(X_i)/n = \theta^2/n$. This fact means that this unbiased estimator attains the Cramer-Rao bound and so it is MVUE.

Example 2.7.6 (Uniform on $(0, \theta)$). Now suppose that X_1, \dots, X_n are i.i.d. observations from the uniform distribution on the interval $[0, \theta]$. First, let us calculate the Fisher Information for this distribution. The density is $1/\theta$ and the score function is $-\log(\theta)$, so by definition:

$$I_X(\theta) = \mathbb{E} \left[\left(\frac{d}{d\theta} s(X, \theta) \right)^2 \right] = \mathbb{E} \left[\frac{1}{\theta^2} \right] = \frac{1}{\theta^2}$$

(If we use identity (2.1), we get:

$$I_X(\theta) = -\mathbb{E} \left[\frac{d^2}{d\theta^2} s(X, \theta) \right] = -\mathbb{E} \left[\frac{d^2}{d\theta^2} s(X, \theta) \right] = \mathbb{E} \left[\frac{1}{\theta^2} \right] = \frac{1}{\theta^2},$$

which is the same. Note, however, that the conditions of Lemma 2.7.3 that justified (2.1) are not satisfied in this example, so we are lucky that (2.1) gives the correct result.) So the Cramer-Rao inequality predicts that every unbiased estimator should have variance $\geq \theta^2/n$. We have calculated earlier the expectation and variance of the ML estimator in this example which is $X_{(n)}$. This estimator is biased but we can correct its bias and consider the unbiased estimator $\hat{\theta} = \frac{n+1}{n} X_{(n)}$. Its variance is

$$\mathbb{V}\text{ar}(\hat{\theta}) = \theta^2 \frac{(n+1)^2}{n^2} \frac{n}{(n+2)(n+1)^2} = \theta^2 \frac{1}{n(n+2)}$$

So this estimator clearly violates the Cramer-Rao bound. The reason is that the conditions of the Theorem 2.7.4 are not satisfied.

We will explain some ideas behind the proof of the Cramer-Rao bound below. Now let us turn to the asymptotic optimality of MLE. The main result here is that under some regularity conditions,

$$n\mathbb{V}\text{ar}(\hat{\theta}_{ML}) \rightarrow \frac{1}{I_X(\theta)}$$

as $n \rightarrow \infty$.

The point is that the MLE in the limit attains the Cramer-Rao bound. In this sense it is an asymptotically MVUE, or in other terminology it is asymptotically efficient.

Ideas of the proof of the Cramer-Rao bound

We are going to prove the bound for the case when X is continuous and its range does not depend on the parameter and when $n = 1$. The proof in the general case is difficult and you can find it in graduate level textbooks.

Lemma 2.7.7. *Assume the range of X does not depend on θ and the density is positive and continuously differentiable in θ . Then,*

$$\mathbb{E}\left[\frac{d}{d\theta}s(X, \theta)\right] = 0.$$

Proof. Note that by chain rule:

$$\frac{d}{d\theta}s(X, \theta) = \frac{\frac{d}{d\theta}f(x, \theta)}{f(x, \theta)}.$$

$$\begin{aligned}\mathbb{E}\left[\frac{d}{d\theta}s(X, \theta)\right] &= \int_a^b \frac{d}{d\theta}s(x, \theta)f(x, \theta)dx = \int_a^b \frac{d}{d\theta}f(x, \theta)dx \\ &= \frac{d}{d\theta} \int_a^b f(x, \theta)dx = \frac{d}{d\theta}1 = 0.\end{aligned}$$

□

Corollary of Lemma 2.7.7: $I(\theta) = \mathbb{V}\text{ar}\left(\frac{d}{d\theta}s(X, \theta)\right)$.

Proof of the Cramer-Rao bound for $n = 1$. Let $\hat{\theta}(X)$ be an unbiased estimator of θ based on just one datapoint X . Let us write $s'(\theta)$ instead of $\frac{d}{d\theta}s(X, \theta)$. By using Lemma 2.7.7 and the Cauchy - Schwarz inequality for covariance:

$$|\mathbb{E}(s'(\theta)\hat{\theta})| = |\text{Cov}(s'(\theta), \hat{\theta})| \leq \sqrt{\mathbb{V}\text{ar}(s'(\theta))\mathbb{V}\text{ar}(\hat{\theta})}.$$

Or,

$$\mathbb{V}\text{ar}(\hat{\theta}) \geq \frac{|\mathbb{E}(s'(\theta)\hat{\theta})|}{I(\theta)}$$

All this would hold if $\hat{\theta}$ was biased. The next step is crucial.

$$\begin{aligned}\mathbb{E}(s'(\theta)\hat{\theta}) &= \int_a^b \frac{d}{d\theta}f(x, \theta)\hat{\theta}(x)dx = \frac{d}{d\theta} \int_a^b f(x, \theta)\hat{\theta}(x)dx \\ &= \frac{d}{d\theta}\mathbb{E}\hat{\theta} = \frac{d}{d\theta}\theta = 1.\end{aligned}$$

□

Example 2.7.8. • $Y_i \sim \text{Bernoulli}(p)$ with PMF $p^{y_i}(1-p)^{1-y_i}$

- We know that $\hat{p} = \sum_{i=1}^n Y_i/n$ is the MLE for p
- Try to derive the Cramer-Rao Lower Bound for \hat{p}

$$\begin{aligned}
 & \frac{1}{n\mathbb{E}\left[-\frac{\partial^2 \log f(Y|\theta)}{\partial \theta^2}\right]} = \frac{1}{n\mathbb{E}\left[-\frac{\partial^2 \log[p^Y(1-p)^{1-Y}]}{\partial p^2}\right]} \\
 &= \frac{1}{n\mathbb{E}\left[-\frac{\partial^2 [Y \log(p) + (1-Y) \log(1-p)]}{\partial p^2}\right]} = \frac{1}{n\mathbb{E}\left[\frac{Y}{p^2} + \frac{1-Y}{(1-p)^2}\right]} \\
 &= \frac{1}{n\left[\frac{p}{p^2} + \frac{1-p}{(1-p)^2}\right]} = \frac{1}{n\left[\frac{1}{p} + \frac{1}{1-p}\right]} = \frac{p(1-p)}{n}
 \end{aligned}$$

- We already know that $\text{Var}(\hat{p}) = \frac{p(1-p)}{n}$. So indeed, \hat{p} is the MVUE.

Chapter 3

Hypothesis testing

3.1 Basic definitions

- Chapter 8: make statistical inference about the population (parameter) by
 - point estimators (with unbiasedness & low variance), and
 - confidence intervals.
- Chapter 9: different mathematical properties of point estimators, and how to find good estimators.
- Chapter 10:
 - Hypothesis tests: answer scientific questions using statistics
 - Different philosophy and goal.
 - Some connection with confidence interval.

Examples of “tests” in real life

- What people may want to know:
 1. Does smoking cause lung cancer?
 2. Is global warming real?

3. Are men more likely to run a stop sign than women?
 4. Does chemotherapy really cure cancer?
 5. Is a new medicine effective in increasing longevity?
- Beyond just scientific interest.
 - Business decisions, military actions, political strategies.

The basic philosophy of statistical testing is “Proof by Contradiction”.

1. Identify the question of interest:
DOES SMOKING LEAD TO LUNG CANCER?
2. Try to prove causality by assuming the **opposite theory** (“does not lead”), which is called **null hypothesis**, and showing that it leads to a contradiction.
3. Namely, if the data looks **very** improbable under the null hypothesis, then you can conclude that the data contradicts the null hypothesis, so it should be rejected and your theory should be accepted instead.
4. However, if the data does not look **very** improbable under null hypothesis, then you cannot reject it and so you don’t have enough evidence in support of your, alternative, point of view.

Terminology

- Hypothesis
 - A statement about a population, usually of the form that a parameter takes a particular numerical value (e.g. $\theta = 2$) or falls in a certain range of values (e.g. $\theta > 2$).
- Null Hypothesis H_0
 - The statement of **no effect**.
 - This is the statement that we will assume as true when we will try to show that it leads to improbable conclusions.

- It is usually denoted by H_0 and it is usually very specific. For example it can state: “the treatment has no effect”.
- Alternative Hypothesis H_a
 - The statement of **some effect**.
 - The statement that we actually want to confirm by showing that H_0 should be **rejected**.
 - Usually it is denoted by H_a ; it can be specific like: the percentage of recoveries after a medicine was used increased by 10%. This is a point hypothesis. Or it can be less specific: the percent of recoveries after the treatment increased by at least 10%. This is called the composite alternative hypothesis.
- The hypotheses must be stated before collecting, viewing or analyzing the data.

Besides the null and alternative hypothesis, the statistical test is defined by a *test statistic* and a *rejection region*.

- Recall: a **statistic** is a function of random observations Y_i (the data). A statistic cannot be defined in terms of the unknown θ .
- A **test statistic (TS)** is a statistic, that is, a function of the data. Its intention is different from an estimator which is also a function a data. The test statistic should help us to answer the question “*how close is the data sample to what we would expect if the null hypothesis H_0 were true?*”
- If the null hypothesis is expressed in terms of a parameter of the model, for example if it has the form like $\theta = \theta_0$, where θ_0 is a specific value, then often you can use a test statistic in the form

$$TS = \frac{\hat{\theta} - \theta_0}{\sigma_{\hat{\theta}}},$$

where $\hat{\theta}$ is an estimator of θ .

- In this case, if $|TS|$ is **large** (so TS is far from 0) then it **might indicate** that the data is not compatible with the null hypothesis.

Definition 3.1.1. The **reject region** RR is a set of possible values of a test statistic $TS(X_1, \dots, X_n)$ so that if value of TS for the observed data is in the rejection region RR , then we reject the null hypothesis H_0 .

- Example: if we reject H_0 when $TS > 1$, then $RR = \{x : x > 1\}$, that is $RR = (1, \infty)$.
- Often, it is more convenient to write the RR as some inequality, such as $RR : TS < a$. This is equivalent to $RR = \{x : x < a\} = (-\infty, a)$.
- If the TS is not in the RR , then we fail to reject the null hypothesis H_0 .
- Other commonly used terms: do not reject, do not have enough evidence to reject, etc.
- Next: How to find the reject region?

Example 3.1.2. Let X_1, \dots, X_n be distributed with Bernoulli distribution with parameter p .

- Null hypothesis: $H_0: p = 0.5$
- Alternative hypothesis $H_a: p > 0.5$.
- Test the hypothesis H_0 against alternative H_a .
- We can define our test statistic to be

$$TS = \frac{\hat{p} - 0.5}{SE(\hat{p})},$$

where \hat{p} is the sample proportion and $SE(\hat{p}) = \sqrt{0.5(1 - 0.5)/n}$

- A suitable reject region seems to be $TS > t$, since large TS indicates that p estimated from data is much larger than 0.5.
- But which t should we choose?

	Does not reject H_0	Reject H_0 in favor of H_a
H_0 is true	Correct decision	Type I Error
H_a is true	Type II Error	Correct decision

Types of errors that a test can make

(a) **Type I Error** or **False Positive** or **False Discovery**

- occurs if we reject H_0 and accept H_a when H_0 is in fact true.
- Probability of making a Type I error, also called (significance)

LEVEL of the test:

$$\text{Level of the test, } \alpha = P(\text{Type I Error}|t) = P(\text{Reject } H_0|H_0, t).$$

(b) **Type II Error** or **False Negative**

- Occurs if we fail to reject H_0 when H_0 is false. We failed to make a discovery
- Probability of making a Type II error. :

$$\beta = P(\text{Type II Error}|t) = P(\text{Does not reject } H_0|\theta \in H_a, t).$$

- The quantity $1 - \beta$ is also called the **POWER of the test**

In the working example:

- As $t > 0$ increases, harder to reject H_0 . Then $\alpha \downarrow$, $\beta \uparrow$
- As $t > 0$ decreases, easier to reject H_0 . Then $\alpha \uparrow$, $\beta \downarrow$
- α and β are always inversely related;
 - It is impossible to minimize both at the same time.
- In scientific practice and in drug development, researchers typically consider a Type I error (“False Discovery”) more serious error than a Type II error (“Failure to make a discovery”). In medical application, such as testing for a disease, however, it is often more important to minimize Type II error. Here we proceed with the assumption that Type I error is more important.

- In this case, a value for α is chosen before initiating a hypothesis test.
- Common values for α are 0.01, 0.05 and 0.10;
- Choose the rejection region so that

$$\mathbb{P}(\text{Type I Error}) = \mathbb{P}(t \in RR|H_0) = \alpha$$

- If $\alpha = 0.05$, this choice of the rejection region means that in 5% of the data samples from a population where H_0 is actually true, the test will reject the H_0 .

Type I against Type II errors

- What is the consequence of a Type I error?
 - Conclude that a drug is effective when in fact that it is not.
 - Conclude that a foreign policy is working when in fact that it is not.
 - Ultimately: huge amount of money spent for nothing
- What is the consequence of a Type II error?
 - Conclude that a drug is ineffective when in fact it is a good drug.
 - Conclude that a potentially working foreign policy is not useful.
 - Ultimately: Lost opportunity

Eventually, the choice of balance between type I and type II error depend on a cost-benefit analysis, which is outside of the area of statistics.

Summary Design of the test:

- Set up H_0
- Set up H_a
- Define a reasonable statistic TS .
- Figure out the distribution of TS under H_0

- Choose a small significance level α (like 5%) and find a reject region (RR) so that $\mathbb{P}(\text{make a type I error}) = \mathbb{P}(TS \in RR | H_0 \text{ is true}) = \alpha$.

Application of the test: If the observed value of the TS is in the RR , then reject H_0 in favor of H_a ; otherwise, decide that the test fails to reject H_0 , and conclude that there is no sufficient evidence at the significance level α that H_a is true.

Quiz 3.1.3. What hypothesis states equality or no difference, or no relationship/effect?

- A. statistical hypothesis
- B. null hypothesis
- C. alternative hypothesis

Quiz 3.1.4. A Type I error is when:

- A. We reject the null hypothesis when it is actually true
- B. We obtain the wrong test statistic
- C. We fail to reject the null hypothesis when it's actually false
- D. We reject the alternate hypothesis when it's actually true

Quiz 3.1.5. A level of significance of 5% means:

- A. There's a 5% chance there is an error in test decision.
- B. There's a 5% chance we'll be wrong if we fail to reject the null hypothesis
- C. There's a 5% chance we'll be wrong if we reject the null hypothesis.
- D. The alternative hypothesis is not significant.

3.2 Calculating the Level and Power of a Test

Given a test statistic and a reject region (RR) of a test, how do we find the probabilities of errors α and β ?

- Type I Error:

- Occurs if we reject H_0 when H_0 is true.
- Probability of making a Type I error:

$$\begin{aligned}\alpha &= P(\text{Type I Error}) = P(\text{rejecting } H_0 \mid H_0 \text{ is true}) \\ &= P(\textcolor{red}{TS} \in \textcolor{red}{RR} \mid H_0 \text{ is true}).\end{aligned}$$

- Type II Error:

- Occurs if we fail to reject H_0 when H_0 is false.
- Probability of making a Type II error:

$$\begin{aligned}\beta &= P(\text{Type II Error}) = P(\text{fail to reject } H_0 \mid \theta \in H_a) \\ &= P(\textcolor{red}{TS} \notin \textcolor{red}{RR} \mid \theta \in H_a).\end{aligned}$$

Example 3.2.1. An experimenter has prepared a drug dosage level that she claims will induce sleep for 80% of people suffering from insomnia. After examining the dosage, we feel that her claims regarding the effectiveness of the dosage are inflated. In an attempt to disprove her claim, we administer her prescribed dosage to 20 insomniacs and we observe Y , the number for whom the drug dose induces sleep. We wish to test the hypothesis $H_0 : p = .8$ versus the alternative, $H_a : p < .8$. Assume that the rejection region $\{y \leq 12\}$ is used.

- What is the probability of type I error (level of the test) α ?
- What is the probability of type II error β if $H_a : p = 0.6$?
- What is the probability of type II error β if $H_a : p = 0.4$?

In this example Y is our test statistic (the complete data consists of observations for each insomniac). This statistic is distributed according to the

binomial distribution with parameters $n = 20$ and p . If we assume H_0 then $p = 0.8$. Then we need to calculate

$$\alpha = \mathbb{P}(Y \leq 12 | H_0 : p = 0.8).$$

We can do it using tables or issuing *R* command `pbinom(12, 20, 0.8)`. which gives us the result $\alpha = 0.03214266$. This is the significance level (or size) of this test.

Under the alternative hypothesis $H_a : p = 0.6$, we have

$$\beta = \mathbb{P}(Y > \leq 12 | H_0 : p = 0.6) = 1 - \mathbb{P}(Y \leq \leq 12 | H_0 : p = 0.6).$$

We can calculate it as $\beta = 1 - \text{pbinom}(12, 20, 0.6) = 0.4158929$. This is a rather large probability of error.

If $H_a : p = 0.4$, then $\beta = 1 - \text{pbinom}(12, 20, 0.4) = 0.02102893$. Here both the probability of type I and type II errors are small because in this case it is easy to detect from the data whether a null hypothesis or an alternative is true. The probability that we encounter the data which would be likely under both alternatives is small.

Typically the researchers are more concerned about the type I error (“false discovery rate”) and the test designed in the following fashion.

Suppose we were interested in testing a statement about the parameter θ of the population;

- The null hypothesis $H_0 : \theta = \theta_0$
- The alternative hypothesis could be one of the following
 - $H_a : \theta > \theta_0$ (one-sided test)
 - $H_a : \theta < \theta_0$ (one-sided test)
 - $H_a : \theta \neq \theta_0$ (two-sided test)
- Using the sample data to find an estimator of θ , denote it by $\hat{\theta}$;
- Define the test statistic

$$TS = \frac{\hat{\theta} - \theta_0}{\sigma_{\hat{\theta}}}.$$

- For the reject region (RR):
 - $\{TS > t\}$ (one-sided test)
 - $\{TS < t\}$ (one-sided test)
 - $\{|TS| > t\}$ (two-sided test)
- Cutoff t is chosen such that

$$P(TS \text{ is in RR} \mid H_0 \text{ is true}) = P(TS \text{ is in RR} \mid \theta = \theta_0) = \alpha$$

Quiz 3.2.2. We are interested in this problem: “Is the proportion of babies born male different from 50%?” In the sample of 200 births, we found that 96 babies born were male. We tested the claim using a test with the level of significance 1% and found that the conclusion is “Fail to reject H_0 .” What could we use as interpretation?

- A. The proportion of babies born male is not 0.50.
- B. There is not enough evidence to say that the proportion of babies born male is different from 0.50.
- C. There is not enough evidence to say that the proportion of babies born male is 0.50.

Quiz 3.2.3. The level of significance is the maximum probability of committing a type II error.

- A. True
- B. False

Example 3.2.4. A machine in a factory must be repaired if it produces more than 10% defectives among the large lot of items that it produces in a day. A random sample of 100 items from the day's production contains 15 defectives, and the supervisor says that the machine must be repaired. Does the sample evidence support his decision? Use a test with level .01.

- The null hypothesis $H_0 : p = p_0$ (where $p_0 = 0.1$ here.)
- The alternative hypothesis $H_a : p > p_0$;
- An estimator of p is \hat{p} the sample proportion;
- Define the test statistic

$$TS = \frac{\hat{p} - p_0}{\sigma_{\hat{p}}} = \frac{\hat{p} - p_0}{\sqrt{p_0(1 - p_0)/n}}.$$

- The reject region (RR): $\{TS > t\}$; with the cutoff t chosen such that $P(TS > t | p = p_0) = \alpha$.

Under the assumption H_0 is true, TS is approximately $N(0, 1)$, because the sample size is large ($n = 100$). So the equation becomes a statement about a normal random variable.

$$P(TS > t | p = p_0) = \alpha \Rightarrow P(Z > t) = \alpha,$$

where Z is a normal random variable. We know that it holds when $t = z_\alpha$. So the reject region RR becomes

$$RR : \{TS > z_\alpha\};$$

The problem asks for level $\alpha = 0.01$, so $z_{0.01} = 2.33$, and the reject region is

$$RR : \{TS > 2.33\};$$

By using the data provided in this problem, we calculate the observed value of the test statistic TS as

$$ts = \frac{\hat{p} - p_0}{\sigma_{\hat{p}}} = \frac{0.15 - 0.10}{\sqrt{0.1(1 - 0.1)/100}} = \frac{5}{3} = 1.667.$$

Since ts is NOT in the reject region, we fail to reject H_0 at the level $\alpha = 0.01$. and come to the conclusion that there is NO sufficient evidence to support the statement that the machine must be repaired, at the significance level $\alpha = 0.01$.

Note that if we used a different α , say $\alpha = 0.05$, then $z_{0.05} = 1.645$, and the reject region would be

$$RR : \{TS > 1.645\};$$

Then, ts would be in the reject region, and we would reject H_0 at the level $\alpha = 0.05$. The conclusion would become: at the significance level $\alpha = 0.05$, there is sufficient evidence to support the statement that the machine must be repaired.

So, the decision of the hypothesis test depends on the value of α – the level of tolerance for the type I error. The report about the decision should always specify the value of α .

Not that the decision of a hypothesis test has a random nature! It depends on the realized data. In particular, if H_0 is true, we incorrectly reject it with probability α .

Now let us look at the calculation of the probability of type II error in the situation when we have a large sample and an estimator of the parameter which is distributed approximately normally. Let us consider the same example as before.

Example 3.2.5. Let X_1, \dots, X_n be distributed according to the Bernoulli distribution with parameter p .

- Null hypothesis: $H_0: p = 0.1$
- Alternative hypothesis $H_a: p > 0.1$.
- Test the hypothesis H_0 against alternative H_a .

As in the previous example, we define the test statistic as

$$TS = \frac{\hat{p} - 0.1}{SE(\hat{p})},$$

where \hat{p} is the sample proportion and $SE(\hat{p}) = \sqrt{0.1(1 - 0.1)/n}$. The rejection region is $TS > z_\alpha$, where α is the pre-specified level (=size) of the test, which is the probability of type I error.

Now, what about the probability of type II error for this test? In order to calculate this probability, we need to have a specific alternative hypothesis H_a about the parameter p . Suppose, for example that $H_a : p = 0.15$. This is a natural choice since we observed 15 defective machine out of 100. If we assume that the alternative hypothesis is true, then we know that for large n , the quantity

$$Z = \frac{\hat{p} - 0.15}{\sqrt{0.15(1 - 0.15)/n}}$$

is distributed as a standard normal random variable. Therefore,

$$\hat{p} = 0.15 + (\sqrt{0.15(1 - 0.15)/n})Z,$$

and we re-write the test statistic as

$$TS = \frac{\hat{p} - 0.1}{\sqrt{0.1(1 - 0.1)/n}} = \frac{0.15 - 0.1 + Z\sqrt{0.15(1 - 0.15)/n}}{\sqrt{0.1(1 - 0.1)/n}}$$

So the probability of type II error for the test with level α is

$$\begin{aligned} \beta = \mathbb{P}(TS \leq z_\alpha | H_a) &= \mathbb{P}\left(\frac{0.05 + Z\sqrt{0.15(1 - 0.15)/n}}{\sqrt{0.1(1 - 0.1)/n}} \leq z_\alpha\right) \\ &= \mathbb{P}\left(Z\sqrt{0.15(1 - 0.15)/n} \leq -0.05 + z_\alpha\sqrt{0.1(1 - 0.1)/n}\right) \\ &= \mathbb{P}\left(Z \leq \frac{-0.05 + z_\alpha\sqrt{0.1(1 - 0.1)/n}}{\sqrt{0.15(1 - 0.15)/n}}\right) \end{aligned}$$

and this quantity is easy to evaluate by using software or by referring to tables. If we use $\alpha = 0.01$ and $n = 100$, then $z_\alpha = 2.33$ we calculate that

$$\frac{-0.05 + z_\alpha\sqrt{0.1(1 - 0.1)/n}}{\sqrt{0.15(1 - 0.15)/n}} = 0.5573$$

and

$$\mathbb{P}(Z \leq 0.5573) = 0.71.$$

So the probability that we make an error of the type II is rather large, 71%

Note that the *power* of the test is by definition $1 - \beta$, so we have a method to calculate both the probability of type II error and the power of the test. In this example the power of the test is $1 - 0.71 = 29\%$.

It is important to note that both depend on the value of the parameter under the alternative hypothesis. For example, if we changed our alternative hypothesis to $H_a : p = 0.2$, then the probability of type II error would be equal to

$$\beta = \mathbb{P}\left(Z \leq \frac{-0.1 + z_\alpha \sqrt{0.1(1-0.1)/n}}{\sqrt{0.2(1-0.2)/n}}\right) = \mathbb{P}(Z \leq -0.7525) = 22.6\%$$

This is certainly much better outcome. Under this alternative hypothesis the test has much more statistical power. The power equals $1 - 0.226 = 77.4\%$.

Example 3.2.6 (Covid-19 in New York and California). The total number of cases of Covid-19 in New York State is $\approx 203,123$ with number of total deaths 10,834 as of April 14. The corresponding numbers for California are 25,536 and 782. Test the hypothesis that the mortality rate in New York is higher than the mortality rate in California.

- The null hypothesis $H_0 : \theta = \theta_0$ (where $\theta = p_1 - p_2$ and $\theta_0 = 0$.)
- The alternative hypothesis $H_a : \theta > \theta_0$;
- An estimator of $\theta = p_1 - p_2$ is $\hat{\theta} = \hat{p}_1 - \hat{p}_2$, the difference in sample proportion;
- Find the test statistic

$$TS = \frac{\hat{\theta} - \theta_0}{\sigma_{\hat{\theta}}} = \frac{\hat{p}_1 - \hat{p}_2 - 0}{\sqrt{\frac{p_1(1-p_1)}{n_1} + \frac{p_2(1-p_2)}{n_2}}} \approx ?.$$

- Form the reject region (RR): $TS > z_\alpha$.

An additional difficulty here is that the null hypothesis does not specify the exact values of p_1 and p_2 . It only says that $p_1 = p_2$. For this reason, we

need to estimate p_1 and p_2 . We use the “**pooled sample proportion**”, suggested by the fact that H_0 claims that $p_1 = p_2$.

$$\tilde{p} := \frac{Y_1 + Y_2}{n_1 + n_2}$$

This is the best guess about p_1 and p_2 we can obtain when $p_1 = p_2$ (that is, under the assumption that H_0 is true.)

Using the data provided in this problem, we can calculate:

$$\begin{aligned}\hat{p}_1 &= \frac{10834}{203123} = 0.05333, \\ \hat{p}_2 &= \frac{782}{25536} = 0.03062, \\ \hat{p}_1 - \hat{p}_2 &= 0.02271\end{aligned}$$

$$\begin{aligned}\tilde{p} &= \frac{10834 + 782}{203123 + 25536} = 0.05080, \\ \sigma_{\hat{p}_1 - \hat{p}_2} &= \sqrt{\tilde{p}(1 - \tilde{p})\left(\frac{1}{n_1} + \frac{1}{n_2}\right)} = 0.001457,\end{aligned}$$

and the observed value of the test statistic TS is

$$ts = \frac{\hat{\theta} - \theta_0}{\sigma_{\hat{\theta}}} = \frac{0.05080 - 0}{0.001457} = 15.5789$$

So, in this case it is obvious that the null hypothesis can be rejected at $\alpha = 0.01$. The data give strong support to the hypothesis that the mortality rate in New York is higher than that in California.

Quiz 3.2.7. Which of the following indicates a right-tailed one-sided test?

- A. $H_a : \mu < 15$
- B. $H_0 : \mu < 15$
- C. $H_a : \mu > 15$
- D. $H_0 : \mu > 15$

Quiz 3.2.8. Statistical power is the probability of

- A. making a Type I error.
- B. rejecting the null when it is true.
- C. making a Type II error.
- D. rejecting the null when it is false.

The previous two examples were about testing hypotheses about population proportions.

Now let us look at the hypotheses about population means. We still maintain the assumption that the sample size is large and therefore we can rely on the normality of the parameter estimator distribution.

Example 3.2.9. A random sample of 37 second graders who participated in sports had manual dexterity scores with mean 32.19 and standard deviation 4.34. An independent sample of 37 second graders who did not participate in sports had manual dexterity scores with mean 31.68 and standard deviation 4.56.

- a. Test to see whether sufficient evidence exists to indicate that second graders who participate in sports have a higher mean dexterity score. Use $\alpha = .05$.
- b. For the rejection region used in part (a), calculate β when $\mu_1 - \mu_2 = 3$.

The null hypothesis is $H_0 : \mu_1 = \mu_2$ and the alternative is $H_a : \mu_1 > \mu_2$.

A suitable estimator for $\theta = \mu_1 - \mu_2$ is $\hat{\theta} = \bar{X} - \bar{Y}$, the difference of the sample means. Its variance is

$$\sigma_{\hat{\theta}} = \sqrt{\sigma_1^2/n_1 + \sigma_2^2/n_2}.$$

Since we do not know the exact values of σ_1^2 and σ_2^2 , we will use estimates for these variances. Then, the test statistic is

$$\begin{aligned} TS &= \frac{\hat{\theta} - \theta_0}{\hat{\sigma}_{\hat{\theta}}} = \frac{\bar{X} - \bar{Y} - 0}{\sqrt{s_1^2/n_1 + s_2^2/n_2}} \\ &= \frac{32.19 - 31.68}{\sqrt{4.34^2/37 + 4.56^2/37}} = 0.4928 \end{aligned}$$

Since the samples are relatively large ($n_1 = n_2 > 30$), the test statistic is distributed as a standard normal random variable. Since $\alpha = 0.05$ and $TS \leq z_{0.05} = 1.645$, the test statistic is not in the rejection region and we are not able to reject the null hypothesis. The data does not give enough evidence to indicate that second graders who participate in sports have a higher mean dexterity score.

Now let us consider the second question. What is β if $\mu_1 - \mu_2 = 3$?

In this case, we know that

$$Z = \frac{\bar{X} - \bar{Y} - 3}{\sqrt{s_1^2/n_1 + s_2^2/n_2}}$$

is approximately standard normal random variable. What we need to calculate is the probability that we do not reject the null hypothesis, that is $\mathbb{P}(TS \leq z_\alpha)$. So, we need to express TS in terms of Z :

$$TS = \frac{\bar{X} - \bar{Y}}{\sqrt{s_1^2/n_1 + s_2^2/n_2}} = \frac{\bar{X} - \bar{Y} - 3 + 3}{\sqrt{s_1^2/n_1 + s_2^2/n_2}} = Z + \frac{3}{\sqrt{s_1^2/n_1 + s_2^2/n_2}}$$

Then, the desired probability is

$$\begin{aligned}\beta &= \mathbb{P}(TS \leq z_\alpha) = \mathbb{P}\left[Z + \frac{3}{\sqrt{s_1^2/n_1 + s_2^2/n_2}} \leq z_\alpha\right] \\ &= \mathbb{P}\left[Z \leq z_\alpha - \frac{3}{\sqrt{s_1^2/n_1 + s_2^2/n_2}}\right]\end{aligned}$$

After plugging in numbers, we get

$$\begin{aligned}\beta &= \mathbb{P}\left[Z \leq 1.645 - \frac{3}{\sqrt{(4.34^2 + 4.56^2)/37}}\right] \\ &= \mathbb{P}[Z \leq -1.2538] = \text{pnorm}(-1.2538) = 0.105...\end{aligned}$$

So, $\beta = 10.5\%$ and the power of this test is $1 - \beta = 89.5\%$.

More generally, if we use the test statistic

$$TS = \frac{\hat{\theta} - \theta_0}{\hat{\sigma}_{\hat{\theta}}},$$

where the sample size is large, the estimator $\hat{\theta}$ has an approximately normal distribution, and the standard deviation of the estimator $\hat{\theta}$ is estimated from

the data (not calculated on the basis of the hypothesis), then we can write simple formulas for β .

If the alternative hypothesis is $\theta_a > \theta_0$, and we use the rejection region $TS > z_\alpha$, then

$$\beta = \mathbb{P}\left[Z \leq z_\alpha - \frac{\theta_a - \theta_0}{\hat{\sigma}_{\hat{\theta}}}\right]$$

If the alternative hypothesis is $\theta_a < \theta_0$ and the rejection region is $TS < -z_\alpha$, then

$$\beta = \mathbb{P}\left[Z \geq -z_\alpha - \frac{\theta_a - \theta_0}{\hat{\sigma}_{\hat{\theta}}}\right],$$

which can also be written as

$$\beta = \mathbb{P}\left[Z \leq z_\alpha - \frac{\theta_0 - \theta_a}{\hat{\sigma}_{\hat{\theta}}}\right]$$

by the symmetry of the distribution of the normal random variable.

Finally, if the alternative hypothesis is $\theta_a \neq \theta_0$ and we decided to use the symmetric rejection region $|TS| \geq z_{\alpha/2}$, then

$$\beta = \mathbb{P}\left[Z \leq z_{\alpha/2} - \frac{\theta_a - \theta_0}{\hat{\sigma}_{\hat{\theta}}}\right] - \mathbb{P}\left[Z \leq -z_{\alpha/2} - \frac{\theta_a - \theta_0}{\hat{\sigma}_{\hat{\theta}}}\right],$$

What happens if $\theta_a = \theta_0$, or for example, if $\theta_a = \theta_0 + \varepsilon$, where ε is very small? This is the situation when the alternative hypothesis is barely distinguishable from the null hypothesis. It is easy to see that in this case $\beta = 1 - \alpha$, and the power = α . This is the worst case scenario for the test and we conclude that the power of a test cannot drop down below its size.

Now what happens if $|\theta_a - \theta_0|$ becomes larger. In the case of one-sided hypotheses, it is easy to see from formulas that β declines. In fact, it is possible to check that this also holds for the two-sided hypothesis. The more the alternative differs from the null hypothesis, the less is the probability of type II error β (for a fixed level α).

3.2.1 Additional examples

Example 3.2.10 (Shear strength of soils). Shear strength of soils is a quantity important in civil engineering. Shear strength measurements derived from unconfined

compression tests for two types of soils gave the results shown in the following table (measurements in tons per square foot). Do the soils appear to differ with respect to average shear strength, at the 1% significance level?

Soil Type I	Soil Type II
$n_1 = 30$	$n_2 = 35$
$\bar{y}_1 = 1.65$	$\bar{y}_2 = 1.43$
$s_1 = 0.26$	$s_2 = 0.22$

The null hypothesis $H_0 : \theta = \theta_0$ (where $\theta = \mu_1 - \mu_2$ and $\theta_0 = 0$.) This is simply a different formulation of the hypothesis $H_0 : \mu_1 = \mu_2$.) The alternative hypothesis $H_a : \theta \neq \theta_0$.

An estimator of $\theta = \mu_1 - \mu_2$ is $\hat{\theta} = \bar{Y}_1 - \bar{Y}_2$, the difference in sample mean. So we take the test statistic

$$TS = \frac{\hat{\theta} - \theta_0}{\sigma_{\hat{\theta}}} = \frac{\hat{\theta} - \theta_0}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}} \approx \frac{\hat{\theta} - \theta_0}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}.$$

The rejection region (RR): $\{|TS| > t\}$ and choose the cutoff t so that

$$P(|TS| > t | \theta = \theta_0) = \alpha$$

When H_0 is true, the TS is approximately $N(0, 1)$, and we can take $t = z_{\alpha/2}$.

$$RR : \{|TS| > z_{\alpha/2}\};$$

For $\alpha = 0.01$, $z_{0.01/2} = 2.575$, which gives $RR : \{|TS| > 2.575\}$.

By using the data,

$$ts = \frac{\hat{\theta} - \theta_0}{\sigma_{\hat{\theta}}} = \frac{1.65 - 1.43 - 0}{\sqrt{\frac{0.26^2}{30} + \frac{0.22^2}{35}}} = 3.65.$$

Since ts in the reject region, our decision: reject H_0 at the level $\alpha = 0.01$. We conclude that at the significance level $\alpha = 0.01$, there is sufficient evidence that the soils appear to differ with respect to average shear strength.

Example 3.2.11. A political researcher believes that the fraction p_1 of Republicans strongly in favor of the death penalty is greater than the fraction p_2 of Democrats strongly in favor of the death penalty. He acquired independent random samples of 200 Republicans and 200 Democrats and found 46 Republicans and 34 Democrats strongly favoring the death penalty. Does this evidence provide statistical support for the researcher's belief? Use $\alpha = .05$.

Using the data provided in this problem, we can calculate:

$$\tilde{p} = \frac{Y_1 + Y_2}{n_1 + n_2} = \frac{46 + 34}{400} = 0.2,$$

$$\sigma_{\hat{p}_1 - \hat{p}_2} = \sqrt{\frac{\tilde{p}(1 - \tilde{p})}{n_1} + \frac{\tilde{p}(1 - \tilde{p})}{n_2}} = \sqrt{\frac{2 \times 0.2 \times 0.8}{200}} = 0.04,$$

and the observed value of the test statistic TS is

$$ts = \frac{\hat{\theta} - \theta_0}{\sigma_{\hat{\theta}}} = \frac{(46/200 - 34/200) - 0}{0.04} = 1.5$$

What should be the conclusion?

3.3 Determining the sample size

Now, consider Example 3.2.9 again. Suppose that we are not satisfied that the probability of type II error β is around 10%. What should the sample size that would give $\beta = 5\%$?

Recall that the formula for β is

$$\beta = \mathbb{P}\left[Z \leq z_\alpha - \frac{3}{\sqrt{s_1^2/n_1 + s_2^2/n_2}}\right].$$

We can use the symmetry of the normal distribution and write it as

$$\beta = \mathbb{P}\left[Z > -z_\alpha + \frac{3}{\sqrt{s_1^2/n_1 + s_2^2/n_2}}\right],$$

which can be re-written as

$$z_\beta = -z_\alpha + \frac{3}{\sqrt{s_1^2/n_1 + s_2^2/n_2}}$$

If we assume that $n_1 = n_2$, then this is the same as

$$z_\beta + z_\alpha = \frac{3\sqrt{n}}{\sqrt{s_1^2 + s_2^2}},$$

and the formula for the appropriate sample size is

$$n = (s_1^2 + s_2^2) \left[\frac{z_\beta + z_\alpha}{3} \right]^2,$$

Note that 3 here represents $\theta_a - \theta_0$, so the general formula is

$$n = (s_1^2 + s_2^2) \left[\frac{z_\beta + z_\alpha}{\theta_a - \theta_0} \right]^2,$$

This formula works for both right-tailed and left-tailed one-sided tests. However, there is no simple formula for two-sided hypotheses.

By plugging in the numbers from the example, we get

$$n = (4.34^2 + 4.56^2) \left[\frac{1.645 + 1.645}{3} \right]^2 \approx 47.66$$

Since we cannot use fraction as a sample size, we conclude that sample size $n = 48$ would be sufficient.

3.4 Relation with confidence intervals

Suppose we are still working with large size samples and we know that the estimator $\hat{\theta}$ for a parameter θ has a normal distribution. In fact let us assume that we calculated the confidence interval, with the confidence level α . For concreteness, let us focus on one-sided lower bound confidence interval

$$CI = (\hat{\theta} - z_\alpha \hat{\sigma}_{\hat{\theta}}, \infty).$$

Does it help us with hypothesis testing? Well, the confidence interval says that the true value of the parameter is likely to be larger than $\hat{\theta} - z_\alpha \hat{\sigma}_{\hat{\theta}}$. So if we test the null hypothesis $H_0: \theta = \theta_0$, and it happens that θ_0 is outside of the confidence interval, that is, if $\theta_0 < \hat{\theta} - z_\alpha \hat{\sigma}_{\hat{\theta}}$, then we should reject the null hypothesis.

The only provision here is that the CI should be in agreement with the alternative hypothesis, that is, the alternative hypothesis should be $H_a: \theta > \theta_0$.

If our alternative hypothesis is that $\theta < \theta_0$, then it is more appropriate to consider the upper bound confidence interval:

$$CI = (-\infty, \hat{\theta} + z_\alpha \hat{\sigma}_{\hat{\theta}}).$$

This confidence interval tells us that the true value of the parameter is likely to be large than $\hat{\theta} + z_{\alpha}\hat{\sigma}_{\hat{\theta}}$, so if θ_0 is greater then this quantity we should reject the null hypothesis.

Similarly, if we use a two-sided alternative hypothesis $H_a : \theta \neq \theta_0$, then it is appropriate to use the two-sided confidence interval

$$CI = (\hat{\theta} - z_{\alpha/2}\hat{\sigma}_{\hat{\theta}}, \hat{\theta} + z_{\alpha/2}\hat{\sigma}_{\hat{\theta}}),$$

and reject the null hypothesis $H_0 : \theta = \theta_0$ if θ_0 is outside of the confidence interval.

On the formal level, consider, say, the first case, when the alternative is $H_a : \theta > \theta_0$.

Then the rejection region is

$$TS = \frac{\hat{\theta} - \theta_0}{\hat{\sigma}_{\hat{\theta}}} > z_{\alpha}.$$

But this condition can be re-written as

$$\theta_0 < \hat{\theta} - z_{\alpha}\hat{\sigma}_{\hat{\theta}},$$

and this is exactly the condition that θ_0 is outside of the low-bound confidence interval, as we claimed above. The other two cases can be done similarly.

3.5 p -values

The p -value of a test is useful if one wants to report how strongly the evidence in the data speaks against the null hypothesis.

Recall that we saw several times the situation when for $\alpha = 0.01$ we could not reject the null hypothesis, the evidence was not strong enough, but for $\alpha = 0.05$, we could reject the null. (This is because when $\alpha = 0.05$ we could allow to make type I error more frequently.)

For any data sample if we consider very large α then the test statistic is likely to land in the rejection region, which is very wide in this case and we are likely to reject the test. However, as we gradually decrease α , we

become more conservative, the rejection region shrinks, and at some point we switch from rejecting H_0 for this data sample to saying that there is not enough evidence in the data to support the rejection. This point is called the p -value of the test.

Note especially that unlike the level and the power of the test, the p -value depends both on the test (that is, on the way to calculate the test statistic and the rejection region) and on the data. If the data sample looks more unlikely for the null hypothesis than another sample, that is, if it has a larger test statistic, then the switch from rejection to non-rejection happens later, for smaller α , and p - value for such data sample is *smaller*!

Definition 3.5.1. The p -value is the smallest significance level α at which the observed data indicates that H_0 should be rejected.

While definition above is easy to use, it is a bit difficult to grasp or to explain to a client who does not know what is the significance level of a test. In this case, the following equivalent definition might be useful.

Definition 3.5.2. The p -value is the probability, – calculated assuming that the null hypothesis is true, – of obtaining a value of the test statistic, which is at least as contradictory to H_0 as the value calculated from the available sample.

It is very easy to calculate the p -value. We just set the threshold in the rejection region equal to the observed value of the test statistic and calculate the probability of this rejection region under the null hypothesis.

Say, let the test have the rejection region $RR : \{TS > t\}$ and let ts be the observed value of the test statistic. Then the p -value is $\Pr\{TS > ts|H_0\}$.

In practice, for large sample tests it often boils down to calculating the cumulative function of the standard normal distribution at the test statistic value ts .

Benefits of the p -value:

- It is a universal measure of the strength of the evidence.
- It describes how extreme the data would be if the H_0 were true.

- It answers the question: “Assuming that the null is true, what is the chance of observing a sample like this, or even worse?”

Example 3.5.3. Urban storm water can be contaminated by many sources, including discarded batteries. When ruptured, these batteries release metals of environmental significance. The paper “Urban Battery Litter” (J. Environ. Engr., 2009: 46–57) presented summary data for characteristics of a variety of batteries found in urban areas around Cleveland. A sample of 51 Panasonic AAA batteries gave a sample mean zinc mass of 2.06 g. and a sample standard deviation of .141 g. Does this data provide compelling evidence for concluding that the population mean zinc mass exceeds 2.0 g.?

With m denoting the true average zinc mass for such batteries, the relevant hypotheses are $H_0 : m = 2.0$ versus $H_a : m > 2.0$. The sample size is large enough so that a z -test can be used without making any specific assumption about the shape of the population distribution. The test statistic value is

$$z = \frac{\bar{x} - 2.0}{s/\sqrt{n}} = \frac{2.06 - 2.0}{.141/\sqrt{51}} = 3.04$$

So, we calculate the p-value:

$$p - \text{value} = \mathbb{P}(Z > 3.04) = 1 - \text{pnorm}(3.04) = 0.118\%$$

This means that the null hypothesis would be rejected by tests with $\alpha = 5\%$, $\alpha = 1\%$, and even with $\alpha = 0.2\%$, although it could not be rejected at the level of $\alpha = 0.1\%$. We would conclude that the sample appears to highly contradict the null hypothesis, and so there is a compelling evidence that the population mean zinc mass exceeds 2.0 g.

p -values for large sample tests (aka z -tests)

1. The parameter θ is one of the following: μ , p , $\mu_1 - \mu_2$ and $p_1 - p_2$;
2. The sample size n is large enough.
3. The null hypothesis is $H_0 : \theta = \theta_0$

4. The test statistic is

$$TS = \frac{\hat{\theta} - \theta_0}{\sigma_{\hat{\theta}}} \sim N(0, 1) \text{ under } H_0$$

and the observed test statistic using the given data is ts ;

If the alternative hypothesis is

- $H_a : \theta > \theta_0$ then $p\text{-value} = P(TS > ts | \theta = \theta_0) = 1 - \Phi(ts)$
- $H_a : \theta < \theta_0$ then $p\text{-value} = P(TS < ts | \theta = \theta_0) = \Phi(ts)$
- $H_a : \theta \neq \theta_0$ then $p\text{-value} = P(|TS| > |ts| | \theta = \theta_0)$
 $= P(TS > |ts| | \theta = \theta_0) + P(TS < -|ts| | \theta = \theta_0) = 2(1 - \Phi(|ts|))$

Quiz

Suppose, we are interested in testing $H_0: \mu = \mu_0$ against $H_a: \mu > \mu_0$. We will reject H_0 at level $\alpha = 0.05$ if μ_0 is

- A. larger than 95% upper confidence bound for μ .
- B. larger than 95% lower confidence bound for μ .
- C. smaller than 95% upper confidence bound for μ .
- D. smaller than 95% lower confidence bound for μ .

Quiz

An educator is interested in determining the number of hours of TV watched by 4-year-old children. She wants to show that the average number of hours watched per day is more than 4 hours. To test her claim she took a random sample of 100 youngsters. Which of the following values for the sample mean would have the largest p -value associated with it.

- A. 2
- B. 3.9
- C. 4
- D. 5

3.6 Small-sample hypothesis tests for population means

If the sample size is small ($n < 30$) then we cannot hope that the Central Limit Theorem will ensure that the test statistic

$$TS = \frac{\hat{\theta} - \theta_0}{\hat{\sigma}_{\hat{\theta}}}$$

has the standard normal distribution. In this case the only way out is to make sure that the data is at least approximately normal, perhaps by applying an appropriate transformation to the data.

From now on, in this section we will assume that the data is normal. Even in this case, the distribution of the test statistic differs significantly from the normal distribution. This means that when we calculate the probabilities of type I and II errors, or when we calculate the p -values, we cannot calculate probabilities like

$$\mathbb{P}(TS > x)$$

as if the TS were a standard normal random variable. This would result in wrong probabilities.

Luckily, the distribution for this test statistics is still known and can be calculated by a computer algorithm. It can also be found in tables.

$$TS = \frac{\hat{\theta} - \theta_0}{\sigma_{\hat{\theta}}} \sim t\text{-distribution if } H_0 \text{ is true } (\theta = \theta_0),$$

The degrees of freedom for t -distributions depends on whether $\theta = \mu$ or $\mu_1 - \mu_2$.

If we are interested in testing of $H_0 : \mu = \mu_0$, then we use the test statistic

$$TS = \frac{\bar{Y} - \mu_0}{S/\sqrt{n}} \sim t_{n-1} \text{ if } H_0 \text{ is true } (\mu = \mu_0)$$

If we have two samples, X_1, \dots, X_{n_1} and Y_1, \dots, Y_{n_2} , with population means μ_1 and μ_2 , respectively, then we are often interested in testing $H_0 : \mu_1 - \mu_2 = \theta_0$.

$$v = \frac{\left(\frac{s_1^2}{m} + \frac{s_2^2}{n}\right)^2}{\frac{(s_1^2/m)^2}{m-1} + \frac{(s_2^2/n)^2}{n-1}} = \frac{[(se_1)^2 + (se_2)^2]^2}{\frac{(se_1)^4}{m-1} + \frac{(se_2)^4}{n-1}}$$

where

$$se_1 = \frac{s_1}{\sqrt{m}} \quad se_2 = \frac{s_2}{\sqrt{n}}$$

(round v down to the nearest integer).

Figure 3.1: Degrees of freedom for the test statistic when the variances are not the same

Here, two different situations are possible. A bit simple situation is when we can assume that the variances in two samples are the same $\sigma_1^2 = \sigma_2^2 = \sigma^2$. (We could check this assumption by an appropriate test!) Then we can use the test statistic

$$TS = \frac{\bar{X} - \bar{Y} - \theta_0}{S_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}} \sim t_{n_1+n_2-2} \text{ if } H_0 \text{ is true } (\theta = \theta_0);$$

where we use the pooled-sample standard deviation as an estimator for σ^2 .

$$S_p = \sqrt{S_p^2} = \sqrt{\frac{(n_1-1)S_1^2 + (n_2-1)S_2^2}{n_1+n_2-2}} = \dots$$

In this case the degrees of freedom of the t -distribution are $df = n_1 + n_2 - 2$.

A more difficult situation arises when we cannot assume that the variances σ_1^2 and σ_2^2 are equal. Then we have to use this test statistic:

$$TS = \frac{\bar{X} - \bar{Y} - \theta_0}{\sqrt{\frac{S_1^2}{n_1} + \frac{S_2^2}{n_2}}},$$

where S_1^2 and S_2^2 are sample variances in the two samples. It turns out that the distribution of this TS is *approximately* a t -distribution, but the formula for the degrees of freedom is quite complicated. See Figure 3.1.

Note, however, that some researchers suggested that this procedure should be used if there are doubts about whether the variances are same.

After the distribution of the test-statistic is determined the rest is simple, one would only need to replace z_α with the t_α that is calculated for the t -distribution with correct number of degrees of freedom.

- $H_a : \theta > \theta_0 \Leftrightarrow RR : \{TS > t_\alpha\}$
- $H_a : \theta < \theta_0 \Leftrightarrow RR : \{TS < -t_\alpha\}$
- $H_a : \theta \neq \theta_0 \Leftrightarrow RR : \{|TS| > t_{\alpha/2}\}$

The quantities t_α can be found from the tables or by using the R command `qt`. In particular t_α for ν degrees of freedom can be calculated as `qt(1 - α , ν)`.

The calculation of the probability of type II error β and the power $1 - \beta$ is in fact very similar to the calculations in the case of the normal distribution. Again, one only needs to use the t -distribution with the correct number of degrees of freedom instead of the standard normal distribution.

This is also true for p -values.

- $H_a : \theta > \theta_0 \Leftrightarrow p\text{-value} = P(TS > ts | \theta = \theta_0)$
- $H_a : \theta < \theta_0 \Leftrightarrow p\text{-value} = P(TS < ts | \theta = \theta_0)$
- $H_a : \theta \neq \theta_0 \Leftrightarrow p\text{-value} = P(|TS| > |ts| | \theta = \theta_0) = 2P(TS > |ts|)$

where the test statistic has the t -distribution with an appropriate number of degrees of freedom. The tables only give a range for p -value. For precise probability, one must use R command `pt(a, df)` whose output is $\mathbb{P}(T < a)$ where $T \sim t(df)$.

Example 3.6.1. An Article in *American Demographics* investigated consumer habits at the mall. We tend to spend the most money when shopping on weekends, particularly on Sundays between 4:00 and 6:00PM, while Wednesday-morning shoppers spend the least.

Independent random samples of weekend and weekday shoppers were selected and the amount spent per trip to the mall was recorded as shown in the following table:

Weekends	Weekdays
$n_1 = 20$	$n_2 = 20$
$\bar{y}_1 = \$78$	$\bar{y}_2 = \$67$
$s_1 = \$22$	$s_2 = \$20$

- Is there sufficient evidence to claim that there is a difference in the average amount spent per trip on weekends and weekdays? Use $\alpha = 0.05$.
- What is the attained significance level (p -value)?
- What if $n_1 = 10$ and $n_2 = 10$?

Let us do this example for $n_1 = n_2 = 20$ and assume that $\sigma_1^2 = \sigma_2^2$.

Our null hypothesis is that $\mu_1 = \mu_2$, so we want to calculate

$$TS = \frac{\bar{X} - \bar{Y}}{S_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}}$$

where

$$\begin{aligned} S_p &= \sqrt{S_p^2} = \sqrt{\frac{(n_1 - 1)s_1^2 + (n_2 - 1)s_2^2}{n_1 + n_2 - 2}} = \sqrt{\frac{19 \times 22^2 + 19 \times 20^2}{20 + 20 - 2}} \\ &= 21.0238 \end{aligned}$$

So,

$$TS = \frac{78 - 67}{21.0238 \sqrt{\frac{1}{20} + \frac{1}{20}}} = 1.654556$$

For the normal distribution $z_{0.05} = 1.645$ so the test would reject the null hypothesis if the sample were large.

If we want to test H_0 at the level $\alpha = 0.05$ and use the t -distribution we want t_α for $\nu = 20 + 20 - 2 = 38$.

We calculate it as $qt(.95, 38) = 1.685954$, so we conclude that the evidence is not sufficient to reject H_0 at the level $\alpha = 0.05$.

The p -value can be calculated as $1 - pt(1.654556, 38) = 5.31\%$.

3.7 Hypothesis testing for population variances

Occasionally, we are interested in testing variances. The most frequent example is when we test equality of variance in two samples, in order to see if the corresponding populations are really different in a certain aspect. Sometimes we might be interested to see that the variance does not exceed a certain threshold. This problem arises in quality control.

Let us consider first testing the hypothesis $H^0 : \sigma^2 = \sigma_0^2$. If the sample is large, then we can use

$$TS = \frac{S^2 - \sigma_0^2}{\sigma_{S^2}},$$

with a suitable estimator for σ_{S^2} .

This approach does not generalize easily to small samples since it is difficult to calculate the exact distribution of the ratio. So we look at an alternative method, which works both for large and small samples.

So let us assume that X_1, \dots, X_n are from a Normal distribution $N(\mu, \sigma^2)$ with unknown mean μ and unknown variance σ^2 .

We use the result that we know from the section about variance estimation.

$$TS = \frac{(n-1)S^2}{\sigma_0^2} \sim \chi^2(n-1) \text{ when } H_0 \text{ is true}$$

For the case when the alternative hypothesis is $H_a : \sigma^2 > \sigma_0^2$, the rejection region is similar to the RR in z - and t -tests:

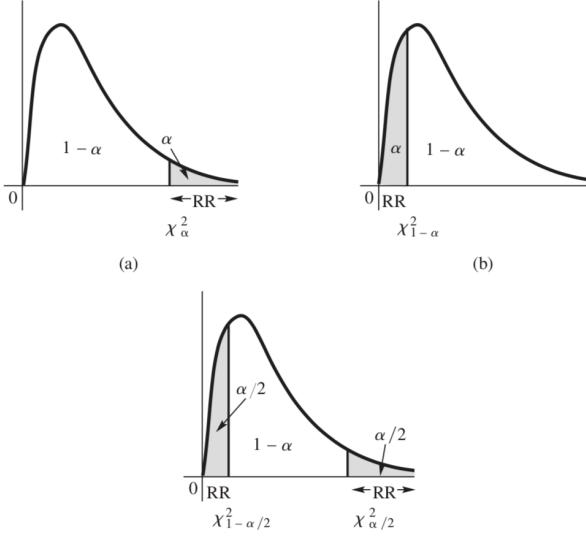
$$RR : TS > \chi_\alpha^2(n-1),$$

where $\chi_\alpha^2(n-1)$ solves the equation $\mathbb{P}(TS > x) = \alpha$ if we know that TS is distributed as a χ^2 random variable with $n-1$ degrees of freedom.

The R command for calculating this quantity is `qchisq(1 - α , $n-1$)`.

FIGURE 10.10

Rejection regions
RR for testing
 $H_0 : \sigma^2 = \sigma_0^2$ versus
(a) $H_a : \sigma^2 > \sigma_0^2$;
(b) $H_a : \sigma^2 < \sigma_0^2$;
and (c) $H_a : \sigma^2 \neq \sigma_0^2$



For the alternative hypothesis $H_a : \sigma^2 < \sigma_0^2$, there is some difference from the case of z - or t -tests because the χ^2 distribution is not symmetric relative to zero. Instead of using $-\chi_\alpha^2(n-1)$ as a threshold, we use $\chi_{1-\alpha}^2(n-1)$. So the rejection region in this case is

$$RR : TS < \chi_{1-\alpha}^2(n-1),$$

Finally, if the alternative hypothesis is $H_a : \sigma^2 \neq \sigma_0^2$, then it is conventional to use the following rejection region:

$$RR : TS < \chi_{1-\alpha/2}^2(n-1) \text{ or } TS > \chi_{\alpha/2}^2(n-1)$$

Correspondingly, the p-values for these alternative hypotheses are as follows. If $H_a : \sigma^2 > \sigma_0^2$, then

$$p\text{-value} = P(TS > ts) = 1 - \text{pchisq}(ts, n-1).$$

If $H_a : \sigma^2 < \sigma_0^2$, then

$$p\text{-value} = P(TS < ts) = \text{pchisq}(ts, n-1)$$

If $H_a : \sigma^2 \neq \sigma_0^2$, then we need to think a bit harder. When we decrease α then at some point either $\chi_{\alpha/2}^2$ or $\chi_{1-\alpha/2}^2$ hits the value of the tests statistic

ts that was realized in the sample. At this moment the test stops rejecting the null hypothesis. So if $\chi_{\alpha/2}^2$ hits ts first, then we conclude that this is the critical α^* which is equal the p-value. Hence in this case the p-value equals

$$\alpha^* = 2\mathbb{P}(TS > \chi_{\alpha^*/2}^2) = 2\mathbb{P}(TS > ts)$$

Note that at that moment $ts > 1 - \alpha^*/2$ so $\mathbb{P}(TS < ts) > 1 - \alpha^*/2$ and $2P(TS < ts) > 2 - \alpha^* > 1$.

If $\chi_{1-\alpha/2}^2$ hits ts first, then by a similar argument we find that p-value equals

$$\alpha^* = 2\mathbb{P}(TS < \chi_{1-\alpha^*/2}^2) = 2\mathbb{P}(TS < ts)$$

and $P(TS > ts) > 1$.

So we have two candidates for the p-value: $2P(TS < ts)$ and $2P(TS > ts)$, and we know that if one of them is indeed the p-value (and so less than 1 as a probability), then the other is greater than 1. So we can simply choose the minimal of these two numbers. In summary,

$$p\text{-value} = 2 \times \min[P(TS > ts), P(TS < ts)]$$

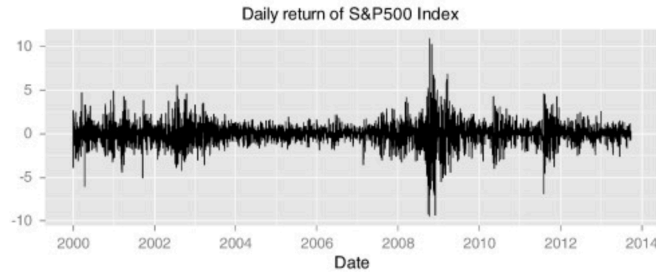
In case you pick up the wrong one between $P(TS > ts)$ or $P(TS < ts)$, your answer will exceed 1, which is an immediate warning sign because probability cannot be greater than 1.

Example 3.7.1. An experimenter was convinced that the variability in his measuring equipment results in a standard deviation of 2. Sixteen measurements yielded $s^2 = 6.1$. Do the data disagree with his claim? Determine the p -value for the test. What would you conclude if you chose $\alpha = 0.05$?

- $H_0 : \sigma^2 = 4$ and $H_a : \sigma^2 \neq 4$
- Test statistic: $\frac{(n-1)S^2}{4}$
- Observed value for test statistic is $ts = \frac{(16-1)6.1}{4} = 22.875$
- R gives $pchisq(22.875, 15) = 0.9132$ and $1 - pchisq(22.875, 15) = 0.0868$.

- Hence $p\text{-value} = 2 \times 0.0868 = 17.36\% > 5\%$.

We conclude that the data do not give enough evidence to disagree with his claim.



Now consider the test for equality of variances in two population. In some situation, a researcher is interested to know whether the data variation in two samples indicated the different variances in corresponding populations. For example:

- comparing precision of two measuring instruments;
- the variation in quality of a product at two locations or at two different time periods;
- variation in scores for two test procedures.
- variation in outcomes for two medical procedures.
- variation in market returns for two time periods or in two different countries.

Suppose that we have two samples with n_1 and n_2 observations, respectively, from the normal distributions with variances σ_1^2 and σ_2^2 respectively. We want to test the hypothesis $H_0 : \sigma_1^2 = \sigma_2^2$ against the alternative $H_a : \sigma_1^2 > \sigma_2^2$.

If S_1^2 and S_2^2 are sample variances, then define the test statistic

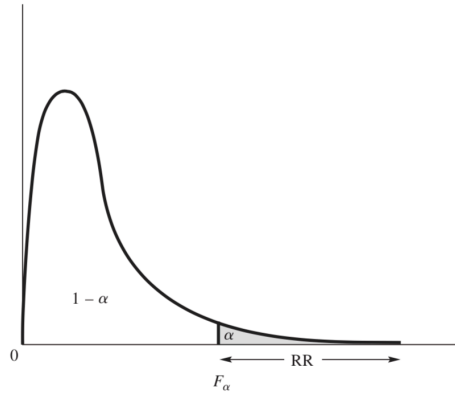
$$TS = \frac{S_1^2}{S_2^2}.$$

Under the null hypothesis this ratio is distributed as so-called **F-distribution** with $n_1 - 1$ and $n_2 - 1$ degrees of freedom. (F is for Fisher, who designed the test.)

So we can use the Rejection Region

$$RR = \{TS > F_\alpha\}$$

FIGURE 10.12
Rejection region
RR for testing
 $H_0 : \sigma_1^2 = \sigma_2^2$ versus
 $H_a : \sigma_1^2 > \sigma_2^2$



If the hypothesis $H_0 : \sigma_1^2 = \sigma_2^2$ but we test against the alternative $H_a : \sigma_1^2 < \sigma_2^2$. (instead of $H_a : \sigma_1^2 > \sigma_2^2$) then we can simply use

$$TS = \frac{S_2^2}{S_1^2}.$$

instead of

$$TS = \frac{S_1^2}{S_2^2}.$$

The new tests statistics is distributed as F -random variable with degrees of freedom $n_2 - 1$ and $n_1 - 1$.

What if we want to test H_0 against the alternative hypothesis $H_a : \sigma_1^2 \neq \sigma_2^2$?

It turns out that in this case we can use the test statistics

$$TS = \frac{S_1^2}{S_2^2},$$

which is distributed as the F -random variable with $n_1 - 1$ and $n_2 - 1$ degrees of freedom and use the rejection region

$$RR = \left\{ \frac{1}{TS} > F_{n_1-1; \alpha/2}^{n_2-1} \text{ or } TS > F_{n_2-1; \alpha/2}^{n_1-1} \right\}$$

Notice also the degrees of freedom in the numerator and denominator of the thresholds!

It is worthwhile to repeat: the test is very sensitive to the assumption that the data are normally distributed.

Example 3.7.2. A study was conducted by the Florida Game and Fish Commission to assess the amounts of chemical residues found in the brain tissue of brown pelicans. In a test for DDT, random samples of $n_1 = 10$ juveniles and $n_2 = 13$ nestlings produced the results shown in the accompanying table (measurements in parts per million, ppm).

Juveniles	Nestlings
$n_1 = 10$	$n_2 = 13$
$\bar{y}_1 = .041$	$\bar{y}_2 = .026$
$s_1 = .017$	$s_2 = .006$

Are you willing to assume that the underlying population variances are equal? Test $\sigma_1^2 = \sigma_2^2$ against $\sigma_1^2 > \sigma_2^2$ at $\alpha = 0.01$. What is the p -value?

The test statistic is

$$TS = \frac{0.017^2}{0.006^2} = 8.027778.$$

$$p\text{-value} = 1 - \text{pf}(8.027778, 9, 12) = 0.07\%.$$

So we would reject the null at 1% level.

3.8 Neyman - Pearson Lemma and Uniformly Most Powerful Tests

So far we talked about specific tests and have not be concerned with evaluating and comparing tests.

Recall that for the estimation problem we had the concept of the Minimal Variance Unbiased Estimator. If we can find such an estimator, then we could agree, that this is the best estimator among available.

For a test the natural measure of its goodness is its power. If we can have two tests with the same α , we will prefer the test that has the large power.

Note however that the power of the test is a function of the parameter under the alternative hypothesis. So if one test has larger power than another under one value of the parameter θ_a , it can actually has smaller power under another value of θ_a .

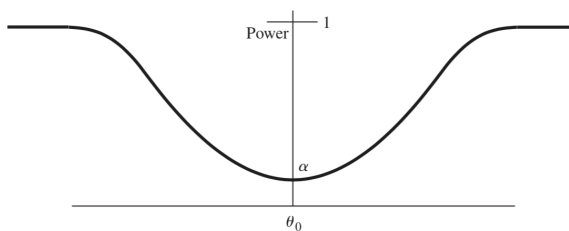
Recall that the power of the test $= 1 - \beta = \Pr(\text{reject } H_0 | \text{if } H_a \text{ is true})$. It can be calculated only if H_a is a simple hypothesis.

Definition 3.8.1. A hypothesis is said to be **simple** if this hypothesis uniquely specifies the distribution of the population from which the sample is taken. Any hypothesis which is not simple is called a **composite** hypothesis.

For example, the hypothesis $\theta = 2$ is simple, the hypothesis $\theta > 2$ is composite.

Power is a curve, it depends on the value of the parameter in the alternative hypothesis. Can we build a test with the “best” power curve?

FIGURE 10.13
A typical power curve for the test of $H_0 : \theta = \theta_0$ against the alternative $H_a : \theta \neq \theta_0$



First, let us not to be too ambitious and try to **find the test with the maximum power** when the significance level is α and when the alternative hypothesis is simple, $H_a : \theta = \theta_a$.

Theorem 3.8.2 (The Neyman-Pearson Lemma). *For testing between $H_0 : \theta = \theta_0$ vs $H_a : \theta = \theta_a$, the test with the reject region*

$$RR = \left\{ \frac{L(\theta_0|X)}{L(\theta_a|X)} < t \right\}$$

where t is chosen so that $\mathbb{P}(\frac{L(\theta_0)}{L(\theta_a)} < t | \theta = \theta_0) = \alpha$, **is the most powerful α -level test** for H_0 versus H_a .

In other words, if the alternative hypothesis is simple: $\theta = \theta_\alpha$, the best test statistic is

$$TS = \frac{L(\theta_0|X)}{L(\theta_\alpha|X)}$$

This TS measures how likely is the data under the null hypothesis compared with its likelihood under the alternative hypothesis. You reject H_0 if the ratio is too small, with the threshold chose so that the level of this test is α .

The theorem says that this test has the largest power to reject H_0 (among the tests with the same α) **provided** the alternative is fixed at θ_α .

In order to construct the Neyman-Pearson test, we need to know the distribution of the test statistics, which is not always easy. Here is an example where the distribution of $\frac{L(\theta_0)}{L(\theta_\alpha)}$ is not too hard to find.

Example 3.8.3. Suppose we have just one observation in the sample, $Y \sim f(y|\theta) = \theta y^{\theta-1} \mathbb{1}_{\{0 < y < 1\}}$. Find the most powerful test for $H_0 : \theta = 2$ against $H_a : \theta = 1$ at significance level $\alpha = 0.05$.

The likelihood here is simply $L(\theta|y) = \theta y^{\theta-1}$. (Only one observation - so no products.)

The ratio is

$$\frac{L(\theta_0|y)}{L(\theta_\alpha|y)} = \frac{\theta_0}{\theta_\alpha} y^{\theta_0 - \theta_\alpha}$$

If we want the size of the test be α , we should have

$$\begin{aligned} \Pr \left[\frac{L(\theta_0|Y)}{L(\theta_\alpha|Y)} < t \middle| H_0 \right] &= \Pr \left[\frac{\theta_0}{\theta_\alpha} Y^{\theta_0 - \theta_\alpha} < t \middle| H_0 \right] = \\ &= \Pr \left[Y < \left(\frac{\theta_\alpha}{\theta_0} t \right)^{1/(\theta_0 - \theta_\alpha)} \middle| H_0 \right] = \alpha. \end{aligned}$$

Under the null hypothesis, Y has density $\theta_0 y^{\theta_0-1}$, so we can calculate the cumulative distribution function as $F_Y(y) = y^{\theta_0}$, and

$$\begin{aligned} \Pr \left[\frac{L(\theta_0|Y)}{L(\theta_\alpha|Y)} < t \middle| H_0 \right] &= \Pr \left[Y < \left(\frac{\theta_\alpha}{\theta_0} t \right)^{1/(\theta_0 - \theta_\alpha)} \middle| H_0 \right] \\ &= \left(\frac{\theta_\alpha}{\theta_0} t \right)^{\theta_0/(\theta_0 - \theta_\alpha)} = \alpha. \end{aligned}$$

Hence, the threshold in the test should be set to

$$t = \frac{\theta_0}{\theta_a} \alpha^{(\theta_0 - \theta_a)/\theta_0}.$$

In our example, the test statistic is $\frac{\theta_0}{\theta_a} Y^{\theta_0 - \theta_a} = 2Y$, and

$$t = \frac{2}{1} 0.05^{(2-1)/2} = 2 \times 0.2236.$$

Equivalently, the most powerful test with $\alpha = 0.05$ in this case has $RR = \{Y < 0.2236\}$

In N-P lemma, the test is guaranteed to be most powerful level- α test against a specific alternative hypothesis. What if we try a different alternative hypothesis?

Consider the previous example: We found that the most powerful test has the rejection region:

$$\left\{ \frac{\theta_0}{\theta_a} Y^{\theta_0 - \theta_a} < \frac{\theta_0}{\theta_a} \alpha^{(\theta_0 - \theta_a)/\theta_0} \right\} = \left\{ Y^{\theta_0 - \theta_a} < \alpha^{(\theta_0 - \theta_a)/\theta_0} \right\}$$

If $\theta_a < \theta_0$, then we can take the power $1/(\theta_0 - \theta_a)$ on both sides and get that the rejection region is

$$\left\{ Y < \alpha^{1/\theta_0} \right\}.$$

It is the same for all $\theta_a < \theta_0$. But if $\theta_a > \theta_0$ then we would get a completely different test:

$$\left\{ Y > (1 - \alpha)^{1/\theta_0} \right\}$$

Say if in the previous example we choose $H_a : \theta = 4$ then we would get test $RR = \{Y > \dots\}$.

When a test (that is a test statistic TS and a rejection region RR) maximizes the power for every value of $\theta \in \Omega_a$, it is said to be a **uniformly most powerful** (or UMP) test for $H_0 : \theta = \theta_0$ versus composite hypothesis $H_a : \theta \in \Omega_a$.

For example, in our previous example, the Neyman-Pearson test is the UMP for $H_a : \theta > \theta_0$. It is also UMP for $H_a : \theta < \theta_0$. (Note that in this case

it is actually a different test. We call it the Neyman-Pearson test because it was obtained by applying the Neyman - Pearson lemma to a specific $\theta_a < \theta$. It is just happened in this example that all these test coincide for all $\theta_a < \theta$.) However, in case of the two-sided hypothesis, $H_a : \theta_a \neq \theta_0$, Neyman-Pearson is not helpful because it gives two different tests depending on whether $\theta_a < \theta_0$ or $\theta_a > \theta_0$. In fact, in this case there is no uniformly most powerful test.

In many cases, even for one sided hypothesis, UMP tests do not exists. However, they are especially rare if the alternative is two sided $H_a : \theta \neq \theta_0$, or if we test a vector parameter and the alternative hypothesis is not simple. (That is, the alternative hypothesis is not just a specific value $\vec{\theta}_a$, for which the Neyman-Pearson lemma would give us a most powerful test.)

Example 3.8.4 (Neyman-Pearson lemma applied to normal data). $Y_1, \dots, Y_n \sim N(\mu, \sigma)$. Consider $H_0 : \mu = \mu_0$ versus $H_a : \mu = \mu_1$. We assume that σ^2 is KNOWN (fixed). Otherwise the N-P lemma is not applicable: the hypothesis is not simple. What is the most powerful test with level α ?

The likelihood is

$$L(\mu) = (2\pi\sigma^2)^{-n/2} \exp \left\{ -\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \mu)^2 \right\}$$

The ratio of the likelihoods defined in the Neyman Pearson lemma is:

$$\begin{aligned} \frac{L(\mu_0)}{L(\mu_1)} &= \frac{(2\pi\sigma^2)^{-n/2} \exp \left\{ -\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \mu_0)^2 \right\}}{(2\pi\sigma^2)^{-n/2} \exp \left\{ -\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \mu_1)^2 \right\}} \\ &= \exp \left\{ -\frac{1}{2\sigma^2} \left[\sum_{i=1}^n (y_i - \mu_0)^2 - \sum_{i=1}^n (y_i - \mu_1)^2 \right] \right\} \\ &= \exp \left\{ -\frac{1}{2\sigma^2} \left[\sum_{i=1}^n 2y_i(\mu_1 - \mu_0) + \mu_0^2 - \mu_1^2 \right] \right\} \end{aligned}$$

The Neyman Pearson Lemma says that the most powerful test is the one

with some appropriate threshold t which rejects H_0 when

$$\begin{aligned} & \exp \left\{ -\frac{1}{2\sigma^2} \left[\sum_{i=1}^n 2y_i(\mu_1 - \mu_0) + \mu_0^2 - \mu_1^2 \right] \right\} < t \\ \iff & -\frac{1}{2\sigma^2} \left[\sum_{i=1}^n 2y_i(\mu_1 - \mu_0) + \mu_0^2 - \mu_1^2 \right] < t' \\ \iff & \sum_{i=1}^n 2y_i(\mu_1 - \mu_0) + \mu_0^2 - \mu_1^2 > t'' \\ & \iff 2n\bar{y}(\mu_1 - \mu_0) > t''' \end{aligned}$$

That is, the test tells us to reject H_0 when

$$\begin{aligned} \bar{y} &> t'''/(2n(\mu_1 - \mu_0)), \text{ if } \mu_1 - \mu_0 > 0, \text{ or,} \\ \bar{y} &< t'''/(2n(\mu_1 - \mu_0)), \text{ if } \mu_1 - \mu_0 < 0. \end{aligned}$$

$$\begin{aligned} \bar{y} &> A, \text{ if } \mu_1 - \mu_0 > 0, \text{ or,} \\ \bar{y} &< B, \text{ if } \mu_1 - \mu_0 < 0, \end{aligned}$$

where the thresholds A and B are chosen so that the level of the test is α . Indeed, such a test would be what our intuition would have driven us.

The NP lemma provides theoretical justification for this. Is this test uniformly most powerful among all the tests against the composite alternative hypothesis $H_a : \mu > \mu_0$? the composite $H_a : \mu < \mu_0$?

Can we construct uniformly most powerful tests for $H_a : \mu \neq \mu_0$?

3.9 Likelihood ratio test

Theoretical Question #2 How do we design a test?

Suppose we have a model with many parameters $\vec{\theta} = (\theta_1, \dots, \theta_k)$, and want to test that one or more of the parameters is 0. Or may be we want to test some relationship between parameters, such as $\theta_1 = \theta_2$, or more generally that $c_1\theta_1 + c_2\theta_2 + \dots + c_k\theta_k = 0$.

How do we test such a hypothesis?

Let Ω_0 be the **set of parameters that satisfy our null hypothesis**. For example, it can be that Ω_0 = the set of all parameters $\vec{\theta}$ such that $\theta_1 = \theta_2$, and all other parameters can be arbitrary. Note that the null hypothesis here can be a composite hypothesis, since there can be many $\vec{\theta}$ that satisfy this condition. All parameters which are not needed to formulate the null hypothesis are often called *nuisance parameters*.

Then, let Ω_a be the set of all possible alternative values for parameters. For example, our alternative can be that $\theta_1 > \theta_2$ and all other θ_i are arbitrary. The alternative hypothesis is typically a composite hypothesis in practical applications.

One approach for the design of a statistical test is to find an estimator for a parameter that encapsulate our hypothesis, and then use our knowledge about the distribution of this estimator. For example, if we test the hypothesis $\theta_1 = \theta_2$, then we can find the ML estimator for the difference $\theta_1 - \theta_2$ and use the fact that in large samples this estimator is normal and that it is possible to calculate its variance. This is a very useful approach. Its deficiency is that we need to obtain the estimator and its variance before we are able to construct the test. In addition, the variance often depends on the true value of parameter θ so if our null hypothesis is not simple but has some nuisance parameters, then we are in trouble.

The second approach is based on the Neyman-Pearson lemma and uses the ratio $\frac{L(\vec{\theta}_0|\vec{Y})}{L(\vec{\theta}_a|\vec{Y})}$ as the test statistic. This will give the most powerful test. The deficiencies is that we need to find the distribution of this test statistic. Additionally, this approach is quite restrictive. It requires both the null and alternative hypotheses to be **simple** and not composite.

In this section we consider the third alternative, which is often very convenient since it works for composite hypotheses, and in case of large samples requires essentially no calculation except the maximization of some likelihood functions.

Define the total feasible parameter set $\Omega = \Omega_0 \cup \Omega_a$. Define the likelihood

ratio statistic by

$$\lambda = \frac{\max_{\vec{\theta} \in \Omega_0} L(\vec{\theta}|\vec{Y})}{\max_{\vec{\theta} \in \Omega} L(\vec{\theta}|\vec{Y})}$$

where $L(\vec{\theta}|\vec{Y})$ is the likelihood of the vector parameter $\vec{\theta}$ given that the observed data is $\vec{Y} = (Y_1, \dots, Y_n)$.

Use the rejection region $RR = \{\lambda < k\}$, where the threshold k is determined by the requirement that the level of the test is α .

This appears to be not an especially useful since we need to do two constrained maximizations and we do not know the distribution of λ , so we cannot calculate the threshold k . In fact, it appears that λ is a quite complicated function of the data: it is the ratio of two constrained maximums of the likelihoods, which are itself complicated functions!

The power of this method is that we can do the maximization numerically and this is a relatively easy given enough computing power. The most important fact, however, is that k can be calculated efficiently when the data sample is large.

Conceptually, the likelihood ratio test makes a lot of sense.

- If H_0 is true, then **with high probability**, the constrained maximum likelihood (with maximum over Ω_0) would be close to unconstrained maximum likelihood (maximum over Ω), the denominator would give the same result as the numerator, and λ would be around 1.
- If H_0 is false while H_a is true, then the unconstrained maximum likelihood would be much larger than the constrained maximum likelihood, the denominator **likely to be** much greater than the numerator, which leads to a **small value of λ**
- Therefore, we may use λ as a test statistic and reject H_0 if $\lambda < k$.

But what threshold k to choose if we want level α test?

The practical applications of the likelihood ratio test are based on the following amazing theorem.

Theorem 3.9.1 (Wilks' Theorem). *Let X_1, X_2, \dots, X_n are i.i.d random observations and $\lambda(X_1, \dots, X_n)$ is the likelihood ratio for the hypothesis Ω_0 against Ω_a . Let r_0 denote the number of free parameters that are specified by $H_0: \vec{\theta} \in \Omega_0$ and let r denote the number of free parameters specified by the statement $\vec{\theta} \in \Omega$. Then, for large n , for all $\theta_0 \in \Omega_0$, the statistic $-2 \log \lambda$ has approximately a χ^2 distribution with $r - r_0$ degree of freedoms.*

The number of free parameters is the total number of parameters minus the number of *equality* constraints, so the difference $r - r_0$ is simply the excess in the number of equality constraints that define Ω_0 over the number of equality constraints that define Ω . (Note that the inequality constraints do not count - they are not important for large sample analysis.)

Therefore, the rejection region for the likelihood ratio test in large samples has a very simple form:

$$RR : \{-2 \log \lambda > \chi^2_\alpha(r - r_0)\},$$

The proof of Wilks' Theorem is not easy and depends on the fact that under assumption that $\vec{\theta} \in \Omega_0$ both constrained and unconstrained maximum likelihood estimators of $\vec{\theta}$ have the approximately Gaussian distributions with the same mean $\vec{\theta}$ and different variances which can be explicitly calculated. Then the asymptotic distribution of the $-\log \lambda$ can be explicitly calculated and it turns out that the exact value of $\theta \in \Omega_0$ disappear from this expression.

So the idea here is the same as in our first method for the design of the hypothesis test. However, the calculation of variances is done behind the scene and the value of the parameter θ magically disappears in the final answer.

Example 3.9.2. Suppose that an engineer wishes to compare the number of complaints per week filed by union stewards for two different shifts at a manufacturing plant. One hundred independent observations on the number of complaints gave means $\bar{x} = 20$ for shift 1 and $\bar{y} = 22$ for shift 2. Assume that the number of complaints per week on the i -th shift has a Poisson

distribution with mean λ_i , for $i = 1, 2$. Use the likelihood ratio method to test $H_0 : \lambda_1 = \lambda_2$ against $H_a : \lambda_1 \neq \lambda_2$ with $\alpha \approx 0.01$.

By taking the product the individual density functions we find the likelihood function:

$$L(\lambda_1, \lambda_2) = \frac{1}{C} e^{-n\lambda_1} (\lambda_1)^{\sum_{i=1}^n x_i} \times e^{-n\lambda_2} (\lambda_2)^{\sum_{i=1}^n y_i},$$

where $C = x_1! \dots x_n! y_1! \dots y_n!$ and $n = 100$.

Here we will be able to do maximizations analytically, although it could also be done numerically.

Log likelihood function:

$$\ell(\lambda_1, \lambda_2) = -\log C + \left(\sum_{i=1}^n x_i \right) \log \lambda_1 - n\lambda_1 + \left(\sum_{i=1}^n y_i \right) \log \lambda_2 - n\lambda_2.$$

If it is assumed that $\lambda_1 = \lambda_2 = \lambda$, then the maximization of the log-likelihood function leads (after a calculation) to the constrained MLE estimator

$$\hat{\lambda}^{ML} = \frac{1}{2}(\bar{x} + \bar{y}) = 21.$$

If we do not assume that $\lambda_1 = \lambda_2$, then the unconstrained maximum likelihood estimator of the vector (λ_1, λ_2) is (after a calculation)

$$\hat{\lambda}_1^{ML} = \bar{x} = 20 \text{ and } \hat{\lambda}_2^{ML} = \bar{y} = 22.$$

Then, log likelihood ratio is the *difference* of the log-likelihoods evaluated at the constrained and unconstrained MLE estimators.

$$-2 \log \text{likelihood ratio} = -2 \left(\ell(\hat{\lambda}^{ML}, \hat{\lambda}^{ML}) - \ell(\hat{\lambda}_1^{ML}, \hat{\lambda}_2^{ML}) \right)$$

Calculation gives:

$$\begin{aligned} -2 \log \text{likelihood ratio} &= -2 \left(\ell(\hat{\lambda}^{ML}, \hat{\lambda}^{ML}) - \ell(\hat{\lambda}_1^{ML}, \hat{\lambda}_2^{ML}) \right) \\ &= -2 \left[-\log k + (n\bar{x} + n\bar{y}) \log \hat{\lambda}^{ML} - 2n\hat{\lambda}^{ML} \right. \\ &\quad \left. + \log k - n\bar{x} \log \hat{\lambda}_1^{ML} + n\hat{\lambda}_1^{ML} - n\bar{y} \log \hat{\lambda}_2^{ML} + n\hat{\lambda}_2^{ML} \right] \end{aligned}$$

Some terms cancel out and we find

$$\begin{aligned} -2 \log \text{likelihood ratio} &= -2 \left[(100 \times 20 + 100 \times 22) \log 21 \right. \\ &\quad \left. - 100 \times 20 \log 20 - 100 \times 22 \log 22 \right] \\ &= 9.5274 \end{aligned}$$

By Wilks theorem we should use the rejection region $RR = \{-2 \log \lambda > \chi_{\alpha=0.01, df=1}^2 = 6.635\}$. Hence we reject $H_0 : \lambda_1 = \lambda_0$ at significance level $\alpha = 0.01$.

In fact, in this example, we can also use the first method based on the estimator of the parameter $\lambda_1 - \lambda_2$. Indeed, since parameter λ is the mean of the Poisson distribution, the problem can be thought as a problem about the equality of means in two samples. The difficulty is that the sample standard deviations are not given. However, we know the distribution of data observations (Poisson).

Note that \bar{x} is an estimator of λ_1 which is approximately normal with distribution $\mathcal{N}(\lambda_1, \lambda_1/n)$. Similarly \bar{y} is approximately normal independent random variable with distribution $\mathcal{N}(\lambda_2, \lambda_2/n)$.

Hence, under the null hypothesis, we have that the test statistic

$$TS = \frac{\bar{y} - \bar{x}}{\sigma_{\bar{y} - \bar{x}}} \sim \mathcal{N}(0, 1).$$

Under the null hypothesis, a good estimator of $\sigma_{\bar{y} - \bar{x}}$ is

$$\sqrt{\frac{\hat{\lambda}}{n} + \frac{\hat{\lambda}}{n}}$$

where $\hat{\lambda} = \frac{1}{2}(\bar{x} + \bar{y}) = 21$. So, we calculate:

$$TS = \frac{22 - 20}{\sqrt{2 \times \frac{21}{100}}} = \frac{2 \times 10}{\sqrt{42}} = 3.086.$$

This is greater than $z_{0.01} = 2.33$, so $H_0 : \lambda_1 = \lambda_2$ should be rejected at level $\alpha = 0.01$.

Quiz 3.9.3. Suppose that I have collected a random sample to test $H_0 : \mu = \mu_0$ v.s. $H_a : \mu > \mu_0$ and I end up rejecting H_0 at level $\alpha = 0.05$ based on my sample. If I decided to change α from 0.05 to 0.01, then based on the same sample that I have in my hand, I would

- A. definitely fail to reject the H_0 at level $\alpha = 0.01$;
- B. definitely reject the H_0 at level $\alpha = 0.01$;
- C. either reject the H_0 at level $\alpha = 0.01$ or fail to reject the H_0 at level $\alpha = 0.01$ depending on the sample;
- D. have to toss a fair coin to decide what to do.

Quiz 3.9.4. Suppose that I am interested in testing $H_0 : \mu = \mu_0$ against $H_a : \mu \neq \mu_0$. I calculate the type II error probability β using the alternative value of parameter μ_a . Then, β will be smaller if I

- A. Decrease the type I error probability α ;
- B. Decrease the sample size n ;
- C. Decrease the distance between μ_a and μ_0 ;
- D. None of the above is correct.

3.9.1 An Additional Example

This section gives an example, in which the likelihood ratio test is designed explicitly without using the approximation provided by Wilks' theorem.

Example 3.9.5. A service station has six gas pumps. When no vehicles are at the station, let p_i denote the probability that the next vehicle will select pump i (where $i = 1, 2, \dots, 6$). We have a sample of size n in which x_i vehicles have chosen pump i . We wish to test $H_0 : p_1 = \dots = p_6 = \frac{1}{6}$ versus the alternative $H_a : p_1 = p_3 = p_5; p_2 = p_4 = p_6 = \theta \neq \frac{1}{6}$.

What is the likelihood function of the parameters p_1, p_2, \dots, p_6 given the sample data if no restriction are imposed on the parameters?

This is a multinomial distribution so

$$L(\vec{p}|\vec{x}) = \binom{n}{x_1, \dots, x_6} \prod_{i=1}^6 p_i^{x_i}$$

What are the likelihood functions under the hypothesis H_0 and H_a , respectively?

Under H_0 , the likelihood is

$$\binom{n}{x_1, \dots, x_6} \left(\frac{1}{6}\right)^n.$$

Under H_a , we calculate that $p_1 = p_2 = p_3 = 1/3 - \theta$, and the likelihood

$$\binom{n}{x_1, \dots, x_6} \left(\frac{1}{3} - \theta\right)^{x_1+x_3+x_5} \theta^{x_2+x_4+x_6}.$$

Suppose that $X = X_2 + X_4 + X_6$ is the number of customers in the sample that select an even numbered pump. What is the maximum likelihood estimator of the parameter θ under the alternative hypothesis H_a ?

Maximization of likelihood under H_a is equivalent to maximization of log-likelihood, that is, of

$$\ell(\theta) = c + (n - X) \log \left(\frac{1}{3} - \theta\right) + X \log \theta,$$

Then,

$$\begin{aligned} \ell'(\theta) &= -(n - X) \frac{1}{\frac{1}{3} - \theta} + X \frac{1}{\theta} = 0, \\ -\theta n + X\theta + X/3 - X\theta &= 0, \\ \hat{\theta}^{MLE} &= \frac{X}{3n}. \end{aligned}$$

Express the likelihood ratio statistic λ in terms of X .

Under H_a , the likelihood

$$\binom{n}{x_1, \dots, x_6} \left(\frac{1}{3} - \theta\right)^{n-X} \theta^X.$$

Substituting the MLE estimate of θ in the definition of λ , we get:

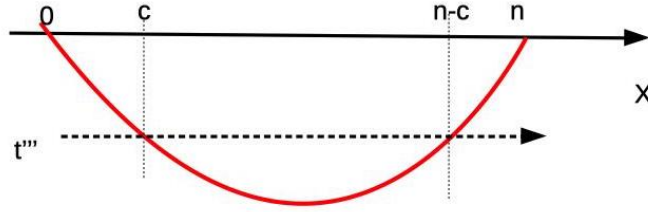
$$\lambda = \frac{(1/6)^n}{\left(\frac{1}{3} - \frac{X}{3n}\right)^{n-X} \left(\frac{X}{3n}\right)^X}$$

The rejection region for likelihood ratio test is $\{\lambda \leq t\}$, where t is a threshold. This is the same as $\{-\log \lambda \geq t'\}$, and we can re-write this region in our case as

$$(n - X) \log \left(\frac{n - X}{3n} \right) + X \log \left(\frac{X}{3n} \right) \geq t'',$$

or

$$(n - X) \log (n - X) + X \log (X) \geq t''',$$



The second derivative of the function on the left is

$$\frac{1}{n - X} + \frac{1}{X} > 0,$$

which means that this function is convex and so:

1. there can be only two solutions to the equality
2. By symmetry, if one of these solutions is c , then the other is $n - c$.
3. The inequality is satisfied only if $X \geq c$ or if $X \leq n - c$.

Let $n = 10$ and $c = 9$. Determine significance level α of the test and its power when $\theta = p_2 = p_4 = p_6 = 1/10$.

Under the H_a , the probability that one of the pumps 2, 4, or 6 is visited equals 3θ . Hence X (the number of visits of these pumps) is distributed as binomial with parameter $3\theta = 0.3$. For H_0 it is the binomial with probability 0.5. Hence,

$$\begin{aligned} \alpha &= \Pr(X \geq 9 | n = 10, p = 0.5) + \Pr(X \leq 1 | n = 10, p = 0.5) \\ &= 1 - \Pr(X \leq 8 | n = 10, p = 0.5) + \Pr(X \leq 1 | n = 10, p = 0.5) \\ &= 1 - 0.989 + 0.011 = 0.022 \end{aligned}$$

and

$$\begin{aligned} \text{power} &= 1 - \beta = \Pr(X \geq 9 | n = 10, p = 0.3) + \Pr(X \leq 1 | n = 10, p = 0.3) \\ &= 1 - \Pr(X \leq 8 | n = 10, p = 0.3) + \Pr(X \leq 1 | n = 10, p = 0.3) \\ &= 1 - 1.000 + 0.149 = 0.149 \end{aligned}$$

Chapter 4

Bayesian Inference

4.1 Estimation

The Bayesian inference is a collection of statistical methods based on a different statistical philosophy. The statistical model is still consist of observations X_1, \dots, X_n which are random with the distribution $f(\vec{X}|\theta)$ that depend on a vector of parameters θ . However, while the classical statistic treats the parameters as fixed and unknown quantities, the Bayesian statistic models researchers' beliefs about the parameters by using probability theory. This adds a second layer of randomness: now the parameters θ of the data-generating distribution $f(\vec{X}|\theta)$ have their own probability distributions which model our beliefs about them.

In fact, the parameters are treated as random variables that have two probability distributions: before and after the data is observed. Their distribution before the data is observed is described by a *prior distribution* with density (or mass) function $p(\theta)$. The *posterior distribution* is the distribution of parameters after the data is observed. It captures our beliefs after they were modified by the observed data. The density (or mass) function of the posterior distribution is the conditional density $p(\theta|x_1, \dots, x_n)$. It can be calculated from the prior distribution and the data by using Bayes' formula:

$$p(\theta|x_1, \dots, x_n) = \frac{f(x_1, \dots, x_n|\theta)p(\theta)}{\int f(x_1, \dots, x_n|\varphi)p(\varphi) d\varphi}.$$

The integral in the denominator is a normalizing constant – it does not depend on θ . Often, it is not written explicitly, and the formula is written as

$$p(\theta|x_1, \dots, x_n) \propto f(x_1, \dots, x_n|\theta)p(\theta).$$

The symbol \propto means “proportional to”.

Example 4.1.1. A biased coin is tossed n times. Let X_i be 1 or 0 as the i -th toss is or is not a head. The probability of a head is θ . Suppose we have no idea how biased the coin is, and we place a uniform prior distribution on θ , to give a so-called “non-informative prior” of

$$p(\theta) = 1, \quad 0 \leq \theta \leq 1.$$

Let t be the number of heads. Then, the posterior distribution of θ is

$$p(\theta|x_1, \dots, x_n) \propto \theta^t(1 - \theta)^{n-t}$$

By inspection we realize that if the appropriate constant on the right-hand side, then we have the density of the Beta distribution with parameters $(t + 1, n - t + 1)$. This is the posterior distribution of θ given x .

After a bit of reflection, we realize that if we start with the prior distribution which is the beta distribution with parameters α_1, α_2 , that is, if $p(\theta) \propto \theta^{\alpha_1-1}(1 - \theta)^{\alpha_2-1}$, then the posterior distribution is

$$p(\theta|x_1, \dots, x_n) \propto \theta^{t+\alpha_1-1}(1 - \theta)^{n-t+\alpha_2-1},$$

that is, it is still the beta distribution but with updated parameters $t+\alpha_1$ and $n-t+\alpha_2$. Note that here α_1 and α_2 are parameters for the prior distribution of the parameter θ ! Sometimes they are called hyper-parameters.

Definition 4.1.2. Let $f(x|\theta)$ be the distribution of data point x given the parameter θ . Let $p(\theta)$ be a prior distribution for the parameter. If the posterior distribution $p(\theta|x_1, \dots, x_n)$ has the same functional form as the prior but with altered parameter values, then the prior $p(\theta)$ is said to be *conjugate* to the distribution $f(x|\theta)$.

The conjugate priors are very convenient in modeling and are often used in practice. Here is another example.

Exercise 4.1.3. Suppose X_1, \dots, X_n are normally distributed $X_i \sim N(\mu, \sigma)$ with unknown parameter μ and known $\sigma = 1$. Let the prior distribution for μ be $N(0, \tau^{-2})$ for known τ^{-2} . Then the posterior distribution for μ is

$$N\left(\frac{\sum_{i=1}^n x_i}{n + \tau^2}, \frac{1}{n + \tau^2}\right).$$

In many cases, in practical applications we need a point estimator and not a posterior distribution of the parameter. Bayesian statistic addresses this concern with the concept of the loss function. A *loss function* $L(\theta, a)$ is a measure of the loss incurred by estimating the value of the parameter to be a when its true value is θ . The estimator $\hat{\theta}$ is chosen to minimize the expected loss $\mathbb{E}L(\theta, \hat{\theta})$, where the expectation is taken over θ with respect to the posterior distribution $p(\theta|\vec{x})$,

$$\hat{\theta} = \arg \min_a \mathbb{E}[L(\theta, a)]$$

Theorem 4.1.4. (a) Suppose that the loss function is quadratic in error: $L(\theta, a) = (\theta - a)^2$. Then the expected loss is minimized by taking $\hat{\theta}$ to be the mean of the posterior distribution:

$$\hat{\theta} = \int \theta p(\theta|x_1, \dots, x_n) d\theta.$$

(b) Suppose the loss function is the absolute value of the error: $L(\theta, a) = |\theta - a|$. Then the expected loss is minimized by taking $\hat{\theta}$ to be the median of the posterior distribution.

Example 4.1.5 (Coin Tosses). Consider the setting of Example 4.1.1. The posterior distribution is Beta distribution with parameters $t+1$ and $n-t+1$. So, by properties of Beta distribution, the *posterior mean* estimator is

$$\hat{\theta} = \frac{t+1}{n+2},$$

and the *posterior median* estimator needs to be calculated numerically. Note that both are different from the standard estimator $\bar{x} = t/n$.

Note that both posterior mean and posterior median estimators depends on our choice of the prior distribution.

The Bayesian analogue of confidence intervals is *credible sets*. A set A is a *credible set* with probability $1 - \alpha$ if the probability that the parameter belongs to the set is α when the probability is calculated with respect to the posterior distribution. That is,

$$\mathbb{P}\theta \in A = \int_A p(\theta|x_1, \dots, x_n) d\theta = 1 - \alpha.$$

Calculation of the credible sets is straightforward when we know the posterior distribution. It is important to understand that confidence sets depend on the choice of prior distribution.

4.2 Hypothesis testing

Bayesian inference is also different in its approach to hypothesis testing from the approach of the classical statistics. In fact, the hypothesis testing plays less significant role in Bayesian inference simply because the idea that a continuous parameter can be for sure equal to a specific value is at odds with main idea of Bayesian inference to model beliefs about parameters with probabilities. Still it is possible to evaluate two competing hypothesis using Bayesian methods.

In the classical case, the setup consists of a pair of hypotheses: the null hypothesis H_0 and the alternative hypotheses H_a . The rejection region is selected as a set of data samples for which we reject of H_0 in favor of H_a . This set is selected in such a way that the probabilities of making the wrong decisions, - the probabilities of type I and type II errors, α and β , - are small.

The deficiency of this method is that it works really well only if both the null and alternative are simple hypothesis, so that we can calculate α and β unambiguously. In this case, the Neyman-Pearson lemma provides us with the most powerful test, that is the test, that has the smallest β for a given α . This test is based on the ratio of likelihood functions, that is on the ratio

of data densities for the parameters θ_0 and θ_a :

$$\frac{L(\theta_0|x_1, \dots, x_n)}{L(\theta_a|x_1, \dots, x_n)} \equiv \frac{f(x_1, \dots, x_n|\theta_0)}{f(x_1, \dots, x_n|\theta_a)}.$$

Sometimes, it is possible to develop the best test (Uniformly Most Powerful Test) even when the alternative hypothesis is composite. However, in many cases, for example, for the two-sided alternative hypothesis there are no UMP tests. In addition, in order to search for UMP tests, we have to assume that the null hypothesis is simple. This is somewhat unsatisfactory since in many practical situations it is difficult to justify a specific value for the null hypothesis.

In practice, statisticians are satisfied with reasonable, although not UMP, tests, which would allow us to do testing even if the null hypothesis is not simple. One of this tests, the likelihood ratio test is based on the test statistic:

$$\lambda(x_1, \dots, x_n) = \frac{\max_{\theta \in \Omega_0} L(\theta|x_1, \dots, x_n)}{\max_{\theta \in \Omega_0 \cup \Omega_a} L(\theta|x_1, \dots, x_n)} \equiv \frac{\max_{\theta \in \Omega_0} f(x_1, \dots, x_n|\theta)}{\max_{\theta \in \Omega_0 \cup \Omega_a} f(x_1, \dots, x_n|\theta)}$$

In other words, we choose the value θ_0 in the null hypothesis set Ω_0 , which gives the largest probability density of the data, and compare this density with the maximum of the data probability density when the parameter is allowed to vary over both the null (Ω_0) and alternative (Ω_a) hypothesis sets. We reject the null if the ratio of these two probabilities is smaller than a threshold.

While this procedure is very reasonable, it does not have a clear probabilistic justification.

In contrast, in the Bayesian inference, the null hypothesis is typically not a single value but a big set of parameters $H_0 : \theta \in \Omega_0$ and the alternative is the complement of this set, $H_a : \theta \in \Omega_a = \Omega_0^c$. For example, we can have the null hypothesis $H_0 : \theta \leq \theta_0$ and the alternative $H_a : \theta > \theta_0$.

The null hypothesis is rejected by the Bayesian test, if the ratio of pos-

terior probabilities of hypotheses:

$$\begin{aligned}\lambda_B(x_1, \dots, x_n) &= \frac{\mathbb{P}[\theta \in \Omega_0]}{\mathbb{P}[\theta \in \Omega_a]} = \frac{\int_{\Omega_0} p(\theta|x_1, \dots, x_n) d\theta}{\int_{\Omega_a} p(\theta|x_1, \dots, x_n) d\theta} \\ &= \frac{\int_{\Omega_0} f(x_1, \dots, x_n|\theta)p(\theta) d\theta}{\int_{\Omega_a} f(x_1, \dots, x_n|\theta)p(\theta) d\theta}\end{aligned}$$

is smaller than a certain threshold t , which measures the degree of our conservatism. For example, if the threshold is set to $1/3$, then we reject the null hypothesis H_0 only if the posterior probability of H_0 is three times smaller than the posterior probability of H_a .

This resembles the likelihood ratio test, except instead of maximizing the data density over the set of parameters Ω_0 and Ω_a , we take the average of the data densities by using the prior probability distribution $p(\theta)$.

It is in principle possible to define α and β of the Bayesian test as the **average** probabilities of making type I and type II errors, where the average is calculated with respect to the prior distribution. However, the definition is more complicated. In addition, the Bayesian analysis is most useful when the amount of the data is not overwhelmingly large compared to our prior beliefs. In this situation, there is no analogue of Wilkes' theorem for the likelihood ratio, and so it is significantly more difficult to develop a test with a given α . For this reason α and β are very rarely used in Bayesian inference.

Here is an example, how a Bayesian test applies in practice.

Example 4.2.1. Let X_1, \dots, X_n be a sample from an exponentially distributed population with density $f(x|\theta) = \theta e^{-\theta x}$. (Note that this is a slightly different parameterization of the exponential distribution. The mean of the distribution is $\mu = 1/\theta$. Suppose the prior distribution is Gamma distribution with parameters α and β . Test the hypothesis $H_0 : \theta \leq \theta_0$ versus $H_a : \theta > \theta_0$.

The density of the data sample is

$$f(x_1, \dots, x_n|\theta) = \theta^n e^{-\theta \sum_{i=1}^n x_i}.$$

The prior is

$$p(\theta) \propto \theta^{\alpha-1} e^{-\theta/\beta}.$$

The posterior distribution is

$$p(\theta|x_1, \dots, x_n) \propto \theta^{n+\alpha-1} e^{-\theta(\sum_{i=1}^n x_i + 1/\beta)}$$

So the posterior distribution is the Gamma distribution with parameters

$$\begin{aligned}\alpha' &= n + \alpha, \\ \beta' &= \frac{1}{\sum_{i=1}^n x_i + 1/\beta} = \frac{\beta}{\sum_{i=1}^n x_i + 1}\end{aligned}$$

In particular we showed that the Gamma distribution is the conjugate prior for the exponential distribution. We reject the null hypothesis only if

$$\mathbb{P}[\theta \leq \theta_0] < \frac{t}{1+t}.$$

In R, this can be solved by checking if

$$\text{pgamma}(\theta_0, \alpha', 1/\beta') < \frac{t}{1+t}.$$

For example, let $n = 10$, $\sum x_i = 1.26$, $\alpha = 3$ and $\beta = 5$. (These are perhaps obtained by reviewing prior studies about θ .) We want to test the null hypothesis that $H_0 : \mu > .12$ against $H_a : \mu \leq .12$ using $t = 1$ (This not a very conservative test. We reject null hypothesis if its probability is smaller than the probability of the alternative.) In terms of the parameter θ , the hypotheses are

$$H_0 : \theta < 1/ (.12) = 8.333 \text{ against } H_a : \theta \geq 8.333.$$

We calculate

$$\begin{aligned}\alpha' &= n + \alpha = 10 + 3 = 13 \\ 1/\beta' &= \sum_{i=1}^n x_i + 1/\beta = 1.26 + 1/5 = 1.46,\end{aligned}$$

then

$$\text{pgamma}(\theta_0, \alpha', 1/\beta') = \text{pgamma}(8.333, 13, 1.46) = 0.4430332$$

Since this is smaller than $1/(1+t) = 1/2$, we can reject the null hypothesis.

One observation about the Bayesian hypothesis testing is that the results of the tests depend on the choice of the prior distribution and this choice should be careful and well-justified. The second observation is that in practice it is sometimes difficult to calculate the probabilities under the posterior distribution. This calculation may involve difficult integrations. In this respect, the classical approach is often computationally easier.