

Skriptum zur Vorlesung “Mathematical Statistics”

Karl Grill¹ (karl.grill@tuwien.ac.at)

Version 0.2020.2

Inhaltsverzeichnis

1	Introduction	3
1.1	Basic Notions	3
1.2	Sample Statistics	4
1.3	Quantiles and Order Statistics	6
1.4	Problems	7
2	Basic Theory of Estimation	9
2.1	Point Estimation	9
2.2	Confidence Intervals	12
2.2.1	The normal distribution	13
2.2.2	Proportions	14
2.3	Problems	15
3	Testing	17
3.1	Basic Notions	17
3.2	The Neyman-Pearson and likelihood ratio tests	19
3.3	Special tests for the normal distribution	20
3.3.1	Tests for μ	20
3.3.2	Tests for σ^2	20
3.3.3	Two-sample tests	21
3.4	Uniformly Optimal Tests	22
3.4.1	Monotone Likelihood Ratios	22
3.4.2	Unbiased Tests.	23
3.5	Tests and confidence intervals	25
3.6	Problems	26
4	Analysis of Variance (ANOVA)	28
4.1	The Fisher-Cochran Theorem	28
4.2	One-way analysis of variance	32
4.3	Two-way analysis of variance	33
4.4	Problems.	34
5	Linear Regression.	36
5.1	Simple linear regression.	36
5.2	Multiple linear regression.	37
5.3	Other functional relations.	38
5.4	Problems	38
6	The χ^2-Family of Tests.	40
6.1	The multinomial distribution and the χ^2 -statistic	40
6.2	The χ^2 goodness-of-fit test.	42
6.3	The χ^2 test of independence.	42
6.4	The χ^2 test for homogeneity.	43
6.5	Problems	43

7	The Kolmogorov-Smirnov Test	45
7.1	The one-sample test	45
7.2	The Two-Sample Test	46
7.3	Problems	47
A	Tables	48
B	Solutions of Problems	57
B.1	Problems from Chapter 1	57
B.2	Problems from Chapter 2	58
B.3	Problems from Chapter 3	61
B.4	Problems from Chapter 4	64
B.5	Problems from Chapter 5	66
B.6	Problems from Chapter 6	68
B.7	Problems from Chapter 7	71
C	German translations of technical terms	73

Kapitel 1

Introduction

Smoking is the principal cause of statistics

Irish Graffiti

1.1 Basic Notions

The field of statistics is usually divided into two parts: first, *descriptive statistics* which deals with the gathering and representation of large sets of data, and *mathematical* or *inductive statistics* whose object is the determination of properties of a large population based on what is called a random sample. The latter is the object of the present lecture.

Usually, statistical questions arise in conjunction with a (presumably large) *population* of individuals (these need not necessarily be human or even animate — the collection of all tv sets produced in a factory within a certain period of time makes a perfectly sound statistical population). These individuals have one or more properties which are going to be studied - e.g., length, height, weight, sex, color of eyes or hair, . . . These data are usually either numerical values (measurements) or some sort of classification (like sex, colours). We can represent the latter by, e.g., integer numbers starting from one, so there is no obvious loss of generality if we assume that our data are real numbers (actually, there is the possibility to allow more general spaces than the real numbers, e.g., some general metric space, but we won't go into that).

The reasoning behind all of mathematical statistics is the observation that if we select an individual from our population randomly, then the value of a certain datum associated with that individual is a random variable whose distribution is the (relative) frequency distribution of that datum within the population. Furthermore, if we draw a sample of n individuals from the population *with replacement* then the values of the datum will be independent random variables with the same distribution. If sampling is done without replacement, there is some dependency between those random variables, some of which will be explored in the problems below. If the size of the population is large with respect to sample size, however, the effect of non-replacement is small. Thus, from now on, we will assume that sampling is done with replacement. Furthermore, although real populations will always generate discrete distributions, we will allow any distribution for our sample random variables (which should make sense at least as a kind of approximation). This gives us

Definition 1.1

A random sample of size n from distribution P is a sequence X_1, \dots, X_n of independent, identically distributed random variables with common distribution P .

The random variables in this definition may be d -dimensional (e.g., if we study height and weight and any possible relation between them).

The job of statistics now is to make some statement about the distribution P based on the sample X_1, \dots, X_n . The following theorem shows that we can expect this to work (for this reason, it is also called the fundamental theorem of mathematical statistics):

Theorem 1.1

(Glivenko-Cantelli) Let X_1, \dots be a sequence of i.i.d.r.v. with common distribution function F and define the *empirical distribution function* as

$$F_n(x) = \frac{\#\{i \leq n : X_i \leq x\}}{n}.$$

Then, with probability one as $n \rightarrow \infty$,

$$\sup_{x \in \mathbf{R}} |F_n(x) - F(x)| \rightarrow 0.$$

Proof. Fix $\epsilon > 0$. Choose $M > 1/\epsilon$, $x_0 = -\infty$, $x_M = \infty$, and for $k = 1, \dots, M-1$ choose x_k in such a way that $F(x_{k-1}) \leq \frac{k}{M} \leq F(x_k)$ and $x_k \geq x_{k-1}$. Then, there is a (random) n_0 such that for $n > n_0$ and any $k = 1, \dots, M-1$

$$|F_n(x_k) - F(x_k)| < \epsilon$$

and

$$|F_n(x_k - 0) - F(x_k - 0)| < \epsilon.$$

by the strong law of large numbers. For any $x \in \mathbf{R}$ we can find a k with $1 \leq k \leq M$ and $x_{k-1} \leq x < x_k$. With this k , we have, for $n > n_0$

$$F_n(x) \geq F_n(x_{k-1}) \geq F(x_{k-1}) - \epsilon \geq F(x_k - 0) - \epsilon - \frac{1}{M} \geq$$

$$F(x_k - 0) - 2\epsilon \geq F(x) - 2\epsilon.$$

Similarly

$$F_n(x) \leq F_n(x_k) \leq F(x_k) + \epsilon \leq F(x_{k-1}) + \epsilon + \frac{1}{M} \leq$$

$$F(x_{k-1}) + 2\epsilon \leq F(x) + 2\epsilon.$$

The last two inequalities, together with the arbitrariness of ϵ prove the theorem.

In many cases we have some information about the shape of the possible distributions P . For example, if we are only interested in whether a voter is in favor of a given party or not, we may put

$$X_i = \begin{cases} 1 & \text{if } i\text{-th person polled is in favor of party A} \\ 0 & \text{otherwise.} \end{cases}$$

In this case, we have

$$\mathbb{P}(X_i = 1) = 1 - \mathbb{P}(X_i = 0) = p,$$

where p is the percentage of voters in favor of party A. In this case, the distribution is determined by the parameter p .

More generally, if Θ is a subset of \mathbf{R}^d and for any $\theta \in \Theta$, P_θ is a probability distribution, then we call $\{P_\theta, \theta \in \Theta\}$ a *parametric family* of probability distributions. In general, we will assume that such a parametric family is dominated by a σ -finite measure μ , i.e., for each $\theta \in \Theta$ the probability measure P_θ is absolutely continuous with respect to μ . In that case, the distribution P_θ is determined by its Radon-Nikodym derivative $f_\theta = \frac{dP_\theta}{d\mu}$. In most cases of practical interest, the dominating measure is either Lebesgue measure or some counting measure.

If the assumption that the distribution of the sample comes from a given parametric family is made, we speak of *parametric statistics*, otherwise of *nonparametric statistics*.

1.2 Sample Statistics

There are lies, damned lies, and statistics.

W. Churchill

A *statistic* is merely something that can be calculated from a sample. A little more mathematically formulated:

Definition 1.2

If f is a function from \mathbb{R}^n to \mathbb{R}^d and (X_1, \dots, X_n) is a random sample, then we call

$$T = f(X_1, \dots, X_n)$$

a statistic.

Important examples of statistics are the sample mean

$$\bar{X}_n = \frac{\sum_{i=1}^n X_i}{n}$$

and the sample variance

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

An important notion is that of a *sufficient statistic*. Loosely speaking, a sufficient statistic is one that already contains all the information about the parameter that can be obtained from the sample.

Definition 1.3

A statistic T is called sufficient for the parameter θ if the conditional distribution of (X_1, \dots, X_n) given T does not depend on θ .

Remark. The statistic T may of course be d -dimensional.

An important criterion is the following:

Theorem 1.2

Let $(P_\theta, \theta \in \Theta)$ be a parametric family of distributions dominated by a measure μ , and $f_\theta = \frac{dP_\theta}{d\mu}$. The statistic $T = T(X_1, \dots, X_n)$ is sufficient for θ if the so-called likelihood function

$$L(x_1, \dots, x_n, \theta) = f_\theta(x_1) \dots f_\theta(x_n)$$

admits a decomposition

$$L(x_1, \dots, x_n, \theta) = g(T(x_1, \dots, x_n), \theta)h(x_1, \dots, x_n),$$

where $h(\cdot)$ does not depend on θ .

Proof. We prove the theorem for the case where μ is a counting measure, i.e., for discrete P_θ . The proof of the general case is not much harder, but one has to be careful about some measure-theoretic technicalities, so we rather stick with the simple case. In this case,

$$\begin{aligned} \mathbb{P}_\theta(T = t) &= \sum_{\{(x_1, \dots, x_n): T(x_1, \dots, x_n) = t\}} L(x_1, \dots, x_n, \theta) = \\ &= \sum_{\{(x_1, \dots, x_n): T(x_1, \dots, x_n) = t\}} g(t, \theta)h(x_1, \dots, x_n) = g(t, \theta)H(x_1, \dots, x_n), \end{aligned}$$

where $H(\cdot)$ does not depend on θ . Thus

$$\mathbb{P}_\theta(X_1 = x_1, \dots, X_n = x_n | T = t) = \begin{cases} \frac{h(x_1, \dots, x_n)}{H(x_1, \dots, x_n)} & \text{if } T(x_1, \dots, x_n) = t, \\ 0 & \text{otherwise,} \end{cases}$$

which does not depend on θ .

One example for a sufficient set of statistics that is always good is the order statistics. These are defined as

$$X_{n:i} = \inf\{x : \#\{j \leq n : X_j \leq x\} \geq i\}.$$

In other words, $(X_{n:1}, \dots, X_{n:n})$ are X_1, \dots, X_n in ascending order. It is obvious from the previous criterion that these are sufficient. We will explore their distribution in the next section.

1.3 Quantiles and Order Statistics

They called him Gautama Buddha
 Long time ago.
 He turned the world to order
 Don't you know?

Yussuf Islam, then (1974) still Cat Stevens, "Jesus"

The order statistics introduced in the last section are closely related to the concept of quantiles:

Definition 1.4

Let X be a random variable with distribution F , $0 \leq p \leq 1$. The number $x = x_p$ is called a p -quantile of X resp. F if

$$F(x_p - 0) \leq p \leq F(x_p).$$

Similarly, x_p is called an empirical p -quantile of the sample X_1, \dots, X_n if it is a p -quantile of its empirical distribution function F_n .

Obviously, if np is not an integer, then the unique empirical p -quantile is $X_{n:[np]}$, whereas if np is an integer, then any value between $X_{n:np}$ and $X_{n:np+1}$ is an empirical p -quantile.

The most important quantiles are the median $x_{0.5}$, which serves as a parameter of location (similar to the expectation; for a symmetric distribution, the expectation and the median coincide, or, more precisely, the expectation is a median) and the quartiles $x_{0.25}$ and $x_{0.75}$ which provide the interquartile range $IQR = x_{0.75} - x_{0.25}$, a parameter of dispersion (like the more common variance resp. standard deviation).

We denote the distribution function of $X_{n:k}$ by $F_{n:k}$. It may be easiest to start with the distributions of $X_{n:1}$ and $X_{n:n}$, which are just the minimum and maximum of the sample. Elementary considerations yield

$$F_{n:1}(x) = 1 - (1 - F(x))^n, F_{n:n}(x) = F(x)^n.$$

From there, it is only a small step to the general case:

$$F_{n:k}(x) = \mathbb{P}(X_{n:k} \leq x) = \mathbb{P}(|\{i \leq n : X_i \leq x\}| \geq k) = \sum_{i=k}^n \binom{n}{i} F(x)^i (1 - F(x))^{n-i}.$$

The joint distribution of two order statistics doesn't raise too many problems either: namely, for $k \leq l$, we can evaluate

$$\mathbb{P}(X_{n:k} > x, X_{n:l} \leq y) = \mathbb{P}(|\{i : X_i \leq x\}| < k, |\{i : X_i > y\}| \leq n - l)$$

as a sum of multinomial probabilities. In the particular case $k = 1, l = n$ This yields

$$F_{n:1,n:n}(x, y) \begin{cases} F(y)^n - (F(y) - F(x))^n & \text{if } x \leq y \\ F(y)^n & \text{otherwise.} \end{cases}$$

In the case that F is absolutely continuous with density f , the joint distribution of the order statistics $(X_{n:i_1}, \dots, X_{n:i_k}), i_1 < \dots < i_k$ can be given by its density:

Theorem 1.3

$$f_{i_1, \dots, i_k}(x_1, \dots, x_k) = \frac{n!}{\prod_{j=0}^k (i_{j+1} - i_j - 1)!} \prod_{j=1}^k f(x_j) \prod_{j=0}^k (F(x_{j+1}) - F(x_j))^{i_{j+1} - i_j - 1} [x_1 < \dots < x_n].$$

In particular, the joint density of all order statistics is

$$f(x_1, \dots, x_n) = \begin{cases} n! f(x_1) \dots f(x_n) & \text{if } x_1 < \dots < x_n, \\ 0 & \text{otherwise.} \end{cases}$$

As F is continuous, X_1, \dots, X_n are all different with probability 1. For $x_1 < \dots < x_n$, we can find an $\epsilon > 0$ such that $x_{i+1} - x_i > 2\epsilon$ for all $i = 1, \dots, n-1$. For $h_i < \epsilon$ the intervals $[x_i - h_i, x_i + h_i]$ are disjoint, so

$$\mathbb{P}(X_{n:i} \in [x_i - h_i, x_i + h_i], i = 1, \dots, n) = n! \mathbb{P}(X_i \in [x_i - h_i, x_i + h_i], i = 1, \dots, n),$$

which proves the last statement of the theorem, and remainder is obtained by calculating the marginal density of this distribution.

For our further investigations, let us consider a sequence (U_1, \dots, U_n) of independent random variables with a uniform distribution on $[0, 1]$, and let $(U_{n:1}, \dots, U_{n:n})$ denote its order statistics. It is a well-known fact that $X_i = F^{-1}(U_i)$ has distribution function F , and by the monotonicity of F^{-1} we have $X_{n:k} = F^{-1}(U_{n:k})$.

1.4 Problems

We can work it out

The Beatles

1. In an urn there are N slips of paper with numbers x_1, \dots, x_N written on them. Let

$$\mu = \frac{x_1 + \dots + x_N}{N},$$

$$\sigma^2 = \frac{(x_1 - \mu)^2 + \dots + (x_N - \mu)^2}{N}.$$

Let (X_1, \dots, X_n) be the results of drawing n items with replacement.

Calculate the expectation and variance of \bar{X}_n .

2. In the previous example, calculate the expectation of S_n^2 .
3. Solve problem 1 for sampling without replacement.
4. Solve problem 2 for sampling without replacement.
5. Show that (\bar{X}_n, S_n^2) is a pair of sufficient statistics for the normal distributions with densities

$$f_{\mu, \sigma^2}(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}.$$

6. Find a sufficient statistic for the family of uniform distributions with densities

$$f_\theta(x) = \begin{cases} 1/\theta & \text{if } 0 \leq x \leq \theta, \\ 0 & \text{otherwise} \end{cases}$$

7. Find a pair of sufficient statistics for the family of uniform distributions with densities

$$f_\theta(x) = \begin{cases} 1/(\theta_2 - \theta_1) & \text{if } \theta_1 \leq x \leq \theta_2, \\ 0 & \text{otherwise} \end{cases}$$

8. Show that \bar{X}_n is sufficient for the family of exponential distributions with densities

$$f_{\theta}(x) = \begin{cases} \frac{1}{\theta} e^{-\theta x} & \text{if } x > 0, \\ 0 & \text{otherwise} \end{cases}$$

9. Find a pair of sufficient statistics for the family of gamma distributions with densities

$$f_{\alpha,\beta}(x) = \begin{cases} \frac{\beta^{\alpha} x^{\alpha-1}}{\Gamma(\alpha)} e^{-\beta x} & \text{if } x > 0, \\ 0 & \text{otherwise.} \end{cases}$$

Kapitel 2

Basic Theory of Estimation

How many raods must a man walk down
before you call him a man?

B. Dylan

2.1 Point Estimation

Throughout this chapter, we shall assume that (X_1, \dots, X_n) is a random sample whose distribution is from some parametric family $(P_\theta, \theta \in \Theta)$. We are trying to find an approximate value for θ from the sample. We define

Definition 2.1

An *estimator* is a sequence $(\hat{\theta}_n, n \in \mathbb{N})$ of statistics, where $\hat{\theta}_n$ is a function of (X_1, \dots, X_n) .

Remark. Obviously, this definition is there so as to allow us to let the sample size tend to infinity. Apart from that, an estimator can be almost everything (identically 0, for example), so we need a few additional definitions to get a little more sense out of it.

Definition 2.2

An estimator $\hat{\theta}_n$ is called

- (weakly) consistent if $\hat{\theta}_n \rightarrow \theta$ in probability with respect to P_θ (i.e., $\mathbb{P}_\theta(|\hat{\theta}_n - \theta| > \epsilon) \rightarrow 0$ for any $\epsilon > 0$),
- strongly consistent if $\hat{\theta}_n \rightarrow \theta$ with probability one,
- unbiased, if $\mathbb{E}_\theta(\hat{\theta}_n) = \theta$,
- efficient, if it is unbiased and has minimal variance among all unbiased estimators,
- asymptotically unbiased, if $\mathbb{E}_\theta(\hat{\theta}_n) \rightarrow \theta$.

The most important notion, of course, is consistency, because it ensures that the estimated value will be close to the actual value of the parameter if we choose the sample size n big enough.

Let us consider a few examples:

1. The sample mean \bar{X}_n is an unbiased estimator of the expectation μ of the underlying distribution. The laws of large numbers imply that it is also consistent.
2. The sample variance S_n^2 is an unbiased consistent estimator of the variance σ^2 - that is why we use $n - 1$ in the denominator.

3. If we replace the denominator $n - 1$ in the definition of S_n^2 by n , the resulting estimator will no longer be unbiased; it will, however, still be consistent and asymptotically unbiased.

Next, we present two methods for calculating estimators. The first one is the *method of moments*. This originates from the fact that by the law of large numbers, the sample mean converges to the expectation of the underlying distribution with probability one. Thus, if the expectation is a function of the parameter that has a continuous inverse f^{-1} , then $f^{-1}(\bar{X}_n)$ will be a consistent estimator of the parameter. If there is more than one parameter (as for the normal distribution), we use the higher moments $\mathbb{E}(X^k)$ which may be approximately calculated from the sample as

$$\frac{1}{n} \sum_{i=1}^n X_i^k.$$

For the normal distribution, this yields the equations

$$\hat{\mu} = \bar{X}_n$$

and

$$\hat{\sigma}^2 + \hat{\mu}^2 = \frac{1}{n} \sum_{i=1}^n X_i^2,$$

so, finally

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

The idea behind the second method is to choose the parameter in such a way that the particular sample has maximal probability of being observed. In other words, we choose $\hat{\theta}$ in such a way that the Likelihood function L is maximal. This is, rather surprisingly, called the maximum likelihood method, and the resulting estimator is called the maximum likelihood (ML) estimator. In most cases, the maximum can be found by using the standard method of zeroing the derivative (or the partial derivatives if there is more than one parameter) of the likelihood function. Usually, it is convenient to first take the logarithm of the likelihood function.

Under rather weak conditions, the ML estimator is consistent; we will not prove this strictly, but rather give a heuristic reasoning. Namely, suppose that θ_0 is the “real” parameter. The logarithm of the likelihood function, evaluated at some θ , is

$$l(\theta) = \sum_{i=1}^n \log(f_\theta(X_i)).$$

By the law of large numbers, we have

$$\frac{1}{n} l(\theta) \rightarrow \mathbb{E}_{\theta_0}(\log f_\theta(X_i)) = \int \log(f_\theta(x)) f_{\theta_0}(x) d\mu(x).$$

Using the inequality $\log x \leq x - 1$, we get

$$\begin{aligned} \int \log(f_\theta(x)) f_{\theta_0}(x) d\mu(x) - \int \log(f_{\theta_0}(x)) f_{\theta_0}(x) d\mu(x) &= \\ \int \log\left(\frac{f_\theta(x)}{f_{\theta_0}(x)}\right) f_{\theta_0}(x) d\mu(x) &\leq \\ \int \left(\frac{f_\theta(x)}{f_{\theta_0}(x)} - 1\right) f_{\theta_0}(x) d\mu(x) &= \int f_\theta(x) d\mu(x) - \int f_{\theta_0}(x) d\mu(x) = 1 - 1 = 0. \end{aligned}$$

Thus, asymptotically, the likelihood function should attain its maximum at θ_0 .

We now have a look at efficient estimators. First, we derive a lower bound for the variance of an unbiased estimator:

Theorem 2.1

(Cramér-Rao) Let X be a random variable with distribution P_θ , where $\theta \in \Theta$ is a real parameter, and Θ is supposed to be an interval. Furthermore, assume that the density $f_\theta(x)$ is twice differentiable with respect to θ and that both $|f'|$ and $|f''|$ (where $'$ denotes differentiation with respect to θ) are bounded above, uniformly in θ , by some integrable function. Furthermore, let $\hat{\theta}$ be an unbiased estimator of θ . Then

$$\mathbf{Var}_\theta(\hat{\theta}) \geq \frac{1}{I(\theta)},$$

where

$$I(\theta) = \mathbb{E}\left(\left(\frac{\partial \log f_\theta(X)}{\partial \theta}\right)^2\right) = -\mathbb{E}\left(\frac{\partial^2 \log f_\theta(X)}{\partial \theta^2}\right)$$

is the so-called Fisher information.

Remark. If we have a random sample of size n , we may consider the whole sample as one n -dimensional random variable; it is readily obtained that the Fisher information of this n -dimensional variable is just n times the Fisher information $I(\theta)$ of one coordinate, thus we have for an unbiased estimator $\hat{\theta}_n$ based on the sample of size n :

$$\mathbf{Var}_\theta(\hat{\theta}_n) \geq \frac{1}{nI(\theta)}.$$

This implies that, in cases where the Cramér-Rao theorem is applicable, we have a lower bound of order $1/n$ for the variance of an unbiased estimator.

Proof. By assumption, we have, for all θ ,

$$\int f_\theta(x) d\mu = 1$$

and

$$\int \hat{\theta} f_\theta(x) d\mu = \theta.$$

We may take derivatives on both sides of these equations and obtain (we let $'$ denote partial differentiation with respect to θ)

$$\int f'_\theta(x) d\mu = 0$$

and

$$\int \hat{\theta} f'_\theta(x) d\mu = 1.$$

Subtracting θ times the first equation from the second yields

$$\int (\hat{\theta} - \theta) \frac{f'_\theta(x)}{f_\theta(x)} f_\theta(x) d\mu = 1.$$

By the Cauchy-Schwarz inequality, we get

$$1 \leq \int (\hat{\theta} - \theta)^2 f_\theta(x) d\mu \int \left(\frac{f'_\theta(x)}{f_\theta(x)}\right)^2 f_\theta(x) d\mu.$$

In the last equation, the first factor is the variance of $\hat{\theta}$, and the second one is the Fisher information.

What remains to be proven is the equality of the first and second expression for the Fisher information. This is left as an exercise to the reader.

There is, of course, the question if the Cramér-Rao bound can be attained. This is not the case in general; we give a characterization of the distributions that attain the bound in the problem section. Under some additional conditions, however, it can be shown that the ML estimator is asymptotically unbiased and that the ratio of its variance and the Cramér-Rao bound tends to 1, i.e., it is “almost” efficient.

A final theorem tells us that, in case there is a sufficient statistic, we can limit our search for “good” estimators to those that are a function of the sufficient statistic.

Theorem 2.2

If $\hat{\theta}$ is an unbiased estimator and T is a sufficient statistic, then

$$\tilde{\theta} = \mathbb{E}_\theta(\hat{\theta}|T)$$

is also an unbiased estimator and has a variance that is not greater than that of $\hat{\theta}$.

Proof. First observe that $\tilde{\theta}$ does not depend on θ .

Then,

$$\mathbb{E}_\theta(\tilde{\theta}) = \mathbb{E}_\theta(\mathbb{E}_\theta(\hat{\theta}|T)) = \mathbb{E}_\theta(\hat{\theta}) = \theta,$$

so $\tilde{\theta}$ is unbiased. Furthermore, applying the Cauchy-Schwarz inequality to the conditional expectation, we have

$$\mathbb{E}_\theta(\tilde{\theta}^2) = \mathbb{E}_\theta((\mathbb{E}_\theta(\hat{\theta}|T))^2) \leq \mathbb{E}_\theta(\mathbb{E}_\theta(\hat{\theta}^2|T)) = \mathbb{E}_\theta(\hat{\theta}^2).$$

This, of course, implies that the variance of $\tilde{\theta}$ is at most that of $\hat{\theta}$.

2.2 Confidence Intervals

In the previous section, we studied the problem of finding reasonable estimates for a parameter. We defined some desirable properties of such estimators and gave some general principles for their construction. One question that is left open is that of the accuracy of our estimates — sometimes we'd like to know how close to the “real” parameter our estimator may be. We are now going to tackle this problem.

We will illustrate some of the ideas of this chapter on the following model:

$$\mathbb{P}(X_i = 1) = \theta, \mathbb{P}(X_i = 0) = 1 - \theta, \quad (0 < \theta < 1).$$

First, it is obvious that, in general, we can't give any absolute bounds on the parameter (based on a sample); in fact, in our simple example, whatever the value of θ may be, any sequence of zeros and ones has a positive probability to be the actual outcome of our sample. Therefore, we have to settle for a little bit less, and this we do by admitting a small probability of missing the right value of θ :

Definition 2.3

A *confidence interval* with *coverage probability* γ is an interval $[a, b]$, where a and b are sample statistics with

$$\mathbb{P}(a \leq \theta \leq b) \geq \gamma.$$

Remarks.

1. Popular choices for γ are 0.95 (most usual) and 0.99.
2. If we have equality in the definition of the confidence interval, then we call it an *exact* confidence interval. In general, we can't expect exact confidence intervals to exist. In our example and for a sample size of one, for instance, we only have four numbers a_0 , b_0 , a_1 , and b_1 (i.e., one interval each for the cases $X_1 = 0$ and $X_1 = 1$) to work with, and it is obvious that we can't get a constant probability γ for θ to be in one of those intervals (apart from trivial cases, of course). One way to get exact confidence intervals in cases like this, too, would be to resort to *randomization*. This means that we carry out an additional random experiment (like rolling a die) and use its outcome to determine the limits of our confidence interval. We won't go into this now, but we will pick this idea up again in the chapter on testing.
3. If there is more than one parameter, we may of course study each one independently and give single confidence intervals for each one of them. In some cases however, this is not quite

adequate, and one is more interested in studying *confidence regions* that have a probability γ of containing the point whose coordinates are the individual parameters. In that case, there is the question of the shape of the confidence region. The most obvious choice is, of course, a cartesian product of one-dimensional intervals, but sometimes other choices, e.g., balls or ellipsoids, are more natural (usually because the associated distributions can be more easily calculated).

4. There are few general rules for the construction of confidence intervals, so we are going to study some cases of practical importance. Of course, if there is a sufficient statistic, we may use this as a starting point; failing that, the ML estimator is sometimes a good starting point (under sufficient regularity, it is not only asymptotically unbiased and efficient, but even has an asymptotic normal distribution with expectation θ and variance $1/nI(\theta)$).

2.2.1 The normal distribution

As the normal distribution has two parameters, there are two types of confidence intervals: one for the mean μ , and one for the variance σ^2 . Starting with the mean, let us first assume that we know σ^2 . In that case, we try to find the limits of our confidence interval as

$$[\bar{X}_n - c, \bar{X}_n + c]$$

with c determined by the equation

$$\mathbb{P}(\bar{X}_n - c \leq \mu \leq \bar{X}_n + c) = \gamma.$$

The latter leads to

$$\begin{aligned} \gamma &= \mathbb{P}(\mu - c \leq \bar{X}_n \leq \mu + c) = \mathbb{P}\left(-\frac{c}{\sqrt{\sigma^2/n}} \leq \frac{\bar{X}_n - \mu}{\sqrt{\sigma^2/n}} \leq \frac{c}{\sqrt{\sigma^2/n}}\right) = \\ &= 2\Phi\left(\frac{c}{\sqrt{\sigma^2/n}}\right) - 1. \end{aligned}$$

So, finally

$$c = \sqrt{\frac{\sigma^2}{n}} \Phi^{-1}\left(\frac{1+\gamma}{2}\right),$$

and our confidence interval becomes

$$\left[\bar{X}_n - z_{\frac{1+\gamma}{2}} \sqrt{\frac{\sigma^2}{n}}, \bar{X}_n + z_{\frac{1+\gamma}{2}} \sqrt{\frac{\sigma^2}{n}} \right],$$

where we have replaced the notation $\Phi^{-1}(\alpha)$ by the more usual z_α .

Now, let us turn to the case where σ^2 is unknown. In that case, we cannot use the confidence interval from above because it still contains the unknown variance. Of course, we have to replace it by the sample variance, and there are actually two ways to do this:

1. We may just replace σ^2 by S_n^2 and treat the resulting confidence interval as an approximate one (suitable if n is large).
2. Or, we may turn the approximate confidence interval into an exact one by making use of the following

Theorem 2.3

Let X_1, X_2, \dots be i.i.d. normal random variables with mean μ and parameter σ^2 . Then

1. \bar{X}_n has a normal distribution with mean μ and variance σ^2/n .

2. $\frac{n-1}{\sigma^2} S_n^2$ has a χ^2 -distribution with $n - 1$ degrees of freedom, which has density

$$\frac{x^{\frac{n-3}{2}}}{\Gamma(\frac{n-1}{2}) 2^{\frac{n-1}{2}}} e^{-\frac{x}{2}} \quad (x > 0).$$

3. \bar{X}_n and S_n^2 are independent.

4. $\frac{\bar{X}_n - \mu}{\sqrt{S_n^2/n}}$ has a t - (or Student-) distribution with $n - 1$ degrees of freedom, which has density

$$\frac{1}{\sqrt{\pi(n-1)}} \frac{\Gamma(\frac{n}{2})}{\Gamma(\frac{n-1}{2})} \left(1 + \frac{x^2}{n-1}\right)^{-\frac{n}{2}}.$$

We omit the proof.

Using this theorem, we get the confidence interval for μ :

$$[\bar{X}_n - t_{n-1, \frac{1+\gamma}{2}} \sqrt{\frac{S_n^2}{n}}, \bar{X}_n + t_{n-1, \frac{1+\gamma}{2}} \sqrt{\frac{S_n^2}{n}}],$$

where $t_{n,\alpha}$ denotes the α -quantile of the t -distribution.

Using the same ideas as above, we calculate the following confidence intervals for σ^2 :

1. If μ is known:

$$\left[\frac{\sum_{i=1}^n (X_i - \mu)^2}{\chi_{n, \frac{1+\gamma}{2}}^2}, \frac{\sum_{i=1}^n (X_i - \mu)^2}{\chi_{n, \frac{1-\gamma}{2}}^2} \right],$$

2. If μ is unknown:

$$\left[\frac{(n-1)S_n^2}{\chi_{n-1, \frac{1+\gamma}{2}}^2}, \frac{(n-1)S_n^2}{\chi_{n-1, \frac{1-\gamma}{2}}^2} \right].$$

2.2.2 Proportions

A question that arises very often in practice is that for the proportion of individuals in the population that have a certain property (e.g., voters of a given party, defective items in some manufacturing process). We already know that this corresponds to finding the parameter θ in our standard model

$$\mathbb{P}(X = 1) = \theta, \mathbb{P}(X = 0) = 1 - \theta, \quad (0 < \theta < 1).$$

We are going to calculate an approximate confidence interval for θ . Our starting point for this is the central limit theorem which tells us that \bar{X}_n has an approximate normal distribution with mean θ and variance $\frac{\theta(1-\theta)}{n}$. This would give us the confidence interval

$$[\bar{X}_n - z_{\frac{1+\gamma}{2}} \sqrt{\frac{\theta(1-\theta)}{n}}, \bar{X}_n + z_{\frac{1+\gamma}{2}} \sqrt{\frac{\theta(1-\theta)}{n}}],$$

but, unfortunately, we don't know the value of θ . There are two ways out of this dilemma: the simpler one would be to simply replace θ by its estimator \bar{X}_n ; if we want to be a little more exact, we can solve the equation

$$\bar{X}_n \pm z_{\frac{1+\gamma}{2}} \sqrt{\frac{\theta(1-\theta)}{n}} = \theta$$

and use the solutions of this equation (which will become quadratic if properly handled) as the limits of our confidence interval.

2.3 Problems

With a little help from my friends

The Beatles

1. Show that the distributions that attain the Cramér-Rao bound have densities of the form

$$f(x) = \exp(c(x) + a(x)A(\theta) + B(\theta)),$$

where

$$B'(\theta) = -\theta A'(\theta).$$

Besides, the efficient estimator is

$$\hat{\theta} = a(X),$$

or, for the case of a sample of size n

$$\hat{\theta}_n = \frac{1}{n} \sum_{i=1}^n a(X_i),$$

and coincides with the ML estimator.

2. Calculate the ML estimator for the parameter λ of the exponential distribution. Calculate the expectation and find a factor that will make it unbiased.
3. Let X_1, \dots, X_n be a sample from a uniform distribution on $[0, \theta]$ ($\theta > 0$).
- Calculate the ML estimator for θ .
 - Calculate an estimator for θ using the moment method.
 - Modify the ML estimator so that it is unbiased.
 - Compare the variances of the estimators in (b) and (c).
4. For the exponential distribution, calculate the Cramér-Rao bound for the variance of an unbiased estimator of the parameter λ . Does the unbiased estimator from problem 2 attain this bound?
5. Let $\mathbb{P}(X_i = 1) = \theta$, $\mathbb{P}(X_i = 0) = 1 - \theta$ ($0 < \theta < 1$). Show that \bar{X}_n is an efficient estimator for θ .
6. Show that the ML estimator for the parameter λ of a Poisson distribution is efficient.
7. Let X_1, \dots, X_n be a sample from a uniform distribution on $[\theta, 2\theta]$ ($\theta > 0$).
- Calculate the ML estimator for θ .
 - Modify the ML estimator so that it becomes unbiased.
 - Show that

$$\tilde{\theta} = \frac{1}{3}(\max(X_1, \dots, X_n) + \min(X_1, \dots, X_n))$$

is another unbiased estimator, and that it has smaller variance than the one in part (b).

8. Let (X_1, \dots, X_n) be a sample from a distribution with density

$$f_{\theta}(x) = \frac{1}{2} \exp(-|x - \theta|).$$

Calculate the ML estimator for θ .

9. Assume that the following sample is from a normal distribution:

1.1 2.1 1.5 2.3 1.4 2.6 1.7 1.3 2.8 1.6

Calculate 95% confidence intervals for μ and σ^2 .

10. In problem 9, calculate a 99% confidence interval for μ if $\sigma^2 = 0.5$ is known.
11. Use the central limit theorem to calculate an approximate confidence interval for the parameter λ of the Poisson distribution.
12. Use the central limit theorem to calculate an approximate confidence interval for the parameter λ of the exponential distribution.
13. Let X_1, X_2, \dots, X_n be a sample from an exponential distribution. Use the fact that $2n\lambda\bar{X}_n$ has a χ^2 -distribution with $2n$ degrees of freedom to calculate a confidence interval for λ .
14. Let X be Poisson with parameter λ , and $Y \chi^2$ with $2n$ degrees of freedom. Use the equality

$$\mathbb{P}(X \geq n) = \mathbb{P}(Y \leq 2\lambda)$$

to calculate a confidence interval for the parameter λ of a Poisson distribution.

(Hint: we try to find $l(x) < u(x)$ with $\mathbb{P}_\lambda(\lambda < l(X)) = \mathbb{P}(\lambda > u(X)) = \frac{1-\gamma}{2}$. Replace λ by $l(k)$ resp. $u(k)$ and use the above relation.)

15. Calculate a confidence interval for the parameter $\theta > 0$ of a uniform distribution on $[0, \theta]$.

Kapitel 3

Testing

A hundred experiments can't prove me right.
Yet, one single experiment can prove me wrong.

Sir Isaac Newton

3.1 Basic Notions

The problems we are studying in this chapter are of the following kind: based on a sample, we have to decide whether the underlying distribution has a given property or not. We may illustrate this with a simple example:

Suppose we want to take part in the statisticians' favorite game — coin tossing — and, in order to win for sure, we have a coin made that shows “heads” with probability $3/4$ and “tails” with probability $1/4$. But, as things go, nobody took care which way around the coin was struck, so we don't know which side is the one with the greater probability.

In this case, what we do, of course, is flip the coin a number of times (100, say) and choose the face that shows more often as the one with higher probability.

Using this way of making our decision, we can go wrong in either of two ways — we may decide for “heads” when “tails” would have been the right choice, and vice versa (In this example, both possibilities are symmetric, but this is not usually the case). The probability of making a wrong decision is

$$\mathbb{P}(X < 50 | \theta = .75) = \Phi\left(\frac{50 - 75}{\sqrt{100 * .75 * .25}}\right) = \Phi(-5.78) = 4 * 10^{-9}.$$

In general, things are less symmetric than in this example, and for either one of the two possibilities there may be more than one choice for the underlying probability distribution. This leads us to

Definition 3.1

A *hypothesis* is any nonempty subset H of the set of probability distributions in our statistical model.

If we have a parametric model, a hypothesis can be expressed in terms of the parameter. In that case, we speak of *parametric* hypotheses.

If the hypothesis contains only one distribution, it is called a *simple hypothesis*, otherwise it is a *composite hypothesis*.

For parametric hypotheses with a single real parameter, we can further distinguish one-sided (like $\theta < a$ or $\theta > b$) and two-sided (like $\theta \neq a$ or $\theta \notin [a, b]$) hypotheses.

With these definitions, we can restate the task of a statistical test as: decide between two hypotheses, using a random sample.

We are going to break symmetry now and call one hypothesis the “null” hypothesis H_0 and the other one the “alternative” (or sometimes the “one-hypothesis”) H_1 . Instead of saying which hypothesis we decide for, we are going to say that we “accept” or “reject” the null hypothesis. As the decision is based on a sample, we can define our test by the set of sample values that will make us reject the null hypothesis — the so-called “rejection region”, or, equivalently, by the indicator function ϕ of the rejection region, which has the nice feature that $\phi(X_1, \dots, X_n)$ is just the index of the hypothesis we choose as our decision. Thus, we arrive at the definition

Definition 3.2

A (non-randomized) test is a function from \mathbb{R}^n to $\{0, 1\}$.

A randomized test is a function from \mathbb{R}^n to $[0, 1]$, with the interpretation that $\phi(X_1, \dots, X_n)$ is the probability that we reject H_0 (e.g., if $\phi = 0.5$, we may flip a coin and reject if it shows “heads” and accept if it shows “tails”).

Now that we have defined our test, we may look at how things can go wrong:

Definition 3.3

An *error of the first kind* occurs if we reject H_0 although it is true.

An *error of the second kind* occurs if we accept H_0 although it is wrong.

Now, the probabilities of incurring an error of the first or second kind are functions of the actual distribution, so there is no fixed value for them if our hypotheses are composite. In fact, if $P \in H_0$,

$$\mathbb{P}(\text{first kind error}|P) = \mathbb{E}_P(\phi),$$

whereas for $P \in H_1$

$$\mathbb{P}(\text{second kind error}|P) = \mathbb{E}_P(1 - \phi).$$

The usual way to cope with this fact is to impose an upper limit on the probability of a first kind error, which is usually denoted by α and called the “*level of significance*” of the test (the most common choice is $\alpha = .05$).

Viewed as a function of the actual distribution (or of the actual parameter) the probability of rejecting H_0 is called the *power function* of the test, and the probability of accepting H_0 is called the *operation characteristic* of the test.

In most cases of practical interest, it is convenient to express the test in terms of a test statistic (often an estimator of the parameter, or some sufficient statistic), and to reject if the value of the test statistic exceeds (or falls short of) some critical value. So, the usual procedure in actually carrying out a test consists of the following steps:

1. State the null hypothesis and alternative.
2. Fix the level α .
3. Choose an appropriate test statistic.
4. Use the distribution of the test statistic under the null hypothesis to calculate the critical value.
5. Calculate the value of the test statistic from the sample and make a decision based on whether this value is above or below the critical value.

To be complete, in order to also include the possibility of using a randomized test, we would have to add a random experiment for deciding whether the null hypothesis should be rejected if the value of the test statistic agrees with the critical value, but this is not really relevant for practical applications: on one hand, a real-world customer would not be too pleased to learn that her question is ultimately decided by tossing a coin, and on the other hand, for all but the smallest sample sizes, the granularity of even a discrete distribution is usually fine enough that any given level of significance can be approximated with satisfying accuracy.

3.2 The Neyman-Pearson and likelihood ratio tests

For a given level of significance, one would like to have the probability of an error of the second kind as small as possible. Of course, the question arises whether this can be achieved. In the case that both H_0 and H_1 are simple, this question is answered by the following

Theorem 3.1

(Neyman-Pearson) If $H_0 = \{P_0\}$ and $H_1 = \{P_1\}$, then the best (randomized) test with level α (based on one observation) is the Neyman-Pearson test which has the form

$$\phi(x) = \begin{cases} 1 & \text{if } f_1(x) > kf_0(x), \\ c & \text{if } f_1(x) = kf_0(x), \\ 0 & \text{if } f_1(x) < kf_0(x), \end{cases}$$

where

$$f_i = \frac{dP_i}{d(P_0 + P_1)},$$

(actually, any density with respect to a common dominating measure works) and $k \in [0, \infty]$ and $c \in [0, 1]$ are calculated from the equation

$$\mathbb{E}_0(\phi) = \alpha.$$

(We write \mathbb{E}_0 and \mathbb{E}_1 to denote expectations with respect to P_0 and P_1 , respectively)

Remark: Although the theorem is stated for sample size one, the general case can be handled by viewing a sample of size n as one observation of an n -dimensional random variable; the Radon-Nikodym derivatives in the theorem then become the respective likelihood functions.

Proof. It is quite obvious that k and c can be chosen in such a way that $\mathbb{E}_0(\phi) = \alpha$. What we have to prove, then, is that the test we construct in this way is optimal, i.e., that any other test of level α has a larger probability for an error of the second kind. In other words, we have to prove that for any test ψ with

$$\mathbb{E}_0(\psi) \leq \alpha$$

we have

$$\mathbb{E}_1(\psi) \leq \mathbb{E}_1(\phi).$$

To this end, consider

$$\mathbb{E}_1(\psi) - k\mathbb{E}_0(\psi).$$

Letting $P = P_0 + P_1$, this becomes

$$\begin{aligned} \int \psi(x)(f_1(x) - kf_0(x))dP(x) &= \int_{f_1(x) < kf_0(x)} \psi(x)(f_1(x) - kf_0(x))dP(x) + \\ &\int_{f_1(x) = kf_0(x)} \psi(x)(f_1(x) - kf_0(x))dP(x) + \int_{f_1(x) > kf_0(x)} \psi(x)(f_1(x) - kf_0(x))dP(x) \leq \\ &\int_{f_1(x) < kf_0(x)} 0(f_1(x) - kf_0(x))dP(x) + \int_{f_1(x) = kf_0(x)} c(f_1(x) - kf_0(x))dP(x) + \\ &\int_{f_1(x) > kf_0(x)} 1(f_1(x) - kf_0(x))dP(x) = \\ &\mathbb{E}_1(\phi) - k\mathbb{E}_0(\phi). \end{aligned}$$

So,

$$\mathbb{E}_1(\psi) \leq \mathbb{E}_1(\phi) + k(\mathbb{E}_0(\psi) - \alpha) \leq \mathbb{E}_1(\phi).$$

Based on the ideas behind the Neyman-Pearson theorem is the following definition which gives us a general principle for constructing tests. Although there is no general theorem like the Neyman-Pearson theorem for it, it gives useful tests in many situations:

Definition 3.4

The likelihood ratio test is given by

$$\phi(X_1, \dots, X_n) = \begin{cases} 0 & \text{if } l > \lambda, \\ c & \text{if } l = \lambda, \\ 1 & \text{if } l < \lambda, \end{cases}$$

where the likelihood ratio statistic l is defined by

$$l = \frac{\sup_{\theta \in H_0} L(X_1, \dots, X_n, \theta)}{\sup_{\theta \in \Theta} L(X_1, \dots, X_n, \theta)}.$$

3.3 Special tests for the normal distribution

3.3.1 Tests for μ

In testing hypotheses for the mean μ of a normal distribution, we have to distinguish two cases, according to whether the variance σ^2 is known or not. In the easier case where σ^2 is known, we can start out from simple hypotheses and the Neyman-Pearson lemma which gives us (see the problem section), assuming $H_0 : \mu = \mu_0$ and $H_1 : \mu = \mu_1 > \mu_0$, rejection for

$$\bar{X}_n > \mu_0 + z_{1-\alpha} \sqrt{\frac{\sigma^2}{n}}.$$

The first observation that may be made is that the actual value of μ_1 does not enter in this test. This means that this test is the best test with level α for all $\mu_1 > \mu_0$. So, we may use the same test for testing $H_0 : \mu = \mu_0$ against the composite alternative $H_1 : \mu > \mu_0$. Further investigation shows that we can replace the simple null hypothesis by $H_0 : \mu \leq \mu_0$ because for $\mu < \mu_0$, the probability of rejecting decreases. Thus, the above test is optimal for testing $H_0 : \mu \leq \mu_0$ against $H_1 : \mu > \mu_0$.

For the two-sided alternative, namely for testing $H_0 : \mu = \mu_0$ against $H_1 : \mu \neq \mu_0$, we have two possible approaches. First, we may use some intuitive reasoning and use a test that rejects if the distance $\bar{X}_n - \mu$ is large in absolute value. This leads us to rejecting if

$$|\bar{X}_n - \mu_0| > z_{1-\frac{\alpha}{2}} \sqrt{\frac{\sigma^2}{n}}.$$

The second approach is to use the likelihood ratio method, which gives the same result.

Now, if σ^2 is unknown, we are going to replace it by its estimator, S_n^2 . As was the case with confidence intervals, we have to use the quantiles of the t -distribution with $n-1$ degrees of freedom and get the rejection regions

$$\bar{X}_n - \mu_0 > t_{n-1; 1-\alpha} \sqrt{\frac{S_n^2}{n}}$$

or

$$\bar{X}_n - \mu_0 < t_{n-1; 1-\alpha} \sqrt{\frac{S_n^2}{n}}$$

for the one-sided problems and

$$|\bar{X}_n - \mu_0| > t_{n-1; 1-\frac{\alpha}{2}} \sqrt{\frac{S_n^2}{n}}$$

in the two-sided case. These are usually called the t -tests.

3.3.2 Tests for σ^2

Using the ideas of the preceding section, we get the following tests for the variance of a normal distribution: If μ is known:

For $H_0 : \sigma^2 \leq \sigma_0^2$ against $H_1 : \sigma^2 > \sigma_0^2$: reject if

$$\sigma_0^{-2} \sum_{i=1}^n (X_i - \mu)^2 > \chi_{n;1-\alpha}^2.$$

For $H_0 : \sigma^2 \geq \sigma_0^2$ against $H_1 : \sigma^2 < \sigma_0^2$: reject if

$$\sigma_0^{-2} \sum_{i=1}^n (X_i - \mu)^2 < \chi_{n;\alpha}^2.$$

For $H_0 : \sigma^2 = \sigma_0^2$ against $H_1 : \sigma^2 \neq \sigma_0^2$: reject if

$$\sigma_0^{-2} \sum_{i=1}^n (X_i - \mu)^2 > \chi_{n;1-\frac{\alpha}{2}}^2.$$

or

$$\sigma_0^{-2} \sum_{i=1}^n (X_i - \mu)^2 < \chi_{n;\frac{\alpha}{2}}^2.$$

If μ is unknown:

For $H_0 : \sigma^2 \leq \sigma_0^2$ against $H_1 : \sigma^2 > \sigma_0^2$: reject if

$$\sigma_0^{-2} \sum_{i=1}^n (X_i - \bar{X}_n)^2 > \chi_{n-1;1-\alpha}^2.$$

For $H_0 : \sigma^2 \geq \sigma_0^2$ against $H_1 : \sigma^2 < \sigma_0^2$: reject if

$$\sigma_0^{-2} \sum_{i=1}^n (X_i - \bar{X}_n)^2 < \chi_{n-1;\alpha}^2.$$

For $H_0 : \sigma^2 = \sigma_0^2$ against $H_1 : \sigma^2 \neq \sigma_0^2$: reject if

$$\sigma_0^{-2} \sum_{i=1}^n (X_i - \bar{X}_n)^2 > \chi_{n-1;1-\frac{\alpha}{2}}^2.$$

or

$$\sigma_0^{-2} \sum_{i=1}^n (X_i - \bar{X}_n)^2 < \chi_{n-1;\frac{\alpha}{2}}^2.$$

3.3.3 Two-sample tests

In some cases, one wants to compare two samples from different populations having a normal distribution. Thus, there are two samples: (X_1, \dots, X_n) from a normal distribution with mean μ_1 and variance σ_1^2 and (Y_1, \dots, Y_m) with mean μ_2 and variance σ_2^2 . The most common question is: are the two means equal? Unfortunately, there is no simple test for this question without a restriction on the variance; if the two variances are equal, however, we can use the test statistic

$$\frac{\bar{X}_n - \bar{Y}_m}{\sqrt{\sum_{i=1}^n (X_i - \bar{X}_n)^2 + \sum_{j=1}^m (Y_j - \bar{Y}_m)^2}} \sqrt{\frac{nm(n+m-2)}{n+m}}$$

which has a t -distribution with $n+m-2$ degrees of freedom.

If we don't have any information on the relative sizes of the two variances, we can use the following test statistic:

$$\frac{\bar{X}_n - \bar{Y}_m}{\sqrt{\frac{1}{n(n-1)} \sum_{i=1}^n (X_i - \bar{X}_n)^2 + \frac{1}{m(m-1)} \sum_{j=1}^m (Y_j - \bar{Y}_m)^2}}$$

which has an approximate t distribution whose number of degrees of freedom f is obtained from the equation

$$\frac{1}{f} = \frac{c^2}{m-1} + \frac{(1-c)^2}{n-1},$$

where

$$c = \frac{\frac{1}{m(m-1)} \sum_{j=1}^m (Y_j - \bar{Y}_m)^2}{\frac{1}{n(n-1)} \sum_{i=1}^n (X_i - \bar{X}_n)^2 + \frac{1}{m(m-1)} \sum_{j=1}^m (Y_j - \bar{Y}_m)^2}.$$

From the above it seems natural that one might want to test for the equality of two variances, too. The case of the highest practical importance is the one where both means are unknown. In that case one uses the test statistic

$$F = \frac{S_X^2}{S_Y^2}.$$

where

$$S_X^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X}_n)^2,$$

and

$$S_Y^2 = \frac{1}{m-1} \sum_{i=1}^m (Y_i - \bar{Y}_m)^2.$$

which has an F -distribution with $(n-1, m-1)$ degrees of freedom. The F -distribution with (a, b) degrees of freedom is the distribution of bX/aY , where X and Y are independent χ^2 -distributed with a resp. b degrees of freedom.

3.4 Uniformly Optimal Tests

3.4.1 Monotone Likelihood Ratios

We return to the one-sided test for the mean of a normal distribution with a given variance that we saw in the previous section. For this particular test, starting from the optimal test for simple hypotheses, it was possible to extend first the alternative (because its actual value did not show in the final test, apart from its position — above or below — relative to the value of the null hypothesis), and then also the null (because as a test for the composite null, it has the same level of significance as for the simple null). Closer inspection of the calculations show that this is owing to the fact that, for $\mu_0 < \mu_1$, the likelihood ratio

$$\frac{L(X_1, \dots, X_n; \mu_1)}{L(X_1, \dots, X_n; \mu_0)} = \exp\left(\frac{n(\mu_0^2 - \mu_1^2) + 2(\mu_1 - \mu_0) \sum_{i=1}^n X_i}{2\sigma^2}\right)$$

is a nondecreasing function of the sufficient statistic

$$T = \sum_{i=1}^n X_i.$$

This property is worth a particular definition:

Definition 3.5

The dominated model $\mathcal{P} = \{P_\theta, \theta \in \Theta \subseteq \mathbb{R}\}$ is said to have *monotone likelihood ratios*, if there is a sufficient statistic T and if, for $\theta_1 < \theta_2$, the likelihood ratio

$$\frac{L(X_1, \dots, X_n; \theta_2)}{L(X_1, \dots, X_n; \theta_1)}$$

is a nondecreasing function of T .

Theorem 3.2: Karlin-Rubin

If the model \mathcal{P} has monotone likelihood ratios, then for any $\theta_0 \in \Theta$, there is a uniformly optimal test of level α for the one-sided hypotheses $H_0 : \theta \leq \theta_0$ against $H_1 : \theta > \theta_0$ of the form

$$\phi = \begin{cases} 1 & \text{if } T > t_c, \\ c & \text{if } T = t_c, \\ 0 & \text{if } T < t_c. \end{cases} \quad (3.1)$$

The proof is just a paraphrase of our observations about the normal distribution: we start from the test for $H_0 : \theta = \theta_0$ against $H_1 : \theta = \theta_1$ ($\theta_1 > \theta_0$). It is easily seen that there is a test of the form (3.1) that is equivalent to the Neyman-Pearson test (in the = part of the definition of the Neyman-Pearson test, the constant c may be replaced by a function, as long as its expectation under the null remains equal to α). Then, of course, the values c and t_c are determined with the help of the distribution of T under the null $\theta = \theta_0$, so the actual value of θ_1 does not enter the definition of ϕ , so this test is optimal for every $\theta_1 > \theta_0$. Thus, we are able to inflate the alternative. For $\theta_2 < \theta_0$, we can interpret ϕ as a test for $H_0 : \theta = \theta_2$ against $H_1 : \theta = \theta_0$ with level $\alpha' = \mathbb{E}_{\theta_2}(\phi)$, and as such, it is optimal. As such, it is at least as good as the trivial test that is constant $= \alpha'$, so

$$\alpha = \mathbb{E}_{\theta_0} \geq \alpha' = \mathbb{E}_{\theta_2}(\phi).$$

Thus, ϕ is of level α for the compound null $H_0 : \theta \leq \theta_0$.

3.4.2 Unbiased Tests.

Above, we have seen that the one-sided test for the mean μ of a normal distribution with known variance is optimal in the sense that any other test of the same level has at least the same probability of a second kind error for any $\mu \in H_1$. One may ask (and that's what we are doing) whether there is a similar optimality criterion for the two-sided test. We cannot expect that without further conditions, as the best test for $H_0 : \mu = \mu_0$ against $H_1 : \mu > \mu_0$ as well as the one for $H_1 : \mu < \mu_0$ can be used as tests for the two-sided alternative. A test that is optimal for all $\mu \neq \mu_0$ should therefore have the same probability for an error of the second kind as the test for $H_1 : \mu < \mu_0$ if $\mu < \mu_0$, and the same as the test for $H_1 : \mu > \mu_0$ if $\mu > \mu_0$. This, however, is not possible. Thus, we have to add some additional restriction to our tests:

Definition 3.6

A test ϕ of level α is called unbiased if for any $P \in H_1$:

$$\mathbb{E}_P(\phi) \geq \alpha.$$

(This is a natural requirement: we wouldn't want any distribution from H_1 to give us a higher probability of accepting H_0 than the distributions from H_0 .)

For a parametric test for $H_0 : \theta = \theta_0$ against $H_1 : \theta \neq \theta_0$, and assuming that the densities $f_\theta(x)$ are sufficiently differentiable, this implies that for an unbiased test

$$\mathbb{E}_{\theta_0}(\phi) = \alpha,$$

and

$$\frac{\partial}{\partial \theta}(\mathbb{E}_\theta(\phi))|_{\theta=\theta_0} = 0.$$

For a fixed θ_1 , we may copy the idea of the Neyman-Pearson lemma and try to find a test ϕ of the following form:

$$\phi(X_1, \dots, X_n) = \begin{cases} 1 & \text{if } L(X_1, \dots, X_n, \theta_1) > kL(X_1, \dots, X_n, \theta_0) + \tilde{k}L'(X_1, \dots, X_n, \theta_0), \\ c & \text{if } L(X_1, \dots, X_n, \theta_1) = kL(X_1, \dots, X_n, \theta_0) + \tilde{k}L'(X_1, \dots, X_n, \theta_0), \\ 0 & \text{if } L(X_1, \dots, X_n, \theta_1) < kL(X_1, \dots, X_n, \theta_0) + \tilde{k}L'(X_1, \dots, X_n, \theta_0). \end{cases}$$

If we can find such a test (with $0 \leq c \leq 1$, $k \geq 0$ and \tilde{k} determined by the equations

$$\mathbb{E}_{\theta_0}(\phi) = \alpha,$$

and

$$\frac{\partial}{\partial \theta}(\mathbb{E}_{\theta}(\phi))|_{\theta=\theta_0} = 0,$$

then we can prove in the same way as for the Neyman-Pearson lemma, that any other unbiased test must have a greater probability of a second kind error (it is not guaranteed, however, that we succeed); of course we hope that the actual choice of θ_1 doesn't enter in the final form of the test, so that our test will be uniformly optimal for all $\theta \neq \theta_0$.

Let us try this for the mean of a normal distribution with known variance: in that case, the likelihood function is

$$L(X_1, \dots, X_n, \mu) = \frac{1}{\sqrt{(2\pi\sigma^2)^n}} \exp\left(-\frac{(n-1)S_n^2 + n(\bar{X}_n - \mu)^2}{2\sigma^2}\right).$$

For its derivative, we have

$$L'(X_1, \dots, X_n, \mu) = \frac{n(\mu - \bar{X}_n)}{\sigma^2} \frac{1}{\sqrt{(2\pi\sigma^2)^n}} \exp\left(-\frac{(n-1)S_n^2 + n(\bar{X}_n - \mu)^2}{2\sigma^2}\right).$$

After some calculation, the inequality

$$L(X_1, \dots, X_n, \mu_1) > kL(X_1, \dots, X_n, \mu_0) + \tilde{k}L'(X_1, \dots, X_n, \theta_0)$$

turns into

$$\exp\left(\frac{n(2\bar{X}_n - \mu_0 - \mu_1)(\mu_1 - \mu_0)}{2\sigma^2}\right) < k + \tilde{k} \frac{n(\mu - \bar{X}_n)}{\sigma^2}.$$

The right hand side is linear in \bar{X}_n , whereas the left hand side is a convex function of \bar{X}_n . Thus, there are (at most) 2 points a and b where the two functions coincide, and the inequality above is equivalent to

$$a < \bar{X}_n < b.$$

The distribution of \bar{X}_n is continuous, so we don't need the number c in the definition of our test, and it becomes

$$\phi = \begin{cases} 0 & \text{if } a \leq \bar{X}_n \leq b, \\ 1 & \text{if } \bar{X}_n < a \text{ or } \bar{X}_n > b. \end{cases}$$

a and b are to be determined from the conditions

$$\mathbb{E}_{\mu_0}(\phi) = \alpha,$$

and

$$\frac{\partial}{\partial \mu}(E_{\mu}(\phi))|_{\mu=\mu_0} = 0.$$

We first calculate

$$\mathbb{E}_{\mu}(\phi) = \mathbb{P}_{\mu}(\bar{X}_n < a) + \mathbb{P}_{\mu}(\bar{X}_n > b) = \Phi\left(\frac{a - \mu}{\sqrt{\sigma^2/n}}\right) + 1 - \Phi\left(\frac{b - \mu}{\sqrt{\sigma^2/n}}\right).$$

taking derivatives with respect to μ , we obtain

$$\frac{\partial}{\partial \mu}(E_{\mu}(\phi)) = \frac{1}{\sqrt{2\sigma^2/n}} e^{-\frac{(b-\mu)^2}{2\sigma^2/n}} - \frac{1}{\sqrt{2\sigma^2/n}} e^{-\frac{(a-\mu)^2}{2\sigma^2/n}}.$$

Thus, the equation

$$\frac{\partial}{\partial \mu}(E_{\mu}(\phi))|_{\mu=\mu_0} = 0$$

is equivalent to

$$(a - \mu)^2 = (b - \mu)^2,$$

and the final test turns out to be identical to the one we already know.

In a more general setting, we may consider

Definition 3.7: Exponential family

We call the dominated model $\mathcal{P} = P_\theta, \theta \in \Theta$ an *exponential family* if its densities admit a representation

$$f_\theta(x) = C(\theta)g(x) \exp\left(\sum_{i=1}^k \tau_i(\theta)\xi_i(x)\right).$$

If $\Theta \subseteq \mathbb{R}^k$ and $\tau_i(\theta) = \theta_i$, we speak of a naturally parametrized exponential family.

In particular, we consider one-parameter exponential families, i.e., the case $k = 1$, so assuming natural parametrization

$$f_\theta(x) = g(x)\exp(\theta\xi(x) - \gamma(\theta)),$$

with a convex function γ (we could consider the case that τ is strictly monotone which is slightly more general than natural parametrization).

It is readily seen that

$$T = \sum_{i=1}^n \xi(X_i)$$

is a sufficient statistic, and that this family has monotone likelihood ratios.

As in the case of the normal distribution that we studied above, we try to find an unbiased test for $H_0 : \theta = \theta_0$ that is optimal for a given $\theta_1 \neq \theta_0$. The remainder of the argument proceeds almost literally as in the special case, leading to

Theorem 3.3

Let $\mathcal{P} = \{P_\theta, \theta \in \Theta\}$ be a naturally parametrized exponential family. Then there is a uniformly optimal unbiased test for $H_0 : \theta = \theta_0$ against the two-sided alternative of the form

$$\phi = \begin{cases} 1 & \text{if } T < a, \\ c_1 & \text{if } T = a, \\ 0 & \text{if } a < T < b, \\ c_2 & \text{if } T = b, \\ 1 & \text{if } T > b, \end{cases}$$

where $a < b$ and $c_1, c_2 \in [0, 1]$ are determined by the equations

$$h(\theta_0) = \mathbb{E}_{\theta_0}(\phi) = \alpha$$

and

$$h'(\theta_0) = \mathbb{E}_{\theta_0}((T - n\gamma'(\theta_0))\phi) = 0.$$

3.5 Tests and confidence intervals

If we take a closer look at the tests for the mean of a normal distribution, it may be noted that we reject the null $H_0 : \mu = \mu_0$ at level α iff μ_0 is contained in the confidence interval with coverage probability $\gamma = 1 - \alpha$.

This is a general feature:

Theorem 3.4

If I is a confidence interval for θ with coverage probability γ , then the rule “reject if θ_0 is not in I ” yields a level $\alpha = 1 - \gamma$ test for the null $H_0 : \theta = \theta_0$.

On the other hand, if, for every $\theta_0 \in \Theta$, we are given a Test ϕ_{θ_0} of level α for the null

$H_0 : \theta = \theta_0$, then

$$I = \{\theta_0 \in \Theta : \phi_{\theta_0} \text{ accpets}\}$$

defines a confidence region for θ with coverage probability $\gamma = 1 - \alpha$.

Remark. The confidence “interval” obtained from the second part of this theorem need not be an interval, in fact, it is not too hard to come up with an example where it is not even measurable.

The proof is just an application of the definitions.

3.6 Problems

We can swing together 'cause we think
we're doin' it right

Lindisfarne

1. For a sample from a normal distribution, calculate the best test for $H_0 : \mu = \mu_0$ against $H_1 : \mu = \mu_1$ if σ^2 is known.
2. For a sample from a normal distribution, calculate the likelihood ratio test for $H_0 : \mu = \mu_0, 0 < \sigma^2 < \infty$ against the two-sided alternative. Show that this test is equivalent to the t -test.
3. For a sample from a Poisson distribution, find the best test for $\lambda = \lambda_0$ against $\lambda = \lambda_1$.
4. For a sample from a uniform distribution on $[0, \theta]$ construct a test for $H_0 : \theta < \theta_0$ against $H_1 : \theta > \theta_0$.
5. If there is a sufficient statistic, show that the likelihood ratio test is always a function of that sufficient statistic.
6. Find the form of the best test for $H_0 : \mu = \mu_0, \sigma^2 = \sigma_0^2$ against $H_1 : \mu = \mu_1, \sigma^2 = \sigma_1^2$ (for a normal distribution).
7. For an exponential distribution, construct a test for testing $H_0 : \lambda = \lambda_0$ against the two-sided alternative.
8. Let (X_1, \dots, X_n) and (Y_1, \dots, Y_m) be two samples from normal distributions with parameters (μ_1, σ_1^2) and (μ_2, σ_2^2) , respectively. Calculate the likelihood ratio test for testing $\sigma_1^2 = \sigma_2^2$ and show that it is equivalent to the F -test.
9. In a sample of 200 products, 15 defective items were found. The manufacturer claims that the percentage of defective items is at most 5%. Test at the 5% level.
10. A sample of 60 bottles of beer showed a mean content of 470 ccm with a sample variance of 2000. Assuming a normal distribution, test $H_0 : \mu = 500$.
11. The lifetime of lightbulbs is supposed to be exponentially distributed. A certain brand is claimed to have an expected lifetime of 2000 hours. A sample showed the following lifetimes:
900 1760 2140 2050 2820 2200 3000 600 1500 1980
Test at the 5% level.
12. For the following sample from a normal distribution
1.1 0.7 1.6 1.3 0.8 1.4 0.7 0.6 1.8 1.4 1.3
test $H_0 : \mu = 1.0$ against the two-sided alternative, if $\sigma^2 = 0.5$
13. Do the previous problem for unknown variance.
14. For the sample in problem 12, test whether $\sigma^2 = 0.5$ ($\alpha = 0.1$).

15. Calculate the density of the F -distribution.

16. Use the central limit theorem to show that for large n

$$\chi_{n,p}^2 \approx n + z_p \sqrt{2n}.$$

17. Use the result of the previous problem to prove that for large a and b

$$F_{a,b,p} \approx 1 + z_p \sqrt{\frac{2(a+b)}{ab}}.$$

(hint: Replace the χ_n^2 -distributed random variables by $n + U\sqrt{2n}$, where U is approximately standard normal.)

18. $((X_1, Y_1), \dots, (X_n, Y_n))$ is a sample from a two-dimensional normal distribution. Devise a test for $H_0 : \mu_X = 2\mu_Y$.

(hint: any linear combination of X_n and Y_n is normally distributed.)

19. Suppose that (X_1, \dots, X_n) and (Y_1, \dots, Y_m) are samples from normal distributions with equal variances. Construct a test for $H_0 : \mu_X = 2\mu_Y$.

(hint: this should look similar to the test for equal means, only the denominator has to be changed, and the scaling factor needs some adjustment.)

20. The following is another sample from a normal distribution:

1.0 0.9 0.7 1.4 1.1 0.8 0.6 0.9 1.1

Assuming equal variances, test whether this has the same mean as the sample in problem 12

21. In the previous problem, use the F -test to check whether the variances are equal ($\alpha = .1$).

Kapitel 4

Analysis of Variance (ANOVA)

Little boxes on the hillside
little boxes made of ticky-tacky
little boxes on the hillside
little boxes all the same

Pete Seeger

4.1 The Fisher-Cochran Theorem

The name “analysis of variance” is slightly misleading; in fact, what this chapter is about is a comparison of the *means* of several samples from different normal distributions. The variances don’t enter the picture - they are even supposed to be the same for all the individual samples. There is, however, a reason why this name persists: namely, we use a decomposition of the sum of the squares of the sample values for our test.

To put things in more precise terms, let k samples $(X_{i1}, \dots, X_{in_i})$ ($i = 1, \dots, k$) be given. assume that X_{ij} has a normal distribution with mean μ_i and variance σ^2 (the same for all samples). We want to test $H_0 : \mu_1 = \mu_2 = \dots = \mu_k$ against the alternative that at least one μ_i is different from the others. This question could of course also be approached by using the likelihood ratio method and you will be pleased to know (left as an exercise) that the method we develop here is equivalent to the likelihood ratio test.

We start out from the sum

$$S_I = \sum_{i=1}^k \sum_{j=1}^{n_i} (X_{ij} - \bar{X}_N)^2,$$

where

$$N = \sum_{i=1}^k n_i$$

and

$$\bar{X}_N = \frac{1}{N} \sum_{i=1}^k \sum_{j=1}^{n_i} X_{ij}$$

is the mean of the combined samples.

By simple calculations,

$$\sum_{j=1}^{n_i} (X_{ij} - \bar{X}_N)^2 = \sum_{j=1}^{n_i} (X_{ij} - \bar{X}_i)^2 + n_i (\bar{X}_i - \bar{X}_N)^2,$$

where

$$\bar{X}_i = \frac{1}{n_i} \sum_{j=1}^{n_i} X_{ij}$$

is the mean of the i th sample. Thus, we arrive at the decomposition

$$S_I = S_{II} + S_{III},$$

where

$$S_{II} = \sum_{i=1}^k \sum_{j=1}^{n_i} (X_{ij} - \bar{X}_i)^2,$$

and

$$S_{III} = \sum_{i=1}^k n_i (\bar{X}_i - \bar{X}_N)^2.$$

The distribution of S_{II} is not affected by the actual values of μ_1, \dots, μ_k . For S_{III} , we can expect it to be smallest if all means are equal, and to increase if we depart from this hypothesis (because \bar{X}_i is close to μ_i). So, we will reject our null hypothesis of equal means if S_{III} is large with respect to S_{II} .

In order to obtain the critical value for such a test, we need the distributions of the sums S_{II} and S_{III} . To obtain those, we will do a little general study of sums of squares. As usual, we start with some definitions:

Definition 4.1

Let X_1, \dots, X_n be independent standard normal random variables. Furthermore, let

$$Y_i = \sum_{j=1}^n a_{ij} X_j \quad (i = 1, \dots, k),$$

where (a_{ij}) are real numbers. Then we call

$$\sum_{i=1}^k Y_i^2$$

a sum of squares.

Furthermore, the maximum number of linearly independent random variables among the Y 's (which is k minus the number of linear relations among the Y 's) is called the number of degrees of freedom of the sum of squares. This is the maximum number of Y 's that can be independently chosen, whereas the rest will follow from the linear relations among them.

The first thing to observe is that we have to prove that the number of degrees of freedom is well defined, i.e., it does not depend on the special choice of a_{ij} . It seems most convenient to use matrix notation for this. So, let

$$X = (X_1, \dots, X_n)^T$$

and

$$Y = (Y_1, \dots, Y_k)^T$$

be the column vectors whose elements are the numbers X_i resp. Y_i , and

$$A = (a_{ij})_{k \times n}$$

be the matrix with k rows and n columns whose elements are a_{ij} .

Then

$$Y = AX$$

and the sum of squares can be written as

$$S = Y^T Y = X^T A^T A X.$$

The number of degrees of freedom of S is readily identified as the rank of A . But this is also the rank of $A^T A$, which follows from the fact that their null spaces coincide. Namely, if

$$Ax = 0,$$

then also

$$A^T Ax = 0;$$

on the other hand

$$A^T Ax = 0$$

implies

$$(Ax)^T(Ax) = x^T A^T Ax = 0,$$

so we also have

$$Ax = 0.$$

Now, it is immediate that, if \tilde{A} is another matrix that generates the same sum of squares, then $\tilde{A}^T \tilde{A} = A^T A$, so the ranks of A and \tilde{A} must be the same.

Examples.

1. $X_1^2 + \dots + X_k^2$ is the simplest sum of squares, and has k degrees of freedom.
2. $(X_1 - \bar{X}_n)^2 + \dots + (X_n - \bar{X}_n)^2$ is another sum of squares. Its number of degrees of freedom is $n - 1$, because if we put

$$Y_i = X_i - \bar{X}_n,$$

then we have the linear relation

$$\sum Y_i = 0.$$

It is obvious that there are no other linear relations.

3. $5X_1^2 + 3X_2^2 + 17(3X_1 - 5X_2)^2$ has 2 degrees of freedom.

The principal tool in our treatment of sums of squares is the following

Theorem 4.1

(FISHER-COCHRAN) If the sum of squares S has a χ^2 -distribution with f degrees of freedom and if, for $i = 1, \dots, k$, S_i is a sum of squares with f_i degrees of freedom satisfying

$$S = S_1 + \dots + S_k,$$

then S_1, \dots, S_k are independent and S_i is χ^2 -distributed with f_i degrees of freedom if and only if

$$f = f_1 + \dots + f_k.$$

Remark.

1. Observe that the independence of S_1, \dots, S_k is not required; instead, it will follow from the theorem if the sum of their numbers of degrees of freedom is the number of degrees of freedom of their sum.
2. It is a surprising fact that if a sum of squares is χ^2 -distributed with f degrees of freedom, then f is the number of degrees of freedom of f (this can be deduced from the proof below).

Proof. The “if” part simply states that a sum of independent χ^2 -distributed random variables is χ^2 -distributed with a number of degrees of freedom that is the sum of the numbers of degrees of freedoms of the summands. So, what remains to be proven is the other direction. First, it is sufficient to prove the theorem for $k = 2$. Namely, it is obvious that the number of degrees of freedom of a sum of sums of squares is at most the sum of the numbers of degrees of freedom of the summand terms. Thus $S_2 + \dots + S_k$ can have at most $f_2 + \dots + f_k$ degrees of freedom. If it

had less than that, then $S_1 + \dots + S_k$ would have less than $f_1 + \dots + f_k$ degrees of freedom, which is in violation of our assumptions. Thus, the general case is obtained from $k = 2$ by induction.

Assume that

$$S = X_1^2 + \dots + X_f^2.$$

This entails no loss of generality, because any sum of squares can be written in the form

$$S = \sum_{i=1}^k a_i Y_i^2,$$

where Y_i are independent standard normal random variables that are linear combinations of X_i , and a_i are some positive constants. Namely, if we use matrix notation again, we have

$$S = X^T A^T A X.$$

$A^T A$ is a symmetric nonnegative definite matrix, so there is an orthogonal matrix B (i.e., $B^T B = I$) such that $B^T A^T A B$ is a diagonal matrix (its diagonal elements are the eigenvalues λ_i of $A^T A$). Let

$$Y = B^T X.$$

It follows that Y is a vector of independent standard normal random variables, and we have

$$S = X^T A^T A X = Y^T B^T A^T A B Y = \sum_{i=1}^f \lambda_i Y_i^2.$$

It is an easy application of the theory of characteristic functions (left as an exercise to the reader) to prove that if that sum has a χ^2 -distribution with f degrees of freedom, then $k = f$ and $\lambda_i = 1$.

Now, S_1 and S_2 can be written as

$$S_i = X^T A_i^T A_i X \quad (i = 1, 2),$$

where A_i are $k_i \times f$ matrices with rank f_i . As $A_i^T A_i$ is a symmetric nonnegative definite matrix, there is again an orthogonal $f \times f$ matrix B (i.e., $B^T B = I$) such that $B^T A_1^T A_1 B$ is a diagonal matrix (and the diagonal elements are the eigenvalues λ_i of $A_1^T A_1$). Let

$$Z = B^T X.$$

It follows that Z is a vector of independent standard normal random variables, and we have

$$\sum_{i=1}^f Z_i^2 = Z^T Z = X^T B B^T X = X^T X = \sum X_i^2 = S$$

and

$$S_1 = X^T A_1^T A_1 X = Z^T B^T A_1^T A_1 B Z = \sum_{i=1}^f \lambda_i Z_i^2.$$

Obviously, the number of degrees of freedom of S_1 is equal to the rank of $B^T A_1^T A_1 B$, which is the number of nonzero elements among $\lambda_1, \dots, \lambda_f$. This gives us

$$S_2 = S - S_1 = \sum_{i=1}^f (1 - \lambda_i) Z_i^2.$$

The number of degrees of freedom of S_2 is the number of nonzero elements among $1 - \lambda_1, \dots, 1 - \lambda_f$. We know that exactly f_1 among the λ_i are nonzero, so $f_2 = f - f_1$ are zero. For those f_2 , $1 - \lambda_i$ is nonzero, so the rest must be zero, or otherwise the number of degrees of freedom of S_2 would be greater than f_2 . Thus, exactly f_1 of the λ_i are 1, and the remaining f_2 are 0. This means that, for $i = 1, 2$, the sum S_i is the sum of the squares of f_i independent standard normal random variables, which proves the theorem.

4.2 One-way analysis of variance

Now we can go about evaluating the null distribution for our test of $H_0 : \mu_1 = \mu_2 = \dots = \mu_k$ against the alternative that the means are different. Remember that our test statistic was

$$\frac{\frac{1}{k-1} \sum_{i=1}^k n_i (\bar{X}_i - \bar{X})^2}{\frac{1}{n-k} \sum_{i=1}^k \sum_{j=1}^{n_i} (X_{ij} - \bar{X}_i)^2}.$$

First of all, it is clear that the distribution of this statistic does not depend on σ^2 and μ (which is the common value of $\mu_1 \dots \mu_k$), so we may assume that the X_{ij} are standard normal.

The sum

$$\sum_{i=1}^k \sum_{j=1}^{n_i} (X_{ij} - \bar{X})^2.$$

Has a χ^2 distribution with $n - 1$ degrees of freedom (in fact, this is just the sum of the squares of the differences of n standard normal random variables and their average). The inner sum can be decomposed in the following way:

$$\sum_{j=1}^{n_i} (X_{ij} - \bar{X})^2 = \sum_{j=1}^{n_i} (X_{ij} - \bar{X}_i)^2 + n_i (\bar{X}_i - \bar{X})^2.$$

Summing up over i we get

$$\sum_{i=1}^k \sum_{j=1}^{n_i} (X_{ij} - \bar{X})^2 = \sum_{i=1}^k \sum_{j=1}^{n_i} (X_{ij} - \bar{X}_i)^2 + \sum_{i=1}^k n_i (\bar{X}_i - \bar{X})^2.$$

The sum on the left-hand side has $n - 1$ degrees of freedom. On the right-hand side, the first sum has $n - k$ degrees of freedom (summing $X_{ij} - \bar{X}_i$ over j gives zero for each i), and the second sum has $k - 1$ degrees of freedom (here, $\sum n_i (\bar{X}_i - \bar{X}) = 0$ is the only linear relation). Thus, the number of degrees of freedom on the left-hand side is the sum of the numbers of degrees of freedom on the right-hand side, so we can apply the Fisher-Cochran theorem and get that

$$S_1 = \sum_{i=1}^k \sum_{j=0}^{n_i} (X_{ij} - \bar{X}_i)^2$$

has a χ^2 -distribution with $n - k$ degrees of freedom, and

$$S_2 = \sum_{i=1}^k n_i (\bar{X}_i - \bar{X})^2$$

has a χ^2 distribution with $k - 1$ degrees of freedom, and S_1 and S_2 are independent. This implies that our test statistic

$$F = \frac{S_2 / (k - 1)}{S_1 / (n - k)}$$

has an F distribution with $k - 1$ and $n - k$ degrees of freedom. So, we reject H_0 if $F > F_{k-1, n-k, 1-\alpha}$.

Once we have found out that there is a difference between the sample (by getting a rejection from the F -test), we might want to find a confidence interval for μ_i (or a difference $\mu_i - \mu_j$). This can be done using the standard formula for sample i , but this would be a little wasteful because it doesn't make use of all the information that is available. In fact

$$\hat{\sigma}^2 = \frac{1}{n - k} \sum_{i=1}^k \sum_{j=0}^{n_i} (X_{ij} - \bar{X}_i)^2$$

is a more accurate estimate of σ^2 than the sample variance of the i th sample alone. This leads us to the confidence interval

$$\left[\bar{X}_i - t_{n-k, \frac{1+\gamma}{2}} \sqrt{\frac{\hat{\sigma}^2}{n_i}}, \bar{X}_i + t_{n-k, \frac{1+\gamma}{2}} \sqrt{\frac{\hat{\sigma}^2}{n_i}} \right].$$

for μ_i , or, for $\mu_i - \mu_j$,

$$[\bar{X}_i - X_j - t_{n-k, \frac{1+\gamma}{2}} \sqrt{\frac{\hat{\sigma}^2(n_i + n_j)}{n_i n_j}}, \bar{X}_i - X_j + t_{n-k, \frac{1+\gamma}{2}} \sqrt{\frac{\hat{\sigma}^2(n_i + n_j)}{n_i n_j}}].$$

4.3 Two-way analysis of variance

In this section, the fun will be increased even more by allowing more than one influence on our data; in the simplest example, we may assume that some industrial product is worked upon by two different types of machines, A and B ; the final length (or weight, or whatever) of the product is supposed to be the sum of the influences of both machines. Assume that we have machines A_1, \dots, A_n resp. B_1, \dots, B_k of each type. For testing purposes, we may produce one product with each possible combination of machines. Let X_{ij} be the size of the product produced with the combination A_i with B_j . Under our assumptions, X_{ij} will have a normal distribution with mean $\mu_i + \nu_j$ and variance σ^2 (yes, again we demand equal variances). We may want to test whether the influences of the machines of type A are all equal, for example. We shall see in a minute how this is done, but first observe that we get the same results for $\mu_i + \nu_j$ if we add a constant to every μ_i and subtract the same constant from every ν_j . We avoid the ambiguities that may arise from this by writing $\mu_{ij} = a_i + b_j + \mu$ and demanding that $\sum a_i = \sum b_j = 0$. With those conventions, we want to test

$$H_0 : a_i = 0, i = 1, \dots, n$$

against the alternative that at least one a_i is different from zero.

We could find the decomposition of the sum of squares

$$\sum_{i=1}^n \sum_{j=1}^k (X_{ij} - \bar{X})^2$$

that fits our needs by simple heuristic arguments, yet it is desirable to have a general principle that could be used in a greater number of cases. This will be the likelihood principle under a new name. For the normal distribution, finding the maximum of the likelihood function is equivalent to minimizing the sum of squares

$$\sum (X_{ij} - \mu_{ij})^2.$$

Thus, if we suppose that μ depends on a parameter θ , then the ML estimator will be the number that minimizes the sum of squares above, hence it is called the least squares estimator (least squares estimators can be used for other distributions than the normal distribution, but they are at their best if there is a normal distribution to work with).

In order to study the general theory, let us replace the double index above by a single one, and let us assume that, for $i = 1 \dots, n$, each X_i has a normal distribution with mean μ_i and variance σ^2 . Furthermore, let μ_i be a linear function of k parameters $\theta_1, \dots, \theta_k$. We can write this down in matrix notation by defining the column vectors

$$X = (X_1, \dots, X_n)^T,$$

$$\mu = (\mu_1, \dots, \mu_n)^T$$

and

$$\theta = (\theta_1, \dots, \theta_k)^T$$

and a $n \times k$ matrix A such that

$$\mu = A\theta.$$

We assume that $k < n$ and that A has rank k .

We wish to test the null hypothesis that some set of linear relations among the parameters pertains, i.e., there is some $r \times k$ matrix B such that we have $H_0 : B\theta = 0$. Again $r < k$, and the

rank of B is supposed to be r . We can first find the estimator for θ without the restriction $B\theta = 0$ by the usual method of zeroing the derivatives with respect to θ_i which gives us

$$\hat{\theta} = (A^T A)^{-1} A^T X.$$

Letting $\hat{\mu} = A\hat{\theta}$, and using some matrix algebra, we find that

$$X^T X = (X - \hat{\mu})^T (X - \hat{\mu}) + \hat{\mu}^T \hat{\mu}.$$

If the X_i have standard normal distributions, then the sum on the left-hand side has a χ^2 -distribution with n degrees of freedom.

The first sum on the right-hand side has $n - k$ degrees of freedom because there are k linear relations between the components of $X - \hat{\mu}$ by virtue of the equation $A^T(X - \hat{\mu}) = 0$. The second sum on the right-hand side has k degrees of freedom, as this is the rank of the matrix $A(A^T A)^{-1} A^T$ that generates this sum of squares. Thus, by the Fisher-Cochran theorem, the distributions of the sums on the right-hand side are χ^2 with $n - k$ and k degrees of freedom, respectively.

If we take the additional condition $B\theta = 0$ into account, the least squares estimator can be obtained by Lagrange's multiplier method, giving an estimator $\tilde{\theta}$ and a vector of estimated means $\tilde{\mu}$. A little more matrix algebra brings us to the decomposition

$$\hat{\mu}^T \hat{\mu} = \tilde{\mu}^T \tilde{\mu} + (\hat{\mu} - \tilde{\mu})^T (\hat{\mu} - \tilde{\mu}).$$

The numbers of degrees of freedom are k for the left-hand side, $k - r$ for the first sum on the right, and r for the third sum on the right. Thus, by the Fisher-Cochran theorem, they have the corresponding χ^2 distributions.

Finally, we arrive at the following: The statistic

$$F = \frac{\sum (\hat{\mu}_i - \tilde{\mu}_i)^2 / r}{\sum (X_i - \hat{\mu}_i)^2 / (n - k)}$$

has a F distribution with r and $n - k$ degrees of freedom, where k is the number of degrees of freedom of our parameter space (=number of parameters- number of linear relations among them), and r is the number of independent equations in our null hypothesis.

In our original question, there are nk Variables X_{ij} , $n+k+1$ parameters $(\mu, a_1, \dots, a_n, b_1, \dots, b_k)$ and 2 linear relations ($\sum a_i = \sum b_i = 0$), giving us $n + k - 1$ degrees of freedom for the parameter space. Our null hypothesis $a_1 = \dots = a_n = 0$ has only $n - 1$ linearly independent equations, because the last one is implied by the condition $\sum a_i = 0$.

We see, that in the general analysis above, we have to replace n by nk , k by $n + k - 1$, and r by $n - 1$.

If we write $X_{.j}$ for $n^{-1} \sum_{i=1}^n X_{ij}$, and use the analogous definitions for $X_{i.}$ and $X_{..}$, we can write our test statistic as

$$F = \frac{\sum_{i,j} (X_{i.} - X_{..})^2 / (n - 1)}{\sum_{i,j} (X_{ij} - X_{i.} - X_{.j} + X_{..})^2 / (n - 1)(k - 1)} = \frac{(k - 1)k \sum_i (X_{i.} - X_{..})^2}{\sum_{i,j} (X_{ij} - X_{i.} - X_{.j} + X_{..})^2},$$

and we reject H_0 if F exceeds $F_{n-1, (n-1)(k-1); 1-\alpha}$.

4.4 Problems.

We shall overcome some day.

Pete Seeger

1. Modify the two-way test for the case that we have samples of sizes n_{ij} instead of single observations for each combination of i and j .

2. (Two-way ANOVA with interaction) In general, we cannot expect the effects of two different influences to just add up. This prompts us to consider the more general model

$$\mu_{ij} = \mu + a_i + b_j + c_{ij}$$

with the conditions

$$\begin{aligned} \sum_i a_i &= \sum_j b_j = 0, \\ \sum_i c_{ij} &= 0 \quad (j = 1, \dots, k), \\ \sum_j c_{ij} &= 0 \quad (i = 1, \dots, n). \end{aligned}$$

Here a_i and b_j are the usual linear influences, and the interaction terms c_{ij} represent the way the joint influence deviates from the simple summation of the individual influences.

In this case, there are too many parameters to work with single observations, so we have to use samples of size $n_{ij} = N$ (i.e., all the samples have equal sizes) for each combination of i and j . Calculate a test for $H_0 : c_{ij} = 0$.

3. Show that for the comparison of two samples, ANOVA is equivalent to the t -test.
 4. For four samples of iron ore from different providers, the metal contents were measured:

i	n_i	\bar{X}_i	S_i^2
1	10	5.3	0.11
2	12	5.1	0.13
3	10	5.0	0.12
4	11	5.4	0.10

Test for equality of the means.

5. in the previous example, calculate a confidence interval for μ_1 .
 6. Compare the confidence interval from the previous example with the one calculated from the first sample alone.
 7. Construct a confidence interval for a_i in the two-way ANOVA.
 8. In a factory, ball bearings are manufactured in two production stages. There are four different units for the first stage, and three different units for the second stage. One ball bearing for each combination is measured for its inner diameter:

j	i			
	1	2	3	4
1	10.5	10.4	10.1	9.9
2	10.3	10.3	10.2	10.0
3	10.2	10.1	10.4	10.3

Test whether there is a difference between the units in stage 1.

9. In the previous problem, assume that two bearings are made for each combination:

i j	1			2			3			4		
	1	2	3	1	2	3	1	2	3	1	2	3
1	10.6	10.4	10.3	10.5	10.4	10.2	10.2	10.3	10.5	10.0	10.1	10.4
2	10.4	10.2	10.1	10.2	10.0	10.4	10.2	10.1	10.3	10.3	9.9	10.2

Again, test whether there is a difference between the units of stage 1.

Kapitel 5

Linear Regression.

And you know something is happening
but you don't know what it is
do you, Mr. Jones?

B. Dylan

5.1 Simple linear regression.

In this chapter, we investigate simple functional relationships between two random variables. We assume that x and y are related through the equation

$$y = ax + b,$$

where a and b are real constants. If we have to find the values of a and b from measurements, and if all measurements are exact, we can calculate a and b from only two observations. But, alas, things usually are not as perfect as they should be, so we have to account for measurement errors in some way. So, we rather write

$$y = ax + b + e,$$

where e is a (random) error (not necessarily $=2.71828\dots$). Thus, our task is to estimate a and b from n observations. We fix numbers x_1, \dots, x_n — those are supposed to be nonrandom; there is another branch of the theory that deals with random x 's (called correlation theory), but we won't touch that here; a great part of regression theory, however, can be considered as correlation theory conditioned on x_1, \dots, x_n .

Now, we have to make some assumptions on the errors e_i . We want those to have expectation zero (which means that there is no systematic error in the measurements) and a variance σ^2 that is the same for all observations. Then, as usual, the errors for different measurements are supposed to be independent. With these assumption alone, one can develop some part of the theory. In order to get confidence intervals, or for testing, we must make some assumptions on the distribution of e_i , and, of course, we pick the normal distribution. So, for us, e_i will have a normal distribution with mean 0 and variance σ^2 . The estimators we use for a and b will be the maximum likelihood estimators which turn, as is customary for the normal distribution, into least squares estimators. The calculations yield

$$\hat{a} = \frac{\sum_{i=1}^n (x_i - \bar{x})(Y_i - \bar{Y})}{\sum_{i=1}^n (x_i - \bar{x})^2}.$$

and

$$\hat{b} = \bar{Y} - \hat{a}\bar{x}.$$

The expectations of \hat{a} and \hat{b} are of course the parameters a and b themselves. With a little work, one calculates the variances

$$\mathbb{V}(\hat{a}) = \frac{1}{\sum_{i=1}^n (x_i - \bar{x})^2} \sigma^2$$

and

$$\mathbb{V}(\hat{b}) = \left(\frac{1}{n} + \frac{\bar{x}^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right) \sigma^2.$$

We still need an estimator for σ^2 . This is achieved by

$$\hat{\sigma}^2 = \frac{1}{n-2} \sum_{i=1}^n (Y_i - \hat{a}x_i - \hat{b})^2.$$

By the Fisher-Cochran theorem, the sum in the above formula, divided by σ^2 , has a χ^2 distribution with $n-2$ degrees of freedom, and is independent of \hat{a} and \hat{b} . This directly leads us to the confidence intervals

$$\left[\hat{a} - t_{n-2; \frac{1+\gamma}{2}} \sqrt{\frac{\hat{\sigma}^2}{\sum_{i=1}^n (x_i - \bar{x})^2}}, \hat{a} + t_{n-2; \frac{1+\gamma}{2}} \sqrt{\frac{\hat{\sigma}^2}{\sum_{i=1}^n (x_i - \bar{x})^2}} \right]$$

for a and (in the following formulae, “ t ” always means “ $t_{n-2; \frac{1+\gamma}{2}}$ ”)

$$\left[\hat{b} - t \sqrt{\hat{\sigma}^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right)}, \hat{b} + t \sqrt{\hat{\sigma}^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right)} \right]$$

for b . In many cases, one wants a confidence interval for $y(x) = ax + b$:

$$\left[\hat{y}(x) - t \sqrt{\hat{\sigma}^2 \left(\frac{1}{n} + \frac{(x - \bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right)}, \hat{y}(x) + t \sqrt{\hat{\sigma}^2 \left(\frac{1}{n} + \frac{(x - \bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right)} \right].$$

Sometimes, one wants to find an interval that contains the value $Y(x) = ax + b + e$ with a given probability γ . This prediction interval is given by

$$\left[\hat{y}(x) - t \sqrt{\hat{\sigma}^2 \left(1 + \frac{1}{n} + \frac{(x - \bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right)}, \hat{y}(x) + t \sqrt{\hat{\sigma}^2 \left(1 + \frac{1}{n} + \frac{(x - \bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right)} \right].$$

The conceptual difference between the two types of intervals is as follows: the simple confidence interval just has to cover the value of the function $ax + b$; the prediction interval, on the other hand, covers the value $ax + b + e$, and the additional error term is the reason that the prediction interval is wider. In other words, the confidence interval is for the “theoretical” value (which, in many cases, cannot be observed), the prediction interval for the actual measurement (which is the one we usually can observe, but only at the cost of an additional measurement error).

5.2 Multiple linear regression.

Now, we consider functions of the form

$$y = a_1x_1 + \dots + a_kx_k.$$

(There is no b in this formula, but it can be treated by letting $x_i = 1$ for some i . We fix n sets of numbers (x_{1l}, \dots, x_{kl}) ($l = 1, \dots, n$) and measure the related

$$Y_l = a_1x_{1l} + \dots + a_kx_{kl} + e_l.$$

With the usual assumptions on our error terms e_i , we get the equations

$$\sum_{i=1}^k a_i b_{ij} = c_j \quad (j = 1, \dots, k),$$

where

$$b_{ij} = \sum_{l=1}^n x_{il}x_{jl}$$

and

$$c_j = \sum_{l=1}^n Y_l x_{jl}.$$

Again, \hat{a}_i is an unbiased estimator of a_i , and the random variables \hat{a}_i have a joint normal distribution with covariance matrix $\sigma^2 B^{-1}$, where B is the $(k \times k)$ -matrix with entries b_{ij} .

Again, we have the estimator

$$\hat{\sigma}^2 = \frac{1}{n-k} \sum_{i=1}^n \left(Y_i - \sum_{j=1}^k \hat{a}_j x_{ji} \right)^2$$

for $\hat{\sigma}^2$ which is independent of $\hat{a}_1, \dots, \hat{a}_k$.

Thus, we arrive at the confidence intervals

$$\left[\hat{a}_i - t_{n-k, \frac{1+\gamma}{2}} \sqrt{\hat{\sigma}^2 (B^{-1})_{ii}}, \hat{a}_i + t_{n-k, \frac{1+\gamma}{2}} \sqrt{\hat{\sigma}^2 (B^{-1})_{ii}} \right].$$

For

$$y(x) = a_1 x_1 + \dots + a_k x_k$$

we get the estimator

$$\hat{y}(x) = \hat{a}_1 x_1 + \dots + \hat{a}_k x_k$$

and the confidence interval

$$\left[\hat{y}(x) - t_{n-k, \frac{1+\gamma}{2}} \sqrt{\hat{\sigma}^2 x^T B^{-1} x}, \hat{y}(x) + t_{n-k, \frac{1+\gamma}{2}} \sqrt{\hat{\sigma}^2 x^T B^{-1} x} \right],$$

where x denotes the column vector with entries x_i .

We also have a prediction interval:

$$\left[\hat{y}(x) - t_{n-k, \frac{1+\gamma}{2}} \sqrt{\hat{\sigma}^2 (1 + x^T B^{-1} x)}, \hat{y}(x) + t_{n-k, \frac{1+\gamma}{2}} \sqrt{\hat{\sigma}^2 (1 + x^T B^{-1} x)} \right],$$

5.3 Other functional relations.

In practice, there is often the question of fitting a nonlinear function to our data. Especially important are the exponential and power functions ($y = AB^x$ or $y = Ax^B$), and the polynomials ($y = a_0 + a_1 x + \dots + a_k x^k$).

The polynomials can easily be treated using the methods of the previous section, by letting $x_{il} = x_i^{l-1}$. The other two types of functions present more problems; actually, we can write down the equations for the least squares estimators, but those cannot be solved in closed form (they can be solved numerically, of course); so we rather take logarithms on both sides to convert this into simple linear regression. We have to be aware, however, that by doing so we have changed our stochastic model. Now the error term does no longer come in as an additive term, but it is added in the exponent. Anyway, this enables us to calculate confidence intervals and prediction intervals, which are not easily obtained by other methods.

5.4 Problems

All together now

The Beatles

1. Suppose that n is even and that we are free to choose x_1, \dots, x_n from the interval $[-A, A]$. How do we have to choose x_1, \dots, x_n in order to estimate

(a) a

(b) b

as accurately as possible?

2. As before, for n odd.
3. Sometimes the variance of the error term is a function $\sigma^2(x)$ of x . Show that in this case the ML estimators for a and b can be obtained by using ordinary linear regression for

$$\frac{y}{\sigma(x)} = a \frac{x}{\sigma(x)} + b \frac{1}{\sigma(x)}.$$

4. For the following data

i	1	2	3	4	5
x	0.5	1.0	2.0	2.5	3.0
Y	1.3	2.7	1.5	2.0	2.2

fit a linear regression function and calculate a prediction interval for $x = 3.5$.

5. In problem 4 , use a regression function of the form $y = ae^{bx}$.
6. In problem 4 , use a regression function of the form $y = \frac{a}{b+x}$.
7. In problem 4 , use a quadratic polynomial as the regression function.
8. In problem 4 , calculate a confidence interval for a .
9. In problem 4 , calculate a confidence interval for b .
10. In problem 4 , calculate a confidence interval for $y(1.5)$.

Kapitel 6

The χ^2 -Family of Tests.

Whereupon the Plumber said in tones of disgust:

“I suggest that we proceed at once to infinity.”

J.L. Synge

6.1 The multinomial distribution and the χ^2 -statistic

In many cases we have to deal with testing whether a certain discrete random variable has a given distribution.

In the simplest setting, assume that X can take k values — without restricting generality, we may assume that these are $1, \dots, k$ — with probabilities p_1, \dots, p_k , and we want to test $H_0 : p_i = p_{i0}$ ($i = 1, \dots, k$) against the alternative that at least one p_i is different from p_{i0} . The standard way to approach this problem is to use the likelihood ratio test, of course, and the likelihood ratio statistic turns out to be

$$\ell = \prod_{i=1}^k \frac{p_{i0}^{Y_i}}{(Y_i/n)^{Y_i}},$$

where Y_i is the number of occurrences of i among X_1, \dots, X_n .

In order to calculate the critical value for this test, we need the distribution of Y_1, \dots, Y_k under the null hypothesis. This is a multinomial distribution:

$$\mathbb{P}(Y_1 = y_1, \dots, Y_k = y_k) = \begin{cases} \frac{n!}{y_1! \dots y_k!} p_{1,0}^{y_1} \dots p_{k,0}^{y_k} & \text{if } y_1 + \dots + y_k = n, \\ 0 & \text{otherwise.} \end{cases}$$

For the sake of simplicity, we drop the subscript “0” from p_{i0} in the following calculations.

Now, at least in theory, we are able to use the multinomial distribution to evaluate the null distribution of ℓ . In general, however, this is too tedious to be carried out in practice. So, we rather assume that n is large and try to use some approximation. We put

$$Y_i = np_i + U_i \sqrt{n},$$

where U_i has mean zero and variance $p_i(1 - p_i)$ and is approximately normal. Furthermore, using the expansion

$$\log(1 + x) = x - \frac{x^2}{2} + o(x^2)$$

we get

$$\log \ell = -\frac{1}{2} \sum_{i=1}^k \frac{U_i^2}{p_i} + o(1) = -\frac{1}{2} \sum_{i=1}^k \frac{(Y_i - np_i)^2}{np_i} + o(1).$$

Thus, for large n , the likelihood ratio test is approximately equivalent to the test that rejects H_0 if

$$\chi^2 = \sum_{i=1}^k \frac{(Y_i - np_i)^2}{np_i}$$

is large.

What we need now is a result on the distribution of χ^2 . First, let us consider the joint distribution of U_1, \dots, U_k . Approximately, they have normal distributions with mean zero; the variances and covariances are

$$\mathbb{V}(U_i) = p_i(1 - p_i)$$

and

$$\mathbf{Cov}(U_i, U_j) = -p_i p_j \quad (i \neq j).$$

Let U_0 be another normal random variable with mean zero and variance one, independent from U_1, \dots, U_k , and let

$$V_i = U_i + p_i U_0.$$

It is readily verified that

$$\mathbb{V}(V_i) = p_i$$

and

$$\mathbf{Cov}(V_i, V_j) = 0.$$

So (at least as an approximation), the random variables V_i are independent normal random variables. Furthermore, we have the decomposition

$$\sum_{i=1}^k \frac{V_i^2}{p_i} = U_0^2 + \sum_{i=1}^k \frac{U_i^2}{p_i},$$

(In the last equation, we use the fact that $\sum U_i = 0$) where the left-hand side has k degrees of freedom, and the two summands on the right-hand side have one and $k - 1$ degrees of freedom, respectively. This, by virtue of the Fisher-Cochran theorem, permits us to state that

$$\chi^2 = \sum_{i=1}^k \frac{(Y_i - np_i)^2}{np_i}$$

has an approximate χ^2 -distribution with $k - 1$ degrees of freedom.

Now, let us turn to the case where H_0 is composite. Assume that H_0 states that p_1, \dots, p_k is of the form $p_i = p_i(\theta_1, \dots, \theta_s)$, where $\theta_1, \dots, \theta_s$ are real parameters. It can be shown, using Taylor expansions and the ideas from analysis of variance, that, if we use the ML estimators $\hat{\theta}_i$,

$$\chi^2 = \sum_{i=1}^k \frac{(Y_i - np_i(\hat{\theta}_1, \dots, \hat{\theta}_s))^2}{np_i(\hat{\theta}_1, \dots, \hat{\theta}_s)}$$

has an approximate χ^2 -distribution with $k - 1 - s$ degrees of freedom (assuming that the parameters are all independent; otherwise we have to use the number of degrees of freedom of the set of parameters instead, which again is the number of parameters that can be independently chosen). Be advised, however, that this result depends on some regularity conditions, similar to the ones in the Cramér-Rao theorem (which in our case amount to saying that p_i is sufficiently often differentiable with respect to θ). These conditions may be violated even in simple cases (the uniform distribution is one example). Anyway, it is very useful and very widely used, even in cases where it is not completely justified.

As the χ^2 -distribution is only an asymptotic one, we need some rule for its application. The usual rule of thumb is that np_i should be at least 5 for each i .

6.2 The χ^2 goodness-of-fit test.

There is not much new in this section. We now want to test $H_0 : X$ has a given distribution. If this distribution is discrete with a finite number k of possible values, we can directly use the first result from the previous section. This means, we calculate

$$\chi^2 = \sum_{i=1}^k \frac{(Y_i - np_i)^2}{np_i},$$

and reject if this exceeds $\chi_{k-1, 1-\alpha}^2$. If the distribution has infinitely many possible values (in particular if it is continuous), or if some of the numbers np_i are less than 5, we can divide the set of possible values into classes, and let Y_i be the number of occurrences of values from class number i among X_1, \dots, X_n .

In some cases, the distribution of X is not completely specified (the most common example is the normal distribution — often one wants to know whether the data are normally distributed, but doesn't know mean and variance). For that case, too, the previous section shows the way: we estimate the parameters by the maximum likelihood method (although this is not exactly what was done in the previous section) and use these to calculate the probabilities for each class (to me, it seems convenient to choose the classes in such a way that the probabilities p_i are all equal; in the case of the normal distribution, this means that the class boundaries are $\bar{X} + z_{i/k} \sqrt{S_n^2}$, ($i = 1, \dots, k$)). Finally, we calculate the χ^2 statistic in the usual way, and use $k - 1 -$ number of estimated parameters as the number of degrees of freedom.

6.3 The χ^2 test of independence.

In this section, we look at paired observations (or two-dimensional ones, if you prefer) (X_i, Y_i) , $i = 1, \dots, n$ and we want to test whether X and Y are independent. We assume that both X and Y have a discrete distribution (otherwise, we have to use a suitable classification); without restricting generality, we may assume that X takes values 1 to k , and that Y takes values 1 to m . If X and Y are in fact independent, then for $1 \leq j \leq k$ and $1 \leq l \leq m$, letting

$$p_j = \mathbb{P}(X = j)$$

and

$$q_l = \mathbb{P}(Y = l)$$

we have

$$\mathbb{P}(X = j, Y = l) = \mathbb{P}(X = j)\mathbb{P}(Y = l) = p_j q_l.$$

Thus, we have expressed the joint distribution of X and Y by the $k + m$ parameters p_1, \dots, p_k and q_1, \dots, q_m . This almost cries out for an application of the χ^2 test. If we let

$$Z_{jl} = \#\{i : X_i = j, Y_i = l\},$$

$$Z_{.l} = \#\{i : Y_i = l\} = \sum_{j=1}^k Z_{jl},$$

and

$$Z_{j.} = \#\{i : X_i = j\} = \sum_{l=1}^m Z_{jl},$$

then the ML estimators for p_j and q_l are $Z_{j.}/n$ and $Z_{.l}/n$, respectively. The χ^2 statistic is

$$\sum_{j=1}^k \sum_{l=1}^m \frac{(Z_{jl} - Z_{j.}Z_{.l}/n)^2}{Z_{j.}Z_{.l}/n}.$$

We have to find the number of degrees of freedom of this. The number of classes, of course, is km . We had to estimate the $k + m$ parameters, but only $k + m - 2$ are independent because we have the relations $\sum p_j = \sum q_l = 1$. Thus, the number of degrees of freedom of our χ^2 statistic is

$$km - 1 - (k + m - 2) = (k - 1)(m - 1).$$

6.4 The χ^2 test for homogeneity.

This time, we have k samples $(X_{i1}, \dots, X_{in_i}), i = 1, \dots, k$. We assume that X_{ij} can take values from 1 to m and want to check if the underlying distributions are the same for all samples. We let

$$n = \sum_{i=1}^k n_i,$$

$$Y_{il} = \#\{j : X_{ij} = l\},$$

$$Z_l = \#\{(i, j) : X_{ij} = l\}.$$

under the null hypothesis that all samples have the same distribution $\mathbb{P}(X_{ij} = l) = p_l$, the ML estimator for p_l is Z_l/n , and we arrive at the χ^2 statistic

$$\sum_{i,l} \frac{(Y_{il} - n_i Z_l/n)^2}{n_i Z_l/n},$$

which has

$$(m-1)(k-1)$$

degrees of freedom. (It may be noted that this is very similar to the test for independence; in fact, it is only a different way of interpreting the same data)

6.5 Problems

All my trials, Lord
soon be over

trad.

1. A die is rolled 600 times showing the following results:

1	2	3	4	5	6
95	111	89	106	88	113

Test whether the die is fair.

2. Test whether the following numbers have a normal distribution with mean 1 and variance 0.5:

1.95 2.13 0.74 0.05 2.70 1.84 1.25 1.00 0.70 1.02

1.71 0.58 0.16 1.12 1.16 0.21 0.70 1.30 1.06 1.75

3. Do the previous problem for unknown mean and variance.
4. Test whether the data from problem 2 have an exponential distribution.
5. Test whether the data from problem 2 have a uniform distribution on $[0, \theta]$.

6. Test whether the following numbers have a Poisson distribution:

0 1 2 1 3 4 3 3 2 4 3 0 3 5 3 2 3 1 2 4 3 3 3 5 2 4 1 0 3 2 2 3 2 3 4 5 6 1 0 6

7. In the previous problem, test whether the data have a binomial distribution with $n = 6$ and $p = 1/2$.
8. Do the previous problem for unknown p .

9. If X_1, \dots, X_n are uniformly distributed on $[0, \theta]$, let Y_1, \dots, Y_{n-1} be X_1, \dots, X_n without the maximum $M = \max_{i \leq n} X_i$. Furthermore, let $Z_i = Y_i/M$. Show that Z_1, \dots, Z_{n-1} is a sample from a uniform distribution on $[0, 1]$. Use this fact to show that the χ^2 -statistic in problem 5 has an approximate χ^2 -distribution with $k-1$ (and not $k-2$!) degrees of freedom (this shows that some rules should not be followed blindly - in this case, we see another anomaly of the uniform distribution that is caused by the fact that the densities $f(x, \theta)$ are not differentiable with respect to θ . The validity of the “subtract the number of parameters estimated from the number of degrees of freedom” depends on the differentiability of the densities, and on the fact that the ML estimator can be obtained by zeroing the derivative.

Now, what if a problem of this kind comes up at a written exam? There are two possibilities: first, just follow the rule, and you’re safe (after all, the exams are there to check if you know the rules; besides, I believe that two out of three lecturers in the field are themselves unaware of this particular problem); second, you might use your superior knowledge, but state *very clearly* that you know the rule, but deliberately break it in *this special case* because it is the more sensible thing to do).

10. A sample of 500 students showed the following marks in mathematics and english:

english	mathematics				
	1	2	3	4	5
1	40	10	0	0	0
2	25	30	15	20	10
3	40	35	45	30	50
4	20	10	20	30	20
5	0	0	10	20	20

Are the marks in the two subjects independent?

11. For $m = k = 2$, give a simple (?) formula for the χ^2 statistic in the test for independence.
12. 460 bolts are classified as “conforming” or “non-conforming” with respect to their measurements and their breaking strength:

meas.	breaking strength	
	conf.	non-c.
conf.	416	23
non-c.	16	5

Are the two classifications independent?

13. Three professors evaluate the same 100 tests with the following results

mark	prof. 1	prof. 2	prof. 3
1	35	30	28
2	22	24	20
3	16	20	22
4	12	20	17
5	15	6	13

Are the three professors equally strict?

Kapitel 7

The Kolmogorov-Smirnov Test

a discussion of order ... has become essential to any understanding of the foundation of mathematics

B.Russell

7.1 The one-sample test

This is another goodness of fit test, i.e., we want to test H_0 : the distribution function of X is F against the alternative that it is different from F . Our point of departure is the Glivenko-Cantelli theorem which states that the empirical distribution function

$$F_n(x) = \frac{1}{n} \#\{j : X_j < x\}$$

converges uniformly to the actual distribution function of X . Thus, we may choose

$$D_n = \|F_n - F\| = \sup_x |F_n(x) - F(x)|$$

as our test statistic. H_0 will be rejected if D_n is greater than a critical value. The following theorem shows that — at least for continuous F — this critical value can be determined independently from F .

Theorem 7.1

If the distribution function F is continuous, then the null distribution of D_n does not depend on F .

For the proof of this theorem, we need the following

Lemma 7.1

1. If the distribution function F of X is continuous, then $F(X)$ is uniformly distributed on $[0, 1]$.
2. For any distribution function F and U uniform on $[0, 1]$, $X = F^{-1}(U)$ has distribution function F , where

$$F^{-1}(x) = \sup\{t : F(t) < x\}.$$

We don't need the second assertion, so we only prove the first one.

By the intermediate value theorem, for any $0 < y < 1$, there is an x such that $F(x) = y$. Thus

$$\mathbb{P}(F(X) \leq y) \geq \mathbb{P}(X \leq x) = F(x) = y$$

and

$$\mathbb{P}(F(X) < y) \leq \mathbb{P}(X < x) = F(x - 0) = y$$

This shows that $F(X)$ is uniformly distributed on $[0, 1]$.

Now, let us turn to the proof of the theorem. We first calculate a formula for D_n which is a little more convenient. To this end, let

$$X_{n:1} < X_{n:2} < \dots < X_{n:n}$$

be the order statistics (i.e., the sample X_1, \dots, X_n in ascending order). Then for $0 \leq k \leq n$ $F_n(x) = \frac{k}{n}$ if $X_{n:k} < x < X_{n:k+1}$ (let $X_{n:0} = -\infty$, $X_{n:n+1} = +\infty$). Thus,

$$\begin{aligned} D_n &= \max_k \sup_{X_{n:k} < x < X_{n:k+1}} |F_n(x) - F(x)| = \\ &= \max_k \sup_{X_{n:k} < x < X_{n:k+1}} \left| \frac{k}{n} - F(x) \right| = \\ &= \max_k \max \left(\left| \frac{k}{n} - F(X_{n:k}) \right|, \left| \frac{k}{n} - F(X_{n:k+1}) \right| \right) = \\ &= \max_k \max \left(\left| \frac{k}{n} - F(X_{n:k}) \right|, \left| \frac{k-1}{n} - F(X_{n:k}) \right| \right). \end{aligned}$$

If we let $U_i = F(X_i)$, then U_1, \dots, U_n will be a sample from a uniform distribution on $[0, 1]$, and $U_{n:k} = F(X_{n:k})$. Thus, D_n is the same whether we calculate it from X or from U . This proves the theorem.

Thus, one can calculate the finite sample distribution of D_n , although the calculations may get quite intricate. For large n , there is an asymptotic result:

Theorem 7.2

Let $\lambda_n = \sqrt{n}D_n$. For $n \rightarrow \infty$, the distribution function of λ_n is approximately

$$K(x) = \begin{cases} \sum_{n=-\infty}^{\infty} (-1)^n e^{-2n^2 x^2} & \text{if } x > 0 \\ 0 & \text{if } x \leq 0. \end{cases}$$

The proof of this theorem is beyond the scope of this lecture.

7.2 The Two-Sample Test

Now, we are going to compare two sample X_1, \dots, X_n and Y_1, \dots, Y_m . We want to test the hypothesis that X and Y have the same distribution. Using the empirical distribution functions F_n and G_m , we define our test statistic as

$$D_{n,m} = \|F_n - G_m\|.$$

If we assume that the common distribution function F of X and Y is continuous, one can prove in the same way as before that the distribution of $D_{n,m}$ does not depend on F . In addition, there is the asymptotic result:

Theorem 7.3

$$\lambda_{n,m} = \sqrt{\frac{nm}{n+m}} D_{n,m}$$

has a limiting distribution with distribution function K .

7.3 Problems

To everything - turn, turn, turn,
there is a season - turn, turn, turn,
and a time for ev'ry purpose under
heaven.

Book of ecclesiastes, adap. P. Seeger

1. Use the Kolmogorov-Smirnov test for problem 2, chapter 6.
2. Test whether the following data are uniformly distributed on $[0, 1]$.
0.23 0.31 0.45 0.67 0.87 0.49 0.93 0.80
3. Use the two-sample version of the Kolmogorov-Smirnov test for problem 20, chapter 3.
4. Show that

$$\begin{aligned}\mathbb{P}(D_n < x) &= \mathbb{P}\left(\frac{k}{n} - x < U_{n:k} < \frac{k-1}{n} + x\right) = \\ &= n! \mathbb{P}\left(\frac{k}{n} - x < U_k < \frac{k-1}{n} + x, U_1 < U_2 < \dots < U_n\right).\end{aligned}$$

5. Use the previous problem to calculate the distribution of D_2 .
6. If both X_i and Y_j are uniformly distributed on $[0, 1]$, calculate the Variance of $\eta_{m,n}(x) = F_n(x) - G_m(x)$ and the covariance of $\eta_{m,n}(x)$ and $\eta_{m,n}(y)$ (This should make it plausible that the normalizing factor in our limit theorems is right).

Anhang A

Tables

The distribution function of the normal distribution:

$$\Phi(x) = \int_{-\infty}^x \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}}$$

	0	1	2	3	4	5	6	7	8	9
0.0	.500	.504	.508	.512	.516	.520	.524	.528	.532	.536
0.1	.540	.544	.548	.552	.556	.560	.564	.567	.571	.575
0.2	.579	.583	.587	.591	.595	.599	.603	.606	.610	.614
0.3	.618	.622	.626	.629	.633	.637	.641	.644	.648	.652
0.4	.655	.659	.663	.666	.670	.674	.677	.681	.684	.688
0.5	.691	.695	.698	.702	.705	.709	.712	.716	.719	.722
0.6	.726	.729	.732	.736	.739	.742	.745	.749	.752	.755
0.7	.758	.761	.764	.767	.770	.773	.776	.779	.782	.785
0.8	.788	.791	.794	.797	.800	.802	.805	.808	.811	.813
0.9	.816	.819	.821	.824	.826	.829	.831	.834	.836	.839
1.0	.841	.844	.846	.848	.851	.853	.855	.858	.860	.862
1.1	.864	.867	.869	.871	.873	.875	.877	.879	.881	.883
1.2	.885	.887	.889	.891	.893	.894	.896	.898	.900	.901
1.3	.903	.905	.907	.908	.910	.911	.913	.915	.916	.918
1.4	.919	.921	.922	.924	.925	.926	.928	.929	.931	.932
1.5	.933	.934	.936	.937	.938	.939	.941	.942	.943	.944
1.6	.945	.946	.947	.948	.949	.951	.952	.953	.954	.954
1.7	.955	.956	.957	.958	.959	.960	.961	.962	.962	.963
1.8	.964	.965	.966	.966	.967	.968	.969	.969	.970	.971
1.9	.971	.972	.973	.973	.974	.974	.975	.976	.976	.977
2.0	.977	.978	.978	.979	.979	.980	.980	.981	.981	.982
2.1	.982	.983	.983	.983	.984	.984	.985	.985	.985	.986
2.2	.986	.986	.987	.987	.987	.988	.988	.988	.989	.989
2.3	.989	.990	.990	.990	.990	.991	.991	.991	.991	.992
2.4	.992	.992	.992	.992	.993	.993	.993	.993	.993	.994
2.5	.994	.994	.994	.994	.994	.995	.995	.995	.995	.995
2.6	.995	.995	.996	.996	.996	.996	.996	.996	.996	.996
2.7	.997	.997	.997	.997	.997	.997	.997	.997	.997	.997
2.8	.997	.998	.998	.998	.998	.998	.998	.998	.998	.998
2.9	.998	.998	.998	.998	.998	.998	.998	.999	.999	.999

Quantiles z_p of the normal distribution:

p	z_p	p	z_p	p	z_p
.51	.025	.71	.553	.91	1.341
.52	.050	.72	.583	.92	1.405
.53	.075	.73	.613	.93	1.476
.54	.100	.74	.643	.94	1.555
.55	.126	.75	.674	.95	1.645
.56	.151	.76	.706	.96	1.751
.57	.176	.77	.739	.97	1.881
.58	.202	.78	.772	.975	1.960
.59	.228	.79	.806	.98	2.054
.60	.253	.80	.842	.99	2.326
.61	.279	.81	.878	.991	2.366
.62	.305	.82	.915	.992	2.409
.63	.332	.83	.954	.993	2.457
.64	.358	.84	.994	.994	2.512
.65	.385	.85	1.036	.995	2.576
.66	.412	.86	1.080	.996	2.652
.67	.440	.87	1.126	.997	2.748
.68	.468	.88	1.175	.998	2.878
.69	.496	.89	1.227	.999	3.090
.70	.524	.90	1.282	.9999	3.719

Quantiles $t_{n;p}$ of Student's t -distribution with n degrees of freedom:

n	.9	.95	.975	.99	.995	n	.9	.95	.975	.99	.995
1	3.078	6.314	12.706	31.821	63.675	26	1.316	1.706	2.056	2.479	2.779
2	1.886	2.920	4.303	6.965	9.725	27	1.314	1.703	2.052	2.473	2.467
3	1.638	2.353	3.183	4.541	5.841	28	1.313	1.701	2.048	2.467	2.763
4	1.533	2.132	2.776	3.747	4.604	29	1.311	1.699	2.045	2.462	2.756
5	1.476	2.015	2.571	3.365	4.032	30	1.310	1.697	2.042	2.457	2.750
6	1.440	1.943	2.447	3.143	3.707	31	1.309	1.696	2.040	2.453	2.744
7	1.415	1.895	2.365	2.998	3.499	32	1.309	1.694	2.037	2.449	2.738
8	1.397	1.860	2.306	2.896	3.355	33	1.308	1.692	2.035	2.445	2.733
9	1.383	1.833	2.262	2.821	3.250	34	1.307	1.691	2.032	2.441	2.728
10	1.372	1.812	2.228	2.764	3.169	35	1.306	1.690	2.030	2.438	2.724
11	1.363	1.796	2.201	2.718	3.106	40	1.303	1.684	2.021	2.423	2.704
12	1.356	1.782	2.179	2.681	3.055	45	1.301	1.679	2.014	2.412	2.690
13	1.350	1.771	2.160	2.650	3.012	50	1.299	1.676	2.009	2.403	2.678
14	1.345	1.761	2.145	2.624	2.977	55	1.297	1.673	2.004	2.396	2.668
15	1.341	1.753	2.131	2.602	2.947	60	1.296	1.671	2.000	2.390	2.660
16	1.337	1.746	2.120	2.583	2.921	65	1.295	1.669	1.997	2.385	2.654
17	1.333	1.740	2.110	2.567	2.898	70	1.294	1.667	1.994	2.381	2.648
18	1.330	1.734	2.101	2.552	2.878	75	1.293	1.665	1.992	2.377	2.643
19	1.328	1.729	2.093	2.539	2.861	80	1.292	1.664	1.990	2.374	2.639
20	1.325	1.725	2.086	2.528	2.845	85	1.292	1.663	1.988	2.371	2.635
21	1.323	1.721	2.080	2.518	2.831	90	1.291	1.662	1.987	2.368	2.632
22	1.321	1.717	2.074	2.508	2.819	95	1.291	1.661	1.985	2.366	2.629
23	1.319	1.714	2.069	2.500	2.807	100	1.290	1.660	1.984	2.364	2.626
24	1.318	1.711	2.064	2.492	2.797	105	1.290	1.659	1.983	2.362	2.623
25	1.316	1.708	2.060	2.485	2.787	∞	1.282	1.645	1.960	2.326	2.576

Quantiles $\chi^2_{n;p}$ of the χ^2 -distribution with n degrees of freedom

n	.005	.01	.02	.025	.05	.1	.5	.9	.95	.975	.98	.99	.995
1	.000	.000	.001	.001	.004	.016	.455	2.706	3.841	5.024	5.412	6.635	7.879
2	.010	.020	.040	.051	.103	.211	1.386	4.605	5.991	7.378	7.824	9.210	10.597
3	.072	.115	.185	.216	.352	.584	2.366	6.251	7.815	9.348	9.837	11.345	12.838
4	.207	.297	.429	.484	.711	1.064	3.357	7.779	9.488	11.143	11.668	13.277	14.860
5	.412	.554	.752	.831	1.145	1.610	4.351	9.236	11.070	12.832	13.308	15.086	16.750
6	.676	.872	1.134	1.237	1.635	2.204	5.348	10.645	12.592	14.449	15.033	16.812	18.548
7	.989	1.239	1.564	1.690	2.167	2.833	6.346	12.017	14.067	16.013	16.622	18.475	20.278
8	1.344	1.646	2.032	2.180	2.733	3.490	7.344	13.362	15.507	17.535	18.168	20.090	21.955
9	1.735	2.088	2.532	2.700	3.325	4.168	8.343	14.684	16.919	19.023	19.679	21.666	23.589
10	2.156	2.558	3.059	3.247	3.940	4.865	9.342	15.987	18.307	20.483	21.161	23.209	25.188
11	2.603	3.053	3.609	3.816	4.575	5.578	10.341	17.275	19.675	21.920	22.618	24.725	26.757
12	3.074	3.571	4.178	4.404	5.226	6.304	11.340	18.549	21.026	23.336	24.054	26.217	28.300
13	3.565	4.107	4.765	5.009	5.892	7.042	12.340	19.812	22.362	24.736	25.472	27.688	29.819
14	4.075	4.660	5.368	5.629	6.571	7.790	13.339	21.064	23.685	26.119	26.873	29.141	31.319
15	4.601	5.229	5.985	6.262	7.261	8.547	14.339	22.307	24.996	27.488	28.259	30.578	32.801
16	5.142	5.812	6.614	6.908	7.962	9.312	15.338	23.542	26.269	28.845	29.633	32.000	34.267
17	5.697	6.408	7.255	7.564	8.672	10.085	16.338	24.769	27.587	30.191	30.995	33.409	35.718
18	6.265	7.015	7.906	8.231	9.390	10.835	17.338	25.909	28.869	31.526	32.346	34.805	37.156
19	6.844	7.633	8.567	8.907	10.117	11.651	18.338	27.204	30.144	32.852	33.687	36.191	38.582
20	7.434	8.260	9.237	9.591	10.851	12.443	19.337	28.412	31.410	34.170	35.020	37.566	39.997
21	8.034	8.897	9.915	10.283	11.591	13.240	20.337	29.615	32.671	35.479	36.343	38.932	41.401
22	8.643	9.542	10.600	10.982	12.338	14.041	21.337	30.813	33.924	36.781	37.659	40.289	42.796
23	9.260	10.196	11.293	11.689	13.091	14.848	22.337	32.007	35.172	38.076	38.968	41.638	44.181
24	9.886	10.856	11.992	12.401	13.848	15.659	23.337	33.196	36.415	39.364	40.270	42.980	45.559
25	10.520	11.524	12.697	13.120	14.611	16.473	24.337	34.382	37.652	40.646	41.566	44.324	46.928
26	11.160	12.198	13.409	13.844	15.379	17.292	25.336	35.563	38.885	41.923	42.856	45.642	48.290
27	11.808	12.879	14.125	14.573	16.151	18.114	26.336	36.741	40.113	43.194	44.140	46.963	49.645
28	12.461	13.565	14.847	15.308	16.928	18.939	27.336	37.916	41.337	44.461	45.419	48.278	50.993
29	13.121	14.256	15.574	16.047	17.708	19.768	28.336	39.087	42.557	45.722	46.693	49.588	52.336
30	13.787	14.953	16.306	16.791	18.493	20.599	29.336	40.256	43.773	46.979	47.962	50.892	53.672
31	14.458	15.655	17.042	17.539	19.281	21.434	30.336	41.422	44.985	48.232	49.226	52.191	55.003
32	15.134	16.362	17.783	18.291	20.072	22.271	31.336	42.585	46.194	49.480	50.487	53.486	56.328
33	15.815	17.074	18.527	19.047	20.867	23.110	32.336	43.745	47.400	50.725	51.743	54.776	57.648
34	16.501	17.789	19.275	19.806	21.664	23.952	33.336	44.903	48.602	51.966	52.995	56.061	58.964
35	17.192	18.509	20.027	20.569	22.465	24.797	34.336	46.059	49.802	53.203	54.244	57.342	60.275
40	20.707	22.164	23.838	24.433	26.509	29.051	39.335	51.805	55.758	59.342	60.436	63.691	66.766
45	24.311	25.901	27.720	28.366	30.612	33.350	44.335	57.505	61.656	65.410	66.555	69.957	73.166
50	27.991	29.707	31.664	32.357	34.764	37.689	49.335	63.167	67.505	71.420	72.613	76.154	79.490
55	31.735	33.570	35.659	36.398	38.958	42.060	54.335	68.796	73.311	77.380	78.619	82.292	85.749
60	35.534	37.485	39.699	40.482	43.188	46.459	59.335	74.397	79.082	83.298	84.580	88.397	91.952
65	39.383	41.444	43.779	44.603	47.450	50.883	64.335	79.973	84.821	89.177	90.501	94.422	98.105
70	43.275	45.442	47.893	48.758	51.739	55.329	69.334	85.527	90.531	95.023	96.388	100.425	104.215
75	47.206	49.475	52.039	52.942	56.054	59.795	74.334	91.061	96.217	100.839	102.243	106.393	110.286
80	51.172	53.540	56.213	57.153	60.391	64.278	79.334	96.578	101.879	106.629	108.069	112.329	116.321
85	55.170	57.634	60.412	61.389	64.749	68.777	84.334	102.079	107.522	112.393	113.871	118.236	122.325
90	59.196	61.754	64.635	65.647	69.126	73.291	89.334	107.565	113.145	118.136	119.649	124.116	128.299
95	63.250	65.898	68.879	69.925	73.520	77.818	94.334	113.038	118.752	123.858	125.405	129.973	134.247
100	67.328	70.065	73.142	74.222	77.929	82.358	99.334	118.498	124.342	129.561	131.142	135.806	140.169

Quantiles $F_{a,b;0.95}$ of the F -distribution with a and b degrees of freedom

b	a																	
	1	2	3	4	5	6	7	8	9	10	15	20	30	40	50	100	∞	
1	161	200	216	225	230	234	237	239	241	242	246	248	250	251	252	253	254	
2	18.51	19.00	19.20	19.20	19.30	19.30	19.40	19.40	19.40	19.40	19.40	19.40	19.50	19.50	19.50	19.50	19.50	
3	10.13	9.55	9.28	9.12	9.01	8.94	8.89	8.85	8.81	8.79	8.69	8.66	8.62	8.59	8.58	8.55	8.55	
4	7.71	6.94	6.59	6.39	6.26	6.16	6.09	6.04	6.00	5.96	5.86	5.80	5.75	5.72	5.70	5.66	5.66	
5	6.61	5.79	5.41	5.19	5.05	4.95	4.88	4.82	4.77	4.74	4.62	4.56	4.50	4.46	4.44	4.41	4.41	
6	5.99	5.14	4.76	4.53	4.39	4.28	4.21	4.15	4.10	4.06	3.94	3.87	3.81	3.77	3.75	3.71	3.66	
7	5.59	4.74	4.35	4.12	3.97	3.87	3.79	3.73	3.68	3.64	3.51	3.44	3.38	3.34	3.32	3.27	3.27	
8	5.32	4.46	4.07	3.84	3.69	3.58	3.50	3.44	3.39	3.35	3.22	3.15	3.08	3.04	3.02	2.97	2.97	
9	5.12	4.26	3.86	3.63	3.48	3.37	3.29	3.23	3.18	3.14	3.01	2.94	2.86	2.83	2.80	2.76	2.77	
10	4.96	4.10	3.71	3.48	3.33	3.22	3.14	3.07	3.02	2.98	2.85	2.77	2.70	2.66	2.64	2.59	2.55	
12	4.75	3.89	3.49	3.26	3.11	3.00	2.91	2.85	2.80	2.75	2.62	2.54	2.47	2.43	2.40	2.35	2.33	
14	4.61	3.75	3.34	3.11	2.96	2.85	2.76	2.70	2.65	2.60	2.46	2.39	2.31	2.27	2.24	2.19	2.19	
16	4.49	3.63	3.22	2.99	2.85	2.74	2.66	2.59	2.54	2.49	2.35	2.28	2.19	2.15	2.12	2.07	2.07	
18	4.41	3.55	3.14	2.91	2.77	2.66	2.58	2.51	2.46	2.41	2.27	2.19	2.11	2.06	2.04	1.98	1.99	
20	4.35	3.49	3.10	2.87	2.71	2.60	2.51	2.45	2.39	2.35	2.20	2.12	2.04	1.99	1.97	1.91	1.88	
25	4.24	3.39	2.99	2.76	2.60	2.49	2.40	2.34	2.28	2.24	2.09	2.01	1.92	1.87	1.84	1.78	1.77	
30	4.17	3.32	2.92	2.69	2.53	2.42	2.33	2.27	2.21	2.16	2.01	1.93	1.84	1.79	1.76	1.70	1.66	
40	4.08	3.23	2.84	2.61	2.45	2.34	2.25	2.18	2.12	2.08	1.92	1.84	1.74	1.69	1.66	1.59	1.55	
50	4.03	3.18	2.79	2.56	2.40	2.29	2.20	2.13	2.07	2.03	1.87	1.78	1.69	1.63	1.60	1.52	1.48	
60	4.00	3.15	2.76	2.53	2.37	2.25	2.17	2.10	2.04	1.99	1.84	1.75	1.65	1.59	1.56	1.48	1.43	
70	3.98	3.13	2.74	2.50	2.35	2.23	2.14	2.07	2.02	1.97	1.81	1.72	1.62	1.57	1.53	1.45	1.39	
80	3.96	3.11	2.72	2.49	2.33	2.21	2.13	2.06	2.00	1.95	1.79	1.70	1.60	1.54	1.51	1.43	1.37	
90	3.95	3.10	2.71	2.47	2.32	2.20	2.11	2.04	1.99	1.94	1.78	1.69	1.59	1.53	1.49	1.41	1.35	
100	3.94	3.09	2.70	2.46	2.31	2.19	2.10	2.03	1.97	1.93	1.77	1.68	1.57	1.52	1.48	1.39	1.32	
200	3.89	3.04	2.65	2.42	2.26	2.14	2.06	1.98	1.93	1.88	1.72	1.62	1.52	1.46	1.41	1.32	1.25	

Remember that $F_{a,b,1-p} = 1/F_{b,a,p}$.

Quantiles $F_{a,b;0.99}$ of the F -distribution with a and b degrees of freedom

b	a																	
	1	2	3	4	5	6	7	8	9	10	15	20	30	40	50	100	∞	
1	4999	5403	5625	5764	5859	5928	5982	6022	6056	6157	6209	6261	6287	6300	6330	6366	6366	
2	99.0	99.2	99.2	99.3	99.3	99.4	99.4	99.4	99.4	99.4	99.4	99.5	99.5	99.5	99.5	99.5	99.5	
3	30.80	29.50	28.70	28.20	27.90	27.70	27.50	27.30	27.20	26.90	26.70	26.50	26.40	26.40	26.20	26.10	26.10	
4	18.00	16.70	16.00	15.50	15.20	15.00	14.80	14.70	14.50	14.20	14.00	13.80	13.70	13.70	13.60	13.50	13.50	
5	13.30	12.10	11.40	11.00	10.70	10.50	10.30	10.20	10.10	9.72	9.55	9.38	9.29	9.24	9.13	9.02	9.02	
6	10.90	9.78	9.15	8.75	8.47	8.26	8.10	7.98	7.87	7.56	7.40	7.23	7.14	7.09	6.99	6.88	6.88	
7	9.55	8.45	7.85	7.46	7.19	6.99	6.84	6.72	6.62	6.31	6.16	5.99	5.91	5.86	5.75	5.65	5.65	
8	8.65	7.59	7.01	6.63	6.37	6.18	6.03	5.91	5.81	5.52	5.36	5.20	5.12	5.07	4.96	4.86	4.86	
9	8.02	6.99	6.42	6.06	5.80	5.61	5.47	5.35	5.26	4.96	4.81	4.65	4.57	4.52	4.42	4.31	4.31	
10	7.56	6.55	5.99	5.64	5.39	5.20	5.06	4.94	4.85	4.56	4.41	4.25	4.17	4.12	4.01	3.91	3.91	
12	6.93	5.95	5.41	5.06	4.82	4.64	4.50	4.39	4.30	4.01	3.86	3.70	3.62	3.57	3.47	3.36	3.36	
14	6.51	5.56	5.04	4.69	4.46	4.28	4.14	4.03	3.94	3.66	3.51	3.35	3.27	3.22	3.11	3.00	3.00	
16	6.23	5.29	4.77	4.44	4.20	4.03	3.89	3.78	3.69	3.41	3.26	3.10	3.02	2.97	2.86	2.75	2.75	
18	6.01	5.09	4.58	4.25	4.01	3.84	3.71	3.60	3.51	3.23	3.08	2.92	2.84	2.78	2.68	2.57	2.57	
20	5.85	4.94	4.43	4.10	3.87	3.70	3.56	3.46	3.37	3.09	2.94	2.78	2.69	2.64	2.54	2.42	2.42	
25	5.57	4.68	4.18	3.86	3.63	3.46	3.32	3.22	3.13	2.85	2.70	2.54	2.45	2.40	2.29	2.17	2.17	
30	5.39	4.51	4.02	3.70	3.47	3.30	3.17	3.07	2.98	2.70	2.55	2.39	2.30	2.25	2.13	2.01	2.01	
40	5.18	4.31	3.83	3.51	3.29	3.12	2.99	2.89	2.80	2.52	2.37	2.20	2.11	2.06	1.94	1.80	1.80	
50	5.06	4.20	3.72	3.41	3.19	3.02	2.89	2.79	2.70	2.42	2.27	2.10	2.01	1.95	1.82	1.68	1.68	
60	4.98	4.13	3.65	3.34	3.12	2.95	2.82	2.72	2.63	2.35	2.20	2.03	1.94	1.88	1.75	1.60	1.60	
70	4.92	4.08	3.60	3.29	3.07	2.91	2.78	2.67	2.59	2.31	2.15	1.98	1.89	1.83	1.70	1.54	1.54	
80	4.88	4.04	3.56	3.26	3.04	2.87	2.74	2.64	2.55	2.27	2.12	1.94	1.85	1.79	1.66	1.49	1.49	
90	4.85	4.01	3.54	3.23	3.01	2.84	2.72	2.61	2.52	2.24	2.09	1.92	1.82	1.76	1.62	1.46	1.46	
100	4.82	3.98	3.51	3.21	2.99	2.82	2.69	2.59	2.50	2.22	2.07	1.89	1.80	1.73	1.60	1.43	1.43	
200	4.71	3.88	3.41	3.11	2.89	2.73	2.60	2.50	2.41	2.13	2.07	1.79	1.69	1.63	1.48	1.28	1.28	

Critical values $D_{n;p}$ of the Kolmogorov-Smirnov one-sample test.

n	0.95	0.99	n	0.95	0.99
			21	0.2827	0.3443
8	0.4543	0.5419	22	0.2809	0.3367
9	0.4300	0.5133	23	0.2749	0.3295
10	0.4093	0.4889	24	0.2693	0.3229
11	0.3912	0.4677	25	0.2640	0.3166
12	0.3754	0.4491	26	0.2591	0.3106
13	0.3614	0.4325	27	0.2544	0.3050
14	0.3489	0.4176	28	0.2499	0.2997
15	0.3376	0.4042	29	0.2457	0.2947
16	0.3273	0.3920	30	0.2417	0.2899
17	0.3180	0.3809	35	0.2243	0.2690
18	0.3094	0.3706	40	0.2101	0.2521
19	0.3014	0.3612	45	0.1984	0.2380
20	0.2941	0.3524	50	0.1884	0.2260

For $n > 50$ use the approximation $D_{n;p} \approx K_p/\sqrt{n}$.

For $p > 1/2$, $K_p \approx \sqrt{\log(\frac{2}{1-p})}/2$. In particular

$$K_{0.95} = 1.358, \quad K_{0.99} = 1.628.$$

critical values $D_{n,m,0.95}$ of the two-sample Kolmogorov-Smirnov test:

	8	9	10	11	12	13	14	15	16	17	18	19	20	22	24	26	28	30	35	40	45	5	
8	.6250																						
9	.6250	.5555																					
10	.5875	.5777	.6000																				
11	.5795	.5858	.5363	.5454																			
12	.5937	.5555	.5333	.5000	.5000																		
13	.5673	.5470	.5230	.5104	.5128	.4615																	
14	.5625	.5396	.5142	.5129	.5000	.4835	.5000																
15	.5500	.5333	.5000	.5090	.5000	.4820	.4619	.4666															
16	.5625	.5277	.5187	.5000	.4843	.4807	.4687	.4750	.4375														
17	.5294	.5294	.5058	.4919	.4803	.4705	.4579	.4470	.4522	.4117													
18	.5416	.5000	.5000	.4848	.4768	.4615	.4523	.4444	.4375	.4313	.4444												
19	.5263	.5146	.4947	.4832	.4649	.4615	.4511	.4421	.4342	.4303	.4122	.4210											
20	.5312	.5055	.5000	.4818	.4666	.4538	.4428	.4333	.4250	.4235	.4194	.4000	.4000										
22	.5284	.5101	.4863	.4545	.4621	.4545	.4415	.4333	.4232	.4144	.4116	.4019	.3954	.3636									
24	.5000	.5000	.4875	.4621	.4583	.4455	.4285	.4250	.4192	.4093	.4027	.3947	.3916	.3844	.3750								
26	.5144	.5000	.4769	.4580	.4487	.4230	.4230	.4153	.4086	.4004	.3931	.3906	.3846	.3758	.3669	.3461							
28	.5000	.4920	.4714	.4577	.4434	.4340	.4285	.4142	.4040	.3991	.3888	.3815	.3803	.3701	.3571	.3516	.3571						
30	.5125	.4888	.4666	.4545	.4333	.4256	.4142	.4000	.4000	.3901	.3888	.3789	.3666	.3636	.3597	.3487	.3416	.3333					
35	.5000	.4698	.4571	.4415	.4333	.4175	.4000	.4000	.3910	.3865	.3777	.3699	.3657	.3519	.3452	.3384	.3285	.3238	.3142				
40	.5000	.4777	.4525	.4363	.4250	.4153	.4000	.3916	.3890	.3750	.3666	.3631	.3512	.3465	.3333	.3288	.3223	.3175	.3035	.3000			
45	.4888	.4666	.4444	.4363	.4185	.4034	.3952	.3777	.3777	.3712	.3666	.3555	.3511	.3414	.3287	.3205	.3142	.3111	.2952	.2838	.2666		
50	.4850	.4666	.4420	.4272	.4150	.4046	.3928	.3813	.3725	.3647	.3588	.3494	.3400	.3327	.3241	.3153	.3092	.3000	.2891	.2800	.2688	.260	

critical values $D_{n,m,0.99}$ of the two-sample Kolmogorov-Smirnov test:

	8	9	10	11	12	13	14	15	16	17	18	19	20	22	24	26	28	30	35	40	45	5		
8	.7500																							
9	.7500	.6666																						
10	.7125	.6888	.7100																					
11	.7045	.6969	.6454	.6446																				
12	.6666	.6666	.6583	.6515	.5902																			
13	.6730	.6581	.6307	.6153	.6089	.6153																		
14	.6696	.6507	.6285	.6103	.6071	.5714	.5714																	
15	.6666	.6444	.6400	.6000	.5888	.5794	.5809	.5333																
16	.6250	.6388	.6125	.5965	.5885	.5673	.5580	.5500	.5625															
17	.6397	.6078	.6117	.5828	.5735	.5701	.5546	.5490	.5220	.5294														
18	.6388	.6172	.5944	.5808	.5555	.5555	.5515	.5333	.5277	.5359	.5000													
19	.6381	.6198	.5947	.5837	.5614	.5587	.5413	.5263	.5230	.5077	.5116	.4736												
20	.6250	.6000	.6000	.5727	.5708	.5500	.5392	.5200	.5156	.5058	.5000	.4921	.5000											
22	.6250	.5959	.5863	.5454	.5530	.5454	.5259	.5121	.5085	.5000	.4924	.4880	.4795	.4545										
24	.6250	.6018	.5791	.5606	.5416	.5288	.5208	.5111	.5000	.4901	.4884	.4758	.4666	.4545	.4583									
26	.6105	.5897	.5692	.5524	.5352	.5384	.5137	.5025	.4927	.4773	.4743	.4676	.4615	.4527	.4407	.4230								
28	.6071	.5714	.5607	.5519	.5238	.5192	.5000	.4952	.4821	.4789	.4682	.4661	.4571	.4415	.4345	.4217	.4285							
30	.6083	.5777	.5700	.5424	.5194	.5128	.4952	.5000	.4812	.4725	.4666	.4578	.4516	.4363	.4250	.4179	.4130	.4000						
35	.5964	.5746	.5428	.5324	.5214	.5010	.4857	.4761	.4714	.4537	.4428	.4421	.4371	.4259	.4130	.4054	.3938	.3904	.3714					
40	.6000	.5638	.5500	.5272	.5104	.4942	.4821	.4666	.4625	.4485	.4430	.4342	.4250	.4147	.4000	.3942	.3857	.3758	.3650	.3506				
45	.5861	.5555	.5333	.5111	.5000	.4888	.4746	.4666	.4541	.4444	.4333	.4257	.4177	.4111	.3953	.3854	.3777	.3666	.3555	.3416	.3333			
50	.5825	.5600	.5400	.5109	.4966	.4830	.4700	.4600	.4487	.4411	.4300	.4210	.4110	.4000	.3891	.3769	.3700	.3600	.3462	.3350	.3248	.320		

Anhang B

Solutions of Problems

B.1 Problems from Chapter 1

1. $\mathbb{E}\bar{X}_n = \frac{1}{n} \sum \mathbb{E}X_i = \mu,$

$$\mathbb{V}\bar{X}_n = \left(\frac{1}{n}\right)^2 \sum \mathbb{V}X_i = \frac{\sigma^2}{n}.$$

(As drawing is done with replacement, X_1, \dots, X_n are independent.)

2.

$$\begin{aligned} \mathbb{E}(S_n^2) &= \frac{1}{n-1} \mathbb{E}\left(\sum_{i=1}^n (X_i - \bar{X}_n)^2\right) = \frac{1}{n-1} \mathbb{E}\left(\sum_{i=1}^n X_i^2 - n\bar{X}_n^2\right) = \\ &= \frac{1}{n-1} (n\mathbb{E}(X_1^2) - n\mathbb{E}(\bar{X}_n^2)) = \frac{1}{n-1} (n((\mathbb{E}X_1)^2 + \mathbb{V}(X_1)) - n((\mathbb{E}\bar{X}_n)^2 + \mathbb{V}(\bar{X}_n))) = \\ &= \frac{n}{n-1} (\mathbb{V}(X_1) - \mathbb{V}(\bar{X}_n)) = \sigma^2. \end{aligned}$$

3. By symmetry, each number x_1, \dots, x_n is equally likely to be the one drawn in the k -th place. This means that the distribution of X_k is the same as in problem 1. So, again, the expectation of the sample mean is μ .

For the variance,

$$\mathbb{V}(\bar{X}_n) = \frac{1}{n^2} (n\mathbb{V}X_1 + n(n-1)\mathbf{Cov}(X_1, X_2)).$$

One can easily calculate the covariance above directly, but the easiest way is to observe that for $n = N$ the variance of the sample mean is 0. Thus, the covariance term is

$$-\frac{\sigma^2}{N-1}$$

and we finally obtain

$$\mathbb{V}(\bar{X}_n) = \frac{\sigma^2}{n} \frac{N-n}{N-1}.$$

4. As in problem 2, we have

$$\mathbb{E}(S_n^2) = \frac{n}{n-1} (\mathbb{V}(X_1) - \mathbb{V}(\bar{X}_n)).$$

Using the result of problem 3, we finally obtain $\frac{N}{N-1}\sigma^2$.

5.

$$\begin{aligned} L(X_1, \dots, X_n, \mu, \sigma^2) &= \frac{1}{\sqrt{(2\pi\sigma^2)^n}} \exp\left(-\frac{\sum_{i=1}^n (X_i - \mu)^2}{2\sigma^2}\right) = \\ &= \frac{1}{\sqrt{(2\pi\sigma^2)^n}} \exp\left(-\frac{(n-1)S_n^2 + n(\bar{X}_n - \mu)^2}{2\sigma^2}\right). \end{aligned}$$

6. $T = \max(X_1, \dots, X_n)$.
7. $T_1 = \max(X_1, \dots, X_n)$, $T_2 = \min(X_1, \dots, X_n)$.
8. The likelihood function is

$$L = \theta^n \exp(-\theta n \bar{X}_n) I_{[0, \infty]^n}(X_1, \dots, X_n).$$

This is of the form $f(\theta, \bar{X}_n)g(X_1, \dots, X_n)$.

9. $T_1 = \prod X_i$, $T_2 = \sum X_i$.

B.2 Problems from Chapter 2

1. In our proof of the Cramér-Rao theorem, we use Cauchy's inequality. For equality to hold, the two factors under the integral must be proportional. Thus,

$$\frac{\partial}{\partial \theta}(\log f_\theta(x)) = (\hat{\theta} - \theta)C(\theta)$$

where C is the proportionality factor which may depend on θ . By integrating over θ we get

$$\log f_\theta(x) = c(x) + \hat{\theta} \int C(\theta) d\theta - \int \theta C(\theta) d\theta.$$

This is already of the form we need. If we look at a sample of size n , the likelihood function must be of the same form, with the only difference that c and $\hat{\theta}$ are functions of n variables X_1, \dots, X_n . $X_1 = X_2 = \dots = X_n$, we see that f_θ has the form

$$\log f_\theta(x) = c_f(x) + a(x)A(\theta) + B(\theta),$$

where

$$c_f(x) = \frac{1}{n} c(x, \dots, x),$$

$$a(x) = \hat{\theta}(x, \dots, x),$$

$$A(\theta) = \frac{1}{n} \int C(\theta) d\theta,$$

and

$$B(\theta) = \frac{1}{n} \int \theta C(\theta) d\theta,$$

and comparing coefficients, we see that

$$\hat{\theta} = \frac{1}{n} \sum a(X_i).$$

On the other hand, by differentiating the equation

$$\int f_\theta(x) d\mu(x) = 1$$

with respect to θ , we find that if the densities are of this form, then $\hat{\theta}$ is an unbiased estimator and attains the Cramér-Rao bound.

- 2.

$$\hat{\lambda} = \frac{1}{\bar{X}_n}.$$

\bar{X}_n has a Gamma distribution with density

$$\frac{x^{n-1} (n\lambda)^n e^{-n\lambda x}}{n!} \quad (x > 0).$$

This yields

$$\mathbb{E}(\hat{\lambda}) = \mathbb{E}\left(\frac{1}{\bar{X}_n}\right) = \int_0^\infty \frac{x^{n-2}(n\lambda)^n e^{-n\lambda x}}{(n-1)!} = \frac{n\lambda}{n-1}.$$

The unbiased estimator is

$$\tilde{\lambda} = \frac{n-1}{n\bar{X}_n}.$$

3. (a) The likelihood function is

$$L(X_1, \dots, X_n, \theta) = \begin{cases} \theta^{-n} & \text{if } \theta \geq \max(X_1, \dots, X_n) \\ 0 & \text{otherwise.} \end{cases}$$

The first branch is decreasing, so the maximum is attained by choosing θ as small as possible. So, the ML estimator is

$$\hat{\theta}_a = \max(X_1, \dots, X_n).$$

- (b)

$$\hat{\theta}_b = 2\bar{X}_n.$$

- (c)

$$\hat{\theta}_c = \frac{n+1}{n}\hat{\theta}_a.$$

- (d)

$$\mathbb{V}(\hat{\theta}_b) = \frac{\theta^2}{3n},$$

$$\mathbb{V}(\hat{\theta}_c) = \frac{\theta^2}{n(n+2)},$$

(Use the fact that $\max(X_1, \dots, X_n)$ has distribution function $(x/\theta)^n$ for $0 \leq x \leq \theta$ and the latter is smaller for $n > 1$ (for $n = 1$ the two estimators are the same).)

4. $\log f_\lambda(x) = \log \lambda - \lambda x$. Differentiating twice with respect to λ gives $-\lambda^{-2}$. Thus the Cramér-Rao bound is λ^{-2} . The estimator from problem 2 does not attain this bound — its variance is larger by a factor $\frac{n-1}{n-2}$.
5. We verify that the estimator attains the Cramér-Rao bound. This is most easily verified by showing that the density

$$f_\theta(x) = \theta^x(1-\theta)^{1-x}$$

satisfies the conditions from problem 1.

6. Again, use problem 1.

7. (a) $\hat{\theta}_a = \frac{1}{2} \max(X_1, \dots, X_n)$.

(b) $\hat{\theta}_b = \frac{2(n+1)}{2n+1}\hat{\theta}_a$.

(c) $\mathbb{V}(\hat{\theta}_b) = \frac{n\theta^2}{(2n+1)^2(n+2)}$.

The joint density of $U = \max(X_1, \dots, X_n)$ and $V = \min(X_1, \dots, X_n)$ is

$$f(u, v) = \begin{cases} \frac{n(n-1)(u-v)^{n-2}}{\theta^n} & \text{if } \theta < v < u < 2\theta, \\ 0 & \text{otherwise,} \end{cases}$$

so

$$\mathbb{V}\tilde{\theta} = \frac{n(n-1)}{9\theta^n} \int_\theta^{2\theta} \int_\theta^u (u+v)^2(u-v)^{n-2} dv du - \theta^2 = \frac{2\theta^2}{9(n+1)(n+2)}.$$

For $n > 1$ this is smaller than the variance of $\hat{\theta}_b$.

8. We have to minimize

$$\sum_{i=1}^n |X_i - \theta|.$$

Unfortunately, the usual way of zeroing the derivative doesn't work perfectly in this case; but if θ is different from X_1, \dots, X_n , then the derivative exists and equals

$$\sum \text{sig}(\theta - X_i) = \#\{i : X_i < \theta\} - \#\{i : X_i > \theta\}.$$

Thus, if more than $n/2$ of the X_i are below θ , the sum above is increasing, otherwise it is decreasing. So, the minimum is attained if

$$\theta = X_{n:\frac{n+1}{2}}$$

if n is odd, and for any

$$X_{n:\frac{n}{2}} \leq \theta \leq X_{n:\frac{n}{2}+1}$$

if n is even, in other words, if θ is a median of X_1, \dots, X_n .

9. $\bar{X}_n = 1.84$, $S_n^2 = .3338$.

Confidence interval for μ : [1.43, 2.25]

Confidence interval for σ^2 : [0.158, 1.11]

10. [1.264, 2.416].

11. $[\bar{X}_n - z_{\frac{1+\gamma}{2}} \sqrt{\frac{\bar{X}_n}{n}}, \bar{X}_n + z_{\frac{1+\gamma}{2}} \sqrt{\frac{\bar{X}_n}{n}}]$

12. $[\frac{1 - z_{\frac{1+\gamma}{2}} n^{-1/2}}{\bar{X}_n}, \frac{1 + z_{\frac{1+\gamma}{2}} n^{-1/2}}{\bar{X}_n}]$

13. $[\frac{\chi_{2n; \frac{1-\gamma}{2}}^2}{2n\bar{X}_n}, \frac{\chi_{2n; \frac{1+\gamma}{2}}^2}{2n\bar{X}_n}]$.

14. If we find $l(x) < u(x)$ such that for all λ

$$\mathbb{P}_\lambda(\lambda < l(X)) \leq \frac{1-\gamma}{2}$$

and

$$\mathbb{P}_\lambda(\lambda > u(X)) \leq \frac{1-\gamma}{2},$$

then $[l(X), u(X)]$ will be a confidence interval with coverage probability γ . Assume that both l and u are nondecreasing. For $l(k-1) < \lambda < l(k)$

$$\mathbb{P}_\lambda(\lambda < l(X)) = \mathbb{P}_\lambda(X \geq k) = \mathbb{P}(Y \leq 2\lambda),$$

where Y is $\chi_i^2 2k$. This probability should be at most $\frac{1-\gamma}{2}$, so we choose

$$l(k) = \frac{1}{2} \chi_{2k; \frac{1-\gamma}{2}}^2.$$

In a similar way we obtain

$$u(k) = \frac{1}{2} \chi_{2k+2; \frac{1+\gamma}{2}}^2.$$

If we have a sample of size n , $S = \sum_{i=1}^n X_i$ has a Poisson distribution with parameter $n\lambda$, so finally, the confidence interval is

$$[\frac{1}{2n} \chi_{2S; \frac{1-\gamma}{2}}^2, \frac{1}{2n} \chi_{2S+2; \frac{1+\gamma}{2}}^2].$$

15. Let $\hat{\theta} = \max(X_1, \dots, X_n)$. Since

$$\mathbb{P}(\hat{\theta} < x\theta) = x^n,$$

we can use the confidence interval

$$[\hat{\theta}, \hat{\theta} \frac{1}{(1-\gamma)^{1/n}}]$$

B.3 Problems from Chapter 3

1. First assume that $\mu_1 > \mu_0$. The ratio f_1/f_0 has the logarithm

$$\frac{\sum(X_i - \mu_0)^2}{2\sigma^2} - \frac{\sum(X_i - \mu_1)^2}{2\sigma^2} = \frac{\mu_1 - \mu_0}{\sigma^2} \sum X_i + C,$$

where C is some constant. Thus, the Neyman-Pearson test is equivalent to

$$\phi = \begin{cases} 1 & \text{if } \sum X_i > K, \\ c & \text{if } \sum X_i = K, \\ 0 & \text{if } \sum X_i < K. \end{cases}$$

As $\sum X_i$ has a continuous distribution, we can let $c = 0$. We calculate K from

$$\alpha = \mathbb{P}_0(\sum X_i > K) = 1 - \Phi\left(\frac{K - n\mu_0}{\sqrt{n\sigma^2}}\right),$$

so

$$K = n\mu_0 + z_{1-\alpha}\sqrt{n\sigma^2}.$$

for $\mu_1 < \mu_0$, we reject if

$$\sum X_i < n\mu_0 + z_\alpha\sqrt{n\sigma^2}.$$

2. The denominator of the likelihood ratio statistic attains its maximum for $\mu = \bar{X}_n$ and

$$\sigma^2 = \frac{1}{n} \sum (X_i - \bar{X}_n)^2.$$

In the numerator, $\mu = \mu_0$ is fixed, and the maximum is attained for

$$\sigma^2 = \frac{1}{n} \sum (X_i - \mu_0)^2.$$

The likelihood ratio becomes

$$\left(\frac{\sum(X_i - \bar{X}_n)^2}{\sum(X_i - \mu_0)^2}\right)^{n/2} = \left(\frac{\sum(X_i - \bar{X}_n)^2}{\sum(X_i - \bar{X}_n)^2 + n(\bar{X}_n - \mu_0)^2}\right)^{n/2} = \left(1 + \frac{T^2}{n-1}\right)^{-n/2},$$

where T is the t statistic. Thus we reject if $|T|$ is larger than a critical value; this is exactly the two-sided t test.

3. For $\lambda_1 > \lambda_0$, the Neyman-Pearson test is

$$\phi = \begin{cases} 1 & \text{if } \sum X_i > K, \\ c & \text{if } \sum X_i = K, \\ 0 & \text{if } \sum X_i < K, \end{cases}$$

where K and $0 \leq c < 1$ are determined from

$$c \frac{(n\lambda_0)^K}{K!} e^{-n\lambda_0} + \sum_{i>K} \frac{(n\lambda_0)^i}{i!} e^{-n\lambda_0} = \alpha.$$

4. For $\theta_1 > \theta_0$, the Neyman-Pearson test is

$$\phi = \begin{cases} 1 & \text{if } \theta_0 < \max(X_1, \dots, X_n) \leq \theta_1, \\ \alpha & \text{if } \max(X_1, \dots, X_n) \leq \theta_0, \end{cases}$$

We can modify ϕ on a set of probability 0:

$$\tilde{\phi} = \begin{cases} 1 & \text{if } \theta_0 < \max(X_1, \dots, X_n), \\ \alpha & \text{if } \max(X_1, \dots, X_n) \leq \theta_0. \end{cases}$$

This does no longer depend on the value of θ_1 , so it is the best test for $H_0 : \theta = \theta_0$ against $H_1 : \theta > \theta_0$. For any $\theta < \theta_0$, The rejection probability is α , so it is also a test for $H_0 : \theta \leq \theta_0$ against $H_1 : \theta > \theta_0$.

5. The likelihood function is

$$L(X_1, \dots, X_n, \theta) = g(X_1, \dots, X_n)f(T, \theta).$$

The factor g cancels in the likelihood ratio statistic, so it only depends on T .

6. Reject if

$$\sum \left(\frac{(X_i - \mu_0)^2}{\sigma_0^2} - \frac{(X_i - \mu_1)^2}{\sigma_1^2} \right) > t_c.$$

Unfortunately, we cannot handle the distribution of this. (Properly scaled, we obtain a so-called non-central χ^2 distribution which is the distribution of $(X_1 - a)^2 + \sum_{i=2}^n X_i^2$, where X_i are independent standard normal).

7. We can easily calculate the likelihood ratio test which rejects if

$$\lambda_0 \bar{X}_n e^{-\lambda_0 \bar{X}_n} < K,$$

which is equivalent to

$$X < \frac{A}{\lambda_0} \text{ or } X > \frac{B}{\lambda_0},$$

where A and B are determined by the equations

$$Ae^{-A} = Be^{-B}$$

and

$$\int_A^B \frac{x^{n-1}}{(n-1)!} e^{-x} dx = 1 - \alpha.$$

Another approach is to use the ideas from section 3.3 to obtain the best unbiased test. This turns out (after some calculations) to be identical to the likelihood ratio test.

8. In the likelihood ratio statistic, the denominator is maximal for $\mu_1 = \bar{X}$, $\mu_2 = \bar{Y}$,

$$\sigma_1^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2,$$

$$\sigma_2^2 = \frac{1}{m} \sum_{i=1}^m (Y_i - \bar{Y})^2.$$

In the numerator, we have to put

$$\sigma^2 = \frac{1}{n+m} \left(\sum_{i=1}^n (X_i - \bar{X})^2 + \sum_{i=1}^m (Y_i - \bar{Y})^2 \right).$$

The likelihood ratio statistic is a function of F .

9. We test $H_0 : p = 0.05$ against $H_1 : p > 0.05$ (we don't care about p being smaller).

$\hat{p} = 0.075$. The critical value is

$$0.05 + 1.645 \sqrt{\frac{.05 * .95}{200}} = .0754$$

We accept H_0 .

10. This time, test $H_0 : \mu = 500$ against $H_1 : \mu < 500$. We reject if

$$\bar{X}_n < 500 - t_{59, .95} \sqrt{\frac{2000}{60}} = 490.4$$

We reject H_0 .

11. Test $H_0 : \lambda = 1/2000$ against $H_1 : \lambda > 1/2000$ (i.e., test against smaller mean lifetime) We reject if $S = \sum X_i$ is smaller than a critical value. This critical value can be obtained by making use of the fact that $2\lambda S$ has a χ^2 distribution with $2n$ degrees of freedom:

$$t_c = \frac{\chi_{2n;0.05}^2}{2\lambda} = 10851,$$

$S = 18950$ is greater than that, so we accept H_0 .

12. $\bar{X}_n = 1.155$. We reject if

$$|\bar{X}_n - 1.0| > 1.96\sqrt{.5/11} = 0.417,$$

and as the left hand side is only 0.155, we accept H_0 .

13. $S_n^2 = 0.1627$, and we use

$$t = \frac{1.155 - 1.0}{\sqrt{0.1627/11}} = 1.274$$

This is smaller in modulus than $t_{10;0.975} = 2.030$, so we accept H_0 .

14. We reject if

$$(n-1)S_n^2 < \sigma_0^2 \chi_{n-1,0.05}^2 = 1.970$$

or

$$(n-1)S_n^2 > \sigma_0^2 \chi_{n-1,0.95}^2 = 24.996.$$

As the left-hand side is only 1.79, we reject H_0 .

- 15.

$$f_{a,b}(x) = \frac{a}{b} \frac{\Gamma(\frac{a+b}{2})}{\Gamma(\frac{a}{2})\Gamma(\frac{b}{2})} \frac{(\frac{a}{b}x)^{\frac{a}{2}-1}}{(1 + \frac{a}{b}x)^{\frac{a+b}{2}}} \quad (x > 0).$$

16. As χ_n^2 is the distribution of a sum of the squares of n independent standard normal random variables, it follows from the central limit theorem that it has an approximate normal distribution. As its expectation and variance are n and $2n$, respectively, we get the desired result.

17. $F_{a,b}$ is the distribution of $\frac{X/a}{Y/b}$, where X and Y are independent χ^2 with a resp. b degrees of freedom. In the light of the previous problem, we may approximately (up to terms of smaller order) replace X by $a + U\sqrt{2a}$, and Y by $b + V\sqrt{2b}$, where U and V are standard normal. This gives

$$F \approx \frac{1 + U\sqrt{2/a}}{1 + V\sqrt{2/b}} \approx 1 + U\sqrt{\frac{2}{a}} - V\sqrt{\frac{2}{b}} = 1 + W\sqrt{\frac{2(a+b)}{ab}},$$

where W is another standard normal random variable.

18. Under H_0 , $Z_i = X_i - 2Y_i$ is normal with mean zero, so we can use the standard t test on Z_1, \dots, Z_n .

19. We consider the general case of testing $H_0 : \mu_Y = a\mu_X$ under the assumption of equal variances. We easily find that

$$\bar{Y}_m - a\bar{X}_n \sim \mathcal{N}(0, \sigma^2(\frac{1}{m} + \frac{a^2}{n})),$$

$$\frac{m-1}{\sigma^2} S_Y^2 \sim \chi_{m-1}^2,$$

and

$$\frac{n-1}{\sigma^2} S_X^2 \sim \chi_{n-1}^2,$$

so we use the test statistic

$$T = \frac{\bar{Y}_m - a\bar{X}_n}{\sqrt{(n-1)S_X^2 + (m-1)S_Y^2}} \sqrt{\frac{nm(n+m-2)}{n+a^2m}},$$

which has a t distribution with $m+n-2$ degrees of freedom.

20. $\bar{Y}_m = 0.9444$, $S_Y^2 = 0.05778$

$$t = \frac{1.155 - 0.944}{\sqrt{10 * .1627 + 8 * 0.05778}} \sqrt{\frac{11 * 9 * 18}{20}} = 1.26$$

This is less than $t_{18,0.975} = 2.101$, so we accept H_0 .

21. $F = \frac{.1627}{.05778} = 2.816 < F_{10,8;0.95} = 3.35$, so we accept H_0 .

B.4 Problems from Chapter 4

1. We denote the variables in the sample for combination (i, j) by X_{ijl} , $l = 1, \dots, n_{ij}$. Our model is

$$\mu_{ij} = \mu + a_i + b_j$$

with

$$\sum_{i=1}^n a_i n_{ij} = \sum_{j=1}^k b_j n_{ij} = 0$$

Let

$$n_{i.} = \sum_{j=1}^k n_{ij},$$

$$n_{.j} = \sum_{i=1}^n n_{ij},$$

$$n_{..} = \sum_{j=1}^k n_{.j},$$

then we get

$$\hat{\mu} = \frac{1}{n_{..}} \sum_{i,j,l} X_{ijl} = \bar{X}_{...},$$

$$\hat{a}_i = \frac{1}{n_{i.}} \sum_{j,l} X_{ijl} - \hat{\mu} = \bar{X}_{i..} - \bar{X}_{...},$$

$$\hat{b}_j = \frac{1}{n_{.j}} \sum_{i,l} X_{ijl} - \hat{\mu} = \bar{X}_{.j.} - \bar{X}_{...}.$$

In addition to the definitions of $X_{...}$ etc. as given above, let

$$\bar{X}_{ij.} = \frac{1}{n_{ij}} \sum_l X_{ijl}.$$

Finally, the test statistic is

$$F = \frac{\sum_i n_{i.} (\bar{X}_{i..} - \bar{X}_{...})^2}{\sum_{i,j,l} (X_{ijl} - X_{i..} - X_{.j.} + X_{...})^2},$$

and this has $k - 1$ and $n_{..} - (n + k - 1)$ degrees of freedom.

2. With the same notations as in the previous problem, our test statistic is

$$F = \frac{\sum_{i,j} n_{ij} (X_{ij.} - X_{i..} - X_{.j.} + X_{...})^2 / (n - 1)(k - 1)}{\sum_{i,j,l} (X_{ijl} - X_{ij.})^2 / (n_{..} - nk)}.$$

This has $(n - 1)(k - 1)$ and $n_{..} - nk$ degrees of freedom.

3. We show that the F statistic is just the square of the t statistic used in the t test for the comparison of two samples of equal variance. The denominator is already of the form we need, so we only have to check the numerator. This becomes

$$n_1(\bar{X}_1 - \bar{X})^2 + n_2(\bar{X}_2 - \bar{X})^2,$$

and as

$$\bar{X} = \frac{1}{n}(n_1\bar{X}_1 + n_2\bar{X}_2),$$

we obtain

$$(\bar{X}_1 - \bar{X}_2)^2 \frac{n_1 n_2^2 + n_2 n_1^2}{n} = (\bar{X}_1 - \bar{X}_2)^2 \frac{n_1 n_2}{n_1 + n_2},$$

which is exactly what we need.

4.

$$F = \frac{\sum_{i=1}^4 n_i (\bar{X}_i - \bar{X})^2 / 3}{\sum_{i=1}^4 (n_i - 1) S_i^2 / (n - k)} = \frac{1.06/3}{4.5/39} = 3.0622 > 2.84 = F_{3,39;.95}$$

We reject H_0 .

5. The confidence interval using all the samples is

$$[5.3 - t_{39,.975} \sqrt{\frac{4.5}{39 * 10}}, 5.3 + t_{39,.975} \sqrt{\frac{4.5}{39 * 10}}] = [5.0827, 5.5173].$$

6. From sample 1 alone we get

$$[5.0628, 5.5372].$$

7. we let

$$S^2 = \frac{1}{(n-1)(k-1)} \sum_{i,j} (X_{ij} - X_{i.} - X_{.j} + X_{..})^2.$$

with this notation, and from the fact that

$$\hat{a}_i = X_{i.} - X_{..}$$

has a normal distribution with mean a_i and variance $\frac{n-1}{nk} \sigma^2$, we get the confidence interval

$$[X_{i.} - X_{..} - t_{(n-1)(k-1), \frac{1+\gamma}{2}} \sqrt{\frac{(n-1)S^2}{nk}}, \\ X_{i.} - X_{..} + t_{(n-1)(k-1), \frac{1+\gamma}{2}} \sqrt{\frac{(n-1)S^2}{nk}}].$$

8. This is a simple two-way question with $n = 4$, $k = 3$; the F statistic is

$$F = 1.6071 < 4.76 = F_{3,6;.95}$$

we accept H_0 .

9.

$$F = 1.25 < 3.16 = F_{3,18;0.95},$$

we accept H_0 .

B.5 Problems from Chapter 5

1. Of course, we want to make the variance of our estimators as small as possible. This means that

- (a) for \hat{a} , we have to make $\sum(x_i - \bar{x})^2$ as large as possible. The derivative with respect to x_i is

$$2(x_i - \bar{x}).$$

Thus, the sum is increasing if $x_i > \bar{x}$, and decreasing if $x_i < \bar{x}$. The maximum, therefore, can only be attained if $x_i = A$ or $x_i = -A$. Let k be the number of x 's that are $= A$. Then the sum becomes

$$A^2(n - \frac{1}{n}(n - 2k)^2).$$

This becomes a minimum if $k = n/2$.

- (b) for \hat{b} , the variance will be minimal if $\bar{x} = 0$. We don't need any other condition, but it might be best to use the set of x 's from part (a), as this also makes the variance of \hat{a} a minimum.
2. (a) The analysis from problem 1 stays valid, but now, we can't choose $k = n/2$ as n is odd. Instead, both $k = (n - 1)/2$ and $k = (n + 1)/2$ will minimize the variance of \hat{a} .
- (b) Again, the variance of \hat{b} is a minimum if $\bar{x} = 0$. This time, however, the solution of part (a) does not satisfy this condition.
3. As always, we assume that the errors are normally distributed. The likelihood function is

$$\frac{1}{\sqrt{(2\pi)^n \prod \sigma(x_i)^2}} \exp\left(-\sum \frac{(Y_i - ax_i - b)^2}{2\sigma(x_i)^2}\right).$$

Maximizing this likelihood function is equivalent to finding the minimum of

$$\sum \frac{(Y_i - ax_i - b)^2}{\sigma(x_i)^2} = \sum \left(\frac{Y_i}{\sigma(x_i)} - a\frac{x_i}{\sigma(x_i)} - b\frac{1}{\sigma(x_i)}\right)^2.$$

This is exactly what we have to minimize to obtain the usual least squares estimators for a and b in the modified problem.

- 4.

x	Y	x^2	xY	Y^2
0.5	1.3	0.25	0.65	1.69
1.0	2.7	1.00	2.70	7.29
2.0	1.5	4.00	3.00	2.25
2.5	2.0	6.25	5.00	4.00
3.0	2.2	9.00	6.60	4.84
9.0	9.7	20.50	17.95	20.07

$$\hat{a} = .114$$

$$\hat{b} = 1.735$$

$$\hat{\sigma}^2 = .398$$

$$\hat{Y}(3.5) = .114 * 3.5 + 1.735 = 2.134$$

prediction interval: $[-2.221, 6.459]$.

5. We let $Z = \log Y$ and compute the regression $Z = \log a + bx = B + Ax$:

x	Z	x^2	xZ	Z^2
0.5	0.2624	0.25	0.1312	0.0688
1.0	0.9933	1.00	0.9933	0.9866
2.0	0.4055	4.00	0.8110	0.1644
2.5	0.6931	6.25	1.7328	0.4804
3.0	0.7885	9.00	2.3655	0.6217
9.0	3.1428	20.50	6.0338	2.3419

$$\hat{A} = 0.08762 = \hat{b}$$

$$\hat{B} = 0.47085 = \log \hat{a}$$

$$\hat{a} = 1.6014$$

$$\hat{\sigma}^2 = 0.1111$$

$$\hat{Z}(3.5) = 0.47085 + .08762 * 3.5 = 0.7775$$

$$\hat{Y}(3.5) = 2.176$$

prediction interval for Z : $[-0.1894, 1.7444]$.

prediction interval for Y : $[0.8275, 5.7224]$.

6. We let $Z = 1/Y$ and compute the regression $Z = \frac{b}{a} + \frac{1}{a}x = B + Ax$:

x	Z	x^2	xZ	Z^2
0.5	0.7692	0.25	0.3846	0.5917
1.0	0.3704	1.00	0.3704	0.1372
2.0	0.6667	4.00	1.3333	0.4444
2.5	0.5000	6.25	1.2500	0.2500
3.0	0.4545	9.00	1.3636	0.2066
9.0	2.7608	20.50	4.7019	1.6299

$$\hat{A} = -0.06222$$

$$\hat{B} = 0.6642$$

$$\hat{a} = -16.07$$

$$\hat{b} = -10.67$$

$$\hat{\sigma}^2 = 0.0295$$

$$\hat{Z}(3.5) = 0.6642 - .06222 * 3.5 = 0.4464$$

$$\hat{Y}(3.5) = 2.2400$$

prediction interval for Z : $[-0.0517, 0.9446]$.

The prediction interval for Y would be the image of the prediction interval for Z under the function $x \rightarrow 1/x$. This doesn't make much sense because the prediction interval for Z contains 0 (the image of this is actually the union of two infinite intervals, which is not what we have in mind).

7. $Y = a + bx + cx^2$

x	Y	x^2	xY	Y^2	x^3	x^2Y	x^4
0.5	1.3	0.25	0.65	1.69	0.125	0.325	0.0625
1.0	2.7	1.00	2.70	7.29	1.000	2.700	1.0000
2.0	1.5	4.00	3.00	2.25	8.000	6.000	8.0000
2.5	2.0	6.25	5.00	4.00	15.625	12.500	39.0625
3.0	2.2	9.00	6.60	4.84	27.000	19.800	81.0000
9.0	9.7	20.50	17.95	20.07	51.75	41.325	129.125

$$\begin{array}{ccc|ccc} 5\hat{a}+ & 9\hat{b}+ & 20.5\hat{c} = & & 9.7 \\ 9\hat{a}+ & 20.5\hat{b}+ & 51.75\hat{c} = & & 17.95 \\ 20.5\hat{a}+ & 51.75\hat{b} & 129.125\hat{c} = & & 41.325 \end{array}$$

$$\hat{a} = 1.782$$

$$\hat{b} = 0.0376$$

$$\hat{c} = 0.0221$$

$$\hat{\sigma}^2 = .598$$

$$\hat{Y}(3.5) = 1.782 + 0.0376 * 3.5 + 0.0221 * 3.5^2 = 2.184$$

The variance of \hat{Y} is $K\sigma^2$, where

$$K = A + 3.5B + 3.5^2C$$

and A, B, C are the solutions of

$$\begin{array}{ccc|ccc} 5\hat{A}+ & 9\hat{B}+ & 20.5\hat{C} = & & 1 \\ 9\hat{A}+ & 20.5\hat{B}+ & 51.75\hat{C} = & & 3.5 \\ 20.5\hat{A}+ & 51.75\hat{B} & 129.125\hat{C} = & & 3.5^2 \end{array}$$

$K = .03558$, so the variance for the prediction is $1.03558\sigma^2$, and we get the prediction interval: $[-1.202, 5.570]$.

8. $[-0.509, 0.736]$.

9. $[-0.400, 3.870]$.

10. $\hat{y}(1.5) = 1.906$

confidence interval: $[0.943, 2.869]$

B.6 Problems from Chapter 6

1.

$$\chi^2 = \frac{(95 - 100)^2}{100} + \frac{(111 - 100)^2}{100} + \dots + \frac{(113 - 100)^2}{100} = 6.16 < 11.070 = \chi_{5;0.95}^2$$

The die is fair.

2.

i	class	Y_i	np_i	$(Y_i - np_i)^2/np_i$
1	$(-\infty, 0.524]$	3	5	0.8
2	$(0.524, 1]$	5	5	0.0
3	$(1, 1, 476]$	6	5	0.2
4	$(1.476, \infty)$	6	5	0.2
Σ	-	20	20	1.2

$$\chi^2 = 1.2 < 7.815 = \chi_{3;.95}^2$$

We accept the null hypothesis.

3. $\hat{\mu} = \bar{X} = 1.156$, $\hat{\sigma}^2 = S^2 = .4844$

i	class	Y_i	np_i	$(Y_i - np_i)^2/np_i$
1	$(-\infty, 0.532]$	3	5	0.8
2	$(0.532, 1]$	5	5	0.0
3	$(1, 1, 468]$	6	5	0.2
4	$(1.468, \infty)$	6	5	0.2
Σ	-	20	20	1.2

$$\chi^2 = 1.2 < 3.841 = \chi_{1;.95}^2$$

We accept the null hypothesis.

4. $\hat{\lambda} = 1/\bar{X} = 0.864$

i	class	Y_i	np_i	$(Y_i - np_i)^2/np_i$
1	(0, 0.333]	3	5	0.8
2	(0.333, 0.802]	4	5	0.2
3	(0.802, 1.604]	7	5	0.8
4	(1.604, ∞)	6	5	0.2
Σ	-	20	20	2.0

$$\chi^2 = 2.0 < 5.991 = \chi_{2;.95}^2$$

We accept the null hypothesis.

5. $\hat{\theta} = \max_i(X_i) = 2.70$

i	class	Y_i	np_i	$(Y_i - np_i)^2/np_i$
1	(0, 0.675]	4	5	0.2
2	(0.675, 1.35]	10	5	5.0
3	(1.35, 2.025]	4	5	0.2
4	(2.025, 2.70]	2	5	1.8
Σ	-	20	20	7.2

$$\chi^2 = 7.2 > 5.991 = \chi_{2;.95}^2$$

We reject the null hypothesis.

6. $\hat{\lambda} = \bar{X} = 2.675$

i	class	Y_i	np_i	$(Y_i - np_i)^2/np_i$
1	0-1	9	10.129	0.078
2	2	8	9.861	0.351
3	3	13	8.793	2.013
4	4	5	5.880	0.132
5	5- ∞	5	5.337	0.021
Σ	-	40	40.000	2.595

$$\chi^2 = 2.595 < \chi_{3;.95}^2 = 7.815$$

We accept H_0 .

7.

i	class	Y_i	np_i	$(Y_i - np_i)^2/np_i$
1	0-1	9	4.375	4.889
2	2	8	9.375	0.202
3	3	13	12.5	0.020
4	4	5	9.375	2.042
5	5-6	5	4.375	0.089
Σ	-	40	40.000	7.422

$$\chi^2 = 7.422 < \chi_{4;.95}^2 = 9.488$$

We accept H_0 .

8. $\hat{p} = \bar{X}/6 = 0.446$

i	class	Y_i	np_i	$(Y_i - np_i)^2/np_i$
1	0-2	17	9.937	5.020
2	3	13	12.065	0.072
3	4	5	11.248	3.471
4	5-6	5	6.750	0.454
\sum	-	40	40.000	9.017

$$\chi^2 = 9.017 > \chi_{2,.95}^2 = 5.991$$

We reject H_0 .

9. By symmetry, for $0 < z_i < 1$,

$$\begin{aligned} \mathbb{P}(Z_1 \leq z_1, \dots, Z_{n-1} \leq z_{n-1} | M = m) &= \\ \mathbb{P}(Y_1 \leq mz_1, \dots, Y_{n-1} \leq mz_{n-1} | M = m) &= \\ \mathbb{P}(Y_1 \leq mz_1, \dots, Y_{n-1} \leq mz_{n-1} | X_1 < m, \dots, X_{n-1} < m, X_n = m) &= \\ \mathbb{P}(X_1 \leq mz_1, \dots, X_{n-1} \leq y_{n-1} | X_1 < m, \dots, X_{n-1} < m) &= z_1 * \dots * z_{n-1}. \end{aligned}$$

Thus, Z_1, \dots, Z_{n-1} are independent, uniformly distributed on $[0, 1]$.

Now, in our usual χ^2 test situation, we would divide the interval $[0, M]$ into k intervals of length M/k . Denote the number of X 's in the i th interval by U_i . Then the χ^2 statistic is

$$\chi_X^2 = \frac{k}{n} \sum_{i=1}^k (U_i - n/k)^2.$$

On the other hand, if we denote the number of Z_i 's in the interval $(\frac{i-1}{k}, \frac{i}{k}]$ by V_i , then

$$U_i = \begin{cases} V_i & \text{if } i < k \\ V_i + 1 & \text{if } i = k. \end{cases}$$

Furthermore, the χ^2 statistic for the sample Z_1, \dots, Z_{n-1} is

$$\chi_Z^2 = \frac{k}{n-1} \sum_{i=1}^k (V_i - (n-1)/k)^2,$$

and a little calculation yields

$$\chi_X^2 = \frac{n-1}{n} \chi_Z^2 + \frac{k}{n} (V_k - \frac{n-1}{k}) + \frac{k-1}{n}.$$

For large n , this implies that χ_X^2 and χ_Z^2 have the same limiting distribution. We know that χ_Z^2 has an approximate χ^2 distribution with $k-1$ degrees of freedom, so the same is true for χ_X^2 .

10. $\chi^2 = 170.7 > 26.269 = \chi_{16,.95}^2$

We reject the null hypothesis of independence.

11.

$$\chi^2 = \frac{n(Z_{11}Z_{22} - Z_{12}Z_{21})^2}{Z_{1.}Z_{.2}Z_{.1}Z_{.2}}$$

12. $\chi^2 = 12.090 > 3.841 = \chi_{1,.95}^2$

We reject the hypothesis of independence.

13. $\chi^2 = 8.4727 < 15.507 = \chi_{8,.95}^2$

We accept the null hypothesis that the three professors are equal.

B.7 Problems from Chapter 7

1.

i	$X_{n:i}$	$F(X_{n:i})$	$F_n(X_{n:i} - 0)$	$F_n(X_{n:i})$	Δ_i
1	0.05	0.090	0.00	0.05	0.090
2	0.16	0.117	0.05	0.10	0.067
3	0.21	0.133	0.10	0.15	0.033
4	0.58	0.278	0.15	0.20	0.128
5	0.70	0.337	0.20	0.25	0.137
6	0.70	0.337	0.25	0.30	0.087
7	0.74	0.356	0.30	0.35	0.056
8	1.00	0.500	0.35	0.40	<u>0.150</u>
9	1.02	0.512	0.40	0.45	0.112
10	1.06	0.532	0.45	0.50	0.087
11	1.12	0.567	0.50	0.55	0.067
12	1.16	0.587	0.55	0.60	0.037
13	1.25	0.637	0.60	0.65	0.037
14	1.30	0.663	0.65	0.70	0.037
15	1.71	0.841	0.70	0.75	0.141
16	1.75	0.855	0.75	0.80	0.105
17	1.84	0.883	0.80	0.85	0.083
18	1.95	0.918	0.85	0.90	0.068
19	2.13	0.945	0.90	0.95	0.045
20	2.70	0.992	0.95	1.00	0.042

$D_{20} = 0.150 < 0.2941$. We accept H_0 .

2.

i	$X_{n:i}$	$F_n(X_{n:i} - 0)$	$F_n(X_{n:i})$	Δ_i
1	0.23	0.000	0.125	<u>0.230</u>
2	0.31	0.125	0.250	0.185
3	0.45	0.250	0.375	0.200
4	0.49	0.375	0.500	0.115
5	0.67	0.500	0.625	0.170
6	0.80	0.625	0.750	0.175
7	0.87	0.750	0.875	0.120
8	0.93	0.875	1.000	0.070

$D_8 = 0.23 < 0.4543$: We accept H_0 .

3. As for the one-sample test, we first look for an easier way to write down the test statistic: it can easily be shown that

$$\sup\{|G_m(x) - F_n(x)|\} = \sup\{|G_m(x) - F_n(x)| : x \in \{X_i, i \leq n\} \cup \{Y_j, j \leq m\}\}.$$

Thus, we use the following table:

x	$F_n(x)$	$G_m(x)$	$\Delta(x)$
0.6	0.111	0.091	0.020
0.7	0.222	0.273	0.051
0.8	0.333	0.364	0.031
0.9	0.556	0.364	0.192
1.0	0.667	0.364	0.303
1.1	0.889	0.455	<u>0.434</u>
1.3	0.889	0.636	0.253
1.4	1.000	0.818	0.182
1.6	1.000	0.909	0.091
1.8	1.000	1.000	0.000

$$D_{nm} = 0.434 < .5959 = D_{9,11,.95}$$

so we accept the null hypothesis.

4. Let S be the set of all permutations of $\{1, \dots, n\}$. Then

$$\begin{aligned} \mathbb{P}(D_n < x) &= \mathbb{P}\left(\left|\frac{k}{n} - U_{n:k}\right| < x, \left|\frac{k-1}{n} - U_{n:k}\right| < x\right) = \mathbb{P}\left(\frac{k}{n} - x < U_{n:k} < \frac{k-1}{n} + x\right) = \\ &= \sum_{s \in S} \mathbb{P}\left(\frac{k}{n} - x < U_{s(k)} < \frac{k-1}{n} + x, U_{n:k} = U_{s(k)}\right) = \\ &= \sum_{s \in S} \mathbb{P}\left(\frac{k}{n} - x < U_{s(k)} < \frac{k-1}{n} + x, U_{s(1)} \leq U_{s(2)} \leq \dots \leq U_{s(n)}\right) = \\ &= n! \mathbb{P}\left(\frac{k}{n} - x < U_k < \frac{k-1}{n} + x, U_1 \leq U_2 \leq \dots \leq U_n\right). \end{aligned}$$

5.

$$\mathbb{P}(D_2 \leq x) = \mathbb{P}\left(\frac{1}{2} - x \leq U_1 \leq x, 1 - x \leq U_2 \leq \frac{1}{2} + x, U_1 < U_2\right)$$

If $x < \frac{1}{4}$, $\mathbb{P}(D_2 \leq x) = 0$.

If $\frac{1}{4} \leq x < \frac{1}{2}$,

$$\mathbb{P}(D_2 \leq x) = 2\mathbb{P}\left(\frac{1}{2} - x \leq U_1 \leq x, 1 - x \leq U_2 \leq \frac{1}{2} + x\right) = 2\left(\frac{1}{2} - 2x\right)^2.$$

If $\frac{1}{2} \leq x \leq 1$,

$$\begin{aligned} \mathbb{P}(D_2 \leq x) &= \mathbb{P}(U_{2:1} \leq x, U_{2:2} \geq 1 - x) = \\ &= 1 - \mathbb{P}(U_{2:1} \geq x) - \mathbb{P}(U_{2:2} \leq 1 - x) + \mathbb{P}(U_{1:1} \geq x, U_{2:2} \leq 1 - x) = \\ &= 1 - 2(1 - x)^2 + 0 = 1 - 2(1 - x)^2. \end{aligned}$$

6.

$$\eta_{mm}(x) = F_n(x) - G_m(x) = \sum_{i=1}^n \sum_{j=1}^m (V_i(x) - W_j(x)),$$

where

$$V_i = \begin{cases} 1 & \text{if } X_i \leq x, \\ 0 & \text{otherwise,} \end{cases}$$

$$W_j = \begin{cases} 1 & \text{if } Y_j \leq x, \\ 0 & \text{otherwise,} \end{cases}$$

if $i_1 \neq i_2$, $V_{i_1}(x)$ and $V_{i_2}(y)$ are independent; also, if $j_1 \neq j_2$, $W_{j_1}(x)$ and $W_{j_2}(y)$ are independent; furthermore, $V_i(a)$ and $W_j(b)$ are independent for any choice of i, j, a , and b . For the remaining cases, we get the covariances

$$\mathbf{Cov}(V_i(x), V_i(y)) = \mathbf{Cov}(W_j(x), W_j(y)) = x \wedge y - xy,$$

where we have put

$$x \wedge y = \min(x, y).$$

Thus

$$\begin{aligned} \mathbf{Cov}(\eta_{mm}(x), \eta_{mm}(y)) &= \sum_{i=1}^n \frac{1}{n^2} \mathbf{Cov}(V_i(x), V_i(y)) + \sum_{j=1}^m \frac{1}{m^2} \mathbf{Cov}(W_j(x), W_j(y)) = \\ &= \frac{1}{n}(x \wedge y - xy) + \frac{1}{m}(x \wedge y - xy) = \frac{n+m}{nm}(x \wedge y - xy). \end{aligned}$$

Letting $x = y$ gives the variance of $\eta_{mn}(x)$.

Anhang C

German translations of technical terms

accept annehmen

alternative hypothesis Alternativhypothese, Gegenhypothese, Eins-Hypothese (brr)

analysis of variance Varianzanalyse

composite hypothesis zusammengesetzte Hypothese

confidence interval Konfidenzintervall

confidence region Konfidenzbereich

consistent konsistent

coverage probability Überdeckungswahrscheinlichkeit

degrees of freedom Freiheitsgrade

descriptive statistics beschreibende Statistik

distribution Verteilung

efficient effizient

empirical distribution function empirische Verteilungsfunktion

estimation Schätzung

estimator Schätzer

frequency Häufigkeit

goodness-of-fit test Anpassungstest

hypothesis Hypothese

inductive statistics schließende Statistik

interquartile range Quartilsdistanz

level of significance Signifikanzniveau (kurz Niveaue)

likelihood Likelihood (bei Leuten, die Fremdwörter nicht so sehr mögen: "Plausibilität")

likelihood ratio test Likelihoodquotiententest

linear regression lineare Regression

mathematical statistics Mathematische Statistik
maximum likelihood estimator Maximum Likelihood Schätzer (auch: plausibler Schätzer)
median Median
method of moments Momentenmethode
multiple linear regression mehrfache lineare Regression
non-parametric statistics nichtparametrische Statistik
null hypothesis Nullhypothese
one-sided hypothesis einseitige Hypothese
one-way analysis of variance einfache Varianzanalyse
operation characteristic Operationscharakteristik
parameter Parameter
parametric family parametrische Familie
parametric statistics parametrische Statistik
population Grundgesamtheit, Population
power function Macht, Mächtigkeit, Schärfe
prediction interval Vorhersageintervall
quantile Quantil, Fraktil
quartile Quartil
randomized randomisiert
reject ablehnen, verwerfen
rejection region Verwerfungsbereich
sample mean Stichprobenmittel
sample variance Stichprobenvarianz
sample Stichprobe
simple hypothesis einfache Hypothese
simple linear regression einfache lineare Regression
statistic Statistik
statistics Statistik
sufficient statistic suffiziente Statistik
test Test, Signifikanztest
two-sided hypothesis zweiseitige Hypothese
two-way analysis of variance doppelte Varianzanalyse
unbiased estimator erwartungstreuer Schätzer, unverzerrter Schätzer
unbiased test unverfälschter Test
uniform distribution Gleichverteilung

Index

- alternative hypothesis, 18
- analysis of variance, 28
 - one-way, 32
 - two-way, 33
- chi-square distribution, 14
- Chi-square statistic, 41
- Chi-square test
 - goodness of fit, 42
 - unknown parameters, 42
 - homogeneity, 43
 - independence, 42
- composite hypothesis, 17
- confidence interval
 - linear regression, 37
 - normal distribution, 13
 - mean, 13, 14
 - variance, 14
 - proportions, 14
 - regression
 - multiple, 38
- confidence regions, 13
- consistent, 9
 - strongly, 9
 - weakly, 9
- coverage probability, 12
- Crámer-Rao theorem, 11
- degrees of freedom, 29
- descriptive statistics, 3
- distribution
 - chi-square, 14
 - F, 22
 - Student, 14
 - t, 14
- empirical distribution function, 4
- error
 - first kind, 18
 - second kind, 18
- estimator, 9
 - asymptotically unbiased, 9
 - consistent, 9
 - efficient, 9
 - least squares, 33
 - maximum likelihood, 10
 - strongly consistent, 9
 - unbiased, 9
 - weakly consistent, 9
- exponential family, 25
- family
 - parametric, 4
- first kind error, 18
- Fischer-Cochran theorem, 30
- Fisher information, 11
- hypothesis, 17
 - alternative, 18
 - composite, 17
 - one-sided, 17
 - two-sided, 17
 - null, 18
 - one-sided, 17
 - parametric, 17
 - simple, 17
 - two-sided, 17
- inductive statistics, 3
- interquartile range, 6
- IQR, 6
- least squares estimator, 33
- least squares method, 33
- level of significance, 18
- likelihood function, 5
- linear regression, 36
 - confidence interval, 37
 - multiple, 37
 - confidence interval, 38
 - prediction interval, 38
 - prediction interval, 37
 - simple, 36
- mathematical statistics, 3
- maximum likelihood estimator, 10
- maximum likelihood method, 10
- median, 60
- method of moments, 10
- monotone likelihood ratios, 22
- Neyman-Pearson test, 19
- Neyman-Pearson theorem, 19
- nonparametric statistics, 4
- normal distribution
 - mean, 20

- null hypothesis, 18
- one-sided hypothesis, 17
- operation characteristic, 18
- order statistics, 6
- parametric family, 4
- parametric hypothesis, 17
- parametric statistics, 4
- population, 3
- power function, 18
- prediction interval, 37
 - multiple regression, 38
- quantile, 6
- randomized test, 18
- rejection region, 18
- sample, 3
- sample mean, 5
- sample variance, 5
- second kind error, 18
- simple hypothesis, 17
- statistic, 5
 - order, 6
 - sufficient, 5
- statistics
 - descriptive, 3
 - inductive, 3
 - mathematical, 3
 - nonparametric, 4
 - parametric, 4
- strongly consistent, 9
- Student's distribution, 14
- sufficient statistic, 5
- sum of squares, 29
 - degrees of freedom, 29
- t distribution, 14
- test, 17, 18
 - Kolmogorov-Smirnov, 45
 - asymptotic distribution, 46
 - one-sample, 45
 - two-sample, 46
 - level of significance, 18
 - Neyman-Pearson, 19
 - non-randomized, 18
 - normal distribution, 20
 - two-sample, 21
 - variance, 20
 - randomized, 18
 - t-test, 20
 - unbiased, 23
- theorem
 - Cramer-Rao, 11
 - Fischer-Cochran, 30
 - Glivenko-Cantelli, 4
 - Karlin-Rubin, 23
 - Neyman-Pearson, 19
 - two-sided hypothesis, 17
 - unbiased test, 23
 - weakly consistent, 9