

Lecture Notes on Acceptance Testing

Yakov Ben-Haim
Faculty of Mechanical Engineering
Technion — Israel Institute of Technology
Haifa 32000 Israel
yakov@technion.ac.il
<http://www.technion.ac.il/yakov>

Primary source material: William W. Hines, and Douglas C. Montgomery, *Probability and Statistics in Engineering and Management Science*. 3rd ed. Sections 11-1–11-4, 11-11, 11-12.

Notes to the Student:

- These lecture notes are not a substitute for the thorough study of books. These notes are no more than an aid in following the lectures.

- Section 11 contains review exercises that will assist the student to master the material in the lecture and are highly recommended for review and self-study. The student is directed to the review exercises at selected places in the notes. They are not homework problems, and they do not entitle the student to extra credit.

Contents

1	Thickness Measurement: Testing a Sample Mean	3
2	Sequential Sampling: Testing a Mean	9
3	Matching Two Dimensions: Testing Two Sample Means	10
4	Engine Warming: A χ^2 Test	13
5	Failure Rate: Poisson Distribution and the χ^2 Test	15
6	χ^2 Test of Independence in a 2-way Table	17
7	Acceptance Sampling of a Large Population	20
8	Sample Size for Detecting a Change: Threshold Tests	24
8.1	Sample Size and Error Probabilities	24
8.2	Uncertain Effect Size and Variance	28
8.3	Uncertain Sample PDF	31
8.3.1	Background	31
8.3.2	Info-gap Approach to Determining the Sample Size	32
8.3.3	An Approximate Robustness for Small Effect Size	34
9	Tests of the Mean with Distributional Uncertainty	35
9.1	Distributional Uncertainty	35
9.2	Info-Gap Representations of Distributional Uncertainty	37
9.3	Robustness Functions with CDF Uncertainty	38
9.3.1	Binary Test: Formulation	38
9.3.2	Robustnesses for Type I Errors	39
9.3.3	Robustnesses for Type II Errors	40
9.3.4	Decisions and Judgments	41

9.4	Robustness Functions with PDF Uncertainty: Definitions	44
9.4.1	Binary Test: Formulation	44
9.4.2	Robustness for Falsely Rejecting H_0 . (Type I Error.)	45
9.4.3	Robustness for Falsely Accepting H_0 . (Type II Error.)	46
9.5	Robustness with PDF Uncertainty: Numerical Examples	47
9.5.1	Robustness for Type-I Error	47
9.5.2	Robustness for Type-II Error	47
10	Accelerated Lifetime Testing: Simple Case	48
10.1	Formulation	48
10.2	Uncertainty and Robustness	48
10.3	Evaluating the Robustness	49
11	Review Exercises	51

1 Thickness Measurement: Testing a Sample Mean

¶ In the first few sections we illustrate various **statistical hypothesis tests** for acceptance testing.

¶ Suppose we measure the thickness of a plate at N widely separated points: x_1, \dots, x_N .

These measurements differ from one another due to random measurement error, as well as fluctuations in the local thickness.

How do we use these measurements to decide if the plate “really” has thickness T ?

What does “really has thickness T ” mean? Perhaps: $E(T) = \mu$.

¶ A **random sample** is a set of independent measurements made on the same population. That is, a random sample is a set of independent and identically distributed (i.i.d.) random variables.

¶ The **sample mean** of a random sample is defined as:

$$\bar{x} = \frac{1}{N} \sum_{i=1}^N x_i \quad (1)$$

Theorem 1. If a random sample of size N is taken from a population with mean μ and variance σ^2 , then the sample mean \bar{x} has mean μ and variance σ^2/N . That is:

$$E(\bar{x}) = \mu \quad (2)$$

$$\text{var}(\bar{x}) = E([\bar{x} - E(\bar{x})]^2) = \frac{\sigma^2}{N} \quad (3)$$

Proof. First consider eq.(2):

$$E(\bar{x}) = E\left(\frac{1}{N} \sum_{i=1}^N x_i\right) = \int_{x_1 \dots x_N} p(x_1, \dots, x_N) \frac{1}{N} \sum_{i=1}^N x_i dx_1 \cdots x_N \quad (4)$$

Because the measurements of the random sample are independent:

$$p(x_1, \dots, x_N) = \prod_{i=1}^N p(x_i) \quad (5)$$

Thus:

$$E(\bar{x}) = \frac{1}{N} \sum_{i=1}^N \int p(x_i) x_i dx_i = \frac{1}{N} \sum_{i=1}^N \mu = \mu \quad (6)$$

Which completes the proof of eq.(2). Note that this is independent of $p(x_i)$.

Now consider eq.(3):

$$\text{var}(\bar{x}) = E[(\bar{x} - \mu)^2] = E\left[\left(\frac{1}{N} \sum_{i=1}^N (x_i - \mu)\right)^2\right] \quad (7)$$

$$= \frac{1}{N^2} \sum_i \sum_j E[(x_i - \mu)(x_j - \mu)] \quad (8)$$

$$= \frac{1}{N^2} \sum_i E(x_i - \mu)^2 = \frac{\sigma^2}{N} \quad (9)$$

Note that this is independent of $p(x_i)$.

Review exercise 1, p. 51.

Review exercise 2, p. 51.

Theorem 2. If a random sample is taken from a normal population, then the sample mean is normal. ■

Combining the last two theorems we can assert:

$$x_i \sim \mathcal{N}(\mu, \sigma^2) \implies \bar{x} \sim \mathcal{N}(\mu, \sigma^2/N) \quad (10)$$

Review exercise 3, p. 51.

“Theorem” 3. An approximate statement of the **central limit theorem**: The sample mean of a random sample will be approximately normal for large sample size ($N > 30$, rough number). ■

Review exercise 4, p. 51.

¶ So, let us suppose that the thickness measurement is normally distributed, or that the number of measurements is large. Thus:

$$\bar{x} \sim \mathcal{N}(\mu, \sigma^2/N) \quad (11)$$

where:

μ = true mean thickness.

σ^2 = variance of the thickness measurements.

N = sample size.

Also, assume that we **know the value of σ^2** .

¶ How do we decide whether:

- The true thickness of the plate is T ?
- To accept or reject the plate?

We use an **hypothesis test**.

¶ The **null hypothesis**:

$$H_0 : \quad \mu = T \quad (12)$$

T = desired thickness: a known value.

μ = true thickness: an unknown value.

¶ The **alternative hypothesis**:

$$H_1 : \quad \mu \neq T \quad (13)$$

This is a **two-tailed test**. The test would be **one-tailed** if the alternative hypothesis were:

$$H_1 : \quad \mu > T \quad (14)$$

Or:

$$H_1 : \quad \mu < T \quad (15)$$

Review exercise 5, p. 51.

¶ **Level of confidence, α :**

- Probability of obtaining a result at least as extreme as the observed result, conditioned upon H_0 .
- Probability of rejecting H_0 erroneously.

For the two-tailed test in question (see fig. 1 on p.5):

$$\alpha = \text{Prob} \left(|\bar{x} - T| \geq |\bar{x}_o - T| \mid H_0 \right) \quad (16)$$

\bar{x} = the random variable “sample mean”.

\bar{x}_o = the observed value of the random variable \bar{x} .

Review exercise 6, p. 51.

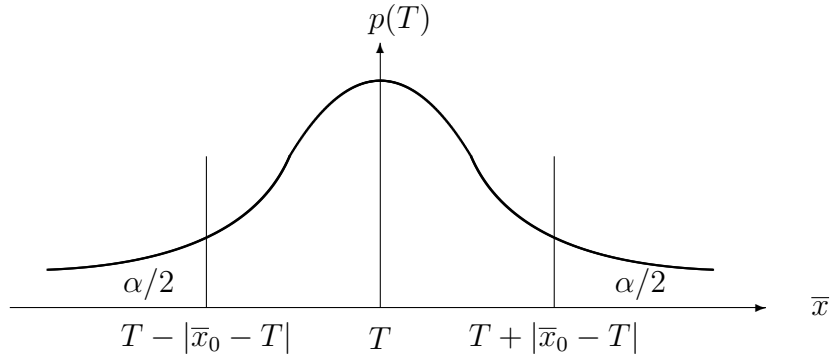


Figure 1: Sketch of probability density, illustrating the level of confidence in eq.(16).

¶ Interpreting the level of confidence:

If α is small: reject H_0 .

If α is large: accept H_0 .

¶ How to evaluate the level of confidence?

Conditioned upon H_0 , we can assert:

$$\bar{x} \sim \mathcal{N}(T, \sigma^2/N) \quad (17)$$

\bar{x} can be standardized as:

$$z = \frac{\bar{x} - T}{\sigma/\sqrt{N}} \sim \mathcal{N}(0, 1) \quad (18)$$

Review exercise 7, p. 51.

Now the level of confidence can be expressed as (see fig. 2 on p.6):

$$\alpha = \text{Prob} \left(|\bar{x} - T| \geq |\bar{x}_o - T| \mid H_0 \right) \quad (19)$$

$$= \text{Prob} \left(\frac{|\bar{x} - T|}{\sigma/\sqrt{N}} \geq \frac{|\bar{x}_o - T|}{\sigma/\sqrt{N}} \mid H_0 \right) \quad (20)$$

$$= 2 \left[1 - \Phi \left(\frac{|\bar{x}_o - T|}{\sigma/\sqrt{N}} \right) \right] \quad (21)$$

where $\Phi(\cdot)$ is the cdf of the standard normal distribution. Define:

$$z_o = \frac{|\bar{x}_o - T|}{\sigma/\sqrt{N}} \quad (22)$$

$$\alpha = 2[1 - \Phi(z_o)] \quad (23)$$

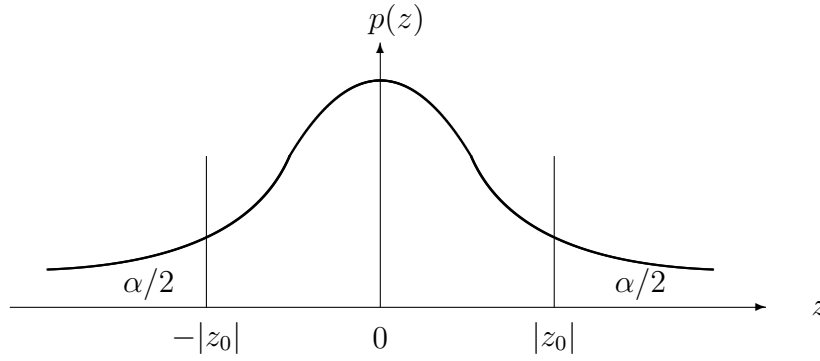


Figure 2: Sketch of probability density illustrating level of confidence in eq.(21).

¶ Numerical example.

Measurements: 1.3, 1.2, 1.4, 1.3, 1.1

Desired thickness: $T = 1.36$.

Known variance: $\sigma^2 = 0.01 \implies \sigma = 0.1$.

Hence: $N = 5$ and $\bar{x}_o = 1.26$.

$$z_o = \frac{1.26 - 1.36}{0.1/\sqrt{5}} = -2.236 \quad (24)$$

$$\alpha = \text{Prob} \left(|z| \geq |z_o| \mid H_0 \right) = 2 [1 - \Phi(|z_o|)] = 0.024 \quad (25)$$

- Note: $\Phi(|z_o|) = 0.988$.

• So, the probability of getting a value as large or larger than the observed standardized sample mean is 0.024.

- This is rather small, so we tend to reject H_0 .
If so, then we reject H_0 at the 0.024 level of confidence.
- Similarly, $0.024 =$ probability of falsely rejecting H_0 .

¶ We tested H_0 on p. 4 under the assumption that σ^2 is known.
What do we do if σ^2 is **unknown**?

¶ Two cases:
 N large.
 N small.

¶ **Case 1: N large.**

If $N \geq 25$ (rough number), then s^2 , the sample variance, is a good estimate of the true population variance.

$$s^2 = \frac{1}{N-1} \sum_{i=1}^N (x_i - \bar{x})^2 \quad (26)$$

We can now assume, conditioned on H_0 , that:

$$t = \frac{\bar{x} - T}{s/\sqrt{N}} \sim \mathcal{N}(0, 1) \quad (27)$$

This assumption would be precise if $s^2 = \sigma^2$.
Now we proceed with the test as before.

Review exercise 8, p. 51.

¶ **Case 2: N small.**

If $N < 25$ (rough number), then the sample variance s^2 is not a good estimate of the population variance.

The statistic:

$$t = \frac{\bar{x} - T}{s/\sqrt{N}} \quad (28)$$

is broader than $\mathcal{N}(0, 1)$ since both \bar{x} and s display variation.

This is a **t statistic** with $N - 1$ degrees of freedom (dofs).

¶ We repeat the numerical example, without knowledge of σ^2 .

Measurements: 1.3, 1.2, 1.4, 1.3, 1.1

Desired thickness: $T = 1.36$.

$N = 5$ and $\bar{x}_o = 1.26$.

Sample variance: $s^2 = 0.013 \implies s = 0.1140$.

The observed t statistic is:

$$t_o = \frac{\bar{x} - T}{s/\sqrt{N}} = \frac{1.26 - 1.36}{0.1140/\sqrt{5}} = -1.961 \quad (29)$$

The dofs: $5 - 1 = 4$.

From the table of the t distribution (transparency AS-p.7.1):

$\alpha =$	0.1	0.05	0.025	0.01
$\nu = 4 :$	1.533	2.132	2.776	3.747

So, with 4 dofs:

The probability of t_4 exceeding 1.533 is 0.1.

The probability of t_4 exceeding 2.132 is 0.05.

Etc.

The level of confidence of this two-tailed test, with $t_o = -1.961$, is:

$$\alpha = \text{Prob} \left(|t| \geq |t_o| \mid H_0 \right) \approx 2 \times 0.07 = 0.14 \quad (30)$$

This is not small, so we cannot reject H_0 .

The probability of falsely rejecting H_0 is 0.14. ■

Review exercise 9, p. 51.

2 Sequential Sampling: Testing a Mean

¶ In the previous example we found that, with 5 measurements, we reject H_0 at 0.14 level of confidence.

This rejection is not very convincing. **Review exercise 10, p. 51.**

If we measured more, we could probably make a better, more confident decision.

How many measurements to make?

One approach is the idea of **sequential sampling**:

Continue adding measurements until the level of confidence is clear cut.

The following table shows an example.

N	x_i	\bar{x}	s^2	s/\sqrt{N}	$ t_o $	α
5	1.3, 1.2, 1.4, 1.3, 1.1	1.26	0.013	0.0510	1.961	$2 \times 0.07 = 0.14$
8	1.1, 1.3, 1.2	1.2375	0.01125	0.0375	3.267	$2 \times 0.007 = 0.014$
11	1.1, 1.2, 1.1	1.2091	0.01091	0.0315	4.7915	< 0.002

Table 1: Data from a sequential test of the mean thickness measurement. (Transparency)

After 11 measurements we can stop:

Our confidence in rejecting H_0 is great.

Review exercise 11, p. 51.

¶ General theory: sequential analysis.¹

¹Abraham Wald, *Sequential Analysis*.

3 Matching Two Dimensions: Testing Two Sample Means

¶ Let us suppose that we are measuring two dimensions, to see if they match. For instance, the inner and outer dimensions of pieces that need to fit together.

These two dimensions are measured repeatedly, with a measuring device which has random errors:

- The random sample of the inner dimension is: x_1, \dots, x_N .
- The random sample of the outer dimension is: y_1, \dots, y_M .

We need not assume that the sample sizes, N and M are the same.

The true mean values of these two samples, which are the true dimensions, are:

μ_1 = true (but unknown) inner dimension.

μ_2 = true (but unknown) outer dimension.

¶ We will assume that these two sets of measurements are **normally distributed**, each with known variance σ^2 .

Will the pieces fit snugly?

How confident are we of the answer?

In other words, we wish to test the null hypothesis:

$$H_0 : \quad \mu_1 = \mu_2 \quad (31)$$

against one of the following alternative hypotheses:

$$H_1 : \quad \mu_1 \neq \mu_2 \quad (32)$$

or:

$$H_1 : \quad \mu_1 > \mu_2 \quad (33)$$

or:

$$H_1 : \quad \mu_1 < \mu_2 \quad (34)$$

We choose an alternative hypothesis depending upon our prior information.

Review exercise 12, p. 51.

Let \bar{x} and \bar{y} be the sample means:

$$\bar{x} = \frac{1}{N} \sum_{i=1}^N x_i, \quad \bar{y} = \frac{1}{M} \sum_{i=1}^M y_i \quad (35)$$

Consider the statistic:

$$\Delta = \bar{x} - \bar{y} \quad (36)$$

What is the mean and variance of Δ , if H_0 holds?

Note that we cannot answer this question under H_1 .

Review exercise 13, p. 51.

$$E(\Delta) = E(\bar{x} - \bar{y}) = E(\bar{x}) - E(\bar{y}) = 0 \quad (37)$$

$$\text{var}(\Delta) = \text{var}(\bar{x} - \bar{y}) = \text{var}(\bar{x}) + \text{var}(\bar{y}) = \underbrace{\frac{\sigma^2}{N} + \frac{\sigma^2}{M}}_{\sigma_\Delta^2} \quad (38)$$

In eq.(38) we have used the fact that \bar{x} and \bar{y} are statistically independent because they are means of separate random samples:

$$\text{var}(\bar{x} - \bar{y}) = E \left[\left(\bar{x} - \bar{y} - [E(\bar{x} - \bar{y})] \right)^2 \right] \quad (39)$$

$$= E \left[\left(\bar{x} - E(\bar{x}) \right)^2 \right] - 2E \left[\left(\bar{x} - E(\bar{x}) \right) \left(\bar{y} - E(\bar{y}) \right) \right] + E \left[\left(\bar{y} - E(\bar{y}) \right)^2 \right] \quad (40)$$

$$= E \left[\left(\bar{x} - E(\bar{x}) \right)^2 \right] + E \left[\left(\bar{y} - E(\bar{y}) \right)^2 \right] \quad (41)$$

$$= \text{var}(\bar{x}) + \text{var}(\bar{y}) \quad (42)$$

How is Δ distributed if H_0 is true?

If we assume both of the following:

- The samples are large.
- The measurements are normally distributed.

Then:

$$\Delta \sim \mathcal{N}(0, \sigma_\Delta^2) \quad (43)$$

Alternatively, if we assume that the measurements are normal but the samples are not large, then:

$$\Delta \sim t_{(N+M-2)} \quad (44)$$

Review exercise 14, p. 51.

Assume both normality and large samples, and define:

$$z = \frac{\Delta}{\sigma_\Delta} \quad (45)$$

If H_0 holds, then:

$$z \sim \mathcal{N}(0, 1) \quad (46)$$

Review exercise 15, p. 51.

Using the two-sided alternative hypothesis of eq.(32) on p.10, we formulate the level of confidence as (fig. 3, p.12):

$$\alpha = \text{Prob} \left(|z| \geq |z_o| \mid H_0 \right) \quad (47)$$

Or, if we use the 1-sided alternative hypothesis of eq.(33), the level of confidence becomes (fig. 4, p.12):

$$\alpha = \text{Prob} \left(z \geq z_o \mid H_0 \right) \quad (48)$$

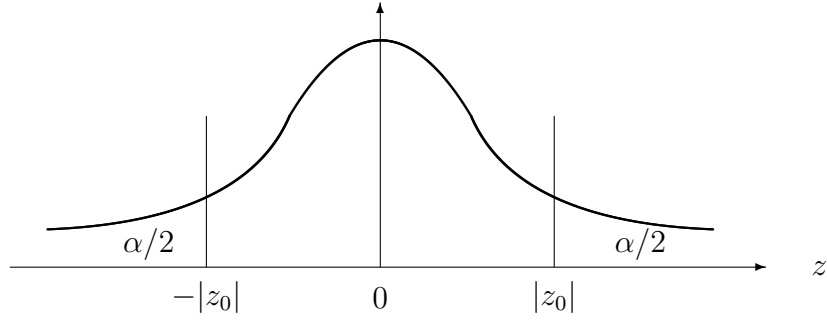


Figure 3: Sketch of probability density illustrating level of confidence in eq.(47).

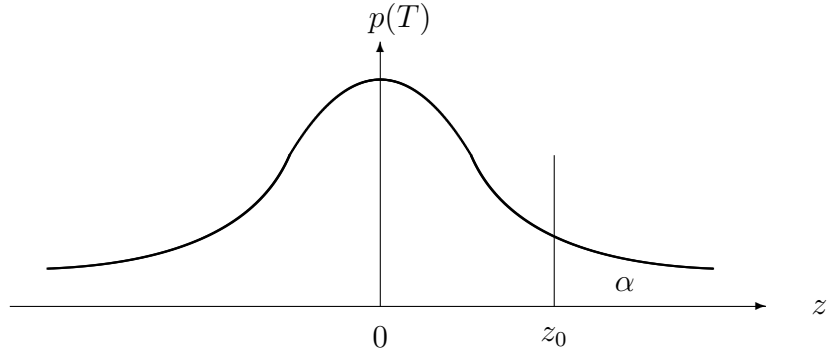


Figure 4: Sketch of probability density illustrating level of confidence in eq.(48).

In the 2-tailed case we obtain:

$$\alpha = 2 [1 - \Phi(z_o)] \quad (49)$$

In the 1-tailed case we obtain:

$$\alpha = 1 - \Phi(z_o) \quad (50)$$

4 Engine Warming: A χ^2 Test

¶ We use an example to introduce the idea of a χ^2 hypothesis test.

¶ The temperature of an operating engine fluctuates between “warm” and “hot”. For normal operation the engine should be “hot” a fraction $p_h = 0.15$ of the time.

Maintenance records since the last overhaul of the engine show that the engine was “warm” at $N_w = 162$ and “hot” at $N_h = 44$ statistically independent sample instants.

$$\frac{N_h}{N_h + N_w} = \frac{44}{206} \approx 0.21 \quad (51)$$

Is the engine operating properly?

Review exercise 16, p. 51.

¶ The χ^2 test for categorical data is suitable for addressing this question.

We formulate the χ^2 test as follows.

Given:

- K types of outcomes of an ‘experiment.’
- n_i outcomes of type i , $i = 1, \dots, K$.
- $N = \sum_{i=1}^K n_i$ = total number of outcomes.

Null Hypothesis:

$$H_0 : \quad p_i = \text{probability of type } i \text{ outcome, } i = 1, \dots, K \quad (52)$$

where p_1, \dots, p_K are known values.

The alternative hypothesis, H_1 , is simply: H_0 is false. That is, at least one of the probabilities in H_0 is wrong.

The χ^2 statistic is:

$$\chi^2 = \sum_{i=1}^K \frac{(n_i - Np_i)^2}{Np_i} \quad (53)$$

Explanation:

- The numerator is a prediction-error for type- i outcomes.
- We expect χ^2 to be small if H_0 is correct.

Theorem 4: If N is large and if H_0 is true, then the statistic in eq.(53) is approximately a χ^2 random variable with $K - 1$ dofs. ■

¶ The χ^2 distribution is shown in transparency AS-p.13.2 for several values of the dof.

¶ How does one calculate the level of confidence, α ?

• Recall:

◦ α is the probability, conditioned upon the null hypothesis, of obtaining a value more extreme than the observed statistic.

◦ α = probability of falsely rejecting H_0 .

Hence, let χ_o^2 be the observed value of the χ^2 statistic of eq.(53). The level of confidence is:

$$\alpha = \text{Prob} \left(\chi^2 \geq \chi_o^2 \mid H_0 \right) \quad (54)$$

$$\begin{aligned} \alpha = \text{'small'} &\implies \chi_o^2 = \text{'extreme'} \implies \text{Reject } H_0. \\ \alpha = \text{'large'} &\implies \chi_o^2 = \text{'not extreme'} \implies \text{Accept } H_0. \end{aligned}$$

‘Small’ and ‘Large’ are understood from the natural calibration of probability: from 0 to 1.

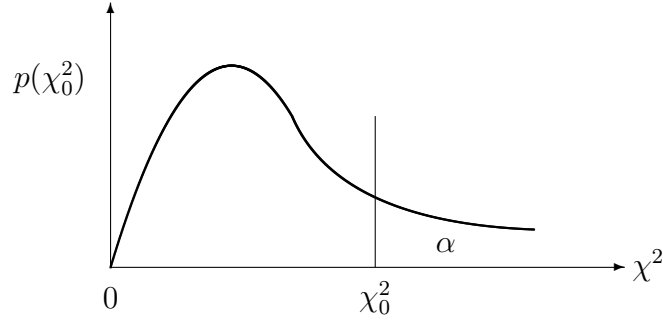


Figure 5: Level of confidence in eq.(54)

¶ Let us apply this theorem to our example.

$K = 2$: 2 possible states: ‘hot’ and ‘warm’.

$N_w = 162$, $N_h = 44$, $N = 206$.

H_0 : $p_h = 0.15$, $p_w = 0.85$.

$$\chi_o^2 = \frac{(162 - 206 \times 0.85)^2}{206 \times 0.85} + \frac{(44 - 206 \times 0.15)^2}{206 \times 0.15} = 6.533 \quad (55)$$

DOFs = $2 - 1 = 1$.

χ^2 table:

$\alpha =$	0.05	0.025	0.01
$\nu = 1 :$	3.84	5.02	6.63

From this we see that $\alpha \approx 0.01$.

So, we reject H_0 at the 1% confidence level.

In other words, the evidence is strong that the engine is running ‘hot’.

Review exercise 17, p. 51.

5 Failure Rate: Poisson Distribution and the χ^2 Test

¶ A computer controlled milling machine operates automatically except when jamming, tool breakage or other failures occur. Under normal circumstances these failures occur at a rate of about 1 or 2 per day. Also, the distribution in time of the failures is **thought to be** a Poisson process: (1) constant average failure rate; (2) events occur independently.

¶ Data have accumulated over a 50-day period for this machine. In this time 75 failures occurred, so the average failure rate is:

$$\lambda = \frac{75 \text{ failures}}{50 \text{ days}} = 1.5 \frac{\text{failures}}{\text{day}} \quad (56)$$

Furthermore, we know how many days had 0, 1, 2 and 3 or more failures:

Failures/day	# days	# failures
0	12	0
1	17	17
2	9	18
3+	12	40
Totals:	50	75

Table 2: Failure data.

¶ We want to test the hypothesis: the distribution over time of failures is described by a Poisson process.

Recall the Poisson distribution:

P_i = probability of exactly i failures in a single day.

$$P_i = \frac{e^{-\lambda} \lambda^i}{i!}, \quad i = 0, 1, 2, \dots \quad (57)$$

¶ We can use a χ^2 test to test this hypothesis.

First define some notation:

N = total number of days.

$p_i = P_i$, the Poisson distribution in H_0 , for $i = 0, \dots, 2$.

$p_3 = \sum_{i=3}^{\infty} P_i$.

n_i = observed number of days with i failures, $i = 0, \dots, 3$.

Np_i = expected number of days with i failures, $i = 0, \dots, 3$.

Now the null hypothesis is:

H_0 : distribution of failures is p_i with $\lambda = 1.5/\text{day}$.

H_1 : H_0 is wrong.

¶ The χ^2 statistic is:

$$\chi_o^2 = \sum_{i=0}^3 \frac{(n_i - Np_i)^2}{Np_i} = 1.695 \quad (58)$$

The DOFs:

$$DOF = \underbrace{4}_{\text{catagories}} - \underbrace{1}_{\text{normalization}} - \underbrace{1}_{\text{estimating } \lambda} = 2 \quad (59)$$

¶ Recall the p th quantile with ν DOFs, $\chi_{(\nu),p}^2$:

$$p = \text{Prob} \left(\chi_{(\nu)}^2 \leq \chi_{(\nu),p}^2 \right) \quad (60)$$

From a χ^2 table, the p -quantiles for $\nu = 2$ are:

$$\begin{array}{rcc} p = & 0.5 & 0.6 \\ \chi_{(2)}^2 = & 1.386 & 1.833 \end{array}$$

From this table we see that the level of confidence, with $\chi_o^2 = 1.695$, is:

$$\alpha = \text{Prob} \left(\chi^2 \geq \chi_o^2 \middle| H_0 \right) \approx 1 - 0.55 = 0.45 \quad (61)$$

This is very large, so we accept H_0 : the distribution of failures is Poisson.

¶ Now consider a **different machine**, for which the data are:

Failures/day	# days	# failures
0	15	0
1	13	13
2	8	16
3+	14	46
Totals:	50	75

Table 3: Failure data.

As before, the average failure rate is:

$$\lambda = 1.5 \frac{\text{failures}}{\text{day}} \quad (62)$$

With these results the observed χ^2 value is:

$$\chi^2 = 5.86 \quad (63)$$

With 2 DOFs, this implies that the level of confidence is:

$$\alpha = \text{Prob} \left(\chi^2 \geq \chi_o^2 \middle| H_0 \right) \approx 0.06 \quad (64)$$

This is rather small, so we reject H_0 : the distribution of failures is not Poisson.

Some factor is causing either:

- Variation in the average failure rate.
- Failure inter-dependence.

In short: clustering of failures.

6 χ^2 Test of Independence in a 2-way Table

¶ Consider the following situation:

Two different systems are used to perform a particular mission. Records indicate the number of failed and successful missions:

	Successful Missions	Failed Missions	Total Missions
System 1	86	35	121
System 2	64	37	101
Total Missions	150	72	

Table 4: Data on successes and failures of two systems.

¶ **Question:**

Is there solid evidence for the contention that system 1 is more reliable than system 2?

In other words, are the rows and columns of the table independent?

That is, by choosing the row (system), do I have sound evidence that I significantly influence the column (success or failure) in which I “fall”?

We can use the χ^2 test to evaluate the evidence.

¶ First we must formulate the **null hypothesis**.

p_{ij} = probability of an event in row i , column j .

$p_{i\bullet}$ = probability of an event in row i .

$p_{\bullet j}$ = probability of an event in column j .

The null hypothesis states:

There is statistical independence between rows and columns.

That is:

$$H_0 : p_{ij} = p_{i\bullet}p_{\bullet j} \quad (65)$$

The alternative hypothesis states that H_0 is false.

¶ Now we can formulate the χ^2 statistic. Define:

n_{ij} = number of events in row i and column j .

r = number of rows.

c = number of columns.

N = total number of events = $\sum_{i=1}^r \sum_{j=1}^c n_{ij}$.

The χ^2 statistic is calculated as:

$$\chi^2 = \sum_{i=1}^r \sum_{j=1}^c \frac{(n_{ij} - Np_{i\bullet}p_{\bullet j})^2}{Np_{i\bullet}p_{\bullet j}} \quad (66)$$

The level of confidence is the probability of χ^2 obtaining a value greater than the observed

value, χ_o^2 , conditioned on H_0 :

$$\alpha = \text{Prob} \left(\chi^2 \geq \chi_o^2 \mid H_0 \right) \quad (67)$$

¶ **A difficulty:** in our example we don't know the values of $p_{i\bullet}$ and $p_{\bullet j}$, so we can't calculate χ_o^2 .

Solution: estimate $p_{i\bullet}$ and $p_{\bullet j}$ from the data: n_{ij} :

$$\hat{p}_{i\bullet} = \frac{n_{i\bullet}}{N}, \quad \hat{p}_{\bullet j} = \frac{n_{\bullet j}}{N} \quad (68)$$

where:

$n_{i\bullet}$ = sum of i th row.

$n_{\bullet j}$ = sum of j th column.

¶ How many DOFs do we have?

DOF = number of categories – number of constraints.

The number of categories is rc .

3 types of constraints:

Type 1. 1 constraint (normalization):

$$\sum_{i=1}^r \sum_{j=1}^c n_{ij} = \sum_{i=1}^r \sum_{j=1}^c N p_{ij} \quad (69)$$

Type 2. $r - 1$ constraints:

Estimate $\hat{p}_{1\bullet}, \dots, \hat{p}_{r-1\bullet}$.

$(\hat{p}_{r\bullet} = 1 - \sum_{i=1}^{r-1} \hat{p}_{i\bullet})$.

Type 3. $c - 1$ constraints:

Estimate $\hat{p}_{\bullet 1}, \dots, \hat{p}_{\bullet c-1}$.

$(\hat{p}_{\bullet c} = 1 - \sum_{j=1}^{c-1} \hat{p}_{\bullet j})$.

So the number of DOFs is:

$$DOF = rc - 1 - (r - 1) - (c - 1) = (r - 1)(c - 1) \quad (70)$$

¶ Now the observed statistic can be calculated as:

$$\chi_o^2 = \sum_{i=1}^r \sum_{j=1}^c \frac{(n_{ij} - N \hat{p}_{i\bullet} \hat{p}_{\bullet j})^2}{N \hat{p}_{i\bullet} \hat{p}_{\bullet j}} \quad (71)$$

¶ In our numerical example we have:

$n_{11} = 86, \quad n_{12} = 35, \quad n_{1\bullet} = 121$

$n_{21} = 64, \quad n_{22} = 37, \quad n_{2\bullet} = 101$

$n_{\bullet 1} = 150, \quad n_{\bullet 2} = 72, \quad N = 222$

Hence the estimated probabilities are:

$$\hat{p}_{1\bullet} = \frac{n_{1\bullet}}{N} = 0.545, \quad \hat{p}_{2\bullet} = \frac{n_{2\bullet}}{N} = 0.455 \quad (72)$$

$$\hat{p}_{\bullet 1} = \frac{n_{\bullet 1}}{N} = 0.676, \quad \hat{p}_{\bullet 2} = \frac{n_{\bullet 2}}{N} = 0.324 \quad (73)$$

The observed statistic becomes:

$$\chi_o^2 = 1.493 \quad (74)$$

And the number of DOFs is:

$$DOF = (2 - 1)(2 - 1) = 1 \quad (75)$$

Is 1.493 large or small? Should we accept or reject H_0 ?

We need to calibrate χ_o^2 using the **quantile** values of the χ^2 distribution with 1 DOF.

Define:

$\chi_{(\nu)}^2 = \chi^2$ random variable with ν DOFs.

$\chi_{(\nu),p}^2 = p$ th quantile of $\chi_{(\nu)}^2$: fig. 6.

$\chi_{(\nu),p}^2$ is defined in:

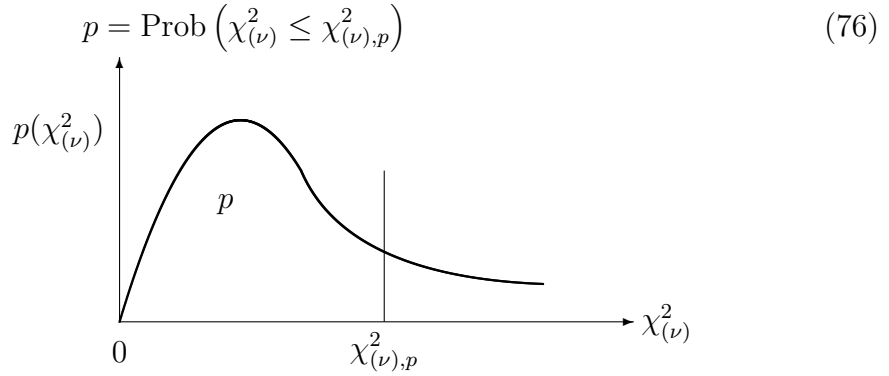


Figure 6: p th quantile of $\chi_{(\nu)}^2$, eq.(76)

From a χ^2 table, the p -quantiles for $\nu = 1$ are:

$$\begin{array}{rcl} p = & 0.75 & 0.80 \\ \chi_{(1),p}^2 = & 1.323 & 1.642 \end{array}$$

From this table we see that:

$$\text{Prob}(\chi_{(1)}^2 \geq 1.323) = 1 - 0.75 = 0.25 \quad (77)$$

$$\text{Prob}(\chi_{(1)}^2 \geq 1.642) = 1 - 0.80 = 0.20 \quad (78)$$

So, since $\chi_o^2 = 1.493$ we see that:

$$\alpha \approx 0.22 \quad (79)$$

We cannot reject H_0 at 0.2 level of confidence.

So we accept H_0 :

- Columns and rows are independent.
- The two systems are not significantly different in reliability.

Review exercise 18, p. 51.

7 Acceptance Sampling of a Large Population

¶ We have a large batch of items, among which an unknown fraction p are defective. We wish to sample this population to decide whether or not to accept the batch.

¶ Define:

N = sample size.

p_a = acceptable proportion of defective items.

p_u = unacceptable proportion of defective items.

$$p_u > p_a \quad (80)$$

p_a and p_u can be interpreted:

- We desire the defective fraction to be no greater than p_a .
- We are willing to live with a defective fraction as large as p_u .

Review exercise 19, p. 51.

p = true but unknown fraction of defective items in the batch.

P_A = probability of accepting the batch.

¶ Our algorithm for accepting or rejecting the batch is:

Accept if and only if:

$$\frac{\text{number of defective items in sample}}{\text{sample size}} \leq p_a \quad (81)$$

¶ Two types of errors can be made:

- I. Acceptance with $p > p_u$. Bad acceptance.
- II. Rejection with $p < p_a$. Bad rejection.

Review exercise 20, p. 52.

¶ How do we calculate P_A , the probability of acceptance?

Assume that $N \ll$ batch size. Thus the sample does not significantly change the composition of the population.

¶ Binomial distribution:

$b(i; p, N)$ = probability of exactly i defectives in a sample of size N .

$$b(i; p, N) = \binom{N}{i} p^i (1-p)^{N-i}, \quad \binom{N}{i} = \frac{N!}{i!(N-i)!} \quad (82)$$

Review exercise 21, p. 52.

¶ Acceptance probability:

$$P_A = \sum_{i=0}^{p_a N} b(i; p, n) \quad (83)$$

$$= \sum_{i=0}^{p_a N} \binom{N}{i} p^i (1-p)^{N-i} \quad (84)$$

¶ **Example.** Suppose:

$$N = 100.$$

$$p_a = 0.01$$

Hence:

$$p_a N = 1$$

$$P_A = \sum_{i=0}^1 b(i; p, 100) = (1-p)^{100} + 100p(1-p)^{99} \quad (85)$$

A plot of P_A vs. p reveals the probabilities of type I and type II errors.

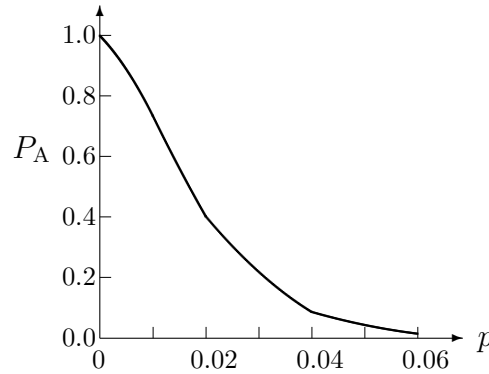


Figure 7: Acceptance probability in eq.(85)

Recall that:

p_a = max “acceptable” fraction of defectives.

p_u = max “tolerable” fraction of defectives.

¶ **Type I errors:**

Acceptance with $p > p_u$. Bad acceptance.

The probability of a type I error is called the **consumer’s risk**:

$$P_I = P_A(p = p_u) \quad (86)$$

¶ **Type II errors:**

Rejection with $p < p_a$. Bad rejection.

The probability of a type II error is called the **producer’s risk**:

$$P_{II} = 1 - P_A(p = p_a) \quad (87)$$

Review exercise 22, p. 52.

When we are designing a sampling scheme, we can evaluate it with 2 pairs of numbers:

$$(p_u, P_I), \quad (p_a, P_{II})$$

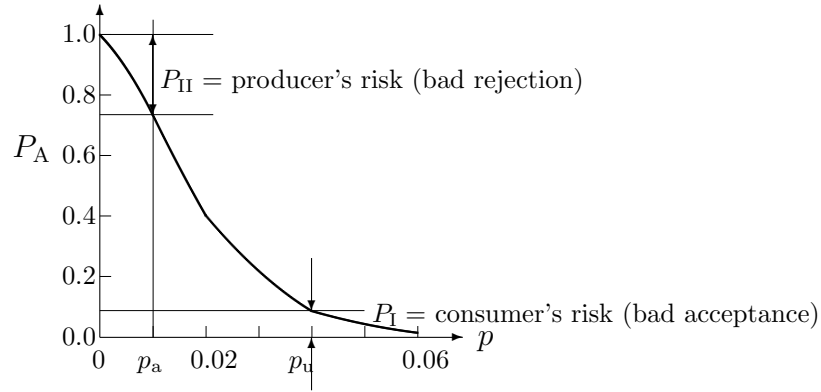


Figure 8: Illustration of type I and type II errors in eqs.(86) and (87) using eq.(85) ($N = 100$). $p_a = 0.01$, $p_u = 0.04$.

¶ **Example.** Consider $P_A(p)$ in eq.(85) on p.21.
Choose $p_a = 0.01$ and $p_u = 0.04$.

With $N = 100$, we have: $p_a N = 1$:

Consumer's risk: $P_I = P_A(p = p_u) = \mathbf{0.0872}$.

Producer's risk: $P_{II} = 1 - P_A(p = p_a) = 1 - 0.736 = \mathbf{0.264}$.

With $N = 200$, we have: $p_a N = 2$:

$$P_A = \sum_{i=0}^2 b(i; p, 100) = (1 - p)^{200} + 200p(1 - p)^{199} + \frac{(200)(199)}{2} p^2 (1 - p)^{198} \quad (88)$$

Consumer's risk: $P_I = P_A(p = p_u) = \mathbf{0.0125}$.

Producer's risk: $P_{II} = 1 - P_A(p = p_a) = 1 - 0.677 = \mathbf{0.323}$.

Review exercise 23, p. 52.

By **increasing** the sample size we:

- **Increased** the producer's risk: $0.264 \rightarrow 0.323$.
- **Decreased** the consumer's risk: $0.0872 \rightarrow 0.0125$.

¶ At $N \rightarrow \infty$ we expect rectangular P_A vs. p :

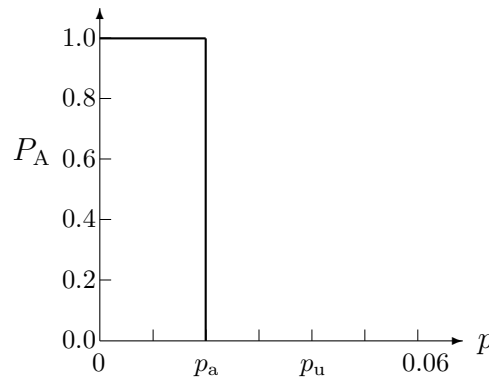


Figure 9: Asymptotical acceptance probability, $N \rightarrow \infty$.

Consumer's risk: $P_I = P_A(p = p_u) = 0$.

Producer's risk: $P_{II} = 1 - P_A(p = p_a) = 0$.

For finite N , $P_A(p)$ oscillates with N , as a result of the discrete, binary nature of the distribution.

Review exercise 24, p. 52.

8 Sample Size for Detecting a Change: Threshold Tests

Source material: David R. Fox, Yakov Ben-Haim, Keith R. Hayes, Michael McCarthy, Brendan Wintle, Piers Dunstan, 2007, An info-gap approach to power and sample size calculations, *Environmetrics*, vol. 18, pp.189–203.

8.1 Sample Size and Error Probabilities

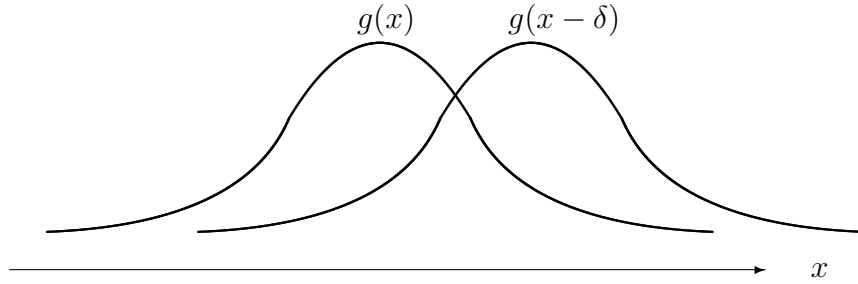


Figure 10: Pdf's $g(x)$ and $g(x - \delta)$ shifted to the right by δ .

¶ The dispute.

- x is the measurement of the output of a system.
- One side argues: x has not changed, its mean is μ_0 and its pdf is $g(x)$.
- The other side argues: x has changed, its mean is $\mu_1 = \mu_0 + \delta$ and its pdf has shifted to $g(x - \delta)$.
- $\delta > 0$ means that $g(x - \delta)$ is shifted to the right of $g(x)$.
- δ is the “effect size”: the shift in the distribution of x .

Review exercise 25, p. 52.

- Null and alternative hypotheses:

$$H_0 : \quad x \sim g(x) \quad (89)$$

$$H_1 : \quad x \sim g(x - \delta) \quad (90)$$

¶ The question: How large a sample is needed to confidently resolve the dispute?

¶ Random sample.

- A random sample of x -values is taken. n = sample size.
- \bar{x} = sample mean. $f_n(\bar{x})$ is the pdf of the sample mean.
- $f_n(\bar{x})$ may differ from $g(x)$. E.g., large n implies $f_n(\bar{x})$ is normal regardless of $g(x)$.
- Null and alternative hypotheses:

$$H_0 : \quad \bar{x} \sim f_n(\bar{x}) \quad (91)$$

$$H_1 : \quad \bar{x} \sim f_n(\bar{x} - \delta) \quad (92)$$

- $\delta > 0$ means that $f_n(\bar{x} - \delta)$ is shifted δ to the right of $f_n(\bar{x})$.

¶ **Example:** Normal distribution:

$$x \sim \mathcal{N}(\mu, \sigma^2) \quad \text{implies} \quad \bar{x} \sim \mathcal{N}(\mu, \sigma^2/n) \quad (93)$$

Review exercise 26, p. 52.

- Hence the null and alternative hypotheses are:

$$H_0 : \quad \bar{x} \sim \mathcal{N}(\mu, \sigma^2/n) \quad (94)$$

$$H_1 : \quad \bar{x} \sim \mathcal{N}(\mu + \delta, \sigma^2/n) \quad (95)$$

- The pdf of the sample mean depends on the sample size.

¶ **Critical value, C :**

- Reject H_0 iff $\bar{x} > C$.
- C depends on the sample size.

¶ **Errors:**

- Type I error: false rejection of H_0 .
- Type II error: false rejection of H_1 .

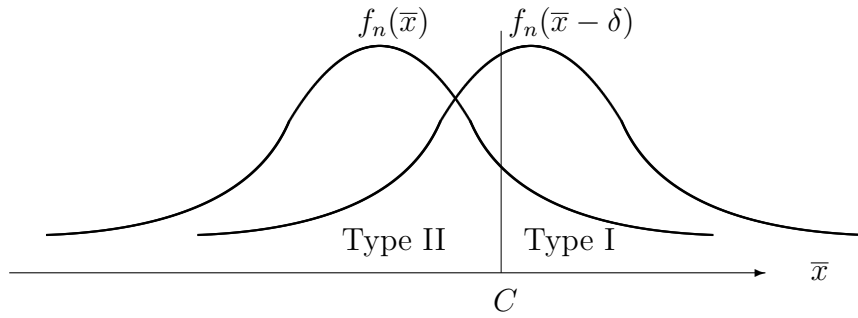


Figure 11: Pdf's for H_0 and H_1 , illustrating type-I and type-II errors.

¶ **Error Probabilities:**

- α = probability of falsely rejecting H_0 :

$$\alpha = \text{Prob}(\bar{x} > C | H_0) \quad (96)$$

$$= \int_C^\infty f_n(\bar{x}) d\bar{x} \quad (97)$$

- β = probability of falsely rejecting H_1 = probability of falsely accepting H_0 :

$$\beta = \text{Prob}(\bar{x} \leq C | H_1) \quad (98)$$

$$= \int_{-\infty}^C f_n(\bar{x} - \delta) d\bar{x} \quad (99)$$

$$= \int_{-\infty}^{C-\delta} f_n(\bar{x}) d\bar{x} \quad (100)$$

- Power = $1 - \beta$ = probability of correctly rejecting H_1 .

¶ **Example: Normal distribution.**

- Null and alternative hypotheses are as in eqs.(94) and (95), p.25.

- Probability of type I error:

$$\alpha = \text{Prob}(\bar{x} > C | H_0) \quad (101)$$

$$= \text{Prob}\left(\frac{\bar{x} - \mu}{\sigma/\sqrt{n}} > \frac{C - \mu}{\sigma/\sqrt{n}}\right) \quad (102)$$

$$= 1 - \Phi\left(\frac{C - \mu}{\sigma/\sqrt{n}}\right) \quad (103)$$

where Φ is the cumulative distribution function of the standard normal distribution.

- Probability of type II error:

$$\beta = \text{Prob}(\bar{x} \leq C | H_1) \quad (104)$$

$$= \text{Prob}\left(\frac{\bar{x} - (\mu + \delta)}{\sigma/\sqrt{n}} \leq \frac{C - (\mu + \delta)}{\sigma/\sqrt{n}}\right) \quad (105)$$

$$= \Phi\left(\frac{C - (\mu + \delta)}{\sigma/\sqrt{n}}\right) \quad (106)$$

- Note: α and β change in opposite directions as n increases.

Review exercise 27, p. 52.

¶ **Choose the sample size, n , so that:**

- $\alpha =$ specified value, e.g. 0.02.
- $\beta \leq$ specified value, e.g. 0.1.

Review exercise 28, p. 52.

¶ **Example: normal distribution.**

- Given: μ , δ and σ .
- Choose type-I error probability, α , say $\alpha = 0.02$.
- For any sample size n , the critical value, C , is found from eq.(103), p.26, as:

$$\Phi\left(\frac{C - \mu}{\sigma/\sqrt{n}}\right) = 1 - \alpha \quad (107)$$

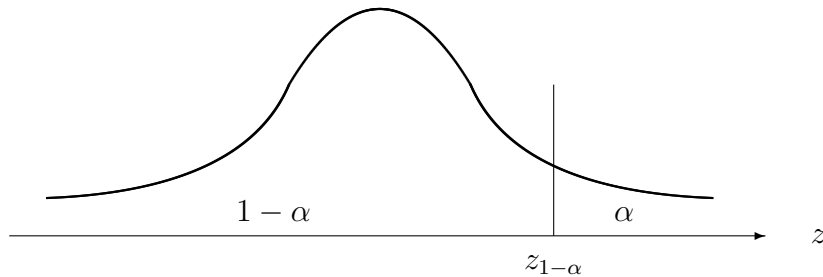


Figure 12: Sketch of probability density illustrating the critical value, eq.(109).

So:

$$\frac{C - \mu}{\sigma/\sqrt{n}} = z_{1-\alpha} = (1 - \alpha)\text{th quantile of } \Phi \quad (108)$$

So:

$$C = \mu + \frac{z_{1-\alpha}\sigma}{\sqrt{n}} \quad (109)$$

- Now the type-II error probability, β , is, from eq.(106), p.26:

$$\beta = \Phi\left(\frac{C - (\mu + \delta)}{\sigma/\sqrt{n}}\right) \quad (110)$$

$$= \Phi\left(z_{1-\alpha} - \frac{\delta}{\sigma/\sqrt{n}}\right) \quad (111)$$

- Numerical example: $\mu = 0$, $\delta = 0.01$, $\sigma = 0.007$.
 - Require $\alpha = 0.02$ so $z_{1-\alpha} = 2.05$.
 - Table 5 shows n , C and β .

n	C , eq.(109)	$\frac{C - (\mu + \delta)}{\sigma/\sqrt{n}}$	β , eq.(110)
2	0.010147	0.029695	0.512
5	0.0064175	-1.14438	$1 - 0.8729 = 0.127$
10	0.0045379	-2.46754	$1 - 0.9934 = 0.0066$

Table 5: Sample size n , critical value C , and type-II error probability β .

- If we require $\beta \leq 0.1$ then:
 - $n = 5$ is too small.
 - $n = 10$ is more than big enough.
- Suppose that $\delta > 0$. Then:

$$\Phi\left(z_{1-\alpha} - \frac{\delta}{\sigma/\sqrt{n}}\right) < \Phi(z_{1-\alpha}) \quad (112)$$

Thus, from eqs.(107) and (111):

$$\beta = \Phi\left(z_{1-\alpha} - \frac{\delta}{\sigma/\sqrt{n}}\right) < \Phi(z_{1-\alpha}) = 1 - \alpha \quad (113)$$

That is:

$$\beta < 1 - \alpha \quad (114)$$

Trade-off: small α means that β may be large.

8.2 Uncertain Effect Size and Variance

¶ Effect size:

- Effect size: $\Delta = \mu_0 - \mu_1$. Note: Δ (here) = $-\delta$ (section 8.1).
- Consider the upper-tail hypothesis: $\Delta < 0$.
- A similar derivation can be formulated for other cases.

¶ The problem:

- Assume normal distribution.
- We have an estimate of the effect size, $\tilde{\Delta}$, but we are unsure how negative it really should be.
- We have an estimate $\tilde{\sigma}$ of the population standard deviation but we are unconfident that this estimate is correct.

¶ Fractional-error info-gap model:

$$\mathcal{U}(h, \tilde{\Delta}, \tilde{\sigma}) = \left\{ (\Delta, \sigma) : \begin{aligned} (1+h)\tilde{\Delta} &\leq \Delta \leq \min[0, (1-h)\tilde{\Delta}] \\ \max[0, (1-h)\tilde{\sigma}] &\leq \sigma \leq (1+h)\tilde{\sigma} \end{aligned} \right\}, \quad h \geq 0 \quad (115)$$

¶ **Power** = $1 - \beta = 1 -$ probability of type-II error, eq.(111), p.27:

$$\text{Power}(\Delta, \sigma, n) = 1 - \Phi\left(\frac{\Delta\sqrt{n}}{\sigma} + z_{1-\alpha}\right) \quad (116)$$

where $z_{1-\alpha}$ is the $(1 - \alpha)$ th quantile of the standard normal distribution.

¶ **Robustness** of sample size n , with requirement that the power be no less than $1 - \beta_c$:

$$\hat{h}(n, \beta_c) = \max \left\{ h : \left(\min_{(\Delta, \sigma) \in \mathcal{U}(h, \tilde{\Delta}, \tilde{\sigma})} \text{Power}(\Delta, \sigma, n) \right) \geq 1 - \beta_c \right\} \quad (117)$$

¶ **Inner minimum** in eq.(117):

$$\mu(h) = \min_{(\Delta, \sigma) \in \mathcal{U}(h, \tilde{\Delta}, \tilde{\sigma})} \text{Power}(\Delta, \sigma, n) \quad (118)$$

- $\mu(h)$ decreases as h increases: nesting of uncertainty sets.
- Robustness: greatest h such that $\mu(h) \geq 1 - \beta_c$.
- $\mu(h)$ is monotonic in h , so robustness is the max h satisfying $\mu(h) = 1 - \beta_c$.
- Plot of $\mu(h)$ vs. h is plot of $1 - \beta_c$ vs. $\hat{\alpha}(n, \beta_c)$.
- See fig. 13.

¶ Derive the robustness function.

- $\mu(h)$ occurs for the greatest allowed value of Δ/σ , which is negative and occurs when $\Delta = \min[0, (1-h)\tilde{\Delta}]$ and when $\sigma = (1+h)\tilde{\sigma}$.

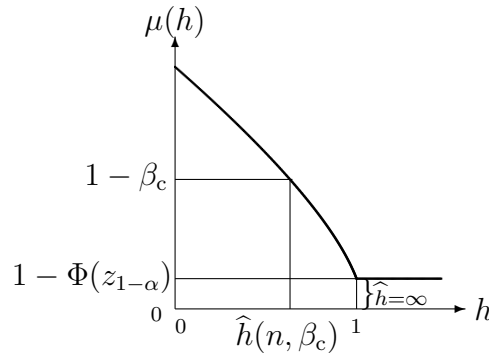


Figure 13: Illustration of the calculation of robustness.

- If $h \geq 1$ then:
 - Power is minimized when $\Delta = 0$ so:

$$\text{Power}(\Delta, \sigma, n) = 1 - \Phi(z_{1-\alpha}) = \mu(h) \quad (119)$$

as in horizontal section of the curve in fig. 13.

- Robustness is infinite when $1 - \beta_c < 1 - \Phi(z_{1-\alpha})$.
- Thus very low demanded power (large β_c) implies very high robustness:

$$\hat{h}(n, \beta_c) = \infty \quad \text{if } \beta_c > \Phi(z_{1-\alpha}) \quad (120)$$

¶ If $h < 1$: the robustness is the greatest value of h satisfying:

$$\Phi \left[\frac{(1-h)\tilde{\Delta}\sqrt{n}}{(1+h)\tilde{\sigma}} + z_{1-\alpha} \right] \leq \beta_c \quad (121)$$

- Let us denote by $q(\beta_c)$ the β_c quantile of the standard normal distribution:

$$\beta_c = \int_{-\infty}^{q(\beta_c)} \phi(x) dx \quad (122)$$

- $q(\beta_c)$ increases from $-\infty$ to $+\infty$ as β_c increases from 0 to 1.
- The robustness is the greatest value of h satisfying:

$$\frac{(1-h)\tilde{\Delta}\sqrt{n}}{(1+h)\tilde{\sigma}} + z_{1-\alpha} \leq q(\beta_c) \quad (123)$$

- If:

$$\frac{\tilde{\Delta}\sqrt{n}}{\tilde{\sigma}} + z_{1-\alpha} > q(\beta_c) \quad (124)$$

then the robustness is zero for this value of β_c , and positive robustness is obtained only for greater values of β_c (lower power).

- Define:

$$\nu = \frac{q(\beta_c) - z_{1-\alpha}}{\tilde{\Delta}\sqrt{n}/\tilde{\sigma}} \quad (125)$$

The robustness is positive only if $\nu < 1$ (recalling that $\tilde{\Delta} < 0$).

• Now, solving eq.(123) (as an equality) for h , in the case that eq.(124) does not hold (that is, $\nu < 1$), we obtain the robustness:

$$\hat{h}(n, \beta_c) = \begin{cases} \frac{1 - \nu}{1 + \nu} & \text{if } \nu < 1 \\ 0 & \text{else} \end{cases}, \quad \text{if } \beta_c \leq \Phi(z_{1-\alpha}) \quad (126)$$

The complete robustness function is eqs.(120) and (126).

¶ **Trade-off: Robustness vs. power.**

Applying the chain rule for differentiation to eq.(126), one finds:

$$\frac{\partial \hat{h}(n, \beta_c)}{\partial (1 - \beta_c)} < 0 \quad (127)$$

The robustness $\hat{h}(n, \beta_c)$ decreases as the demanded power, $1 - \beta_c$, increases: high aspirations are vulnerable to uncertainty.

¶ **Trade-off: Robustness vs. sample size.**

One finds:

$$\frac{\partial \hat{h}(n, \beta_c)}{\partial n} > 0 \quad (128)$$

Thus the robustness increases as the sample size increases.

8.3 Uncertain Sample PDF

8.3.1 Background

¶ Binary statistic test.

- x = decision statistic.
- Distribution of x under H_1 equals distribution under H_0 shifted up by δ :

$$H_0 : \quad x \sim f(x) \quad (129)$$

$$H_1 : \quad x \sim f(x - \delta) \quad (130)$$

- We accept the null hypothesis iff $x \leq C$.
- Determine:
 - Sample size, n .
 - Critical value, C .

¶ Definitions.

- α = level of significance = probability of type-I error (falsely reject H_0).
- $\beta(f)$ = probability of type-II error (falsely reject H_1) = $1 - \text{power}$.
- δ = non-negative effect size.
- $f(x)$ = pdf of decision statistic.
- C = critical value.

$$1 - \alpha = \int_{-\infty}^C f(x) dx \quad (131)$$

$$\beta(f) = \int_{-\infty}^C f(x - \delta) dx = \int_{-\infty}^{C-\delta} f(x) dx = 1 - \alpha - \int_{C-\delta}^C f(x) dx \quad (132)$$

¶ Standard statistical approach.

- Known sampling distribution: $\tilde{f}(x)$.
- $\tilde{f}(x)$ depends on sample size (number of measurements.)
- Specify α and δ .
- Determine C and β from eqs.(131) and (132).
- Increase sample size until the power is adequate.

8.3.2 Info-gap Approach to Determining the Sample Size

¶ **Approach.**

- Sampling distribution is uncertain.
- Evaluate info-gap robustness of the estimated power.
- Determine sample size, n , so that adequate power is adequately robust.

¶ **Info-gap model for pdf uncertainty: fractional-error.**

$$\mathcal{U}(h, \tilde{f}) = \left\{ f(x) : f \in \mathcal{P}, |f(x) - \tilde{f}(x)| \leq h\tilde{f}(x) \right\}, \quad h \geq 0 \quad (133)$$

\mathcal{P} is the set of all non-negative and normalized pdfs on the domain of x .

¶ **Performance requirement.**

- Power = $1 - \beta$. Require large power; small β .
- $1 - \beta_d$ = demanded power.
- Analyst requires $\beta \leq \beta_d$.

¶ **Robustness:**

$$\hat{h}(N, \beta_d) = \max \left\{ h : \left(\max_{f \in \mathcal{U}(h, \tilde{f})} \beta(f) \right) \leq \beta_d \right\} \quad (134)$$

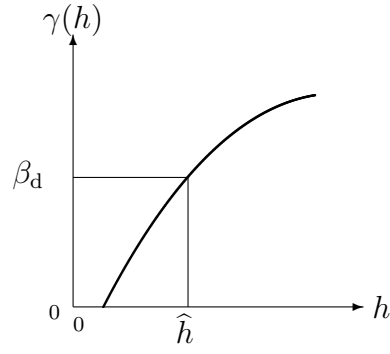


Figure 14: Illustration of the calculation of robustness.

¶ **Evaluating robustness.**

- Denote inner maximum in eq.(134) by $\gamma(h)$.
- Robustness is $\max h$ such that $\gamma(h) \leq \beta_d$.
- Uncertainty sets $\mathcal{U}(h, \tilde{f})$ are nested with respect to h .
- Thus $\gamma(h)$ increases as h increases.
- Thus robustness is $\max h$ at which $\gamma(h) = \beta_d$.
- $\gamma(h)$ is inverse of $\hat{h}(N, \beta_d)$:

$$\gamma(h) = \beta_d \quad \text{if and only if} \quad \hat{h}(N, \beta_d) = h \quad (135)$$

8.3.3 An Approximate Robustness for Small Effect Size

¶ **Special case**, very small effect size:

$$\delta \ll 1 \quad (136)$$

Derive approximate expression for robustness.

¶ Now eq.(132), p.31, can be approximated as:

$$\beta(f) = 1 - \alpha - f(C)\delta \quad (137)$$

¶ **Nominal critical value, \tilde{C} :**

- $\tilde{C} = (1 - \alpha)$ th quantile of the best-estimated pdf $\tilde{f}(x)$.
- \tilde{C} depends on sample size n .

¶ **Maximizing pdf.** The pdf in $\mathcal{U}(h, \tilde{f})$ which maximizes β is very nearly:

$$\hat{f}(x) = \begin{cases} \tilde{f}(x) & \text{if } x < \tilde{C} - \delta \\ (1 - h)\tilde{f}(x) & \text{if } x \in [\tilde{C} - \delta, \tilde{C}] \\ (1 + wh)\tilde{f}(x) & \text{if } x > \tilde{C} \end{cases} \quad (138)$$

where w is a very small positive number which normalizes $\hat{f}(x)$. That is, w is determined so that the decrement in \hat{f} in $[\tilde{C} - \delta, \tilde{C}]$ is compensated by the increment in (\tilde{C}, ∞) :

$$wh[1 - \tilde{F}(\tilde{C})] = h\tilde{f}(\tilde{C})\delta \quad (139)$$

where \tilde{F} is the cumulative distribution function of \tilde{f} .

¶ **Evaluating $\gamma(h)$.** $\beta(\hat{f})$ in eq.(137) becomes:

$$\gamma(h) = \beta(\hat{f}) = 1 - \alpha - (1 - h)\tilde{f}(\tilde{C})\delta \quad (140)$$

Note: $\gamma(h)$ increases as h increases.

¶ **Evaluating robustness.**

Equate $\gamma(h)$ in eq.(140) to β_d and solve for h :

$$\hat{h}(n, \beta_d) = \begin{cases} 0 & \text{if } \beta_d < 1 - \alpha - \tilde{f}(\tilde{C})\delta \\ \frac{\beta_d - 1 + \alpha + \tilde{f}(\tilde{C})\delta}{\tilde{f}(\tilde{C})\delta} & \text{else} \end{cases} \quad (141)$$

- Robustness increases as β_d increases. Trade-off:
high power \iff low robustness.
- Robustness is zero when β_d equals the nominal value, $\beta(\tilde{f})$.
- Robustness depends on sample size through nominal critical value \tilde{C} .
- This derivation is contingent on the assumption in eq.(136), p.34.

9 Tests of the Mean with Distributional Uncertainty

§ **Source:** Yakov Ben-Haim, 2008, Tests of the Mean with Distributional Uncertainty: An Info-Gap Approach, working paper.²

9.1 Distributional Uncertainty

§ **Statistical tests of the mean** depend on various assumptions about the data and population, such as:

- Normality.
- Random sampling: independent measurements with same instrument from same population which is unaffected by the measurement process.
- Stationarity of the sampled population.

§ **Distributional uncertainty:**

- Violations of assumptions about data and population, *unknown to the analyst*.
- E.g.:
 - Non-normality.
 - Sampling protocol varies. E.g., some observers are experts, some are not.
 - Population evolves during the sample.
 - Population is influenced by the sample.
- Examples:
 - Franklin³ uses a range of observational data from many different sources over the past 150 years—of varying and uncertain accuracy and reliability—to evaluate change in bird assemblages in northern Australia.
 - McCarthy⁴ uses museum collections to evaluate trends in marsupials and monotremes, recognizing that variable and uncertain collection efforts introduce uncertainties.
 - Burgman *et al*⁵ recognize that “collection frequencies will reflect changing trends in museum and herbarium collections”, which introduces uncertainties in evaluating extinction threats based on historical development of collections.
 - Stewart-Oaten *et al*⁶ study tests of changes of a mean population property, before and after an impact, where the impact cannot be replicated (e.g., construction of a power plant). They note that data from such measurements “do not necessarily satisfy” the assumptions of standard tests. They state that “there is no panacea” for violation of test assumptions, and if the assumptions “are seriously wrong, alternative analyses are needed. This will often require a long time series of data.” These authors discuss many sources of violation of test assumptions, stressing the importance of unknown skewness of distributions or correlations among measurements.

§ **The problem:**

When violations are unknown and uncharacterized, the analyst cannot correct for them.

²Files: \papers\T-Test\ct03.tex and ttest07.tex.

³Franklin, Donald C., 1999, Evidence of disarray amongst granivorous bird assemblages in the savannas of northern Australia, a region of sparse human settlement, *Biological Conservation*, 90: 53–68.

⁴McCarthy, Michael A., 1998, Identifying declining and threatened species with museum data, *Biological Conservation*, 83: 9–17.

⁵Burgman, Mark A., Roger C. Grimson and Scott Ferson, 1995, Inferring threat from scientific collections, *Conservation Biology*, 9: 923–928.

⁶Stewart-Oaten, Allan, James R. Bence, and Craig W. Osenberg, 1992, Assessing effects of unreplicated perturbations: No simple solutions, *Ecology*, vol. 73, #4, pp.1396–1404.

§ **Statistical tools exist for managing distributional uncertainty.**

- Careful test design.
- Non-parametric methods weaken some assumptions, e.g. normality.
 - These tests do assume random sampling, and usually are asymptotic.
 - They can be very sensitive to outliers.
- Given adequate data, one can model the data as a mixture of populations.
- Outliers can be managed using Jackknife or trimmed-means techniques.
- Method of M -estimates.

9.2 Info-Gap Representations of Distributional Uncertainty

§ θ is the test statistic. It may be a t statistic, but not necessarily.

§ Tests of the mean:

$$H_0 : \quad x \sim g(x) \quad (142)$$

$$H_1 : \quad x \sim g(x - \delta) \quad (143)$$

§ Estimated pdfs.

- Let $\tilde{f}_i(\theta)$ denote the best guess of the pdf of the test statistic t , under hypothesis H_i .
- For instance:
 - If θ is the t statistic then $\tilde{f}_0(\theta)$ is the t distribution with $n - 1$ degrees of freedom.
 - $\tilde{f}_1(\theta) = \tilde{f}_0(\theta - \delta)$ where $\delta = (T_1 - T_0)/(s/\sqrt{n})$ is the shift between the two hypotheses.
 - Thus $\tilde{f}_1(\theta)$ is formed by shifting $\tilde{f}_0(\theta)$ to the right by δ .

§ A fractional-error info-gap model:

$$\mathcal{U}_i(h, \tilde{f}_i) = \left\{ f(\theta) : f(\theta) \in \mathcal{P}, |f(\theta) - \tilde{f}_i(\theta)| \leq h f_t^*, \forall \theta \right\}, \quad h \geq 0 \quad (144)$$

- \mathcal{P} is the set of all normalized non-negative pdf's.
- f_t^* is a normalization constant with units of probability density. For instance the mode:

$$f_t^* = \max_{\theta} \tilde{f}(\theta) \quad (145)$$

If $\tilde{f}(\theta)$ is a t distribution then $f_t^* = \tilde{f}(0)$.

§ A more restrictive fractional-error info-gap model:

$$\mathcal{U}_i(h, \tilde{f}_i) = \left\{ f(\theta) : f(\theta) \in \mathcal{P}, |f(\theta) - \tilde{f}_i(\theta)| \leq h \tilde{f}_i(\theta), \forall \theta \right\}, \quad h \geq 0 \quad (146)$$

- The variation on the tails dies out if $\tilde{f}_i(\theta)$ becomes small on the tails, unlike eq.(144).

§ Estimated cdfs.

- Let $\tilde{F}_i(\theta)$ denote the best guess of the cdf of the test statistic t , under hypothesis H_i .
- For instance:
 - If θ is the t statistic then $\tilde{F}_0(\theta)$ is the t distribution with $n - 1$ degrees of freedom for the statistic in eq.(167).
 - $\tilde{F}_1(\theta) = \tilde{F}_0(\theta - \delta)$ where $\delta = (T_1 - T_0)/(s/\sqrt{n})$.
 - Thus $\tilde{F}_1(\theta)$ is formed by shifting $\tilde{F}_0(\theta)$ to the right by δ .

§ Uniform-bound info-gap model:

$$\mathcal{U}_i(h) = \left\{ F(\theta) : F(\theta) \in \mathcal{P}, |F(\theta) - \tilde{F}_i(\theta)| \leq h, \forall \theta \right\}, \quad h \geq 0 \quad (147)$$

where \mathcal{P} is the set of all normalized non-negative cdf's.

9.3 Robustness Functions with CDF Uncertainty

§ This section is based on file \papers\T-Test\ct03.tex.

9.3.1 Binary Test: Formulation

§ **Data.**

- $X = \{x_1, \dots, x_n\}$
- Not necessarily a random sample of any known distribution.

§ **Decision.** Two simple hypotheses about the population mean:

$$H_0 : \quad \mu = T_0 \quad (148)$$

$$H_1 : \quad \mu = T_1 \quad (149)$$

where each T_i is a specified number, and $T_1 > T_0$.

§ **Size and power.**

- θ is a statistic, for instance the t statistic.
- $F_i(\theta)$ is the cdf of θ under H_i .
- For any distribution $F(\theta)$, $q_\alpha(F)$ is the $(1 - \alpha)$ th quantile of $F(\theta)$:

$$1 - \alpha = F[q_\alpha(F)] \quad (150)$$

- We reject H_0 with significance α if:

$$\theta \geq q_\alpha(F_0) \quad (151)$$

- The size α , and power, $1 - \beta$, are defined in:

$$1 - \alpha = F_0[q_\alpha(F_0)] \quad (152)$$

$$\beta = F_1[q_\alpha(F_0)] \quad (153)$$

- The size, α , is the probability of *falsely rejecting* the null hypothesis, H_0 .
- $1 - \alpha$ is the probability of correctly accepting H_0 .
- The power, $1 - \beta$, is the probability of *correctly rejecting* H_0 .
- β is the probability of falsely rejecting H_1 .

9.3.2 Robustnesses for Type I Errors

§ **Decision threshold.** Test of size α^* which rejects H_0 when:

$$\theta \geq q_{\alpha^*}(\tilde{F}_0) \quad (154)$$

α^* : “nominal” size of the test, based on best-estimate of cdf under H_0 , \tilde{F}_0 .

§ Note that:

$$\tilde{F}_0[q_{\alpha^*}(\tilde{F}_0)] = 1 - \alpha^* \quad (155)$$

§ **Robustness for falsely rejecting H_0 :**

• Maximum horizon of uncertainty, h , at which the test at nominal size α^* falsely rejects H_0 with probability no greater than α :

$$\hat{h}_0(\alpha^*, \alpha) = \max \left\{ h : \left(\min_{F \in \mathcal{U}_0(h)} F[q_{\alpha^*}(\tilde{F}_0)] \right) \geq 1 - \alpha \right\} \quad (156)$$

• We use the quantile $q_{\alpha^*}(\tilde{F}_0)$ because the test is implemented with the quantile of the best-guess distribution under H_0 , $\tilde{F}_0(\theta)$, and is of nominal size α^* .

• The actual size (probability of falsely rejecting H_0) is determined by the unknown true distribution under H_0 , $F(\theta)$, which is info-gap-uncertain.

§ **Relation to type I error** (falsely rejecting H_0):

• $\hat{h}_0(\alpha^*, \alpha)$ is the greatest horizon of uncertainty up to which the probability of type I error is no greater than α .

§ **The Robustness**, $\hat{h}_0(\alpha^*, \alpha)$, for the info-gap model in eq.(147), is:

$$\hat{h}_0(\alpha^*, \alpha) = \alpha - \alpha^* \quad (157)$$

or zero if this is negative.

- α is the **effective size**, while α^* is the **nominal size**.
- For any choice of α^* , the robustness curve for type-I error, $\hat{h}_0(\alpha^*, \alpha)$ vs. α , is independent of the form of the test: t test, Wilcoxon signed-ranks test, etc.
- The implementation of the test, eq.(154), does depend on the type of test, through the value of the quantile $q_{\alpha^*}(\tilde{F}_0)$.

§ **Derivation of eq.(157).**

- Define the following step function:

$$V(x) = \begin{cases} 0, & \text{if } x < 0 \\ x, & \text{if } 0 \leq x \leq 1 \\ 1, & \text{else} \end{cases} \quad (158)$$

- Let $m_0(h)$ denote the inner minimum in eq.(156).
- The robustness, $\hat{h}_0(\alpha^*, \alpha)$, is the greatest non-negative h for which $m_0(h) = 1 - \alpha$.
- If there is no such h , then the robustness is zero.
- The inner min results when $F(\theta)$ is minimal at $q_{\alpha^*}(\tilde{F}_0)$, subject to membership in $\mathcal{U}_0(h)$.
- From the info-gap model in eq.(147) we find:

$$m_0(h) = V(\tilde{F}_0[q_{\alpha^*}(\tilde{F}_0)] - h) = V(1 - \alpha^* - h) \quad (159)$$

- Recall that $\tilde{F}_0[q_{\alpha^*}(\tilde{F}_0)] = 1 - \alpha^*$.
- The greatest value of h at which $m_0(h) = 1 - \alpha$ is the robustness, eq.(157).

9.3.3 Robustnesses for Type II Errors

§ Robustness for falsely accepting H_0 .

• $\hat{h}_1(\alpha^*, \beta)$ is the greatest horizon of uncertainty up to which the probability of falsely accepting H_0 , with a test of nominal size α^* , is no greater than β :

$$\hat{h}_1(\alpha^*, \beta) = \max \left\{ h : \left(\max_{F \in \mathcal{U}_1(h)} F[q_{\alpha^*}(\tilde{F}_0)] \right) \leq \beta \right\} \quad (160)$$

§ $1 - \beta^*$ is the nominal power:

$$1 - \beta^* = 1 - \tilde{F}_1[q_{\alpha^*}(\tilde{F}_0)] \quad (161)$$

• $\hat{h}_1(\alpha^*, \beta)$, for the info-gap model in eq.(147), is:

$$\hat{h}_1(\alpha^*, \beta) = 1 - \beta^* - (1 - \beta) \quad (162)$$

or zero if this is negative.

• $1 - \beta$ as the **effective power**. $1 - \beta^*$ is the **nominal power**.

• For any choice of α^* , $\hat{h}_1(\alpha^*, \beta)$ vs. β , depends on the form of the test, unlike for the type-I robustness. This is because the value of β^* depends on α^* through the cdf's of the test statistic, \tilde{F}_0 and \tilde{F}_1 .

§ Derivation of eq.(162).

- $m_1(h)$ denotes the inner maximum in eq.(160).
- The robustness, $\hat{h}_1(\alpha^*, \beta)$, is the greatest h at which $m_1(h) = \beta$.
- From the info-gap model in eq.(147), and using $V(x)$ in eq.(158):

$$m_1(h) = V(\tilde{F}_1[q_{\alpha^*}(\tilde{F}_0)] + h) \quad (163)$$

Equating this to β and solving for h we find the robustness in eq.(162) with the aid of the expression for the nominal power in eq.(161).

9.3.4 Decisions and Judgments

§ Two decisions, two judgments:

- *Decide* on the nominal test size α^* and the sample size n .
 - Together these decisions determine the decision threshold $q_{\alpha^*}(\tilde{F}_0)$ in eq.(154), p.39.
- *Judge* what are reliable and acceptable values of effective size α and effective power $1 - \beta$.
 - Do this by considering $\hat{h}_0(\alpha^*, \alpha)$ and $\hat{h}_1(\alpha^*, \beta)$.
 - α (size or level of significance) is the probability of falsely rejecting H_0 .
 - $1 - \beta$ (power) is the probability of correctly rejecting H_0 .

§ Example: t test.

- Test statistic, $\theta = (\bar{x} - T_0)(s/\sqrt{n})$. \bar{x} is sample mean, s^2 is sample variance, and n is sample size.
- Estimated distribution under H_0 , $\tilde{F}_0(\theta)$, is the cdf of the t statistic with $n - 1$ degrees of freedom.
- Estimated distribution under H_1 is $\tilde{F}_1(\theta) = \tilde{F}_0(\theta - \delta)$ where $\delta = (T_1 - T_0)/(s/\sqrt{n})$.
- True distributions under H_0 and H_1 are unknown; uncertainty is represented by info-gap model in eq.(147), p.37.

§ No distributional uncertainty: no need for judgments:

- α^* is the actual size.
- Actual power, $1 - \beta^*$, is entirely determined by α^* and n .
- Values of α^* and $1 - \beta^*$ are shown in table 6.
- Power increases with increasing n at fixed α^* .
- Power increases with increasing α^* at fixed n .

$\alpha^* = 0.01$		$\alpha^* = 0.05$	
n	$1 - \beta^*$	n	$1 - \beta^*$
5	0.1027	3	0.1784
7	0.3185	4	0.3736
9	0.5400	5	0.5390
12	0.7644	7	0.7457
31	0.9980	31	0.9997

Table 6: Size and power in the absence of distributional uncertainty.

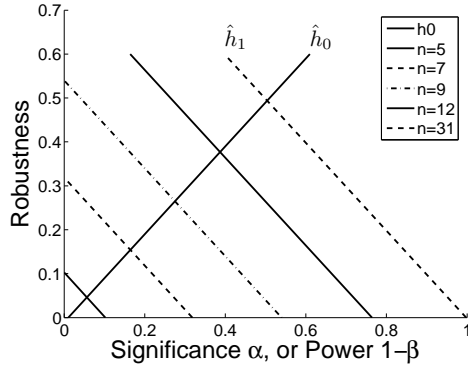


Figure 15: Robustness curves for the t test, $\hat{h}_0(\alpha^*, \alpha)$ for falsely rejecting H_0 , and $\hat{h}_1(\alpha^*, \alpha)$ for falsely rejecting H_1 . Nominal size is $\alpha^* = 0.01$. $\hat{h}_1(\alpha^*, \alpha)$ calculated at 5 different sample sizes: $n = 5, 7, 9, 12$ and 31 . $\delta = 1$.

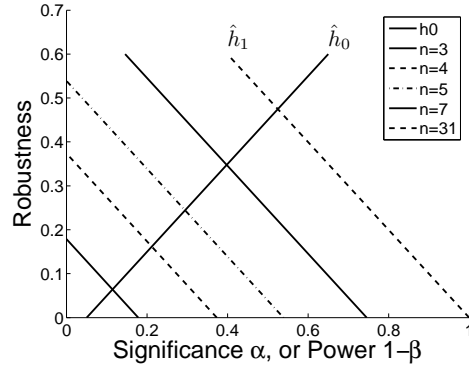


Figure 16: Robustness curves for the t test, $\hat{h}_0(\alpha^*, \alpha)$ for falsely rejecting H_0 , and $\hat{h}_1(\alpha^*, \alpha)$ for falsely rejecting H_1 . Nominal size is $\alpha^* = 0.05$. $\hat{h}_1(\alpha^*, \alpha)$ calculated at 5 different sample sizes: $n = 3, 4, 5, 7$ and 31 . $\delta = 1$.

§ **Robustness curves.** Figs. 15 and 16:

- $\hat{h}_0(\alpha^*, \alpha)$ vs. α (positive slope).
 - No robustness for nominal size: $\hat{h}_0(\alpha^*, \alpha^*) = 0$.
 - Positive slope: Trade-off: robustness is exchanged for significance.
- $\hat{h}_1(\alpha^*, \beta)$ vs. $1 - \beta$ (negative slope).
 - No robustness for nominal power: $\hat{h}_1(\alpha^*, \beta^*) = 0$.
 - Negative slope: Trade-off: robustness is exchanged for power.

§ **Judging reliable effective size, α :**

- The test designed for $\alpha^* = 0.01$ will falsely reject H_0 with probability ≤ 0.05 if $F(\theta)$ differs from $\tilde{F}_0(\theta)$ by no more than 0.04 in cumulative probability.
 - E.g., tails no heavier than 0.04 of total distribution.
 - E.g., outlying sub-population no larger than 0.04 of total distribution.

§ **Judging effective power, $1 - \beta$:**

- A test designed for size $\alpha^* = 0.01$ with sample size $n = 9$ (dot-dash in fig. 15), has no robustness for power 0.54 (the horizontal intercept and nominal power).
- This test will falsely accept H_0 with probability of 0.44 if the actual cdf differs from the estimated cdf by no more than 0.1.

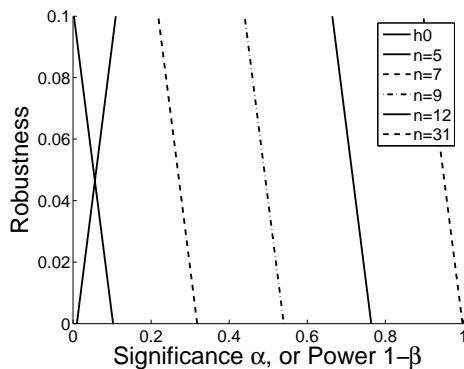


Figure 17: Same as fig. 15.

§ Choosing the sample size, n .

- Only type-II robustness is influenced by the sample size.
- The nominal and effective power:
 - increase with increasing sample size,
 - are influenced by the nominal size α^* .
- Choose n in light of the effective power and robustness which are needed.
- See fig. 17, which is expanded from fig. 15.

§ Choosing the sample size, n , continued.

- In fig. 17 consider nominal size $\alpha^* = 0.01$.
- Judgment: effective size $\alpha = 0.05$ is adequate and reliable because $\hat{h}_0(0.01, 0.05) = 0.04$.
- Apply this robustness to type II: Require $\hat{h}_1(\alpha^*, \beta) = 0.04$.
- From fig. 17: effective powers of 0.50, 0.72 and 0.96 for sample sizes 9, 12 and 31.
- Judgment: power of 0.50 is too small, so we require a sample larger than $n = 9$.
- Judgment: if power of 0.72 is adequate then we adopt a sample of size 12.
- Choosing a sample of size 31 would result in power of 0.96.

§ Judgments of robustness: how much robustness is needed?

- Robustness has units of probability (in this example).
- Thus judge adequate robustness probabilistically.
- This not necessary: analogical inference.

§ Choosing the sample size, n , continued.

- Previously we required $\hat{h}_0(\alpha^*, \alpha) = \hat{h}_1(\alpha^*, \beta)$.
- This is not necessary. We can make separate judgments for type I and type II robustnesses.

9.4 Robustness Functions with PDF Uncertainty: Definitions

§ This section and the next are based on file \papers\T-Test\ttest07.tex.

9.4.1 Binary Test: Formulation

§ **Data.**

- $X = \{x_1, \dots, x_n\}$
- Not necessarily a random sample of any known distribution.

§ **Decision.** Two simple hypotheses about the population mean:

$$H_0 : \quad \mu = T_0 \quad (164)$$

$$H_1 : \quad \mu = T_1 \quad (165)$$

where each T_i is a specified number, and $T_1 > T_0$.

§ **Sample mean and variance:**

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i, \quad s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2 \quad (166)$$

§ **The t statistic for testing H_0 is:**

$$t = \frac{\bar{x} - T_0}{s/\sqrt{n}} \quad (167)$$

which has a t distribution with $n - 1$ degrees of freedom under H_0 (in the absence of distributional uncertainty).

§ **Size and power of the test.**

- Let $f_i(t)$ denote the probability density of t under H_i .
- For any density $f(t)$, let $q_\alpha(f)$ denote the $(1 - \alpha)$ th quantile of $f(t)$:

$$\int_{-\infty}^{q_\alpha(f)} f(t) dt = 1 - \alpha \quad (168)$$

- We reject H_0 with significance α if:

$$t \geq q_\alpha(f_0) \quad (169)$$

- The size α , and power, $1 - \beta$, are defined in:

$$1 - \alpha = \int_{-\infty}^{q_\alpha(f_0)} f_0(t) dt \quad (170)$$

$$\beta = \int_{-\infty}^{q_\alpha(f_0)} f_1(t) dt \quad (171)$$

- α is the probability of falsely rejecting the null hypothesis, H_0 .
- $1 - \beta$ is the probability of correctly rejecting H_0 .

9.4.2 Robustness for Falsely Rejecting H_0 . (Type I Error.)

- Consider a t test of size α^* , which rejects H_0 when:

$$t > q_{\alpha^*}(\tilde{f}_0) \quad (172)$$

- The robustness is the maximum horizon of uncertainty, h , up to which the t test at size α^* falsely rejects H_0 with probability no greater than α :

$$\hat{h}_0(t, \alpha^*, \alpha) = \max \left\{ h : \left(\min_{f \in \mathcal{U}_0(h, \tilde{f}_0)} \int_{-\infty}^{q_{\alpha^*}(\tilde{f}_0)} f(t) dt \right) \geq 1 - \alpha \right\} \quad (173)$$

- $q_{\alpha^*}(\tilde{f}_0)$: the test is implemented with the quantile of the best-guess distribution under H_0 , \tilde{f}_0 , and is of nominal size α^* .
- Actual size (probability of falsely rejecting H_0) is determined by the unknown true distribution under H_0 , f .
- The inverse of $\hat{h}_0(t, \alpha^*, \alpha)$ is defined as:

$$m_0^t(h, \alpha^*) = 1 - \alpha \quad \text{if and only if} \quad \hat{h}_0(t, \alpha^*, \alpha) = h \quad (174)$$

- An explicit expression for the inverse of $\hat{h}_0(t, \alpha^*, \alpha)$ is:⁷

$$m_0^t(h, \alpha) = [c_1(h) - c_2(h)]h f_t^* + \tilde{F}_0[c_2(h)] - \tilde{F}_0[c_1(h)] \quad (175)$$

where:

$$c_1(h) = -\tilde{f}_0^{-1}(h f_t^*) \quad (176)$$

$$c_2(h) = \min[\tilde{f}_0^{-1}(h f_t^*), q_{\alpha^*}(\tilde{f}_0)] \quad (177)$$

- $\tilde{f}_0(t)$ is the pdf of the t variate with $n - 1$ dofs.
- $\tilde{F}_0(t)$ is the cumulative distribution function of the t variate with $n - 1$ dofs.
- $\tilde{f}_0^{-1}(h)$ is the inverse of $\tilde{f}_0(t)$ for $t \geq 0$.
Thus $-\tilde{f}_0^{-1}(h)$ is the smallest value of t at which $\tilde{f}_0(t) = h$.
So $-\tilde{f}_0^{-1}(0) = -\infty$ and $\tilde{f}_0^{-1}[\tilde{f}_0(0)] = 0$.

§ Type I error (falsely rejecting H_0).

- $\hat{h}_0(t, \alpha^*, \alpha)$ is the greatest horizon of uncertainty at which:
the probability of type I error is no greater than α .
- The test is implemented so that the probability of type I error is no greater than α^* assuming no distributional uncertainty.

⁷Yakov Ben-Haim, 2008, Tests of the Mean with Distributional Uncertainty: An Info-Gap Approach, working paper. Appendix A.

9.4.3 Robustness for Falsely Accepting H_0 . (Type II Error.)

- Consider, as before, a t test of size α^* , which rejects H_0 when:

$$t > q_{\alpha^*}(\tilde{f}_0) \quad (178)$$

- The robustness is the greatest horizon of uncertainty up to which the probability of falsely accepting H_0 , with a t test of size α^* , is no greater than β :

$$\hat{h}_1(t, \alpha^*, \beta) = \max \left\{ h : \left(\max_{f \in \mathcal{U}_1(h, \tilde{f}_1)} \int_{-\infty}^{q_{\alpha^*}(\tilde{f}_0)} f(t) dt \right) \leq \beta \right\} \quad (179)$$

- An explicit expression for the inverse of $\hat{h}_1(t, \alpha^*, \beta)$ is:⁸

$$M_1^t(h, \alpha^*) = 1 + (c_4 - c_3)hf_t^* - \tilde{F}_1(c_4) + \tilde{F}_1(c_3) \quad (180)$$

where:

$$c_3(h) = q_{\alpha^*}(\tilde{f}_0) \quad (181)$$

$$c_4(h) = \max[\tilde{f}_1^{-1}(hf_t^*), q_{\alpha^*}(\tilde{f}_0)] \quad (182)$$

- We have assumed that $q_{\alpha^*}(\tilde{f}_0) \geq \delta$, which in practice will always hold.
- $\tilde{f}_1(t) = \tilde{f}_0(t - \delta)$.
- $\tilde{F}_1(t)$ is the cumulative distribution function of the $\tilde{f}_1(t)$.
- $\tilde{f}_1^{-1}(h)$ is the inverse of $\tilde{f}_1(t)$ for $t \geq \delta$.
- A plot of h vs. $M_1^t(h, \alpha^*)$ is identical to a plot of $\hat{h}_1(t, \alpha^*, \beta)$ vs. β .

§ Nominal power.

- Let $1 - \beta^*$ be the nominal power:

$$\beta^* = \int_{-\infty}^{q_{\alpha^*}(\tilde{f}_0)} \tilde{f}_1(t) dt \quad (183)$$

From the contraction and nesting axioms we recognize that $\hat{h}_1(t, \alpha^*, \beta^*) = 0$ and $\hat{h}_1(t, \alpha^*, \beta) > 0$ only for $\beta > \beta^*$.

⁸Yakov Ben-Haim, 2008, Tests of the Mean with Distributional Uncertainty: An Info-Gap Approach, working paper. Appendix B.

9.5 Robustness with PDF Uncertainty: Numerical Examples

9.5.1 Robustness for Type-I Error

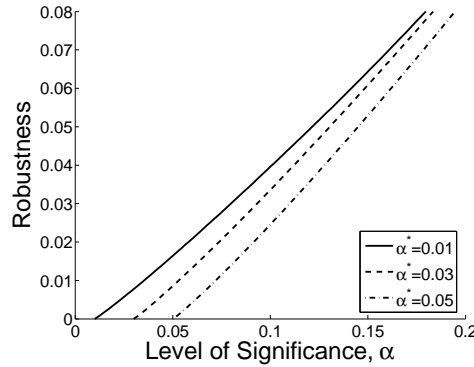


Figure 18: Robustness curves for the t test, $\hat{h}_0(t, \alpha^*, \alpha)$, for falsely rejecting H_0 , at three design sizes, $\alpha^* = 0.01, 0.03$ and 0.05 . Sample size $n = 17$. $f_t^* = \max_t \tilde{f}_0(t)$.

§ Trade-off:

- Robustness vs. level of significance.
- Zero robustness at nominal level of significance.

§ What does $\hat{h}_0(t, \alpha^*, \alpha) = 0.02$ mean?

- The true pdf, $f(t)$, can deviate from $\tilde{f}_0(t)$ by a ‘bump’ (or dimple) no larger than $0.02f_t^*$ if size α is not to be exceeded.
- This is a small bump: the tail of $\tilde{f}_0(t)$ becomes as thin as $0.02f_t^*$ at about 3σ ’s from the mean.
- So, $\hat{h}(t, \alpha^*, \alpha) = 0.02$ might be sufficient robustness only if immunity to small deviations on the far tails is sufficient.

9.5.2 Robustness for Type-II Error

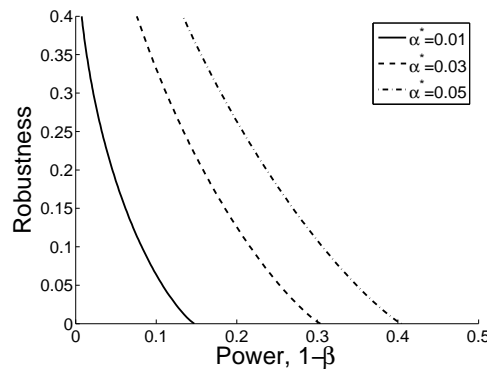


Figure 19: Robustness curves for the t test, $\hat{h}_1(t, \alpha^*, \beta)$, for correctly rejecting H_0 , at three design sizes, $\alpha^* = 0.01, 0.03$ and 0.05 . Sample size $n = 17$. $f_t^* = \max_t \tilde{f}_0(t)$.

§ Trade-off:

- Robustness vs. power of the test.
- Zero robustness at nominal power.

10 Accelerated Lifetime Testing: Simple Case

10.1 Formulation

§ **Lifetime testing:** Measure MTTF or other statistical characterization of a system under operating conditions.⁹

§ **Accelerated lifetime testing:** Measure MTTF or other statistical characterization of a system under conditions which are *more stressful* than ordinary operating conditions. Then *deduce lifetime* under ordinary conditions.

§ **Lifetime** of a device is denoted ℓ , which depends on the “stress” which the system is subject to: $\ell(s)$.

§ **Linear lifetime-stress model:**

- We adopt a piece-wise linear model:

$$\ell_m(s, c) = \begin{cases} (s - s_0)c & \text{if } s \leq s_0 \\ 0 & \text{if } s \geq s_0 \end{cases} \quad (184)$$

- $c < 0$.
- s_0 is known.
- That is, $\ell = 0$ for stress greater than s_0 .
- Lifetime increases as stress decreases below s_0 .

§ **Data:** we have measured (estimated) ℓ at stress $s_1 < s_0$: $\ell(s_1)$ is known.

§ **Best-estimated model:** Given the data, we estimate c :

$$\hat{c} = \frac{\ell(s_1)}{s_1 - s_0} \quad (185)$$

§ **Requirement:** Estimate $\ell(s)$ for $s_2 < s_1$.

10.2 Uncertainty and Robustness

§ **Uncertain information:** We expect that:

$$\ell(s_2) > \ell_m(s_2, \hat{c}) \quad (186)$$

- Lifetime at low stress, s_2 , should be greater than linear prediction.
- We don't know how much greater.

§ **Info-gap model of uncertainty:**

$$\mathcal{U}(h) = \{\ell(s_2) : 0 \leq \ell(s_2) - \ell_m(s_2, \hat{c}) \leq h\}, \quad h \geq 0 \quad (187)$$

§ **Performance function:**

- Consider linear model with coefficient c . Squared error is:

$$E^2(c) = [\ell(s_1) - \ell_m(s_1, c)]^2 + [\ell(s_2) - \ell_m(s_2, c)]^2 \quad (188)$$

⁹This example is programmed in GapZapper: Domain: Statistics, Application: Accel-Lifetime-Test-Simple.

§ Performance requirement:

$$E^2(c) \leq E_c^2 \quad (189)$$

§ Robustness of linear model $\ell_m(s, c)$:

$$\hat{h}(c, E_c) = \max \left\{ h : \left(\max_{\ell(s_2) \in \mathcal{U}(h)} E^2(c) \right) \leq E_c^2 \right\} \quad (190)$$

10.3 Evaluating the Robustness

§ We begin by evaluating the inverse of the robustness. We then invert this.

§ Let $\mu(h)$ denote the inner maximum in the robustness, eq.(190). This is the **inverse of the robustness**:

$$\mu(h) = E_c^2 \quad \text{implies} \quad \hat{h}(c, E_c) = h \quad (191)$$

Equivalently:

$$\sqrt{\mu(h)} = E_c \quad \text{implies} \quad \hat{h}(c, E_c) = h \quad (192)$$

- Plot of h vs $\sqrt{\mu(h)}$ is the same as $\hat{h}(c, E_c)$ vs E_c .

§ **Extreme values of $\ell(s_2)$ at horizon of uncertainty h .** From info-gap model of eq.(187):

$$\ell_m(s_2, \hat{c}) \leq \ell(s_2) \leq \ell_m(s_2, \hat{c}) + h \quad (193)$$

§ Evaluating $\mu(h)$:

- $\mu(h)$ occurs at one of the extreme values of $\ell(s_2)$.
- That is, $\mu(h)$ is the greater of the following two values:

$$\mu_1 = [\ell(s_1) - \ell_m(s_1, c)]^2 + [\ell_m(s_2, \hat{c}) - \ell_m(s_2, c)]^2 \quad (194)$$

$$\mu_2(h) = [\ell(s_1) - \ell_m(s_1, c)]^2 + [\ell_m(s_2, \hat{c}) + h - \ell_m(s_2, c)]^2 \quad (195)$$

• **Assumption:** Given our uncertain information, eq.(186), we will only consider slopes c which are steeper than \hat{c} :

$$c \leq \hat{c} < 0 \quad (196)$$

Recall that s_0 is known with certainty.

- Thus:

$$\ell_m(s_2, c) \geq \ell_m(s_2, \hat{c}) \quad (197)$$

- Hence the inverse robustness is:

$$\mu(h) = \begin{cases} \mu_1 & \text{if } h < 2[\ell_m(s_2, c) - \ell_m(s_2, \hat{c})] \\ \mu_2(h) & \text{else} \end{cases} \quad (198)$$

• When $c = \hat{c}$ then $\mu(h) = \mu_2(h)$. Also note that from the definition of \hat{c} one finds that $\ell(s_1) = \ell_m(s_1, \hat{c})$. Hence:

$$\mu_2(h) = h^2 \quad (199)$$

So:

$$\hat{h}(\hat{c}, E_c) = E_c \quad (200)$$

Generally, we see that eq.(198) can be inverted and, together with eq.(194) and (195), we obtain an explicit expression for the robustness. Set $\mu_2(h) = E_c^2$ in eq.(198) and solve for h to obtain:

$$\hat{h}(c, E_c) = \begin{cases} 0 & \text{if } E_c \leq \sqrt{\mu_1} \\ \sqrt{E_c^2 - [\ell(s_1) - \ell_m(s_1, c)]^2} - \ell_m(s_2, \hat{c}) + \ell_m(s_2, c) & \text{else} \end{cases} \quad (201)$$

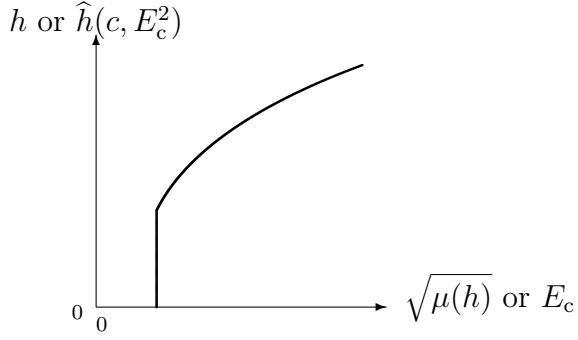


Figure 20: Illustration of the robustness based on eq.(198).

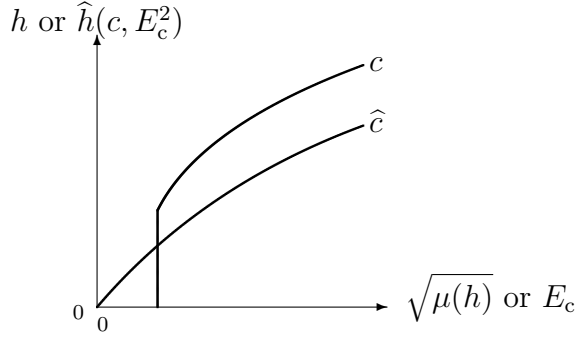


Figure 21: Illustration of the robustness based on eq.(198).

11 Review Exercises

§ The exercises in this section are not homework problems, and they do not entitle the student to credit. They will assist the student to master the material in the lecture and are highly recommended for review and self-study.

1. Explain that the variance, eq.(7), p.3, can be written as:

$$\text{var}(\bar{x}) = E(x^2) - E(x)^2 \quad (202)$$

2. What is it significant and important that theorem 1, p.3, does not depend on the probability distribution of the observations?
3. Does eq.(10), p.4, depend on the probability distribution of the observations?
4. The central limit theorem, p.4, explains the widespread (but not universal) occurrence of normal distributions. (Is this why the normal distribution is called “normal”?) Think up some examples of natural, social, or other phenomena that are described by normal distributions.
5. Why is the hypothesis test in eqs.(12) and (13), p.4, called a “two-tailed” test? Why is the hypothesis test in eqs.(14) and (15), p.5, called a “one-tailed” test?
6. Why is the level of confidence, eq.(16), p.5, defined as “the probability of obtaining a result *at least as extreme as* the observed result”, rather than as “the probability of obtaining a result *equal to* the observed result”?
7. Explain eq.(18), p.5: where does the “standard normal” distribution, $\mathcal{N}(0,1)$, come from?
8. Why is $N \geq 25$ a “rough number” on p.7? Why can’t we state a precise value for N above which the normal distribution applies exactly? For what type of distribution would you need $N \gg 25$ or $N \ll 25$?
9. The probability of 0.14, in eq.(30), p.8, may not sound small to everybody. How do you decide what is small, medium, huge, big enough, etc?
10. Do you agree with the assertion, on p.9, that $\alpha = 0.14$ is “not very convincing”? Why?
11. Will the sequential procedure in table 1, p.9, always result in α *decreasing* and thus leading to rejection of H_0 ? What would the table look like if really H_0 should be rejected?
12. When would you choose the alternative hypothesis in eq.(33) or (34) rather than (32), p.10?
13. Why can’t we answer the question following eq.(36), p.10, under hypothesis H_1 ?
14. Explain eqs.(43) and (44), p.11.
15. Why is eq.(46), p.11, true?
16. $0.21 > 0.15$, so why isn’t the answer to the question following eq.(51), p.13, obvious?
17. Regarding the conclusion that the engine is running hot, following eq.(55), p.14: Didn’t we already know this from the fact that $0.21 > 0.15$? What additional insight does the χ^2 test provide (if any)?
18. The χ^2 test is used to test “categorical” data: events that occur in different types, classes, or categories, as distinct from events that result in a real number. The tests in sections 4, p.13, 5, p.15, and 6, p.17, are very different, though they are all χ^2 tests. What are the categories in these tests?
19. What is the difference between p_u and p_a in eq.(80), p.20? Why isn’t “unacceptable” anything greater than “acceptable”? What deeper problem of meaning and judgment is this distinction trying to grapple with?

20. Are the two types of errors on p.20 the *only* errors that one can make? Are they equally severe?
21. What assumptions underlie the binomial distribution in eq.(82), p.20?
22. Why is P_I in eq.(86), p.21, call the consumer's risk? Why is P_{II} in eq.(87), p.21, call the producer's risk?
23. In the example on p.22, is a sample size $N = 200$ better or worse than a sample size of $N = 100$? Why?
24. An infinite sample size is ideal, as shown in fig. 9, p.22, but this is usually not feasible. What sample size is big enough? Why does it matter who decides?
25. In fig. 10, p.24, explain why $\delta > 0$ implies (1) pdf shifts to the right and (2) is represented by $g(t - \delta)$ rather than $g(t + \delta)$.
26. Explain eq.(93), p.25. (Recall theorems in section 1).
27. Why do α and β in eqs.(101)–(106), p.26, change in opposite directions as n increases? What is the significance of this?
28. Why the different treatment of α and β on p.26: $\alpha = 0.02$ and $\beta \leq 0.1$?