

## Problem Set on Acceptance Tests

### List of problems:

1. Sampling until the first failure, page 1 (p.7).
2. Fixed sample size, page 1 (p.8).
3. Test for equivalent performance, page 2 (p.9).
4. Run-time faults or programming errors? page 2 (p.10).
5. Milled surface, page 2.
6. Thermal expansion, page 3.
7. Vibration amplitude, page 3.
8. Dispute about the mean, page 3.
9. Robustness of the significance level, with info-gap-uncertain pdf, page 3 (p.13).
10. Robustness of the  $t$  test to distributional uncertainty, page 4 (p.13).
11. Maximum likelihood estimation, p.4 (p.14).
12. Accelerated lifetime testing: simple case, p.4 (p.16).
13. Single hypothesis test, p.5 (p.18).

1. **Sampling until the first failure** (p.7). We will sample a very large population of items and test each item for integrity. Each item will either pass (P) or fail (F). We will stop sampling when the first F occurs.

- (a) What is the probability distribution of the sample size? What is the average and standard deviation of the sample size if the fraction of F's in the population is  $p = 0.01$ ?
- (b) We have sampled  $N = 215$  items without finding an F. Test the null hypothesis:

$$H_0 : p = 0.01 \quad (1)$$

against the alternative hypothesis:

$$H_1 : p < 0.01 \quad (2)$$

- (c) What is the smallest sample size at which we reject  $H_0$  in eq.(1) at a 0.025 level of confidence, if no F has yet been obtained?
- (d) In a different batch the first F occurred at  $N = 70$  items. Test the null hypothesis:

$$H_0 : p = 0.01 \quad (3)$$

against the alternative hypothesis:

$$H_1 : p > 0.01 \quad (4)$$

Repeat this test if the first failure occurs at  $N = 10$ , and again at  $N = 2$ . What do you conclude from this?

2. **Fixed sample size.** (p.8) We will sample  $N$  items from a very large population, where each item either passes (P) or fails (F). The value of  $N$  is chosen before we begin. The probability of an F is denoted  $p$ , but its value is unknown.

- (a) What is the probability distribution of the number of F's among the  $N$  tests?
- (b) If we observe  $m$  F's among the  $N$  tests then our estimate of  $p$  would be  $m/N$ . Specifically, if we had observed 1 F among the  $N$  tests our estimate of  $p$  would be  $p_1 = 1/N$ . However, suppose that we observe no F's among the  $N$  tests. We now wish to choose between the following two hypotheses:

$$H_0 : \quad p = p_1 \quad (5)$$

$$H_1 : \quad p < p_1 \quad (6)$$

- (c) Suppose that  $p$  truly equals zero: F's do not and cannot occur. How many tests would be needed to demonstrate this? What is the implication, more generally, when testing for very small  $p$ ?

3. **Test for equivalent performance.** The performance of a system is evaluated as 'Low', 'Moderate' or 'High'. The performances of two versions of this system are summarized in table 1. Do these two systems have equivalent performance?

Version	Low	Moderate	High	Total
I	160	140	40	340
II	40	60	60	160
Totals	200	200	100	500

Table 1: Data for problem 3.

4. **Run-time faults or programming errors?** A large computer program has been run repeatedly with different input streams. Inputs occur from many different sources and at various different times during operation. Run-time faults can arise either from input errors or from programming errors. In the former case we expect a Poisson distribution in time of run-time faults. In the latter case we expect fairly consistent and recurrent run-time faults, though not completely so since some inputs may not activate some programming faults.

13 runs completed without any run-time faults; 13 runs incurred a single fault; 4 runs incurred two faults each, and 2 runs had 3 faults. These data are summarized in table 2.

# of faults/run	0	1	2	3	
# of runs	13	13	4	2	Total = 32
# of faults	0	13	8	6	Total = 27

Table 2: Data for problem 4.

With what confidence do you accept or reject the contention that the faults arose from input errors? (Hint: Estimate the coefficient of the Poisson distribution from the data, which removes an additional degree of freedom.)

5. **Milled surface.** The thickness of a milled surface is assessed as the deviation above a reference level. This deviation is measured in microns using a profilometer. In a particular sample the deviation was measured 8 times at a single point, with the following results: 12.30, 12.37, 12.37, 11.79, 12.17, 11.86, 12.31, 11.99. The desired deviation above the reference level is

11.90 microns. Assuming that the measurement errors are normally distributed, is this milled surface acceptably smooth?

6. **Thermal expansion.** The coefficient of thermal expansion of two welded metals must match, in order to prevent cracking during thermal cycling. The expansion coefficients are measured for random samples of the two metals, with the following results, in units of  $10^{-6}/\text{C}$  (strain/degree C):

Metal 1: 21, 24, 22, 23, 25, 22

Metal 2: 24, 23, 22, 25, 25, 24, 25, 23

Would you recommend welding these two metals for a thermal-cycling application? (Assume normal distribution of errors.)

7. **Vibration amplitude.** The amplitude of vibration in a milling machine increases as the tool bit wears out. An automatic monitoring system detects increased vibration in order to warn the operator of tool bit wear. Vibration levels in the milling machine are automatically categorized by the monitoring system as 'low', 'moderate' or 'high'. The vibration level is sampled periodically. During normal use of a sharp new bit, the fraction of 'low', 'moderate' and 'high' vibrations are  $p_1 = 0.85$ ,  $p_2 = 0.10$  and  $p_3 = 0.05$ , respectively.

Of the past 100 samples, 72 samples are at the 'low' vibration level, 17 are 'moderate' and 11 samples are 'high'. Should the operator change the tool bit?

8. **Dispute about the mean.** One person claims that the random variable  $x$  is uniformly distributed on the interval  $[0, 1]$ . Another person claims that  $x$  has an average greater than  $1/2$ . A large random sample of size  $N = 50$  has been made, and the sample mean is  $\bar{x} = 0.53$ . Construct and implement a statistical hypothesis test to test the first person's claim as opposed to the second person's claim.

9. **Robustness of the significance level, with info-gap-uncertain pdf.** (p.13) Consider a test between two hypotheses, based on a sample of size  $n$ . Under hypothesis  $H_0$  the sample mean is thought to be normally distributed with known mean and variance,  $\mu_0$  and  $\sigma^2/n$ . Under  $H_1$  the distribution of the sample mean is shifted to the right by a known positive quantity,  $\delta$ . We wish to distinguish between the hypotheses:

$$H_0 : \quad \mu = \mu_0 \quad (7)$$

$$H_1 : \quad \mu = \mu_0 + \delta \quad (8)$$

While the actual distribution under  $H_1$  is known to be the distribution under  $H_0$  shifted by  $\delta$ , the actual shape of the distribution in both cases is uncertain. Let  $\tilde{f}(\bar{x})$  denote the best estimate of the pdf of the sample mean under  $H_0$ , which is  $\mathcal{N}(\mu_0, \sigma^2/n)$ . An info-gap model for uncertainty in the actual distribution under  $H_0$  is:

$$\mathcal{U}(h, \tilde{f}) = \left\{ f(\bar{x}) : f(\bar{x}) \geq 0, \int_{-\infty}^{\infty} f(\bar{x}) d\bar{x} = 1, |f(\bar{x}) - \tilde{f}(\bar{x})| \leq h\tilde{f}(\bar{x}), \forall \bar{x} \right\}, \quad h \geq 0 \quad (9)$$

The critical value,  $C$ , is the rejection threshold: reject  $H_0$  if and only if  $\bar{x} > C$ . The significance level, given pdf  $f(\bar{x})$ , is the probability of falsely rejecting  $H_0$ :

$$SL(f) = \text{Prob}(\bar{x} > C | H_0) \quad (10)$$

$$= \int_C^{\infty} f(\bar{x}) d\bar{x} \quad (11)$$

The robustness of the test, for significance level  $\alpha$ , is the greatest horizon of uncertainty in the pdf, up to which the significance level does not exceed  $\alpha$ . Derive an explicit expression for the robustness, for small  $\alpha$ .

10. **Robustness of the  $t$  test to distributional uncertainty.** (p.13) Use matlab program du03.m<sup>1</sup> to explore several properties of the robustnesses to distributional uncertainty for type I and type II errors when using a  $t$  test.

(a) The robustnesses to distributional uncertainty for type I and type II errors are  $\hat{h}_0(t, \alpha^*, \alpha)$  and  $\hat{h}_1(t, \alpha^*, \beta)$ . Explore how these robustnesses vary with variation of the number of degrees of freedom of the  $t$  test. Explain your results intuitively.

(b) Level of significance,  $\alpha$ , and power,  $1 - \beta$ , trade off against each other. For instance, for an ordinary  $t$  test with 17 DoFs, levels of significance of 0.01, 0.03 and 0.05 have associated powers of 0.15, 0.31 and 0.41, respectively. Explore this trade-off at different levels of robustness. Specifically, when:

$$\hat{h}_0(t, \alpha^*, \alpha) = \hat{h}_1(t, \alpha^*, \beta) = \text{constant} \quad (12)$$

what are the values for  $\alpha$  and  $1 - \beta$  for  $\alpha^* = 0.01, 0.03$  and  $0.05$ ?

11. **Maximum likelihood estimate.** (p.14). Given a random sample,  $x = (x_1, \dots, x_n)$ , what is the maximum likelihood estimate of the parameter  $\lambda$  for each of the following distributions:

(a) Uniform distribution:

$$p(x) = \begin{cases} \frac{1}{\lambda}, & \text{if } 0 \leq x \leq \lambda \\ 0, & \text{else} \end{cases} \quad (13)$$

(b) Triangular distribution:

$$p(x) = \begin{cases} -\frac{2x}{\lambda^2} + \frac{2}{\lambda}, & \text{if } 0 \leq x \leq \lambda \\ 0, & \text{else} \end{cases} \quad (14)$$

(c) Exponential distribution:

$$p(x) = \lambda e^{-\lambda x}, \quad x \geq 0 \quad (15)$$

12. **Accelerated lifetime testing: simple case,**<sup>2</sup> p.4 (p.16). The lifetime of a device is denoted  $\ell$ , which depends on the “stress”,  $s$ , which the system is subject to:  $\ell(s)$ . We know that the lifetime is zero for any stress exceeding the value  $s_0$ . Also, it is believed that the lifetime-stress relation is roughly linear for lower stress:

$$\ell_m(s, c) = \begin{cases} (s - s_0)c & \text{if } s \leq s_0 \\ 0 & \text{if } s \geq s_0 \end{cases} \quad (16)$$

where  $c < 0$ .

We have measured the lifetime,  $\ell(s_1)$ , at stress  $s_1 < s_0$ . We wish to estimate the lifetime at lower stress  $s_2 < s_1$ . Consider the following specific case.  $s_0 = 1$ ,  $s_1 = 0.8$ ,  $\ell(s_1) = 10$  and  $s_2 = 0.4$ .

<sup>1</sup>\papers\T-Test\du03.m

<sup>2</sup>See Lecture Notes on Acceptance Testing, section 10 (acctes.tex).

(a) Use the known lifetimes to evaluate the coefficient  $\hat{c}$  in eq.(16) and to predict the lifetime at stress  $s_2$ .

(b) Our understanding indicates that the lifetime at low stress  $s_2$  will be longer than the lifetime which is predicted based on the measurement at stress  $s_1$ . However, we do not know how much longer. Using the estimated coefficient  $\hat{c}$ , calculate the robustness to lifetime uncertainty at critical errors 0, 0.2, 0.35 and 0.65. Discuss the meaning of these results.

(c) Now consider a more negative coefficient (steeper slope)  $c = 1.01\hat{c}$ . What is the lifetime prediction with this value of  $c$ ? Calculate the robustness for this  $c$  at the critical errors 0, 0.2, 0.35 and 0.65. Discuss the meaning of these results.

(d) Now consider an even more negative coefficient (steeper slope)  $c = 1.02\hat{c}$ . What is the lifetime prediction with this value of  $c$ ? Calculate the robustness for this  $c$  at the critical errors 0, 0.2, 0.35 and 0.65. Discuss the meaning of these results.

13. **Single hypothesis test**, p.5 (p.18). A particular property,  $x$ , (e.g. height, temperature, longevity, etc.) of a healthy population is a random variable with known mean and variance  $\mu$  and  $\sigma^2$ . The central limit theorem asserts that, regardless of the distribution of  $x$ , the mean,  $\bar{x}$ , of a large random sample of size  $N$  is normally distributed with mean  $\mu$  and variance  $\sigma^2/N$ . Given an observed value of the sample mean,  $\bar{x}_{\text{obs}}$ , we want to test the hypothesis that the population is healthy:

$$H_0 : \bar{x} \sim \mathcal{N}(\mu, \sigma^2/N) \quad (17)$$

We will reject  $H_0$  if  $\bar{x}_{\text{obs}}$  is an implausible value, conditioned on  $H_0$ . Specifically, we reject  $H_0$  if, conditioned on  $H_0$ :

$$\text{Prob} \left( \left| \frac{\bar{x} - \mu}{\sigma/\sqrt{N}} \right| > \left| \frac{\bar{x}_{\text{obs}} - \mu}{\sigma/\sqrt{N}} \right| ; H_0 \right) \leq \alpha \quad (18)$$

where  $\alpha$  is a 'level of significance'. If  $H_0$  holds and if the sample is statistically random, then the distribution of  $\frac{\bar{x} - \mu}{\sigma/\sqrt{N}}$  would be standard normal,  $\mathcal{N}(0, 1)$ . Let  $\Phi(z)$  denote the cumulative probability distribution (CPD) for  $\mathcal{N}(0, 1)$ . The problem is that we are unsure that the sample is truly random: statistically independent measurements from the same population. Thus we are unsure that the true distribution of  $\frac{\bar{x} - \mu}{\sigma/\sqrt{N}}$ , call it  $F(\cdot)$ , is actually  $\Phi(\cdot)$ . We represent this uncertainty with the following info-gap model in which we introduce a simplifying assumption that the distributions are symmetric around the origin:

$$\mathcal{U}(h) = \left\{ F(z) : F(-\infty) = 0, F(\infty) = 1, F(z) = 1 - F(-z), \frac{dF}{dz} \geq 0, |F(z) - \Phi(z)| \leq h \right\}, \quad h \geq 0 \quad (19)$$

- (a) We have an observed value of the sample mean,  $\bar{x}_{\text{obs}}$ . Suppose that eq.(18) holds based on this observation, implying that we should reject  $H_0$ . How much can the distribution of  $\frac{\bar{x} - \mu}{\sigma/\sqrt{N}}$  deviate from  $\mathcal{N}(0, 1)$  without changing this decision? That is, derive the robustness as a function of the rejection threshold,  $\alpha$ .
- (b) In contrast to part 13a, suppose that we have observed a sample mean,  $\bar{x}_{\text{obs}}$ , for which eq.(18) does not hold, implying that we should accept  $H_0$ . Derive the robustness as a function of the rejection threshold,  $\alpha$ .

## Solutions to Homework on Acceptance Tests

### List of problem solutions:

1. Sampling until the first failure, page 7.
2. Fixed sample size, page 1.
3. Test for equivalent performance, page 9.
4. Run-time faults or programming errors? page 10.
5. Milled surface, page 11.
6. Thermal expansion, page 12.
7. Vibration amplitude, page 12.
8. Dispute about the mean, page 13.
9. Robustness of the significance level, with info-gap-uncertain pdf, page 13 (p.3).
10. Robustness of the  $t$  test to distributional uncertainty, page 13 (p.4).
11. Maximum likelihood estimation, p.14 (p.4).
12. Accelerated lifetime testing: simple case, p.16 (p.4).
13. Single hypothesis test, p.18 (p.5).

**Solution to problem 1.**

**1a**  $x$  = number of samples to the first failure, and has a geometric distribution, where  $p$  is the probability of an F:

$$f(x) = (1 - p)^{x-1}p, \quad x = 1, 2, \dots \quad (20)$$

Check normalization:

$$\sum_{x=1}^{\infty} f(x) = p \sum_{x=1}^{\infty} (1 - p)^{x-1} = p \sum_{x=0}^{\infty} (1 - p)^x = p \frac{1}{1 - (1 - p)} = 1 \quad (21)$$

The mean and variance of  $x$  are:

$$E(x) = \frac{1}{p}, \quad \text{var}(x) = \frac{1 - p}{p^2} \quad (22)$$

So:

$$E(x) = \frac{1}{p} = \frac{1}{0.01} = 100 \quad (23)$$

$$\sqrt{\text{var}(x)} = \sqrt{\frac{1 - p}{p^2}} \approx \frac{1}{p} = 100 \quad (24)$$

**1b** Our observation is  $x = N = 215$  samples without observing an “F”. A large  $x$  is evidence against  $H_0$ . The level of significance is the probability of an equally or more extreme result (impugning  $H_0$ ) conditioned on  $H_0$ . Hence the level of significance is:

$$\alpha = \Pr[x \geq N | H_0] = \sum_{x=N}^{\infty} f(x) = p \sum_{x=N}^{\infty} (1 - p)^{x-1} = p(1 - p)^{N-1} \underbrace{\sum_{x=0}^{\infty} (1 - p)^x}_{\frac{1}{1 - (1 - p)} = \frac{1}{p}} = (1 - p)^{N-1} \quad (25)$$

Under  $H_0$ :  $p = 0.01$  so  $\alpha = (0.99)^{N-1} = (0.99)^{214} = 0.116$ . This is not small so we do not (yet) reject  $H_0$ .

**1c** We require  $\alpha = (0.99)^{N-1} = 0.025$  hence:

$$N = 1 + \frac{\ln 0.025}{\ln 0.99} = 368.04 \quad (26)$$

So, at  $N = 369$  we reject  $H_0$  at  $\alpha = 0.025$  level of significance.

**1d** The first F occurred at item  $N = 70$ . A small  $x$  impugns  $H_0$ . (This is the reverse of part 1b because  $H_1$  in eq.(4) is the reverse of  $H_1$  in eq.(2).) Hence the level of significance is:

$$\alpha = \Pr[x \leq N | H_0] = \sum_{x=1}^N f(x) = p \sum_{x=1}^N (1-p)^{x-1} = p \frac{(1-p)^N - 1}{(1-p) - 1} = 1 - (1-p)^N \quad (27)$$

Under  $H_0$ :  $p = 0.01$  so  $\alpha = 1 - (0.99)^N = 1 - (0.99)^{70} = 0.505$ . This is not small so we do not reject  $H_0$ .

Repeat this with the first F at  $N = 10$ , and we find that the level of significance is  $\alpha = 1 - (0.99)^N = 1 - (0.99)^{10} = 0.0956$ . This is not small so we still do not reject  $H_0$ .

Repeat this with the first F at  $N = 2$ , and we find that the level of significance is  $\alpha = 1 - (0.99)^N = 1 - (0.99)^2 = 0.0199$ . This is small (but certainly not tiny) so we reject  $H_0$ .

The general conclusion is that it is “difficult” to reject  $H_0$ . The “problem” is that the standard deviation is nearly equal to the mean, as seen in eqs.(23) and (24).

### Solution to problem 2.

**2a** The probability of exactly  $m$  F's among  $N$  tests is the binomial distribution:

$$f(m|N) = \binom{N}{m} p^m (1-p)^{N-m} \quad (28)$$

where the binomial coefficient is:

$$\binom{N}{m} = \frac{N!}{m!(N-m)!} \quad (29)$$

The mean and variance of the binomial distribution are:

$$E(m) = Np, \quad \text{var}(m) = Np(1-p) \quad (30)$$

Thus, for small  $p$ , we see:

$$\frac{\sqrt{\text{var}(m)}}{E(m)} = \frac{\sqrt{Np(1-p)}}{Np} = \sqrt{\frac{(1-p)}{Np}} \approx \frac{1}{\sqrt{Np}} \quad (31)$$

For instance, this ratio equals 1 for  $N = 100$  and  $p = 0.01$ . In other words, the binomial distribution tends to be wide.

**2b** The level of significance,  $\alpha$ , is the probability of an equally or more extreme result than observed (less than or equal to 0). A small value of  $\alpha$  impugns  $H_0$ . Thus, with  $p_1 = 1/N$ :

$$\alpha = \Pr(x = 0 | H_0) = f(0|N) = \binom{N}{0} p_1^0 (1-p_1)^N = \left(1 - \frac{1}{N}\right)^N \quad (32)$$

See results in table 3. The confidence converges very slowly (and is not even monotonic at very large  $N$ ) as the sample size rises.

**2c** If you are looking for something that isn't there, you will never find it. And even if  $p$  is positive but very very tiny, you are very very unlikely to ever find an F. The implication is that estimating  $p$  based on fitness tests is infeasible when  $p$  is tiny. The best one can hope for is to establish confidence that  $p$  is no greater than a specified value. We explored this in problem 1d. We saw that this also is often inconclusive. If  $p$  is very small then fitness testing should probably be augmented with (or replaced by) other methods of system analysis. For example, attempting to identify possible failure mechanisms, and verifying their rarity.



$N$	$\alpha$
2	0.25
5	0.3277
10	0.3487
20	0.3585
30	0.3617
50	0.3642
100	0.3666
$10^3$	0.3677
$10^6$	0.3679
$10^{15}$	0.3682
$\infty$	1

Table 3: Levels of significance, eq.(32).

**Solution to problem 3.** The systems are equivalent if the rows and columns of the data, table 1, are statistically independent. We will use the  $\chi^2$  test to test this hypothesis.

$c$  = number of columns = 3.

$r$  = number of rows = 2.

$p_{ij}$  = probability of an outcome in row  $i$  and column  $j$ . We can estimate this probability as:

$$p_{ij} = \frac{n_{ij}}{N} \quad (33)$$

where  $n_{ij}$  = number of outcomes in row  $i$  and column  $j$ .

$p_{\bullet i}$  = probability of an outcome in column  $i$ .

$p_{j\bullet}$  = probability of an outcome in row  $j$ .

The hypothesis of statistical independence of rows and columns is:

$$H_0 : p_{ij} = p_{i\bullet}p_{\bullet j} \quad (34)$$

The alternative hypothesis is:

$$H_1 : H_0 \text{ is false} \quad (35)$$

$$p_{\bullet 1} = \frac{200}{500} = 0.4, \quad p_{\bullet 2} = \frac{200}{500} = 0.4, \quad p_{\bullet 3} = \frac{100}{500} = 0.2, \quad p_{1\bullet} = \frac{340}{500} = 0.68, \quad p_{2\bullet} = \frac{160}{500} = 0.32 \quad (36)$$

$$\chi^2 = \sum_{i=1}^c \sum_{j=1}^r \frac{(n_{ij} - Np_{\bullet i}p_{j\bullet})^2}{Np_{\bullet i}p_{j\bullet}} = 49.63 \quad (37)$$

DOF = number of categories – number of constraints:

$$\text{DOF} = (r - 1)(c - 1) = 2 \quad (38)$$

$$\chi_{2,0.005}^2 = 10.6 \implies \Pr(\chi_2^2 \geq 10.6) = 0.005 \implies \Pr(\chi_2^2 \geq 49.63) \ll 0.005 \quad (39)$$

So we strongly reject  $H_0$ , with level of confidence greater than  $1 - 0.005$ .

**Solution to problem 4.** We must test the following hypotheses:

$H_0$ : Faults are due to input errors.

$H_1$ : Faults are due to programming errors.

These hypotheses are **roughly** equivalent to:

$H_0$ : The number of faults per run is a Poisson random variable.

$H_1$ : The number of faults per run is **not** a Poisson random variable.

If  $H_0$  is true, then  $n$ , the number of faults per run, will be distributed as:

$$P_n = \frac{e^{-\lambda} \lambda^n}{n!} \quad (40)$$

where  $\lambda$  = average number of faults per run.

We do not know  $\lambda$  but we can estimate it from the data:

$$\lambda = \frac{27}{32} \approx 0.84 \quad (41)$$

We can test  $H_0$  against  $H_1$  with a  $\chi^2$  test.

We have four categories:

Faults per run	Probability conditioned on $H_0$
0	$P_0 = e^{-\lambda} = 0.432$
1	$P_1 = \lambda e^{-\lambda} = 0.363$
2	$P_2 = \frac{\lambda^2 e^{-\lambda}}{2!} = 0.152$
3 or more	$\sum_{n=3}^{\infty} P_n = 0.053$

Table 4: Categories for  $\chi^2$  test for problem 4.

DOF = number of categories - number of constraints:

$$4 - 2 = 2 \quad (42)$$

1 constraint: normalization. 1 constraint: estimate of  $\lambda$ .

Define:

$n_k$  = number of runs with  $k$  faults.

$N$  = total number of runs.

$P_k$  = probability of run with  $k$  faults.

So we calculate the  $\chi^2$  statistic as:

$$\chi^2 = \sum_{k=1}^4 \frac{(n_k - NP_k)^2}{NP_k} = 0.422 \quad (43)$$

The level of significance in our test is  $\alpha$  = probability of a more extreme result, conditioned upon  $H_0$ :

$$\alpha = \Pr(\chi_2^2 \geq \chi_{\text{obs}}^2 | H_0) \quad (44)$$

$$= \Pr(\chi_2^2 \geq 0.422 | H_0) \quad (45)$$

$$> 0.5 \quad (46)$$

This is **not small**, so we cannot reject  $H_0$ .

**Solution to problem 5.** Use the  $t$  statistic to test between the following two hypotheses:

$$H_0: x = 11.90 \quad (47)$$

$$H_1: x \neq 11.90 \quad (48)$$

The data result in:

$$N = 8, \quad \bar{x} = 12.145, \quad s^2 = 0.054943, \quad s = 0.23440 \quad (49)$$

Thus the  $t$  statistic, conditioned on  $H_0$ , is:

$$t = \frac{\bar{x} - 11.90}{s/\sqrt{8}} \sim t_{(7)} \quad (50)$$

The observed value of the statistic is:

$$t_o = \frac{12.145 - 11.90}{0.2344/\sqrt{8}} = 2.956 \quad (51)$$

If  $t_o$  is “greatly different” from zero, then reject  $H_0$ . We use the level of confidence to evaluate the degree of deviation of  $t_o$  from zero:

$$\alpha = \Pr(|t| \geq |t_o| \mid H_0) \quad (52)$$

$$= 2 \times [\Pr(t_{(7)} \geq 2.956)] \quad (53)$$

$$= 2 \times (1 - 0.99) = 0.02 \quad (54)$$

This is small so we reject  $H_0$ .

**Solution to problem 6.** Use the  $t$  test to compare two means. The null- and alternative hypotheses are:

$$H_0 : \mu_1 = \mu_2 \quad (55)$$

$$H_1 : \mu_1 \neq \mu_2 \quad (56)$$

The basic results are:

$$\text{Metal 1: } N_1 = 6 \quad \bar{x}_1 = 22.833 \quad s_1^2 = 2.1666 \quad s_1 = 1.47194$$

$$\text{Metal 2: } N_2 = 8 \quad \bar{x}_2 = 23.875 \quad s_2^2 = 1.2679 \quad s_2 = 1.1260$$

Define:

$$\Delta = \bar{x}_1 - \bar{x}_2 \quad (57)$$

for which the observed value is:

$$\Delta_o = -1.042 \quad (58)$$

Under hypothesis  $H_0$ :

$$E(\Delta) = 0, \quad \text{var}(\Delta) = \text{var}(\bar{x}_1) + \text{var}(\bar{x}_2) \approx \frac{s_1^2}{N_1} + \frac{s_2^2}{N_2} = \sigma_\Delta^2 = 0.5196 \quad (59)$$

Thus  $\sigma_\Delta = 0.7208$ .

Since the errors have a normal distribution we conclude, under  $H_0$ , that:

$$\Delta \sim t_{(N_1+N_2-2)} \quad (60)$$

A large value of  $|\Delta|$  is evidence against  $H_0$ . The level of significance is:

$$\alpha = \Pr\left(|\Delta| \geq |\Delta_o| \mid H_0\right) \quad (61)$$

$$= \Pr\left(\frac{|\Delta|}{\sigma_\Delta} \geq \frac{|\Delta_o|}{\sigma_\Delta} \mid H_0\right) \quad (62)$$

$$\approx 2 \times 0.08 = 0.16 \quad (63)$$

$\alpha$  is not small so we do not reject  $H_0$ .

**Solution to problem 7.** Use the  $\chi^2$  test with  $K = 3$  categories. In normal operation the probabilities of outcomes in each of the three categories are:

$$p_1^o = 0.85, \quad p_2^o = 0.10, \quad p_3^o = 0.05 \quad (64)$$

We wish to test between the following hypotheses:

$$H_0 : \quad p_i = p_i^o, \quad i = 1, \dots, 3 \quad (65)$$

$$H_1 : \quad H_0 \text{ is false} \quad (66)$$

The  $\chi^2$  statistic is:

$$\chi^2 = \sum_{i=1}^K \frac{(N_i - Np_i^o)^2}{Np_i^o} \quad (67)$$

where  $N_i$  is the number of outcomes in category  $i$  and  $N$  is the total number of outcomes. This statistic is distributed as  $\chi_{(3-1)}^2$  if  $H_0$  holds.

The observed statistic is:

$$\chi_o^2 = 14.08 \quad (68)$$

The level of significance is:

$$\alpha = \Pr(\chi_{(2)}^2 \geq \chi_o^2 \mid H_0) = \Pr(\chi_{(2)}^2 \geq 14.08) \leq 1 - 0.999 = 0.001 \quad (69)$$

This is very small so we reject  $H_0$  and change the tool bit.

**Solution to problem 8.** The two hypotheses are:

$$H_0 \quad x \sim U[0, 1] \quad (70)$$

$$H_1 \quad E(x) > 1/2 \quad (71)$$

Under  $H_0$  the mean and variance of  $x$  are  $1/2$  and  $1/12$ . Since the sample is large, the central limit theorem implies that, under  $H_0$ ,  $\bar{x} \sim \mathcal{N}(\frac{1}{2}, \frac{1/12}{N})$ . Thus the level of significance is:

$$\alpha = \text{Prob}(\bar{x} \geq \bar{x}_o | H_0) \quad (72)$$

$$= \text{Prob} \left( \frac{\bar{x} - \mu_{\bar{x}}}{\sigma_{\bar{x}}} \geq \frac{\bar{x}_o - \mu_{\bar{x}}}{\sigma_{\bar{x}}} \mid H_0 \right) \quad (73)$$

$$= 1 - \Phi \left( \frac{\bar{x}_o - \mu_{\bar{x}}}{\sigma_{\bar{x}}} \right) \quad (74)$$

$$= 1 - \Phi \left( \frac{0.53 - 0.5}{1/\sqrt{12} \times 50} \right) = 1 - \Phi(0.735) = 1 - 0.7673 = 0.2327 \quad (75)$$

This is not large so we cannot reject  $H_0$ .

**Solution to problem 9.** (p.3) The robustness is the greatest horizon of uncertainty up to which all pdf's result in significance level no less than  $\alpha$ :

$$\hat{h}(\alpha) = \max \left\{ h : \left( \max_{f \in \mathcal{U}(h, \tilde{f})} SL(f) \right) \leq \alpha \right\} \quad (76)$$

Let  $M(h)$  denote the inner maximum in eq.(76). For small significance level,  $\alpha \ll 1$ , this maximum occurs when the upper tail is as fat as possible at horizon of uncertainty  $h$ :

$$f(\bar{x}) = (1 + h)\tilde{f}(\bar{x}) \quad (77)$$

Thus:

$$M(h) = (1 + h) \int_C^\infty \tilde{f}(\bar{x}) d\bar{x} \quad (78)$$

The robustness is the greatest value of  $h$  for which:

$$M(h) = \alpha \quad (79)$$

Let  $\tilde{\alpha}$  denote the best estimate of the significance level:

$$\tilde{\alpha} = \int_C^\infty \tilde{f}(\bar{x}) d\bar{x} \quad (80)$$

Combining eqs.(78)–(80) we find the robustness to be:

$$\hat{h}(\alpha) = \frac{\alpha}{\tilde{\alpha}} - 1 \quad (81)$$

or zero if this expression is negative. Note:

- The robustness depends on the sample size via  $\tilde{\alpha}$ , which depends on the sample size via the estimated distribution  $\tilde{f}$ .
- The robustness is zero for the estimated significance,  $\tilde{\alpha}$ .
- The robustness is positive for the significance greater (lower confidence) than  $\tilde{\alpha}$ .

**Solution to problem 10.** (p.4)

(a) The robustnesses to distributional uncertainty for type I and type II errors are  $\hat{h}_0(t, \alpha^*, \alpha)$  and  $\hat{h}_1(t, \alpha^*, \beta)$ . The  $t$  test is constructed with the nominal pdfs for level of significance  $\alpha^*$ .  $\hat{h}_0(t, \alpha^*, \alpha)$

DoF	$\alpha^* = 0.01$		$\alpha^* = 0.03$		$\alpha^* = 0.05$	
	$\alpha$ for $\widehat{h}_0 = 0.04$	$1 - \beta$ for $\widehat{h}_1 = 0.2$	$\alpha$ for $\widehat{h}_0 = 0.04$	$1 - \beta$ for $\widehat{h}_1 = 0.2$	$\alpha$ for $\widehat{h}_0 = 0.04$	$1 - \beta$ for $\widehat{h}_1 = 0.2$
5	0.1210	0	0.1273	0.0674	0.1411	0.1510
17	0.1004	0.0480	0.1117	0.1608	0.1274	0.2392
50	0.0956	0.0752	0.1080	0.1884	0.1240	0.2627

Table 5: Results for problem 10(a).

is the greatest horizon of uncertainty up to which the probability of falsely rejecting  $H_0$  is no greater than  $\alpha$ .  $\widehat{h}_1(t, \alpha^*, \beta)$  is the greatest horizon of uncertainty up to which the probability of falsely accepting  $H_0$  is no greater than  $\beta$ . The “power” is defined as  $1 - \beta$ , the probability of falsely rejecting  $H_1$ .

Table 5 shows values  $\alpha$  for which  $\widehat{h}_0(t, \alpha^*, \alpha) = 0.04$ , and values of  $1 - \beta$  for which  $\widehat{h}_1(t, \alpha^*, \beta) = 0.2$ , for various degrees of freedom and various values of  $\alpha^*$ . From this table we observe the following two facts:

1. At fixed  $\alpha^*$  (the nominal level of significance),  $\widehat{h}_0(t, \alpha^*, \alpha)$  reaches robustness of 0.04 at lower  $\alpha$  (more significant rejection) as the DoF increases. This means that the robustness increases as the DoF increases.
2. At fixed  $\alpha^*$  (the nominal level of significance),  $\widehat{h}_1(t, \alpha^*, \beta)$  reaches robustness of 0.2 at larger  $1 - \beta$  (greater power) as the DoF increases. This means that the robustness increases as the DoF increases.

(b) We seek the values of  $\alpha$  and  $1 - \beta$  at which:

$$\widehat{h}_0(t, \alpha^*, \alpha) = \widehat{h}_1(t, \alpha^*, \beta) = 0, 0.04, \text{ and } 0.08 \quad (82)$$

$\widehat{h}_0 = \widehat{h}_1$	$\alpha^* = 0.01$	0.03	0.05	$\alpha^* = 0.01$	0.03	0.05
0	$\alpha = 0.0100$	0.0300	0.0500	$1 - \beta = 0.1505$	0.3066	0.4068
0.04	0.1004	0.1117	0.1274	0.1172	0.2646	0.3605
0.08	0.1789	0.1828	0.1942	0.0946	0.2333	0.3248

Table 6: Results for problem 10(b).

Results are shown in table 6 for 17 DoFs. We observe:

1. Robust level of significance ( $\alpha$ ) trades-off against robust power ( $1 - \beta$ ) at all levels of robustness.
2. The trade-off is less severe at high robustness than at low robustness. For instance, at robustness of 0.04, the level of significance deteriorates less than at robustness of zero, for approximately the same relative improvement in power.

**Solution to problem 11.** (p.4). Define  $s(x) = 1$  if  $x \geq 0$  and zero otherwise.

(a) The likelihood function is the product of pdf's:

$$L(x|\lambda) = \prod_{i=1}^n p(x_i) = \prod_{i=1}^n \frac{1}{\lambda} s(\lambda - x_i) s(x_i) \quad (83)$$

$$= \frac{1}{\lambda^n} \text{ if } x_i \in [0, \lambda] \forall i \quad (84)$$

and zero otherwise.

Thus the MLE is:

$$\hat{\lambda} = \max_i x_i \quad (85)$$

The MLE is not defined if there is a negative measurement, which simply refutes this pdf.

(b) The likelihood function is the product of pdf's:

$$L(x|\lambda) = \prod_{i=1}^n p(x_i) = \prod_{i=1}^n \left( -\frac{2x_i}{\lambda^2} + \frac{2}{\lambda} \right) s(\lambda - x_i) s(x_i) \quad (86)$$

$$\frac{\partial L}{\partial \lambda} = \frac{2}{\lambda^2} \prod_{i=1}^n \left( \frac{2x_i}{\lambda} - 1 \right) s(\lambda - x_i) s(x_i) \quad (87)$$

Thus the MLE is for  $\lambda$  to equal  $2x_i$  for some measurement  $i$ . Also, we require  $\lambda \geq \max_i x_i$ . This allows several possible solutions. Consider the 2nd derivative:

$$\frac{\partial^2 L}{\partial \lambda^2} = \frac{4}{\lambda^3} \prod_{i=1}^n \left( -\frac{3x_i}{\lambda} + 1 \right) s(\lambda - x_i) s(x_i) \quad (88)$$

For a maximum we require that an odd number of terms be negative.

The MLE is not defined if there is a negative measurement, which simply refutes this pdf.

For instance:

Suppose  $n = 1$  where  $0 \leq x_1$ . Then  $\lambda = 2x_1$  implies  $L' = 0$  and  $L'' < 0$ . So this is the MLE.

Suppose  $n = 2$  where  $0 \leq x_1 \leq x_2$ , where  $-\frac{3x_1}{2x_2} + 1 > 0$ . Then  $\lambda = 2x_2$  implies  $L' = 0$  and  $L'' < 0$ . So this is the MLE.

Suppose  $n = 2$  where  $0 \leq x_1 \leq x_2$ , where  $-\frac{3x_1}{2x_2} + 1 < 0$ . Then  $\lambda = 2x_2$  implies  $L' = 0$  and  $L'' > 0$ . So this is *not* the MLE.

(c) The likelihood function is the product of pdf's:

$$L(x|\lambda) = \prod_{i=1}^n p(x_i) = \lambda^n \prod_{i=1}^n e^{-\lambda x_i} s(x_i) = \lambda^n e^{-\lambda X} \prod_{i=1}^n s(x_i) \quad (89)$$

where  $X = \sum_{i=1}^n x_i$ .

$$\frac{\partial L}{\partial \lambda} = \lambda^{n-1} e^{-\lambda X} (n - \lambda X) \prod_{i=1}^n s(x_i) \quad (90)$$

Thus the MLE is:

$$\hat{\lambda} = \frac{n}{X} \quad (91)$$

The MLE is not defined if there is a negative measurement, which simply refutes this pdf.

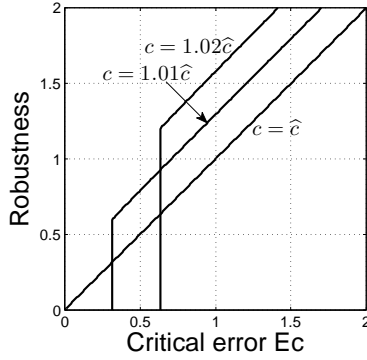
**Solution to problem 12. (p.4).**

Figure 1: Robustness curves for problem 12.

**(a)** Given the data, we estimate  $c$ :

$$\hat{c} = \frac{\ell(s_1)}{s_1 - s_0} = \frac{10}{0.8 - 1} = -50 \quad (92)$$

Thus the lifetime prediction is:

$$\ell(s_2) = (s_2 - s_0)\hat{c} = (0.4 - 1)(-50) = 30 \quad (93)$$

**(b)** We know that the robustness curve for the estimated slope is a straight line at  $45^\circ$ :

$$\hat{h}(\hat{c}, E_c) = E_c \quad (94)$$

Thus the robustness based on  $\hat{c}$  is shown in the 2nd column of table 7

$E_c$	$\hat{h}(\hat{c}, E_c)$	$\hat{h}(1.01\hat{c}, E_c)$	$\hat{h}(1.02\hat{c}, E_c)$
0	0	0	0
0.2	0.2	0	0
0.35	0.35	0.64	0
0.65	0.65	0.94	1.2

Table 7: Data for problem 3.

**(c)** The lifetime prediction with  $c = 1.01\hat{c} = -50.5$  is:

$$\ell(s_2) = (s_2 - s_0)c = (0.4 - 1)(-50.5) = 30.3 \quad (95)$$

A 1% increase in predicted lifetime.

To calculate the robustness we use:

$$\hat{h}(c, E_c) = \begin{cases} 0 & \text{if } E_c \leq \sqrt{\mu_1} \\ \sqrt{E_c^2 - [\ell(s_1) - \ell_m(s_1, c)]^2} - \ell_m(s_2, \hat{c}) + \ell_m(s_2, c) & \text{else} \end{cases} \quad (96)$$

where:

$$\mu_1 = [\ell(s_1) - \ell_m(s_1, c)]^2 + [\ell_m(s_2, \hat{c}) - \ell_m(s_2, c)]^2 \quad (97)$$



We find  $\sqrt{\mu_1} = 0.32$ , which is the value at which the robustness curve lifts off the  $E_c$  axis. Thus  $\widehat{h}(c, 0) = \widehat{h}(c, 0.2) = 0$  while  $\widehat{h}(c, 0.35)$  and  $\widehat{h}(c, 0.65)$  will be positive:

$$\widehat{h}(c, 0.35) = \sqrt{0.35^2 - [10 - 10.1]^2} - 30.0 + 30.3 = 0.635 \quad (98)$$

$$\widehat{h}(c, 0.65) = \sqrt{0.65^2 - [10 - 10.1]^2} - 30.0 + 30.3 = 0.942 \quad (99)$$

See column 3 of table 7.

**(d)** The lifetime prediction with  $c = 1.02\widehat{c} = -51$  is:

$$\ell(s_2) = (s_2 - s_0)c = (0.4 - 1)(-51) = 30.6 \quad (100)$$

A 2% increase in predicted lifetime.

We find  $\sqrt{\mu_1} = 0.63$ , which is the value at which the robustness curve lifts off the  $E_c$  axis. Thus  $\widehat{h}(c, 0) = \widehat{h}(c, 0.2) = \widehat{h}(c, 0.35) = 0$  while  $\widehat{h}(c, 0.65)$  will be positive:

$$\widehat{h}(c, 0.65) = \sqrt{0.65^2 - [10 - 10.2]^2} - 30.0 + 30.6 = 1.22 \quad (101)$$

See column 4 of table 7.

**Solution to problem 13.** (p.5).

**(13a)** Define the random variable  $z = \frac{\bar{x} - \mu}{\sigma/\sqrt{N}}$  and denote its CPD as  $F(z)$ . Let  $P_{\text{rej}}(F)$  denote the probability of rejection in eq.(18) on p.5. Define  $z_{\text{obs}} = \frac{\bar{x}_{\text{obs}} - \mu}{\sigma/\sqrt{N}}$ . Thus:

$$P_{\text{rej}}(F) = \text{Prob}(|z| \geq |z_{\text{obs}}|) = F(-|z_{\text{obs}}|) + [1 - F(|z_{\text{obs}}|)] \quad (102)$$

$$= 2[1 - F(|z_{\text{obs}}|)] \quad (103)$$

where eq.(103) exploits the symmetry around the origin of the CPD's in the info-gap model.

The definition of the robustness is:

$$\hat{h}_1(\alpha) = \max \left\{ h : \left( \max_{F \in \mathcal{U}(h)} 2[1 - F(|z_{\text{obs}}|)] \right) \leq \alpha \right\} \quad (104)$$

Let  $m(h)$  denote the inner maximum, which occurs when  $F(|z_{\text{obs}}|)$  is as small as possible at horizon of uncertainty  $h$ :

$$F(|z_{\text{obs}}|) = \max \left[ \frac{1}{2}, \Phi(|z_{\text{obs}}|) - h \right] \quad (105)$$

where the ' $\frac{1}{2}$ ' comes from the symmetry of the CPD's in the info-gap model. Thus, for  $h \leq \Phi(|z_{\text{obs}}|) - \frac{1}{2}$ :

$$m(h) = 2[1 - \Phi(|z_{\text{obs}}|) + h] \leq \alpha \implies \hat{h}_1(\alpha) = \frac{\alpha - 2[1 - \Phi(|z_{\text{obs}}|)]}{2} \quad (106)$$

or zero if this is negative.<sup>3</sup> Note that, for all  $0 \leq \alpha \leq 1$ ,  $\hat{h}$  does not exceed  $\Phi(|z_{\text{obs}}|) - \frac{1}{2}$  so we need not evaluate  $m(h)$  for larger values of  $h$ .

The robustness curve in eq.(106) is plotted in fig. 2 for the following parameter values:  $\mu = 1$ ,  $\sigma = 1.3$  and  $N = 30$ . The observed sample mean is  $\bar{x}_{\text{obs}} = 1.5$  for which  $z_{\text{obs}} = 2.1066$  and  $\Phi(|z_{\text{obs}}|) = 0.9824$ .

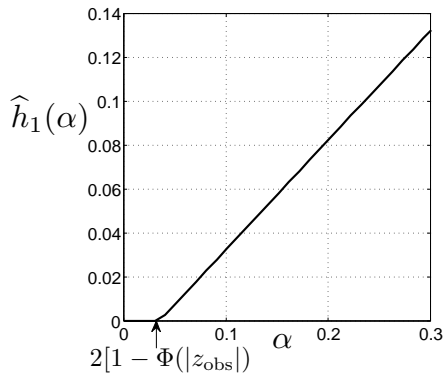


Figure 2: Robustness curve for problem 13a.

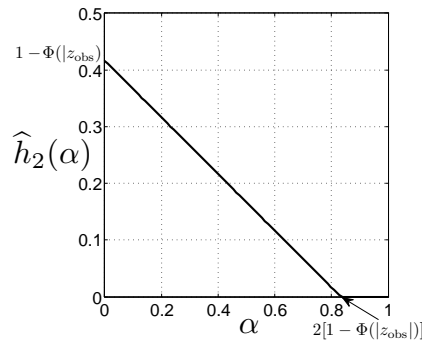


Figure 3: Robustness curve for problem 13b.

**(13b)** In analogy to eq.(104), the definition of the robustness is:

$$\hat{h}_2(\alpha) = \max \left\{ h : \left( \min_{F \in \mathcal{U}(h)} 2[1 - F(|z_{\text{obs}}|)] \right) \geq \alpha \right\} \quad (107)$$

Let  $m(h)$  denote the inner minimum, which occurs when  $F(|z_{\text{obs}}|)$  is as large as possible at horizon of uncertainty  $h$ :

$$F(|z_{\text{obs}}|) = \min [1, \Phi(|z_{\text{obs}}|) + h] \quad (108)$$

<sup>3</sup>Use matlab program \lectures\reltest\prob12.m to calculate robustness curves for both parts 13a and 13b

Thus, for  $h \leq 1 - \Phi(|z_{\text{obs}}|)$ :

$$m(h) = 2[1 - \Phi(|z_{\text{obs}}|) - h] \geq \alpha \implies \boxed{\hat{h}_2(\alpha) = \frac{2[1 - \Phi(|z_{\text{obs}}|)] - \alpha}{2}} \quad (109)$$

or zero if this is negative. Note that, for all  $0 \leq \alpha \leq 1$ ,  $\hat{h}$  does not exceed  $1 - \Phi(|z_{\text{obs}}|)$  so we need not evaluate  $m(h)$  for larger values of  $h$ .

The robustness curve in eq.(109) is plotted in fig. 2 for the following parameter values:  $\mu = 1$ ,  $\sigma = 1.3$  and  $N = 30$ . The observed sample mean is  $\bar{x}_{\text{obs}} = 1.05$  for which  $z_{\text{obs}} = 0.2107$  and  $\Phi(|z_{\text{obs}}|) = 0.5834$ .