# Statistical Techniques for CS/BIT 2021-1B

## Practice test #1

Time: 2hrs 15min

**Solutions.**

**1.** **a.** Highest weekly sale is 82.0 and the corresponding $z$-score is $\frac{82-52.56}{9.84} \approx 2.99$.
According to the empirical rule, $z$-scores which are larger than 2 in absolute value ($|z| > 2$) occur with a probability of 5%, if the distribution is approximately normal. Furthermore, $z$-scores which are larger than 3 in absolute value occur with a probability of 0.3%. So for $z$-scores of at least 2.99, the probability of occurrence lies between 0.3% and 5% (more precise: 1.4%), not exceptional in case of 45 observations.

**b.** $20^{\text{th}}$ percentile is $\frac{x_{(9)}+x_{(10)}}{2} = 43.75$.
$80^{\text{th}}$ percentile is $\frac{x_{(36)}+x_{(37)}}{2} = 61.00$.

$n = 45$
Minimum: 34.0
$Q_1 = x_{(12)} = 45.7$
$m = x_{(23)} = 50.6$
$Q_3 = x_{(34)} = 58.9$
Maximum: 82.0

The $1.5\times$IQR range is: $(25.9, 78.7)$. All observations except for the highest weekly sale $(x_{(45)} = 82)$ fall in this range. So there is only one outlier acoording to the $1.5\times$IQR rule.

**c.** –The skewness coefficient 0.601 and the kurtosis 0.573 are close enough to the reference values of the normal distribution (0 and 0). (Note here that SPSS reports *excess kurtosis*, which is kurtosis$-3$, so the reference value is indeed zero).
–The histogram roughly follows the symmetric bell shaped curve of the normal distribution.
–The of Q-Q plot does not show any systematic deviations from the line $y = x$.

Each of these three aspects point in the direction of the normal distribution as a correct model for the weekly book sales.

**d.** –Hypotheses are $H_0 \colon F(x) = \Phi\left(\frac{x-\mu}{\sigma}\right)$ versus $H_1 \colon F(x) \neq \Phi\left(\frac{x-\mu}{\sigma}\right)$, with $\alpha = 0.10$.
–Reject $H_0$ if $W \leq c = 0.953$.
–Observed value $w = 0.974$ does not lie in the rejection region. The distribution is not demonstrably non-normal at 10% level of significance.

**2.** Only option (b) is true. Option (d) is not true because we do not know whether the estimators $T_1$ and $T_2$ are unbiased.

**3.** The only correct expression is (c).

**4.**  **a.** We are interested in constructing a confidence interval for a population mean. We have covered a procedure to do this if either: the population is normal, or the sample size is considerably large.

–From the classical numerical summary we notice that the skewness of 1.443 is more than two (even three) times the standard error of 0.419.
–From the histogram we can also see that the data is skewed to the right, and the normal curve does not provide a good fit.

We conclude that the data should not be assumed normal. Since the sample size of $n = 20$ is too small, we cannot invoke the Central Limit Theorem either.

Therefore, there is *no method* covered in this course to construct a confidence interval for the population mean.

**b.** This interpretation is wrong! The *sample* mean is always included in the confidence interval. It is the *population* mean that is trapped only 95 out of 100 times in the long run.

**5.**  **a.** The first observation is that the two samples (scores of Test 1 and Test 2) are not independent! They were measured for the same set of students; meaning each student took both tests. Therefore, we need to pair the samples and focus on the *difference* of the scores.

The next observation is that the normal Q-Q plot of the difference casts a lot of doubt on the assumption that this variable is normally distributed. The sample size $n = 7$ is also too small, so the only alternative we have is to perform a Sign Test on the median.

For a sign test we use as test statistic $X$: "The number of students who scored higher on Test 1". The observed value is thus $x = 6$.
Under the null hypothesis we have $X \sim \text{Binom}(7, 0.5)$. The p-value is given by:

$$\mathbf{P}(X \geq 6 \,|\, H_0 \text{ true}) = 1 - \mathbf{P}(X \leq 5 \,|\, H_0 \text{ true}) = 1 - 0.938 = 0.062.$$

**b.** The interpretation is not correct. The null hypothesis is either true or false (even if we don't know which). Its truth value is fixed, so not a random variable, and we cannot attach a probability to it.

**6.**  **a.** The null hypothesis is that the *population* variances (for the response grade of each interface) are equal. The alternative hypothesis is that these population variances are different.

**b.** In such case we fail to reject the null hypothesis. This does *not* prove the null hypothesis. It simply means that from the data collected there is not enough evidence to support the alternative hypothesis (at some given significance level).

**c.** The only test covered in this course that compares population variances is the $F$-test. An important assumption in this test is that both populations are normally distributed. Looking at the nummerical summaries, we see that skewness and kurtosis are within two standard errors from the reference values. Moreover, the normal Q-Q plots do not show a

severe deviation form the diagonal line. From the information available we see no reason to doubt the normality assumption. Therefore, we may apply here the $F$-test on the equality of variances.

**7.** The statement of the problem explains that two separate surveys were organized: one among lecturers of technical studies, and one among lecturers of social sciences. Since we have two independent populations which we want to compare, the correct test is a test for homogeneity.

(1) We have 2 independent multinomial distributions for the variable "Opinion". We will say that sample 1 is for social sciences, and sample 2 is technical studies. The samples have size 95 and 120, respectively (total is $n = 215$). We will denote by $N_{ij}$ and $p_{ij}$ the observed counts and the (unknown) cell probabilities. Here $i = 1, 2$ denotes the group index and $j = 1, 2, 3$ the opinion index.

(2) Test $H_0 : p_{1j} = p_{2j}$ for all $j = 1, 2, 3$ versus $H_1 : p_{1j} \neq p_{2j}$ for at least one of $j = 1, 2, 3$.

(3) The test statistic is $\chi^2 = \sum_{j=1}^{3} \sum_{i=1}^{2} \frac{(N_{ij} - \widehat{E}_0 N_{ij})^2}{\widehat{E}_0 N_{ij}}$, with $\widehat{E}_0 N_{ij} = \frac{\text{row total} \times \text{column total}}{n}$.

(4) Under $H_0$, $\chi^2$ is $\chi^2$-distributed with $(r-1)(c-1) = 2$ degrees of freedom.

(5) Observed value is $\chi^2_{\text{obs}} = \frac{(40 - 33.1)^2}{33.1} + \ldots + \frac{(25 - 25.1)^2}{25.1} \approx 4.67$.

(6) We reject $H_0$ if $\chi^2 \geq c = 5.99$ (upper tail probability of 5% in $\chi^2_2$ table).

(7) $\chi^2_{\text{obs}} = 4.67 < 5.99$ , so we fail to reject $H_0$.

(8) At a 5% level of significance the survey does not provide sufficient evidence to claim that lecturers in technical and social sciences hold different opinions with respect to online testing.

$$\text{Grade} = 1 + \frac{\text{\# points}}{36} \times 9$$

Rounded to 1 decimal

| question | 1 | 2 | 3 | 4 | 5 | 6 | 7 | Total |
|---|---|---|---|---|---|---|---|---|
| points possible | 11 | 2 | 2 | 5 | 5 | 5 | 6 | 36 |