

EXERCISE SET 6,
MS-A0503, FIRST COURSE IN PROBABILITY AND STATISTICS

EXPLORATIVE EXERCISES

I will expect that you study the explorative problems before the first lecture of the week. It is very strongly recommended that you work on them in groups.

Problem 1. Eight patients have their blood pressure measured before and after trying a new medicine. The results are listed in the table below:

| | | | | | | | | |
|--------|-----|-----|-----|-----|-----|-----|-----|-----|
| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
| Before | 134 | 174 | 118 | 152 | 187 | 136 | 125 | 168 |
| After | 128 | 176 | 110 | 149 | 183 | 136 | 118 | 158 |

We want to test whether the medicine has any effect on the blood pressure.

- (1) Suggest a test statistic and a null hypothesis. (Hint: We are interested in whether the difference $Y_i = X_{i,\text{after}} - X_{i,\text{before}}$ of the blood pressure of the i^{th} patient can be explained by measuring errors or random variations from one day to the other. What distribution would such random fluctuations have?)
- (2) Is the change statistically significant with confidence level 5%? With confidence level 1%?

Problem 2. Recall that the *covariance* of two random variables X and Y is

$$\text{Cov}(X, Y) = E((X - \mu)(Y - \nu)) = E(XY) - \mu\nu.$$

Suggest a way to estimate $\text{Cov}(X, Y)$ from n observations

$$(x_1, y_1), \dots, (x_n, y_n).$$

(Hint: Compute the expected value of $\sum_i (x_i - \bar{x})(y_i - \bar{y})$.)

Problem 3. The following is the height (in inches) of ten pairs of a father and his son.

| | | | | | | | | | | |
|-----------------|------|------|----|------|------|------|------|------|------|----|
| Fathers' height | 60 | 62 | 64 | 65 | 66 | 67 | 68 | 70 | 72 | 74 |
| Sons' height | 63.6 | 65.2 | 66 | 65.5 | 66.9 | 67.1 | 67.4 | 68.3 | 70.1 | 70 |

- (1) Estimate (using your answer in Problem 2) the covariance of the height of a father and his son.
- (2) (Challenging:) Discuss how (and whether) this could be used to test the hypothesis that the height of a father is independent from that of his son.

HOMEWORK PROBLEMS

The homework problems are reported during the second exercise session of the week. You are allowed and encouraged to work in groups, but every student should be prepared to present the solutions individually. During the last exercise session of the week, the teacher will ask you to mark what problems you have solved, and you get points according to how many problems you marked as solved. If you mark a problem as solved, however, you should also be prepared to present your solution in front of the class.

Homework 1. A factory manufactures nails with target length 10 cm. However, the length of manufactured nails varies randomly according to normal distribution. The quality of nails is controlled such that on each full hour, 30 nails are selected randomly and measured.

In a sample the average length of nails is 10.05 and the sample variance is 0.16cm^2 . Test the null hypothesis that the length of these nails is on average 10 cm using 5% significance level, under the alternative hypothesis that the average length differs from the target length.

Homework 2. There are two machines, M_1 and M_2 , which manufacture screws in a screw factory. The thickness of screws manufactured by these machines vary randomly and independently according to a normal distribution. We pick two independent random samples of screws manufactured by each machine and compute the sample variances of their thickness. Data from the samples is shown on the table below. Test the null hypothesis that the machines manufacture equally thick screws on average, when the alternative hypothesis is that there is a difference between the machines. Use 5% significance level.

| Machine | Average (mm) | Sample variance (mm^2) | Sample size |
|---------|--------------|-----------------------------------|-------------|
| M_1 | 9.9 | 0.25 | 31 |
| M_2 | 10.3 | 0.16 | 21 |

Homework 3. In an opinion poll, 3433 random Finns are asked which party they intend to vote for in the next general election. In the September poll, 17.6% claimed that they would vote for Keskusta, whereas only 16.5% claimed they would in the November poll. On confidence level 95%, is it true that the support for Keskusta has fallen?

Week 6, Homework 1

X_1, \dots, X_{30} length of nails.

$H_0: X_1, \dots, X_{30}$ i.i.d. $\mathcal{N}(10, \sigma)$

Test statistic: $T(X) = \frac{\bar{X} - 10}{s/\sqrt{30}} \sim t_{29}$
assuming H_0 .

Plugging in $\bar{X} = 10.05$ $s^2 = 0.16$, we get

$$T(X) = \frac{0.05}{0.4/\sqrt{30}} \approx 0.68.$$

If $T \sim t_{29}$, then

$$P[|T| \geq 0.68] > P[|T| \geq 2.045] = 0.05, \text{ so}$$

we accept the null hypothesis with significance level ~~5%~~ 5%.

Week 6, Homework 2

X_1, \dots, X_{31} screws from M_1

Y_1, \dots, Y_{21} — " — M_2

~~Hypothesis~~ X_1, \dots, X_{31} iid, ~~μ_1, σ_1^2~~ mean μ_1 , var σ_1^2
 Y_1, \dots, Y_{21} iid, mean μ_2 , var σ_2^2 .

$$H_0 : \mu_1 = \mu_2.$$

$\bar{X} - \bar{Y}$ is an unbiased estimator of $\mu_1 - \mu_2$.

Unfortunately, we do not know the exact distribution of $\bar{X} - \bar{Y}$, even after normalizing with S_1 and S_2 .

$$\text{But } \text{Var}(\bar{X} - \bar{Y}) = \text{Var} \bar{X} + \text{Var} \bar{Y} = \frac{\sigma_1^2}{31} + \frac{\sigma_2^2}{21},$$

so by CLT, $\bar{X} - \bar{Y}$ can be approximated

$$\text{by } \mathcal{N}\left(\mu_1 - \mu_2, \frac{\sigma_1^2}{31} + \frac{\sigma_2^2}{21}\right).$$

Assuming H_0 , we would have

$$\frac{\bar{X} - \bar{Y}}{\sqrt{\frac{\sigma_1^2}{31} + \frac{\sigma_2^2}{21}}} \underset{\text{approx}}{\sim} \mathcal{N}(0, 1).$$

Week 6, Homework 2, Continued

(31 & 21 are large enough)

Assuming large sample sizes, we have

$$\frac{\bar{X} - \bar{Y}}{\sqrt{\frac{\sigma_1^2}{31} + \frac{\sigma_2^2}{21}}} \approx \frac{\bar{X} - \bar{Y}}{\sqrt{\frac{s_1^2}{31} + \frac{s_2^2}{21}}} = \frac{9.9 - 10.3}{\sqrt{\frac{0.25}{31} + \frac{0.16}{21}}} \approx 3.19$$

But if $Z \sim \mathcal{N}(0,1)$, then

$$P[|Z| > 3.19] = ~~2~~ 2\Phi(-3.19)$$

≈ 0.0014 , so we reject H_0 .

Week 6, Homework 3

Define: $X_i = \mathbb{I}_{\{i^{\text{th}} \text{ voter in September poll votes Kesk}\}}$

$Y_i = \mathbb{I}_{\{i^{\text{th}} \text{ voter in November poll votes Kesk}\}}$

$$E[X_i] = \mu_{\text{Sept}} \quad E[Y_i] = \mu_{\text{Nov}}$$

$$H_0: \mu_{\text{Sept}} = \mu_{\text{Nov}}, \quad X_i, Y_j \text{ i.i.d.}$$

Assuming H_0 , ^{approximately} an χ^2 -distributed unbiased estimator for $\text{Var}(X_i - Y_j) = \text{Var}(X_i + Y_j)$ is

$$s^2 = \frac{\sum X_i^2 + \sum Y_i^2 - \frac{(\sum X_i + \sum Y_i)^2}{n}}{2 \cdot 3433 - 1} = \frac{3433 \cdot 0.176^2 + 3433 \cdot 0.165^2 - 0.1705^2 \cdot 3433}{2 \cdot 3433 - 1} \approx \frac{1.705 - 0.175^2}{2} \approx 0.1414$$

So we observed $\frac{\bar{X} - \bar{Y}}{s/\sqrt{n}} = \frac{0.176 - 0.165}{\sqrt{0.1414}/\sqrt{3433}} \approx 1.713$

Assuming H_0 , $\frac{\bar{X} - \bar{Y}}{s/\sqrt{n}} \stackrel{\text{approx}}{\sim} \mathcal{N}(0,1)$, and if $Z \sim \mathcal{N}(0,1)$,

then $P[|Z| \geq 1.713] \approx 0.086$, so we cannot reject H_0 .