

ASSIGNMENT SHEET 6

October 25, 2017

Assignment 1. Let X_1, X_2, \dots, X_n be an i.i.d. sample from a continuous distribution with continuous density function f . Let $h > 0$ and define a partition $\{I_j\}_{j \in \mathbb{Z}}$ of \mathbb{R} , where $I_j = [\kappa + (j-1)h, \kappa + jh)$ for a fixed real number κ .

The *histogram* of X_1, X_2, \dots, X_n with *bin-width* h and *origin* κ is defined as the function $x \mapsto \text{Hist}_{X_1, X_2, \dots, X_n}(x)$, where

$$\text{Hist}_{X_1, X_2, \dots, X_n}(x) = \frac{1}{nh} \sum_{i=1}^n \mathbf{1}(X_i \in I_j) \quad \text{if } x \in I_j$$

for each $j \in \mathbb{Z}$.

(a) Show that $\int_{-\infty}^{\infty} \text{Hist}_{X_1, X_2, \dots, X_n}(x) dx = 1$.

(b) Find the distribution of $nh \text{Hist}_{X_1, X_2, \dots, X_n}(x)$ for each x . Hence, find its mean and variance.

(c) What happens to $\mathbb{E}[\text{Hist}_{X_1, X_2, \dots, X_n}(x)]$ when $h \rightarrow 0$.

(Note : This limit indicates what $\text{Hist}_{X_1, X_2, \dots, X_n}(x)$ is estimating for each x for sufficiently small h .)

(d) Using part (b), find $\mathbb{E}\{[\text{Hist}_{X_1, X_2, \dots, X_n}(x) - f(x)]^2\}$.

(e) What happens to the value of the mean squared error in part (d) when $h \rightarrow 0$ and $nh \rightarrow \infty$?

(f) Interpret the limits $h \rightarrow 0$ and $nh \rightarrow \infty$.

(g) Is $\text{Hist}_{X_1, X_2, \dots, X_n}(x)$ a consistent estimator of $f(x)$? Justify your answer.

Assignment 2. (a) Let $X_1, \dots, X_n \sim \text{Exp}(\lambda)$ be independent. We are interested in testing the hypothesis $H_0 : \lambda = 4$ versus $H_1 : \lambda > 4$. Suppose that $n = 16$ and $\sum_{i=1}^{16} X_i = 4$. Find the p -value obtained for this sample for the test constructed last week. *Hint : you cannot do this in closed form ; use R to plot a function and evaluate the p -value graphically from that plot.*

(b* bonus question) Use the R function `uniroot` to obtain the answer numerically (not graphically). *Hint : you will need to write a new function, that is a translation of the function used in (a). For numerical reasons use `interval = c(0.00001, 1)` instead of `c(0, 1)`.*

(c) Do the same for the case where the sum is 2.1.

(d) In which of the two cases would you reject H_0 at significance level $\alpha = 0.05$?

(e) Find a value x for which the p -value when $\sum X_i = x$ is 0.05.

Assignment 3. The p -value, as a function of the sample, is a random variable taking values in $[0, 1]$. We will see the distribution of this random variable using simulations. Let $X_1, \dots, X_n \sim N(\mu, 1)$ be independent, and consider testing the null hypothesis $H_0 : \mu = 0$ versus $H_1 : \mu \neq 0$.

(a) What is the distribution of $\bar{X} = n^{-1} \sum_{i=1}^n X_i$?

(b) Using part (a), find a number v_α such that the test function $\mathbf{1}\{|\bar{X}| > v_\alpha\}$ has significance level α .

(c) Find a formula, as explicit as possible, to the p -value of the test, as a function of the sample X_1, \dots, X_n .

(d) Use the formula in (c) to empirically find the distribution of the p -value under the null : fix n , generate X_1, \dots, X_n , and compute the p -value. Repeat this `REP` times. Store all the p -values in a numerical vector `p` of length `REP`. Use the command `hist(p)` to plot a histogram of `p`. What do you observe? What happens when you change n ? What happens to this distribution under the alternative?

Assignment 4. Let $X_1, \dots, X_n \sim f(x; \theta)$, where f is 1-parameter exponential family. Another idea for building tests for bilateral hypotheses $H_0 : \theta = \theta_0$ versus $H_1 : \theta \neq \theta_0$ is to directly compare θ_0 to the maximum likelihood estimator $\hat{\theta}$, which we know enjoys good properties. We cannot simply use $(\hat{\theta} - \theta_0)^2$, since such measure does not take into account the variability of $\hat{\theta}$. We would like to use a standardised version $(\hat{\theta} - \theta_0)^2 / \text{var} \hat{\theta}$. Alas, the variance will typically depend on the unknown value of θ , and needs to be estimated as well. This results in what is known as a *Wald test*.

- (a) What is the asymptotic variance of $\hat{\theta}$? Call it $v(\theta)$ (recall that it depends on θ !)
 (b) Since $\hat{\theta}$ is close to θ , it makes sense to estimate $v(\theta)$ by $v(\hat{\theta})$. We then obtain the *Wald test statistic*

$$T = \frac{(\hat{\theta} - \theta_0)^2}{v(\hat{\theta})}$$

Write a formula for T by plugging in the function v .

- (c) Assuming that v is continuous and $v(\theta) > 0$, find the asymptotic distribution of T . *Hint* : $v(\hat{\theta})/v(\theta) \rightarrow 1$ in probability; use Slutsky's theorem.
 (d) Let X_1, \dots, X_n be independent $N(0, \sigma^2)$ random variables and let $\sigma_0^2 > 0$. Use an approximate Wald test to test $H_0 : \sigma^2 = \sigma_0^2$ versus $H_1 : \sigma^2 \neq \sigma_0^2$ at significance level α .
 (e) What is the likelihood ratio test in this case? Is it the same?

Assignment 5. (a) Let X_1, X_2, \dots, X_n e an i.i.d. sample from the $N(\theta, 1)$ distribution. Find $100(1 - \alpha)\%$ confidence regions for θ by

- (i) inverting the likelihood ratio test, and
 (ii) inverting the Wald test.

Are the two confidence intervals same?

(b) Let X_1, X_2, \dots, X_n e an i.i.d. sample from the $Ber(p)$ distribution. Find $100(1 - \alpha)\%$ confidence regions for p by

- (i) inverting the asymptotic likelihood ratio test obtained using Wilks' theorem,
 (ii) inverting the Wald test, and
 (iii) inverting the asymptotic test obtained using the approximation $\sqrt{n}(\bar{X} - p) / \sqrt{p(1 - p)} \xrightarrow{d} N(0, 1)$.

(iv) Are these confidence intervals the same? Compare the lengths of these intervals for $n = 10, 25$ and 100 .

Assignment 6. Let X_1, X_2, \dots, X_n be an i.i.d. sample from $N(\mu, \sigma^2)$ distribution, where $\mu \in \mathbb{R}$ and $\sigma \in (0, \infty)$ are unknown parameters. In this problem, we will find different types of confidence regions for (μ, σ^2) .

(a) Find $100(1 - \alpha)\%$ confidence interval for μ (when σ is known) and σ^2 (when μ is unknown) separately using \bar{X} and S^2 . Denote these by $R_{1,\alpha}(\mathbf{X})$ and $R_{2,\alpha}(\mathbf{X})$.

(b) Is the region $R_{1,\alpha}(\mathbf{X}) \times R_{2,\alpha}(\mathbf{X})$ a $100(1 - \alpha)\%$ confidence region for (μ, σ^2) ? Otherwise, find β (depending on α) such that $R_{1,\beta}(\mathbf{X}) \times R_{2,\beta}(\mathbf{X})$ is a $100(1 - \alpha)\%$ confidence region for (μ, σ^2) using the Bonferroni method.

(c) Use the independence of \bar{X} and S^2 to find β such that $R_{1,\beta}(\mathbf{X}) \times R_{2,\beta}(\mathbf{X})$ is a $100(1 - \alpha)\%$ confidence region for (μ, σ^2) . Denote this region by $R_A(\mathbf{X})$.

(Note : $R_A(\mathbf{X})$ is called the Mood exact region.)

(d) Which one of the above confidence regions obtained in (b) and (c) do you think is preferable? Justify your answer.

(e) Write down the likelihood ratio statistic for testing the hypothesis $H_0 : \mu = \mu_0, \sigma^2 = \sigma_0^2$

vs $H_1 : \mu \neq \mu_0, \sigma^2 \neq \sigma_0^2$.

(f) Use Wilks' theorem and the expression of the likelihood ratio statistic to derive an asymptotic $100(1 - \alpha)\%$ confidence region for (μ, σ^2) . Denote this region by $R_B(\mathbf{X})$.

FACT : It can be proved that $\sqrt{n}\{(\bar{X}, S^2)^\top - (\mu, \sigma^2)^\top\}$ converges in distribution to $(Z_1, Z_2)^\top$, where $Z_1 \sim N(0, \sigma^2)$, $Z_2 \sim N(0, 2\sigma^4)$, and they are independent.

(g) Use the above fact to find the asymptotic distribution of $U_n = n(\bar{X} - \mu)^2/\sigma^2 + n(S^2 - \sigma^2)/(2\sigma^4)$.

(h) Use part (g) to find an asymptotic $100(1 - \alpha)\%$ confidence region for (μ, σ^2) . Denote this region by $R_C(\mathbf{X})$.

(i) What is the asymptotic distribution of $V_n = n(\bar{X} - \mu)^2/S^2 + n(S^2 - \sigma^2)/(2S^4)$?

(j) Use part (i) to find an asymptotic $100(1 - \alpha)\%$ confidence region for (μ, σ^2) . Denote this region by $R_D(\mathbf{X})$.

(k) Suppose that $n = 10$, $\bar{x} = 0$, $s^2 = 1$ and $\alpha = 0.05$. Write a code in **R** to understand how the above four 95% confidence regions, namely, $R_A(\mathbf{X})$, $R_B(\mathbf{X})$, $R_C(\mathbf{X})$ and $R_D(\mathbf{X})$ look like. (*Hint : Each confidence region will be a set of the form $\{(\mu, \sigma^2) : H(\mu, \sigma^2) \leq h\}$, where H is a real-valued function and h is a real number. You will get the function H after you simplify and put the values of n , \bar{x} , s^2 and α . To draw this set, you can use the following code)*

```
f <- function(a,b) H(a,b)
mu_vals <- seq(from=-1,to=1,length=100)
sig_vals <- seq(from=0.5,to=1.5,length=100)
z <- outer(mu_vals,sig_vals,f)
contour(mu_vals,sig_vals,z,levels=h,drawlabels=FALSE)
abline(h=1,v=0,col="red")
```

What do you observe? What do you observe if you take $n = 25$ and $n = 100$?

(Note : It can be proved that the likelihood based region $R_B(\mathbf{X})$ asymptotically has the smallest expected area.)

Assignment 7. A car manufacturing company wishes to publish information on the oil usage of a new car model. For 12 cars of this new model they measure the liters necessary to run for a 100 km, with the following results

14.60, 11.21, 11.56, 11.37, 13.68, 15.07, 11.06, 16.58, 13.37, 15.98, 12.07, 13.22.

The empirical mean and variance of the sample are $\bar{x}_{12} = 13.31$ and $s_{12}^2 = \frac{1}{11} \sum_{i=1}^{12} (x_i - \bar{x})^2 = 3.69$, respectively. Suppose the sample is normally distributed. We want to test whether the mean consumption is equal to 12.2 litres vs the alternative hypothesis that the mean is different from 12.2 litres.

(a) Write the model, the null and the alternative hypothesis.

(b) Which test statistics would you use?

(c) Which values of the test statistics would be considered “extremes”?

(d) Use the test statistics to test at a 5% significance level.

(e) Repeat for a 10% significance level. Did you find any difference?

(f) Compute the p-value p_{obs} .

(g) Repeat (d) et (e) using an approach based on the p-value p_{obs} .

Assignment 8. In this exercise we will explore how, in testing many hypotheses simultaneously, compiling a list on tests based on a small p-values cut-off might result in many false positives with high probability.

Mars Reconnaissance Orbiter is a NASA orbiter that aims to prove that water persisted on the surface on Mars for long period of times. Assume Europe sent its own orbiter to check for presence of liquid water on Mars. The orbiter will snap a picture at 100 pre-sampled location. Under the null, each sampled location is assumed to have a probability of 5% to host water. To have more certainty, the orbiter will gravitate around the planet until it has taken 200 pictures of each location.

- (i). Run the following code and comment it.

```
setseed(25102017)
positions <- 100;
trials <- 200;
true.p <- 0.05;
alpha <- 0.01;
nrep <- 1000

p <- matrix(replicate(positions*nrep,prop.test(rbinom(1,trials,true.p),
trials,true.p)$p.value),nrep)

dim(p)
mean(apply(p,1,min)<alpha)
```

- (ii). What happens if $\alpha = 0.05$? Comment the result.
- (iii). Use the function `p.adjust` to adjust the p-values using Bonferroni, Holms and Hochberg's correction.
- (iv). * What happens when you change the numbers? For example, if trials "small"? How could you overcome this?