

OZONMASTERS

STATISTICS

Home work

Kozhemyak Vitaly

November 13, 2020

Contents

Task 1	2
T1 (i)	2
T1 (ii)	3
N1	6
1 Task 2	7
1.1 T2	7
1.2 T3	8
2 Task 3	8
2.1 T4	8
2.2 T5	10
3 Task 4	11
3.1 T6*	11
3.1.1 (I)	11
3.1.2 (ii)	12
Appendix	13
N1	13

Task 1

T1 (i)

Statement 1. Consider X_1, \dots, X_n — i.i.d. random variables. Let $F(x)$ be a CDF of random variable $X_i, i = \{1, \dots, n\}$. Then

$$-\sum_{i=1}^n \log(1 - F(X_i)) \sim \Gamma(n, 1).$$

Proof. The CDF of random variable $1 - F(X_i)$ is

$$\mathbb{P}(1 - F(X_i) \leq x) = \mathbb{P}(1 - x \leq F(X_i)) = 1 - \mathbb{P}(1 - x > F(X_i)) = 1 - F(F^{-1}(1 - x)) = x.$$

Thus,

$$1 - F(X_i) \sim U[0, 1] \Rightarrow -\log(1 - F(X_i)) \sim \text{Exp}(1).$$

Now we are going to use a moment-generating function $M_X(t)$ which uniquely determines a distribution. If we have $Y_1, \dots, Y_n \sim \text{Exp}(1)$ and $S_n = Y_1 + \dots + Y_n$ then

$$\begin{aligned} M_{S_n}(t) &= \mathbb{E}[e^{tS_n}] = \mathbb{E}[e^{t(Y_1 + \dots + Y_n)}] = \{i.i.d.\} = \mathbb{E}[e^{tY_1}] \cdot \dots \cdot \mathbb{E}[e^{tY_n}] = \\ &= M_{Y_1}(t) \cdot \dots \cdot M_{Y_n}(t) = \frac{1}{(1 - t)^n}, \quad t < 1. \end{aligned}$$

This is a moment-generating function of Gamma distribution $\Gamma(n, 1)$. If we put $Y_i = -\log(1 - F(X_i))$ we will get the statement. □

Plotting exact confidence intervals

Consider X_1, \dots, X_n — i.i.d. random variables with Weibull CDF

$$F(x) = 1 - e^{-(x/\lambda)^\tau}, \quad x > 0.$$

Choose the next central statistic $Z = Z(X_1, \dots, X_n; \lambda) = \sum_{i=1}^n \left(\frac{X_i}{\lambda} \right)^\tau$. A central statistics satisfies the next rules:

1. $F_Z(z)$ does not depend on λ . Denote $Z_i = \left(\frac{X_i}{\lambda}\right)^\tau$. Then

$$F_{Z_i}(z) = \mathbb{P}(Z_i \leq z) = \mathbb{P}(X_i \leq \lambda z^{1/\tau}) = 1 - e^{-z} \Rightarrow Z_i \sim \text{Exp}(1).$$

Moreover, the result above tells us $Z = \sum_{i=1}^n Z_i \sim \Gamma(n, 1)$.

2. Z is continuous and strictly monotone in λ .

According to the definition of the exact confidence interval

$$\mathbb{P}(q_{\alpha/2} \leq Z \leq q_{1-\alpha/2}) = 1 - \alpha,$$

$$\mathbb{P}\left(\frac{q_{\alpha/2}}{\sum_{i=1}^n X_i^\tau} \leq \frac{1}{\lambda^\tau} \leq \frac{q_{1-\alpha/2}}{\sum_{i=1}^n X_i^\tau}\right) = 1 - \alpha,$$

$$\mathbb{P}\left(\left(\frac{q_{\alpha/2}}{\sum_{i=1}^n X_i^\tau}\right)^{1/\tau} \geq \lambda \geq \left(\frac{q_{1-\alpha/2}}{\sum_{i=1}^n X_i^\tau}\right)^{1/\tau}\right) = 1 - \alpha$$

where $q_{\alpha/2}$ and $q_{1-\alpha/2}$ are left and right $\frac{\alpha}{2}$ -quantiles of Gamma distribution $\Gamma(n, 1)$ with the significance level α respectively.

T1 (ii)

Consider $X_1, \dots, X_n \sim \text{Exp}(\lambda)$. As we know $\mathbb{E}[X_1] = \frac{1}{\lambda}$, $\text{Var}[X_1] = \frac{1}{\lambda^2}$. According to the CLT we have

$$\lambda\sqrt{n}\left(\bar{X} - \frac{1}{\lambda}\right) \rightarrow \mathcal{N}(0, 1), \quad n \rightarrow +\infty.$$

Consider the next two statistics

$$S_1(X_1, \dots, X_n) = \frac{1}{n} \sum_{i=1}^n X_i,$$

$$S_2(X_1, \dots, X_n) = \sqrt{\frac{1}{2n} \sum_{i=1}^n X_i^2}.$$

Plotting asymptotic confidence intervals using S_1 statistics

Using the S_1 statistics we get

$$\mathbb{P} \left(-z_{1-\alpha/2} \leq \lambda \sqrt{n} \left(\bar{X} - \frac{1}{\lambda} \right) \leq z_{1-\alpha/2} \right) = 1 - \alpha,$$

$$\mathbb{P} \left(\frac{1}{\bar{X}} \left(1 - \frac{z_{1-\alpha/2}}{\sqrt{n}} \right) \leq \lambda \leq \frac{1}{\bar{X}} \left(1 + \frac{z_{1-\alpha/2}}{\sqrt{n}} \right) \right) = 1 - \alpha.$$

The length of the confidence interval is

$$l_1 = \frac{2}{\bar{X}} \frac{z_{1-\alpha/2}}{\sqrt{n}}.$$

Plotting asymptotic confidence intervals using S_2 statistics

Using the next formula $\left(\sigma^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \mu)^2 \right)$ we can derive the next equation

$$\frac{1}{\lambda^2} = \frac{1}{n} \sum_{i=1}^n \left(X_i - \frac{1}{\lambda} \right)^2 \Rightarrow \bar{X} = \frac{\lambda}{2} \frac{1}{n} \sum_{i=1}^n X_i^2.$$

Then

$$\mathbb{P} \left(-z_{1-\alpha/2} \leq \lambda \sqrt{n} \left(\bar{X} - \frac{1}{\lambda} \right) \leq z_{1-\alpha/2} \right) = 1 - \alpha,$$

$$\mathbb{P} \left(\sqrt{\frac{2n}{\sum_{i=1}^n X_i^2} \left(1 - \frac{z_{1-\alpha/2}}{\sqrt{n}} \right)} \leq \lambda \leq \sqrt{\frac{2n}{\sum_{i=1}^n X_i^2} \left(1 + \frac{z_{1-\alpha/2}}{\sqrt{n}} \right)} \right) = 1 - \alpha.$$

The length of the confidence interval is

$$l_2 = \sqrt{\frac{2n}{\sum_{i=1}^n X_i^2}} \cdot \left(\sqrt{\left(1 + \frac{z_{1-\alpha/2}}{\sqrt{n}} \right)} - \sqrt{\left(1 - \frac{z_{1-\alpha/2}}{\sqrt{n}} \right)} \right).$$

Let's take a look at the plot below

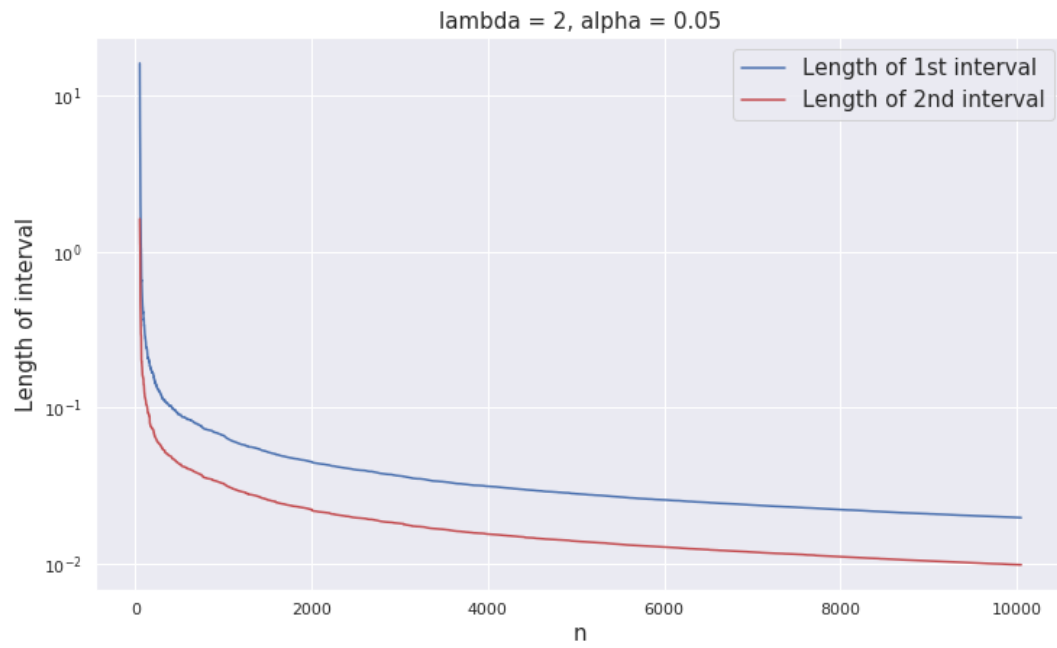


Figure 1: The length of the second interval is less than the length of the first interval ($l_2 < l_1$).

N1

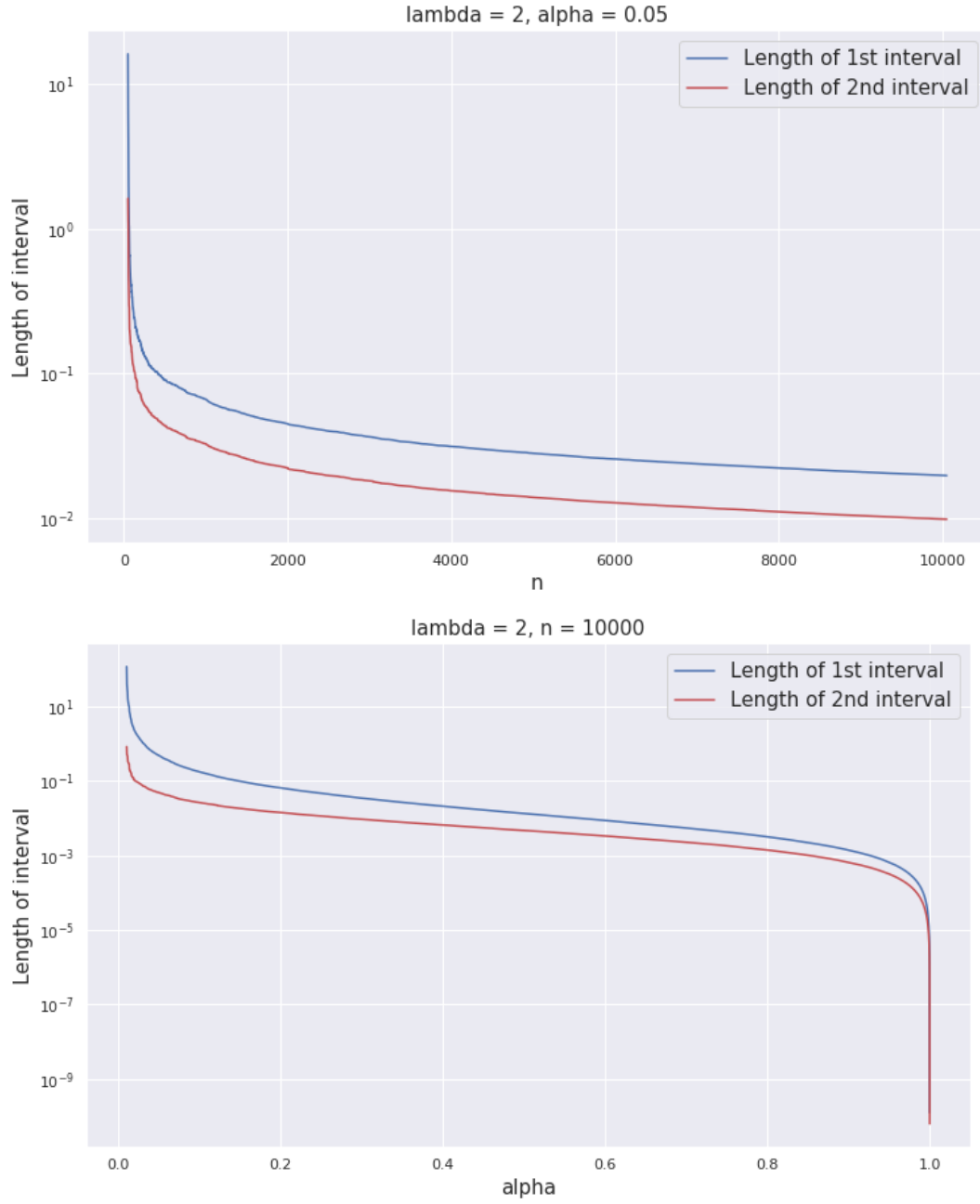


Figure 2: The length of the second interval is less than the length of the first interval ($l_2 < l_1$).

1 Task 2

1 T2

Let X_i, Y_i be random variables such that

$$X_i = \begin{cases} 1, & p, \\ 0, & 1-p, \end{cases} \quad Y_i = \begin{cases} 1, & q, \\ 0, & 1-q, \end{cases}$$

where p is probability to answer "YES" in a city, q is probability to answer "YES" in a village. Using Central Limit Theorem we get

$$\begin{aligned} \frac{\bar{X} - \bar{Y} - (p - q)}{\sqrt{\frac{p(1-p)}{n} + \frac{q(1-q)}{n}}} &= \{\text{Law of Large Numbers}\} = \\ &= \frac{\bar{X} - \bar{Y} - (p - q)}{\sqrt{\frac{\bar{X}(1-\bar{X})}{n} + \frac{\bar{Y}(1-\bar{Y})}{n}}} \rightarrow \mathcal{N}(0, 1), \quad n \rightarrow +\infty. \end{aligned}$$

Denote

$$S_n = \sqrt{\frac{\bar{X}(1-\bar{X})}{n} + \frac{\bar{Y}(1-\bar{Y})}{n}}.$$

The confidence intervals can be calculated as following

$$\begin{aligned} \mathbb{P} \left(-z_{1-\alpha/2} \leq \frac{\bar{X} - \bar{Y} - (p - q)}{S_n} \leq z_{1-\alpha/2} \right) &= 1 - \alpha, \\ \mathbb{P} \left(\bar{X} - \bar{Y} - z_{1-\alpha/2} S_n \leq p - q \leq \bar{X} - \bar{Y} + z_{1-\alpha/2} S_n \right) &= 1 - \alpha. \end{aligned}$$

We can notice that

$$\begin{aligned} S_n &= \sqrt{\frac{\bar{X}(1-\bar{X})}{n} + \frac{\bar{Y}(1-\bar{Y})}{n}} \leq \\ &\leq \left\{ \bar{X}(1-\bar{X}) \leq \frac{1}{4} \right\} \leq \frac{1}{\sqrt{2n}}. \end{aligned}$$

Then

$$\mathbb{P} \left(\bar{X} - \bar{Y} - z_{1-\alpha/2} \frac{1}{\sqrt{2n}} \leq p - q \leq \bar{X} - \bar{Y} + z_{1-\alpha/2} \frac{1}{\sqrt{2n}} \right) = 1 - \alpha.$$

Finally, n_{min} is a solution of the next equation

$$\begin{aligned} \frac{z_{1-\alpha/2}}{\sqrt{2n}} &= 0.05, \\ n_{min} &= \frac{1}{2} \left(\frac{z_{1-\alpha/2}}{0.05} \right)^2, \quad \alpha = 0.05. \end{aligned}$$

1 T3

Consider a null hypothesis $H_0 : p - q = 0$ and an alternative $H_1 : p - q = 0.03$. Let's calculate the Test Statistics

$$T = \frac{\bar{X} - \bar{Y} - (p - q)}{\sqrt{\frac{p(1-p)}{n} + \frac{q(1-q)}{n}}} = \frac{\bar{X} - \bar{Y} - (p - q)}{\sqrt{\frac{\bar{X}(1-\bar{X})}{n} + \frac{\bar{Y}(1-\bar{Y})}{n}}} = \sqrt{n} \frac{0.05 - (p - q)}{\sqrt{0.2023}}.$$

Reject the null hypothesis if

$$T > z_{1-\alpha}$$

where $z_{1-\alpha}$ is a $1 - \alpha$ -quantile of $\mathcal{N}(0, 1)$.

Let's calculate type 1 error as following

$$\mathbb{P}(\text{error 1 type}) = \mathbb{P}(\text{reject } H_0 \mid H_0 \text{ is true}) = \mathbb{P}(T > z_{1-\alpha} \mid p = q).$$

Let's calculate type 2 error as following

$$\mathbb{P}(\text{error 2 type}) = \mathbb{P}(\text{don't reject } H_0 \mid H_0 \text{ is false}) = \mathbb{P}(T \leq z_{1-\alpha} \mid p = q + 0.03).$$

Since we want that

$$\mathbb{P}(\text{error 1 type}) = 0.05,$$

$$\mathbb{P}(\text{error 2 type}) = 0.04,$$

we have to solve the next system of inequalities

$$\begin{cases} \sqrt{n} \frac{0.05}{\sqrt{0.2023}} > z_{0.95} = 1.65, \\ \sqrt{n} \frac{0.02}{\sqrt{0.2023}} \leq z_{0.96} = 1.76, \end{cases} \Leftrightarrow \begin{cases} n > 250.65, \\ n < 1566.61. \end{cases}$$

Finally, we have

$$n_{min} = 251.$$

2 Task 3

2 T4

Consider a null hypothesis $H_0 : \theta = \theta_0$, an alternative $H_1 : \theta > \theta_0$ and LR (Likelihood Ratio) test

$$\lambda(\vec{x}) = \frac{\mathcal{L}(\vec{x}, \theta_0)}{\sup_{\theta \geq \theta_0} \mathcal{L}(\vec{x}, \theta)}.$$

Let $\hat{\theta}$ be the solution for the next problems $\sup_{\theta \geq \theta_0} \mathcal{L}(\vec{x}, \theta)$. Now the LR can be written as following

$$\lambda(\vec{x}) = \left(\frac{\theta_0}{\hat{\theta}}\right)^n \prod_{i=1}^n (1 - x_i)^{\theta_0 - \hat{\theta}}.$$

Thus, we accept H_0 if

$$\lambda(\vec{x}) \geq c \Leftrightarrow n \log \left(\frac{\theta_0}{\hat{\theta}}\right) + (\theta_0 - \hat{\theta}) \sum_{i=1}^n \log(1 - x_i) \geq \log c.$$

Denote

$$\tilde{c} = \frac{\log c - n \log \left(\frac{\theta_0}{\hat{\theta}}\right)}{\theta_0 - \hat{\theta}}.$$

To find the constant \tilde{c} we are going to calculate a type 1 error

$$\mathbb{P}(\text{type 1 error}) = \mathbb{P}(\text{reject } H_0 | H_0 \text{ is true}) = \mathbb{P}\left(\sum_{i=1}^n \log(1 - x_i) < \tilde{c} \mid H_0 \text{ is true}\right).$$

Statement 2. Consider $X_1, \dots, X_n \sim p(x; \theta)$. Then

$$\sum_{i=1}^n \log(1 - X_i) \sim \Gamma\left(n, \frac{1}{\theta}\right).$$

Proof. Consider exponential random variable $X \sim p(x; \theta)$. Then

$$F_{\log(1-X)}(x) = \mathbb{P}(\log(1 - X) \leq x) = 1 - F_X(1 - e^x).$$

Recall that

$$F_X(x) = \begin{cases} 0, & x < 0, \\ 1 - (1 - x)^\theta, & 0 \leq x \leq 1, \\ 1, & x > 1. \end{cases}$$

Thus, we have

$$F_{\log(1-X)}(x) = e^{x\theta} \Leftrightarrow f_{\log(1-X)}(x) = \theta e^{x\theta} \sim \text{Exp}(\theta), \quad x \in (-\infty, 0].$$

As we know

$$\sum_{i=1}^n \log(1 - X_i) \sim \Gamma\left(n, \frac{1}{\theta}\right).$$

□

Back to the type 1 error

$$\mathbb{P}(\text{type 1 error}) = \mathbb{P}\left(\sum_{i=1}^n \log(1 - x_i) < \tilde{c} \mid H_0 \text{ is true}\right) = F_{(n, \theta_0)}(\tilde{c}) = \alpha,$$

$$\tilde{c} = F_{(n, \theta_0)}^{-1}(\alpha) = q_{(n, \theta_0)}(\alpha) \Rightarrow c = \exp\left((\theta_0 - \hat{\theta}) \cdot q_{(n, \theta_0)}(\alpha) + n \log\left(\frac{\theta_0}{\hat{\theta}}\right)\right)$$

where $F_{(n, \theta_0)}$ is a CDF of Gamma distribution $\Gamma\left(n, \frac{1}{\theta_0}\right)$, and $q_{(n, \theta_0)}(\alpha)$ is a α -quantile of Gamma distribution $\Gamma\left(n, \frac{1}{\theta_0}\right)$ (if we consider the left tailed test).

According to the definition of the power of statistic test

$$\begin{aligned} W(\theta) &= \mathbb{P}(\text{reject } H_0 \mid H_1 \text{ is true}) = \\ &= \mathbb{P}\left(\sum_{i=1}^n \log(1 - x_i) < \tilde{c} \mid H_1 \text{ is true}\right) = F_{(n, \theta)}(q_{(n, \theta_0)}(\alpha)), \end{aligned}$$

where $F_{(n, \theta)}$ is a CDF of Gamma distribution $\Gamma\left(n, \frac{1}{\theta}\right)$ and $\theta > \theta_0$. Finally, we have

$$\begin{aligned} W(\theta) &= F_{(n, \theta)}(q_{(n, \theta_0)}(\alpha)) = \int_0^{q_{(n, \theta_0)}(\alpha)} \frac{\theta^n t^{n-1} e^{-t\theta}}{\Gamma(n)} dt = \\ &= \left\{z = \frac{\theta}{\theta_0} t\right\} = \frac{\theta_0}{\theta} \int_0^{\frac{\theta}{\theta_0} \cdot q_{(n, \theta_0)}(\alpha)} \frac{\theta^n \theta_0^{n-1}}{\theta^{n-1}} \frac{z^{n-1} e^{-t\theta_0}}{\Gamma(n)} dz = F_{(n, \theta_0)}\left(\frac{\theta}{\theta_0} \cdot q_{(n, \theta_0)}(\alpha)\right). \end{aligned}$$

2 T5

Note that $\bar{X} \sim \mathcal{N}\left(\mu, \frac{\sigma^2}{n}\right)$. Let

$$E_1(c) = \mathbb{P}(\bar{X} > c \mid \mu = \mu_0) = 1 - \mathbb{P}(\bar{X} \leq c \mid \mu = \mu_0) = 1 - \int_{-\infty}^c \frac{\sqrt{n}}{\sqrt{2\pi}\sigma^2} e^{-\frac{n(x - \mu_0)^2}{2\sigma^2}} dx$$

be a type 1 error and

$$E_2(c) = \mathbb{P}(\bar{X} \leq c | \mu = \mu_1) = \int_{-\infty}^c \frac{\sqrt{n}}{\sqrt{2\pi\sigma^2}} e^{-\frac{n(x - \mu_1)^2}{2\sigma^2}} dx$$

be a type 2 error.

Using the optimality condition we get

$$\begin{aligned} \frac{d(E_1(c) + E_2(c))}{dc} &= 0, \\ \frac{\sqrt{n}}{\sqrt{2\pi\sigma^2}} e^{-\frac{n(c - \mu_1)^2}{2\sigma^2}} &= \frac{\sqrt{n}}{\sqrt{2\pi\sigma^2}} e^{-\frac{n(c - \mu_0)^2}{2\sigma^2}}, \\ c &= \frac{\mu_1 + \mu_0}{2}. \end{aligned}$$

3 Task 4

3 T6*

3.1.1 (I)

Let's denote

$$\hat{\theta} = \operatorname{argmax}_{\theta \geq \theta_0} \mathcal{L}(\vec{x}, \theta) = \operatorname{argmax}_{\theta \geq \theta_0} \prod_{i=1}^n p(x_i, \theta).$$

Definition 1. The Kullback-Leibler divergence between the distribution $p_0(x) = p(x, \theta_0)$ of null hypothesis H_0 , and the distribution $p(x) = p(x, \theta)$ of the alternative H_1 is the expected value of the log-likelihood ratio in favor of the null hypothesis. Then

$$KL(p_0 || p) = \mathbb{E}_{p_0} \left[\log \frac{p(x, \theta_0)}{p(x, \theta)} \right] = \int_{-\infty}^{+\infty} p(x, \theta_0) \log \frac{p(x, \theta_0)}{p(x, \theta)} dx.$$

Let

$$\Lambda(\vec{x}) = \frac{\mathcal{L}(\vec{x}, \theta_0)}{\mathcal{L}(\vec{x}, \hat{\theta})} \tag{1}$$

be a likelihood ratio. Consider the log-likelihood ratio, normalized dividing by n :

$$\Lambda_n(\vec{x}) = \frac{1}{n} \log \frac{\mathcal{L}(\vec{x}, \theta_0)}{\mathcal{L}(\vec{x}, \hat{\theta})} = \frac{1}{n} \sum_{i=1}^n \log \frac{p(x_i, \theta_0)}{p(x_i, \hat{\theta})}.$$

Note X_1, \dots, X_n are i.i.d. and $L_i = \log \frac{p(x_i, \theta_0)}{p(x_i, \hat{\theta})}$ is a random variable. In addition, we know from the strong Law of Large Numbers that

$$\frac{1}{n} \sum_{i=1}^n \log \frac{p(x_i, \theta_0)}{p(x_i, \hat{\theta})} \rightarrow \mathbb{E}[L_1], \quad n \rightarrow +\infty.$$

By definition of expected value

$$\mathbb{E}[L_1] = \int_{-\infty}^{+\infty} q(x) \log \frac{p(x, \theta_0)}{p(x, \hat{\theta})} dx$$

where we can put $q(x) = p(x, \theta_0)$ and then

$$\mathbb{E}[L_1] = \int_{-\infty}^{+\infty} p(x, \theta_0) \log \frac{p(x, \theta_0)}{p(x, \hat{\theta})} dx = KL(p_0 || \hat{p}).$$

Since $KL(p_0 || \hat{p}) \geq 0$ ($KL(p_0 || \hat{p}) = 0 \Leftrightarrow p_0 = \hat{p}$) we have two cases.

- **Case 1.** If $\hat{\theta} > \theta_0$ then $\mathbb{E}[L_1] > 0$.
- **Case 2.** If $\hat{\theta} = \theta_0 \Rightarrow \mathbb{E}[L_1] = 0$.

Finally,

$$\log \Lambda(\vec{x}) = n \mathbb{E}[L_1] = n KL(p_0 || \hat{p}) \geq 0, \quad \hat{\theta} \geq \theta_0.$$

3.1.2 (ii)

Using (1) we can rewrite LR test as following

$$\begin{aligned} \mathcal{L}(\vec{x}, \theta_0) &\geq c_\alpha \mathcal{L}(\vec{x}, \hat{\theta}), \\ \prod_{i=1}^n e^{x_i \theta_0 - d(\theta_0)} &\geq c_\alpha \prod_{i=1}^n e^{x_i \hat{\theta} - d(\hat{\theta})}, \\ \sum_{i=1}^n (x_i \theta_0 - d(\theta_0)) &\geq \log c_\alpha + \sum_{i=1}^n (x_i \hat{\theta} - d(\hat{\theta})), \\ \theta_0 &\geq \hat{\theta} + \frac{\log c_\alpha + nd(\theta_0) - nd(\hat{\theta})}{\sum_{i=1}^n x_i}. \end{aligned}$$

If we denote $\tilde{c}_\alpha = \frac{\log c_\alpha + nd(\theta_0) - nd(\hat{\theta})}{\sum_{i=1}^n x_i}$ then

$$\theta_0 \geq \hat{\theta} + \tilde{c}_\alpha.$$

Appendix

N1

asdf