

"Needle in a Haystack"

What is the limiting power:

$$P_{H_A}(\max y_i > |z(\frac{\alpha}{n})|), n \rightarrow \infty$$

in the case when H_A : one $\mu_i = \mu > 0$

1. Asymptotic full power above threshold:

$$\mu = \mu^{(n)} > (1+\varepsilon)\sqrt{2\log n}$$

assume that $\mu_1 = \mu^{(n)}$

$$\Rightarrow P_{H_A}(\max y_i > |z(\frac{\alpha}{n})|) \geq$$

$$\geq P(y_1 > |z(\frac{\alpha}{n})|) = P(z + \mu^{(n)} > |z(\frac{\alpha}{n})|)$$

$$\stackrel{\sim N(0,1)}{=} \left| \begin{array}{l} |z(\frac{\alpha}{n})| \approx \sqrt{2\log n} \\ \mu^{(n)} - |z(\frac{\alpha}{n})| > \varepsilon\sqrt{2\log n} \end{array} \right| \approx P(z > \frac{\varepsilon\sqrt{2\log n}}{\mu^{(n)}}) \xrightarrow{\mu^{(n)} \rightarrow \infty} 1$$

2. Asymptotic powerlessness below threshold:

$$\mu^{(n)} < (1-\varepsilon)\sqrt{2\log n}$$

assume that $\mu_i = \mu^{(n)}$

$$P_{H_A}(\max y_i > |z(\frac{\alpha}{n})|) \leq \underbrace{P(y_1 > |z(\frac{\alpha}{n})|)}_{\textcircled{1}} + \underbrace{P(\max_{2 \leq i \leq n} y_i > |z(\frac{\alpha}{n})|)}_{\textcircled{2}}$$

$$\textcircled{1}: \underbrace{P(z + \mu^{(n)} > |z(\frac{\alpha}{n})|)}_{\sim N(0,1)} \approx P(z > \varepsilon\sqrt{2\log n}) \xrightarrow{n \rightarrow \infty} 0$$

$$\textcircled{2}: P(\max_{2 \leq i \leq n} y_i > |z(\frac{\alpha}{n})|) = P(\min_{2 \leq i \leq n} p_i \leq \frac{\alpha}{n}) =$$

from lecture 1

$$1 - (1 - \frac{\alpha}{n})^n \rightarrow 1 - e^{-\alpha} \approx \alpha$$

Is there any procedure that can do better?

To answer this question we will analyze the optimal test \Leftrightarrow Neyman-Pearson test;

Theorem (Neyman-Pearson)

y_1, \dots, y_n - iid with pdf $f(y)$ (random sample)

$H_0: f(y) = f_0(y)$ VS $H_1: f(y) = f_1(y)$,
 f_0, f_1 are given functions (H_0, H_1 are simple hypothesis)

\Rightarrow the optimal test (\Leftrightarrow it has the largest power for a given significance level α) is based on the statistic:

$$L(y) = \frac{\prod_{i=1}^n f_1(y_i)}{\prod_{i=1}^n f_0(y_i)} = \frac{f_1(y)}{f_0(y)} \leftarrow \begin{array}{l} \text{likelihood} \\ \text{functions,} \\ y = (y_1, \dots, y_n) \end{array}$$

likelihood ratio.

We reject H_0 for large value of $L(y) \Leftrightarrow$
 $\Leftrightarrow L(y) \geq T_n(\alpha)$, where
 $P_0(L(y) \geq T_n(\alpha)) = \alpha$

We don't have a simple alternative:

$H_0: \mu = (\mu_1, \dots, \mu_n) = 0$ VS $H_1: \text{there exists } i \text{ s.t. } \mu_i \neq 0$

\Rightarrow We restrict H_1 to the "Needle in the haystack" problem:

$\mu^j := (0, \dots, 0, \underset{\substack{\uparrow \\ \text{jth place}}}{\mu}, 0, \dots, 0)$, μ is known,
 (but we don't know where the needle is)

We can model the alternative hypothesis in the following way:

The needle can appear with the same ⁻² probability at each of n places:

1) there exists random variable π that does not depend on $y_1 \dots y_n$ s.t.

$$P(\pi = j) = \frac{1}{n}, \quad j = 1, \dots, n$$

$$2) \mu_1 = \mu^\pi \Leftrightarrow \mu_1 = \mu^j \text{ with probability } \frac{1}{n} \\ j = 1, \dots, n$$

$$H_0: \mu = (0 \dots 0) \quad (\text{vs}) \quad H_1: \mu = \mu^\pi$$

let us calculate likelihood functions:

$$H_0: \mu = (0 \dots 0) \Leftrightarrow y_i \sim N(0, 1) \Leftrightarrow$$

$$f_0(y) = \left(\frac{1}{\sqrt{2\pi}}\right)^n \prod_{i=1}^n e^{-\frac{y_i^2}{2}}$$

$$H_1: \mu = \mu^\pi \Leftrightarrow \mu = \mu^j \text{ with probability } \frac{1}{n} \\ \Leftrightarrow y_i \sim N(0, 1), i \neq j$$

$$y_j \sim N(\mu_j, 1) \quad \text{when } \pi = j$$

$$\Rightarrow f_1(y) = \frac{1}{n} \sum_{j=1}^n \frac{1}{(\sqrt{2\pi})^n} e^{-\frac{(y_j - \mu_j)^2}{2}} \prod_{j \neq i} e^{-\frac{y_i^2}{2}}$$

$$P((y_1, \dots, y_n) \in A) = \left| \text{Law of Total Probability} \right| =$$

$$= \sum_{j=1}^n P((y_1, \dots, y_n) \in A / \pi = j) \cdot P(\pi = j) =$$

$$= \frac{1}{n} \sum_{j=1}^n P(\underbrace{(y_1, \dots, y_j, \dots, y_n)}_{\substack{N(0,1) \quad N(\mu_j,1) \quad N(0,1)}} \in A) =$$

$$= \frac{1}{n} \sum_{j=1}^n \int_A \frac{1}{(\sqrt{2\pi})^n} e^{-\frac{(y_j - \mu_j)^2}{2}} \prod_{j \neq i} e^{-\frac{y_i^2}{2}} dy_1 \dots dy_n$$

independent

Let us calculate Neyman-Pearson statistic

$$L(Y) = \frac{f_1(Y)}{f_0(Y)} = \frac{\frac{1}{n} \sum_{j=1}^n \frac{1}{(\sqrt{2\pi})^n} e^{-\frac{(y_j - \mu)^2}{2}}}{\frac{1}{(\sqrt{2\pi})^n} \prod_{j=1}^n e^{-\frac{y_j^2}{2}}} =$$

$$= \frac{1}{n} \sum_{j=1}^n e^{-\frac{(y_j - \mu)^2}{2} + \frac{y_j^2}{2}} = \frac{1}{n} \sum_{j=1}^n e^{y_j \mu - \frac{\mu^2}{2}}$$

Note that μ is known.

let us prove that this optimal test does not identify $\mu = (1 - \varepsilon) \sqrt{2 \log n}$

Proposition 1: Under H_0 : $L(Y) \xrightarrow{P} 1$
 (\Leftrightarrow) the test does not see the difference between H_0, H_1)

Proof will be further

Proposition 2 (based on Prop 1 \Leftrightarrow we assume that prop 1 is true)

$$\lim_{n \rightarrow \infty} P(\text{Type II Error}) = 1 - \alpha,$$

$$\text{where } P_0(L > T_n(\alpha)) = \alpha$$

$$P(\text{Type II Error}) = P_{H_1}(L \leq T_n(\alpha)) =$$

$$= \int \mathbb{1}_{\{L \leq T_n(\alpha)\}} dP_1 = \left| L = \frac{dP_1}{dP_0} \right| =$$

$$= \int \mathbb{1}_{\{L \leq T_n(\alpha)\}} L dP_0 = \left| L = 1 + (L - 1) \right| =$$

$$= \underbrace{\int \mathbb{1}_{\{L \leq T_n(\alpha)\}} dP_0}_{P_0(L \leq T_n(\alpha)) = 1 - \alpha} + \underbrace{\int \mathbb{1}_{\{L \leq T_n(\alpha)\}} (L - 1) dP_0}_{\xrightarrow{L \rightarrow 1} \text{bounded convergence theorem} \rightarrow 0} \rightarrow 1 - \alpha$$

$\Rightarrow T_n(\alpha)$ is bounded uniformly

Corollary: power of the test = $P_{H_2}(L \geq T_n(\alpha)) = 1 - \alpha$

Proof of Proposition 1.

$$L = \frac{1}{n} \sum_{j=1}^n e^{y_j M - \frac{M^2}{2}}, \quad y_j \sim N(0, 1)$$

$$L \xrightarrow{P} 1, n \rightarrow \infty, \text{ when } \mu = (1-\varepsilon)\sqrt{2\log n}$$

$$L = \frac{1}{n} \sum_{j=1}^n X_j, \text{ where } X_j = \frac{e^{y_j M}}{e^{M^2/2}}$$

$$y_j \sim N(0, 1) \Rightarrow y_j M \sim N(0, \mu^2) \Rightarrow$$

$$\Rightarrow e^{y_j M} \sim \text{lognormal}(0, \mu^2) \Rightarrow E e^{y_j M} = e^{\frac{\mu^2}{2}}$$

$$\Rightarrow E X_j = 1 \quad \xrightarrow{\text{Weak Law of Large Numbers}} \quad L \xrightarrow{P} 1, n \rightarrow \infty ???$$

we cannot use it
because X_j depends on n .

$$L = \frac{1}{n} \sum_{j=1}^n X_j; \quad X_j = e^{y_j M - \frac{M^2}{2}}; \quad \mu = (1-\varepsilon)\sqrt{2\log n}$$

$$T_n = \sqrt{2\log n}$$

$$\tilde{L} = \frac{1}{n} \sum_{j=1}^n X_j \mathbb{1}_{\{y_j \leq T_n\}}$$

$$\text{Properties: } \tilde{L} \leq \sum_{j=1}^n \frac{1}{T_n} \frac{1}{\sqrt{2\pi}} e^{-\frac{T_n^2}{2}} = \frac{1}{\sqrt{2\pi}} \frac{n \cdot \frac{1}{n}}{T_n} \xrightarrow{n \rightarrow \infty} 0$$

$$1) P(L \neq \tilde{L}) \leq P(\max_{1 \leq j \leq n} y_j \geq T_n) \leq \sum_{j=1}^n P(y_j \geq T_n)$$

$$2) E_0 \tilde{L} = P(\varepsilon \sqrt{2\log n}) \quad \varphi \text{ cdf of } N(0,1)$$

$$\Gamma E_0 \tilde{L} = \frac{1}{n} \sum_{j=1}^n E X_j \mathbb{1}_{\{y_j \leq T_n\}} =$$

$$= E X_1 \mathbb{1}_{\{y_1 \leq T_n\}} = \int_{-\infty}^{T_n} e^{y\mu - \frac{\mu^2}{2}} \frac{1}{\sqrt{2\pi}} e^{-\frac{y^2}{2}} dy =$$

$$= \int_{-\infty}^{T_n} \frac{1}{\sqrt{2\pi}} e^{-\frac{(y-\mu)^2}{2}} dy = \Phi(T_n - \mu) = \Phi(\varepsilon \sqrt{2 \log n})$$

$$3) \text{Var}_0(\tilde{L}) = o(1)$$

$$\Gamma \text{Var}_0(\tilde{L}) = \frac{1}{n^2} \text{Var}_0 \left(\sum_{j=1}^n X_j \mathbb{1}_{\{y_j \leq T_n\}} \right) =$$

$$= \frac{1}{n} \text{Var}_0(X_1 \mathbb{1}_{\{y_1 \leq T_n\}}) \stackrel{\text{i.i.d.}}{=} \frac{1}{n} E_0 X_1^2 \mathbb{1}_{\{y_1 \leq T_n\}} =$$

$$= \frac{1}{n} \int_{-\infty}^{T_n} e^{2y\mu - \mu^2} \frac{1}{\sqrt{2\pi}} e^{-\frac{y^2}{2}} dy =$$

$$= \frac{1}{n} e^{\mu^2} \int_{-\infty}^{T_n} \frac{1}{\sqrt{2\pi}} e^{-\left(2\mu^2 - 2y\mu + \frac{y^2}{2}\right)} dy =$$

$$= \frac{1}{n} e^{\mu^2} \int_{-\infty}^{T_n} \frac{1}{\sqrt{2\pi}} e^{-\frac{(y-2\mu)^2}{2}} dy = \frac{1}{n} e^{\mu^2} \Phi(T_n - 2\mu)$$

$$T_n - 2\mu = \sqrt{2 \log n} - 2(1-\varepsilon)\sqrt{2 \log n} = (2\varepsilon - 1)\sqrt{2 \log n}$$

$$a) 2\varepsilon - 1 > 0 \Rightarrow \frac{1}{2} < \varepsilon < 1 \Rightarrow$$

$$\frac{e^{\mu^2}}{n} = \frac{e^{(1-\varepsilon)^2 \cdot 2 \log n}}{n} = n^{\frac{2(1-\varepsilon)^2 - 1}{2}} \xrightarrow{< 0} 0, n \rightarrow \infty$$

$$b) 2\varepsilon - 1 < 0$$

$$\Phi(T_n - 2\mu) \stackrel{\rightarrow -\infty}{\leq} \frac{\varphi(2\mu - T_n)}{2\mu - T_n} = \frac{e^{-\frac{(2\mu - T_n)^2}{2}}}{(2\mu - T_n)\sqrt{2\pi}}$$



$$\Rightarrow \text{Var}_0(\tilde{L}) \leq \frac{1}{n} \frac{e^{\mu^2 - \frac{(2\mu - T_n)^2}{2}}}{(2\mu - T_n)\sqrt{2\pi}} \stackrel{=}{=}$$

$$\mu^2 - \frac{1}{2}(2\mu - T_n)^2 = (1-\varepsilon)^2 \cdot 2 \log n - \frac{1}{2}(1-2\varepsilon)^2 \cdot 2 \log n =$$

$$/ = \log n (2(1-\varepsilon)^2 - (1-2\varepsilon)^2) = \log n (1-2\varepsilon^2) \quad -6-$$

$$\Rightarrow \frac{1}{n} \cdot \frac{1}{(2\mu - \bar{T}_n) \sqrt{2\pi}} e^{\log n - 2 \log n \cdot \varepsilon^2} = \frac{e^{-2 \log n \varepsilon^2}}{(2\mu - \bar{T}_n) \sqrt{2\pi}} \rightarrow 0, n \rightarrow \infty$$

Then from 2) and 3) and Chebyshev's inequality:

$$\tilde{L} = \mathcal{P}(\varepsilon \sqrt{2 \log n}) + O_{P_0}(1)$$

$$\Rightarrow \tilde{L} \rightarrow 1, n \rightarrow \infty$$