Lecture 8 (Continued from Lecture 7)

Goodness of Fit Testing in $L_2$ for the density Model

Let $X_1, X_2, \ldots, X_n$ be iid samples from a density $\eta$ on $([0,1], B[0,1])$. Consider the hypothesis testing problem

$$H_0: \eta \equiv 1 \quad vs \quad H_1 := \{\eta \in \mathcal{P}(\alpha, M, B): \|\eta - 1\|_2 \geqslant r_n\}$$

where $\|\cdot\|_2$ stands for the $L_2$ norm on $[0,1]$. Problems of this kind is often called a goodness of fit problem since we are trying to understand whether the uniform distribution on $[0,1]$ fits the data better compared to set of "rougher" densities.

Let us put ourselves in the context of the general discussion above. Here $Z^{(n)} = (X_1, \ldots, X_n)$, $\Omega_n = [0,1]^n$ and $A_n = B[0,1]^n$. Take $\Theta = \mathcal{P}(\alpha, M, B)$ with $B > 1$, and for any $\theta \in \Theta$ identify $\theta$ with $\eta$. Then $P_{\theta \eta}^{(n)}(A) = \int_A \prod_{i=1}^n \eta(x_i) \, dx_i$ for any $\eta \in \Theta$.

Also consider $\Theta \subseteq L_2[0,1]$ and thus $(\Theta, \|\cdot\|_2)$ is a metric space

The testing problem is then,

$$H_0 = \{\eta_0 \equiv 1\} \quad vs \quad H_1(d, r_n) := \{\eta \in \Theta : \inf_{\theta \eta' \in H_0} d(\eta, \eta')$$

$$= \|\eta - 1\|_2 \geqslant r_n\}$$

Therefore to solve this problem we need to address both (i) (Upper Bound) and (ii) (Lower Bound) of our Definition of minimax optimal $d$-separation.

⑥

(i) Upper Bound

The crux of the argument relies on noting that
$$\eta \in H_0 \iff \|\eta - 1\|_2^2 = 0 \quad \text{and} \quad \eta \in H_1 \iff \eta \in P(\alpha, M, B) \text{ \& } \|\eta - 1\|_2^2 \geqslant r_n^2$$

$\therefore$ If we can estimate $\|\eta - 1\|_2^2$ for $\eta \in P(\alpha, M, B)$, we can hope to tell from it whether the value is "close to 0 or not".

Let $x^*(\eta) = \|\eta - 1\|_2^2 = \|\eta\|_2^2 - 1 = x(\eta) - 1$ where
$$x(\eta) = \int \eta^2(x)\,dx$$

Now note that we have already considered the estimation of $x(\eta)$ with our candidate estimator

$$\hat{x}_n = \frac{1}{n(n-1)} \sum_{i_1 \neq i_2} K_{v_j}(x_{i_1}, x_{i_2})$$

$$= \frac{1}{n(n-1)} \sum_{i_1 \neq i_2} \left[ \varphi_{00}(x_{i_1}) \varphi_{00}(x_{i_2}) + \sum_{l=0}^{j} \sum_{k=0}^{2^l - 1} \psi_{lk}(x_{i_1}) \psi_{lk}(x_{i_2}) \right]$$

$$= \frac{1}{n(n-1)} \sum_{i_1 \neq i_2} \sum_{l=-1}^{j} \sum_k \psi_{lk}(x_{i_1}) \psi_{lk}(x_{i_2}) \qquad \text{where } \psi_{-10} = \varphi_{00}$$

Therefore our estimator of $x^*(\eta)$ is

$$\hat{x}_n = \frac{1}{n(n-1)} \sum_{i_1 \neq i_2} K_{v_j}(x_{i_1}, x_{i_2}) - 1$$

$$= \frac{1}{n(n-1)} \sum_{i_1 \neq i_2} \sum_{l=-1}^{j} \sum_{k=0}^{(2^l - 1) \vee 0} \left( \psi_{lk}(x_{i_1}) - \langle \psi_{lk}, \eta_0 \rangle \right) \left( \psi_{lk}(x_{i_2}) - \langle \psi_{lk}, \eta_0 \rangle \right)$$

where $\eta_0 \equiv 1$ corresponds to the uniform density

$\boxed{7}$

- Since $\hat{x}_n^*$ is an minimax optimal estimator of $x^*(\eta)$ over $\mathcal{P}(\alpha, M, B)$ for properly chosen $j$, we expect $\hat{x}^*$ to be "small" under the null hypothesis and "comparatively larger" under the alternative hypothesis. Therefore we will want to reject $H_0$ for larger values of $\hat{x}_n^*$.

<u>Question</u>: How large a value of $\hat{x}_n^*$ should makes us reject $H_0$?

<u>Idea</u>: The answer to this question is of course guided by the necessity to control the type I error (such that under the alternative one still beats the threshold)

Fix $\alpha > 0$.

$\therefore$ Let our test be $\Psi_n \equiv \mathbb{I}\left(|\hat{x}_n^*| \geq t_n\right)$ and our job is to determine the "smallest $t_n$" s.t $\sup_{\eta \in H_0} P_\eta(\Psi_n = 1) \to 0$ as $n \to \infty$.

By chebychev's inequality, $P_{\eta_0}(\Psi_n = 1) = P_{\eta_0}(|\hat{x}_n^*| \geq t_n) \leq \dfrac{\mathbb{E}_{\eta_0}(\hat{x}_n^{*2})}{t_n^2}$

So, we way to find $t_n$ s.t. $\mathbb{E}_{\eta_0}(\hat{x}_n^{*2})/t_n \leq \alpha/2$ for $0 < \alpha' < 1/2$.

Note that under $\eta_0$ i.e. the null hypothesis

$$\hat{x}_n^* = \frac{1}{n(n-1)} \sum_{i_1 \neq i_2} \sum_{l=-1}^{j} \sum_{k=0}^{2^l \vee 0} (\Psi_{lk}(x_{i_1}) - \langle \Psi_{lk}, \eta_0 \rangle)(\Psi_{lk}(x_{i_2}) - \langle \Psi_{lk}, \eta_0 \rangle)$$

is degenerate.

therefore by a result in Lecture 4,

$$\mathbb{E}_{\eta_0}(\hat{x}_n^{*2}) = \frac{2}{n(n-1)} \mathbb{E}_{\eta_0}\left\{ \sum_{\ell=-1}^{\dot{j}} \sum_{k=0}^{\overset{\circ}{2}-1\vee 0} (\psi_{\ell k}(x_{i_1}) - \langle \psi_{\ell k}, \eta_0 \rangle) \right.$$
$$\left. (\psi_{\ell k}(x_{i_2}) - \langle \psi_{\ell k}, \eta_0 \rangle) \right\}^2 \quad (*)$$

To analyze the above, note that, if $e_1, e_2, \dots$ is
an o.n.b. of $L_2[0,1]$ then (using degeneracy)

$$\mathbb{E}_{\eta_0}\left( \sum_{\ell=1}^{M} (e_\ell(x_1) - \langle e_\ell, \eta_0 \rangle)(e_\ell(x_2) - \langle e_\ell, \eta_0 \rangle) \right)^2$$

$$\leq \mathbb{E}_{\eta_0}\left( \sum_{\ell=1}^{M} e_\ell(x_1) e_\ell(x_2) \right)^2 \qquad \text{(check this \=)}$$

Therefore, from (*) above,

$$\mathbb{E}_{\eta_0}(\hat{x}_n^{*}) \overset{\sim}{=} \frac{2}{n(n-1)} \mathbb{E}_{\eta_0}\left\{ \sum_{\ell=-1}^{\dot{j}} \sum_{k=0}^{\overset{\circ}{2}-1\vee 0} (\psi_{\ell k}(x_1) - \langle \psi_{\ell k}, \eta_0 \rangle) \right.$$
$$\left. (\psi_{\ell k}(x_2) - \langle \psi_{\ell k}, \eta_0 \rangle) \right\}^2$$

$$\leq \frac{2}{n(n-1)} \mathbb{E}_{\eta_0}\left[ \sum_{\ell=-1}^{\dot{j}} \sum_{k=0}^{\overset{\circ}{2}-1\vee 0} \psi_{\ell k}(x_1) \psi_{\ell k}(x_2) \right]^2$$

$$= \frac{2}{n(n-1)} \sum_{\ell=-1,\,k=0}^{\dot{j},\,\overset{\circ}{2}-1} \left\{ \mathbb{E}_{\eta_0}\left( \psi_{\ell k}^2(x) \right) \right\}^2 + \frac{2}{n(n-1)} \sum_{\ell \neq \ell'} \sum_{k \neq k'} \left( \mathbb{E}_{\eta_0} \psi_{\ell k}(x) \psi_{\ell' k'}(x) \right)^2$$

$$\leq \frac{2}{n(n-1)} \iint \left( \sum_{\ell=-1}^{\dot{j}} \sum_{k=0}^{\overset{\circ}{2}-1\vee 0} \psi_{\ell k}(x_1) \psi_{\ell k}(x_2) \right)^2 \eta_0(x_1)\, \eta_0(x_2)\, dx_1 dx_2$$

$$\leq \frac{2 \|\eta_0\|_\infty^2}{n(n-1)} \iint \left( \sum_{\ell=-1}^{\dot{j}} \sum_{k=0}^{\overset{\circ}{2}-1\vee 0} \psi_{\ell k}(x_1) \psi_{\ell k}(x_2) \right)^2 dx_1 dx_2 \quad \text{(if } \eta_0 \in \gamma(\alpha, M, B)$$
$$\text{Here } \eta_0 \leq 1$$
$$\text{so can use}$$
$$B = 1 \text{)}$$

$$= \frac{2 \|\eta_0\|_\infty^2}{n(n-1)} \times 2^{\dot{j}} = \frac{2^{\dot{j}+1}}{n(n-1)} \|\eta_0\|_\infty^2$$

⑨

Therefore we have,

$$\mathbb{P}_{\eta_0}\left(|\hat{x}_n^*| \geq t_n\right) \leq \mathbb{E}_{\eta_0}\left(\hat{x}_\eta^{*2}\right)/t_n^2$$

$$= \frac{2^{j+1}\,\|\eta_0\|_\infty^2}{n(n-1)\,t_n^2} \qquad (**)$$

$\Rightarrow$ Letting $t_n = \sqrt{\dfrac{2^{j+2}\,\|\eta_0\|_\infty^2}{n(n-1)\,\alpha'}}$ we get

$$P_{\eta_0}\left(|\hat{x}_n^*| \geq t_n\right) \leq \alpha'/2$$

$\therefore$ We reject when $|\hat{x}_n^*| \geq \sqrt{\dfrac{2^{j+2}\,\|\eta_0\|_\infty^2}{n(n-1)\,\alpha'}}$

$H_0 : \{\eta_0\}$

We now need to show that if $r_n^2 \geq C(\alpha,M,B)\, n^{-\frac{4\alpha}{4\alpha+1}}$

then,

$$\sup_{\substack{\eta \in P(\alpha,M,B),\\ \|\eta-\eta_0\|_2^2 \geq r_n^2}} \mathbb{P}_\eta\left(\psi_n = 0\right) \leq \alpha'/2$$

To do this we need to analyze $\hat{x}_n^*$ under each $\eta \in P(\alpha,M,B)$, $\|\eta-\eta_0\|_2^2 \geq r_n^2$ separately.

Fix $\eta \in \mathscr{P}(\alpha, M, B)$ such that $\|\eta - \eta_0\|_2^2 \geq c(M, B, \alpha) n^{-\frac{4\alpha}{4\alpha+1}}$

$$\mathbb{P}_\eta(\psi_n = 0) = \mathbb{P}_\eta(|\hat{\chi}_n^*| \leq t_n)$$

$$= \mathbb{P}_\eta(|\hat{\chi}_n^* - \mathbb{E}_\eta \hat{\chi}_n^* + \mathbb{E}_\eta \hat{\chi}_n^*| \leq t_n)$$

$$\leq \mathbb{P}_\eta(|\hat{\chi}_n^* - \mathbb{E}_\eta \hat{\chi}_n^*| \geq -t_n + |\mathbb{E}_\eta \hat{\chi}_n^*|) \qquad (*)$$

Now, $t_n = \sqrt{\dfrac{2^{j+2} \|\eta_0\|_\infty^2}{n(n-1)\alpha}}$

$$|\mathbb{E}_\eta \hat{\chi}_n^*| = \mathbb{E}_\eta\left(\frac{1}{n(n-1)} \sum_{i_1 \neq i_2} K_{V_j}(x_{i_1}, x_{i_2}) - 1\right)$$

$$= \mathbb{E}_\eta\left(\frac{1}{n(n-1)} \sum_{i_1 \neq i_2} \sum_{l=-1}^{j} \sum_{k=0}^{2^l-1\vee 0} (\psi_{lk}(x_{i_1}) - \langle \psi_{lk}, \eta_0\rangle)(\psi_{lk}(x_{i_2}) - \langle \psi_{lk}, \eta_0\rangle)\right)$$

$$= \sum_{l=-1}^{j} \sum_{k=0}^{2^l-1\vee 0} \left\{\mathbb{E}_\eta(\psi_{lk}(x) - \langle \psi_{lk}, \eta_0\rangle)\right\}^2$$

$$= \sum_{l=-1}^{j} \sum_{k=0}^{2^l-1\vee 0} \left\{\langle \psi_{lk}, \eta - \eta_0\rangle\right\}^2 = \|\Pi_{V_j}(\eta - \eta_0)\|_2^2$$

$$= \|\eta - \eta_0\|_2^2 - \|\Pi_{V_j^\perp}(\eta - \eta_0)\|_2^2$$

$$\geq c(M, B, \alpha) n^{-\frac{4\alpha}{4\alpha+1}} - (2M) 2^{-2j\alpha} \qquad \text{as } \eta - \eta_0 \in \mathscr{P}(\alpha, 2M, 2B)$$

$\Rightarrow$ If $2^j \approx n^{\frac{2}{4\alpha+1}}$ i.e. $j \sim \frac{2}{4\alpha+1} \log_2 n$ we have that for sufficiently large $c(M, B, \alpha)$

$\mathbb{E}_\eta|\hat{\chi}_n^*| \geq t_n$. Indeed by taking $c(M, B, \alpha)$ large enough, $\mathbb{E}_\eta|\hat{\chi}_n^*| \geq \dfrac{c(M, B, \alpha)}{3} n^{-\frac{4\alpha}{4\alpha+1}} = \dfrac{c(M, B, \alpha)}{3} 2^{-2j\alpha}$

$$= \dfrac{c(M, B, \alpha)}{3} \sqrt{\dfrac{2^j}{n^2}}$$

Note that the crux of the above choice of $j$ is driven by the fact $2^{-2j\alpha} = \sqrt{\frac{2^j}{n^2}}$ is solved by $2^j \sim n^{\frac{2}{4\alpha+1}}$. This is exactly similar to the bias variance tradeoff while estimating a quadratic functional. The difference is that the $\sqrt{1/n}$ term does not arise in the std. deviation due to first order degeneracy of $\widehat{x}_n^*$.

Collecting the above calculations in (*) implies,

$$\mathbb{P}_\eta(\psi_n = 0) \leq \mathbb{P}_\eta\left(|\widehat{x}_n^* - \mathbb{E}_\eta \widehat{x}_n^*| \geq |\mathbb{E}_\eta \widehat{x}_n^*| - t_n\right)$$

$$= \mathbb{P}_\eta\left(|\widehat{x}_n^* - \mathbb{E}_\eta \widehat{x}_n^*| \geq \underbrace{\|\Pi_{V_j}(\eta - \eta_0)\|_2^2}_{\geq \frac{C(M,B,\alpha)}{3} n^{-\frac{4\alpha}{4\alpha+1}}} - t_n\right)$$

Now, in order to apply chebychev's inequality, we need to understand the variance of $\widehat{x}_n^*$, which can be done by a Hoeffding's Decomposition of $\widehat{x}_n^* - \mathbb{E}_\eta \widehat{x}_n^*$ under $\mathbb{P}_\eta$.

Check: (See chapter 11 & 12, And van Der Vaart : Asymptotic Statistics for reference, on Hoeffding's Decomposition )

$$\widehat{x}_n^* - \mathbb{E}_\eta \widehat{x}_n^* = \frac{2}{n} \sum_{j=1}^n \sum_{\ell=0}^j \sum_k \left(\psi_{\ell k}(x_i) - \langle \psi_{\ell k}, \eta \rangle\right)\langle \psi_{\ell k}, \eta - \eta_0 \rangle$$

$$+ \frac{1}{n(n-1)} \sum_{i_1 \neq i_2} \sum_{\ell=-1}^j \sum_{k=0}^{2^\ell - 1 \vee 0} \left(\psi_{\ell k}(x_i) - \langle \psi_{\ell k}, \eta \rangle\right)\left(\psi_{\ell k}(x_j) - \langle \psi_{\ell k}, \eta \rangle\right)$$

$$= T_1(\eta) + T_2(\eta) \qquad \textcircled{12}$$

$$\Rightarrow \mathbb{P}(\psi_n = 0) \leq \mathbb{P}_\eta\left(|T_1(\eta) + T_2(\eta)| \geq \|\Pi_{V_j^\circ}(\eta - \eta_0)\|_2^2 - t_n\right)$$

$$\bullet \leq \mathbb{P}_\eta\left(|T_1(\eta)| \geq \|\Pi_{V_j^\circ}(\eta - \eta_0)\|_2^2/4\right)$$

$$+ \mathbb{P}_\eta\left(|T_2(\eta)| \geq \frac{3}{4}\|\Pi_{V_j}(\eta - \eta_0)\|_2^2 - t_n\right)$$

$$= I_1 + I_2$$

Controlling $I_1$: $\mathbb{P}_\eta\left(|T_1(\eta)| \geq \|\Pi_{V_j^\circ}(\eta - \eta_0)\|_2^2/4\right)$

$$\leq \frac{\text{Var}_\eta T_1^\bullet(\eta)}{\|\Pi_{V_j^\circ}(\eta - \eta_0)\|_2^4 / 16}$$





$$\text{Var}_\eta(T_1^\bullet(\eta)) = \text{Var}_\eta\left(\frac{2}{n}\sum_{i=1}^{n}\sum_{\ell,k}\left(\psi_{\ell k}^{(x_i)} \mathbb{E}_\eta(\psi_{\ell k})\right)\right)\langle\psi_{\ell k}, \eta - \eta_0\rangle$$

$$= \text{Var}_\eta\left(\frac{2}{n}\sum_{i=1}^{n}\left\{\left(\sum_{\ell,k}\psi_{\ell k}^{(x_i)}\langle\psi_{\ell k}, \eta - \eta_0\rangle\right) - \left(\mathbb{E}_\eta\sum_{\ell,k}\psi_{\ell k}\langle\psi_{\ell k}, \eta - \eta_0\rangle\right)\right\}\right)$$

$$= \frac{4}{n}\text{Var}_\eta\left(\sum_{\ell,k}\psi_{\ell k}(x_i)\langle\psi_{\ell k}, \eta - \eta_0\rangle - \mathbb{E}_\eta\sum_{\ell,k}\psi_{\ell k}^{(x_i)}\langle\psi_{\ell k}, \eta - \eta_0\rangle\right)$$

$$\leq \frac{4}{n}\mathbb{E}_\eta\left(\sum_{\ell,k}\psi_{\ell k}(x_i)\langle\psi_{\ell k}, \eta - \eta_0\rangle\right)^2$$

$$= \frac{4}{n}\int\left(\sum_{\ell,k}\psi_{\ell k}(x)\langle\psi_{\ell k}, \eta - \eta_0\rangle\right)^2\eta(x)\,dx$$

⑬

$$\leq \frac{4 \|\eta\|_\infty^2}{n} \int \left( \sum_{\ell,k} \psi_{\ell k}(x) \langle \psi_{\ell k}, \eta - \eta_0 \rangle \right)^2 dx$$

$$= \frac{4 \|\eta \cdot\|_\infty^2}{n} \sum_{\ell,k} \langle \psi_{\ell k}, \eta - \eta_0 \rangle^2 = \frac{4 \|\eta\|_\infty^2}{n} \| \Pi_{v_j}(\eta - \eta_0) \|_2^2$$

$$\Rightarrow \mathbb{P}_\eta \left( |T_1(\eta)| \geq \| \Pi_{v_j}(\eta - \eta_0) \|_2^2 / 4 \right) \leq \; 216 \|\eta\|_\infty^2 / n \| \Pi_{v_j}(\eta - \eta_0) \|_2^2$$

$$\leq \frac{216 \|\eta\|_2^2}{n^{\frac{1}{4\alpha+1}}}$$

<u>Controlling $I_2$</u>:  First note that

$$\frac{3}{4} \| \Pi_{v_j}(\eta - \eta_0) \|_2^2 - t_n$$

$$= \quad \geq \frac{3}{4} \left( \| \eta - \eta_0 \|_2^2 - (2M) 2^{-2j\alpha} \right) - \sqrt{\frac{2^{j+2} \|\eta_0\|_\infty^2}{n(n-1)\alpha'}}$$

$$\geq \frac{3}{4} \left( c(M,B,\alpha) n^{-\frac{4\alpha}{4\alpha+1}} - (2M) 2^{-2j\alpha} \right) - \sqrt{\frac{2^{j+2} B^2}{n(n-1)\alpha'}}$$

$$\geq \sqrt{\frac{2^{j+3} B^2}{n(n-1)\alpha'}} \qquad \text{if } 2^j \sim n^{\frac{2}{4\alpha+1}} \text{ and } \quad c(M,B,\alpha) \text{ is}$$
$$\text{chosen large enough.}$$

Therefore
$$\mathbb{P}_\eta \left( |T_2(\eta)| \geq \frac{3}{4} \| \Pi_{v_j}(\eta - \eta_0) \|_2^2 - t_n \right)$$

$$\leq \mathbb{P}_\eta \left( \left| \frac{1}{n(n-1)} \sum_{i_1 \neq i_2} \sum_{\ell,k} (\psi_{\ell k}(x_{i_1}) - \langle \psi_{\ell k}, \eta \rangle)(\psi_{\ell k}(x_{i_2}) - \langle \psi_{\ell k}, \eta \rangle) \right| \right.$$

$$\left. \geq \sqrt{\frac{2^{j+3} B^2}{n(n-1)\alpha}} \right)$$

which by calculations similar to (**) above is
$$\leq \alpha/4$$

Therefore we have proved that ~~for~~, for $C(M,B,\alpha)$ large enough, if $r_n^2 = C(M,B,\alpha) \, n^{-\frac{4\alpha}{4\alpha+1}}$, then

$$\sup_{\eta \in H_\perp(d, r_n)} \mathbb{P}_\eta(\psi_n = 0) \le \alpha/2$$

This combined with our bound on type I error proves the upper bound of the theorem.

## (ii) Lower Bound

The proof of the lower bound is ~~so~~ exactly the same as the proof of lower bound for estimation of quadratic functional.

Following our strategy, take $\theta_0 \in H_0$ as $\theta_0 = \eta_0$ and $\oplus_M \subseteq H_\perp(d, \rho_n)$ with $\rho_n \ll r_n = n^{-4\alpha/4\alpha+1}$ as

$\bullet \; \eta_\lambda(x) = \eta_0(x) + c \, 2^{-j\alpha} \dfrac{2^{-j}}{a_n} \sum_k \lambda_k \, \psi_{jk}(x)$

$\lambda \in \{-1,1\}^{2^j}$, $M = 2^{2^j}$

~~where~~ $2^j \sim n^{2/4\alpha+1}$

The rest is an exactly similar second moment argument.
(Do it by yourself to be sure).

and $a_n \to \infty$ is s.t. $\dfrac{2^{-2j\alpha}}{a_n} = n^{-4\alpha/4\alpha+1}$

(this can be done since $\rho_n \ll n^{-4\alpha/4\alpha+1}$ and therefore these $\eta_\lambda$'s $\in H_\perp(d, \rho_n)$)