

Density Estimation:

①

R.V. X has density f on \mathbb{R}^d when

$$P\{X \in A\} = \int_A f(x) dx \quad \text{for all } A \in \mathcal{B}(\mathbb{R}^d), \text{ Borel sets in } \mathbb{R}^d$$

[ie, $P\{X \in A\} = f(x) \cdot \text{Vol}(A)$ if $A = \{x' \in \mathbb{R}^d : d(x', x) \leq r\}$ for small $r > 0$.

Purpose

Estimate unknown density f from an i.i.d. sample X_1, X_2, \dots, X_n drawn from f .

density estimate

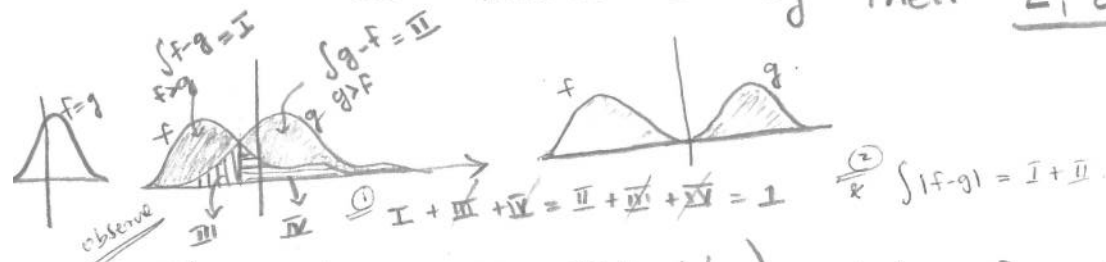
$$f_n(x) = f_n(x; X_1, \dots, X_n) : (\mathbb{R}^d)^{n+1} \rightarrow \mathbb{R}$$

quality of f_n measured by T.V. distance.

$$\text{error} = \sup_{B \in \mathcal{B}(\mathbb{R}^d)} \left| \int_B f_n - \int_B f \right|$$

if this is $< \epsilon$ then all probabilities will be estimated with errors not exceeding ϵ .

The distance between two densities f & g can be measured by their L_1 distance $\int |f-g|$



Thm 1 (Scheffé's Identity). Let f and g be two functions defined on \mathbb{R}^d satisfying $\int f = \int g = 1$.

$$\sup_{B \in \mathcal{B}(\mathbb{R}^d)} \left| \int_B f - \int_B g \right| = \frac{1}{2} \int |f-g|$$

proof

$$\sup_{B \in \mathcal{B}(\mathbb{R}^d)} \left| \int_B f - \int_B g \right| = \int_{f>g} (f-g) = \int_{g>f} (g-f) = \frac{1}{2} \int |f-g|$$

□

why L_1 (TV) ① absolute scale ② physical interpretation

$$0 \leq \frac{1}{2} \int |f-g| \leq 1, \quad \int |f-g| < 0.02$$

\Rightarrow differences in probabilities are at most $\frac{1}{2} \times 0.02 = 0.01$.

③ scale invariance

$$\sup_B |P\{X \in B\} - P\{Y \in B\}| = \sup_B |P\{T(X) \in T(B)\} - P\{T(Y) \in T(B)\}|$$

if T is a bijection and $\{T(B) : B \in \mathcal{B}(\mathbb{R}^d)\} = \mathcal{B}(T(\mathbb{R}^d))$.

④ TV distance decreases on any Borel measurable mapping

$T: \mathbb{R}^d \rightarrow \mathcal{S} \subseteq \mathbb{R}^k$, i.e. for any R.V.s X and Y on $\mathcal{B}(\mathbb{R}^d)$

$$\sup_{A \in \mathcal{B}(\mathbb{R}^d)} |P\{X \in A\} - P\{Y \in A\}| \geq \sup_{A \in \mathcal{B}(\mathbb{R}^k)} |P\{T(X) \in A\} - P\{T(Y) \in A\}|$$

skip

proof:

$$\begin{aligned} & \sup_{A \in \mathcal{B}(\mathbb{R}^k)} |P\{T(X) \in A\} - P\{T(Y) \in A\}| \\ &= \sup_{A \in \mathcal{B}(\mathbb{R}^k)} |P\{X \in T^{-1}(A)\} - P\{Y \in T^{-1}(A)\}| \\ &\leq \sup_{A \in \mathcal{B}(\mathbb{R}^d)} |P\{X \in A\} - P\{Y \in A\}| \end{aligned}$$

since $\{T^{-1}(A) : A \in \mathcal{B}(\mathbb{R}^k)\} \subseteq \mathcal{B}(\mathbb{R}^d)$. □

Minimum Distance Estimate

(3)

setting:

$$X_1, \dots, X_n \stackrel{iid}{\sim} f$$

objective:

non parametric density estimate f_n from X_1, \dots, X_n

with

universal performance guarantees ($\int |f_n - f|$ is small for any $f \in \mathcal{L}_1$ as n gets large).

consider the simpler problem of choosing between two densities f_n and g_n , i.e. construct ϕ_n such that,

$$\int |\phi_n - f| \approx \min(\int |f_n - f|, \int |g_n - f|)$$

Idea

use empirical measure $\mu_n(A) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{[X_i \in A]} = \frac{\text{\# of data point in } A}{n}$

* Scheffé set of ordered pair (F_n, g_n) is:

$$A = A(F_n, g_n) = \{x : F_n(x) > g_n(x)\}$$

observe:

IF $\int f_n = \int g_n = 1$ then by Scheffé's identity,

$$\int |f_n - g_n| = 2 \int_{F_n > g_n} (F_n - g_n) = 2 \int_A f_n - 2 \int_A g_n$$

to get

Scheffé Estimate

$$f_n^* = \begin{cases} f_n & \text{if } \left| \int_A f_n - \mu_n(A) \right| < \left| \int_A g_n - \mu_n(A) \right| \\ g_n & \text{otherwise} \end{cases}$$



eg 1 given data $\text{iii} \dots \text{i}$
 IF $\mu_n(A) = 8/10$, $\int_A f_n = 7/10$, $\int_A g_n = 4/10$
 then $f_n^* = f_n$.

eg 2 := different data $\text{xx} \dots \text{x}$ from another f
 IF $\mu_n(A) = 5/10$
 then $f_n^* = g_n$

Thm (DL6-1) Let f_n and g_n be two density estimates (4)
 with $\int f_n = \int g_n = 1$. For the Scheffé estimate f_n^* , we have

$$\int |f_n^* - f| \leq 3 \min \left(\int |f_n - f|, \int |g_n - f| \right) + 4 \max_{A \in \mathcal{A}} \left| \int_A f - \mu_n(A) \right|,$$

where $\mathcal{A} = \{ \{f_n > g_n\}, \{g_n > f_n\} \}$

Proof: (let's prove a slightly more general thm)

Note: f_n^* has an error that is within $E_n = 4 \sup_{A \in \mathcal{A}} \left| \int_A f - \mu_n(A) \right|$
 of 3 times the best possible error among f_n and g_n .

Now consider the problem of selecting from k densities. (our main problem here)

f_{ni} , $1 \leq i \leq k$, $\int f_{ni} = 1$ for all i

let A_{ij} denote the Scheffé set $A_{ij} = A(f_{ni}, f_{nj}) = \{x : f_{ni}(x) > f_{nj}(x)\}$.

Let the Yatacos class of such Scheffé sets be:

$$\mathcal{A} = \{ A_{ij}, A_{ji} : 1 \leq i < j \leq k \}$$

The Minimum Distance Estimate (MDE) Ψ_n is the f_{ni} of smallest index that minimizes

$$\Delta_i = \sup_{A \in \mathcal{A}} \left| \int_A f_{ni} - \mu_n(A) \right|$$

$$\Psi_n = f_{ni^*}, \quad i^* = \min \{ \operatorname{argmin}_{1 \leq i \leq k} \Delta_i \}$$

Thm 3 (Universal Performance Bound of MDE) [DL 6.3]
 from a finite set of densities.

(5)

For the MDE ψ_n , we have (for each n).

$$\int |\psi_n - f| \leq 3 \min_i \int |f_{n_i} - f| + 4\Delta, \quad \Delta = \sup_{A \in \mathcal{A}} \left| \int f - M_n(A) \right|$$

Proof:

Let $\psi_n = f_{n_i}$ and let $f_{n_j} = \operatorname{argmin}_{1 \leq l \leq k} \int |f_{n_l} - f|$.

Assume $j \neq i$, then

$$(*) \quad \int |\psi_n - f| = \int |f_{n_i} - f| \leq \int |f_{n_j} - f| + \int |f_{n_i} - f_{n_j}| \quad (\Delta_i \neq \Delta_j \text{ of } \mathcal{L}_1)$$

Now, assuming wlog $i < j$:

$$\int |f_{n_i} - f_{n_j}| = 2 \sup_{A \in \{A_{ij}, A_{ji}\}} \left| \int_A f_{n_i} - \int_A f_{n_j} \right| \quad (\text{by Scheffé's identity, Thm 1})$$

$$\leq 2 \sup_{A \in \mathcal{A}} \left| \int_A f_{n_i} - \int_A f_{n_j} \right| \quad (\{A_{ij}, A_{ji}\} \subseteq \mathcal{A})$$

$$\leq 2 \sup_{A \in \mathcal{A}} \left| \int_A f_{n_i} - M_n(A) \right| + 2 \sup_{A \in \mathcal{A}} \left| \int_A f_{n_j} - M_n(A) \right|$$

$$= 2 \sup_{A \in \mathcal{A}} \left| \int_A \psi_n - M_n(A) \right| + 2 \sup_{A \in \mathcal{A}} \left| \int_A f_{n_j} - M_n(A) \right| \quad (\psi_n = f_{n_i} \text{ by assumption})$$

$$\leq 4 \sup_{A \in \mathcal{A}} \left| \int_A f_{n_j} - M_n(A) \right| \quad (\Delta_i \leq \Delta_j)$$

$$\leq 4 \sup_{A \in \mathcal{A}} \left| \int_A f_{n_j} - \int_A f \right| + 4 \sup_{A \in \mathcal{A}} \left| \int_A f - M_n(A) \right| = 4\Delta$$

$$\leq 4 \sup_{B \in \mathcal{B}(\mathbb{R}^d)} \left| \int_B f_{n_j} - \int_B f \right| + 4\Delta \quad (\mathcal{A} \subseteq \mathcal{B}(\mathbb{R}^d))$$

$$(**) \quad = 2 \int |f_{n_j} - f| + 4\Delta$$

$$*)(***) \quad \int |\psi_n - f| \leq 3 \int |f_{n_i} - f| + 4\Delta = 3 \min_i \int |f_{n_i} - f| + 4\Delta \quad \square$$

has to obtain error bounds for MDE we need

(6)

bounds on $\Delta = \sup_{A \in \mathcal{A}} |\mu(A) - M_n(A)|$

no matter what $f \in L_1$ is generating data X_1, \dots, X_n .

$$\mu(A) = \int_A f$$

We are interested in Δ , the maximal deviation of M_n from μ over \mathcal{A} :

$$g(X_1, \dots, X_n) = \Delta = \sup_{A \in \mathcal{A}} |M_n(A) - \mu(A)|$$

We will first show (using concentration \neq s) that

$$P\{|g - \mathbb{E}g| \geq \varepsilon\} \leq 2e^{-2n\varepsilon^2}$$

and thus the maximal deviation is sharply concentrated around its mean and then we will show (using uniform deviation \neq) that this mean ($\mathbb{E}g = \mathbb{E}\Delta$) can be bounded using combinatorial characteristics of \mathcal{A} .

Using concentration #s i.e., #s of the form:

(7)

$$P\{|X_n - EX_n| \geq \varepsilon\} \leq c_1 e^{-nc_2 \varepsilon^2}$$

Recall Markov's # For any non-negative RV X and any

$t > 0$:

$$P\{X \geq t\} \leq \frac{EX}{t}$$

proof:

$$EX = \int_0^{\infty} x dF(x) = \int_0^t x dF(x) + \int_t^{\infty} x dF(x) \geq \int_t^{\infty} x dF(x) \geq t \int_t^{\infty} dF(x) = t P\{X \geq t\}$$

Let $s \in (0, \infty)$, then for any RV X and any $t > 0$, by

Markov's #, $P\{X \geq t\} = P\{e^{sX} \geq e^{st}\} \leq \frac{E\{e^{sX}\}}{e^{st}}$

In Chernoff's bounding method we find an $s > 0$ that minimizes $E\{e^{sX}\}/e^{st}$, the upper bound.

IF $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} X_1$ and $S_n = \sum_{i=1}^n X_i$, we get

$$P\{S_n - ES_n \geq t\} \leq e^{-st} E\left\{\exp\left(s \sum_{i=1}^n (X_i - EX_i)\right)\right\} \\ = e^{-st} \prod_{i=1}^n E\{e^{s(X_i - EX_i)}\} \quad (\text{by independent-ence})$$

So finding tight bounds \Leftrightarrow good upper bounds for M.G.F. of $X_i - EX_i$.

Lemma 4 (upper bounding the MGF of a bounded RV.)

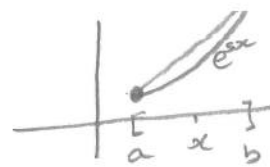
Let X be RV with $EX=0$, $a \leq X \leq b$.

Then for $s > 0$:

$$E\{e^{sX}\} \leq e^{s^2(b-a)^2/8}$$

proof: by convexity

$$e^{sx} \leq \frac{x-a}{b-a} e^{sb} + \frac{b-x}{b-a} e^{sa}$$



for $a \leq x \leq b$.

(8)
 $\frac{1}{2}$
 $\frac{1}{3}$
 $\frac{1}{4}$
 $\frac{1}{5}$
 $\frac{1}{6}$
 $\frac{1}{7}$
 $\frac{1}{8}$
 $\frac{1}{9}$
 $\frac{1}{10}$
 $\frac{1}{11}$
 $\frac{1}{12}$
 $\frac{1}{13}$
 $\frac{1}{14}$
 $\frac{1}{15}$
 $\frac{1}{16}$
 $\frac{1}{17}$
 $\frac{1}{18}$
 $\frac{1}{19}$
 $\frac{1}{20}$

Since $E(X) = 0$,

$$\begin{aligned} E\{e^{sX}\} &\leq E\left\{\frac{X-a}{b-a} e^{sb} + \frac{b-X}{b-a} e^{sa}\right\} = \underbrace{\left(\frac{-a}{b-a}\right)}_p e^{sb} + \underbrace{\left(\frac{b}{b-a}\right)}_{1-p} e^{sa} \\ &= p e^{sb} + (1-p) e^{sa} = \left(\frac{p e^{sb}}{e^{sa}} + \frac{(1-p) e^{sa}}{e^{sa}}\right) e^{sa} \\ &= (1-p + p e^{s(b-a)}) e^{sa} = (1-p + p e^{s(b-a)}) e^{-ps(b-a)} \\ &\stackrel{\text{def}}{=} e^{\phi(u)} \end{aligned}$$

where $u = s(b-a)$ and $\phi(u) = -pu + \log(1-p + p e^u)$

now, $\phi'(u) = \frac{d}{du} \phi(u) = -p + \frac{p e^u}{(1-p + p e^u)} = -p + \frac{p}{p + (1-p) e^{-u}}$

Thus, $\phi(0) = \phi'(0) = 0$

and $\phi''(u) = \frac{p(1-p)e^{-u}}{(p + (1-p)e^{-u})^2} \leq \frac{1}{4}$, since $0 < p < 1$ and $e^{-u} < 1$

Then, by Taylor's theorem, for some $\xi \in [0, u]$

$$\phi(u) = \phi(0) + u\phi'(0) + \frac{u^2}{2} \phi''(\xi) \leq \frac{u^2}{2} = \frac{s^2(b-a)^2}{2} \quad \square$$

Thm 5 (Hoeffding's #). Let X_1, \dots, X_n be indep. RVs such that $X_i \in [a_i, b_i]$ w.p. 1. Then for any $t > 0$,

$$P\{S_n - ES_n \geq t\} \leq e^{-2t^2 / \sum_{i=1}^n (b_i - a_i)^2}$$

and $P\{S_n - ES_n \leq -t\} \leq e^{-2t^2 / \sum_{i=1}^n (b_i - a_i)^2}$

proof:

(9)

From Chernov's bounding method, for any $s > 0, t > 0$

$$P\{S_n - ES_n \geq t\} \leq e^{-st} \prod_{i=1}^n E\{e^{s(X_i - EX_i)}\}$$

$$\leq e^{-st} \prod_{i=1}^n e^{s^2(b_i - a_i)^2/8}$$

by Lemma 4.
since $a_i \leq X_i \leq b_i$
and $E(X_i - EX_i) = 0$

$$= e^{-st} e^{s^2 \sum_{i=1}^n (b_i - a_i)^2/8}$$

$$= e^{-2t^2 / \sum_{i=1}^n (b_i - a_i)^2}$$

by choosing $s = \arg \min_s \left(-st + \frac{s^2 \sum_{i=1}^n (b_i - a_i)^2}{8} \right)$
 $= 4t / \left(\sum_{i=1}^n (b_i - a_i)^2 \right)$

||| by for $P\{S_n - ES_n \leq -t\} = P\{ES_n - S_n \geq t\}$

Example $X_1, \dots, X_n \stackrel{iid}{\sim} \text{Bernoulli}(p) \Rightarrow P\{S_n/n - p \geq \varepsilon\} \leq e^{-2n\varepsilon^2}$

Lemma 6 (Expected Maximal Deviation \neq)

Let $\sigma > 0, n \geq 2$, and let Y_1, \dots, Y_n be real-valued RVs such that for all $s > 0$ and $1 \leq i \leq n$, $E\{e^{sY_i}\} \leq e^{s^2\sigma^2/2}$

Then

$$E\left\{ \max_{i \leq n} Y_i \right\} \leq \sigma \sqrt{2 \ln n}$$

If, additionally, $E\{e^{s(-Y_i)}\} \leq e^{s^2\sigma^2/2}$ for every $s > 0$

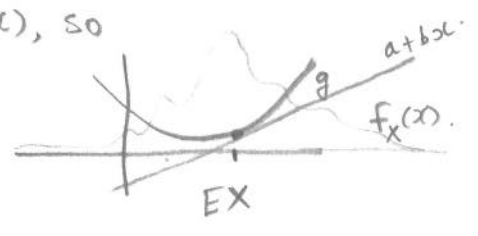
and $1 \leq i \leq n$, then for any $n \geq 1$,

$$E\left\{ \max_{i \leq n} |Y_i| \right\} \leq \sigma \sqrt{2 \ln(2n)}$$

Proof: First let's recall Jensen's \neq $\begin{cases} g \text{ is convex} \Rightarrow E g(X) \geq g(EX) \\ g \text{ is concave} \Rightarrow E g(X) \leq g(EX) \end{cases}$

Let $L(x) = a + bx$ be a line tangent to $g(x)$ at the point EX . Since g is convex it lies above $L(x)$, so

$$\begin{aligned}
 E g(X) &\geq E L(X) = E(a + bX) = a + bEX \\
 &= L(EX) \\
 &= g(EX)
 \end{aligned}$$



Proof

By Jensen's \neq and convexity of $\exp(\cdot)$, for all $s > 0$,

$$\begin{aligned}
 e^{s E \{ \max_{i \leq n} Y_i \}} &\leq E \{ e^{s \max_{i \leq n} Y_i} \} \\
 &= E \{ \max_{i \leq n} e^{s Y_i} \} \\
 &\leq \sum_{i=1}^n E \{ e^{s Y_i} \} \\
 &\leq n e^{s^2 \sigma^2 / 2}
 \end{aligned}$$

Thus, $E \{ \max_{i \leq n} Y_i \} \leq \frac{\ln n}{s} + \frac{s \sigma^2}{2}$

for any $s > 0$, and taking

$$s = \operatorname{argmin}_s \frac{\ln n}{s} + \frac{s \sigma^2}{2} = \sqrt{2 \ln n / \sigma^2}$$

gives

$$\begin{aligned}
 E \{ \max_{i \leq n} Y_i \} &\leq \frac{\ln n}{\sqrt{2 \ln n} / \sigma} + \frac{\sqrt{2 \ln n} \cdot \sigma^2}{2} \\
 &= \frac{2(\ln n) \sigma + (\sqrt{2 \ln n})^2 \sigma}{2 \sqrt{2 \ln n}} = \frac{2(2 \ln n) \sigma}{2 \sqrt{2 \ln n}} = \sigma \sqrt{2 \ln n}
 \end{aligned}$$

$$\begin{aligned}
 \operatorname{argmin}_s \frac{\ln n}{s} + \frac{s \sigma^2}{2} &= \sqrt{2 \ln n / \sigma^2} \\
 \frac{d}{ds} ((\ln n) s^{-1} + \frac{\sigma^2}{2} s) &= 0 \\
 -\ln n s^{-2} + \frac{\sigma^2}{2} &= 0 \\
 s^2 &= \sigma^2 / 2 \ln n \\
 s &= \sqrt{2 \ln n / \sigma^2} \\
 \text{and } \frac{d^2}{ds^2} &= \ln n s^{-3} > 0
 \end{aligned}$$

$$E \{ \max_{i \leq n} Y_i \} \leq \sigma \sqrt{2 \ln n} \quad (\star)$$

Finally, $\max_{i \leq n} |Y_i| = \max(Y_1, -Y_1, Y_2, -Y_2, \dots, Y_n, -Y_n)$

and applying (\star) to $2n$ terms we get

$$E \{ \max_{i \leq n} |Y_i| \} \leq \sigma \sqrt{2 \ln(2n)}$$



Lemma 7: [a conditional straightforward extension of Lemma 4]

(11)

Let V and Z be RVs such that $E\{V|Z\} = 0$ with prob. 1, and for some function h and constant $c \geq 0$:

$$h(Z) \leq V \leq h(Z) + c.$$

Then for all $s > 0$:

$$E\{e^{sV} | Z\} \leq e^{s^2 c^2 / 8}$$

Now, we get to an extension of Hoeffding's \neq .

Thm 8 (Bounded Difference \neq)

Let A be some set and suppose the function $g: A^n \rightarrow \mathbb{R}$ satisfies the bounded difference assumption

$$\sup_{\substack{x_1, \dots, x_n, \\ x_i' \in A}} |g(x_1, \dots, x_n) - g(x_1, \dots, x_{i-1}, x_i', x_{i+1}, \dots, x_n)| \leq c_i, \\ 1 \leq i \leq n.$$

(ie., changing the i^{th} variable of g while fixing all others does not change the value of g by more than c_i)

and let X_1, \dots, X_n be independent RVs. taking values in A .

Then, for all $t > 0$,

$$P\{g(X_1, \dots, X_n) - E g(X_1, \dots, X_n) \geq t\} \leq e^{-2t^2 / \sum_{i=1}^n c_i^2}$$

and $P\{E g(X_1, \dots, X_n) - g(X_1, \dots, X_n) \geq t\} \leq e^{-2t^2 / \sum_{i=1}^n c_i^2}$

proof [by martingale technique of McDiarmid (1989)]

proof

(12)

$$\text{Let } V = g(X_1, \dots, X_n) - \mathbb{E}g(X_1, \dots, X_n) = g - \mathbb{E}g$$

$$\text{let } V_i = \mathbb{E}\{g | X_1, \dots, X_i\} - \mathbb{E}\{g | X_1, \dots, X_{i-1}\}, \quad i=1, \dots, n$$

$$\text{Then } V = \sum_{i=1}^n V_i.$$

$$\text{Let } H_i(X_1, \dots, X_i) = \mathbb{E}\{g(X_1, \dots, X_n) | X_1, \dots, X_i\}$$

and F_i be the distribution of X_i for $i=1, \dots, n$

Then,

$$V_i = H_i(X_1, \dots, X_i) - \int H_i(X_1, \dots, X_{i-1}, x) F_i(dx)$$

$$\text{Let } W_i = \sup_u \left(H_i(X_1, \dots, X_{i-1}, u) - \int H_i(X_1, \dots, X_{i-1}, x) F_i(dx) \right)$$

$$\text{and } Z_i = \inf_v \left(H_i(X_1, \dots, X_{i-1}, v) - \int H_i(X_1, \dots, X_{i-1}, x) F_i(dx) \right)$$

then $Z_i \leq V_i \leq W_i$ for $i=1, \dots, n$ with prob 1

and by bounded difference assumption:

$$W_i - Z_i = \sup_u \sup_v \left(H_i(X_1, \dots, X_{i-1}, u) - H_i(X_1, \dots, X_{i-1}, v) \right) \leq c_i$$

Therefore by Lemma 7, for $i=1, \dots, n$:

$$\mathbb{E}\{e^{sV_i} | X_1, \dots, X_{i-1}\} \leq e^{s^2 c_i^2 / 8}$$

Using the fact that if X, Y are arbitrary bounded RVs

$$\text{then } \mathbb{E}\{XY\} = \mathbb{E}\{\mathbb{E}\{XY | Y\}\} = \mathbb{E}\{Y \mathbb{E}\{X | Y\}\}$$

and Chernov's bounding method; for any $s > 0$:

$$P\{g - \mathbb{E}g > t\} = P\{V > t\} = P\left\{\sum_{i=1}^n V_i > t\right\}$$

$$\leq \frac{\mathbb{E}\{e^{s \sum_{i=1}^n V_i}\}}{e^{st}} = \mathbb{E}\left\{e^{s \sum_{i=1}^{n-1} V_i} e^{sV_n}\right\}$$

$$= \frac{\mathbb{E}\left\{e^{s \sum_{i=1}^{n-1} V_i} \mathbb{E}\{e^{sV_n} | X_1, \dots, X_{n-1}\}\right\}}{e^{st}}$$

$$\leq e^{s^2 c_n / 8} \frac{E \left\{ e^{s \sum_{i=1}^{n+1} v_i} \right\}}{e^{st}}$$

∴ (by repeating n times)

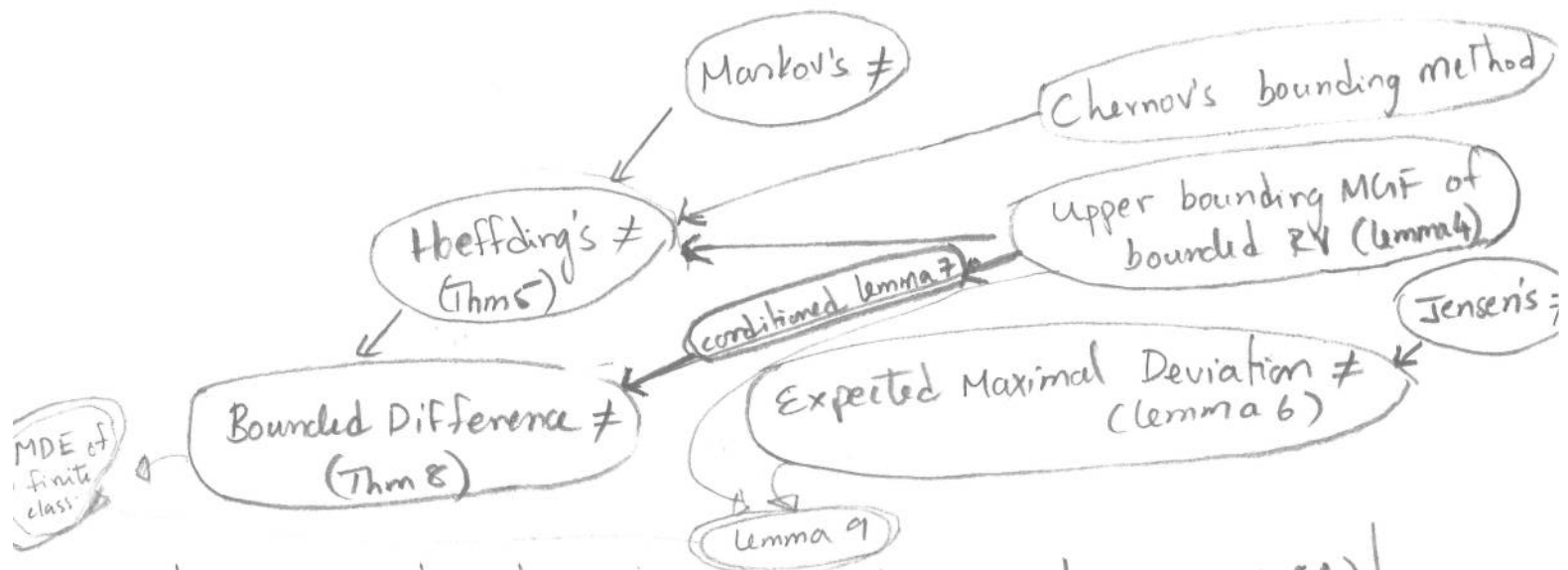
$$\leq e^{-st} e^{s^2 \sum_{i=1}^n c_i^2 / 8}$$

Now choosing the $s = 4t / \sum_{i=1}^n c_i^2$ gives

$$P\{g - Eg \geq t\} \leq e^{-\frac{4t^2}{\sum_{i=1}^n c_i^2} + \frac{2}{8} \frac{16t^2}{\left(\sum_{i=1}^n c_i^2\right)^2} \sum_{i=1}^n c_i^2} = e^{-\frac{2t^2}{\sum_{i=1}^n c_i^2}}$$

$$P\{Eg - g \geq t\} \leq e^{-\frac{2t^2}{\sum_{i=1}^n c_i^2}} \text{ is proved similarly } \blacksquare$$

Recall our main reason for tour of concentration ≠ s



is to upper bound $\Delta = g(X_1, \dots, X_n) = \sup_{A \in \mathcal{A}} |\mu(A) - M_n(A)|$

where X_1, \dots, X_n are iid RVs taking values in \mathcal{X} ,
 \mathcal{A} is a collection of subsets of \mathcal{X} with $\mu(A) = P\{X_i \in A\}$
 and $M_n(A) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{\{X_i \in A\}}$, for any $A \in \mathcal{A}$.

Regardless of the nature of \mathcal{A} , by changing one X_i ,
 g can change by at most $1/n$.
 Therefore by bounded difference ≠ :

$$P \left\{ \left| \sup_{A \in \mathcal{A}} |\mu(A) - \mu_n(A)| - E \left\{ \sup_{A \in \mathcal{A}} |\mu(A) - \mu_n(A)| \right\} \right| > t \right\}$$

$$\leq 2e^{-2t^2 / \sum_{i=1}^n (\frac{1}{n})^2} = 2e^{-2t^2 / \frac{n}{n^2}} = 2e^{-2nt^2}$$

for any n and $t > 0$.

This shows that for any class \mathcal{A} , the maximal deviation of empirical measure μ_n from true measure μ is sharply concentrated around its expected value. ($\Delta = g \rightarrow 0$ in probability if and only if $g \rightarrow 0$ almost surely).

Thus, we only worry about bounding the expected value of the maximal deviation, i.e.

$$E \left\{ \sup_{A \in \mathcal{A}} |\mu(A) - \mu_n(A)| \right\}$$

To do this ^{in generality} we need uniform deviation $\neq s$

but when $|\mathcal{A}| < \infty$ as in our Thm 2's MDE setting

of selecting from k densities f_{n_i} , $1 \leq i \leq k$, $\int f_{n_i} = 1$

for all i , with Yatracos class $\mathcal{A} = \{A_{ij}, A_{ji} : 1 \leq i < j \leq k\}$ made up of

Scheffé sets $A_{ij} = \{x : f_{n_i}(x) > f_{n_j}(x)\}$ then all

we need are lemmas 4 and 6.

Lemma 9 If $\mathcal{A} = \{A_{ij} = \{x : f_{n_i}(x) > f_{n_j}(x)\} ; i, j \in \{1, \dots, k\}, i \neq j\}$ with $k < \infty$

$$\text{Then } E \Delta = E \left\{ \sup_{A \in \mathcal{A}} |\mu(A) - \mu_n(A)| \right\} \leq \sqrt{\frac{\ln(2k^2)}{2n}}$$

proof

Note that $\mathcal{A} = \{A_{ij}, A_{ji} : 1 \leq i < j \leq k\}$ has at most k^2 sets.

let $X_A = \mu(A) - \mu_n(A)$ for each $A \in \mathcal{A}$. For $s > 0$,

$$\begin{aligned}
E\{e^{sX_A}\} &= E\left\{e^{s(\mu(A) - \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{\{X_i \in A\}})}\right\} \\
&= E\left\{e^{s\left(\frac{1}{n} \sum_{i=1}^n (\mu(A) - \mathbb{1}_{\{X_i \in A\}})\right)}\right\} \\
&= \prod_{i=1}^n E\left\{e^{s(\mu(A) - \mathbb{1}_{\{X_i \in A\}})/n}\right\} \quad \text{by independence of } X_1, \dots, X_n. \\
&\leq \prod_{i=1}^n e^{s^2(1/n)^2/8} \\
&= e^{s^2/8 \sum_{i=1}^n 1/n^2} \\
&= e^{s^2(\frac{1}{2\sqrt{n}})^2/2}
\end{aligned}$$

by lemma 4 since $E\left\{\frac{\mu(A) - \mathbb{1}_{\{X_i \in A\}}}{n}\right\} = 0$
and $\mu(A) \leq \mu(A) - \mathbb{1}_{\{X_i \in A\}} \leq \mu(A) + \frac{1}{n}$

Note that we also have $E\{e^{-sX_A}\} \leq e^{s^2(\frac{1}{2\sqrt{n}})^2/2}$.
 Thus for each $A \in \mathcal{A}$, $X_A = \mu(A) - \mu_n(A)$, $E\{e^{sX_A}\} \leq e^{s^2\sigma^2/2}$,
 where $\sigma = \frac{1}{2\sqrt{n}}$. Therefore by lemma 6 (expected maximal deviation \neq)

$$\begin{aligned}
E\left\{\max_{A \in \mathcal{A}} |\mu(A) - \mu_n(A)|\right\} &= E\left\{\max_{A \in \mathcal{A}} X_A\right\} \\
&\leq \frac{1}{2\sqrt{n}} \sqrt{2 \ln(2|\mathcal{A}|)} \leq \frac{1}{2\sqrt{n}} \sqrt{2 \ln(2k^2)} \\
&= \sqrt{\frac{\ln(2k^2)}{2n}} \quad \square
\end{aligned}$$

Finally from Lemma 9 and Thm 2 taking Expectations in we get our full asymptotic story for MDE for the problem of selecting from k densities.

Corollary 10

$$E(|\Psi_n - F|) \leq 3 \min_{i \leq b} \int |F_{ni} - F| + \sqrt{\frac{8 \ln(2k^2)}{n}}$$

Corollary 10 has applications in situations where k is allowed to grow with sample size n , say as k_n , such that $\min_{1 \leq i \leq k_n} |f_{ni} - f| \xrightarrow[n \rightarrow \infty]{} 0$ as $k_n \rightarrow \infty$

and $\sqrt{\frac{8 \ln(2k_n^2)}{n}} \rightarrow 0$, for e.g. $\ln(k_n) \leq \sqrt{n} \iff k_n \leq \exp(n^{1/20})$

concretely,

n	$k_n = \lfloor \exp(n^{1/20}) \rfloor$	$\sqrt{\frac{8 \ln(2k_n^2)}{n}}$	$k_n = \lfloor \exp(n^{1/20}) \rfloor$	$\sqrt{\frac{8 \ln(2k_n^2)}{n}}$
100	150	0.926	23	0.750
1,000	7,4621	0.430	276	0.310
10,000	$\approx 8.1 \times 10^{10}$	0.202	22,026	0.129
100,000	$\approx 2.6 \times 10^{24}$	0.095	$\approx 5.2 \times 10^7$	0.054
1,000,000	$\approx 4.7 \times 10^{54}$	0.045	$\approx 5.4 \times 10^{13}$	0.023

How to control first error term $\min_{1 \leq i \leq k} |f_{ni} - f|$?

Situation 1 *

IF $f \in \{f_{ni} : 1 \leq i \leq k\}$ then $\min_{1 \leq i \leq k} |f_{ni} - f| = 0$
 (this is useful in some simulation settings only)

Situation 2 *

Suppose $f \in \tilde{\mathcal{F}}$, a class of totally bounded densities i.e., $\forall \epsilon > 0, \exists G_\epsilon = \{g_1, g_2, \dots, g_{N_\epsilon}\} \subseteq \tilde{\mathcal{F}}$, s.t.

$$\tilde{\mathcal{F}} \subseteq \bigcup_{i=1}^{N_\epsilon} B_{g_i, \epsilon}, \quad B_{g, r} = \{f : \int |f - g| \leq r\}$$

The smallest such N_ϵ is the complexity or covering number of $\tilde{\mathcal{F}}$ \circ $\lfloor \log_2 N_\epsilon \rfloor$ is the Kolmogorov entropy of $\tilde{\mathcal{F}}$ and G_ϵ is the skeleton of $\tilde{\mathcal{F}}$.

Data Compression
 # bits to store an integer between 1, ..., N_ϵ

Define The Yatracos class.

$$A_\epsilon = \{ \{x: g_i(x) > g_j(x)\} ; g_i, g_j \in G_\epsilon \}$$

with $|A_\epsilon| \leq N_\epsilon^2$

Now, The MDE of Yatracos (skeleton estimate) is:

$$\Psi_n = \operatorname{argmin}_{g_i \in G_\epsilon} \sup_{A \in A_\epsilon} \left| \int_A g_i - M_n(A) \right|$$

Thus, by corollary 10 and Dfn of N_ϵ we have:

$$E \left(\int |\Psi_n - f| \right) \leq 3 \min_{g_i \in G_\epsilon} \int |g_i - f| + \sqrt{\frac{8 \ln(2N_\epsilon^2)}{n}} = 3\epsilon + \sqrt{\frac{8 \ln(2N_\epsilon^2)}{n}} \quad \left[\begin{array}{l} \epsilon > 0 \text{ s.t.} \\ N_\epsilon \geq 2 \end{array} \right]$$

Remark. Regardless of how quickly $N_\epsilon \rightarrow \infty$ as $\epsilon \rightarrow 0$, we can choose an $\epsilon = \epsilon_n = \inf \{ \epsilon : \log N_\epsilon \leq \sqrt{n} \}$ so that the expected L_1 error of the MDE Ψ_n converges to 0 uniformly, provided $f \in \tilde{\mathcal{F}}$ and $\tilde{\mathcal{F}}$ is totally bounded.

Situation 3
(realistic).

But if $f \notin \tilde{\mathcal{F}}$ and $f \in \mathcal{F}$, the set of all L_1 densities.

then by Corollary 10

$$E \left(\int |\Psi_n - f| \right) \leq 3 \min_{g \in G_\epsilon} \int |g - f| + \sqrt{\frac{8 \ln(2N_\epsilon^2)}{n}} \leq 3 \min_{g \in \tilde{\mathcal{F}}} \int |g - f| + 3\epsilon + \sqrt{\frac{8 \ln(2N_\epsilon^2)}{n}}$$

MDE is robust and accounts for the distance between true f and the best estimate in the class $\tilde{\mathcal{F}}$. (by dfn. N_ϵ an $\Delta^k \neq$) If we are looking in $\tilde{\mathcal{F}}$ for one of the best estimates, we are looking in $\tilde{\mathcal{F}}$ for one of the best estimates.

(18)

Situation 4 $f \in L_1$ and we want to construct from data X_1, \dots, X_n
 a set of densities $\{f_{ni} : 1 \leq i \leq k_n\}$ s.t.

$$3 \min_{1 \leq i \leq k_n} \int |f_{ni} - f| \rightarrow 0 \quad \text{and} \quad 4 \sup_{A \in \mathcal{A}} \left| \int_A f - M_n(A) \right| \rightarrow 0$$

Want data-adaptive strategies for constructing $\{f_{ni} : 1 \leq i \leq k_n\}$.

Want better bounds for this that that based on $|\mathcal{A}|$.

Thm 11 (Vapnik-Chervonenkis #)

$$E \left\{ \sup_{A \in \mathcal{A}} |M_n(A) - \mu(A)| \right\} \leq 2 \sqrt{\frac{\ln(2 S_{\mathcal{A}}(n))}{n}}$$

where $S_{\mathcal{A}}(n) = \max_{x_1, \dots, x_n \in \mathbb{R}^d} |\{ \{x_1, \dots, x_n\} \cap A ; A \in \mathcal{A} \}|$.

$S_{\mathcal{A}}(n)$ is V-C shatter coefficient and gives the maximal number of different subsets of a set of n points that can be obtained by intersecting with elements of \mathcal{A} .

Proof (Giné & Zinn (1984))

Introduce X'_1, \dots, X'_n , as an independent copy of X_1, \dots, X_n

$\sigma_1, \dots, \sigma_n$, as n iid. sign variables with $P\{\sigma = -1\} = P\{\sigma = +1\} = \frac{1}{2}$.

That are also independent of $X_1, X'_1, X_2, X'_2, \dots, X_n, X'_n$.

Let
$$\mu'_n(A) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}[X'_i \in A]$$

Then,

$$\begin{aligned}
 & E \left\{ \sup_{A \in \mathcal{A}} |M_n(A) - M(A)| \right\} \\
 &= E \left\{ \sup_{A \in \mathcal{A}} |E \{ M_n(A) - M'_n(A) \mid X_1, \dots, X_n \}| \right\} \\
 &\leq E \left\{ \sup_{A \in \mathcal{A}} E \{ |M_n(A) - M'_n(A)| \mid X_1, \dots, X_n \} \right\} \quad (\text{by Jensen's } \neq \\
 &\hspace{15em} \text{and convexity of } |\cdot|) \\
 &\leq E \left\{ \sup_{A \in \mathcal{A}} |M_n(A) - M'_n(A)| \right\} \quad \left(\begin{array}{l} \sup E(\cdot) \leq E \sup(\cdot) \\ * E(E(\cdot | \mathcal{G})) = E(\cdot) \end{array} \right) \\
 &= \frac{1}{n} E \left\{ \sup_{A \in \mathcal{A}} \left| \sum_{i=1}^n \sigma_i (\mathbf{1}_{[X_i \in A]} - \mathbf{1}_{[X'_i \in A]}) \right| \right\} \quad (\text{since } X_1, X'_1, \dots, X_n, X'_n \text{ are} \\
 &\hspace{15em} \text{i.i.d.}) \\
 &= \frac{1}{n} E \left\{ E \left\{ \sup_{A \in \mathcal{A}} \left| \sum_{i=1}^n \sigma_i (\mathbf{1}_{[X_i \in A]} - \mathbf{1}_{[X'_i \in A]}) \right| \mid X_1, X'_1, \dots, X_n, X'_n \right\} \right\}
 \end{aligned}$$

Since σ_i 's are independent of $X_1, X'_1, \dots, X_n, X'_n$, let us fix $X_1 = x_1, X'_1 = x'_1, \dots, X_n = x_n, X'_n = x'_n$ and investigate

$$E \left\{ \sup_{A \in \mathcal{A}} \left| \sum_{i=1}^n \sigma_i (\mathbf{1}_{[x_i \in A]} - \mathbf{1}_{[x'_i \in A]}) \right| \right\}.$$

Let $\hat{\mathcal{A}} \subseteq \mathcal{A}$ be a collection of sets s.t. any two sets in $\hat{\mathcal{A}}$ have different intersections with the set $\{x_1, x'_1, \dots, x_n, x'_n\}$.

and every possible intersection is represented once

Thus, $|\hat{\mathcal{A}}| \leq S_{\mathcal{A}}(2n)$, and

$$\begin{aligned}
 & E \left\{ \sup_{A \in \mathcal{A}} \left| \sum_{i=1}^n \sigma_i (\mathbf{1}_{[x_i \in A]} - \mathbf{1}_{[x'_i \in A]}) \right| \right\} \\
 &= E \left\{ \max_{A \in \hat{\mathcal{A}}} \left| \sum_{i=1}^n \sigma_i (\mathbf{1}_{[x_i \in A]} - \mathbf{1}_{[x'_i \in A]}) \right| \right\}
 \end{aligned}$$

Now, since each $\sigma_i(1_{[x_i \in A]} - 1_{[x'_i \in A]})$ has mean zero and range $[-1, 1]$, by Lemma 4 (upper bounding MGF of a bounded RV), we have for any $s > 0$:

$$\begin{aligned} E \left\{ e^{s \sum_{i=1}^n \sigma_i(1_{[x_i \in A]} - 1_{[x'_i \in A]})} \right\} \\ = \prod_{i=1}^n E \left\{ e^{s \sigma_i(1_{[x_i \in A]} - 1_{[x'_i \in A]})} \right\} \leq \prod_{i=1}^n e^{s^2/8} = e^{ns^2/8}. \end{aligned}$$

Since the distn. of $\sigma_i(1_{[x_i \in A]} - 1_{[x'_i \in A]})$ is symmetric, Lemma 6 (expected Maximal Deviation) implies that

$$\begin{aligned} E \left\{ \max_{A \in \hat{\mathcal{A}}} \left| \sum_{i=1}^n \sigma_i(1_{[x_i \in A]} - 1_{[x'_i \in A]}) \right| \right\} &\leq \sqrt{2n \log 2 S_{\mathcal{A}}(2n)} \\ &\leq \sqrt{2n \log 2 S_{\mathcal{A}}(n)^2} \\ &\quad (\text{since } S_{\mathcal{A}}(2n) \leq S_{\mathcal{A}}(n)^2) \end{aligned}$$

by dividing by $\frac{1}{n}$ on both sides we get V.C. \neq :

$$E \left\{ \sup_{A \in \mathcal{A}} |M_n(A) - M(A)| \right\} \leq 2 \sqrt{\frac{\log 2 S_{\mathcal{A}}(n)}{n}} \quad \square$$