Chapter 5: Minimaxity · Suppose we observe WEW drawn from a dist. P belonging to a statistical model P. > We'll be moss interested in setting P=Qⁿ with Qⁿ denoting the n-told product measure of dieso Q. Based on data realization w, we take action aeA $\begin{array}{c} & \mathcal{W}e'\mathcal{U} \ consider \ a \ collection \ of \ decision \ rules \\ T: \mathcal{W} \longrightarrow \\ \end{array}$ The quality of an action a is judged by loss fune. L: & · P -> R. The quality of a decision rule, the risk which corresponds to the expected loss: $\mathcal{R}(\mathcal{T},\mathcal{P}) = \int L(\mathcal{T}(\mathbf{n}),\mathcal{P}) d\mathcal{P}(\mathbf{n})$ Sx 1 = Point estimation The objective is estimate $\Psi(P) \in \mathbb{R}$. 4 E.g.: P=Q", and I(P)= SxdQ(x) If $L(a, P) = \{a - \overline{\Psi}(P)\}^2$ is the squared enor loss

then $R(T, P) = \int \{T(w) - \overline{\Psi}(P)\}^2 dP(w)$ is the mean-squared enor. EX2: Estimating a regression function - We observe n'id copies of (X.Y)~Q with support on X × R and we wish to estimate the fune. $f_{\alpha} : x \mapsto \mathbb{E}_{\alpha}(Y|X=x)$ The action space is a collection of funcs. mapping N to R. We may consider an intergrated squared enn loss: $L(a, P) = \int \{a(x) - fa(x)\}^2 dv(x)$ The corresponding risk is known as the mean integrated squared enor. Recall from 581. performance of decision rule Can be judged by its maximal risk: Sup R(T,P) PEOP A minimum estimator minimises the maximal risk

<u>Challenge</u>: Deriving a minimax estimator is eften an intractable problem. Minimax rate optimality Because ne often can't derive the minimax estimator, we'll instead look to find estimator that achieves minimas optimal rate. For now, consider the case P=Q" For each n, let Tn denote the collection of allowed decision rules, In many problems, the risk is always non-negative, and inf sup $R(T, Q^{*}) \xrightarrow{n \to \infty} 0$ TETA QEQ Focus on these devision problems here. A seq. of decision rules $\{T_n\}_{n=1}^{\infty}$ is called minimax rate optimal if the following: (*) limint Tern ace R(T.O") $sup R(T_n, Q^n)$ $Q \in Q$

inf sup R(T,Q") TEJn Q sup R(Tn, Q") sup R(Tn, On) sup R(Tn, Qn) If the numerator is of the form denominator is of the form Cina and C2n^{-B} and (\bigstar) holds off $\alpha = \beta$. Problem: We can't typically derivene the form of the numetor or denometer of (\$)! Sol'n: Find bounds on them Later in this course, we'll derivan upper bounds on the maximal risk of particular estimaten see 87n Sh=,. > That is, we'll find Ella Jaz, s.E. $\sup_{\Omega} R(T_n, Q^n) \leq U_n$ In this chapter, we'll desire limer bounds on the minimax risk, of the form int sup $R(T, Q^n) \ge L_n$ $T \in T_n \quad Q$

Once we have ELn), EUns, K asymp $(\bigstar) \frac{int}{q} \sup_{\mathcal{R}(T_n, \mathcal{R}^n)} R(T, \mathcal{R}^n) \geq \frac{L_n}{u_n} > 0$ Hence, to show (A), v lissing $\frac{L\eta}{U\eta} > 0$ va suffices to show Lower bounds ELnsn=1 int sup $R(T, P) \ge L$ Recall from 581: For any prior TI on P. the Bayes risk of a decision rule T is given by $r(T,TT) = \int R(T,P) dTT(P)$ Lemma: For any decision prob. with rule T, $sup R(T, P) \ge sup r(T, \pi)$ π Proof: VTT, sup over op is lower bounded by $\sup_{\mathcal{P}} R(T, \mathcal{P}) \geq r(T, \pi)$ Taking sup gives the result.

Thm: For any decision prob., mf sup R(T, P) > sup mf r(T, π) T P minimax risk Bayos risk under LFP. Proof: The max-min ineq. yields inf sup $r(T, \pi) \ge \sup_{T} \inf_{T} r(T, \pi)$ Combine this w/ the preceeding Som. inf sup $R(T,P) > \sup_{T} \inf_{T} r(T,T)$ Note: The above implies that, for any prive T, $\inf_{P} \sup_{P} R(T, P) \ge \inf_{T} r(T, T)$ Mechod # 1 for deriving minimax lower bounds; Le Cam's method.

Thm. (Le Cam): Let R be a risk func. defined
according to some non-negative loss L.
$For any r_1, r_2 \subset \mathcal{T},$
inf sup $R(T,P) \ge \frac{1}{2} d(P_1,P_2) P_1 \wedge P_2 _1$
$\geq \frac{1}{4} d(P_1, P_2) \exp[-kL(P_1, P_2)]$
nhere
"discrepancy" $d(P_1, P_2) = \inf_{a \in S} [L(a, P_1) + L(a, P_2)]$
"testing IIP, $\Lambda P_2 II_1 = \int min \left\{ \frac{dP_1}{dv}(w), \frac{dP_2}{dv}(w) \right\} dv(w)$
$KL(P_1, P_2) = \int \int log \left[\frac{dP_1}{dP_2} (m) \right] dP_1(m) P_1 << P_2$
(+ x)
Remark: Understanding discoveryency.
Ne'll look at the discrepency in 2 cases.
a) Point estimation w/ squared own loss.
We want to estimate I(P) E R and the lass
$L(a,P) = \{a - \overline{L}(P)\}^{2}$
A little calculus shows that
$d(P_1, P_2) = \inf_{a \in \mathbb{R}} \left[L(a, P_1) + L(a, P_2) \right]$

 $= \frac{1}{2} \left[\overline{\Psi}(P_i) - \overline{\Psi}(P_2) \right]^2$ b) Estimating a funce w. / integrated squared enor loss soy P=Qⁿ, the action space is a convex subset of the space of funcs mapping from X to R, and $L(a,P)=\int [a(x)-fa(x)]^{2}dy(x)$ In this case, $d(P_1, P_2) = \frac{1}{2} \int [f_0(x) - f_{02}(x)]^2 dv(x)$ where $P_1 = Q_1^n$, $P_2 = Q_2^n$. Remark: Understanding the testing affinity Statistically, we can show that given a draw from P; , j = [1,2] is unknown, the difficulty of identifying the value of j is quantified by the testing affinity. Pr(n) P2(W) Pilm

11PINP211= /min {PIP23 dv - Consider a test $T: \mathcal{W} \to \{1, 2\}$ We'll quantify performance of the decision rule T via misclassificiens enn : $\mathcal{E}(T) = \mathbb{E}_{P_i} \left[\mathbb{1} \{ T(w) = 2 \} \right]$ $+ E_{P_2} [1 \{ T(w) = 1 \}]$ Observe that $\Sigma(T) = \int \left[P_{i}(w) \pounds \{ T(w) = 2 \} \right]$ $+ P_2(n) 1 \{ T(n) = 1 \} dv(n)$ $\geq \int min \{P, (w), P_2(w)\}$ $\begin{bmatrix} \int \int \langle T(w) \rangle = 2 \\ \langle T(w) \rangle = 2$ +1[T(m)=1]]dv(m) $= \int' min(P_1, P_2) dr$ $= l(P_1 \Lambda P_2) l_1$ The lower bound is achieved by TX(w) = augmax P; (w)

Remark: 11P1 1P2111 = 1- TV(P1, P2), where TV(P1,P2)= sup [P,(A)-P2(A)] Proof of main Thm: (Le Cam) Let TT denote a uniform prior over $\{P_1, P_2\}, \text{ that is } TT(P_1) = TT(P_2) = 1/2$ For any rule T, $r(T,\pi) = \sum_{i=1}^{2} \int L[T(w), P_{i}] P_{i}(w) dv(w) \pi i P_{i})$ $= \frac{1}{2} \int \sum_{j=1}^{2} L(T(w), P_j) P_j(w)] dv(w)$ $2 \frac{1}{2} \int mins \mathcal{E} \mathcal{P}_i(w), \mathcal{P}_2(w) \mathcal{E}$ $\left[\sum_{j=1}^{2} L(T(w), P_{j})\right] dr(w)$ $\geq \frac{1}{2} \int min \{ P_1(w), P_2(w) \}$ $\left[\inf_{z \in L(a, P_{z})}dv(w)\right]$ = = & (P1, P2) ||P1/P2/11

Why we also need other lower bounding techniques? - In certain cases, Le Cam will provide a rate-optimal lower bound Is In fact, this is the case for prob. 2 HW1. which concerns estimation of a smooth density at a point to In other cases, Le Cam can't give a rateoptimal LB is This is the case when the goal is to estimate a smooth regression funcin terms of the mean integrated squared error.

Thm (Fam's method): Let Rbe a risk func, defined awarding to some non-negacine loss L. For N ? 3, P1, P2, ..., Px EP and define $\begin{aligned} \mathcal{L} &= \min \ d(P), P_k) = \min \ \inf \left[L(a, P_j) \\ \vec{J} \neq k \\ \vec{P} &= \frac{1}{N} \sum_{j=1}^{N} P_j \\ (uniform measure over \\ P_1, \dots, P_N) \end{aligned}$ $\inf_{T} \sup_{P} R(T,P) \ge \frac{2}{2} \left[1 - \frac{\log 2 + \frac{1}{N} \sum_{j=1}^{N} kL(P_j,\bar{P})}{\log N} \right]$ $\geq \frac{1}{2} \left[1 - \frac{\log 2 + \max k L(P_j, P_k)}{\sum_{j \neq k} log N} \right]$ $=\frac{1}{2}\left[1-\frac{\log 2+k}{\log N}\right],$ where $K = \max_{\substack{i \neq k}} kL(P_i, P_k)$

Remark . The bound increases when : 4 N increases and l.K remain the same L> K decreases and N.1 remain the same L> 1 increases and N, K remain the same Moreover, adding a dist, to a set {P1, P2, ..., PN } · N++ 2--· K++ Proof: (Fano) Let TT denote a uniform prior on P1,..., Prv. For any rule T, we have that $r(T, TT) = \frac{1}{N} \sum_{j=1}^{N} \int L(T(w), P_j) P_j(w) dv(w)$ $\geq \frac{1}{2N} \sum_{i=1}^{N} \int 4 \{L(T(w), P_{i}) \geq \frac{1}{2}\} P_{i}(w) dv(w)$ (2) where we've used that $l = \frac{1}{2} 1 \{l \} = \frac{1}{2} \}$ for all l. n/2 1/2

Continuing Eqn. (2), (3) $r(T, \pi) = \frac{1}{2} \left[1 - \frac{1}{N} \int_{1}^{N} \frac{1}{2} \left[(T(n), P_{j}) + \frac{1}{2} \right] \right]$ $P_{i}(n) dr(n)$ E MAX PS(W) For each w, $\sum_{i=1}^{n} \underline{A} \{ L(T(n), P_j) < \frac{n}{2} \} P_j(n)$ $\leq [\max P_j(n)] \sum_{n=1}^{\infty} \mathbb{1}\left\{L(T(n), P_j) - \frac{1}{2}\right\}$ It suffices to show T = 1 $(max P_j(w)) \sup_{a \to -i} \sum_{j=1}^{n} 1 \{ L(a, P_j) < \frac{1}{2} \}$

To show $\sum_{a} \sum_{y=1}^{N} \mathbb{I}\{L(a, P_{y}) < \frac{1}{2}\} \leq 1$ it will be useful to define $\eta_{a} = \min[L(a, P_{y}) + L(a, P_{y})]$ Noce: n= infla ... 2 = 2a Va. Fix an action R, We'll show that $\sum_{i=1}^{N} \mathbb{1}\{L(a, P_i) < \frac{2a}{2}\} \leq 1$ How many points are at the lefe ? Na/2 $\frac{}{L(a,P_4)} \times L(a,P_2)$ $L(a, P_3)$ $L(a, P_i)$ $L(a, P_2)$ 0 L(a,P4) $L(a, P_3)$ $L(a, P_i)$ Also, because n = la, $\frac{2}{2} \mathbb{1} \{ L(a, P_j) < \frac{1}{2} \} \leq \frac{2}{2} \mathbb{1} \{ L(a, P_j) < \frac{1}{2} \}$ Taking a sup oner all actions on the left

 $P(ug into (3), (4)) + r(T, \pi) \ge \frac{n}{2} [$ Tr Amax P; (n)] dr(w) D_n slide 23, we show that $\max P_j(n) \leq N \overline{p}(n) \log_2 + \sum_{j=1}^{n} P_j(n) \log \frac{P_j(n)}{\overline{p}(n)}$ log N $= \int max P_{J}(w) dr(w) \\ = N \log 2 + \tilde{\Sigma} K L(P_{J}, \bar{P})$ logN log 1 + Tr JE KL(PS, P) log 1 $\frac{1}{N}\int [mass P_{J}(n)] dr(n) \leq$ Plug is L4), log 2 + tr Z KL(P; - p) log N $r(\tau,\pi) \geq \frac{1}{2}$

9x. Minimax lower bound for estimating a smooth regression func. (Xn, Yn) ~ QEQ (χ_{ℓ}, Y_{ℓ}), - We observe - For each Q, the marginal diss. of X is Umi[0,1] Moreover, $Y|\chi = x \sim N(f_0(x), 1)$ where for is smooth - in parcicular, fa & F(B,L), where F(B,L) is the following Hölder class: F(B,L) := { f: f is l-times differentible fixed and $|f^{(l)}(x_i) - f^{(l)}(x_2)| \le L|x_i - x_i|^{\beta-l} \le$ nonneg. conves. for all X1, X2 E[0, 1] S where I is the greatest integer that is Strictly less than p. Generalized Lipschots Class Noce: A suff cond for a fune. I to belong to F(p,L) is that f be (l+1)-times differentible and $\sup |f^{(l+i)}(x)| \in L$

Our objective is to estimate $f_{\alpha}(x) = E_{\alpha}[Y|X=x]$ -> Performance will be quancified by MISE, arising from the loss $L(a, \tilde{P}) = \int \left[a(x) - fa(x)\right]^2 dx$ We'll derive a minimax lowor bound ria Fano method. > To do this, we must select a seg/ collection Po:= {P1..., PNS of dists, in our model. > We're letting Pj = Qj. For any $w \in \{0, 1\}^m$, we define $f_w(x) = \sum_{j=1}^m w_j \phi_j(x)$, where here, for a fixed const. h>0, we have $m \in [8, \frac{1}{h}-1]$ and $\varphi_{\mathcal{J}}(x) = L h^{\beta} K \left(\frac{x - \partial m + 1}{R} \right)$ Here, for suff small a >0, K is defined as

 $K(x) = \operatorname{Rexp}(-\frac{1}{1-4x^2}) \operatorname{I}[1x] < 1/2$ 1/2 The quantity a is chosen that for EF(B,L) for all w E[0,1]^m. Recall $f_{\omega}(x) = \sum_{j=1}^{\infty} \omega_{\sigma} \phi_{\sigma}(x)$ $\mathcal{W}_{1} =$ $\omega_2 = 1$ $\omega_3 = 0$ $W_{4} = 1$ 2 3 m+1 m+1 mti . m+1.

M, for some 72 ≤ 80, 13", we applied Fano's method based on the collection of dists. $\{P_w := O_w : w \in \widehat{\Sigma}\},\$ where Quo E Q has regression for, then we'd find that int sup $R(T, P) \ge \min_{\substack{w \neq v \in \widehat{\Lambda}}} d(P_w, P_v)$ 2 [1- log2+max KL(Pw, Py w=r log IRI] $d(P_{\omega}, P_{\nu}) = \frac{1}{2} \int \left[f_{\omega}(x) - f_{\nu}(x) \right]^2 dx$ Also, from slide 35, $KL(Pw,Py) = \frac{n}{2} \int_{0}^{1} \left[f_{w}(x) - f_{v}(x) \right]^{2} dx$ Since none bump fune, overlag, $\phi_{\mathcal{J}}(x)\phi_{\mathcal{K}}(x)=0$ for $\forall X$ Only diagonal terms left.

 $\int [f_w(x) - f_v(x)]^2 da$ $= \sum_{j=1}^{m} \left[w_{j} - v_{j} \right]^{2} \int \phi_{j}(x)^{*} dx$ $= \sum_{j=1}^{m} \left[w_{j} - v_{j} \right]^{2} L^{2} h^{+j} \int K(u)^{2} du$ $= \sum_{j=1}^{m} \left[w_{j} - v_{j} \right]^{2} L^{2} h^{-\beta+j} \int K(u)^{2} du$ $C_2 L^2 h^{2\beta+1} \sum_{j=1}^{m} [w_j - v_j]^{\nu}$ $=: H(W, V) \qquad Hamming$ dissancePlug in and e3: c2l2 we find $d(P_w, P_v) = {}^{c_3} k^{2\beta + i} H(w, v)$ and $KL(Pw,Pv) = C_3nh^{2\beta+1}H(w, 0)$ Also, since $H(W, V) \leq n$, we have the further lower bound $KL(PW, P_{y}) \leq C_{3}nh^{2\beta+1}m$ $\max_{\substack{w\neq v}} KL(P_w, P_v) \leq C_3 n h^{2\beta+1} m$

Plug in above, $\frac{c_{s}h^{2\beta+1}}{2} \underset{w\neq v}{\overset{\mu}{\longrightarrow}} \mathcal{H}(w,v)$ inf sup R(T,P) > [1 - log2+C3nh2p+1 log151] Following, think about two choices of the index set $\Sigma \subseteq \{0,1\}^m$ Choice 1: What if $\pi = \{0, 1\}^m$? (Envire thing) In this case, $min_{w\neq y} \mathcal{H}(w, y) = 1$ logIn 1 = mlog 2 Plug in, we have $\inf \sup R(T, P) \ge \frac{C_3 h^{2\beta+1}}{2} \left(1 - \frac{\log_2 + (3nh^{2\beta+1})}{m \log_2 2}\right)$ Ino obs 1. If $nh^{2\beta+1}$ is large then RHS is negative \rightarrow Not useful at all Hence, we need $h = O(n^{-\frac{1}{2\beta+1}})$

2. In the best case, the RHS is no more than C3 h 2 Combining these two, we see at best, the RHS would give a lower bound on order of 1/n From 581, you might expect this lover bound is loose, In particular, consider the parametric prob. fa e F where $\mathcal{F} := \{ x \mapsto \beta x : \beta \in [0, L] \}$ (14) $\rightarrow \mathcal{F} : small class$ f m) Holden Smooth From last quarter, you know ME sarcifies $\sqrt{n}(\beta - \beta_0) \longrightarrow \mathcal{N}(0, \theta^2)$ Hence, $n(\hat{\beta} - \beta_0)^2 \longrightarrow \chi^2(1)$

and $(\beta - \beta_0)^2 = O_p(1/n)$ Hence, it's reasonable to expect that $\mathbb{E}_{o_{n}^{n}}\left[\int (\widehat{\beta}_{x} - \beta_{o} \times)^{2} dx\right] = O\left(\frac{1}{n}\right)$ Chore 2: Question: Is there a subset i of 50,1" for which In large and Hamming dissonce of distinct elements is also large? (1,1,0) (1, 0, 0)(D,0,0) Answer: Yes! Varshamov- Gilbert Lemma If m 7, 8, then there exists a subset I of {0,1}" s.t. [I] is large ? 2^{m/8} and min H(w, v) > m/8 $w \neq v \in \Omega$

Choice 2: 52=52 In this case, our lower bound gives that inf sup $R(T, P) \ge \frac{c_3 h^{-p+1} m}{16} \left(1 - \frac{\log 2 + c_3 n h^{-p+1} m}{m \log 2/8}\right)$ $= \frac{c_{3}h^{2\beta+1}}{16} \left(1 - \frac{8}{m} - \frac{8c_{3}nh^{2\beta+1}}{\log 2}\right)$ Choose largests m possible. $(m = \lfloor \frac{1}{R} - 1 \rfloor)$ We'll approximate this lower band m = i/L yielding $LB \triangleq \frac{c_s h^{2\beta}}{16} \left(1 - 8h - \frac{8C_{3n} h^{2\beta+1}}{log 2} \right)$ A non-negative? We now choose h. For $(A) \ge 0$, We need $h \le C_4 n^{-\frac{1}{2\beta+1}}$ for some C_4 . Letting h= (4n^{-2p+1}, shows that there exists const C5 s.t. inf sup R(T,P) ? Cs n 2B+1 For example, if B=2, we derived n⁴⁵ rate lower bound on the minimax risk.

•	•	•	•		t	h	ìs		bn	עע	in	d	•	Ù	Š	C	y	A	a	K. 3	2	ì	n	7	e	rh	25		0	f	r	a	ć	25)	1	• •
•	•	•	•	•	•		•	•	•	•	•	•	•	•		•	•	:	•	•	•	•	•	•	•	•	:	:	•		•	•	•	•		•	
														٠																				•			• •
•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•	
												0							0		0																
					۰	•							•	٠				•							•	•	•	•	•	•				•	•		• •
					•									•																							
																																		•			
•	•	•	•	•	•				•	•	•	•	•	•	•	•	•		•		•		•	•	•	•	•	•	•	•	•	•	•	•	•	•	
•						•								•											•	•	•		•	•				•	•		
		•			•	•					•	•		•				•	•		•				•	•	•		•	•				•	•		
•		•				•								•			•	•				•	•	•	•	•			•					•	•	•	• •
														•													•			•				•			
•		•		•	•	•	•	•	•	•	•	•		•			•	•	•	•	•	•	•	•	•	•	•		•	•				•	•	•	• •
•	•											•		•	•		•	•	•		•		•		•	•	•		•	•				•	•	•	• •
																	•						•		•	•			•								
•		•				•								•			•	•				•	•	•	•	•			•					•	•	•	• •
														•																				•			
•		•		•	•	•	•	•	•	•	•	•		•			•	•	•	•	•	•	•	•	•	•	•		•	•				•	•	•	• •
						•								•			•	•				•	•	•										•	•	•	• •
																																		•			
•	•	•				•			•	•				•	•	•	•					•	•	•	•	•	•	•	•	•	•	•	•	•	•		• •
						•					•			•	•												•			•				•	•		
•	•	•	•						•	•		•	•	•	•		•	•	•		•		•	•			•			•				•			
																																		•			