

# Stat 991 Lecture 3 - Prediction regions & Conformal inference

(given  $\underbrace{z_1, z_2, z_3, \dots, z_n}_{\text{map}} \in \mathbb{Z}$ , predict  $\underline{z_{n+1}}$  (or some part of it))  
 e.g. ~~supervised setting~~  $\underline{z_i} = (x_i, y_i)$   $\underline{z} \rightarrow$  training data. test datapoint  $\underline{z_t}$

- Prediction region: map  $T: \mathbb{Z}^n \rightarrow \mathcal{B}$   $\mathcal{B}$   $\rightarrow$  "Subsets of  $\mathbb{Z}$ " sigma-algebra over  $\mathbb{Z}$

$$[z_t \in T(z)] \quad \{ \dots [ \dots ] \dots \} \subset \mathbb{Z}$$

$$z_{n+1} \in T(z_1, z_2, \dots, z_n)$$

- marginal coverage:  $P_{z, z_t} (z_t \in T(z)) \geq 1 - \alpha$

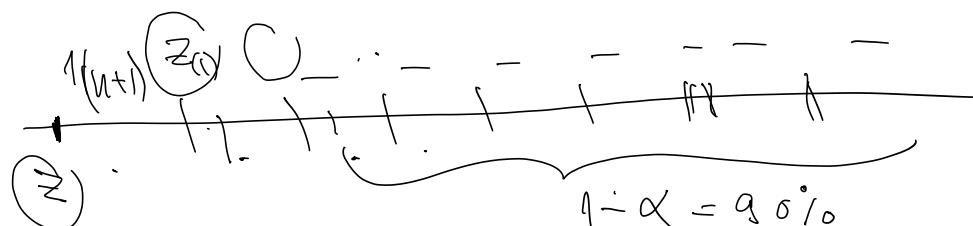
- contrast with point prediction:  $f: \mathbb{Z}^n \rightarrow \mathbb{Z}$

- Intuition 1-D:  $\frac{1}{(n+1)} \frac{1}{(n+1)} \dots \frac{1}{(n+1)}$   $\xrightarrow{\text{statistically equivalent blocks}}$

$z_{(i)}$  -  $i^{\text{th}}$  order statistic of sample  $\mathbb{Z}$

$$P(z_t \in [z_{(j)}, z_{(j+1)}]) \geq \frac{1}{n+1} \quad \left[ \begin{array}{l} \text{assuming ties occur w/p } 0 \\ \text{with } n+1 \end{array} \right]$$

$z_t, z_1, \dots, z_n \stackrel{\text{iid}}{\sim} f \text{ density on } \mathbb{R}$



- ranks:  $\text{rank}(x, S) = \#\{s \in S : s \leq x\}$

$$\text{rank}(1, \{2, 0\}) = 1, \text{ b/c } 0 \leq 1 \rightarrow 2 \neq 1$$

$$P_k(z) = \frac{1}{n+1} \cdot \underbrace{z}_{\text{stat. equiv. blocks}}$$

$$\geq 1 \rightarrow n+1 \rightarrow \frac{n+1}{n+1} = 1$$

$$\geq 2 \rightarrow n \rightarrow \frac{n}{n+1}$$

$$\geq k+1 \rightarrow n+1-k \rightarrow \frac{n+1-k}{n+1} = 1 - \frac{k}{n+1}$$

want  $\geq 1-\alpha$

$$1 - \frac{k}{n+1} \geq 1 - \alpha$$

$$\alpha \geq \frac{k}{n+1}$$

$$\underline{\alpha}(n+1) \geq k \rightarrow k = \lfloor \underline{\alpha}(n+1) \rfloor$$

floor of  $0.5 - 0$   
 $\lfloor 0.5 \rfloor = 0$   
 $\lfloor 1 \rfloor = 1$

Conclusion :  $T(z) = \left\{ z : \underline{\alpha}(z, z_1, \dots, z_n) \geq \lfloor \underline{\alpha}(n+1) \rfloor \right\}$

has  $1 - \alpha$  <sup>marginal</sup> coverage

$T(z) = \left\{ z : \underline{\alpha}(z, z_1, \dots, z_n) \geq \lfloor \underline{\alpha}(n+1) \rfloor \right\}$

Example :  $n=5$  (

$1-\alpha=0.9$   $90\%$   $z \quad \lfloor 0.1 \cdot 5 \rfloor = \lfloor 0.5 \rfloor = 0$

$\Rightarrow T(z) = \mathbb{R} = (-\infty, \infty)$

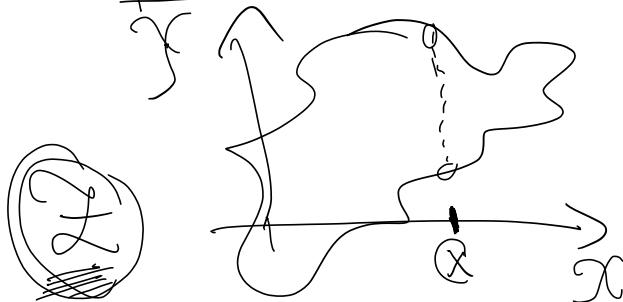
$80\%$   $\lfloor 0.2 \cdot 5 \rfloor = \lfloor 1.0 \rfloor = 1$

$T(z) = \left[ z_{(1)}, \infty \right)$

Then Fraser '50s

all "non-parametric" prediction regions need to be based on order statistics.

Supervised setting  $\underline{z} = (x, y)$



$$T(x; \underline{z})_{\underline{z}} = \{y : (x, y) \in T(z)\}$$

Claim :  $P(\underbrace{y_t \in T(x_t, z)}_{(x_1, y_1), \dots, (x_n, y_n)}) = P(\underbrace{z_t \in T(z)}_{\geq 1 - \alpha})$

4:18

Principle : given "exchangeable" data  $z_1, \dots, z_n$

consider "symmetric" fns  $\left[ \begin{array}{l} m(z_1, z) \in \mathbb{R} \\ m(z_2, z) \\ \vdots \\ m(z_n, z) \in \mathbb{R} \end{array} \right]$  preserve "exch."

& then use  $\mapsto$  construction.

$$z_1, \dots, z_n \rightarrow \frac{\sum_{i=1}^n z_i}{n}$$

$$m(z_1, \dots, z_n) = m(\underbrace{z_{\pi_1}, z_{\pi_2}, \dots, z_{\pi_n}}_{\text{permutation } \pi \text{ on } \{1, \dots, n\}}) \quad \text{for any } \pi \in S_n$$

Def:  $m(z) = m(z_\pi) \quad \forall \pi \in S_n$   
Symm function

$$\overline{(z_1, \dots, z_n)}$$

$$\overline{(z_n)}$$

$$m(a, b_1, \dots, b_n) = \left| a - \frac{b_1 + \dots + b_n}{n} \right|$$

$$m((z_1, z_2))$$

$$m(z_2, \underline{z})$$

$$m(z_n, \underline{z})$$

$$\left| z_1 - \frac{\sum z_i}{n} \right|$$

$$m(a)$$

If  $m(a, \cdot)$  is symmetric (in the last  $n$  variables)

&  $z_1, \dots, z_n$  are exchangeable

then  $\underline{m}(z_1, z), \dots, m(z_n, z)$  are exchangeable

Z

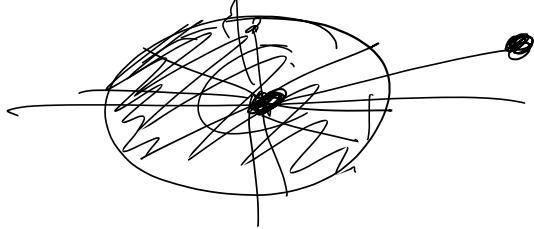
J

R

$$\boxed{m(a, b_1, \dots, b_n) = \|a - \sum_i b_i\|} \quad \|\cdot\| \text{ norm on } Z$$

$$\boxed{\|z_1 - \bar{z}\|, \|z_2 - \bar{z}\|, \dots, \|z_n - \bar{z}\| \in}$$

$$T(z_{n+1}, \bar{z}) = \tan^{-1} \left( \frac{\|z_{n+1} - \bar{z}^+\|}{\|z_j - \bar{z}^+\|} \right) \quad j=1, \dots, n$$

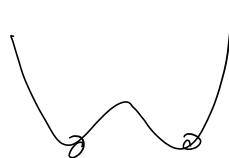


$$\bar{z}^+ = \{z\} \cup \{z_{n+1}\}$$

$$= \{z_1, \dots, z_{n+1}\}$$

$$\begin{aligned} &\geq L(n+1)\alpha \\ &\leq R(n+1)(1-\alpha) \end{aligned}$$

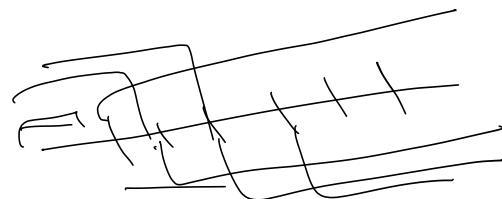
$$\text{E.g.: } Z_i = (x_i, y_i)$$



$$m(z, z^+) = |y - \hat{f}(x, z^+)|$$

$\hat{f}(x, z^+)$  is symmetric in the 2<sup>nd</sup> variable  $z^+$   
(prediction fn. does not depend on order)

$$\text{OLS: } \hat{f}(x, z^+) = \underbrace{x^\top (\cancel{\hat{X}^\top \hat{X}})^{-1} \cancel{\hat{X}^\top} y}_{\text{OLS.}} \quad \min_{\beta} \|y - X\beta\|_2^2$$



$$x : (n+1) \times \dim(x)$$

$$y : (n+1) \times 1 \rightarrow \text{quantiles.}$$

$$|y_{n+1} - \hat{f}(x_{n+1}, z^+)| \leq \underbrace{\alpha \left| y_i - \hat{f}(x_i, z^+) \right|}_{i=1 \dots n}$$

$$y_i - \underbrace{x_i^\top (\hat{X}^\top \hat{X})^{-1} \hat{X}^\top}_{u_i} \begin{pmatrix} y_1 \\ \vdots \\ y_{n+1} \end{pmatrix}$$

$$= y_i - u_i^\top y = y_i - \sum_{j=1}^{n+1} u_{ij} y_j$$

$$= v_i^\top y = \underbrace{a_i y_{n+1}}_{= a_i} + \underbrace{b_i y_i}_{= b_i}$$

need to have sufficiently many inequalities of the form

$$|\bar{a}_{n+1} \underbrace{y_{n+1}}_{= y} + \bar{b}_{n+1}| \leq |\bar{a}_i y_{n+1} + \bar{b}_i| \quad i=1 \dots n$$

$$y \rightarrow |a_i y + b_i|$$

$$y \rightarrow |ay + b|$$

$$rk \leq k$$

11

$y_{nt}$

2 1, 2 3 2 3 2

1

2 3 2

1

disconnected

1

Exchangeability :  $(z_1, z_2, z_3 \dots, z_n \dots)$

each sequence  
if for any permutation  $\pi \in S_n$

$(\underbrace{z_1, z_2, \dots, z_n, \dots}_n)_{n>1}$  has the same distribution as

$$(z_{\overline{1}}, z_{\overline{1}_2} \dots z_{\overline{n}_n}, z_{n+1} \dots)$$

e.g., if  $(z_i)_{i \geq 1}$  are iid.  $\Rightarrow$  exch.

de Finetti's thm 1937 : excl  $\Rightarrow$  mixture of iids  
(?) latent param  $T$

(7) latent param  $\Gamma$

st. cond. on  $\Gamma$ ,  $(Z_i)_{i \in I}$  are iid

## Pf of tank construction

$$\left( \text{Th} \left( \overbrace{\underline{z_{n+1}}, \underline{z_1} \dots \underline{z_n}}^{x_{n+1}} \right) \right) \geq \dots$$

↑ Claim:  $(z_1, z_2, \dots, z_{n+1})$  are exch, (then) is uniformly dist.  $h[1, \dots, n+1]$  & this occurs w.p. = 0

rank ( ~~$\sum_{n+1}$~~ ,  $\underline{z_1 \dots z_{n+1}}$ )

= rank ( ~~$\sum$~~ ,  $\underline{z_1 \dots z_m}$ )

=  $\text{rn}(\underline{z_2 \dots -})$

$X = dY$

$X$  has same dist.  
as  $Y$ .

$U \sim \text{Unif}\{1 \dots n+1\}$

rank ( ~~$\sum_{n+1}$~~ ,  $\underline{\{z_1 \dots z_m\}})$

condition

$\underline{z_1 \dots z_m}$

~~$\sum_{n+1}$~~

~~$\sum_{n+1}$~~

$\sum_{n+1}^n \approx \text{Uniform}\{1 \dots n+1\} \Rightarrow \text{rk}(\underline{z_{n+1}}, \underline{\{z_1 \dots z_{n+1}\}})$

$\sim U\{1 \dots n+1\}$

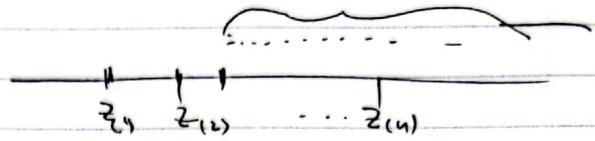


STAT 991

Lec #:

Recap.

$$\bar{z} = (z_1, \dots, z_n)^{1-D}$$



$$T(\bar{z}) = \{z : \text{rank}(z; \bar{z}) \geq \lfloor \alpha(n) \rfloor\}$$

Claim: if  $(z_1, \dots, z_n, z_{n+1})$  are exchangeable, then

$$P(z_{n+1} \in T(\bar{z})) \geq 1 - \alpha. \quad \leftarrow \text{marginal coverage (validity)}$$

How to use it? - deductions

(1) Functions preserving exchangeability.)

(2) Supervised setting: fiber / cross-section.

$$T(x, \bar{z}) = \{y : (x, y) \in T(\bar{z})\}$$

~~1-D const.~~

> Let's consider maps  $\bar{z}_i \rightarrow z'_i$

$$z_1 \quad m_1(\bar{z}) = z'_1$$

$$\vdots \quad \rightarrow \quad \vdots \quad \text{etc. for } \bar{z}' = m(\bar{z})$$

$$z_n \quad m_n(\bar{z}) = z'_{n+k}$$

then are  $(z'_1, \dots, z'_{n+k})$  are exch. [

When are  $z'_1, \dots, z'_{n+k}$  exch?

Claim: if for each permutation  $\pi \in S_n$ ,

(Deza & Verducci) there is a permutation  $\pi' \in S_n$

such that  $m(\bar{z}_{\pi}) = [\bar{m}(\bar{z})]_{\pi'}$

(permute  $\bar{z}$  by  $\pi \Rightarrow$  permute  $m(\bar{z})$  by  $\pi'$ )

Then  $\bar{z}' = m(\bar{z})$  is exch. ~~of coordinates~~.

$z_1 \quad z_2 \quad z_3$

Example

$$z_1 - \frac{z_1 + z_3}{2} \sim z_2 - \frac{z_1 + z_3}{2} \sim z_3 - \frac{z_1 + z_2}{2}$$

more generally:

$$m_i(z_1, \dots, z_n) = \tilde{m}(z_i; \{z\})$$

permutation  
inv. w.r.t.  
2<sup>nd</sup> argument.

$$\text{or } m_i(z_1, \dots, z_n) = \tilde{m}(z_i, \{z_{-i}\})$$

$\{z\}, \{z_{-i}\}$

used to check how well  $z_i$  "conforms" to  $z$ .

Remark:  $m(z_{\pi}) = [m(z)]_{\pi}$  ~~is eq. mean~~:  $m$  is permutation-equivariant. [or co-variant].

- this is an instance of more general group equivariance

$$\forall g \in G, \exists g' \in G': m(gz) = g'm(z),$$

- studied in

- representation theory of finite groups [Serre '77]

- modern perspective: Deep Network architecture

design - Deep Sets [Fisher et al '17]  
PointNet.

Research questions:

- what are the implications?  
"deep conformal prediction"

Key point: preserved under compositions on feature

$$z \rightarrow z' \rightarrow z''$$

↓      ↓  
pres. excl pres. excl

full algo

input:

$$\text{given: } (x_i, y_i) \quad i=1 \dots n.$$

$\vdash (x_{n+1}, \cdot)$

Output:

Prediction set for  $y_{n+1}$

$$T(x_{n+1}, z)$$

- for each  $y_{n+1}$  in [ how discretize? ].
- calculate  $w_i = \Delta(y_i, f(x_i, z^+))$ ,  $i=1 \dots n+1$ .
  - : preserve exch.  $\cancel{z^+} \rightarrow z^+$  - with  $\pi$  if you permute
  - $\cancel{z^+} \rightarrow z_{\pi}^+$  : if  $f$  is perm. inv.
  - then,  $w(z_{\pi}^+) = [w(z^+)]_{\pi}$  so. over.
- if.  $\text{rank}(w_{n+1}, w) \leq \bar{l}(n+1)(1-\alpha)$ , include  $y_{n+1}$  in  $T(x_{n+1}, z)$
- how obtain  $f$ ?  $\rightarrow f(\cdot)$  can be any function of  $z^+$ , if it depends symmetrically.

- e.g. ERM:  $\min_{f \in \mathcal{H}} \frac{1}{n+1} \sum_{i=1}^{n+1} L(y_i, f(z_i))$  | if # alg. = M L

(1) run any alg. after randomly permuting data.

What we really need is  $\tilde{w} =_d w_{\pi}$ .  $\cancel{\neq \pi}$ . (exch. of  $w$ ).

enough if  $\cancel{\exists^+} \cancel{z^+} =_d z_{\pi}^+ \cancel{+ \pi}$ .

$$w_i = \Delta(y_i, f(x_i, z_{\pi}^+)), \quad \tilde{\pi} \sim \text{Unit}(S_{n+1})$$

More generally

$$\bar{z} - [\bar{z}_\pi]_\pi$$

$$\begin{aligned} w(\bar{z}_\pi)_{\pi^{-1}(i)} &= \Delta \left( y_{\pi \cdot \pi^{-1}(i)}, f(x_{\pi \cdot \pi^{-1}(i)}, \bar{z}_{\pi \cdot \pi}) \right) \\ &= \Delta \left( y_i, f(x_i, \bar{z}_{\pi \cdot \pi}) \right) \\ &= \Delta(\bar{z}) = d. \end{aligned}$$

$$w(z)_i = (y_i + (x_i, \bar{z}_\pi^*))$$

So, we find  $w(\bar{z}_\pi)_{\pi^{-1}(i)} = d = w(z)_{\pi(i)}$

Moreover/similarly :  $w(\bar{z}_\pi)_{\pi^{-1}} =_1 w(z)$

$$w(\bar{z}_\pi) =_1 [w(z)]_\pi.$$

Now, if  $\bar{z}_\pi =_d \bar{z}$ , i.e.  $\bar{z}$  is exchangeable

$$\Rightarrow \text{So. } w =_1 w_\pi, \text{ so } w \text{ is mech.}$$

needed to  
be symmetric

$\bar{z}$  is inv exch.

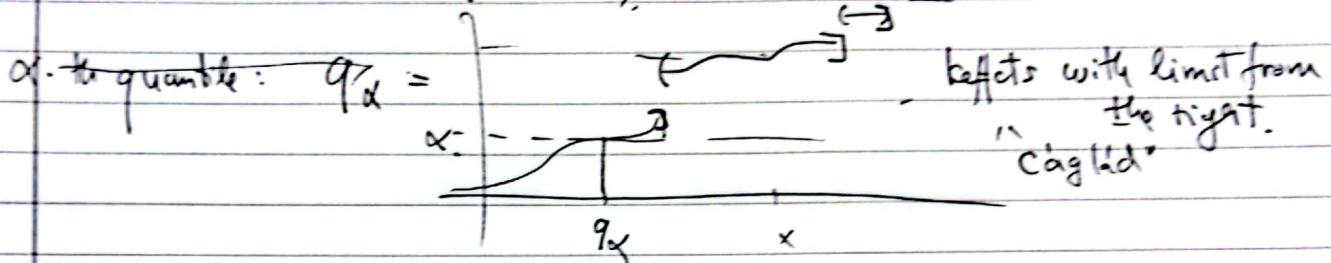
Conclusion : if  $\pi \sim \text{Unif(Surj.)}$ , &  $f(\cdot, \bar{z}^*)$  is any fn.

then  $\Delta(y_i, f(x_i, \bar{z}^*)), i \in [n]$  are exch.

Alternative presentation: in terms of quantiles.

Quantile: Let  $F$  be a cumulative distribution function (cdf)  $\alpha \in \mathbb{R}$ .

$$F(x) = P(X \leq x) \text{ for some } x.$$



$$q_\alpha = F^{-1}(\alpha) = \inf \{x : F(x) \geq \alpha\}.$$

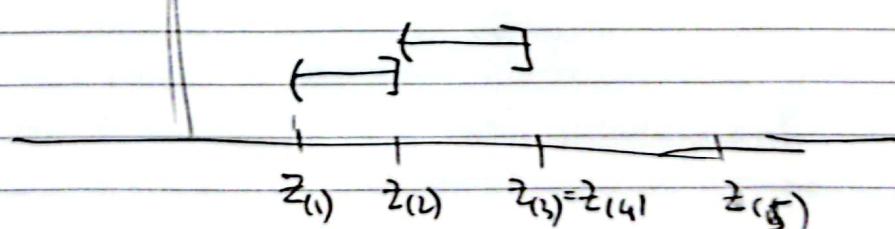
two cases: if there is a value  $\tilde{q}_\alpha$  st.  $F(\tilde{q}_\alpha) = \alpha$ .

then  $q_\alpha$  is the smallest such value.  
or else  $F(q_\alpha) > \alpha$ . but  $\lim F(p) < \alpha$

iff.  $F$  has a point mass at  $q_\alpha$ , i.e.  $q_\alpha$  is a point in range of  $F$ .

spec example?

$$z_1, \dots, z_n \rightarrow \text{empirical cdf} : F_n(x) = \frac{\#\{z_i \leq x\}}{n}.$$



$$\text{turns out: } q_\alpha = z_{(\lceil n\alpha \rceil)}$$

Now remember this:

$$\alpha \in (0, 1) \rightarrow F(x) \geq \alpha \Leftrightarrow F_n(x) \geq \frac{\lceil n\alpha \rceil}{n} \Leftrightarrow x \geq z_{(1)} \Rightarrow q_\alpha = x = z_{(1)}$$

back to CI:  $\text{rank}(\mathbf{z}, h\mathbf{z}_1, \dots, h\mathbf{z}_n) \geq L(u_{t+1})\alpha$

$$( \Rightarrow ) \quad \mathbf{z} \geq q_{\beta}(\{\mathbf{z}\}), \quad \beta = \frac{L(u_{t+1})\alpha}{n}$$

slight abuse of  
notation.

$$( \Rightarrow ) \quad \text{rk } (\mathbf{z}, \mathbf{z}^+) \geq L(u_{t+1})\alpha + 1 \quad ( \Rightarrow ) \quad \mathbf{z} \geq q_{\beta'}(h\mathbf{z}, \mathbf{z}^+), \quad \beta' = \frac{(u_{t+1})\alpha + 1}{n+1}$$

- Computational burden: need to re-train Alg 1 for all  $y \in \mathcal{Y}$ .

$\downarrow$   
big drawback.

- efficient for certain special ML models
  - ridge & lasso regularization.

Examples:  $\leftarrow$  due to comput. burden, ~~not~~ not easy to find good example  
[datir]

Alternative perspective, vi. Statistical decision theory. (Takeuchi '75+)

- Decision theory: general unifying framework in statistics
  - usually for estimation problems (i.e. fixed param)
  - loss, risk  $\rightarrow$  admissible, minimax optimal etc.
- Extension to prediction [of a random quantity]  
by Takeuchi & collaborators.

$$\mathbb{E}_{(x,y) \sim P_\theta} \sigma(x,y)$$

$$P_{(x,y) \sim P_\theta, Y} (y \in T_u(x)) \geq 1 - \alpha$$

$(X, Y) \sim P_\theta, \theta \in \Theta$

randomized pred. rul'n

$$T_u(x) = \{y \mid f(x, y) > u\}, \quad f: X \times Y \rightarrow [0, 1]$$

$$\text{marginal cov: } \forall \theta: P_{(x,y) \sim P_\theta, Y} \cup \cup \cup [0, 1].$$

$T_u$  is a locally best prediction region at some  $\theta_0 \in \Theta$ ,

for some loss  $L: Y \rightarrow [-\infty, \infty]$ , marginally valid  
& Lebesgue measure  $\mu$  over  $Y$ . if it minimizes  $\int_L(y) d\mu(y)$ .

$$\begin{aligned} \mathbb{E}_{X \sim P_{\theta_0}, Y} \int_Y I(y \in T_u(x)) L(y) d\mu(y) \\ = \int_Y \mathbb{E}_{X \sim P_{\theta_0}} \sigma(X, y) L(y) d\mu(y). \end{aligned}$$

~~Exact coverage~~:  $P_{\theta_0}(y \in T_u(x)) = 1 - \alpha \quad \forall \theta \in \Theta$

Suppose ( $\exists$ ) complete sufficient statistic  $W = w(x, y)$

$$\Rightarrow \mathbb{E}_\theta W = 0, \quad \forall \theta \in \Theta \Rightarrow W = 0 \text{ - a.s. } \forall \theta$$

$\Rightarrow \mathbb{P}_{(x,y)|W} (y \in T_u(x))$  does not depend on  $\theta$ .  
law of

example:  $Z_1, \dots, Z_n$  iid  $\sim f$  - density  $R$ .

$W = \{Z_1, \dots, Z_n\}$  = order statistics - complete suff.

$Z|W = \text{rank}_S$  - dist unif over  $S_n$ .

Then: Exact cov ( $\Rightarrow P_{\theta_0}(y \in T_u(x) | W = w) = 1 - \alpha$ )

↓ doesn't depend on  $w$ .

~~map~~  ~~$\phi = \phi(x)$~~ ,

comp. suff.

Suppose:  $\exists$  1-1 map  $(x, y) \mapsto (x, v, w)$

st.  $y = \phi(x, v, w)$  has jacobian  $|\mathcal{J}| = |\mathcal{A}|(x, v, w)$

Theorem (Takeuchi): The locally best prediction reg. is.

$$f(x, y) = \begin{cases} 1 & : P_\theta(x, v | w) > c_w \cdot P_{f_0}(x) \cdot L(y)(\mathcal{J}) \\ \dots & \\ 0 & \end{cases} = \begin{cases} & \\ & \leftarrow \text{fulfill} \\ & \text{constraint.} \end{cases}$$

Pf: Follows from Neyman-Pearson lemma. w/

"null density",  $\sim P_{f_0}(x) L(y) \sim P_{f_0}(x) L(\phi(x, v, w) | \mathcal{J})$

"alternative":  $\sim$  density of  $(x, y)$  give  $w = w$ .  
 $P_\theta(x, v | w)$ .

reject null ( $\Rightarrow y \in \text{set}$ ).

Example:  $Z_1, \dots, Z_{n+1}$  iid f. : &  ~~$\mathcal{A}$~~

Exact. coverage  $P\left(Z_{n+1} \in \mathcal{A}(Z) \mid \{Z_1, \dots, Z_{n+1}\}\right) = 1-\alpha$

$$\frac{1}{n+1} \sum_{i=1}^{n+1} F(Z_i \mid \{Z_1, \dots, Z_{i-1}\}) = 1-\alpha.$$

Lec 5

$$Z \in \mathcal{L}$$

\*  $Z \sim P_\theta, \theta \in \Theta$  - statistical model

\* a statistic  $W: \mathcal{Z} \rightarrow \mathbb{R}$  is complete (for this model) if for all  $\psi: \mathcal{Z} \rightarrow \mathbb{R}$

if  $\mathbb{E}_\theta^\psi(W) = 0$ . Then  $\forall \theta \Rightarrow \psi(W) = 0$ .  $P_\theta$  a.s.  $\forall \theta \in \Theta$

next, introduce suffi. (P-3) see "Theory of Point Estimation" book, p. 42

→ examples & non-examples

(A)  $Z_1, Z_2 \sim N(\theta, 1)$ ,  $\theta \in \mathbb{R}$ .

(1)  $W = Z_1 - Z_2 \rightarrow$  not complete  
 $\mathbb{E}_\theta W = 0$ .  $\nexists \theta$ , but  $W \neq 0$

suffi?

x

(2)  $W = Z_1 \rightarrow$

$$\mathbb{E}_\theta f(w) = 0 \quad \forall \theta$$

$$\mathbb{E}_\theta f(w + \theta) = 0 \quad \forall \theta$$

$$\int f(w) \varphi(w + \theta) dw = 0 \quad . \quad d/d\theta$$

$$\int f(w) \varphi'(w + \theta) dw = 0 \quad \vdots \quad \forall \theta$$

$$\int f(w) \varphi^{(k)}(w + \theta) dw = 0 \quad \vdots \quad \forall \theta$$

$\{\varphi^{(k)}\}_{k \geq 0}$  form basis for  $L^2(\mathbb{R}) \rightarrow f = 0$  in  $L^2(\mathbb{R})$

$\Rightarrow f = 0$  wrt. Lebesgue measure

2

$$(3) \quad W = \frac{Z_1 + Z_2}{2}$$

$W \sim N\left(\mu, \frac{1}{2}\right) \Rightarrow$  follows from (2). ✓

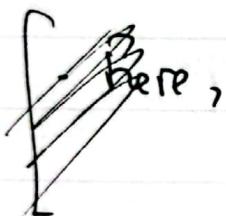
suffi

$$(B) \quad Z_1, Z_2, Z_3 \sim U[\theta, \theta+1], \quad \theta \in \mathbb{R}.$$

$W = Z_{(1)}, Z_{(2)}$  not complete,

$$\therefore E_{\theta}(W = \frac{1}{3}) = 0. \neq \theta.$$

✓



$$(C) \quad Z_1, \dots, Z_n \stackrel{iid}{\sim} f, \text{ where } f \text{ any density on } \mathbb{R}. \quad \checkmark$$

$$W = \{Z_1, \dots, Z_n\} = \mathbb{R}$$

so complete

construct rich parametric subfamily.

$$\propto \exp\left(-\sum_{i=1}^n \theta_i \left(\sum_i z_i^i\right) - \sum_i z_i^{2n}\right)$$



\* A statistic  $W : \mathbb{Z} \rightarrow \mathbb{Z}'$  is sufficient. (for  $\mathbb{Z}$  in this model)

if the distribution of  $\mathbb{Z} | W = w$   
does not depend on  $\theta$ ,  $w$ -a.s.  $P_\theta$ -a.s.  
"all info about  $\theta$  is in  $W$ "

- Sufficient & complete

"has info about  $\theta$ , & no ancillary info"

Neyman-Pearson Lemma (30's)

- Given two prob' densities  $P_0, P_1$ ,

an ~~fixed~~ optimal test :  $\phi : \mathbb{Z} \rightarrow [0, 1]$ .

of the null hypothesis

$$H_0 : \mathbb{Z} \sim P_0.$$

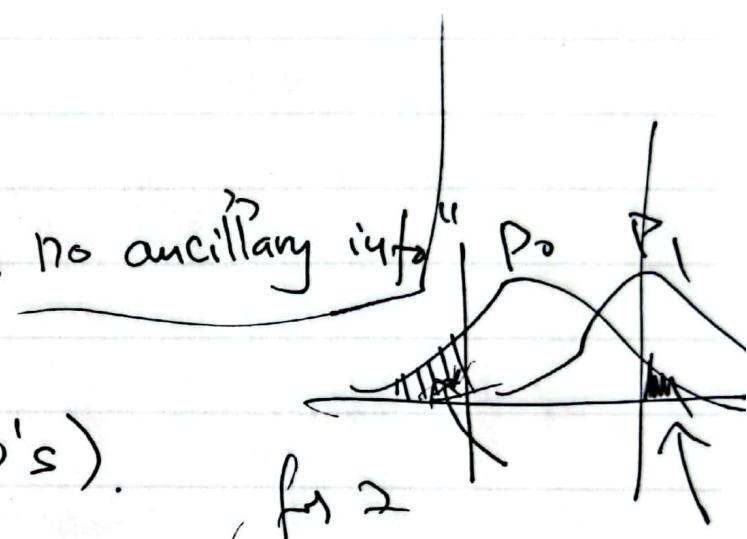
of level  $\alpha$  :

$$\mathbb{E}_{\mathbb{Z} \sim P_0} \phi \leq \alpha$$

Maximizing the power

$$\mathbb{E}_{\mathbb{Z} \sim P_1} \phi$$

$$P_{H_1}(\text{reject})$$



- interpreted as,  
given  $\mathbb{Z}$ ,  
reject  $H_0$  w/  $\phi(z)$

$$\Rightarrow P_{H_1}(\text{reject}) \leq \alpha$$

has the form.

$$\phi(z) = \begin{cases} 1 & P_1(z) > c \cdot P_0(z) \\ \tilde{\gamma} & = \\ 0 & P_1(z) < c \cdot P_0(z). \end{cases}$$

~~either  
for  
elsewhere~~

where  $c$  is the unique constant.

s.t.  $\exists \tilde{\gamma}$  s.t.  $E_{P_0} \phi = \alpha$ .

Details of Takeuchi's Theorem:

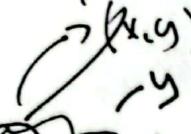
constraint:  $\int_{x,y} f(x,y) \underbrace{P(x,y|w)}_{\text{null}} = 1 - \alpha.$

max:  $\int_{x,y} f(x,y) \underbrace{P(x|w) L(y)}_{\text{alternative}} d\pi(y).$

let,  $f = f_{\text{alt}}$ .

$f = 1$  if  $\text{alt} > \text{null}$  & change variables  
density density D.

E.g.)



-  $z_1, \dots, z_{n+1}$  -  $\rho$  density on R

-  $w(z) = [z_1, \dots, z_{n+1}]$  - complete; suff.

-  $Z^+ \leftrightarrow (w, R)$  one-to-one.  $\begin{cases} \text{a.s. } \\ \text{assume in} \end{cases}$  distinct. ranks. everything below.

-  $Z^+ | w \sim \text{uniform over } W$ .

- exact coverage:  $P(\forall j \in T(x) \mid w=w) = 1-\alpha$

$P(z_{n+1} \in T(Z) \mid w=w) = 1-\alpha$ .

$$= \frac{1}{(n+1)!} \cdot \underbrace{\sum_{\substack{w \\ Z}}}_{\text{ways}} \underbrace{\sum_{\substack{S \\ \text{subset}}} \sum_{\substack{I \\ \text{subset}}} \sum_{\substack{J \\ \text{subset}}} \sum_{\substack{K \\ \text{subset}}} \dots}_{\text{ways}}$$

$$T(Z) = \left\{ z_{n+1} \mid \begin{array}{l} + (z_{n+1}, z_{-n+1}) > z_i \\ + (z_{n+1}, z_i) > z_{-n+1} \end{array} \right\},$$

$$\begin{aligned} & \frac{1}{(n+1)!} \sum_{I \in S_{n+1}} P(z_{\bar{I}(n+1)} \in T(Z^+_{-I(n+1)})) \quad \left( \text{ways to choose } I \right) \\ & \Rightarrow \frac{1}{(n+1)!} \sum_{I \in S_{n+1}} f(I_{\bar{I}(n+1)}; Z^+_{-I(n+1)}) \end{aligned}$$

Suppose  $f$  is symmetric in last  $n$  vars

[or, average it above]

$$= \frac{1}{n+1} \sum_{i=1}^{n+1} f(Z_i; Z^+_{-i}) = 1-\alpha.$$



can make it depend on  $\zeta^+$

$$= \frac{1}{n+1} \sum_{i=1}^{n+1} f(z_i, \zeta^+) = 1 - \alpha. \quad (\star)$$

- optimal prod region w/  $1-\alpha$  coverage, if for some  $z_1, \dots, z_{n+1}$  iid  $P_0$ :

$f_j(\text{iven } w = \zeta^+)$

$P_{\theta_0}$ . & for some  $f(\cdot, z^+)$  satisfying  $(\star)$ .

$f = 1$  iff.

$$P_{\theta}(x, v | w) > c_w P_{\theta_0}(x) L(v) \quad \begin{cases} & \\ & \end{cases}$$

$\underbrace{\zeta^+}_{\text{choose const.}}$

$\underbrace{P_0(z_1) \dots P_0(z_{n+1})}_{\text{const.}}$

unif. over all ranks,  
does not dep. on  $\zeta^+$

$$= C \quad >$$

$$(2) \quad P_0(z_1) \dots P_0(z_{n+1}) < C (\zeta^+)$$

$$\Leftrightarrow P_0(z_{n+1}) > c. (\zeta^+)$$

Thus: optimal prod region:

$$T_n(z_{n+1}) = \overline{\{P_0(z_{n+1}) > c(\zeta^+)\}}$$

$$\text{s.t.} \quad + \overline{\{P_0(z_{n+1}) = c(\zeta^+)\}}$$

7.

$$f(z_{n+1}, \zeta z) = \begin{cases} 1, & P_0(z_{n+1}) > c(\zeta z) \\ \gamma, & \\ 0, & \end{cases} = \begin{cases} <, & \\ & \end{cases}$$

$$\text{s.t. } \frac{1}{n+1} \sum_{i=1}^{n+1} f(z_{n+1}, \zeta z) = 1 - \alpha.$$

$$\Rightarrow \text{Sort } P_0(z_{i_1}) \leq P_0(z_{i_2}) \leq \dots \leq P_0(z_{i_{n+1}})$$

$$0 \dots 0 + \underbrace{\gamma + \dots + \gamma}_{\text{if equal}} + 1 + \dots + 1 = \underbrace{1}_{\leq (n+1)(1-\alpha)}$$

$\Leftrightarrow$  Conformal inference w.l. Conformity measure  $P_0$ .

Conclusion :  $\bar{C}_I$  is opti

when  $z_i \sim \text{iid } p$ .  $p$  density on  $\mathbb{R}$ .

- among all exactly valid pred. sets.

$$P(z_{n+1} \in T_u(z)) = 1 - \alpha \quad \text{if } p.$$

$$T_u(z) = \{ z_{n+1} \mid z + (z_{n+1}(z)) \geq u \}.$$

- $\bar{C}_I$  minimizes average ~~length~~ mis-class prob.

$$\int_{z_{n+1}} P(z_{n+1} \notin T_u(z)) dz_{n+1}, \quad \text{under } p = P_0$$

with. Conformity measure  
 $P_0(z)$

A lot to explore

- other loss fns  $L$  ?
- predicting multiple outputs.  
e.g.,  $z_{n+1}, z_{n+2}$ ? shrinkage
- supervised setting ?
- optimality against more than a fixed density ?
- implicit optimality ~~for~~ in ML & existing methods ?
- minimax optimality ?

Split (inductive) CI [Vark'02, Papadopoulos et al'02]

$$\mathcal{Z} = (\underbrace{\mathcal{Z}_1, \dots, \mathcal{Z}_{n_1}}_{\mathcal{Z}_T}, \underbrace{\mathcal{Z}_{n_1+1}, \dots, \mathcal{Z}_{n_0}}_{\mathcal{Z}_C})^T$$

"training set"

"Calibration set!"

- "proper training set"  $\rightarrow$  indep. of  $\mathcal{Z}$ ,  
use it to train  $f_M : \mathcal{Z}^{n_1+1} \rightarrow \mathbb{R}$

- construct

$$M = \left( m(\mathcal{Z}_{n_1+1}, \mathcal{Z}_T), m(\mathcal{Z}_{n_1+2}, \mathcal{Z}_T), \dots, m(\mathcal{Z}_{n_0}, \mathcal{Z}_T), m(\mathcal{Z}_T, \mathcal{Z}_T) \right)$$

Claim for a conditional on  $\mathcal{Z}_T$ .

if  $\mathcal{Z}_C \dots (\mathcal{Z}_C, \mathcal{Z}_T)$  exch,

then, entries of  $M$  are exch.

- so this also holds unconditionally over  $\mathcal{Z}_T$

if  $\mathcal{Z}_T \perp\!\!\!\perp \mathcal{Z}_C$  (doesn't have to be iid).

- So, can. use 1-D rank construction

for  $M_t$  :  $\text{rank}(m(\mathcal{Z}_t, \mathcal{Z}_T); M_t)$

- no need for retraining

$$\begin{aligned} & \Rightarrow \left[ \underbrace{(n_0 - n)}_{\text{(or. } \leq \Gamma n^{\frac{n}{(1-\alpha)}})} + 1 \right] \alpha \end{aligned}$$

good if small  $\rightarrow \langle \cdot \rangle$

10



\* e.g.  $z = (x, y)$ ,  $M(z, z_T) = |y - f(x, z_T)|$

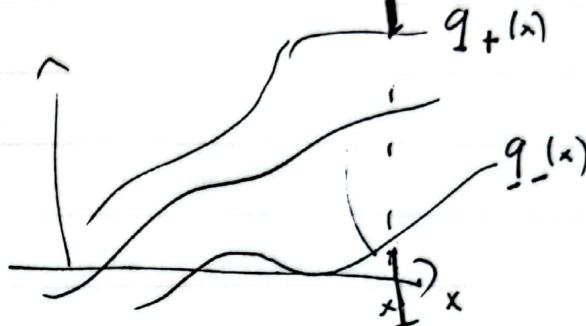
any predictor of  $y$

Trained on  $z_T$ , no sym. needed

Ramanu, Patterson,  
Candes, 2019

examples quantile reg.

\* e.g.  $z = (x, y)$



$$M(z, z_T) = \max \{ q_+(x, z_T) - y, y - q_-(x, z_T) \}$$

→

$$-q_- = q_+ = q \rightarrow |q(x, z_T) - y|$$

-  $q_- \leq q_+$  → "distance" for to region  $[q_-, q_+]$

→ negative inside the region

minimized at  $y = \frac{q_- + q_+}{2}$

pred neg:  $M(z_{t+1}, z_T) \leq \varphi_F(M(z_t, z_T))$

$$\Gamma_K \leq \overline{\Gamma}(n\tau)(\Gamma_K) \quad (n+1)(1-\alpha) = n\beta$$

$$q_\beta = \mathbb{E}(\Gamma_{n\beta})$$

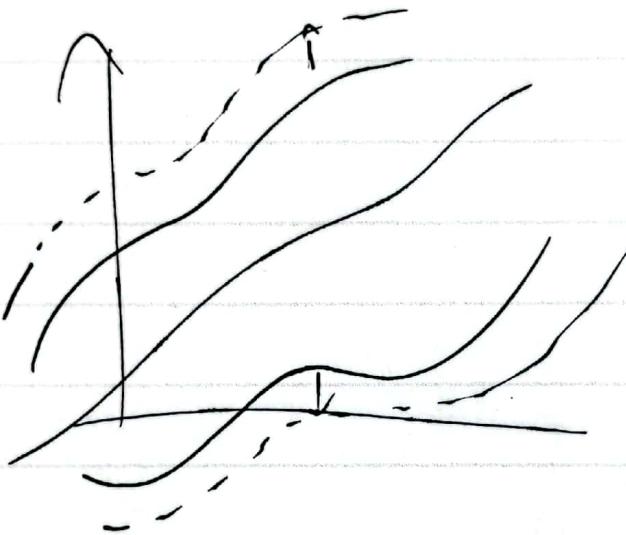
$$\beta = (1 + \frac{1}{n})(1 - \alpha)$$

$$q_-(x) - y \leq q_\beta$$

$$y - q_+(x) \leq q_\beta$$

$$y_t \in [q_-(x_t) - q_\beta, q_+(x_t) + q_\beta].$$

$$f_t \in [q_-(x_t) - q_\beta, q_+(x_t) + q_\beta]$$



→ Look at their paper for experiments.

## STAT 991. Lect.

- Recall Inductive Conformal Prediction

- $Z = \{z_1, \dots, z_{n_1}, z_{n_1+1}, \dots, z_{n_r}\}$

- $M = (m(z_{n_1+1}, z_T), \dots, m(z_{n_r}, z_T))$ .

- $T(z, Z) = \{z_t \mid \text{rank}(z, M) \geq \lfloor \frac{(n+1)(1-\alpha)}{\Gamma(n+1)(\ln n)} \rfloor\}$ .

Claim: • Condition for any fixed  $z_T = z_T$ .

• if  $(z_c, z_t)$  exch., then  $P(z_t \in T(z)) > 1 - \alpha$ .

- example .. quantile reg  $m(z, Z_T) = \max \{q_{-}(x, Z_T)^{-y}, 1 - q_{+}(x, Z_T)\}$   
[Roman et al '19] → see paper for details

- Beyond marginal validity = - standard guarantee.  $P(Y \in C(x)) \geq 1 - \alpha$
- conditional guarantee  $P(Y \in C(x) | A(x) = a) \geq 1 - \alpha$   
 $\rightarrow$  almost always better.

"Cluster-conditionals" / Mondrian CP [Vovk et al '03]



- K pre-defined clusters
- $Z[k]$ ,  $k \in [K] := \{1, \dots, k\}$  - data inclus. k.
- Mondrian score:  $A : \mathbb{Z}^4 \rightarrow \mathbb{R}^4$ 
  - permutation equivariant. within each cluster.



$$\text{# perm } \bar{\pi}_k \text{ of } Z[k] \quad A(Z_1, \dots, Z_k)_{\bar{\pi}_k} = A(Z)_{\bar{\pi}_k}^{(k)}$$



- Algo: CI in the cluster of  $Z_t = (x_t, y_t)$  has  $\frac{1}{t}$ -cluster validity

- Systematic approach [Vovk '13], validity conditional on
  - Training set  $Z$ ,  $\rightarrow$  a.k.a PAC & connected to feature regions
  - Test features  $X_t$   $\rightarrow$  mostly impossibility in finite sample
  - Test labels  $y_t$  some easy results Lei & Wasserman

} Apply CI per class "Mondrian"
 

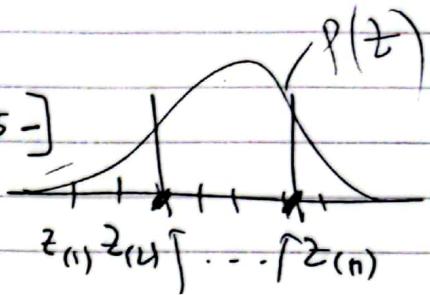
- Vovk, Lei

 $\rightarrow$  confidence-interval / test duality
 

- Guan & Tibshirani.

## Training set cond. validity

- Tolerance regions [Wilks '41, Wald '43, Tukey '45 -]



$Z = (z_1, \dots, z_n) \sim i.i.d. P.$  density on  $\mathbb{R}^n$ .

find  $T$ . s.t.

:  $(\mathbb{R}^n + \mathbb{R})$

needs id,  
each not enough  
 $z_t \sim P$ .

$$\mathbb{E}_Z \underbrace{P_{Z_t}(z_t \in T(z))}_{\text{P}} > 1 - \alpha$$

$$\mathbb{P}_Z \left[ \underbrace{P_{Z_t}(z_t \in T(z))}_{\text{P}} \in [b, c] \right] > 1 - \alpha$$

- Claim :  $P_{Z_t}(z_{(r)} \leq z_t \leq z_{(s)}) \sim \text{Beta}(s-r, n+1-(s-r))$

if  $p_i < \frac{1}{2}$ .

continuous

Proof sketch

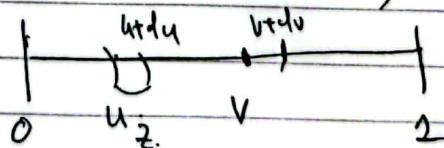
- Statement unchanged if transform  $\tilde{z}_i := F(z_i)$ ,  $F$  - cdf of  $f$ .  
( $\because F$  is strictly increasing)  
thus, assume wlog.  $z_i \sim U[0, 1]$ .

$$\text{then} : P_{Z_t}(z_{(r)} \leq z_t \leq z_{(s)}) = z_{(s)} - z_{(r)}$$

- Fact : joint distribution of  $(z_{(r)}, z_{(s)})$  has pdf.

$$\propto u^{r-1} (v-u)^{s-r-1} (1-v)^{n-s}.$$

Heuristically,



$$P(z_{(r)} \in [u, u+du], z_{(s)} \in [v, v+dv]) \sim u^{r-1} dv (v-u-du)^{s-r-1} du / (1-v-du)^n$$

Multinomial pmf

- \* obtain marginal density of  $z_{(s)} - z_{(r)}$  by integration, get desired claim

$$Q \sim \text{Beta}(s-r, n+r-s)$$

If  $T(z) = [z_{(r)}, z_{(s)}]$ , need to choose  $r, s$  such that

$$\mathbb{E} \text{Beta}(\dots) Q \geq a \quad (\Rightarrow \frac{s-r}{n+r} \geq a)$$

$$\mathbb{P}(Q \in [b, c]) > 1-\alpha.$$

$$\underline{z_{(r)}} \quad \overline{z_{(s)}} \quad \overline{\dots}$$

More general perspectives:

- statistically equivalent blocks (Turkey)<sup>(1944)</sup>

$$\cdot T(z) = \bigcup_i [z_{(\tau_{2i-1})}, z_{(\tau_{2i})}], \text{ for any non-decreasing sequence } (\tau_i)_{i \geq 1}$$

$$\text{Same coverage as } [z_{(r)}, z_{(s)}] : b-a = \sum \tau_{2i} - \sum \tau_{2i-1}$$

 ~~Sequential~~ ~~(sequential)~~  $\rightarrow$  construct sequence of total stat. equiv. blocks  
[Turkey, Fraser]  $\leftarrow$  assume all values distinct

- start w/  $g_1: \mathcal{Z} \rightarrow \mathbb{R}$ ,
- iteratively find  $i_1 = \arg \max_i g_1(z_i)$
- construct function  $g_j$ , dep. possibly depending on  $\{z_{i_1}, \dots, z_{i_{j-1}}\}$
- $i_j = \arg \max_{i \notin \{i_1, \dots, i_{j-1}\}} g_j(z_i)$



"partition space

according to the functions  $(g_i)$

then

$$\text{let } A_{ij} = \{z : g_j(z) > g_i(z_{ij})\}$$

$$B_{ij} = (\mathcal{Z} \setminus \bigcup_{i' < j} A_{i'}) \cap A_{ij}$$

stat. equiv. block

assumes iid, stronger than excl

but more flexible than standard see Conformal.

(stronger assumption  $\Rightarrow$  more flexibility)

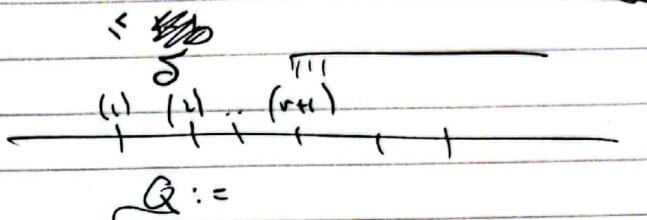
Copper-Pearson - CI

[Patra et al' 2020]

Equivalence of Tolerance Regions, Inductive CP, PAC Prod. Sets

For Aim :  $P_Z \{ P_{Z_t} (Z_t \in T(z)) \geq 1-\delta \} \geq 1-\alpha$

if



$T(z) = [z_{(r+1)}, +\infty)$ .

Wilks :  $P_{Z_t} (Z_t \in T(z)) = P_{Z_t} (Z_{(r+1)} \leq Z_t) \sim \text{Beta}(n-r, r+1)$

Let  $F_{\text{Beta}(u, v)}$  be cdf of  $\text{Beta}(u, v)$

need  $P(Q z_{(r+1)}) \geq 1-\alpha$

$P(Q \leq z_{(r+1)}) \leq \alpha$   $F_{\text{Beta}(n-r, r+1)}(z_{(r+1)}) \leq \alpha$ . (x)

Claim :  $F_{\text{Binomial}(n, p)}(r) = F_{\text{Beta}(n-r, r+1)}(u(r))$

check by taking derivative w.r.t.  $P$ .

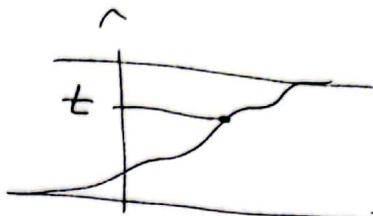
(\*)  $\Leftrightarrow F_{\text{Binom}}(n, p)(r) \leq \alpha$  (\*\*)

Given,  $(n, r)$ , find  $u(r) = u(r)$ . s.t.  $(**)$  holds (uds)

This  $u(r)$  is known as the Copper-Pearson upper confidence bound for  $P_0$

$P_0$ , when  $R \sim \text{Bin}(n, p_0)$ .

Claim : we have  $P_R \{ P_0 \in [0, u(R)] \} \geq 1-\alpha$



Fact : If R.V.  $R$  w/ cdf  $F$ .  $P(F(R) \leq t) \leq t$ ,  $\forall t \in \mathbb{R}$ .

Goal :  $P(p_0 \in [0, u(\cdot)]) > 1 - \alpha$

$$P_0 < U(R)$$

but  $p \rightarrow F(r, p)$  is

$$\left( \leftarrow \right) F_{\text{Bin}(n, p)}(R) > F_{\text{Bin}(n, p)}(r) = \alpha. \quad \begin{cases} r \leq n-1 \\ r = n, \text{constant} \end{cases}$$

$$P(F(R) > \alpha) = 1 - P(F(R) \leq \alpha) > 1 - \alpha \quad \square$$

Note: Clopper-Pearson CI relies on general principle of "inverting the CDF" for constructing CIs.

Conclusion : One-Sided Tolerance region  $\Leftrightarrow$  Clopper-Pearson CI.

Alternative way to present:

Inductive Confamal prediction produces PAC confidence sets. [Vovk '13]

$$T(z) = \left\{ z : f(z) \geq \frac{q_{t+1}}{n} (f(z_i), i \in [n]) \right\}$$

$$\hookrightarrow \text{ICP w/ CM } f = \{z : \text{rank}(f(z), \{f(z_i)\}_{i \in [n]}) \geq t+1\}$$

We have shown that if  $q > F_{\text{Bin}(n, \delta)}(r)$ , then

$$\text{if } (z, z_t) \text{ are iid}, \quad P_z \{ P_{z_t} \{ z_t \in T(z) \} \geq 1 - \delta \} \geq 1 - \alpha. \quad \square$$

[This follows b/c.  $T(z) = [f(z)]_{(r+1), \infty} =$  our previous construction  
& (\*\*\*) holds]

[so, Vovk's result already follows from Wilks' 1941 result]

Another perspective

Nested pred. sets / "PAC pred sets". [Park et al '20]

- Consider  $(T_\tau), \tau \in \mathbb{R}$ ;  $T_\tau : \mathcal{X} \rightarrow \mathcal{Y}$
- $T_\tau(z) = \{z : f(z) > \tau\}$  — nested:  $T_{\tau_1} \subset T_{\tau_2}, \text{ if } \tau_1 > \tau_2$
- $R_\tau = \sum_{i=1}^n I(z_i \notin T_\tau) = \#\{i : f(x_{i+1}) < \tau\}$
- choose  $\hat{\tau}$  s.t.  $F_{\text{Beta}(n, \delta)}(R_{\hat{\tau}}) \leq \alpha$ . [+]
- Claim:  $P_Z \left\{ P_{Z_i} \left\{ z_i \in T_{\hat{f}(z_i)} \right\} \geq k - \delta \right\} \geq 1 - \alpha$ .

Equivalence to 1-sided tolerance region, etc:

- Let  $r$  be the largest integer s.t.  $F_{\text{Beta}(n, \delta)}(r) \leq \alpha$ .  
Solve [+] s.t.  $\hat{\tau}$  is as large as possible. Then

$$R_{\hat{\tau}} = r = \#\{i : f(z_i) < \hat{\tau}\}$$

- Let  $f_{(i)}$  be order stats. of  $\{f(x_{i+1})\}$

then  $f_{(r)} < \hat{\tau} < f_{(r+1)}$

$\hat{\tau}$  maximal  $\Rightarrow \hat{\tau} = f_{(r+1)}$

- $T_{\hat{\tau}} = \{z : f(z) > f_{(r+1)}\} \Rightarrow$  same as Wilks's construction.

Conclusion: Saw equivalence of

- 1-sided tolerance regions
- 1-sided Clopper-Pearson confidence interval
- Conditional validity of IC P
- PAC Prediction Sets

→ note: All this can be also extended directly to the supervised setting

## Feature-conditional validity

- $\mathcal{Z}_i = (x_i, y_i)$ ,  $\sim_{iid} P$ ,  $\mathcal{Z}_t = (x_t, y_t) \sim P$

for  $P$ -almost all  $x$ :  $P_{\mathcal{Z}, y_t} \{ y_t \in \bar{I}(\mathcal{Z}, x_t) \mid X_t = x \} > 1 - \alpha$ .

$$\mathcal{X} = \mathbb{R}^q, Y = \mathbb{R}$$

Thm [Levy-Wasserman 14] if  $T$  has  $1 - \alpha$ -cond. validity, for any  $P$ , then, given any  $P$  and ~~a non-atom~~ non-atom  $x$  of  $P$ .

$$P \left( \lim_{\delta \rightarrow 0} \text{ess sup}_{x: \|x - x\| \leq \delta} \mu(T(x)) = \infty \right) = 1.$$

Possible remedy: asymptotic conditional coverage

$$P\text{-a.e. } x : \sup_{x \in \text{supp}(x)} \left[ P_{\mathcal{Z}, y_t} (y_t \in \bar{I}(\mathcal{Z}, x_t)) - \alpha \right] \xrightarrow[P]{+} 0.$$

algorithm:

- partition  $X$  into hypercubes
- within each region, estimate marginal density of  $y$
- if  $p(x, y)$  smooth  $p(x, y) \approx p(y) \text{-const}$  in small reg's
- apply C.P w.l.c.m  $\hat{P}$

## Label-conditional validity

$$P(y \in T(x) | y = y) >_{1-\alpha} \forall y \in Y.$$

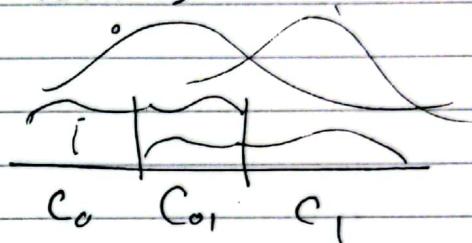
- special case of Mondrian CP [Vovk et al '03]

• in the inductive case, class-conditional marginal guarantee stated by Vovk '13

- Leili4 : -  $\mathbb{R}$ -class classification :

• min. ambiguity :

$$\begin{aligned} & \min P(C_0) \rightarrow P(-1|y=j) \\ \text{s.t. } & C_0 \cup C_1 = X, P_j(C_j) >_{1-\alpha}, \forall j=0,1 \end{aligned}$$



• using Neyman-Pearson lemma :

$$C_0 = \{x : P(y=1|x) \leq t_0\}$$

$$C_1 = \{x : P(y=0|x) \geq t_1\} \cup C_0$$

- estimation by plug-in, under margin conditions, conv. rate
- finite-sample coverage via ICP.

- Guad Tibshirani '97 : - multiclass

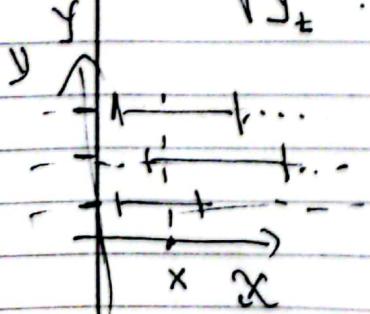
- split data over classes  $j=1, \dots, J$  for  $x_t$
- apply standard or iCP.  $\xrightarrow{\text{gt}}$

$$\forall y_t : P_{Z, x_t | y_t = y_t} (x_t \in T_{y_t}(z)) >_{1-\alpha}.$$

$$T(z, x_t) = \{y_t : x_t \in T_{y_t}(z)\}$$

[similar to construction of

confidence intervals based on tests]  
"duality" between hypothesis testing & conf. intervals



by definition:  $P \{ z, x_t | u_i = y_t \mid y_t \in T(z, x_t) \} \geq 1 - \alpha$

- Same guarantee as before , diff construction

Conformal prediction

- other topics:
  - between ICP & CP: cross-CP, jackknife +, sub-sampling
  - more general losses: ICP-type analysis via concentration inequalities
- papers to look at
  - Park et al '20: "PAC prod. sets" → good DeepNN examples.

Calibration

idea: if we predict a probability for an event, then that should be close to the actual probability of observing it.

Def:  $\rightarrow A(\text{outcomes})$  - events in a space  $X$ . e.g. weather today  
 $\rightarrow$  random. evnt. e.g. - rain  $\in \{0, 1\}$ , tomorrow - level of rain  $\in [0, \infty)$

· probability predictor / forecaster. e.g., Human expert  
 $p: X \rightarrow [0, 1]$ . predicts prob. of rain.  
 neutral net.

def:  $p$  is calibrated if.  $\forall q \in [0, 1]$

$$\Pr(A_{\{q\}} | p(x) = q) = q$$

Notes :- if  $p(x)$  can take any value, in general "unattainable"  
 need to relax

· in ML :  $A = \{f(x) = y\}$

$$p(x) = \frac{\text{top } k \text{ softmax scores (sum)}}{\text{classifi} \uparrow \text{observed class}}$$

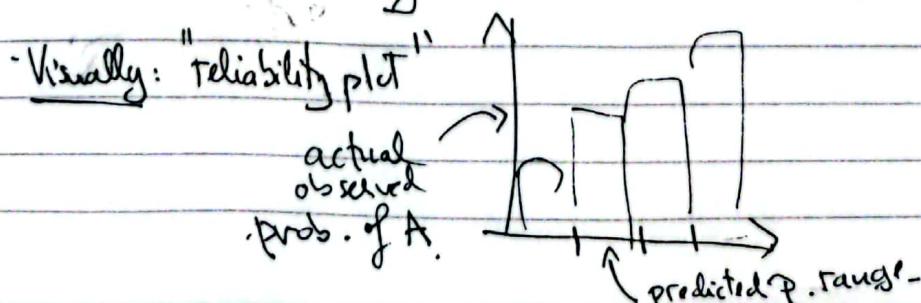
What about  $p(x)$ ? consider:

- output of typical classifier / NN.
- $s(x) \in \mathbb{R}^K, \Delta_K - K\text{-dimensional simplex}$ .
- "predicted probability" of each class  
(has nothing to do with probability)
- Prediction  $f(x) \in \arg \max_i s(x)_i$
- Confidence  $p(x) = \max_i s(x)_i$

History / Context: a rich notion; see Lichtenstein et al '80.

- 1906: Australian meteorologist Cooke
- 1958: British statistician D.R. Cox [test Bernoulli parameters]
- 1950: US meteor. Brier  $\rightarrow$  Brier Score
- 1960s: psychology, social sciences  $\rightarrow$  overconfidence
- 1980s: statistics: decision-theory, Bayesian, online re-calibration [Dawid, DeGroot, Fienberg, Foster, Volgra...]
- 1962: US meteor. R. Miller [test several batches of Bernoullis]
- 2012+: Deep learning:  $\rightarrow$  over-confidence
  - ↳ many approaches, mostly heuristic
  - ↳ guarantees: e.g. fn binning [Park et al. '20]

Studied under many names: reliability, external validity, testism of confidence, appropriateness of confidence, secondary validity, reliability



## A few vignettes/perspectives

\* Testing :  $P(A | p(x) = q) = q + \epsilon q$ .

- Suppose we have ~~be~~ completely unstructured,  $p(\cdot)$ , ~~last of values~~  $q_i$
- consider the Rvs  $\mathbb{X}_i = \mathbb{1}_{A_i} | p(x) = q_i$ .
- Then, it reduces to:
  - .  $\mathbb{X}_i \sim \text{Bernoulli}(p_i)$ ,  $p_i \in (0, 1)$   
& some dependence structure
  - Calibration null :  $p_i = q_i, \forall i$
- Simplest case,  $\mathbb{X}_i$  - indep.
- in this form, ~~is~~ a multiple hypothesis testing problem
- Miller'62 :  $X_{ij} \sim \text{Ber}(p_{ij})$ .  $\forall i$ , indep across  $j$ .
  - test  $H_0 : p_{ij} = q_{ij}, \forall i, j$
  - to test  $H_i$ :  $\underbrace{p_{ij}}_{H_i: \text{fixed } i, j} = q_{ij}$
- Chi-squared test statistic :  $T_i = \frac{\left( \sum_j X_{ij} - \sum_j p_{ij} \right)^2}{\sum_j p_{ij} (1-p_{ij})}$
- by CLT, expect that as  $\sum_j p_{ij} (1-p_{ij}) \rightarrow \infty$ .  
 $T_i \Rightarrow_d \chi^2(1)$
- to test  $H_0$  : use  $T^\top \hat{\Sigma}^{-1} T$ ,  $\chi^2(1)$  ~~not~~ null dist.
- .  $\hat{\Sigma}$  : estimated correlation matrix of entries of  $T$
- . [telep works?]

## \* Scoring rules

- $q \in [0,1]$  - prob. pred  $\Rightarrow; Q$  : r.v. of probability prediction

- $r(p) = P(A \mid p(x)=p)$

- . Brier score:  $E((p - I(A=1))^2)$  (\*) / over joint dist. of prediction  $P$  event  $A$

predicted prob.      true event.

$$\begin{aligned} &= E \left\{ (p - r(p))^2 + r(p)(1 - r(p)) \right\} \\ &\quad \text{(bias-var. decomp)} \\ &\quad \text{predicted prob.      true prob.} \\ &\quad \text{calibration} \end{aligned}$$

Sharpness  
↳  $r(p)$  "sharp" when this is small  $r(p) \approx 0, 1$ .

- more generally, for a  $K$ -class prediction problem,

bt ..

$$s : h_{1..K} \times D_K \rightarrow \mathbb{R} \text{ be a scoring rule.}$$

$s(y, p)$  - Poss when obs. class  $y \in 1..K$  & predict  $p \in P$

- let  $q$  - true pmf of  $y$ , & TIBK. [again stat. decision theory]

$$R(p, q) = E_{y \sim q} s(y, p).$$

e.g., Brier Score:  $s(y, p) = \|p_y - p\|^2 = \left\| \begin{pmatrix} 0 \\ \vdots \\ p_y \\ \vdots \\ 0 \end{pmatrix} - p \right\|^2$

in binary case,

enough to look at one coordinate.

so get  $(y - p)^2$ , as before (\*) :  $y = I(A=1)$

gen:  $\sum_y \|p_y - p\|^2 = \sum_y q(y) \left[ (1 - p(y))^2 + \|p\|^2 - p(y)^2 \right] = \|p\|^2 2(p, q) + 1 - 2p(q) + \|p\|^2 = \|p - q\|^2 - \|q\|^2 + 1$

def: proper scoring rule:  $\exists q \in \Delta_K : q \in \arg \min_{p \in \Delta_K} R(p; q)$

def. Strictly  $\rightarrow$   $q$  is unique minimizer.

e.g., Brier score is a strictly proper scoring rule.

$$\arg \min_p \left\{ \|p - q\|^2 - \|q\|^2 + 1 \right\} = p.$$

S-entropy:  $R(p, p) = \mathbb{E}_{y \sim p} s(y, p)$ .

e.g. Br. sc.:  $1 - \|p\|^2$

More examples

logarithm

- $s(y, p) = -\log p(y)$

$$R(p, q) = \mathbb{E}_{y \sim q} - \log p(y) = \underbrace{\mathbb{E}_{y \sim q} \log \left( \frac{q(y)}{p(y)} \right)}_{KL(q, p)} - \underbrace{\mathbb{E}_{y \sim q} \log(p(y))}_{H(q)}$$

is a str. prop. sc. rule.

KL div

Shannon entropy

normalized prob.

spherical. sc. rule:  $s(y, p) = -\frac{p(y)}{\|p\|}$  · strictly prop.

$$R(p, q) = \mathbb{E}_{y \sim q} S(y, p)$$

$$= \mathbb{E}_{y \sim q} \left\{ S(y, p) - S(y, q) \right\}$$

$\geq 0$  for proper Bayes rule.

$> 0$ , unless  $p = q$  for S.p.l.r.

$$+ \mathbb{E}_{y \sim q} S(y, q)$$

$R(q, q)$  · entropy

- this is just loss fn.

- stat. decision problem:

$$\min_{f \in \mathcal{F}} R(P_{f(x)}, q(y))$$

- prob dist  
w/ softmax scores f.

(could say  $\max_{p \in \mathcal{P}} \dots p(x) \dots$ )

- we prod error pushes us to 0-1. (over-conf)

↳ especially if overfit

↳ esp. if. use log loss

$$\mathbb{E}_{y \sim q} - \log p(y).$$

- conn. to pred sets / t

	. pred st. ch. ugh	here	
What?	$x \rightarrow C(x), \underset{\text{subset}}{\cancel{s}}$	$x \rightarrow \vec{q}(x) \in \Delta_{ Y -1}$ : prob. dist. array	
aim:	$y \in C(x)$	$P(y x) \approx q(x)$	

## \* Finite-sample guarantees for binning.

- Suppose  $(A_i, P)$  have a joint distribution event. e.g. ML:  $\{f(x) = y\}, P(x), (x^y)$  data prob. pred.

- Let  $\bigcup_{j=1}^k B_j = [0, 1]$  be a partition of  $[0, 1]$ .  $b_1, b_2, \dots, b_k$

- let  $P_i = P(A_i | P \in B_j)$ .

- Can act!

- Given a sample  $\mathcal{Z} = \{(A_i, P_i), i=1..n\}$

- achieve conditional coverage for  $(P_j)$ ,  $j \in [k]$ :

by constructing  $\text{CI}_j$  (e.g. Clopper-Pearson).

using  $\mathcal{Z}_j = \{(A_i, P_i) : A_i \in B_j\}$

$$P(P_j \in \text{CI}_j | \mathcal{Z}_j) \geq 1 - \alpha.$$

- Can also control ensure coverage of discretization when  $P \in B_j$ .

# of  $\hat{P}_j$ : equal to  $P_j$  w.p.  $P(P \in B_j)$

- e.g. ML:  $b(x)$  - bin of  $x$

$$c(x) = P_{(X, Y)}(f(x) = y | P(x) \in B_{b(x)})$$

- per-bin accuracy

- quantile  $P_{X^y}(c(x) \in \text{CI}(x)) \geq 1 - \alpha$ .



- Summary: - calibration

- perspectives: - testing

- scoring rules

- CIs after binning

others: re-calibration.