

$$K_2(\pi^{\leq s+1}) - K_2(\pi^{\leq s}) = \Delta_{s+1}(p_{s+1}, y_{s+1}) \leq \frac{2V^{p_{s+1}}}{\sqrt{n(\pi^{\leq s}, p_{s+1})}} \cdot (y_{s+1} - p_{s+1})$$

$$+ \frac{1}{n(\pi^{\leq s}, p_{s+1})}$$

$$\left(\frac{\sum_{t=1}^s \mathbb{1}(p_t = p_{s+1}) (y_t - p_t)}{\sqrt{n(\pi^{\leq s}, p_{s+1})}} \right)^2$$

Want to find a dist over p_{s+1} to make ∇ small $\forall y_{s+1}$

$V_S^p \geq 0$ means our preds are weakly too small
 ≤ 0 (too big)

IP $V_S^1(\pi^{\leq s}) \geq 0$ [y_t 's are $\geq p_t$ on days where $p_t = 1$]
 Let $p_{s+1} = 1$ w.p. 1

$Q^{\geq 0} \cdot (y_{s+1} - p_{s+1})$
 setting $p_{s+1} = 1$ makes this as small as possible ≤ 0

IF $V_S^0(\pi^{\leq s}) \leq 0$ [y_t 's are $\leq p_t$ on $p_t = 0$ days]
 Let $p_{s+1} = 0$ w.p. 1

$Q^{\leq 0} \cdot (y_{s+1} - p_{s+1})$
 setting $p_{s+1} = 0$ makes $y_{s+1} - p_{s+1}$ max (and $\neq 0$)
 So prod is min (≤ 0)

One dumb obs/reminder

$$V_S^p(\pi) = \sum \pi(p_t = p) (y_t - p)$$

$$n(\pi, p) = 0 \Rightarrow$$

$$\frac{\sqrt{n(\pi^{\leq s}, p)}}{\sqrt{n(\pi^{\leq s}, p)}} \leq 1$$

so $|V_S^p| \leq \sqrt{n(\pi, p)}$

$$\text{Else } (V'_S(\pi \leq s) < 0 \quad \& \quad V^0_S(\pi \leq s) > 0) \\ \exists p \in \{0, \frac{1}{m}, \frac{2}{m}, \dots, 1\} \text{ s.t. } \begin{aligned} &V^p_S(\pi \leq s) \geq 0 \\ &V^{p+\frac{1}{m}}_S(\pi \leq s) \leq 0 \end{aligned}$$

Compute q :

$$\frac{q \cdot V^p_S(\pi \leq s)}{\sqrt{n(\pi \leq s, p)}} + \frac{(1-q) V^{p+\frac{1}{m}}_S(\pi \leq s)}{\sqrt{n(\pi \leq s, p+\frac{1}{m})}} = 0$$

Play p w.p. q
 $p+\frac{1}{m}$ w.p. $1-q$

This means
the

$$\mathbb{E}[\Delta'_{s+1}(p_{s+1}, y_{s+1})] \leq q \Delta^p_S(\pi \leq s)(y_s - p) + (1-q) \Delta^{p+\frac{1}{m}}_S(\pi \leq s)(y_s - p - \frac{1}{m}) \\ = \frac{1}{q} \cdot -\frac{1}{m} \Delta^{p+\frac{1}{m}}_S(\pi \leq s)$$

Ends up $\leq \frac{2}{m}$

So, this strategy makes the 1st term $O(\frac{1}{m})$
Won't show but $\sum_{s=1}^T \frac{1}{n(\pi \leq s, p_{s+1})} \leq O(m \log(\frac{T}{m}))$

$$\text{So } \mathbb{E}[K_2] = \mathbb{E}[K_2 / T] \leq O\left(\frac{1}{m} + \frac{m}{T} \log \frac{T}{m}\right)$$

Now, we argued one can get online calibration.
It's also more meaningful than online MMC since
@ the least it won't allow y $\begin{smallmatrix} 01010101 \\ p10101010 \end{smallmatrix}$

But it still is mostly arguing that y 's & p 's
aren't too anticorrelated, still pretty global

If your predictions are constant.

Eg. $P = \bar{y}$ is calibrated.

More meaningful when your predictor
has more variability

What about asking for "good" pred on sets defined
in a way \perp of predictions?

Eg, are p 's similarly interpretable for
 $x \in \text{Tall}$ as $x \in \text{Short}$?

Also kind of weird that we haven't touched x 's
yet, other than that predictions might come from them.

• Attempt 1 : Marginal mean consistency for
(quantile) groups

2 : Calibration | groups

Let $G \subseteq 2^X$ be a collection of group fns,
 $g \in G$ is a subset of X , $g(x) \in \{0, 1\}$
 $g(x) = 1$ "x belongs to g"

Then, $\mu(g) = \Pr_{x \sim D_X}[g(x) = 1]$ is g's frequency/mass

Def f is α -approx group MMC for G if
 $\forall g \in G$

$$\mu(g) \left[\mathbb{E}[f(x) | g(x) = 1] - \mathbb{E}[f] \right]^2 \leq \alpha$$

If f isn't α -approx group MMC for G ,
 patching will decrease 2 ed error, by

$$\mu(g) \cdot \Delta^2 \geq \alpha$$

So $\frac{1}{\alpha}$ rounds will lead to termination.

This is more meaningful than MMC for
 expressive G : eg $G = 2^X$.

However, this should lead us to wondering
 how much additional data we need to satisfy
 it... to estimate MMC for each group, for
 example.

[Hint: one way to do this is to treat G like \mathcal{A} !] $|G|$ or $VCC(G)$

Let's talk now about group conditional calibration.
We call this multicalibration:

Def Fix $f: \mathcal{X} \rightarrow [0, 1]$ and group $g: \mathcal{X} \rightarrow \{0, 1\}$. Avg² cal error of f on g is

$$K_2(f, g, \mathcal{D}) = \sum_v \Pr[f(x)=v \mid g(x)=1] (v - \mathbb{E}[y \mid f(x)=v, g(x)=1])^2$$

Ideally, want α avg² ^{multi} calibration error $\forall g \in G$,
but weighted by group mass

$$\frac{K_2(f, g, \mathcal{D})}{\mu(g)} \leq \alpha.$$

well, if your f isn't α -multicalibrated wot G ,
we can patch it!

$$h(x, f; v \rightarrow v', g) = \begin{cases} v' & \text{if } f(x)=v, g(x)=1 \\ f(x) & \text{o/w} \end{cases}$$

One subtlety here. originally calibration patching
mapped values $v \in R(f)$ to $v' \in C(R(f))$,
no more than $|R(f)|$ outputs.

w this, $v \in R(f) \nleftrightarrow g \in G \Rightarrow v'$, could increase
to $|R(f)| \cdot |G|$.

is even worse since we're adding to $R(f)$

Our Brier score \downarrow is $B(\hat{f}_{t+1}) - B(f_t) = \mu_t(v_t, g_t)(v'_t - v_t)^2$

A natural fix is to define our grid of prediction values in advance,

$$m = \lceil \frac{1}{\alpha} \rceil, \quad R(f) = \sum_{i=1}^m \frac{1}{m} \dots \{ \}$$

$$\approx \mu_t(v_t, g_t)(v'_t - v_t)^2 - \frac{1}{4m^2}$$