

Lemma 6.1.3.a

Let  $\mathcal{H}$  be closed under affine transformations. Suppose  $f$  is calibrated,

and  $\exists h \in \mathcal{H}$  s.t.

$$\mathbb{E}[(h(x) - y)^2 - (f(x) - y)^2 \mid f(x) = v] \geq \alpha \Big] A$$

Then  $\exists h' \in \mathcal{H}$  s.t.

$$\mathbb{E}[h'(x)(y - v) \mid f(x) = v] \geq \frac{\alpha}{2} \Big] B$$

and also if  $B$ , then  $A$   
w. different constants.

Def 1 Fix  $\mathcal{D} \in \Delta \mathcal{Z}$  & function class  $\mathcal{H}$ . Let  $f^*(x) = \mathbb{E}[y]$   
 $y \sim \mathcal{D}(x)$

We say  $\mathcal{H}$  satisfies the weak learner cond. wrt  $\mathcal{D}$   
if  $\forall S \subset \mathcal{X}$  w.  $\Pr[x \in S] > 0$ ,

$$\mathbb{E}_{(x,y) \sim \mathcal{D}} [(f^*(x) - y)^2 | x \in S] < \min_{c \in \mathbb{R}} \mathbb{E}_{(x,y) \sim \mathcal{D}} [(c - y)^2 | x \in S]$$

then  $\exists h \in \mathcal{H}$  s.t.

$$\mathbb{E}_{(x,y) \sim \mathcal{D}} [(h(x) - y)^2] < \min_{c \in \mathbb{R}} \mathbb{E}_{(x,y) \sim \mathcal{D}} [(c - y)^2]$$

---

For any subset of  $\mathcal{X}$ , if  $f^*$  is better than the best constant,  $\exists h \in \mathcal{H}$  which is better too (though  $f^*$  might be  $\gg h$ )

---

When is it the case that any multicalibrated  $f$  wrt  $\mathcal{H}$  is Bayes optimal?

---

Theorem 34 Fix  $\mathcal{D}$ . Let  $\mathcal{H}$  be a set of fns closed under affine transforms:  $h \in \mathcal{H} \Rightarrow \alpha h(x) + b \in \mathcal{H}$

MC wrt  $\mathcal{H}$  implies Bayes-OPT over  $\mathcal{D}$  iff

$\mathcal{H}$  satisfies the weak learning condition.

[assume  $\mathcal{X}, \mathcal{Y}$  countable]

$\Rightarrow$  (Show weak learning  $\Rightarrow$  Every MC is Bayes opt.)  
 if not,  $\exists$  MC  $f$  which isn't Bayes opt.:

$$\mathbb{E}[(y - f^*(x))^2] < \mathbb{E}[(y - f(x))^2]$$

$$\Leftrightarrow \sum_v \Pr[f(x)=v] \mathbb{E}[(y - f(x))^2 | f(x)=v] - (y - f^*(x))^2 \Pr[f(x)=v] > 0$$

$$\Rightarrow \exists v \text{ s.t. } \mathbb{E}[(y - v)^2 | f(x)=v] > (y - f^*(x))^2$$

(where  $f(x)=v$ )

that is a subset of  $\mathcal{X}$ , and since  $f$  is calibrated, it's prediction is right on average there

$$\mathbb{E}[(v - y)^2 | x \in S] = \min_{c \in \mathcal{R}} \mathbb{E}[(c - y)^2 | x \in S]$$

$$\Rightarrow \mathbb{E}[(f^*(x) - y)^2 | x \in S] < \mathbb{E}[(v - y)^2 | x \in S]$$

So, weak learning implies  $\exists h \in \mathcal{H}$  s.t.  
 $\mathbb{E}[(h(x) - y)^2 | x \in S] < \min_{c \in \mathcal{R}} \mathbb{E}[(c - y)^2 | x \in S]$

(lemma)  $\Rightarrow \exists h' \in \mathcal{H}$  s.t.  $\mathbb{E}[h'(x)(y - v) | f(x)=v] > 0$   
 A violation of MC wrt  $\mathcal{H}$  ~~Contradiction~~

$\leftarrow$  (For any  $\mathcal{H}$  not satisfying WLC wrt  $\mathcal{D}$ , MC wrt  $\mathcal{H}$  &  $\mathcal{D}$  doesn't imply Bayes-OPT over  $\mathcal{D}$ .  
 In particular  $\exists f$  MC wrt  $\mathcal{H}$  &  $\mathcal{D}$  which isn't Bayes-OPT:  $\mathbb{E}[(f(x)-y)^2] > \mathbb{E}[(f^*(x)-y)^2]$

---

Since  $\mathcal{H}$  doesn't satisfy WLC over  $\mathcal{D}$ ,  $\exists S \subseteq \mathcal{X}$  w  $\Pr[x \in S] > 0$  st.

$$\min_h \mathbb{E}[(h(x)-y)^2 | x \in S] > \min_{c \in \mathbb{R}} [(c-y)^2 | x \in S]$$

Let  $c(S) = \mathbb{E}[y | x \in S]$

Define  $f(x) = \begin{cases} h^*(x) & x \notin S \\ c(S) & x \in S \end{cases}$

---

Then  $\mathbb{E}[(f(x)-y)^2]$

$$\begin{aligned}
 &= \Pr[x \in S] \mathbb{E}[(c(S)-y)^2 | x \in S] + \Pr[x \notin S] \mathbb{E}[(f^*(x)-y)^2 | x \notin S] \\
 &> \Pr[x \in S] \mathbb{E}[(f^*(x)-y)^2 | x \in S] + \Pr[x \notin S] \mathbb{E}[(f^*(x)-y)^2 | x \notin S] \\
 &= \mathbb{E}[(f^*(x)-y)^2]
 \end{aligned}$$

So  $f$  isn't Bayes OPT. Is it MC wrt  $\mathcal{H}, \mathcal{D}$ ?

Informally, it should be: it's Bayes OPT everywhere but  $S$ , and on  $S$ , it's better than the best  $h \in \mathcal{H}$ .

Suppose  $f$  isn't MC wrt  $\mathcal{A}$ , then  
 $\exists h$  and  $v \in \mathcal{R}(f)$  s.t.

$$\mathbb{E}[h(x)(y-v) \mid f(x)=v] > 0$$

So Lemma  $\Rightarrow \exists h'$ :

$$\star \mathbb{E}[(h'(x)-y)^2 - (f(x)-y)^2 \mid f(x)=v] > 0$$

Notice  $v$  must be  $c(S)$ , since  $f(x)$  is  
 Bayes OPT  
 elsewhere.

$$\mathbb{E}[(h'(x)-y)^2 \mid f(x)=c(S)]$$

$$= \underbrace{\Pr[x \in S]}_A \mathbb{E}[(h'(x)-y)^2 \mid f(x)=c(S), x \in S]$$

$$+ \Pr[x \notin S] \mathbb{E}[(h'(x)-y)^2 \mid f(x)=c(S), x \notin S]$$

$$\geq A + \Pr[x \notin S] \mathbb{E}[(f-y)^2 \mid f(x)=c(S), x \notin S]$$

Since  $f = f^*$  here

So for  $\mathcal{H}$ , it must be that

$$\mathbb{E}[(h'(x) - y)^2 | x \in S, f(x) = c(S)] \\ < \mathbb{E}[(f(x) - y)^2 | f(x) = c(S)]$$

$$\Rightarrow \mathbb{E}[(h'(x) - y)^2 | x \in S] < \mathbb{E}[(c(S) - y)^2 | x \in S]$$

A contradiction to  $\mathcal{H}$  violating the WLC.

---

So, if  $f$  is exactly MC wrt  $\mathcal{H}, \mathcal{D}$ ,  
and  $\mathcal{H}$  satisfies WLC

$$\Rightarrow f \text{ is Bayes OPT}$$

[and if  $f$  is Bayes OPT & exactly  
MC wrt  $\mathcal{H}, \mathcal{D}$ ,  $\mathcal{H}$  must satisfy  
the WLC].

Approximate variants hold too.

# Conformal prediction

Let's talk about classification!

$$\mathcal{X} = \text{---} \quad \mathcal{Y} = \{1, \dots, K\}$$

$$f: \mathcal{X} \rightarrow \mathcal{Y}$$

How good  
is a given  
prediction?

Most of our  
contemporary models  
actually produce  
 $f(x) \in \Delta \mathcal{Y}$ , "probability  
of labels"

→ A1: you could train a  
calibrated regressor to predict  
 $\Pr[\text{correct}]$

A2: you could output sets  
 $C_x$  s.t.  $y_x \in C_x$  w.p.  $\overline{1-\delta}$

[Note, these might feel like confidence intervals, that is a prediction set, but these coverage guarantees generally hold only if the world follows the parametric model class]

---

$T: \mathcal{X} \rightarrow 2^{\mathcal{Y}}$  s.t.  $\Pr[Y \in T(x)] \approx 1 - \delta$   
= High dimensional.

Idea: learn a nonconformity score fn  $s: \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$

---

We have some model  $f$

Given  $x, f(x), y$ , how surprising is  $y$ ?

[example  $s(x, y) = |y - f(x)|$  for regression]

---

Parametrizes a  $\mathbb{I}$ -d set  $T(x, \tau) =$

$$\{y : s(x, y) \leq \tau\}$$



## Weak Marginal guarantee

Given a sample  $D \sim \mathcal{D}^n$ , a calibration set  
produce  $\mathcal{T}(x)$  that have this set  
coverage guarantee: for new  $(x, y) \sim \mathcal{D}$

$$1 - \delta \leq \Pr_{\substack{D \sim \mathcal{D}^n \\ (x, y) \sim \mathcal{D}}} [y \in \mathcal{T}(x)] \leq 1 - \delta + \frac{1}{n+1}$$

[can be made w.p. 1]

Let  $\tau$  be the smallest val s.t.

$$\sum_{i=1}^n \mathbb{1}(s(x_i, y_i) \leq \tau) \geq (1 - \delta)(n+1)$$

[the empirical  $\frac{(1 - \delta)(n+1)}{n}$  quantile of  $\mathcal{D}$ ]

$$\text{Output } \mathcal{T}_D(x) = \{ \hat{y} : s(x, \hat{y}) \leq \tau \}$$

We might also want guarantees  
that aren't just marginal,

eg

Given  $G \subseteq 2^X$ , group conditional covy:

$$\Pr_{(x,y) \sim D} [y \in \tau(x) \mid g(x) = \underline{1}] = 1 - \delta$$

---

To get such guarantees is a little trickier, our coverage/pred sets will now look like

$$\tau_D^f(x) = \{ \bar{y} : s(\bar{y}, x) \leq f(x) \}$$

was  $\tau$   $\swarrow$

$\tau_D^f(x)$  will have group covy

$\Leftrightarrow f(x)$  has group conditional quantile guarantees!

[Which we can obtain]

---

Unfortunately, these  $f(x)$ 's instead of a fixed  $\tau$  add complexities.

w.p.  $\delta$ , output  $\emptyset$  or  $1 - \delta$  output  $\tau$ .

↳ This has the right coverage Prob. for every  $x, g, \dots$  but it's completely uninformative.

---

To fix, we ask for coverage to hold [group  $\phi$  value of thresh, eg  $\forall g \in G, v \in \mathcal{R}(f)$ ]

$$\Pr_{(X,Y)}[y \in \mathcal{T}_g^f(x) \mid y(x)=1, f(x)=v] = 1 - \delta$$

---

This will hold  $\Leftrightarrow f(x)$  is a multicalibrated quantile predictor for  $q = 1 - \delta$

---

This stuff is doable in online settings too  $\dots$