

## 5: Information and Statistical Inference-2

(1)

**A** How much "information" is revealed by one coin toss?

Let us revisit our coin-bias estimation problem.

Suppose that we observe independent tosses of a coin with

bias  $\frac{1}{2}$  or  $\frac{1}{2} + \varepsilon$ . We assume a Bayesian setup with a

uniform prior on two coins. In our language,  $X \sim \text{unif}\{\theta_0, \theta_1\}$

and  $Y|X \sim \text{iid } \text{Ber}(p_X)$  with  $p_0 = \frac{1}{2}$  and  $p_1 = \frac{1}{2} + \varepsilon$ .

Initial uncertainty of  $X$  is  $H(X) = 1$ . How much is this uncertainty reduced when we toss the coin, namely

how much information is revealed per coin toss?

One answer to this question can be seen from our

formula:  $P_e^*(\text{unif}) = \frac{1}{2} (1 - d(W_0^n, W_1^n))$

We can look at how fast  $d(W_0^n, W_1^n)$  increases per

sample. Specifically, how large is  $d(W_0^n, W_1^n)/n$ ?

Earlier we saw that  $\frac{d(W_0^n, W_1^n)}{n} \leq d(W_0, W_1) = \varepsilon$ .

But this bound is weak. Indeed, by Pinsker's inequality

$$\frac{d(W_0^n, W_1^n)}{n} \leq \frac{1}{n} \sqrt{\frac{\ln 2}{2} D(W_0^n \| W_1^n)}$$

$$= \frac{1}{n} \sqrt{\frac{n \ln 2}{2} D(W_0 \| W_1)} \quad (\text{additivity of divergence})$$

$$= \sqrt{\frac{\ln 2}{2n}} \log \frac{1}{1-\varepsilon^2} \approx \frac{\varepsilon}{\sqrt{n}} \text{ when } \varepsilon \text{ is small} \quad (2)$$

Thus, the "information" revealed per toss is much smaller than what was suggested by our earlier bound.

In particular, we need at least (roughly)  $\frac{1}{\varepsilon^2}$  tosses to get  $P_e^*(\text{unif}) \leq 1/3$ .

→ This notion of information is not very formal, but gives us an idea how  $D(P||Q)$  better captures information than  $d(P, Q)$ .

### B M-any hypothesis testing : Maximum Likelihood (ML) and Maximum A posteriori Probability (MAP) tests

We now move to M-any hypothesis testing. We consider a Bayesian setting where it is assumed that the unknown  $X$  is generated from a fixed distribution  $P_X$  on  $\mathcal{X} = \{1, \dots, M\}$ .

The observation  $Y$  is generated from  $W_x$  when  $X=x$ . A test now is a mapping  $g: \mathcal{Y} \rightarrow \mathcal{X}$  and its probability of error

is given by  $P(g(Y) \neq X) = \sum_x P_X(x) \sum_{y: g(y) \neq x} W(y|x)$ .

We denote this error by  $P_e(g|P_X)$ , and denote by  $P_e^*(P_X)$  the minimum of  $P_e(g|P_X)$  over  $g$ .

\*  $P_x$  is uniform

(3)

For  $|\mathcal{X}|=2$ , we saw the LLRT given by

$$\log \frac{W_0(y)}{W_1(y)} \stackrel{0}{\gtrless} 0.$$

We can reinterpret this test as

$$g(y) = \operatorname{argmax}_{x \in \{0,1\}} W_x(y) = \operatorname{argmax}_{x \in \{0,1\}} \log W_x(y).$$

This test is much more general and is called the Maximum Likelihood (ML) test:

$$\underline{g_{ML}(y) = \operatorname{argmax}_{x \in \mathcal{X}} \log W_x(y)}.$$

\* General  $P_x$

We can consider a natural generalization of the ML test where discount  $\log W_x(y)$  by  $\lambda_x \geq 0$  and use

$$g(y) = \max_{x \in \mathcal{X}} \log W_x(y) - \lambda_x.$$

This  $\lambda$  is often called a regularizer and its role is to make the test favor certain hypotheses  $\mathcal{X}_0 \subseteq \mathcal{X}$ . Indeed, regularized ML is a very popular algorithm in statistics and machine learning.

Note that  $\log W_x(y) - \lambda_x = \log 2^{-\lambda_x} W_x(y)$ , where  $2^{-\lambda_x} \in [0, 1]$ . Thus,

$$g(y) = \operatorname{argmax}_{x \in \mathcal{X}} 2^{-\lambda_x} W_x(y) = \operatorname{argmax}_{x \in \mathcal{X}} 2^{-\lambda_x} W_x(y) / \sum_x 2^{-\lambda_x}.$$

We can interpret  $2^{-\lambda_x} / \sum_{x'} 2^{-\lambda_{x'}}$  as a distribution on  $\mathcal{X}$ . (4)

A natural choice is  $P_x(x)$ , for which

$$P_x(x) W(y|x) = P_{x|y}(x|y) P_y(y).$$

Thus,

$$g(y) = \underset{x \in \mathcal{X}}{\operatorname{argmax}} P_{x|y}(x|y) P_y(y) = \underset{x \in \mathcal{X}}{\operatorname{argmax}} \underline{P_{x|y}(x|y)}.$$

Since  $P_{x|y}(x|y)$  denote the a posteriori prob. of  $X$  given  $y$ ,  
the test above is called the MAP test, denoted

$$g_{\text{MAP}}(y) = \underset{x \in \mathcal{X}}{\operatorname{argmax}} P_{x|y}(x|y).$$

In fact, this test attains  $P_e^*(P_x)$ .

Theorem (Optimality MAP)

$$P_e(g_{\text{MAP}} | P_x) = P_e^*$$

Proof. For any  $g: \mathcal{Y} \rightarrow \mathcal{X}$ ,

$$\begin{aligned} P_e(g | P_x) &= \mathbb{P}(g(Y) \neq X) = P_y(y) \mathbb{P}(X \neq g(y) | Y=y) \\ &\geq \sum_y P_y(y) \mathbb{P}(X \neq g(y) | Y=y) \\ &\geq \sum_y P_y(y) \min_g \mathbb{P}(X \neq g(y) | Y=y) \quad \begin{matrix} (\text{in fact, this} \\ \text{is an equality}) \\ \text{why?} \end{matrix} \\ &= \sum_y P_y(y) \min_{\tilde{x}} \mathbb{P}(X \neq \tilde{x} | Y=y) \end{aligned}$$

$$\begin{aligned} \text{Note that } \min_{\tilde{x}} P(X \neq \tilde{x} | Y=y) &= 1 - \max_{\tilde{x}} P(X = \tilde{x} | Y=y) \\ &= 1 - \max_{\tilde{x}} P_{X|Y}(\tilde{x}|y). \end{aligned} \quad (5)$$

Thus, the right-side of the previous inequality coincides with  $P_e(g_{MAP} | P_x)$ , giving

$$P_e(g | P_x) \geq P_e(g_{MAP} | P_x).$$

□

**C** Which distance governs the minimum probability of error for Many HT?

We saw two variants of binary hypothesis testing: for the first, the min. prob. of error was determined by  $d(W_0, W_1)$ , and for the second, the min prob. of error (of type II) was determined by  $D(W_0 || W_1)$ .

Is there a similar geometric "notion" of distance that governs  $P_e^*(P_x)$ ?

We build using an heuristic argument.

$$\begin{aligned} P_e(g | P_x) &= \sum_x P_x(x) \sum_{x' \neq x} P(g(y)=x' | X=x) \\ &\leq (M-1) \max_{x, x': x \neq x'} P(g(y)=x' | X=x) \end{aligned}$$

Consider

$$\begin{aligned} \max_{x \neq x'} P(g(y)=x' | X=x) &= 2 \max_{x \neq x'} \frac{1}{2} P(g(y)=x' | X=x) \\ &\quad + \frac{1}{2} P(g(y)=x | X=x') \end{aligned}$$

(6)

We can roughly approximate this expression as

$$1 - \min_{x \neq x'} d(W_x, W_{x'}).$$

Alternatively, we can use heuristics from Stein's lemma and approximate  $P(g(Y)=x' | X=x) \approx 2^{-D(W_x || W_{x'})}$ , saying that

$$P_e^*(P_X) \approx M 2^{-\min_{x \neq x'} D(W_x || W_{x'})}.$$

Thus:

- Some notion of how spread out  $W_x$  are in a statistical distance determines the error
- We can't distinguish more than  $\min_{x \neq x'} D(W_x || W_{x'})$  bits of uncertainty ( $\log M \leq \min_{x \neq x'} D(W_x || W_{x'})$ )

Of course, this answer is only heuristic. But it gives us an idea of how to characterize the "information revealed by  $Y$  about  $X$ ".

### D Information revealed by $Y$ about $X$

We now return to the question we have been pursuing. First, we define a notion of uncertainty remaining in  $X$  given  $Y$ .

(7)

### \* Conditional entropy of $X$ given $Y$

$H(X|Y) \equiv H(P_{XY} | P_Y)$  is defined as

$$:= \sum_y P_Y(y) H(P_{X|Y=y}),$$

namely the average entropy of the conditional distribution.

At this point, this definition appears out of the blue.

But soon we will see that this is a good definition.

Note

$$\begin{aligned}
 H(X|Y) &= \sum_y P_Y(y) \sum_x P_{X|Y}(x|y) \log \frac{1}{P_{X|Y}(x|y)} \\
 &= \sum_{y,x} P_{XY}(x,y) \log \frac{1}{P_{X|Y}(x|y)} = \mathbb{E} \left[ \log \frac{1}{P_{X|Y}(x|y)} \right] \\
 (\text{here } P_{X|Y}(\cdot|\cdot) \text{ is a fixed function of } (x,y)) \quad &\quad \swarrow \\
 &= \sum_{x,y} P_{XY}(x,y) \log \frac{P_Y(y)}{P_{XY}(x,y)} \\
 &= \sum_{x,y} P_{XY}(x,y) \log \frac{1}{P_{XY}(x,y)} \\
 &- \underbrace{\sum_{x,y} P_{XY}(x,y) \log \frac{1}{P_Y(y)}}_{\sum_y P_Y(y) \log \frac{1}{P_Y(y)}} = \underbrace{H(X,Y) - H(Y)}_{\text{uncertainty in } (X,Y) - \text{uncertainty in } H(Y)}
 \end{aligned}$$

In the previous calculation, we have used the Bayes rule several times : (8)

$$\begin{aligned} P_{X,Y}(x,y) &= P_{X|Y}(x|y) P_Y(y) \\ &= P_{Y|X}(y|x) P_X(x) \end{aligned}$$

### \* Information

Now that we have defined the notion of residual uncertainty, the conditional entropy of  $X$  given  $Y$ , we can define information revealed by  $Y$  about  $X$  as

(uncertainty of  $X$ ) - (uncertainty remaining in  $X$  after  $Y$  is revealed)

$$= H(X) - H(X|Y) =: \underline{\underline{I(X \wedge Y)}} \quad \text{mutual information b/w } X \text{ and } Y$$

$$= H(X) + H(Y) - H(X,Y)$$

$$= H(Y) - H(Y|X) \quad \text{Thus, information revealed by } Y \text{ about } X \text{ equals information revealed by } X \text{ about } Y$$

$$= H(X,Y) - H(X|Y) - H(Y|X)$$

$$\begin{aligned} \text{Note } I(X \wedge Y) &= H(X) - H(X|Y) = \mathbb{E} \left[ \log \frac{1}{P_X(x)} - \log \frac{1}{P_{X|Y}(x|y)} \right] \\ &= \mathbb{E} \left[ \log \frac{P_{X|Y}(x|y)}{P_X(x)} \right] \\ &= \mathbb{E} \left[ \log \frac{P_{X,Y}(x,y)}{P_X(x) P_Y(y)} \right] \\ &= D(P_{X,Y} || P_X P_Y) \quad \text{divergence b/w joint and independent dist.} \\ &= \mathbb{E}_Y [D(P_{X|Y} || P_X)] \end{aligned}$$

→ Given a channel  $W: \mathcal{X} \rightarrow \mathcal{Y}$  and a distribution  $P_X$  on  $\mathcal{X}$ , ⑨

let  $Y$  be the output of the channel when the input  $X \sim P_X$ .

Then,

$$\begin{aligned} I(X \wedge Y) &= \mathbb{E}_X [D(W_X \| W \circ P_X)] \\ &= \sum_x P_X(x) D(W_x \| \sum_{x'} P_X(x') W_{x'}) . \end{aligned}$$

In particular, for  $P_X = \text{Unif}(\{1, \dots, M\})$ ,

$$I(X \wedge Y) = \underbrace{\frac{1}{M} \sum_x D(W_x \| \frac{1}{M} \sum_{x'} W_{x'})}_{\text{avg. distance from the centroid}}$$

It is this avg. distance that will be seen to govern the difficulty of M-ary HT

### E Information-error bound: Fano's inequality

Fano's inequality (without proof)

For  $P_X = \text{unif}(\{1, \dots, M\})$ ,  $\rightarrow h(q) = q \log \frac{1}{q} - (1-q) \log \frac{1}{1-q}$

$$\begin{aligned} P_e^*(\text{unif}) &\geq 1 - \frac{I(X \wedge Y) + h(P_e^*(\text{unif}))}{\log M} \\ &\geq 1 - \frac{I(X \wedge Y) + 1}{\log M} \end{aligned}$$

Thus, if  $P_e^*(\text{unif}) \leq \frac{1}{3}$ ,

$$\frac{3}{2} (I(X \wedge Y) + 1) \geq \log M \quad \begin{array}{l} \text{mutual information limits} \\ \text{the number of values of } X \\ \text{that can be resolved by } Y \end{array}$$