

Lecture 7: October 18, 2017

Lecturer: Madhur Tulsiani

1 Binary hypothesis testing

In this lecture, we apply the tools developed in the past few lectures to understand the problem of distinguishing two distributions (special cases of which have been discussed in the previous lectures). This problem is also known as the hypothesis testing. Suppose we have two distributions P_0 and P_1 on a finite universe U . The “universe” chooses one of the two distributions and generates the data, which consists of a sequence $\bar{x} \in U^n$ chosen either from P_0^n or P_1^n . The true distribution is unknown to us, but we are guaranteed that once P_0 or P_1 is chosen, all n samples in the sequence \bar{x} are sampled independently from the chosen distribution. The goal is to distinguish between the following two hypotheses:

- H_0 : The true distribution is P_0 .
- H_1 : The true distribution is P_1 .

Sometimes H_0 is also referred to as the null (default) hypothesis. We will consider (deterministic) tests $T : U^n \rightarrow \{0, 1\}$, which take the sequence of samples \bar{x} as input and select one of the hypotheses. There are two types of errors we will be concerned with

$$\alpha(T) := \mathbb{P}_{\bar{x} \sim P_0^n} [T(\bar{x}) = 1] \quad (\text{False Positive})$$

$$\beta(T) := \mathbb{P}_{\bar{x} \sim P_1^n} [T(\bar{x}) = 0] \quad (\text{False Negative}).$$

The following claim is easy to prove based on the properties of total-variation distance considered earlier.

Claim 1.1 $\min_T \{\alpha(T) + \beta(T)\} = 1 - \delta_{TV}(P_0^n, P_1^n)$.

Recall that optimal test for the above claim should be of the form

$$T(\bar{x}) = \begin{cases} 1 & \text{if } P_1^n(\bar{x}) \geq P_0^n(\bar{x}) \\ 0 & \text{if } P_1^n(\bar{x}) < P_0^n(\bar{x}) \end{cases}.$$

One may ask why should we only consider the optimal tests for minimizing the sum $\alpha(T) + \beta(T)$. We may care more about a false positive than a false negative, and may want to minimize a weighted sum (or some other monotone function) of the errors. The following lemma shows that all optimal tests should be of the form above, which make a decision only based on the *ratio* $P_0^n(\bar{x})/P_1^n(\bar{x})$

Lemma 1.2 (Neyman-Pearson Lemma) *Let T be a test of the form*

$$T(\bar{x}) = \begin{cases} 1 & \text{if } P_1^n(\bar{x})/P_0^n(\bar{x}) \geq \Delta \\ 0 & \text{if } P_0^n(\bar{x})/P_1^n(\bar{x}) < \Delta, \end{cases}$$

for some constant $\Delta \geq 0$. Let T' be any other test. Then,

$$\alpha(T') \geq \alpha(T) \quad \text{or} \quad \beta(T') \geq \beta(T).$$

Proof: The proof follows simply from the observation that for all $\bar{x} \in U^n$

$$(T(\bar{x}) - T'(\bar{x})) \cdot (P_1^n(\bar{x}) - \Delta \cdot P_0^n(\bar{x})) \geq 0.$$

This is true because if $P_1^n(\bar{x}) - \Delta \cdot P_0^n(\bar{x})$, then $T(\bar{x}) = 1$ and the first quantity is non-negative. Similarly, when $P_1^n(\bar{x}) - \Delta \cdot P_0^n(\bar{x})$ is negative, $T(\bar{x}) = 0$ and $T(\bar{x}) - T'(\bar{x}) \leq 0$. Summing over all $\bar{x} \in U^n$ on both sides gives

$$\begin{aligned} & \mathbb{E}_{\bar{x} \sim P_1^n} [T(\bar{x}) - T'(\bar{x})] - \Delta \cdot \mathbb{E}_{\bar{x} \sim P_0^n} [T(\bar{x}) - T'(\bar{x})] \geq 0 \\ \Rightarrow & ((1 - \beta(T)) - (1 - \beta(T'))) - \Delta \cdot (\alpha(T) - \alpha(T')) \geq 0 \\ \Rightarrow & \frac{\beta(T') - \beta(T)}{\alpha(T) - \alpha(T')} \geq \Delta \geq 0. \end{aligned}$$

Thus, $\alpha(T) - \alpha(T') \geq 0$ implies $\beta(T') - \beta(T) \geq 0$. ■

We now discuss how to analyze the error probabilities for the optimal tests as characterized by the Neyman-Pearson lemma. As before, let $P_{\bar{x}}$ denote the type (empirical distribution on U) of the sequence \bar{x} . Check that the test $T(\bar{x})$ considered above can be written in the following form

$$\frac{P_1^n(\bar{x})}{P_0^n(\bar{x})} \geq \Delta \quad \Leftrightarrow \quad D(P_{\bar{x}} \| P_0) - D(P_{\bar{x}} \| P_1) \geq \frac{1}{n} \cdot \log \Delta.$$

We define the following sets of probability distributions.

$$\begin{aligned} \Pi & := \left\{ P \mid D(P \| P_0) - D(P \| P_1) \geq \frac{1}{n} \cdot \log \Delta \right\} \\ \Pi^c & := \left\{ P \mid D(P \| P_0) - D(P \| P_1) < \frac{1}{n} \cdot \log \Delta \right\} \end{aligned}$$

Check the following property of the sets Π and Π^c .

Exercise 1.3 Check that both the sets Π and Π^c are convex (and are in fact defined by linear inequalities in the distributions P). Also, check that Π is a closed set.

We know from Sanov's theorem that

$$\begin{aligned}\alpha(T) &= \mathbb{P}_{\bar{x} \sim P_0^n} [P_{\bar{x}} \in \Pi] \approx 2^{-n \cdot D(P_0^* \| P_0)} \\ \beta(T) &= \mathbb{P}_{\bar{x} \sim P_1^n} [P_{\bar{x}} \in \Pi^c] \approx 2^{-n \cdot D(P_1^* \| P_1)},\end{aligned}$$

where $P_0 = \arg \min_{P \in \Pi} \{D(P \| P_0)\}$. Also, since Π^c is not a closed set, we define P_1^* with respect to the closure of Π^c of Π^c i.e., $P_1 = \arg \min_{P \in \overline{\Pi^c}} \{D(P \| P_1)\}$. We will see later how to compute the distributions which minimize the KL-divergence (known as I-projections) as in the bounds above. The distributions P_0^* and P_1^* in the above bounds turn out to be of the form

$$P_0^*(x) = P_1^*(x) = \frac{P_0^\lambda(x) \cdot P_1^{1-\lambda}(x)}{\sum_{y \in U} P_0^\lambda(y) \cdot P_1^{1-\lambda}(y)},$$

where λ is the solution to an optimization problem. While the above analysis gives the optimal bounds for optimal all tests characterized by the Neyman-Pearson lemma, the bound we will use the most is the lower bound in terms of the total variation distance i.e.,

$$\min_T \{\alpha(T) + \beta(T)\} \geq 1 - \delta_{TV}(P_0, P_1).$$

We will now develop such a bound for the case of multiple hypotheses.

2 Fano's inequality and multiple hypothesis testing

Fano's inequality is concerned with Markov chains, which we saw before in the context of data processing inequality. We will denote the Markov chain as $Z \rightarrow Y \rightarrow \hat{Z}$. In the context of hypothesis testing, we can think of Z as the choice of an unknown hypothesis from some finite set (hypothesis class) U_Z . We think of Y as the "data" generated from this hypothesis, say a sequence \bar{x} of n independent samples. Finally, we think of \hat{Z} as a "guess" for Z , which depends only on the data. Fano's inequality is concerned with the probability of error in the guess, defined as $p_e = \mathbb{P}[\hat{Z} \neq Z]$. We have the following statement

Lemma 2.1 (Fano's inequality) Let $Z \rightarrow Y \rightarrow \hat{Z}$ be a Markov chain, and let $p_e = \mathbb{P}[\hat{Z} \neq Z]$. Let $H(p_e)$ denote the binary entropy function computed at p_e . Then,

$$H(p_e) + p_e \cdot \log(|U_Z| - 1) \geq H(Z|\hat{Z}) \geq H(Z|Y).$$

Proof: We define a binary random variable, which indicates an error i.e

$$E := \begin{cases} 1 & \text{if } \hat{Z} \neq Z \\ 0 & \text{if } \hat{Z} = Z \end{cases}$$

The bound in the inequality then follows from considering the entropy $H(Z, E|\hat{Z})$.

$$H(Z, E|\hat{Z}) = H(Z|\hat{Z}) + H(E|Z, \hat{Z}) = H(Z|\hat{Z}),$$

since $H(E|Z, \hat{Z}) = 0$ (why?) Another way of computing this entropy is

$$\begin{aligned} H(Z, E|\hat{Z}) &= H(E|\hat{Z}) + H(Z|E, \hat{Z}) \\ &= H(E|\hat{Z}) + p_e \cdot H(Z|E=1, \hat{Z}) + (1-p_e) \cdot H(Z|E=0, \hat{Z}) \\ &\leq H(E) + p_e \cdot H(Z|E=1, \hat{Z}) \\ &\leq H(p_e) + p_e \cdot \log(|U_Z| - 1). \end{aligned}$$

Comparing the two expressions then proves the claim. ■

We can use Fano's inequality to derive a convenient way of obtaining a lower bound for testing multiple hypotheses. However, we need the following property of KL-divergence.

Exercise 2.2 Prove that KL-divergence is (strictly) convex in both its arguments i.e., $\forall \alpha \in (0, 1)$ and all $P_1 \neq P_2, Q_1 \neq Q_2$,

$$\begin{aligned} D(\alpha \cdot P_1 + (1-\alpha) \cdot P_2 \| Q) &< \alpha \cdot D(P_1 \| Q) + (1-\alpha) \cdot D(P_2 \| Q) \\ D(P \| \alpha \cdot Q_1 + (1-\alpha) \cdot Q_2) &< \alpha \cdot D(P \| Q_1) + (1-\alpha) \cdot D(P \| Q_2) \end{aligned}$$

In fact, KL-divergence is jointly convex in both its arguments but we will need this property.

Let $\{P_v\}_{v \in \mathcal{V}}$ be a collection of hypotheses. Let the environment choose one of the hypotheses uniformly at random (denoted by a random variable V) and let $\bar{x} \sim P_v^n$ be a sequence of independent samples from a chosen distribution P_v (denoted by the random variable \bar{X}). We will now bound the probability of error for a classifier \hat{V} for V . Note that $V \rightarrow \bar{X} \rightarrow \hat{V}$ is a Markov chain.

Proposition 2.3 Let $V \rightarrow \bar{X} \rightarrow \hat{V}$ be the Markov chain as above. Then,

$$p_e = \mathbb{P}[V \neq \hat{V}] \geq 1 - \frac{n \cdot \mathbb{E}_{v_1, v_2 \in \mathcal{V}} [D(P_{v_1} \| P_{v_2})] + 1}{\log |\mathcal{V}|}.$$

Proof: From Fano's inequality, we have that

$$1 + p_e \cdot \log |\mathcal{V}| \geq H(p_e) + p_e \cdot \log |\mathcal{V}| \geq H(V|\bar{\mathbf{X}}) = \log |\mathcal{V}| - I(V; \bar{\mathbf{X}}).$$

We can now analyze the mutual information between V and $\bar{\mathbf{x}}$ using the equivalent expression in terms of KL-divergence.

$$\begin{aligned} I(V; \bar{\mathbf{x}}) &= D(P(V, \bar{\mathbf{X}}) \| P(V)P(\bar{\mathbf{X}})) \\ &= D(P(V) \| P(V)) + \mathbb{E}_{v \in \mathcal{V}} [D(P(\bar{\mathbf{X}}|V=v) \| P(\bar{\mathbf{X}}))] \\ &= \mathbb{E}_{v \in \mathcal{V}} [D(P_v^n \| \bar{P})], \end{aligned}$$

where $\bar{P} = \mathbb{E}_{v \in \mathcal{V}} [P_v^n]$ denotes the marginal distribution of $\bar{\mathbf{X}}$. Using the convexity of KL-divergence in the second argument, Jensen's inequality and the chain rule for KL-divergence, we get

$$\mathbb{E}_{v \in \mathcal{V}} [D(P_v^n \| \bar{P})] \leq \mathbb{E}_{v_1, v_2 \in \mathcal{V}} [D(P_{v_1}^n \| P_{v_2}^n)] = n \cdot \mathbb{E}_{v_1, v_2 \in \mathcal{V}} [D(P_{v_1} \| P_{v_2})].$$

Combining the bounds gives

$$1 + p_e \cdot \log |\mathcal{V}| \geq \log |\mathcal{V}| - n \cdot \mathbb{E}_{v_1, v_2 \in \mathcal{V}} [D(P_{v_1} \| P_{v_2})],$$

which proves the claim. ■