

Lecture-26: Fano's method

1 Fano's method

We discuss another method for proving minimax lower bound by reduction to multiple hypothesis testing. We call this program *Fano's method*, based on the Fano's inequality to show the impossibility result in one of the steps.

Step 1. We assume that the loss function $L : \Theta \times \Theta \rightarrow \mathbb{R}_+$ is a metric.

Step 2. Consider an ϵ -packing of the parameter space Θ , namely, a finite collection of parameters $T \triangleq \{\theta_i \in \Theta : i \in [M]\} \subset \Theta$ whose minimum separation is $\min_{i \neq j \in [M]} L(\theta_i, \theta_j) \geq \epsilon$.

Step 3. Suppose we can show that given observation X one cannot reliably distinguish these hypotheses $\Theta' \triangleq \{\theta_1, \dots, \theta_M\}$. That is, $P_{\theta_i} \{\tilde{\theta}(X) \neq \theta_i\} > 0$.

Step 4. Then the best estimation error $\mathbb{E}L(\hat{\theta}, \theta)$ is at least proportional to ϵ .

Step 5. The impossibility of testing is often shown by applying Fano's inequality in Corollary A.3, which bounds the probability of error of testing in terms of the mutual information.

Theorem 1.1 (Fano). Let $L : \Theta \times \Theta \rightarrow \mathbb{R}_+$ be a metric on parameter space Θ . Fix an estimator $\hat{\theta}$. For any $T \subseteq \Theta$ and $\epsilon > 0$,

$$P \left\{ L(\theta, \hat{\theta}) \geq \frac{\epsilon}{2} \right\} \geq 1 - \frac{C(T) + \ln 2}{\ln M(T, L, \epsilon)}, \quad (1)$$

where $C(T) \triangleq \sup_{\pi \in \mathcal{M}(T)} I(\theta; X)$ is the capacity of the channel $\theta \rightarrow X$ with input space T . Consequently,

$$\inf_{\hat{\theta}} \sup_{\theta \in \Theta} \mathbb{E}_{\theta} L(\theta, \hat{\theta})^r \geq \sup_{T \subseteq \Theta, \epsilon > 0} \left(\frac{\epsilon}{2} \right)^r \left(1 - \frac{C(T) + \ln 2}{\ln M(T, L, \epsilon)} \right). \quad (31.11)$$

Proof. It suffices to show (1), since the second result follows from the first applying Markov inequality for increasing function $f : \mathbb{R}_+ \rightarrow \mathbb{R}_+$ defined as $f(x) = x^r$ for any $x \in \mathbb{R}_+$ and $r \geq 1$. Fix $T \subseteq \Theta$. Consider an ϵ -packing $T' \triangleq \{\theta_1, \dots, \theta_M\} \subset T$ such that $\min_{i \neq j \in [M]} L(\theta_i, \theta_j) > \epsilon$. For each $\theta \in T'$, we define $\frac{\epsilon}{2}$ balls $B(\theta, \frac{\epsilon}{2}) \triangleq \{\theta' \in \Theta : L(\theta, \theta') \leq \frac{\epsilon}{2}\}$, and observe that $\{B(\theta, \frac{\epsilon}{2}) : \theta \in T'\}$ is a set of disjoint balls. Let $\theta : \Omega \rightarrow T'$ be uniformly distributed and $X \sim P_{\theta}$ conditioned on parameter θ . Given any estimator $\hat{\theta} : \mathcal{X} \rightarrow \Theta$, construct a test by rounding estimate $\hat{\theta}(X)$ to output $\tilde{\theta}(X) \triangleq \arg \min_{\theta \in T'} L(\theta, \hat{\theta}(X))$. Let $\theta \in T'$. From the triangle inequality for metric L and definition of $\tilde{\theta}(X)$, we get $L(\theta, \tilde{\theta}(X)) \leq L(\theta, \hat{\theta}(X)) + L(\hat{\theta}(X), \tilde{\theta}(X)) \leq 2L(\theta, \hat{\theta}(X))$, and thus

$$P \{ \theta \neq \tilde{\theta}(X) \} = \mathbb{E} \mathbb{1}_{\{\theta \neq \tilde{\theta}(X)\}} \leq \mathbb{E} \mathbb{1}_{\{L(\theta, \tilde{\theta}(X)) > \epsilon\}} \leq P \left\{ L(\theta, \hat{\theta}(X)) > \frac{\epsilon}{2} \right\}.$$

Recall that $|T'| = M$. The result follows from the application of Fano's inequality from Corollary A.3 to the Markov chain $\theta \rightarrow X \rightarrow \hat{\theta} \rightarrow \tilde{\theta}$ with uniformly distributed $\theta \in T'$ to lower bound $P \{ \theta \neq \tilde{\theta}(X) \}$ and using the definition of $C(T) \geq I(\theta; X)$. \square

In applying Fano's method, since it is often difficult to evaluate the capacity $C(T)$, it is useful to recall from Theorem 5.9 that $C(T)$ coincides with the KL radius of the set of distributions $\{P_{\theta} : \theta \in T\}$, namely, $C(T) \triangleq \inf_Q \sup_{\theta \in T} D(P_{\theta} \| Q)$. As such, choosing any Q leads to an upper bound on the capacity. As an application, we revisit the d -dimensional GLM in Corollary 28.8 under the ℓ_q loss for $1 \leq q \leq \infty$, with the particular focus on the dependency on the dimension. For a different application in sparse setting see Exercise VI.12.

Example 1.2. Consider GLM with *i.i.d.* sample size m , where $P_\theta = \mathcal{N}(\theta, I_d)^{\otimes m}$. Taking natural logarithms here and below, we have

$$D(P_\theta \| P_{\theta'}) = \frac{m}{2} \|\theta - \theta'\|_2^2.$$

In other words, KL-neighborhoods are ℓ_2 balls. As such, let us apply Theorem 1.1 to $T = B_2(\rho)$ for some $\rho > 0$ to be specified. Then $C(T) \leq \sup_{\theta \in T} D(P_\theta \| P_0) = m \frac{\rho^2}{2}$. To bound the packing number from below, we applying the volume bound in Theorem 27.3,

$$M(B_2(\rho), \|\cdot\|_q, \epsilon) \geq \frac{\rho^d \text{vol}(B_2)}{\epsilon^d \text{vol}(B_q)} \geq \left(\frac{c_q \rho d^{\frac{1}{q}}}{\epsilon \sqrt{d}} \right)^d.$$

for some constant c_q , where the last step follows the volume formula (27.13) for ℓ_q balls. Choosing $\rho = \sqrt{\frac{d}{n}}$ and $\epsilon = \frac{c_q}{2} \rho d^{\frac{1}{q} - \frac{1}{2}}$, an application of Theorem 1.1 yields the minimax lower bound

$$R_q \triangleq \inf_{\hat{\theta}} \sup_{\theta \in \mathbb{R}^d} \mathbb{E}_{X \sim \theta} \|\hat{\theta} - \theta\|_q \geq C_q \frac{d^{\frac{1}{q}}}{m} \quad (31.12)$$

for some constant C_q depending on q . This is the same lower bound as that in (30.9) obtained via the mutual information method plus the Shannon lower bound which is also volume-based.

Remark 1. For any $q \geq 1$, (31.12) is rate-optimal since we can apply the MLE $\hat{\theta} = X$. Note that at $q = \infty$, the constant C_q is still finite since $\text{vol}(B_\infty) = 2d$. However, for the special case of $q = \infty$, (31.12) does not depend on the dimension at all, as opposed to the correct dependency $\sqrt{\ln d}$ shown in Corollary 28.8. In fact, previously in Example 31.4 the application of Assouad's lemma yields the same suboptimal result. So is it possible to fix this looseness with Fano's method? It turns out that the answer is yes and the suboptimality is due to the volume bound on the metric entropy, which, as we have seen in Section 27.3, can be ineffective if ϵ scales with dimension. Indeed, if we apply the tight bound of $M(B_2, \|\cdot\|_\infty, \epsilon)$ in (27.18)¹, with $\epsilon = \sqrt{\frac{c \ln d}{m}}$ and $\rho = \sqrt{\frac{c' \ln d}{m}}$ for some absolute constants c, c' , we do get $R_\infty \geq \approx \sqrt{\frac{\ln d}{m}}$ as desired.

Remark 2. It is sometimes convenient to further bound the KL radius by the KL diameter, since $C(T) \leq \text{diam}_{\text{KL}}(T) \triangleq \sup_{\theta, \theta' \in T} D(P_{\theta'} \| P_\theta)$ (cf. Corollary 5.8). This suffices for Example 31.5.

Remark 3. In Theorem 1.1 we actually lower bound the global minimax risk by that restricted on a parameter subspace $T \subset \Theta$ for the purpose of controlling the mutual information, which is often difficult to compute. For the GLM considered in Example 31.5, the KL divergence is proportional to squared ℓ_2 distance and T is naturally chosen to be a Euclidean ball. For other models such as the covariance model (Exercise VI.16) wherein the KL divergence is more complicated, the KL neighborhood T needs to be chosen carefully. Later in Section 32.4 we will apply the same Fano's method to the infinite-dimensional problem of estimating smooth density.

A Fano's inequality

Definition A.1 (Binary entropy and KL divergence). Consider binary random variables $X, Y : \Omega \rightarrow \mathcal{X} \triangleq \{0, 1\}$ with respective probability mass functions $(p, 1-p), (q, 1-q) \in \mathcal{M}(\mathcal{X})$ for any $p, q \in [0, 1]$. Then binary entropy $h : [0, 1] \rightarrow [0, 1]$ is defined as $h(p) \triangleq H(X) = -p \ln p - (1-p) \ln(1-p)$ for all $p \in [0, 1]$, and the binary KL divergence is defined as $d : [0, 1] \times [0, 1] \rightarrow \mathbb{R}_+$ as $d(p, q) \triangleq D(P_X \| P_Y) = p \ln \frac{p}{q} + (1-p) \ln \frac{(1-p)}{(1-q)}$ for all $p, q \in [0, 1]$.

Theorem A.2 (Fano's inequality). Let $|\mathcal{X}| = M < \infty$, $X \rightarrow Y \rightarrow \hat{X}$ be a Markov chain, and $P_e \triangleq P\{X \neq \hat{X}\}$. Then the following are true

(a) $H(X | Y) \leq F_M(1 - P_e) \triangleq P_e \ln(M - 1) + h(P_e)$.

¹In fact, in this case we can also choose the explicit packing $\{\epsilon e_1, \dots, \epsilon e_d\}$

(b) If $P_{\max} \triangleq \max_{x \in \mathcal{X}} P_X(x) > 0$, then $I(X; Y) \geq (1 - P_e) \ln \frac{1}{P_{\max}} - h(P_e)$ regardless of $|\mathcal{X}|$.

Proof. Consider two joint distributions $P_{X,Y,\hat{X}} = P_X P_{Y|X} P_{\hat{X}|Y}$ and $Q_{X,Y,\hat{X}} = Q_X P_Y P_{\hat{X}|Y}$, and the data processor (kernel) $(X, Y, \hat{X}) \mapsto \mathbb{1}_{\{X \neq \hat{X}\}}$. We note that X and Y are independent under Q , and the observation Y has identical marginal $P_Y = Q_Y$ under both P and Q . Further, the kernel $P_{\hat{X}|Y} = Q_{\hat{X}|Y}$, i.e. the estimator for X is same for both distributions based on the observation Y . We recall the KL divergence data processing inequality for Markov chain $X \rightarrow Y \rightarrow \hat{X}$, such that

$$D(P_{X,Y,\hat{X}} \| Q_{X,Y,\hat{X}}) \geq D(P_{X,\hat{X}} \| Q_{X,\hat{X}}) \geq d(P\{X = \hat{X}\} \| Q\{X = \hat{X}\}).$$

From the definition of KL divergence, mutual information, and joint distributions P, Q , we get

$$D(P_{X,Y,\hat{X}} \| Q_{X,Y,\hat{X}}) = \mathbb{E}_P \ln \frac{dP_X dP_{Y|X}}{dQ_X dP_Y} = D(P_X \| Q_X) + I(X; Y) \geq d(P\{X = \hat{X}\} \| Q\{X = \hat{X}\}).$$

(a) Let $U_X \in \mathcal{M}(X)$ be a uniform distribution over \mathcal{X} . Then for $Q_X = U_X$, we obtain $D(P_X \| Q_X) = \ln M - H(X)$ and since $I(X; Y) = H(X) - H(X|Y)$, we get $D(P_X \| Q_X) + I(X; Y) = \ln M - H(X|Y)$. Further, since X and Y are independent under Q , we get $Q\{X = \hat{X}\} = \sum_{x \in \mathcal{X}} Q\{X = \hat{X} = x\} = \sum_{x \in \mathcal{X}} Q_X(x) Q_{\hat{X}}(x) = \frac{1}{M}$. It follows that

$$\ln M - H(X|Y) \geq d\left(P_e \| 1 - \frac{1}{M}\right) = -h(P_e) + (1 - P_e) \ln M + P_e \ln \frac{M}{M-1}.$$

(b) When $P_X = Q_X$, we get $D(P_X \| Q_X) = 0$ and $Q\{X = \hat{X}\} = \sum_x P_X(x) Q_{\hat{X}}(x) \leq P_{\max}$. Therefore,

$$\begin{aligned} I(X; Y) &\geq d(1 - P_e \| Q\{X = \hat{X}\}) = -h(P_e) - (1 - P_e) \ln Q\{X = \hat{X}\} - P_e \ln Q\{X \neq \hat{X}\} \\ &\geq -h(P_e) + (1 - P_e) \ln \frac{1}{Q\{X = \hat{X}\}} \geq -h(P_e) - (1 - P_e) \ln P_{\max}. \end{aligned}$$

□

The following corollary of the previous result emphasizes its role in providing converses or impossibility results for statistics and data transmission.

Corollary A.3 (Lower bound on average probability of error). Consider Markov chain $W \rightarrow X \rightarrow Y \rightarrow \hat{W}$, where W is uniform on $[M] \triangleq \{1, \dots, M\}$. Then

$$P_e \triangleq P\{W \neq \hat{W}\} \geq 1 - \frac{I(X; Y) + h(P_e)}{\ln M} \geq 1 - \frac{I(X; Y) + \ln 2}{\ln M}.$$

Proof. Apply Theorem A.2 and the data processing for mutual information $I(W; \hat{W}) \leq I(X; Y)$. □