# Lecture 11

<u>Review</u>  * Distribution learning

- minimax lower bounds: Fano's method

- Scheffé Selector

* Relating estimation to M-ary hypothesis testing

<u>Agenda</u>:  → Hypothesis testing / Distribution Testing

→ Uniformity testing

$\boxed{A}$ <u>Distribution Testing Formulation</u>

Let $\mathcal{C}$ denote a class of distributions.

Denote $d(P, \mathcal{C}) = \min\limits_{Q \in \mathcal{C}} d(P, Q)$.

Let $\overline{\mathcal{C}}_\varepsilon$ denote the class $\{P : d(P, \mathcal{C}) \geq \varepsilon\}$.

By observing samples $X_1, \ldots, X_n$, <u>test</u> if the samples
were generated from $P \in \mathcal{C}$ or $P \in \overline{\mathcal{C}}_\varepsilon$.

Specifically, a test $T : \mathcal{X}^n \to \{0, 1\}$ constitutes a $(\delta_1, \delta_2)$-test

if $\forall P \in \mathcal{C}, \quad P\left(T(x^n) = 0\right) \geq 1 - \delta_1$,

and
$\forall P \in \overline{\mathcal{C}}_\varepsilon, \quad P\left(T(x^n) = 1\right) \geq 1 - \delta_2$.

Denote by $n_{\delta_1, \delta_2}$ the least $n$ such that we can find
such a test $T$. For simplicity, we fix $\delta_1 = \delta_2 = \frac{1}{3}$.

Denote $n^* = n_{\frac{1}{3}, \frac{1}{3}}$.

A few popular formulations:

Consider $\mathcal{C} \subseteq \mathcal{P}_k \equiv (k-1)$-dimensional prob. simplex.

For all our examples, $\mathcal{C}$ will be closed ($\Rightarrow$ compact).

(1) <u>Uniformity Testing</u>: $\mathcal{C} = \{ \text{unif}[k] \}$

$$n^* = \Theta\left( \frac{\sqrt{k}}{\varepsilon^2} \right)$$

(2) <u>Identity Testing</u>: For $P \in \mathcal{P}_k$, $\mathcal{C} = \{P\}$.

$$n^* = \Theta\left( \frac{\sqrt{k}}{\varepsilon^2} \right).$$

(3) <u>Closeness Testing</u>: $\mathcal{X} = [k] \times [k]$

Are two sequences $(X_1, ..., X_n)$ and $(Y_1, ..., Y_m)$ generated from the same distributions or $(P, Q)$ s.t. $d(P, Q) \geq \varepsilon$.

[a slight modification of the general formulation]

$$n^* = \Theta\left( \max\left\{ \frac{n^{2/3}}{\varepsilon^{4/3}}, \frac{\sqrt{n}}{\varepsilon^2} \right\} \right)$$

(4) <u>Independence testing</u>: Observe sequence $(X^n, Y^n)$:

$\mathcal{C} = \{ P_X P_Y : P_X P_Y \in \mathcal{P}_k \}$

$$n^* = \Theta\left( \frac{n}{\varepsilon^2} \right)$$

B  <u>Uniformity Testing</u>: Collision-based tester

→ Goldreich and Ron, "On testing expansion in bounded-degree graphs", 2000.

→ Batu, Fischer, Fortnow,... "Testing random variables for independence and identity", FOCS 2001.

→ Diakonikolas, Gouleakis, Peebles, Price, "Collision-based testers are optimal for uniformity and closeness", 2016.

<u>Key idea</u>: Under uniform distribution, the no. of samples required to see collisions is roughly $\sqrt{k}$. (Birthday Paradox) In fact, a careful analysis will show that the uniform distribution takes the "longest" for collisions to appear.

Thus, we can make a test based on no. of collisions for uniformity.

<u>The Test</u>: $S \equiv S(x^n) = \sum_{i<j} \mathbb{1}_{\{X_i = X_j\}}$

$$T(x^n) = \begin{cases} \text{unif}, & S \leq \tau \\ \varepsilon\text{-away from unif}, & S > \tau, \end{cases}$$

where $\tau = \binom{n}{2} \dfrac{1 + c \cdot \varepsilon^2}{k}$. <span style="color:red">(N.B. → requires the knowledge of $k$)</span>

<u>Analysis</u>: Notation $\|P\|_\alpha = \sum_x P(x)^\alpha$

$\mathbb{E}[S] = \binom{n}{2} P(X_1 = X_2) = \binom{n}{2} \|P\|_2^2$

$\mathbb{E}[S^2] = \mathbb{E}\left[ \sum_{i_1 < j_1} \sum_{i_2 < j_2} \mathbb{1}_{\{X_{i_1} = X_{j_1}, X_{i_2} = X_{j_2}\}} \right]$

$$= \mathbb{E}\left[ \sum_{\substack{\text{two} \\ \text{distinct}}} \mathbb{1}_{\{x_i = x_j\}} + \sum_{\substack{\text{three} \\ \text{distinct}}} \mathbb{1}_{\{x_i = x_j = x_\ell\}} \right.$$

$$\left. + \sum_{\substack{i_1 < j_1, \ i_2 < j_2 \\ \text{all distinct}}} \mathbb{1}_{\{x_{i_1} = x_{j_1}\}} \mathbb{1}_{\{x_{i_2} = x_{j_2}\}} \right]$$

$$= \binom{n}{2} \|P\|_2^2 + 6\binom{n}{3} \|P\|_3^3 + \binom{n}{2}\binom{n-2}{2} \|P\|_2^4$$

$$\Rightarrow \mathrm{Var}(S) = \binom{n}{2} \|P\|_2^2 + n(n-1)(n-2) \|P\|_3^3$$

$$- \binom{n}{2} \left[ \underbrace{\frac{n!}{2!(n-2)!} - \frac{(n-2)!}{2!(n-4)!}}_{} \right] \|P\|_2^4$$

$$\frac{n(n-1) - (n-2)(n-3)}{2} = \frac{4n-6}{2} = 2n-3$$

$$= \binom{n}{2} \left( \|P\|_2^2 - \|P\|_2^4 \right) + n(n-1)(n-2) \left( \|P\|_3^3 - \|P\|_2^4 \right)$$

$$\leq \binom{n}{2} \|P\|_2^2 + n^3 \left( \|P\|_3^3 - \|P\|_2^4 \right).$$

<u>Under uniform distribution</u> $\left( P = \mathrm{unif}\,[k] \right)$

$$\|P\|_2^2 = \frac{1}{k}, \quad \|P\|_3^3 = \frac{1}{k^2}, \quad \mathbb{E}[S] = \binom{n}{2}\frac{1}{k}$$

Therefore, $P\left( \left| S - \binom{n}{2}\frac{1}{k} \right| > \tau \right) \leq \binom{n}{2}\frac{1}{k\tau^2}$

$$\Rightarrow P\left( \left| S - \binom{n}{2}\frac{1}{k} \right| > c\sqrt{3\binom{n}{2}}\frac{1}{k} \right) \leq \frac{1}{3}, \text{ for every } c > 1.$$

(1) $\quad P\left( S \leq \binom{n}{2}\frac{1}{k} + c\sqrt{3} \cdot \sqrt{\binom{n}{2}\frac{1}{k}} \right) \geq \frac{2}{3}, \quad \forall \ c > 1.$

**For $P$ s.t. $d(P, \text{unif}) > \varepsilon$**

$$\frac{1}{2} \sum_{i=1}^{k} \left| P_i - \frac{1}{k} \right| > \varepsilon \Rightarrow \sum_{i=1}^{k} P_i^2 - \frac{1}{k} > \frac{4\varepsilon^2}{k}.$$

$$\Rightarrow \|P\|_2^2 \geq \frac{1 + 4\varepsilon^2}{k}$$

**Case 1:** $\binom{n}{2} \|P\|_2^2 > n^3 \left( \|P\|_3^3 - \|P\|_2^4 \right)$

$$P\left( \left| S - \binom{n}{2} \|P\|_2^2 \right| > \tau \right) \leq \frac{n^2 \|P\|_2^2}{\tau^2}$$

$$\Rightarrow P\left( S > \binom{n}{2} \|P\|_2^2 - c\sqrt{3}\, n \|P\|_2 \right) \geq \frac{2}{3}, \quad \forall\, c \geq 1.$$

$(2) \Rightarrow P\left( S > \binom{n}{2} \frac{1}{k} + 4\binom{n}{2} \frac{\varepsilon^2}{k} - c\sqrt{3} \cdot \frac{n}{\sqrt{k}} \sqrt{1 + 4\varepsilon^2} \right) \geq \frac{2}{3}, \quad \forall\, c > 1.$

**Case 2:** $\binom{n}{2} \|P\|_2^2 \leq n^3 \left( \|P\|_3^3 - \|P\|_2^4 \right)$

We will use $-\|P\|_2^4$ to remove a $\frac{1}{k^2}$ term.

$$\|P\|_3^3 - \|P\|_2^4 \leq \sum_{i=1}^{k} P_i^3 - \frac{1}{k^2}$$

$$= \sum_{i=1}^{k} \left[ \left( P_i - \frac{1}{k} \right)^3 + \frac{1}{k^3} + 3\left( P_i - \frac{1}{k} \right)\frac{1}{k^2} + 3\left( P_i - \frac{1}{k} \right)^2 \frac{1}{k} \right] - \frac{1}{k^2}$$

$$= \sum_{i=1}^{k} \left( P_i - \frac{1}{k} \right)^3 + \frac{3}{k} \sum_{i=1}^{k} \left( P_i - \frac{1}{k} \right)^2$$

$$\leq \|P - P_{\text{unif}}\|_2^3 + \frac{3}{k} \|P - P_{\text{unif}}\|_2^2 = g_k^2 \left( \|P - P_{\text{unif}}\|_2 \right)$$

(To be continued in the next lecture)