

We are interested in bounding

$$\mathbb{P} \left( \sup_{C \in \mathcal{C}} \left| \frac{1}{n} \sum_{i=1}^n I(X_i \in C) - \mathbb{P}(C) \right| \geq t \right)$$

In Lecture 7 we hinted at Symmetrization as a way to deal with the unknown  $\mathbb{P}(C)$ .

**Lemma 10.1**(Symmetrization). *If  $t \geq \sqrt{\frac{2}{n}}$ , then*

$$\mathbb{P} \left( \sup_{C \in \mathcal{C}} \left| \frac{1}{n} \sum_{i=1}^n I(X_i \in C) - \mathbb{P}(C) \right| \geq t \right) \leq 2\mathbb{P} \left( \sup_{C \in \mathcal{C}} \left| \frac{1}{n} \sum_{i=1}^n I(X_i \in C) - \frac{1}{n} \sum_{i=1}^n I(X'_i \in C) \right| \geq t/2 \right).$$

*Proof.* Suppose the event

$$\sup_{C \in \mathcal{C}} \left| \frac{1}{n} \sum_{i=1}^n I(X_i \in C) - \mathbb{P}(C) \right| \geq t$$

occurs. Let  $X = (X_1, \dots, X_n) \in \{\sup_{C \in \mathcal{C}} |\frac{1}{n} \sum_{i=1}^n I(X_i \in C) - \mathbb{P}(C)| \geq t\}$ . Then

$$\exists C_X \text{ such that } \left| \frac{1}{n} \sum_{i=1}^n I(X_i \in C_X) - \mathbb{P}(C_X) \right| \geq t.$$

For a fixed  $C$ ,

$$\begin{aligned} \mathbb{P}_{X'} \left( \left| \frac{1}{n} \sum_{i=1}^n I(X'_i \in C) - \mathbb{P}(C) \right| \geq t/2 \right) &= \mathbb{P} \left( \left( \frac{1}{n} \sum_{i=1}^n I(X'_i \in C) - \mathbb{P}(C) \right)^2 \geq t^2/4 \right) \\ &\leq \text{(by Chebyshev's Ineq)} \frac{4\mathbb{E} \left( \frac{1}{n} \sum_{i=1}^n I(X'_i \in C) - \mathbb{P}(C) \right)^2}{t^2} \\ &= \frac{4}{n^2 t^2} \sum_{i,j} \mathbb{E}(I(X'_i \in C) - \mathbb{P}(C))(I(X'_j \in C) - \mathbb{P}(C)) \\ &= \frac{4}{n^2 t^2} \sum_{i=1}^n \mathbb{E}(I(X'_i \in C) - \mathbb{P}(C))^2 = \frac{4n\mathbb{P}(C)(1 - \mathbb{P}(C))}{n^2 t^2} \leq \frac{1}{nt^2} \leq \frac{1}{2} \end{aligned}$$

since we chose  $t \geq \sqrt{\frac{2}{n}}$ .

So,

$$\mathbb{P}_{X'} \left( \left| \frac{1}{n} \sum_{i=1}^n I(X'_i \in C_X) - \mathbb{P}(C_X) \right| \leq t/2 \right) \geq 1/2$$

because  $C_X$  does not depend on  $X'$ . Assume that the event

$$\left| \frac{1}{n} \sum_{i=1}^n I(X'_i \in C_X) - \mathbb{P}(C_X) \right| \leq t/2$$

occurs. Recall that

$$\left| \frac{1}{n} \sum_{i=1}^n I(X_i \in C_X) - \mathbb{P}(C_X) \right| \geq t.$$

Hence, it must be that

$$\left| \frac{1}{n} \sum_{i=1}^n I(X_i \in C_X) - \frac{1}{n} \sum_{i=1}^n I(X'_i \in C_X) \right| \geq t/2.$$

We conclude

$$\begin{aligned} \frac{1}{2} &\leq \mathbb{P}_{X'} \left( \left| \frac{1}{n} \sum_{i=1}^n I(X'_i \in C_X) - \mathbb{P}(C_X) \right| \leq t/2 \right) \\ &\leq \mathbb{P}_{X'} \left( \left| \frac{1}{n} \sum_{i=1}^n I(X_i \in C_X) - \frac{1}{n} \sum_{i=1}^n I(X'_i \in C_X) \right| \geq t/2 \right). \end{aligned}$$

Clearly,

$$\begin{aligned} &\mathbb{P}_{X'} \left( \left| \frac{1}{n} \sum_{i=1}^n I(X_i \in C_X) - \frac{1}{n} \sum_{i=1}^n I(X'_i \in C_X) \right| \geq t/2 \right) \\ &\leq \mathbb{P}_{X'} \left( \sup_{C \in \mathcal{C}} \left| \frac{1}{n} \sum_{i=1}^n I(X_i \in C) - \frac{1}{n} \sum_{i=1}^n I(X'_i \in C) \right| \geq t/2 \right). \end{aligned}$$

Since indicators are 0, 1-valued,

$$\begin{aligned} &\frac{1}{2} I \left( \sup_{C \in \mathcal{C}} \left| \frac{1}{n} \sum_{i=1}^n I(X_i \in C) - \mathbb{P}(C) \right| \geq t \right) \\ &\leq \mathbb{P}_{X'} \left( \sup_{C \in \mathcal{C}} \left| \frac{1}{n} \sum_{i=1}^n I(X_i \in C) - \frac{1}{n} \sum_{i=1}^n I(X'_i \in C) \right| \geq t/2 \right). \end{aligned}$$

Now, take expectation with respect to  $X_i$ 's to obtain

$$\begin{aligned} &\mathbb{P}_X \left( \sup_{C \in \mathcal{C}} \left| \frac{1}{n} \sum_{i=1}^n I(X_i \in C) - \mathbb{P}(C) \right| \geq t \right) \\ &\leq \mathbb{P}_{X, X'} \left( \sup_{C \in \mathcal{C}} \left| \frac{1}{n} \sum_{i=1}^n I(X_i \in C) - \frac{1}{n} \sum_{i=1}^n I(X'_i \in C) \right| \geq t/2 \right). \end{aligned}$$

**Theorem 10.1.** *If  $VC(\mathcal{C}) = V$ , then*

$$\mathbb{P} \left( \sup_{C \in \mathcal{C}} \left| \frac{1}{n} \sum_{i=1}^n I(X_i \in C) - \mathbb{P}(C) \right| \geq t \right) \leq 4 \left( \frac{2en}{V} \right)^V e^{-\frac{nt^2}{8}}.$$

*Proof.*

$$\begin{aligned} & 2\mathbb{P} \left( \sup_{C \in \mathcal{C}} \left| \frac{1}{n} \sum_{i=1}^n I(X_i \in C) - \frac{1}{n} \sum_{i=1}^n I(X'_i \in C) \right| \geq t/2 \right) \\ &= 2\mathbb{P} \left( \sup_{C \in \mathcal{C}} \left| \frac{1}{n} \sum_{i=1}^n \varepsilon_i (I(X_i \in C) - I(X'_i \in C)) \right| \geq t/2 \right) \\ &= 2\mathbb{E}_{X, X'} \mathbb{P}_\varepsilon \left( \sup_{C \in \mathcal{C}} \left| \frac{1}{n} \sum_{i=1}^n \varepsilon_i (I(X_i \in C) - I(X'_i \in C)) \right| \geq t/2 \right). \end{aligned}$$

The first equality is due to the fact that  $X_i$  and  $X'_i$  are i.i.d., and so switching their names (i.e. introducing random signs  $\varepsilon_i$ ,  $\mathbb{P}(\varepsilon_i = \pm 1) = 1/2$ ) does not have any effect. In the last line, it's important to see that the probability is taken with respect to  $\varepsilon_i$ 's, while  $X_i$  and  $X'_i$ 's are fixed.

Let

$$a(C) = (I(X_1 \in C) - I(X'_1 \in C), \dots, I(X_n \in C) - I(X'_n \in C)).$$

By Sauer's lemma,

$$\Delta_{2n}(\mathcal{C}, X_1, \dots, X_n, X'_1, \dots, X'_n) \leq \left( \frac{2en}{V} \right)^V.$$

In other words, any class will be equivalent to one of  $C_1, \dots, C_N$  on the data, where  $N \leq \left(\frac{2en}{V}\right)^V$ . Hence,

$$\begin{aligned}
& 2\mathbb{E}_{X, X'} \mathbb{P}_\varepsilon \left( \sup_{C \in \mathcal{C}} \left| \frac{1}{n} \sum_{i=1}^n \varepsilon_i (I(X_i \in C) - I(X'_i \in C)) \right| \geq t/2 \right) \\
&= 2\mathbb{E}_{X, X'} \mathbb{P}_\varepsilon \left( \sup_{1 \leq k \leq N} \left| \frac{1}{n} \sum_{i=1}^n \varepsilon_i (I(X_i \in C_k) - I(X'_i \in C_k)) \right| \geq t/2 \right) \\
&= 2\mathbb{E}_{X, X'} \mathbb{P}_\varepsilon \left( \bigcup_{k=1}^N \left| \frac{1}{n} \sum_{i=1}^n \varepsilon_i (I(X_i \in C_k) - I(X'_i \in C_k)) \right| \geq t/2 \right) \\
&\leq 2\mathbb{E} \sum_{k=1}^n \mathbb{P}_\varepsilon \left( \left| \frac{1}{n} \sum_{i=1}^n \varepsilon_i (I(X_i \in C_k) - I(X'_i \in C_k)) \right| \geq t/2 \right) \\
&\leq 2\mathbb{E} \sum_{k=1}^n 2 \exp \left( -\frac{-n^2 t^2}{8 \sum_{i=1}^n (I(X_i \in C) - I(X'_i \in C))^2} \right) \\
&\leq 2\mathbb{E} \sum_{k=1}^n 2 \exp \left( -\frac{-n^2 t^2}{8n} \right) \leq 2 \left( \frac{2en}{V} \right)^V 2e^{-\frac{nt^2}{8}},
\end{aligned}$$

□

where the first inequality above follows from the Hoeffding's inequality (see Lecture 7).