

ECE 6980
Algorithmic and Information-Theoretic Methods in Data Science

Instructor: Jayadev Acharya
Scribe: Ziteng Sun

Lecture #0
25th September, 2017

1 Scheffe Estimator

- 1 [Recap] Given M distributions, with $\frac{\log M}{\epsilon^2}$ samples from p , we can find a distribution close to p .
- 2 \mathcal{P} is a collection of distributions (e.g. all distributions over $[k]$). Let N_ϵ be the covering number of \mathcal{P} , which is the minimum value of M such that $\exists p_1, p_2, p_3, \dots, p_M \in \mathcal{P}$ such that $\forall p \in \mathcal{P}, \exists p_j$ such that $d(p, p_j) \leq \epsilon$ (d is a metric). Then we can learn a distribution from \mathcal{P} with $\frac{\log N_\epsilon}{\epsilon^2}$ samples using Scheffe estimators.
- 3 $\log N_\epsilon$ is called the metric entropy of \mathcal{P} with respect to metric d .
- 4 For distributions over $[k]$, $N_\epsilon \sim \left(\frac{k}{\epsilon}\right)^k$, so the complexity of the learner is $\frac{\log N_\epsilon}{\epsilon^2} \sim \frac{k \log(k/\epsilon)}{\epsilon^2}$.

2 Learning a distribution with exponentially tiny error

As from previous lectures, we can learn a distribution over $[k]$ using $O\left(\frac{k}{\epsilon^2}\right)$ samples with probability larger than 0.9. Using boosting technique, we can make the failure probability an arbitrary δ by adding a multiplicative factor of $\log \frac{1}{\delta}$. However, with the following analysis, we can show the complexity is $O\left(\frac{k + \log(1/\delta)}{\epsilon^2}\right)$.

Theorem 1. Mcdiarmid's Inequality

Suppose $f : \mathcal{X}^n \rightarrow \mathbb{R}$ is c -bounded difference, which means $|f(x_1, x_2, \dots, x_n) - f(x_1, \dots, x'_i, \dots, x_n)| \leq c, \forall x_1, \dots, x_n, x'_i \in \mathcal{X}$. Further suppose X_1, X_2, \dots, X_n are independent, then we have:

$$\Pr(f(X_1, X_2, \dots, X_n) - \mathbb{E}[f(X_1, X_2, \dots, X_n)] > \epsilon) \leq \exp\left(-\frac{2\epsilon^2}{nc^2}\right) \quad (1)$$

For reference, see <http://cs.nyu.edu/~rostami/ml/2007/ashish-mcdiarmid.pdf>.

Let $f(x_1, x_2, \dots, x_n) = \sum_{x=1}^k \left| \frac{N_x}{n} - \Pr(x) \right|$. Then $f(x_1, x_2, \dots, x_n)$ is $\frac{2}{n}$ -bounded difference. So we have:

$$\Pr(f(X_1, X_2, \dots, X_n) - \mathbb{E}[f(X_1, X_2, \dots, X_n)] > \epsilon) \leq \exp\left(-\frac{2(\epsilon/2)^2}{n(2/n)^2}\right) = \exp\left(-\frac{n\epsilon^2}{8}\right) \quad (2)$$

- 1 From previous lectures, with $O\left(\frac{k}{\epsilon^2}\right)$ samples, $\mathbb{E}[f(X_1, X_2, \dots, X_n)] \leq \epsilon/2$.
- 2 From Mcdiarmid's Inequality, with $O\left(\frac{\log(1/\delta)}{\epsilon^2}\right)$ samples, we can make the failure probability less than δ .

Hence the total sample complexity is $O\left(\frac{\log(1/\delta) + k}{\epsilon^2}\right)$. So for $\delta = e^{-k}$, the sample complexity will remain $O\left(\frac{k}{\epsilon^2}\right)$.

3 Property Estimation

Let $f : \mathcal{P} \rightarrow R$ be a property of a distribution, which includes:

- 1 $\mathbb{E}[X]$
- 2 Entropy $\sum_x p(x) \log \frac{1}{p(x)}$
- 3 Mode
- 4 Support size
- 5 Distance to uniformity $\sum_x |p(x) - \frac{1}{k}|$
- 6 Number of heavy hitters
- 7 $\sum p_x^2$.

Definition 1. Learning the property of a distribution

Let p be a distribution over $[k]$ and $x_1, x_2, x_3, \dots, x_n \sim p$ are independent samples. The goal is to find an estimator of the property \hat{f} such that with probability > 0.9 , we have $|\hat{f}(x_1, x_2, \dots, x_n) - f(p)| < \epsilon$.

Definition 2. Symmetric Property

A property f is symmetric if $f(p_\sigma) = f(p) \forall \sigma \in S_k$ where S_k is all the permutation over $[k]$.

Next we consider the problem of entropy estimation. $H(p) = \sum_x p(x) \log(\frac{1}{p(x)})$ and the performance of the empirical estimator.

Empirical estimator of entropy:

$$\hat{H}(p) = H(\hat{p}_n) = \sum_x \frac{N_x}{n} \log\left(\frac{n}{N_x}\right) \quad (3)$$

It can be shown that $H(p) \geq \mathbb{E}[H(\hat{p}_n)]$. Now we consider the expectation of the bias and the variance of \hat{H} .

$$H(p) - E[H(\hat{p}_n)] = \sum_x \left[p(x) \log \frac{1}{p(x)} - \mathbb{E}[\hat{p}_n(x)] \log \frac{1}{\hat{p}_n(x)} \right] \quad (4)$$

$$= \sum_x \mathbb{E} \left[p(x) \log \frac{1}{p(x)} - \hat{p}_n(x) \log \frac{1}{\hat{p}_n(x)} \right] \quad (5)$$

$$= \sum_x \mathbb{E} \left[(p(x) - \hat{p}_n(x)) \log \frac{1}{p(x)} \right] + \sum_x \mathbb{E} \left[\hat{p}_n(x) \log \frac{\hat{p}_n(x)}{p(x)} \right] \quad (6)$$

$$= \sum_x \mathbb{E} \left[\hat{p}_n(x) \log \frac{\hat{p}_n(x)}{p(x)} \right] \quad (7)$$

$$= \sum_x \mathbb{E} [D_{KL}(\hat{p}_n(x), p(x))] \quad (8)$$

$$\leq \mathbb{E} \left[\sum_x \frac{(\hat{p}_n(x) - p(x))^2}{p(x)} \right] \quad (9)$$

$$= \sum_x \frac{\mathbb{E} [(\hat{p}_n(x) - p(x))^2]}{p(x)} \quad (10)$$

$$= \sum_x \frac{p(x)(1 - p(x))}{np(x)} \quad (11)$$

$$= \frac{k - 1}{n} \quad (12)$$

We can also show that $Var(\hat{p}_n(x)) < \frac{\log^2 n}{n}$.

Definition 3. Bias-Variance Decomposition of an Estimator

Suppose \hat{z} is an estimator for a random variable z , then we have:

$$\mathbb{E}[(z - \hat{z})^2] = (z - \mathbb{E}[\hat{z}])^2 + \mathbb{E}[(\hat{z} - \mathbb{E}[\hat{z}])^2] \quad (13)$$

Hence if we want to estimate the entropy with high probability, we need $H(p) - E[H(\hat{p}_n)] = O(\epsilon)$ and $Var[H(\hat{p}_n)] = O(\epsilon^2)$. Hence we need $O(\frac{k}{\epsilon} + \frac{\log^2 k}{\epsilon^2})$ samples