

Lecture 5: February 12

Lecturer: Alessandro Rinaldo

Scribe: David Zhao

Note: *LaTeX template courtesy of UC Berkeley EECS dept.*

Disclaimer: *These notes have not been subjected to the usual scrutiny reserved for formal publications. They may be distributed outside this class only with the permission of the Instructor.*

In this lecture, we prove the bounded differences inequality and give an example of its application. We then begin our discussion of density estimation.

5.1 The bounded differences inequality

Recall the bounded differences property we defined in the previous lecture.

Definition: We say that a function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ satisfies the bounded differences property (BDP) if $\exists L_1, \dots, L_n > 0$ s.t. for any (X_1, \dots, X_n) in the domain of f , for any coordinate k ,

$$\sup_{x', y'} |f(x_1, \dots, x_{k-1}, y', x_{k+1}, \dots, x_n) - f(x_1, \dots, x_{k-1}, x', x_{k+1}, \dots, x_n)| \leq L_k$$

Now, we show that functions of random variables that satisfy the BDP concentrate around their expectations.

Theorem (Bounded Differences Inequality): Let X_1, \dots, X_n be independent random variables, and let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be a function that satisfies the BDP as defined above, with constants L_1, \dots, L_n . Then

$$\mathbb{P}(|Z - \mathbb{E}[Z]| > t) \leq 2 \exp\left(-\frac{2t^2}{\sum_{k=1}^n L_k^2}\right)$$

where $Z = f(X_1, \dots, X_n)$.

Proof: Suppose D_1, \dots, D_n is a martingale difference sequence, and suppose that $a_k \leq D_k \leq b_k$, $\forall k$ a.e. for sequences of constants $\{a_k\}$ and $\{b_k\}$. Then by the Azuma-Hoeffding inequality, we have

$$\mathbb{P}\left(\left|\sum_{k=1}^n D_k\right| > t\right) \leq 2 \exp\left(-\frac{2t^2}{\sum_{k=1}^n (b_k - a_k)^2}\right), \forall t > 0$$

Let

$$D_k = \mathbb{E}[Z|X_1, \dots, X_k] - \mathbb{E}[Z|X_1, \dots, X_{k-1}], \quad k > 1$$

$$D_0 = \mathbb{E}[Z]$$

Then $Z - \mathbb{E}[Z] = \sum_{k=1}^n D_k$. Note that D_1, \dots, D_n form a martingale difference sequence.

Now, for each $k = 1, \dots, n$, define

$$A_n = \inf_x \left\{ \mathbb{E}[Z|X_1, \dots, X_{k-1}, x] - \mathbb{E}[Z|X_1, \dots, X_{k-1}] \right\}$$

$$B_n = \sup_x \left\{ \mathbb{E}[Z|X_1, \dots, X_{k-1}, x] - \mathbb{E}[Z|X_1, \dots, X_{k-1}] \right\}$$

Observe that

$$D_k - A_k = \mathbb{E}[Z|X_1, \dots, X_k] - \inf_x \left\{ \mathbb{E}[Z|X_1, \dots, X_{k-1}, x] \right\} \geq 0 \text{ a.e.}$$

Similarly,

$$B_k - D_k = \sup_x \left\{ \mathbb{E}[Z|X_1, \dots, X_{k-1}, x] \right\} - \mathbb{E}[Z|X_1, \dots, X_k] \geq 0 \text{ a.e.}$$

In other words, for sequences of random variables $\{A_k\}$ and $\{B_k\}$, it holds that $A_k \leq D_k \leq B_k$, $\forall k$ a.e.

Finally, we unpack $B_k - A_k$ for each k , and show that it is $\leq L_k$, using the independence assumption.

$$\begin{aligned} B_k - A_k &= \sup_x \left\{ \mathbb{E}[Z|X_1, \dots, X_{k-1}, x] \right\} - \inf_y \left\{ \mathbb{E}[Z|X_1, \dots, X_{k-1}, y] \right\} \\ &= \sup_x \left\{ \int f(X_1, \dots, X_{k-1}, x, x_{k+1}, \dots, x_n) dP(x_{k+1}, \dots, x_n | X_1, \dots, X_{k-1}, x) \right\} \\ &\quad - \inf_y \left\{ \int f(X_1, \dots, X_{k-1}, y, x_{k+1}, \dots, x_n) dP(x_{k+1}, \dots, x_n | X_1, \dots, X_{k-1}, y) \right\} \end{aligned}$$

Let x^* and y^* be the (random) numbers that achieve the sup and inf over z , respectively, of

$$\int f(X_1, \dots, X_{k-1}, z, x_{k+1}, \dots, x_n) dP(x_{k+1}, \dots, x_n | X_1, \dots, X_{k-1}, z)$$

Note that in general,

$$dP(x_{k+1}, \dots, x_n | X_1, \dots, X_{k-1}, x^*) \neq dP(x_{k+1}, \dots, x_n | X_1, \dots, X_{k-1}, y^*)$$

unless X_1, \dots, X_n are independent, in which case both expressions are equal to

$$dP(x_{k+1}) \dots dP(x_n)$$

But the theorem does assume that X_1, \dots, X_n are independent, so we have that

$$\begin{aligned} B_k - A_k &= \sup_x \left\{ \int f(X_1, \dots, X_{k-1}, x, x_{k+1}, \dots, x_n) dP(x_{k+1}, \dots, x_n | X_1, \dots, X_{k-1}, x) \right\} \\ &\quad - \inf_y \left\{ \int f(X_1, \dots, X_{k-1}, y, x_{k+1}, \dots, x_n) dP(x_{k+1}, \dots, x_n | X_1, \dots, X_{k-1}, y) \right\} \\ &= \sup_x \left\{ \int f(X_1, \dots, X_{k-1}, x, x_{k+1}, \dots, x_n) dP(x_{k+1}) \dots dP(x_n) \right\} \\ &\quad - \inf_y \left\{ \int f(X_1, \dots, X_{k-1}, y, x_{k+1}, \dots, x_n) dP(x_{k+1}) \dots dP(x_n) \right\} \\ &\leq \sup_{x_1, \dots, x_{k-1}, x_{k+1}, \dots, x_n, x, y} \left\{ \left| f(x_1, \dots, x_{k-1}, x, x_{k+1}, \dots, x_n) - f(x_1, \dots, x_{k-1}, y, x_{k+1}, \dots, x_n) \right| \right\} \\ &\leq L_k, \quad \forall k \text{ a.e.} \end{aligned}$$

and the claim follows. ■

5.2 Application of the bounded differences inequality

Example: Let G_n be an Erdos-Renyi graph on n nodes. Let the $\binom{n}{2}$ potential edges be i.i.d. Bernoulli random variables, where 1 indicates the presence of an edge and 0 indicates its absence. Let C_n be the clique number, i.e. the size of the largest complete subgraph.

We show that C_n concentrates around its expectation. Note that C_n can be thought of as the result of a function applied to all of the $\binom{n}{2}$ potential edges of the graph. This function satisfies the BDP, because if we only change one edge at a time while keeping all other edges fixed, then C_n changes by at most 1. In other words, $L_k = 1, \forall k$. It follows that

$$\mathbb{P}\left(\left|\frac{C_n}{n} - E\left[\frac{C_n}{n}\right]\right| \geq t\right) = \mathbb{P}\left(|C_n - E[C_n]| \geq nt\right) \leq 2 \exp\left(-\frac{2n^2 t^2}{\sum_{k=1}^n 1}\right) = 2 \exp(-2nt^2), \forall t > 0$$

Example: We could also consider T_n , the number of triangles in the graph. Note that if T_n is the result of a function applied to all of the $\binom{n}{2}$ potential edges, then changing one edge at a time can change the value of T_n by as much as $n - 2$. So here, we cannot apply the BDP in a meaningful way.

5.3 Introduction to density estimation

We now turn our attention to a different type of problem. Let $X_1, \dots, X_n \stackrel{i.i.d.}{\sim} f$ where f is an unknown density. We seek to estimate f using our sample. Note that f is a point not in Euclidean space, but in a functional space, which makes density estimation a complex problem.

One type of density estimator is called the kernel density estimator.

Definition: Let $K : \mathbb{R} \rightarrow \mathbb{R}_+$ be a function such that $\int_{\mathbb{R}} K(x) dx = 1$. Then for a parameter $h > 0$ known as the bandwidth, define the kernel density estimator as

$$\hat{f}_h(x) = \frac{1}{nh} \sum_{k=1}^n K\left(\frac{X_k - x}{h}\right), x \in \mathbb{R}$$

In other words, we are simply convolving the empirical measure (consisting of point masses at the observed data points) with a kernel function to “smooth it out”.

Let us define a function $x \rightarrow f_h(x)$ as the pointwise expectation $\mathbb{E}[\hat{f}_h(x)]$ at every x . Then $f_h(x)$ is a deterministic probability density, determined by the kernel function K and the choice of bandwidth h .

Note that there are two levels of approximation here: the deviation $\hat{f}_h(x) - f_h(x)$ representing the variance, and the deviation $f_h - f$ representing the bias. As $h \rightarrow 0$, the bias $\rightarrow 0$ and the variance $\rightarrow \infty$.

We never explicitly observe the bias because we obviously do not know the true underlying function f . If we make certain smoothness assumptions about f belonging to a particular class of functions (e.g. a Holder class), we can derive expressions for both the bias and variance, which lead to well-known results on the optimal choice of bandwidth h as a function of n , and the resulting optimal rate of convergence.

Note that the above well-known results usually consider convergence rates in the sense of mean squared error, or L_2 error defined in terms of the L_2 norm:

$$\|\hat{f}_h(x) - f_h(x)\|_2 = \sqrt{\int (\hat{f}_h(x) - f_h(x))^2 dx}$$

The L_2 norm is popular because it is convenient to work with computationally, but a more natural way of thinking statistically about how “well” \hat{f}_h approximates f_h is the L_1 norm:

$$\|\hat{f}_h(x) - f_h(x)\|_1 = \int |\hat{f}_h(x) - f_h(x)| dx = 2d_{TV}(\hat{P}_h, P_h) = 2 \sup_{A \in \mathcal{A}} |\hat{P}_h(A) - P_h(A)|$$

where $d_{TV}(P, Q)$ is the TV distance between probability measures P and Q . Note that TV distance defines a very strong notion of similarity, as it considers the subset of the domain over which the two measures differ the most. So two measures that are close in TV distance are very similar over all possible subsets we might take over the domain.

Unfortunately, the L_1 norm is inconvenient to work with in practice because it is difficult to compute. For reference, see [DL01] for various methods to compute and use the L_1 norm in density estimation.

We conclude today by noting that the L_1 distance satisfies the BDP, with $L_k = \frac{2}{n}$, $\forall k$, since

$$\begin{aligned} f(x_1, \dots, x_{k-1}, x, x_{k+1}, \dots, x_n) - f(x_1, \dots, x_{k-1}, y, x_{k+1}, \dots, x_n) &\leq \frac{1}{nh} \int \left| K\left(\frac{x-z}{h}\right) - K\left(\frac{y-z}{h}\right) \right| dz \\ &\leq \frac{1}{nh} \left[h \int K(w) dw + h \int K(v) dv \right] = \frac{2}{n} \end{aligned}$$

using the substitutions $w = \frac{x-z}{h}$ and $v = \frac{y-z}{h}$.

It follows that

$$\mathbb{P}\left(\left|L_1(\hat{f}_h, f_h) - \mathbb{E}[L_1(\hat{f}_h, f_h)]\right| > t\right) \leq 2 \exp\left(-\frac{nt^2}{2}\right), \forall t > 0$$

The remarkable thing is that this concentration property holds regardless of the choice of bandwidth h .

Of course, in this case, it is easy to show that $L_1(\hat{f}_h, f_h)$ concentrates around its expectation, but, as we have noted, difficult to actually compute what that expectation is.

References

- [DL01] L. DEVROYE and G. LUGOSI, “Combinatorial Methods in Density Estimation,” *Springer Series in Statistics*, 2001.