

1 Martingales

We have seen that if $X = X_1 + \dots + X_n$ is a sum of independent $\{0, 1\}$ random variables, then X is tightly concentrated around its expected value $\mathbb{E}[X]$. The fact that the random variables were $\{0, 1\}$ -valued was not essential; similar concentration results hold if we simply assume that they are in some bounded range $[-L, L]$. One can also relax the independence assumption, as we will see next.

Consider a sequence of random variables X_0, X_1, X_2, \dots . The sequence $\{X_i\}$ is called a *discrete-time martingale* if it holds that

$$\mathbb{E}[X_{i+1} \mid X_0, X_1, \dots, X_i] = X_i$$

for every $i = 0, 1, 2, \dots$. More generally, the sequence $\{X_i\}$ is a martingale with respect to another sequence of random variables $\{Y_i\}$ if for every i , it holds that

$$\mathbb{E}[X_{i+1} \mid Y_0, Y_1, \dots, Y_i] = X_i.$$

Note that this is equivalent to $\mathbb{E}[X_{i+1} - X_i \mid Y_0, Y_1, \dots, Y_i] = 0$. If one thinks of $\{Y_0, Y_1, \dots, Y_i\}$ as all the “information” up to time i , then this says that the difference $X_{i+1} - X_i$ is unbiased conditioned on the past up to time i . Observe that for any i , we have

$$\mathbb{E}[X_i] = \mathbb{E}[\mathbb{E}[X_i \mid X_0, \dots, X_{i-1}]] = \mathbb{E}[X_{i-1}] = \dots = \mathbb{E}[X_0].$$

Martingales form an extremely useful class of random processes that appear in a vast array of settings (e.g., finance, machine learning, information theory, statistical physics, etc.). The classic example is that of a Gambler whose bank roll is X_0 . At each time, she chooses to play some game in the casino at some stakes. If we assume that every game is fair (that is, the expected utility from playing the game is 0), then the sequence $\{X_0, X_1, \dots\}$ forms a martingale, where X_i is the amount of money she has at time i .

Remark 1.1. The correct level of generality at which to define martingales involves a filtration. Formally, this is an increasing sequence of σ -algebras on our measure space $(\Omega, \mu, \mathcal{F})$: $\mathcal{F}_0 \subseteq \mathcal{F}_1 \subseteq \dots \subseteq \mathcal{F}$. A sequence of random variables $\{X_i\}$ is a martingale with respect to the filtration $\{\mathcal{F}_i\}$ if $\mathbb{E}[X_{i+1} \mid \mathcal{F}_i] = X_i$ for every $i \geq 0$.

1.1 Doob martingales

One reason martingales are so powerful is that they model a situation where one gains progressively more information over time. Suppose that \mathcal{U} is a set of objects, and $f : \mathcal{U} \rightarrow \mathbb{R}$. Let X be a random variable taking values in \mathcal{U} , and let $\{Y_i\}$ be another sequence of random variables. The associated *Doob martingale* is given by

$$X_i = \mathbb{E}[f(X) \mid Y_0, Y_1, \dots, Y_i].$$

In words, this is our “estimate” for the value of $f(X)$ given the information contained in $\{Y_0, \dots, Y_i\}$. To see that this is always a martingale with respect to $\{Y_i\}$, observe that

$$\mathbb{E}[X_{i+1} \mid Y_0, \dots, Y_i] = \mathbb{E}[\mathbb{E}[f(X) \mid Y_0, \dots, Y_{i+1}] \mid Y_0, \dots, Y_i] = \mathbb{E}[f(X) \mid Y_0, \dots, Y_i] = X_i,$$

where we have used the tower rule of conditional expectations.

Balls in bins. Suppose we throw m balls into n bins one at a time. At step i , we place ball i in a uniformly random bin. Let C_1, C_2, \dots, C_m be the sequence of (random) choices, and let C denote the final configuration of the system, i.e. exactly which balls end up in which bins.

Now we can consider a functional like $f(C) = \#$ of empty bins. If $X_i = \mathbb{E}[f(C) \mid C_1, \dots, C_i]$, then $\{X_i\}$ is a (Doob) martingale. It is straightforward to calculate that

$$\mathbb{E}[X_m] = \mathbb{E}[X_0] = \mathbb{E}[f(C)] = n \cdot \left(1 - \frac{1}{n}\right)^m.$$

Suppose we are interested the concentration of $X_m = f(C)$ around its mean value. Of course, we can write $X_m = Z_1 + \dots + Z_m$ where Z_i is the indicator of whether the i th been is empty after all the balls have been thrown. But note that the $\{Z_i\}$ variables are not independent—in particular, if I tell you that $Z_1 = 1$ (bin 1 is empty), it decreases slightly the likelihood that other bins are empty.

The vertex exposure filtration. Recall that $\mathcal{G}_{n,p}$ denotes the random graph model where an undirected graph on n vertices is chosen by including every edge independently with probability p . Suppose the vertices are numbered $\{1, 2, \dots, n\}$. Let $G \sim \mathcal{G}_{n,p}$ and denote by G_i the induced subgraph on the vertices $\{1, \dots, i\}$. G_0 denotes the empty graph.

Let $\chi(G)$ denote the chromatic number of G , and consider the Doob martingale

$$X_i = \mathbb{E}[\chi(G) \mid G_0, \dots, G_i].$$

If we wanted to understand concentration properties of $X_n = \chi(G)$, this seems even more daunting. The chromatic number is a very complicated parameter of a graph! Nevertheless, we will now see that martingale concentration inequalities allow us to achieve tight concentration using very limited information about a sequence of random variables.

2 The Hoeffding-Azuma inequality

Say that a martingale $\{X_i\}$ has L -bounded increments if

$$|X_{i+1} - X_i| \leq L$$

for all $i \geq 0$. (The preceding inequality is meant to hold with probability 1.)

Theorem 2.1. *For every $L > 0$, if $\{X_i\}$ is a martingale with L -bounded increments, then for every $\lambda > 0$ and $n \geq 0$, we have*

$$\mathbb{P}[X_n \geq X_0 + \lambda] \leq e^{-\frac{\lambda^2}{2L^2n}}$$

$$\mathbb{P}[X_n \leq X_0 - \lambda] \leq e^{-\frac{\lambda^2}{2L^2n}}$$

We will prove this in the next lecture. It's useful to note the following special case of the theorem.

Corollary 2.2. *Suppose that Z_1, Z_2, \dots, Z_n are independent random variables taking values in the interval $[-L, L]$. Put $Z = Z_1 + \dots + Z_n$ and $\mu = \mathbb{E}[Z]$. Then for every $\lambda > 0$, we have*

$$\mathbb{P}[Z \geq \mu + \lambda] \leq e^{-\lambda^2/(2L^2n)}$$

$$\mathbb{P}[Z \leq \mu - \lambda] \leq e^{-\lambda^2/(2L^2n)}$$

The Lipschitz condition. Recall the setting of Doob martingales, where \mathcal{U} is a set. Suppose that we can describe every element $u \in \mathcal{U}$ by a sequence of values $u = (u_1, u_2, \dots, u_n)$. (For instance, every configuration of m balls in n bins can be described by the sequence of which balls go into which bins.)

Say that f is L -Lipschitz if it holds that for every $i = 1, \dots, n$ and for every two elements $u = (u_1, u_2, \dots, u_i, \dots, u_n) \in \mathcal{U}$ and $u' = (u_1, u_2, \dots, u'_i, \dots, u_n) \in \mathcal{U}$ that differ only in the i th coordinate, we have

$$|f(u) - f(u')| \leq L.$$

Let $Z = (Z_1, \dots, Z_n)$ be a \mathcal{U} -valued random variable such that the random variables $\{Z_i\}$ are independent. We now confirm that the Doob martingale $X_i = \mathbb{E}[f(Z) \mid Z_1, \dots, Z_i]$ has L -bounded increments.

Let Z'_{i+1} be an independent copy of Z_{i+1} conditioned on Z_1, \dots, Z_i , and let $Z' = (Z_1, \dots, Z_i, Z'_{i+1}, \dots, Z_n)$. Then:

$$\begin{aligned} |X_{i+1} - X_i| &= \left| \mathbb{E}[f(Z) \mid Z_1, \dots, Z_{i+1}] - \mathbb{E}[f(Z) \mid Z_1, \dots, Z_i] \right| \\ &= \left| \mathbb{E} [f(Z) - f(Z') \mid Z_1, \dots, Z_{i+1}] \right| \\ &\leq \mathbb{E} [|f(Z) - f(Z')| \mid Z_1, \dots, Z_{i+1}] \\ &\leq L, \end{aligned}$$

where in the last step we have used the fact that the term inside the absolute value signs is always at most L by the L -Lipschitz property of f , and the fact that Z and Z' differ in at most one coordinate.

Remark 2.3. Note the power of [Theorem 2.1](#) combined with this construction of Doob marginales. It means that if we have any random variable $Z = (Z_1, \dots, Z_n)$ that is built out of independent pieces of information $\{Z_i\}$ and some quantity $f(Z)$ that we care about does not depend too much on changing any single piece of information, then $f(Z)$ is tightly concentrated about its mean. This is a vast generalization of the fact that sums of independent, bounded random variables are highly concentrated (cf. [Corollary 2.2](#)).

The number of empty bins. First let's apply this to balls and bins. Recall that for a sequence of choices C_1, \dots, C_m (where C_i is the bin that the i th ball is thrown into), we put $f(C_1, \dots, C_m)$ to be the number of empty bins. Then clearly f is 1-Lipschitz: Changing the fate of ball i can only change the number of empty bins by 1. Therefore the corresponding martingale $X_i = \mathbb{E}[f(C_1, \dots, C_m) \mid C_1, \dots, C_i]$ has 1-bounded increments, and Azuma's inequality implies that

$$\mathbb{P}[X_n \geq X_0 + \lambda] \leq e^{-\frac{\lambda^2}{2m}}.$$

Recall that $X_0 = \mathbb{E}[X_n] = n(1 - \frac{1}{m})^n$. Consider the situation where $m = n$ and thus $X_0 \asymp \frac{n}{e}$. If we put $\lambda = C\sqrt{n}$, we see that with high probability we expect the number of empty bins to be in the interval $\frac{n}{e} \pm O(\sqrt{n})$.

The chromatic number. Similarly, consider the vertex exposure martingale. We have to be a little more careful here to describe a graph G by a sequence (Z_1, \dots, Z_n) of *independent* random variables. The key is to think about Z_i containing the information on edges from vertex i to the vertices $\{1, \dots, i-1\}$ so that we have independence.

Since we can identify a graph G with the vector (Z_1, \dots, Z_n) , we can think of the chromatic number as a function $\chi(Z_1, \dots, Z_n)$. The function χ satisfies the 1-Lipschitz property because

changing the edges adjacent to some vertex i can only change the chromatic number by 1. The chromatic number cannot increase by more than one because we could always color i a new color; it cannot decrease by more than one because if we could color the graph without vertex i with c colors, then we can color the whole graph with $c + 1$ colors.

So the martingale $X_i = \mathbb{E}[\chi(G) \mid Z_1, \dots, Z_i] = \mathbb{E}[\chi(G) \mid G_1, \dots, G_i]$ has 1-bounded increments and Azuma's inequality tells us that

$$\mathbb{P}[\chi(G) \geq \mathbb{E}[\chi(G)] + \lambda] \leq e^{-\frac{\lambda^2}{2n}}.$$

Even without having any idea how to compute $\mathbb{E}[\chi(G)]$, we are able to say something significant about its concentration properties.

Remark 2.4. By the way, if $G \sim \mathcal{G}_{n,1/2}$, then $\mathbb{E}[\chi(G)] = n/(2 \log_2 n)$, so the concentration window here—which is $O(\sqrt{n})$ —is again quite small with respect to the expectation. In the next lecture, we will see how a more clever use of Azuma's inequality can achieve even better concentration of $\chi(G)$.