

Balls and Bins II : Poisson Random Variables and Poisson Approximation¹

- Recall the balls-and-bins setting: m balls are independently thrown into n bins. $L_i^{(m)}$ is the random variable indicating the number of balls in the i th bin. These are identical but not independent random variables, whose expectation is $\frac{m}{n}$.

In this lecture, we connect these random loads with *Poisson* random variables which are a powerful class of discrete random variables. In some sense, they form the discrete analog of the famous Gaussian random variables. Of note will be the following “approximation theorem”: to argue about events involving the random load vector $\vec{L}^{(m)} := (L_1^{(m)}, L_2^{(m)}, \dots, L_n^{(m)})$, it suffices to argue about a vector of *independent* Poissons, which is a much easier thing to do.

- To show the connection, let us figure out the probability $L_i^{(m)}$ is exactly r for some non-negative integer r . We see that

$$\Pr[L_i^{(m)} = r] = \underbrace{\binom{m}{r}}_{\text{ways to select } r \text{ balls}} \cdot \underbrace{\left(\frac{1}{n}\right)^r}_{\text{which all fall in bin } i} \cdot \underbrace{\left(1 - \frac{1}{n}\right)^{m-r}}_{\text{and the rest don't.}} \quad (1)$$

$$\underbrace{\approx}_{\text{when } r \ll n} \frac{m^r}{r!} \cdot \left(\frac{1}{n}\right)^r \cdot e^{-\frac{m}{n}} \quad (2)$$

Let’s list out the approximations: we have approximated $m(m-1)\dots(m-r+1) \approx m^r$, we have approximated $(1 - \frac{1}{n}) \approx e^{-\frac{1}{n}}$, and $m-r \approx m$. All of these are “ok”, when $n \gg 1$ and $r \ll n$. But the point is actually to show the connection with Poisson random variables which we describe next.

- Poisson Random Variables.** A Poisson random variable Z with parameter μ , denoted as $Z \sim \text{Pois}(\mu)$, is a non-negative integer valued random variable with pdf defined as

$$\text{For non-negative integer } r, \quad \Pr[Z = r] = \frac{e^{-\mu} \mu^r}{r!} \quad (\text{Poisson Random Variable})$$

Note that (2) is exactly the RHS of (Poisson Random Variable) when $\mu = \frac{m}{n}$. Let’s verify a couple of things, and then look at some magical properties of these variables.

Claim 1. Z defined in (Poisson Random Variable) is a valid probability distribution.

Proof. The RHS in (Poisson Random Variable) is indeed > 0 for any r . We need to check that it sums to 1. Indeed,

$$\sum_{r=0}^{\infty} \Pr[Z = r] = e^{-\mu} \cdot \underbrace{\sum_{r=0}^{\infty} \frac{\mu^r}{r!}}_{\text{This is } e^{\mu}} = 1$$

□

¹Lecture notes by Deeparnab Chakrabarty. Last modified : 21st April, 2023
 These have not gone through scrutiny and may contain errors. If you find any, or have any other comments, please email me at deeparnab@dartmouth.edu. Highly appreciated!

Claim 2. The expectation of $Z \sim \text{Pois}(\mu)$ is μ .

Proof.

$$\mathbf{Exp}[Z] = e^{-\mu} \sum_{r=1}^{\infty} \frac{r \cdot \mu^r}{r!} = \mu e^{-\mu} \sum_{r=1}^{\infty} \frac{\mu^{r-1}}{(r-1)!} = \mu \cdot \underbrace{e^{-\mu} \sum_{s=0}^{\infty} \frac{\mu^s}{s!}}_{=1 \text{ Claim 1}} = \mu$$

□

Exercise: Calculate the variance of $Z \sim \text{Pois}(\mu)$. Surprised?

Before we dive into the deeper connection with balls and bins, let's cover a powerful fact about Poisson random variables.

Theorem 1 (Sum of independent Poissons is Poisson). Let Z_1, \dots, Z_n be n independent Poisson random variables with $Z_i \sim \text{Pois}(\mu_i)$. Then, $Z := \sum_{i=1}^n Z_i$ is $\sim \text{Pois}(\mu)$ where $\mu := \sum_{i=1}^n \mu_i$.

Proof. Let's prove this for $n = 2$ and the rest follows inductively. Let $Z = Z_1 + Z_2$. Then,

$$\begin{aligned} \Pr[Z = r] &= \sum_{s=0}^r \Pr[Z_1 = s \wedge Z_2 = r - s] \stackrel{\text{independence}}{=} \sum_{s=0}^r \Pr[Z_1 = s] \cdot \Pr[Z_2 = r - s] \\ &= \sum_{s=0}^r \left(\frac{e^{-\mu_1} \mu_1^s}{s!} \right) \cdot \left(\frac{e^{-\mu_2} \mu_2^{r-s}}{(r-s)!} \right) \\ &= e^{-(\mu_1 + \mu_2)} \sum_{s=0}^r \frac{\mu_1^s \mu_2^{r-s}}{s!(r-s)!} = \frac{e^{-\mu}}{r!} \sum_{s=0}^r \underbrace{\frac{r!}{s!(r-s)!}}_{\text{observe this is } \binom{r}{s}} \mu_1^s \mu_2^{r-s} \\ &= \frac{e^{-\mu} \mu^r}{r!} \text{ by the Binomial Theorem} \quad \square \end{aligned}$$

This above facts allow us to prove exactly the same Chernoff bounds for sums of Poisson variables (which, recall, are very different from Bernoulli variables; in particular, these Poisson random variables are unbounded.)

Theorem 2 (Chernoff Bounds for Sums of Independent Poissons.). Let X be a Poisson random variable with parameter μ . Then for any $t > 0$, we have

$$\Pr[X \geq (1+t)\mu] \leq e^{-\mu \cdot g(t)} \quad \text{and} \quad \Pr[X \leq (1-t)\mu] \leq e^{-\mu \cdot h(t)} \quad (3)$$

where $g(t) := (1+t) \ln(1+t) - t$ and $h(t) := (1-t) \ln(1-t) + t$.

Remark: Consequently, using [Theorem 1](#) one gets the following. Suppose X_1, \dots, X_n are *independent* Poisson random variables and $X = \sum_{i=1}^n X_i$. Then for any $\varepsilon \in (0, 1)$,

$$\Pr[X \leq (1 - \varepsilon) \mathbf{Exp}[X]] \leq e^{-\frac{\varepsilon^2 \mathbf{Exp}[X]}{2}} \quad (\text{LT})$$

and

$$\Pr[X \geq (1 + \varepsilon) \mathbf{Exp}[X]] \leq e^{-\frac{\varepsilon^2 \mathbf{Exp}[X]}{3}} \quad (\text{UT1})$$

For the “upper tail”, that is for “larger” deviations, we have when $1 \leq t \leq 4$, we have the following (changing ε to t so as to underscore that the deviation is big)

$$\Pr[X \geq (1 + t) \mathbf{Exp}[X]] \leq e^{-\frac{t^2 \mathbf{Exp}[X]}{4}} \quad (\text{UT2})$$

and for $t > 4$ (really large), we have

$$\Pr[X \geq (1 + t) \mathbf{Exp}[X]] \leq e^{-\frac{t \ln t \mathbf{Exp}[X]}{2}} \quad (\text{UT3})$$

- **The Poisson Approximation : Connection with Balls and Bins.** Till now, the connection between balls-and-bins and Poisson random variables seems a bit tenuous: (2) is after all an approximation. Is thinking of the $L_i^{(m)}$'s as Poisson random variables correct? Is it useful? The following theorem captures this connection rigorously, and is called the *Poisson Approximation*.

Theorem 3 (Poisson Approximation for Balls and Bins.).

Suppose you throw m balls into n bins, each ball independently landing on a bin uniformly at random. Let $\mathcal{E}^{(m)}$ be an event of interest whose indicator random variable is a function of $f(L_1^{(m)}, \dots, L_n^{(m)})$. Consider a second experiment where we choose n *independent and identical* Poisson random variables (Z_1, \dots, Z_n) where each $Z_i \sim \text{Pois}(\frac{m}{n})$. Then,

$$\Pr[\mathcal{E}^{(m)}] := \Pr[f(L_1^{(m)}, \dots, L_n^{(m)}) = 1] \leq e\sqrt{m} \cdot \Pr[f(Z_1, \dots, Z_n) = 1] \quad (\text{Gen-PA})$$

One can state a stronger statement if the events $\mathcal{E} := \mathcal{E}^{(m)}$ have probability *monotone* in m . More precisely, the events $\mathcal{E}^{(m)}$ as a function of m are monotonically non-decreasing, if $m \leq m'$ implies $\Pr[\mathcal{E}^{(m)}] \leq \Pr[\mathcal{E}^{(m')}]$. The events $\mathcal{E}^{(m)}$ as a function of m are monotonically non-increasing if $m \leq m'$ implies $\Pr[\mathcal{E}^{(m)}] \geq \Pr[\mathcal{E}^{(m')}]$. The event $\mathcal{E} := \mathcal{E}^{(m)}$ are called monotone in m if they are either monotonically non-decreasing or monotonically non-increasing. Most events of interest are monotone.

Theorem 4 (Poisson Approximation for Balls and Bins: Monotone Events.).

Suppose you throw m balls into n bins, each ball independently landing on a bin uniformly at random. Let $\mathcal{E}^{(m)}$ be a monotone event of interest whose indicator random variable is a function of $f(L_1^{(m)}, \dots, L_n^{(m)})$. Consider a second experiment where we choose n *independent and*

identical Poisson random variables (Z_1, \dots, Z_n) where each $Z_i \sim \text{Pois}(\frac{m}{n})$. Then,

$$\Pr[\mathcal{E}^{(m)}] := \Pr[f(L_1^{(m)}, \dots, L_n^{(m)}) = 1] \leq 2 \cdot \Pr[f(Z_1, \dots, Z_n) = 1] \quad (\text{Mon-PA})$$

- *Lower bound on the maximum load.* It should be clear how [Theorem 3](#) can be useful : we now have *independence* over the various bins which was missing in the normal balls-and-bins setting. Let us illustrate this by showing a converse to a theorem we showed in a previous lecture : when we throw n balls independently into n different bins, the maximum load is in fact $\Omega(\frac{\ln n}{\ln \ln n})$ with high probability.

Theorem 5. For large enough n , if we throw n balls into n bins, then the probability the maximum load is $\leq \frac{\ln n}{2 \ln \ln n}$ is at most $2e^{-\sqrt{n}}$.

Proof. We are interested in upper-bounding the bad event \mathcal{E} which occurs if all loads $L_i^{(m)} \leq \frac{\ln n}{2 \ln \ln n}$. First note that once we fix n , these events are monotonically decreasing in m ; the more balls we throw the less the likelihood of all balls being small. Therefore, we can apply [Theorem 4](#).

Define $f(x_1, \dots, x_n) = 1$ if all $x_i \leq \frac{\ln n}{2 \ln \ln n}$, and 0 otherwise. We upper bound the probability $\Pr[f(Z_1, \dots, Z_n) = 1]$, where $Z_i \sim \text{Pois}(1)$ (note that $m = n$ and therefore, $m/n = 1$), and by [\(Mon-PA\)](#), $\Pr[\mathcal{E}] \leq 2 \Pr[f(Z_1, \dots, Z_n) = 1]$.

First, fix an $Z_i \sim \text{Pois}(1)$ and let us calculate the probability this is less than $L := \lfloor \frac{\ln n}{2 \ln \ln n} \rfloor$.

$$\Pr[Z_i \leq L] = e^{-1} \sum_{j \leq L} \frac{1}{j!} = e^{-1} \cdot \left(e - \underbrace{\sum_{j > L} \frac{1}{j!}}_{\geq \frac{1}{(L+1)!}} \right) \leq 1 - \frac{1}{e(L+1)!}$$

Now, since the Z_i 's are *independent*, we get that $\Pr[f(Z_1, \dots, Z_n) = 1] = \Pr[\bigwedge_{i=1}^n \{Z_i \leq L\}] = (\Pr[Z_i \leq L])^n$. Using [\(Mon-PA\)](#), we get

$$\Pr[\mathcal{E}] = \Pr[f(L_1^{(n)}, \dots, L_n^{(n)}) = 1] \leq 2 \cdot \left(1 - \frac{1}{e(L+1)!} \right)^n \quad (4)$$

What remains is a calculation similar to the upper bound proof. We get that for large enough n ,

$$\ln(e(L+1)!) \leq \ln L^L = L \ln L \leq \frac{\ln n}{2 \ln \ln n} \cdot (\ln \ln n) = \frac{\ln n}{2} \Rightarrow e(L+1)! \leq \sqrt{n}$$

Substituting in [\(4\)](#), we get

$$\Pr[f(L_1^{(n)}, \dots, L_n^{(n)}) = 1] \leq 2 \left(1 - \frac{1}{\sqrt{n}} \right)^n \leq 2e^{-\sqrt{n}} \quad \square$$

Use : $(1 - t) \leq e^{-t}$ to see this

- **The Proof of the Poisson Approximation Theorem.** The main observation is the following elementary lemma which states that if we throw m balls into n bins, then the **distribution** of the load vector is *precisely* the same as the distribution of n **independent** Poisson random variables with parameter $\mu := \frac{m}{n}$ *conditioned* on the event that their sum is m . That is, if we sample n independent Poisson random variables with parameter $\frac{m}{n}$ and reject anything whose sum is not m , then the resulting distribution of vectors is the same as the distribution of the loads on the n bins when m balls are thrown.

Lemma 1. For any tuple of non-negative integers (m_1, m_2, \dots, m_n) such that $\sum_{i=1}^n m_i = m$,

$$\Pr[(L_1^{(m)}, L_2^{(m)}, \dots, L_n^{(m)}) = (m_1, \dots, m_n)] = \Pr\left[(Z_1, Z_2, \dots, Z_n) = (m_1, \dots, m_n) \mid \sum_{i=1}^n Z_i = m\right]$$

where each $Z_i \sim \text{Pois}(\frac{m}{n})$ and are mutually independent.

Proof. There is not much to this lemma rather than a calculation. Let us calculate the LHS. How many ways can we split m balls into n sets such that set i has m_i balls? This is precisely the *multinomial coefficient*, and equals

$$\binom{m}{m_1, m_2, \dots, m_n} = \frac{m!}{m_1! m_2! \dots m_n!}$$

Given such a split, what is the probability that the first specified m_1 balls go into bin 1? The answer is $(\frac{1}{n})^{m_1}$. Similarly for the other bins. And therefore,

$$\Pr[(L_1^{(m)}, L_2^{(m)}, \dots, L_n^{(m)}) = (m_1, \dots, m_n)] = \frac{m!}{m_1! m_2! \dots m_n!} \cdot \left(\frac{1}{n}\right)^m \quad (\text{LHS})$$

Now let's compute the RHS. We get,

$$\Pr\left[(Z_1, Z_2, \dots, Z_n) = (m_1, \dots, m_n) \mid \sum_{i=1}^n Z_i = m\right] = \frac{\Pr[(Z_1, Z_2, \dots, Z_n) = (m_1, \dots, m_n)]}{\sum_{i=1}^n \Pr[Z_i = m]} \quad (5)$$

Note that the numerator event implies the denominator event and therefore we don't include it as an "and" in the numerator. Now, the $\Pr[Z_i = m_i] = \frac{e^{-\mu} \mu^{m_i}}{m_i!}$, and the Z_i 's are independent. Therefore,

$$\Pr[(Z_1, Z_2, \dots, Z_n) = (m_1, \dots, m_n)] = \frac{e^{-n\mu} \mu^m}{m_1! m_2! \dots m_n!}$$

Finally, by [Theorem 1](#), $\sum_{i=1}^n Z_i$ is also a Poisson random variable with parameter $n\mu$. Therefore, $\Pr[\sum_{i=1}^n Z_i = m] = \frac{e^{-n\mu} (n\mu)^m}{m!}$. Plugging these into (5), we get

$$\begin{aligned} \Pr\left[(Z_1, Z_2, \dots, Z_n) = (m_1, \dots, m_n) \mid \sum_{i=1}^n Z_i = m\right] &= \frac{e^{-\mu} \mu^{m_i} \cdot m!}{e^{-n\mu} (n\mu)^m \cdot m_1! m_2! \dots m_n!} \\ &= \frac{m!}{m_1! m_2! \dots m_n!} \cdot \left(\frac{1}{n}\right)^m \\ &\stackrel{(\text{LHS})}{=} \Pr[(L_1^{(m)}, L_2^{(m)}, \dots, L_n^{(m)}) = (m_1, \dots, m_n)] \quad \square \end{aligned}$$

- *Completing the proof.* Now we can prove [Theorem 3](#) and [Theorem 4](#). In fact, one can establish more general statements than in ([Gen-PA](#)) and ([Mon-PA](#)). One can show that for *non-negative* function $f : \mathbb{Z}^n \rightarrow \mathbb{R}_{\geq 0}$, one has

$$\mathbf{Exp}[f(L_1^{(m)}, L_2^{(m)}, \dots, L_n^{(m)})] \leq e\sqrt{m} \cdot \mathbf{Exp}[f(Z_1, Z_2, \dots, Z_n)]$$

This implies the theorem since the expectation is the same as probability of occurrence for an indicator random variable. We start with the RHS:

$$\begin{aligned} \mathbf{Exp}[f(Z_1, \dots, Z_n)] &= \sum_{k=0}^{\infty} \mathbf{Exp}[f(Z_1, \dots, Z_n) | \sum_i Z_i = k] \cdot \Pr[\sum_{i=1}^n Z_i = k] \\ &\geq \mathbf{Exp}[f(Z_1, \dots, Z_n) | \sum_{i=1}^n Z_i = m] \cdot \Pr[\sum_{i=1}^n Z_i = m] \\ &\stackrel{\text{Lemma 1}}{=} \mathbf{Exp}[f(L_1^{(m)}, L_2^{(m)}, \dots, L_n^{(m)})] \cdot \frac{e^{-m} m^m}{m!} \end{aligned}$$

where the inequality used the non-negativity of f . The proof of [Theorem 3](#) follows since $m! < e\sqrt{m}(m/e)^m$.

Exercise: Prove the above. That is, for all integer $m \geq 1$, we have $m! < e\sqrt{m}(m/e)^m$. Hint: take (natural) logs of both sides, and replace summation via integration (in the correct direction).

To replace the $e\sqrt{m}$ by 2 for *monotone* events, one is a bit more careful with the inequality. We first note that [Lemma 1](#) gives that for *any* k , we have

$$\mathbf{Exp}[f(Z_1, \dots, Z_n) | \sum_i Z_i = k] = \mathbf{Exp}[f(L_1^{(k)}, \dots, L_n^{(k)})] = \Pr[\mathcal{E}^{(k)}]$$

because f is an indicator function. Now if monotonically non-decreasing, we get that for $k \geq m$, the above RHS is $\geq \Pr[\mathcal{E}^{(m)}]$ which in turn is $\mathbf{Exp}[f(Z_1, \dots, Z_n) | \sum_{i=1}^n Z_i = m]$. That is,

$$\forall k \geq m, \quad \mathbf{Exp}[f(Z_1, \dots, Z_n) | \sum_i Z_i = k] \geq \mathbf{Exp}[f(Z_1, \dots, Z_n) | \sum_{i=1}^n Z_i = m] \quad (6)$$

We can plug this into $\mathbf{Exp}[f(Z_1, \dots, Z_n)]$ as follows

$$\begin{aligned} \mathbf{Exp}[f(Z_1, \dots, Z_n)] &= \sum_{k=0}^{\infty} \mathbf{Exp}[f(Z_1, \dots, Z_n) | \sum_i Z_i = k] \cdot \Pr[\sum_{i=1}^n Z_i = k] \\ &\geq \sum_{k=m}^{\infty} \mathbf{Exp}[f(Z_1, \dots, Z_n) | \sum_i Z_i = k] \cdot \Pr[\sum_{i=1}^n Z_i = k] \\ &\geq \mathbf{Exp}[f(Z_1, \dots, Z_n) | \sum_{i=1}^n Z_i = m] \cdot \sum_{k=m}^{\infty} \Pr[\sum_{i=1}^n Z_i = k] \\ &= \mathbf{Exp}[f(Z_1, \dots, Z_n) | \sum_{i=1}^n Z_i = m] \cdot \Pr[\sum_{i=1}^n Z_i \geq m] \end{aligned}$$

where the first inequality used non-negativity of f , the second used (6). Now one uses another pretty fact about Poisson random variables.

Fact 1. Let $Z \sim \text{Pois}(m)$ where m is an integer. Then $\text{Med}(Z) = m$. That is, $\Pr[Z \geq m] \geq \frac{1}{2}$ and $\Pr[Z \leq m] \geq \frac{1}{2}$.

Plugging this fact into the above chain of inequalities gives

$$\mathbf{Exp}[f(Z_1, \dots, Z_n)] \geq \frac{1}{2} \cdot \mathbf{Exp}[f(Z_1, \dots, Z_n) | \sum_{i=1}^n Z_i = m] = \mathbf{Exp}[f(L_1^{(m)}, \dots, L_n^{(m)}) = 1] = \mathbf{Pr}[\mathcal{E}]$$

thus proving [Theorem 4](#) for monotonically non-decreasing events. Do you see how to get the 2 when f is monotonically non-increasing events? You should be.