

29 Jan 2025

The Coupon Collector Problem

Def The geometric distribution with parameter p is the distribution of the # of coin tosses until (and including) the first time heads is tossed, when tossing a coin with $\Pr(\text{heads}) = p$, and independent coin-toss outcomes.

$$\Pr(X=k) = (1-p)^{k-1} \cdot p$$

\uparrow $k-1$ tosses yield tails... and k^{th} toss is heads.

Expectation of geometric rand var

$$E[X] = \sum_{k=1}^{\infty} k \cdot \Pr(X=k)$$

$$= \sum_{k=1}^{\infty} k (1-p)^{k-1} \cdot p$$

Let $q = 1-p$

$$= p \sum_{k=1}^{\infty} k q^{k-1}$$

$$= p \frac{d}{dx} \left[\sum_{k=1}^{\infty} x^k \right]_{x=q} = p \frac{d}{dx} \left[\frac{x}{1-x} \right]_q$$

$$= p \frac{d}{dx} \left[1 + \frac{x}{1-x} \right]_{x=q}$$

$$= p \frac{d}{dx} \left[\frac{1}{1-x} \right]_{x=q}$$

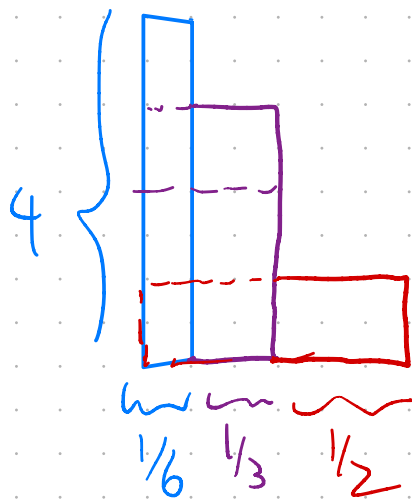
$$= \frac{p}{(1-q)^2} = \frac{p}{p^2} = \frac{1}{p}$$

Fact. If X is a random variable taking non-negative integer values,

$$E[X] = \sum_{j=0}^{\infty} \Pr(X > j)$$

Proof by picture. Suppose $X = \begin{cases} 1 & \text{w. prob. } 1/2 \\ 3 & \text{w. prob. } 1/3 \\ 4 & \text{w. prob. } 1/6 \end{cases}$

$$E[X] = \sum k \cdot \Pr(X = k)$$



For geometric RV,

$$\begin{aligned} \Pr(X > j) &= \Pr(\text{first } j \text{ coin tosses are tails}) \\ &= (1-p)^j \end{aligned}$$

$$E[X] = \sum_{j=0}^{\infty} \Pr(X > j) = \sum_{j=0}^{\infty} (1-p)^j = \frac{1}{1-(1-p)} = \frac{1}{p}$$

Last derivation. (principle of deferred decisions)

If first coin toss is H, $X = 1$.

If $\sim \sim \sim \sim \sim$ T, $X = 1 + Y$

where Y is $\text{Geom}(p)$.

$$\begin{aligned} E[X] &= E[X \mid \text{first toss H}] \cdot \Pr(H) \\ &\quad + E[X \mid \text{first toss T}] \cdot \Pr(T) \end{aligned}$$

$$= 1 \cdot p + (1 + E(X)) \cdot (1-p)$$

$$p \cdot E(X) = 1 \cdot p + 1 \cdot (1-p) = 1$$

$$E(X) = \frac{1}{p}$$

Variance of geometric RV

$$\text{Var}(X) = E[X^2] - (E(X))^2 = \frac{1}{p^2}$$

$$\begin{aligned} E[X^2] &= E[X^2 \mid \text{first toss is H}] \cdot p \\ &\quad + E[X^2 \mid \text{first toss is T}] \cdot (1-p) \end{aligned}$$

$$= 1 \cdot p + E[(1+Y)^2] \cdot (1-p) \quad \sim \text{Geom}(p)$$

$$= 1 \cdot p + E[(1+X)^2] \cdot (1-p)$$

$$= p + \left(\underbrace{E[1]}_1 + \underbrace{E[2x]}_{\frac{2}{p}} + E[x^2] \right) (1-p)$$

$$E[x^2] = p + (1-p) + \frac{2(1-p)}{p} + (1-p)E[x^2]$$

$$p \cdot E[x^2] = 1 + \frac{2-2p}{p}$$

$$E[x^2] = \frac{1}{p} + \frac{2}{p^2} - \frac{2}{p} = \frac{2-p}{p^2}$$

$$\text{Var}(X) = E[x^2] - (EX)^2 = \frac{1-p}{p^2}$$

COUPON COLLECTOR!

Phase 1 ends at $\tau_1 =$ first time at least 1 bin is occupied
 $= 1$

⋮

Phase k ends at $\tau_k =$ first time at least k bins are occupied.

If $\tau_k \leq t < \tau_{k+1}$ it means

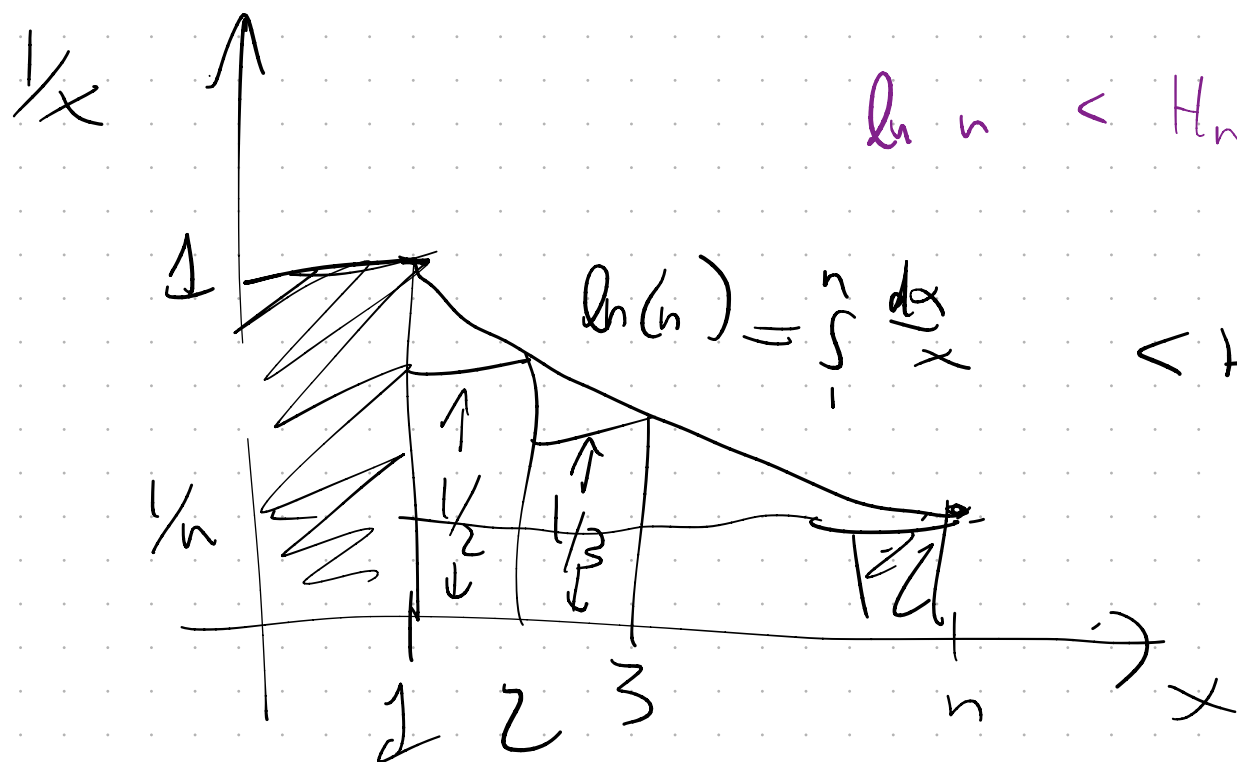
exactly k bins are occupied at time t

\implies Pr (ball thrown at $t+1$ occupies new bin)

$$= \frac{n-k}{n}$$

The rand var $X_{k+1} = \tau_{k+1} - \tau_k$ is
 geom. distr. with parameter $p = \frac{n-k}{n}$,
 independent of $X_1, X_2, \dots, X_{k-1}, X_{k+1}, \dots, X_n$.

$$\begin{aligned}
 E[\tau_n] &= E[\tau_1 + (\tau_2 - \tau_1) + (\tau_3 - \tau_2) + \dots + (\tau_n - \tau_{n-1})] \\
 &= E(X_1 + X_2 + \dots + X_n) \\
 &= \sum_{k=1}^n E\left(\text{Geom}\left(\frac{n-(k-1)}{n}\right)\right) \\
 &= \sum_{k=1}^n \frac{n}{n-(k-1)} = \frac{n}{n} + \frac{n}{n-1} + \dots + \frac{n}{1} \\
 &= n \cdot \left[1 + \frac{1}{2} + \dots + \frac{1}{n}\right] \\
 &= n \cdot H_n \leftarrow n^{\text{th}} \text{ harmonic number.}
 \end{aligned}$$



$$\ln n < H_n < 1 + \ln n.$$

$$\ln(n) = \int_1^n \frac{dx}{x} < H_n < 1 + \int_1^n \frac{dx}{x} = 1 + \ln(n)$$

Suppose $\frac{1}{p} \in \mathbb{N}$.

$$\Pr(\text{Geom}(p) > \frac{1}{p}) = (1-p)^{1/p}$$

$$< (e^{-p})^{1/p} = e^{-1} \approx 0.36\dots$$

Let $m_{\text{coupon}}(n)$ be the smallest m s.t.

$$\Pr(\tau_n > m) \leq \frac{1}{2}.$$

Estimating $m_{\text{coupon}}(n) \dots$

To show $\Pr(\tau_n \text{ far from } \mathbb{E}\tau_n \text{ is small})$

we use Chebyshev's Ineq.

$$\text{Var}(\tau_n) = \text{Var}(X_1 + X_2 + \dots + X_n)$$

$$= \text{Var}(X_1) + \dots + \text{Var}(X_n)$$

$$= \sum_{k=1}^n \frac{1 - \left(\frac{n-k+1}{n}\right)^2}{\left(\frac{n-k+1}{n}\right)^2}$$

$$< \sum_{k=1}^n \frac{1}{\left(\frac{n-k+1}{n}\right)^2} = n^2 \left[1 + \frac{1}{2^2} + \frac{1}{3^2} + \dots + \frac{1}{n^2} \right]$$

$$< 2n^2$$

$\forall x$

$$\Pr(|T_n - \mathbb{E}T_n| > x)$$

$$= \Pr\left(\left(T_n - \mathbb{E}T_n\right)^2 > x^2\right)$$

$$< \frac{\mathbb{E}\left(\left(T_n - \mathbb{E}T_n\right)^2\right)}{x^2} < \frac{2n^2}{x^2}$$

$$\Pr\left(|T_n - \mathbb{E}T_n| > 2n\right) < \frac{2n^2}{4n^2} = \frac{1}{2}$$

The event $T_n > n \cdot H_n + 2n$

cannot have probability $> \frac{1}{2}$.

$$\therefore m_{\text{coupon}}(n) < n H_n + 2n$$

$$< n \ln(n) + 3n$$

Also

$$m_{\text{coupon}}(n) > n \ln(n) - 2n$$