

From prev lecture: Coupon collector problem, Markov and Chebyshov's inequality.

Recalling Chebyshov's inequality:

Let X be a r.v with expectation $\mathbb{E}[X]$ and standard deviation σ_X . Then for any $t \in \mathbb{R}_{\geq 0}$

$$\Pr[|X - \mathbb{E}[X]| \geq t\sigma_X] \leq \frac{1}{t^2}$$

$$\Pr[X \geq 2\mathbb{E}[X]]$$

Let $\vartheta = t\sigma_X$. Then $t^2 = \frac{\vartheta^2}{\sigma_X^2}$ and $\sigma_X^2 = \text{Var}[X]$.

$$\Rightarrow \Pr[|X - \mathbb{E}[X]| \geq \vartheta] \leq \frac{\text{Var}[X]}{\vartheta^2}.$$

$$\Pr[|X - \mathbb{E}[X]| \geq \mathbb{E}[X]] \leq \frac{\text{Var}[X]}{\mathbb{E}[X]^2}$$

We wanted to find out prob that $X \geq 2\mathbb{E}[X]$.

Recall that this was at most $\frac{1}{2}$ from Markov ineq.

By setting $\vartheta = \mathbb{E}[X]$, we get that

$$\Pr[|X - \mathbb{E}[X]| \geq \mathbb{E}[X]] \leq \frac{\text{Var}[X]}{(\mathbb{E}[X])^2}.$$

$$\begin{aligned} |X - \mu| &\geq t \\ \downarrow \\ X &> \mu + t \\ \text{and} \\ X &\leq \mu - t. \end{aligned}$$

For coupon collector problem, $\text{Var}[X] = O(n^2)$ and $\mathbb{E}[X] = \Theta(n \ln n)$

Thus, we get that $\Pr[|X - \mathbb{E}[X]| \geq \mathbb{E}[X]] \leq O\left(\frac{1}{n^2}\right)$.

(In comparison, Markov gives only $\frac{1}{2}$.)

Remark: $\text{Var}[X+Y] = \text{Var}[X]+\text{Var}[Y]+2\text{Covar}[X,Y]$.

$$\mathbb{E}[(X-\mathbb{E}[X])(Y-\mathbb{E}[Y])].$$

In coupon collector problem, $X=\sum X_i$ and X_i 's were independent r.v.s and thus $\text{Var}[X]=\sum_{i=1}^n \text{Var}[X_i]$.

Other arguments to obtain tail probabilities:

Probability of not obtaining i^{th} coupon after s steps for some $s > n \ln n + \Theta(n)$:

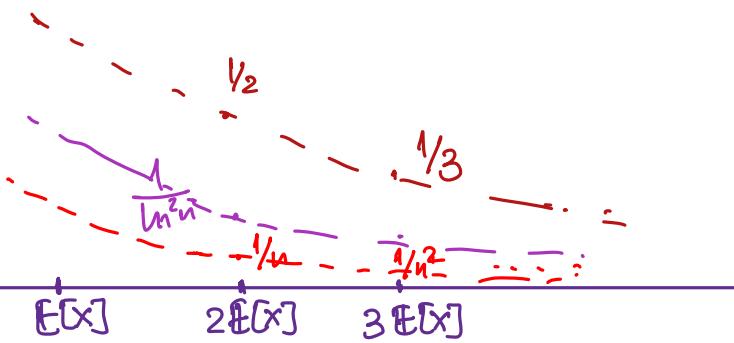
$$\left(1 - \frac{1}{n}\right)^s = \left(1 - \frac{1}{n}\right)^{n \cdot \frac{s}{n}} \sim e^{-\frac{s}{n}}.$$

Probability that there exists some coupon that is not yet picked up after s trials, by union bound is \leq sum of probabilities that each coupon is not picked after s trials.

$$\approx \sum_{i=1}^n e^{-\frac{s}{n}} = \frac{n}{e^{s/n}}. \quad \text{w.p. } \geq 1 - \frac{1}{n}; \text{ all coupons are collected in } s = 2n \ln n \text{ trials.}$$

Set $s = 2n \ln n$. Then the tail prob is at most $\frac{1}{n}$.

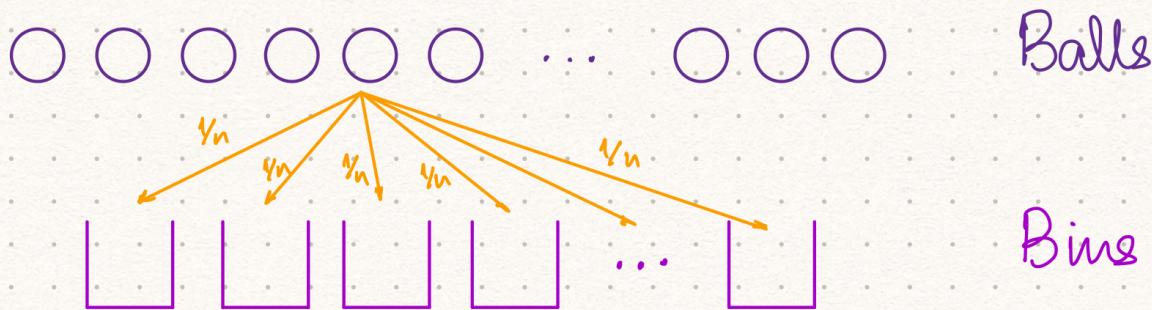
Tail probabilities:



2.2 Balls and Bins.

We are given m "indistinguishable" balls and n bins. Each ball is placed in a bin that is chosen uniformly at random.

Question: How are these balls distributed amongst the bins.



Question 1: What is the expected no. of balls in each bin?

Question 2: What is the prob that no bin gets more than $k-1$ balls?

Let E_j be the event that bin j has $\geq k$ balls.

$$\Pr[E_j] = \sum_{l=k}^m \Pr[\text{bin } j \text{ gets exactly } l \text{ balls}]$$

$$= \sum_{l=k}^m \binom{m}{l} \cdot \left(\frac{1}{n}\right)^l \cdot \left(1 - \frac{1}{n}\right)^{m-l}$$

$$\left(\frac{a}{b}\right)^b \leq \binom{a}{b} \leq \left(\frac{ea}{b}\right)^b$$

$$\leq \sum_{l=k}^m \binom{m}{l} \cdot \left(\frac{1}{n}\right)^l \leq \sum_{l=k}^m \left(\frac{em}{l}\right)^l \cdot \left(\frac{1}{n}\right)^l = \sum_{l=k}^m \left(\frac{em}{ln}\right)^l.$$

Say $m=n$. Then $\Pr[E_j] \leq \sum_{l=k}^n \left(\frac{e}{k}\right)^l \leq \sum_{l=k}^n \left(\frac{e}{k}\right)^l$

$$\begin{array}{ll}
 b_1 & b_2 \\
 \text{trial 1} & B_1 \quad B_2 \xrightarrow{\text{2 bins}} \\
 \text{trial 2} & B_1 \quad B_1 \\
 \text{trial 3} & B_2 \quad B_2 \\
 \text{trial 4} & B_2 \quad B_1 \leftarrow
 \end{array}
 \quad = \left(\frac{e}{k}\right)^k \left[1 + \left(\frac{e}{k}\right) + \dots + \left(\frac{e}{k}\right)^{n-k} \right] \\
 \leq \left(\frac{e}{k}\right)^k \cdot \frac{1}{1 - \left(\frac{e}{k}\right)}$$

Let $k = \left\lceil \frac{3 \ln n}{m \ln n} \right\rceil$. Then $\Pr[E_j] \leq \frac{1}{n^2}$.

$$\begin{aligned}
 & \left(\frac{e \ln \ln(n)}{3 \ln(n)} \right)^{\frac{3 \ln n}{\ln \ln n}} \cdot \frac{3 \ln n}{3 \ln n - \ln \ln(n)} \\
 & \leq \left(\frac{\ln \ln(n)}{\ln(n)} \right)^{\frac{3 \ln(n)}{\ln \ln(n)}} \approx e^{(\ln \ln n - \ln \ln n) \frac{3 \ln n}{\ln \ln n}} \\
 & = e^{-\left(1 - \frac{\ln \ln \ln(n)}{\ln \ln(n)}\right) \cdot 3 \ln(n)}
 \end{aligned}$$

$$\leq \frac{1}{n^2} \quad \text{For sufficiently large } n, 1 - \frac{\ln \ln \ln(n)}{\ln \ln(n)} \geq \frac{2}{3}.$$

$$\mathbb{E}[X_j] = \Pr[X_j = 1] = \Pr[E_j] \leq \frac{1}{n^2}.$$

X_j be a 0-1 r.v s.t

$$X_j = \begin{cases} 1 & \text{if bin } j \text{ contains more than } k \\ & \text{balls.} \\ 0 & \text{otherwise.} \end{cases}$$

Prob that no bin contains more than $k-1$ balls

$$\begin{aligned}
 &= 1 - \Pr[\exists \text{ a bin with at least } k \text{ balls}] \\
 &= 1 - \Pr\left[\sum_{i=1}^n x_i \geq 1\right] \quad \text{Application of Markov ineq.} \\
 &\geq 1 - \left(\frac{\mathbb{E}\left[\sum_{i=1}^n x_i\right]}{1}\right) \\
 &= 1 - \sum_{i=1}^n \mathbb{E}[x_i] \\
 &\leq 1 - n \cdot \frac{1}{n^2} \\
 &= 1 - \frac{1}{n}.
 \end{aligned}$$

Theorem: With probability of at least $1 - \frac{1}{n}$, no bin has more than $\frac{3 \ln(n)}{\ln \ln(n)}$ balls in it.
 (Max load is $\overset{\uparrow}{k}$ w.p. $1 - o(1)$).

Question: Does there exist a bin with $\frac{\ln n}{3 \ln \ln n}$ balls?

Recall that by Chebyshov's inequality,

$$\Pr[|x - \mathbb{E}[x]| \geq t] \leq \frac{\text{Var}[x]}{t^2}.$$

Remark: Prob that no bin contains more than $k-1$ balls is at most $\Pr[|x - \mathbb{E}[x]| \geq \mathbb{E}[x]]$ and by Chebyshov's this qty is at most $\frac{\text{Var}[x]}{(\mathbb{E}[x])^2}$ where $X = \sum x_i$.

Thus, we need an upper bound on $\text{Var}[x]$ and lower bound on $\mathbb{E}[x]$.

$\mathbb{E}[x_i] =$ prob that bin i contains at least k balls
 \geq prob that bin i contains exactly k balls.

$$= \binom{n}{k} \left(\frac{1}{n}\right)^k \left(1 - \frac{1}{n}\right)^k$$

Set $k = \frac{\ln(n)}{3\ln\ln(n)}$.

$$\geq \left(\frac{n}{k}\right)^k \cdot \frac{1}{n^k} \cdot e^{-\frac{k}{n}}$$

$$\geq \frac{1}{e} \cdot \left(\frac{1}{k}\right)^k$$

$$\geq \frac{1}{e} \left(\frac{3\ln\ln(n)}{\ln n}\right)^{\frac{\ln n}{3\ln\ln n}}$$

$$= \frac{1}{e} \cdot e^{\frac{\ln n}{3\ln\ln n}} \left(\ln(3\ln\ln(n)) - \ln(\ln n)\right)$$

$$= \frac{1}{e} \cdot e^{-\frac{\ln n}{3}} \left(\frac{\ln\ln n - \ln 3 - \ln\ln\ln(n)}{\ln\ln n}\right)$$

$$\approx \frac{1}{e} \cdot \frac{1}{n^{2/3}} \quad \left. \right\}$$

$$\mathbb{E}[x] = \sum_{i=1}^n \mathbb{E}[x_i] \geq \frac{n^{2/3}}{e} \quad \left. \right\}.$$

Claim: x_i and x_j for $i \neq j$ are negatively correlated.

Thus, $\text{Cov}(x_i, x_j) \leq 0$.

Thus, $\text{Var}[x] = \text{Var}\left[\sum_{i=1}^n x_i\right] \leq \sum_{i=1}^n \text{Var}[x_i]$.

$\text{Var}[x_i] = \mathbb{E}[(x_i - \mathbb{E}[x_i])^2] \leq 1$ as $x_i \in \{0, 1\}$. Fact

Thus, $\text{Var}[X] \leq n$. $X = \text{no. of bins with at least } k \text{ balls.}$

$$\begin{aligned} \text{Thus, } \Pr[X=0] &= \\ &\leq \Pr[|X - \mathbb{E}[X]| \geq \mathbb{E}[X]] \\ &\leq \frac{\text{Var}[X]}{(\mathbb{E}[X])^2} \\ &\leq \frac{n}{\left(\frac{n^{4/3}}{e^2}\right)} = \frac{e^2}{n^{1/3}} \end{aligned}$$

$\Pr[X=0]$ is a bad event when you want to count \exists of bin w/ $\geq k$ balls.

Theorem: With a prob of at least $1 - \frac{e^2}{n^{4/3}}$

a bin with $\frac{\ln n}{3 \ln \ln n}$ balls

(Max load is at least k w.p. $1 - o(1)$).

Question: Markov vs Chebyshov (when to use them?)
(Can you get $1 - \frac{e^2}{n^{4/3}}$ kind of a bound from Markov?)

Let $Y = n - X$ be a r.v.

$$\begin{aligned} \Pr[X=0] &= \Pr[Y=n] \\ &= \Pr[Y \geq n] \end{aligned}$$

- Markov states that a non-negative r.v takes values much larger than its expectation with small prob.

- Chebyshov states that a r.v with finite and bounded variance is concentrated around its expectation. Smaller the var \Rightarrow better concentration.

$$\begin{aligned} \text{Success prob} &\leq \frac{\mathbb{E}[Y]}{n} \\ &= \frac{\mathbb{E}[n-X]}{n} \\ &\sim \frac{1}{2 \cdot n^{1/3}} \geq n - \frac{n^{2/3}}{e} \\ \left(1 - \frac{e^2}{n^{4/3}}\right) &= 1 - \frac{1}{e \cdot n^{1/3}}. \end{aligned}$$