

Problem 3:

m balls into n bins

1. X_i : r.v. for bin i -th ball

Show $\Pr[X_i = X_j] = \frac{1}{n}$ for $i \neq j$

$$\Pr[X_i = X_j] \underset{\substack{\uparrow \\ \text{law of total probability}}}{=} \sum_{k=1}^n \Pr[X_i = X_j = k] = \sum_{k=1}^n \Pr[X_i = k \wedge X_j = k] \underset{\substack{\uparrow \quad \uparrow \\ \text{independent events}}}{=} \sum_{k=1}^n \Pr[X_i = k] \cdot \Pr[X_j = k] \\ = \sum_{k=1}^n \frac{1}{n} \cdot \frac{1}{n} = \sum_{k=1}^n \frac{1}{n^2} = \frac{1}{n}$$

2. C r.v. for collisions (collision: two balls are in the same bin)

View this as a Bernoulli trial:

- For every distinct (unordered) pair of balls: "success" if in the same bin
"not success" otherwise

- Success probability $p = \frac{1}{n}$ (see 1.)

- Number of pairs (trials): $N = \binom{m}{2}$

$$\Rightarrow E[C] = pN = \binom{m}{2} \frac{1}{n} \quad (\text{exp. of binomial distribution})$$

$$\text{Var}[C] = N \cdot p \cdot (1-p) = \binom{m}{2} \frac{1}{n} \cdot \left(1 - \frac{1}{n}\right)$$

But: still need to check independence of trials!

(Definition of independence)

Claim: For all distinct pairs of ball (a, b) and (c, d) : $\Pr[X_a = X_b \wedge X_c = X_d] = \Pr[X_a = X_b] \cdot \Pr[X_c = X_d]$

\rightarrow easy to check, similar to 1.

3. Assumption: $m \leq \sqrt{n}$

$$\Pr[C \geq 1] = \Pr\left[C \geq d \cdot \underbrace{\binom{m}{2} \cdot \frac{1}{n}}_{=E[C]}\right] \leq \frac{1}{d} = \frac{\binom{m}{2}}{n} = \frac{m \cdot (m-1)}{2n} \leq \frac{\frac{1}{2}m^2}{n} \leq \frac{\frac{n}{2}}{n} = \frac{1}{2}$$

for $d = \frac{n}{\binom{m}{2}}$ ↑
Markov

4. Assumption $m > 3\sqrt{n}$

$$\binom{m}{2} = \frac{m \cdot (m-1)}{2} \geq \frac{(3\sqrt{n})^2}{2} = \frac{9}{2}n$$

$$\Pr[C \leq 1] = \Pr[-C \geq -1] = \Pr\left[\underbrace{\binom{m}{2} \cdot \frac{1}{n}}_{=E[C]} - C \geq \binom{m}{2} \frac{1}{n} - 1\right]$$

$$\leq \Pr\left[| \underbrace{\binom{m}{2} \frac{1}{n} - C }_{=|C - E[C]|} | \geq \underbrace{\binom{m}{2} \frac{1}{n} - 1}_{d \cdot \text{Var}[C]} \right] \leq \frac{1}{d^2} = \frac{\left(1 - \frac{1}{n}\right)^2}{\left(1 - \frac{n}{\binom{m}{2}}\right)^2}$$

with $d = \frac{1 - \frac{n}{\binom{m}{2}}}{1 - \frac{1}{n}}$

$$\stackrel{n \geq 2}{\leq} \frac{\left(1 - \frac{1}{2}\right)^2}{\left(1 - \frac{n}{\frac{9}{2}n}\right)^2} = \frac{1}{4} \cdot \frac{1}{\left(\frac{7}{9}\right)^2} = \frac{81}{4 \cdot 49} \leq \frac{1}{2}$$

2. Modification of reservoir sampling

Assumption: $i > k$

We prove by induction that $\Pr[e_i \in S \text{ after step } i+t] = p \left(1 - \frac{p}{k}\right)^t$

Base case: $t=0$ $\Pr[e_i \in S \text{ after step } i+0] = p = p \cdot \left(1 - \frac{p}{k}\right)^0$

Inductive step ($t \rightarrow t+1$): $\Pr[E_i \in S \text{ after step } i+t+1] = \Pr[E_i \in S \text{ after step } i+t] \cdot \Pr[E_i \text{ survives}]$

$$\Pr[E_i \text{ survives}] = \underbrace{\Pr[E_{i+t+1} \text{ not kept}]}_{1-p} + \underbrace{\Pr[E_{i+t+1} \text{ kept}]}_p \cdot \underbrace{\Pr[E_i \text{ not replaced when } E_{i+t+1} \text{ kept}]}_{\left(1 - \frac{1}{k}\right) = \frac{k-1}{k}}$$

$$= 1-p + p \cdot \frac{k-1}{k} = 1 - \frac{p}{k}$$

$$\textcircled{*} = \underbrace{p \cdot \left(1 - \frac{p}{k}\right)^t}_{\text{by IH}} \cdot \left(1 - \frac{p}{k}\right) = p \cdot \left(1 - \frac{p}{k}\right)^{t+1}$$

2. exp. number of steps = $\frac{1}{1 - \text{survival prob.}}$

(geometrically distributed)