

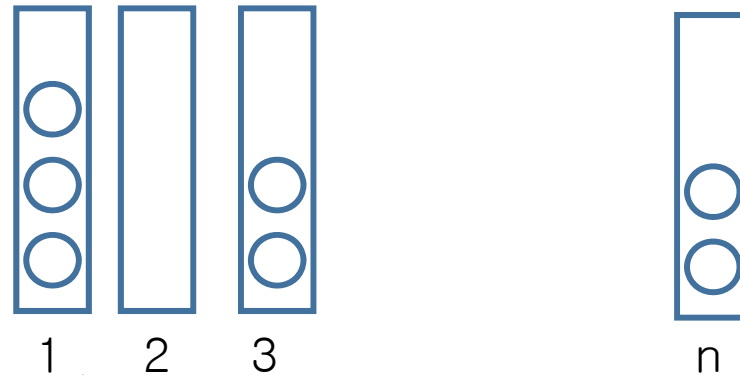
Balls into Bins

Name: Chong-kwon Kim

SCONE
Lab.

Balls & Bins

- Put m balls into n bins uniformly at random



- Same (or similar) problems
 - Birthday paradox
 - Hash table
 - Coupon collection
 - Random allocation of requests to servers
 - Bucket sort,
- We are interested in various statistics such as
 - Number of (non)–empty bins
 - # balls in the most crowded bin
 - Max. # balls w/o collisions, ..

Difficult to obtain
exact values
→ Approximation

Load & Max. Load

- (Max) Load

- Average (Maximum) load: Average (Maximum) # balls in a bin
- If $m=n$, average load = $m/n = 1$

- Assume n balls into n bins. Then

$\Pr(\text{Max. load} \geq (3 \ln n / \ln \ln n)) \leq 1/n$ for large n

- Proof

- Let E_1 : Event that bin1 receives at least M balls

$$\Pr(E_1) = \Pr(X_1 \geq M) < \binom{n}{M} \cdot \left(\frac{1}{n}\right)^M$$

- Use inequalities and obtain $\binom{n}{M} \cdot \left(\frac{1}{n}\right)^M < \frac{1}{M!} < \left(\frac{e}{M}\right)^M$

- Probability that any bin has at least M balls

$$= \Pr\left(\bigcup_{i=1}^n E_i\right) < n \left(\frac{e}{M}\right)^M \quad \leftarrow \text{For } M \geq 3 \ln n / \ln \ln n$$

$$\leq n \left(\frac{e \ln \ln n}{3 \ln n}\right)^{3 \ln n / \ln \ln n}$$

Union bound

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

$$= e^{-2 \ln n + 3 \ln n \cdot \frac{\ln \ln \ln n}{\ln \ln n}}$$

$$\leq 1/n \text{ (for large } n \text{ such as } n = e^{e^{e^e}})$$

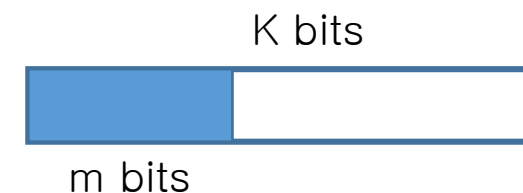
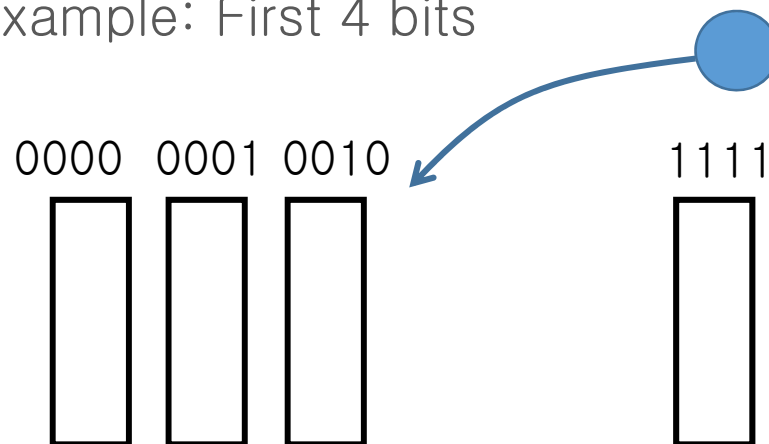
$$\frac{k^k}{k!} < \sum_{i=0}^{\infty} \frac{k^i}{i!} = e^k$$

$$k! > \left(\frac{k}{e}\right)^k$$

Compare to the prior approach that uses (Chernoff) inequalities

Application – Bucket Sort

- Lower bound for average complexity of comparison based sorting is $\Omega(n \log n)$
 - Example: Quick Sort
- Complexity of bucket sort is $O(n)$
- Input: $n = 2^m$ elements each chosen uniformly at random from $[0, 2^K)$ where $K \geq m$
- First, sort based on left-most m bits
 - Example: First 4 bits



Assume national IDs consist of 13 digits
Each ID can be expressed as a 47 bit long number
There are $n \approx 2^{26}$ IDs
Sort based on the first 26 bits

Application – Bucket Sort

1. Prepare n buckets (bins) and put input elements based on first m bits
 - All elements in j th bucket is larger than those in i th bucket if $j > i$
 - If we can put an element into bucket in $O(1) \rightarrow O(n)$ for this stage
2. Sort elements in each bin and combine all bins
 - \rightarrow Total complexity = $O(n)$

• Complexity of stage 2

- X_j : # elements in bucket j

$$X_j: B(n, 1/n) \rightarrow E[X_j] = 1, E[X_j^2] = \frac{n(n-1)}{n^2} + 1 = 2 - 1/n < 2$$

Var[X_j]

- Complexity of sorting X_j elements = $c X_j^2 < c X_j \log X_j$
- $E[\sum_{i=1}^n X_j^2] = 2 \cdot c \cdot n$

BB and Poisson Distribution



- Again m balls into n bins
- X : Number of empty bins

– Let $X_j = \begin{cases} 1, & \text{if bin } j \text{ is empty} \\ 0, & \text{o.w} \end{cases}$

– $E[X_j] = (1 - \frac{1}{n})^m$

– $E[X] = E[\sum_{i=1}^n X_i] = n \cdot (1 - \frac{1}{n})^m \approx n \cdot e^{-m/n}$

Poisson was a prolific French mathematician
First paper at 19 and more than 300 papers

Fraction of empty bins $= e^{-m/n}$

- Generalize \rightarrow Prob. that a bin has exactly r balls, p_r

$$\begin{aligned} p_r &= \binom{m}{r} \cdot \left(\frac{1}{n}\right)^r \cdot \left(1 - \frac{1}{n}\right)^{m-r} \\ &= \frac{1}{r!} \frac{m(m-1) \cdots (m-r+1)}{n^r} \cdot \left(1 - \frac{1}{n}\right)^{m-r} \\ &\approx \frac{(m/n)^r e^{-m/n}}{r!} \quad (\text{For } r \ll m) \end{aligned}$$

Poisson Distribution

- A discrete Poisson random variable X with parameter μ , $\text{Poi}(\mu)$, is

$$\Pr(X=j) = \frac{\mu^j e^{-\mu}}{j!}$$

- Properties

- $\sum_{j=0}^{\infty} \Pr(X = j) = \sum_{j=0}^{\infty} \frac{\mu^j e^{-\mu}}{j!}$
- $E[X] = \sum_{j=0}^{\infty} j \cdot \frac{\mu^j e^{-\mu}}{j!} = \mu$

- Lemma: The sum of finite number of independent Poisson random variables is Poisson

- Proof

- Let X and Y be independent Poisson with means μ_1 and μ_2
- $\Pr(X+Y=j) = \sum_{k=0}^j \Pr(X = k) \cap (Y = j - k) \rightarrow \frac{(\mu_1 + \mu_2)^j e^{-(\mu_1 + \mu_2)}}{j!}$

Any other methods?

Poisson Distribution

- Alternatively, we can show $\text{Poi}(\mu_1) + \text{Poi}(\mu_2) = \text{Poi}(\mu_1 + \mu_2)$ using MGFs

Uniqueness of MGF

- First, show that MGF of a Poisson is

$$M_X(t) = e^{\mu(e^t - 1)}$$

- Let X and Y be two Poisson with means μ_1 and μ_2

$$M_{X+Y}(t) = M_X(t) M_Y(t) = e^{\mu_1(e^t - 1)} e^{\mu_2(e^t - 1)}$$

Bounds on Poisson Distribution

- Let X be a Poisson r.v. with mean μ

- If $x > \mu$, then $\Pr(X \geq x) \leq \frac{e^{-\mu}(e\mu)^x}{x^x}$
- If $x < \mu$, then $\Pr(X \leq x) \leq \frac{e^{-\mu}(e\mu)^x}{x^x}$

- Proof

What method would you use?
Chernoff inequality. & find a proper t

- For any $t > 0$, and $x > \mu$

$$\Pr(X \geq x) = \Pr(e^{tX} \geq e^{tx}) \leq \frac{\mathbf{E}[e^{tX}]}{e^{tx}} = e^{\mu(e^t - 1) - xt}$$

$$\Rightarrow \Pr(X \geq x) \leq e^{x - \mu - x \ln(x/\mu)}$$

$$= \frac{e^{-\mu}(e\mu)^x}{x^x}.$$

Min. at $t = \ln(x/\mu) > 0$

- For any $t < 0$ and $x < \mu$

$$\Pr(X \leq x) = \Pr(e^{tX} \geq e^{tx}) \leq \frac{\mathbf{E}[e^{tX}]}{e^{tx}}$$

Again, let $t = \ln(x/\mu) < 0$

Poisson as Limit of Binomial

- Already showed that Bins and Balls model can be approximated with a Poisson distribution
- Let X_n be a Binomial w/ parameters n and p , where p is function of n and $\lambda = \lim_{n \rightarrow \infty} np$ is a constant and is independent of n . Then for any k

$$\lim_{n \rightarrow \infty} \Pr(X_n = k) = \frac{\lambda^k e^{-\lambda}}{k!}$$

More rigorous proof

- Proof

$$\Pr(X_n = k) = \binom{n}{k} p^k (1-p)^{n-k}$$

$$\Pr(X_n = k) \leq \frac{n^k}{k!} p^k \frac{(1-p)^n}{(1-p)^k}$$

$$\leq \frac{(np)^k}{k!} \frac{e^{-pn}}{1-pk}$$

$$= \frac{e^{-pn} (np)^k}{k!} \frac{1}{1-pk}$$

From Taylor expansion, for $|x| \leq 1$
 $e^x (1-x^2) \leq 1+x \leq e^x$ ($-p$ as x)

For $k > 0$, $(1-p)^k \geq 1-pk$

Poisson as Limit of Binomial

– Also, $\Pr(X_n = k) \geq \frac{(n-k+1)^k}{k!} p^k (1-p)^n$

$e^x(1-x^2) \leq 1+x \leq e^x$

$\geq \frac{((n-k+1)p)^k}{k!} e^{-pn} (1-p^2)^n$

For $k > 0$, $(1-p)^k \geq 1-pk$

$\geq \frac{e^{-pn}((n-k+1)p)^k}{k!} (1-p^2n).$

– Combining, we have

$$\frac{e^{-pn}(np)^k}{k!} \frac{1}{1-pk} \leq \Pr(X_n = k) \leq \frac{e^{-pn}((n-k+1)p)^k}{k!} (1-p^2n)$$

$$\lim_{n \rightarrow \infty} \frac{e^{-pn}(np)^k}{k!} \frac{1}{1-pk} = \frac{e^{-\lambda} \lambda^k}{k!}$$

$$\lim_{n \rightarrow \infty} \frac{e^{-pn}((n-k+1)p)^k}{k!} (1-p^2n) = \frac{e^{-\lambda} \lambda^k}{k!}$$

New Inequality

- We already learnt three bounds
- (Most) Analysis of a single bin can be done by applying the three inequalities
- Analysis of multiple bins requires additional bound
 - Poisson approximation

Poisson Approximation

- In Balls into Bins model, one difficulty is the *dependency* between bins

- If bin i is empty, then the probability that other bins are empty decreases

- Again throw m balls into n bins

- AND, Consider two sets of random variables

$$\{X_1^{(m)}, X_2^{(m)}, \dots, X_n^{(m)}\} \text{ and } \{Y_1^{(m)}, Y_2^{(m)}, \dots, Y_n^{(m)}\}$$

- $X_i^{(m)}$ be the # balls in i -th bin
- $Y_1^{(m)}, Y_2^{(m)}, \dots, Y_n^{(m)}$ are independent Poisson with mean m/n .

Exact Case: # balls in a bin when m balls are thrown to n bins

Poisson Case: # balls in a bin is Poisson with mean m/n

Poisson Approximation

- Theorem: The distribution of $\{Y_1^{(m)}, Y_2^{(m)}, \dots, Y_n^{(m)}\}$ conditioned on $\sum_i Y_i^{(m)} = k$ is the same as $\{X_1^{(k)}, X_2^{(k)}, \dots, X_n^{(k)}\}$, regardless of m .

$Y_i^{(m)}$ is Poisson with mean m/n

$$\forall i, Y_i^{(m)} = k_i \iff X_i^{(k)} = k_i$$

- Proof

- Throwing k balls into n bins, the probability that $\{X_1^{(k)}, X_2^{(k)}, \dots, X_n^{(k)}\} = (k_1, k_2, \dots, k_n)$, $\sum_i k_i = k$ is given by
- $$\frac{\binom{k}{k_1, k_2, \dots, k_n}}{n^k} = \frac{k!}{(k_1!)(k_2!) \cdots (k_n!) n^k}$$
- For any k_1, k_2, \dots, k_n with $\sum_i k_i = k$, consider the probability that $\{Y_1^{(m)}, Y_2^{(m)}, \dots, Y_n^{(m)}\} = (k_1, k_2, \dots, k_n)$ conditioned on $\sum_i Y_i^{(m)} = k$

Poisson Approximation

$$\Pr\left((Y_1^{(m)}, \dots, Y_n^{(m)}) = (k_1, \dots, k_n) \mid \sum_{i=1}^n Y_i^{(m)} = k\right) \\ = \frac{\Pr((Y_1^{(m)} = k_1) \cap (Y_1^{(m)} = k_2) \cap \dots \cap (Y_n^{(m)} = k_n))}{\Pr(\sum_{i=1}^n Y_i^{(m)} = k)}.$$

$$Y_i^{(m)} \sim \text{Poi}(m/n) \quad \Rightarrow \quad \sum_i Y_i^{(m)} \sim \text{Poi}(m)$$

$$\frac{\Pr((Y_1^{(m)} = k_1) \cap (Y_1^{(m)} = k_2) \cap \dots \cap (Y_n^{(m)} = k_n))}{\Pr(\sum_{i=1}^n Y_i^{(m)} = k)} = \frac{\prod_{i=1}^n e^{-m/n} (m/n)^{k_i} / k_i!}{e^{-m} m^k / k!} \\ = \frac{k!}{(k_1!)(k_2!) \dots (k_n!) n^k},$$

Poisson Approximation

Proved that Joint distributions of **Exact Case** and **conditioned Poisson Case** are the same. How about individual random variables? Or more generally

- Let $f(x_1, \dots, x_n)$ be a nonnegative function. Then

$$\mathbf{E}[f(X_1^{(m)}, X_2^{(m)}, \dots, X_n^{(m)})] \leq e\sqrt{m} \cdot \mathbf{E}[f(Y_1^{(m)}, Y_2^{(m)}, \dots, Y_n^{(m)})]$$

- Proof

$$\mathbf{E}[f(Y_1^{(m)}, \dots, Y_n^{(m)})] = \sum_{k=0}^{\infty} \mathbf{E}\left[f(Y_1^{(m)}, \dots, Y_n^{(m)}) \mid \sum_{i=1}^n Y_i^{(m)} = k\right] \Pr\left(\sum_{i=1}^n Y_i^{(m)} = k\right)$$

Count $k=m$ case only

$$\geq \mathbf{E}\left[f(Y_1^{(m)}, \dots, Y_n^{(m)}) \mid \sum_{i=1}^n Y_i^{(m)} = m\right] \Pr\left(\sum_{i=1}^n Y_i^{(m)} = m\right)$$

Previous theorem

$$= \mathbf{E}[f(X_1^{(m)}, \dots, X_n^{(m)})] \Pr\left(\sum_{i=1}^n Y_i^{(m)} = m\right),$$

$$= \mathbf{E}[f(X_1^{(m)}, \dots, X_n^{(m)})] \frac{m^m e^{-m}}{m!}$$

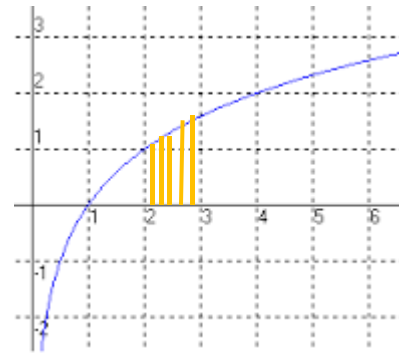
Note: $m! < e\sqrt{m} \left(\frac{m}{e}\right)^m$

$$\mathbf{E}[f(Y_1^{(m)}, \dots, Y_n^{(m)})] \geq \mathbf{E}[f(X_1^{(m)}, \dots, X_n^{(m)})] \frac{1}{e\sqrt{m}}.$$

Poisson Approximation

- To prove that

$$m! < e\sqrt{m}\left(\frac{m}{e}\right)^m$$



- Proof

- Note that $\ln x$ is concave

$$\Rightarrow \int_{i-1}^i \ln x \, dx \geq \frac{\ln(i-1) + \ln i}{2}.$$

$$\int_1^n \ln x \, dx \geq \sum_{i=1}^n \ln i - \frac{\ln n}{2}$$

$$n \ln n - n + 1 \geq \ln(n!) - \frac{\ln n}{2}.$$

- Taking exponentiation to both sides, we obtain the result

Poisson Approximation

- Theorem: Any event that takes place with probability p in the **PC** takes place with probability at most $pe^{\sqrt{m}}$ in the **EC**
- Proof
 - Let f be the indicator function of the event (i.e $f = 1$ if the event occurs, $f = 0$ ow). Then $E[f]$ is the probability that the event occurs.

Poisson Approximation

- Theorem: Let $f(x_1, \dots, x_n)$ be a nonnegative function such that $E[f(X_1^{(m)}, X_2^{(m)}, \dots, X_n^{(m)})]$ is either monotonically increasing or decreasing in m . Then

$$E[f(X_1^{(m)}, X_2^{(m)}, \dots, X_n^{(m)})] \leq 2 E[f(Y_1^{(m)}, Y_2^{(m)}, \dots, Y_n^{(m)})]$$

$$E[f(X_1^{(m)}, X_2^{(m)}, \dots, X_n^{(m)})] \leq E[f(X_1^{(m+1)}, X_2^{(m+1)}, \dots, X_n^{(m+1)})]$$

Exercise 5.13 & 5.14

- Corollary: Let \mathcal{E} be an event whose probability is monotonically increasing (or decreasing) in the number of balls. If \mathcal{E} has probability p in the **PC**, then \mathcal{E} has probability at most $2p$ in the **EC**.

Lemma 5.12

- Claim: Assume n balls into n bins. The maximum load is at least $\ln n / \ln \ln n$ with probability at least $1 - 1/n$.

- Proof:

- $\Pr(\text{Max. load is at least } M) \geq 1 - 1/n$

- $\Pr(\text{Max. load} < M) \leq 1/n$

- Consider Poisson case

- Prob. that bin 1 has at least M balls $\geq 1/eM!$

- Prob. that no bins has $\geq M$ balls $\Rightarrow (1 - 1/eM!)^n \leq e^{-n/eM!}$

Event that Max. load $< M$

- Now, from $\Pr(\text{EC}) \leq e\sqrt{n} \Pr(\text{PC})$,

need to prove that $e\sqrt{n} e^{-n/eM!} < 1/n$

- Sufficient to prove that $e^{-n/eM!} < n^{-2}$

- $M! < n / 2e \ln n$

$M! \leq e\sqrt{M} \left(\frac{M}{e}\right)^M \leq M\left(\frac{M}{e}\right)^M$ when n are suitably large (also M are quite large)

- $\ln M! \leq M \ln M - M + \ln M$

Lemma 5.12

$$\begin{aligned} -\ln M! &\leq M \ln M - M + \ln M \\ &= \frac{\ln n}{\ln \ln n} (\ln \ln n - \ln \ln \ln n) - \frac{\ln n}{\ln \ln n} + (\ln \ln n - \ln \ln \ln n) \\ &\leq \ln n - \frac{\ln n}{\ln \ln n} \\ &\leq \ln n - \ln \ln n - \ln(2e) \end{aligned}$$

Appendix: Revisit Birthday Paradox

- Sequential selection

- First person selects a birthday out of 365 days
- Second person selects a birthday out of (365–1) days
- ..
- i–th person selects a birthday out of (365–i+1) days

- Pr(All m people have different birthdays) (Assume n different days)

$$= \prod_{j=1}^{m-1} \left(1 - \frac{j}{n}\right)$$

$$= \prod_{j=1}^{m-1} e^{-\frac{j}{n}}$$

$$m \ll n$$

$$= \exp\left\{-\sum_{j=1}^{m-1} \frac{j}{n}\right\}$$

$$\approx e^{-m^2/2n}$$

- Let $\text{Pr}(\text{No birthday match}) = \frac{1}{2} \Rightarrow m^2/2n = \ln 2$

$$\Rightarrow m = \sqrt{2n \cdot \ln 2} = \Omega(\sqrt{n})$$

- For $n=365$, $m \approx 23$