

COMMONWEALTH OF AUSTRALIA

Copyright Regulations 1969

WARNING

This material has been reproduced and communicated to you by or on behalf of the University of Sydney pursuant to Part VB of the Copyright Act 1968 (**the Act**). The material in this communication may be subject to copyright under the Act. Any further copying or communication of this material by you may be the subject of copyright protection under the Act.

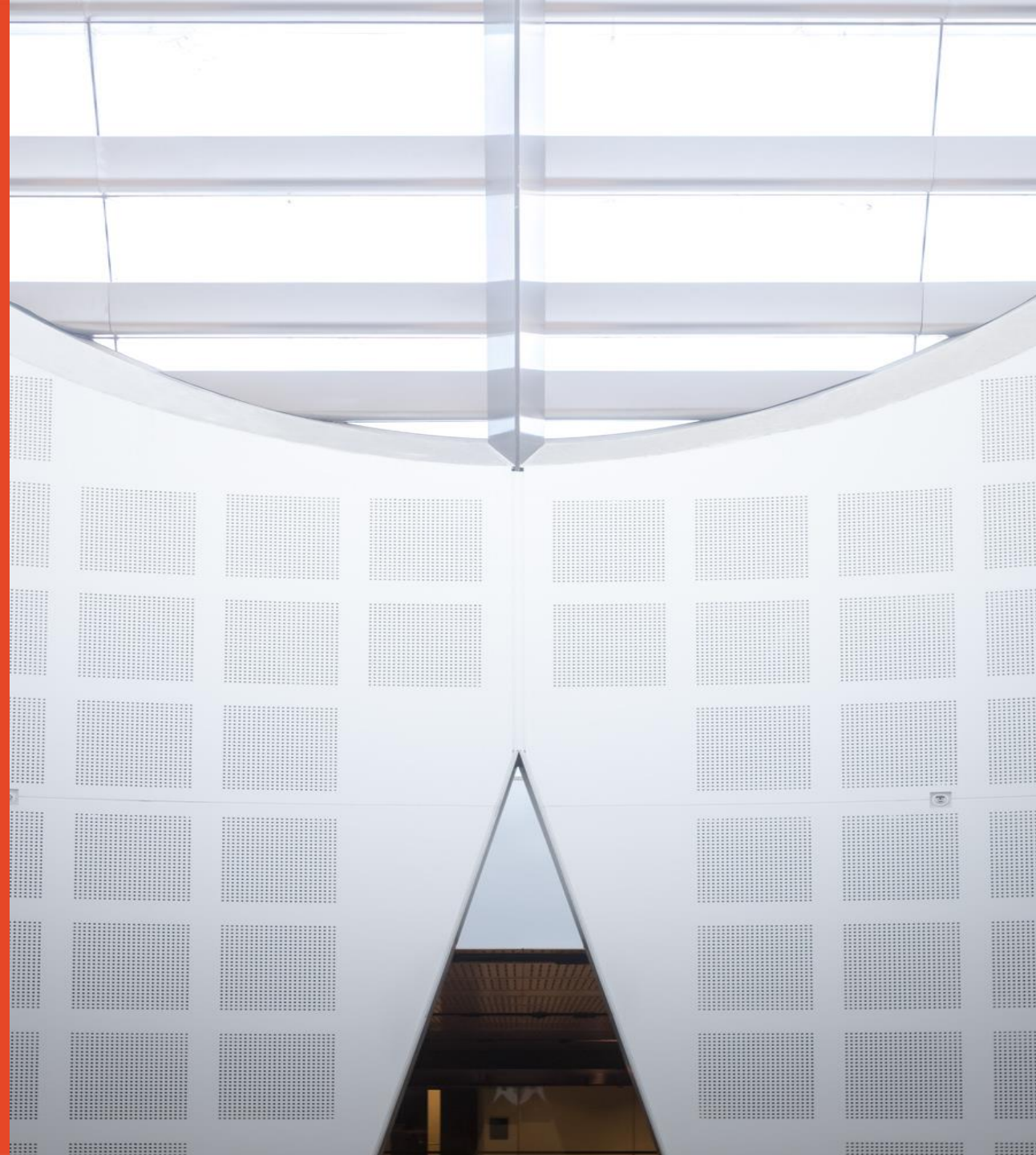
Do not remove this notice.

COMPx270: Randomised and
Advanced Algorithms
Lecture 3: Balls in Bins

Clément Canonne
School of Computer Science




THE UNIVERSITY OF
SYDNEY



A question 🎂

There are quite a few people in the classroom right now. What are the odds two of you (at least) have the same birthday?

A question

Theorem. (The  paradox) If you gather 23 people in a room, then with probability at least 50% a pair will sharing their birthday.

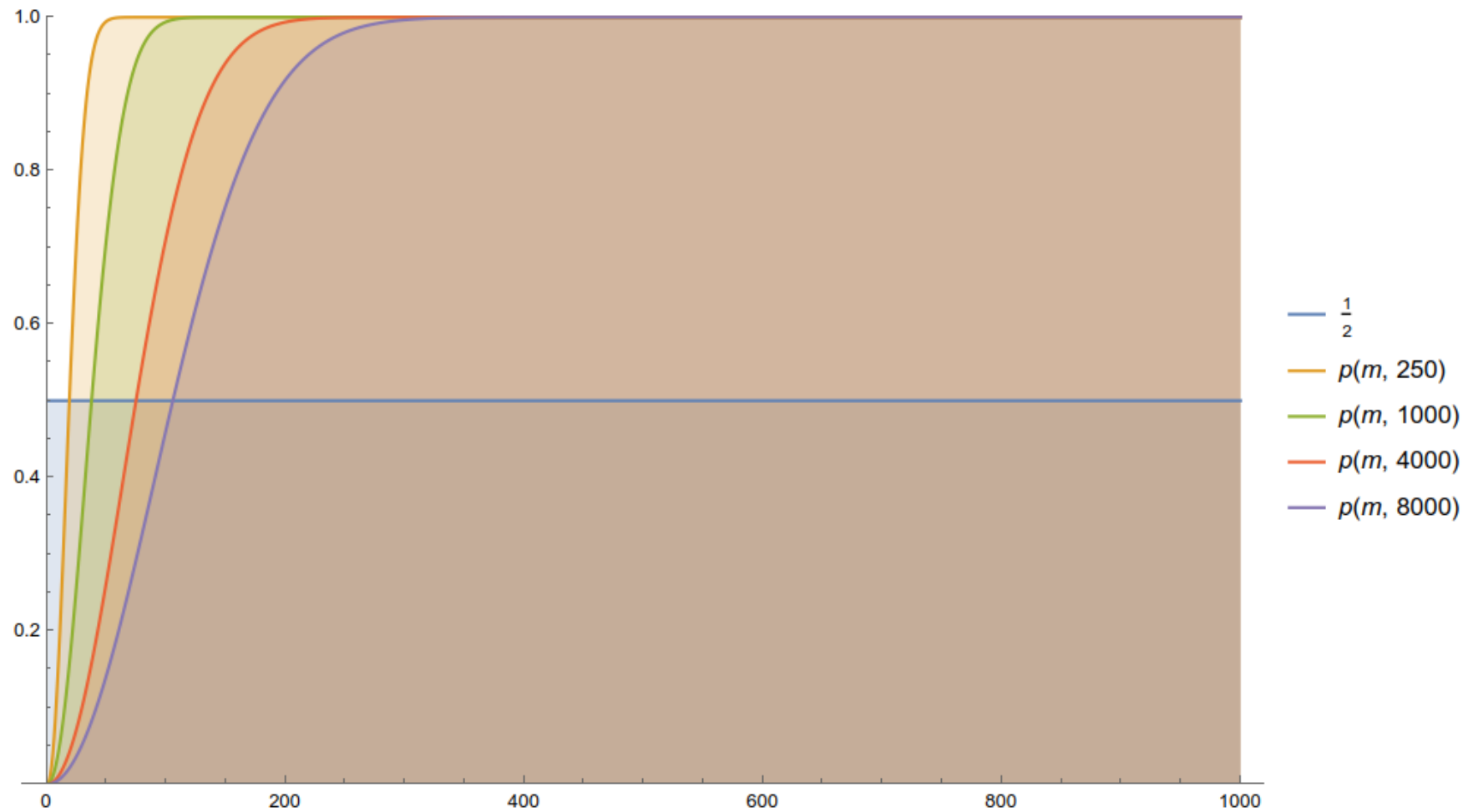
An answer? 🎂

Theorem. If you gather m people and give each a number uniform between 1 and n , then the probability $p(m, n)$ that at least two have the same number is...

$$p_{m,n} = 1 - \frac{n!}{n^m (n - m)!} = 1 - \frac{m!}{n^m} \binom{n}{m}$$

Proof.

```
DiscretePlot[{1/2, p[m, 250], p[m, 1000], p[m, 4000], p[m, 8000]}, {m, 1, 1000}, PlotRange -> {0, 1}, PlotLegends -> "Expressions"]
```



Let's start simple: $m=2$  

Now, for large values of "2"...

C: number of collisions when throwing m  into n . What is

$$c(m, n) = E[C]?$$

Number of collisions

... and what is $\text{Var}[C]$?



Number of collisions

... and what is $\text{Var}[C]$?

If C was Binomial
($C \sim \text{Bin}(\binom{m}{2}, \frac{1}{n})$)

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

C is not Binomial, b/c Y_1, \dots, Y_ℓ not indep^t.

so ...

$$\text{Var} C = \underset{\uparrow ?}{\mathbb{E}[C^2]} - \underset{\uparrow \text{known}}{\mathbb{E}[C]^2}$$

$$C = \sum_{\ell=1}^{\binom{m}{2}} Y_\ell$$

0 if no collision on ℓ^{th} pair, 1 o/w

$$\mathbb{E} C = \frac{\binom{m}{2}}{n}$$



$$M = \binom{m}{2}$$

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

Number of collisions

$$\mathbb{E}[C^2] = \mathbb{E}\left[\left(\sum_{e=1}^M Y_e\right)^2\right] = \mathbb{E}\left[\sum_{e=1}^M \sum_{h=1}^M Y_e Y_h\right]$$

(yay)

$$\sum_{e=1}^M \sum_{h=1}^M \mathbb{E}[Y_e Y_h]$$

$$\mathbb{E}[Y_e Y_h]$$

↑

$$\binom{m}{2} \frac{1}{n}$$

$$+ \binom{m}{2} \binom{m-2}{2} \frac{1}{n^2}$$

$$+ 6 \binom{m}{3} \frac{1}{n^2}$$

$$= \binom{m}{2} \frac{1}{n} \left(1 - \frac{1}{n}\right) + \binom{m}{2} \frac{1}{n^2} \left. \begin{array}{l} \text{omitting} \\ \text{distinct} \\ \text{of cases} \end{array} \right\}$$

$$\text{Var } C = \mathbb{E}[C^2] - \mathbb{E}[C]^2$$

$$= \binom{m}{2} \frac{1}{n} \left(1 - \frac{1}{n}\right) \quad \checkmark$$

Number of collisions

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

Now, we can use Chebyshev:

$$\Pr[X = 0] \leq 1/2$$

for $m = \Omega(\sqrt{n})$.

Number of collisions

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

Now, we can use Chebyshev:

$$\Pr[X = 0] \leq 1/2$$

for $m = \Omega(\sqrt{n})$. **Is it tight?**

Number of collisions

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

Now, we can use Chebyshev:

$$\Pr[X = 0] \leq 1/2$$

for $m = \Omega(\sqrt{n})$.

By Markov, we also have

$$\Pr[X \neq 0] = \Pr[X \geq 1] \leq E[X] \leq 1/2$$

for $m = O(\sqrt{n})$.

Applications?

Bounding the variance: is it always that bad?

Two tricks (and even 3).

① $\text{Var} \sum_i X_i \neq \sum_i \text{Var} X_i$ in general
But "often" $\text{Var} \sum_i X_i \leq \sum_i \text{Var} X_i$
 X_1, \dots, X_n "negatively correlated"

② $\text{Var} X \leq \mathbb{E}[X^2]$ ($\text{Var} X = \mathbb{E}[X^2] - \underbrace{\mathbb{E}[X]^2}_{\geq 0} \leq \mathbb{E}[X^2]$)

(③ "Pairwise independence" is enough
(will get back to it))

Coverage (Coupon collector)



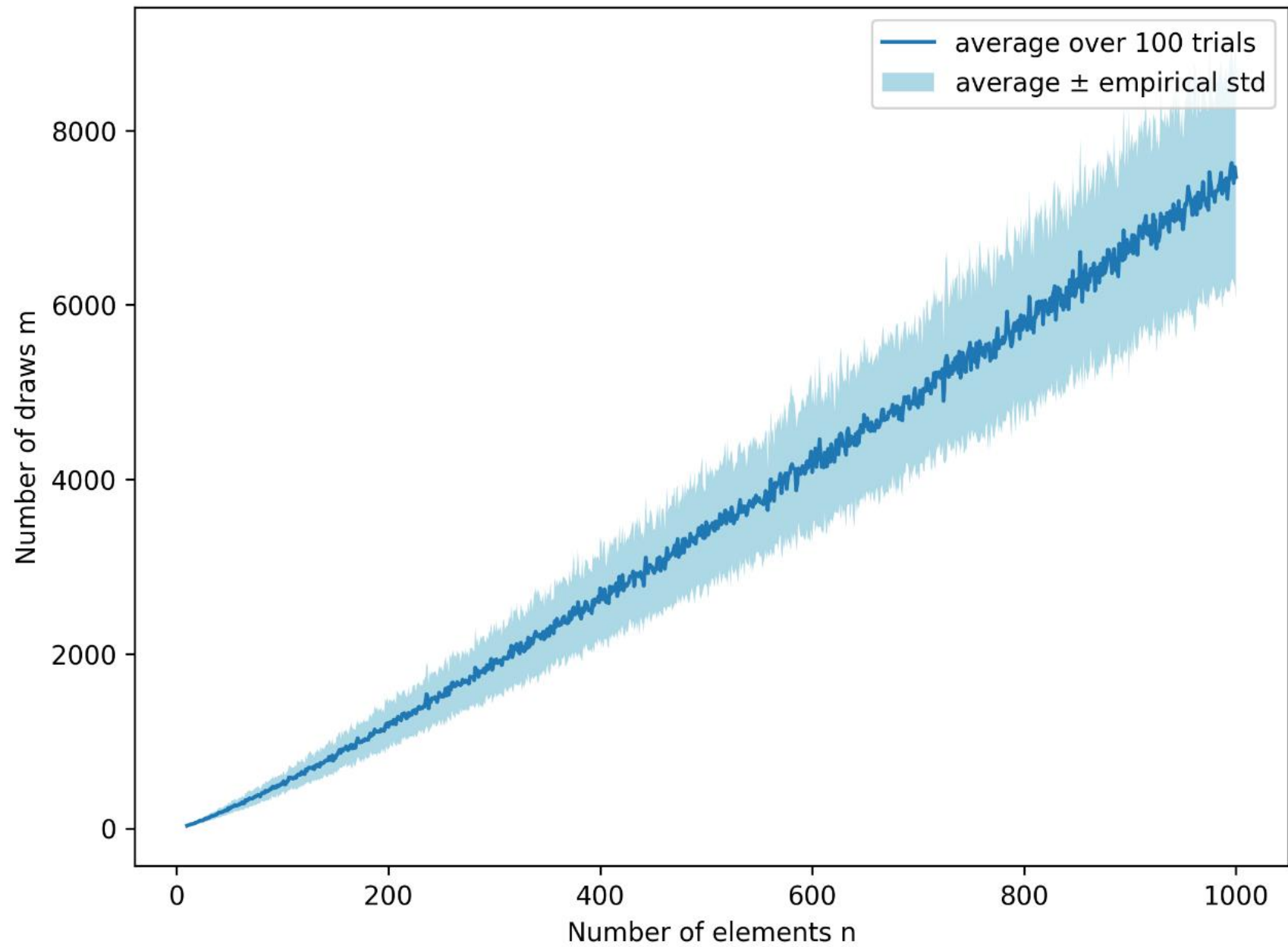
"What is the expected number of balls $M(n)$ we need to throw before each of the n bins contains **at least one** ball?"

Coverage (Coupon collector)



"What is the expected number of balls $M(n)$ we need to throw before each of the n bins contains **at least one** ball?"

- $\Theta(n)$?
- $\Theta(n \log n)$?
- $\Theta(n^2)$?
- Something else?



Coverage (Coupon collector)



"What is the expected number of balls $M(n)$ we need to throw before each of the n bins contains **at least one** ball?"

Theorem. In expectation, $M(n) = \Theta(n \log n)$ balls. (Even more precisely: $n \ln n + O(n)$.)

Intuition: throw n balls. $\Pr[\text{no Pikachu}] = \left(1 - \frac{1}{n}\right)^n \approx \frac{1}{e}$

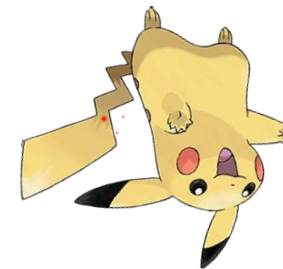
($\Pr[\exists \text{ Pokemon not seen}] \leq_{\text{union}} n \cdot \left(1 - \frac{1}{n}\right)^n \approx \frac{n}{e}$)

$\mathbb{E}[\text{Pokemon not seen}] = n \left(1 - \frac{1}{n}\right)^n \approx \frac{n}{e}$

n more: $\mathbb{E}[\] = \frac{n}{e^2}$. n more: $\frac{n}{e^3}$

Proof.

T_i : conditioned on having $i-1$ Pokémon,
what is the # of Pokéballs until I get
an i^{th} Pokémon (new)



$$T_1 = 1$$

$$M(n) = T_1 + \dots + T_n$$

$$\mathbb{E}[M(n)] = \sum_{i=1}^n \mathbb{E}[T_i]$$

$\Pr[\text{getting } i^{\text{th}} \text{ new} \mid \text{I have } i-1 \text{ distinct Pokémon}] =$

$$\frac{n - (i-1)}{n}$$

Geometric!

g.v.

$$T_i \sim \text{Geom}\left(\frac{n-i+1}{n}\right)$$

$$\mathbb{E}[T_i] = \frac{n}{n-i+1}$$

$$= \sum_{i=1}^n \frac{n}{n-i+1}$$

$$= n \sum_{i=1}^n \frac{1}{n-i+1}$$

$$\stackrel{\text{"}l=n-i+1\text{"}}{=} n \cdot \sum_{l=1}^n \frac{1}{l} = n H_n$$

$$= n(\ln n + O(1))$$

What about the variance?



$$\text{Var } M(n) = \text{Var}(T_1 + \dots + T_n)$$

$$= \sum_{i=1}^n \text{Var } T_i$$

independent

(magic)

$$= \sum_{i=1}^n \frac{n^2}{(n-i+1)^2} \cdot \underbrace{\left(1 - \frac{n-i+1}{n}\right)}_{\leq 1}$$

$$\leq n^2 \sum_{i=1}^n \frac{1}{(n-i+1)^2}$$

$$= n^2 \sum_{e=1}^n \frac{1}{e^2}$$

(n ∞)

$$\leq n^2 \sum_{e=1}^{\infty} \frac{1}{e^2} = \frac{\pi^2}{6} n^2$$

Geom(p)
ratio
variance
 $\frac{1-p}{p^2}$

Chapman

$$M(n) = n \log n \pm O(n)$$

w.h.p.

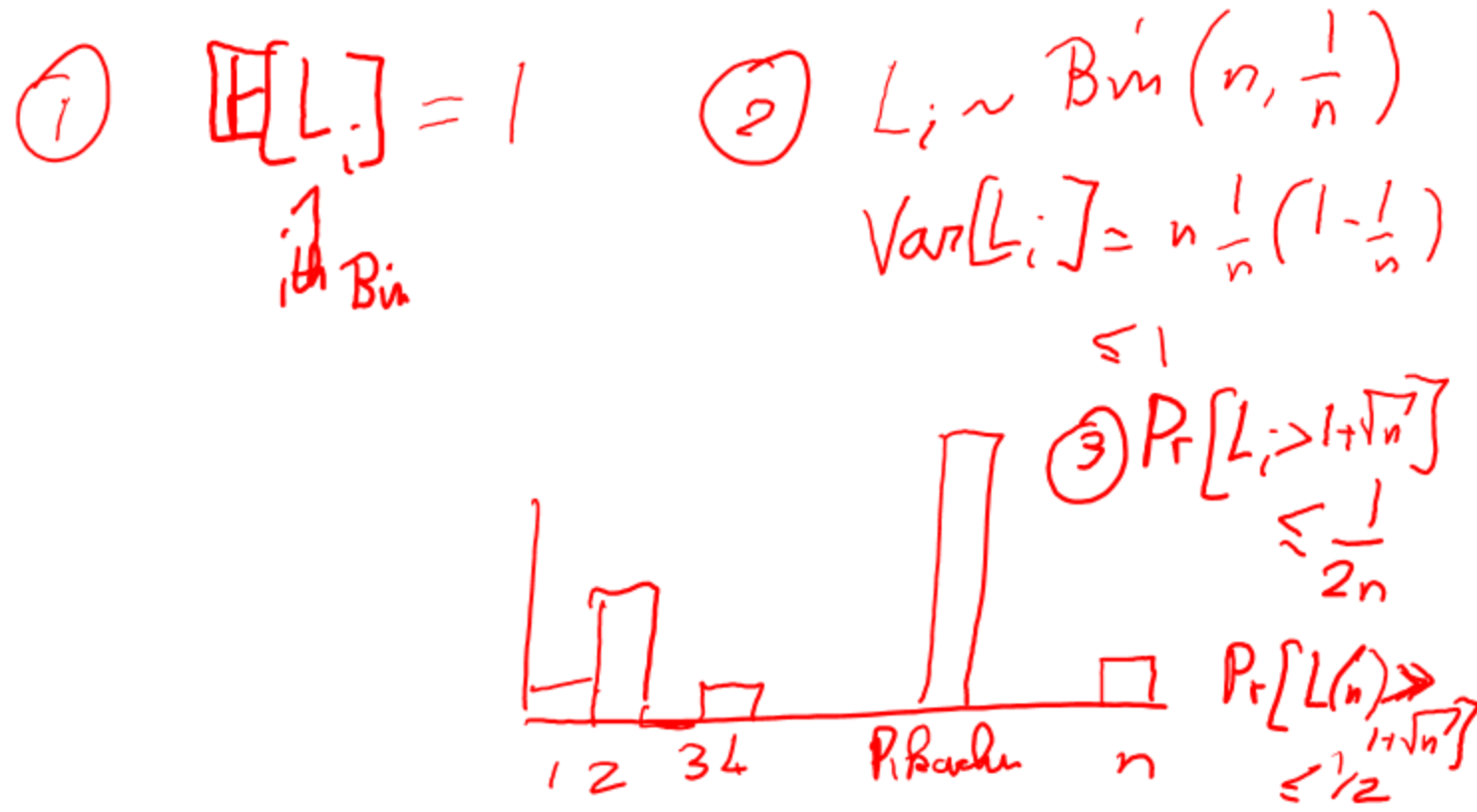
Load balancing

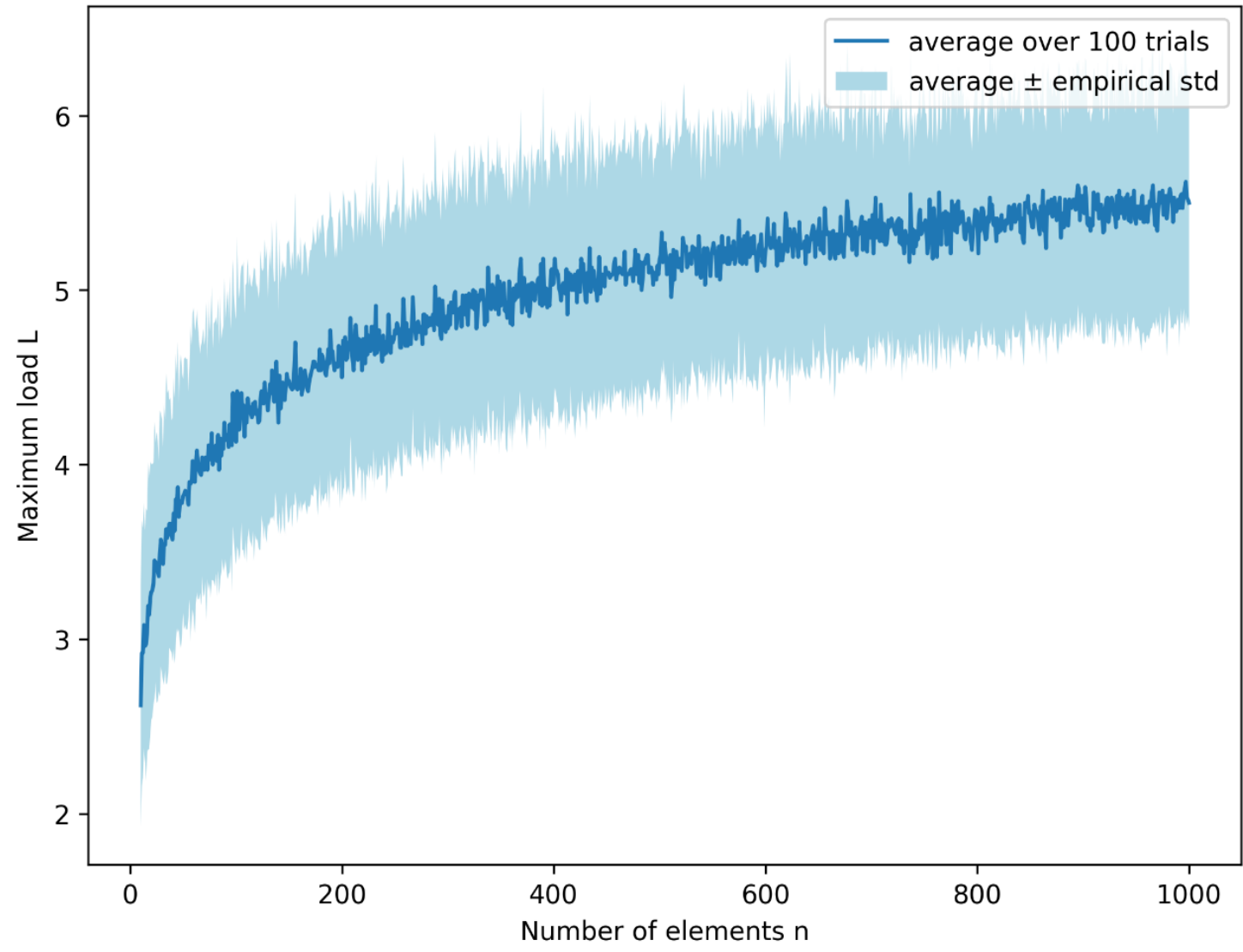
"What is the expected number of balls $L(n)$ the fullest of the n bins contains after throwing n balls?"

Load balancing

"What is the expected number of balls $L(n)$ the fullest of the n bins contains after throwing n balls?"

- $\Theta(1)$?
- $\Theta(\log n)$?
- $\Theta(\sqrt{n})$?
- Something else?





Load balancing

"What is the expected number of balls $L(n)$ the fullest of the n bins contains after throwing n balls?"

Theorem. The expected maximum load is $L(n) = \Theta(\log n / \log \log n)$.

$$\begin{aligned} \mathbb{E}[L] &= \sum_{k=1}^n \Pr[L \geq k] \leq n \sum_{k=1}^n \Pr[L_1 \geq k] \\ \uparrow & \quad \uparrow \\ \max_{1 \leq i \leq n} L_i & \quad \end{aligned}$$
$$\Pr[L \geq k] = \Pr[\exists i, L_i \geq k] \leq \sum_{i=1}^n \Pr[L_i \geq k]$$

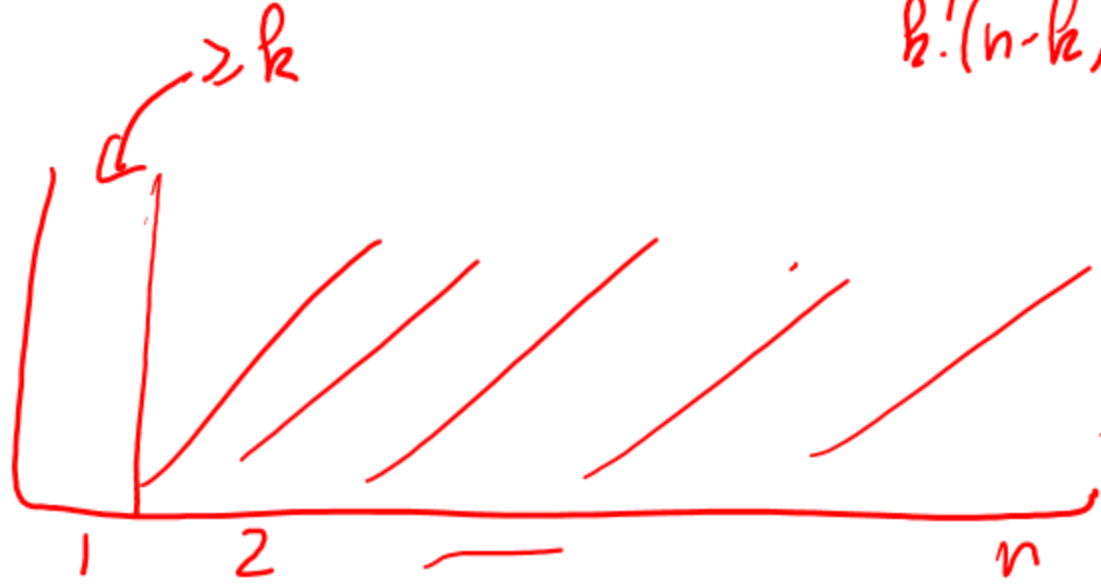
union

Proof.

$1 \leq k \leq n$

$$P[L_1 \geq k] \leq \binom{n}{k} \left(\frac{1}{n}\right)^k$$

$$= \frac{\binom{n}{k} n!}{k!(n-k)! n^k}$$



Proof.

$$\left(\frac{n}{k}\right)^k \leq \binom{n}{k} \leq \left(\frac{e \cdot n}{k}\right)^k$$

$$\Pr[L_1 \geq k] \leq \left(\frac{en}{k}\right)^k \cdot \frac{1}{n^k} = \frac{e^k}{k^k} \quad \checkmark$$

$$\mathbb{E}[L] \leq \sum_{k=1}^n \Pr[L \geq k] \leq \sum_{k=1}^n \min\left(\frac{ne^k}{k^k}, 1\right)$$

$$\leq \sum_{k=1}^l 1 + \sum_{k=l+1}^n \frac{ne^k}{k^k} \leq l + \sum_{k=l+1}^{\infty} \frac{ne^k}{k^k}$$

Pick l s.t

$$\frac{ne^k}{k^k} < \frac{1}{2^k} \text{ for } k \geq l+1$$

$$\leq l + o(1)$$

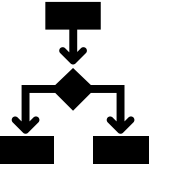
$$\leq O\left(\frac{\log n}{\log \log n}\right)$$

$$\text{Need } 2^l ne^l \leq l^l$$

$$l = O\left(\frac{\log n}{\log \log n}\right) \text{ works}$$

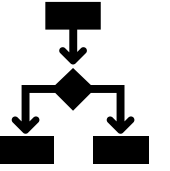
$$\left. \begin{array}{l} \Pr[L \geq k] \\ \leq \frac{ne^k}{k^k} \\ \Pr[L \geq k] \leq 1 \end{array} \right\}$$

Load balancing (a twist)



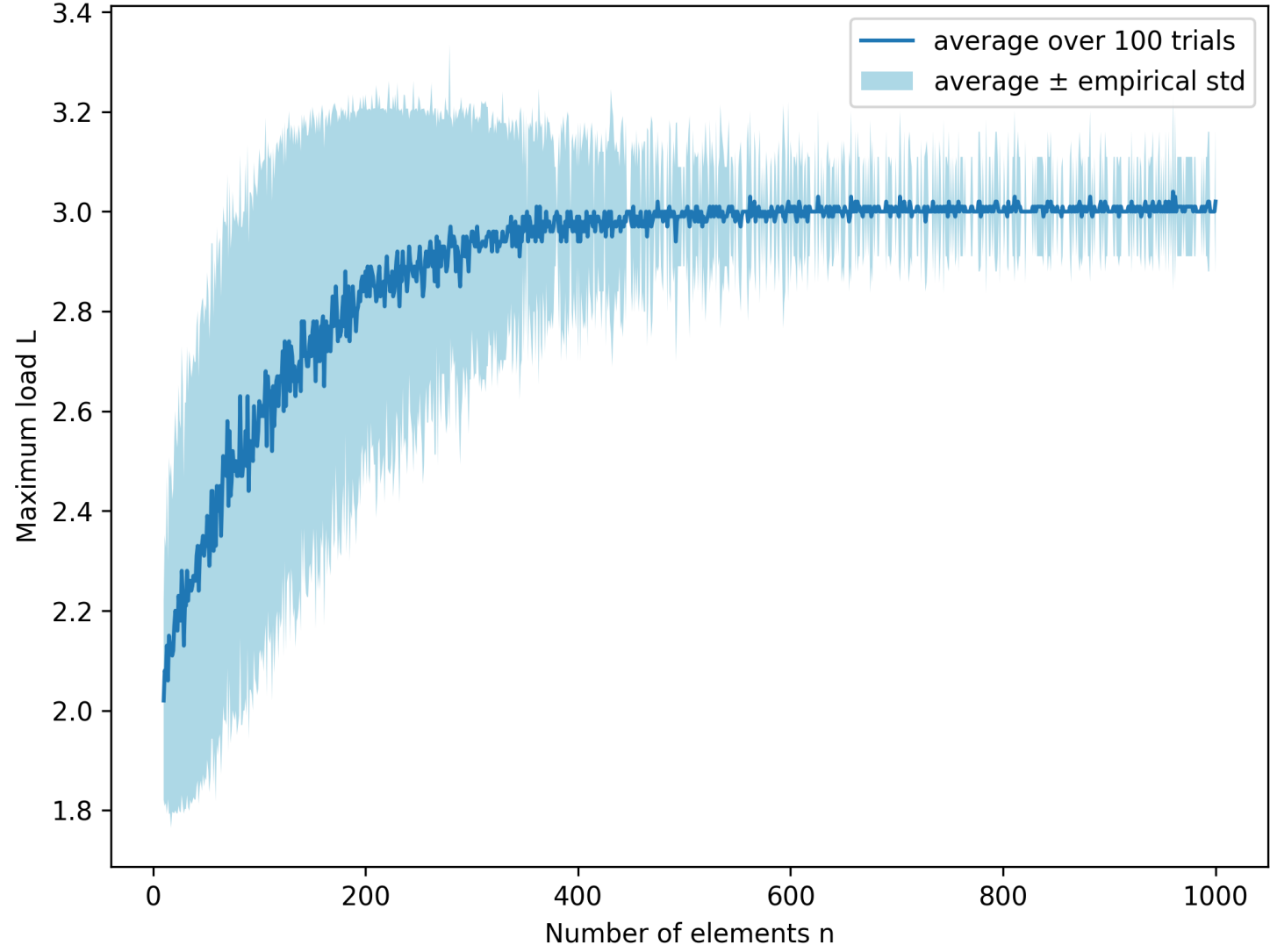
"Now, every time you throw a ball, it selects **two** bins at random, and goes to the least full of the two. What is the maximum expected load?"

Load balancing (a twist)

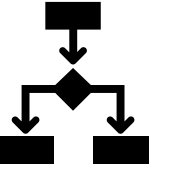


"Now, every time you throw a ball, it selects **two** bins at random, and goes to the least full of the two. What is the maximum expected load?"

- $\Theta(1)$?
- $\Theta(\log \log n)$?
- $\Theta(\sqrt{\log n})$?
- Something else?



Load balancing (a twist)



"Now, every time you throw a ball, it selects **two** bins at random, and goes to the least full of the two. What is the maximum expected load?"

Theorem. The expected maximum load now $\Theta(\log \log n)$.

