Lecture 4: Balls & Bins

Rafael Oliveira

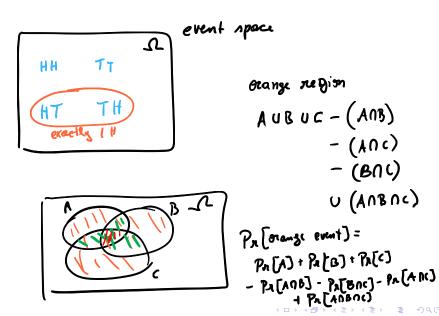
University of Waterloo Cheriton School of Computer Science rafael.oliveira.teaching@gmail.com

May 20, 2021

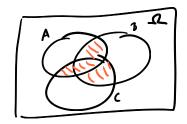
Overview

- Introduction
 - Probability basic notions
 - Balls and Bins
 - Analyses
- Coupon Collector and Power of Two Choices
 - Coupon Collector
 - Power of Two Choices
- Acknowledgements

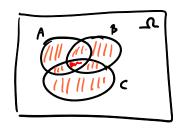
Event Spaces and Inclusion-Exclusion



Union Bound and Inclusion-Exclusion



Union Bound and Inclusion-Exclusion



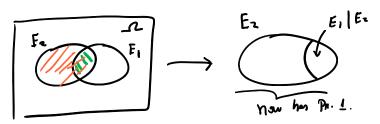
agriculties this to any level of intersections.

• The conditional probability of E_1 given E_2 is

$$\Pr[E_1 \mid E_2] := \frac{\Pr[E_1 \cap E_2]}{\Pr[E_2]}$$

probability

that E_1 happens
given that E_2 happened

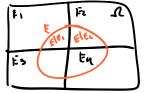


• The *conditional probability* of E_1 given E_2 is

$$\Pr[E_1 \mid E_2] := \frac{\Pr[E_1 \cap E_2]}{\Pr[E_2]}$$

• If E_1, \ldots, E_k partition our sample space, then for any event E

$$Pr[E] = \sum_{i=1}^{K} Pr[E \mid E_i] \cdot Pr[E_i]$$



• The *conditional probability* of E_1 given E_2 is

$$\Pr[E_1 \mid E_2] := \frac{\Pr[E_1 \cap E_2]}{\Pr[E_2]}$$

• If E_1, \ldots, E_k partition our sample space, then for any event E

$$Pr[E] = \sum_{i=1}^{k} Pr[E \mid E_i] \cdot Pr[E_i]$$
le:

• Simple Bayes' rule:

$$Pr[E_1 \mid E_2] = \frac{Pr[E_2 \mid E_1] \cdot Pr[E_1]}{Pr[E_2]}$$

• The *conditional probability* of E_1 given E_2 is

$$\Pr[E_1 \mid E_2] := \frac{\Pr[E_1 \cap E_2]}{\Pr[E_2]}$$

• If E_1, \ldots, E_k partition our sample space, then for any event E

$$\Pr[E] = \sum_{i=1}^{k} \Pr[E \mid E_i] \cdot \Pr[E_i]$$

Simple Bayes' rule:

$$Pr[E_1 \mid E_2] = \frac{Pr[E_2 \mid E_1] \cdot Pr[E_1]}{Pr[E_2]}$$

• Bayes' rule: E_1, \ldots, E_k partition our sample space then for event E

$$\Pr[E_i \mid E] = \frac{\Pr[E \cap E_i]}{\Pr[E]} = \frac{\Pr[E \mid E_i] \cdot \Pr[E_i]}{\sum_{i=1}^k \Pr[E \mid E_j] \cdot \Pr[E_j]}$$

Setup: we have m balls and we want to put them in n bins.

As it is NBA playoff season, we will do this by throwing each ball into a *uniformly random* bin *independently*.

Setup: we have m balls and we want to put them in n bins.

As it is NBA playoff season, we will do this by throwing each ball into a *uniformly random* bin *independently*.

Setup: we have m balls and we want to put them in n bins.

As it is NBA playoff season, we will do this by throwing each ball into a *uniformly random* bin *independently*.

We are interested in the following questions:

• What is the expected number of balls in a bin?

Setup: we have m balls and we want to put them in n bins.

As it is NBA playoff season, we will do this by throwing each ball into a *uniformly random* bin *independently*.

- What is the expected number of balls in a bin?
- What is the expected number of empty bins?

Setup: we have m balls and we want to put them in n bins.

As it is NBA playoff season, we will do this by throwing each ball into a *uniformly random* bin *independently*.

- What is the expected number of balls in a bin?
- What is the expected number of empty bins?
- What is "typically" the maximum number of balls in any bin?

Setup: we have m balls and we want to put them in n bins.

As it is NBA playoff season, we will do this by throwing each ball into a *uniformly random* bin *independently*.

- What is the expected number of balls in a bin?
- What is the *expected* number of empty bins?
- What is "typically" the maximum number of balls in any bin?
- What is the *expected* number of bins with *k* balls in them?

Setup: we have m balls and we want to put them in n bins.

As it is NBA playoff season, we will do this by throwing each ball into a *uniformly random* bin *independently*.

- What is the expected number of balls in a bin?
- What is the expected number of empty bins?
- What is "typically" the maximum number of balls in any bin?
- What is the expected number of bins with k balls in them?
- For what values of *m* do we expect to have *no empty bins*? (coupon collector)

Why Learn About Balls and Bins?

In next lectures, we are going to learn about and analyse *randomized algorithms*. While we will usually analyse the *expected running times* of the algorithms, we would also like to know if the algorithm runs in time close to its expected running time *most of the time*.

Why Learn About Balls and Bins?

In next lectures, we are going to learn about and analyse *randomized algorithms*. While we will usually analyse the *expected running times* of the algorithms, we would also like to know if the algorithm runs in time close to its expected running time *most of the time*.

Running time *small* with *high probability better than* small expected running time.

Why Learn About Balls and Bins?

In next lectures, we are going to learn about and analyse *randomized algorithms*. While we will usually analyse the *expected running times* of the algorithms, we would also like to know if the algorithm runs in time close to its expected running time *most of the time*.

Running time *small* with *high probability better than* small expected running time.

In **this lecture**, we will analyse random processes (*balls & bins*) which underlie several randomized algorithms!

Applications ranging from:

- data structures
- outing in parallel computers
- many more!

$$\mathbb{E}[\# \text{ balls in bin } j] = \mathbb{E}\left[\sum_{i=1}^m B_{ij}\right]$$

$$\mathbb{E}[\# \text{ balls in bin } j] = \mathbb{E}\left[\sum_{i=1}^m B_{ij}\right]$$

$$= \sum_{i=1}^m \mathbb{E}\left[B_{ij}\right] \qquad \qquad \text{(linearity of expectation)}$$

$$\mathbb{E}[\# \text{ balls in bin } j] = \mathbb{E}\left[\sum_{i=1}^{m} B_{ij}\right]$$

$$= \sum_{i=1}^{m} \mathbb{E}\left[B_{ij}\right] \qquad \text{(linearity of expectation)}$$

$$= \sum_{i=1}^{m} \Pr\left[\text{ball } i \text{ in bin } j\right]$$

$$\Pr\left[B_{ij} = i\right] \cdot 4 + O \cdot \Pr\left[B_{ij} = 0\right]$$

$$\mathbb{E}[\# \text{ balls in bin } j] = \mathbb{E}\left[\sum_{i=1}^{m} B_{ij}\right]$$

$$= \sum_{i=1}^{m} \mathbb{E}\left[B_{ij}\right] \qquad \text{(linearity of expectation)}$$

$$= \sum_{i=1}^{m} \Pr\left[\text{ball } i \text{ in bin } j\right]$$

$$= \sum_{i=1}^{m} \frac{1}{n} = \frac{m}{n} \qquad \text{(uniformly at random)}$$

Let us label the m balls $1, \ldots, m$, and the n bins $1, 2, \ldots, n$. Let B_{ij} be the indicator variable that ball i was thrown into bin j.

$$\mathbb{E}[\# \text{ balls in bin } j] = \mathbb{E}\left[\sum_{i=1}^m B_{ij}\right]$$

$$= \sum_{i=1}^m \mathbb{E}\left[B_{ij}\right] \qquad \text{(linearity of expectation)}$$

$$= \sum_{i=1}^m \Pr\left[\text{ball } i \text{ in bin } j\right]$$

$$= \sum_{i=1}^m \frac{1}{n} = \frac{m}{n} \qquad \text{(uniformly at random)}$$

When m = n, expectation of one ball per bin. How often will this actually happen?

$$N_i = \begin{cases} 1 & \text{if bin } i \text{ is empty} \\ 0 & \text{otherwise} \end{cases}$$

$$\mathbb{E}[\# ext{ empty bins}] = \mathbb{E}\left[\sum_{i=1}^n N_i
ight]$$

$$\mathbb{E}[\# \text{ empty bins}] = \mathbb{E}\left[\sum_{i=1}^n N_i\right]$$

$$= \sum_{i=1}^n \mathbb{E}\left[N_i\right] \qquad \qquad \text{(linearity of expectation)}$$

Let N_i be the indicator variable that bin i is empty.

$$\mathbb{E}[\# \text{ empty bins}] = \mathbb{E}\left[\sum_{i=1}^{n} N_{i}\right]$$

$$= \sum_{i=1}^{n} \mathbb{E}[N_{i}] \qquad \text{(linearity of expectation)}$$

$$= \sum_{i=1}^{n} \Pr[\text{bin } i \text{ is empty}]$$

$$\text{No ball landed in bin i}$$

$$\text{Pr}\left[B_{ki} = 0\right] = 1 - \frac{1}{n}$$
by independence of ball throwing

29 / 90

$$\mathbb{E}[\# \text{ empty bins}] = \mathbb{E}\left[\sum_{i=1}^{n} N_i\right]$$

$$= \sum_{i=1}^{n} \mathbb{E}[N_i] \qquad \text{(linearity of expectation)}$$

$$= \sum_{i=1}^{n} \Pr[\text{bin } i \text{ is empty}]$$

$$= \sum_{i=1}^{n} (1 - 1/n)^m$$

$$= n \cdot (1 - 1/n)^m \approx n \cdot e^{-m/n}$$

Let N_i be the indicator variable that bin i is empty.

$$\mathbb{E}[\# \text{ empty bins}] = \mathbb{E}\left[\sum_{i=1}^n N_i\right]$$

$$= \sum_{i=1}^n \mathbb{E}\left[N_i\right] \qquad \text{(linearity of expectation)}$$

$$= \sum_{i=1}^n \Pr\left[\text{bin } i \text{ is empty}\right]$$

$$= \sum_{i=1}^n \left(1 - 1/n\right)^m$$

$$= n \cdot \left(1 - 1/n\right)^m \approx n \cdot e^{-m/n}$$

When m = n, expected fraction of empty bins is $\frac{1}{e}$.

When m = n, first calculation had expectation of one ball per bin.

When m = n, first calculation had expectation of one ball per bin.

When m = n, second calculation had expectation of 1/e fraction of empty bins.

When m = n, first calculation had expectation of one ball per bin.

When m = n, second calculation had expectation of 1/e fraction of empty bins.

Which expectation should I actually "expect"?

When m = n, first calculation had expectation of one ball per bin.

When m = n, second calculation had expectation of 1/e fraction of empty bins.

Which expectation should I actually "expect"?

As we mentioned earlier, this is where *concentration of probability measure* tries to address. It turns out that the *second random variable* (and thus second calculation) is concentrated around the mean (i.e., expectation).

So we "expect" (or it is "typical") to see around 1/e-fraction of empty bins when m=n

Maximum load in a bin

What is the "typical" maximum number of balls in a bin?

As we saw in the previous slide, "typical" is related to concentration of probability measure.

Maximum load in a bin

What is the "typical" maximum number of balls in a bin?

As we saw in the previous slide, "typical" is related to concentration of probability measure.

Let us first see a simpler problem, which is known as the *birthday paradox*: for what value of m do we expect to see two balls in one bin?

Birthday Paradox

The probability that there are no collisions after we have thrown m balls is:

$$1 \cdot \left(1 - \frac{1}{n}\right) \cdot \left(1 - \frac{2}{n}\right) \cdot \cdots \cdot \left(1 - \frac{m-1}{n}\right) \leq e^{-1/n} \cdot \cdots \cdot e^{-\frac{m-1}{n}} \approx e^{\frac{-m^2}{2n}}$$

$$\int_{n}^{n} \int_{n}^{n} dx \, dx \, dx \, dx \, dx \, dx$$

In a call (six)

half

Birthday Paradox

The probability that there are no collisions after we have thrown m balls is:

$$1 \cdot \left(1 - \frac{1}{n}\right) \cdot \left(1 - \frac{2}{n}\right) \cdot \dots \cdot \left(1 - \frac{m-1}{n}\right) \le e^{-1/n} \cdot \dots \cdot e^{-\frac{m-1}{n}} \approx e^{\frac{-m^2}{2n}}$$

This is $\leq 1/2$ when $m = \sqrt{2n \ln(2)}$. For n = 365, this is $m \approx 22.4$ for the probability that two people *(balls)* have birthday on the same date *(bins)* to become $\geq 1/2$.

Birthday Paradox

The probability that there are no collisions after we have thrown m balls is:

$$1 \cdot \left(1 - \frac{1}{n}\right) \cdot \left(1 - \frac{2}{n}\right) \cdot \dots \cdot \left(1 - \frac{m-1}{n}\right) \leq e^{-1/n} \cdot \dots \cdot e^{-\frac{m-1}{n}} \approx e^{\frac{-m^2}{2n}}$$

This is $\leq 1/2$ when $m = \sqrt{2n \ln(2)}$. For n = 365, this is $m \approx 22.4$ for the probability that two people *(balls)* have birthday on the same date *(bins)* to become $\geq 1/2$.

Thus, expect to see collision (two balls in the same bin) when $m = \Theta(\sqrt{n})$. This appears in several places:

- hashing
- factoring
- many more

```
\Pr[\text{bin 1 has } \geq k \text{ balls}] \leq \sum_{\substack{S \text{ subset}[n] \ i \in S \\ |S| = k}} \prod_{i \in S} \Pr[\text{ball } i \text{ in bin 1}]
```

$$\begin{aligned} \Pr[\text{bin 1 has } \geq k \text{ balls}] &\leq \sum_{\substack{S \text{ } \textit{subset}[n] \\ |S| = k}} \prod_{i \in S} \Pr[\text{ball } i \text{ in bin 1}] \\ &= \sum_{\substack{S \text{ } \textit{subset}[n] \\ |S| = k}} \prod_{i \in S} \frac{1}{n} \end{aligned}$$

$$\begin{aligned} \Pr[\text{bin 1 has } \geq k \text{ balls}] &\leq \sum_{\substack{S \text{ subset}[n] \\ |S| = k}} \prod_{i \in S} \Pr[\text{ball } i \text{ in bin 1}] \\ &= \sum_{\substack{S \text{ subset}[n] \\ |S| = k}} \prod_{i \in S} \frac{1}{n} \\ &= \binom{n}{k} \cdot \frac{1}{n^k} \leq \left(\frac{ne}{k}\right)^k \cdot \frac{1}{n^k} = \frac{e^k}{k^k} \end{aligned}$$

What is the probability that a particular bin (say bin 1) has $\geq k$ balls in it?

$$\begin{aligned} \Pr[\text{bin 1 has } \geq k \text{ balls}] &\leq \sum_{\substack{S \text{ } \textit{subset}[n] \\ |S| = k}} \prod_{i \in S} \Pr[\text{ball } i \text{ in bin 1}] \\ &= \sum_{\substack{S \text{ } \textit{subset}[n] \\ |S| = k}} \prod_{i \in S} \frac{1}{n} \\ &= \binom{n}{k} \cdot \frac{1}{n^k} \leq \left(\frac{ne}{k}\right)^k \cdot \frac{1}{n^k} = \frac{e^k}{k^k} \end{aligned}$$

By union bound

$$\Pr[\text{some bin has } \ge k \text{ balls}] \le \sum_{i=1}^n \Pr[\text{bin i has } \ge k \text{ balls}] \le n \cdot \frac{e^k}{k^k}$$

$$\Pr[\text{some bin has } \ge k \text{ balls}] \le n \cdot \frac{e^k}{k^k} = e^{\ln n + k - k \ln k}$$

$$\Pr[\text{some bin has } \ge k \text{ balls}] \le n \cdot \frac{e^k}{k^k} = e^{\ln n + k - k \ln k}$$

$$\Pr[\max | \text{load is } \leq k] = 1 - \Pr[\text{some bin has } > k \text{ balls}] \geq 1 - e^{\ln n + k - k \ln k}$$
 all bias have $\leq k$ balls

$$\Pr[\text{some bin has } \ge k \text{ balls}] \le n \cdot \frac{e^k}{k^k} = e^{\ln n + k - k \ln k}$$

$$\Pr[\max \text{ load is } \leq k] = 1 - \Pr[\text{some bin has } > k \text{ balls}] \geq 1 - e^{\ln n + k - k \ln k}$$

When will the above probability be large (say >> 1/2)?

$$\Pr[\text{some bin has } \ge k \text{ balls}] \le n \cdot \frac{e^k}{k^k} = e^{\ln n + k - k \ln k}$$

$$Pr[\max load is \le k] = 1 - Pr[\text{some bin has } > k \text{ balls}] \ge 1 - e^{\ln n + k - k \ln k}$$

When will the above probability be large (say >> 1/2)?

When
$$k \ln k > \ln n$$
. Setting $k = 3 \frac{\ln n}{\ln \ln n}$ does it.

$$\Pr[\text{some bin has } \ge k \text{ balls}] \le n \cdot \frac{e^k}{k^k} = e^{\ln n + k - k \ln k}$$

$$\Pr[\max \text{ load is } \leq k] = 1 - \Pr[\text{some bin has } > k \text{ balls}] \geq 1 - e^{\ln n + k - k \ln k}$$

When will the above probability be large (say >> 1/2)?

When
$$k \ln k > \ln n$$
. Setting $k = 3 \frac{\ln n}{\ln \ln n}$ does it.

With high probability, max load is $O\left(\frac{\ln n}{\ln \ln n}\right)$.

$$\Pr[\text{some bin has } \ge k \text{ balls}] \le n \cdot \frac{e^k}{k^k} = e^{\ln n + k - k \ln k}$$

$$\Pr[\max \text{ load is } \leq k] = 1 - \Pr[\text{some bin has } > k \text{ balls}] \geq 1 - e^{\ln n + k - k \ln k}$$

When will the above probability be large (say >> 1/2)?

When
$$k \ln k > \ln n$$
. Setting $k = 3 \frac{\ln n}{\ln \ln n}$ does it.

With high probability, max load is
$$O\left(\frac{\ln n}{\ln \ln n}\right)$$
.

This comes up in hashing and in analysis of approximation algorithms (for instance, best known approximation ratio for congestion minimization).

- Introduction
 - Probability basic notions
 - Balls and Bins
 - Analyses

- Coupon Collector and Power of Two Choices
 - Coupon Collector
 - Power of Two Choices

Acknowledgements

For what value of m do we expect to have no empty bins?

For what value of m do we expect to have no empty bins?

Why is this problem called the coupon collector problem?

Because we can formulate it in the following way:

- suppose each bin is a different coupon
- we buy one coupon at random (like kinder eggs/pack action cards)
- what is the number of coupons that we need to buy to collect all of them?

For what value of m do we expect to have no empty bins?

Why is this problem called the coupon collector problem?

Because we can formulate it in the following way:

- suppose each bin is a different coupon
- we buy one coupon at random (like kinder eggs/pack action cards)
- what is the number of coupons that we need to buy to collect all of them?

Let X_i be the number of balls thrown to get from i+1 empty bins to i empty bins. Let X be the number of balls thrown until we have no empty bins.

$$X = \sum_{i=0}^{n-1} X_i$$

- $X_i \leftarrow \#$ balls thrown to get from i empty bins to i-1 empty bins
- $X \leftarrow \#$ balls thrown until we have no empty bins

- ullet $X_i \leftarrow \#$ balls thrown to get from i empty bins to i-1 empty bins
- X ← # balls thrown until we have no empty bins

$$\mathbb{E}[X] = \mathbb{E}\left[\sum_{i=1}^{n} X_{i}\right] = \sum_{i=1}^{n} \mathbb{E}[X_{i}]$$
Inearity of expectation

- $X_i \leftarrow \#$ balls thrown to get from i empty bins to i-1 empty bins
- $X \leftarrow \#$ balls thrown until we have no empty bins

$$\mathbb{E}[X] = \mathbb{E}\left[\sum_{i=1}^{n} X_{i}\right] = \sum_{i=1}^{n} \mathbb{E}\left[X_{i}\right]$$

What is $\mathbb{E}[X_i]$?

- $X_i \leftarrow \#$ balls thrown to get from i empty bins to i-1 empty bins
- $X \leftarrow \#$ balls thrown until we have no empty bins

$$\mathbb{E}[X] = \mathbb{E}\left[\sum_{i=1}^{n} X_{i}\right] = \sum_{i=1}^{n} \mathbb{E}\left[X_{i}\right]$$

What is $\mathbb{E}[X_i]$?

 X_i geometric random variable with parameter $p = \frac{i}{n}$.

Number of trials until the first success, where success probability p.

$$\Pr[X_i = k] = (1-p)^{k-1} \cdot p$$
when we thinker ball into empty bin i empty bins
attempts is 196-1

Coupon Collector - Computing $\mathbb{E}[X]$

$$X_{i} + k \text{ less volum over } N$$

$$E[X_{i}] = \sum_{k=1}^{\infty} k \cdot P_{n}[X_{i} = k] = \sum_{k=1}^{\infty} P_{n}[X_{i} \ge k] = \sum_{k=1}^{\infty} P_{n}[X_{i} \ge k] = \sum_{k=1}^{\infty} P_{n}[X_{i} \ge k] = \sum_{k=1}^{\infty} (1-p)^{k-1} = \frac{1}{1-(1-p)} = \frac{1}{p} = \frac{n}{i}$$

$$E[X] = \sum_{i=1}^{n} E[X_{i}] = \sum_{i=1}^{n} \frac{1}{1-(1-p)} = n \sum_{i=1}^{n} \frac{1}{1-(1-p)} \approx n \log n$$

$$E[X] = \sum_{i=1}^{n} E[X_{i}] = \sum_{i=1}^{n} \frac{1}{1-(1-p)} \approx n \log n$$

Coupon Collector - Computing $\mathbb{E}[X]$

This $n \ln n$ bound shows up in:

- cover time of random walks in complete graph
- number of edges needed in graph sparsification
- many more places

We now know that when n balls are thrown into n bins, the maximum load is $\Theta(\ln n / \ln \ln n)$ with constant probability.¹



¹we'll maybe see lower bound later

We now know that when n balls are thrown into n bins, the maximum load is $\Theta(\ln n / \ln \ln n)$ with constant probability.¹

Consider following variant: what if when throwing a ball in a bin, before we throw the ball we choose two bins uniformly at random and put the ball in the bin with fewer balls?

(Independently if you think about it request kelly)

¹we'll maybe see lower bound later

We now know that when n balls are thrown into n bins, the maximum load is $\Theta(\ln n / \ln \ln n)$ with constant probability.¹

Consider following variant: what if when throwing a ball in a bin, before we throw the ball we choose two bins uniformly at random and put the ball in the bin with fewer balls?

This simple modification reduces maximum load to $O(\ln \ln n)!$



¹we'll maybe see lower bound later

We now know that when n balls are thrown into n bins, the maximum load is $\Theta(\ln n / \ln \ln n)$ with constant probability.¹

Consider following variant: what if when throwing a ball in a bin, before we throw the ball we choose two bins uniformly at random and put the ball in the bin with fewer balls?

This simple modification reduces maximum load to $O(\ln \ln n)!$

Intuition/idea: let the height of a bin be the # balls in that bin. This process tells us that to get one bin with height h+1 we must have at least two bins of height h.

We can bound # bins with height at least h (because this will tell us how likely it is to get to height h+1).



¹we'll maybe see lower bound later

 $N_h :=$ number of bins with height at least h

Pr[at least one bin of height h+1] $\leq \left(\frac{N_h}{n}\right)^2$ we much chose 2 bins of height h

 $N_h :=$ number of bins with height at least h

$$\Pr[\text{at least one bin of height } h+1] \leq \left(\frac{N_h}{n}\right)^2$$

• Say we have only n/4 bins with 4 items (i.e. height 4)

$$\Pr[\text{at least one bin of height } h+1] \leq \left(\frac{N_h}{n}\right)^2$$

- Say we have only n/4 bins with 4 items (i.e. height 4)
- \bullet Probability of selecting two such bins is 1/16

$$\Pr[\text{at least one bin of height } h+1] \leq \left(\frac{N_h}{n}\right)^2$$

- Say we have only n/4 bins with 4 items (i.e. height 4)
- ullet Probability of selecting two such bins is 1/16
- So we should expect only n/16 bins with height 5
- And only $n/256 = n/16^2 = n/2^{2^3}$ bins with height 6

$$\Pr[\text{at least one bin of height } h+1] \leq \left(\frac{N_h}{n}\right)^2$$

- Say we have only n/4 bins with 4 items (i.e. height 4)
- \bullet Probability of selecting two such bins is 1/16
- So we should expect only n/16 bins with height 5
- And only $n/256 = n/16^2 = n/2^{2^3}$ bins with height 6
- Repeating this, we should expect $\frac{n}{2^{2^{h-3}}}$ bins of height h

$$\Pr[\text{at least one bin of height } h+1] \leq \left(\frac{N_h}{n}\right)^2$$

- Say we have only n/4 bins with 4 items (i.e. height 4)
- \bullet Probability of selecting two such bins is 1/16
- So we should expect only n/16 bins with height 5
- And only $n/256 = n/16^2 = n/2^{2^3}$ bins with height 6
- Repeating this, we should expect $\frac{n}{2^{2^{h-3}}}$ bins of height h
- So expect $\log \log n$ maximum height after throwing n balls.

A bit more intuition

 $N_h :=$ number of bins with height at least h

$$\Pr[\text{at least one bin of height } h+1] \leq \left(\frac{N_h}{n}\right)^2$$

- Say we have only n/4 bins with 4 items (i.e. height 4)
- \bullet Probability of selecting two such bins is 1/16
- So we should expect only n/16 bins with height 5
- And only $n/256 = n/16^2 = n/2^{2^3}$ bins with height 6
- Repeating this, we should expect $\frac{n}{2^{2^{h-3}}}$ bins of height h
- So expect $\log \log n$ maximum height after throwing n balls.

How do we turn this into a proof?

Use following Chernoff bound on binomial random variable B(n, p) with ntrials and success probability p.²

$$\Pr[B(n,p) \ge 2np] \le e^{-np/3}$$

Use following Chernoff bound on binomial random variable B(n, p) with n trials and success probability p.²

$$\Pr[B(n,p) \ge 2np] \le e^{-np/3}$$

• $\beta_4 := n/4$ and $\beta_{i+1} = 2\beta_i^2/n$.

²That is, $\Pr[B(n,p)=k]=\binom{n}{k}\cdot p^k\cdot (1-p)^{n-k}$.

Use following Chernoff bound on binomial random variable B(n, p) with n trials and success probability $p.^2$

$$\Pr[B(n,p) \ge 2np] \le e^{-np/3}$$

- $\beta_4 := n/4$ and $\beta_{i+1} = 2\beta_i^2/n$.
- E(h,t) := event that after all t balls are thrown, $N_h \le \beta_h$

²That is, $\Pr[B(n,p)=k]=\binom{n}{k}\cdot p^k\cdot (1-p)^{n-k}$.

Use following Chernoff bound on binomial random variable B(n,p) with n trials and success probability p.²

$$\Pr[B(n,p) \ge 2np] \le e^{-np/3}$$

- $\beta_4 := n/4$ and $\beta_{i+1} = 2\beta_i^2/n$.
- E(h,t):= event that after all t balls are thrown, $N_h \leq \beta_h$
- Pr[E(4, n)] = 1 (why?)

$$E(4,n) = \text{event that after all } n \text{ bells thrown}$$

$$N_4 \leq \frac{n}{4} \qquad N_4 := \# \text{ bins } w_i \text{ th} \qquad \frac{4 \text{ bells}}{4 \text{ bells}}$$

²That is, $\Pr[B(n,p)=k]=\binom{n}{k}\cdot p^k\cdot (1-p)^{n-k}$.

Use following Chernoff bound on binomial random variable B(n, p) with n trials and success probability $p.^2$

$$\Pr[B(n,p) \ge 2np] \le e^{-np/3}$$

- $\beta_4 := n/4$ and $\beta_{i+1} = 2\beta_i^2/n$.
- E(h, t) := event that after all t balls are thrown, $N_h \le \beta_h$
- Pr[E(4, n)] = 1
- We will prove that if E(h, n) holds with high probability then so does E(h+1, n) (so long as h is "small enough")

• $Y_t(h)$ be the indicator variable that t^{th} ball has height $\geq h+1$ (i.e., was placed in a bin that had height h)

•
$$\Pr[Y_t(h) = 1 \mid \underline{E(h, t)}] \le \left(\frac{N_h}{n}\right)^2 \le \frac{\beta_h^2}{n^2}$$

• If $p_h := \frac{\beta_h^2}{n^2}$ then happens than $N_h \leq \beta_h$

$$\Pr\left[\sum_{t=1}^{n} Y_{t}(h) > k \mid E(h, n)\right] \leq \Pr\left[B(n, p_{h}) > k \mid E(h, n)\right]$$
balls held height had

$$\Pr[\widetilde{N_{h+1}} > k \mid E(h, n)] \leq \Pr\left[\sum_{t=1}^{n} Y_t(h) > k \mid E(h, n)\right]$$

$$\leq \Pr\left[B(n, p_h) > k \mid E(h, n)\right]$$

$$\Pr[N_{h+1} > k \mid E(h, n)] \le \Pr[B(n, p_h) > k \mid E(h, n)]$$

$$\Pr[N_{h+1} > k \mid E(h, n)] \le \Pr[B(n, p_h) > k \mid E(h, n)]$$

Setting $k = \beta_{h+1} = 2np_h$ above, we get

$$\Pr[N_{h+1} > k \mid E(h, n)] \le \Pr[B(n, p_h) > k \mid E(h, n)]$$

Setting $k = \beta_{h+1} = 2np_h$ above, we get

$$\begin{split} \Pr[N_{h+1} > \beta_{h+1} \mid E(h,n)] &\leq \Pr\left[B(n,p_h) > \beta_{h+1} \mid E(h,n)\right] \\ & \leq \frac{\Pr\left[B(n,p_h) > \beta_{h+1}\right]}{\Pr[E(h,n)]} \\ &\leq \frac{1}{\Pr[E(h,n)] \cdot e^{np_h/3}} \end{split} \tag{Chernoff}$$

$$\Pr[N_{h+1} > k \mid E(h, n)] \le \Pr[B(n, p_h) > k \mid E(h, n)]$$

Setting $k = \beta_{h+1} = 2np_h$ above, we get

$$\begin{split} \Pr[N_{h+1} > \beta_{h+1} \mid E(h,n)] &\leq \Pr\left[B(n,p_h) > \beta_{h+1} \mid E(h,n)\right] \\ &\leq \frac{\Pr\left[B(n,p_h) > \beta_{h+1}\right]}{\Pr[E(h,n)]} \\ &\leq \frac{1}{\Pr[E(h,n)] \cdot e^{np_h/3}} \end{split} \tag{Chernoff}$$

Thus, setting $p_h \cdot n \ge 6 \ln n$ we get

$$\Pr[\text{not } E(h+1,n) \mid E(h,n)] = \Pr[N_{h+1} > \beta_{h+1} \mid E(h,n)] \le \frac{1}{n^2 \Pr[E(h,n)]}$$

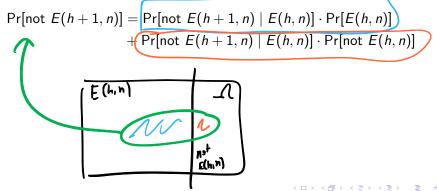
$$\Pr[\text{not } E(h+1,n) \mid E(h,n)] = \Pr[N_{h+1} > \beta_{h+1} \mid E(h,n)] \le \frac{1}{n^2 \Pr[E(h,n)]}$$

$$\Pr[\text{not } E(h+1,n) \mid E(h,n)] = \Pr[N_{h+1} > \beta_{h+1} \mid E(h,n)] \le \frac{1}{n^2 \Pr[E(h,n)]}$$

Now, to bound the final probability, we have:

$$\Pr[\text{not } E(h+1,n) \mid E(h,n)] = \Pr[N_{h+1} > \beta_{h+1} \mid E(h,n)] \le \frac{1}{n^2 \Pr[E(h,n)]}$$

Now, to bound the final probability, we have:



$$\Pr[\text{not } E(h+1,n) \mid E(h,n)] = \Pr[N_{h+1} > \beta_{h+1} \mid E(h,n)] \le \frac{1}{n^2 \Pr[E(h,n)]}$$

Now, to bound the final probability, we have:

$$\begin{split} \Pr[\mathsf{not}\ E(h+1,n)] &= \Pr[\mathsf{not}\ E(h+1,n)\mid E(h,n)] \cdot \Pr[E(h,n)] \\ &+ \Pr[\mathsf{not}\ E(h+1,n)\mid E(h,n)] \cdot \Pr[\mathsf{not}\ E(h,n)] \\ &\leq \frac{1}{n^2} + \Pr[\mathsf{not}\ E(h,n)] \quad (\mathsf{so}\ \mathsf{long}\ \mathsf{as}\ p_h n \geq 6\,\mathsf{ln}\ n) \end{split}$$

$$\Pr[\text{not } E(h+1,n) \mid E(h,n)] = \Pr[N_{h+1} > \beta_{h+1} \mid E(h,n)] \le \frac{1}{n^2 \Pr[E(h,n)]}$$

Now, to bound the final probability, we have:

$$\begin{aligned} \Pr[\mathsf{not}\ E(h+1,n)] &= \Pr[\mathsf{not}\ E(h+1,n)\mid E(h,n)] \cdot \Pr[E(h,n)] \\ &+ \Pr[\mathsf{not}\ E(h+1,n)\mid E(h,n)] \cdot \Pr[\mathsf{not}\ E(h,n)] \\ &\leq \frac{1}{n^2} + \Pr[\mathsf{not}\ E(h,n)] \quad (\mathsf{so}\ \mathsf{long}\ \mathsf{as}\ p_h n \geq 6 \,\mathsf{ln}\ n) \end{aligned}$$

To finish the proof, need to show:

- $p_h \cdot n \ge 6 \ln n$ for $h = O(\ln \ln n)$ (easy calculation we did it)
- Handle the case where $p_h \cdot n < 6 \ln n$. (another Chernoff bound see Lap Chi's notes)

Acknowledgement

- Lecture based largely on Lap Chi's notes and on [Motwani & Raghavan 2007, Chapter 3].
- See Lap Chi's notes at https://cs.uwaterloo.ca/~lapchi/cs466/notes/L04.pdf

References I



Motwani, Rajeev and Raghavan, Prabhakar (2007)

Randomized Algorithms



Mitzenmacher, Michael, and Eli Upfal (2017)

Probability and computing: Randomization and probabilistic techniques in algorithms and data analysis.

Cambridge university press, 2017.