

COMPSCI 614: Randomized Algorithms with Applications to Data Science

Prof. Cameron Musco

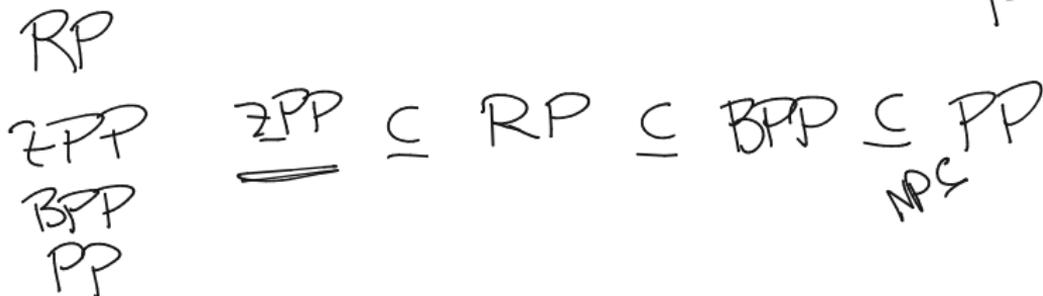
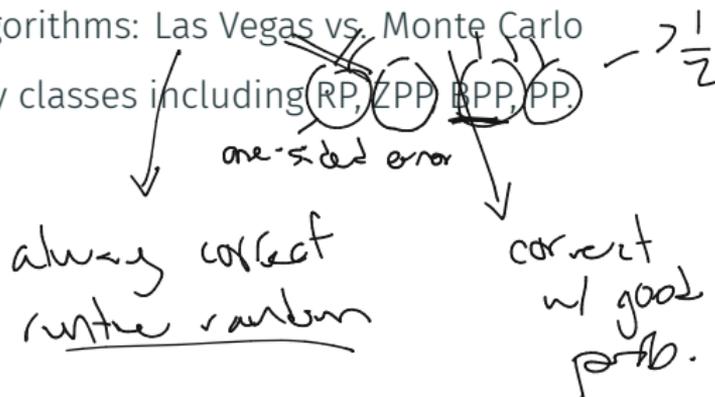
University of Massachusetts Amherst. Spring 2024.

Lecture 2

Summary

Last Class:

- Course logistics/overview of planned content.
- Intro to randomized algorithms: Las Vegas vs. Monte Carlo
- Randomized complexity classes including RP, ZPP, BPP, PP



Summary

Last Class:

- Course logistics/overview of planned content.
- Intro to randomized algorithms: Las Vegas vs. Monte Carlo
- Randomized complexity classes including RP, ZPP, BPP, PP.

This Class: Basic probability review with algorithmic applications.

- Conditional probability, Baye's theorem, and independence.
Application to polynomial identity testing.
- Linearity of expectation and variance. Application to randomized Quicksort analysis.
- Maybe start on concentration inequalities (Markov's and Chebyshev's).

Basic Probability Review

Conditional Probability Review

Consider two random events A and B.

B = die roll is even
A = die roll is "2"

- Conditional Probability:

$$P(A|B) \cdot P(B) = Pr(A \cap B) \quad \left| \quad Pr(A|B) = \frac{Pr(A \cap B)}{Pr(B)} \rightarrow \frac{1/6}{1/2} = \frac{1}{3} \right.$$

- Baye's Theorem:

$$Pr(B|A) = \frac{Pr(A|B) \cdot P(B)}{P(A)}$$

- Independence: A and B are independent if:

$$Pr(A|B) = Pr(A).$$

Using the definition of conditional probability, independence means:

$$Pr(A|B) = \frac{Pr(A \cap B)}{Pr(B)} = Pr(A) \implies Pr(A \cap B) = Pr(A) \cdot Pr(B)$$

Independence

Sets of events: For a set of n events, A_1, \dots, A_n , the events are k -wise independent if for any subset S of at most k events,

$$\Pr\left(\bigcap_{i \in S} A_i\right) = \prod_{i \in S} \Pr(A_i).$$

For $k = n$ we just say the events 'are independent'.

$$\begin{aligned} D_1 &= 1 \\ D_2 &= 1 \\ D_3 &= 1 \end{aligned}$$

$$D_1 \quad D_2$$

$$D_3$$

$$A: D_1 = D_2$$

$$P(A) = \frac{1}{6}$$

$$P(B) = \frac{1}{6}$$

$$P(C)$$

$$\begin{aligned} P(A \cap B) &= \frac{1}{36} \\ &= P(A) \cdot P(B) \end{aligned}$$

$$B: D_2 = D_3$$

$$C: D_1 = D_3$$

$$\begin{aligned} P(A \cap B \cap C) &= \frac{1}{36} \\ &\neq P(A) \cdot P(B) \cdot P(C) \end{aligned}$$

Independence

Sets of events: For a set of n events, A_1, \dots, A_n , the events are **k -wise independent** if for any subset S of at most k events,

$$\Pr\left(\bigcap_{i \in S} A_i\right) = \prod_{i \in S} \Pr(A_i).$$

For $k = n$ we just say the events 'are independent'.

Random Variables: Two random variables X , Y are independent if for all s, t , $X = s$ and $Y = t$ are independent events. In other words:

$$\Pr(X = s \cap Y = t) = \Pr(X = s) \cdot \Pr(Y = t).$$

Application 1: Polynomial Identity Testing

Polynomial Identity Testing

Given an n -variable degree- d polynomial $p(x_1, x_2, \dots, x_n)$, determine if the polynomial is **identically zero**. I.e., if $p(x_1, x_2, \dots, x_n) = 0$ for all x_1, \dots, x_n .

Polynomial Identity Testing

Given an n -variable degree- d polynomial $p(x_1, x_2, \dots, x_n)$, determine if the polynomial is **identically zero**. I.e., if $p(x_1, x_2, \dots, x_n) = 0$ for all x_1, \dots, x_n . E.g., you are given:

$$p(x_1, x_2, \dots, x_3) = x_3(x_1 - x_2)^3 + (x_1 + 2x_2 - x_3)^2 - x_1(x_2 + x_3)^2.$$

$$\underline{x_3 x_1^2 \cdot x_3 \cdot 2x_1 x_2 + \dots} \stackrel{?}{=} 0$$

Polynomial Identity Testing

Given an n -variable degree- d polynomial $p(x_1, x_2, \dots, x_n)$, determine if the polynomial is **identically zero**. I.e., if $p(x_1, x_2, \dots, x_n) = 0$ for all x_1, \dots, x_n . E.g., you are given:

$$p(x_1, x_2, \dots, x_3) = x_3(x_1 - x_2)^3 + (x_1 + 2x_2 - x_3)^2 - x_1(x_2 + x_3)^2.$$

$x_1 \quad x_3 \quad x_2^3$

- Can expand out all the terms and check if they cancel. But the number of terms can be as large as $\binom{n+d}{d}$ – i.e., exponential in the number of variables n and the degree d . $\approx n^d$

Polynomial Identity Testing

Given an n -variable degree- d polynomial $p(x_1, x_2, \dots, x_n)$, determine if the polynomial is **identically zero**. I.e., if $p(x_1, x_2, \dots, x_n) = 0$ for all x_1, \dots, x_n . E.g., you are given:

$$p(x_1, x_2, \dots, x_3) = x_3(x_1 - x_2)^3 + (x_1 + 2x_2 - x_3)^2 - x_1(x_2 + x_3)^2.$$

- Can expand out all the terms and check if they cancel. But the number of terms can be as large as $\binom{n+d}{d}$ – i.e., exponential in the number of variables n and the degree d .

Extremely Simple Randomized Algorithm: Just pick random values for x_1, \dots, x_n and evaluate the polynomial at these values. With high probability, if $p(x_1, \dots, x_n) = 0$, the polynomial is identically 0!

$$p(\underline{5, 2, \dots, -1}) = -1(5 - 2)^3 + (5 + 2 \cdot 2 + 1)^2 - 5(2 - 1)^2 = \underline{68}.$$

0
[CO-RP]

if $p \neq 0$
(may) witness incorrectly
if $p = 0$
always correctly

IPP ?
RP
BPP

Polynomial Identity Testing

Given an n -variable degree- d polynomial $p(x_1, x_2, \dots, x_n)$, determine if the polynomial is **identically zero**. I.e., if $p(x_1, x_2, \dots, x_n) = 0$ for all x_1, \dots, x_n . E.g., you are given:

$$p(x_1, x_2, \dots, x_3) = x_3(x_1 - x_2)^3 + (x_1 + 2x_2 - x_3)^2 - x_1(x_2 + x_3)^2.$$

- Can expand out all the terms and check if they cancel. But the number of terms can be as large as $\binom{n+d}{d}$ – i.e., exponential in the number of variables n and the degree d .

Extremely Simple Randomized Algorithm: Just pick random values for x_1, \dots, x_n and evaluate the polynomial at these values. With high probability, if $p(x_1, \dots, x_n) = 0$, the polynomial is identically 0!

$$p(5, 2, \dots, -1) = -1(5 - 2)^3 + (5 + 2 \cdot 2 + 1)^2 - 5(2 - 1)^2 = 68.$$

What style algorithm is this? BPP, ZPP, RP, something else?

co-RP

Polynomial Identity Testing Proof

Schwartz-Zippel Lemma: For any n -variable degree- d polynomial $p(x_1, \dots, x_n)$ and any set S , if z_1, \dots, z_n are selected independently and uniformly at random from S , then $\Pr[p(z_1, \dots, z_n) \neq 0] \geq 1 - \frac{d}{|S|}$.

$$d = 10$$

$$1 - \frac{d}{|S|} = 1 - \frac{1}{10} = 90\%$$

$$S = \{0, 1, 2, 3, \dots, 100\}$$

$$\{0, .5, 1, 1.5, \dots\}$$

$$x_1 - x_2$$

Polynomial Identity Testing Proof

Schwartz-Zippel Lemma: For any n -variable degree- d polynomial $p(x_1, \dots, x_n)$ and any set S , if z_1, \dots, z_n are selected independently and uniformly at random from S , then $\Pr[p(z_1, \dots, z_n) \neq 0] \geq 1 - \frac{d}{|S|}$.

Proof: Via induction on the number of variables n

Base Case $n = 1$:

$$p(x_1) = x_1^2 + 3x_1^3 + 4x_1^4 + \dots$$

p has at most d roots
 $p(x) = 0$ for at most d values of x

worst case, all roots are in S

probability that z_1 is a root so $p(z_1) = 0$

$$\Pr(p(z_1) = 0) \leq \frac{d}{|S|} \geq 1 - \frac{d}{|S|}$$

Polynomial Identity Testing Proof

Schwartz-Zippel Lemma: For any n -variable degree- d polynomial $p(x_1, \dots, x_n)$ and any set S , if z_1, \dots, z_n are selected independently and uniformly at random from S , then $\Pr[p(z_1, \dots, z_n) \neq 0] \geq 1 - \frac{d}{|S|}$.

Proof: Via induction on the number of variables n

Induction Step $n > 1$:

Polynomial Identity Testing Proof

Schwartz-Zippel Lemma: For any n -variable degree- d polynomial $p(x_1, \dots, x_n)$ and any set S , if z_1, \dots, z_n are selected independently and uniformly at random from S , then $\Pr[p(z_1, \dots, z_n) \neq 0] \geq 1 - \frac{d}{|S|}$.

Proof: Via induction on the number of variables n

Induction Step $n > 1$:

- Let k be the max degree of x_1 in $p(\dots)$. Assume w.l.o.g. that $k > 0$. Write $p(x_1, \dots, x_n) = x_1^k \cdot q(x_2, \dots, x_n) + r(x_1, \dots, x_n)$. E.g.,

$$k = 2 \quad \underline{x_1^2 x_2 + x_1^2 x_3 + x_1 x_2 x_3 + x_2 x_3} = \underline{x_1^2 \cdot \underbrace{(x_2 + x_3)}_{q(\dots)}} + \underbrace{x_1 x_2 x_3 + x_2 x_3}_{r(\dots)}$$

Polynomial Identity Testing Proof

Schwartz-Zippel Lemma: For any n -variable degree- d polynomial $p(x_1, \dots, x_n)$ and any set S , if z_1, \dots, z_n are selected independently and uniformly at random from S , then $\Pr[p(z_1, \dots, z_n) \neq 0] \geq 1 - \frac{d}{|S|}$.

Proof: Via induction on the number of variables n

Induction Step $n > 1$:

- Let k be the max degree of x_1 in $p(\dots)$. Assume w.l.o.g. that $k > 0$. Write $p(x_1, \dots, x_n) = x_1^k \cdot q(x_2, \dots, x_n) + r(x_1, \dots, x_n)$. E.g.,

$$x_1^2 x_2 + x_1^2 x_3 + x_1 x_2 x_3 + x_2 x_3 = x_1^2 \cdot \underbrace{(x_2 + x_3)}_{q(\dots)} + \underbrace{x_1 x_2 x_3 + x_2 x_3}_{r(\dots)}$$

- Observe: $q(\dots)$ is non-zero, with $\boxed{n-1}$ variables and degree $d-k$.

Polynomial Identity Testing Proof

Schwartz-Zippel Lemma: For any n -variable degree- d polynomial $p(x_1, \dots, x_n)$ and any set S , if z_1, \dots, z_n are selected independently and uniformly at random from S , then $\Pr[p(z_1, \dots, z_n) \neq 0] \geq 1 - \frac{d}{|S|}$.

Proof: Via induction on the number of variables n

Induction Step $n > 1$:

- Let k be the max degree of x_1 in $p(\dots)$. Assume w.l.o.g. that $k > 0$. Write $p(x_1, \dots, x_n) = x_1^k \cdot q(x_2, \dots, x_n) + r(x_1, \dots, x_n)$. E.g.,

$$x_1^2 x_2 + x_1^2 x_3 + x_1 x_2 x_3 + x_2 x_3 = x_1^2 \cdot \underbrace{(x_2 + x_3)}_{q(\dots)} + \underbrace{x_1 x_2 x_3 + x_2 x_3}_{r(\dots)}$$

- Observe: $q(\cdot)$ is non-zero, with $n - 1$ variables and degree $d - k$.
- So, by inductive assumption, $\Pr[\underbrace{q(z_2, \dots, z_n)}_{\neq 0}] \geq 1 - \frac{d-k}{|S|}$.

Polynomial Identity Testing Proof

Schwartz-Zippel Lemma: For any n -variable degree- d polynomial $p(x_1, \dots, x_n)$ and any set S , if z_1, \dots, z_n are selected independently and uniformly at random from S , then $\Pr[p(z_1, \dots, z_n) \neq 0] \geq 1 - \frac{d}{|S|}$.

Proof: Via induction on the number of variables n

Induction Step $n > 1$:

- Let k be the max degree of x_1 in $p(\dots)$. Assume w.l.o.g. that $k > 0$. Write $p(x_1, \dots, x_n) = x_1^k \cdot q(x_2, \dots, x_n) + r(x_1, \dots, x_n)$. E.g.,

$$x_1^2 x_2 + x_1^2 x_3 + x_1 x_2 x_3 + x_2 x_3 = x_1^2 \cdot \underbrace{(x_2 + x_3)}_{q(\dots)} + \underbrace{x_1 x_2 x_3 + x_2 x_3}_{r(\dots)}$$

- Observe: $q(\cdot)$ is non-zero, with $n - 1$ variables and degree $d - k$.
- So, by inductive assumption, $\Pr[q(z_2, \dots, z_n) \neq 0] \geq 1 - \frac{d-k}{|S|}$.
- Assuming $q(z_2, \dots, z_n) \neq 0$, then $p(\underline{x_1}, \underline{z_2}, \dots, \underline{z_n})$ is a degree k non-zero univariate polynomial in x_1 .

Polynomial Identity Testing Proof

Assuming $q(z_2, \dots, z_n) \neq 0$, then $p(x_1, z_2, \dots, z_n)$ is a degree k non-zero univariate polynomial in x_1 .

Example:

$$p(x_1, x_2, x_3) = x_1^2 x_2 + x_1^2 x_3 + x_1 x_2 x_3 + x_2 x_3 = x_1^2 \underbrace{(x_2 + x_3)}_{q(\dots)} + \underbrace{x_1 x_2 x_3 + x_2 x_3}_{r(\dots)}$$

$2+1$ $x_1 \cdot 2 \cdot 1$ $2 \cdot 1$

$$p(x_1, z_2, z_3) = p(x_1, \underline{2}, \underline{1}) = \underline{x_1^2 \cdot 3 + 2x_1 + 2}$$

$$q(z_2, \dots, z_n)$$

Polynomial Identity Testing Proof

Assuming $q(z_2, \dots, z_n) \neq 0$, then $p(x_1, z_2, \dots, z_n)$ is a degree k non-zero univariate polynomial in x_1 .

Example:

$$p(x_1, x_2, x_3) = x_1^2 x_2 + x_1^2 x_3 + x_1 x_2 x_3 + x_2 x_3 = x_1^2 \cdot \underbrace{(x_2 + x_3)}_{q(\dots)} + \underbrace{x_1 x_2 x_3 + x_2 x_3}_{r(\dots)}$$

$$p(x_1, z_2, z_3) = \underbrace{p(x_1, 2, 1)} = x_1^2 \cdot 3 + 2x_1 + 2.$$

- degree k poly in 1-variate

Next Step: Again applying the inductive hypothesis,

$$\Pr[\underbrace{p(z_1, \dots, z_n)} \neq 0 \mid \underbrace{q(z_2, \dots, z_n)} \neq 0] \geq \underbrace{1 - \frac{k}{|S|}}.$$

Polynomial Identity Testing Proof

Assuming $q(z_2, \dots, z_n) \neq 0$, then $p(x_1, z_2, \dots, z_n)$ is a degree k non-zero univariate polynomial in x_1 .

Example:

$$p(x_1, x_2, x_3) = x_1^2 x_2 + x_1^2 x_3 + x_1 x_2 x_3 + x_2 x_3 = x_1^2 \cdot \underbrace{(x_2 + x_3)}_{q(\dots)} + \underbrace{x_1 x_2 x_3 + x_2 x_3}_{r(\dots)}$$

$$p(x_1, z_2, z_3) = p(x_1, 2, 1) = x_1^2 \cdot 3 + 2x_1 + 2.$$

Next Step: Again applying the inductive hypothesis,

$$\Pr[p(z_1, \dots, z_n) \neq 0 | q(z_2, \dots, z_n) \neq 0] \geq 1 - \frac{k}{|S|}.$$

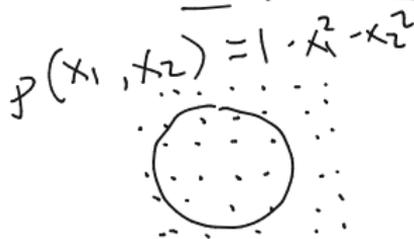
Overall:

$$\Pr[p(z_1, \dots, z_n) \neq 0] \geq \Pr[p(z_1, \dots, z_n) \neq 0 \cap q(z_2, \dots, z_n) \neq 0]$$

$$= \Pr[p(\dots) \neq 0 | q(\dots) \neq 0] \cdot \Pr[q(\dots) \neq 0]$$

$$\geq \left(1 - \frac{k}{|S|}\right) \cdot \left(1 - \frac{d-k}{|S|}\right) \geq 1 - \frac{d}{|S|}$$

$$\geq 1 - \frac{d}{|S|^2}$$



\int^n

Polynomial Identity Testing Proof

Assuming $q(z_2, \dots, z_n) \neq 0$, then $p(x_1, z_2, \dots, z_n)$ is a degree k non-zero univariate polynomial in x_1 .

Example:

$$p(x_1, x_2, x_3) = x_1^2 x_2 + x_1^2 x_3 + x_1 x_2 x_3 + x_2 x_3 = x_1^2 \cdot \underbrace{(x_2 + x_3)}_{q(\dots)} + \underbrace{x_1 x_2 x_3 + x_2 x_3}_{r(\dots)}$$

$$p(x_1, z_2, z_3) = p(x_1, 2, 1) = x_1^2 \cdot 3 + 2x_1 + 2.$$

Next Step: Again applying the inductive hypothesis,

$$\Pr[p(z_1, \dots, z_n) \neq 0 | q(z_2, \dots, z_n) \neq 0] \geq 1 - \frac{k}{|S|}.$$

Overall: $\Pr[p(z_1, \dots, z_n) \neq 0] \geq \Pr[p(z_1, \dots, z_n) \neq 0 \cap q(z_2, \dots, z_n) \neq 0]$

$p(z_1) = 1 - x_2^2 - x_1^2$ $|S| = 2$

$\frac{1}{|S|} = 1$ $= \Pr[p(\dots) \neq 0 | q(\dots) \neq 0] \cdot \Pr[q(\dots) \neq 0]$

$1 - \frac{2}{|S|} = 0$ $\geq \left(1 - \frac{k}{|S|}\right) \cdot \left(1 - \frac{d-k}{|S|}\right) \geq 1 - \frac{d}{|S|}$



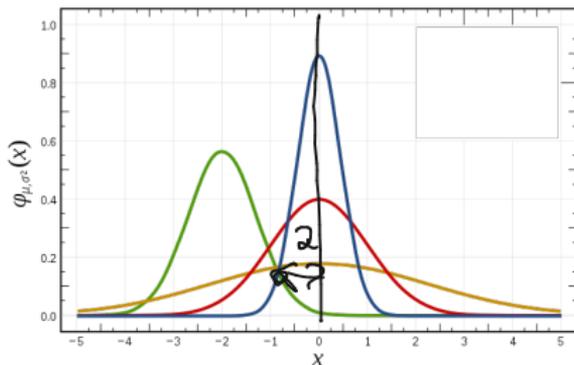
This completes the proof of Schwartz-Zippel.

Expectation and Variance Review

Expectation and Variance

Consider a random X variable taking values in some finite set $S \subset \mathbb{R}$. E.g., for a random dice roll, $S = \{1, 2, 3, 4, 5, 6\}$. $\curvearrowright \{2, 4, 6, 10, 12\}$

- **Expectation:** $\mathbb{E}[X] = \sum_{s \in S} \underline{\text{Pr}(X = s)} \cdot \underline{s}$.
- **Variance:** $\text{Var}[X] = \mathbb{E}[(X - \mathbb{E}[X])^2]$.



Exercise: Verify that for any scalar α , $\mathbb{E}[\alpha \cdot X] = \alpha \cdot \mathbb{E}[X]$ and $\text{Var}[\alpha \cdot X] = \alpha^2 \cdot \text{Var}[X]$.

Linearity of Expectation

$\mathbb{E}[X + Y] = \mathbb{E}[X] + \mathbb{E}[Y]$ for any random variables X and Y . No matter how correlated they may be!



Linearity of Expectation

$\mathbb{E}[X + Y] = \mathbb{E}[X] + \mathbb{E}[Y]$ for any random variables X and Y . No matter how correlated they may be!

Proof:

$$X+Y \text{ when } X=s + Y=t$$

$$\mathbb{E}[X + Y] = \sum_{s \in S} \sum_{t \in T} \Pr(X = s \cap Y = t) \cdot (s + t)$$

Linearity of Expectation

$\mathbb{E}[X + Y] = \mathbb{E}[X] + \mathbb{E}[Y]$ for any random variables X and Y . No matter how correlated they may be!

Proof:

$$\begin{aligned}\mathbb{E}[X + Y] &= \sum_{s \in S} \sum_{t \in T} \Pr(X = s \cap Y = t) \cdot (s + t) \\ &= \sum_{s \in S} \underbrace{\left(\sum_{t \in T} \Pr(X = s \cap Y = t) \right)}_{\Pr(X = s)} \cdot s + \sum_{t \in T} \sum_{s \in S} \Pr(X = s \cap Y = t) \cdot t \\ &\quad \underbrace{\sum_{s \in S} \Pr(X = s) \cdot s}_{\mathbb{E}[X]} \quad \quad \quad \downarrow \\ &\quad \quad \quad \mathbb{E}[Y]\end{aligned}$$

Linearity of Expectation

$\mathbb{E}[X + Y] = \mathbb{E}[X] + \mathbb{E}[Y]$ for any random variables X and Y . No matter how correlated they may be!

Proof:

$$\begin{aligned}\mathbb{E}[X + Y] &= \sum_{s \in S} \sum_{t \in T} \Pr(X = s \cap Y = t) \cdot (s + t) \\ &= \sum_{s \in S} \sum_{t \in T} \Pr(X = s \cap Y = t) \cdot s + \sum_{t \in T} \sum_{s \in S} \Pr(X = s \cap Y = t) \cdot t \\ &= \sum_{s \in S} \Pr(X = s) \cdot s + \sum_{t \in T} \Pr(Y = t) \cdot t\end{aligned}$$

Linearity of Expectation

$\mathbb{E}[X + Y] = \mathbb{E}[X] + \mathbb{E}[Y]$ for any random variables X and Y . No matter how correlated they may be!

Proof:

$$\begin{aligned}\mathbb{E}[X + Y] &= \sum_{s \in S} \sum_{t \in T} \Pr(X = s \cap Y = t) \cdot (s + t) \\ &= \sum_{s \in S} \sum_{t \in T} \Pr(X = s \cap Y = t) \cdot s + \sum_{t \in T} \sum_{s \in S} \Pr(X = s \cap Y = t) \cdot t \\ &= \sum_{s \in S} \Pr(X = s) \cdot s + \sum_{t \in T} \Pr(Y = t) \cdot t \\ &= \mathbb{E}[X] + \mathbb{E}[Y].\end{aligned}$$

Linearity of Expectation

$\mathbb{E}[X + Y] = \mathbb{E}[X] + \mathbb{E}[Y]$ for any random variables X and Y . No matter how correlated they may be!

Proof:

$$\begin{aligned}\mathbb{E}[X + Y] &= \sum_{s \in S} \sum_{t \in T} \Pr(X = s \cap Y = t) \cdot (s + t) \\ &= \sum_{s \in S} \sum_{t \in T} \Pr(X = s \cap Y = t) \cdot s + \sum_{t \in T} \sum_{s \in S} \Pr(X = s \cap Y = t) \cdot t \\ &= \sum_{s \in S} \Pr(X = s) \cdot s + \sum_{t \in T} \Pr(Y = t) \cdot t \\ &= \mathbb{E}[X] + \mathbb{E}[Y].\end{aligned}$$

Maybe the single most powerful tool in the analysis of randomized algorithms.

Linearity of Variance

Var[X + Y] = Var[X] + Var[Y] when X and Y are independent.

Linearity of Variance

$\text{Var}[X + Y] = \text{Var}[X] + \text{Var}[Y]$ when X and Y are independent.

Claim 1: (exercise) $\text{Var}[X] = \mathbb{E}[X^2] - \mathbb{E}[X]^2$ (via linearity of expectation)

$\mathbb{E}[(X - \mathbb{E}[X])^2]$

Claim 2: (exercise) $\mathbb{E}[XY] = \mathbb{E}[X] \cdot \mathbb{E}[Y]$ (i.e., X and Y are uncorrelated) when X, Y are independent.

Linearity of Variance

$\text{Var}[X + Y] = \text{Var}[X] + \text{Var}[Y]$ when X and Y are independent.

Claim 1: (exercise) $\text{Var}[X] = \mathbb{E}[X^2] - \mathbb{E}[X]^2$ (via linearity of expectation)

Claim 2: (exercise) $\mathbb{E}[XY] = \mathbb{E}[X] \cdot \mathbb{E}[Y]$ (i.e., X and Y are uncorrelated) when X, Y are independent.

Together give:

$$\text{Var}[X + Y] = \mathbb{E}[(X + Y)^2] - \mathbb{E}[X + Y]^2$$

Linearity of Variance

$\text{Var}[X + Y] = \text{Var}[X] + \text{Var}[Y]$ when X and Y are independent.

Claim 1: (exercise) $\text{Var}[X] = \mathbb{E}[X^2] - \mathbb{E}[X]^2$ (via linearity of expectation)

Claim 2: (exercise) $\mathbb{E}[XY] = \mathbb{E}[X] \cdot \mathbb{E}[Y]$ (i.e., X and Y are uncorrelated) when X, Y are independent.

Together give: $\mathbb{E}[X^2 + 2XY + Y^2]$

$$\begin{aligned}\text{Var}[X + Y] &= \mathbb{E}[(X + Y)^2] - \mathbb{E}[X + Y]^2 \\ &= \underbrace{\mathbb{E}[X^2]} + \underbrace{2\mathbb{E}[XY]} + \underbrace{\mathbb{E}[Y^2]} - (\mathbb{E}[X] + \mathbb{E}[Y])^2\end{aligned}$$

Linearity of Variance

$\text{Var}[X + Y] = \text{Var}[X] + \text{Var}[Y]$ when X and Y are independent.

Claim 1: (exercise) $\text{Var}[X] = \mathbb{E}[X^2] - \mathbb{E}[X]^2$ (via linearity of expectation)

Claim 2: (exercise) $\mathbb{E}[XY] = \mathbb{E}[X] \cdot \mathbb{E}[Y]$ (i.e., X and Y are uncorrelated) when X, Y are independent.

Together give:

$$\begin{aligned}\text{Var}[X + Y] &= \mathbb{E}[(X + Y)^2] - \mathbb{E}[X + Y]^2 \\ &= \mathbb{E}[X^2] + 2\mathbb{E}[XY] + \mathbb{E}[Y^2] - (\mathbb{E}[X] + \mathbb{E}[Y])^2 \\ &= \mathbb{E}[X^2] + 2\mathbb{E}[XY] + \mathbb{E}[Y^2] - \mathbb{E}[X]^2 - 2\mathbb{E}[X] \cdot \mathbb{E}[Y] - \mathbb{E}[Y]^2 \\ &= \underbrace{\mathbb{E}[X^2] - \mathbb{E}[X]^2}_{\text{Var}(X)} + \underbrace{\mathbb{E}[Y^2] - \mathbb{E}[Y]^2}_{\text{Var}(Y)} + \underbrace{2\mathbb{E}[XY] - 2\mathbb{E}[X]\mathbb{E}[Y]}_{0 \text{ by claim 2}}\end{aligned}$$

Linearity of Variance

$\text{Var}[X + Y] = \text{Var}[X] + \text{Var}[Y]$ when X and Y are independent.

Claim 1: (exercise) $\text{Var}[X] = \mathbb{E}[X^2] - \mathbb{E}[X]^2$ (via linearity of expectation)

Claim 2: (exercise) $\mathbb{E}[XY] = \mathbb{E}[X] \cdot \mathbb{E}[Y]$ (i.e., X and Y are uncorrelated) when X, Y are independent.

Together give:

$$\begin{aligned}\text{Var}[X + Y] &= \mathbb{E}[(X + Y)^2] - \mathbb{E}[X + Y]^2 \\ &= \mathbb{E}[X^2] + 2\mathbb{E}[XY] + \mathbb{E}[Y^2] - (\mathbb{E}[X] + \mathbb{E}[Y])^2 \\ &= \mathbb{E}[X^2] + 2\mathbb{E}[XY] + \mathbb{E}[Y^2] - \mathbb{E}[X]^2 - 2\mathbb{E}[X] \cdot \mathbb{E}[Y] - \mathbb{E}[Y]^2 \\ &= \mathbb{E}[X^2] + \mathbb{E}[Y^2] - \mathbb{E}[X]^2 - \mathbb{E}[Y]^2\end{aligned}$$

Linearity of Variance

$\text{Var}[X + Y] = \text{Var}[X] + \text{Var}[Y]$ when X and Y are independent.

Claim 1: (exercise) $\text{Var}[X] = \mathbb{E}[X^2] - \mathbb{E}[X]^2$ (via linearity of expectation)

Claim 2: (exercise) $\mathbb{E}[XY] = \mathbb{E}[X] \cdot \mathbb{E}[Y]$ (i.e., X and Y are uncorrelated) when X, Y are independent.

Together give:

$$\begin{aligned}\text{Var}[X + Y] &= \mathbb{E}[(X + Y)^2] - \mathbb{E}[X + Y]^2 \\ &= \mathbb{E}[X^2] + 2\mathbb{E}[XY] + \mathbb{E}[Y^2] - (\mathbb{E}[X] + \mathbb{E}[Y])^2 \\ &= \mathbb{E}[X^2] + 2\mathbb{E}[XY] + \mathbb{E}[Y^2] - \mathbb{E}[X]^2 - 2\mathbb{E}[X] \cdot \mathbb{E}[Y] - \mathbb{E}[Y]^2 \\ &= \mathbb{E}[X^2] + \mathbb{E}[Y^2] - \mathbb{E}[X]^2 - \mathbb{E}[Y]^2 \\ &= \text{Var}[X] + \text{Var}[Y].\end{aligned}$$

Linearity of Variance

Exercise: Verify that for random variables X_1, \dots, X_n ,

$$\text{Var} \left(\sum_{i=1}^n X_i \right) = \sum_{i=1}^n \text{Var}(X_i),$$

whenever the variables are 2-wise independent (also called pairwise independent).

Application 2: Quicksort with Random Pivots

Quicksort

Quicksort(X): where $X = (x_1, \dots, x_n)$ is a list of numbers.

1. If X is empty: return X .
2. Else: select pivot p uniformly at random from $\{1, \dots, n\}$.
3. Let $X_{lo} = \{i \in X : x_i < x_p\}$ and $X_{hi} = \{i \in X : x_i \geq x_p\}$ (requires $n - 1$ comparisons with x_p to determine).
4. Return the concatenation of the lists $[\text{Quicksort}(X_{lo}), (x_p), \text{Quicksort}(X_{hi})]$.

Quicksort

Quicksort(X): where $X = (x_1, \dots, x_n)$ is a list of numbers.

1. If X is empty: return X .
2. Else: select pivot p uniformly at random from $\{1, \dots, n\}$.
3. Let $X_{lo} = \{i \in X : x_i < x_p\}$ and $X_{hi} = \{i \in X : x_i \geq x_p\}$ (requires $n - 1$ comparisons with x_p to determine).
4. Return the concatenation of the lists
 $[\text{Quicksort}(X_{lo}), (x_p), \text{Quicksort}(X_{hi})]$.

4	5	2	8	1	3	6	9	7	0
---	---	---	---	---	---	---	---	---	---

Quicksort

Quicksort(X): where $X = (x_1, \dots, x_n)$ is a list of numbers.

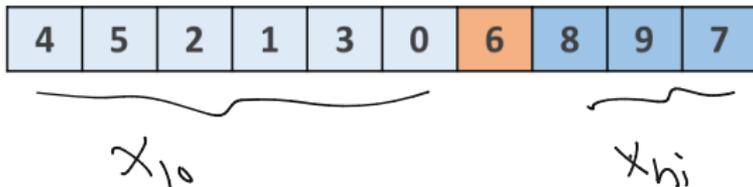
1. If X is empty: return X .
2. Else: select pivot p uniformly at random from $\{1, \dots, n\}$.
3. Let $X_{lo} = \{i \in X : x_i < x_p\}$ and $X_{hi} = \{i \in X : x_i \geq x_p\}$ (requires $n - 1$ comparisons with x_p to determine).
4. Return the concatenation of the lists
 $[\text{Quicksort}(X_{lo}), (x_p), \text{Quicksort}(X_{hi})]$.

4	5	2	8	1	3	6	9	7	0
---	---	---	---	---	---	---	---	---	---

Quicksort

Quicksort(X): where $X = (x_1, \dots, x_n)$ is a list of numbers.

1. If X is empty: return X .
2. Else: select pivot p uniformly at random from $\{1, \dots, n\}$.
3. Let $X_{lo} = \{i \in X : x_i < x_p\}$ and $X_{hi} = \{i \in X : x_i \geq x_p\}$ (requires $n - 1$ comparisons with x_p to determine).
4. Return the concatenation of the lists
[Quicksort(X_{lo}), (x_p), Quicksort(X_{hi})].



Quicksort

Quicksort(X): where $X = (x_1, \dots, x_n)$ is a list of numbers.

1. If X is empty: return X .
2. Else: select pivot p uniformly at random from $\{1, \dots, n\}$.
3. Let $X_{lo} = \{i \in X : x_i < x_p\}$ and $X_{hi} = \{i \in X : x_i \geq x_p\}$ (requires $n - 1$ comparisons with x_p to determine).
4. Return the concatenation of the lists
[Quicksort(X_{lo}), (x_p), Quicksort(X_{hi})].

0	1	2	3	4	5	6	7	8	9
---	---	---	---	---	---	---	---	---	---

$$\# [T] = \underline{\underline{O(n \log n)}}$$

Quicksort

Quicksort(X): where $X = (x_1, \dots, x_n)$ is a list of numbers.

1. If X is empty: return X .
2. Else: select pivot p uniformly at random from $\{1, \dots, n\}$.
3. Let $X_{lo} = \{i \in X : x_i < x_p\}$ and $X_{hi} = \{i \in X : x_i \geq x_p\}$ (requires $n - 1$ comparisons with x_p to determine).
4. Return the concatenation of the lists
[Quicksort(X_{lo}), (x_p), Quicksort(X_{hi})].



What is the worst case running time of this algorithm?

$O(n^2)$

Randomized Quicksort Analysis

Theorem: Let T be the number of comparisons performed by $\text{Quicksort}(X)$. Then $\mathbb{E}[T] = O(n \log n)$.

Randomized Quicksort Analysis

Theorem: Let T be the number of comparisons performed by $\text{Quicksort}(X)$. Then $\mathbb{E}[T] = O(n \log n)$.

- For any $i, j \in [n]$ with $i < j$, let $\mathbf{I}_{ij} = 1$ if x_i, x_j are compared at some point during the algorithm, and $\mathbf{I}_{ij} = 0$ if they are not. An **indicator random variable**.

Randomized Quicksort Analysis

Theorem: Let T be the number of comparisons performed by $\text{Quicksort}(X)$. Then $\mathbb{E}[T] = O(n \log n)$.

- For any $i, j \in [n]$ with $i < j$, let $I_{ij} = 1$ if x_i, x_j are compared at some point during the algorithm, and $I_{ij} = 0$ if they are not. An **indicator random variable**.
- We can write $\underline{T} = \underline{\sum_{i=1}^{n-1} \sum_{j=i+1}^n I_{ij}}$.

Randomized Quicksort Analysis

Theorem: Let T be the number of comparisons performed by $\text{Quicksort}(X)$. Then $\mathbb{E}[T] = O(n \log n)$.

- For any $i, j \in [n]$ with $i < j$, let $I_{ij} = 1$ if x_i, x_j are compared at some point during the algorithm, and $I_{ij} = 0$ if they are not. An **indicator random variable**. *(are these independent?)*
- We can write $T = \sum_{i=1}^{n-1} \sum_{j=i+1}^n I_{ij}$. Thus, via **linearity of expectation**

$$\underline{\mathbb{E}[T]} = \sum_{i=1}^{n-1} \sum_{j=i+1}^n \mathbb{E}[I_{ij}]$$

Randomized Quicksort Analysis

Theorem: Let T be the number of comparisons performed by $\text{Quicksort}(X)$. Then $\mathbb{E}[T] = O(n \log n)$.

- For any $i, j \in [n]$ with $i < j$, let $I_{ij} = 1$ if x_i, x_j are compared at some point during the algorithm, and $I_{ij} = 0$ if they are not. An **indicator random variable**.
- We can write $T = \sum_{i=1}^{n-1} \sum_{j=i+1}^n I_{ij}$. Thus, via **linearity of expectation**

$$\mathbb{E}[T] = \sum_{i=1}^{n-1} \sum_{j=i+1}^n \mathbb{E}[I_{ij}] = \sum_{i=1}^{n-1} \sum_{j=i+1}^n \Pr[x_i, x_j \text{ are compared}]$$

Randomized Quicksort Analysis

Theorem: Let T be the number of comparisons performed by $\text{Quicksort}(X)$. Then $\mathbb{E}[T] = O(n \log n)$.

- For any $i, j \in [n]$ with $i < j$, let $I_{ij} = 1$ if x_i, x_j are compared at some point during the algorithm, and $I_{ij} = 0$ if they are not. An **indicator random variable**.
- We can write $T = \sum_{i=1}^{n-1} \sum_{j=i+1}^n I_{ij}$. Thus, via **linearity of expectation**

$$\mathbb{E}[T] = \sum_{i=1}^{n-1} \sum_{j=i+1}^n \mathbb{E}[I_{ij}] = \sum_{i=1}^{n-1} \sum_{j=i+1}^n \Pr[x_i, x_j \text{ are compared}]$$

So we need to upper bound $\Pr[x_i, x_j \text{ are compared}]$.

Randomized Quicksort Analysis

Upper bounding $\Pr[x_i, x_j \text{ are compared}]$:

Randomized Quicksort Analysis

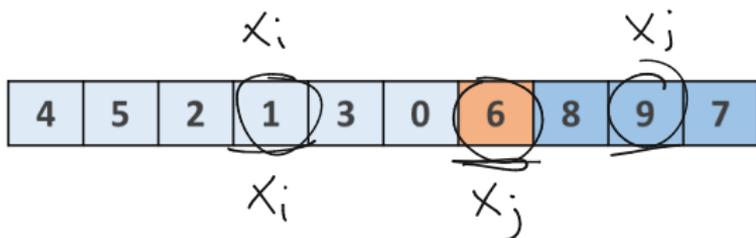
Upper bounding $\Pr[x_i, x_j \text{ are compared}]$: $x_i < x_j$

- Assume without loss of generality that $\underline{x_1} \leq \underline{x_2} \leq \dots \leq \underline{x_n}$. This is just 'renaming' the elements of our list. Also recall that $i < j$.

Randomized Quicksort Analysis

Upper bounding $\Pr[x_i, x_j \text{ are compared}]$:

- Assume without loss of generality that $x_1 \leq x_2 \leq \dots \leq x_n$. This is just 'renaming' the elements of our list. Also recall that $i < j$.
- At **exactly one step of the recursion**, x_i, x_j will be 'split up' with one landing in X_{hi} and the other landing in X_{lo} , or one being chosen as the pivot. x_i, x_j are only ever compared in this later case – if one is chosen as the pivot when they are split up.



Randomized Quicksort Analysis

Upper bounding $\Pr[x_i, x_j \text{ are compared}]$:

- Assume without loss of generality that $x_1 \leq x_2 \leq \dots \leq x_n$. This is just 'renaming' the elements of our list. Also recall that $i < j$.
- At **exactly one step of the recursion**, x_i, x_j will be 'split up' with one landing in X_{hi} and the other landing in X_{lo} , or one being chosen as the pivot. x_i, x_j are only ever compared in this later case – if one is chosen as the pivot when they are split up.
- The split occurs when some element between x_i and x_j is chosen as the pivot. The possible elements are x_i, x_{i+1}, \dots, x_j .

4	5	2	1	3	0	6	8	9	7
---	---	---	---	---	---	---	---	---	---

$j - i + 1$

$$\Pr(x_i \text{ \& } x_j \text{ are compared}) = \frac{2}{j - i + 1}$$

Randomized Quicksort Analysis

Upper bounding $\Pr[x_i, x_j \text{ are compared}]$:

- Assume without loss of generality that $x_1 \leq x_2 \leq \dots \leq x_n$. This is just 'renaming' the elements of our list. Also recall that $i < j$.
- At **exactly one step of the recursion**, x_i, x_j will be 'split up' with one landing in X_{hi} and the other landing in X_{lo} , or one being chosen as the pivot. x_i, x_j are only ever compared in this later case – if one is chosen as the pivot when they are split up.
- The split occurs when some element between x_i and x_j is chosen as the pivot. The possible elements are x_i, x_{i+1}, \dots, x_j .

4	5	2	1	3	0	6	8	9	7
---	---	---	---	---	---	---	---	---	---

- $\Pr[x_i, x_j \text{ are compared}]$ is equal to the probability that either x_i or x_j are chosen as the splitting pivot from this list. Thus,

$$\Pr[x_i, x_j \text{ are compared}] = \frac{2}{j-i+1}$$

Randomized Quicksort Analysis

So Far: Expected number of comparisons is given as:

$$\mathbb{E}[T] = \sum_{i=1}^{n-1} \sum_{j=i+1}^n \Pr[x_i, x_j \text{ are compared}].$$

And we computed $\Pr[x_i, x_j \text{ are compared}] = \frac{2}{j-i+1}$.

Randomized Quicksort Analysis

So Far: Expected number of comparisons is given as:

$$\mathbb{E}[T] = \sum_{i=1}^{n-1} \sum_{j=i+1}^n \Pr[x_i, x_j \text{ are compared}].$$

And we computed $\Pr[x_i, x_j \text{ are compared}] = \frac{2}{j-i+1}$. Plugging in:

$$\mathbb{E}[T] = \sum_{i=1}^{n-1} \sum_{j=i+1}^n \frac{2}{j-i+1}$$

Randomized Quicksort Analysis

So Far: Expected number of comparisons is given as:

$$\mathbb{E}[T] = \sum_{i=1}^{n-1} \sum_{j=i+1}^n \Pr[x_i, x_j \text{ are compared}].$$

And we computed $\Pr[x_i, x_j \text{ are compared}] = \frac{2}{j-i+1}$. Plugging in:

$$\mathbb{E}[T] = \sum_{i=1}^{n-1} \sum_{\substack{j=i+1 \\ \underline{j=i+1}}}^n \frac{2}{j-i+1} = \sum_{i=1}^{n-1} \sum_{\substack{k=2 \\ \underline{k=j-i+1}}}^{n-i+1} \frac{2}{k}$$

Randomized Quicksort Analysis

So Far: Expected number of comparisons is given as:

$$\mathbb{E}[T] = \sum_{i=1}^{n-1} \sum_{j=i+1}^n \Pr[x_i, x_j \text{ are compared}].$$

And we computed $\Pr[x_i, x_j \text{ are compared}] = \frac{2}{j-i+1}$. Plugging in:

$$\begin{aligned} \mathbb{E}[T] &= \sum_{i=1}^{n-1} \sum_{j=i+1}^n \frac{2}{j-i+1} = \sum_{i=1}^{n-1} \sum_{k=2}^{n-i+1} \frac{2}{k} \\ &\leq \sum_{i=1}^{n-1} \sum_{k=1}^n \frac{2}{k} \leq \underbrace{2 \cdot (n-1)}_{=} \cdot \underbrace{\sum_{k=1}^n \frac{1}{k}}_{=} \end{aligned}$$

Randomized Quicksort Analysis

So Far: Expected number of comparisons is given as:

$$\mathbb{E}[T] = \sum_{i=1}^{n-1} \sum_{j=i+1}^n \Pr[x_i, x_j \text{ are compared}].$$

And we computed $\Pr[x_i, x_j \text{ are compared}] = \frac{2}{j-i+1}$. Plugging in:

$$\begin{aligned} \mathbb{E}[T] &= \sum_{i=1}^{n-1} \sum_{j=i+1}^n \frac{2}{j-i+1} = \sum_{i=1}^{n-1} \sum_{k=2}^{n-i+1} \frac{2}{k} \\ &\leq \sum_{i=1}^{n-1} \sum_{k=1}^n \frac{2}{k} \leq 2 \cdot (n-1) \cdot \sum_{k=1}^n \frac{1}{k} = 2n \cdot \underline{H_n} = \underline{O(n \log n)}. \end{aligned}$$