

THE PROBABILISTIC METHOD

WEEK 11: APPLICATIONS, RANDOMIZED ALGORITHMS



JOSHUA BRODY
CS49/MATH59
FALL 2015

READING QUIZ

Which Frequency Moment can be approximated by a randomized algorithm that uses logarithmic space?

- (A) F_0
- (B) F_2
- (C) F_6
- (D) F_{11}
- (E) Multiple Answers Correct

READING QUIZ

Which Frequency Moment can be approximated by a randomized algorithm that uses logarithmic space?

(A) F_0

(B) F_2

(C) F_6

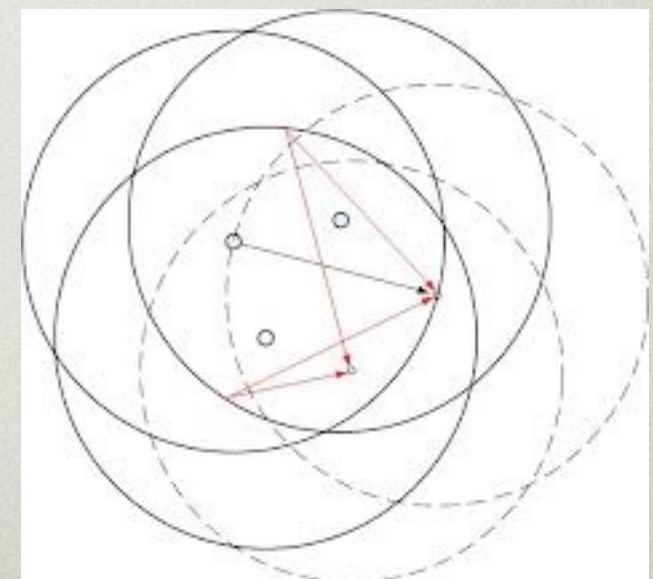
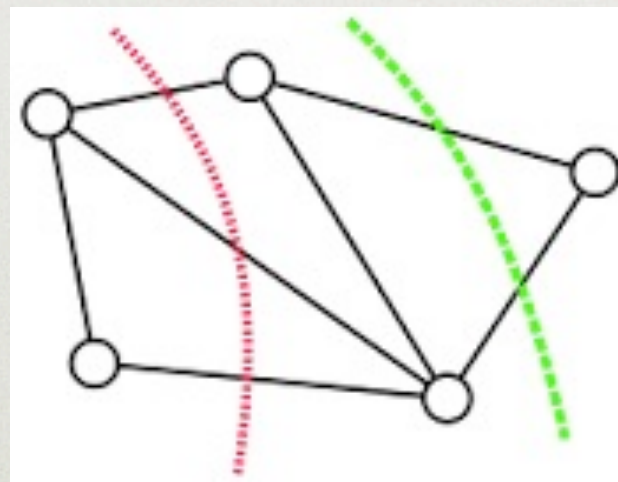
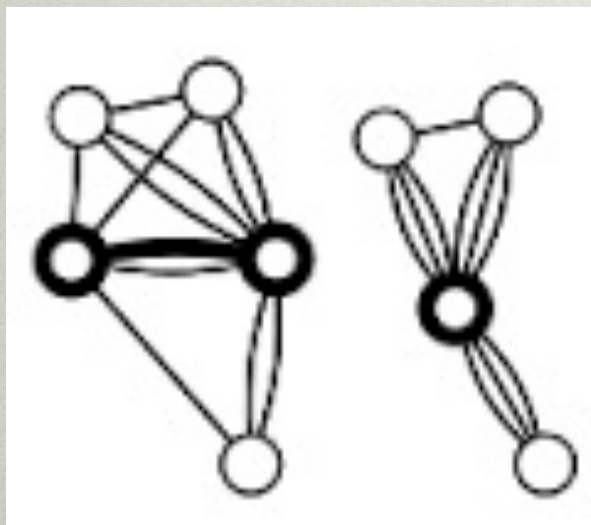
(D) F_{11}

(E) Multiple Answers Correct

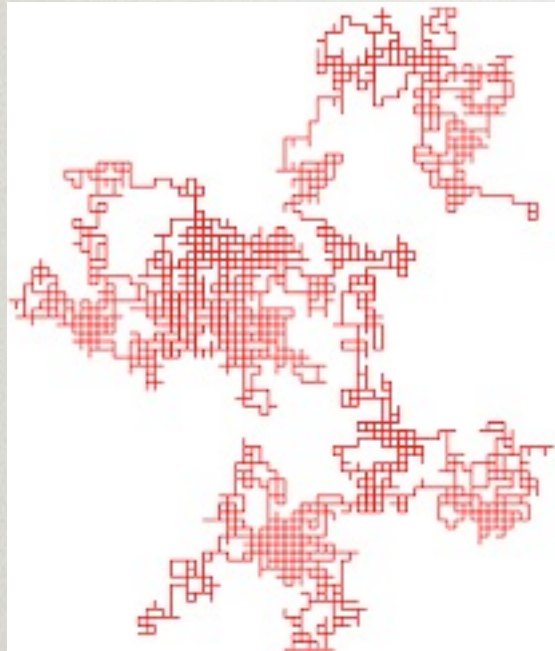
RANDOMIZED ALGORITHMS

Types of randomization:

- **algorithm** uses randomness
 - ▶ Las Vegas: *no error, expected runtime*
 - ▶ Monte Carlo: *small error, worst-case runtime*
- **input** is randomized (algorithm often deterministic)



RANDOMIZED ALGORITHMS

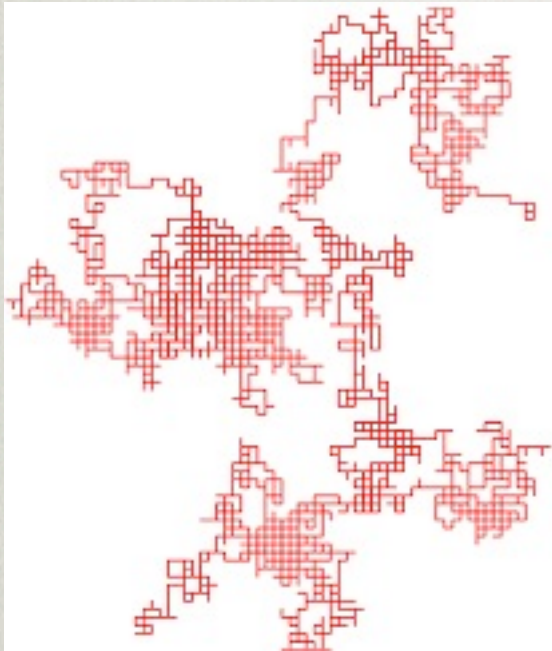


Why randomized algorithms?

- reduce runtime (or space or ...)
- simpler algorithm
- deterministic not known or impossible



RANDOMIZED ALGORITHMS



Why randomized algorithms?

- reduce runtime (or space or ...)
- simpler algorithm
- deterministic not known or impossible

History

- developed late 1970's
- *computational geometry*
- *primality testing* (deterministic algorithm not known until 2003)
- many applications (machine learning, streaming, ...) assume randomness **by default**



CLICKER QUESTION

You see a series of integers $a_1, a_2, \dots, a_m \in \{1, \dots, n\}$ that appear one at a time.

Design a randomized algorithm that outputs a_K for a uniform $1 \leq K \leq m$, using as little space as possible.

How much space is used?

- (A) $O(m)$
- (B) $O(\sqrt{m})$
- (C) $O(\log(m))$
- (D) $O(\log(\log(m)))$
- (E) None of the above

CLICKER QUESTION

You see a series of integers $a_1, a_2, \dots, a_m \in \{1, \dots, n\}$ that appear one at a time.

Design a randomized algorithm that outputs a_K for a uniform $1 \leq K \leq m$, using as little space as possible.

How much space is used?

(A) $O(m)$

(B) $O(\sqrt{m})$

(C) $O(\log(m))$

(D) $O(\log(\log(m)))$

(E) None of the above

DATA STREAM MODEL

DATA STREAM MODEL

- Stream: m elements from universe of size n
 - e.g., 6,9,4,2,1,3,8,8,6,5,7,2,9,3,3,4,7,12, ...
- Goal: Compute function F of stream e.g., # *distinct elements*, *frequency moments*, *entropy*, ...
- Example: $F_0(S)$ = number of *distinct elements* in stream S

DATA STREAM MODEL

- Stream: m elements from universe of size n
 - e.g., $6, 9, 4, 2, 1, 3, 8, 8, 6, 5, 7, 2, 9, 3, 3, 4, 7, 12, \dots$
- Goal: Compute function F of stream e.g., # *distinct elements*, *frequency moments*, *entropy*, ...
- Example: $F_0(S)$ = number of *distinct elements* in stream S
- Caveats: limited space (e.g., $\log(mn)$ or $\log^2(mn)$), sequential access to S

DATA STREAM MODEL

- Stream: m elements from universe of size n
 - e.g., 6,9,4,2,1,3,8,8,6,5,7,2,9,3,3,4,7,12, ...
- Goal: Compute function F of stream e.g., # *distinct elements*, *frequency moments*, *entropy*, ...
- Example: $F_0(S)$ = number of *distinct elements* in stream S
- Caveats: limited space (e.g., $\log(mn)$ or $\log^2(mn)$), sequential access to S
- Stream sources:
 - Network traffic: one pass, data not saved is lost forever
 - Data from hard disk: multiple passes possible, but costly

DATA STREAM MODEL

- Stream: m elements from universe of size n
 - e.g., 6,9,4,2,1,3,8,8,6,5,7,2,9,3,3,4,7,12, ...
- Goal: Compute function F of stream e.g., # *distinct elements*, *frequency moments*, *entropy*, ...
- Example: $F_0(S)$ = number of *distinct elements* in stream S
- Caveats: limited space (e.g., $\log(mn)$ or $\log^2(mn)$), sequential access to S
- Stream sources:
 - Network traffic: one pass, data not saved is lost forever
 - Data from hard disk: multiple passes possible, but costly

Alon, Matias, Szegedy '96 take home messages:

- (i) Randomization usually needed
- (ii) Approximation usually needed

DATA STREAM MODEL

- Stream: m elements from universe of size n
 - e.g., $6, 9, 4, 2, 1, 3, 8, 8, 6, 5, 7, 2, 9, 3, 3, 4, 7, 12, \dots$
- Goal: Compute function F of stream e.g., # *distinct elements*, *frequency moments*, *entropy*, ...
- Example: $F_0(S)$ = number of *distinct elements* in stream S
- Caveats: limited space (e.g., $\log(mn)$ or $\log^2(mn)$), sequential access to S
- Stream sources:
 - Network traffic: one pass, data not saved is lost forever
 - Data from hard disk: multiple passes possible, but costly

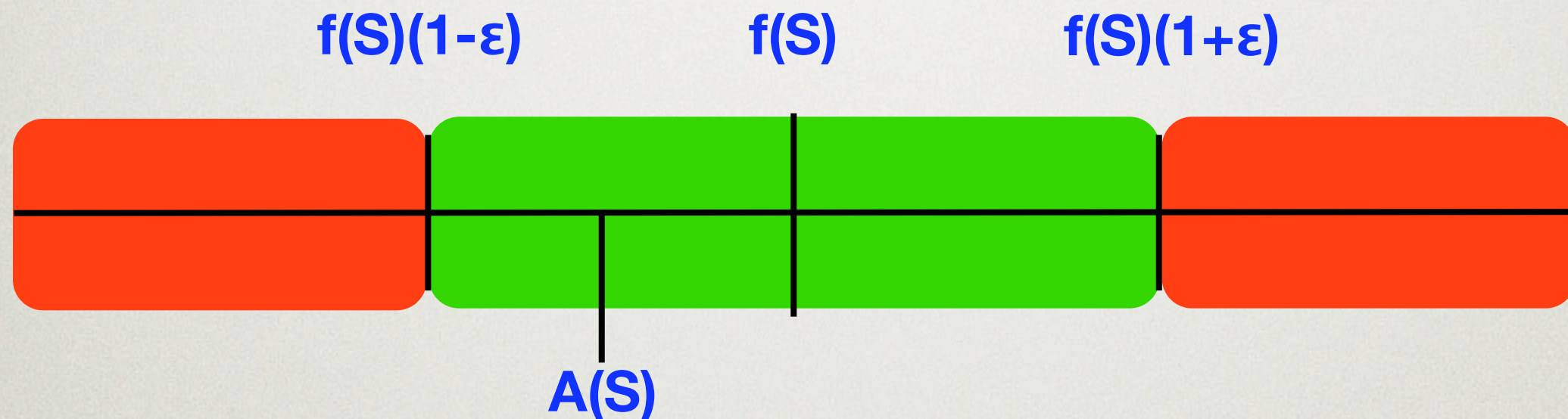
Theorem:

[Alon-Matias-Szegedy '96]

Approximating $F_2(s)$ possible in $O(\log(mn)/\epsilon^2)$ space.

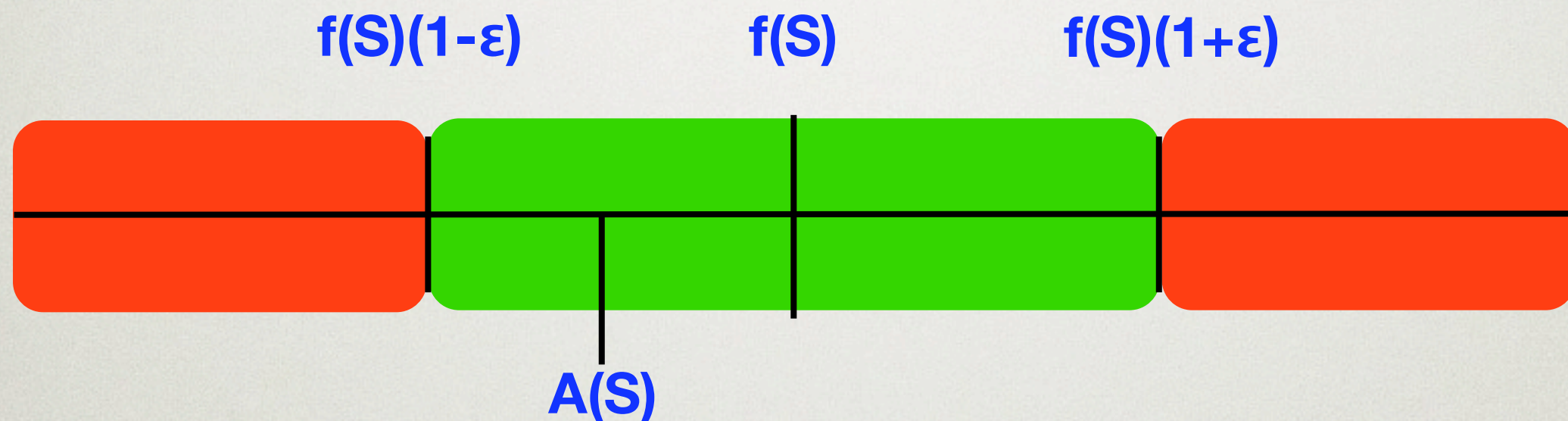
PAC ALGORITHMS

Probably Approximately Correct (PAC) model:
“most of the time we’re close enough”



PAC ALGORITHMS

Probably Approximately Correct (PAC) model:
“most of the time we’re close enough”



Definition: A streaming algorithm **A** is an (ϵ, δ) -approximation for **f** if for any stream **S**

$$\Pr[f(s)(1-\epsilon) \leq A(S) \leq f(s)(1+\epsilon)] \geq 1-\delta \square$$

CLICKER QUESTION

What is $E[H_j]$?

- (A) $E[H_j] = -1$
- (B) $E[H_j] = 0$
- (C) $E[H_j] = 1/2$
- (D) $E[H_j] = f_j$
- (E) None of the above

CLICKER QUESTION

What is $E[H_j]$?

(A) $E[H_j] = -1$

(B) $E[H_j] = 0$

(C) $E[H_j] = 1/2$

(D) $E[H_j] = f_j$

(E) None of the above

CLICKER QUESTION

What is $E[H_i H_j]$?

- (A) $E[H_i H_j] = -1$
- (B) $E[H_i H_j] = 0$
- (C) $E[H_i H_j] = 1$
- (D) It Depends
- (E) Multiple Answers Correct

CLICKER QUESTION

What is $E[H_i H_j]$?

(A) $E[H_i H_j] = -1$

(B) $E[H_i H_j] = 0$

(C) $E[H_i H_j] = 1$

(D) It Depends

(E) Multiple Answers Correct

THE PROBABILISTIC METHOD



Some of us see the world in terms of expected value. We are very different from the rest of you.

www.chalkboardmanifesto.com