

Plan

Logistics

Review

Frequent Items

Wednesday 6pm Games!

Warner 210

Projects: 1 or 2 ppl,

Check out topics  
on home page

Problem set due Friday 6pm

↳ overleaf

↳ colab

➤ combine  
into 7 pdf

Advice: aim to finish  
day assigned

check in form after  
problem solving session

Recommend: read written notes  
day before class

# Review

$X$  random variable

$$X = \begin{cases} 1 & \text{wp } 1/6 \\ 2 & \text{wp } 1/6 \\ 3 & \\ 4 & \\ 5 & \\ 6 & \end{cases}$$

$$\Pr(X=x) \stackrel{\text{dice}}{=} 1/6$$

$$\mathbb{E}[X] = \sum_x x \cdot \Pr(X=x)$$

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

$X, Y$  r.v. independent iff

$$\forall x, y \quad \Pr(X=x | Y=y) = \Pr(X=x)$$

$$\Leftrightarrow \Pr(X=x \wedge Y=y) = \Pr(X=x) \Pr(Y=y)$$

uniform samples

$$D_{i,j} = \begin{cases} 1 & \text{if } i^{\text{th}}, j^{\text{th}} \text{ same} \\ 0 & \text{else} \end{cases}$$

$D_{1,2}, D_{3,4}$  independent!

$D_{1,2}, D_{2,3}$  independent!

$D_{1,2}, D_{2,1}$  not independent

## Facts

$$E[cX] = c E[X]$$

↙ constant

$$\text{Var}(cX) = c^2 \text{Var}(X)$$

$$E[X+Y] = E[X] + E[Y] \text{ always}$$

$$E[XY] = E[X]E[Y] \quad \text{iff uncorrelated}$$

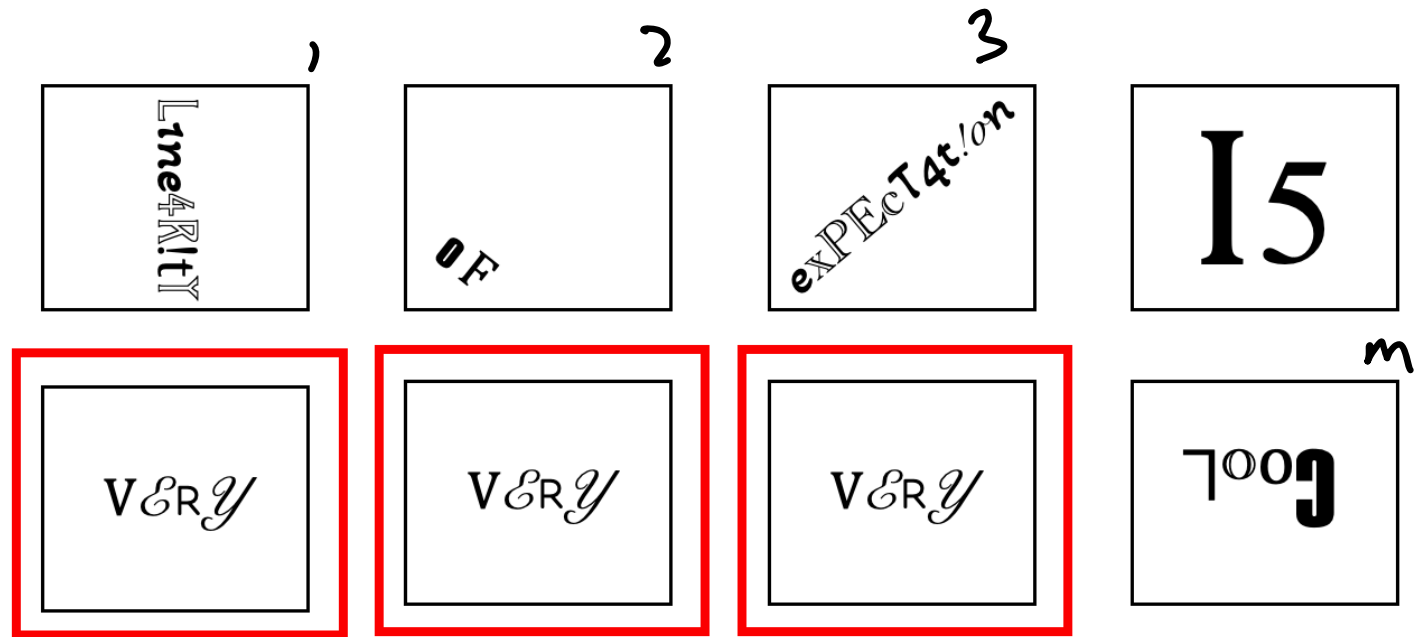
↑ independence

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

$$\text{Var}(X+Y) = \text{Var}(X) + \text{Var}(Y) \quad \text{iff uncorrelated}$$

↑ independence

# Set Size Estimation



$$D = \sum_{i=1}^m \sum_{j=i+1}^m D_{i,j} \quad \leftarrow \# \text{ samples}$$

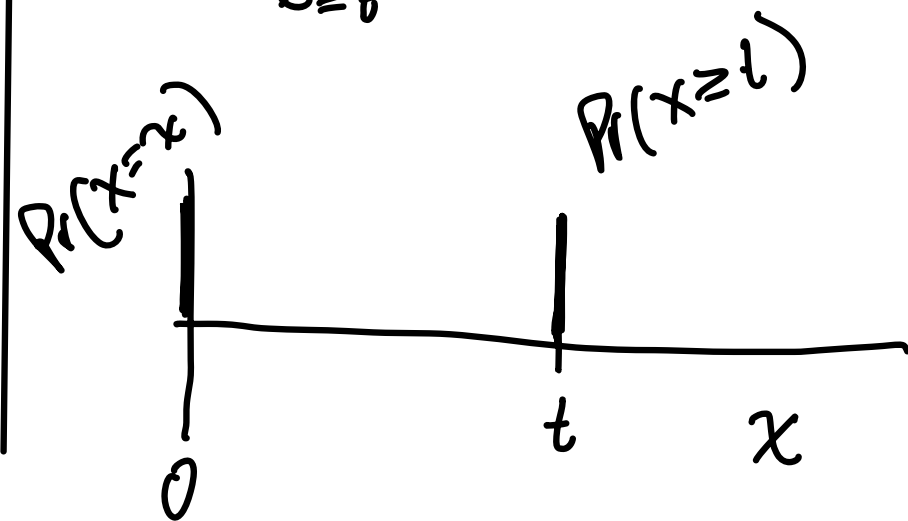
$$E[D] = \frac{m(m-1)}{2 \cdot n} \quad \leftarrow \begin{matrix} \text{set} \\ \text{size} \end{matrix}$$

If we  $\hat{D} \geq E[D]$ ,  
 then  $n \geq \frac{m(m-1)}{2 \hat{D}}$

Markov's  $X$  non-negative,  
 $t > 0$   

$$\Pr(X \geq t) \leq \frac{E[X]}{t}$$

$$\begin{aligned} E[X] &= \sum_x x \cdot \Pr(X=x) \\ &= \sum_{x \geq t} x \Pr(X=x) + \sum_{x < t} x \Pr(X=x) \\ &\geq t \sum_{x \geq t} \Pr(X=x) + 0 = t \Pr(X \geq t) \end{aligned}$$



# Frequent Items Estimation

$x_1, x_2, x_3, \dots, x_n$  ← products, searches, videos

Which are most frequent?

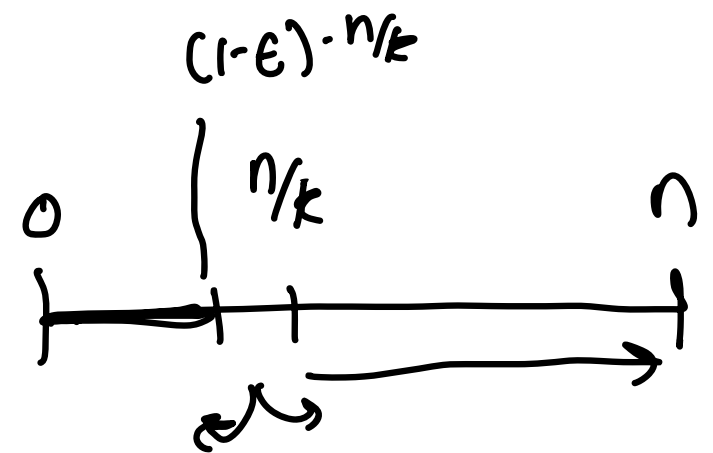
Storing counts takes too much space (esp. if items are pairs)

Params:  $k$  integer,  $0 < \epsilon < 1$  error

Return:

(1) every item that appears  $\frac{n}{k}$  times

(2) only items that appear at least  $(1-\epsilon) \frac{n}{k}$



We'll estimate frequency:

how often item appears

$$f(v) = \sum_{i=1}^n \mathbb{I}[x_i = v]$$

$\begin{cases} 1 & \text{if } x_i = v \\ 0 & \text{else} \end{cases}$

We'll return estimate  $\hat{f}(v)$

$$f(v) \leq \hat{f}(v) \leq f(v) + \frac{\epsilon}{k} \cdot n$$

↑  
one-sided  
error

with probability  $9/10$

then return  $v$  s.t.  $\hat{f}(v) \geq \frac{n}{k}$

(1) If  $f(v) \geq \frac{n}{k}$ ,

$$\hat{f}(v) \geq f(v) \geq \frac{n}{k}$$

(2)  $\frac{n}{k} \leq \hat{f}(v) \leq f(v) + \frac{\epsilon}{k} \cdot n$

$$\Rightarrow \frac{n}{k} \leq f(v) + \frac{\epsilon}{k} \cdot n$$

$$\frac{n}{k} - \frac{\epsilon}{k} \cdot n \leq f(v)$$

$$\Rightarrow (1-\epsilon) \frac{n}{k} \leq f(v)$$

# Hash Functions!

$h: U \rightarrow \{1, \dots, m\}$  ← consistently maps to a random value

- $\Pr(h(x) = i) = \frac{1}{m}$

- $h(x), h(y)$  independent r.v.

$$\Rightarrow \Pr(h(x) = h(y)) \leq \frac{1}{m}$$

# Count-Min Sketch

Choose  $h$  hash-function

Initialize  $m$ -length array  $A$

For every  $x_i$ ,

$$A[h(x_i)] = A[h(x_i)] + 1$$

$$\tilde{f}(v) = A[h(v)] \geq f(v)$$

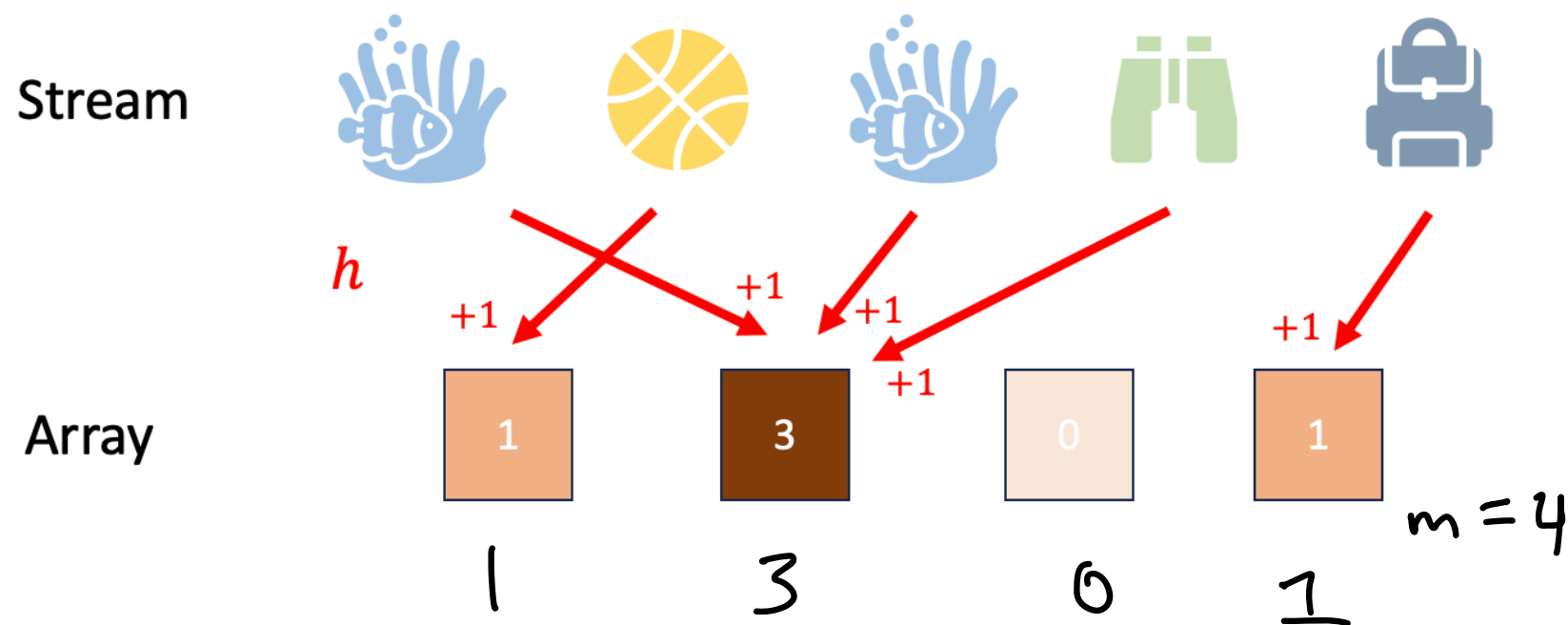
$$\hat{f}(v) = f(v) + \sum_{y \in U \setminus v} f(y) \mathbb{I}[h(y) = h(v)]$$

all other error items

$$(1) \mathbb{E}[\text{error}] \leq \frac{n}{m} \checkmark ;$$

$$(2) \Pr(\text{error} \geq t) \leq 1/2$$

↑  
what is  $t$ ?





$$E \left[ \sum_{y \in U_v} f(y) \mathbb{1}[h(y) = h(v)] \right]$$

$$= \sum_{y \in U_v} E \left[ f(y) \mathbb{1}[h(y) = h(v)] \right]$$

by linearity of expectation

$$= \sum_{y \in U_v} f(y) E \left[ \mathbb{1}[h(y) = h(v)] \right]$$

$$= \sum_{y \in U_v} f(y) \cdot \frac{1}{m}$$

$$= \frac{1}{m} \sum_{y \in U_v} f(y)$$

$$\leq \frac{1}{m} \cdot n = \frac{n}{m}$$

$$1 \cdot \Pr(h(y) = h(v))$$

$$+ 0 \cdot (1 - \Pr(h(y) = h(v)))$$

$$= \Pr(h(y) = h(v))$$

$$\leq \frac{1}{m}$$

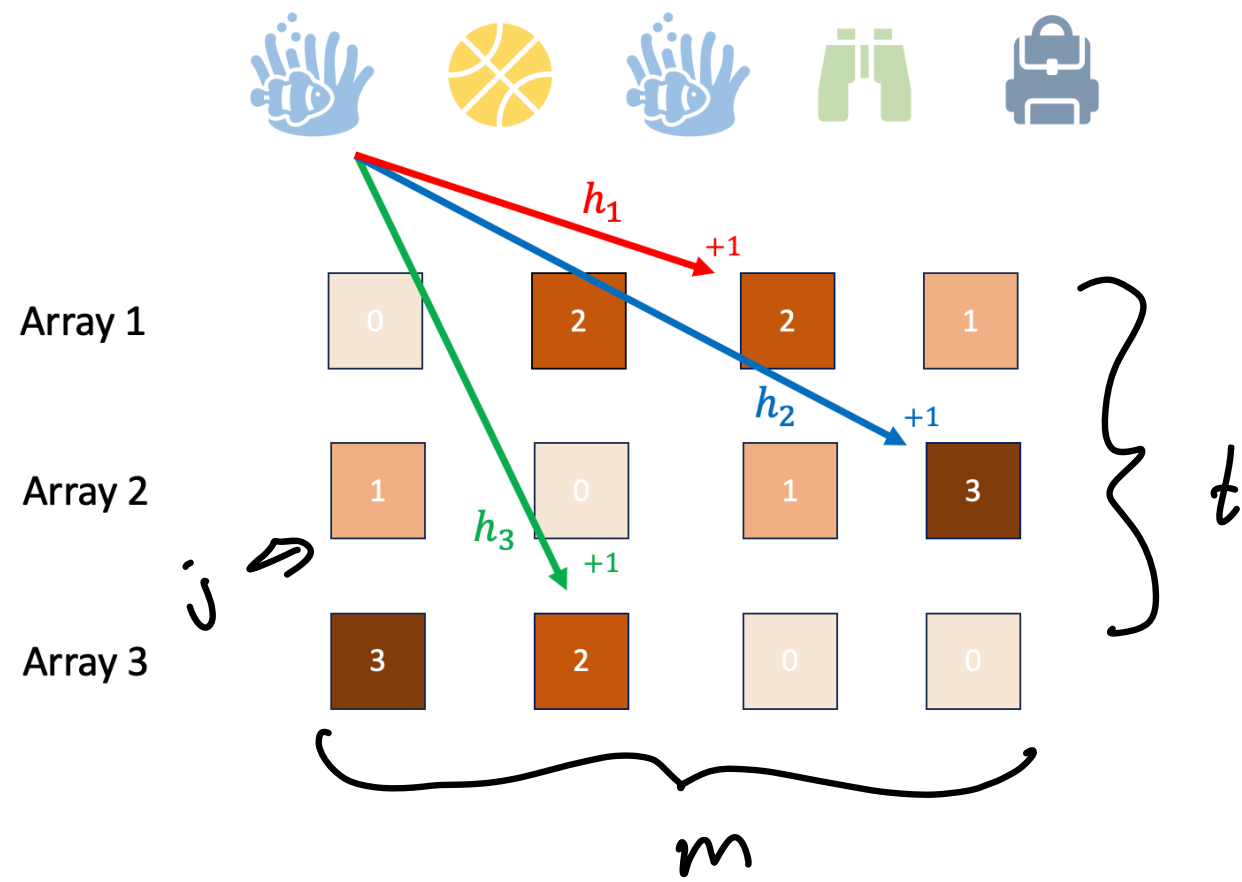
Want:  $\Pr(\text{error} \geq t) \leq 1/2$   
           $\uparrow$   
          what is  $t$ ?

$$\Pr(x \geq t) \leq \frac{\mathbb{E}[x]}{t}$$

$$\Pr(\text{error} \geq t) \leq \frac{\mathbb{E}[\text{error}]}{t} \leq \frac{1/2}{t} = 1/2$$

$$t = \frac{2n}{m}$$

Boost with repetition!



$$\hat{f}(v) = \min \{ A[h_j(v)] : j \in \{1, \dots, t\} \}$$

$$\geq f(v)$$

$$\Pr(\text{error} \geq \frac{2n}{m}) \leq 1/2$$

$$m = \frac{2k}{\epsilon} \Rightarrow \Pr(\text{error} \geq \frac{n\epsilon}{k}) \leq 1/2$$

For every  $j$  w.p  $1/2$

$$f(v) \leq A[h_j(v)] \leq f(v) + \frac{\epsilon n}{k}$$

$$\Pr(\hat{f}(v) \geq f(v) + \frac{\epsilon n}{k})$$

$$\stackrel{\text{indep}}{=} \prod_{j=1}^t \Pr(A[h_j(v)] \geq f(v) + \frac{\epsilon n}{k})$$

$$\leq (1/2)^t = \delta$$

$$\frac{1}{2^t} = \delta \quad \frac{1}{\delta} = 2^t$$

$$\log_2(1/\delta) = t$$

So we have  $f(v) \leq \hat{f}(v) \leq f(v) + \frac{\epsilon n}{k}$  w.p.  $1-\delta$

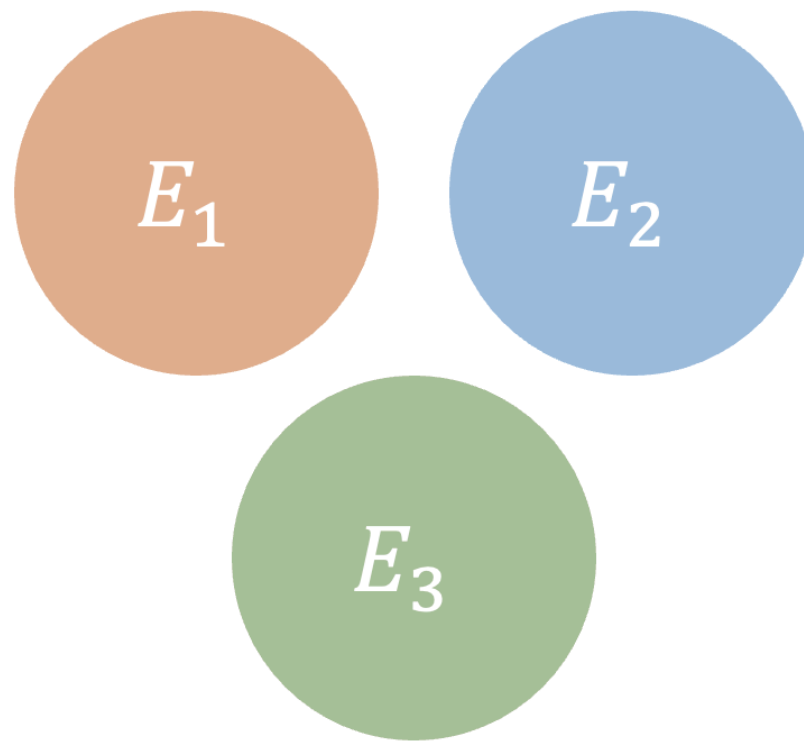
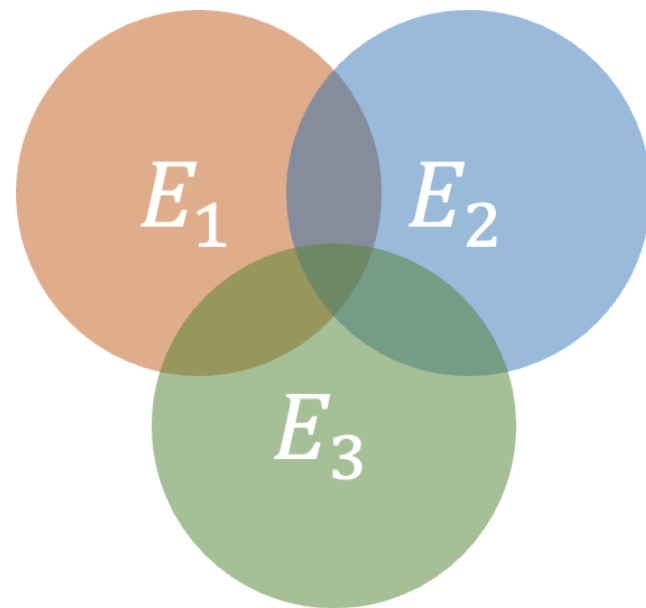
Space  $O(mt) = O\left(\frac{k}{\epsilon} \cdot \log(1/\delta)\right)$

but only holds for one item  $v$

Union Bound

Events  $E_1, \dots, E_n$

$$\Pr(E_1 \cup E_2 \cup E_3 \dots \cup E_n) \leq \Pr(E_1) + \Pr(E_2) + \dots + \Pr(E_n)$$



$$\Pr(\text{fail for } v_1 \cup \text{fail for } v_2 \cup \dots \cup \text{fail for } v_{|U|})$$

$$\stackrel{\text{union}}{\leq} \Pr(\text{fail for } v_1) + \Pr(\text{fail for } v_2) + \dots + \Pr(\text{fail for } v_{|U|})$$

$$\leq \delta + \delta + \dots + \delta = \delta \cdot |U| \leq \delta \cdot n \stackrel{\text{want}}{=} \frac{1}{10}$$

$$\delta^* = \frac{1}{10n}$$

$$O\left(\frac{k}{\epsilon} \log\left(\frac{1}{\delta}\right)\right) \stackrel{*}{=} O\left(\frac{k}{\epsilon} \log 10n\right)$$