

①

Thu 3/21

Last time: Graph Entropy

$$G = (V, E), \quad H(G) = \min_{\substack{X \text{ uniform on } V \\ Y \subseteq V \text{ independent set} \\ X \in Y}} I(X; Y)$$

Subadditivity:  $H(G_1 \cup G_2) \leq H(G_1) + H(G_2)$

Monotonicity:  $H(G_1) \leq H(G_1 \cup G_2)$

Disjoint union:  $G_1, \dots, G_k$  connected components of  $G$ ,

$$H(G) = \sum_{i \in [k]} p_i H(G_i), \quad p_i \triangleq \frac{|V(G_i)|}{|V(G)|}$$

Today: Applications in Combinatorics.  
(for lower bounds)

[ Intuition: A combinatorial process that ~~can~~ produces a "high entropy" graph from "small entropy" ones needs a lot of steps by subadditivity.]

Quick example: What is the smallest number of bipartite graphs that is needed to cover the complete graph  $K_n$ ?

(2)

\* Construction (upper bound):

Represent each vertex by an  $l$ -bit string,  $l := \lceil \log n \rceil$ .

The  $i^{\text{th}}$  bipartite graph connects every two vertices whose binary representation differs ~~at~~ at the  $i^{\text{th}}$  coordinate.

$\Rightarrow l$  graphs are sufficient.

\* Lower bound: ~~Let~~ Let  $G_1, \dots, G_l$  be the bipartite graphs.

$$H(G_1 \cup \dots \cup G_l) = H(K_n) = \log n$$

$$\text{But } \begin{cases} H(G_i) \leq 1 \end{cases}$$

$$\begin{cases} H(G_1 \cup \dots \cup G_l) \leq \sum_{i \in [l]} H(G_i) \leq l. \end{cases}$$

$$\Rightarrow l \geq \lceil \log n \rceil.$$

□

[similar bound can be shown by a chromatic number argument].

3

## Perfect hash function family:

Setting: A database where each "file" is an element of  $[N]$  (eg, a  $\lfloor \log N \rfloor$ -bit string).

A hash function maps a file to a much smaller domain, i.e.,  $h: [N] \rightarrow [b]$ ,  
 $b \ll N$ .

~~Instead of storing each file~~  $x \in [N]$   
[ Instead of storing each file, an indexing database can store the hash  $h(x)$ . ]

In order to search for  $x$  in the database, search for  $h(x)$  in the hash table.

If  $b \ll n$ , a lot of "collision" must occur, i.e.,  $x, y$  s.t.  $h(x) = h(y)$ .

To work around this, we use a "family" of hash functions  $\mathcal{H} = \{ h_1, \dots, h_t \}$ ,  $h_i: [N] \rightarrow [b]$ .

We can hope to design a "nice" family that can differentiate every set of up to  $k$  files.

That is,  $\forall S \subseteq [N], |S| = k, \exists h \in \mathcal{H}$  s.t.

$h$  is injective on  $S$ .

$$(k \leq b)$$

(4)

\* This is called a " $k$ -perfect hash family".

\* Question: How small can  $t = |T|$  be?

\* Upper bound:  $t = O\left(k \log\left(\frac{N}{k}\right)\right)$  is possible,  
where  $b \geq k^2$ .

Proof sketch: Pick each function  $h_i: [N] \rightarrow [b]$   
uniformly randomly and independently.

Fix  $S \subseteq [N]$ ,  $|S| = k$ .

$$\Pr[h_i \text{ is injective on } S] = 1 \cdot \frac{b-1}{b} \cdot \frac{b-2}{b} \cdots \frac{b-k+1}{b}$$

$$\geq \left(1 - \frac{k}{b}\right)^k \geq \left(1 - \frac{1}{k}\right)^k \geq \frac{1}{4}$$

$$\Rightarrow \Pr[\forall i, h_i \text{ is not injective on } S] \leq \left(\frac{3}{4}\right)^t$$

Union bound on all  $S$ :

$$\Pr(\text{family is not perfect}) \leq \binom{N}{k} \left(\frac{3}{4}\right)^t$$

$$\leq \left(\frac{Ne}{k}\right)^k \left(\frac{3}{4}\right)^t = 2^{O\left(k \log\left(\frac{N}{k}\right)\right) - \Omega(t)}$$

$$< 1 \text{ for some } t = O\left(k \log \frac{N}{k}\right).$$

□

⑤

Lower bounds:

$$* t \geq \frac{\log N}{\log b} \quad (\forall k \geq 2)$$

distinct

Proof: For every  $x_1, x_2 \in [N]$ , we must have

$$(h_1(x_1), \dots, h_t(x_1)) \neq (h_1(x_2), \dots, h_t(x_2)).$$

By the pigeonhole principle,

$$N \leq b^t \Rightarrow t \geq \frac{\log N}{\log b}. \quad \square$$

\* Stronger lower bound via graph entropy:

Theorem:  $t \geq \frac{b^{k-1}}{b(b-1)\dots(b-k+2)} \cdot \frac{\log(N-k+2)}{\log(b-k+2)}$

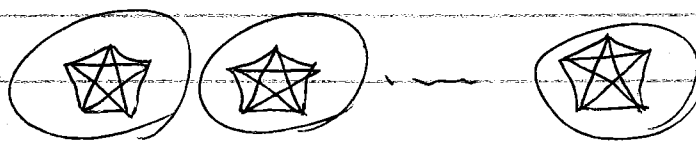
(Fredman, Komlós '84)

Proof (assuming  $b|N$ ).

Define  $G = (V, E)$ .

$$V = \{ (D, x) : D \subseteq [N], |D| = k-2, x \in [N] \setminus D \}$$

$$E = \{ (D, x_1), (D, x_2) : \forall D, x_1 \neq x_2 \}$$



→ a clique of size  $N-k+2$  for each  $D$ .

⑥

$$\Rightarrow H(G) = H(\text{each component}) = \log(N-k+2).$$

Now we construct  $\{G_h\}_{h \in \mathcal{H}}$ , s.t.  $G = \bigcup_{h \in \mathcal{H}} G_h$ .

$$\begin{cases} V(G_h) = V(G) \end{cases}$$

$$\begin{cases} E(G_h) = \{ (D, x_1), (D, x_2) : h \text{ is injective on } D \cup \{x_1, x_2\} \} \end{cases}$$

[For every  $(D, x_1, x_2)$ ,  $\exists h \in \mathcal{H}$  injective on  $D \cup \{x_1, x_2\}$ ]

$$\Downarrow$$
$$G = \bigcup_{h \in \mathcal{H}} G_h$$

all we need to do is to upper bound  $H(G_h)$ .

\* If  $h$  is not injective on  $D$

$$\Rightarrow G_h = \text{empty} \Rightarrow H(G_h) = 0.$$

\* If not,  ~~$G_h$  is  $(b-k+2)$ -partite~~ ~~information only depends on~~

$G_h$  is  $(b-k+2)$ -partite with parts

$$A_i := \{ (D, x) : h(x) = i \}$$

$(i \in h(D))$

subadditivity of  $H(G)$

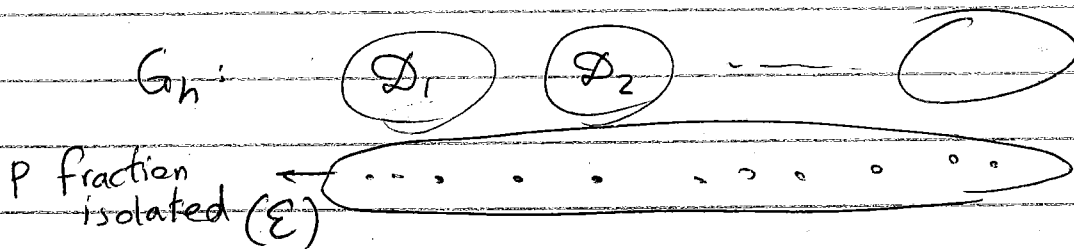
$$\Rightarrow H(G_h) \leq \log(b-k+2) \stackrel{\uparrow}{\Rightarrow} \frac{\log(N-k+2)}{\log(b-k+2)}$$

7

### Improving the bounds:

Observation: each  $G_h$  has a large number of isolated vertices.

~~uniformly random~~  
 $p := \Pr(\text{uniformly random vertex } v \text{ of } G_h \text{ is isolated})$



$$\text{Disjoint union} \Rightarrow H(G_h) = \underbrace{p H(E)}_0 + (1-p) H(D_1 \cup D_2 \cup \dots)$$

$$\leq (1-p) \log(b-k+2).$$

### Bounding p:

$(D, x)$  is isolated iff  $h$  is not injective on  $D \cup \{x\}$

$$\Rightarrow p = \Pr_S(h \text{ is not injective on } S)$$

$|S|=k-1$

Claim: To lower bound  $p$ , wlog we can

assume  $|h^{-1}(1)| = |h^{-1}(2)| = \dots = |h^{-1}(b)| = N/b$ .

8

(of Claim)

Proof Suppose  $|h^{-1}(1)| > |h^{-1}(2)|$ ,

and  $h(x)=1$ . Show that by making  $h(x)=2$ ,

$P$  only gets larger.

increase by  $|h^{-1}(h(x))|$ .  
can only decrease

$$\Pr_S(h \text{ injective on } S) = \Pr(x \in S) \Pr(h \text{ inj. on } S \mid x \in S) + \Pr(x \notin S) \Pr(\text{''} \mid x \notin S)$$

remain unchanged.

$$\text{Now, } \Pr_S(h \text{ injective on } S) = 1 \cdot \frac{b-1}{b} \frac{b-2}{b} \dots \frac{b-k+2}{b}$$

~~$H(G_h)$~~

$$\Rightarrow t \geq \frac{H(G)}{\max_h H(G_h)} \geq \frac{\log(N-k+2)}{(1-p) \log(b-k+2)}$$

$$= \frac{b^{k-1}}{b(b-1) \dots (b-k+2)} \frac{\log(N-k+2)}{\log(b-k+2)}$$

□