

ALGORITHM DESIGN

TL;DR: ONLY WHAT MATTERS, IN DETAIL.

1) KARP-RABIN FINGERPRINT

- HOW TO CHECK EQUALITY OF FILES?

BINARY STRINGS = NUMBERS
 s

DEFINE A FINGERPRINT FUNCTION

$$F(s) = s \bmod p$$

RANDOMLY CHOSEN
PRIME NUMBER

WE HAVE THE PROBLEM OF THE COLLISION

GIVEN TWO NUMBERS (FILES) WE HAVE TWO
CASES :

$$k_1, k_2$$

1) $F(k_1) \neq F(k_2) \Rightarrow k_1 \neq k_2$, OK

2) $F(k_1) = F(k_2) \Rightarrow \begin{cases} k_1 = k_2 & \text{OK} \\ k_1 \neq k_2 & \text{FALSE POSITIVE} \end{cases}$

- 1-SIDED ERROR: WHAT IS THE PROBABILITY

THAT THIS HAPPEN?

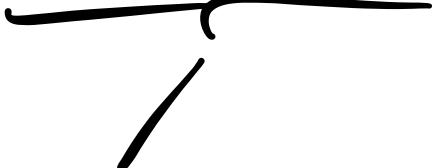
BEFORE WE COMPUTE THIS PROBABILITY, WE NEED
TO:

MORE ON THIS LATER

- a) CHOOSE A PARAMETER τ "SUFFICIENTLY LARGE"
- b) CHOOSE A RANDOM PRIME p SO THAT THE ERROR
PROBABILITY IS LOW: $p \in [\varepsilon, \dots, \tau]$

• LET'S NOW CONSIDER THE FALSE POSITIVE

$$F(k_1) = F(k_2) \wedge k_1 \neq k_2$$



THIS MEANS THAT $k_1 \bmod p = k_2 \bmod p$

OR EQUIVALENTLY $\underbrace{k_1 \equiv k_2 \bmod p}$

CONGRUENCE

WE CALL SUCH p A "BAD PRIME"

WHAT IS THE PROBABILITY THAT WE PICK
A BAD PRIME?



$$P(\text{ERROR}) =$$

$$= P_{p \leq \tau} (k_1 \neq k_2 \wedge F(k_1) = F(k_2)) =$$

$$= \frac{\# \text{BAD PRIMES IN } [2, \tau]}{\# \text{PRIMES IN } [2, \tau]} =$$

$$= \frac{?}{\tau / \log \tau}$$

THIS IS A FACT, STRAIGHT
FROM THE HEAVENS

1) OBSERVATION: k_1, k_2 REQUIRE N BITS FOR BEING
REPRESENTED

2) THEOREM: #DISTINCT PRIME DIVISORS OF ANY
 $\overbrace{\tau}^{\text{NUMBER}} < 2^N$ IS AT MOST N
SKIP THE PROOF.

1) + 2) \Rightarrow AT MOST #BAD PRIMES $\leq N$

/

WE CAN CONCLUDE FROM

$$\leq \frac{N}{\tau / \log(\tau)} = \frac{N}{N^{C+1}} = \frac{1}{N^C} \quad \text{$$

$$\text{WE KNOW THAT } \frac{\tau}{\log \tau} = N^{C+1}$$

WHERE C IS AN ARBITRARY CONSTANT

✳: C IS AN ARBITRARY CONSTANT

WE CAN CHOOSE τ "SUFFICIENTLY LARGE"
TO MAKE C LARGE AND HAVE A VERY
LOW ERROR PROBABILITY

2) CUCKOO HASHING

A DICTIONARY IS A DATA STRUCTURE FOR
STORING ITEMS THAT SUPPORT 3 OPS.

- 1) LOOKUP
 - 2) INSERT (IF NOT PRESENT)
 - 3) DELETE (IF PRESENT)
- N.B. THE PRESENTED VERSION
DONT SUPPORT THE DELETION OP.

WITH CUCKOO HASHING : $O(1)$ WORST-CASE
(WITH THE PROBLEM OF COLLISIONS)



IDEA:

TAKE Two HASH FUNCTIONS FROM THE UNIVERSAL FAMILY \mathcal{H}

h_1, h_2

WE WANT TO INSERT x :

- a) CHECK $T[h_1(x)]$; IF FREE INSERT x
- b) IF NOT FREE, CHECK $T[h_2(x)]$ AND IF FREE STORE IT THERE
- c) IF NOT FREE: PUT x IN $T[h_1(x)]$, THROWS OUT THE CURRENT ELEMENT y
- d) FIND A PLACE FOR y , STARTING WITH $h_2(y)$



THE OPTIONS:

- a) AT SOME POINT YOU FIND A PLACE FOR y
- b) IF THERE IS NO PLACE (OR THE OPERATION IS TAKING TOO LONG) WE REBUILD T

• INSERTION COST

WE TAKE AT RANDOM h_1, h_2 FROM \mathcal{H} .



GIVEN A KEY $x \in U$ AND $i, j \in [M]$, WE HAVE

$$P(h_1(x) = i \wedge h_2(x) = j) \approx \frac{1}{M^2}$$

TAKE IT AS A THEOREM: TO REMEMBER IT
THINK AS

$$P(h_1(x) = i \wedge h_2(x) = j) = \underbrace{P(h_1(x) = i)}_{1/M} \cdot \underbrace{P(h_2(x) = j)}_{1/M} \Rightarrow \frac{1}{M^2}$$

NOW CONSIDER THE SET OF KEYS S AND h_1, h_2
SET WE WANT
TO STORE

WE CAN BUILD AN UNDIRECTED GRAPH $G = (V, E)$
WHERE:

- $V = \{\emptyset, 1, \dots, M-1\}$; THE TABLE POSITIONS
- $E = \{(h_1(x), h_2(x)) : x \in S\}$; THE EDGES
ARE RANDOM

INSERTING A NEW ELEMENT x MEANS ADDING
AN EDGE $(h_1(x), h_2(x))$ IN E .

AT THIS POINT WE HAVE:



- a) EITHER $T[h_1(x)]$ OR $T[h_2(x)]$ IS FREE, OR
- b) WE FOLLOW A PATH IN G UNTIL THE INSERTION COMPLETE
- c) AS b) BUT WE TRANSVERSE A CYCLE: WE HAVE TO REHASH THE WHOLE TABLE

WE NOW ANALYZE THE CASE b) AND c).

WE ASSUME THAT M (THE SIZE OF THE TABLE)

$$IS > \Sigma \cdot C \cdot N$$

$$\text{---} |IS| = N$$

CONSTANT
 $> \Sigma$

b) A FREE POSITION IS ALREADY FOUND

LET i BE THE STARTING POSITION AND j THE FINAL POSITION IN THE PATH.

THE PATH $i \rightsquigarrow j$ HAS LENGTH k IN G WITH THE NOTATION $i \xrightarrow{k} j$

- LEMMA: $P(\exists i \xrightarrow{k} j) \leq \frac{1}{C^k \cdot M}$



PROOF (BY INDUCTION)

- BASE CASE: $K = 1 \Rightarrow i \rightsquigarrow j$ IS AN EDGE

$$P(\exists (i,j) \in E) =$$

$$= \sum_{x \in S} P \left(\begin{array}{l} h_1(x) = i \wedge h_2(y) = j \\ h_1(x) = j \wedge h_2(y) = i \end{array} \right) = \frac{1/M^2}{N} \text{ THEOREM}$$

$|S| = N$

$$= N \cdot 2 \cdot \frac{1/M^2}{N} \leq \frac{M}{2C} \cdot 2 \cdot \frac{1/M^2}{N} = \frac{1}{C^k M} \text{ with } k=1$$

WE ASSUMED $M > 2CN \Rightarrow N < \frac{M}{2C}$

- INDUCTIVE CASE: $K > 1$

WE THEN HAVE: $i \rightsquigarrow_r r \rightarrow j$
 PATH EDGE

$$P(\exists i \rightsquigarrow j) \leq$$

$$\leq \sum_{r \in V \setminus \{i, j\}} P \left(\underbrace{\exists i \rightsquigarrow_r r}_{B} \wedge \underbrace{\exists r \rightarrow j}_{A} \right) =$$

$$= \underbrace{P(\exists i \rightsquigarrow_r r)}_{\text{P}} \cdot \underbrace{P(\exists r \rightarrow j | \exists i \rightsquigarrow_r r)}_{\text{P}}$$

$$P(A \wedge B) = P(A|B) \cdot P(B)$$

INDUCTIVE STEP

$$\leq \frac{1}{C^{k-1} M}$$

$\leq \frac{1}{C^k M}$
BASE CASE

$$\leq \frac{1}{C^{k-1} M} \cdot \frac{1}{C^k M} = \frac{1}{C^{2k} M^2} < \frac{1}{C^{k+1} M} \blacksquare$$

Now, we know that the insertion cost in this case is

$O(1+k)$
COMPUTE THE 1st POSITION $i \mapsto j$
FOLLOW THE PATH $i \mapsto j$

WE WANT TO COMPUTE THE AVERAGE COST!

$$O(1 + E[k]) =$$

$$= O\left(1 + \underbrace{\sum_{k=1}^N k \cdot \frac{1}{C^k M}}_{\text{BY DEF OF EXPECTED VALUE}}\right) =$$

BY DEF OF EXPECTED VALUE

$$= O(1 + 1/M) = O(1), \underline{\text{OK}}$$

$$\sum_{K=1}^H \frac{1}{C^K H} = \frac{1}{H}$$

IS A KNOWN SUMMATION

c) NO FREE POSITION, REHASHING

IN THIS CASE A CYCLE APPEARS AND THUS A REHASHING IS PERFORMED IN $O(N)$ TIME.

WE CHOOSE 2 NEW HASH FUNCTIONS $h'_1, h'_2 \in \mathcal{H}$ AND WE REINSERT EVERYTHING.

WITH WHICH PROBABILITY THIS HAPPENS?

$$\begin{aligned}
 P(\exists \text{ CYCLE in } G) &\leq \underbrace{P(A \cup B)}_{\text{UNION BOUND}} \leq P(A) + P(B) \\
 &\leq \sum_{i=0}^{H-1} \sum_{K \geq 1} P(\exists i \xrightarrow{h'_1} j, j=i) \leq \\
 &\leq \sum_{i=0}^{H-1} \sum_{K \geq 1} \frac{1}{C^K H} = \\
 &= \underbrace{\frac{1}{H} \sum_{i=0}^{H-1} \sum_{K \geq 1} \frac{1}{C^K}}_{\text{TAKEN OUT}}
 \end{aligned}$$

$$\left| < \frac{1}{M} \sum_{i=0}^{M-1} \frac{1}{c-i} = \right.$$

$$\sum_{k \geq 1} \frac{1}{ck} < \frac{1}{c-1}; \text{ know SUM.}$$

$$= \cancel{\frac{1}{M}} \cancel{\frac{1}{c-1}} \cancel{\sum_{i=0}^{M-1}} 1 = \underbrace{\frac{1}{c-1}}$$

WE ASSUMED $c > 2$
 \Rightarrow THE REHASHING OCCURS
 WITH PRB. $p = \frac{1}{c-1} < 1$
 AND TAKES $O(n)$

HOW MANY REHASHING?

$$\begin{aligned}
 & 1 \text{ REHASH WITH PRB. } p \\
 & 2 = = = = p^2 \\
 & \vdots \\
 & t = = = = p^t
 \end{aligned}
 \left. \right\} \text{ THE EXPECTED NUMBER OF REHASHING IS } \sum_{t \geq 1} t \cdot p^t = O(1)$$

$p < 1$

TO CALCULATE: WE ARE INSERTING x THAT WE ASSUME "STICKS" FROM i !

① $P(\exists \text{ CYCLE IN } G, \underline{\text{STARTING FROM } i})$

$$= P(\exists i \rightsquigarrow j, j=i)$$

$$= \sum_{k \geq 1} \frac{1}{C^k M}$$

$$= \frac{1}{M} \cdot \underbrace{\sum_{k \geq 1} \frac{1}{C^k}}_{< \frac{1}{C-1}} < \frac{1}{C-1} \text{ known}$$

$$< \frac{1}{M} \cdot \frac{1}{C-1}$$

$$= \frac{1}{M} \cdot p \text{ as } p = \frac{1}{C-1}$$

$$= \frac{p}{M}$$

WE PAY THE COST OF RELABELING ($O(n)$) WITH PROBABILITY P_M :

THE EXPECTED COST is $O(n \cdot P_M) = O(1)$
OF A SINGLE INSERTION

3) BLOOM FILTERS

WE HAVE A SET S OF KEYS, $S \subseteq U$.

WE DO NOT WANT TO STORE THEM, AS THEY ARE TOO MANY.

THE GOAL IS TO CHECK $x \in S$, without storing S .



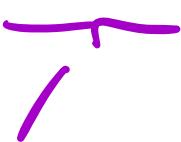
WE TAKE AT RANDOM K FUNCTIONS $h_1, \dots, h_k \in \mathcal{H}$, WHERE $h_i : U \rightarrow [M]$.

THEN WE USE A BITVECTOR B OF SIZE M INITIALIZED TO ALL \emptyset s, AND WE DO

$$x \in S \Rightarrow \forall i \in [M]. B[h_i(x)] = 1$$

THIS TELLS US HOW WE DO THE INSERTION

AND THE QUERY $\overline{x \in S}$



$$\text{RETURN } B[h_1(x)] \wedge \dots \wedge B[h_K(x)]$$

AS BEFORE: THERE IS THE PROBLEM OF 1-SIDED ERROR

$$\{ x \notin S \wedge \bigwedge_{i \in [M]} B[h_i(x)] = 1 \}$$

LET'S SAY THAT THE PROBABILITY OF ERROR IS f :

$$f = P(\text{ERROR}) = P((x \notin S) \wedge (\bigwedge_{i \in [k]} B[h_i(x)]))$$

• TWO POINTS TO CONSIDER:

- 1) WHAT IS f ?
- 2) HOW CAN WE CONTAIN f ?

1) PROBABILITY ERROR f

SUPPOSE THAT THE BLOOM FILTER B IS BUILT, AND TAKE A GENERIC POSITION q AND AN ELEMENT x

STEP a)

DEF OF PROB.

$$P(B[q] = 1) = 1 - P(B[q] = \emptyset)$$

STEP b)

LETS FIX A FUNCTION h_i .

- $h_i \in \mathcal{H} : P(h_i(x) = q) = 1/M$
- $B[q] = \emptyset$ WHEN $h_i(x) \neq q$

$$\begin{aligned} \text{THEN WE HAVE: } P(h_i(x) \neq q) &= 1 - P(h_i(x) = q) \\ &= 1 - 1/M \end{aligned}$$

EQUIVALENTLY WE CAN WRITE

$$P(B[q] = \emptyset \mid h_i \text{ is chosen}) = 1 - \frac{1}{N}$$

STEP C)

LET'S RELAX THE STEP b) : DON'T FIX h_i BUT
CONSIDER THEM ALL :

$$\begin{aligned} P(B[q] = \emptyset) &= P(\forall i \in [k]. h_i(x) \neq q) = \\ &= \left(1 - \frac{1}{N}\right)^k \end{aligned}$$

STEP d)

LET'S RELAX THE STEP c) : DON'T CONSIDER THE
ELEMENT x BUT CONSIDER THEM ALL

$$\underbrace{P(B[q] = \emptyset)}_{\text{PROB } B[q] = \emptyset} = \left(\left(1 - \frac{1}{N}\right)^k \right)^N \quad |S| = N$$

FOR ALL ELEMENTS $x \in S$

AND FOR ALL h_i

STEP e)

WE CAN NOW APPLY STEP d) TO THE STEP a)

$$\begin{aligned}
 P(B[q]=1) &= 1 - P(B[q]=\emptyset) \\
 &= 1 - \underbrace{\left(1 - \frac{1}{M}\right)^{NK}}_{\text{LET'S CALL IT } p'} \\
 &= 1 - p'
 \end{aligned}$$

WE HAVE THEN COMPUTED f :

$$\begin{aligned}
 f = P(\text{ERROR}) &= P((x \notin S) \wedge (\bigwedge_{i \in [k]} B[h_i(x)])) \\
 &= (1 - p')^K \quad \text{EVER POSITION } q \in B
 \end{aligned}$$

2) THE BEST PARAMETERS TO MINIMIZE f

WE WANT TO MINIMIZE $f = (1 - p')^K$ WHERE $p' = \left(1 - \frac{1}{M}\right)^{NK}$

• REMEMBER THE TAYLOR EXPANSION:

$$(1 + r) \approx e^r \text{ FOR SMALL } r$$

$$\begin{aligned}
 p' &= \left(1 - \frac{1}{M}\right)^{NK} \underset{r}{\approx} \left(e^{-\frac{1}{M}}\right)^{NK} = e^{-\frac{NK}{M}} \\
 &\text{LET'S CALL THIS } p \\
 &\text{AND } p' \approx p
 \end{aligned}$$

WE NOW TRY TO MINIMIZE $(1-p)^k$

LET'S EXPLOIT THE LOGS:

$$P = e^{\frac{-NK}{M}} \Rightarrow \underline{\log(p)} = \underline{\log(e^{-NK/M})} =$$

PROP.
LOGS

$$\underline{\log_e(e)} = -\frac{NK}{M} \log(e)$$

$$\underline{\log_e(e)} = 1 = -\frac{NK}{M}$$

$$\underline{\log(p)} = -\frac{NK}{M} \Rightarrow K = -\frac{M}{N} \log(p)$$

- TAKE LOGS:

WE HAVE
FOUND K!

$$f \approx (1-p)^k \Rightarrow \underline{\log(f)} = \underline{\log((1-p)^k)} =$$

$= k \log(1-p) =$

$= -\frac{M}{N} \log(p) \log(1-p)$

K

$$\text{red} \quad \text{Log}(f) = -\frac{M}{N} \text{Log}(p) \text{Log}(1-p)$$

↓
THIS TELLS US "HOW MINIMIZE" STEP:

WE MINIMIZE OVER P, AND THE BEST CHOICE FOR P IS $\frac{1}{2}$; THEN:

$$f \approx (1-p)^k = \left(1 - \frac{1}{2}\right)^k = \left(\frac{1}{2}\right)^k = \frac{1}{2^k} =$$

$$\begin{aligned} & \frac{1}{\sum \frac{M}{N} \text{Log}\left(\frac{1}{2}\right)} = \frac{1}{\sum \frac{M}{N} \text{Log}(2)} = \\ & \text{Log}(a/b) \leftarrow \text{SEE LATER} \\ & \text{Log}(a) = \left(\sum \frac{-\text{Log}(2)}{N}\right)^{\frac{M}{N}} \approx 0.618^{\frac{M}{N}} \quad \left(\begin{array}{l} \text{YOU HAVE} \\ \text{TO ADJUST} \\ \text{M ACCORDINGLY} \end{array} \right) \\ & \text{Log}(b) \\ & + \text{Log}\left(\frac{1}{2}\right) = -1 = 0.618 \end{aligned}$$

A PRETTY LOW ERROR
PROBABILITY

$$\text{Log}(1 - \text{Log}(2))$$

SPACE COMPLEXITY OF BF

THE SPACE COMPLEXITY IS

$M + \underbrace{O(k) \text{ WORDS}}$

BITS TO STORE
THE BINARY

STORE K HASH FUNCTIONS:
STORING A PAIR (a, b)
FOR EACH h_i

CONSIDER f:

$$\textcircled{1} \quad f = \frac{1}{\sum_{i=1}^N \log(2)} \Rightarrow \log(f) = \frac{N}{\sum_{i=1}^N \log(2)} \cdot \log\left(\frac{1}{\sum_{i=1}^N \log(2)}\right) \Rightarrow$$

$\overbrace{\quad \quad \quad}^{= -1}$

$$\Rightarrow N \approx \frac{N \log(f)}{\log\left(\frac{1}{\sum_{i=1}^N \log(2)}\right) \log(2)} = \frac{N \log\left(\frac{1}{M}\right)}{\log(2)} = 1.44 \cdot N \log\left(\frac{1}{M}\right)$$

BITS

4) LOAD BALANCING

M SERVERS, N JOBS

- FAIR LOAD: $\frac{M}{N}$ JOBS PER SERVER

RANDOMLY CHOOSE AN HASH FUNCTION $h \in \mathcal{H}$
USING h SO THAT WE HAVE JOB i ASSIGNED TO THE
MACHINE $h(i)$

↳ THIS GIVES A FAIR LOAD

PROOF:

X_j = LOAD FOR THE MACHINE J = #JOBS ASSIGNED TO J

LET'S DEFINE X_{ij} :

$$X_{ji} = \begin{cases} 1 & \text{IF THE JOB } i \text{ IS ASSIGNED TO } J = h(i) \\ 0 & \text{OTHERWISE} \end{cases}$$

with prob
 $P = 1/M$

SO WE CAN COMPUTE

$$\bullet X_j = \sum_{i=0}^{N-1} X_{ji}$$
$$\bullet E[X_j] = \sum_{i=0}^{N-1} E[X_{ji}] = \sum_{i=0}^{N-1} P(X_{ji} = 1) = \frac{N}{M}$$

$$E[X_j] = E\left[\sum_{i=0}^{N-1} X_{ji}\right]$$

$$\text{DEF OF } X_j = \sum_{i=0}^{N-1} E[X_{ji}]$$

LINERARITY OF E

WHAT ABOUT THE MAX LOAD?

TO COMPUTE THE MAX LOAD WE FIRST NEED
TO REMEMBER TWO RULES:

a) $P(X_{j_i} = 1 \wedge X_{j_{i'}} = 1) = P(X_{j_i} = 1)P(X_{j_{i'}} = 1)$

$\underbrace{\quad}_{\substack{\text{Two DIFFERENT TASKS} \\ i, i' BOTH ASSIGNED TO \\ THE SAME MACHINE j}}$ $\underbrace{\quad}_{\substack{| \\ \text{As } h \in H \text{ IS} \\ \text{Z-WAY INDEPENDENT!}}}$

$$= P(h(i) = j)P(h(i') = j)$$

b) X, Y, Z PAIRWISE INDEPENDENT

INDICATOR
VARIABLES

$$\Rightarrow \sigma^2(X+Y+Z)$$

$$= \sigma^2(X) + \sigma^2(Y) + \sigma^2(Z)$$

VARIANCE:

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

HOW FAR FROM THE MEAN A R.V.'S VALUES
ARE LIKELY TO BE.

FROM a) AND b) WE CAN SAY THAT

$$\sigma^2(X_J) = \sum_{i=0}^{N-1} \sigma^2(X_{J_i})$$

a) TELLS US THAT OUR RANDOM VARS ARE PAIRWISE INDEPENDENT!

THE FACT THAT $X_{J_i}=1$ TELLS US NOTHING ABOUT THE VALUE OF $X_{J_{ii}}$. THIS FOR EACH i, j POSSIBLE. Hence we can use b) with our RVs.

THEN, WE KNOW THAT

$$\sigma^2(X_J) = \sum_{i=0}^{N-1} \sigma^2(X_{J_i}) =$$

$$= \sum_{i=0}^{N-1} \underbrace{\left(E[X_{J_i}^2] - \overbrace{E^2[X_{J_i}]} \right)}_{\text{DEF OF } \sigma^2} =$$

$$= \underbrace{E[X_{J_i}]}_{1/M} : X_{J_i}^2 = \begin{cases} 1^2 & \dots \\ 0^2 & \dots \end{cases}$$

$$\left(\frac{1}{M}\right)^2 = \frac{1}{M^2}$$

$$= N \left(\frac{1}{M} - \frac{1}{M^2} \right)$$

HENCE WE HAVE

$$\sigma = \sqrt{N\left(\frac{1}{M} - \frac{1}{M^2}\right)}$$

NOW, TO COMPUTE (OR BETTER SAID, BOUND) THE MAX LOAD OF A MACHINE WE USE THE **CHEBYSHEV'S INEQUALITY**

$$P\left(\underbrace{|X - E[X]|}_{\text{DISTANCE FROM}} \geq k\sigma\right) \leq \frac{1}{k^2}$$

DISTANCE FROM

X AND ITS EXP. VALUE

A CONSTANT OF
OUR CHOOSING

NOW WE SET $k = \sqrt{2M}$. $\underbrace{\sqrt{2M}}_{\text{K}} \cdot \underbrace{N\left(\frac{1}{M} - \frac{1}{M^2}\right)}_{\sigma} = \sqrt{2} \cdot N\left(1 - \frac{1}{M}\right) \approx \sqrt{2N(1)} = \sqrt{2N}$

THEN $\sigma_k = \sqrt{2M} \cdot \underbrace{N\left(\frac{1}{M} - \frac{1}{M^2}\right)}_{\sigma} \approx \sqrt{2N}$

$$\sqrt{a} \cdot \sqrt{b} = \sqrt{ab}$$

NUMBER OF
SERVERS

AND WE GET

$$E[X_j]$$

$$P\left(|X_j - \underbrace{\frac{N}{M}}_{E[X_j]}| \geq \underbrace{\sqrt{2N}}_{\sigma_k}\right) \leq \frac{1}{\underbrace{2M}_{K^2}}$$

WE CAN THEN COMPUTE THE FOLLOWS PROBABILITY

$$P\left(\max_j : \left|X_j - \frac{N}{M}\right| \leq \sqrt{2N}\right) =$$

$$= 1 - P\left(\exists j : \left|X_j - \frac{N}{M}\right| \geq \sqrt{2N}\right)$$

$$\geq 1 - \sum_{j=\emptyset}^{M-1} P\left(\left|X_j - \frac{N}{M}\right| \geq \sqrt{2N}\right)$$

UNION BOUND

a) \leq b)

$$= 1 - \underbrace{\sum_{j=\emptyset}^{M-1} \frac{1}{2M}}_{\text{CHEBYSHEV}}$$

$$= 1 - \frac{1}{2} \quad / \text{SIMPLIFY THE } \sum_{j=\emptyset}^{M-1} 1 = M$$

$$\geq \frac{1}{2}$$

HENCE:

$$\text{P}\left(\max_j : \left|X_j - \frac{N}{M}\right| \leq \sqrt{2N}\right) \geq \frac{1}{2}$$

THE MAX LOAD (THE LOAD OF THE MOST LOADED MACHINE j) IS AT MOST $N/M + \sqrt{2N}$, WITH PROB.

GREATER THAN $\frac{1}{2}$

5) COUNT-MIN SKETCHES FOR FREQ. ITEMS

CONSIDER A STREAM (POSSIBLY INFINITE) $a_0 \dots$
OF ELEMENTS $\in U$

WE WANT TO KNOW HOW MANY TIMES A GIVEN ELEMENT
 a HAS APPEARED SO FAR IN THE STREAM.

LET'S DEFINE A FREQUENCY ARRAY:

$F[a] = \# \text{TIMES } a \text{ APPEARED SO FAR}$

THE POINT IS THAT $\|F\| = \sum_{a \in U} F[a]$ IS NOT
COMPUTABLE WHEN U IS LARGE (OR WHEN THE STREAM
IS INFINITE)

NON

COMPUTABLE WHEN U IS LARGE (OR WHEN THE STREAM
IS INFINITE) NOT ENOUGH MEMORY

FOR SOLVING THIS PROBLEM WE FIRST SET 2
OBJECTIVES:

- 1) SMALL SPACE OCCUPATION
- 2) 2 WELL DEFINED PARAMETERS:

- δ : THE ERROR PROBABILITY

- ϵ : WE CAN'T COMPUTE $F[\cdot]$ EXACTLY,
SO WE COMPUTE AN ϵ -APPROX

THE ALGORITHM THAT SOLVE THE PROBLEM
IS CALLED COUNT-MIN SKETCH.

WE SET: KINDA THE ALPHABET OF THE STREAM

- $|U| = N$
- UPDATE OP. FOR F : $F[i] \leftarrow F[i] + f_{a,i}$
EVERY $i \in [N]$
- QUERY OP.: RETURN $F[i]$
- $\|F\| = \sum_{a \in U} F[a]$ { DEF }

AND GIVEN ϵ, δ WE COMPUTE AN APPROXIMATED
 $\tilde{F}[i]$ WHERE

$$\tilde{F}[i] \leq \tilde{F}[i] \leq F[i] + \epsilon \|F\|$$

This is TRUE WITH PROB. $1 - \delta$

THE PROBLEM IS THAT IF $F[i]$ IS SMALL COMPARED TO $\epsilon \|F\|$ THE ESTIMATION IS USELESS

• INSIGHT: CM-SKETCH IS A SORT OF BLOOM FILTERS WITH COUNTERS.

GIVEN THE INPUT ϵ AND δ WE BUILD A TABLE $T = R \times C$, WHERE

- R : Rows, $R = \text{WG}(\frac{1}{\delta})$
- C : Columns, $C = \frac{\epsilon}{\delta}$

AND EACH ENTRY OF THE TABLE IS A COUNTER INITIALIZED TO \emptyset .

WE THEN CHOOSE UNIFORMLY AND RANDOMLY

$h_0, \dots, h_{R-1} \in \mathcal{H}$ HASH FUNCTIONS WITH \mathcal{H} A Two-way independent universal hash function, WHERE $M = C$.

↳ SEEN BEFORE IN LOAD BALANCING

THE ITEM $i \in U$ OF THE STREAM IS ASSOCIATED WITH

THE ENTRIES AT POSITIONS

$$\{ \langle i, h_\phi(i) \rangle, \dots, \langle i, h_{R-1}(i) \rangle \}$$

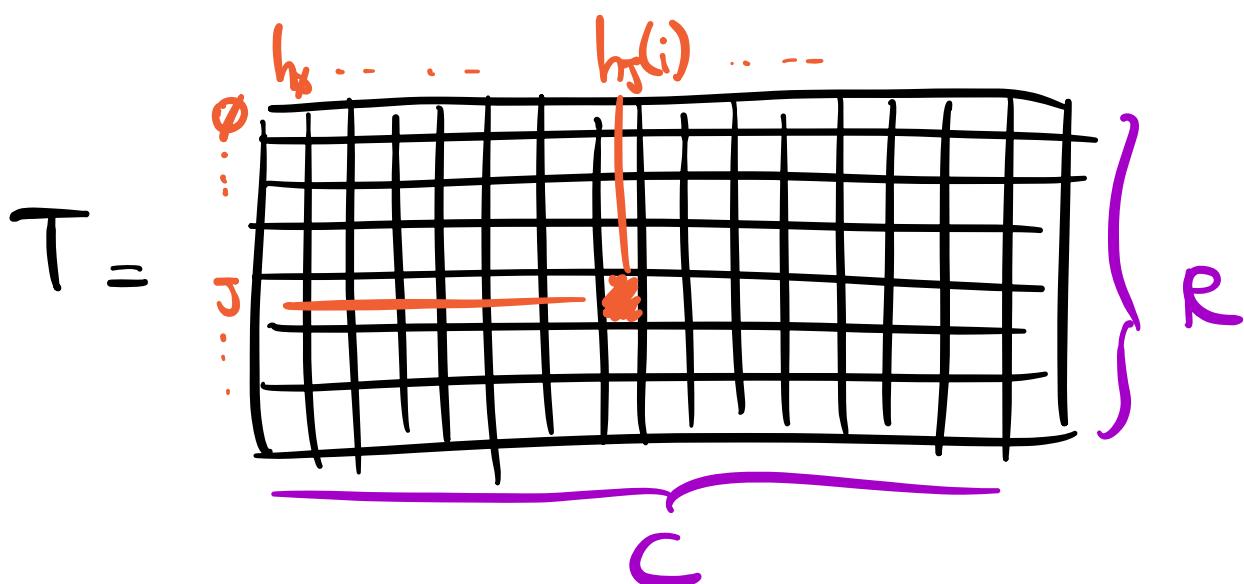
WE THEN HAVE THE TWO OPERATIONS:

- UPDATE: A NEW ELEMENT i IS SEEN IN THE STREAM:

$F[i]++ \equiv$ INCREMENT BY 1 THE CELLS
AT POSITIONS $\langle j, h_j(i) \rangle$
FOR $j \in [\phi, R-1]$

- QUERY:

RETURN $\min_{\phi \leq j \leq R-1} T[j][h_j(i)]$



IT'S CLEAR THAT THE OPERATIONS COST $O(R)$

* ISSUE: $\tilde{F}[i] > F[i]$ BY CONSTRUCTION,
BUT IT COULD BE TOO LARGE!

LET j BE THE MIN OF *

THEN

$$\tilde{F}[i] = T[j][h_j(i)] = F[i] + X_{ji}$$

WHERE X_{ji} IS "GARBAGE" DUE TO THE UPDATES
 $F[k]$ OF ELEMENTS $k \neq i$ (collisions!)

X_{ji} IS INCREASING (W.R.T $F[i]$) THAT THE $T[j][h_j(i)]$ HAD DUE TO THE COLLISIONS OF
ELEMENTS $k \neq i$.

IF ↑ THIS COLLISION HAPPENS ($h_j(i) = h_j(k), i \neq k$),
IT HAPPENS N TIMES WE SEE k IN THE STREAM,
HENCE $F[k]$ TIMES.

MORE SPECIFICALLY:

$$X_{ji} = \sum_{\substack{k \neq i \\ h_j(i) = h_j(k)}} F[k] = \sum_{k=1}^N I_{jik} \cdot F[k]$$

WHERE I_{jik} IS AN INDICATOR VARIABLE!

$$I_{jik} = \begin{cases} 1 & \text{if } h_j(i) = h_j(k), \text{ with } i \neq k \\ 0 & \text{otherwise} \end{cases}$$

WE WANT TO SHOW: $P(X_{ji} > \epsilon | F) < \delta$ ★



AT THE END

WE COMPUTE

DEF OF IND. VAR

$$c = \frac{e}{\epsilon}$$

$$\{ E[I_{jik}] = P(I_{jik} = 1) = \frac{1}{c} = \frac{\epsilon}{e}$$

$h_j: U \rightarrow [c]$ μ IS UNIVERSAL:

THE PROB A COLLISION IS AT MOST $\frac{1}{c}$

WE USE IT TO COMPUTE THE EXPECTED VALUE
OF X_{ji} :

$$E[X_{ji}] =$$

mmmmmmmm

$$\text{DEF OF } - = \sum_{k=1}^N E[I_{jik}] \cdot F[k]$$

X_{ji} + LIKELIHOOD
OF Σ

JUST COMPUTED ↑

$$= \sum_{k=1}^N \frac{\epsilon}{e} \cdot F[k]$$

$$\begin{aligned}
 & \text{TAKEN OUT} = \frac{e}{e} \sum_{k=1}^n F[k] \\
 & \text{NOT APPEAR } k \quad \text{DEF OF } \|F\| \\
 & = \frac{e}{e} \|F\| \quad \text{~~~~~}
 \end{aligned}$$



WE SHOWED

$$\bar{E}[X_{j_i}] = \frac{e}{e} \|F\| \Rightarrow e \bar{E}[X_{j_i}] = e \|F\|$$

WE CAN APPLY THE MARKOV'S INEQUALITY:

$$P(X > z) \leq \frac{\bar{E}[X]}{z}$$

IN THIS CASE WE SET

- $z = e \|F\|$

- $X = X_{j_i}$

AND WE GET: M.I

$$P(X_{j_i} > e \|F\|) \leq \frac{\bar{E}[X_{j_i}]}{e \|F\|} = \frac{\bar{E}[X_{j_i}]}{e \bar{E}[X_{j_i}]} = \frac{1}{e}$$

AND WE CAN FINALLY CAPTURE:

$$P(\underbrace{\tilde{F}[i] > F[i] + \epsilon ||F||}_{T}) = \\ = P(\forall j: \underbrace{P(X_{j,i} > \epsilon ||F||)}_{})$$

$\tilde{F}[i]$ IS A MINIMUM OPERATION OVER THE ROWS OF T : IT HAS TO BE VERIFIED FOR EVERY ROW!

$$\leq \prod_{j=1}^R \frac{1}{e}, \text{ AS } h_j \in \mathcal{H}$$

$\underbrace{\quad}_{\text{COMPUTED BEFORE}} \uparrow$

$$= \frac{1}{e^R}$$

$$= \delta$$

$$R = \log(1/\delta)$$

\Rightarrow

$$\frac{1}{e^{\log(1/\delta)}} = \frac{1}{\frac{1}{\delta}} = \delta$$

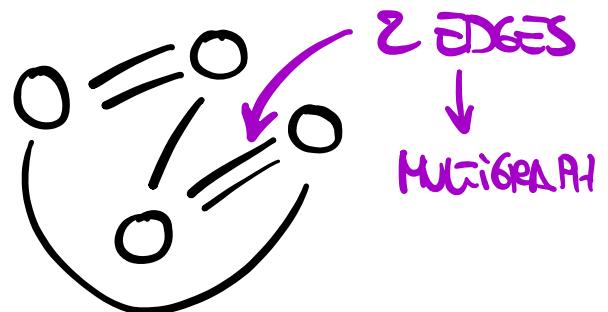
WHICH PROVES $\star \blacksquare$

6) (RANDOMIZED) MIN-CUT ALGORITHM

CONSIDER AN UNDIRECTED (MULTI)GRAPH
 $G = (V, E)$, WHERE

- $|V| = N$
- $|E| = M$

MULTIPLE EDGES BETWEEN
THE SAME NODES



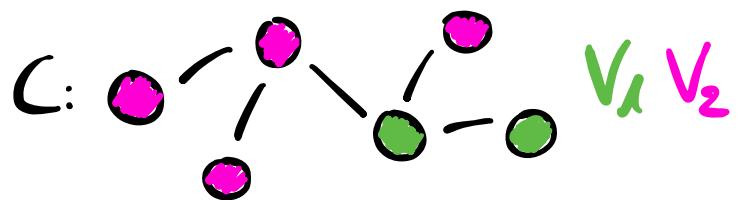
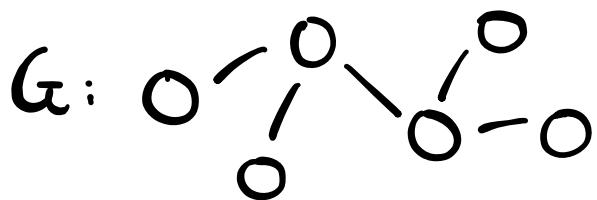
WE DEFINE A CUT OF THE GRAPH

$$C = (V_1, V_2)$$

WHERE

- $V_1 \cap V_2 = \emptyset$
 - $V_1 \cup V_2 = V$
- } WE CALL V_1, V_2 PARTITIONS OF V

PARTITIONS CAN BE SEEN AS COLORINGS



WE THEN DEFINE A CUT-SET OF G

$$\underline{E(V_1, V_2)} = \left\{ \underbrace{uv \in E :}_{\text{UV}} \begin{array}{l} (u \in V_1 \wedge v \in V_2) \\ \vee \\ (u \in V_2 \wedge v \in V_1) \end{array} \right\}$$

SOID EASY: SET OF EDGES THAT CROSS THE PARTITIONS OF A CUT.

ALSO NOTATED WITH $E(C)$

THE PREVIOUS DEFINITIONS ARE USED TO DEFINE SUBGRAPHS OF G:

- NODE INDUCED (SUB) GRAPHS:

$$G[V_i] = G(V_i, E_i) = G(u, v \in V_i, uv \in E)$$

- SUBGRAPH INDUCED BY THE CUT-SET

$$G[E(V_1, V_2)] = (V_{1,2}, E(V_1, V_2))$$

$$\text{WITH } \underline{V_{1,2}} = \left\{ u \in V : uv \in E(V_1, V_2) \right\}$$

"NODES OF THE CUTSET"

COMPUTING THE SIZE OF $G[E(C)]$ IS HARD

- COMPUTING THE LOWER BOUND TAKE POOR TIME

- COMPUTING THE UPP. BOUND IS NP-HARD

WE SAY THAT A CUT $C = (V_1, V_2)$ IS A MIN-CUT IF IS TRUE THAT

$$\forall C'. |E(C)| \leq |E(C')|$$

THE CUT C IS A MIN-CUT IF THE ASSOCIATED CUT-SET IS THE SMALLEST CUT-SET.

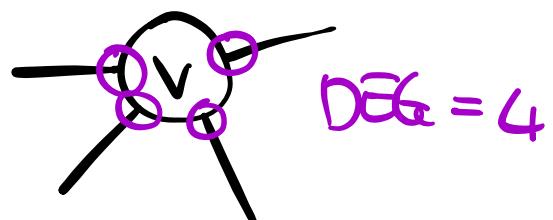
THE CUT WITH THE LEAST AMOUNT OF EDGES

• OBSERVATION: THE MIN-CUT IS NOT UNIQUE

• OBSERVATION: $\forall v \in V. |E(v)| \leq \underbrace{\text{DEG}(v)}$

WITH C MIN-CUT

DEGREE OF V :
EDGES OF V



EDGE CONTRACTION

CONSIDER AN EDGE $uv \in E$

WE DEFINE THE GRAPH $G / \{uv\} = (V', E')$, WHERE

- $V' = (V - \{u, v\}) \cup \{\bar{uv}\}$

- $\tilde{E}' = (E - \{uv\}) \cup \{\bar{u}\bar{v}z : u \in \bar{E} \vee v \in \bar{E}\}$

Contracting AN EDGE uv MEANS TO "MERGE" THE NODES u AND v IN A SINGLE NODE $\bar{u}v$ AND THEN SUBSTITUTE ALL THE EDGES WHERE EITHER u OR v APPEAR WITH EDGES TO/FROM $\bar{u}v$

- INTUITION: WHEN PERFORMING THE EDGE CONTRACTION WITH $u, v \in V$, WE ARE SAYING THAT u, v BELONG TO THE SAME "SIDE" OF THE CUT.

CONSIDER A CONNECTED GRAPH G :

EVERY TWO NODES ARE CONNECTED BY AN EDGE

THEOREM A)

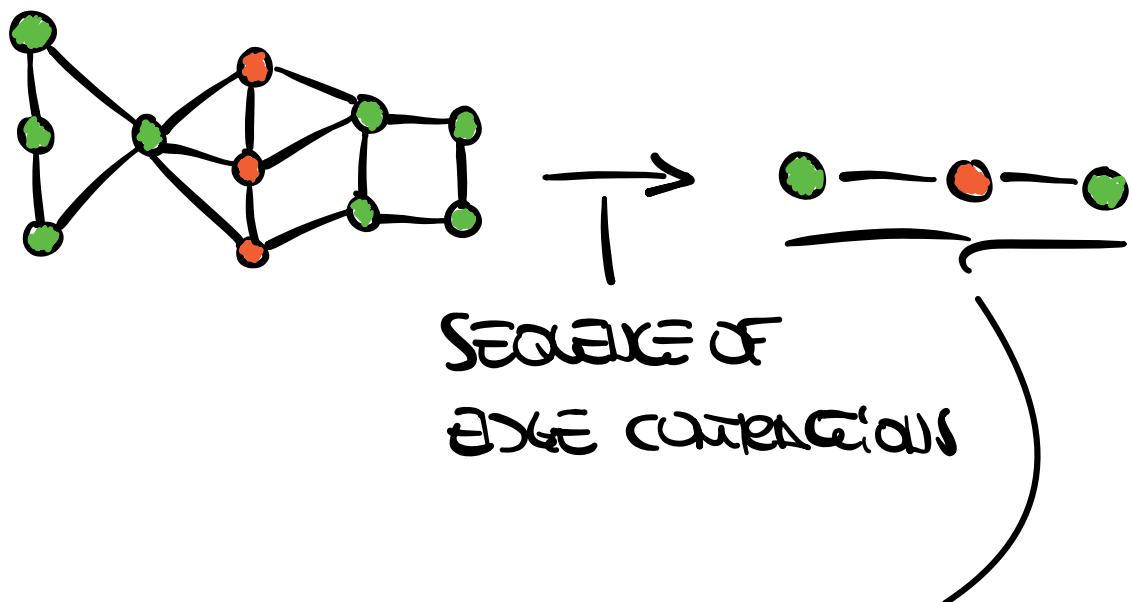
ANY SEQUENCE OF EDGE CONTRACTION LEADING JUST 2 NODES GIVES A CUT $C(V_1, V_2)$

PROOF BY INDUCTION ON #NODES

THEOREM B)

GIVEN A CUT $C = (V_1, V_2)$, IT DOES NOT NECESSARILY MEAN THAT EXISTS A SEQUENCE OF EDGE CONTRACTION THAT LEADS TO JUST 2 NODES

PROOF)



THE CUT $\neq \text{ } \textcolor{red}{\text{ }} \textcolor{green}{\text{ }} \text{ } \text{ } \text{ } \text{ }$ IS NOT OBTAINABLE
BY EDGE CONTRACTION

WE HAVE NOW ALL THE INGREDIENTS TO DEFINE
AN ALGORITHM THAT "GUESSES" THE MIN-CUT OF G

GUESS MIN CUT ($G = (V, E)$) {

WHILE ($|V| > 2$)

{ 1) CHOOSE RANDOMLY $uv \in E$

2) $G = G - uv$ { EDGE CONTRACTION

RETURN $|E|$ { NUMBER OF EDGES THAT REMAIN
WITH JUST 2 NODES { MULTIGRAPH!

HOW WE CHOOSE AN EDGE $uv \in E$?

- CHOOSE u WITH PRWB. $\frac{1}{|V|}$

- CHOOSE $v \in \frac{V(G)}{\overline{T}}$ WITH PROB $\frac{1}{deg(v)}$

NOTATION: SET OF VERTEXES CONNECTED TO v

THANKS TO THEOREM A) THE ALGORITHM CAN PROVIDE THE SIZE OF A MIN-CUT

WHEN THE PREVIOUS ALGORITHM FAIL?

IT FAILS WHEN AN EDGE uv IS A BAD EDGE, WHICH IS AN EDGE THAT BELONGS TO EVERY POSSIBLE MIN-CUT.

IN SYMBOLS:

$\{ uv \in E \text{ IS BAD WHEN } \forall \text{MIN-CUT } C. uv \in E(C) \}$

UV BELONGS TO EVERY MIN-CUT SETS!

IN WORDS: GUESSMINCUT CAN'T FIND A MIN-CUT WHEN IT REMOVES ALL EDGE THAT WOULD BELONG TO EVERY MIN-CUT. → AKA EDGE CONTRACTION

LET $k = |E(C)|$ BE THE SIZE OF A MIN-CUT SET.

PROPOSITION:

THERE ARE AT MOST K BAD EDGES

THE PROOF IS THE INTUITIVE ONE

SO WE HAVE THAT

$$\text{⊗ } \rightarrow P(\text{UV is BAD}) = \frac{k}{M} - |E| = n$$

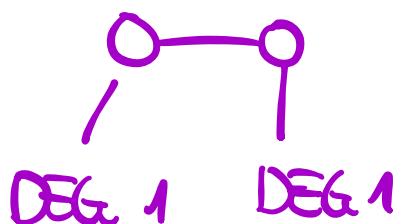
SEE LATER

WE KNOW THAT

$$1) \sum_{u \in V} \text{DEG}(u) = 2M$$

UNDIRECTED GRAPH

EVERY EDGE COUNTED
TWICE



$$2) \forall u \in V: \text{DEG}(u) \geq k$$

OTHERWISE THE MIN-CUT SIZE WOULD BE
SMALLER, BUT k IS THE SIZE OF $E(C)$

IN A GESSMIN CUT CONTEXT: ONLY TWO NODES:

DUMB INSIGHT



$$1) \wedge 2) \Rightarrow \frac{\varepsilon M}{T} > \frac{NK}{M} \Leftrightarrow M \geq \frac{NK}{\varepsilon}$$

1) M 2)

WE CAN USE IT IN \otimes (SEE NSOE)

$$P(\text{uv is bad}) \leq \frac{K}{M} \leq \frac{K}{\frac{NK}{\varepsilon}} = \frac{\varepsilon}{NK} \cdot K = \frac{\varepsilon}{N}$$

WE THEN HAVE

- $P(N) = \text{GUESSMIN FINDS } |ECC| \text{ FOR THE GRAPH } G = (V, E), |V| = N$

$$P(N) \geq \begin{cases} 1 & \text{WITH } N = 2 \\ \left(1 - \frac{\varepsilon}{N}\right) P(N-1) & \text{WITH } N > 2 \end{cases}$$

PRB. THE CHOSEN
EDGE IS NOT A BAD EDGE
GUESS MIN CUT
WITHOUT THE
MERGED NODE

WE EXPAND THAT RECURSIVELY

- NOTE: $1 - \frac{\varepsilon}{N} = \frac{N-\varepsilon}{N}$

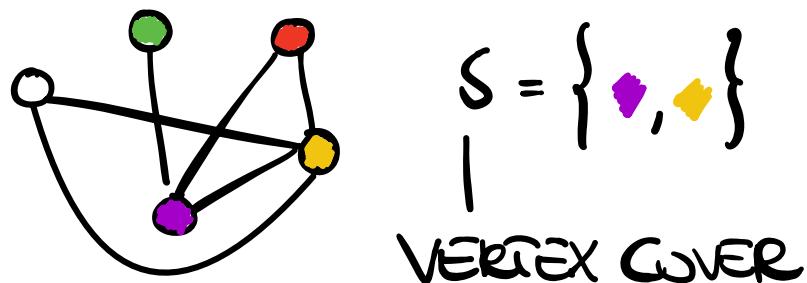
$$P(N) \geq \left(\frac{N-2}{N}\right) \cdot \left(\frac{N-3}{N-1}\right) \cdots \frac{1}{3} \cdot 1 = \frac{(N-2)!}{\frac{N!}{2}} = \frac{2}{N(N-1)} \blacksquare$$

7) VERTEX COVER

WE NOW CONSIDER THE VERTEX COVER PROBLEM, WHICH IS NP-COMPLETE.

GIVEN AN UNDIRECTED GRAPH $G = (V, E)$, WE SAY THAT:

$S \subseteq V$ IS A VERTEX COVER IF $\forall uv \in E : (u \in S) \vee (v \in S)$



WE CAN NOW DISTINGUISH BETWEEN THE VERTEX COVER PROBLEM AND ITS DECISIONAL FORM

1) MINIMIZATION: FIND K_{\min} : $|S| = k_{\min}$ S.T. S IS A VERTEX COVER OF MINIMAL SIZE:

$\forall S' : |S'| < |S| \Rightarrow S'$ IS NOT A V.C.

2) DECISION: CONSIDER VC_K , A VERTEX COVER OF SIZE K , DOES EXISTS A SMALLER VC?

↳ AKA: IS VC_K MINIMAL?

WE FOCUS ON DECISION

- BASELINE FOR VERTEX COVER

THERE ARE $\binom{N}{K} \approx N^K$ SUBSETS S OF V

WE CAN CHECK EACH SUBSET IN POLYTIME,

HENCE THE COST IS $O(N^K \cdot \text{Poly}(n))$

???

$\boxed{(O(N^{K+O(1)})}$

SUBSETS PRICE FOR EACH SUBSET GROSSI NOTES

PARAMETERIZED ALGORITHMS: { EXACT
EXponential
ALGORITHMS }

WE CAN SOLVE THE PROBLEM IN EXACTLY $O(f(k) \text{Poly}(n))$ TIME, WHERE $f(k)$ CAN BE LARGE.

THIS IS USEFUL WHEN $K \ll n$.

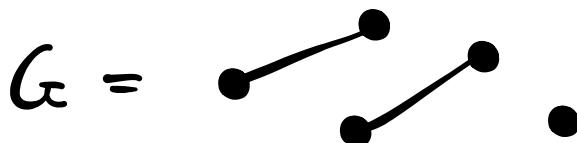


A) KERNELIZATION

CONSIDER AN UNDIRECTED GRAPH $G = (V, E)$.

LET'S CONSIDER THE EASY CASES:

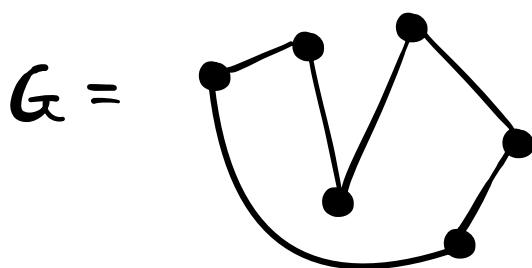
CASE 1:



THE MAXIMUM DEGREE OF A NODE IS 1:

$$\text{MAX_DEGREE}(G) = 1$$

CASE 2:



IN THIS CASE

$$\text{MAX_DEGREE}(G) = 2$$

CONSIDER $G = (V, E)$. WITH $\langle G, k \rangle$ WE MEAN THAT WE WANT TO CHECK WHETHER A VERTEX COVER OF SIZE $\leq k$ DOES EXIST IN G .

WE THEN APPLY THE FOLLOWING TWO RULES AS LONG AS POSSIBLE:

R1) IF v IS AN ISOLATED NODE ($\text{DEG}(v)=\phi$) THEN
SOLVE $\langle G - \{v\}, k \rangle$

R2) IF $\text{DEG}(v) > k$ THEN SOLVE $\langle G - v, k-1 \rangle$

IT'S CLEAR THAT R1) AND R2) TAKES $\text{POLY}(n)$ TIME.
WHEN R1) AND R2) ARE NOT APPLICABLE ANYMORE WE
APPLY THE RULE

R3) • IF $(k < \phi) \vee (G \text{ HAS } > k+k^2 \text{ NODES}) \vee$
|
 $(G > k^2 \text{ EDGES})$



• THEN:

REPORT NO VERTEX COVER

• ELSE:

G_K BE THE CURRENT GRAPH (G REDUCED)

WITH R1 AND R2), SOLVE VERTEX COVER

ON $\langle G_K, k \rangle$ AND RETURN THE ANSWER

THE REDUCED GRAPH WITH R1) AND R2) IS
SAID THE KERNEL OF G .

THAT'S THE KERNELIZATION.

LUDR. MY SPECULATIONS

R3)

IF $\langle G_K, k \rangle$ HAS A V.C. OF SIZE $\leq k$ THEN G_K MUST HAVE $\underbrace{\leq k^2 + k \text{ NODES}}_1$ AND $\underbrace{\leq k^2 \text{ EDGES}}_2$

LET S BE SUCH V.C.

• $|S| \leq k$ — DEF OF $\langle G, k \rangle$
1) • $\forall u \in S. \deg(u) \leq k$ } $\Rightarrow |E| \leq |S| \cdot k = k^2$
• S IS A VC
DEF OF $\langle G, k \rangle$ PERHAPS? TODO

2) $\underbrace{|S|}_{\leq k} + |G_K - S| \leq k + k^2$

THE APPLICATION OF THE RULES TAKE POLY TIME, SO:

- 1) WE TAKE A GRAPH
- 2) APPLY THE RULES AND OBTAIN $\langle G_K, k \rangle$
- 3) RUN THE BASELINE ON $\langle G_K, k \rangle$

$$O(\underbrace{((k^2)^k + \text{POLY}(n))}_{\text{FUNCTION IN } k, \text{BASELINE COST}}) = O(\underbrace{f(k) + \text{POLY}(n)}_{\text{F}})$$

B) BRANCHING TECHNIQUE

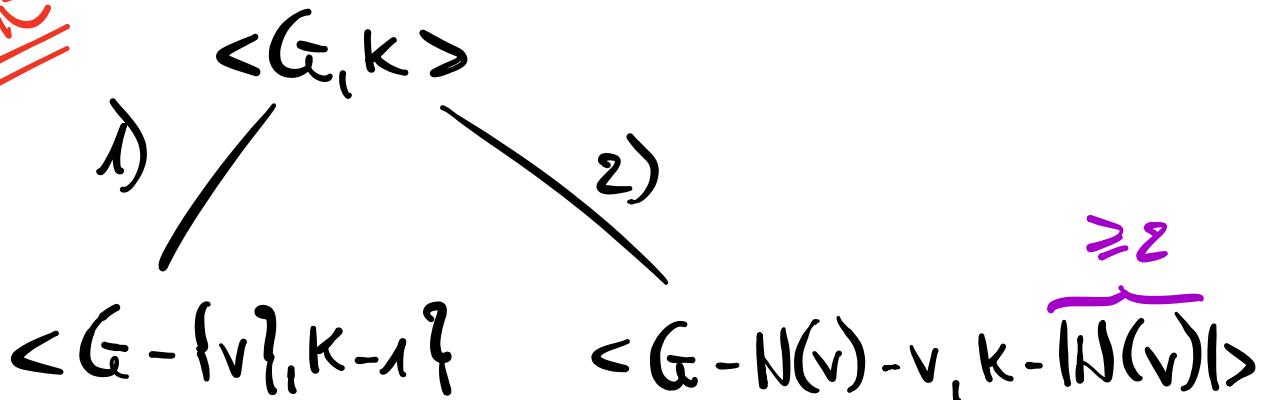
LET'S CONSIDER $\langle G, k \rangle$ AS BEFORE

LET'S DEFINE v AS THE NODE IN G WITH MAXIMUM DEGREE, WE WILL HAVE RECURSION ON THIS DEGREE.

- BASE CASE: $\deg(v) = 1$
- RECURSIVE CASE: $\deg(v) > 1$

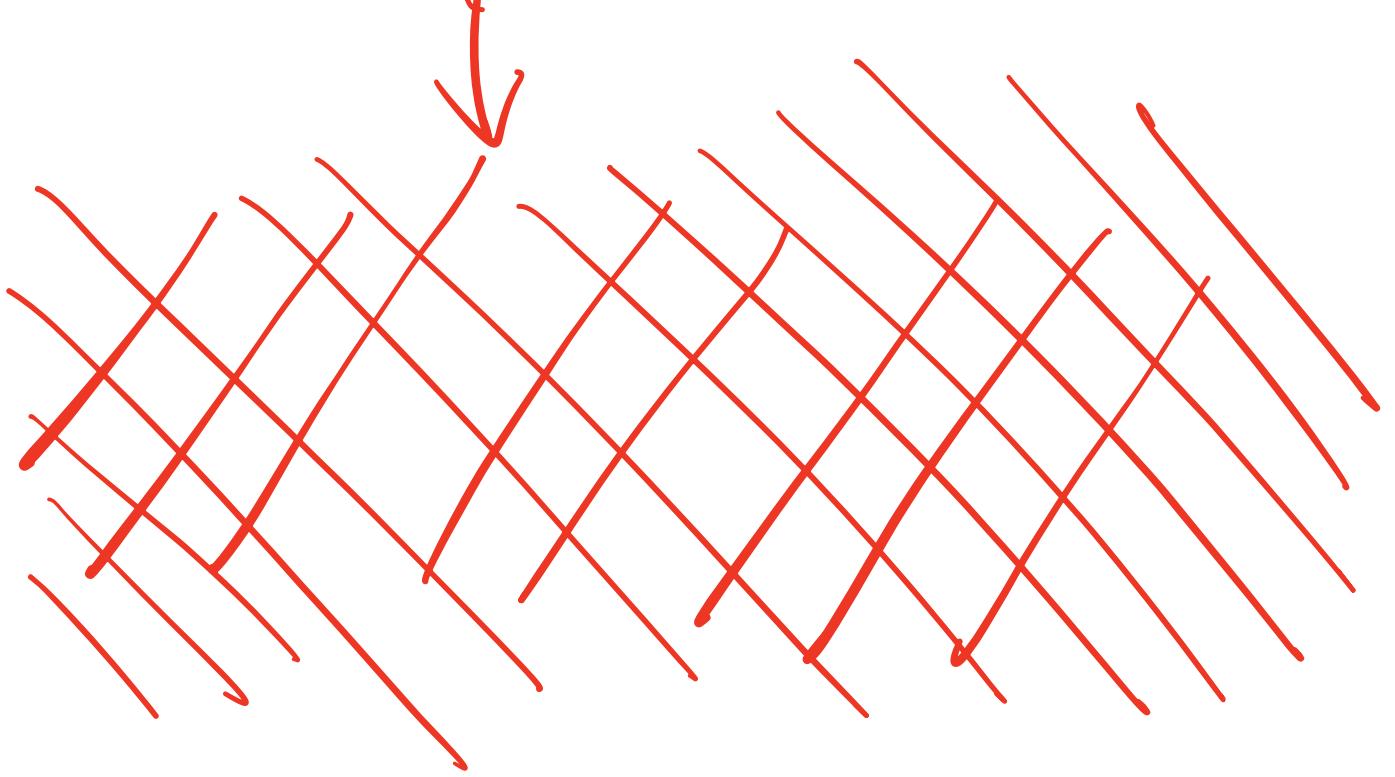
- (1) EITHER $v \in VC$, AND WE CAN'T SAY ANYTHING FOR THE NODES $\in N(v)$, OR
(2) EVERY NODE IN $N(v)$ BELONGS TO VC

COME



LEAVES:
BASE CASE

TODO: COME FUNZIONA?



THE TOTAL COST IS GIVEN BY

$$\# \text{RECURSIVE CALLS} \cdot \text{POLY}(n)$$

EACH CALL GENERATES TWO CALLS, HECKE

$$\# \text{CALLS} \leq 2 \cdot \# \text{LEAVES}$$

AND THE LEAVES ARE THE BASE CASE(S): $L(k)$

$$L(k) = \underbrace{L(k-1)}_{1)} + \underbrace{L(k - \lceil \log_2 k \rceil)}_{2)} \leq \geq 2$$

$\leq L(k-1) + L(k-2)$

IT'S FIBONACCI: $L(k) \sim \Phi^k$, $\Phi = 1.6\dots$

THE TOTAL COST IS THEN $O(\Phi^k \text{POLY}(n))$

BRANCHING + KERNELIZATION

WE CAN PUT THINGS TOGETHER.

VC COSTS:

$$O(n^k)$$

↓ KERNELIZATION

$$O\left(\underbrace{((k^2)^k}_{f(k)} + \text{POLY}(n))\right)$$

? — ↓ BRANCHING

$$O(\Phi^{k^2} + \text{POLY}(n))$$

=

$$O\left(\underbrace{1.6^{k^2}}_{T} + \text{POLY}(n)\right)$$

EVEN THOUGH IT CAN BE EXPONENTIAL
IN n , IN MOST CASES k IS SMALL

TODO / \equiv

QUESTO $\in K, n$?

PS: STO ROBO ↑ E' W CASO DI FPT?
PENSARE K?

A SNEAK-PEEK TO APPROXIMATING VC

MORE DETAILS IN A FOLLOWING SECTION

CONSIDER A GRAPH $G = (V, E)$.

BY DEF S IS A VERTEX COVER IF

$$\forall v \in V. \exists s \in S. v \in s$$

AND WE THEN HAVE THE DECISIONAL PROBLEM VC_k

$$\exists? S \subseteq V. |S| \leq k \wedge S \text{ is a V.C.}$$

WE NOW DEFINE VC^* , THE PROBLEM OF
FINDING THE SMALLEST VERTEX COVER OF G

$$VC^* = \underline{k_{\min}}. \forall S \subseteq V. |S| < k$$

$$\Rightarrow$$

S IS NOT A VERTEX COVER

HOW WE FIND k_{\min} ?

NOT BY DUMB BINARY SEARCH : k IS AT THE EXPONENT.

WE SIMPLY START WITH $k=1$ AND SOLVE, AS SEEN BEFORE, $\langle G, k \rangle$.

IT WORKS, BUT IT IS SLOW.

WE CAN APPROXIMATE THE SOLUTION.

WE DO A 2-APPROXIMATION, A CASE OF A TECHNIQUE CALLED R-APPROXIMATION

WE RELAX THE OPTIMALITY: INSTEAD OF LOOKING FOR S^* , WE LOOK FOR A VERTEX COVER $\tilde{S} \subseteq V$ SUCH THAT $|\tilde{S}| \leq \frac{2}{R} \cdot K_{\min}$

THE PROCEDURE FOR FINDING \tilde{S} IS:

$$\tilde{S} = \{\}$$

FOR EACH uv IN \bar{E} :

IF ($u \notin \tilde{S}$) \wedge ($v \notin \tilde{S}$)

$$\tilde{S} = \tilde{S} \cup \{u, v\}$$

RETURN \tilde{S}

AND WE SEE THAT

1) \tilde{S} IS A VC BR DESIGN

2) $|\tilde{S}| \leq 2 \cdot K_{\min} = 2 \cdot |S^*|$

PROOF:

S^* : BY DEF $\forall u, v \in E, u \in S^* \text{ or } v \in S^*$

CONSIDER \tilde{S} : FOR EVERY TWO EDGES THAT "CONTRIBUTES" TO \tilde{S}

a) NO ENDPOINT IN COMMON

ALL THE EDGES THAT CONTRIBUTE TO \tilde{S}
ARE ABOUT DIFFERENT NODES.

IT CAN'T BE THAT $uv \wedge uy$ OR ANY OTHER
SIMILAR VARIATION.

b) THE #EDGES ARE $\frac{|\tilde{S}|}{2}$, UNDIRECTED
GRAPH
THAT CONTRIBUTE TO \tilde{S}

c) EACH OF SUCH EDGE HAS AT LEAST A NODE IN \tilde{S}

d) b) + c) $\Rightarrow \frac{|\tilde{S}|}{2} \leq |S^*| \equiv |\tilde{S}| \leq 2|S^*| \blacksquare$

8) RANDOM FPT ALGORITHMS

FPT: FIXED-PARAMETER TRACTABLE

A DESIGN TECHNIQUE FOR SOLVING COMBINATORIALLY HARD (MOSTLY NP-HARD) PROBLEMS.

FOR SOME OF THIS PROBLEMS FPT ACHIEVE BOTH EFFICIENCY AND OPTIMACY OF THE SOLUTION

$O(f(k) \cdot n^c)$ WITH c A CONSTANT

UNDER FPT FALLS TO APPROACHES:

A) COLOR CODING

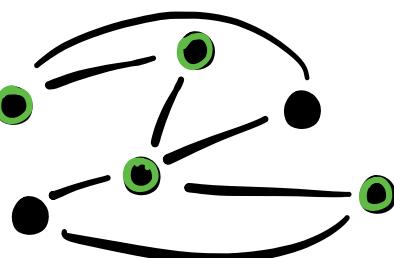
B) RANDOMIZED SEPARATION } SKIPPED

A) Color Coding

CONSIDER AN UNDIRECTED GRAPH $G = (V, E)$

WE DEFINE A K-PATH A PATH OF K NODES WHICH ARE ALL DISTINCT

$k=4$



THE PROBLEM OF FINDING THE K-PATH IS NPC
AS IF $K = N = |V|$ THIS IS HAM.

WE ATTACK THE K-PATH PROBLEM BY COLOR
CODING.

CONSIDER A K-COLORING $X: V \rightarrow [K]$ THAT
ASSOCIATE TO EVERY VERTEX A COLOR.

RANDOMLY!

$$\hookrightarrow \forall v \in V. \forall c \in [K]. P(X(v) = c) = \frac{1}{K}$$

OBS. BASICALLY $X \in \mathcal{H}$

THEN WE SAY THAT:

A SET $V' \subseteq V$ IS COLORFUL



$\forall u, v \in V'$ WITH $u \neq v. X(u) \neq X(v)$

SAY EASY: V' IS COLORFUL IF EVERY NODE IN
 V' IS MAPPED TO AN "UNIQUE" COLOR.

WE THEN HAVE THE FOLLOWING PROPOSITIONS:

P.a) A COLORFUL WALK OF K NODES IS A K-PATH

T

- WALK: SEQUENCE OF NODES (AND EDGES) WHERE REPETITION IS ADMITTED
- PATH: A WALK IN WHICH NEITHER THE NODES NOR EDGES ARE REPEATED
- THE OPPOSITE IS NOT ALWAYS TRUE (OVR, COLORS ARE RANDOMLY ASSIGNED)

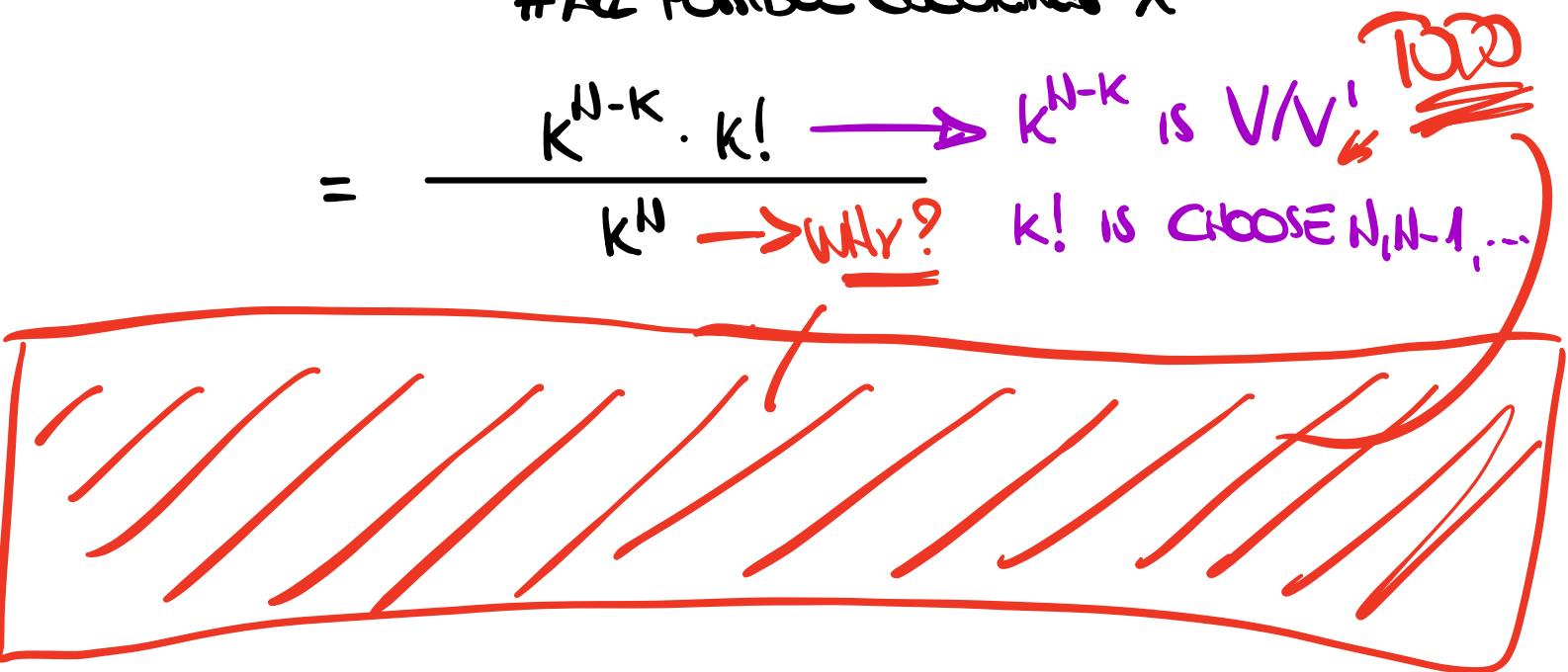
P.b) A WALK OF k NODES IS A k -PATH IF AND ONLY IF EXISTS A COLORING X S.T. THE WALK IS COLORFUL.

Consider $V' \subseteq V$ s.t. $|V'| = k$

$P(V' \text{ is colorful}) =$

$$= \frac{\# \text{colorings } X \text{ that makes } V' \text{ colorful}}{\# \text{all possible colorings } X}$$

$$= \frac{k^{N-k} \cdot k!}{k^N} \rightarrow \underline{\text{Why?}} \quad k! \text{ is choose } N, N-1, \dots$$



$$= \frac{k!}{k^k} - \frac{k^{N-k}}{k^N} = \frac{\frac{k^N}{k^k}}{\frac{k^N}{k^N}} = \frac{k^N}{k^k} \cdot \frac{1}{k^{kN}} = \frac{1}{k^k}$$

$\geq e^{-k}$ $k! \geq \left(\frac{k}{e}\right)^k$ known fact

$$\Rightarrow \frac{k!}{k^k} \geq \frac{\left(\frac{k}{e}\right)^k}{k^k} = \frac{k!}{k^k} \geq \frac{\frac{k^k}{e^k}}{k^k}$$

$\overbrace{\phantom{\frac{k^k}{e^k}}}$

$$= \frac{1}{e^k} = e^{-k} \blacksquare$$

QUESTION: GIVEN A RANDOM COLORING X , HOW DO WE FIND COLORFUL PATHS (k -PATHS)?

DYNAMIC PROGRAMMING

- a) ANY SUBPATH OF A COLORFUL PATH IS A COLORFUL PATH.
- b) A COLORFUL PATH OF LENGTH J ENDS IN A NODE U IF AND ONLY IF EXISTS A NODE $V \in N(U)$ SUCH THAT A COLORFUL PATH OF LENGTH $J-1$ ENDS IN V AND NONE OF THE NODES IN

THAT PATH HAVE A COLOR $X(v)$

c) IN THE COLORFUL PATH OF LENGTH $J-1$ ONLY THE COLORS MATTER, NOT THEIR ORDER OR PATHS.

a) + b) + c) : WE CAN USE DP AND BUILD A TABLE (AKA MATRIX) PATH AS FOLLOWS

TRUE

$$\text{PATH}[S][v] = \begin{cases} 1 & \text{IF } S = \{X(v)\} \\ \text{OR}_{v \in N(v)} \text{PATH}[S \setminus X(v)][v] & \text{IF } |S| \geq 2 \\ \emptyset & \text{otherwise} \end{cases}$$

IS THERE A COLORFUL PATH THAT ENDS IN v AND THAT USES THE COLORS IN S

POWER SET OF k COLORS (S)

THE TABLE PATH HAS $\sum_{k=1}^{k-1}$ ROWS AND $N = |V|$ COLUMNS

HENCE THE FINAL RESULTS WILL BE IN THE LAST ROW, WHERE $|S| = k$

EACH ENTRY $\text{PATH}[S][v]$ TAKES $O(|N(v)|)$ TIME

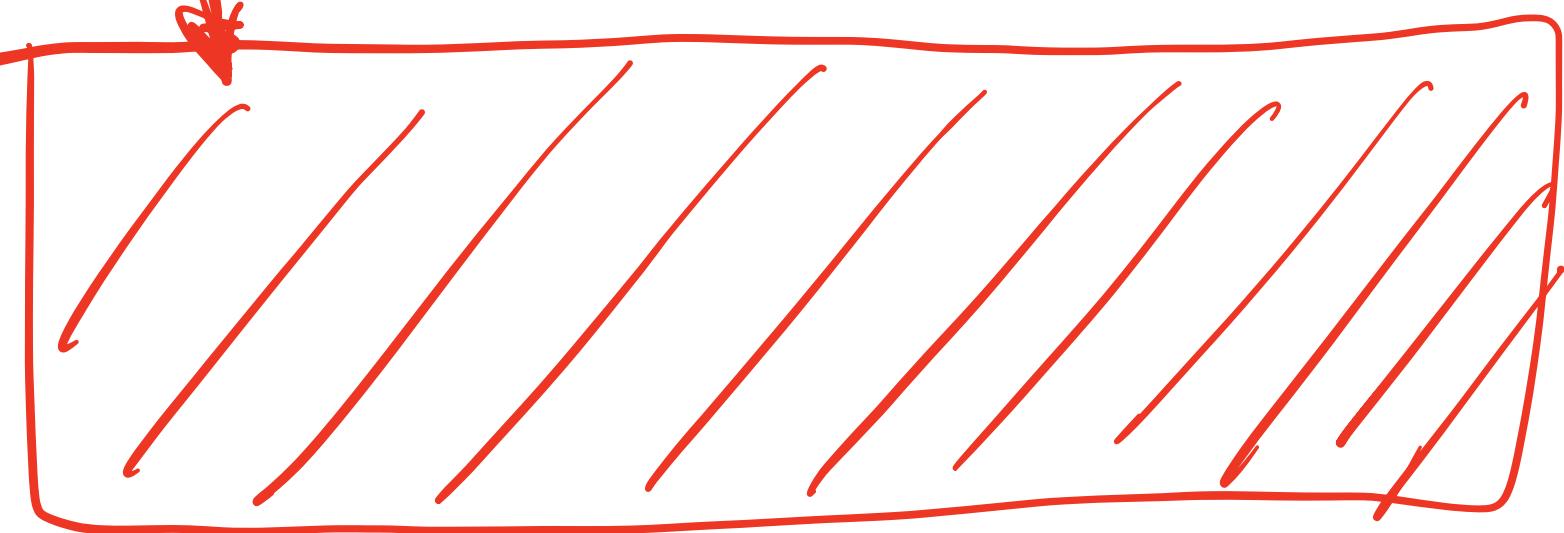
TO FILL.

THEN THE TOTAL COST IS:

$$\sum_{v \in V} |N(v)| =$$

TODO 

$$\left. \begin{array}{l} = O(M) \text{ PER ROW, } M = |E| \\ = O(2^k \cdot M) \\ = O(2^k \cdot N^2) \end{array} \right\}$$



FINAL TOUCH

PRB. OF FIND 1 COLORFUL, COMPUTED ↑

WE REPEAT THE ABOVE FOR \bar{e}^{-k} TIMES, EACH TIME RANDOMLY CHOOSING A COLORING X .

WE STOP WHEN WE FIND A 1 IN THE LAST ROW.

WE RECALL THE UNION BOUND: CONSTANT PROBABILITY OF SUCCESS, THE RUNNING TIME IS

$$O(\bar{e}^{-k} 2^k N^2) = O\left(\frac{2^k}{\bar{e}^k} N^2\right)$$

$O(f(k) \cdot n^c)$, FPT ■

WHY THE COLOR CODING IS POWERFUL?

TAKE A NODE U AND CONSIDER ALL THE k -PATHS ENDING IN U .

POTENTIALLY, WE MAY HAVE UP TO n^k SUCH PATHS,
AND A SIGNIFICANT FRACTION OF THEM ARE COLORFUL
IN X . — PROPOSITION P.a)

FOR EACH COLORFUL k -PATH, THE SET OF COULE IS
THE SAME, $[k]$

WITH JUST ONE SET WE REPRESENT ALL THE COLORFUL
K-PATH.

IN GENERAL: FOR A j -PATH, WHERE $j \in [k]$, WE
HAVE $\binom{k}{j}$ SET OF COULES TO REMEMBER, INSTEAD
OF $O(n^j)$ j -PATHS

g) R-APPROXIMATION

HARDNESS

- P: PROBLEMS THAT CAN BE SOLVED IN POLY TIME

- NP: PROBLEMS THAT ADMIT A POLY-TIME CERTIFICATE
- NP-HARD: PROBLEMS THAT ARE AT LEAST AS HARD AS THE HARDEST PROBLEM IN NP
- NP-COMPLETE: PROBLEMS BOTH IN NP AND NP-HARD

TO ATTACK SUCH HARD PROBLEMS:

a) EXACT SOLUTIONS: FPT

WE HAVE SEEN K-PATH

b) APPROXIMATE SOLUTIONS (FOR OPTIMIZATION PROBLEMS)

- MIN-COST
- MAX-COST

A TECHNIQUE CALLED R-APPROXIMATION,
THE APPROXIMATED SOLUTION \tilde{S} IS SUCH THAT

$$\bullet \text{MIN} = \frac{\text{COST}(\tilde{S})}{\text{OPT}} \leq R$$

$$\bullet \text{MAX} = \frac{\text{OPT}}{\text{COST}(\tilde{S})} \leq R$$

WITH $R > 1$.

WE HAVE ALREADY APPLIED THIS TECHNIQUE WITH VERTEX COVER.

10) NP-COMPLETE PROBLEMS: TSP

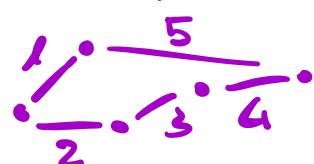
THE TRAVEL SALESMAN PROBLEM

- N CITIES
- D_{ij} : THE COST/DISTANCE FROM THE CITY i TO THE CITY j
- WE DEFINE A TOUR A PERMUTATION OF ALL THE CITIES

IT'S CLEAR THAT

$$\text{COST OF THE TOUR} = \sum_{i=1}^{N-1} (D_{c_i, c_{i+1}}) + D_{N,1}$$

Tour: $i, i+1, \dots, i+N-1, i$



AND WE WANT TO MINIMIZE THE COST.

THEOREM: TSP ∈ NPC

1) THERE IS NO R-APPROX WITH $R < 1$ FOR TSP,
UNLESS $P = NP$.

ALTERNATIVELY SAID: APPROXIMATING TSP IS
AS HARD AS SOLVING TSP

2) IT CAN BE REDUCED IN HAM IN POLY-TIME
(SKIPPED)

• METRIC TSP

WE CONSIDER A SUB-VERSION OF TSP AND
TO FIND IT WE USE R-APPROX.

IN THIS VERSION OF THE PROBLEM THE GRAPH
 $G = (V, E)$, IS COMPLETE BY HYP. AND THE EDGES
RESPECT THE

TRIANGLE INEQUALITY: T.I

$$D_{ij} + D_{jk} \geq D_{ik} \quad D_{ik} \leq D_{ij} + D_{jk}$$

$G = (V, E, D)$ $\forall v \in V, \exists u \in E$

| SET OF COSTS FOR EVERY EDGE

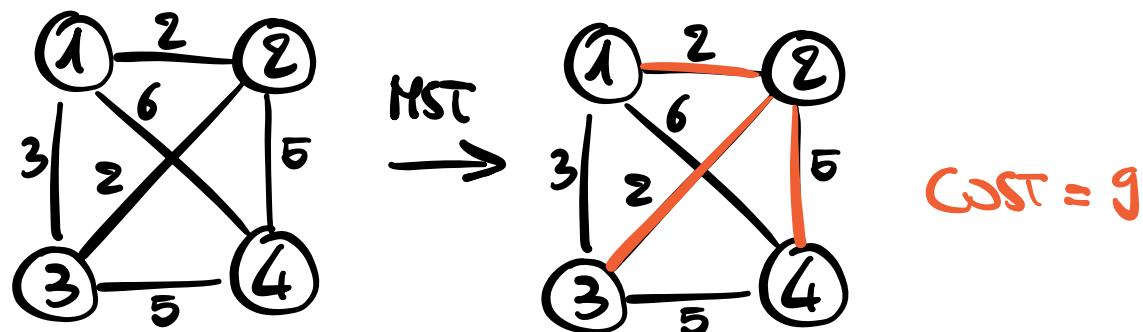
G IS WEIGHTED

THE ALGORITHM TO FIND THE SOLUTION IS THIS:

- 1) CONSIDER $G = (V, E, D)$ COMPLETE, AND THE NODE U FOR WHICH WE WANT THE TSP
- 2) COMPUTE $S = \text{MINIMUM SPANNING TREE OF } G \{ \otimes$
 $= M\text{ST}(G) \} \text{ WHERE } U \text{ IS THE ROOT FOR US.}$
- 3) COMPUTE THE TOUR T : PERFORM A ^{PREORDER} TRAVERSAL OF S : EACH TIME WE SEE A NODE IN S FOR THE FIRST TIME, WE ADD IT TO T !

MINIMUM SPANNING TREE

A MINIMUM SPANNING TREE IS A SUBSET OF THE EDGES OF A CONNECTED AND WEIGHTED UNDIRECTED GRAPH, THAT CONNECTS ALL THE VERTICES, WITHOUT AN CYCLE AND WITH THE MINIMAL POSSIBLE EDGE WEIGHT.



④ A MST IS A TREE IN ONLY THE SENSE OF LACK OF CYCLES AND BEING CONNECTED.

PICK A ROOT AND SET THE ORDER FOR THE SONS,
THEN A DFS WILL GIVE A PREORDER TRAVERSAL

EG: WE CHOOSE ① AS A ROOT AND THE ORDER IS THE NUMBERS:

1 → 2 → 3 → 4 → 1 IS THE TSP : T
NOT IN MST: HERE THERE IS THE NECESSITY OF T.I.
AND THIS IS THE POINT!

- 1) THE COST OF THE OPTIMAL SOLUTION IS GREATER THAN THE COST OF THE MST BY DEF, AS A TOUR WITHOUT AN EDGE IS SURELY A MST BY DEF. THE OPPOSITE IS NOT AUTOMATIC
- 2) THE COMPUTED TOUR IS AN APPROXIMATED SOLUTION:

$$\underbrace{\text{cost}(T)}_{\text{COST OF OUR APPROX. SOLUTION}} \leq \varepsilon \underbrace{\text{cost}(S)}_{\text{COST OF THE SPANNING TREE}}$$

T. i + EVERY EDGE IS SUMMED TWICE

1) + 2)

$$\text{COST}(T) \leq 2 \cdot \text{COST}(S) \leq 2 \cdot \text{OPT} \blacksquare$$

11) NP-COMPLETE PROBLEMS: KNAPSACK

THE KNAPSACK PROBLEM

- N ELEMENTS
- W IS THE KNAPSACK CAPACITY
- w_i IS THE OCCUPANCY OF THE i-TH ELEMENT
- v_i IS THE VALUE OF THE i-TH ELEMENT
- A SOLUTION $S \subseteq [n]$ IS FEASIBLE IF

$$\sum_{i \in S} w_i \leq W$$

NOTE:

$$\forall i \in N, w_i \leq W$$

WE WANT TO MAXIMIZE

$$\sum_{i \in S} v_i$$

THE GREEDY APPROACH

WE DEFINE $\frac{v_i}{w_i}$ AS THE VALUE PER UNIT.

WE DO THE FOLLOWING HYPOTHESIS, WITHOUT LOSS OF GENERALITY:

$$\left\{ \frac{v_1}{w_1} \geq \dots \geq \frac{v_n}{w_n} \right.$$

TAKE THE i -TH ELEMENT, WITH i AS SMALL AS POSSIBLE, AS SOON AS POSSIBLE

WE CAN FOLLOW THIS GREEDY APPROACH, WHICH GIVES US AN R -APPROXIMATION ($R = 2$):

- $\tilde{S} = \emptyset$
- $W' = W$, W' REPRESENT THE RESIDUAL CAPACITY OF THE KNOTPSACK
- $\forall i = 1, 2, \dots, N:$
 - IF $w_i \leq W'$ THEN:
 - $\tilde{S} = \tilde{S} \cup \{i\}$
 - $W' = W' - w_i$
- RETURN \tilde{S}

THE PREVIOUS ALGORITHM IS A GOOD APPROXIMATION
BUT IT HAS A GLITCH.

CONSIDER THE FOLLOWING SCENARIO



- a) $v_1 = v_2 = \dots = v_{N-1} = 1 \wedge v_N = W - 1$
 b) $w_1 = w_2 = \dots = w_{N-1} = 1 \wedge w_N = W$

$$a) + b) \Rightarrow \begin{cases} \bullet \forall i \in [1, N-1], \frac{v_i}{w_i} \\ \bullet \frac{v_N}{w_N} < 1 \end{cases}$$

AND WE THEN HAVE

- $\tilde{S} = \{1, \dots, N-1\}, \text{cost}(\tilde{S}) = \sum_{i \in S} v_i = N-1$
- THE OPTIMAL IS $S^* = \{N\}, \text{cost}(S^*) = W - 1$

THE APPROXIMATION RATIO IS

$$\frac{\text{OPT}}{\text{cost}(\tilde{S})} = \frac{W-1}{N-1} = k$$

→

K ARBITRARILY LARGE

LET'S FIX THE GLITCH:

WE DEFINE M_G AS THE VALUE RETURNED BY THE GREEDY APPROACH.

LET THEN $V_{\text{MAX}} = \max_{i \in [N]} v_i$

AS A SOLUTION TO THE KNAAPSACK WE RETURN THE MAXIMUM BETWEEN v_g AND v_{\max}

THE FIX GIVES US A $\frac{1}{2}$ -APPROXIMATION

PROOF:

LET'S DEFINE AN $\frac{\text{UPPER BOUND}}{\text{UB}}$ TO THE $\frac{\text{OPT. SOL.}}{\text{OPT}}$

$$\text{UB} \geq \text{OPT}$$

LET'S THEN DEFINE J AS THE SMALLEST ELEMENT IN $[n]$ THAT DOES NOT FIT IN W' ($J \geq 2$)

IN SYMBOLS:

$$\sum_{i=1}^{J-1} w_i \leq W \text{ BUT } \sum_{i=1}^J w_i > W$$

WE THEN USE THE FOLLOWING NOTATION:

$$\bullet \overline{w}_J = \sum_{i=1}^{J-1} w_i$$

$$\bullet \overline{v}_J = \sum_{i=1}^{J-1} v_i$$

THEN WE SAY THAT

$$UB = \sum_{i=1}^J v_i = \bar{V}_J + v_J > OPT$$

MIND THAT UB IS NOT A SOLUTION, THE OBJECT J CANT FIT THE KHAAPSACK BY DEF.

SINCE $OPT < \bar{V}_J + v_J$ WE HAVE TWO CASES

$$1) \bar{V}_J > v_J \Rightarrow OPT < \sum \bar{V}_j \leq \sum M_G$$

$$\begin{aligned} \bar{V}_J + v_J &> OPT + \bar{V}_J > v_J \\ \Rightarrow & \\ OPT &< \sum \bar{V}_j \end{aligned}$$

AS S IS THE GREEDY
SOLUTION IN $\{x_1, \dots, x_{T-1}\} \subseteq S$

$$2) \bar{V}_J \leq v_J \Rightarrow OPT < \sum v_j \leq \sum V_{MAX}$$

$$v_J \leq V_{MAX} \text{ BY DEF}$$

SUMMING UP:

$$OPT \leq \sum (\max(M_G, V_{MAX})) = 2 \cdot \text{cost}(\tilde{S})$$

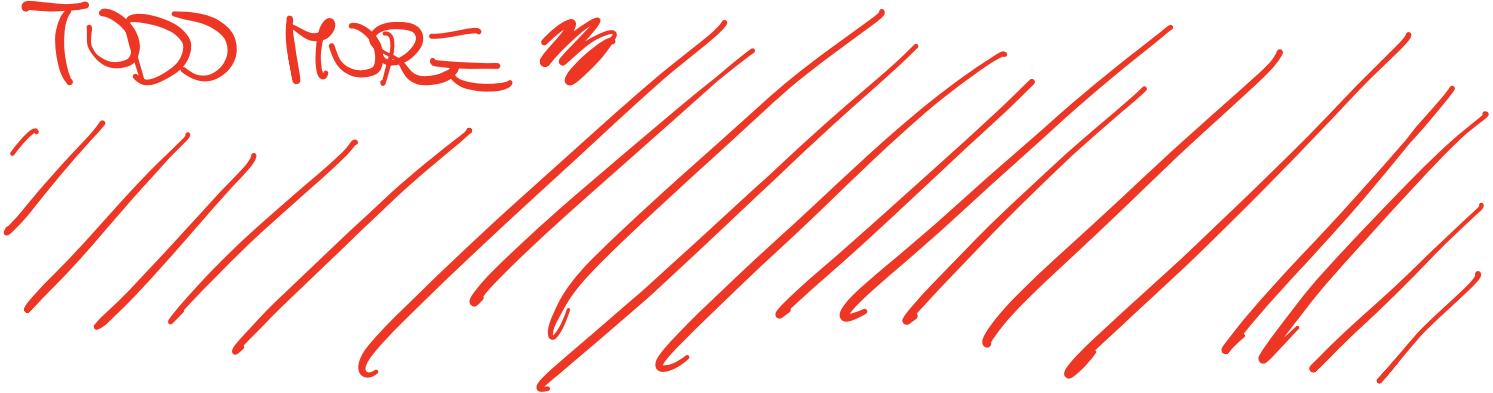
$$\Rightarrow$$

$$\frac{OPT}{\text{cost}(\tilde{S})} \leq 2 \quad \blacksquare$$

12) FPTAS: FULLY POLYNOMIAL APPROX. SOLUTION

IT'S A ALGORITHM DESIGN TECHNIQUE TO OBTAIN
AN R-APPROX FOR ANY $R > 1$

TODD MORE m



TODD FPTAS



KNAPSACK FPTAS: DP ε — CAN FA? SCAL $v_i \in$
UN GREEDY? BKH?

SOLVE AN INSTANCE OF THE KNAPSACK IN

| $O(n^2 \cdot v_{\max})$ TIME \otimes EXACTLY

IDEA: SCALE THE VALUES v_i TO $\tilde{v}_i = \left\lfloor \frac{v_i}{k} \right\rfloor$ FOR
SOME SUITABLE FACTOR k TO BE FIXED ASAP.

WE THEN HAVE TWO "INSTANCES" OF THE PROBLEM

- { a) ORIGINAL: v_i, w_i, W
b) SCALED: \tilde{v}_i, w_i, W

S IS FEASIBLE IN a) \Leftrightarrow S IS FEASIBLE IN b

LET'S FIX K SO THAT \tilde{v}_i ARE POLY(N) AS
WE WILL PAY $O(n^2 \tilde{v}_{\max})$

A GOOD CHOICE: \leftarrow PERCENT ?

$$\tilde{v}_i = \left\lfloor \frac{v_i}{k} \right\rfloor \leq \left\lfloor \frac{N}{\epsilon} \right\rfloor \Rightarrow \frac{v_i}{k} \leq \frac{N}{\epsilon} \Rightarrow k = \frac{v_{\max} \cdot \epsilon}{N}$$

- RUN DP ε ON b) IN TIME $O(n^2 \tilde{v}_{\max}) = O\left(\frac{n^3}{\epsilon}\right)$

OBSERVATION: WE FIND THE EXACT
OPTIMAL SOLUTION \tilde{s} FOR b)

PEPE?

How good is it for a)?

- ## • CLEARLY IT IS FEASIBLE

$$\text{WHAT ABOUT } \sum_{i \in S} v_i \geq (1-\epsilon) \sum_{i \in S^*} v_i$$

~~CLEARLY:~~

$$\left\{ \sum_{i \in S} v_i^2 \geq \sum_{i \in S^*} v_i^2 \right\}$$

PERCHÉ?

WE THEN HAVE 3 FACTS:

1) $\forall i. \quad v_i \geq r_i^2$

$$2) k \cdot v_i \geq v_i - k \quad \{ \text{wie vor Di?}$$

$$3) \sum_{i \in S} \tilde{v}_i \geq \sum_{i \in S^*} \tilde{v}_i \quad \left\{ \begin{array}{l} \text{perché?} \\ \text{Cosa vale } \tilde{v}_i? \end{array} \right.$$

THEN: IN THE CONTEXT OF THE VERSION a) OF
THE PROBLEM

$$\text{COST}(\tilde{S}) =$$

$$= \sum_{i \in \tilde{S}} v_i$$

$$1) - \geq \sum_{i \in \tilde{S}} k \tilde{v}_i$$

$$3) - \geq \sum_{i \in S^*} k \tilde{v}_i$$

$$2) - \geq \sum_{i \in S^*} (v_i - k)$$

⑧

$$|S^*| \leq N$$

$$= \sum_{i \in S^*} v_i - \sum_{i \in S^*} k$$

$$\geq \text{OPT} - NK$$

$$\lceil NK = \left(\frac{\epsilon \cdot v_{\max}}{N} \right) N = \epsilon \cdot v_{\max} \leq \epsilon \cdot \text{OPT}$$

$$\geq \text{OPT} - \epsilon \text{OPT}$$

$$= (1 - \epsilon) \text{OPT} \blacksquare$$