

# Approximate Counting

Decision problem in NP:

$\Pi$  is in NP if for any YES instance  $I$ ,  $\exists$  a proof that  $I$  is a YES instance that can be verified in poly time.

Ex: "Is there a Hamiltonian cycle in the given graph?"

"Is there a tour of length  $\leq k$  for a travelling Salesperson instance".

"Is the given Boolean formula satisfiable".

Corresponding to these decision problems, we can define counting problems.

Counting problem for a decision problem  $\Pi$  in NP:

Given an instance  $I$  of  $\Pi$ , produce as output a non-negative integer that is the no. of solutions (or proofs) for the instance  $I$ .

$$\#P = \left\{ F_n : \{0,1\}^{n^c} \rightarrow \mathbb{Z}_{\geq 0} \mid \begin{array}{l} F_n \text{ counts the no. of accepting} \\ \text{paths in a NTM, and } n > 0 \end{array} \right\}$$

- Ex: "No. of Hamiltonian cycles in the given graph"  
 "No of satisfying assignments to a Boolean formula"  
 "No of k-cliques in a given graph".

These are at least as hard as their decision versions:

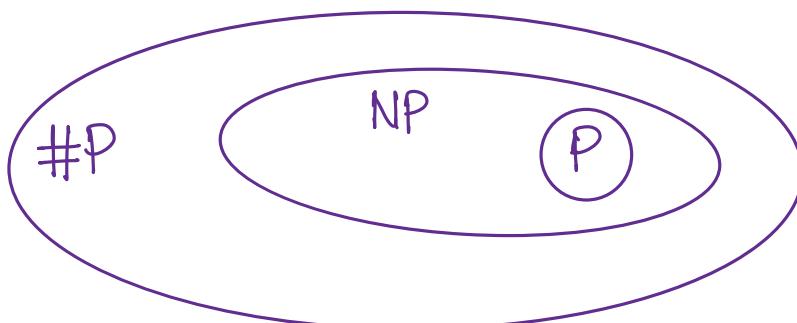


Figure: Hierarchy  
of functional  
Complexity classes.

Qn: Can we approximately count them?

- We do not know of det. approx. schemes
- what about randomized approx schemes?

Defn: A polynomial approximation scheme (PAS) for a counting problem  $P$  is a deterministic algo that takes as input an instance  $I$  and  $\epsilon \in \mathbb{R}_{>0}$ , and in time polynomial in  $n = |I|$  produces an output  $A(I)$  such that

$$(1 - \epsilon) \#(I) \leq A(I) \leq (1 + \epsilon) \#(I).$$

A fully polynomial approximation scheme (FPAS) is a PAS that runs in time  $\text{poly}(n, 1/\epsilon)$ .

Defn: A polynomial randomized approximation scheme (PRAS) for a counting problem  $P$  is a randomized algorithm  $A$  that takes as input an instance  $I$  and a real no.  $\epsilon > 0$ , an in time  $\text{poly}(n)$  produces  $A(I)$  s.t

$$\Pr \left[ (1-\epsilon) \#(I) \leq A(I) \leq (1+\epsilon) \#(I) \right] \geq \frac{3}{4}.$$

Fully polynomial randomized approximation scheme (FPRAS) is a PRAS that runs in time  $\text{poly}(n, 1/\epsilon)$ .

With prob  $< \frac{1}{4}$ , we assume nothing about the dist. of  $A(I)$  from  $\#(I)$ .

Defn: An  $(\epsilon, \delta)$ -FPRAS for a counting problem is a FPRAS that takes as input an instance  $I$  and computes an  $\epsilon$ -approximation to  $\#(I)$  w.p  $\geq 1-\delta$  in time  $\text{poly}(n, \frac{1}{\epsilon}, \frac{1}{\delta})$ .

Abstract problem:

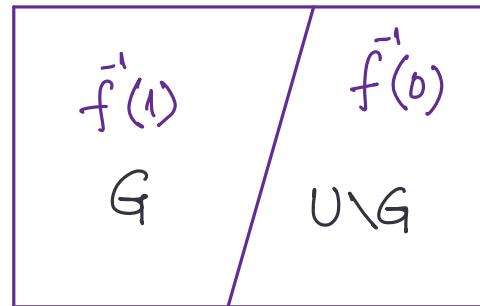
Let  $U$  be a finite set of known size. Let  $f: U \rightarrow \{0, 1\}$  be a Boolean function over  $U$ . Let  $G = \{u \in U \mid f(u) = 1\}$ .

Assumptions:  $u$  can be sampled uniformly and

randomly from  $U$ . Given  $u \in U$ ,  $f(u)$  can be computed "quickly".

Question: Estimate  $|G|$ .

Attempt 1:



Sample  $u_1, \dots, u_N$  from  $U$  independently from  $U$ .

For all  $i \in [N]$ , let

$$Y_i = \begin{cases} 1 & \text{if } f(u_i) = 1, \\ 0 & \text{otherwise.} \end{cases} \quad Y_i = 1 \text{ iff } u_i \in G.$$

Let  $Z$  be a r.v s.t

$$Z = \frac{|U|}{N} \cdot \sum_{i=1}^N Y_i.$$

We call it estimator random variable.

Claim:  $E[Z] = |G|$

$$\begin{aligned} E[Z] &= E\left[\frac{|U|}{N} \cdot \sum_{i=1}^N Y_i\right] = \frac{|U|}{N} \cdot \sum_{i=1}^N E[Y_i] \cdot \frac{|U|}{N} \\ &= \frac{|U|}{N} \sum_{i=1}^N \Pr[Y_i = 1] = \frac{|G|}{N} = \frac{|G|}{|U|} \end{aligned}$$

[Estimator theorem]

Theorem: Let  $\rho = \frac{|G|}{|U|}$ . Then the Monte Carlo method yields an  $\epsilon$ -approximation to  $|G|$  w. prob  $> 1 - \delta$  provided

$$N \geq \frac{4}{\epsilon^2 \rho} \cdot \ln \frac{2}{\delta}.$$

Proof: Note that  $\Pr[Y_i = 1] = \frac{|G|}{|U|}$ . Let  $Y = \sum_{i=1}^N Y_i$ .

$$\begin{aligned}
 & \Pr[(1-\varepsilon)|G| \leq Z \leq (1+\varepsilon)|G|] \\
 &= \Pr[(1-\varepsilon)Np \leq Y \leq (1+\varepsilon)Np] \quad \text{where } p = \frac{|G|}{|U|} \\
 &\geq 1 - \Pr[Y > (1+\varepsilon)Np] - \Pr[Y < (1-\varepsilon)Np] \\
 &= 1 - 2 \cdot e^{-Np\varepsilon^2/4} \quad \Rightarrow \quad \frac{\delta}{2} \geq e^{-\frac{Np\varepsilon^2}{4}} \\
 &\geq \frac{1-\delta}{1-\delta} \quad \Rightarrow \quad \frac{Np\varepsilon^2}{4} \geq \ln \frac{2}{\delta} \quad \text{N = poly}\left(\frac{1}{\varepsilon}, \ln \frac{1}{\delta}\right) \\
 &\quad \text{and } \frac{1}{p} \text{ is exponential.} \\
 &\quad \Rightarrow N \geq \frac{4}{p\varepsilon^2} \ln \frac{2}{\delta}.
 \end{aligned}$$

Question: But,  $N$  is dependent on  $p = \frac{|G|}{|U|}$ .  
 Problem is not solved yet.

$\hookrightarrow N$  could be exponential.

Ex:  $F$  is a Boolean formula in Disjunctive Normal Form. That is OR of ANDs.

$$F = \underbrace{T_1 \vee T_2 \vee \dots \vee T_m}_{\text{OR of terms}} \quad \text{and} \quad T_i = \underbrace{L_1 \wedge L_2 \wedge \dots \wedge L_k}_{\text{AND of literals}}$$

$$U = \{T, F\}^n \quad = \text{conjunctive clauses} \quad \text{conjunction} = \wedge \quad \text{disjunction} = \vee$$

$\#F :=$  No. of satisfying solutions of  $F$ .

$$G = \{u \in U \mid F(u) = T\} \quad \frac{|G|}{|U|} \text{ could be inverse exp.}$$

# Fix: Coverage argument // Biased sampling.

Space of assignments:  
 $H_1, \dots, H_m \subseteq V$

$H_i$  = subset of assignments that satisfy term  $T_i$ .

$\left| \bigcup_{i=1}^m H_i \right|$  is exactly the count we are interested in  
 $\rightsquigarrow H$ .

$$|H_i| = 2^{n-r_i} \quad x_1 \wedge x_2 \dots \wedge x_{r_i} \quad T_i = L_{i,1} \wedge \dots \wedge L_{i,r_i}$$

$\uparrow \quad \uparrow$

$T_i$  has  $r_i$  literals.

$U = H_1 \uplus H_2 \uplus \dots \uplus H_m \quad \{ \text{multiset union.} \}$

$$|U| = \sum_{i=1}^m |H_i| \geq |H| \text{ where } H = \bigcup_{i=1}^m H_i$$

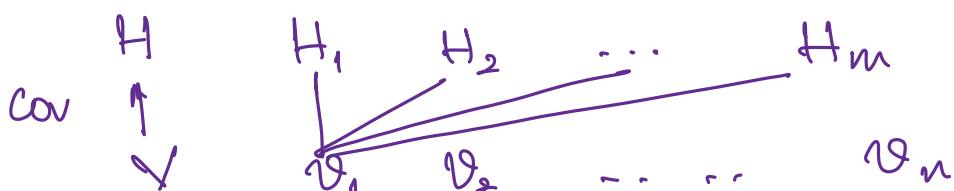
$\forall \quad T_1, T_2, T_3, \dots$

$\downarrow \quad \downarrow \quad \downarrow$

$(v,1) \quad (v,2) \quad (v,3)$

Coverage of an assignment  $v$ :

$$\text{cov}(v) = \underbrace{\{ (v,i) \mid (v,i) \in U \}}_{\text{ }} \quad || \quad \underbrace{\{ i \mid (v,i) \in U \}}$$



$$-\lvert \text{cov}(\vartheta) \rvert \leq m. \text{ and } U = \bigcup_{\vartheta \in H} \text{cov}(\vartheta) \\ |U| = \sum_{\vartheta \in H} |\text{cov}(\vartheta)|.$$

Defn:  $f: U \rightarrow \{0,1\}$  be defined as follows:

$$f((\vartheta, i)) = \begin{cases} 1 & \text{if } i = \min \{j \mid (\vartheta, j) \in U\} \\ 0 & \text{otherwise} \end{cases}$$

and  $G = \{(\vartheta, j) \mid f((\vartheta, j)) = 1\}$ .

Obs:  $|G|$  is exactly  $|H|$

It is easy to compute  $H_i$  and its cardinality.

Lemma:  $\rho = \frac{|G|}{|U|} \geq \frac{1}{m}$ .

{Similarly,  $|H|$  can be computed easily.

Proof:  $|U| = \sum_{\vartheta \in H} |\text{cov}(\vartheta)|$

$$\leq \sum_{\vartheta \in H} m$$

$$= m|H|.$$

$$= m \cdot |G|. \Rightarrow \frac{|G|}{|U|} \geq \frac{1}{m}.$$

Not done is to tell you how to sample from  $U$ .

## [Estimator theorem]

Theorem: Let  $p \geq 1/m$ . Then the Monte Carlo method yields an  $\epsilon$ -approximation to  $|G|$  w. prob  $\geq 1 - \delta$  provided

$$N \geq \frac{4m}{\epsilon^2} \cdot \ln \frac{2}{\delta}.$$

$$\begin{array}{cc} T_1 & T_2 \\ \downarrow & \downarrow \\ (\underline{\vartheta}, 1) & (\underline{\vartheta}, 2) \end{array}$$

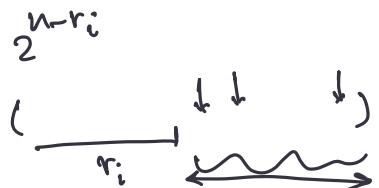
Recall that  $u_1, \dots, u_N$  were sampled uniformly and ind.  
 $\downarrow$   
 $(\underline{\vartheta}_1, i_1), (\underline{\vartheta}_2, i_2), \dots, (\underline{\vartheta}_N, i_N)$ .

For  $(\underline{\vartheta}, i)$  we define a way to sample it uniformly.

→ First sample  $i \in [m]$  w.p.  $\frac{|H_i|}{|U|} = \frac{|H_i|}{\sum_{i=1}^m |H_i|}$  ↗ no. of assignments satisfying term  $i$ .

$$U = H_1 \uplus H_2 \uplus \dots \uplus H_m.$$

$$\phi_i = \Pr[\text{Sampling } i] = \frac{|H_i|}{|U|} \text{ and } \sum_{i=1}^m \phi_i = 1$$

→ Sample  $\cong$  from  $H_i$  uniformly ↗ 

$$\Pr[(\underline{\vartheta}, i) \text{ is sampled}] = \Pr[\underline{\vartheta} \text{ is sampled from } H_i] \Big|_{\substack{i \text{ is sampled} \\ \text{from } [m]}}$$

$$\cdot \Pr[i \text{ is sampled}].$$

$$= \frac{1}{|H_i|} \times \frac{|H_i|}{|U|} = \frac{1}{|U|}.$$

$$N = \text{poly}(m, \frac{1}{\epsilon}, \frac{1}{\delta}).$$

$$Z' = \frac{|U'|}{N} \cdot \sum X_i$$

↙

$$\mathbb{E}[Z'] = |G|.$$

$$T_1 \vee T_2 \vee \dots \vee T_m . \quad |H_i| = 2^{n-r_i}$$

→ Define a distribution

$$i, \quad \Pr_r[i] = \frac{2^{n-r_i}}{\sum_{i=1}^m 2^{n-r_i}} . \quad \checkmark$$

→ Sample  $\underset{i=1}{\overset{m}{\in}} H_i$

$$N \geq \frac{4m}{\epsilon^2} \ln \frac{2}{\delta}$$

$$\underset{1-\delta}{=}$$

$$\epsilon.$$

