

Today:

- overview of other types of distribution testing (see previous notes)

- new theme: "Querying Big Data Sets"

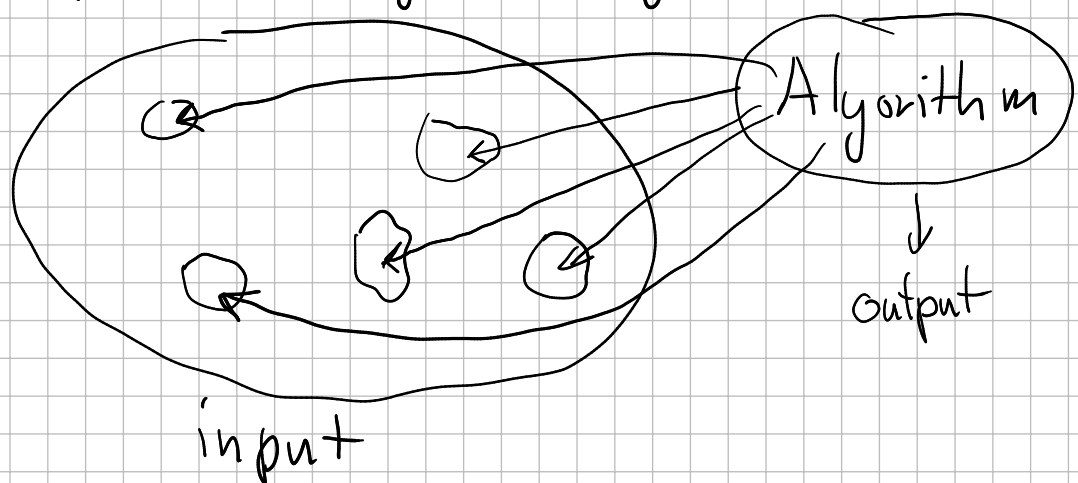
example: estimating maximum matching or minimum vertex cover size

New theme: Querying Big Data Sets

- input: big data set

- don't want to read all of it

- what can we say about its properties after looking at tiny fraction?



Goal: sublinear-time algorithms

Note: results very dependent on type of queries/samples allowed

---

Simplest example: estimate fraction of elements with specific property, e.g., fraction of numbers that are prime

Input: sequence of numbers  $s_1, s_2, \dots, s_n$

Queries: for any  $i \in [n]$ , can get  $s_i$

How many queries needed to estimate fraction of prime numbers up to  $\pm \epsilon$ ?

Solution: - sample  $O(1/\epsilon^2)$  numbers

- return the fraction of primes in the sample

---

Now: graphs  $G = (V, E)$

Allowed queries:

- get uniformly random  $v \in V$

- for any  $v \in V$ , get degree of  $v$

- for any  $v \in V$  &  $i \in \mathbb{Z}_+$ , get  $i$ -th neighbor of  $v$

18-2

# Maximum matching & minimum vertex cover

Matching = subset of edges that share  
no endpoints

Maximum matching = highest cardinality matching

Maximal matching = no edge can be added  
or it won't be a matching

$$\textcircled{*} \stackrel{=}{=} \left( \begin{array}{l} \# \text{ edges in maximal matching} \\ \geq \frac{1}{2} \left( \# \text{ edges in maximum matching} \right) \end{array} \right)$$

Vertex cover = subset of vertices such that at  
least one endpoint of each edge  
in the set

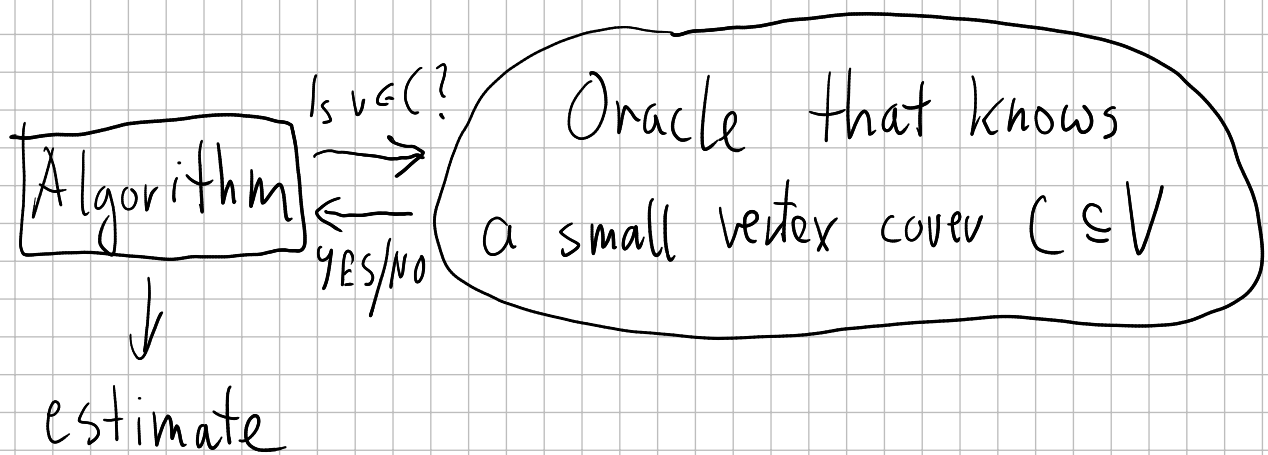
Minimum vertex cover = minimum cardinality  
vertex cover

$$\textcircled{*} \leq \underbrace{\left( \text{size of minimum vertex cover} \right)}_{\leftarrow} \leq 2 \textcircled{*}$$

Our main goal: estimate this

Ideal situation: "Our lucky day"

Someone gave us an oracle that provides access to a small vertex cover



$O(1/\epsilon^2)$  queries about randomly selected vertices suffice to estimate  $|C|$  up to  $\pm \epsilon n$  with probability 99/100

Bad news: today is not our lucky day  
and Tuesday won't be either

We'll have to construct our oracle ourselves