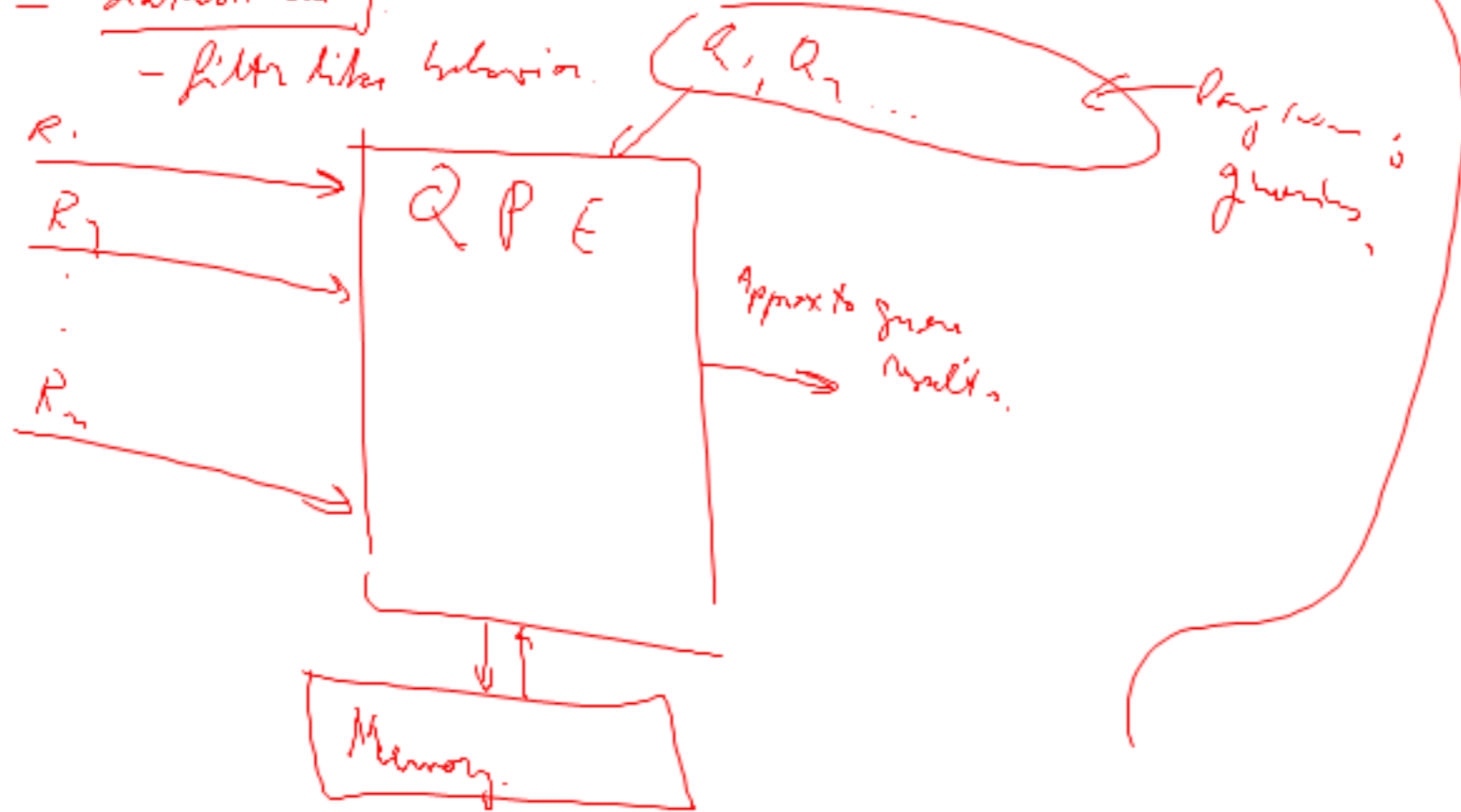


Lecture 14

Sketches

- look at all the data
- use small amount of memory.
- datastreaming.
- filter like behavior.



- no assumptions
- not making predictions.

Given D want to compute

$f(D)$

could require a lot of space.

Compute $\tilde{f}(D)$

$\tilde{f}(D)$ approx $f(D)$

$\tilde{f}(D)$ requires little memory

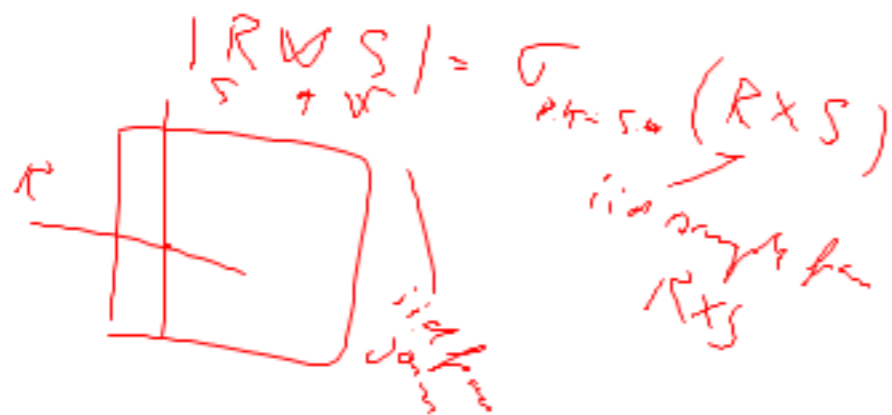
History of Sampling

- survey sampling.

Pearson, Fisher 1930

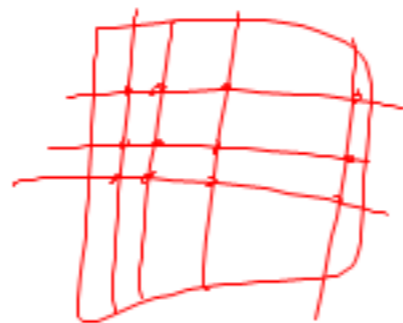
- in DB

- Frankel 1981
 selectivity estimation
 joins with joins



Before (Ripped Join)

Get R' , S' samples
 consider $|R' \bowtie S'| \dots$



$|R'| \cdot |S'|$

↓

w , anchored join

$V_q \sim \frac{1}{|R'| \cdot |S'|}$



$V_q \sim \frac{1}{|R'|}$

sample $R \times S$

easy to analyze

$|R'| = |S'|$

Linting of Sketches

Fajot - Martin algo for

counting distinct elements
- space $O(\# \text{ of distinct elements})$

- can you count in less space?

- need approx

Let $h(x)$ be a min. hash of x .

$h_S(x)$

• for every item x on the stream

- compute $h(x)$

- compute the ~~sum~~ position of the first 1 in the binary rep. of $h(x)$

1 0 1 1 0 1 1 0 1 0

↑
maintain map such that

Claim:

$\frac{1}{m}$

2

approx to number of distinct values