

Algorithms for Big Data – no. 22934

Problems Set 11

General instructions: Please keep your answers short and easy to read. You can use algorithms, results, calculations or notation that appear in our course material without repeating them, unless asked explicitly to redo them.

1. Show that given any algorithm that runs in time $T(n)$ on inputs of size n with probability of error ϵ , one can convert it into a new algorithm that runs in time $O(T(n) \log(1/\epsilon))$ with probability of error at most ϵ . Hint: run $\log(1/\epsilon)$ times and take the majority answer. Use Chernoff bounds.
2. You are given an approximation scheme for SUM such that $\text{SUM}(S)$ is approximated by $\text{SUM}'(S)$, and SUM' runs in time polynomial in n and $1/\epsilon$. Construct an approximation scheme for COUNT such that $\text{COUNT}(S)$ is approximated by $\text{COUNT}'(S)$, and COUNT' runs in time polynomial in n and $1/\epsilon$.
and $\text{COUNT}'(S)$.
3. (Coupon Collector Problem). Given a die with s sides. What is the expected number of times you need to roll the die in order to see each of the s sides? Hint: Given that you saw k sides, how many times do you need to roll the die to see the $(k+1)$ side? Then use linearity of expectation.
4. Recall that n stands for the length of the input stream (i.e. the number of the tokens) and k is the number of different possible tokens. Assume that the value of every token is between 1 and k , i.e. $1 \leq \text{token} \leq k$. Design a streaming algorithm that at any point (not known in advance) receives a query Q (i.e. a subset of tokens values) and outputs and estimate what fraction of tokens in the stream belong to Q within additive error ϵ . Note that k is given only at query time (not in advance). Hint: Maintain $1/\epsilon$ random samples and use them to estimate the fraction in Q .
5. a. Suppose we are guaranteed that some token in the stream appears more than half the time, i.e., there exists (unknown) token t with frequency $> n/2$. Design a streaming algorithm with space complexity $O(\log n)$ bits that finds this item t . Hint: Store only two items.

b. Next, extend your algorithm to output also a $(1 \pm \epsilon)$ approximation to its frequency f_t . Make sure to clearly state the space complexity of your algorithms. Hint: Use ideas of question 4.