## **Homework 8 Solutions**

Note: These solutions are not necessarily model answers. Rather, they are designed to be tutorial in nature, and sometimes contain a little more explanation than an ideal solution. Also, bear in mind that there may be more than one correct solution. The maximum total number of points available is 35.

- 1. (a) The random variables in the family are indexed by  $x \in \mathbb{Z}_p$ , so there are p of them.
  - (b) As stated in the hint, our goal is to prove that Pr[f<sub>a</sub>(x) = y] = <sup>1</sup>/<sub>p</sub> for all x, y ∈ Z<sub>p</sub>, for a chosen u.a.r. 4pts This event is equivalent to ∑<sup>d-1</sup><sub>i=0</sub> a<sub>i</sub>x<sup>i</sup> = y. Using the principle of deferred decisions, we can first choose the coefficients a<sub>1</sub>,..., a<sub>d-1</sub>, at which point the event becomes a<sub>0</sub> = y - ∑<sup>d-1</sup><sub>i=1</sub> a<sub>i</sub>x<sup>i</sup>. Thus there is a unique value of a<sub>0</sub> ∈ Z<sub>p</sub> that makes the event true, so since a<sub>0</sub> is chosen u.a.r. from Z<sub>p</sub>, independently of the other a<sub>i</sub>, the probability of the event is exactly <sup>1</sup>/<sub>p</sub>, as required.

1pt

- (c) Again following the hint, our goal is to prove that  $\Pr[(f_a(x_1) = y_1) \cap (f_a(x_2) = y_2) \cap \ldots \cap (f_a(x_d) = 4pts y_d)] = \frac{1}{p^d}$  for any distinct  $x_1, \ldots, x_d \in \mathbb{Z}_p$  and any  $y_1 \ldots, y_d \in \mathbb{Z}_p$ , where again the probability is over the random choice of a. But this event defines d points on the degree-(d-1) polynomial  $f_a$ , so the polynomial, and hence its coefficients  $a_0, a_1, \ldots, a_{d-1}$ , are uniquely determined. Since each such polynomial is picked with uniform probability from the set of  $p^d$  polynomials, the probability of this event is exactly  $\frac{1}{p^d}$ , as required.
- (a) Let C be any set of k vertices in K<sub>n</sub>. The probability that C is monochromatic is equal to two times 2pts Pr[all edges in C are red]. But this probability in turn is equal to 2<sup>-(k/2)</sup>, since the edge colors are d-wise independent. Hence the probability that any C is monochromatic is (by a union bound) at most (<sup>n</sup><sub>k</sub>)2<sup>-(k/2)+1</sup>, which is less than 1 when n ≤ 2<sup>k/2</sup>, as we saw in class using a simple computation.
  - (b) We need at least m d-wise independent random variables (one for each edge of  $K_n$ ). Using the con-*lpt* struction of part (a), this entails that  $p \ge m$ .
  - (c) Using the parameters  $m = \binom{n}{2}$  and  $d = \binom{k}{2}$ , we see that the size of the sample space is  $p^d \le n^{k^2}$ . 4pts (Each point in the sample space is a choice of d coefficients from  $\mathbb{Z}_p$ .) Since this is polynomial in n for any fixed k, we can therefore deterministically cycle through every sample point a: for each such point, we then use the construction of the previous question to generate m d-wise independent colors for the edges, and for each such 2-coloring we check each of the  $\binom{n}{k} = O(n^k)$  cliques in  $K_n$  to see if it is monochromatic. We stop when we find a 2-coloring that is k-good (which must happen, since by part (a) some such coloring in our sample space exists and we are trying all of them). The resulting algorithm is deterministic and requires time  $O(n^{k^2+k})$ .
  - (d) We can allocate one group of  $\binom{n}{k}$  processors to each of the  $O(n^{k^2})$  sample points. We first use  $p \leq 4pts$   $2m \leq n^2$  of these processors to evaluate, in parallel, the polynomial  $f_a(x)$  at m points x, for this particular sample point a, thus generating the colors for the edges of  $K_n$ . Each of these evaluations of a degree-(d-1) polynomial can be done in constant time (since  $d = \binom{k}{2}$  is a constant<sup>1</sup>). Then, each of the  $\binom{n}{k}$  processors checks one of the k-cliques of  $K_n$  to see if it is monochromatic: this takes constant time since k is constant. All processors finding a monochromatic clique report back to one single "lead" processor in the group: if none have reported, then the lead processor announces to the world that a k-good coloring has been found. The total (parallel) time required is  $O(\log n)$  since the final reporting can be done in logarithmic time.

Most students forgot to parallelize the generation of the edge colors for each set of coefficients  $a_0, a_1, \ldots, a_{d-1}$ . Without parallelizing this step, the generation of the edge colors will take  $O(n^2)$  time for each set of coefficients.

<sup>&</sup>lt;sup>1</sup>Technically we should charge  $O(\log n)$  time for this since we are working with integers with  $O(\log n)$  bits. But this does not affect the overall time complexity.

- 3. (a) The running time is O(n). For each of the two sets  $S_1, S_2$ , we run *n* hashing operations and increment 2*pts* counters *n* times. Finally, we need to make *n* comparisons.
  - (b) If  $S_1$  and  $S_2$  are identical, then clearly the *i*th counters for both tables will match for all *i*, regardless *1pt* of the choice of hash function. (This follows from the definition of a hash function: it always hashes the same element to the same position.)
  - (c) Suppose  $S_1$  and  $S_2$  are not identical, and furthermore  $S_1, S_2$  are disjoint. Fix some  $x \in S_1$ , and 4pts note from our assumptions that  $x \notin S_2$ . For each  $y \in S_2$ ,  $y \neq x$ , we have by the property of universal hash functions that  $\Pr_h[h(x) = h(y)] \leq 1/cn$ . Taking a union bound over all  $y \in S_2$ , we have  $\Pr_h[\exists y \in S_2 : h(x) = h(y)] \leq 1/c$ . Hence, with probability 1 1/c over h, we have that  $h(y) \neq h(x)$  for all  $y \in S_2$ . In that case, the h(x)th counter for  $S_2$  is 0, whereas that for  $S_1$  is at least 1, so the algorithm outputs "no".
- 4. (a) It is clear that the new algorithm still outputs "no" with probability 1 on input  $x \notin L$ . For input 4pts  $x \in L$ , let  $Y = \sum_{i=1}^{t} Y_i$ , where  $Y_i$  is the indicator r.v. for the event that  $\mathcal{A}(x, r_i)$  outputs "yes". Then,  $E[Y_i] \leq 1/2$  and  $Var[Y_i] \leq 1/4$ . Hence,  $E[Y] \leq s/2$  and (since the  $Y_i$  are pairwise independent)  $Var[Y] = \sum_{i=1}^{t} Var[Y_i] \leq s/4$ . The probability that the new algorithm outputs "no" is given by Pr[Y = 0]. Applying Chebyshev's inequality, we have

$$\Pr[Y = 0] \le \Pr[|Y - E[Y]| \ge s/2] \le \frac{\operatorname{Var}[Y]}{(s/2)^2} = 1/s.$$

This yields the required bound on the error probability.

(b) To implement part (a), we need to generate s = 1/δ ≤ 2<sup>t</sup> pairwise independent uniform random 2pts strings in {0,1}<sup>t</sup>. We can do this using the construction discussed in class and in [MU, Lemma 15.2]. Namely, let 2<sup>t</sup> < q < 2<sup>t+1</sup> be a prime and use the construction to generate q pairwise independent samples from Z<sub>q</sub>, using only two independent random integers in Z<sub>q</sub>. This requires only O(t) random bits, as required.

Most students forgot to explain why  $s \leq 2^t$ . Note that the constructions we know (e.g., from lecture and Problem 1) for generating pairwise independent bit strings of length t using O(t) mutually independent bits can only generate up to  $2^t$  strings.

(c) We need to run  $\mathcal{A} s = 1/\delta$  times in the new scheme, versus  $O(\log(\delta^{-1}))$  times in the standard 2pts approach. Hence the running times are  $O(t/\delta)$  and  $O(t \log \delta^{-1})$  respectively. Note, therefore, that the new scheme uses significantly fewer random bits, but at the expense of a significantly increased running time.