



## Min-Hashing.

In **min-hashing** we created an estimator with “one-sided error”: our estimate was always an *overestimate*. I.e., for the target value  $v$  we created a random variable  $X$  such that  $\Pr[X \geq v] = 1$ . Suppose  $\mathbf{E}[X] = \mu$ .

1. Show that  $\Pr[X \geq 2\mu] \leq \frac{1}{2}$ . (“The probability of one estimate being too large is at most 50%.”)
2. Use this to show that if we take  $k$  independent copies  $X_1, X_2, \dots, X_k$  of the r.v.  $X$ , then  $\Pr[\min_{i=1}^k(X_i) \geq 2\mu] \leq 2^{-k}$ .
3. Show that  $k = \lg(1/\delta)$  gives  $2^{-k} = \delta$ . (If we want error probability  $2^{-100}$ , take the minimum of 100 independent estimates.)

## Fingerprinting.

**Many Patterns:** You are given a set of patterns  $P_1, P_2, \dots, P_k$  of equal length (all of them having length  $n$ ) and a text  $T$  of length  $m$ . Give an algorithm to find all the locations  $i$  such that some pattern  $P_j$  occurs as a substring of  $T$  starting at location  $i$ . The expected runtime should be  $O(kn + m)$ , and the probability of error is at most 0.01.<sup>1</sup>

---

<sup>1</sup>Assume you can do arithmetic operations on numbers of size  $O(\log(kmn))$  in constant time, even modulo a prime.