

A Variance of LogEstimator

We now bound the variance of our estimator by $O(\log^2 k)$. Recall that the output of `LogEstimator` is given by $\log(\mathbf{X}/t) - g(\mathbf{B}_1, \dots, \mathbf{B}_r)$, where the function g is bounded. Since the variance we seek is $O(\log^2 k)$, it suffices to show that the variance of $\log(\mathbf{X}/t)$ is $O(\log^2 k)$ with $i \sim \mathcal{D}$, since subtracting g changes the estimate by at most a constant (see Lemma 2.3).

Lemma A.1. *Let $i \sim \mathcal{D}$ and \mathbf{X} denote the number of independent trials from $\text{Ber}(p_i)$ before we see t successes. Then, $\text{Var}[\log(\mathbf{X}/t)] = O(\log^2 k)$.*

Proof. Let $X_{\max} = 2kt$, and consider the random variable $\mathbf{X}' = \min\{\mathbf{X}, X_{\max}\}$. Then

$$\begin{aligned} \text{Var}[\log(\mathbf{X}/t)] &\leq \mathbf{E}\left[\left(\log(\mathbf{X}/t) - \log(\mathbf{X}'/t) + \log(\mathbf{X}'/t)\right)^2\right] \\ &\leq 2 \cdot \mathbf{E}\left[\left(\log(\mathbf{X}/t) - \log(\mathbf{X}'/t)\right)^2\right] + 2 \cdot \mathbf{E}\left[\log^2(\mathbf{X}'/t)\right] \\ &\leq 2 \cdot \mathbf{E}\left[\log^2(\mathbf{X}/\mathbf{X}')\right] + 2\log^2(2k) \\ &\leq \frac{4}{\ln^2(2)} \cdot \mathbf{E}\left[\left(\sqrt{\frac{\mathbf{X}}{\mathbf{X}'}} - 1\right)^2\right] + 2\log^2(2k), \end{aligned}$$

where we used that $\log(\mathbf{X}'/t) \leq \log(2k)$ always, and that $\log(z) \leq \sqrt{z-1}/\ln(2)$ for all $z \geq 1$. Then,

$$\mathbf{E}\left[\frac{\mathbf{X}}{\mathbf{X}'} - 1\right] \leq \mathbf{E}\left[\frac{\mathbf{X}}{X_{\max}}\right] = \frac{1}{X_{\max}} \sum_{i=1}^k p_i \cdot t = \frac{tk}{X_{\max}} = 2.$$

□

B Omitted Details from Section 2

Proof of Claim 2.5. Notice that \mathbf{X} is the number of trials from $\text{Ber}(p_i)$ until we see t successes. We now have the following string of equalities:

$$\begin{aligned} \mathbf{E}_{\mathbf{X}, \mathbf{B}_1, \dots, \mathbf{B}_r} \left[\eta - \log\left(\frac{1}{p_i}\right) \right] &= \mathbf{E}_{\mathbf{X}}[\log \mathbf{Y}] - \mathbf{E}_{\mathbf{B}_1, \dots, \mathbf{B}_r}[g(\mathbf{B}_1, \mathbf{B}_2, \dots, \mathbf{B}_r)] \\ &= \mathbf{E}_{\mathbf{X}}[f(\mathbf{Y}) + h(\mathbf{Y})] - g(p_i, p_i^2, \dots, p_i^r) = \mathbf{E}_{\mathbf{X}}[h(\mathbf{Y})], \end{aligned}$$

where we used the fact that g is a linear function, and that $\mathbf{E}[\mathbf{B}_\ell] = p_i^\ell$ in order to substitute

$$\mathbf{E}_{\mathbf{B}_1, \dots, \mathbf{B}_r}[g(\mathbf{B}_1, \dots, \mathbf{B}_r)] = g(p_i, p_i^2, \dots, p_i^r).$$

Furthermore, we divide $\log \mathbf{Y} = f(\mathbf{Y}) + h(\mathbf{Y})$, where $f(z)$ is the degree- r Taylor expansion of $\log z$ at 1, and $h(z) = \log z - f(z)$ is the error in the degree- r Taylor expansion of $\log(z)$, i.e.,

$$h(z) = \log(z) - f(z).$$

Finally, by construction of g , $\mathbf{E}[f(\mathbf{Y})] = g(p_i, p_i^2, \dots, p_i^r)$, which gives the desired equality. □

Verifying \mathbf{Y} is subgamma. Recall that \mathbf{X} is the number of independent draws from a $\text{Ber}(p)$ distribution until we see t successes. In other words, we may express $\mathbf{X} = \mathbf{X}_1 + \dots + \mathbf{X}_t$, where \mathbf{X}_i is the number of draws of $\text{Ber}(p)$ before we get a single success. Then, we always satisfy

$$\mathbf{E}[\mathbf{X}_i] = \frac{1}{p} \quad \Pr[\mathbf{X}_i > \ell] = (1-p)^{\lceil \ell \rceil} < e^{-p\ell}.$$

This, in turn, implies that for any $r \geq 1$

$$(\mathbf{E}[|\mathbf{X}_i - 1/p|^r])^{1/r} \leq (\mathbf{E}_{\mathbf{X}_i, \mathbf{X}'_i}[|\mathbf{X}_i - \mathbf{X}'_i|^r])^{1/r} \leq 2(\mathbf{E}[|\mathbf{X}_i|^r])^{1/r} = O(r/p),$$

where the first line is by Jensen's inequality, and the second is by the triangle inequality and Hölder inequality. Finally, we use the tail bound on \mathbf{X}_i to upper bound the expectation of $|\mathbf{X}_i|^r$. Then, we have

$$\begin{aligned} \mathbf{E}\left[e^{\lambda(\mathbf{X}_i - 1/p)}\right] &= 1 + \lambda \mathbf{E}[\mathbf{X}_i - 1/p] + \sum_{k=2}^{\infty} \frac{\lambda^k}{k!} \cdot \mathbf{E}[|\mathbf{X}_i - 1/p|^k] \\ &= 1 + \sum_{k=2}^{\infty} \frac{\lambda^k}{k!} (O(k/p))^k \leq 1 + O(\lambda^2/p^2), \quad \text{when } |\lambda| \text{ sufficiently smaller than } p \\ &\leq \exp(O(\lambda^2/p^2)) \end{aligned}$$

Then, since $\mathbf{X}_1, \dots, \mathbf{X}_t$ are all independent, we have

$$\mathbf{E}\left[e^{\lambda(\mathbf{X} - t/p)}\right] \leq \exp(O(\lambda^2 t/p^2)) \implies \mathbf{E}\left[e^{\lambda(\mathbf{Y} - 1)}\right] \leq \exp(O(\lambda^2/t)),$$

and this bound is valid whenever $|\lambda|$ is sufficiently smaller than t .

C Omitted Proofs from Section 3

Proof of Lemma 3.1. The approach is to estimate

$$\mathbf{E}_{i \sim \mathcal{D}}[h_t(p_i)] = \mathbf{E}_{i \sim \mathcal{D}}[g(p_i, p_i^2, \dots, p_i^t)]. \quad (10)$$

There exists an algorithm using $O(\log(1/\epsilon)/\epsilon^2)$ samples to estimate the above quantity: for $j \in \{0, \dots, O(1/\epsilon^2)\}$, one takes a sample $i_j \sim \mathcal{D}$ and uses $r = O(\log(1/\epsilon))$ additional samples $\mathbf{s}_1, \dots, \mathbf{s}_r \sim \mathcal{D}$ to define

$$\mathbf{B}_m^{(j)} \stackrel{\text{def}}{=} \mathbb{1}\{\mathbf{s}_1 = \dots = \mathbf{s}_m = i_j\} \quad \text{and} \quad \mathbf{Z}_j = g(\mathbf{B}_1^{(j)}, \dots, \mathbf{B}_r^{(j)}).$$

Then, let \mathbf{Z} be the average of all \mathbf{Z}_j 's, which is an unbiased estimate to $\mathbf{E}_{i \sim \mathcal{D}}[g(p_i, p_i^2, \dots, p_i^t)]$. Since g is bounded (from Lemma 2.3), the variance of $O(1/\epsilon^2)$ such values is a large constant factor smaller than ϵ^2 . By Chebyshev's inequality, we estimate (10) to error $\pm\epsilon$ with probability at least 0.9. With that estimate, we will now use Lemma 2.4. Specifically, the entropy of \mathcal{D} is exactly $\mathbf{E}_{i \sim \mathcal{D}}[\log(1/p_i)]$, and we have

$$\begin{aligned} \left| \mathbf{E}_{i \sim \mathcal{D}}[\log(1/p_i)] - (\hat{H} - \mathbf{Z}) \right| &\leq \epsilon + \left| \mathbf{E}_{i \sim \mathcal{D}}[\log(1/p_i)] - (\hat{H} - \mathbf{Z}) \right| \\ &\leq \epsilon + \mathbf{E}_{i \sim \mathcal{D}} \left[\left| \log\left(\frac{1}{p_i}\right) - \mathbf{E}[\eta_i] \right| \right] \leq 2\epsilon, \end{aligned}$$

where η_i is the result of running $\text{LogEstimator}(\mathcal{D}, i)$. □

Proof of Lemma 3.2. We note that since $\log(\cdot)$ is monotone increasing, we must have $H \geq \tilde{H}$. To see that it is not much larger, note that we always have $\log z = \ln(z)/\ln(2) \leq (z-1)/\ln(2)$, which means

$$\begin{aligned} H - \tilde{H} &= \mathbf{E}_{i, \mathbf{X}}[\log(\mathbf{X}/\mathbf{X}')] \leq \frac{1}{\ln(2)} \mathbf{E}_{i, \mathbf{X}} \left[\frac{\mathbf{X}}{\min\{\mathbf{X}, X_{\max}\}} - 1 \right] \leq \frac{1}{\ln(2)} \mathbf{E}_{i, \mathbf{X}} \left[\frac{\mathbf{X}}{X_{\max}} \right] \\ &= \frac{1}{X_{\max} \cdot \ln(2)} \sum_{i=1}^k p_i \cdot \frac{t}{p_i} = \frac{tk}{X_{\max} \cdot \ln(2)} = \epsilon. \end{aligned}$$

□

Proof of Lemma 3.4. Substituting the r_ℓ values into Lemma 3.3 ensures $\mathbf{E}[\text{Error}^2] \leq \epsilon^2/10$. Hence the estimator is within $\pm\epsilon$ of \tilde{H} with probability 0.9 by Chebyshev's inequality.

For the intervals $\ell = \{1, \dots, L-1\}$, we always spend r_ℓ tries to determine whether a sample falls within a particular interval. Note that we take one sample to determine $i \sim \mathcal{D}$, and then we take at

most b_ℓ samples. Therefore, the sample complexity for these is

$$\begin{aligned} \sum_{\ell=1}^{L-1} r_\ell \cdot b_\ell &= \frac{80tk}{\epsilon^2} \cdot \sum_{\ell=1}^{L-1} \frac{\log^2(\log^{(\ell-1)}(k)/\epsilon)}{(\log^{(\ell)} k)^3} = \frac{80tk}{\epsilon^2} \cdot \sum_{\ell=1}^{L-1} \frac{(3 \log^{(\ell)}(k) + \log(1/\epsilon))^2}{(\log^{(\ell)} k)^3} \\ &\leq kt \cdot O(\log^2(1/\epsilon)/\epsilon^2), \end{aligned}$$

where we used the fact that

$$\sum_{\ell=1}^{L-1} \frac{1}{(\log^{(\ell)} k)} \leq \frac{1}{1} + \frac{1}{\exp(1)} + \frac{1}{\exp(\exp(1))} + \frac{1}{\exp(\exp(\exp(1)))} + \dots = O(1).$$

Finally, it remains to bound the expected sample complexity of the bucket L . Here, we note

$$r_L = \frac{O(1)}{\epsilon^2} \cdot \log^2\left(\frac{\log^{(L-1)} k}{\epsilon}\right) \leq O\left(\frac{\log^2(1/\epsilon)}{\epsilon^2}\right).$$

Therefore, the expected sample complexity for interval L is $r_L \cdot \sum_{i=1}^k p_i \cdot \frac{t}{p_i} = O(k \log^4(1/\epsilon)/\epsilon^2)$. \square

D Conjectured Lower Bound

Recall that without a memory constraint the sample complexity is known to be $n = \Theta(\max\{\epsilon^{-1} \cdot k/\log(k/\epsilon), \epsilon^{-2} \log^2 k\})$ [VV17, VV11, JVHW15, WY16]. To prove a $\Omega(k/\epsilon^2)$ lower bound for the memory constrained version, we conjecture the following randomized process can be used to generate distributions over $[2k]$ that look alike to any constant space algorithm that uses $o(k/\epsilon^2)$ samples but they have *different* entropies.

Suppose we have k Bernoulli random variables with parameter α : Y_1, \dots, Y_k . And, we have k Rademacher random variables Z_1, \dots, Z_k (that are $+1$ or -1 with probability $1/2$). We construct distribution p in such a way that it is uniform over k pairs of elements $(1, 2), (3, 4), \dots, (2k-1, 2k)$. However, conditioning on pair $(2i-1, 2i)$, we may have a constant bias based on the random variable Y_i . And, we decide about the direction of the bias based on Z_i . More precisely, we set the probabilities in p as follows:

$$p_{2i-1} = \frac{1 + Y_i \cdot Z_i/4}{2k}, \quad p_{2i} = \frac{1 - Y_i \cdot Z_i/4}{2k} \quad \forall i \in [k].$$

Now, it is not hard to show that if we generate two distributions as above with $\alpha = (1 + \epsilon)/2$ and $\alpha = (1 - \epsilon)/2$, then their entropies are $\Theta(\epsilon)$ separated with a constant probability. Thus, any algorithm that can estimate the entropy has to *distinguish* $\alpha = (1 + \epsilon)/2$ from $\alpha = (1 - \epsilon)/2$. Intuitively, to learn α , we would require to *determine* $\Omega(1/\epsilon^2)$ many of Y_i 's. Since we have only a constant words of memory, we cannot perform the estimation of the Y_i 's in parallels. Thus, any natural algorithm would require to draw $\Omega(k/\epsilon^2)$ samples.