How Much Randomness is Needed to Sample from a Discrete Distribution?

Eklavya Sharma

There is a discrete distribution \mathcal{D} of support $S \subseteq \mathbb{Z}$. We want to sample from \mathcal{D} , and our only source of randomness is a (potentially unfair) k-sided die ($k \geq 2$), where the probability of each face is known.

For $i \in \{1, \ldots, k\}$, let q_i be the probability that the die shows i. Without loss of generality assume $0 < q_1 \le q_2 \le \ldots \le q_k$. For $i \in \{0, \ldots, k\}$, let $Q_i := \sum_{j=1}^i q_j$.

1 Algorithm and Notation

Let Z be a random variable from \mathcal{D} . For any $i \in S$, let $p_i := \Pr(Z = i)$ and let I_i be the half-open interval $(\Pr(Z < i), \Pr(Z \le i)]$.

Lemma 1. $(0,1) \subseteq \bigcup_{i \in S} I_i \subseteq [0,1].$

Proof. (TODO)

 Algorithm 1 sampleFromDie: Sample from \mathcal{D} , where throwDie ~ unif($\{1, 2, \dots, k\}$).

 1: Set $\ell = 0$ and r = 1.

 2: for $t \in \mathbb{Z}_{\geq 0}$ do

 3: if $(\ell, r] \subseteq I_i$ for some $i \in S$ then

 4: return (i, t)

 5: end if

 6: v = throwDie()

 7: $\ell, r = \ell + Q_{v-1}(r-\ell), \ell + Q_v(r-\ell)$

 8: end for

Let (Y,T) sampleFromDie's output. Here Y is the random sample from \mathcal{D} we are looking for and T is the number of times throwDie is called.

For $t \in \mathbb{Z}_{\geq 1}$, let v_t be the value of the t^{th} die roll. For $t \in \mathbb{Z}_{\geq 0}$, let ℓ_t and r_t be the value of ℓ and r in sampleFromDie after t iterations have completed. Let $J_t := (\ell_t, r_t]$. (The algorithm may terminate before the t^{th} iteration completes. However, we look at the hypothetical scenario where we continue running the algorithm instead of terminating at return statements. Hence, v_t , ℓ_t , and r_t are well-defined for all t.)

Let \mathcal{J}_t be the support of J_t .

2 Correctness

For any interval I, let len(I) be its length.

Lemma 2. For all $t \in \mathbb{Z}_{\geq 0}$ and $W_t \in \mathcal{J}_t$, we have $\Pr(J_t = W_t) = \operatorname{len}(W_t)$.

Proof. (TODO)

Theorem 3. Let Y be sampleFromDie's output. Then $Pr(Y = i) = p_i$ for all $i \in S$.

Proof. (TODO)

3 Running Time for Bounded |S|

Without loss of generality, assume $S = \{1, 2, ..., n\}$. For $i \in \{0, 1, ..., n\}$, let $s_i := \Pr(Z \leq i)$. For any proposition P, let := (P) be 1 if P is true and 0 if P is false.

Lemma 4. For all $t \in \mathbb{Z}_{\geq 0}$ and $W_t \in \mathcal{J}_t$, we have $\operatorname{len}(W_t) \in [q_1^t, q_k^t]$.

Proof. (TODO)

Lemma 5. $Pr(T > t) \le \min(q_k^t(n-1), 1).$

Proof. Let $C := \{s_1, \ldots, s_{n-1}\}$. Then $T > t \iff C \cap J_t \neq \emptyset$. For any $W \in \mathcal{J}_t$, let $n_W := |C \cap W|$. Then

$$Pr(T > t) = Pr(C \cap J_t \neq \emptyset)$$

= $\sum_{W \in \mathcal{J}_t} Pr(J_t = W) := (n_W \ge 1)$
 $\leq \sum_{W \in \mathcal{J}_t} q_k^t n_W$
= $q_k^t \sum_{W \in \mathcal{J}_t} |C \cap W|$
= $q_k^t (n-1).$

Corollary 5.1. Let $\delta \in (0,1)$ and $t := \left\lceil \frac{\log(n-1) + \log(1/\delta)}{\log(1/q_k)} \right\rceil$. Then $\Pr(T > t) \le \delta$.

Theorem 6. Let $h := \left\lceil \frac{\log(n-1)}{\log(1/q_k)} \right\rceil$. Then $E(T) \le h + \frac{(n-1)q_k^h}{1-q_k} \le h + \frac{1}{1-q_k}$.

Proof. For any $t \in \mathbb{Z}_{\geq 0}$,

$$t \ge h \iff t \ge \frac{\log(n-1)}{\log(1/q_k)} \iff t \log(1/q_k) \ge \log(n-1)$$
$$\iff q_k^{-t} \ge n-1 \iff (n-1)q_k^t \le 1.$$

Hence,

$$\begin{split} \mathbf{E}(T) &= \sum_{t=0}^{\infty} \Pr(T > t) \le \sum_{t=0}^{\infty} \min(1, q_k^t (n-1)) \\ &\le h + \sum_{t=h}^{\infty} q_k^t (n-1) = h + (n-1) q_k^h \sum_{t=0}^{\infty} q_k^t \\ &= h + \frac{(n-1) q_k^h}{1 - q_k} \le h + \frac{1}{1 - q_k}. \end{split}$$

If we have a fair *n*-sided die, then we get $E(T) \leq 2$.

4 Entropy-Based Bound

All intervals in \mathcal{J}_t are said to have *depth* t.

Definition 1. Let $\mathcal{J} := \bigcup_{t=0}^{\infty} J_t$. An interval $W \in \mathcal{J}$ is said to be maximal in I_i if $W \subseteq I_i$ and for all $W' \in \mathcal{J}$, we have $W' \supseteq W \implies W' \nsubseteq I_i$.

A die-decomposition of I_i is the set of all sets in \mathcal{J} that are maximal in I_i .

Lemma 7. For any $i \in S$, let \mathcal{K}_i be a die-decomposition of I_i . Then $\bigcup_{W \in \mathcal{K}_i} W = I_i$.

Proof. (TODO)

Lemma 8. Let (Y,T) be the output of sampleFromDie. Then $(\ell_T, r_T) \in \mathcal{K}_Y$.

Proof. (TODO)

Lemma 9. For any $i \in S$ and $t \in \mathbb{Z}_{\geq 0}$, we get

$$\sum_{W \in \mathcal{K}_i \cap \mathcal{J}_t} \operatorname{len}(W) \le (2 - q_1 - q_k) q_k^{t-1}.$$

Proof. (TODO)

Lemma 10.

$$\mathcal{E}(T) \leq \frac{2 - q_1 - q_k}{1 - q_k} \sum_{i \in S} p_i^{\frac{\log(1/q_k)}{\log(1/q_1)}} \left(\left\lceil \frac{\log(1/p_i)}{\log(1/q_1)} \right\rceil + \frac{q_k}{1 - q_k} \right).$$

Proof. Let $\beta_{i,t} := \sum_{W \in \mathcal{K}_i \cap \mathcal{J}_t} \operatorname{len}(W)$. Then $\beta_{i,t} \leq \min(p_i, (2 - q_1 - q_k)q_k^{t-1})$. $\beta_{i,t} > 0 \iff |\mathcal{K}_i \cap \mathcal{J}_t| > 0 \implies q_1^t \leq p_i$. Let $h_i := \left\lceil \frac{\log(1/p_i)}{\log(1/q_1)} \right\rceil$. Then

$$q_1^t \le p_i \iff t \log(1/q_1) \ge \log(1/p_i) \iff t \ge h_i.$$

Hence, $t < h_i \implies \beta_{i,t} = 0$.

$$\begin{split} \mathbf{E}(T) &= \sum_{i \in S} \sum_{t=0}^{\infty} \mathbf{E}(T|J_T \in \mathcal{K}_i \cap \mathcal{J}_t) \operatorname{Pr}(J_T \in \mathcal{K}_i \cap \mathcal{J}_t) \\ &= \sum_{i \in S} \sum_{t=0}^{\infty} t \operatorname{Pr}(J_T \in \mathcal{K}_i \cap \mathcal{J}_t) \\ &= \sum_{i \in S} \sum_{t=0}^{\infty} t \beta_{i,t} = \sum_{i \in S} \sum_{t=h_i}^{\infty} t \beta_{i,t} \\ &\leq (2 - q_1 - q_k) \sum_{i \in S} \sum_{t=h_i}^{\infty} t q_k^{t-1} \\ &= \frac{2 - q_1 - q_k}{1 - q_k} \sum_{i \in S} q_k^{h_i - 1} \left(h_i + \frac{q_k}{1 - q_k}\right) \\ &= \frac{2 - q_1 - q_k}{1 - q_k} \sum_{i \in S} p_i^{\frac{\log(1/q_k)}{\log(1/q_1)}} \left(h_i + \frac{q_k}{1 - q_k}\right) \end{split}$$

Suppose $q_1 = q_k = 1/k$. Then we get

$$E(T) \leq \frac{2(1-1/k)}{1-1/k} \sum_{i \in S} p_i \left(\lceil \log_k(1/p_i) \rceil + \frac{k}{k-1} \right)$$
$$= 2 \sum_{i \in S} p_i \left\lceil \log_k(1/p_i) \rceil + \frac{2k}{k-1} \right.$$