# 3: Randomness and Entropy

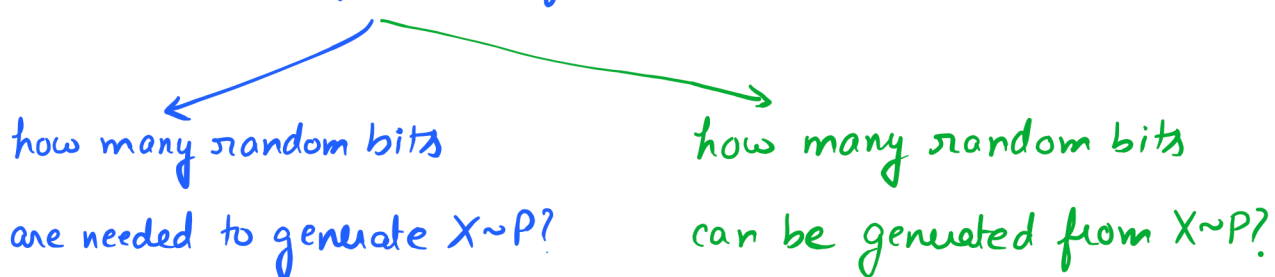[A] Uncertainty and randomness

In the previous two lectures, we identified entropy as a measure of uncertainty and therefore, the information revealed on observing $X \sim P$ is $H(P)$.

Information = reduction of uncertainty

But does uncertainty equal randomness?

How much randomness is there in $X \sim P$?

(two possible definitions)

how many random bits are needed to generate $X \sim P$?

how many random bits can be generated from $X \sim P$?

Some examples:

Example 1. How many independent unbiased coin flips are needed to generate samples from:

(a) $P = \left( \frac{1}{2}, \frac{1}{4}, \frac{1}{4} \right)$     (b) $P = \left( \frac{2}{3}, \frac{1}{3} \right)$

avg. no.          worst case
                  on avg.

Example 2. How many independent samples from $Ber(p)$ are needed to generate one sample from $Ber(\frac{1}{2})$?

**B** Measuring distance between distributions:

Total variation distance $d(P, Q)$

$$d(P, Q) := \frac{1}{2} \sum_x |P(x) - Q(x)|$$

Properties

1. $d(P, Q) = \sup_A P(A) - Q(A) = \sup_A |P(A) - Q(A)|$

$$= \sup_B Q(B) - P(B) = \sup_B |Q(B) - P(B)|$$

where the sup in the first line is attained by

$$A^* = \{x : P(x) > Q(x)\}$$

and that in the second line by

$$B^* = \{x : Q(x) > P(x)\}$$

2. $0 \leq d(P, Q) \leq 1,$

with equality on the left-side iff $P = Q$

and equality on the right-side iff $supp(P) \cap supp(Q) = \emptyset$.

3. We can extend the definition to distributions with densities:

$$d(P, Q) = \frac{1}{2} \int |f(x) - g(x)| \, dx$$

$\hookrightarrow$ density of P $\quad \longrightarrow$ density of Q

$$= \sup_A P(A) - Q(A)$$

where the equality in the second line attained by

$$A^* = \{x : f(x) > g(x)\}.$$

## C Generating almost random bits

Given a sample $X \sim P$, we want to generate $U \sim \text{unif}(\{0,1\}^\ell)$.

Specifically, for $0 < \varepsilon < 1$, what is the largest $\ell$ s.t. we can find $f : \mathcal{X} \to \{0,1\}^\ell$ s.t.

$$d\left(P_{f(X)}, P_U\right) \leq \varepsilon ?$$

### A scheme

Let $\mathcal{J}_\geqslant = \{x : -\log p(x) > \lambda\}$ and suppose that
$$P(\mathcal{J}_\geqslant) \geqslant 1 - \varepsilon.$$

Consider a partition of $\mathcal{J}_\geqslant$ into parts $\mathcal{X}_1, \ldots, \mathcal{X}_M, \mathcal{X}_{M+1}$ that each part $\mathcal{X}_i$ satisfies

(#1)     $2^{-\lambda} \cdot N < P(\mathcal{X}_i) \leq 2^{-\lambda} \cdot (N+1)$,    for $1 \leq i \leq M$.

and

(#2)                 $P(\mathcal{X}_{M+1}) \leq 2^{-\lambda} \cdot N$

<span style="color:red">→ we can only find such a partition because $p(x) \leq 2^{-\lambda} \; \forall x \in \mathcal{J}_\geqslant$.</span>

Thus,

$$\left(\frac{1 - \varepsilon - 2^{-\lambda} \cdot N}{N+1}\right) 2^\lambda \;\leq\; M \;<\; \frac{2^\lambda}{N}.$$

Let $U \sim \text{unif}(\{1, \ldots, M'\})$ where $M' = \dfrac{2^\lambda}{N}$ and

$$f(x) = \begin{cases} i, & \text{if } x \in \mathcal{X}_i \text{ for some } i, \\ \bot, & \text{if } x \notin \mathcal{J}_\geqslant. \end{cases}$$

We now analyse our scheme.

$$2\, d\left(P_{f(x)}, P_U\right) = \sum_{m=1}^{M} \left| P_{f(x)}(m) - \frac{1}{M'} \right| + \sum_{m=M+1}^{M'} \left| P_{f(x)}(m) - \frac{1}{M'} \right|$$

$$+ \ P_{f(x)}(\perp)$$

$$P_{f(x)}(\perp) = P\left(J_\lambda^c\right) \leq \varepsilon$$

$$P_{f(x)}(M+1) \leq \frac{2^{-\lambda}}{N}, \qquad P_U(M+1) = \frac{2^{-\lambda}}{N} = \frac{1}{M'}$$

$$P_{f(x)}(m) = 0, \qquad P_U(m) = \frac{2^{-\lambda}}{N}, \quad \text{for } m \in \{M+2, \ldots, M'\}.$$

$$\Rightarrow \sum_{m=M+1}^{M'} \left| P_{f(x)}(m) - \frac{1}{M'} \right| = \frac{M'-M}{M'} - P_{f(x)}(M+1)$$

$$\leq 1 - \frac{M}{M'} \leq 1 - \frac{(1 - \varepsilon - 2^{-\lambda} \cdot N) \cdot N}{N+1}$$

$$\sum_{m=1}^{M} \left| P_{f(x)}(m) - \frac{1}{M'} \right| = \sum_{m=1}^{M} \left| P(x_m) - \frac{1}{M'} \right|$$

$$\leq M \cdot 2^{-\lambda}$$

$$\leq \frac{1}{N},$$

where we used (#1) and (#2).

On combining the bounds above, we get

$$2\, d\left(P_{f(x)}, P_U\right) \leq \frac{1}{N} + \frac{1}{N+1} + 2\varepsilon + \varepsilon \cdot 2^{-\lambda}$$

Setting $N = \frac{1}{\varepsilon} - 1$,

$$2d(P_{f(x)}, P_U) \leq 2\left(\varepsilon + \varepsilon + \frac{2^{-\lambda}}{\varepsilon}\right) + \varepsilon$$

$$\Rightarrow d(P_{f(x)}, P_U) \leq \frac{5\varepsilon}{2} + \frac{2 \cdot 2^{-\lambda}}{\varepsilon}.$$

Namely, we have the following theorem:

Theorem. For $\lambda > 0$ and $0 < \varepsilon < 1$, suppose that

$$P\left(\{x : -\log p(x) > \lambda\}\right) \geq 1 - \varepsilon.$$

Then, for $U \sim$ uniformly over at least $\boxed{r = \lambda - \log \frac{1}{\varepsilon}}$ bits,

we have

$$d(P_{f(x)}, P_U) \leq \frac{5}{2}\varepsilon + 2 \cdot 2^{-r}.$$

In particular, if $\lambda \geq 2\log\frac{1}{\varepsilon}$, we get $\Big\}$ We have generated

$\lambda - \log\frac{1}{\varepsilon}$ bits.

$$d(P_{f(x)}, P_U) \leq \frac{7}{2}\varepsilon.$$

So, what is the largest $\lambda$ we can have?

It's roughly the entropy!

You can use the Chebyshev's inequality to show:

$$P\left(-\log P(x) > H(P) - \sqrt{\frac{\text{Var}\left(-\log P(x)\right)}{\varepsilon}}\right)$$

$$\geq 1 - \varepsilon.$$

## D  Generating a sample $X \sim P$ using uniform bits

What is the least number of bits $l$ such that for $U \sim \text{unif}(\{0, \beta^l\})$ we can find a function $f$ s.t. $d(P_{f(u)}, P) \leq \varepsilon$?

### A scheme

Consider $U \sim \text{unif}(\{1, \ldots, M\})$ and suppose that

$$P\left(\{x : -\log p(x) \leq \lambda\}\right) \geq 1 - \varepsilon.$$

Denoting $\mathcal{I}_\lambda = \{x : -\log p(x) \leq \lambda\}$, each $x \in \mathcal{I}_\lambda$ satisfies $p(x) \geq 2^{-\lambda}$. Suppose that $M > 2^{-\lambda}$. Consider a partition $\mathcal{Y}_1, \ldots, \mathcal{Y}_{M'+1}$ of $\{1, \ldots, M\}$ such that $M' = |\mathcal{I}_\lambda|$ and

$$p(x_i) < P_U(\mathcal{Y}_i) \leq p(x_i) + \frac{1}{M}, \quad 1 \leq i \leq M',$$

where $x_1, \ldots, x_{M'}$ denote the elements of $\mathcal{I}_\lambda$. Further, let $x_{M'+1}$ denote an arbitrary element outside $\mathcal{I}_\lambda$.

Define

$$f(y) = x_i \quad \text{if } y \in \mathcal{Y}_i, \quad 1 \leq i \leq M'+1.$$

Then,

$$2d(P_{f(u)}, P) \leq \sum_{x \in \mathcal{I}_\lambda} |P_{f(u)}(x) - p(x)|$$
$$+ P(\mathcal{I}_\lambda^c) + P_U(\mathcal{Y}_{M'+1})$$

For the first term, by our construction,

$$\left| P_{f(v)}(x_i) - p(x_i) \right| = \left| P_v(y_i) - p(x_i) \right|$$
$$\leq \frac{1}{M} \quad \text{for } i \in \{1, \dots, M'\}.$$

Thus,

$$\sum_{x \in J} \left| P_{f(v)}(x) - p(x) \right| \leq \frac{M'}{M}.$$

For the third term, for each $i \in \{1, \dots, M'\}$,

$$p(x_i) \leq P_v(y_i).$$

which on summing over $i \in \{1, \dots, M\}$ gives

$$P(J_\lambda) \leq P_v\left( \bigcup_{i=1}^{M'} y_i \right) = 1 - P_v(y_{M'+1})$$

$$\Rightarrow P_v(y_{M'+1}) \leq P(J_\lambda^c).$$

Thus, noting that $P(J_\lambda^c) \leq \varepsilon$,

$$d(P_{f(v)}, P) \leq \frac{M'}{2M} + \varepsilon = \frac{|J_\lambda|}{2M} + \varepsilon$$

$$\leq \frac{2^\lambda}{2 \cdot M} + \varepsilon \qquad \text{(why?)}$$

In particular, upon choosing $M = \frac{2^\lambda}{\varepsilon}$, we get

$$d(P_{f(v)}, P) \leq \frac{3\varepsilon}{2}.$$ We have shown the following result:

**Theorem** Suppose that
$$P\left(\{x: -\log p(x) \leq \lambda\}\right) \geq 1-\varepsilon.$$
Then, for $\boxed{\ell = \lambda + \log \frac{1}{\varepsilon}}$, using $U \sim \text{unif}(\{0,1\}^\ell)$
we can find a function $f$ s.t.
$$d\left(P_{f(U)}, P\right) \leq \frac{3}{2}\varepsilon.$$

We saw in the previous lecture that a good estimate
for $\lambda$ required in the theorem above is $\dfrac{H(P)}{\varepsilon}$.

**E** _Typical sets and Entropy_

In all our proofs, the following sets appeared:
$$\mathcal{J}_\lambda^{(1)} = \{x: -\log p(x) \leq \lambda\}$$
$$\mathcal{J}_\lambda^{(2)} = \{x: -\log \quad x) > \lambda\}$$
We want $\lambda$ both the sets have probability exceeding $1-\varepsilon$.
We saw that a good choice is roughly $H(P)$.

* $|\mathcal{J}_\lambda^{(1)}| \leq 2^\lambda$

* If $P\left(\mathcal{J}_\lambda^{(2)}\right) \geq 1-\varepsilon$, $\quad |\mathcal{J}_\lambda^{(2)}| \geq 2^\lambda(1-\varepsilon)$

These sets are called _typical sets_ since they contain
elements which occur with large probability.

Furthermore, probabilities of elements in these sets are "close"

to a uniform distribution on this set.

* <u>The case of iid source</u>

$X \equiv (X_1, \ldots, X_n) \sim$ iid $P$

Then, for $\lambda = n(H(P) + \eta)$ and $Z_i = -\log P(X_i)$

$$P(J_\lambda^{(1)}) = P\left(\sum_{i=1}^{n} Z_i \leq n(H(P) + \eta)\right)$$

$$\geq 1 - \sqrt{\frac{Var(Z_1)}{n\eta}}.$$

and for $\lambda = n(H(P) - \eta)$,

$$P(J_\lambda^{(2)}) = P\left(\sum_{i=1}^{n} Z_i > n(H(P) - \eta)\right)$$

$$\geq 1 - \sqrt{\frac{Var(Z_1)}{n\eta}}.$$

* <u>Entropy is additive</u>

For independent $X_1, X_2$,

$$H(X_1, X_2) = H(X_1) + H(X_2)$$

This is perhaps the most important property of Shannon entropy.