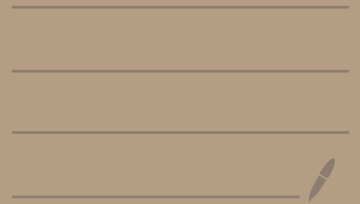


Information Theory & Coding

Sept 29, 2020



Recall:

$$H(u) \leq E[\text{code word length}^{\text{best}}] \leq H(u) + 1$$

$$\Rightarrow \frac{1}{n} H(u_1 \dots u_n) \leq E\left[\frac{1}{n} \text{length } \hat{c}_n(u_1 \dots u_n)\right] \leq \frac{1}{n} H(u_1 \dots u_n) + \frac{1}{n}$$

of bits/letter

Def: given a source u_1, u_2, u_3, \dots
(a stochastic process), we say the entropy-rate of the source is

$$\lim_{n \rightarrow \infty} \frac{1}{n} H(u_1 \dots u_n) =: \mathcal{H}(\{u_i\})$$

if the limit exists.

with this definition:

$$\lim_{n \rightarrow \infty} \frac{1}{n} E(\text{length } \hat{c}_n(u_1 \dots u_n)) = \mathcal{H}(\{u_i\})$$

Ex: if u_1, u_2, u_3 are i.i.d then

$$H(\{u_i : i \in \mathbb{N}\}) = H(u_1)$$

Pf: $H(u_1 \dots u_n) \stackrel{\text{chain rule}}{=} H(u_1) + H(u_2|u_1) + \dots + H(u_n|u_1 \dots u_{n-1})$

$$= H(u_1) + H(u_2) + \dots + H(u_n)$$

\nearrow independence ident dist
 $= \sum_n H(u_1)$

$$\Rightarrow \frac{1}{n} H(u_1 \dots u_n) = H(u_1) //$$

Def: A stochastic process u_1, u_2, \dots is

said to be stationary if for every $n \geq 1$

↳ every $k \geq 1$ and every $u_1 \dots u_n$

$$\Pr(u_1 \dots u_n = u_1 \dots u_n)$$

$$= \Pr(u_{1+k}, u_{2+k}, \dots, u_{n+k} = u_1 \dots u_n).$$

Thm: If U_1, U_2, U_3, \dots is a stationary process then

$\lim_{n \rightarrow \infty} \frac{1}{n} H(U_1 \dots U_n)$ exists and

equals $\lim_{n \rightarrow \infty} H(U_n | U_1 \dots U_{n-1})$

Pf: Let $a_n = H(U_n | U_1 \dots U_{n-1})$. 1st claim

$0 \leq a_{n+1} \leq a_n$ which will imply the

$a = \lim_n a_n$ exists. To show the claim:

$$a_{n+1} = H(U_{n+1} | U_n, U_{n-1}, \dots, U_1)$$

$$\leq H(U_{n+1} | U_n, \dots, U_2) \quad (\text{cond. reduction entropy})$$

$$\Rightarrow H(U_n | U_{n-1}, \dots, U_1) = \underline{a_n}$$

$$(U_1 \dots U_n) \sim (U_2 \dots U_{n+1})$$

has the same stats as (stationarity).

Now observe

$$S_n \triangleq \frac{1}{n} H(u_1 \dots u_n)$$

$$= \frac{1}{n} \left[H(u_1) + H(u_2(u_1)) + \dots + H(u_n(u_1 \dots u_{n-1})) \right]$$

$$= \frac{1}{n} \left(\underbrace{a_1} + \underbrace{a_2} + \dots + \underbrace{a_n} \right)$$

Fact: (Cesàro): if x_1, x_2, x_3, \dots is a \mathbb{R} -

valued sequence with $x = \lim_{n \rightarrow \infty} x_n$ then

$\gamma_n \triangleq \frac{1}{n} (x_1 + \dots + x_n)$ also has a limit &

$$\lim_{n \rightarrow \infty} \gamma_n = x.$$

All we need to do is to prove Cesàro:

$$\left(\lim_{n \rightarrow \infty} x_n = x \right) \iff \forall \epsilon > 0 \exists N_0(\epsilon) \forall n \geq N_0(\epsilon) \left(|x_n - x| < \epsilon \right)$$

We need to show that $\lim_{n \rightarrow \infty} \gamma_n = x$

$$\gamma_n - x = \frac{1}{n} \left[(x_1 - x) + \dots + (x_n - x) \right]$$

$$|\gamma_n - x| \leq \frac{1}{n} \left[|x_1 - x| + \dots + |x_n - x| \right]$$

$$= \frac{1}{n} \left[\sum_{i=1}^{n_0(\varepsilon)} |x_i - x| + \sum_{i=n_0(\varepsilon)+1}^n |x_i - x| \right] \quad n \geq n_0(\varepsilon)$$

$$\leq \varepsilon + \frac{1}{n} \sum_{i=1}^{n_0(\varepsilon)} |x_i - x|$$

for $\forall x$ choose $n \geq \max\{n_0(\varepsilon), \dots\}$

$$\frac{\sum_{i=1}^{n_0(\varepsilon)} |x_i - x|}{n} < \varepsilon$$

$$=: n_1(\varepsilon).$$

thus for we have \implies

$$|y_n - x| < \varepsilon + \varepsilon = 2\varepsilon.$$

$$\implies \lim y_n = x. \quad //$$

Ex: Suppose u_1, u_2, u_3, \dots is a

Markov process, & stationary; i.e.

$$\left. \begin{aligned} &P(u_n = u_n | u_{n-1} = u_{n-1}, \dots, u_1 = u_1) \\ &= P(u_n = u_n | u_{n-1} = u_{n-1}) \end{aligned} \right\} \text{Markovity.}$$

Then

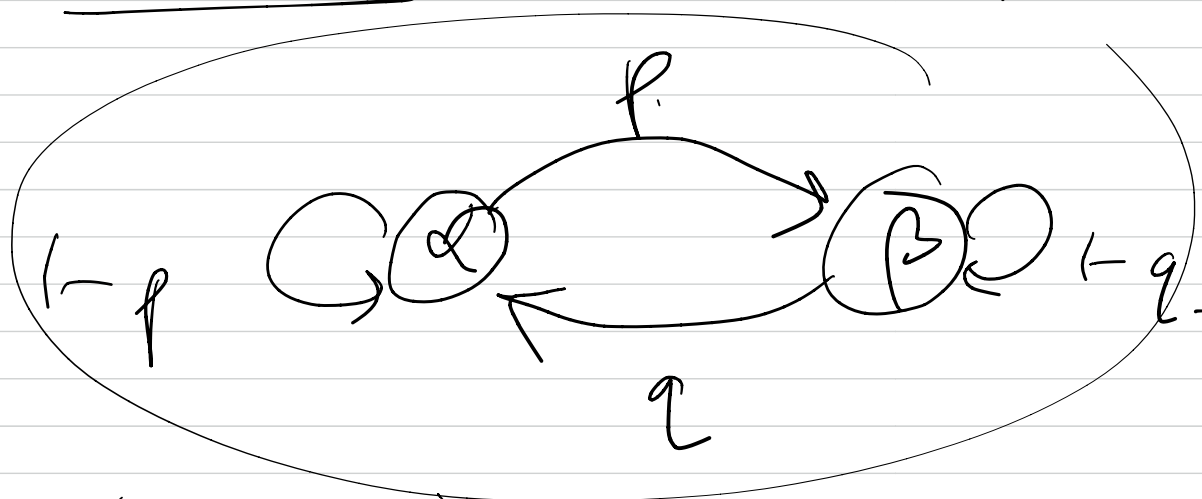
$$H(\{u_i : i \in \mathbb{N}\}) = \lim_{n \rightarrow \infty} H(u_n | u_{n-1}, \dots, u_1)$$

$$= \lim_{n \rightarrow \infty} H(u_n | u_{n-1})$$

$$= \lim_{n \rightarrow \infty} H(u_2 | u_1)$$

$$= H(u_2 | u_1).$$

Example: $u_i \in \{\alpha, \beta\}$.



$$P(u_{n+1} = \beta | u_n = \alpha) = p$$

$$P(u_{n+1} = \alpha | u_n = \beta) = q$$

$$\pi(\alpha) = \underbrace{P(u_{n+1} = \alpha)}_{\pi(\alpha)} = \underbrace{P(u_n = \alpha)}_{\pi(\alpha)}(1-p) + \underbrace{P(u_n = \beta)}_{(1 - \pi(\alpha))} q$$

$$\pi(\alpha) = \pi(\alpha)(1-p) + (1-\pi(\alpha))q$$

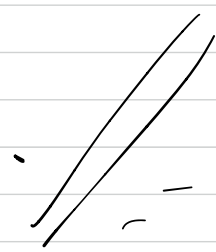
$$\pi(\alpha) = \Pr(U_n = \alpha) = \Pr(U_{n+1} = \alpha) = \dots = \Pr(U_1 = \alpha)$$

$$\pi(\alpha) = \frac{q}{p+q}, \quad \pi(\beta) = \frac{p}{p+q}$$

$$\Rightarrow H(U_2 | U_1) = \frac{q}{p+q} \underbrace{H(U_2 | U_1 = \alpha)} + \frac{p}{p+q} H(U_2 | U_1 = \beta)$$

$$= \frac{q}{p+q} \left[p \log_2 \frac{1}{p} + (1-p) \log_2 \frac{1}{1-p} \right]$$

$$+ \frac{p}{p+q} \left[q \log_2 \frac{1}{q} + (1-q) \log_2 \frac{1}{1-q} \right]$$



Back to where we left off yesterday: Remember

we were considering iid processes we had defined the notion of typicality:

Def: given a distribution p on \mathcal{U} , $\epsilon > 0$, $n \in \{1, 2, 3, \dots\}$

we say a sequence u_1, u_2, \dots, u_n to be ϵ -typical w.r.t p if

$$\forall u \in \mathcal{U} \quad \frac{\#\{i : u_i = u\}}{n} \in [(1-\epsilon)p(u), (1+\epsilon)p(u)].$$

We let $T(n, \epsilon, p) =$ set of all ϵ -typical sequences w.r.t p of length n .

① We have seen that if $(u_1, \dots, u_n) \in T(n, p, \epsilon)$

$$\Pr(\underbrace{u_1, \dots, u_n}_{\text{i.i.d.} \sim p} = \underbrace{u_1, \dots, u_n}) = \underline{\underline{2^{-n H(p)} (1 \pm \epsilon)}}$$

② $1 \geq \Pr(\underbrace{(u_1, \dots, u_n)}_{\text{i.i.d.} \sim p} \in T(n, p, \epsilon))$

$$= \sum_{(u_1, \dots, u_n) \in T} \Pr(u_1, \dots, u_n = u_1, \dots, u_n) \geq \sum_{(u_1, \dots, u_n) \in T} 2^{-n H(p)} (1 \pm \epsilon)$$

$$= |T(n, p, \epsilon)| \cdot 2^{-n H(p)} (1 \pm \epsilon)$$

$$\Rightarrow \left(|T(n, p, \epsilon)| \leq 2^{n H(p)} (1 \pm \epsilon) \right)$$

$\underbrace{\quad}_{=}$ $\underbrace{\quad}_{\sim p}$

③. we will show

$$\Pr((U_1, \dots, U_n) \in T(n, p, \epsilon)) \approx 1.$$

Lemma: Fix $p, \epsilon > 0$, then

$$\lim_{n \rightarrow \infty} \Pr(\underbrace{(U_1, \dots, U_n)}_{\text{i.i.d. } \sim p} \in T(n, p, \epsilon)) = 1.$$

Pf: $\{(U_1, \dots, U_n) \notin T(n, p, \epsilon)\}$

$$= \bigcup_{u \in \mathcal{U}} \left\{ \underbrace{|\{i : U_i = u\}|}_{\frac{1}{n}} \notin [(1-\epsilon)p(u), (1+\epsilon)p(u)] \right\}$$

$$\Pr((U_1, \dots, U_n) \notin T) \leq \sum_{u \in \mathcal{U}} \Pr(\downarrow)$$

we will show that for every $u \in \mathcal{U}$

$$\Pr\left(\frac{1}{n} |\{i : U_i = u\}| \notin [(1-\epsilon)p(u), (1+\epsilon)p(u)]\right)$$

$\rightarrow 0$ as n gets large

To do so, let $X_i = \begin{cases} 1, & U_i = u \\ 0, & U_i \neq u \end{cases}$

$$\rightarrow \frac{1}{n} |\{i : U_i = u\}| = \frac{1}{n} \sum_{i=1}^n X_i$$

observe X_1, X_2, \dots, X_n are i.i.d.

$$E(X_i) = p(u), \quad \text{Var}(X_i) = p(u) - p(u)^2 = p(u)(1-p(u)).$$

$$\Pr\left(\left|\frac{1}{n} \sum_{i=1}^n X_i - p(u)\right| > \varepsilon p(u)\right)$$

Recall Chebyshev's inequality:

$$\Pr(|Y - E(Y)| > \alpha) \leq \frac{\text{Var}(Y)}{\alpha^2}$$

$$\leq \frac{n \text{Var}(X_1)}{n^2 \varepsilon^2 p(u)^2} = \frac{\text{Var}(X_1)}{n \varepsilon^2 p(u)^2}$$

$$\Rightarrow \lim_{n \rightarrow \infty} \Pr\left(\frac{1}{n} |\{i: u_i = u\}| \notin [(1-\varepsilon)p(u), (1+\varepsilon)p(u)]\right) = 0$$

we have proved $\forall \varepsilon > 0, \exists n_0(\varepsilon)$ s.t. $\forall n > n_0(\varepsilon)$

$$\Pr\left(\underbrace{(u_1, \dots, u_n)}_{\sim \text{iid } p} \in T(n, p, \varepsilon)\right) > 1 - \varepsilon$$

Corollary: for n large enough

$$(1 - \varepsilon) < \Pr\left(\underbrace{(u_1, \dots, u_n)}_{\text{iid } p} \in T(n, p, \varepsilon)\right) = \sum_{(u_1, \dots, u_n) \in T} \Pr(u_1, \dots, u_n)$$

$$\leq \sum_{(u_1, \dots, u_n) \in T} 2^{-n H(u)} (1 - \varepsilon) = |T(n, p, \varepsilon)| \cdot 2^{-n H(u)} (1 - \varepsilon)$$

$$\Rightarrow |T(n, p, \varepsilon)| > \underbrace{(1 - \varepsilon) 2^{n H(u)}}_{(1 - \varepsilon) 2^{n H(u)}}$$

Summary: $T(n, p, \epsilon)$ has the following properties:

①. $(u_1 \dots u_n) \in T(n, p, \epsilon)$ then

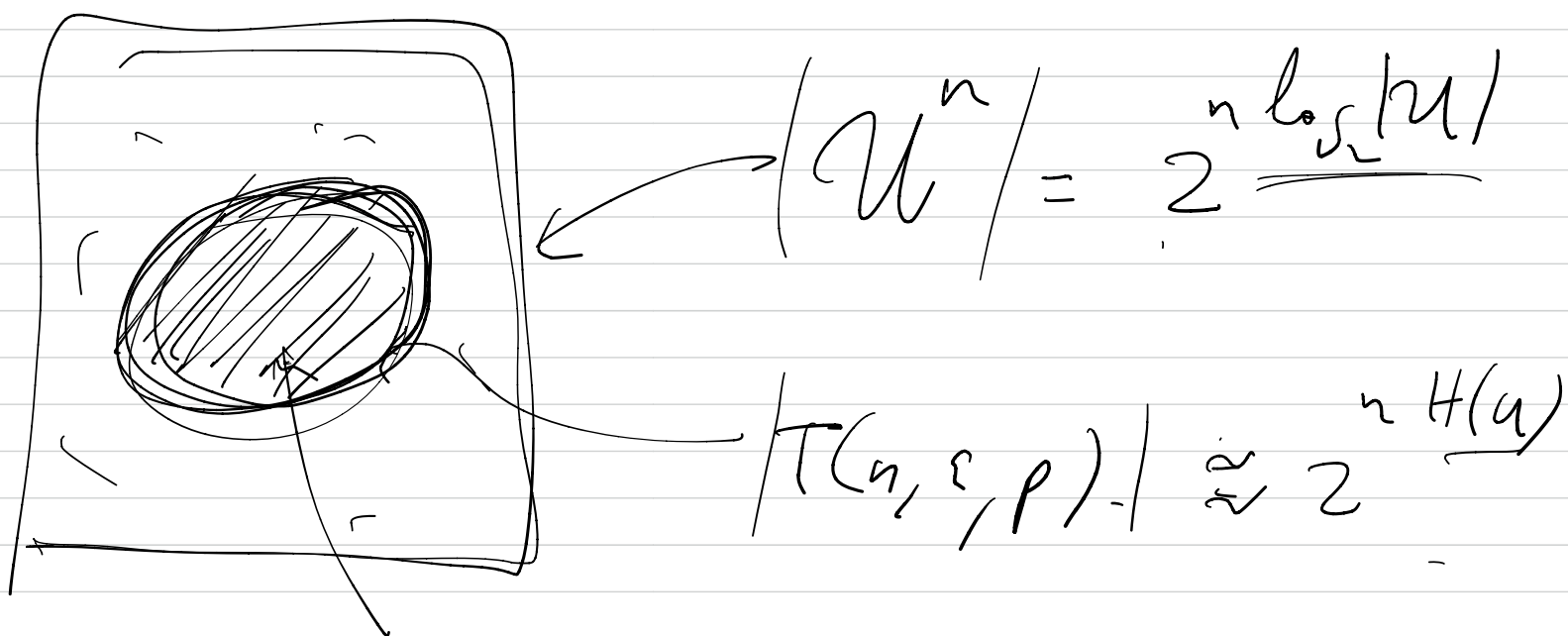
$$\Pr(U_1 \dots U_n = u_1 \dots u_n) = 2^{-n H(u) (1 \pm \epsilon)}$$

iid $\sim p$

②. $|T(n, p, \epsilon)| \leq 2^{n H(u) (1 + \epsilon)}$

③. for large n
 $|T(n, p, \epsilon)| > (1 - \epsilon) 2^{n H(u) (1 - \epsilon)}$

[Asymptotic equipartition Property].



each element has the \approx same probability of being produced by the iid $\sim p$ source.

So an "almost correct" view of the iid source is a probabilistic device that picks uniformly at random an element from T .

Conceptual way for compressing iid sources,

- give each element of $T(n, \epsilon, p)$ a binary representation. Since $|T| \leq 2^{nH(u)(1+\epsilon)}$

$\lceil n(1+\epsilon)H(u) \rceil$ bits is enough.

- when the source produces $u_1 \dots u_n$ check if $(u_1 \dots u_n) \in T$, and if so emit the binary representation prefixed by 0 else ~~end universe~~

emit 1 followed by $\lceil n \log_2 |T| \rceil$ bit representation of $(u_1 \dots u_n)$.

So we describe n letters by

$$\leq \begin{cases} \frac{nH(u)(1+\epsilon) + 1 + 1}{1} & \text{if } u_1 \dots u_n \in T \\ \frac{n \log_2 |T| + 2}{1} & \text{else} \end{cases}$$

bits

In bits/letter we have

$$\leq \left\{ \begin{array}{l} \underbrace{H(u) + \epsilon}_1 + \frac{2}{n} \quad u_1 \dots u_n \in T \\ \underbrace{L(u)}_1 + \frac{2}{n} \quad \notin T \end{array} \right\}$$

This gives a way to represent an iid source with $\approx H(u)$ bits/letter almost all the time