

$$\begin{aligned}
 1 \quad H_b(X) &= -\sum_{x \in X} p(x) \log p(x) \quad D(p||q) = \sum_x p(x) \log \frac{p(x)}{q(x)} = -\mathbb{E}[\log \frac{p(x)}{q(x)}] - H(X) \\
 2 \quad &= -\mathbb{E}[\log \frac{p(x)}{p(x)}] \\
 3 \quad I(X;Y) &= D(p_{XY} || p_X p_Y) = H(X) + H(Y) - H(X,Y) \quad H(X,Y) = -\sum_{x,y} p(x,y) \log p(x,y) \\
 4 \quad &= -\sum_{x,y} p(x,y) \log p(x,y) \\
 5 \quad D(p_{Y|X} || q_{Y|X} | p_X) &= \sum_x p(x) D(p_{Y|X=x} || q_{Y|X=x}) \quad H(X|Y) = \mathbb{E}[H(X|Y)] = \mathbb{E}[-\sum_{x,y} p(x,y) \log p(x|y)] \\
 6 \quad I(X;Y|Z) &= I(Y;X|Z) = H(X|Z) - H(X|Y,Z) = H(Y|Z) - H(Y|X,Z)
 \end{aligned}$$

Inequalities (Jobs): $-\sum p(x) \log p(x) \leq -\sum p(x) \log \frac{p(x)}{\sum a_i} = \Leftrightarrow p=q$ [Jensen: $a_i > 0, f$ convex]

(\log - sum): $\sum_i a_i \log \frac{a_i}{b_i} \geq (\sum_i a_i) \log \frac{\sum a_i}{\sum b_i}$; for: $a_i, b_i > 0, 0 < \sum a_i < \sum b_i \Rightarrow \frac{a_i}{b_i} < 1$, const. $f(\frac{\sum a_i}{\sum b_i}) \leq \frac{\sum a_i f(a_i)}{\sum b_i}$

Fano: $H(X|Y) \leq H(1_{X \neq Y}) + I(X \neq Y)$ ($|X|=D \leq 1 + I(X \neq Y) \leq 1 + I(X \neq Y)$) OR: $f(\mathbb{E}X) \leq \mathbb{E}(f(X))$

entropy $D \geq 0 \leq H(X) \leq \log(|X|)$ information: $I(X;Y) \geq 0 \Leftrightarrow X \perp Y$

$\exists D \geq 0 \leq H(X|Y) \leq H(X)$ $I(X;Y) = H(X) - H(X|Y) = H(Y) - H(Y|X)$

$X=f(Y)$ $I(X;Y) = \sum_i I(X_i;Y|X_1, \dots, X_{i-1})$

$H(X_1, \dots, X_n) = \sum_i H(X_i|X_1, \dots, X_{i-1}) \leq \sum_i H(X_i)$ $I(X;Y) \geq I(X;f(Y))$ $\forall f$

$H(f(X)) \leq H(X) \Leftrightarrow f$ bijective $I(X;Y) | Z \Rightarrow I(X|Z) \geq I(X;Y)$ data processing

$X, Y \text{ iid} \Rightarrow I(X=Y) \geq 2 - H(X)$ $H(X)$ concave w.r.t. p_X

Divergence: $D(p||q) \geq 0 \Leftrightarrow p=q$ $D(p_{XY} || p_{X|Y} p_Y) = D(p_{Y|X} || p_{Y|X} p_X) + D(p_X || p_Y)$

$D(p_{XY} || p_{X|Y} p_Y) > D(p_X || p_X) \geq D(p_{Y|X} || p_{Y|X} p_X) = D(p_X p_{Y|X} || p_X p_{Y|X})$

Convexity: $D(\lambda p_1 + (1-\lambda) p_2 || \lambda p_1 + (1-\lambda) p_2) \leq \lambda D(p_1 || p_1) + (1-\lambda) D(p_2 || p_2)$

Typical sets: $T_n^\epsilon := \{(x_1, \dots, x_n) \in X^n : |\sum_i p_{x_1, \dots, x_n}(x_1, \dots, x_n) - H(X)| \leq \epsilon\} = \frac{1}{n} \log p_{x_1, \dots, x_n}(x_1, \dots, x_n) \rightarrow H(X)$

$\forall \epsilon > 0 \exists N \forall n \geq N$: Weak AEP1 in prob

$\forall \bar{x} \in T_n^\epsilon : 2^{-n(H(X)+\epsilon)} \leq p_{\bar{x}} \leq 2^{-n(H(X)-\epsilon)}$ + same for $S_n^\epsilon := \text{smallest set w/ } \mathbb{P}(X \in S_n^\epsilon) \geq 1 - \epsilon$ prob: info in ϵ

$2) (1-\epsilon) 2^{n(H(X)-\epsilon)} \leq |T_n^\epsilon| \leq 2^{n(H(X)+\epsilon)}$ except $(1-2\epsilon)$ is $\leq |S_n^\epsilon|$. $\mathbb{E}[-\log p_{\bar{x}}] = H$, WLN

3) $\mathbb{P}(X \in T_n^\epsilon) \geq 1 - \epsilon$. **Codes**

$c: X \rightarrow Y$ is a code, is d-ary if $|Y| = d$. $c(x) = c(x_\infty)$

unambiguous: c is injective, unambig. decodable: $c: X \rightarrow Y$ by $c(x_1, \dots, x_n) \mapsto c(x_1) \dots c(x_n)$ is injective

prefix: no codeword $c(x)$ is a prefix of another codeword $c(y)$. $\forall \epsilon > 0 \exists n : X^n \rightarrow S_n^\epsilon$ it is optimal for X^n : $\text{uniqu. dec.} + \text{other uniq. dec.} \Leftrightarrow \mathbb{E}[I(c(X))] \leq \mathbb{E}[I(c(X))]$.

$H_d(X) \leq \mathbb{E}[I(c(X))]$ for c unambig. + d-ary. Kraft-McMillan: $\sum_x |c(x)| \leq 1$. if $(c(x))_{x \in X} \text{ s.t. } \sum_x |c(x)| \leq 1$, then \exists a prefix code c w/ length $L(x)$.

$\forall X: X \rightarrow Y \exists$ opt code c^* st. $H_d(X) \leq \mathbb{E}[I(c(X))] \leq H_d(X) + 1$. Shannon's code: sort probs $p_1 \geq \dots \geq p_n$, $L(x) = -\log(p_x)$

Elias: same setup as Shannon, except $n_r := \sum_{i < r} p_i + p_{r+1}, \ell_r := 1$ \leftarrow not optimal $c_s(x) = \text{shift } \ell_r \text{ digits of } n_r$ $\leftarrow 0(\ell_r \log k)$ to sort d-ary expansion of n_r $L = |Y|$

Huffman: sort $p_1 \geq \dots \geq p_n$, give each a node labelled by prob $\rightarrow H_d(X) + D(p||p^*) \leq \mathbb{E}[I(c(X))] \leq H_d(X) + D(p||p^*) + 1$

- optimal: repeat: join 2 least likely nodes to ℓ_r red list \leftarrow code based on c .

red of code: code is path down tree.

optimality proof: 1) a canonical code always exists

a canonical code: 2) induction on $|X|$: 1 for 2 - optimal for 2^k : given $p_i, p_i' = v/2$ smallest merged $- p_i > p_i' \Rightarrow |c(x_2)| \leq |c(x_2')|$ c^P : canonical form $c^P(c^P) = \text{merged}$ - 2 (optimal) codewords have same len e^P : opt for $p_i, e^P = \text{unmerged}$ \leftarrow $\forall i$ differ only in last digit \leftarrow prove $c^P = e^P$

\Rightarrow expansion gives opt code \Rightarrow Huffman opt as based on expansion

Channel codes or $(M-n)$ code for a DMC (X, M, Y) is (c, d) : $c: S_1 \dots S_m \rightarrow X^n$ encoder $d: Y^n \rightarrow S_1 \dots S_m$ decoder

ECC: i = code book, message set, $c(i)$ = codeword for i . rate $r(c, d) := \frac{1}{n} (\log |X| - C_n)$

$\epsilon_i := \mathbb{P}(d(Y) \neq i | c(i) = \bar{x})$ = error on message i . $\epsilon_{\max} := \max_i \epsilon_i$, $\bar{\epsilon} := \sum \epsilon_i / m$.

a rate $R > 0$ is achievable if $\forall \epsilon > 0 \exists m, n$ a (m, n) code (c, d) with $r(c, d) > R - \epsilon$ and $\epsilon_{\max} < \epsilon$.

DMC: (X, M, Y) realizes a DMC if $M = (p_{Y|X}(y|x))_{x,y}$ in stochastic matrix form, possibly empty

Lossless if Y splits into disjoint Y_i for $i \in \{1, \dots, n\}$ st. $\forall i \in \{1, \dots, n\} \mathbb{P}(Y_i | X=x) = 1 \Rightarrow H(X|Y) = 0$ vs. unless $X \rightarrow Y$ onto info $X \rightarrow H(X|Y) = H(X)$

channel capacity: $C := \sup_{p \in \Delta(X)} I(X;Y) = \sup_{p \in \Delta(X)} (H(X) - H(X|Y)) = \sup_{p \in \Delta(X)} (H(Y) - H(Y|X))$

$0 \leq C \leq \min_{p \in \Delta(X)} H(Y|X)$ calculating C : code $H(X|Y)$ w.r.t. $\sum_{x,y} p_{X,Y}(x,y) \mathbb{P}(X=x)$ then want to use $H(Y) \leq H(\text{unif}) = 1$: find Y in form of X 's dist. find dist of X to make Y uniform. Do everything in terms of X 's dist.

Shannon 2 (noisy channel coding) for (X, M, Y) with capacity $C, R > 0$ is achievable $\Leftrightarrow R \leq C$

Joint AEP: $J_n^\epsilon := \{(\bar{x}, \bar{y}) \in X^n \times Y^n : \max_i \left| \frac{-\log p_{Y|X}(y_i|x_i)}{\log p_X(x_i)} - H(X) \right| \leq \epsilon\}$

- $\lim_{n \rightarrow \infty} \mathbb{P}(X, Y \in J_n^\epsilon) = 1$
- $|J_n^\epsilon| \leq 2^{n(H(X) + \epsilon)}$
- $\left| \frac{-\log p_{Y|X}(y_i|x_i)}{\log p_X(x_i)} - H(X) \right| \leq \epsilon \Leftrightarrow -\log p_{Y|X}(y_i|x_i) \leq H(X) + \epsilon \log p_X(x_i)$
- $X', Y' \text{ same marginals as } X, Y \text{ (but not indep.)} \Rightarrow \exists n_0, \forall n \geq n_0 \mathbb{P}(X, Y \in J_n^\epsilon) \leq 2^{-n(I(X;Y) - \epsilon)}$

1 Finding good channel codes: if n fixed, can try all poss codes: then are $\sim |X|^m$ info code & for rate close to C need $m \sim |X|^n$.

2 randomly generate them - looks like codewords, 2^{nK+1} , so exponential storage & decoding time.

3 non-IID input: disc. stochastic proc: seq of RVs, stationary if $P(X_1=x_1 \dots X_n=x_n) = P(X_{1+j}=x_{1+j} \dots X_{n+j}=x_{n+j})$ for all j

4 entropy rate: $H(X) = \lim_{n \rightarrow \infty} H(X_1, \dots, X_n) = H(X_1) + \text{info}$

5 for a stationary process, $H(X)$ exists = $\lim_{n \rightarrow \infty} H(X_n | X_{n-1}, \dots, X_1)$. If $n \rightarrow \infty$ $H(X_n | X_{n-1}, \dots, X_1)$ is (weakly) decr & limit exists.

6 entropy rate of Markov chain: $H(X) = \lim_{n \rightarrow \infty} H(X_n | X_{n-1}, \dots, X_1) = \sum_i p(x_i) I(x_i)$. Cesaro mean: $a_n \rightarrow a \Rightarrow \frac{1}{n} \sum a_i \rightarrow a$

7 X stationary markov chain, $Y = f(X)$ function of M . chain. note Y not Markov, but stationary

8 $H(Y | Y_{n-1}, \dots, Y_1) \leq H(Y) \leq H(Y_n | Y_{n-1}, \dots, Y_1)$ & equality all the same. \leftarrow we $I(X_1; Y_n | Y_{n-1}, \dots, Y_1) \leq H(X) \leftarrow \sum_{i=1}^n I(X_i; Y_i | Y_{i-1}, \dots, Y_1)$

9 extend backwords to X_0, Y_0 stationary \leftarrow is it converging

10 Proofs - Gibbs: Jensen: f convex $\Rightarrow E_{\pi} f(a) \geq f(E_{\pi} a) + b(a-x)$ for some $b \geq 0 \Rightarrow f(D) \geq f(E[X]) + b(X - E[X])$, take E of both sides

11 Gibbs combine to $\sum p_i = E[\log \frac{p_i}{E}] \geq -\log E[\frac{p_i}{E}] = -\log \sum p_i = 0$. Jensen w/ log convex. \leftarrow see Jensen [discrete]

12 Divergence: 2) chain rule: split prob in logs, take out terms, expand, etc. \leftarrow log sum over i w/ $a_i = p_i, a_2 = (1-p_1)p_2, a_3 = (1-p_1)(1-p_2)p_3$

13 mutual info: $D(p||q) = \sum_i p_i \log \frac{p_i}{q_i}$ for $H(X) = I(Y; X) = H(X) - H(X|Y)$. chain rule (3) on Y or Z \leftarrow $q = f(Y)$ then $(X|Z) \sim Z$ known.

14 entropy $I(X): \sum p_i \log p_i \leq \sum p_i \log q_i = I(X|Y) \leq H(X) \leftarrow$ info all the same. \leftarrow we $I(X_1; Y_n | Y_{n-1}, \dots, Y_1) \leq H(X) \leftarrow \sum_{i=1}^n I(X_i; Y_i | Y_{i-1}, \dots, Y_1)$

15 Fano: $Z = X \oplus Y \Rightarrow H(Z) = H(X) + H(Z|X) = H(X) + H(Y|X) \leq H(X) + H(Y) \leftarrow$ split by pos val of Z . $Z=1 \Rightarrow \leq |X|-1$ pos X vals

16 Work AEP: conv: info + MLLN $I(X_1, \dots, X_n) = \sum p_i I(X_i)$. Conv: $p_i = \frac{1}{n} \sum p_i(x_i)$ $\Rightarrow \sum p_i(x_i) \geq \sum p_i(x_i) \geq \sum p_i(x_i) \geq \sum p_i(x_i)$. Conv: $1 - \epsilon \leq H(X) \leq \sum p_i(x_i) \leq 1$

17 Shannon 1: Let ϵ : $c \in \mathbb{R}_+$ s.t. $c \cdot \log(\frac{1}{c}) \leq 1$. Then get $c \cdot \log(\frac{1}{c}) \leq n(H(X) + \epsilon)$. \leftarrow choose const c s.t. $c \cdot \log(\frac{1}{c}) \leq \epsilon$

18 unif dec $\Rightarrow \sum p_i(x_i) \cdot \log(\frac{1}{p_i}) \leq n(H(X) + \epsilon)$. \leftarrow # of words in X which are const of codewords. \leftarrow $\sum p_i(x_i) \cdot \log(\frac{1}{p_i}) \leq n(H(X) + \epsilon) \Rightarrow \sum p_i(x_i) \leq n(H(X) + \epsilon) \Rightarrow \sum p_i(x_i) \leq n(H(X)) + n\epsilon \leq \frac{n}{2}$

19 given ϵ : $\epsilon = \sum p_i(x_i) \leq 1$, $c(m) = \min \{n \text{ digits in } d \text{ s.t. } \sum p_i(x_i) \leq 1\}$ \leftarrow is a prefix code $\leftarrow X \Rightarrow r_m < d^{m-1}$ but $X \in \{r_1, \dots, r_m\} \Rightarrow d^{m-1}$

20 $E[I(X|Z)] \geq H(X)$: unif prob $q = \frac{1}{d} \leq \frac{1}{2}$, sub into $E[I(X|Z)]$, move q across to M term to make $D(p||q)$, get $\sum p_i \log \frac{p_i}{q} \geq \sum p_i \log \frac{p_i}{q} \geq 0$.

21 opt code exists: $\epsilon_c := -\log D(p||q) \leq \frac{1}{2}$. K-M pref code exists, $E[I(X|Z)]$ right: countably many pref codes with $E[I(X|Z)] < \text{fin. num.}$, so sort and pick min.

22 Shannon code bound: $R \leq C$. \leftarrow Huffman code: canonical code exists: c const $c = w/2$ words x_i, z_n simple, check $p_i < p_j \Rightarrow |c(x_i)| - |c(x_j)| \geq 1 \Rightarrow 0 \leq E[I(X|Z)] \leq C$

23 Shannon 2: take prod $p, m, n \in \mathbb{N}$, T^n sizeable typical set of $P_{XY} = P_X P_Y$ \leftarrow just choose

24 1) generate m random codewords from X^n by randomly sampling P_X n times \leftarrow the channel

2) randomly assign each message a codeword

3) the decoder: a typical-set decoder: on \tilde{Y} , if $\exists ! \tilde{x}$ st. $(\tilde{x}, \tilde{y}) \in T^n$ return \tilde{x} 's message from 2). else return m .

4) call this (m,n) -code. \leftarrow it is RANDOM!!

5) sample from (C, D) \leftarrow $R, m = 2^{n(R-2\epsilon/3)+1}$: (C, D) has rate $R - \frac{2}{3}\epsilon + \frac{1}{n}$

6) sample a message $W \sim \text{Unif}\{1, \dots, m\}$ chose nst $I(W \neq \tilde{W}) < \epsilon/2$, by (lemma): W, \tilde{W} is diff, $I(W \neq \tilde{W}) \geq 0 \text{ as } n \rightarrow \infty$: no prob

7) send $\tilde{X} = C(W)$ on the channel

8) decode to \tilde{W} w/ D .

9) con [only] when $R \leq C$: \leftarrow $\tilde{\epsilon} = \frac{1}{m} \sum \epsilon_i < \frac{1}{2}$. order ϵ_i 's least half have $\epsilon_i < \epsilon$ \leftarrow there always not: now a $(\frac{m}{2}, n)$ code w/ $R - \frac{2}{3}\epsilon > R - \epsilon$, $\epsilon_{\text{min}} \leq \epsilon$.

10) entropy rate exists: \leftarrow non-incr: $H(X_n | X_{n-1}, \dots, X_1) \leq H(X_n | X_{n-1}, \dots, X_2) \leq \dots \leq H(X_n | X_{n-1}, \dots, X_1)$ by stationarity

11) limit exists is $n \rightarrow \infty$ a bit below $\epsilon \geq 0$.

12) Cesaro mean: $\frac{1}{n} \sum \epsilon_i \leq \frac{1}{n} \sum \epsilon_i \leq \frac{1}{n} \sum \epsilon_i + \frac{n-\epsilon}{n} \epsilon \leq \epsilon \Rightarrow \frac{1}{n} \sum \epsilon_i \leq \epsilon$ as $n \rightarrow \infty$

13) $H(X_1, \dots, X_n) = \frac{1}{n} \sum H(X_i | X_{i-1}, \dots, X_1)$

14) work exps converge, + Cesaro mean so conv to lim of code exps

15) Distributions

16) Uniform $\{1, \dots, n\}: f(x) = \frac{1}{n}, \mu = \frac{n+1}{2}, \sigma^2 = \frac{n^2-1}{12}$

17) $Ber(p): P(X=1) = p, P(X=0) = 1-p, \mu = p, \sigma^2 = p(1-p)$

18) $Bin(n, p): P(X=k) = \binom{n}{k} p^k (1-p)^{n-k}, \mu = np, \sigma^2 = np(1-p)$

19) Geometric(p): $P(X=k) = (1-p)^{k-1} p, \mu = \frac{1}{p}, \sigma^2 = \frac{1-p}{p^2}$, for both

20) Normal (μ, σ^2): $f(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp(-\frac{(x-\mu)^2}{2\sigma^2})$

21) Poisson (λ): $f(k) = \lambda^k e^{-\lambda} / k!, \mu = \lambda, \sigma^2 = \lambda$. discrete

22) Exp (λ): $f(x) = \lambda e^{-\lambda x}, P(X > x) = e^{-\lambda x}, \mu = \lambda, \sigma^2 = \lambda$

23) Gamma (r, λ): $f(x) = \frac{1}{\Gamma(r)} \lambda^r x^{r-1} e^{-\lambda x} [x^r e^{-\lambda x} \ln \Gamma(r)] \sim \text{Gamma}(r, \lambda)$

24) Gamma func: In general complicated, $\Gamma(n) = (n-1)!$, $\Gamma(1) = 1$

25) Taylor series: \leftarrow approximation: \leftarrow Simple notes

26) $\chi_{(0, \infty)} = \sum_{n=0}^{\infty} x^n$

27) $\Lambda(x, \lambda, \mu) = f(x) + \lambda(x - \mu) + \mu(h_i - h_i)$

28) $e^x = \sum_{n=0}^{\infty} x^n / n!$

29) $\ln(1+x) = \sum_{n=0}^{\infty} (-1)^n x^n / n$

30) $\chi_{(0, \infty)} = \sum_{n=0}^{\infty} x^n / n!$

31) \leftarrow set to 0. suff suff: $e^x \rightarrow 0$

32) maximizing entropy: if have target dist p^* , w/ it: $\max_{p \in \text{class}, 2^{nK+1}} I(p || p^*)$

33) $H(Y) = -\sum p_i \log p_i = -D(p^* || p) \leftarrow$ $\sum p_i \log p_i^* = \dots$ w/ def of p_i^*

34) also bound: e.g.: $E[Y] \leq \alpha \leftarrow$ $E[Y] = \sum x_i p_i$

35) BSG (capacity): $[\text{con. quota}] = 1 - H(Y|X) \leftarrow$ $\alpha \in [\text{sat. quota}]$, concave $H(C, \alpha) = H(X)$