

\log DENOTES LOGARITHM TO THE BASE 2

\ln DENOTES LOGARITHM TO THE BASE e

ELEMENTARY STATEMENT

$$\log x \leq \frac{x-1}{\ln 2} \quad * \text{ FOR ALL } x > 0, \text{ EQUALITY}$$

HOLDS IFF $x=1$.

PROOF

WE CONSIDER $f(x) = \log x - \frac{x-1}{\ln 2}$

$$f'(x) = \frac{1}{x \ln 2} - \frac{1}{\ln 2} \quad : \text{ IT EQUALS 0 AT } x=1$$

$$f''(x) = -\frac{1}{x^2 \ln 2} < 0 \quad : f(x) \text{ IS STRICTLY CONCAVE HENCE } x=1 \text{ IS THE UNIQUE GLOBAL MAXIMUM}$$

LEMMA

(LOG-10M INEQUALITY)

FOR ARBITRARY NONNEGATIVE NUMBERS

$a_i, b_i, i=1, \dots, n$ WE HAVE

$$\sum_{i=1}^n a_i \log \frac{a_i}{b_i} \geq a \cdot \log \frac{a}{b}$$

WHERE $a = \sum_{i=1}^n a_i, b = \sum_{i=1}^n b_i$

$a_i b = b_i a \text{ FOR } i=1, \dots, n$

THE EQUALITY HOLDS IFF

~~$a_i = c b_i \text{ FOR } i=1, \dots, n$~~

(HERE $a \cdot \log \frac{a}{b}$ IS DEFINED TO BE 0 IF $a=0$ AND $+\infty$ IF $a > 0 = 0$: CONCLUSION BASED ON CONTINUITY)

PROOF

• WE MAY ASSUME THAT THE a_i 'S ARE POSITIVE SINCE DELETING THE PAIRS (a_i, b_i) WITH $a_i = 0$ (IF ANY), THE LEFT-HAND SIDE REMAINS UNCHANGED WHILE THE RIGHT-HAND SIDE DOES NOT DECREASE

• WE CAN ALSO ASSUME THAT THE b_i 'S ARE ALSO POSITIVE (OTHERWISE THE LEFT-HAND SIDE IS TRIVIAL)

• FURTHER, IT SUFFICES TO PROVE THE LEMMA FOR $a=b$, SINCE MULTIPLYING THE b_i 'S BY A CONSTANT DOES NOT AFFECT THE INEQUALITY

THE CASE IDENTIFYING $X = \frac{b_i}{a_i}$ IN (X) WE GET THAT

$$-\sum_{i=1}^n a_i \log \frac{b_i}{a_i} \geq -\sum_{i=1}^n a_i \frac{\left(\frac{b_i}{a_i} - 1\right)}{\ln 2} = \frac{-\sum_{i=1}^n b_i + \sum_{i=1}^n a_i}{\ln 2} = 0$$

$A = \{a_1, a_2, \dots, a_n\}$: FINITE SET OF CARDINALITY $|A|$
 LET $P = \{P(a), a \in A\}$ AND $Q = \{Q(a), a \in A\}$ BE PROB. DISTRIBUTIONS ON THE SET A . THEY HAVE BACH-WEISE DISTANCE (INFORMATION DIVERGENCE, J-DIVERGENCE, RELATIVE ENTROPY) IT IS DEFINED BY

$$D(P||Q) = \sum_{a \in A} P(a) \log \frac{P(a)}{Q(a)} \quad (\text{DUE TO THE CONVENTION IT IS FINITE IFF } P \ll Q, \text{ i.e., IF } Q(a) = 0 \text{ IMPLIES } P(a) = 0)$$

HOW MUCH WE CAN DISTINGUISH P AND Q WITH STATISTICAL TEST

LOG-LIKELIHOOD INEQUALITY IMPLIES THAT $D(P||Q) \geq 0$ WITH EQUALITY IFF $P=Q$.

(CONVERGENCE OF PROB DISTRIBUTIONS) $P_n \rightarrow P$ MEANS POINTWISE CONVERGENCE, THAT IS, $P_n(a) \rightarrow P(a)$ FOR EACH $a \in A$. TOPOLOGICAL CONCEPTS FOR PROB. DISTRIBUTIONS (CONTINUITY, OPEN AND CLOSED SETS, ETC.) ARE MEANT FOR THE TOPOLOGY INDUCED BY POINTWISE CONVERGENCE

KL DIVERGENCE IS LOWER SEMI CONTINUOUS IN THE PAIR (P, Q) , I.E., $(P_n, Q_n) \xrightarrow{n \rightarrow \infty} (P, Q)$ IMPLIES THAT $\liminf_{n \rightarrow \infty} D(P_n || Q_n) \geq D(P || Q)$. IT IS CONTINUOUS AT EACH (P, Q) WITH STRICTLY POSITIVE Q .

DEFINITION $D(P||Q)$ IS A CONVEX FUNCTION OF THE PAIR (P, Q) PROVED IT IS STRICTLY CONVEX IN P AT $P=Q$ ALSO

PROOF SUPPOSE THAT $P = \alpha \cdot P_1 + (1-\alpha)P_2$, $Q = \alpha \cdot Q_1 + (1-\alpha)Q_2$, I.E., $P(a) = \alpha \cdot P_1(a) + (1-\alpha)P_2(a)$ FOR EVERY $a \in A$ AND SIMILARLY FOR Q , WHERE $0 < \alpha < 1$.

the log-sum inequality implies that for all $a \in X$

$$\alpha \cdot p_1(a) \cdot \log \frac{p_1(a)}{q_1(a)} + (1-\alpha) \cdot p_2(a) \cdot \log \frac{p_2(a)}{q_2(a)} =$$

$$= \alpha \cdot p_1(a) \cdot \log \frac{\alpha p_1(a)}{\alpha q_1(a)} + (1-\alpha) \cdot p_2(a) \cdot \log \frac{(1-\alpha) p_2(a)}{(1-\alpha) q_2(a)} \geq p(a) \cdot \log \frac{p(a)}{q(a)}$$

Similarly for $a \in X$ gives a bound $D(p_1 || q_1) + (1-\alpha) \cdot D(p_2 || q_2) \geq D(p || q)$

LEMMA

The above lemma states also that $D(p || q)$ is strictly convex in p if $q(a) > 0$ for all a .

Variational distance (one avoids use the term for the sake of this):

$$|p - q| = \sum_a |p(a) - q(a)|$$

Pinsker inequality (also called Wilfong-Guennemann-Wulfsberg-Pinsker inequality):

$$D(p || q) \geq \frac{1}{2 \ln 2} |p - q|^2$$

Non uniform convexity: if $D(p || q) < \infty$ then

$$D(p || q) \leq \frac{|p - q|^2}{\alpha_0}$$

where $\alpha_0 = \min_{a: q(a) \neq 0} q(a)$

NUMERICALS WITH APPROXIMATE HINTS

THE ENTROPY OF A REPR. DIST $p = \{ p(a), a \in X \}$ IS DEFINED BY

$$H(p) = - \sum_{a \in X} p(a) \log p(a)$$

$H(p) \geq 0$ BY DEFINITION

IF $q_1 = q_2 = q$ \Leftrightarrow $p_1 = p_2 = p$

$$p_1(a) = \alpha \cdot p_1(a) + (1-\alpha) p_2(a)$$

$$p_2(a) = \alpha \cdot p_1(a) + (1-\alpha) p_2(a)$$

$$\Leftrightarrow p_1(a) = p_2(a)$$

(*) IS AN CONVEXITY IF $\forall a \in X$

$$\alpha \cdot p_1(a) \cdot (\alpha \cdot q_1(a) + (1-\alpha) q_2(a)) = \alpha \cdot q_1(a) (\alpha \cdot p_1(a) + (1-\alpha) p_2(a))$$

$$\text{AND } (1-\alpha) p_2(a) \cdot (-u) = (1-\alpha) (q_2(a)) \cdot (-u)$$

$$H(\text{P II UNIFORM MIT ON } A) = \sum_{a \in A} p(a) \cdot \log \frac{p(a)}{\frac{1}{|A|}} =$$

$$= \sum_{a \in A} p(a) \cdot \log p(a) + \sum_{a \in A} p(a) \log |A| = \log |A| - H(p)$$

HENCE $H(p) = \log |A| - H(\text{P II UNIFORM MIT ON } A)$

- $H(p) \leq \log |A|$, = holds iff p is uniform dist
- $H(p)$: CONTINUOUS AND STRICTLY CONCAVE IN P

LET X BE ~~THE~~ RV TAKING VALUES IN A

$$H(X) = H(p_X) = E(-\log p_X(X))$$

MEASURE OF INFORMATION CONTENT OR MEASURE OF AVERAGE UNCERTAINTY ABOUT THE OUTCOME OF X

SEVERAL AXIOMATIC AND OPERATIONAL JUSTIFICATION EXIST FOR $H(X)$

A^∞ DENOTES THE SET OF ALL INFINITE SEQUENCE $x = x_1^\infty$ WITH $x_i \in A, i \geq 1$.

A^* DENOTES THE SET OF ALL FINITE SEQUENCE DRAWN FROM A (THE EMPTY SET IS ALSO INCLUDED)

THE CONCATENATION OF $u \in A^*$ AND $v \in A^* \cup A^\infty$ IS DENOTED BY uv . A FINITE SEQUENCE u IS A PREFIX OF A FINITE OR INFINITE SEQUENCE w , AND WE WRITE $u \prec w$ IF $w = uv$ FOR SOME v

DATA COMPRESSION

EXAMPLE LET X BE A RV TAKING VALUES IN $A = \{1, 2, 3, 4\}$ HAVING THE FOLLOWING DISTRIBUTION:

$P(X=1) = \frac{1}{2}, P(X=2) = \frac{1}{4}, P(X=3) = \frac{1}{8}, P(X=4) = \frac{1}{8}$

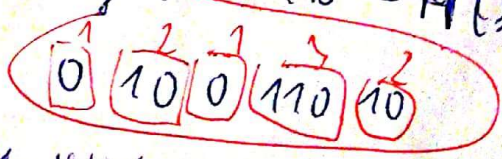
WE WOULD LIKE TO COMPRESS EFFICIENTLY ITS OUTCOME

IDEA ASSIGNING SHORT DESCRIPTIONS TO THE MOST FREQUENT OUTCOMES AND NECESSARILY LONGER DESCRIPTIONS TO THE LESS FREQUENT SYMBOLS
NO CODEWORD IS A PREFIX OF ANOTHER ONE

$C(1) = 0, C(2) = 10, C(3) = 110, C(4) = 111$

THEN $E(L) = \frac{1}{2} \cdot 1 + \frac{1}{4} \cdot 2 + \frac{1}{8} \cdot 3 + \frac{1}{8} \cdot 3 = 1.75 = H(X)$

ADVANTAGE OF THE PREFIX CODE: IF WE USE THE CODE REPRESENTATIVELY WE CAN DECODE THE SYMBOLS UNIVOCALLY AND INSTANTANEOUSLY (THE SYMBOL a CAN BE DECODED AS SOON AS WE COME TO THE END OF THE CODEWORD CORRESPONDING TO IT)



A CODE FOR SYMBOLS IN A WITH IMAGE ALPHABET B IS A MAPPING $C: A \rightarrow B^*$, (ii) LENGTH FUNCTION $L: A \rightarrow \mathbb{N}$ IS DEFINED BY THE FORMULA $C(a) = b_1^{L(a)}$

IN THIS CASE WE ASSUME: $B = \{0, 1\}$, THE CODEWORDS ARE DISTINCT AND DIFFERENT FROM ~~THE~~ EMPTY STRING.

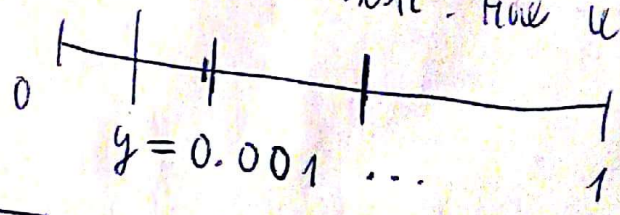
C IS PREFIX IF $C(a) \not\subset C(\tilde{a})$ NEVER HOLDS FOR $a \neq \tilde{a}$ IN A

THE DYADIC REPRESENTATIONS FOR NUMBERS $y \in [0, 1)$ ARE OF THE FORM $y = \sum_{k=1}^{\infty} y_k 2^{-k}$ WITH $y_k \in \{0, 1\}$.

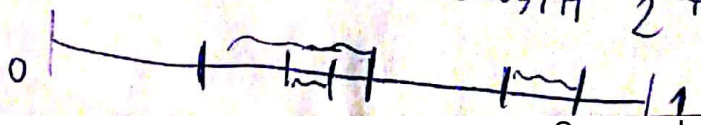
THE NUMBERS OF FORM $y = \frac{m}{2^n}$ WHERE n AND m ARE INTEGERS WITH $n \geq 1, 0 \leq m < 2^n$, ARE CALLED DYADIC RATIONALS.

THEY HAVE TWO DYADIC REPRESENTATIONS: ONE WITH INFINITE 0-S AND ONE WITH INFINITE 1-S (E.G. $\frac{1}{4} = 0.01 = 0.001\dots$). WE WILL USE THE REPRESENTATION WITH THE INFINITE 0-S. ALL OTHER NUMBERS IN $[0, 1)$ HAVE A UNIQUE REPRESENTATION.

NOTE THAT y IS NOT A DYADIC RATIONAL. HOW WE CAN FIND ITS REPRESENTATION



WHILE ALL THE NUMBERS WHOSE DYADIC REPRESENTATION START WITH 0110 : IN THE INTERVAL OF LENGTH $2^{-4} = \frac{1}{16}$ FOLLOWING



LEMMA (KRAFT INEQUALITY)

A FUNCTION $L: A \rightarrow \mathbb{N}$ IS THE LENGTH FUNCTION OF SOME PREFIX CODE IFF IT SATISFIES THE SO CALLED KRAFT INEQUALITY:

$$\sum_{a \in A} 2^{-L(a)} \leq 1$$

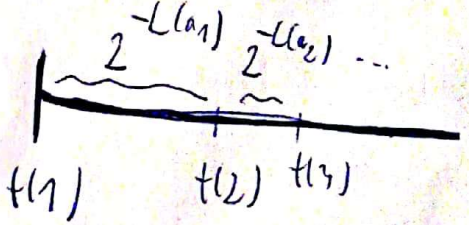
PROOF

ASSUME THAT $(: A \rightarrow B^*$ IS A PREFIX CODE. ASSOCIATE WITH EACH $a \in A$ THE NUMBER $f(a)$ WHOSE BINARY IS THE CONCORD $f(a) = 0.b_1^{L(a)}$ THAT IS $f(a) = 0.b_1 \dots b_{L(a)}$

THE PREFIX CONDITION IMPLIES THAT $f(\tilde{a}) \notin [f(a), f(a) + 2^{-L(a)})$ IF $\tilde{a} \neq a$. HENCE THE INTERVALS $[f(a), f(a) + 2^{-L(a)})$, $a \in A$, ARE DISJOINT. THE TOTAL LENGTH OF THESE DISJOINT INTERVALS ARE LESS THAN 1 $\Rightarrow \sum_{a \in A} 2^{-L(a)} \leq 1$.

CONVERSELY, SUPPOSE A FUNCTION $L: A \rightarrow \mathbb{N}$ SATISFIES $\sum_{a \in A} 2^{-L(a)} \leq 1$. LABEL A SO THAT $L(a_i) \leq L(a_{i+1})$, $i = 1, \dots, k-1$ (OR EQUIVALENTLY $2^{-L(a_i)} \geq 2^{-L(a_{i+1})}$)

$$\text{LET } f(i) = \sum_{j < i} 2^{-L(a_j)}$$



$f(i)$ CAN BE BINARYLY REPRESENTED AS $f(i) = 0.b_1 \dots b_{L(a_i)}$ THEN $(: (a_i) = 0.b_1^{L(a_i)}$ DEFINES A PREFIX CODE WITH LENGTH FUNCTION L .

!!! PREFIXED LENGTHS CAN VARY !!!

KEY CONSEQUENCE OF THE LEMMA IS SHANNON'S NOISELESS CODING THEOREM VIII

THEOREM LET P BE A PROB DIST ON A . THEN EACH PREFIX CODE HAS EXPECTED LENGTH $E(L) = \sum_{a \in A} P(a) L(a) \geq H(P)$. FURTHERMORE, THERE IS A PREFIX CODE WITH LENGTH FUNCTION $L(a) = \lceil -\log P(a) \rceil$ (AND HENCE ITS EXPECTED LENGTH SATISFIES $E(L) < H(P) + 1$)

PROOF OF THE FIRST ASSERTION

$$E(L) - H(P) = - \sum_{a \in A} P(a) \log_2 2^{-L(a)} + \sum_{a \in A} P(a) \log_2 P(a) = \sum_{a \in A} P(a) \log_2 \frac{P(a)}{2^{-L(a)}} \geq \sum_{a \in A} P(a) \log_2 \frac{1}{\sum_{a \in A} 2^{-L(a)}} \geq 0$$

PROOF OF THE SECOND ASSERTION

$$\sum_{a \in A} 2^{-\lceil -\log P(a) \rceil} \leq \sum_{a \in A} 2^{-\log P(a)} = \sum_{a \in A} P(a) \leq 1$$

$L(a) = \lceil -\log P(a) \rceil$ SATISFIES THE WEAK INEQUALITY

REMARK 1 THE ABOVE CONSTRUCTION ACHIEVING $E(L) < H(P) + 1$ IS NOT PRACTICAL WHEN $|A|$ IS LARGE, SINCE IT RELIES ON THE CONSTRUCTION VIEW IN THE PROOF OF THE ~~LEMMA~~ LEMMA ABOUT THE WEAK INEQUALITY WHERE THE ORDERING OF THE LENGTHS ~~IS~~ IS NECESSARY

REMARK 2 IT WILL BE IMPORTANT LATER TO DISTINGUISH THE DISTRIBUTION P ACCORDING TO WHICH WE CALCULATE EXPECTED LENGTH AND THE LOSSY DISTRIBUTION DENOTED BY Q

SINCE THAT $P \ll Q$. WE KNOW THAT THERE EXISTS PREFIX CODE WITH $L(a) = \lceil -\log Q(a) \rceil$. FOR THIS CODE

$$H(P) + D(P||Q) \leq E_p(L) \leq 1 + H(P) + D(P||Q)$$

SINCE $\sum_{a \in A} p(a) (-\log Q(a)) = \sum_{a \in A} p(a) \frac{1}{\log Q(a)} + \sum_{a \in A} p(a) \log p(a)$
 $-\sum_{a \in A} p(a) \cdot \log p(a) = D(P||Q) + H(P)$

THE FOLLOWING THEOREM SHOWS THAT EVEN UN-PREFIX CODES CAN NOT INSTANTLY BEAT THE ENTIRELY LOWER BOUNDS

THEOREM THE LENGTH FUNCTION OF A NOT NECESSARILY PREFIX CODE $C: A \rightarrow B^*$ SATISFIES $\sum_{a \in A} 2^{-L(a)} \leq \log |A|$ AND FOR ANY PROB. DIST P ON A , THE CODE HAS EXPECTED LENGTH $E_p(L) \geq H(P) - \log \log |A|$

~~Handwritten scribbles~~

X RV TAKING VALUES IN A

2020

(1)

$$H(X) = H(P_X) = - \sum_{a \in A} P_X(a) \log P_X(a) = \underline{E(-\log P_X(X))}$$

MEASURE OF INFORMATION CONTENT, MEASURE OF AVERAGE UNCERTAINTY ABOUT THE OUTCOME OF X

LET (X, Y) BE A 2-DIM R.V. WITH VALUES IN $A \times B$

THEN $H(X, Y)$ IS ALSO WELL DEFINED

$$H(X, Y) \triangleq - \sum_{a \in A} \sum_{b \in B} P_{XY}(a, b) \log P_{XY}(a, b) = \underline{E(-\log P_{XY}(X, Y))}$$

COND ENTROPY: $H(Y|X) = \sum_{a \in A} P_X(a) \cdot H(Y|X=a) =$
 $= - \sum_{a \in A} P_X(a) \sum_{b \in B} P_{Y|X}(b|a) \log P_{Y|X}(b|a) =$
 $= - \sum_{a \in A} \sum_{b \in B} P_{XY}(a, b) \log P_{Y|X}(b|a) = \underline{E(-\log P_{Y|X}(Y|X))}$

CHAIN RULE: $H(X, Y) = H(X) + H(Y|X)$

MUTUAL INFORMATION OF X AND Y IS DEFINED BY

$$J(X \wedge Y) \triangleq H(X) - H(X|Y) =$$

$$= - \sum_{a \in A} \sum_{b \in B} P_{XY}(a, b) \log P_X(a) + \sum_{a \in A} \sum_{b \in B} P_{XY}(a, b) \log P_{X|Y}(a|b) =$$

$$= \sum_{a \in A} \sum_{b \in B} P_{XY}(a, b) \cdot \log \frac{P_{XY}(a, b) \cdot P_Y(b)}{P_X(a) \cdot P_Y(b)} = D(P_{XY} || P_X \times P_Y)$$

WHERE $P_X \times P_Y(a, b) = P_X(a) \cdot P_Y(b)$

IT FOLLOWS THAT $J(X \wedge Y) \geq 0, = 0$ IFF X AND Y ARE INDEPENDENT
 MORE THAN THAT $J(X \wedge Y) = H(Y) - H(Y|X) = H(X) + H(Y) - H(X, Y)$

IMPORTANT CONSEQUENCE: $H(X|Y) \leq H(X)$ (CONDITIONALY REDUCES ENTROPY)

$$J(X \wedge Y | Z) \triangleq H(X|Z) - H(X|Y, Z) \geq 0$$

$$\Rightarrow H(X|Y, Z) \leq H(X|Z)$$

$$H(X, Y | Z) = H(X|Z) + H(Y|X, Z)$$

SIMILAR CALCULATIONS AS BEFORE

AT THIS POINT WE USE THE FOLLOWING DEFINITION OF STOCH. PROCESS:
IT IS AN INDEXED SEQUENCE OF RANDOM VARIABLES DEFINED ON A COMMON PROBABILITY SPACE.

WE ASSUME THAT THE INDEX SET IS \mathbb{Z}^+ , AND THAT X_1, X_2, \dots TAKING VALUES IN A FINITE SET A

P_n DENOTES THE JOINT DIST OF $(X_1, \dots, X_n) : P_n(x_1^n) = P_{\text{rob}} \{ X_1=x_1, \dots, X_n=x_n \}$
 $x_1^n \in A^n$

NOTE THAT FOR THESE DISTRIBUTIONS THE CONSISTENCY CONDITIONS
 $P_n(x_1^n) = \sum_{a \in A} P_{n+1}(x_1^n a)$ HOLDS

$X_1, X_2, \dots, X_n, \dots$, IS STATIONARY IF $P_{\text{rob}} \{ X_1=x_1, \dots, X_n=x_n \} =$
 $= P_{\text{rob}} \{ X_{1+l}=x_1, \dots, X_{n+l}=x_n \}$ FOR EVERY n , EVERY $l \geq 1$
AND EVERY $x_1, \dots, x_n \in A$.

SIMPLEST WELL KNOWN EXAMPLE: i.i.d. SEQUENCES

THEOREM ASSUME THAT X_1, X_2, \dots IS A STATIONARY SEQUENCE THEN
BOTH $\lim_{n \rightarrow \infty} \frac{H(X_1, \dots, X_n)}{n}$ AND $\lim_{n \rightarrow \infty} H(X_n | X_1, \dots, X_{n-1})$ EXIST
AND THEY ARE EQUAL. THE LIMIT IS CALLED THE ENTROPY RATE OF
THE PROCESS, AND IT IS DENOTED BY \bar{H} .

LIMIT THEOREM (CESARO MEAN): IF $a_n \rightarrow a$ AND $b_n = \frac{1}{n} \sum_{i=1}^n a_i$
THEN $b_n \rightarrow a$

PROOF ϵ FIXED, $\exists N \in \mathbb{N}$ S.T. $|a_n - a| < \epsilon$ $\forall n \geq N$. THEN IF n IS LARGE
 $|b_n - a| = \left| \frac{1}{n} \sum_{i=1}^n (a_i - a) \right| \leq \frac{1}{n} \sum_{i=1}^{N-1} |a_i - a| + \frac{n - N + 1}{n} \epsilon \leq \frac{2\epsilon}{n} + \epsilon \leq 2\epsilon$

OF THE THEOREM

$$H(x_{n+1} | x_1, \dots, x_n) \leq H(x_{n+1} | x_2, \dots, x_n) =$$

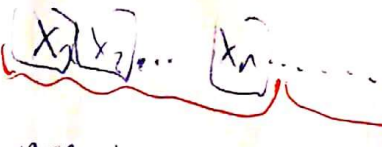
STATIONARITY

$$= H(x_n | x_1, \dots, x_{n-1}) \Rightarrow \lim_{n \rightarrow \infty} H(x_n | x_1, \dots, x_{n-1}) \text{ exists}$$

Let \bar{H} denote this limit. THEN

$$\frac{1}{n} H(x_1, \dots, x_n) = \frac{1}{n} \sum_{i=1}^n H(x_i | x_1, \dots, x_{i-1}) \xrightarrow[n \rightarrow \infty]{\text{Cesàro}} \bar{H}$$

IN REAL-WORLD APPLICATIONS WE HAVE TO COMPRESS STOCHASTIC PROCESSES (DATA SOURCE) $x_1, x_2, \dots, x_n, \dots$. WE CAN ENCODE EACH RANDOM VARIABLE SEPARATELY BUT DUE TO THE 1-bit OVERHEAD IT IS USUALLY NOT THE BEST EVEN FOR I.I.D. SOURCE



n-CODES: WE CAN REDUCE THE OVERHEAD PER SYMBOL BY SPREADING IT OUT OVER MANY SYMBOLS

IN THIS CASE THE CODE C_n IS A MAPPING FROM A^n TO B^* .

WE WILL PROVE SHANNON'S NOISELESS CODING THEOREM AND THE ABOVE CONSIDERATIONS IMPLY

THEOREM

THE MINIMUM EXPECTED CODEWORD LENGTH PER SYMBOL

SATISFIES:

$$\frac{H(x_1, \dots, x_n)}{n} \leq \frac{\min_{C_n} E(L_{C_n})}{n} \leq \frac{H(x_1, \dots, x_n)}{n} + \frac{1}{n}$$

MOREOVER, IF x_1, \dots, x_n, \dots IS A STATIONARY PROCESS THEN

$$\frac{\min_{C_n} E(L_{C_n})}{n} \xrightarrow[n \rightarrow \infty]{} \bar{H}$$

REMARK

THE GAIN OF JOINT CODING CAN BE LARGE IF THE SYMBOLS OF THE DATA SOURCE ARE HIGHLY DEPENDENT.

THE ONLY ALGORITHM PROVIDED BY SHANNON'S NOISELESS CODING THEOREM IS NOT PRACTICAL (WE HAVE TO ORDER THE ELEMENTS ASSIGNED TO THE ELEMENTS OF A^n : NOW

• IT USES THE DISTRIBUTION OF THE DATA SOURCE : LATER

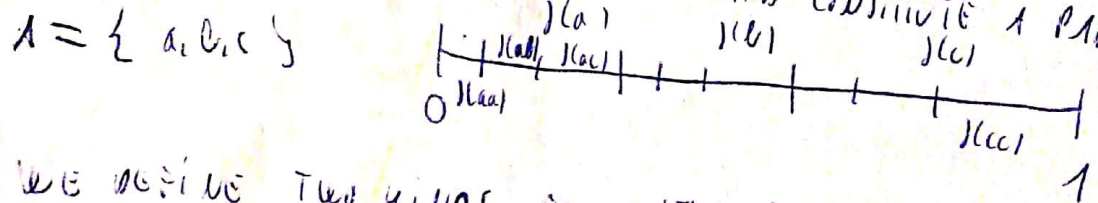
LET $Q_n, n=1,2,\dots$ BE PROBABILITY DISTRIBUTIONS ON THE SET A^n SATISFYING THE CONSISTENCY CONDITIONS $Q_n(x_1^n) = \sum_a Q_{n+1}(x_1^{n+1}a)$. AN ARITHMETIC CODE WITH CODING DISTRIBUTION SEQUENCE $\{Q_n, n \geq 1\}$ IS A SEQUENCE OF n-CODE DEFINED AS FOLLOWS.

FIX AN ORDERING ~~ON THE ELEMENTS OF A~~ ON THE ELEMENTS OF $A: a_1 < a_2 < \dots < a_{|A|}$

FOR EACH n PARTITION THE UNIT INTERVAL $[0, 1)$ INTO SUBINTERVALS $J(x_1^n) = [l(x_1^n), r(x_1^n))$ OF LENGTH $r(x_1^n) - l(x_1^n) = Q_n(x_1^n)$

IN A NESTED MANNER, I.E., SUCH THAT $J(x_1^{n+1}a_1), J(x_1^{n+1}a_2), \dots, J(x_1^{n+1}a_{|A|})$

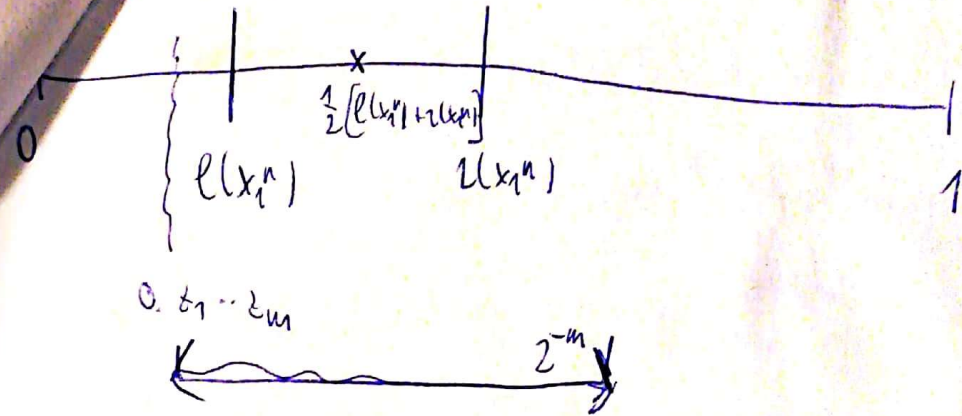
FOLLOW EACH OTHER IN THIS ORDER AND CONSTITUTE A PARTITION OF $J(x_1^n)$.



WE DEFINE TWO KINDS OF ARITHMETIC CODES :

1 $l(x_1^n) = z_1^m$ IF THE ENDPNTS OF $J(x_1^n)$ HAVE BINARY EXPANSION $l(x_1^n) = 0.z_1z_2 \dots z_m 0$, $r(x_1^n) = 0.z_1z_2 \dots z_m 1$

2 $\tilde{l}(x_1^n) = z_1^{\tilde{m}}$ IF THE MIDPOINT OF $J(x_1^n)$ HAS BINARY EXPANSION $0.z_1z_2 \dots z_{\tilde{m}}$, $\tilde{m} = \lceil -\log Q_n(x_1^n) \rceil + 1$



- $(x_1^n) \prec \tilde{(x_1^n)}$
- $L(x_1^n) < \tilde{m} = \tilde{L}(x_1^n)$
- (x_1^n) is 1-1 (THE WORDS ARE DISTINCT)
- $\tilde{(x_1^n)}$ is PREFIX

(i) UNCLE : (x_1^n) is ~~NOT~~ A PREFIX OF $(x_1^n a)$

IN ORDER TO DETERMINE CODEWORDS (x_1^n) OR $\tilde{(x_1^n)}$ THE NESTED PARTITIONS NEED NOT BE ACTUALLY COMPUTED, IT IS ENOUGH TO FIND $J(x_1^n)$. THIS CAN BE DONE IN STEPS, THE i TH STEP IS TO PARTITION THE INTERVAL $J(x_1^{i-1})$ INTO $|A|$ SUBINTERVALS OF LENGTH PROPORTIONAL TO THE CONDITIONAL PROBABILITIES

$$Q(a | x_1^{i-1}) = Q_i(x_1^{i-1} a) / Q_{i-1}(x_1^{i-1}), a \in A$$

LET x_1, x_2, \dots BE A STochastic PROCESS AND LET P_n DENOTES THE JOINT DISTRIBUTION OF (x_1, \dots, x_n) , $n \geq 1$. CONSIDER ARITHMETIC CODE WITH CODING DISTRIBUTION SEQUENCE $\{P_n, n \geq 1\}$. FOR EACH n , ITS EXPECTED CODEWORD LENGTH PER SYMBOL CAN BE UPPER BOUND BY $\frac{H(x_1, \dots, x_n) + 2}{n}$. HENCE, IT CAN BE APPROXIMATED BY A PRACTICAL METHOD WHICH ACHIEVES H ASYMPTOTICALLY.

EMPIRICAL DISTRIBUTIONS VIA TYPES

THE TYPE OF A SEQUENCE $x_1^n \in A^n$ IS ITS EMPIRICAL DISTRIBUTION

$\hat{p} = \hat{p}_{x_1^n}$, THAT IS THE DISTRIBUTION DEFINED BY

$$\hat{p}(a) = \frac{|\{i : x_i = a\}|}{n}, \quad a \in A$$

EXAMPLE

$x_1^5 = 01001, \quad \hat{p}_{x_1^5} = \left(\frac{3}{5}, \frac{2}{5} \right)$

A DISTRIBUTION p ON A IS CALLED n-TYPE IF $p = \hat{p}_{x_1^n}$ FOR SOME $x_1^n \in A^n$

\mathcal{P}_n DENOTES THE SET OF n-TYPES

$A = \{0,1\} \quad \mathcal{P}_5 = \left\{ \binom{5}{0,1}, \binom{4}{1,1}, \binom{3}{2,1}, \binom{2}{2,2}, \binom{1}{3,1}, \binom{0}{4,0} \right\}$

\mathcal{Y}_p^n : TYPE CLASS OF n-TYPE p : IT IS THE SET OF THOSE $x_1^n \in A^n$ FOR WHICH FOR WHICH $\hat{p}_{x_1^n} = p$.

$\mathcal{Y}_{\binom{5}{1,4}} = \{01111, 10111, 11011, 11101\}$

LET $A = \{a_1, \dots, a_t\}$ WHERE $t = |A|$. LET $p \in \mathcal{P}_n$ BE FIXED.

LET ℓ_i DENOTE THE OCCURRENCE OF a_i IN x_1^n FOR ANY FIXED $x_1^n \in \mathcal{Y}_p^n$ (IT DOES NOT MATTER WHICH ELEMENT OF \mathcal{Y}_p^n WE CHOOSE)

THEN THE t-TUPLE $(\ell_1, \ell_2, \dots, \ell_t)$ OF NON-NEGATIVE INTEGERS WITH $\ell_1 + \dots + \ell_t = n$ UNIVOCALLY DETERMINES THE n-TYPE p .

THEN $|\mathcal{P}_n| \leq (n+1)^{|A|}$ TRIVIAALLY HOLDS AND THE FORMULA OF COMBINATION WITH REPETITION GIVES

LEMMA $|\mathcal{P}_n| = \binom{n+|A|-1}{|A|-1}$

NEXT WE WILL USE THAT

$$\binom{n+|A|-1}{|A|-1} = \frac{(n+|A|-1) \dots (n+1)}{(|A|-1)(|A|-2) \dots 1} \leq \frac{(n+|A|-1)}{2} \dots \frac{n+2}{2} (n+1) \leq (n+1)^{|A|-1}$$

MAIN FACT

The number of n -types is polynomial in n

Lemma

For any n -type p :
$$\binom{n+(t-1)}{t-1} 2^{nH(p)} \leq |Y_p^n| \leq 2^{nH(p)}$$

Proof

Let (k_1, \dots, k_t) be the t -tuple which corresponds to the n -type p . We have

$$|Y_p^n| = \frac{n!}{k_1! k_2! \dots k_t!} \quad (\text{permutation with repetition})$$

$$n^n = (k_1 + \dots + k_t)^n = (k_1 + \dots + k_t)(k_1 + \dots + k_t) \dots (k_1 + \dots + k_t) =$$

$$= \sum \frac{n!}{j_1! \dots j_t!} k_1^{j_1} \dots k_t^{j_t}$$

(j_1, \dots, j_t): t tuple of nonnegative integers with $j_1 + \dots + j_t = n$)

Now we show that the largest term in
$$\frac{n!}{k_1! \dots k_t!} k_1^{k_1} \dots k_t^{k_t} \quad (*)$$
 is
$$\frac{n!}{k_1! \dots k_t!} k_1^{k_1} \dots k_t^{k_t}$$

is not the case. Denote the largest term by
$$\frac{n!}{l_1! \dots l_t!} l_1^{l_1} \dots l_t^{l_t}$$
. Then there exist r and s with $l_r > k_r$ and $l_s < k_s$. Decreasing l_r by 1 and increasing l_s by 1 multiplies the corresponding term by $\frac{l_r}{l_r - 1} \frac{l_s}{l_s + 1} > 1$. It gives a contradiction. Hence, we proved $(*)$. Then

$$x_1^{z_1} \dots x_t^{z_t} \leq n^n = \sum_{\substack{z_1 + \dots + z_t = n \\ z_i \geq 0}} \frac{n!}{z_1! \dots z_t!} x_1^{z_1} \dots x_t^{z_t} \leq \binom{n+k-1}{k-1} \frac{n!}{z_1! \dots z_t!} x_1^{z_1} \dots x_t^{z_t}$$

TOTAL OF NUMERICAL INTEGERS WITH $z_1 + \dots + z_t = n$

ANS $\frac{n^n}{x_1^{z_1} \dots x_t^{z_t}} = \prod_{i=1}^t \left(\frac{x_i}{n}\right)^{-z_i} = \prod_{i=1}^t (p(a_i))^{-n p(a_i)}$

$$= 2^{\log_2 \prod_{i=1}^t (p(a_i))^{-n p(a_i)}} = 2^{n H(p)} \quad \square$$

THE NEXT RESULT CONNECTS THE THEORY OF TYPES WITH GENERAL ENTROPY THEORY FOR ANY DIST P ON A. LET P^n DENOTES THE DISTRIBUION OF n INDEPENDENT REALIZING FROM P, THAT IS $P^n(x_1^n) = \prod_{i=1}^n p(x_i)$, $x_1^n \in A^n$

LEMMA FOR ANY DISTRIBUION P ON A AND ANY n-TUPLE Q :

$$P^n(x_1^n) = 2^{-n [D(Q||P) + H(Q)]}$$

IF $x_1^n \in Y_Q^n$ AND

$$\binom{n+k-1}{k-1} 2^{-n D(Q||P)} \leq P^n(Y_Q^n) \leq 2^{-n D(Q||P)}$$

PROOF

$$P^n(x_1^n) = \prod_{a \in A} p(a)^{n \cdot \omega(a)} = 2^{n \sum_a \omega(a) \log_2 p(a)}$$

$$= 2^{-n \sum_a \omega(a) \log_2 \frac{1}{p(a)} - n \sum_a \omega(a) \log_2 \omega(a) + n \sum_a \omega(a) \log_2 \omega(a)}$$

$$= 2^{-n [D(Q||P) + H(Q)]}$$

FINALLY, $P^n(Y_Q^n) = |Y_Q^n| \cdot P^n(x_1^n)$ FOR ANY FIXED $x_1^n \in Y_Q^n$.
 HENCE, THE SECOND CLAIM FOLLOWS BY THE FIRST CLAIM AND THE PREVIOUS LEMMA.