



EXAMPLE OF H
FAHAD ENCODING

ENTROPY AND NUMBER OF BITS

LABEL INVARIANCE

$H(F(x))$

ALTERNATIVE ENTROPIES FOR ANOMALY DETECTION

PERMUTATION ENTROPY

SAMPLE ENTROPY

EXERCISE 3

SUM ENTROPY

$H(X, Y) = H(X) + H(Y)$ FOR S.I. X, Y

$H(X, Y) \leq H(X) + H(Y)$

CONDITIONAL ENTROPY

CHAIN RULE

MEANING OF $H(X|Y)$

$0 \leq H(X|Y) \leq H(X)$

INFORMATION GAIN

INTERPRETATION - FIGURE

FORMULAS

KULLBACK · LEIBLER DIVERGENCE

EXAMPLE

$x_0 \quad x_1$

	P_0	P_1	R_0	R_1	H
→	0.5	0.5	1	1	1
→	0.1	0.9	3.32	0.15	0.46
→	0.01	0.99	6.64	0.015	0.08

FANO CODING

R.V. $X \in \Omega_X$ $P(X)$

WE WANT TO ASSIGN TO EACH OUTCOME

$n \in \Omega_X$ A BINARY VECTOR

$$\Omega_X = \{x_0, x_1, x_2, x_3\}$$

00 01 10 11

IDEA: INSTEAD OF USING FIXED

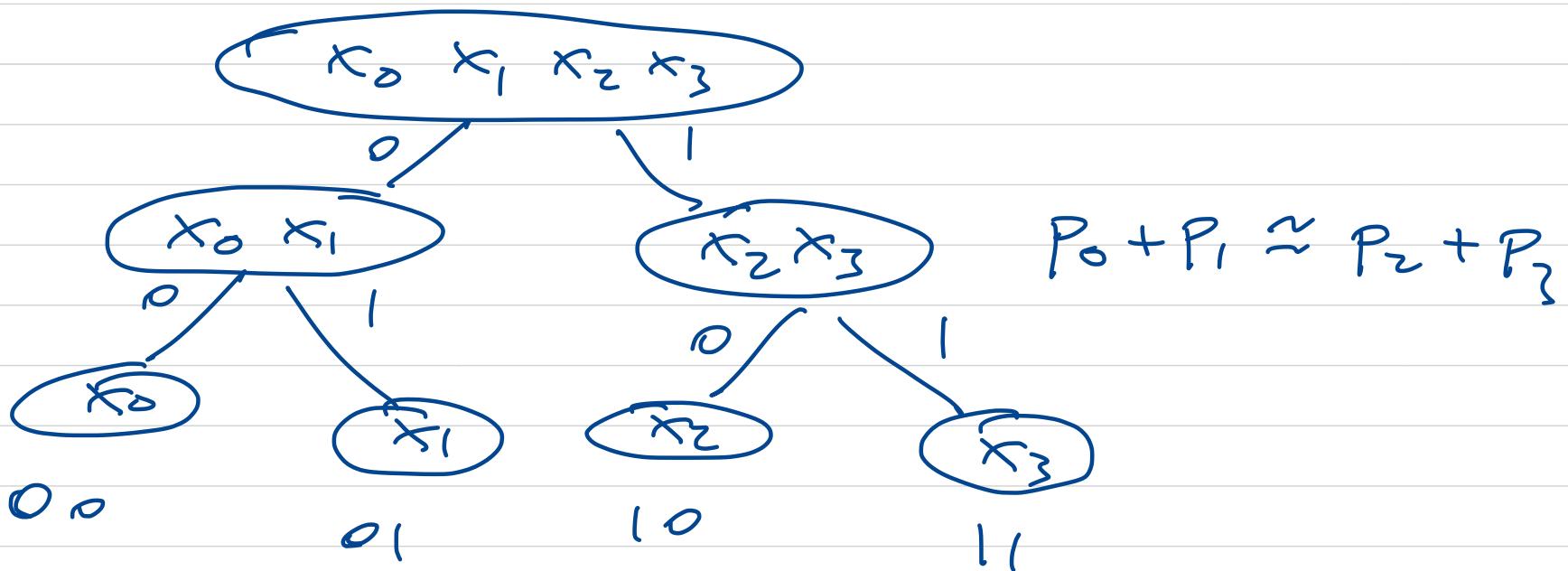
LENGTH CODES WE COULD USE

CODES WITH DIFF. LENGTH TRYING
TO ASSIGN SMALL CODES (FEW BITS)
TO MOST PROBABLE SYMBOLS

FAHAD ENCODER

CHOT OPTIMAL, HUFFMAN ENCODER

IS OPTIMAL \rightarrow STUDIED (IN SECTION 2)



WE DIVIDE EACH SET INTO TWO SETS
SUCH THAT THE SUM OF THEIR PROB.
IS AS CLOSE AS POSSIBLE

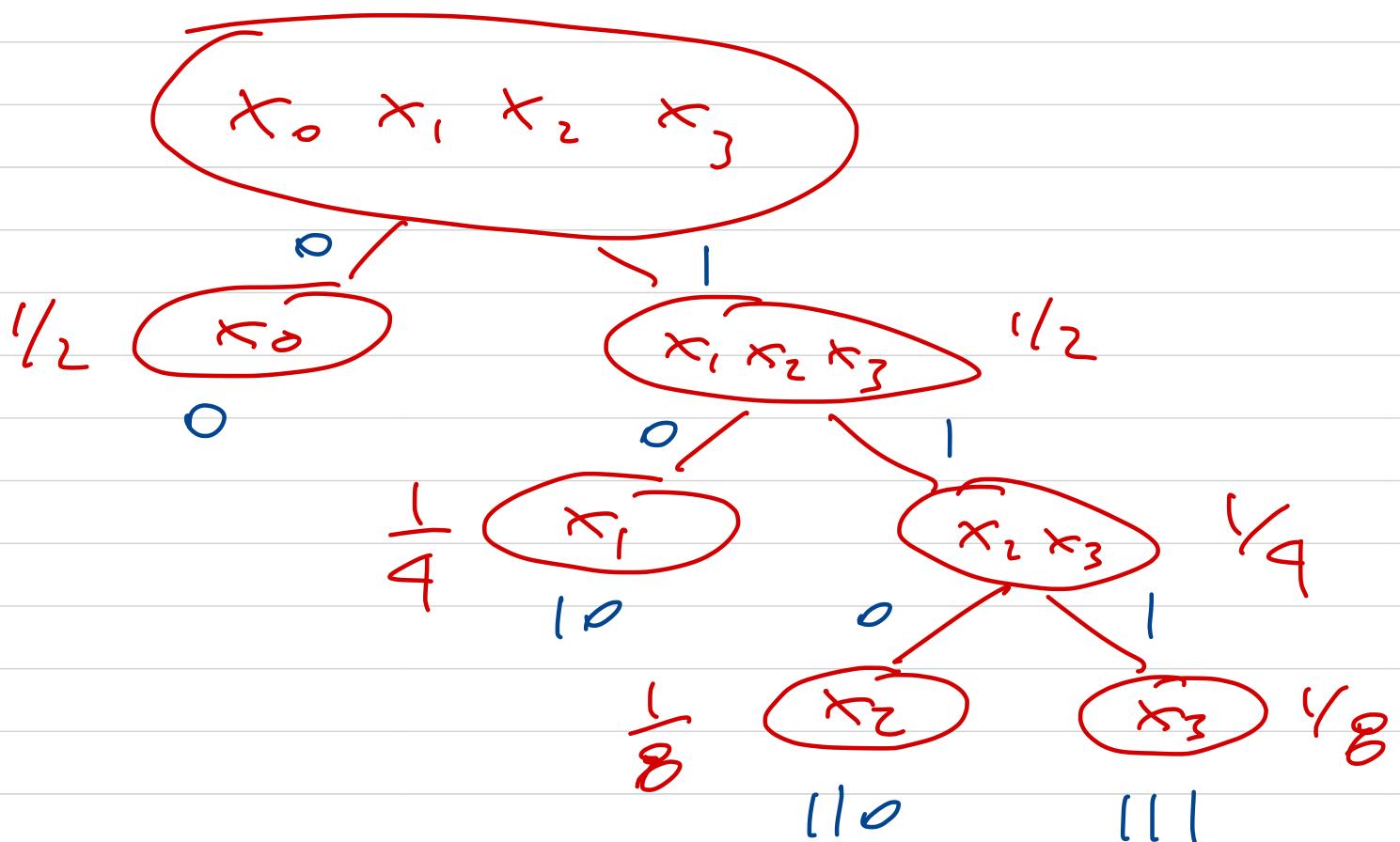
$$x_0 \quad x_1 \quad x_2 \quad x_3$$

$$\frac{1}{2} \quad \frac{1}{4} \quad \frac{1}{8} \quad \frac{1}{8}$$

$$x_0 \quad x_1 \quad x_2 \quad x_3$$

$$\frac{1}{2} \quad \frac{1}{4} \quad \frac{1}{8} \quad \frac{1}{8}$$

0 10 110 111



x_1	x_2	x_3	x_4
$\frac{1}{2}$	$\frac{1}{4}$	$\frac{1}{8}$	$\frac{1}{8}$
0	10	110	111
—	—	—	—
1	2	3	3
m_1	m_2	m_3	m_4
—	—	—	—

$$H(x) = \sum_{i=1}^4 p_i \log_2 \frac{1}{p_i} =$$

$$= \frac{1}{2} \log_2 2 + \frac{1}{4} \log_2 4 +$$

$$+ \frac{1}{8} \log_2 8 + \frac{1}{8} \log_2 8 = 1.75$$

$$\overline{m} = P_1 m_1 + P_2 m_2 + P_3 m_3 + P_4 m_4$$

AVERAGE
OF BITS

$$= \frac{1}{2} + \frac{1}{4} \cdot 2 + \frac{1}{8} \cdot 3 + \frac{1}{8} \cdot 3$$

$$= 1.75 \leq 2$$

OBTAINED
WITH FIXED
LENGTH
CODE

IN GENERAL THIS PROPERTY
HOLDS (SECTION 2)

$$H(x) \leq \bar{m} \leq H(x) + 1$$

FANo ALGORITHM PRODUCES CODES. SATISFYING THE
PREFIX CONDITION

$x_0 \quad x_1 \quad x_2 \quad x_3 \leftarrow$ NO. PREFIX
CODE

10 | 101 100
— — —

101 ?? ↘ 101 101 $x_0 x_1$
 x_2

PREFIX CONDITION

ANY CODE CANNOT BE

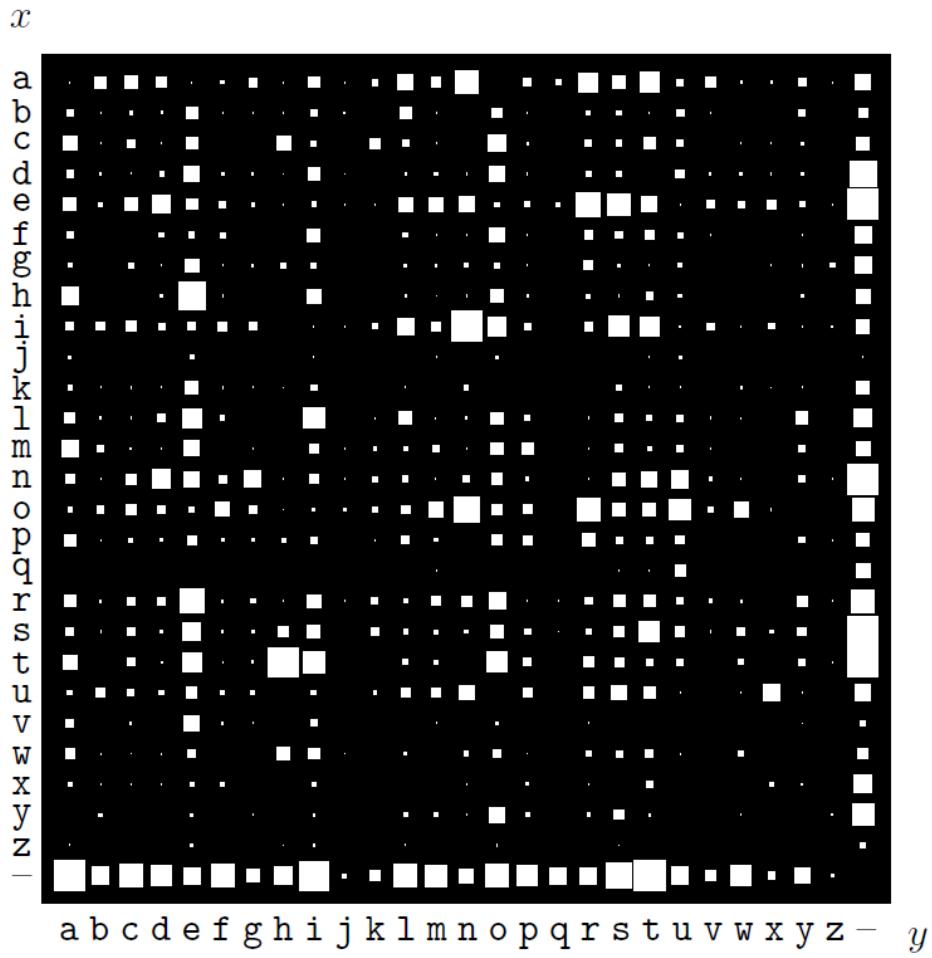
INITIAL PORTION OF

ANOTHER CODE

EXAMPLE OF PROBABILITY DISTRIBUTIONS

i	a_i	p_i	
1	a	0.0575	a
2	b	0.0128	b
3	c	0.0263	c
4	d	0.0285	d
5	e	0.0913	e
6	f	0.0173	f
7	g	0.0133	g
8	h	0.0313	h
9	i	0.0599	i
10	j	0.0006	j
11	k	0.0084	k
12	l	0.0335	l
13	m	0.0235	m
14	n	0.0596	n
15	o	0.0689	o
16	p	0.0192	p
17	q	0.0008	q
18	r	0.0508	r
19	s	0.0567	s
20	t	0.0706	t
21	u	0.0334	u
22	v	0.0069	v
23	w	0.0119	w
24	x	0.0073	x
25	y	0.0164	y
26	z	0.0007	z
27	-	0.1928	-

SINGLE LETTER



TWO CONSECUTIVE LETTERS



LABEL INVARIANCE

$$\begin{matrix} X \\ \{ z_1 & z_i & z_n \} \\ \{ p_1 & p_i & p_n \} \end{matrix}$$

$$\begin{matrix} Y \\ \{ y_1 & y_i & y_n \} \\ \{ p_1 & p_i & p_n \} \end{matrix}$$

TWO RANDOM VARIABLE WITH DIFFERENT OUTCOMES BUT SAME PROBABILITY DISTRIBUTION

$$H = \sum_{i=1}^n p_i \log_2 \frac{1}{p_i}$$

IS THE SAME

LABEL INVARIANCE \rightarrow ENTROPY DOES NOT DEPEND ON OUTCOMES BUT ONLY ON PROB.

FUNCTION OF A RANDOM VAR.

$$X \quad \Omega_x = \{x_0, x_i, x_n\}$$

F
FUNCTION

$$\Omega_y = \{f(x_0), f(x_i), f(x_n)\}$$

$$Y = F(X)$$

$$H(X)$$

$$H(Y)$$

If f is one-to-one

p_1

p_i

p_n

x_1

x_j

x_n

\downarrow

\downarrow

\downarrow

$f(x_1)$

$f(x_i)$

$F(x_n)$

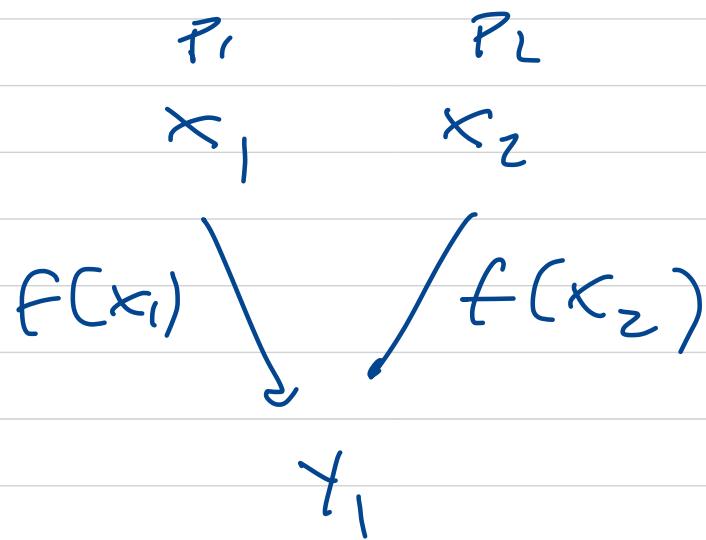
p_1

p_i

p_n

$$H(x) = H(y)$$

IF $F \subseteq \text{NOT ONE - TO - ONE}$



$$P_1 + P_L$$

$$H(F(x)) \leq H(x)$$

PROOF

$$P_1 \log_2 \frac{1}{P_1} + P_2 \log_2 \frac{1}{P_2}$$

$$(P_1 + P_2) \log_2 \frac{1}{P_1 + P_2}$$

$$P_1 \log_2 \frac{1}{P_1} + P_2 \log_2 \frac{1}{P_2} - (P_1 + P_2) \log_2 \frac{1}{P_1 + P_2}$$

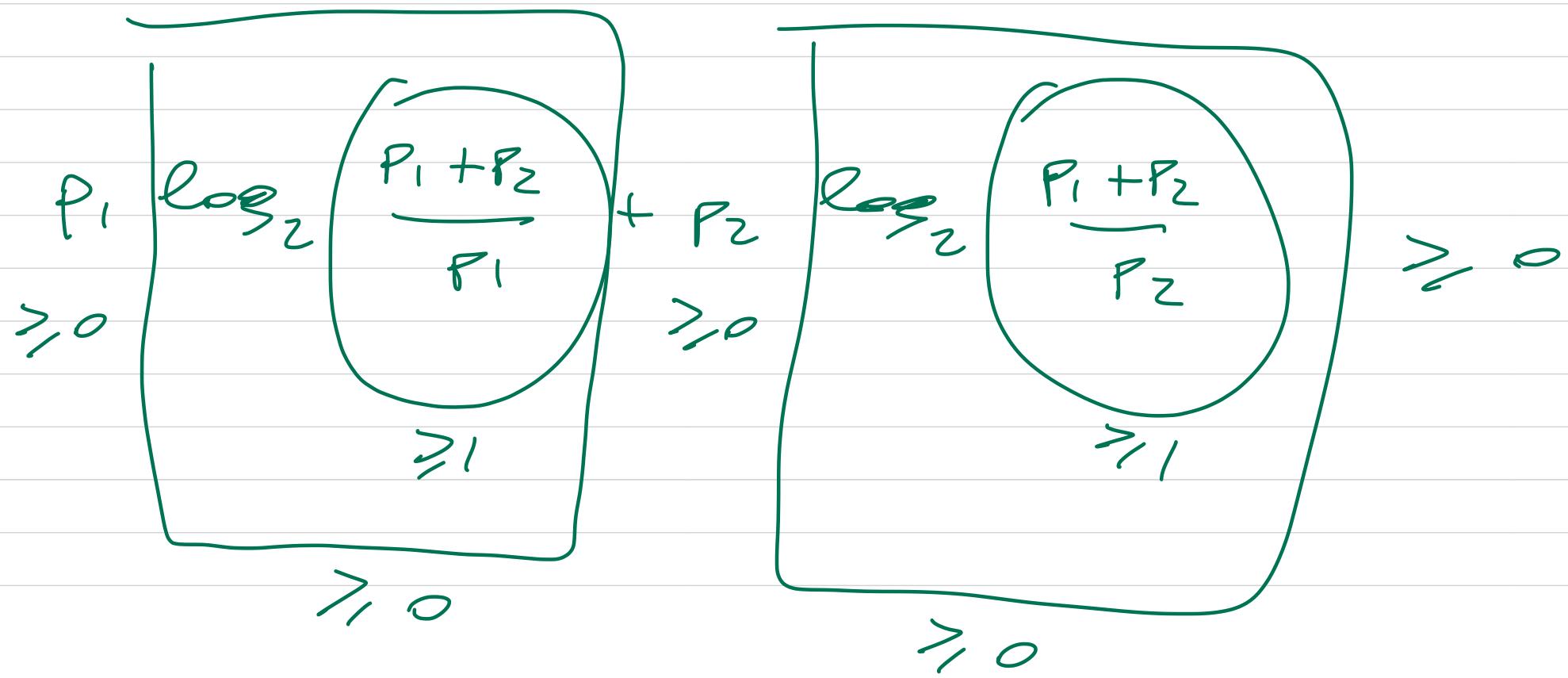
$$- P_1 \log_2 \frac{1}{P_1 + P_2}$$

$$- P_2 \log_2 \frac{1}{P_1 + P_2}$$

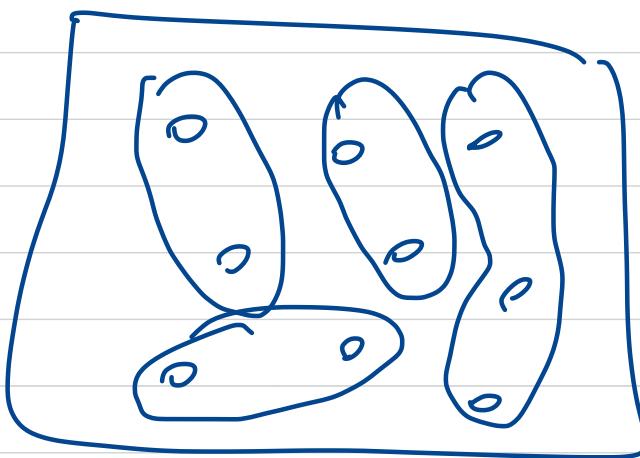
$$P_1 \log_2 \frac{1}{P_1} + P_2 \log_2 \frac{1}{P_2} - (P_1 + P_2) \log_2 \frac{1}{P_1 + P_2}$$

$$- P_1 \log_2 \frac{1}{P_1 + P_2}$$

$$- P_2 \log_2 \frac{1}{P_1 + P_2}$$



EXAMPLE



CLUSTERING IS

A FUNCTION
WHICH IS NOT
ONE-TO-ONE

CLUSTERING ALWAYS DECREASES

THE ENTROPY (I.E., THE AMOUNT
OF INFORMATION)

ALTERNATIVE ENTROPY DEFINITIONS FOR ANOMALY DETECTION

DATA SERIES

RANDOM DATA OR CONTAINS
SOME REGULARITY?

10 10 10 10 10

0110111010

KOLMOGOROV COMPLEXITY OF A SEQUENCE

MINIMUM SIZE OF A MESSAGE (PROGRAM)

WHICH ALLOWS TO REBUILD THE SEQUENCE

IF THE SEQUENCE IS COMPLETELY RANDOM
THE SITE OF THE MESSAGE
IS EQUAL TO THE SITE OF THE
SEQUENCE

PERMUTATION ENTROPY

2 8 9 7 5 2 3 5 4 1 6

ORDER $m = 3$

WE ANALYZE VECTORS OF LENGTH 3
INSIDE THE SEQUENCE

2 8 9 7 5 2 3 5 4
8 9 7 5 2 3 5 4 1
9 7 5 2 3 5 4 1 6

2 8 9 7 5 2 3 5 4
 8 9 + 5 2 3 5 4 1
 9 7 5 2 3 5 4 1 6

1 2 3 3 3 1 1 3 2
 2 3 2 2 1 2 3 2 1
 3 1 1 1 2 3 2 1 3
 - - - / / - - - -

$$m = 3$$

1 1 2 2 3 3
 2 3 1 3 1 2
 3 2 3 1 2 1

2 1 1 1 1 3 | 9

$\frac{2}{P_1}$ $\frac{1}{P_2}$ $\frac{1}{P_3}$ $\frac{1}{P_4}$ $\frac{1}{P_5}$ $\frac{3}{P_6}$

$$H = \sum_{i=1}^6 P_i \log_2 \frac{1}{P_i}$$

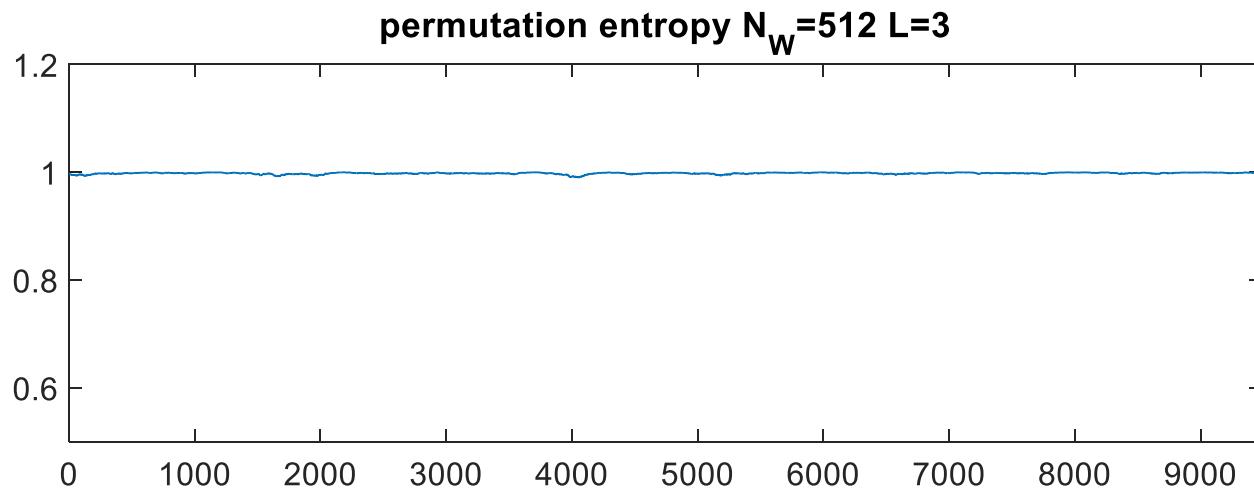
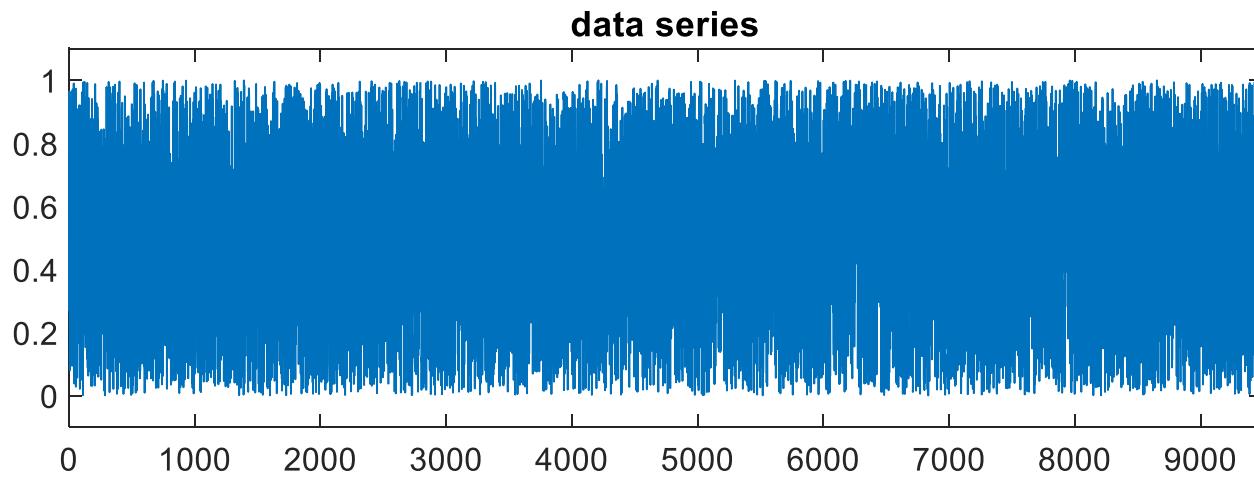
$$0 \leq H \leq \log_2 m!$$

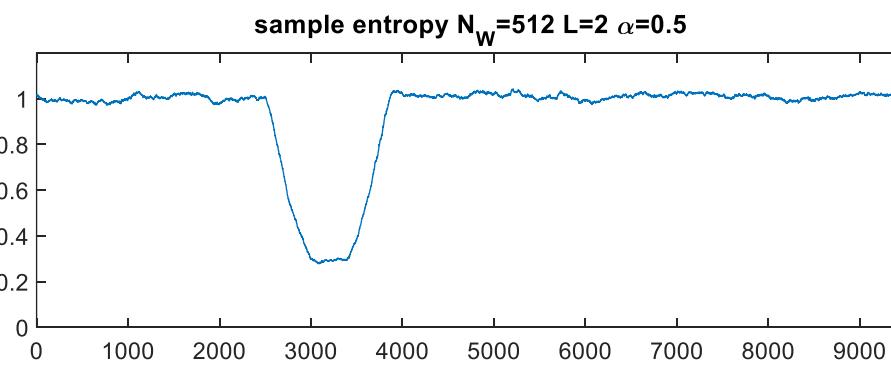
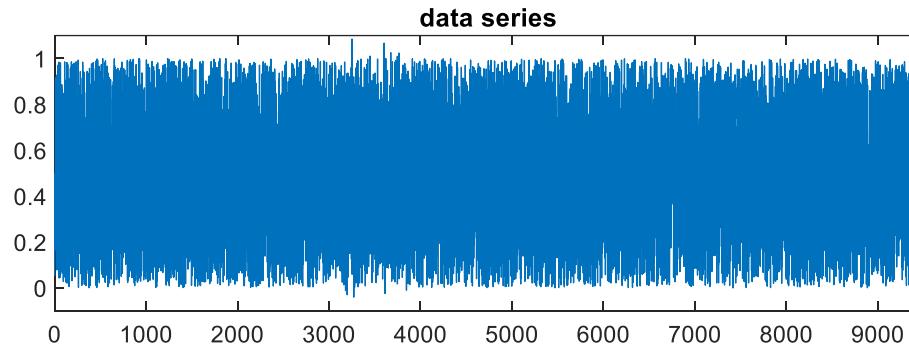
$$0 \leq \bar{H} = \frac{\sum p_i \log_2 \frac{1}{p_i}}{\log_2 m!} \leq 1$$

RANDOM DATA $\rightarrow \bar{H}$ CLOSE TO 1

CORRELATED DATA \rightarrow WE OBSERVE A
DROP OF \bar{H}

EXAMPLE : RANDOM DATA





ALTERNATIVE: SAMPLE ENTRY

$m = 3$

x	x	x	x	x
x	x	x	x	x
x	x	x	x	x

A B OF ORDER m

CHEBYCHEV DISTANCE

$$d(\underline{A}, \underline{B}) = \max_j |A_j - B_j|$$

WE COUNT HOW MANY TIME

$$d(\underline{A}, \underline{B}) \leq t = \alpha \sigma_n$$

ST. DEV.
PARAMETER SEQ,
 ≤ 0.5

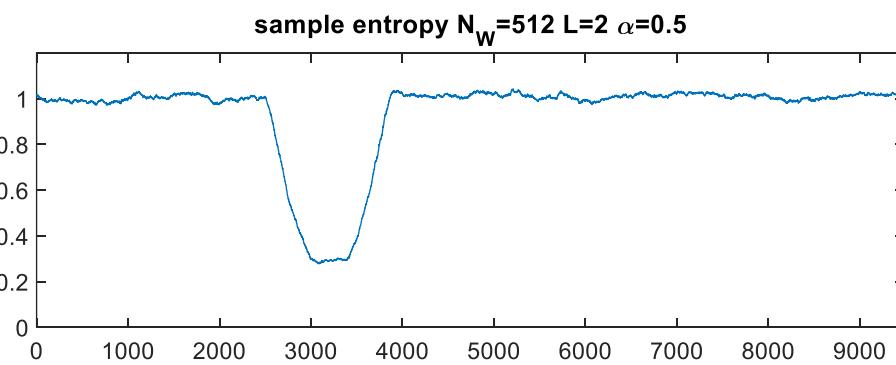
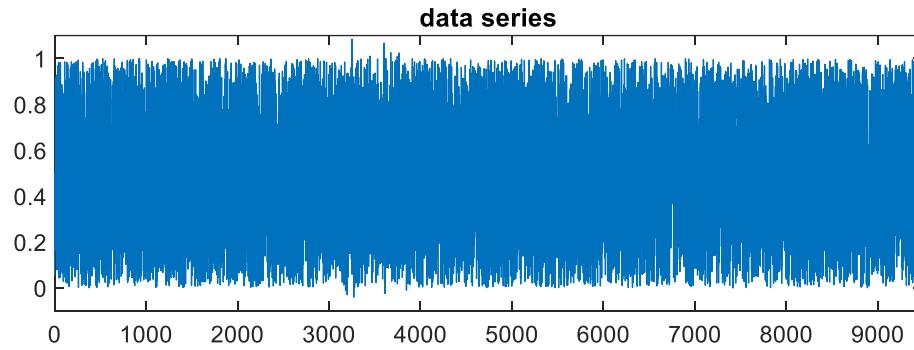
c_m

WE REPEAT THIS OPERATION

WITH ORDER $M+1$

c_{M+1}

$$\text{SAMPLE ENTROPY} = \log_2 \frac{c_m}{c_{m+1}}$$



JOINT ENTROPY

2 RANDOM VARIABLES

$X \in \Sigma_X$

$Y \in \Sigma_Y$

(X, Y)

		Temp < 25	Temp ≥ 25	Y
		Sunny	0.4	0.2
		Cloudy	0.35	0.05
X	Y			

$$H(X, Y) = \sum_{x,y} p(x, y) \log_2 \frac{1}{p(x, y)}$$

$$0 \leq H(X, Y) \leq \log_2 [|\Sigma_X| |\Sigma_Y|]$$

WE WANT TO STUDY RELATIONSHIP

BETWEEN $H(X, Y)$ AND $H(X) + H(Y)$

X Y ARE STATISTICALLY INDEPENDENT

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

FOR ALL THE OTHER CASES

$$H(X, Y) \leq H(X) + H(Y)$$

PROOF

IMPORTANT
PROPERTY
BASED ON
MARGINALIZATION

$$\sum_y p(n, y) = p(n)$$

$\sum_{n,y} p(n, y) f(n)$

$$= \sum_n f(n) \sum_y p(n, y)$$

$$= \sum_n f(n) p(n)$$

$$S. I. \quad p(n, y) = p(n)p(y)$$

$$\sum_{n,y} p(n,y) \log_2 \frac{1}{p(n,y)}$$

$$= \sum_{n,y} p(\rightarrow) p(y) \log_2 \frac{1}{p(\rightarrow)p(y)} =$$

$$= \sum_{n,y} p(\rightarrow) p(y) \log_2 \frac{1}{p(\rightarrow)} +$$

$$+ \sum_{n,y} p(\rightarrow) p(y) \log_2 \frac{1}{p(y)}$$

$$= \sum_{x,y} p(x) p(y) \log_2 \frac{1}{p(x)}$$

$$+ \sum_{x,y} p(x) p(y) \log_2 \frac{1}{p(y)}$$

$$= \sum_x p(x) \log_2 \frac{1}{p(x)}$$

$$+ \sum_y p(y) \log_2 \frac{1}{p(y)}$$

$$= H(X) + H(Y)$$

PROOF 2

$$\sum p(ny) \log_2 \frac{1}{p(ny)} - \sum_n p(n) \log_2 \frac{1}{p(n)}$$

$$- \sum_y p(y) \log_2 \frac{1}{p(y)}$$

$$= \sum_p p(ny) \log_2 \frac{1}{p(ny)} - \sum_{ny} p(ny) \log_2 \frac{1}{p(n)}$$

$$- \sum_{ny} p(ny) \log_2 \frac{1}{p(y)}$$

$$= \sum_{ny} p(ny) \log_2 \frac{1}{p(ny)} - \sum_{ny} p(ny) \log_2 \frac{1}{p(n)}$$

$$- \sum_{ny} p(ny) \log_2 \frac{1}{p(y)}$$

$$= \sum_{ny} p(ny) \log_2 \frac{p(n)p(y)}{p(ny)}$$

$$\leq \log_2 e \sum_{ny} p(ny) \left[\frac{p(n)p(y)}{p(ny)} - 1 \right]$$

$$\leq \log_2 e \sum p(ny) \left[\frac{p(n)p(y)}{p(ny)} - 1 \right]$$

$$= \log_2 e \left[\sum_{ny} p(n)p(y) - \sum p(ny) \right]$$

$$\sum_n p(n) \sum_y p(y)$$

1 . 1

$$1 - 1 = 0$$

$$H(x,y) - H(x) - H(y) \leq 0$$

CONDITIONAL ENTROPY

$$H(x,y) = \sum_{n,y} p(n,y) \log_2 \frac{1}{p(n,y)} = (*)$$

$$p(n,y) = p(n|y)p(y)$$

$$(*) = \sum p(n,y) \log_2 \frac{1}{p(n|y)p(y)} =$$

$$= \sum p(x|y) \log_2 \frac{1}{p(x|y)p(y)} =$$

$$= \boxed{\sum_{x|y} p(x|y) \log_2 \frac{1}{p(x|y)}} H(x|y)$$

$$+ \sum_{x|y} p(x|y) \log_2 \frac{1}{p(y)}$$

 $\sum_y p(y) \log_2 \frac{1}{p(y)} \rightarrow H(Y)$

$$H(x, y) = H(x|y) + H(y)$$

$H(x) =$ UNCERT. ABOUT x

$$H(x|y) = \log_2 p(z|y) \log_2 \frac{1}{p(z|y)}$$

IMPORTANT

$H(x|y)$ IS THE RESIDUAL UNCERT. ABOUT x
WHEN THE OUTCOME OF y
IS REVEALED

IF X AND Y ARE ST. INDEP.

$$H(X|Y) = H(X)$$

EVEN IF WE OBSERVE Y

WE GAIN NO INFO ABOUT X

PROOF

$$H(x, y) = H(x|y) + H(y)$$

IF x AND y ARE IN S !

$$H(x, y) = H(x) + H(y)$$

$$H(x) + H(y) = H(x|y) + H(y)$$

$$0 \leq H(x|y) \leq H(x)$$

PROOF

IF x AND y NOT $\perp\!\!\!\perp$,

$$H(x,y) \leq H(x) + H(y)$$

$$H(x,y) = H(x|y) + H(y)$$

$$H(x|y) \leq H(x)$$

$$P \leq H(x|y) \leq H(x)$$

(If y contains all info

about x $H(x|y) = 0$

(example : y is a one-to-one
function of x)

CHAIN RULE OF ENTROPY

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

$$H(XYZ) = H(X|YZ) + H(Y|Z) + H(Z)$$

IMPORTANT: INFORMATION GAIN

$$H(x)$$

$$H(x|y)$$

$$I(x,y) = H(x) - H(x|y)$$

HOW MUCH WE GAIN ABOUT X

KNOWLEDGE BY OBSERVING Y OUTCOME

$$0 \leq I(x,y) \leq H(x)$$

①

$$I(x,y) = H(x) - H(x|y)$$

$$H(xy) = H(x|y) + H(y)$$

②

$$I(xy) = H(x) + H(y) - H(x,y)$$

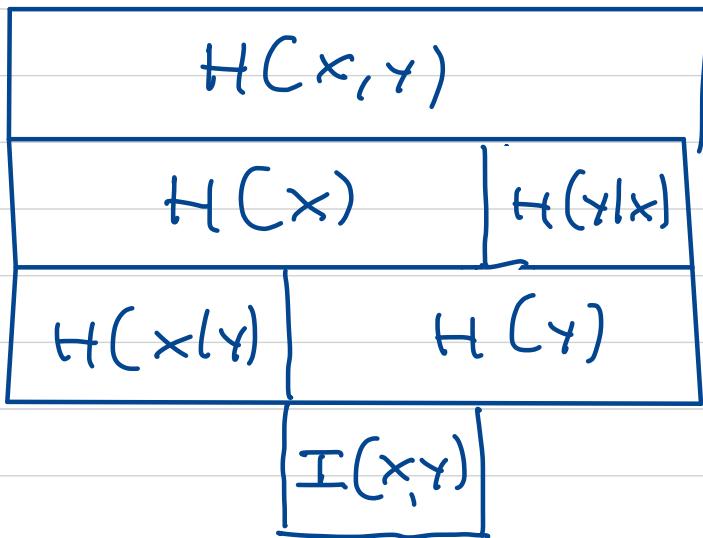
$$H(yx) = H(y|x) + H(x)$$

③

$$I(xy) = H(y) - H(y|x)$$

IHF. Gain is symmetric

IMPORTANT



$$H(X,Y) \leq H(X) + H(Y)$$

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

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

$$I(X;Y) = H(X) - H(X|Y)$$

$$I(X;Y) = H(Y) - H(Y|X)$$

$$I(X;Y) = H(X) + H(Y) - H(X,Y)$$

EXAMPLE 1

$P(X,Y) =$

0.400 0.200

0.350 0.050

$P(X) =$

0.600 0.400

$P(Y) =$

0.750 0.250

$$H(X,Y) = 1.739 \quad H(X) = 0.971 \quad H(X|Y) = 0.928 \quad I(X;Y) = 0.043$$

$$H(X,Y) = 1.739 \quad H(Y) = 0.811 \quad H(Y|X) = 0.768 \quad I(Y;X) = 0.043$$

$\text{pmi} =$

-0.170 0.415

0.222 -1.000

$$\underline{E_Z} = 0.043$$

EXAMPLE 2

P(X, Y) =

0.700 0.020

0.080 0.200

P(X) =

0.720 0.280

P(Y) =

0.780 0.220

H(X, Y) = 1.229

H(X) = 0.855

H(X|Y) = 0.469

I(X;Y) = 0.387

H(X, Y) = 1.229

H(Y) = 0.760

H(Y|X) = 0.374

I(Y;X) = 0.387

pmi =

0.318 -2.986

-1.449 1.699

E_Z = 0.387

FULL BACK- LEIBLER DIVERGENCE

$$\begin{aligned} \mathcal{L} &= \{x_0 \ x_i \ x_n\} && \text{OBSERVED} \\ P &= \{p_1 \ p_i \ p_n\} && \text{DISTRIBUTION} \\ Q &= \{q_1 \ q_i \ q_n\} && \text{DISTRIBUTION} \\ &&& \text{MODEL} \end{aligned}$$

$$D_{KL}(P||Q) = \sum_i p_i \log_2 \frac{p_i}{q_i}$$

$$\text{IF } p_i = q_i \quad \text{then} \quad D(P||Q) = 0$$

PROPERTIES SIMILAR TO DISTANCE

BUT NOT A DISTANCE

$$D(P|Q) \neq D(Q|P)$$

$$D_{KL}(P|Q) \geq 0$$

PROOF

$$D = \sum_i p_i \log_2 \frac{p_i}{q_i}$$

$$-D = \sum_i p_i \log_2 \frac{q_i}{p_i}$$

$$\in \log_2 e \sum_i p_i \left(\frac{q_i}{p_i} - 1 \right)$$

$$= \log_2 e \left[\sum_i q_i - \sum_i p_i \right] = 0$$

$$-D \leq 0 \rightarrow D \geq 0$$

LINK BETWEEN D AND INF.
GAIN

$$I(X;Y) = H(X) + H(Y) - H(X,Y)$$

$$= \sum_{x,y} p(x,y) \log_2 \frac{p(x,y)}{p(x)p(y)}$$

$$= D_{KL}(P(X,Y) \mid\mid P(X)P(Y))$$

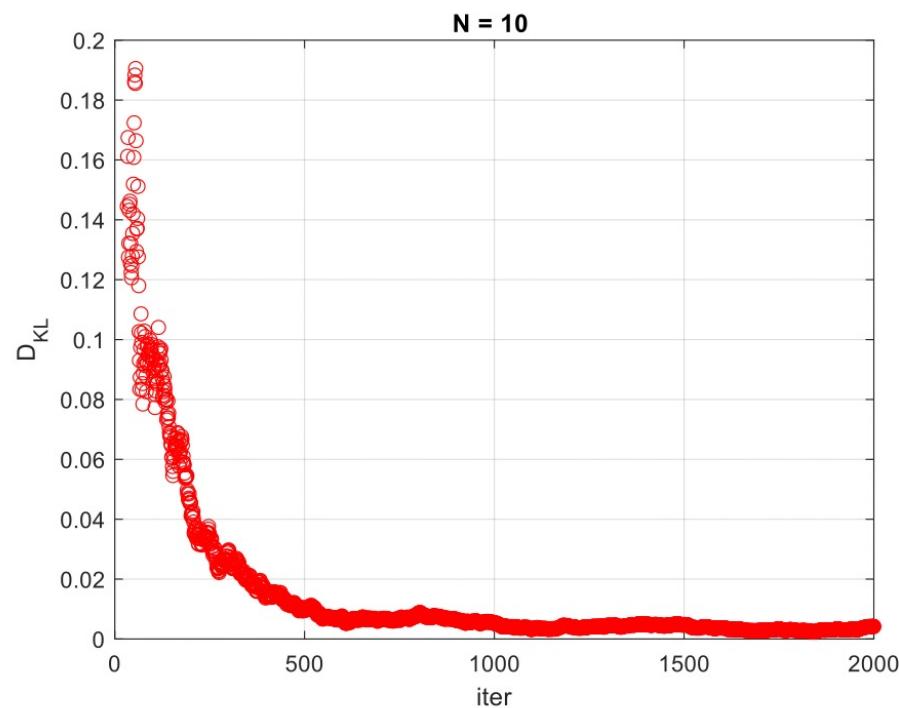
$$D_{KL}(P||Q) = \sum_i P_i \log_2 \frac{P_i}{Q_i}$$

P OBSERVED
Q MODEL

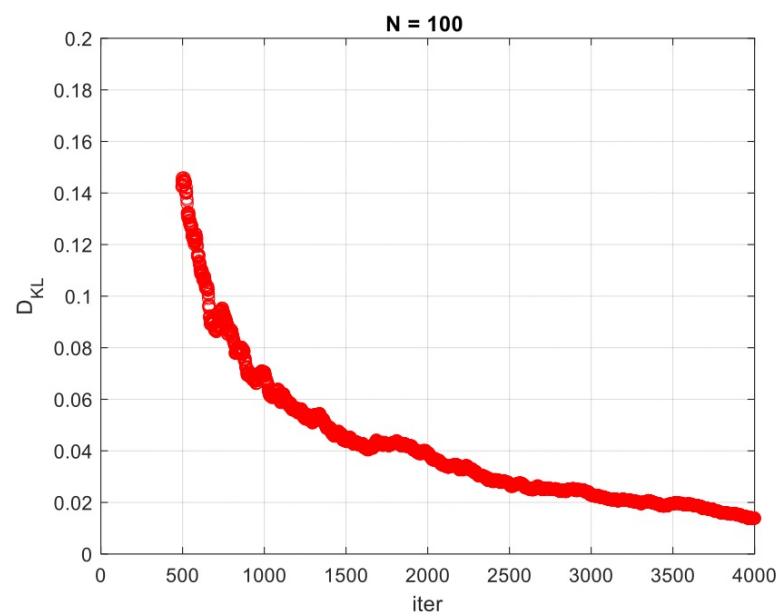
D_{KL} REPRESENTS THE INFO. GAIN
WHEN WE USE THE OBSERVED P
INSTEAD OF THE MODEL DISTR.

\equiv THE LOSS WHEN
WE USE MODEL Q INSTEAD OF
OBSERVED P.

Example 1: OBSERVED DATA = RANDOMLY GENERATED
WITH UNIFORM DISTRIBUTION
BETWEEN 1 AND N



VS. IDEAL
UNIFORM
DISTRIBUTION
BETWEEN
1 AND N

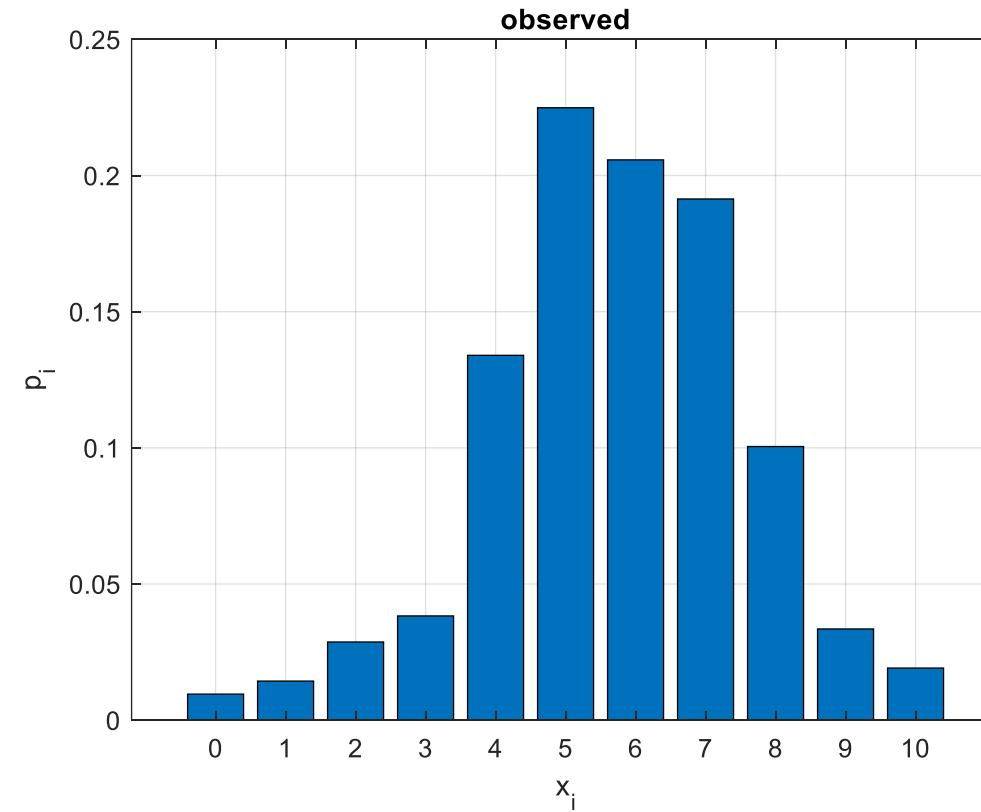


EXAMPLE 2

OBSERVED DISTRIBUTION

WE LOOK FOR THE BEST FIT

WITH BINOMIAL DISTRIBUTION



Information Theory for Data Science

Assignment 1

Introduction to Information Theory and application to Classifiers

Draft version 0.2

Exercises:

1. Entropy of a binary random variable with 3 outcomes (pt. X)
2. Application of the principle of maximum entropy (pt. X)
3. Permutation Entropy for Time Series Anomaly Detection (pt. X)
4. Kullback Leibler divergence from empirical distribution (pt. X)

Exercise 1 - Entropy of a binary random variable

Exercise 1.A

1. Plot the entropy of a binary random variable
2. Discuss the result

Exercise 1.B – Renyi entropy

$$H_\alpha(X) = \frac{1}{1-\alpha} \log_2 \left[\sum_i p_i^\alpha \right] \quad \alpha \geq 0 \quad \alpha \neq 1$$

3. Prove that the limit for $\alpha \rightarrow 1$ is the Shannon entropy.
4. Repeat point 1 with different values of α (smaller and bigger than 1)
5. Comment the results

Exercise 1.C

- Prove that this averaging operation always increases the entropy

$$\{p, p, p_3\} \quad p = \frac{p_1 + p_2}{2}$$

Hint: use the log inequality

Exercise 2 - Application of the principle of maximum entropy

Exercise 2

1. Invent an exercise where you have a random variable X with alphabet Ω_X where each outcome has a given “cost”.
2. Fix the mean value bigger than the arithmetic average of the costs, and apply the principle of maximum entropy to find the probability distribution $P(X)$
3. Plot $P(X)$
4. Repeat with a mean value equal to the arithmetic average and plot the result
5. Repeat with other values of the mean value and plot the results
6. Comment the results

You must numerically solve the equation generated by the Lagrange optimization.

As an example , for Matlab you can use

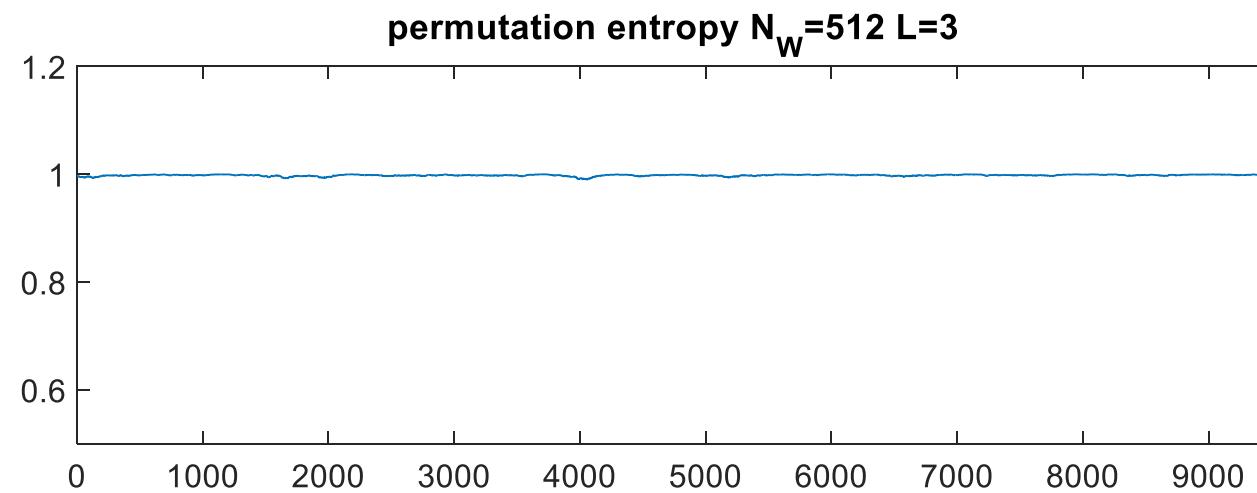
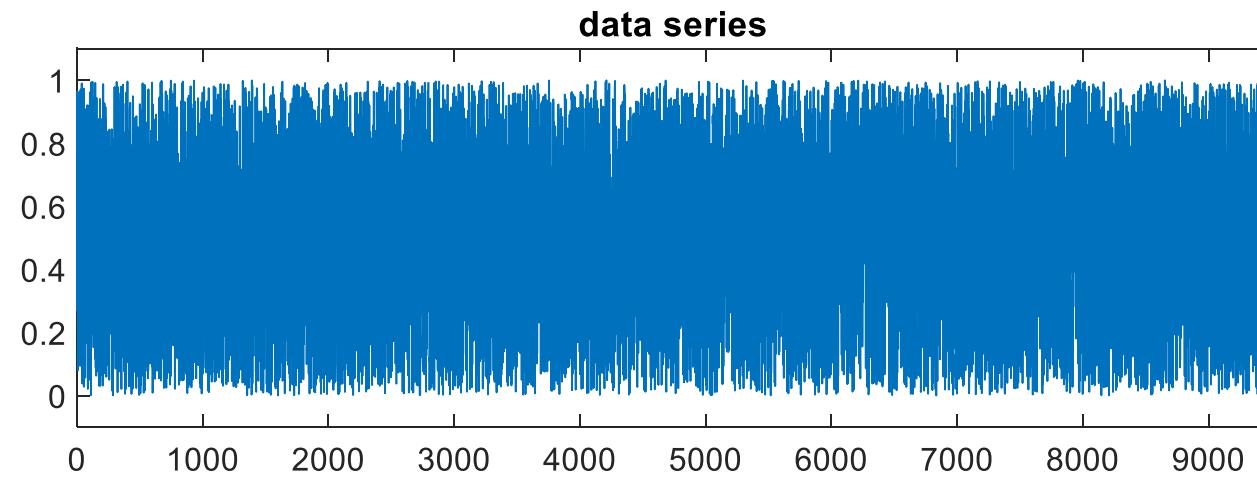
```
syms x
eqn = ( . . . ) *mu == ( . . . );
V = vpasolve(eqn, x, [0 10])
```

Exercise 3 - Permutation Entropy for Time Series Anomaly Detection

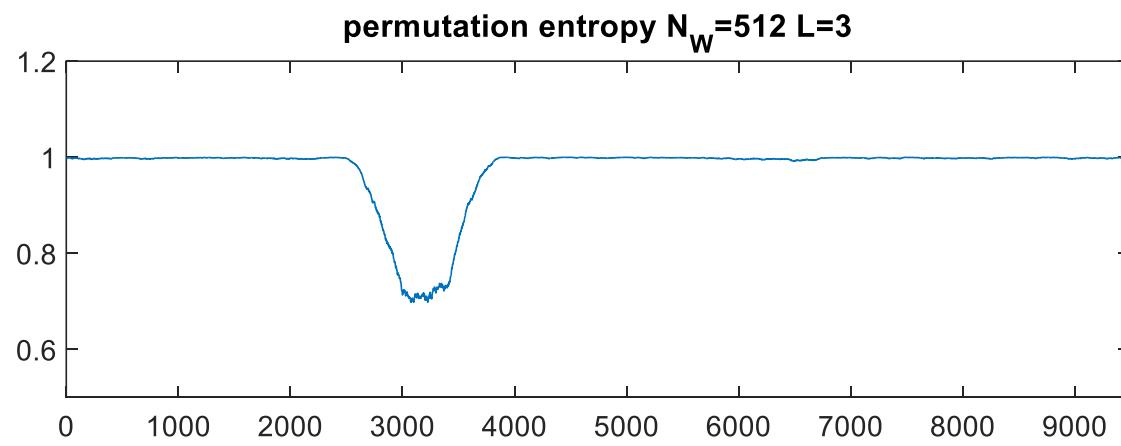
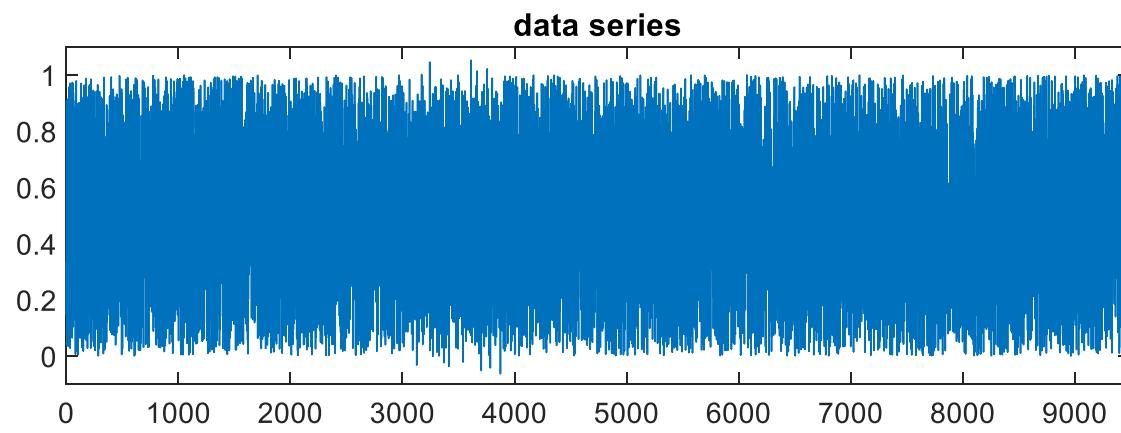
Exercise 3

Write a program that:

- Generates a data series made by 10,000 random symbols.
- Computes the permutation entropy with a sliding window of $N_w = 512$ symbols. (For example, choose an order $m = 3$.)
- Plot the data and the entropy.



- Insert a pattern between 3000 and 4000 made by correlated data (with about the same mean value and variance of the original random data).
- Describe how you generated the correlated data
- Compute the permutation entropy with the same N_w and N_R used before.
- Plot the data and the entropy.
- Comment the results



Some Matlab functions that might be useful

sort
perms

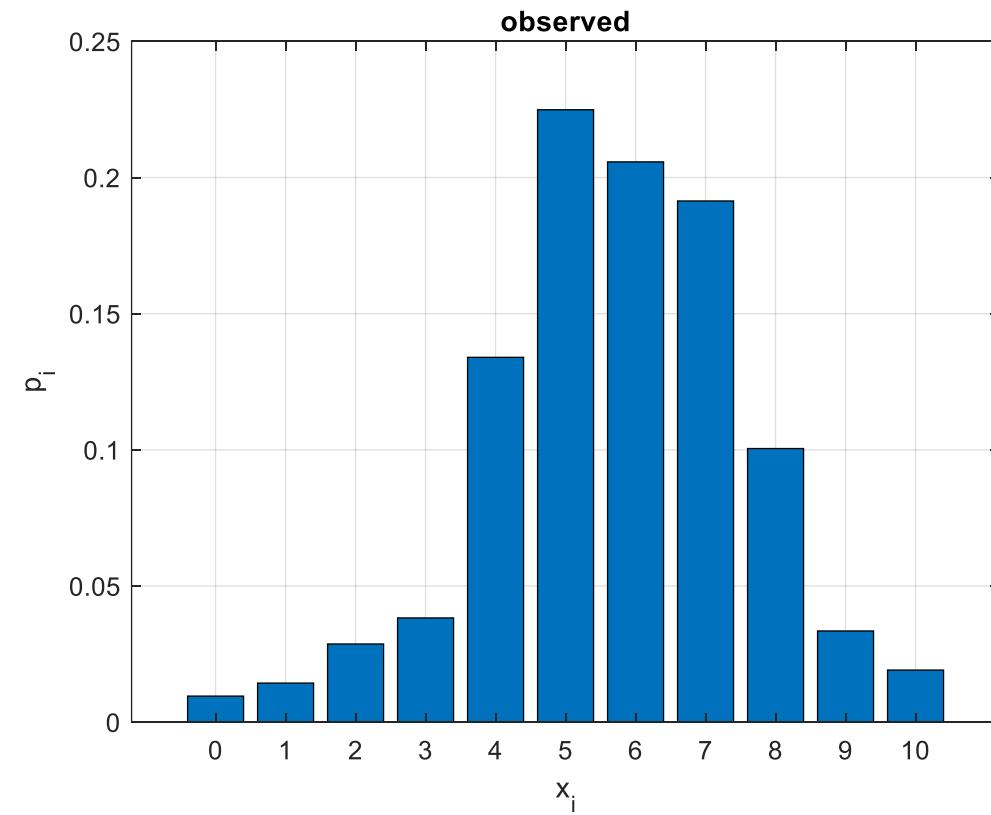
Exercise 4 - Kullback-Leibler distance from empirical distribution

Exercise 4

Consider the observed data

$$x_i = [0 \ 1 \ 2 \ 3 \ 4 \ 5 \ 6 \ 7 \ 8 \ 9 \ 10]$$

$$\text{NUM}(x_i) = [1 \ 3 \ 2 \ 8 \ 22 \ 45 \ 44 \ 42 \ 24 \ 8 \ 3]$$



1. Compare with uniform pmf
2. Write in the title the KL divergence value

Compare with binomial pmf with $0 < p < 1$ (step=0.001)

3. Identify the pmf at minimum KL divergence
4. Plot the pmf
5. In the title write the value of p and KL divergence