

LECTURE 02

01 - 10 - 2024

FANO ENCODING

ENTROPY AND NUMBER OF BITS

LABEL INVARIANCE

$$H(F(x))$$

SPLIT ENTROPY

$$H(X, Y) = H(X) + H(Y) \text{ FOR S.I. } X, Y$$

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

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

## MEANING OF ENTROPY

THE AMOUNT OF INFORMATION

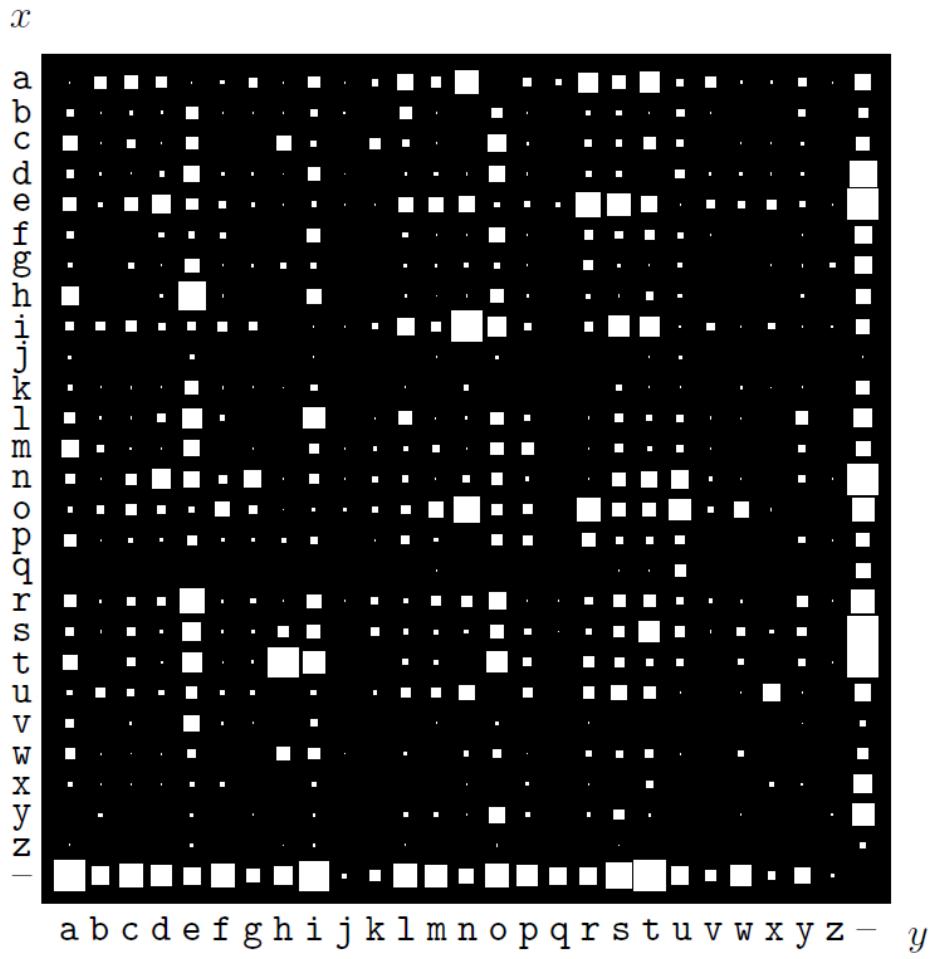
WE RECEIVE WHEN

WE OBSERVE AN EVENT

# EXAMPLES OF PROBABILITY DISTRIBUTIONS

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

SINGLE LETTER



TWO CONSECUTIVE LETTERS

## FAHO ENCODING

$$X \quad \mathcal{R}_X = \{ n_0, n_1, n_n \}$$

$$P_X = \{ p_0, p_1, p_n \}$$

WE COULD USE FIXED LENGTH BINARY CODES TO  
REPRESENT THE DIFFERENT OUTCOMES

$$\text{EX} \quad \mathcal{R}_X = \{ n_0, n_1, n_2, n_3 \}$$
$$00 \quad 01 \quad 10 \quad 11$$

IN THIS WAY WE DO NOT TAKE  
INTO ACCOUNT THEIR DIFFERENT  
PROBABILITIES.

IDEA: USE VARIABLE LENGTH CODES  
(BINARY CODES ASSOCIATED TO DIFFERENT  
OUTCOMES MAY HAVE DIFFERENT LENGTH)  
THIS WAY WE COULD ASSOCIATE  
SHORT CODES TO OUTCOMES WITH HIGH  
PROBABILITIES AND LONG CODES TO THOSE  
WITH LOW PROBABILITIES.

GOAL : REDUCE THE AVERAGE

NUMBER  $\bar{n}$  OF BITS

NECESSARY TO REPRESENT THE OUTCOMES → ADVANTAGE :

WE HAVE LESS BITS TO STORE,  
TRANSMIT, PROCESS.

THE OPTIMAL ENCODING ALGORITHM IS  
THE HUFFMAN ALGORITHM (SECTION 2).

WE PRESENT THE SIMPLER FANO  
ALGORITHM

$$\mathcal{N}_x = \{ n_0, n_1, n_2, n_3 \}$$

$$P_x = \left\{ \frac{1}{2}, \frac{1}{4}, \frac{1}{8}, \frac{1}{8} \right\}$$

FANO'S ALGORITHM

ASSIGN THE BINARY CODES BY USING  $P_x$

WE BUILD A TREE BY DIVIDING  
AT EACH STEP

THE SET OF OUTCOMES INTO  
2 SETS WITH CLOSE  
PROBABILITIES

$$\{ n_0 \quad n_1 \quad n_2 \quad n_3 \}$$
$$\left\{ \frac{1}{2} \quad \frac{1}{4} \quad \frac{1}{8} \quad \frac{1}{8} \right\}$$
$$\{ n_0 \}$$
$$1/2$$
$$0 \quad \quad \quad 1$$
$$\{ n_1 \quad n_2 \quad n_3 \}$$

$n_0$	0
$n_1$	1 0
$n_2$	1 1 0
$n_3$	1 1 1

$$\{ n_1 \}$$
$$1/4$$
$$0 \quad \quad \quad 1$$
$$1/2$$
$$\{ n_2 \quad n_3 \}$$
$$\begin{array}{c} 0 \quad \quad \quad 1 \\ \diagup \quad \quad \quad \diagdown \\ \{ n_2 \} \quad \quad \quad \{ n_3 \} \end{array}$$
$$0 \quad \quad \quad 1$$
$$\begin{array}{c} 0 \quad \quad \quad 1 \\ \diagup \quad \quad \quad \diagdown \\ \{ n_3 \} \end{array}$$
$$1/4$$

$n_0 = 0$ 

$P_0 = 1/2$

 $n_1 = 10$ 

$P_1 = 1/4$

 $n_2 = 110$ 

$P_2 = 1/8$

 $n_3 = 111$ 

$P_3 = 1/8$

THE CODES SATISFY THE PCEFCX

RULE : A CODE IS XOR THE FIRST PART OF ANY OTHER CODE  
CONDITION

$0|10|$   
 $n_0 \quad n_1$

AT THE RX SIDE GIVEN THE BIT SEQUENCE IT IS POSSIBLE TO RECOVER THE SYMBOL SEQUENCE WITHOUT AMBIGUITY

$n_0 = 0$

$$P_0 = 1/2$$

$n_1 = 10$

$$P_1 = 1/4$$

$n_2 = 110$

$$P_2 = 1/8$$

$n_3 = 111$

$$P_3 = 1/8$$

## AVERAGE NUMBER OF BITS

$$\bar{n} = \frac{1}{2}(1) + \frac{1}{4}(2) + \frac{1}{8}(3) + \frac{1}{8}(3)$$

$$1/2$$

$$1/2$$

$$0.75$$

$$= 1.75$$

## ENTROPY

$n_0 \quad n_1 \quad n_2 \quad n_3$

$1/2 \quad 1/4 \quad 1/8 \quad 1/8$

$$H(x) = \sum_i p_i \log_2 \frac{1}{p_i} =$$

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

$$+ \underbrace{\frac{1}{8} \log_2 8}_{3} + \underbrace{\frac{1}{8} \log_2 8}_{3}$$

$$= 1.75$$

IN THIS SPECIFIC EXAMPLE  
( WHERE THE PROBABILITIES HAVE THESE  
PARTICULAR VALUES ) WE OBTAIN

THAT

$$\bar{m} = h(x)$$

IN GENERAL WE CAN  
OBTAIN A CODE WITH

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

IMPORTANT

---

SO THE ENTROPY IS VERY CLOSE  
TO THE AVERAGE # OF BITS  
WE NEED TO REPRESENT  
THE OUTCOMES OF OUR  
RANDOM VARIABLE

## LABEL INVARIANCE

$X$

$$\Sigma_X = \{x_1, x_2, \dots, x_n\}$$

$Y$

$$\Sigma_Y = \{y_1, y_2, \dots, y_m\}$$

$$P_X = \{p_1, p_2, \dots, p_M\}$$

$$P_Y = \{p_1, p_2, \dots, p_N\}$$

$$I_X = I_Y$$

$$H = \sum_i p_i \log_2 \frac{1}{p_i}$$

TWO RANDOM VARIABLES WITH DIFF. ALPHABET  
BUT SAME PROB. HAVE SAME ENTROPY

## FUNCTION OF A RANDOM VARIABLE

$$X \quad \Omega_x = \{ n_1, n_2, \dots, n_n \}$$

$$P_x = \{ p_1, p_2, \dots, p_M \}$$

$$H(x) = \sum_i p_i \log_2 \frac{1}{p_i}$$

$$Y = f(X) \quad \Omega_Y = \{ f(n_1), f(n_2), \dots, f(n_n) \}$$

$$H(Y) \stackrel{?}{=} H(X)$$

$h_1$ ,  $h_2$

$$F(h_1) = F(h_2)$$

IF  $f$  IS INJECTIVE (ONE-TO-ONE,  
NO COLLISIONS)

$$H(f(x)) = H(x)$$

OTHERWISE

$$H(f(x)) \leq H(x)$$

IF  $F$  IS NOT INJECTIVE

$$H(x) \geq H(F(x))$$

PROOF

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

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

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

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

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

$$= p_1 \log_2 \frac{p_1 + p_2}{p_1} + p_2 \log_2 \frac{p_1 + p_2}{p_2} \geq 0$$

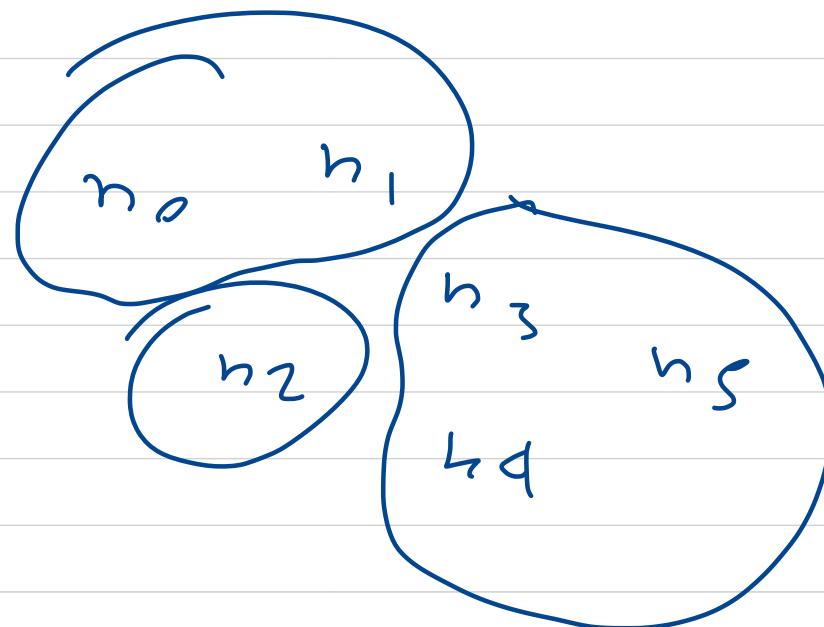
$\overbrace{\quad}$        $\overbrace{\quad}$        $\overbrace{\quad}$        $\overbrace{\quad}$   
 $\geq 0$        $\geq 1$        $\geq 0$        $\geq 1$   
 $\geq 0$        $\geq 1$

THEH

$$p_1 \log_2 \frac{1}{p_1} + p_2 \log_2 \frac{1}{p_2} \geq (p_1 + p_2) \log_2 \frac{1}{p_1 + p_2}$$

## CONSEQUENCES

IF THE FUNCTION IS NOT INJECTIVE  
THE ENTROPY (AMOUNT OF INFO.)  
DECREASES



EXAMPLE:

CLUSTERING  
ALWAYS  
DECREASES  
AMOUNT ON  
INFO.

## JOINT ENTROPY

$X \quad \mathcal{R}_X$

$Y \quad \mathcal{R}_Y$

$(x, y) \quad \mathcal{R}_X \times \mathcal{R}_Y$

		$\text{Temp} < 25$	$\text{Temp} \geq 25$
$X$	Sunny	0.4	0.2
	Cloudy	0.35	0.05

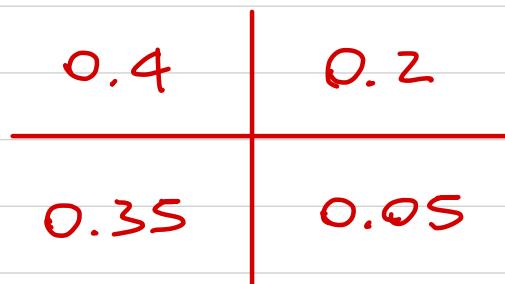
$$H(X, Y) = \sum_{n,y} p(n, y) \log_2 \frac{1}{p(n, y)}$$

RANGE

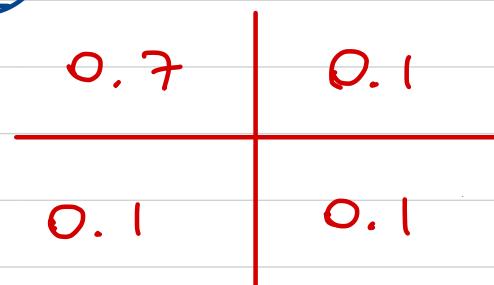
$$0 \leq H(x, y) \leq \log_2 |\mathcal{R}_x| |\mathcal{R}_y|$$

## EXAMPLES

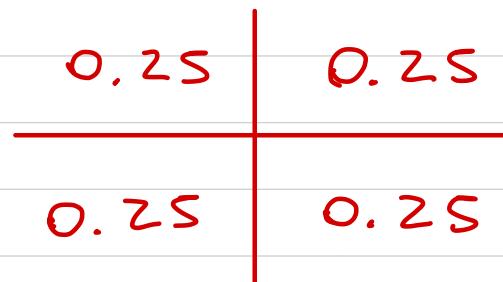
①



②



③

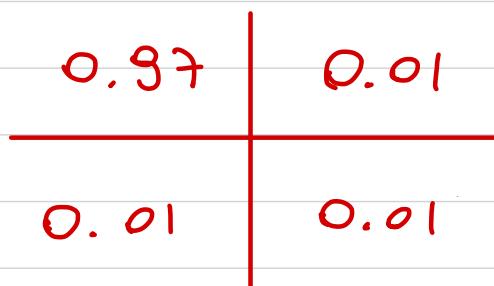


$$H(x, y) =$$

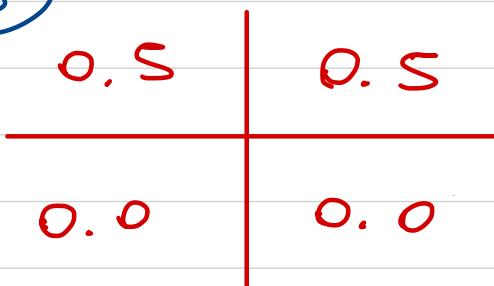
$$H(x, y) =$$

$$H(x, y) =$$

④



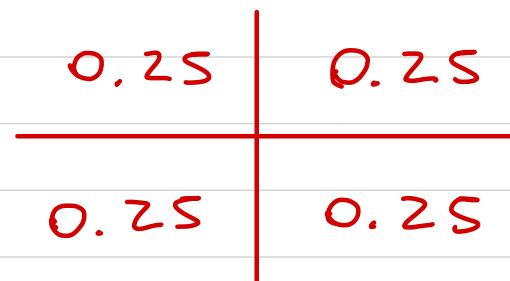
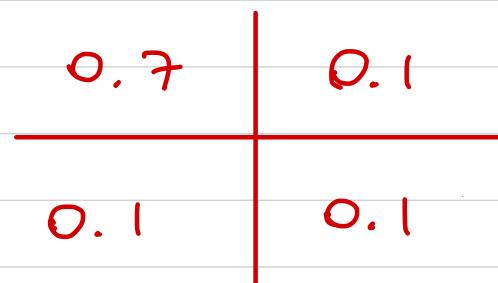
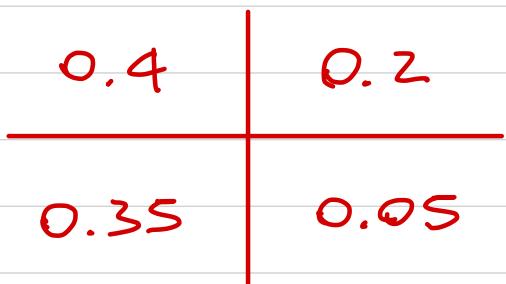
⑤



$$H(x, y) =$$

$$H(x, y) =$$

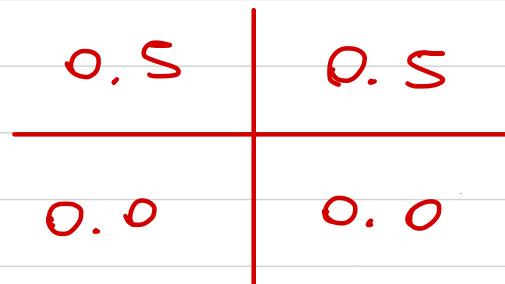
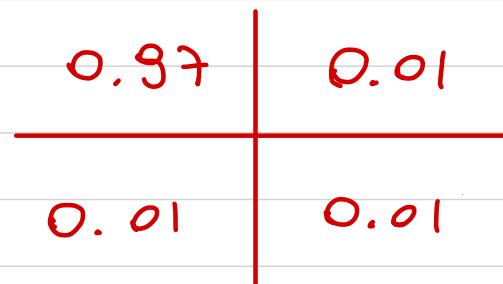
## EXAMPLES



$$H(x, y) = 1.74$$

$$H(x, y) = 1.35$$

$$H(x, y) = 2$$



$$H(x, y) = 0.24$$

$$H(x, y) = 1$$

LINK BETWEEN  $H(x, y)$  AND  $H(x) + H(y)$

$X$  AND  $Y$   
ARE STAT. IND.

$$\text{IFF } P(n, y) = P(n) P(y) \text{ FOR PAIR } (n, y)$$

---

$$\text{IF } X \text{ AND } Y \text{ ARE S.I. } H(x, y) = H(x) + H(y)$$

$X$  AND  $Y$   
ARE NOT S.I.

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

## MEANING

IF  $X$  AND  $Y$  ARE NOT ST. IND.

THEY ARE LINKED SOMEHOW

SO WHEN WE OBSERVE

THEN TOGETHER THE

AMOUNT OF INF. IS

REDUCED

# PROOF : MARGINALIZATION

$$\sum_y P(x, y) = P(x)$$

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

$$\begin{aligned}
 P(\text{SUNNY}) &= P(\text{SUNNY}, \text{TEMP} < 25) + P(\text{SUNNY}, \text{TEMP} \geq 25) \\
 &= 0.4 + 0.2 = 0.6
 \end{aligned}$$

$$P(\text{CLOUDY}) = P(\text{CLOUDY}, \text{TEMP} < 25) + P(\text{CLOUDY}, \text{TEMP} \geq 25) = 0.4$$

## MARGINALIZATION

## PROPERTY

$$\sum_{x,y} p(x,y) \cdot f(x) = \sum_x f(x) \sum_y p(x,y)$$

$$= \sum_x f(x) p(x)$$

$$p(S, T < 25) \cdot \log_2 \frac{1}{p(S)} + P(S, T \geq 25) \log_2 \frac{1}{p(S)} +$$

		$\gamma$	$\text{Temp} \geq 25$
		Sunny	Cloudy
$\times$	Temp < 25	0.4	0.2
	Cloudy	0.35	0.05

$$+ P(C, T < 25) \log_2 \frac{1}{p(C)} + P(C, T \geq 25) \log_2 \frac{1}{p(C)}$$

$$= \log_2 \frac{1}{p(S)} \cdot P(S) + \log_2 \frac{1}{p(C)} \cdot P(C)$$

X AND Y

STATISTICALLY

INDEPENDENT

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

PROOF

$$P(x, y) = P(x) \cdot P(y)$$

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

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

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

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

→ MARGINALIZATION

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

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

→  $H(x, y) = H(x) + H(y)$

NOT STATISTICALLY INDEPENDENT  
 $H(x, y) \leq H(x) + H(y)$

PROOF

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

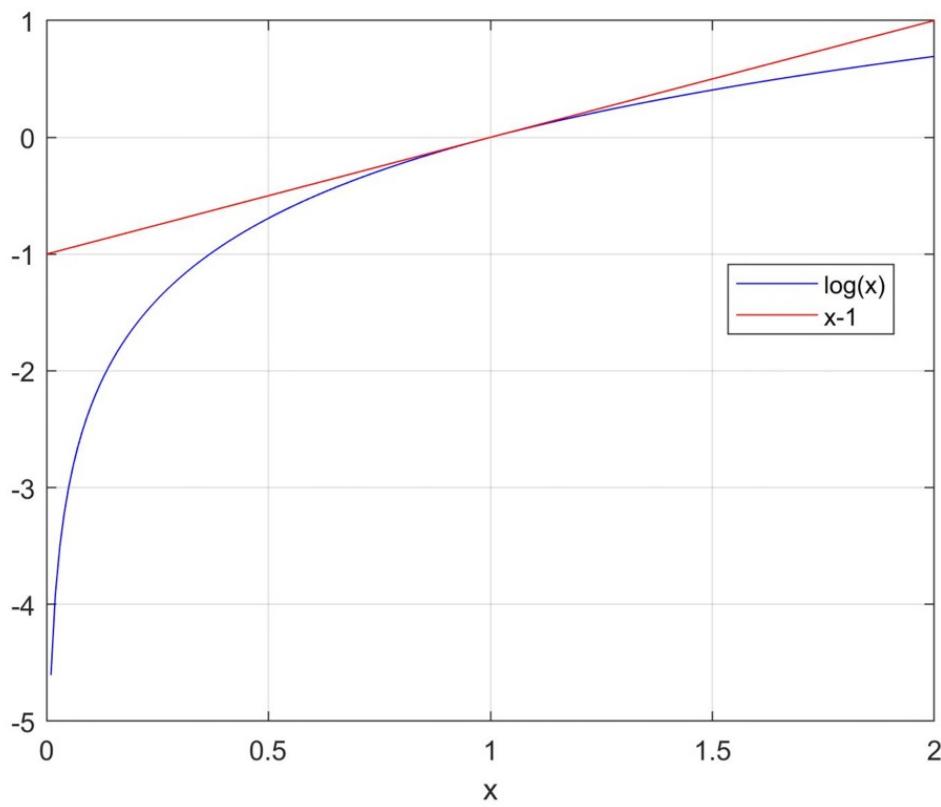
$$= \sum_{ny} p(n, y) \log_2 \frac{1}{p(n, y)} - \sum_n p(n) \log_2 \frac{1}{p(n)} \\ - \sum_y p(y) \log_2 \frac{1}{p(y)}$$

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

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

# LOG INEQUALITY

$$\log_e x \leq x - 1$$



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

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

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

$$= \log_2 e \left[ \sum_n p(n) \sum_y p(y) - \sum_{n,y} p(n,y) \right]$$

$$= \log_2 e \left[ 1 \cdot 1 - 1 \right] = 0$$

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

$(x, y)$

EXAMPLES

	$m_1$	$m_2$
$n_1$	0.4	0.2
$n_2$	0.35	0.05

	$x$	$y$
$n_1$	0.6	$m_1$ 0.75
$n_2$	0.4	$m_2$ 0.25

$$H(x, y) = 1.74$$

$$H(x) = 0.97$$

$$H(y) = 0.81$$

$$H(x, y) \leq H(x) + H(y) = 1.78$$

STATISTICAL IND:  $P(n, y) = P(n) P(y)$

	0.45	0.15
0.30	0.10	

$$H(x, y) = 1.78 = H(x) + H(y)$$

RECAP:

$$H(x)$$

$$H(x, y)$$

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

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

(FF       $x$      $y$ )

ARE ST. IND.

## CONDITIONAL ENTROPY

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

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

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

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

$$= \boxed{\sum_{ny} p(n,y) \log_2 \frac{1}{p(n|y)}} \quad \begin{array}{l} H(x|y) \\ \xleftarrow{\text{CONDITIONAL}} \\ \text{ENTROPY} \end{array}$$

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

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

MEANING (IMPORTANT)

$$H(x,y) = H(x|y) + \underline{H(y)}$$

RESIDUAL UNCERTAINTY

ABOUT  $x$  WHEN

$y$  outcome is

REVEALED

$H(x)$  = ORIGINAL UNCERTAINTY

ABOUT  $x$

CASE 1:  $X$  AND  $Y$  STATISTICALLY  
INDEPENDENT

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

WE OBSERVE  $Y$  BUT

$X$  AND  $Y$  ARE NOT LINKED

$\rightarrow$  WE GAIN NOTHING

$X \quad Y$  STATISTICALY INDEPENDENT

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

PROOF

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

IF  $X$  AND  $Y$  ARE STATISTICALLY IND.

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

$$\rightarrow H(X|Y) + H(Y) = H(X) + H(Y)$$

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

CASE 2:  $X$  AND  $Y$  NOT STATISTICALLY INDEPENDENT

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

WHEN WE OBSERVE  $Y$

THE AMOUNT OF UNCERTAINTY

ABOUT  $X$  IS REDUCED

$X \quad Y$  NOT STATISTICALLY INDEPEND.

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

PROOF

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

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

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

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

CASE 3:  $H(x|y) = 0$

$y$  contains all info about  $x$   
IF  $y = f(x)$  WITH  $f$  INJECTIVE

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

$$x \quad \mathcal{R}_x = \{n_1, n_i, n_n\}$$

$$y = f(x) \quad \mathcal{R}_y = \{f(n_1), f(n_i), f(n_n)\}$$

IF WE OBSERVE  $f(n_i)$

WE AUTOMATICALLY

KNOW  $n_i$

## SUMMARY AND MEANING

$$o \leq h(x|y) \leq h(x)$$

## CHAIN RULE OF ENTROPY

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

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

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

$$+ H(Z)$$

VERY IMPORTANT:

INFORMATION GAIN

X Y

$$I(x; y) = H(x) - \underbrace{H(x|y)}_{\text{RESIDUAL UNCERT. ABOUT } X \text{ AFTER OBSERVING } Y}$$

ORIGINAL  
UNCERT.  
ABOUT X

RESIDUAL  
UNCERT.  
ABOUT X  
AFTER  
OBSERVING  
Y

REDUCTION OF  
UNCERTAINTY ABOUT X  
OBTAINED BY OBSERVING Y

$\equiv$   
INFORMATION GAIN

ABOUT X OBTAINED

BY OBSERVING Y

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

P

A

WHEN

X AND Y

ARE ST. (IND.)

WHEN Y

CONTAINS

ALL INFO

ABOUT X

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

$$\text{GAIN} = 0$$

$$\text{GAIN} =$$

$$\text{ENTIRE } H(x)$$

TWO ALTERNATIVE FORMULATIONS  
FOR  $I(x; y)$

$$I(x; y) = H(x) + H(y) - \underline{H(x, y)}$$

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

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

$$\rightarrow H(x|y) = H(x, y) - H(y)$$

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

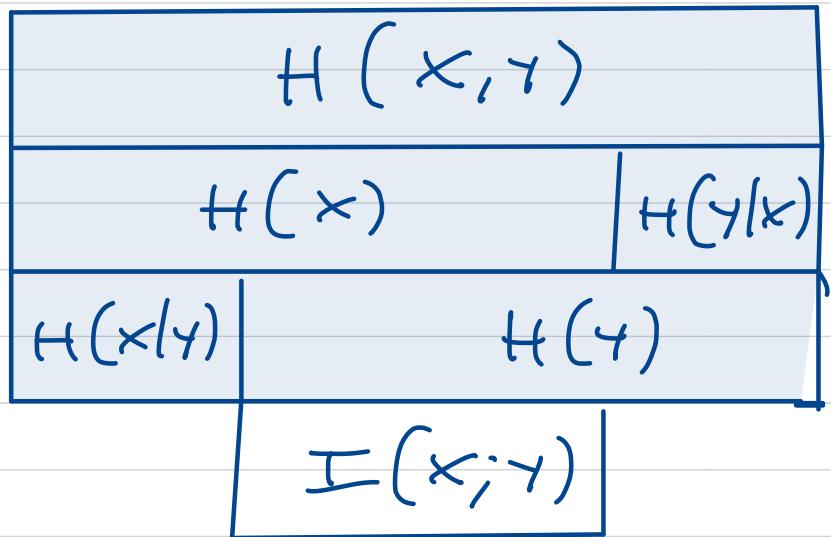
$$I(x; y) = H(Y) - H(Y|x)$$

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

SYMMETRIC !!

## IMPORTANT

## FIGURE



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

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

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

## INFORMATION GAIN RATIO

$$0 \leq I(x; y) = H(x) - H(x|y) \leq H(x)$$

$$0 \leq I'(x; y) = \frac{I(x; y)}{H(x)} \leq 1$$

FOR THE DATA SET OF ASSIGNMENT 2

$x_1 \quad x_2 \quad x_3 \quad | \quad x_R$  TARGET

$$\underset{1 \leq i \leq 3}{\overline{I}(x; x_i)} = \frac{\overline{I}(x; x_i)}{H(x_i)} = \frac{H(x) - H(x|x_i)}{H(x_i)}$$

# KULLBACK - LEIBLER DIVERGENCE

$$X \quad \mathcal{N}_X = \{ n_1, n_2, \dots, n_n \}$$

$$P = \{ p_1, p_2, \dots, p_n \}$$

$$Q = \{ q_1, q_2, \dots, q_n \}$$

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

IT'S USED TO COMPARE 2 PROB. DISTRIBUTIONS

TYPICALLY  $P \equiv$  OBSERVED, MEASURED

$Q \equiv$  ANALYTICAL FIT

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

IF  $p_i = q_i$  THEN  $D_{KL}(P||Q) = 0$

HOT A DISTANCE

$$D_{KL}(P||Q) = D_{KL}(Q||P)$$

$$D(P|Q) \geq 0$$

PROOF

$D \geq 0$

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

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

$$\leq \sum_i p_i \left( \frac{q_i}{p_i} - 1 \right) = \sum_i q_i - \sum_i p_i = \\ = 1 - 1 = 0$$

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

UNIFORM DISTR.

OVER 10 ELEM

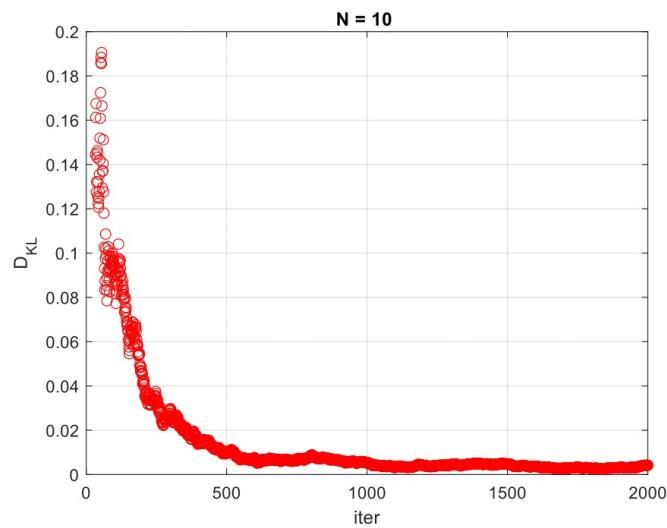
$$q_i = \frac{1}{10} \sum_{j=1}^{10} \theta_{ij}$$

COMPUTATION

OF  $D_{KL}(P||Q)$

FOR INCREASING

# OF ITERATIONS



LINK BETWEEN  $I(x; y)$  AND D

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

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

$$= D(P(x, y) \mid P(x)P(y))$$

## COMMON INTERPRETATION

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

$P_i$  = OBSERVED PROB. PSTR.

$Q_i$  = MODEL/ANALYTICAL PSTR.

$D$  = INF. GAIN WHEN WE USE MEASURED

$D$  = INF. LOSS  
WHEN WE USE  $P_i$  INSTEAD OF MODEL  $Q_i$

MODEL  $Q_i$  INSTEAD OF  
MEASURED  $P_i$

# Information Theory for Data Science

Assignment 1

Introduction to Information Theory and application to Classifiers

Draft version 0.2

## **Exercises:**

1. Entropy of a random variable with 3 outcomes (pt. 3)
2. Entropy of a random variable from a data series (pt. 4)
3. Application of the principle of maximum entropy (pt. 4)
  
4. Exercise on Information Gain Ratio (pt. 4)
5. Kullback Leibler divergence from empirical distribution (pt. 3)
6. Permutation entropy for anomaly detection (pt. 5)
  
7. Information Gain Ratio and Classifier (pt. 7) (not yet available)

## **Exercise 1 - Entropy of a random variable with 3 outcomes (pt. 3)**

## **Exercise 1 - Entropy of a random variable with 3 outcomes**

1. Given a random variable with 3 outcomes, write a program to plot the entropy as a function of all possible probability vectors
2. Start with a probability vector where one of the elements is significantly higher than the others. Apply an iterative averaging procedure (for example, replace each element with the average of itself and its neighbors, followed by normalization). For each updated vector, compute the entropy and plot its value on the figure generated in step 1. Show that as the probability distribution approaches the uniform distribution, the entropy approaches its maximum value. Finally, discuss the results.

## **Exercise 2 - Entropy of a random variable from a data series (pt. 4)**

## **Exercise 2 - Entropy of a random variable from a data series**

1. Identify a data series and estimate the probabilities of the outcomes based on their occurrences, updating the probabilities at each time step.
2. At each time step, compute the entropy, plot its behavior, and discuss the results

Note: In the presentation, include a link to the source of the data series

## **Exercise 3 - Application of the principle of maximum entropy (pt. 4)**

### **Exercise 3.a**

1. Invent an exercise where you have a random variable  $X$  with an alphabet  $\Omega_X$  with 2 outcomes with integer values.
2. Show some examples of the probability distribution  $P(X)$  for different values of the mean value.
3. Discuss the results

## Exercise 3.b

1. Invent an exercise where you have a random variable  $X$  with alphabet  $\Omega_X$  with at least 4 outcomes, where each outcome has an integer value (“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 4 – Exercise on Information Gain Ratio (pt. 4)**

## Exercise 4 – Exercise on Information Gain Ratio

1. Invent an exercise based on a data set like this one.

$x_1$

$x_2$

$x_3$



Athlete	Training Hours	Rest Hours	Gym Workouts	Performance
A	High	Low	Low	Lose
B	Medium	High	Medium	Win
C	Low	Medium	High	Lose
D	High	Medium	Medium	Win
E	Medium	High	Low	Win
F	Low	Low	High	Lose
G	High	Low	Medium	Win
H	Medium	Low	Low	Lose
I	Low	High	Medium	Win
J	High	High	High	Win

The last column represents your target variable X.

$x_1 \quad x_2 \quad x_j$

2. Compute the Information Gain Ratio with respect to all the other variables and select the one that provides more information about X.

$$\text{IGR}(x; x_i) = \frac{I(x; x_i)}{H(x_i)} = \frac{H(x) - H(x|x_i)}{H(x_i)}$$

## **Exercise 5 - Kullback-Leibler distance from empirical distribution (pt. 3)**

## Exercise 5 - Kullback-Leibler distance from empirical distribution

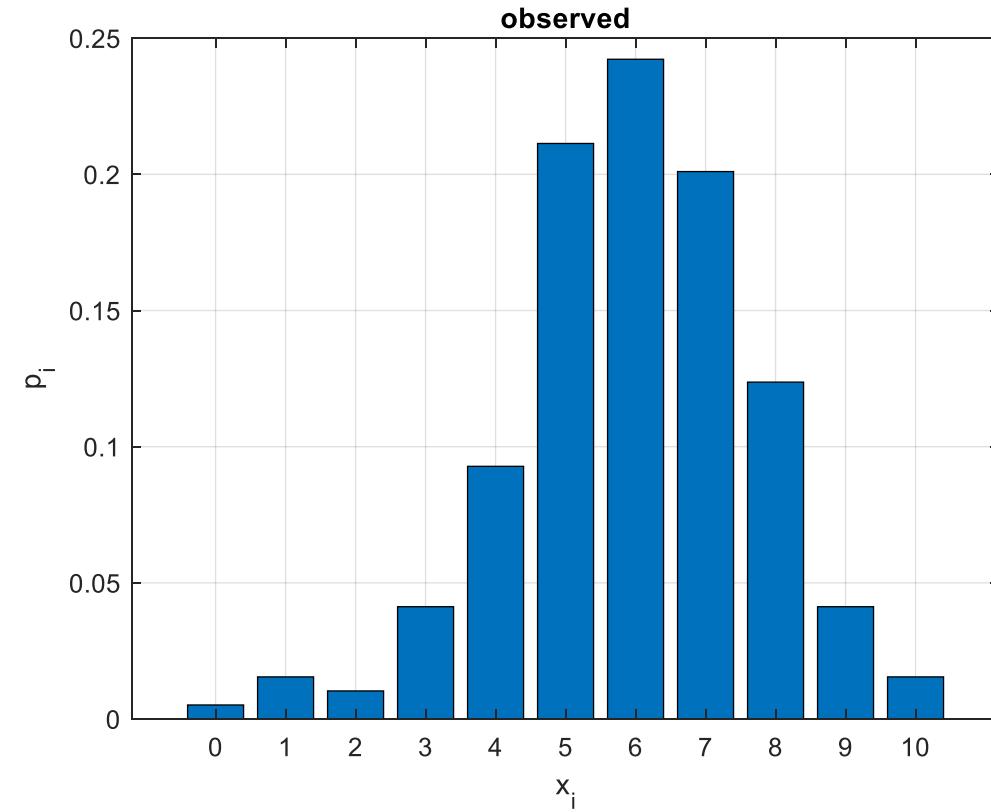
Consider the observed data

$$l \sum x_i l = 11$$



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

$$\text{Number of observed } x_i = [1 \ 3 \ 2 \ 8 \ 18 \ 41 \ 47 \ 39 \ 24 \ 8 \ 3]$$



$$\overbrace{\pi}^l$$

Compare with uniform pmf

1. Plot the two observed and the uniform pmfs
2. In the title write the KL divergence value

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

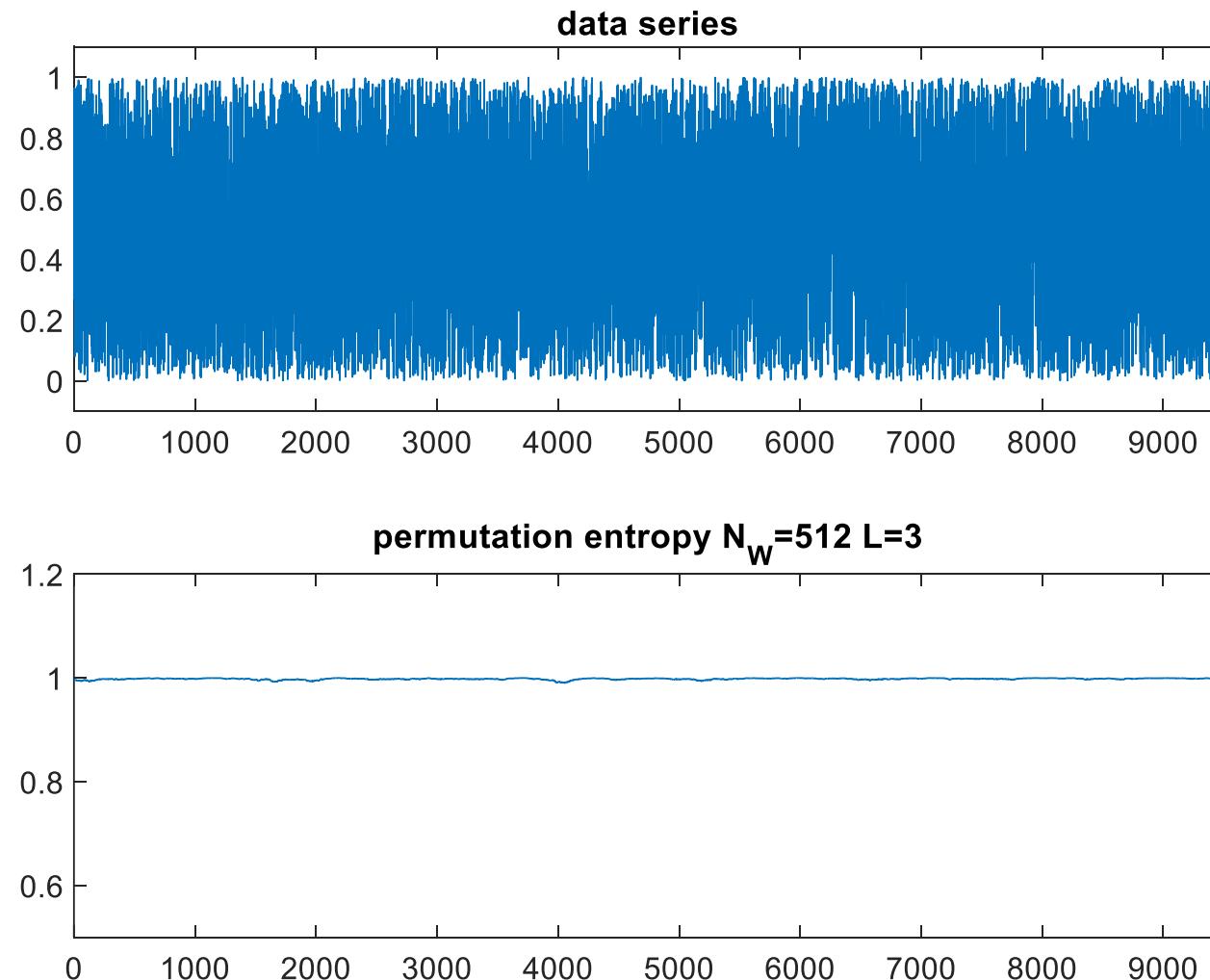
3. Identify the pmf at minimum KL divergence
4. Plot the observed and the best binomial pmfs
5. In the title write the value of  $p$  and KL divergence

## **Exercise 6 - Permutation Entropy for Time Series Anomaly Detection (pt. 5)**

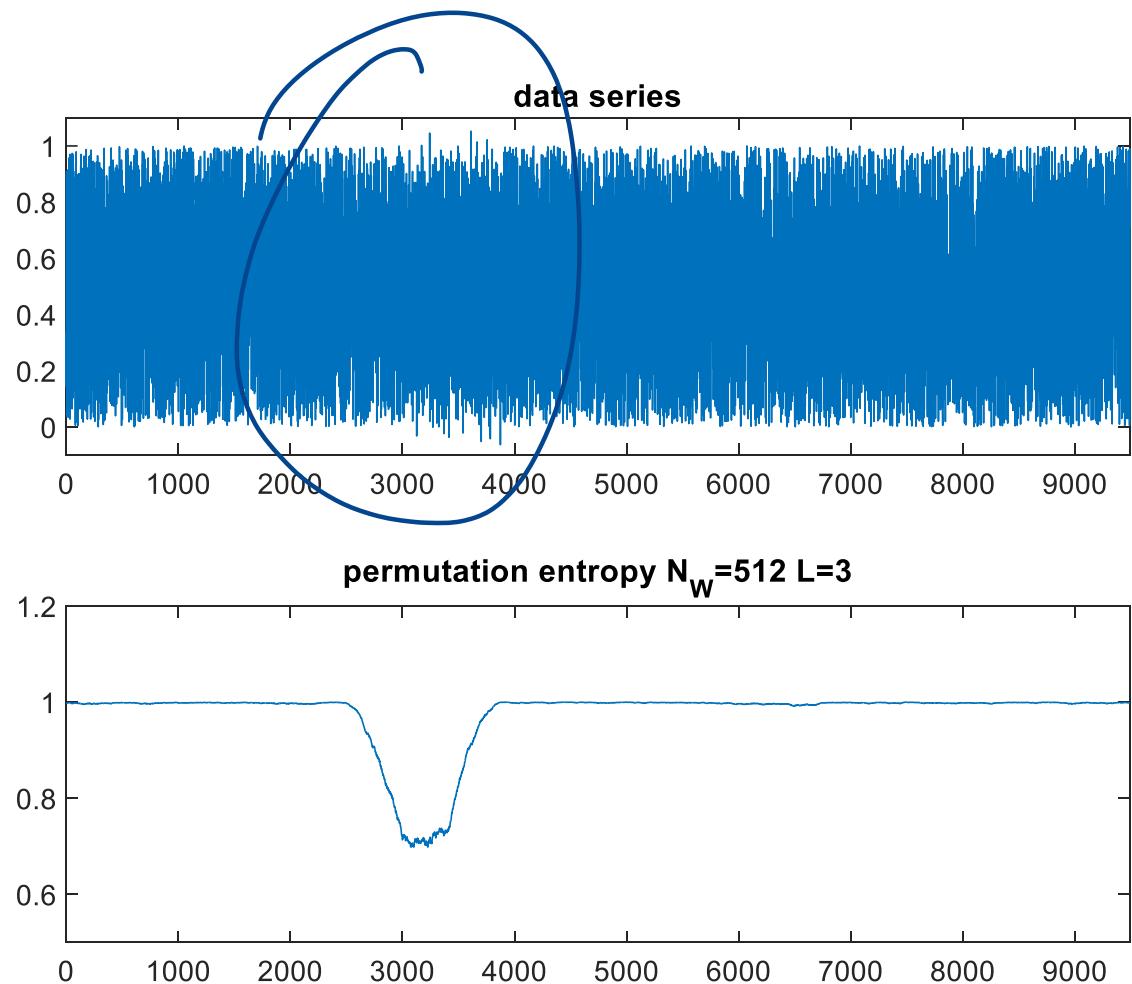
## Exercise 6 - Permutation Entropy for Time Series Anomaly Detection

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  $m$  used before.
- Plot the data and the entropy.
- Comment the results



## Some Matlab functions that might be useful

sort  
perms

## **Exercise 7 - Information Gain Ratio and Classifier (pt. 7)**

# Important

Final version assigned on 08/10/2024 (**to be confirmed**)

Delivery by

- 23/10/2024, 11.59 PM: +2 points
- 30/10/2024, 11.59 PM: +1 point
- 21/02/2025, 11.59 PM: 0 points
- **Later: not accepted**

## 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	7	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
---	---	---	---	---	---

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

$P_1 \ P_2 \ P_3 \ P_4 \ P_5 \ P_6$

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

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

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

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

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

$\alpha_1 \alpha_2 \alpha_3 \alpha_4 \alpha_5 \alpha_6 \alpha_7 \alpha_8 \alpha_9 \alpha_{10}$

•  $N = 10$

$\alpha_1 \alpha_2 \alpha_3 \alpha_4 \alpha_5 \alpha_6 \alpha_7 \alpha_8 \alpha_9 \alpha_{10}$

•  $N_w = 5$

$\alpha_1 \alpha_2 \alpha_3 \alpha_4 \alpha_5 \alpha_6 \alpha_7 \alpha_8 \alpha_9 \alpha_{10}$

•  $m = 3$

$\alpha_1 \alpha_2 \alpha_3 \alpha_4 \alpha_5 \alpha_6 \alpha_7 \alpha_8 \alpha_9 \alpha_{10}$

$\alpha_1 \alpha_2 \alpha_3 \alpha_4 \alpha_5 \alpha_6 \alpha_7 \alpha_8 \alpha_9 \alpha_{10}$

$\alpha_1 \alpha_2 \alpha_3$

$\alpha_2 \alpha_3 \alpha_4$

$\alpha_3 \alpha_4 \alpha_5$

$\rightarrow 1 \text{ VALUE OF P.E.}$

$\alpha_2 \alpha_3 \alpha_4$

$\alpha_3 \alpha_4 \alpha_5$

$\alpha_4 \alpha_5 \alpha_6$

$\rightarrow 1 \text{ VALUE OF P.E.}$

:

:

:

$N - N_w + 1$  VALUES

OF P.E.

$\alpha_6 \alpha_7 \alpha_8$

$\alpha_7 \alpha_8 \alpha_9$

$\alpha_8 \alpha_9 \alpha_{10}$

Some Matlab functions that might be useful

sort  
perms

$M = \text{MATRIX}$

$\begin{matrix} z_0 & z_1 & z_2 & \dots \\ z_1 & z_2 & z_3 & \dots \\ z_2 & z_3 & z_4 & \dots \end{matrix}$

$[B, P] = \text{sort}(M)$

$\uparrow$

CONTAINS PERMUTAT.