# Information and Entropy

- Finite-alphabet random variable
- Information and entropy
- Joint and conditional entropy
- Mutual information
- Relative entropy
- Data processing inequality
- Asymptotic equipartition property
- Typical sets



## Entropy of a Random Variable

Consider a discrete, finite-alphabet random variable X

$$\mathcal{A}_X = \{\alpha_0, \alpha_1, \dots, \alpha_{K-1}\}$$
$$f_X(x) = P(X = x) \quad \forall x \in \mathcal{A}_X$$

"Information" associated with the event X=x

$$h_X(x) = -\log_2 f_X(x)$$

• "Entropy of X" is the <u>expected value</u> of that information  $H(X) = E\{h_X(X)\} = -\sum_{x \in \mathcal{A}_X} f_X(x) \log_2 f_X(x)$ 

Unit: bits

Markus Flierl: EQ2845 Information Theory and Source Coding Information & Entropy no. 2

# Information and Entropy: Properties

- Information  $h_X(x) \ge 0$
- Information  $h_X(x)$  strictly increases with decreasing probability  $f_X(x)$
- Boundedness of entropy

$$\begin{array}{c|c} 0 \leq H(X) \leq \log_2\left(|\mathcal{A}_X|\right) \\ \hline \\ \text{equality if only one} \\ \text{outcome can occur} \end{array} \qquad \begin{array}{c} 0 \leq H(X) \leq \log_2\left(|\mathcal{A}_X|\right) \\ \hline \\ \text{equality if all outcomes} \\ \text{are equally likely} \end{array}$$

 Very likely and very unlikely events do not substantially change entropy

$$-p \log_2 p \rightarrow 0 \quad \text{for } p \rightarrow 0 \text{ or } p \rightarrow 1$$



#### **Example: Binary Random Variable**

$$H(X) = -p \log_2 p - (1-p) \log_2 (1-p)$$





Markus FlierI: EQ2845 Information Theory and Source Coding Information & Entropy no. 4

## Joint Entropy

Consider random vectors (with discrete, finite-alphabet components)

$$\mathbf{X} = (X_0, X_1, \dots, X_{m-1})$$

Entropy

$$H(\mathbf{X}) = E\left[-\log_2 f_{\mathbf{X}}(\mathbf{X})\right] = E\left[h_{\mathbf{X}}(\mathbf{X})\right]$$

Long-hand

$$H(\mathbf{X}) = H(X_0, X_1, ..., X_{m-1})$$
  
=  $-\sum_{x_0} \sum_{x_1} ... \sum_{x_{m-1}} f_{\mathbf{X}}(x_0, x_1, ..., x_{m-1}) \log_2 f_{\mathbf{X}}(x_0, x_1, ..., x_{m-1})$ 



## **Conditional Entropy**

Consider two discrete finite-alphabet r.v. X and Y

$$H(X|Y) = E\left[-\log_2 f_{X|Y}(x,y)\right] = -\sum_y \sum_x f_{X,Y}(x,y)\log_2 f_{X|Y}(x,y)$$
$$= -\sum_y f_Y(y)\sum_x f_{X|Y}(x,y)\log_2 f_{X|Y}(x,y)$$

- Conditional entropy H(X|Y) is average additional information, if Y is already known
- Joint entropy:  $H(X,Y) = E\left[-\log_2 f_{X,Y}(X,Y)\right]$   $= E\left[-\log_2 \left(f_Y(y)f_{X|Y}(X,Y)\right)\right]$   $= E\left[-\log_2 f_Y(y)\right] + E\left[-\log_2 f_{X|Y}(X,Y)\right]$ = H(Y) + H(X|Y) (Chain rule)



Markus FlierI: EQ2845 Information Theory and Source Coding Information & Entropy no. 6

## **Mutual Information**

- "Mutual information" is the average information that random variables X and Y convey about each other
  - Reduction in uncertainty about x, if y is observed
  - Reduction in uncertainty about y, if x is observed

$$I(X;Y) = H(X) - H(X | Y) = H(Y) - H(Y | X)$$
  
=  $\sum_{x} \sum_{y} f_{X,Y}(x, y) \log_2 \frac{f_{X,Y}(x, y)}{f_X(x) f_Y(y)}$ 

• Properties  $0 \le I(X;Y) = I(Y;X)$  $I(X;Y) \le H(X)$  $I(X;Y) \le H(Y)$ 



Markus FlierI: EQ2845 Information Theory and Source Coding Informat

Information & Entropy no. 7

### **Mutual Information**



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



Information & Entropy no. 8

# **Conditioning Reduces Entropy**

Property of mutual information

$$0 \le I(X;Y) = H(X) - H(X|Y)$$

 Knowing another random variable Y can only reduce the uncertainty in X

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

Equality if and only if X and Y are independent

## Chain Rule for Entropy

$$H(X_1, X_2) = H(X_1) + H(X_2|X_1)$$

#### $H(X_1, X_2, X_3) = H(X_1) + H(X_2, X_3 | X_1)$ = $H(X_1) + H(X_2 | X_1) + H(X_3 | X_2, X_1)$ :

$$H(X_1, X_2, \dots, X_n) = \sum_{i=1}^n H(X_i | X_{i-1}, \dots, X_1)$$

Markus FlierI: EQ2845 Information Theory and Source Coding Information & Entropy no. 10

# Joint Entropy and Statistical Dependency

Theorem (Independence bound on entropy)

$$H(X_0, X_1, X_2, ..., X_{m-1}) \leq H(X_0) + H(X_1) + ... + H(X_{m-1})$$

$$\boxed{\qquad}$$
Equality for statistical independence of  $X_0, X_1, X_2, ..., X_m$ 

- Exploiting statistical dependencies can reduce bit-rate
- Statistically independent components can be compressed and decompressed separately without loss



## **Conditional Mutual Information**

Mutual information between X and Y given Z

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

Chain rule for information

$$I(X_1, X_2, \dots, X_n; Y) = \sum_{i=1}^n I(X_i; Y | X_{i-1}, \dots, X_1)$$



## **Relative Entropy**

Kullback Leibler distance between two PMFs p(x) and q(x)

$$H(p \parallel q) = \sum_{x} p(x) \log_2 \frac{p(x)}{q(x)}$$

In general

$$H(p \parallel q) \neq H(q \parallel p)$$

Information inequality

 $H(p \parallel q) \geq 0$ 

With equality if and only if p(x)=q(x) for all x.



## **Relative Entropy and Mutual Information**

 Mutual information is the relative entropy between the joint distribution and the product distribution

$$I(X;Y) = H(f_{X,Y} \parallel f_X f_Y)$$

Non-negativity of mutual information

 $I(X;Y) \geq 0$ 

• With equality if and only if X and Y are independent.

## **Data Processing Inequality**

- No clever manipulation of the data can improve the inferences that can be made from the data
- X, Y, Z are said to form a Markov chain in that order

 $X \to Y \to Z$ 

if the joint PMF can be written as

$$f_{X,Y,Z} = f_X f_{Y|X} f_{Z|Y}$$

Note, X and Z are conditionally independent given Y

$$f_{X,Z|Y} = f_{X|Y} f_{Z|Y}$$

## Data Processing Inequality

• Chain rule for mutual information I(X;Y,Z) = I(X;Z) + I(X;Y|Z) = I(X;Y) + I(X;Z|Y)

• If  $X \to Y \to Z$ , X and Z are conditionally independent given Y

$$I(X;Z|Y) = 0$$

• Non-negativity of mutual information:  $I(X; Y|Z) \ge 0$ 

• Data processing inequality: If 
$$X \to Y \to Z$$
  
 $I(X;Y) > I(X;Z)$ 

## Remark: Weak Law of Large Numbers

 Law of Large Numbers: The average of the results obtained from a large number of trials is close to the expected value

$$\frac{1}{n}(U_1 + U_2 + \dots + U_n) \xrightarrow{p} E[U] \quad \text{for} \quad n \to \infty$$

Convergence in probability (weak law)

$$\lim_{n \to \infty} \Pr\left( \left| \frac{1}{n} \sum_{i=1}^{n} U_i - E[U] \right| < \epsilon \right) = 1 \quad \text{for} \quad \epsilon > 0$$



## Asymptotic Equipartition Property

• Let  $X_1$ ,  $X_2$ , ... be i.i.d. with  $f_X$ 

$$-\frac{1}{n}\log_2 f_{X_1,X_2,\dots,X_n} = -\frac{1}{n}\sum_{i=1}^n \log_2 f_{X_i}$$
$$\stackrel{p}{\to} E[-\log_2 f_X]$$
$$= H(X)$$

Asymptotic Equipartition Property (AEP): If X<sub>1</sub>, X<sub>2</sub>, ... are i.i.d. with f<sub>X</sub>, then

$$-\frac{1}{n}\log_2 f_{X_1,X_2,\dots,X_n} \xrightarrow{p} H(X)$$



# **Typical Set**

The typical set A<sub>ε</sub><sup>(n)</sup> with respect to f<sub>X</sub> is the set of sequences (x<sub>1</sub>, x<sub>2</sub>, ..., x<sub>n</sub>) ∈ X<sup>n</sup> with the property

$$2^{-n[H(X)+\epsilon]} \le f_{X_1,X_2,\dots,X_n} \le 2^{-n[H(X)-\epsilon]}$$

The typical set has probability nearly 1

$$\Pr\left\{A_{\epsilon}^{(n)}\right\} > 1 - \epsilon$$

- All elements of the typical set are nearly equiprobable
- The number of elements in the typical set is nearly 2<sup>nH</sup>

$$\left| A_{\epsilon}^{(n)} \right| \le 2^{n[H(X) + \epsilon]}$$



Represent sequences in  $\mathcal{X}^n$  using nH(X) bits on average.



Markus Flierl: EQ2845 Information Theory and Source Coding Information & Entropy no. 20