

ECE253/CSE208 Introduction to Information Theory

Lecture 6: Entropy Rate

Dr. Yu Zhang

ECE Department

University of California, Santa Cruz

- Chap 4 of *Elements of Information Theory (2nd Edition)* by Thomas Cover & Joy Thomas

Markov Chains

Consider a discrete-time Markov chain X_1, X_2, \dots, X_{n+1} , we have

$$\begin{aligned}\Pr(X_{n+1} = x_{n+1} | X_n = x_n, \dots, X_1 = x_1) &= \Pr(X_{n+1} = x_{n+1} | X_n = x_n) \\ p(x_{n+1}, \dots, x_1) &= p(x_1)p(x_2|x_1) \dots p(x_{n+1}|x_n)\end{aligned}$$

All knowledge of the past states is embedded in the current state.

A Markov chain of order k (the future state depends on the past k states):

$$p(x_{n+1} | x_n, \dots, x_1) = p(x_{n+1} | x_n, \dots, x_{n-k+1})$$

Time-homogeneous Markov chain: The transition probability is independent of time.

That is, for all n

$$\Pr(X_{n+1} = j | X_n = i) = \Pr(X_n = j | X_{n-1} = i), \quad \forall i, j \in \mathcal{S} := \{1, 2, \dots, m\}.$$

Let matrix $\mathbf{P}_{m \times m}$ contain all transition probabilities, whose (i, j) -th entry is

$$P_{ij} = \Pr(X_{n+1} = j | X_n = i), \quad \forall i, j \in \mathcal{S} := \{1, 2, \dots, m\}. \quad (1)$$

Three things in a Markov chain: A sequence of random variables (the chain), a state space (from which the random variables take values), and the rules for transition (the transition probability matrix).

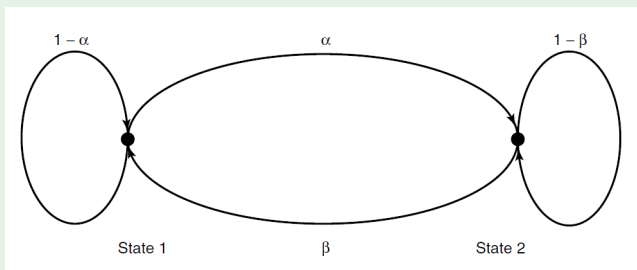
Two-state Markov chain

Example

(Two-state Markov chain).

$$\mathbf{P} = \begin{bmatrix} 1 - \alpha & \alpha \\ \beta & 1 - \beta \end{bmatrix}$$

As shown in the figure below:



Q: Given $\Pr(X_n = i)$, find $\Pr(X_{n+1} = j)$, $\forall i, j \in \{1, 2, \dots, m\}$.

A: $\Pr(X_{n+1} = j) = \sum_i \Pr(X_n = i, X_{n+1} = j) = \sum_i \Pr(X_n = i) \Pr(X_{n+1} = j | X_n = i)$.

Stationary Distribution

Define a row vector to collect all state probabilities at time $n + 1$

$$\boldsymbol{\pi}^{(n+1)} = \left[\pi_1^{(n+1)}, \pi_2^{(n+1)}, \dots, \pi_m^{(n+1)} \right],$$

where $\pi_j^{(n+1)} = \Pr(X_{n+1} = j)$, $\forall j \in \{1, 2, \dots, m\}$. Hence, we have

$$\boldsymbol{\pi}^{(n+1)} = \boldsymbol{\pi}^{(n)} \mathbf{P}.$$

Definition (Stationary/steady-state/invariant/equilibrium distribution)

Stationary distribution of a Markov chain: $\boldsymbol{\pi} \mathbf{P} = \boldsymbol{\pi}$ and $\mathbf{0} \leq \boldsymbol{\pi} \leq \mathbf{1}$, $\boldsymbol{\pi} \mathbf{1} = 1$, where $\mathbf{1}$ is the all-ones column vector with an appropriate dimension.

Hence, stationary distribution $\boldsymbol{\pi}$ is a *fixed point* of the transformation represented by \mathbf{P} , which is a *left eigenvector* of \mathbf{P} corresponding to eigenvalue 1.

Example

Find the stationary distribution of the aforementioned example of the two-state MC.

$$\begin{cases} \boldsymbol{\pi} \mathbf{P} = \boldsymbol{\pi} \\ \boldsymbol{\pi} \mathbf{1} = 1 \end{cases} \Rightarrow \boldsymbol{\pi} = \left[\frac{\beta}{\alpha + \beta}, \frac{\alpha}{\alpha + \beta} \right].$$

Classification of States

Definition (Accessible and Communicate)

- State j is said to be **accessible** from state i if $P_{ij}^{(n)} > 0$ for some $n \geq 0$, which is denoted as $i \rightarrow j$.
- Two states i and j are said to **communicate** if they are accessible from each other, which is denoted as $i \leftrightarrow j$. (i.e., there are directed paths between i and j).

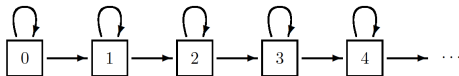


Figure: Each state is accessible from all its previous states, but not vice versa.

Communicate is an *equivalence relation*, meaning that

- Reflexivity: $i \leftrightarrow i$
- Symmetry: If $i \leftrightarrow j$, then $j \leftrightarrow i$
- Transitivity: If $i \leftrightarrow j$ and $j \leftrightarrow k$, then $i \leftrightarrow k$

An equivalence relation divides the state space into disjoint classes of equivalent states that is called **communication classes**.

Irreducible MC

Definition (Irreducible MC: every state can be reached from every other state)

It is possible to go with positive probability from any state to any other states in a finite number of steps. That is, $\exists n < \infty$, $\Pr(X_n = j | X_0 = i) = P_{ij}^{(n)} > 0$, $\forall i, j \in \mathcal{S}$.

- A Markov chain is irreducible iff all states belong to one communication class; i.e., all states communicate with each other.
- A Markov chain is reducible iff if there are two or more communication classes.
- A finite Markov chain is irreducible iff its transition graph is strongly connected (there is a path between any pair of two vertices).

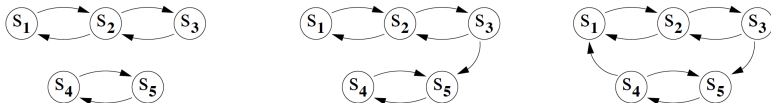


Figure: The first two Markov chains are reducible while the last one is irreducible.

Aperiodic MC

Definition (Aperiodic)

The period of a state i is defined as $k = \gcd\{n > 0, P_{ii}^{(n)} > 0\}^a$. That is, if $P_{ii}^{(n)} = 0$ when n is not a multiple of k and k is the greatest integer with this property. If $k = 1$, the state is said to be aperiodic. **A Markov chain is aperiodic if every state is aperiodic.**

^a**gcd** is the greatest common divisor; e.g., $\gcd\{6, 8, 10, \dots\} = 2$; $\gcd\{3, 5, 7, \dots\} = 1$.

- All states in the same communication class have the same period.
- An irreducible MC only needs one aperiodic state to imply the chain is aperiodic.

Consider a finite irreducible Markov chain:

- If there is a self transition in the diagram, then the chain is aperiodic.
- Suppose $P_{ii}^{(\ell)} > 0$ and $P_{ii}^{(m)} > 0$. If ℓ and m are co-prime ($\gcd(\ell, m) = 1$), then state i is aperiodic.

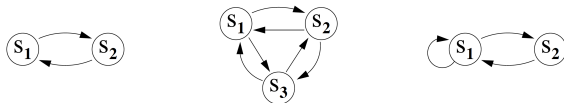


Figure: The first one has period 2 while the last two are aperiodic.

An Exercise¹

Question: Find all communication classes and their periods.

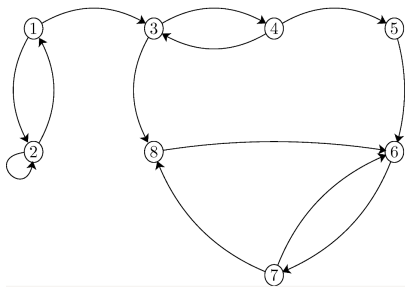


Figure: A state transition diagram.

Answer:

- Class 1 = {1, 2}, aperiodic
- Class 2 = {3, 4}, period = 2
- Class 3 = {5}, period = 0 (transient state)
- Class 4 = {6, 7, 8}, aperiodic

¹https://www.probabilitycourse.com/chapter11/11_2_4_classification_of_states.php

Unique Stationary Distribution and Limiting Distribution

Theorem

For an irreducible, aperiodic, and finite-state Markov chain, there exists a finite integer N such that $P_{ij}^{(n)} > 0$, for all $i, j \in S$ and all $n \geq N$.

Theorem

An irreducible and aperiodic finite-state Markov chain has a unique stationary distribution.

Lemma

For an irreducible, aperiodic, and finite-state Markov chain, any initial distribution converges to the unique stationary distribution as $n \rightarrow \infty$.

- Assuming irreducibility, the stationary distribution is always unique if it exists, and its existence can be implied by positive recurrence of all states.
- The stationary distribution has the interpretation of the limiting distribution when the chain is irreducible and aperiodic.

Motivating Question

Shannon used Markov Chain to describe English texts:

- From 0th-order to 3rd-order letters; 1st-order, 2nd-order words

3. THE SERIES OF APPROXIMATIONS TO ENGLISH

To give a visual idea of how this series of processes approaches a language, typical sequences in the approximations to English have been constructed and are given below. In all cases we have assumed a 27-symbol 'alphabet,' the 26 letters and a space.

1. Zero-order approximation (symbols independent and equiprobable).

XFOML RXKHRJFFJUJ ZLPWCFWKCYJ FFJEYVKQSGHYD QPAAMKBZAACIBZLHJQD.

2. First-order approximation (symbols independent but with frequencies of English text).

OCRO HLI RGWR NMIELWIS EU LL NBNESEBYA TH EEI ALHENHTTPA OOBTTVA NAH BRL.

3. Second-order approximation (digram structure as in English).

ON IE ANTSOUTINYS ARE T INCTORE ST BE S DEAMY ACHIN D ILONASIVE TU- COOWE AT TEASONARE FUSO TIZIN ANDY TOBE SEACE CTISBE.

4. Third-order approximation (trigram structure as in English).

IN NO IST LAT WHEY CRATICT FROURE BIRS GROCID PONDENOME OF DEMONSTURES OF THE REPTAGIN IS REGOACTONA OF CRE.

5. First-order word approximation. Rather than continue with tetragram, . . . , n -gram structure it is easier and better to jump at this point to word units. Here words are chosen independently but with their appropriate frequencies.

REPRESENTING AND SPEEDILY IS AN GOOD APT OR COME CAN DIFFERENT NATURAL HERE HE THE A IN CAME THE TO OF TO EXPERT GRAY COME TO FURNISHES THE LINE MESSAGE HAD BE THESE.

6. Second-order word approximation. The word transition probabilities are correct but no further structure is included.

THE HEAD AND IN FRONTAL ATTACK ON AN ENGLISH WRITER THAT THE CHARACTER OF THIS POINT IS THEREFORE ANOTHER METHOD FOR THE LETTERS THAT THE TIME OF WHO EVER TOLD THE PROBLEM FOR AN UNEXPECTED.

Figure: In Section 3 of the paper “A Mathematical Theory of Communication (1948)”: Shannon asked “Can we define a quantity which will measure, in some sense, how much information is “produced” by such a process, or better, at what rate information is produced?”

Stationary Stochastic Process

Discrete-time Information Sources:

- Communications take place continually rather than a finite period of time; e.g., Internet cellular networks, radio stations, TV programs, etc.
- The info source can be modeled as a discrete-time stochastic process $\{X_k\}_{k=1}^{\infty}$

Definition (Stationary stochastic process)

A stochastic process $\{X_k\}$ is *strongly stationary* if the joint distribution of any subset of the sequence is invariant w.r.t. any time shifts. That is,

$$\Pr\{X_1 = x_1, X_2 = x_2, \dots, X_n = x_n\} = \Pr\{X_{1+l} = x_1, X_{2+l} = x_2, \dots, X_{n+l} = x_n\}$$

for every n , every shift l , and for all $x_1, \dots, x_n \in \mathcal{X}$.

Time-Homogeneous vs Stationary Markov Chains²

Q: Is a time-homogeneous Markov chain a stationary process? **A:** Not necessarily.

Example

Consider an MC with $\mathbf{P} = \begin{pmatrix} 1/2 & 1/2 \\ 1/2 & 1/2 \end{pmatrix}$ and the initial probability state $\pi^{(0)} = (1, 0)$, then $\pi^{(1)} = \pi^{(0)} \cdot \mathbf{P} = (1/2, 1/2) \neq \pi^{(0)}$. So, the chain is not stationary.

- *Time-homogeneous MC:*

$$\Pr(X_{n+1} = j | X_n = i) = \Pr(X_n = j | X_{n-1} = i), \forall n.$$

- *Stationary MC:* $\Pr(X_0 = x_0, \dots, X_n = x_n) = \Pr(X_l = x_0, \dots, X_{n+l} = x_n), \forall n, l.$

1. Every stationary chain is time-homogeneous (proved by Bayes' rule).
2. A time-homogeneous MC is stationary iff the distribution of X_0 is a stationary distribution of the MC.

²https://en.wikipedia.org/wiki/Markov_chain

Entropy Rate for Stochastic Processes

Q: How does the entropy of a sequence grow with n ?

A: We define the entropy rate as the rate of growth:

Definition (Per symbol entropy of n random variables)

$$H(\mathcal{X}) := \lim_{n \rightarrow \infty} \frac{1}{n} H(X_1, X_2, \dots, X_n)$$

when the limit exists.

Special cases:

- $\{X_i\}$ are i.i.d.: $H(\mathcal{X}) = H(X_1)$.
- $\{X_i\}$ are independent but not identical: $H(\mathcal{X}) = \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n H(X_i)$. But, the limit may not even exist; e.g, $H(X_i) = i$.

Definition (Conditional entropy of the last random variable given the past.)

$$H'(\mathcal{X}) := \lim_{n \rightarrow \infty} H(X_n | X_{n-1}, \dots, X_1)$$

Entropy Rate for Stationary Stochastic Processes

Theorem (Entropy rate)

For a stationary stochastic process, $H(\mathcal{X}) = H'(\mathcal{X})$.

Proof: We first show that $H'(\mathcal{X})$ is well-defined. Due to stationarity, we have

$$0 \leq \underbrace{H(X_n | X_{n-1}, \dots, X_1)}_{a_n} \leq H(X_n | X_{n-1}, \dots, X_2) = \underbrace{H(X_{n-1} | X_{n-2}, \dots, X_1)}_{a_{n-1}}$$

Hence, sequence $\{a_n\}$ is monotonically non-increasing and lower bounded (by zero). It must converge $\lim_{n \rightarrow \infty} H(X_n | X_{n-1}, \dots, X_1) =: H'(\mathcal{X})$. Recall that

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

By the lemma below, the RHS converges to $H'(\mathcal{X})$. So does the LHS, which is $H(\mathcal{X})$ in the limit.

Lemma (Cesàro mean)

If a sequence $\{a_n\} \rightarrow c$, its running average $\{b_n := \frac{1}{n} \sum_{i=1}^n a_i\} \rightarrow c$.

Entropy Rate for Stationary Stochastic Processes (cont'd)

Lemma

For a stationary Markov chain, $H'(\mathcal{X}) = H(X_2|X_1)$.

Proof: $H'(\mathcal{X}) = \lim_{n \rightarrow \infty} H(X_n|X_{n-1}, X_{n-2}, \dots, X_1) = \lim_{n \rightarrow \infty} H(X_n|X_{n-1}) = H(X_2|X_1)$.

Theorem

Let $\{X_i\}$ be a stationary Markov chain with stationary distribution μ and transition probability matrix \mathbf{P} . If $X_1 \sim \mu$, then the entropy rate $H(\mathcal{X}) = -\sum_{ij} \mu_i P_{ij} \log P_{ij}$

Proof: $H(\mathcal{X}) = H(X_2|X_1) = \sum_i \Pr(X_1 = i) H(X_2|X_1 = i) = \sum_i \mu_i \sum_j P_{ij} \log P_{ij}^{-1}$.

Note: For an irreducible, aperiodic, and finite MC, any initial distribution converges to the stationary distribution. So, even though we may not start from the stationary distribution, the entropy rate $H(\mathcal{X})$ given above is still correct.

The entropy rate of a stationary Markov chain is not dependent on its initial distribution, but only on the transitions between the states and the stationary distribution.

Entropy Rate of Random Walk over Graph

Example (Random walk over a weighted graph)

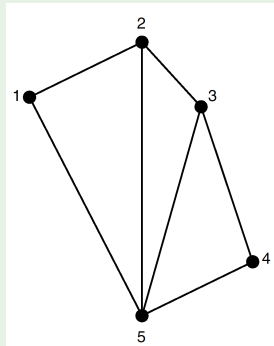


Figure: A weighted graph $G(\mathcal{N}, \mathcal{E}, \mathcal{W})$.

- $w_{ij} = w_{ji}$ denotes the edge weight between nodes i and j (0 if no edges).
- Given $X_n = i$, the probability of moving from node i to j is $P_{ij} = \frac{w_{ij}}{\sum_k w_{ik}} = \frac{w_{ij}}{w_i}$, where $w_i := \sum_k w_{ik}$ is the total weight of all edges connecting with node i .
- Intuitively, the stationary distribution of any node $i \in \mathcal{N}$ should be proportional to its degree w_i , which can be derived as $\pi_i = \frac{w_i}{2w}$, where $w \triangleq \sum_{i,j:j>i} w_{ij}$ is the total weight of all edges.
- Sanity check of the stationary distribution:
$$\sum_i \pi_i P_{ij} = \sum_i \frac{w_i}{2w} \frac{w_{ij}}{w_i} = \frac{w_j}{2w} = \pi_j, \quad \forall j \in \mathcal{N}.$$
- Locality property of this stationary distribution: it depends only on the total weight and the weight of edges connected to the node.

Entropy Rate of Random Walk over Graph (cont'd)

Example (cont.)

Hence, the entropy rate is

$$H(\mathcal{X}) = H(X_2|X_1) = - \sum_{ij} \mu_i P_{ij} \log P_{ij} \quad (2)$$

$$= - \sum_{ij} \frac{w_{ij}}{2w} \log \frac{w_{ij}}{w_i} \quad (3)$$

$$= - \sum_{ij} \frac{w_{ij}}{2w} \log \left(\frac{w_{ij}}{2w} \times \frac{2w}{w_i} \right) \quad (4)$$

$$= - \sum_{ij} \frac{w_{ij}}{2w} \log \frac{w_{ij}}{2w} + \sum_i \frac{w_i}{2w} \log \frac{w_i}{2w} \quad (5)$$

$$= H \underbrace{\left(\dots, \frac{w_{ij}}{2w}, \dots \right)}_{|\mathcal{N}|^2 \text{ terms}} - H \underbrace{\left(\dots, \frac{w_i}{2w}, \dots \right)}_{|\mathcal{N}| \text{ terms}} \quad (6)$$

If all the edges have equal weight, the stationary distribution becomes $\pi_i = \frac{D_i}{2D}$, where D_i is the degree of node i and D is the total degree of the graph. The entropy rate becomes $H(\mathcal{X}) = \log(2D) - H\left(\frac{D_1}{2D}, \dots, \frac{D_{|\mathcal{N}|}}{2D}\right)$.

Functions of Markov Chain

Theorem

Consider a stationary Markov chain $\{X_i\}$ and $Y_i = \phi(X_i)$ for all i . We have:

$$H(Y_n|Y_{n-1}, \dots, Y_1, X_1) \leq H(\mathcal{Y}) \leq H(Y_n|Y_{n-1}, \dots, Y_1)$$
$$\lim_{n \rightarrow \infty} H(Y_n|Y_{n-1}, \dots, Y_1, X_1) = H(\mathcal{Y}) = \lim_{n \rightarrow \infty} H(Y_n|Y_{n-1}, \dots, Y_1).$$

Claim: For a stationary MC $\{X_i\}$, $\{Y_i = \phi(X_i)\}$ is stationary, but not necessarily a MC (unless ϕ is injective or constant)³.

Example (Process $\{Y_i\}$ is not MC)

Consider a mapping $\phi : \mathcal{X} = \{x_1, x_2, x_3\} \mapsto \mathcal{Y} = \{y, z\}$ s.t. $\phi(x_1) = \phi(x_2) = y, \phi(x_3) = z$. Transitions on \mathcal{X} are $\Pr(x_1 \rightarrow x_3) = \Pr(x_3 \rightarrow x_1) = \Pr(x_2 \rightarrow x_2) = 1$.

Assuming X_0 is uniformly distributed on \mathcal{X} , we have

$$P(Y_2 = z | Y_1 = y) = \frac{P(Y_1 = y, Y_2 = z)}{P(Y_1 = y)} = \frac{P(X_0 = x_3)}{P(X_0 \in \{x_2, x_3\})} = 0.5$$
$$P(Y_2 = z | Y_1 = y, Y_0 = z) = \frac{P(Y_0 = z, Y_1 = y, Y_2 = z)}{P(Y_0 = z, Y_1 = y)} = \frac{P(X_0 = x_3)}{P(X_0 = x_3)} = 1.$$

³<https://math.stackexchange.com/questions/2262424/>

Injective vs Surjective Functions⁴

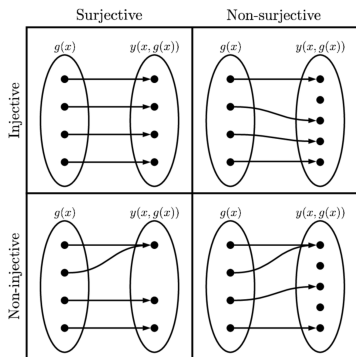


Figure: Injective/one-to-one/left

invertible: each element of the codomain is mapped to by *at most* one element of the domain ($\text{image} \subseteq \text{codomain}$).

Surjective/onto/right invertible: each element of the codomain is mapped to by *at least* one element of the domain ($\text{image} = \text{codomain}$).

Bijjective/one-to-one

correspondence/invertible: each element of the codomain is mapped to by *exactly one* element of the domain.

- Function $f : \mathcal{A} \mapsto \mathcal{B}$ is left invertible if $\exists g : \mathcal{B} \mapsto \mathcal{A}$ such that $g(f(x)) = x, \forall x \in \mathcal{A}$
- Function $f : \mathcal{A} \mapsto \mathcal{B}$ is right invertible if $\exists g : \mathcal{B} \mapsto \mathcal{A}$ such that $f(g(y)) = y, \forall y \in \mathcal{B}$

⁴https://en.wikipedia.org/wiki/Bijection,_injection_and_surjection

Functions of Markov Chain (cont'd)

Proof.

- If ϕ is a constant mapping, it is trivial to show stationarity and Markovianity of $\{Y_i\}$.
- If $\phi : \mathcal{X} \mapsto \mathcal{Y}$ is injective, then $\phi^{-1} : \mathcal{Y} \mapsto \mathcal{X}$ is well defined, where \mathcal{Y} is the image of ϕ . Thus, for Markovianity we have:

$$\Pr(Y_{n+1} = y_{n+1} | \{Y_k = y_k\}_{k \leq n}) = \Pr(X_{n+1} = \phi^{-1}(y_{n+1}) | \{X_k = \phi^{-1}(y_k)\}_{k \leq n}) \quad (7a)$$

$$= \Pr(X_{n+1} = \phi^{-1}(y_{n+1}) | X_n = \phi^{-1}(y_n)) \quad (7b)$$

$$= \Pr(Y_{n+1} = y_{n+1} | Y_n = y_n). \quad (7c)$$

Similarly, for stationarity we have

$$\Pr(Y_0 = y_0, \dots, Y_n = y_n) = \Pr(X_0 = \phi^{-1}(y_0), \dots, X_n = \phi^{-1}(y_n)) \quad (8a)$$

$$= \Pr(X_l = \phi^{-1}(y_0), \dots, X_{n+l} = \phi^{-1}(y_n)) \quad (8b)$$

$$= \Pr(Y_l = y_0, \dots, Y_{n+l} = y_n), \forall n, l. \quad (8c)$$

- If ϕ is surjective, but not injective. Then, in (8) by replacing those “=” with “ \in ”, we prove the stationarity of $\{Y_i\}$.

Function of Markov Chain (cont'd)

Proof: We have proved the upper bound. For the lower bound,

$$H(Y_n|Y_{n-1}, \dots, Y_1, X_1) = H(Y_n|Y_{n-1}, \dots, Y_1, X_1, X_0, \dots, X_{-k}) \quad (9)$$

$$= H(Y_n|Y_{n-1}, \dots, Y_1, X_1, X_0, \dots, X_{-k}, Y_0, \dots, Y_{-k}) \quad (10)$$

$$\leq H(Y_n|Y_{n-1}, \dots, Y_1, Y_0, \dots, Y_{-k}) \quad (11)$$

$$= H(Y_{n+k+1}|Y_{n+k}, \dots, Y_1), \quad (12)$$

where

- (9) follows from Markovity of X ;
- (10) is due to the fact that Y is a function of X ;
- (11) follows from conditioning reduces entropy;
- (12) follows from stationarity of Y .

The inequality is true for all k , it is also true in the limit:

$$H(Y_n|Y_{n-1}, \dots, Y_1, X_1) \leq \lim_{k \rightarrow \infty} H(Y_{n+k+1}|Y_{n+k}, \dots, Y_1) = H(\mathcal{Y}).$$

Function of Markov Chain (cont'd)

Next, we show that the upper and lower bounds converge to the same value:

$$\lim_{n \rightarrow \infty} H(Y_n | Y_{n-1}, \dots, Y_1, X_1) = \lim_{n \rightarrow \infty} H(Y_n | Y_{n-1}, \dots, Y_1),$$

which is equivalent to $\lim_{n \rightarrow \infty} I(X_1; Y_n | Y_{n-1}, \dots, Y_1) = 0$.

Note that $\underbrace{I(X_1; Y_n, \dots, Y_1)}_{\text{increases in } n} = H(X_1) - H(X_1 | Y_n, \dots, Y_1) \leq H(X_1)$. Hence, we have

$$H(X_1) \geq \lim_{n \rightarrow \infty} I(X_1; Y_n, Y_{n-1}, \dots, Y_1) \quad (13)$$

$$= \lim_{n \rightarrow \infty} \sum_{i=1}^n I(X_1; Y_i | Y_{i-1}, \dots, Y_1) \quad (14)$$

$$= \sum_{i=1}^{\infty} I(X_1; Y_i | Y_{i-1}, \dots, Y_1) \quad (15)$$

The infinite sum of nonnegative terms is finite \implies the terms must tend to 0.

Hidden Markov Model (HMM)⁵

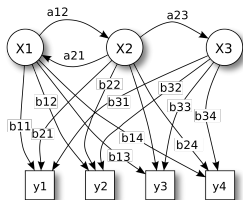


Figure: HMM diagram:

$X \rightarrow$ states;

$y \rightarrow$ possible observations;

$a \rightarrow$ state transition probabilities;

$b \rightarrow$ output (or emission) probabilities.

Given a Markov process $\{X_n\}$, each Y_i is drawn according to $p(y_i|x_i)$, conditionally independent of all the other X_j , $j \neq i$; i.e.,

$$p(x^n, y^n) = p(x^n)p(y^n|x^n) = p(x_1) \prod_{i=1}^{n-1} p(x_{i+1}|x_i) \prod_{i=1}^n p(y_i|x_i)$$

- The same argument used for functions of Markov chain carries over to HMMs; i.e., lower bounding the entropy rate by conditioning it on the underlying Markov state.

⁵Wiki: HMM is widely used in many real applications such as speech recognition, handwriting recognition, musical score following, bioinformatics, etc.

Thank You!

Email: <zhangy@ucsc.edu>

Homepage: <https://people.ucsc.edu/~yzhan419/>