

An introduction to hidden Markov models (HMMs)

25/03/2025

Hidden Markov Models (HMMs) are used to compress, infer & predict behaviours of stochastic processes.

Stochastic processes

A stochastic process is a bi-infinite sequence of random variables in time.

$$\dots X_{-3} X_{-2} X_{-1}, X_0 X_1 X_2 \dots \quad \text{where subscripts denote time.}$$

past $X_{-\infty:0}$ present future $X_{1:\infty}$

In the most general case the random variable X_t is correlated with the entire past, i.e.,

$$\Pr(X_t) = \sum_{X_{-\infty:t}} \Pr(X_t | X_{-\infty:t}) \Pr(X_{-\infty:t} = x_{-\infty:t}). \quad (\text{the sum becomes an integral in the continuous case})$$

Some *special* kinds of stochastic processes

1. Independently & identically distributed (IID) processes. //

$$\Pr(X_t | X_{-\infty:t}) = \Pr(X_t) ; \Pr(X_t, X_{t'}) = \Pr(X_t) \Pr(X_{t'}) ; \Pr(X_t) = \Pr(X_t) \forall t$$

An example of an IID process is a coin flip.

Every coinflip has the same statistics over outputs and is independent from the past coin flips.

2. Processes with finite Markov order. //

$$\Pr(X_t) = f(X_{t-k:t}) \text{ or equivalently } \Pr(X_t | X_{-\infty:t}) = \Pr(X_t | X_{t-k:t})$$

This means that we can discard the history of the past prior to k time steps prior

$$\dots X_{-k-2} X_{-k-1} \dots X_{-k} X_{-k+1} \dots X_0$$

(we can disregard this part of the process.)

3. Markov order-1 process when the process at time t depends only on one time step prior.

$$\Pr(X_t | X_{-\infty:t}) = \underbrace{\Pr(X_t | X_{t-1:t})}_{\text{We normally write } \Pr(X_{t+1}|X_t)}.$$

Markov models

If a process has Markov order 1 then the process is fully captured by the transition matrix

$M(X_{t+1}|X_t)$

e.g. if X takes values in $\mathcal{X} := \{x^{(0)}, x^{(1)}\}$, then

$$M(X_{t+1}|X_t) = \begin{pmatrix} \Pr(x_{t+1}^{(0)}|x_t^{(0)}) & \Pr(x_{t+1}^{(0)}|x_t^{(1)}) \\ \Pr(x_{t+1}^{(1)}|x_t^{(0)}) & \Pr(x_{t+1}^{(1)}|x_t^{(1)}) \end{pmatrix} \Leftrightarrow \begin{array}{ccc} x^{(0)} & & x^{(1)} \\ \xrightarrow{\Pr(x_{t+1}^{(0)}|x_t^{(0)})} & & \xleftarrow{\Pr(x_{t+1}^{(1)}|x_t^{(0)})} \\ & & \xrightarrow{\Pr(x_{t+1}^{(0)}|x_t^{(1)})} \end{array}$$

Then the process at any time-step $t=L$ is given by,

$$\Pr(X_L) = \underbrace{\Pr(X_L | X_{0:L-1})}_{\text{final distribution}} \Pr(X_{0:L-1}) = M^L \Pr(X_{0:L-1})$$

$= M^L M^L \Pr(X_{0:L-1}) = M^{2L} \Pr(X_{0:L-1})$

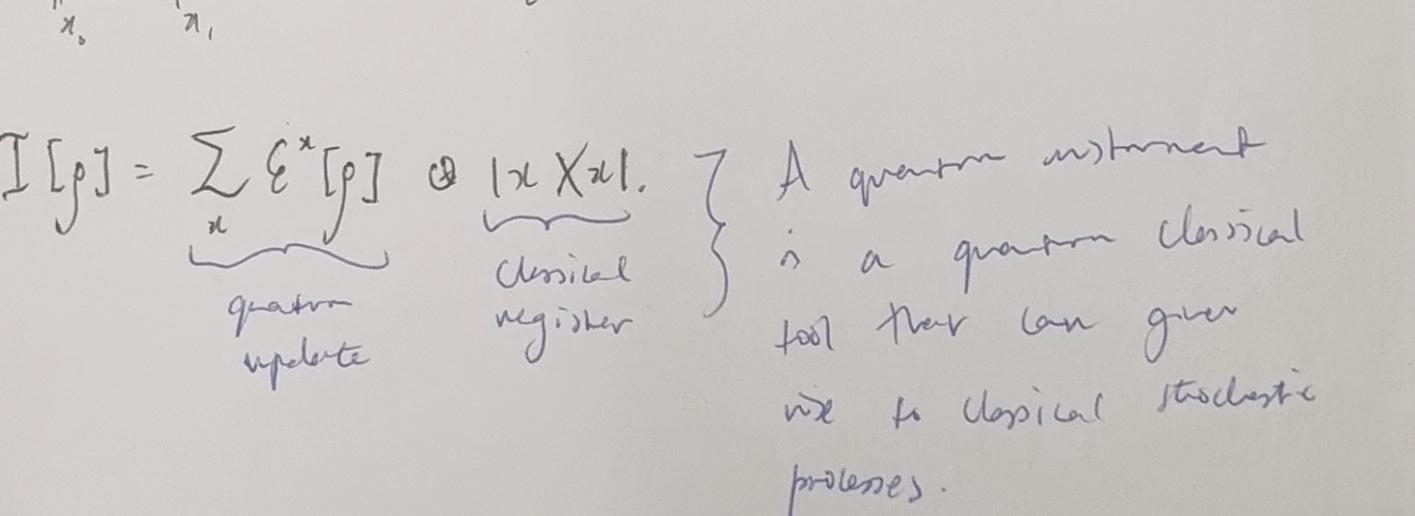
$= M^{2L} \Pr(X_0) \underbrace{\Pr(X_0)}_{\text{initial distribution}}$

If a process has Markov order k , then it can still be expressed as a Markov order 1 process by grouping k variables into a new random variable.

$$\dots X_{-3} X_{-2} X_{-1}, X_0 X_1 X_2 \dots \rightarrow \dots Y_{-3} Y_{-2} Y_{-1}, Y_0 Y_1 Y_2 \dots$$

where $Y_t = X_{t:k}$ and the new random variable takes values from an alphabet with $|X|^k$ symbols.

e.g. for $k=2$:



But say that we have a process that outputs σ symbols from an alphabet with 2 symbols, and has Markov order 10, then the Markov model will require $2^{10} = 1024$ states to express this.

We can do much better with *Hidden Markov models (HMMs)

Hidden Markov Models

Markov models fail in cases where there is infinite Markov order or the internal dynamics of the process are unknown (we instead only have a model of the Markov model).

Def: An HMM is composed of

1. $X := \{x_0, x_1, \dots, x_N\}$ a set of N observables

2. $S := \{s_0, s_1, \dots, s_M\}$ a set of M internal states

3. An $M \times M$ matrix A , $A_{lm} = \Pr(s_{t+1}^{(l)} | s_t^{(m)})$ the probability of transitioning between states of the model

4. An $N \times N$ matrix B , $B_{lm} = \Pr(x_{t+1}^{(l)} | s_t^{(m)})$ the emission probabilities for each state.

Example

$$A = \begin{pmatrix} \Pr(H|H) & \Pr(G|H) \\ \Pr(H|G) & \Pr(G|G) \end{pmatrix}$$

$$B = \begin{pmatrix} \Pr(B|H) & \Pr(R|H) & \Pr(G|R) & \Pr(B|R) \\ \Pr(B|G) & \Pr(R|G) & \Pr(G|R) & \Pr(B|R) \end{pmatrix}$$

$$= \begin{pmatrix} p(H) & p(R) & p(G) & p(B) \\ p(G) & p(R) & p(G) & p(B) \end{pmatrix}$$

$$T^x = \begin{pmatrix} \Pr(x^{(0)}|H) & \Pr(x^{(1)}|H) \\ \Pr(x^{(0)}|G) & \Pr(x^{(1)}|G) \end{pmatrix}$$

$$\sum T^x = T \quad (\text{stochastic})$$

Generally we are not guaranteed that the model representation is unique. Many parameterizations can lead to the same output process!!

We can extract probabilities of sequences by once again considering an initial state η_0 which is a distribution over hidden states.

state edge emitting: $\Pr(X_{0:L} = x_{0:L}) = \prod_{t=0}^{L-1} \Pr((AT(\dots((AT(\eta_0 \circ B_{t,1})) \circ B_{t,2}) \dots) \circ B_{t,L})$

edge emitting: $\Pr(X_{0:L} = x_{0:L}) = \prod_{t=0}^{L-1} \prod_{i=1}^N T^{x_i} \cdot \eta_0$

tensor framework: $\Pr(X_{0:L}) = \underbrace{\langle \dots \rangle}_{\text{vector}} \underbrace{T \dots T}_{\text{matrix}} \underbrace{T_{x_{0:L}}}_{\text{vector}}$

We can define π_s stationary states $\sum_x T^{x,s} = \pi_s$

\rightarrow entropy rates $\sum_s \pi_s \log(\pi_s)$ where $\pi_s = (\pi_1, \pi_2, \dots)$

\rightarrow complexity measures.

The three basic problems of HMMs = hidden tutorial

1) Given an observation sequence $x_{0:L}$ & model $\lambda = (A, B, \eta_0) = (\sum T^x, S, \eta_0)$ how do we efficiently compute $\Pr(x_{0:L} | \lambda)$, the probability of the sequence given by the model

(this is where tensor network formulations are particularly useful).

2) Given the observation sequence $x_{0:L}$ & model λ , how can we find the corresponding state sequence $s_{0:L}$ which best explains the data.

$\max_{s_{0:L}} \Pr(x_{0:L} | s_{0:L}, \lambda)$ - e.g. speech recognition, finding patterns in data etc.

3) How can we adjust model parameters to maximize $\Pr(x_{0:L} | \lambda)$ or even $\Pr(x_{0:L} | x_{0:M}, \lambda)$?

This is a model finding problem to maximize the model's ability to explain / fit the data or to predict future sequences.

(forward/backward method, tensor method, Baum Welch algorithm).

Quantum Hidden Markov Models

Motivation: Quantum models often require less memory than their classical counterparts to generate a model of a given process

arXiv: 2412.12812 for our next work which partly answers why this is

Def: A quantum hidden Markov model (qHMM) is given by a quantum instrument $\mathcal{I} = \sum \mathcal{E}^x$ with $\sum_x \mathcal{E}^x = \mathcal{E}$ CPTP.

Together with an initial state ρ_0 , the instrument gives rise to a classical stochastic process

$$\mathcal{I}[\rho] = \sum_x \mathcal{E}^x[\rho] \otimes \underbrace{|x\rangle\langle x|}_{\text{quantum update}} \quad \text{A quantum instrument is a quantum classical tool that can give rise to classical stochastic processes.}$$

$$\Pr(X_{0:L} = x_{0:L}) = \text{tr} \left(\mathcal{E}^{x_0} \mathcal{E}^{x_1} \dots \mathcal{E}^{x_L} [\rho_0] \right)$$