

Models for Structured Data

Linear Chains

- If we take a person's BP every five minutes over a 24 hour period then there is some significant dependence between successive values
- How to model this dependence?
- Occurs in protein sequences, time series (measurements ordered in time) image data (measurements defined on a spatial grid)

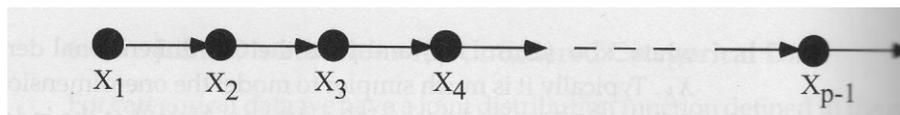
First Order Markov Model

- Structure of the data suggests a natural structuring of models we will build
- T data points observed sequentially y_1, \dots, y_T

$$p(y_1, \dots, y_T) = p_1(y_1) \prod_{t=2}^T p_t(y_t | y_{t-1})$$

Generative interpretation of Markov Model

$$p(y_1, \dots, y_T) = p_1(y_1) \prod_{t=2}^T p_t(y_t | y_{t-1})$$



y 's instead of x 's

First value chosen by drawing a y_1 value randomly according to initial distribution $p(y_1)$

Value at time $t = 2$ chosen according to the conditional density function $p(y_2/y_1)$

y_3 is generated according to $p(y_3/y_2)$

Markov model limitation

- Influence of the past is completely summarized by the value of Y at time $t-1$
- Y does not have any long-range dependencies
- This model may not be accurate in many situations
 - In modeling English text, where Y takes on values such as verb, adjective, noun, etc, deciding whether a verb is singular or plural depends on the subject of the verb which may be much further back than just one word back

Real-valued Y

- Markov model is specified as a conditional Normal distribution

Deterministic function
Linking the past y_{t-1}
to present y_t

$$p(y_t | y_{t-1}) = \frac{1}{\sqrt{2\pi\sigma}} \exp - \frac{1}{2} \left(\frac{y_t - g(y_{t-1})}{\sigma} \right)^2$$

Noise in the model

If g is chosen such that it is a linear function of y_{t-1} : $g(y_{t-1}) = \alpha_0 + \alpha_1 y_{t-1}$

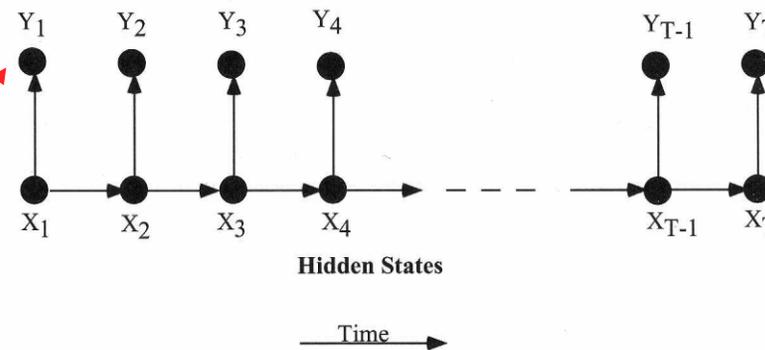
It leads to *first-order autoregressive model*

$$y_t = \alpha_0 + \alpha_1 y_{t-1} + e$$

Hidden State Variable

- Notion of hidden state for sequential and spatial models is prevalent in engineering and the sciences
- Examples include HMMs and Kalman filters

Graphical Model of HMM



Observation
Variable

Hidden
State
Variable

Generative view of HMM

- Observations are generated by moving from left to right along the chain
- Hidden state variable X is categorical (corresponding to m discrete states) and is first order Markov
- Thus x_t is generated by sampling a value from the conditional distribution $p(x_t/x_{t-1})$
- Where $p(x_t/x_{t-1})$ is an $m \times m$ matrix
- Once the state at time t is generated (with value x_t) an observation is generated with probability $p(y_t/x_t)$

View of HMM as a Mixture Model

- m different density functions for the Y variable with added Markov dependence between “adjacent” mixture components x_t and x_{t+1}
- Joint probability of an observed sequence and any particular state sequence is

$$p(y_1, \dots, y_T, x_1, \dots, x_T) = p(x_1) p(y_1 | x_1) \prod_{t=2}^T p(y_t | x_t) p(x_t | x_{t-1})$$

- To calculate $p(y_1, \dots, y_T)$, the likelihood of the observed data, one has to sum the LHS terms over the m^T possible state sequences.
 - Appears to involve a sum over an exponential number of terms
 - Viterbi algorithm performs the calculation in time proportional to $O(m^2T)$

Generalizations of HMMs

- k^{th} order Markov model
 - x_t depends on the previous k states
- Dependence of y s can be generalized
 - y_t depends on the previous k previous y s

Generalizations of HMMs

- Kalman Filters
 - Hidden states are real-valued
 - E.g., unknown velocity or momentum of a vehicle
 - Independence structure is the same as for HMM

Relationship to Finite State Machines

- First order HMM is directly equivalent to a stochastic finite state machine (FSM) with m states
 - Choice of the next state is governed by $p(x_t/x_{t+1})$
- FSMs are simple forms of regular grammars
- Next level up are context-free grammars
 - Augmenting FSM with a stack
 - To remember long-range dependencies such as closing parentheses
 - Models become more expressive but much more difficult to fit to data
- Although simple in structure HMMs have dominated due to difficulties of fitting such data

Markov Random Fields

- Instead of Y s existing in an ordered sequence more general data dependencies
- Such as data on a two-dimensional grid
- MRFs are multidimensional analogs of Markov chains (in two-dimensions a grid structure instead of chains)