

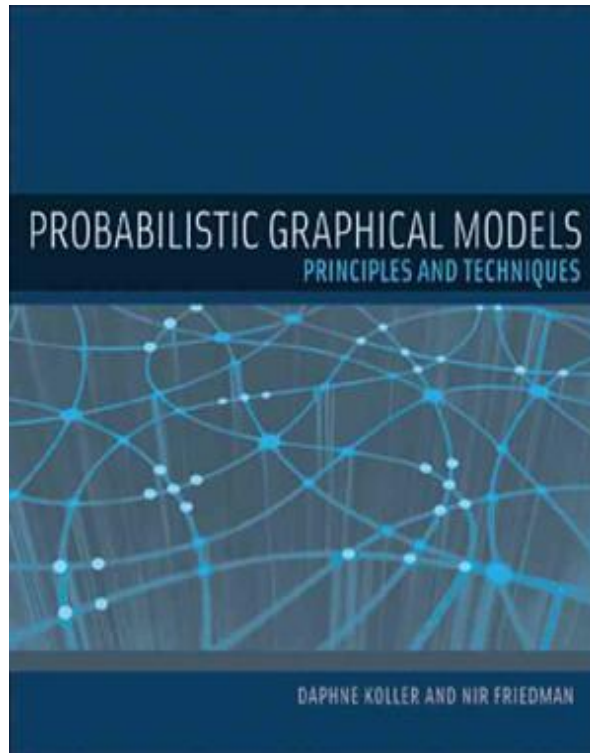
Hidden Markov Model

-- Probabilistic Graphical Model Perspective

Rui Li

Resources

- Textbook and Tutorial



A Tutorial on Hidden Markov Models and Selected Applications in Speech Recognition

LAWRENCE R. RABINER, FELLOW, IEEE

Although initially introduced and studied in the late 1960s and early 1970s, statistical methods of Markov source or hidden Markov modeling have become increasingly popular in the last several years. There are two strong reasons why this has occurred. First the models are very rich in mathematical structure and hence can form the theoretical basis for use in a wide range of applications. Second the models, when applied properly, work very well in practice for several important applications. In this paper we attempt to carefully and methodically review the theoretical aspects of this type of statistical modeling and show how they have been applied to selected problems in machine recognition of speech.

I. INTRODUCTION

Real-world processes generally produce observable outputs which can be characterized as signals. The signals can be discrete in nature (e.g., characters from a finite alphabet, quantized vectors from a codebook, etc.), or continuous in nature (e.g., speech samples, temperature measurements, music, etc.). The signal source can be stationary (i.e., its statistical properties do not vary with time), or nonstationary (i.e., the signal properties vary over time). The signals can be pure (i.e., coming strictly from a single source), or can be corrupted from other signal sources (e.g., noise) or by transmission distortions, reverberation, etc.

A problem of fundamental interest is characterizing such real-world signals in terms of signal models. There are several reasons why one is interested in applying signal models. First of all, a signal model can provide the basis for a theoretical description of a signal processing system which can be used to process the signal so as to provide a desired output. For example if we are interested in enhancing a speech signal corrupted by noise and transmission distortion, we can use the signal model to design a system which will optimally remove the noise and undo the transmission distortion. A second reason why signal models are important is that they are potentially capable of letting us learn a great deal about the signal source (i.e., the real-world process which produced the signal) without having to have the source available. This property is especially important when the cost of getting signals from the actual source is high.

In this case, with a good signal model, we can simulate the source and learn as much as possible via simulations. Finally, the most important reason why signal models are important is that they often work extremely well in practice, and enable us to realize important practical systems—e.g., prediction systems, recognition systems, identification systems, etc., in a very efficient manner.

There are several possible choices for what type of signal model is used for characterizing the properties of a given signal. Broadly one can dichotomize the types of signal models into the class of deterministic models, and the class of statistical models. Deterministic models generally exploit some known specific properties of the signal, e.g., that the signal is a sine wave, or a sum of exponentials, etc. In these cases, specification of the signal model is generally straightforward; all that is required is to determine (estimated) values of the parameters of the signal model (e.g., amplitude, frequency, phase of a sine wave, amplitudes and rates of exponentials, etc.). The second broad class of signal models is the set of statistical models in which one tries to characterize only the statistical properties of the signal. Examples of such statistical models include Gaussian processes, Poisson processes, Markov processes, and hidden Markov processes, among others. The underlying assumption of the statistical model is that the signal can be well characterized as a parametric random process, and that the parameters of the stochastic process can be determined (estimated) in a precise, well-defined manner.

For the applications of interest, namely speech processing, both deterministic and stochastic signal models have had good success. In this paper we will concern ourselves strictly with one type of stochastic signal model, namely the hidden Markov model (HMM). (These models are referred to as Markov sources or probabilistic functions of Markov chains in the communications literature.) We will first review the theory of Markov chains and then extend the ideas to the class of hidden Markov models using several simple examples. We will then focus our attention on the three fundamental problems¹ for HMM design, namely the

Manuscript received January 15, 1988; revised October 4, 1988. The author is with AT&T Bell Laboratories, Murray Hill, NJ 07974-2070, USA.
IEEE Log Number 8825949.

¹The idea of characterizing the theoretical aspects of hidden Markov modeling in terms of solving three fundamental problems is due to Jack Ferguson of IDA (Institute for Defense Analysis) who introduced it in lectures and writing.

0018-9219/89/030200-02\$7.50/0 © 1989 IEEE

Resources

- **Software**

- Hidden Markov Model (HMM) Matlab Toolbox

- By Kevin Murphy

- GraphLab

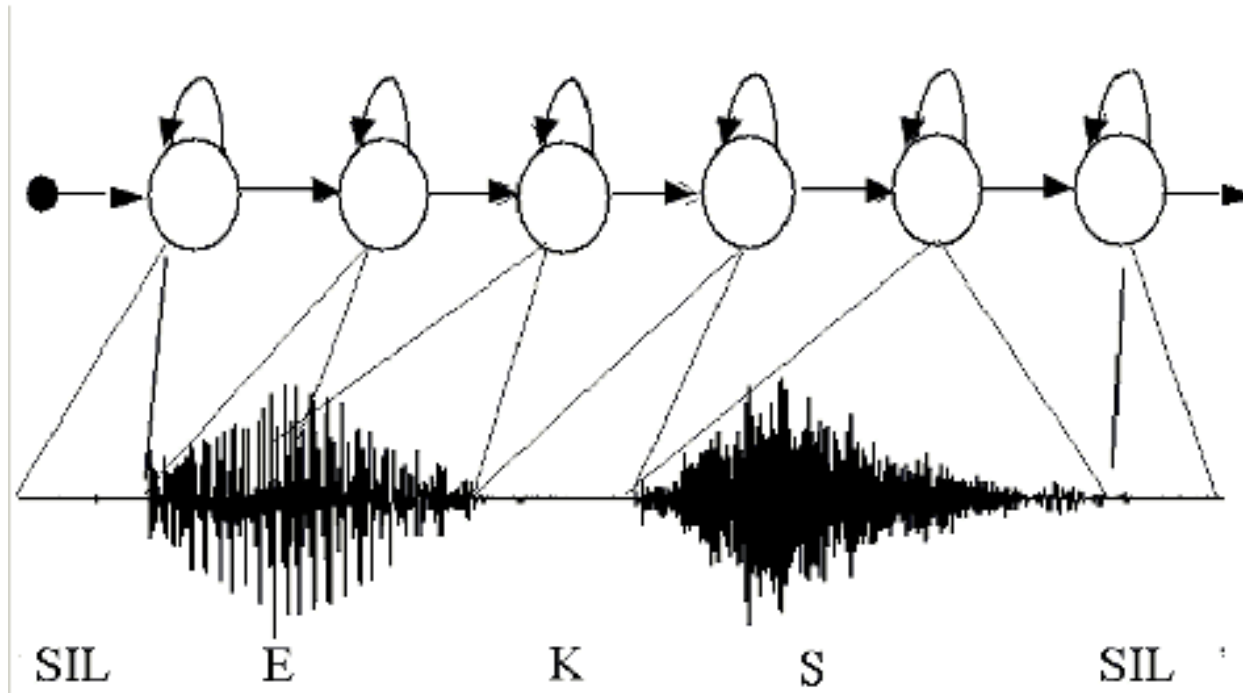
- By CMU

- Hidden Markov Model Toolkit (HTK)

- C Libraries

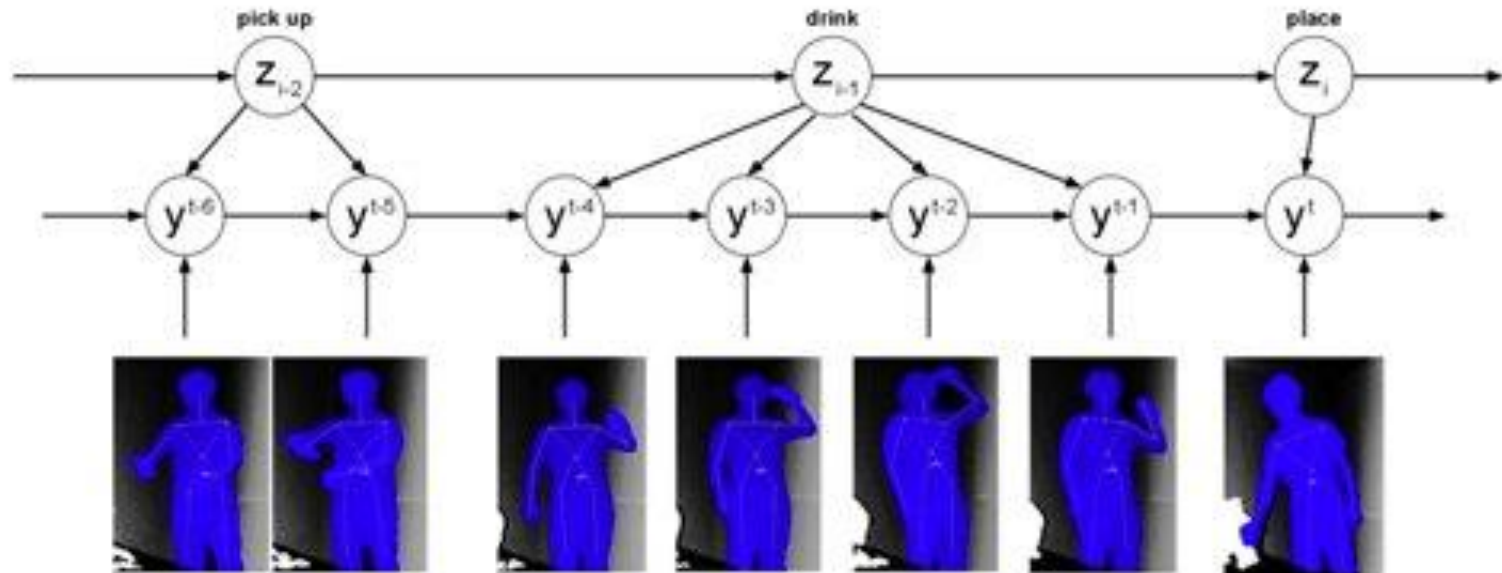
Dynamic Phenomena

- Speech Recognition



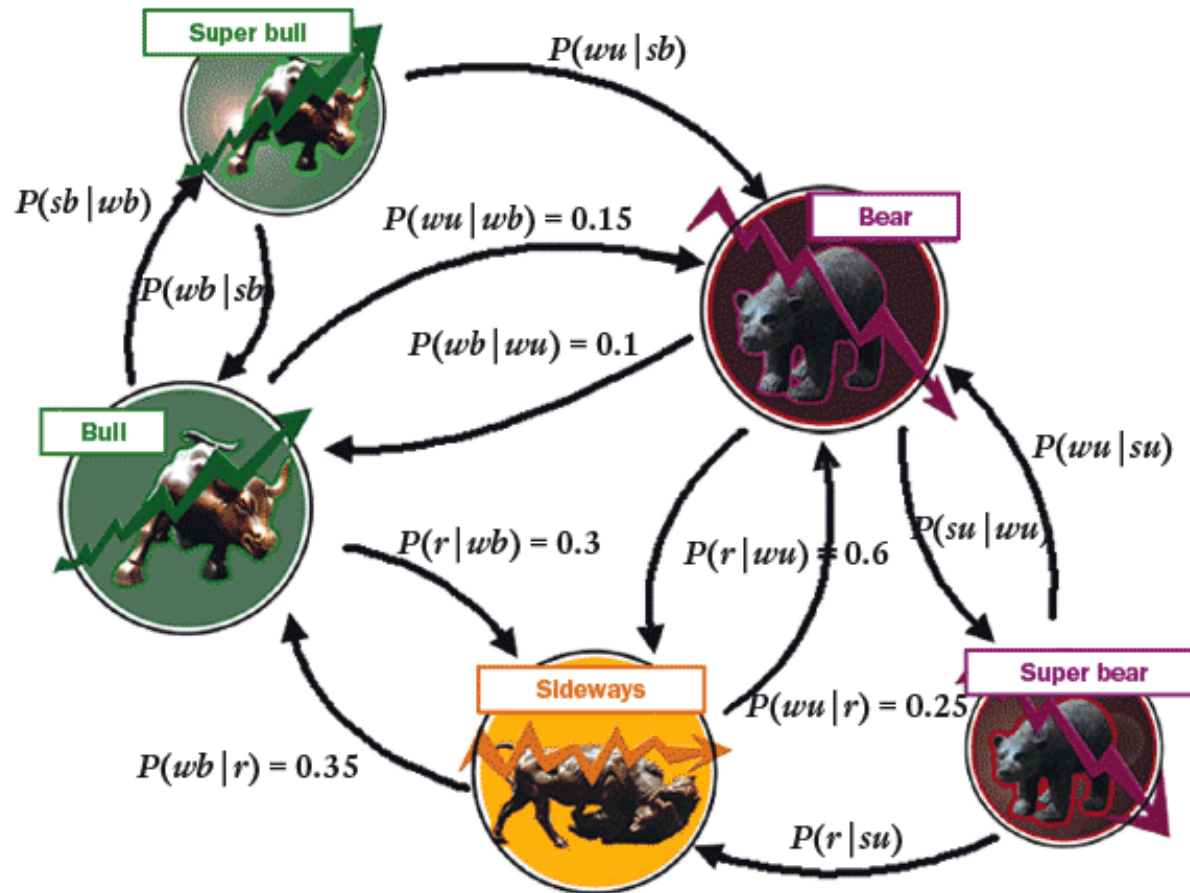
Dynamic Phenomena

- Body Motion Tracking



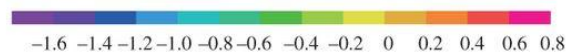
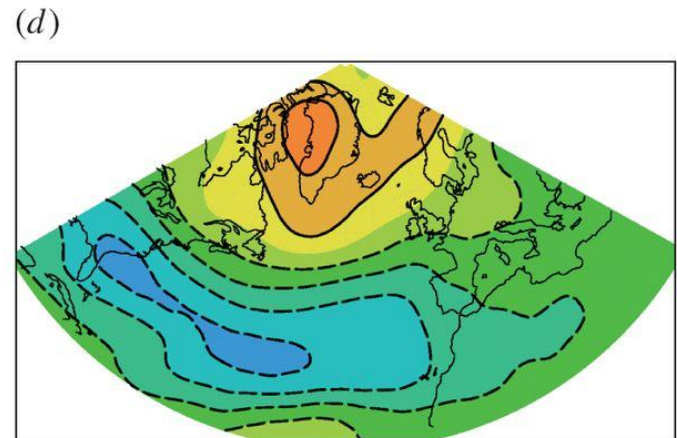
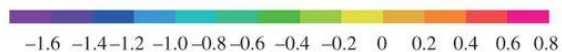
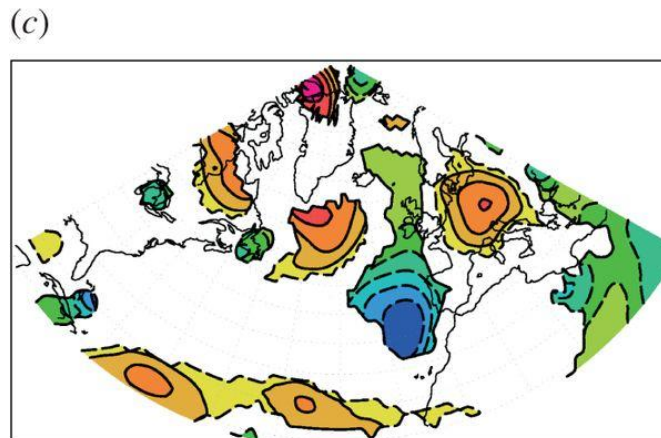
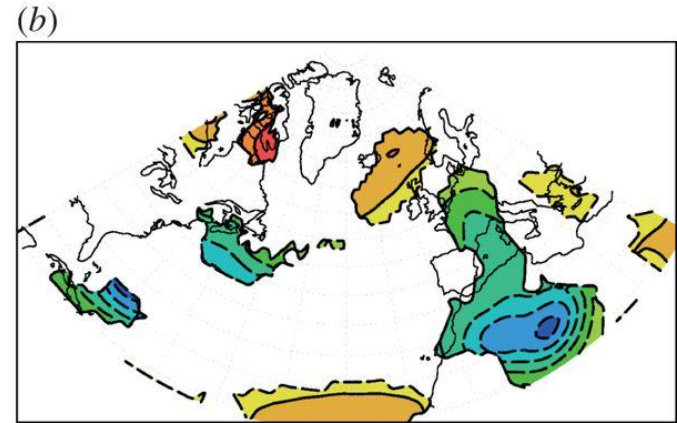
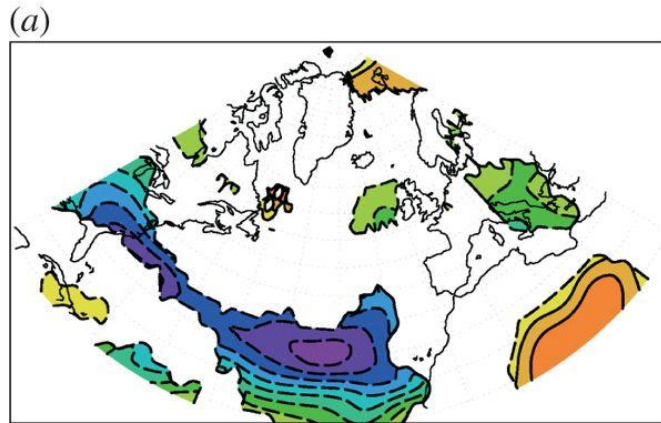
Dynamic Phenomena

- Stock Prediction



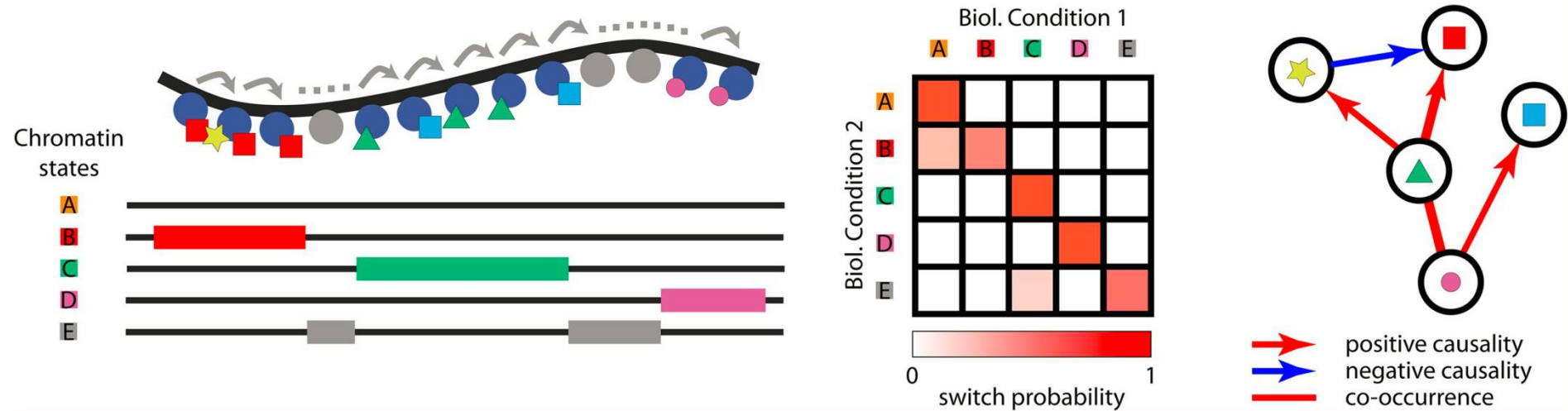
Dynamic Phenomena

- Climate Change



Bioinformatics

- DNA Sequences



Outline

- **Lecture One**

- HMMs as Probabilistic Graphical Models

- Motivation
 - Algebraic representation
 - Graphical representation

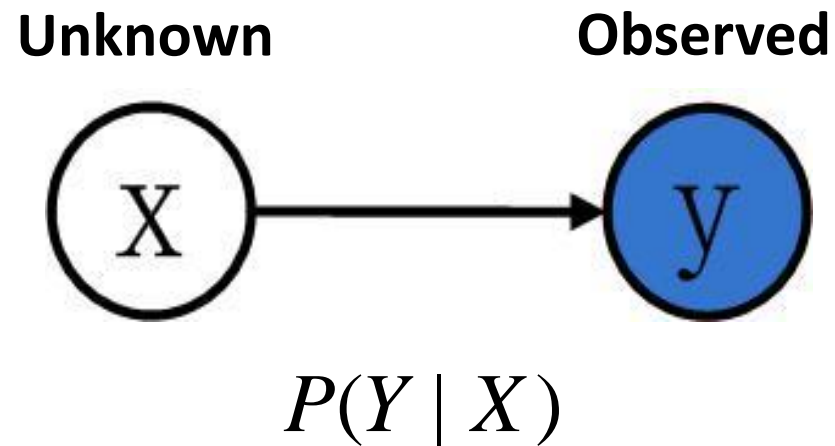
- **Lecture Two**

- HMMs with Inference and Learning

- Message Passing (Forward-Backward)
 - Expectation-Maximization (Baum-Welch)
 - Application Demos

Motivation

- A simple graphical model



Motivation

- An Example

x



y



Motivation

- Inference



Posterior
probability

Prior probability

Noise model

$$P(X | Y) = \frac{P(X, Y)}{P(Y)} = \frac{P(X)P(Y | X)}{P(Y)}$$

Constant

Curse of Dimensionality

X



Size of the lookup table of $P(X)$

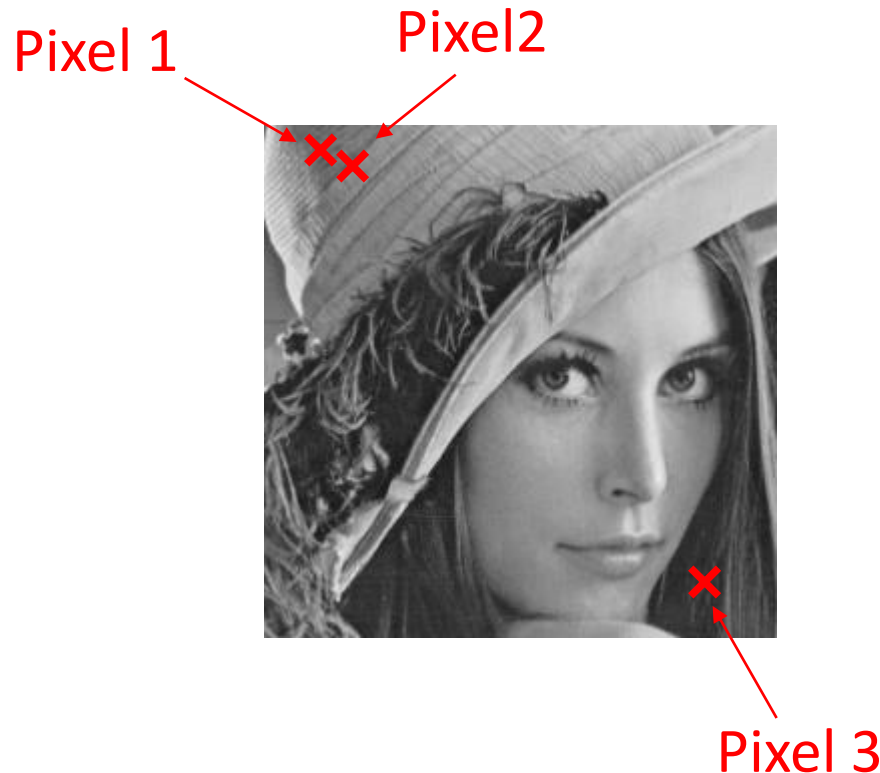
$$|P(X)| = 256^{100 \times 100} = 2^{80000}$$

171 169 167 167 166 165 166 164 167 171 171 174 174 175 173 171
168 168 168 167 166 167 167 165 169 168 174 176 175 175 172
168 167 167 165 166 166 167 167 168 170 178 177 176 174 174 173
168 168 165 169 167 168 167 165 168 175 177 177 175 175 172 171
169 170 167 169 169 168 163 166 172 169 174 173 175 178 173 173
171 169 170 168 169 168 169 168 168 170 175 175 177 178 176
172 171 170 168 169 169 167 168 173 172 173 177 174 175 178 176
172 174 171 170 166 168 167 168 172 172 177 179 172 175 175
171 167 176 169 170 169 168 169 171 172 174 174 175 173 174 178
174 172 173 173 173 174 171 171 172 174 172 172 169 173 173
173 173 173 176 178 172 171 174 174 173 175 175 173 173 171
173 175 178 173 173 171 171 175 175 177 178 175 174 173 175 178
178 175 174 169 173 175 177 175 177 177 174 175 176 177 177 174
173 175 173 174 172 173 174 175 174 171 173 174 175 174 172 171
177 174 175 175 172 171 172 176 172 173 172 172 173 170 170 175
173 171 174 168 176 172 173 173 175 174 171 174 175 173 174 174
175 173 171 172 170 171 176 175 178 172 174 175 175 175 172
181 179 177 172 170 170 169 179 175 174 175 174 172 175 174 175
188 184 179 178 176 176 176 174 172 178 172 174 173 172 174 173
195 191 188 186 185 183 180 177 178 175 174 176 175 174 176 176
200 199 197 193 190 187 185 180 176 175 180 177 175 175 176 177
202 202 199 202 199 194 187 180 175 179 177 176 174 175 176 173

Probabilistic Graphical Model

- **The basic idea**

- $P(X)$ has some locality properties encoded by graphs

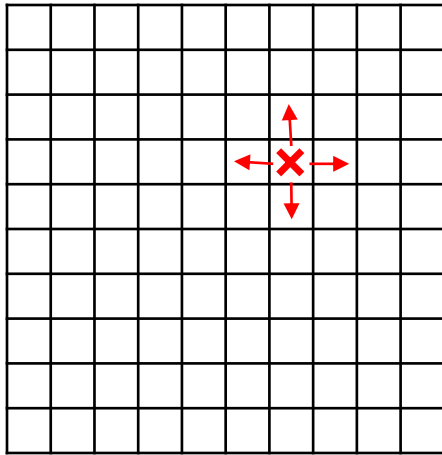


Object Tracking

$X = \{x_1, x_2, \dots, x_{T-1}, x_T\}$ x_t location at time t

$Y = \{y_1, y_2, \dots, y_{T-1}, y_T\}$ y_t sensor measurement at time t

$T = 1000$



Inference: $P(X | Y)$

Computation Complexity:

$$|P(X | Y)| = (10 \times 10)^{1000}$$

Probabilistic Graphical Model

- **PDF Representation**

- $P(X, Y)$

- **Inference**

- Given Y'

- Use $P(X | Y')$ to solve problems

- **Learning**

- Given $\{(x_t, y_t)\}$

- Fit $P(X | Y)$

Hidden Markov Models

- **Definition**

- $\{x_t, y_t\}_{1..T}$ are a HMM, if

- X is a Markov process

$$P(x_t \mid x_1, x_2, \dots, x_{t-1}, x_{t+1}, \dots, x_{T-1}, x_T) = P(x_t \mid x_{t-1})$$

- y_t only depends on x_t

$$P(y_t \mid x_1, x_2, \dots, x_{t-1}, x_t, x_{t+1}, \dots, x_{T-1}, x_T) = P(y_t \mid x_t)$$

Hidden Markov Models

- **Representation**

- Claim: for $\{x_t, y_t\}_{1..T}$ as a HMM

$$P(X, Y) = \prod_{t=1}^T P(x_t | x_{t-1}) P(y_t | x_t)$$

Proof: $P(X, Y) = P(X)P(Y | X)$

$$\begin{aligned} P(X) &= P(x_1, x_2, \dots, x_T) \\ &= P(x_T | x_{T-1}, x_{T-2}, \dots, x_1) P(x_{T-1} | x_{T-2}, x_{T-3}, \dots, x_1) \dots P(x_2 | x_1) P(x_1) \\ &= P(x_T | x_{T-1}) P(x_{T-1} | x_{T-2}) \dots P(x_2 | x_1) P(x_1) \\ &= \prod_{t=1}^T P(x_t | x_{t-1}) \end{aligned}$$

$$\begin{aligned} P(Y | X) &= P(y_1, y_2, \dots, y_T | X) \\ &= P(y_T | y_{T-1}, y_{T-2}, \dots, y_1, X) P(y_{T-1} | y_{T-2}, y_{T-3}, \dots, y_1, X) \dots P(y_2 | y_1, X) P(y_1 | X) \\ &= P(y_T | x_T) P(y_{T-1} | x_{T-1}) \dots P(y_2 | x_2) P(y_1 | x_1) \\ &= \prod_{t=1}^T P(y_t | x_t) \end{aligned}$$

Hidden Markov Models

- **Representation**

- Claim: for $\{x_t, y_t\}_{1..T}$ as a HMM

$$P(X, Y) = \prod_{t=1}^T P(x_t \mid x_{t-1}) P(y_t \mid x_t)$$

Computational Complexity:

before claim: $|P(X, Y)| = 100^{1000} \times 100^{1000}$

after claim: $|P(X, Y)| = 2000 \times 100^2$

Hidden Markov Models

- **Representation**

- Claim: for $\{x_t, y_t\}_{1..T}$ as a HMM

$$P(X, Y) = \prod_{t=1}^T P(x_t \mid x_{t-1}) P(y_t \mid x_t)$$

Statistical queries:

$$P(x_t \mid x_{t-1}, x_{t-2}) = P(x_t \mid x_{t-1})$$

$$P(x_t \mid x_{t-2}, x_{t-3})$$

$$= \sum_{x_{t-1}} P(x_t, x_{t-1} \mid x_{t-2}, x_{t-3})$$

$$= \sum_{x_{t-1}} P(x_t \mid x_{t-1}, x_{t-2}, x_{t-3}) P(x_{t-1} \mid x_{t-2}, x_{t-3})$$

$$= \sum_{x_{t-1}} P(x_t \mid x_{t-1}) P(x_{t-1} \mid x_{t-2})$$

$$= \sum_{x_{t-1}} P(x_t, x_{t-1} \mid x_{t-2})$$

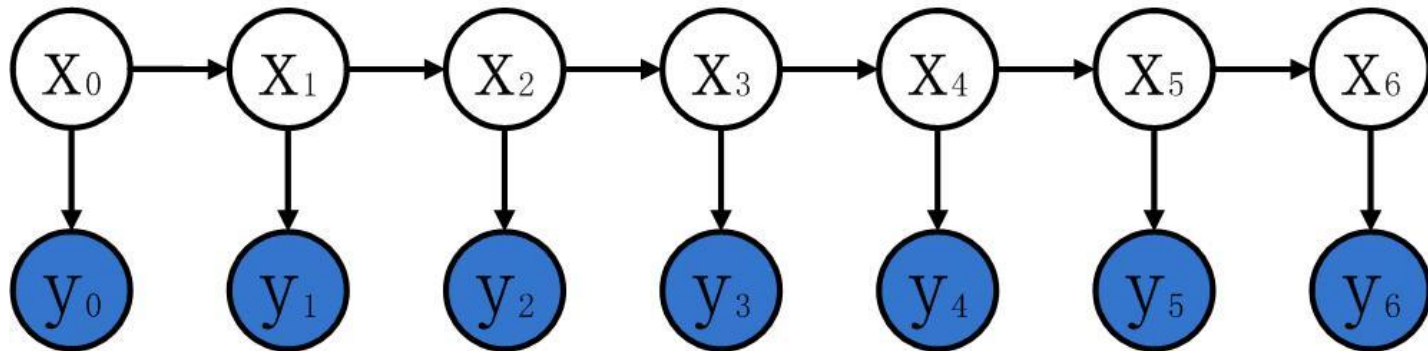
$$= P(x_t \mid x_{t-2})$$

HMMs and Graphical Models

- **Definition**

- The graph represents a HMM is

- a chain of x_t
 - y_t connects to x_t

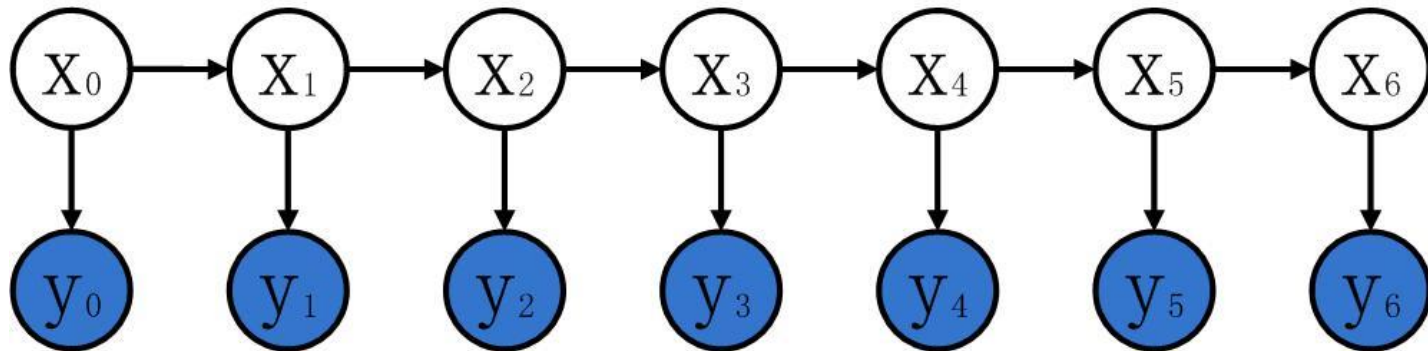


HMMs and Graphical Models

- **Theorem (Hammersley-Clifford)**
 - Given any random variables A, B and C

$$P(C \mid A, B) = P(C \mid B)$$

iff B separates A and C in the graph



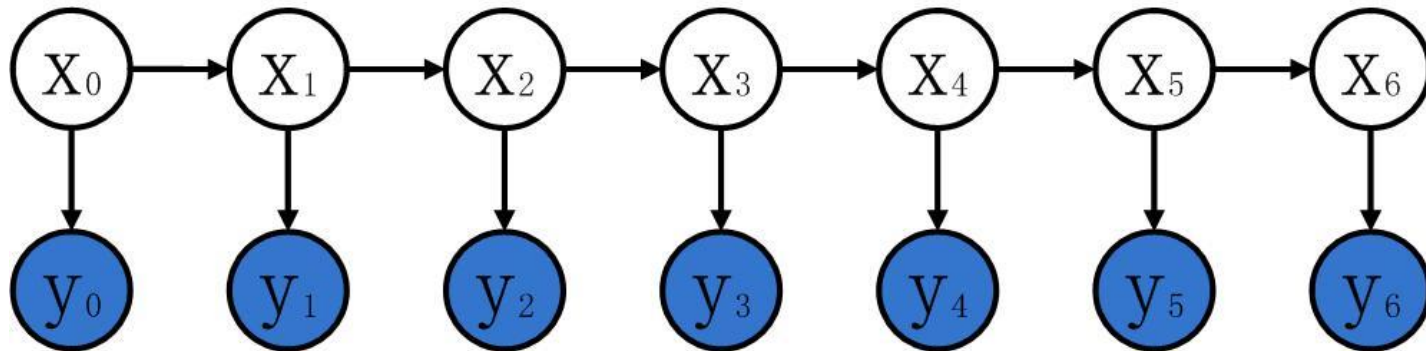
$$P(x_6 \mid x_1, x_2) = P(x_6 \mid x_2)$$

$$P(y_6 \mid y_4, y_5)$$

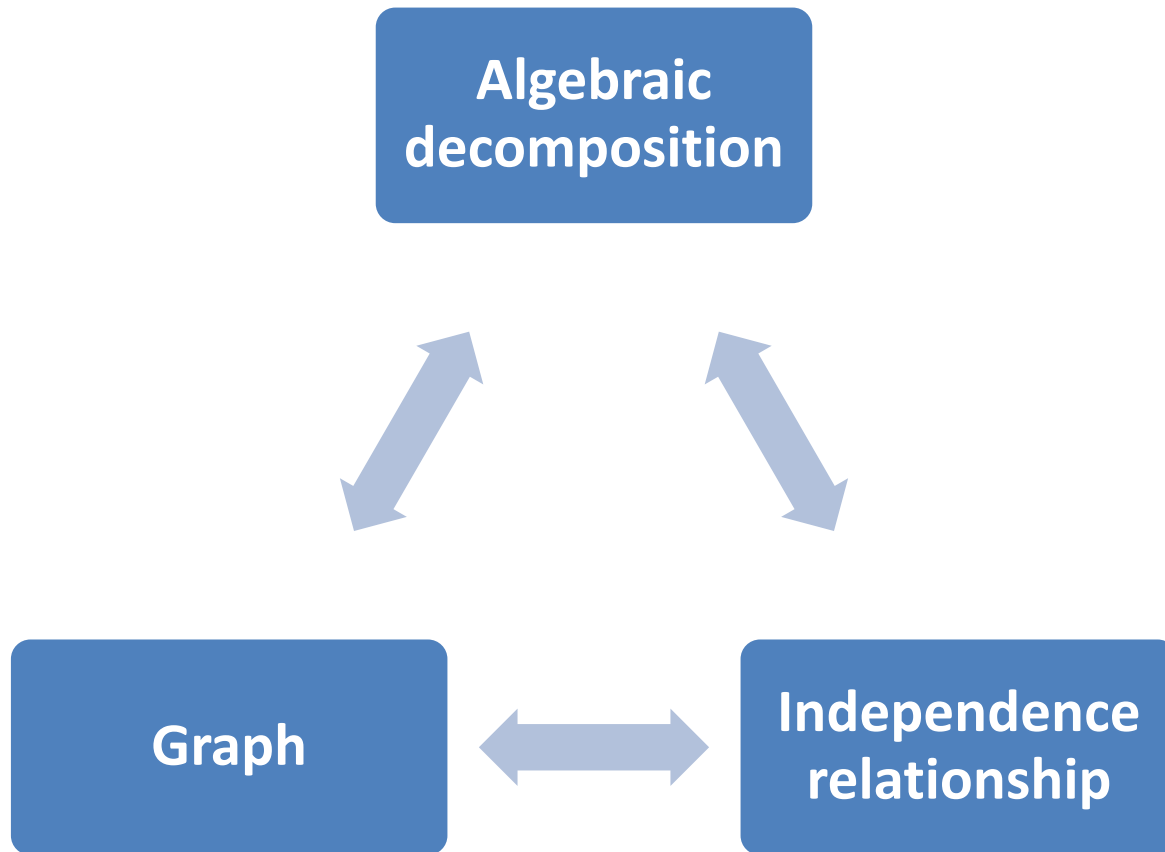
HMMs and Graphical Models

- Graph & Factorization

$$P(X, Y) = \prod_{t=1}^T P(x_t \mid x_{t-1}) P(y_t \mid x_t)$$



HMMs and Graphical Models



Probabilistic Graphical Model

- **PDF Representation**

- $P(X, Y)$

- **Inference**

- Given Y'

- Use $P(X | Y')$ to solve problems

- **Learning**

- Given $\{(x_t, y_t)\}$

- Fit $P(X | Y)$

Inference

- **MAP (Maximum A Posteriori)**

$$X^* = \arg \max_X P(X | Y')$$

- **Marginalization**

$$P(x_t | Y') \quad \forall t$$

Henceforth, “inference”== marginalization