

# Hidden Markov Models

Theoretical Aspects. Octave Implementation. Applications

Alexandru Sorici, Tudor Berariu

Romanian Asociation for Artificial Intelligence

October, 27<sup>th</sup>, 2012

PART 1.

## **Intro**

PART 2.

## **Theory of HMMs**

PART 3.

## **Demo & Discussions**

# Outline

# Intro

## 1 ARIA Education Workshops

- ARIA's Mission
- ARIA Education
- Workshop Program

## 1 ARIA Education Workshops

- ARIA's Mission
- ARIA Education
- Workshop Program

## 1 ARIA Education Workshops

- ARIA's Mission
- ARIA Education
- Workshop Program

# ARIA EDU

:)



## 1 ARIA Education Workshops

- ARIA's Mission
- ARIA Education
- Workshop Program

# Today's Program

9:00	Registration
10:00	Rahaturi despre ARIA
11:00	HMM Theory

# Theory of HMMs

## 2 Machine Learning Applications for HMM

- Machine Learning
- Where do HMMs fit into Machine Learning?

## 3 Theory of HMMs

- The 3 things you want from an HMM
- Mathematical Foundations for HMMs
- Notation Conventions & Framework Description

## 4 Implementing HMMs

- Using the Model for Estimations: the Forward-Backward algorithm
- Learning from Observations: Baum-Welch algorithm
- Uncovering Hidden states: Viterbi algorithm

## 2 Machine Learning Applications for HMM

- Machine Learning
- Where do HMMs fit into Machine Learning?

## 3 Theory of HMMs

- The 3 things you want from an HMM
- Mathematical Foundations for HMMs
- Notation Conventions & Framework Description

## 4 Implementing HMMs

- Using the Model for Estimations: the Forward-Backward algorithm
- Learning from Observations: Baum-Welch algorithm
- Uncovering Hidden states: Viterbi algorithm

# What is Machine Learning?

Machine Learning

..**Trascau** is a beautiful horse.

# Machine Learning Applications

- Computer Vision: Google Car
- Machine Translation
- Speech Recognition
- Recommender Systems
- Intelligent Advertising

## 2 Machine Learning Applications for HMM

- Machine Learning
- Where do HMMs fit into Machine Learning?

## 3 Theory of HMMs

- The 3 things you want from an HMM
- Mathematical Foundations for HMMs
- Notation Conventions & Framework Description

## 4 Implementing HMMs

- Using the Model for Estimations: the Forward-Backward algorithm
- Learning from Observations: Baum-Welch algorithm
- Uncovering Hidden states: Viterbi algorithm



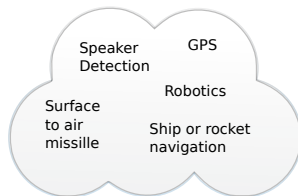
# Machine Learning Classification

## Types of Machine Learning Problems

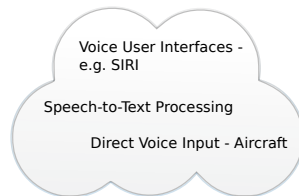
- Regression
- Classification
- Reinforcement Learning
- supervised learning (eg. ..)
- unsupervised

# Sequence / Temporal problems (I)

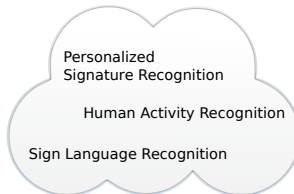
## OBJECT TRACKING



## SPEECH RECOGNITION



## GESTURE RECOGNITION



# Sequence / Temporal problems (II)

## BIOINFORMATICS

Protein Sequencing

Modeling of a Gene  
Regulatory Network

## ECONOMICS

Stock Price Prediction

Econometrics

- estimate a country's economic indicators across time -

# Probabilistic Reasoning over Time - Models

Consider some of the previously presented problems ...

# Probabilistic Reasoning over Time - Models

Consider some of the previously presented problems ...

How do we model such dynamic situations?

# Probabilistic Reasoning over Time - Models

Consider some of the previously presented problems ...

How do we model such dynamic situations?

## States and Observations

- The process of change is viewed as a series of **time slices (snapshots)**
- Each time slice contains a set of random variables
  - $\mathbf{O}_t$  - set of all **observable** evidence variables at time  $t$
  - $\mathbf{Q}_t$  - set of all **unobservable / hidden** state variables at time  $t$

# Probabilistic Reasoning over Time - Assumptions

Consider some of the previously presented problems ...

# Probabilistic Reasoning over Time - Assumptions

Consider some of the previously presented problems ...

What **assumptions** (if any) do we make?



# Probabilistic Reasoning over Time - Assumptions

Consider some of the previously presented problems ...

What **assumptions** (if any) do we make?

## Stationary Process

The process of change is governed by laws **that do not themselves change over time**.

**Implication:** we need to specify conditional distributions only for the variables within a *representative* timeslice.

# Probabilistic Reasoning over Time - Assumptions

Consider some of the previously presented problems ...

What **assumptions** (if any) do we make?

## Stationary Process

The process of change is governed by laws **that do not themselves change over time**.

**Implication:** we need to specify conditional distributions only for the variables within a *representative* timeslice.

## Markov Assumption

The current state in a process of change depends only on a **finite history** of previous states.

**Implication:** there is a **bounded** number of “parents” for the variables in each time slice.

$$P(Q_t | Q_{1:t-1}) = P(Q_t | Q_{t-1}) \quad P(O_t | Q_{1:t}, Q_{1:t-1}) = P(O_t | Q_t)$$

# Probabilistic Reasoning over Time - Inference

What are the basic inference tasks that must be solved?

# Probabilistic Reasoning over Time - Inference

What are the basic inference tasks that must be solved?

## Filtering (monitoring)

The task of computing the **belief state** - the posterior distribution over the **current state**, given all evidence to date.

$$P(\mathbf{Q}_t | \mathbf{o}_{1:t})$$

# Probabilistic Reasoning over Time - Inference

What are the basic inference tasks that must be solved?

## Filtering (monitoring)

The task of computing the **belief state** - the posterior distribution over the **current state**, given all evidence to date.

$$P(\mathbf{Q}_t | \mathbf{o}_{1:t})$$

## Evaluation (likelihood)

The task of computing the **likelihood** of the evidence up to present.

$$P(\mathbf{o}_{1:t})$$

# Probabilistic Reasoning over Time - Inference

## Prediction

The task of computing the posterior distribution over the **future state**, given all evidence to date.

$P(\mathbf{Q}_{t+k} | \mathbf{o}_{1:t})$ , for some  $k > 0$

# Probabilistic Reasoning over Time - Inference

## Prediction

The task of computing the posterior distribution over the **future state**, given all evidence to date.

$$P(\mathbf{Q}_{t+k} | \mathbf{o}_{1:t}), \text{ for some } k > 0$$

## Smoothing (hindsight)

The task of computing the posterior distribution over a **past state**, given all evidence to the present.

$$P(\mathbf{Q}_k | \mathbf{o}_{1:t}), \text{ for some } 1 \leq k < t$$

Provides a better estimate of the state than was available at the time.

# Probabilistic Reasoning over Time - Inference

## Most likely explanation

Given a *sequence of observations*, find the **sequence of states** that is **most likely** to have generated those observations.  $\operatorname{argmax}_{q_{1:t}} \mathbf{P}(\mathbf{q}_{t+k} | \mathbf{o}_{1:t})$ , for some  $k > 0$



# Probabilistic Reasoning over Time - Inference

## Most likely explanation

Given a *sequence of observations*, find the **sequence of states** that is **most likely** to have generated those observations.  $\operatorname{argmax}_{q_{1:t}} \mathbf{P}(\mathbf{q}_{t+k} | \mathbf{o}_{1:t})$ , for some  $k > 0$

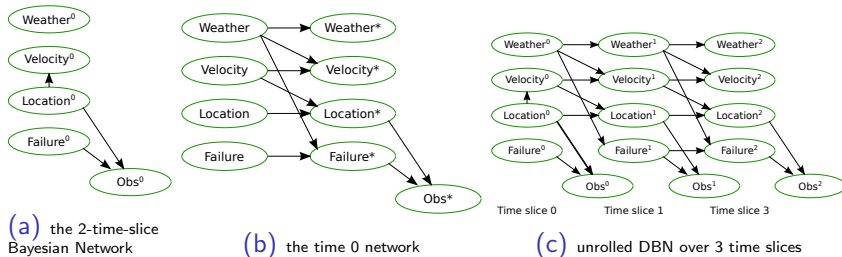
## Learning

Given a set of *observation sequences*, find a method to learn the **transition** (e.g.  $\mathbf{P}(\mathbf{q}_{t+1} = s_j | \mathbf{q}_t = s_i)$ ,  $1 \leq i, j < N$ ) and **sensor** ( $\mathbf{P}(\mathbf{o}_t | \mathbf{q}_t)$ ) **models** from the observations.

# Probabilistic Reasoning over Time - Known Methods

## Dynamic Bayesian Networks (DBN)

A DBN is Bayesian network that represents a temporal probability model.



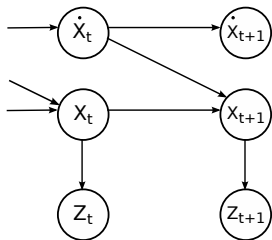
**Figure:** A highly simplified DBN for monitoring a vehicle (Koller and Friedman 2009)

Applied in problems like: object tracking, human activity recognition, protein sequencing etc.

# Probabilistic Reasoning over Time - Known Methods

## Kalman Filters (Linear Dynamical Systems)

A temporal model of one or more real-valued variables that **evolve linearly** over time, with some **Gaussian noise**.



- can be viewed as DBNs where all variables are continuous and all dependencies are linear gaussian
- wide application in **object tracking**

**Figure:** BN structure for a linear dynamical system with position  $X_t$ , velocity  $\dot{X}_t$ , and position measurement  $Z_t$

# Probabilistic Reasoning over Time - Known Methods

## Hidden Markov Models (HMM)

An HMM is a temporal probabilistic model in which the state of the process is described by a **single discrete** random variable. The possible values of the variable are the possible states of the world.

Used successfully in applications like:

- Handwriting Recognition
- Gesture Recognition
- Speech Recognition
- Part-of-Speech Tagging
- DNA Sequencing

## 2 Machine Learning Applications for HMM

- Machine Learning
- Where do HMMs fit into Machine Learning?

## 3 Theory of HMMs

- The 3 things you want from an HMM
- Mathematical Foundations for HMMs
- Notation Conventions & Framework Description

## 4 Implementing HMMs

- Using the Model for Estimations: the Forward-Backward algorithm
- Learning from Observations: Baum-Welch algorithm
- Uncovering Hidden states: Viterbi algorithm

## The 3 fundamental problems (Rabiner 1989)

- Particularization of temporal inference problems to the HMM case
- The restricted structure of the HMM allows for elegant implementations of all the basic algorithms

## The 3 fundamental problems (Rabiner 1989)

- Particularization of temporal inference problems to the HMM case
- The restricted structure of the HMM allows for elegant implementations of all the basic algorithms

### Evaluation Problem

Given a model and a sequence of observations, how do we compute the probability that the **observed sequence** was produced by the model?

## The 3 fundamental problems (Rabiner 1989)

- Particularization of temporal inference problems to the HMM case
- The restricted structure of the HMM allows for elegant implementations of all the basic algorithms

### Evaluation Problem

Given a model and a sequence of observations, how do we compute the probability that the **observed sequence** was produced by the model?

### Best Explanation of Observations Problem

Given a model and a sequence of observations how do we choose a corresponding sequence of **states** which *gives meaning* to the observations?  
How do we *uncover* the hidden part of the model?



## The 3 fundamental problems (Rabiner 1989)

- Particularization of temporal inference problems to the HMM case
- The restricted structure of the HMM allows for elegant implementations of all the basic algorithms

### Evaluation Problem

Given a model and a sequence of observations, how do we compute the probability that the **observed sequence** was produced by the model?

### Best Explanation of Observations Problem

Given a model and a sequence of observations how do we choose a corresponding sequence of **states** which *gives meaning* to the observations?  
How do we *uncover* the hidden part of the model?

### Model Estimation (Training) Problem

Given some observed sequences, how do we adjust the **parameters** of an HMM model that best tries to explain the observations?

# Outline

## 2 Machine Learning Applications for HMM

- Machine Learning
- Where do HMMs fit into Machine Learning?

## 3 Theory of HMMs

- The 3 things you want from an HMM
- **Mathematical Foundations for HMMs**
- Notation Conventions & Framework Description

## 4 Implementing HMMs

- Using the Model for Estimations: the Forward-Backward algorithm
- Learning from Observations: Baum-Welch algorithm
- Uncovering Hidden states: Viterbi algorithm

# Elements of an HMM

- $N$  Hidden States :  $S_1, S_2, \dots S_N$
- $M$  Observable Variables :  $O_1, O_2, \dots O_M$

Parameters:

- Transition Function / Matrix between states
- Emission probabilities
- Initial state probabilities

# Formalisation of the estimation problem



# Formalisation of problem # 2

- $P(Q_1) = \sum_{x=\{1,2\}}^N P(Q_2)\theta\Pi$
- $P(Q_i|q_i = s_x) = i \times x \cdot i \dots$

# Formalisation of parameters estimation problem



## 2 Machine Learning Applications for HMM

- Machine Learning
- Where do HMMs fit into Machine Learning?

## 3 Theory of HMMs

- The 3 things you want from an HMM
- Mathematical Foundations for HMMs
- **Notation Conventions & Framework Description**

## 4 Implementing HMMs

- Using the Model for Estimations: the Forward-Backward algorithm
- Learning from Observations: Baum-Welch algorithm
- Uncovering Hidden states: Viterbi algorithm

# Notation Conventions



# Variables in Octave

# Outline

## 2 Machine Learning Applications for HMM

- Machine Learning
- Where do HMMs fit into Machine Learning?

## 3 Theory of HMMs

- The 3 things you want from an HMM
- Mathematical Foundations for HMMs
- Notation Conventions & Framework Description

## 4 Implementing HMMs

- Using the Model for Estimations: the Forward-Backward algorithm
- Learning from Observations: Baum-Welch algorithm
- Uncovering Hidden states: Viterbi algorithm

# Outline

## 2 Machine Learning Applications for HMM

- Machine Learning
- Where do HMMs fit into Machine Learning?

## 3 Theory of HMMs

- The 3 things you want from an HMM
- Mathematical Foundations for HMMs
- Notation Conventions & Framework Description

## 4 Implementing HMMs

- Using the Model for Estimations: the Forward-Backward algorithm
- Learning from Observations: Baum-Welch algorithm
- Uncovering Hidden states: Viterbi algorithm

# Learning from observations - Reminder

## Model Estimation (Training) Problem

Given some observed sequences, how do we adjust the **parameters** of an HMM model that best tries to explain the observations?

# Learning from observations - Reminder

## Model Estimation (Training) Problem

Given some observed sequences, how do we adjust the **parameters** of an HMM model that best tries to explain the observations?

Adjust the model parameters  $\lambda = (A, B, \Pi)$  to obtain  $\max_{\lambda} P(O|\lambda)$

# Learning from observations - Reminder

## Model Estimation (Training) Problem

Given some observed sequences, how do we adjust the **parameters** of an HMM model that best tries to explain the observations?

Adjust the model parameters  $\lambda = (A, B, \Pi)$  to obtain  $\max_{\lambda} P(O|\lambda)$

The observation sequence used to adjust the model parameters is called a **training** sequence.

Training problem is crucial - allows to create best models for real phenomena.

# Learning from observations - Aspects of the approach

# Learning from observations - Aspects of the approach

## Problem

There is no known way to analytically solve for the model which maximizes the probability of the observation sequence.



# Learning from observations - Aspects of the approach

## Problem

There is no known way to analytically solve for the model which maximizes the probability of the observation sequence.

## Solution

We can choose  $\lambda = (A, B, \Pi)$  such that  $\max_{\lambda} P(O|\lambda)$  is **locally maximized** using an **iterative procedure** such as *Baum-Welch*.

The method is an instance of the *EM algorithm* (Dempster, Laird, and Rubin 1977) for HMMs.

# Baum-Welch algorithm (I)

# Baum-Welch algorithm (I)

We first define some auxiliary variables:

$$\xi_{t,i,j} = \xi_t(i,j) = P(q_t = s_i, q_{t+1} = s_j | O, \lambda)$$

The probability of being in state  $s_i$  at time  $t$  and in state  $s_j$  at time  $t + 1$ , given the model and the observation sequence.

# Baum-Welch algorithm (I)

We first define some auxiliary variables:

$$\xi_{t,i,j} = \xi_t(i,j) = P(q_t = s_i, q_{t+1} = s_j | O, \lambda)$$

The probability of being in state  $s_i$  at time  $t$  and in state  $s_j$  at time  $t + 1$ , given the model and the observation sequence.

$$\gamma_{t,i} = \gamma_t(i) = P(q_t = s_i | O, \lambda)$$

The probability of being in state  $s_i$  at time  $t$ , given the model and the observation sequence.

# Baum-Welch algorithm (I)

We first define some auxiliary variables:

$$\xi_{t,i,j} = \xi_t(i,j) = P(q_t = s_i, q_{t+1} = s_j | O, \lambda)$$

The probability of being in state  $s_i$  at time  $t$  and in state  $s_j$  at time  $t + 1$ , given the model and the observation sequence.

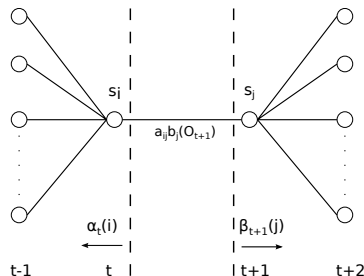
$$\gamma_{t,i} = \gamma_t(i) = P(q_t = s_i | O, \lambda)$$

The probability of being in state  $s_i$  at time  $t$ , given the model and the observation sequence.

From the definitions it follows that:

$$\gamma_t(i) = \sum_{j=1}^N \xi_t(i,j)$$

# Baum-Welch algorithm (II)



**Figure:** Sequence of operations required for the computation of the joint event that the system is in state  $S_i$  at time  $t$  and state  $S_j$  at time  $t+1$  (Rabiner 1989)

$$\alpha_{t,i} = P(o_1, o_2, \dots, o_t, q_t = S_i | \lambda)$$

$$\beta_{t,i} = P(o_{t+1} o_{t+2} \dots o_T | q_t = S_i, \lambda)$$

$$\begin{aligned} \xi_t(i, j) &= \frac{\alpha_{t,i} \cdot a_{i,j} \cdot b_j(o_{t+1}) \cdot \beta_{t+1,j}}{P(O | \lambda)} \\ &= \frac{\alpha_{t,i} \cdot a_{i,j} \cdot b_j(o_{t+1}) \cdot \beta_{t+1,j}}{\sum_{k=1}^N \sum_{l=1}^N \alpha_{t,k} \cdot a_{k,l} \cdot b_l(o_{t+1}) \cdot \beta_{t+1,l}} \end{aligned}$$

## Baum-Welch algorithm (III)

How do these auxiliary variables help?

$$\sum_{t=1}^{T-1} \gamma_t(i) = \text{expected number of transitions from } S_i$$

$$\sum_{t=1}^{T-1} \xi_t(i, j) = \text{expected number of transitions from } S_i \text{ to } S_j$$

# Baum-Welch algorithm (IV)

$\bar{\pi}_i =$  expected no. of times in state  $S_i$  at time  $(t = 1) = \gamma_t(i)$



# Baum-Welch algorithm (IV)

$\bar{\pi}_i =$  expected no. of times in state  $S_i$  at time  $(t = 1) = \gamma_t(i)$

$$\begin{aligned} a_{i,j} &= \frac{\text{expected no. of transitions from } S_i \text{ to } S_j}{\text{expected no. of transition from } S_i} \\ &= \frac{\sum_{t=1}^{T-1} \xi_t(i,j)}{\sum_{t=1}^{T-1} \gamma_t(i)} \end{aligned}$$

# Baum-Welch algorithm (IV)

$\bar{\pi}_i =$  expected no. of times in state  $S_i$  at time  $(t = 1) = \gamma_t(i)$

$$\begin{aligned} \bar{a}_{i,j} &= \frac{\text{expected no. of transitions from } S_i \text{ to } S_j}{\text{expected no. of transition from } S_i} \\ &= \frac{\sum_{t=1}^{T-1} \xi_t(i,j)}{\sum_{t=1}^{T-1} \gamma_t(i)} \end{aligned}$$

$$\begin{aligned} \bar{b}_{j,k}^- &= \frac{\text{expected no. of times in } S_j \text{ observing symbol } v_k}{\text{expected no. of times in } S_j} \\ &= \frac{\sum_{t=1, O_t=v_k}^T \gamma_t(j)}{T} \end{aligned}$$

# Baum-Welch algorithm (V)

The routine for the general case:

```
1  Initialize uniform  $\pi_i$  for  $1 \leq i \leq N$ 
2  Initialize random (stochastic)  $a_{i,j}$ 
3  Initialize uniform  $b_{j,k}$  for  $1 \leq k \leq M$ 
4
5  Repeat until convergence
6      E step:
7          compute auxiliary variables  $\xi_t(i,j)$  and  $\gamma_t(i)$ 
8          using current  $\pi_i$ ,  $a_{i,j}$  and  $b_{j,k}$ 
9
10     M step:
11     compute updated parameter models  $\bar{\pi}_i$ ,  $\bar{a}_{i,j}$ ,  $\bar{b}_{j,k}$ 
```

Baum-Welch Iterative Update

# Baum-Welch - Let's write some code

LET'S WRITE SOME CODE :-)

# Outline

## 2 Machine Learning Applications for HMM

- Machine Learning
- Where do HMMs fit into Machine Learning?

## 3 Theory of HMMs

- The 3 things you want from an HMM
- Mathematical Foundations for HMMs
- Notation Conventions & Framework Description

## 4 Implementing HMMs

- Using the Model for Estimations: the Forward-Backward algorithm
- Learning from Observations: Baum-Welch algorithm
- Uncovering Hidden states: Viterbi algorithm

# Viterbi s-a nascut in ...

# Demo & Discussions

## 5 A Case for HMMs in Symbol Recognition



# A simple symbol recognition application

Features:

# A simple symbol recognition application

Features:

## Define

Define, organize and visualize a dataset of symbols defined with mouse movements.

# A simple symbol recognition application

Features:

## Define

Define, organize and visualize a dataset of symbols defined with mouse movements.

## Train

Train a HMM-based recognition engine on a symbol dataset.

# A simple symbol recognition application

Features:

## Define

Define, organize and visualize a dataset of symbols defined with mouse movements.

## Train

Train a HMM-based recognition engine on a symbol dataset.

## Recognize

Recognize new symbols and view classification metrics.

# A simple symbol recognition application

Features:

## Define

Define, organize and visualize a dataset of symbols defined with mouse movements.

## Train

Train a HMM-based recognition engine on a symbol dataset.

## Recognize

Recognize new symbols and view classification metrics.

Default included symbols: **left arrow**, **right arrow**, **circle**, **square**, **infinity**

# A simple symbol recognition application - A View

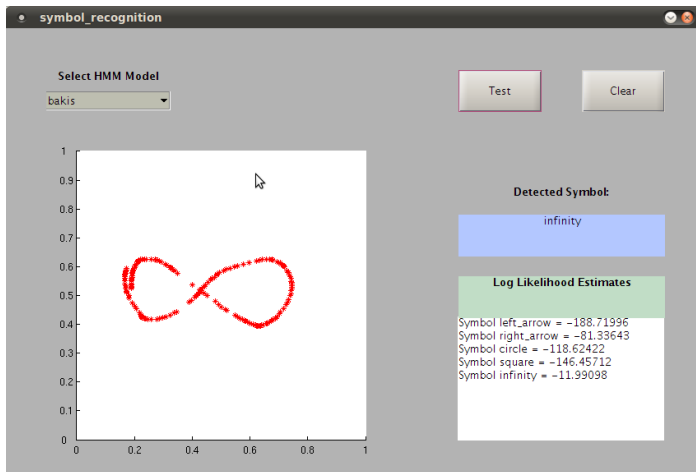
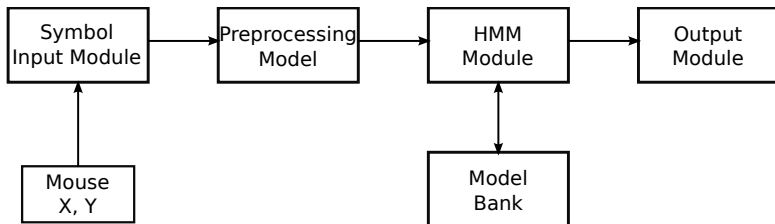


Figure: A view of the symbol recognition application GUI

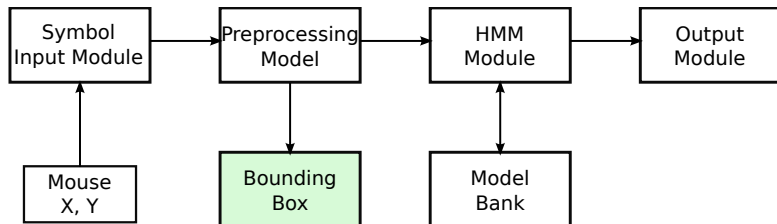
# A simple symbol recognition application - Approach (I)

Adapted from (Yang and Xu 1994).



# A simple symbol recognition application - Approach (I)

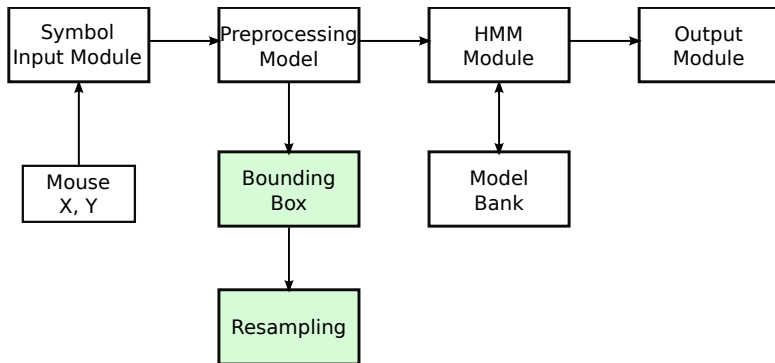
Adapted from (Yang and Xu 1994).





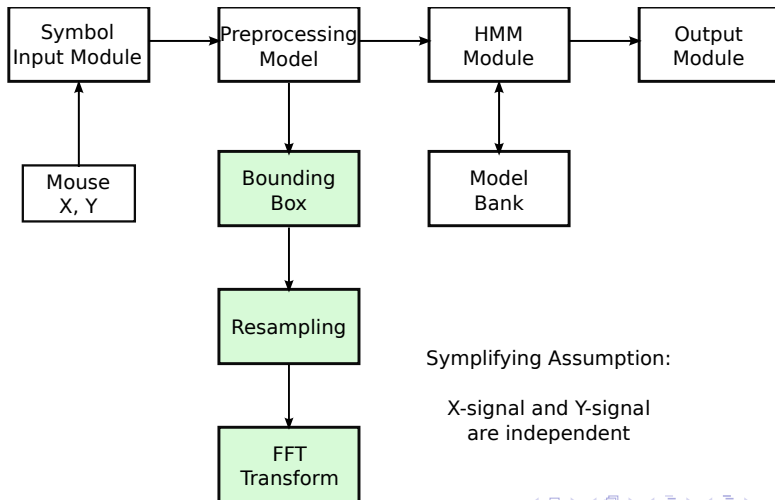
# A simple symbol recognition application - Approach (I)

Adapted from (Yang and Xu 1994).



# A simple symbol recognition application - Approach (I)

Adapted from (Yang and Xu 1994).



Simplifying Assumption:

X-signal and Y-signal  
are independent

# A simple symbol recognition application - Approach (II)

Adapted from (Yang and Xu 1994).

## HMM Structure

$N(\text{number of states}) = 8$

2 discrete observable variables per state -  $\text{coef}_{FFT}(x)$ ,  $\text{coef}_{FFT}(y)$

$M(\text{number of values for each observable variable}) = 256$

Transition model:

- Bakis
- Ergodic

# A simple symbol recognition application - Results

## Dataset size

**5 symbols: left arrow, right arrow, circle, square, infinity**

**100 samples per symbol: 50 training, 10 validation, 40 testing**

```
>> symbol_performance_test('ergodic')
===== Testing trained HMM models =====
## Results for the model of symbol "left_arrow":
Accuracy: 0.97500
Precision: 1.00000
Recall: 0.97500
Confusion matrix line: 39 0 1 0 0 0

## Results for the model of symbol "right_arrow":
Accuracy: 1.00000
Precision: 1.00000
Recall: 1.00000
Confusion matrix line: 0 40 0 0 0 0

## Results for the model of symbol "circle":
Accuracy: 0.90244
Precision: 0.97368
Recall: 0.92500
Confusion matrix line: 0 0 37 2 1 0

## Results for the model of symbol "square":
Accuracy: 0.95238
Precision: 0.95238
Recall: 1.00000
Confusion matrix line: 0 0 0 40 0 0

## Results for the model of symbol "infinity":
Accuracy: 0.97561
Precision: 0.97561
Recall: 1.00000
Confusion matrix line: 0 0 0 0 40 0
```

```
>> symbol_performance_test('bakis')
===== Testing trained HMM models =====
## Results for the model of symbol "left_arrow":
Accuracy: 0.90000
Precision: 1.00000
Recall: 0.90000
Confusion matrix line: 36 0 1 0 0 3

## Results for the model of symbol "right_arrow":
Accuracy: 1.00000
Precision: 1.00000
Recall: 1.00000
Confusion matrix line: 0 40 0 0 0 0

## Results for the model of symbol "circle":
Accuracy: 0.97561
Precision: 0.97561
Recall: 1.00000
Confusion matrix line: 0 0 40 0 0 0

## Results for the model of symbol "square":
Accuracy: 0.97500
Precision: 1.00000
Recall: 0.97500
Confusion matrix line: 0 0 0 39 0 1

## Results for the model of symbol "infinity":
Accuracy: 1.00000
Precision: 1.00000
Recall: 1.00000
Confusion matrix line: 0 0 0 0 40 0
```