# Hidden Markov Models

NIKOLAY YAKOVETS

# A Markov System

$$s_1, .., s_N$$

$S_2$

$S_1$

$S_3$

# A Markov System

$N$ **states**

$$s_1, .., s_N$$

$S_2$

$S_1$

$S_3$

**modeling weather**

# A Markov System

**state changes over time..**

$$S_1 \quad S_2 \quad S_2 \quad S_3 \quad S_2 \quad S_1$$

$q_t$
$time$

$$q_t \in \{s_1, \ldots, s_N\}$$

# A Markov System

**state changes over time..**



$$q_t \in \{s_1, \ldots, s_N\}$$

**modeling weather**

# A Markov Property

**system is memory less..**



$$P(q_{t+1} = S_j | q_t = S_i) = P(q_{t+1} = S_j | q_t = S_i, \text{any earlier history})$$

# A Markov System

**Directed Graph**



$$P(q_{t+1} = S_j | q_t = S_i)$$

# Weather Prediction

## Initial P

| ☀️ | ⛅ | 🌧️ |
|-----|-----|-----|
| 0.5 | 0.2 | 0.1 |

## Transitional P

|        | ☀️  | ⛅  | 🌧️ |
|--------|-----|-----|-----|
| ☀️     | 0   | 0   | 1   |
| ⛅     | 0.5 | 0.5 | 0   |
| 🌧️    | 0.3 | 0.7 | 0   |

Probability of
3-day forecast?: 🌧️ ⛅ ☀️

P(🌧️)P(⛅|🌧️)P(☀️|⛅)=

0.1 * 0.7 * 0.3 = 0.021

# Towards Hidden Markov

what if can't observe the current state?

for example…

# CRAZY VENDING MACHINE

**Prefers dispensing either Coke or Iced Tea**

# CRAZY VENDING MACHINE

Prefers dispensing either Coke or Iced Tea

Changes its mind all the time

# CRAZY VENDING MACHINE

Prefers dispensing either **Coke** or **Iced Tea**

Changes its mind all the time

We don't know its preference at a given moment

# CRAZY VENDING MACHINE

**observations**

**hidden states**

# CRAZY VENDING MACHINE

observation | state

state(t+1) | state(t)

# e.g.

## Probability of vending?:

# e.g.

**Probability of vending?:**

**Consider all HMM paths:**

$$T(\text{CocaCola} \mid \text{CocaCola})\, O(\text{🥤} \mid \text{CocaCola})\ T(\text{CocaCola} \mid \text{CocaCola})\, O(\text{🧃} \mid \text{CocaCola})\ +$$

# e.g.

**Probability of vending?:**

**Consider all HMM paths:**

$$T(\text{Coca-Cola} \mid \text{Coca-Cola})O(\text{Minute Maid} \mid \text{Coca-Cola}) \ T(\text{Coca-Cola} \mid \text{Coca-Cola})O(\text{Lipton} \mid \text{Coca-Cola}) \ +$$

$$T(\text{Coca-Cola} \mid \text{Coca-Cola})O(\text{Minute Maid} \mid \text{Coca-Cola}) \ T(\text{Lipton} \mid \text{Coca-Cola})O(\text{Lipton} \mid \text{Lipton}) \ +$$

# e.g.

## Probability of vending?:

## Consider all HMM paths:

T( 🟥CocaCola | 🟥CocaCola )O( 🥤MinuteMaid | 🟥CocaCola ) T( 🟥CocaCola | 🟥CocaCola )O( 🥤Lipton | 🟥CocaCola ) +

T( 🟥CocaCola | 🟥CocaCola )O( 🥤MinuteMaid | 🟥CocaCola ) T( 🔴Lipton | 🟥CocaCola )O( 🥤Lipton | 🔴Lipton ) +

T( 🔴Lipton | 🟥CocaCola )O( 🥤MinuteMaid | 🔴Lipton ) T( 🟥CocaCola | 🔴Lipton )O( 🥤Lipton | 🟥CocaCola ) +

# e.g.

**Probability of vending?:**

**Consider all HMM paths:**

$T(\text{Coca-Cola}|\text{Coca-Cola})O(\text{Minute Maid}|\text{Coca-Cola})\ T(\text{Coca-Cola}|\text{Coca-Cola})O(\text{Lipton}|\text{Coca-Cola})\ +$

$T(\text{Coca-Cola}|\text{Coca-Cola})O(\text{Minute Maid}|\text{Coca-Cola})\ T(\text{Lipton}|\text{Coca-Cola})O(\text{Lipton}|\text{Lipton})\ +$

$T(\text{Lipton}|\text{Coca-Cola})O(\text{Minute Maid}|\text{Lipton})\ T(\text{Coca-Cola}|\text{Lipton})O(\text{Lipton}|\text{Coca-Cola})\ +$

$T(\text{Lipton}|\text{Coca-Cola})O(\text{Minute Maid}|\text{Lipton})\ T(\text{Lipton}|\text{Lipton})O(\text{Lipton}|\text{Lipton})\ = \ \ldots$

# Hidden Markov

Set of **states S**:

$$S = \{s_1, .., s_N\}$$

# Hidden Markov

Set of **states S**:

$$S = \{s_1, .., s_N\}$$

Output **alphabet K**:

$$K = \{k_1, \ldots, k_M\} = \{1, \ldots, M\}$$

# Hidden Markov

**Initial** state probabilities Π:

$$\Pi = \{\pi_i\}, i \in S$$

# Hidden Markov

**Initial** state probabilities **Π**:

$$\Pi = \{\pi_i\}, i \in S$$

State **transition** probabilities **A**:

$$A = \{a_{ij}\}, i, j \in S$$

# Hidden Markov

**Initial** state probabilities **Π**:

$$\Pi = \{\pi_i\}, i \in S$$

State **transition** probabilities **A**:

$$A = \{a_{ij}\}, i, j \in S$$

Symbol **emission** probabilities **B**:

$$B = \{b_{ijk}\}, i, j \in S, k \in K$$

# Hidden Markov

**State** sequence **X**:

$$X = (X_1, .., X_{T+1})$$

# Hidden Markov

**State** sequence **X**:

$$X = (X_1, .., X_{T+1})$$

**Output** sequence **O**:

$$O = (o_1, .., o_T)$$

# Fundamental Problems

**Evaluation:**
how **likely** is certain observation **O**?

Given:
μ = (**A**, **B**, **Π**)
**O**

Find:
P(**O**|μ)?

# Naïve Evaluation

$$\boxed{P(O|X,\mu)} = \prod_{t=1}^{T} P(o_t|X_t, X_{t+1}, \mu)$$

$$= b_{X_1 X_2 o_1} b_{X_2 X_3 o_2} \cdots b_{X_T X_{T+1} o_T}$$

# Naïve Evaluation

$$\boxed{P(O|X,\mu)} \;=\; \prod_{t=1}^{T} P(o_t|X_t, X_{t+1}, \mu)$$

$$=\; b_{X_1 X_2 o_1} b_{X_2 X_3 o_2} \cdots b_{X_T X_{T+1} o_T}$$

$$\boxed{P(X|\mu)} = \pi_{X_1} a_{X_1 X_2} a_{X_2 X_3} \cdots a_{X_T X_{T+1}}$$

# Naïve Evaluation

$$\boxed{P(O|X,\mu)} = \prod_{t=1}^{T} P(o_t|X_t, X_{t+1}, \mu)$$

$$= b_{X_1 X_2 o_1} b_{X_2 X_3 o_2} \cdots b_{X_T X_{T+1} o_T}$$

$$\boxed{P(X|\mu)} = \pi_{X_1} a_{X_1 X_2} a_{X_2 X_3} \cdots a_{X_T X_{T+1}}$$

$$\boxed{P(O, X|\mu)} = P(O|X, \mu) P(X|\mu)$$

# Naïve Evaluation

$$\boxed{P(O|X,\mu)} \;=\; \prod_{t=1}^{T} P(o_t|X_t, X_{t+1}, \mu)$$

$$=\; b_{X_1 X_2 o_1} b_{X_2 X_3 o_2} \cdots b_{X_T X_{T+1} o_T}$$

$$\boxed{P(X|\mu)} = \pi_{X_1} a_{X_1 X_2} a_{X_2 X_3} \cdots a_{X_T X_{T+1}}$$

$$\boxed{P(O, X|\mu)} = P(O|X,\mu) P(X|\mu)$$

$$\boxed{P(O|\mu)} \;=\; \sum_{X} \boxed{P(O|X,\mu)} \boxed{P(X|\mu)}$$

$$=\; \sum_{X_1 \cdots X_{T+1}} \pi_{X_1} \prod_{t=1}^{T} a_{X_t X_{t+1}} b_{X_t X_{t+1} o_t}$$

# Naïve Evaluation

$$\boxed{P(O|X,\mu)} = \prod_{t=1}^{T} P(o_t | X_t, X_{t+1}, \mu)$$

$$= b_{X_1 X_2 o_1} b_{X_2 X_3 o_2} \cdots b_{X_T X_{T+1} o_T}$$

$$\boxed{P(X|\mu)} = \pi_{X_1} a_{X_1 X_2} a_{X_2 X_3} \cdots a_{X_T X_{T+1}}$$

$$\boxed{P(O, X|\mu)} = P(O|X,\mu) P(X|\mu)$$

$$\boxed{P(O|\mu)} = \sum_X \boxed{P(O|X,\mu)} \boxed{P(X|\mu)}$$

$$= \sum_{X_1 \cdots X_{T+1}} \pi_{X_1} \prod_{t=1}^{T} a_{X_t X_{t+1}} b_{X_t X_{t+1} o_t}$$

$$(2T+1) \cdot N^{T+1}$$

**calculations!**

# Smarter Evaluation

Use DP! FW-BW Alg.

# Smarter Evaluation

**Use DP! FW-BW Alg.**

**DP Table:**

**state over time**

# Smarter Evaluation

**Use DP! FW-BW Alg.**

**DP Table:**

**state over time**



**store forward variables:**

$$\alpha_i(t) = P(o_1 o_2 \cdots o_{t-1}, X_t = i | \mu)$$

# Smarter Evaluation

compute **forward** variables:

**1. initialization:**

$$\alpha_i(1) = \pi_i$$

# Smarter Evaluation

**compute forward variables:**

**1. initialization:**

$$\alpha_i(1) = \pi_i$$

**2. induction:**

$$a_j(t+1) = \sum_{i=1}^{N} \alpha_i(t) a_{ij} b_{ijo_t}$$

# Smarter Evaluation

**compute forward variables:**

**1. initialization:**

$$\alpha_i(1) = \pi_i$$

**2. induction:**

$$a_j(t+1) = \sum_{i=1}^{N} \alpha_i(t) a_{ij} b_{ijo_t}$$

**3. total:**

$$P(O|\mu) = \sum_{i=1}^{N} \alpha_i(T+1)$$

# Smarter Evaluation

**much lower complexity than naïve:**

$$2N^2T$$

**calculations!**

**vs.**

$$(2T + 1) \cdot N^{T+1}$$

**calculations!**

# Smarter Evaluation

**much lower complexity than naïve:**

$$2N^2T$$ **calculations!** vs. $$(2T + 1) \cdot N^{T+1}$$ **calculations!**

**similarly, can work backwards:**

$$\beta_i(t) = P(o_t \cdots o_T | X_t = i, \mu)$$

# Fundamental Problems

**Inference:**

finding **X** that best explains **O**?

Given:

$\mu = ($**A**, **B**, **Π**$)$

**O**

Find:

$\underset{X}{\text{argmax}}\ P($**X**|**O**$,\mu)$

# Smarter Inference

Again, use DP! Viterbi Algorithm

# Smarter Inference

**Again, use DP! Viterbi Algorithm**

**Store:**

    **probability of the most probable path that leads to a node**

$$\delta_j(t) = \max_{X_1 \cdots X_{t-1}} P(X_1 \cdots X_{t-1}, o_1 \cdots o_{t-1}, X_t = j | \mu)$$

# Smarter Inference

**Again, use DP! Viterbi Algorithm**

**Store:**

probability of the most probable path that leads to a node

$$\delta_j(t) = \max_{X_1 \cdots X_{t-1}} P(X_1 \cdots X_{t-1}, o_1 \cdots o_{t-1}, X_t = j | \mu)$$

backtrack through max solution to find the path

# Smarter Evaluation

**compute the variables (fill in the DP table):**

**1 initialization:**

$$\delta_i(1) = \pi_i$$

# Smarter Evaluation

**compute the variables (fill in the DP table):**

**1 initialization:**

$$\delta_i(1) = \pi_i$$

**2.2 induction:**

$$\delta_j(t+1) = \max_{1 \le i \le N} \delta_i(t) a_{ij} b_{ijo_t}$$

# Smarter Evaluation

**compute the variables (fill in the DP table):**

**1 initialization:**

$$\delta_i(1) = \pi_i$$

**2.2 <span style="color:blue">induction</span>:**

$$\delta_j(t+1) = \max_{1 \leq i \leq N} \delta_i(t) a_{ij} b_{ij o_t}$$

**2.2 <span style="color:red">store backtrace</span>:**

$$\psi_j(t+1) = arg \max_{1 \leq i \leq N} \delta_i(t) a_{ij} b_{ij o_t}$$

# Smarter Evaluation

**3 termination and path readout:**

$$\hat{X}_{T+1} = \arg\max_{1 \le i \le N} \delta_i(T+1)$$

$$\hat{X}_t = \psi_{\hat{X}_{t+1}}(t+1)$$

$$P(\hat{X}) = \max_{1 \le i \le N} \delta_i(T+1)$$

# Fundamental Problems

**Estimation:**

finding **μ** that best explains **O**?

Given:

**O**$_{training}$

Find:

$$\underset{\mu}{argmax}\ P(\mathbf{O}_{training}, \mu)$$

# Estimation: MLE

no known analytic method

# Estimation: MLE

no known **analytic method**
find local max using **iterative hill-climb**

# Estimation: MLE

no known **analytic method**
find local max using **iterative hill-climb**
<span style="color:red">**Baum-Welch: (outline)**</span>
**1  choose a model $\mu$ (perhaps randomly)**

# Estimation: MLE

**no known analytic method**
**find local max using iterative hill-climb**
**Baum-Welch: (outline)**
**1 choose a model** $\mu$ **(perhaps randomly)**
**2 estimate** $P(0|\mu)$

# Estimation: MLE

no known **analytic method**
find local max using **iterative hill-climb**
**Baum-Welch: (outline)**
1  choose a model $\mu$ (perhaps randomly)
2  estimate $P(0|\mu)$
3  choose a revised model $\mu$ to maximize the
   values of the paths used a lot…

# Estimation: MLE

no known **analytic method**
find local max using **iterative hill-climb**
**Baum-Welch: (outline)**
1  choose a model μ (perhaps randomly)
2  estimate P(O|μ)
3  choose a revised model μ to maximize the
   values of the paths used a lot…
4  repeat 1-3, hope to converge on values of μ

# When HMMs are good..

Observations are <span style="color:red">ordered</span>

Random process can be represented by a <span style="color:red">stochastic finite state machine</span> with emitting states

# Why HMMs are good..

1. Statistical Grounding
2. Modularity
3. Transparency of a Model
4. Incorporation of Prior Knowledge

# Why HMMs are bad..

1. Markov Chains
2. Local Maxima/Over Fitting
3. Slower Speed

# Speech Recognition



**given an audio waveform, would like to robustly extract & recognize any spoken words**

# Target Tracking



*Radar-based tracking of multiple targets*

*Visual tracking of articulated objects*

**estimate motion of targets in 3D world from indirect, potentially noisy measurements**

# Robot Navigation



CAD
Map

*(S. Thrun,
San Jose Tech Museum)*

Estimated
Map

*Landmark
SLAM
(E. Nebot,
Victoria Park)*

**as robot moves, estimate its world geometry**

# Financial Forecasting



**predict future market behavior from historical data, news reports, expert opinions,..**

# Bioinformatics



**multiple sequence alignment, gene finding, motif/promoter region finding..**

# HMM Applications

HMM can be applied in many more fields where the goal is to recover sequence that is not immediately observable:

- cryptoanalysis
- POS tagging
- MT
- activity recognition
- etc.

# Thank You