

**Hidden Markov Models:
Explanation and Model Learning**



Brian C. Williams
16.410/16.413
Session 21

courtesy of JPL Brian C. Williams, copyright 2000

Reading Assignments

AIMA (Russell and Norvig)

- Ch 15.1-3, 20.3 State Estimation and Hidden Markov Models

From last Monday:

- Ch 13 Review of Probabilities
- Ch 14.1-4 Probabilistic Reasoning


9/13/00 copyright Brian Williams, 2000 2

Outline

- Review
- Explanation and Learning in Statistical Natural Language
 - Decoding using the Viterbi Algorithm
 - Evaluation via Forward and Backward Algorithms.
 - Model learning via the Baum-Welch Algorithm.

9/13/00 copyright Brian Williams, 2000 3


HMM Estimation is Pervasive




courtesy of NASA

Engineering
Operations

Dialogue
Management



courtesy of NASA



Robot
Localization

9/13/00 copyright Brian Williams, 2000 4

Courtesy of Kanna Rajan, NASA Ames. Used with permission.

Posterior Probability, after Observations $X_{1:n} = \mathbf{x}_{1:n}$

$$P(M) = \prod_{M_i \in M} P(M_i)$$

Assume:
 • Apriori mode independence.
 • Consistent obs equally likely

$$P(M | x_{1:n}) = a P(x_n | M) P(M | x_{1:n-1})$$

$P(x_i | M)$ is estimated using model, F , according to:

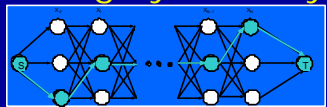
If previous observations $X_{1:i-1} = \mathbf{x}_{1:i-1}$, M and F entails $X_i = x_i$
Then $P(x_i | M) = 1$

If previous observations $X_{1:i-1} = \mathbf{x}_{1:i-1}$, M and F entails $X_i \neq x_i$? v_i
Then $P(x_i | M) = 0$

Otherwise, Assume all consistent assignments to X_i are equally likely observations:
 let $D_i = \{x_c \in D_{X_i} | c, F \text{ is consistent with } X_i = x_c\}$
Then $P(x_i | M) = 1/|D_i|$

9/13/00 copyright Brian Williams, 2000 5

Estimating Dynamic Systems



Given a sequence of observations and commands:

- What is the likelihood of a particular state?
 ⇒ **Belief State Update**: (filtering and smoothing)
- What is the most likely sequence of states that got me here?
 ⇒ **Decoding**: (Viterbi Algorithm)
- What is the most likely sequence of observations generated?
 ⇒ **Evaluation/Prediction**:
- What HMM most likely generated these observations?
 ⇒ **Learning**: (Baum-Welch Algorithm, Expectation-Maximization Algorithm)

9/13/00 copyright Brian Williams, 2000 6

What is the likelihood of a state?

The diagram shows a horizontal timeline from 1 to t. A shaded bar represents the state over time. Above the bar, 'Smoothing' is indicated by a double-headed arrow from 1 to t. 'Filtering' is indicated by a vertical line at time t. 'Prediction' is indicated by a double-headed arrow from t to the end of the timeline.

- Filtering
 - Probabilities of current states
- Prediction
 - Probabilities of future states
- Smoothing
 - Probabilities of past states

9/13/00 copyright Brian Williams, 2000 7

Notation

- S^{t+1} : set of hidden variables in the t+1 time slice
- s^{t+1} : set of values for those hidden variables at t+1
- x^{t+1} : set of observations at time t+1
- $x^{1:t}$: set of observations from all times from 1 to t
- a: normalization constant

9/13/00 copyright Brian Williams, 2000 8

Hidden Markov Models

- Finite States S, Actions A & Observations Ω
- State transition function
 $T(S^i, A^i, S^{i+1}) \equiv P(S^{i+1} | S^i, A^i)$
- Observation function
 $O(S^i, \Omega^i) \equiv P(\Omega^i | S^i)$
- Initial state distribution
 $\Theta(S): P(S^1)$

Notation:
 $\Pi(S)$ denotes
 all subsets of S

9/13/00 copyright Brian Williams, 2000 9

Markov Assumptions

Given a distribution over the current state, the future states and current and future observations are independent of the past.

- **First-order Markov process**
 - $P(S^i | S^{0:t-1}) = P(S^i | S^{t-1})$
- **Markov assumption of evidence**
 - $P(X^i | S^{0:t}, X^{0:t-1}) = P(X^i | S^i)$

9/13/00 copyright Brian Williams, 2000 10

Belief Update Example

$$b(S^{i+1}) = aO(x^{i+1}, S^{i+1}) \sum_{s^i \in S^i} T(S^i, a^i, S^{i+1}) b(S^i)$$

Observed sequence: H T H H H H T H

C_1 0.5 $a \times 0.7 \times [0.9 \times 0.5 + 0.1 \times 0.5] = 0.35a = 0.64$
 C_2 0.5 $a \times 0.4 \times [0.1 \times 0.5 + 0.9 \times 0.5] = 0.20a = 0.36$

$\sum_{s_i \in S^i} P(s_i) = 1$, hence $0.35a + 0.20a = 1 \Rightarrow a = 1.82$

9/13/00 copyright Brian Williams, 2000 11

Diagnosing Dynamic Systems: Via Probabilistic Constraint Automata

- Devices modes
- Probabilistic transitions between modes
- State constraints for each mode
- One automata per component

Open: vlv=open => Outflow = M_2^+ (inflow);
 Stuck open: vlv=stuck open => Outflow = M_2^+ (inflow);
 Closed: vlv=closed => Outflow = 0;
 Stuck closed: vlv=stuck closed=> Outflow = 0;

9/13/00 copyright Brian Williams, 2000 12

Outline

- Review
- Explanation and Learning in Statistical Natural Language
 - Decoding using the Viterbi Algorithm
 - Evaluation via Forward and Backward Algorithms
 - Model learning via the Baum-Welch Algorithm

9/13/00 copyright Brian Williams, 2000 13

Variant on a Hidden Markov Model

A HMM is defined as $\langle S, S, W, E \rangle$

- S is the set of states,
- $s_1 \in S$ is the **single start state**,
- W is the set of **observable symbols**,
- E is the set of **"emitting" transitions**.

Differences from earlier HMM:

- Observations are "words"
- Transitions "emit" observations
- A unique start state

A transition is a four tuple such as $\langle s_2, "had", s_3, 0.3 \rangle$
 denoting: $P(s_2 \xrightarrow{"had"} s_3) = 0.3 \equiv P(s_3, "had" | s_2)$

E combines the transition and observation function in the previous HMM def.

Example:

$s_1 = a$
 $S = \{ a, b \}$
 $W = \{ 0, 1 \}$
 $E = \langle a, "1", a, 0.48 \rangle$
 $\langle a, "0", a, 0.48 \rangle$
 $\langle a, "0", b, 0.04 \rangle$
 $\langle b, "1", a, 1.0 \rangle$

Sentence Parsing Example

HMM = $\langle S, s_1, W, E \rangle$ where $S = \{s_1, s_2, s_3, s_4, s_5, s_6, s_7, s_8\}$; $W = \{ "Roger", \dots \}$;
 $E = \{ \text{transition}_1, \dots \}$ and $\text{transition}_1 = \langle s_2, "had", s_3, 0.3 \rangle \dots$

Observed Word Sequences:

S_1 : Mary had a little Lamb and a big dog.
 S_2 : Roger ordered a lamb curry and a hot dog.
 S_3 : John cooked a hot dog curry.

$P(S_3) = 0.3 * 0.3 * 0.5 * 0.5 * 0.3 * 0.5 = 0.003375$

Problems and Algorithms

Decoding:
 Given a sequence of observations, what is the most likely sequence of hidden states?
 - Solution: Viterbi algorithm

Evaluation: What is the probability of a given sequence of observations?
 - Solution: forward/backward algorithm
 • (These are used in the learning algorithm).

Learning:
 Given a sequence of observations, what HMM transition probabilities maximize the likelihood of the sequence?
 - Solution: Baum-Welch algorithm (form of Expectation-Maximization)

9/13/00 copyright Brian Williams, 2000 16

Outline

- Overview
- Belief State Update
- Explanation and Learning in Statistical Natural Language
 - Decoding using the Viterbi Algorithm
 - Evaluation via Forward and Backward Algorithms
 - Model learning via the Baum -Welch Algorithm.

9/13/00 copyright Brian Williams, 2000 17

HMM Decoding

Finding the most likely state trajectory

Problem: Given HMM = $\langle S, s_i, W, E \rangle$, and observation sequence $w^{1:t-1}$, find the most likely state sequence (denoted $\sigma_i(t)$), ending in s_i at time t :

$$\sigma_i(t) = \arg \max_{s^{1:t} \text{ s.t. } s^t = s_i} P(s^{1:t} | w^{1:t-1})$$

Observe $w^{1:4} = \langle 1, 1, 1, 0 \rangle$

	1	1	1	0
a	.2	.2	.2	.4
b	.1	.1	.1	.3

$\sigma_a^5 = \langle a, b, b, b, a \rangle$

9/13/00 copyright Brian Williams, 2000 18

Probability of the Most Likely State Trajectory

Probability of the most likely sequence $\mathbf{s}_i(t)$ ending in s_i given observations $w^{1:t}$:

$$P(\mathbf{s}_i(t+1) | w^{1:t}) = \frac{P(\mathbf{s}_i(t+1), w^{1:t})}{P(w^{1:t})}$$

$$= \mathbf{a} \max_{s_j \in S} P(s_i^{t+1}, w^t | s_j^t) P(s_j^{t-1}, s_j^t | w^{1:t-1})$$

$$P(\mathbf{s}_i(t+1) | w^{1:t}) = \mathbf{a} \max_{s_j \in S} P(s_i^{t+1}, w^t | s_j^t) P(\mathbf{s}_j(t) | w^{1:t-1})$$

$P(\mathbf{s}_i(1)) = 1$, if $s_i^1 = s_i$, and 0, otherwise.

9/13/00 copyright Brian Williams, 2000 19

Viterbi Algorithm

Problem: Compute $\mathbf{s}_i(t) = \arg \max_{s^t, s^t, s^t = s_i} P(s^t | w^{1:t-1})$

Solution: For n from 1 to t,

- Compute the most likely paths of length n that end at each $s_k \in S$.
- Extend to the most likely paths of length n+1 that end at each $s_p \in S$

$$\mathbf{s}_i(1) = \langle s_i \rangle$$

$$\mathbf{s}_i(t+1) = \mathbf{s}_{j_{\max}}(t) \circ s_i,$$

where $j_{\max} = \arg \max_{k=1}^{|S|} P(\mathbf{s}_k(t) | w^{1:t-1}) P(s_k \rightarrow s_i)$

9/13/00 copyright Brian Williams, 2000 20

Example: Viterbi Algorithm

$\mathbf{s}_i(1) = \langle s_i \rangle$

$\mathbf{s}_i(t+1) = \mathbf{s}_{j_{\max}}(t) \circ s_i$, where $j_{\max} = \arg \max_{k=1}^{|S|} P(\mathbf{s}_k(t) | w^{1:t-1}) P(s_k \rightarrow s_i)$

States\ Obs	Start	1	11	111	1110	
a	Sequence	a	aa	aaa	aaaa	abba
	Probability	1.0	0.2	0.04	0.008	0.005
b	Sequence	b	ab	abb	abbb	abbbb
	Probability	0.0	0.1	0.05	0.025	0.005

Observe "1110"

9/13/00 copyright Brian Williams, 2000 21

Viterbi Pseudo Code

```

Viterbi(<S,s1,W,E>, w1:T, T) // param αi(t) denotes P(αi(t)|w1:t-1)
1. begin
2. for (i ← 1; i ≤ |S|; i ← i+1) {
3.   initialize αi(1) ← {si}, αi(1) ← {1 if si is start state s1, else 0}
4.   for (t ← 1; t < T; t ← t+1) { // for each obs wt ∈ w1:T
5.     for (i ← 1; i ≤ |S|; i ← i+1) { // each state sit+1 at time t+1
6.       initialize jmax ← 0, Pmax ← -1;
7.       for (k ← 1; k ≤ |S|; k ← k+1) // compute arg max over skt
8.         if (αk(t) * P(sit+1, wt|skt) > Pmax) {
9.           Pmax ← αk(t) * P(sit+1, wt|skt); jmax ← k ;
10.        αi(t+1) ← Pmax ;
11.        αi(t+1) ← σjmax(t) ◦ si } // postpend the next state
12.   return αi(T) for 0 < i ≤ |S| that maximizes αi(T)
13. end
    
```

Outline

- Review
- Explanation and Learning in Statistical Natural Language
 - Decoding using the Viterbi Algorithm
 - Evaluation via Forward and Backward Algorithms
 - Model learning via the Baum-Welch Algorithm
- Appendix: Monitoring and Diagnosis via Probabilistic Constraint Automata

9/13/00 copyright Brian Williams, 2000 23

Probability of Observation Sequence

Forward probability $\mathbf{a}_i(t) = P(w^{1:t-1}, s_i^t)$ Similar to belief state, given earlier.

$$\mathbf{a}_i(1) = P(w^{1:0}, s_i^1) = \begin{cases} i = 1 \rightarrow 1.0 \\ \text{otherwise} \rightarrow 0 \end{cases}$$

$$\mathbf{a}_i(t+1) = \sum_{j=1}^c P(s_j \xrightarrow{w^t} s_i) \mathbf{a}_j(t) = \sum_{j=1}^c P(s_j \xrightarrow{w^t} s_i) P(w^{1:t-1}, s_j^t)$$

Observation probability $P(w^{1:T}) = \sum_{i=1}^{|S|} \mathbf{a}_i(T)$

HMM Evaluation

Observation and Forward probabilities

“1110”

t	1	2	3	4	5
ϵ	1	1	1	0	
$\alpha_a(t)$	1.0	0.2	0.05	0.017	0.0148
$\alpha_b(t)$	0.0	0.1	0.07	0.04	0.0131
$P(w^{1:t})$	1.0	0.3	0.12	0.057	0.0279

$$a_i(t+1) = \sum_{j=1}^c P(s_j \rightarrow \xi) a_j(t)$$

$$P(w^{1:t}) = \sum_{i=1}^{|S|} a_i(t+1)$$

$0.2 * 0.1 = 0.02$
 $+$
 $0.1 * 0.5 = 0.05$

Forward Algorithm Pseudo Code

Forward($\langle S, s_1, W, E \rangle, w^{1:T}$) // param $\alpha_i(t)$ denotes $P(s_i | w^{1:t-1})$

1. **begin**
2. **for** $1 \leq i \leq |S|$
3. **initialize** $\alpha_i(1) \leftarrow \{1 \text{ if } s_i \text{ is start state } s_1, \text{ else } 0\}$;
4. **for** $1 \leq t \leq T$ {
5. **for** $1 \leq j \leq |S|$ {
6. **initialize** $\alpha_j(t+1) \leftarrow 0$;
7. **for** $1 \leq i \leq |S|$
8. $\alpha_j(t+1) \leftarrow \alpha_j(t+1) + P(s_j, w^t | s_i) * \alpha_i(t)$;
9. }
10. }
11. **return** ? $\alpha_i(T+1)$ for all $s_i \in S$
12. **end**

9/13/00
copyright Brian Williams, 2000
26

What is the likelihood of a state?

- **Filtering**
 - Probabilities of current states $P(S^t | w^{1:t})$
- **Prediction**
 - Probabilities of future states $P(S^k | w^{1:t})$ for $k > t$
- **Smoothing**
 - Probabilities of past states $P(S^k | w^{1:t})$ for $k < t$

9/13/00
copyright Brian Williams, 2000
27

Smoothing

$$P(S^k | w^{1:T}) = P(S^k | w^{1:k}, w^{k+1:T}) \quad \text{Divide obs}$$

$$= \mathbf{a} P(S^k | w^{1:k}) P(w^{k+1:T} | S^k, w^{1:k}) \quad \text{Bayes}$$

$$= \mathbf{a} P(S^k | w^{1:k}) P(w^{k+1:T} | S^k) \quad \text{Markov}$$

$P(S^k | w^{1:T}) = \mathbf{a} \mathbf{a}_i(k) \mathbf{b}_i(k)$

9/13/00 copyright Brian Williams, 2000 28

Backward Probabilities

Backward probability $\mathbf{b}_i(t) \equiv P(w^{t:T} | s_i^t)$

$$\mathbf{b}_i(T+1) = P(\mathbf{e} | s_i^{T+1}) = 1$$

$\beta_i(t)$ is similar to $\alpha_i(t)$ but starts from the end.

$$\mathbf{b}_j(t-1) = \sum_{j=1}^{|S|} P(s_i \rightarrow s_j) \mathbf{b}_j(t) = \sum_{j=1}^{|S|} P(s_i \rightarrow s_j) \mathcal{P}(w^{t:T} | s_j^t)$$

Observation probability $P(w^{1:T}) = \mathbf{b}_i(1) = P(w^{1:T} | s^1 = s_i^1)$

9/13/00 copyright Brian Williams, 2000 29

Outline

- Review
- Explanation and Learning in Statistical Natural Language
 - Decoding using the Viterbi Algorithm
 - Evaluation via Forward and Backward Algorithms
 - Model learning via the Baum-Welch Algorithm

9/13/00 copyright Brian Williams, 2000 30

HMM Training (Baum-Welch Algorithm)

Approach: Given a training sequence $w^{1:T}$, adjust the HMM state transition probabilities to make the observation sequence as likely as possible.

Training Sequence: $w^{1:8} = 01010210$

9/13/00 copyright Brian Williams, 2000 31

Dealing with Hidden States

Intuitively...

Problem: States are not known.
Solution: Estimate states from Model.

Problem: Transitions no longer deterministic.
Solution: Compute expected # of transitions.

Problem: Model not known.(chicken & egg)
Solution: Bootstrap the model:

1. Guess a model (transition probabilities).
2. Use model to estimate states.
3. Count estimated transitions to get model.

When counting transitions "prorate" each transition by its probability.

9/13/00 copyright Brian Williams, 2000 32

Expectation-Maximization (Baum-Welch)

1. Guess a set of transition probabilities.
2. while (transition probabilities improving) {
 - a. **Expectation:** Use transition P to estimate states $P(S^t | w^{1:T})$.
 - b. **Maximization:** Estimate new transition probabilities by counting expected # of transitions, given state estimates.

- "improvement" measured by comparing *cross-entropy* after each iteration:

$$-\frac{1}{n} \sum_{w^{1:T}} P_{M-1}(w^{1:T}) \log_2 P_M(w^{1:T})$$
- Terminate when change in cross-entropy is less than some θ .

9/13/00 copyright Brian Williams, 2000 33

Estimating The Transition Probability

- $C(s_i, w_k, s_j)$: The expected "count" of transitions $s_i \rightarrow s_j$, during observation sequence $w^{1:T}$:

$$C(s_i \rightarrow s_j) = \sum_{t=1}^T \mathbf{a}_i(t) \mathcal{P}(s_i \rightarrow s_j) \mathbf{b}_j(t+1)$$
- P_e : Estimated transition probability for $s_i \rightarrow s_j$ estimated from observation sequence $w^{1:T}$:

$$P_e(s_i \rightarrow s_j) = \frac{C(s_i \rightarrow s_j)}{\sum_{l=1}^{|S|} \sum_{m=1}^{|W|} C(s_i \rightarrow s_j)}$$

} $C(s_i)$: Expected count of transitions out of s_i

9/13/00 copyright Brian Williams, 2000 34

Baum-Welch Pseudo Code

Baum-Welch ($P_{new}, w^{1:T}, \theta$) // P_{new} = estimated $P(s_i, w_k | s_j)$

1. **do** (// $w^{1:T}$ is a training sequence and θ is a convergence criteria)
2. **for** $l = 1, j = |S|, l = k = |W|$
3. $P_{old}(s_i, w_k | s_j) = P_{new}(s_i, w_k | s_j)$; // remember old probability estimate
4. **compute** $\alpha_i(t), \beta_j(t)$, for all values of $l = |S|$ and $1 = t = T$;
5. **for** $l = 1, j = |S|, l = k = |W|$, and $1 = t = T$ {
6. **initialize** $C(s_i, w_k, s_j) \leftarrow 0$;
7. **for** $l = t = T$
8. $C(s_i, w_k, s_j) \leftarrow \alpha_i(t) P_{old}(s_i, w_k | s_j) \beta_j(t)$;
9. }
10. **for** $l = 1 = |S|$ {
11. **initialize** $C(s_i) \leftarrow 0$;
12. **for** $l = j = |S|, l = k = |W|$
13. $C(s_i) \leftarrow C(s_i) + C(s_i, w_k, s_j)$;
14. **for** $l = j = |S|, l = k = |W|$
15. $P_{new}(s_i, w_k | s_j) = C(s_i, w_k, s_j) / C(s_i)$;
16. }
17. } **while** ($\max \text{Changed}(P_{new}, P_{old}) > \theta$)

Baum-Welch example

Transition probabilities are initially guessed.
Training sequence is "01011"

	1	2	3	4	5	6
ϵ	0	1	0	1	1	
$\alpha_i(t)$	1	0.48	0.27	0.13	0.072	0.035
$\alpha_i(t)$	0	0.04	0	0.01	0	0

	1	2	3	4	5	6
ϵ	0	1	0	1	1	
$\beta_j(t)$	0.035	0.062	0.13	0.23	0.48	1
$\beta_j(t)$	0	0.13	0	0.28	1	1

	0	1	0	1	1	Total	New P
$T(a,0,b)$	0.0052	0	0.0052	0	0	0.01	0.06
$T(b,1,a)$	0	0.0052	0	0.0048	0	0.01	1
$T(a,0,a)$	0.030	0	0.03	0	0	0.06	0.26
$T(a,1,a)$	0	0.03	0	0.03	0.035	0.095	0.58

} $C(a,0,b) + C(a,0,a) + C(a,1,a)$

Notation Summary

Probabilistic transitions:

- written as $P(s_3 \text{ "had" } | s_2) = 0.3$ or as $P(s_2 \xrightarrow{\text{"had"}} s_3) = 0.3$

Observation sequences:

- $w^{1:T}$ denotes the entire sequence of observations.
- ϵ denotes the empty sequence (no observations).

States S

- States are subscripted, $s_i \in S$, where $1 = i = |S|$.
- Superscripts indicate time, for example, s^k is the k^{th} state in a state sequence.

State sequences:

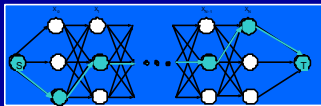
- $\alpha(t)$ denotes the most likely sequence of t states that ends in state s_i .
- $\alpha(t-1) \circ s^t$ concatenates state to the end of the sequence.
- $\alpha(t) \equiv P(s_i | w^{1:t-1})$ denotes the forward probability at time step t of $s^t = s_i$.
- $\beta_i(t) \equiv P(w^{t+1} | s_t^i)$ denotes the backward probability at time step t of $s^t = s_i$.

9/13/00

copyright Brian Williams, 2000

37

Estimating Dynamic Systems



Given a sequence of observations and commands:

- What is the likelihood of a particular state?
 - ⇒ **Belief State Update:** (filtering, smoothing, prediction)
- What is the most likely sequence of states that got me here?
 - ⇒ **Decoding:** (Viterbi Algorithm)
- What is the most likely sequence of observations generated?
 - ⇒ **Evaluation:**
- What HMM most likely generated these observations?
 - ⇒ **Learning:** (Baum-Welch Algorithm, Expectation-Maximization Algorithm)

9/13/00

copyright Brian Williams, 2000

38
