

Applications of HMMs

Any sequence tagging problem

- Named Entity Recognition

Alan Spoon , recently named Newsweek president ,
said Newsweek ad rates would increase in January .

- Morpheme Boundary Detection

Signs of a slowing economy are increasing

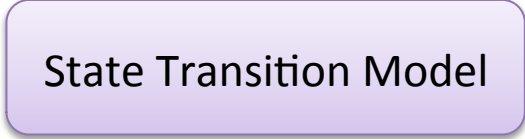

- Simple Machine Translation (no reordering)

das blaue Haus
the blue house

- Speech Recognition testing one



Generative Story for HMMs

- To create observation sequence $o_1 o_2 \dots o_n$
 - For each time step i , generate a latent state s_i conditioned on s_{i-1} A light purple rounded rectangular button with the text "State Transition Model" centered inside.
 - For each time step i , generate an observation symbol o_i conditioned on s_i A light purple rounded rectangular button with the text "Channel Model" centered inside.

HMM Decoding

- Given observation, find best sequence of latent states

State Transition Model

$$\operatorname{argmax}_{\text{state sequences } s} P(s | o) = \operatorname{argmax}_{\text{state sequences } s} P(s)P(o | s)$$

Channel Model

- What is the size of the search space (number of state sequences) for an observation of length n ?

Dynamic Programming!

HMM Decoding Example

- I see this word from some language

Ω Ω Ω

- Which symbols are vowels and which are consonants?
 - Given: state (consonant-vowel) transition probabilities, channel "emission" probabilities

Viterbi Algorithm

V			
C			
#	Ω	⊖	⊠

Prob of state sequence that maximizes $P(s)P(o|s)$ where $o = \underline{\Omega} \ominus$ and s ends with C

Given these probabilities:
 $P(C | \#) = 0.7$
 $P(C | C) = 0.4$
 $P(V | V) = 0.1$

$P(\underline{\Omega} | C) = 0.08$
 $P(\underline{\Omega} | V) = 0.01$
 $P(\ominus | C) = 0.02$
 $P(\ominus | V) = 0.14$
 $P(\boxplus | C) = 0.07$
 $P(\boxplus | V) = 0.20$

Viterbi Algorithm

V			
C			
#	Ω	Ϡ	⊠

$$P(C|\#)P(\underline{\Omega}|C) \\ =0.7*0.08=0.056$$

Given these probabilities:

$$P(C|\#)=0.7$$

$$P(C|C)=0.4$$

$$P(V|V)=0.1$$

$$P(\underline{\Omega}|C)=0.08$$

$$P(\underline{\Omega}|V)=0.01$$

$$P(\underline{\rho}|C)=0.02$$

$$P(\underline{\rho}|V)=0.14$$

$$P(\boxplus|C)=0.07$$

$$P(\boxplus|V)=0.20$$

Viterbi Algorithm

V			
C	0.056		
#	Ω	Ϡ	⊠

$$P(V|\#)P(\underline{\Omega}|V) = 0.3 * 0.01 = 0.003$$

Given these probabilities:

$$P(C|\#)=0.7$$

$$P(C|C)=0.4$$

$$P(V|V)=0.1$$

$$P(\underline{\Omega}|C)=0.08$$

$$P(\underline{\Omega}|V)=0.01$$

$$P(\underline{\rho}|C)=0.02$$

$$P(\underline{\rho}|V)=0.14$$

$$P(\boxplus|C)=0.07$$

$$P(\boxplus|V)=0.20$$

Viterbi Algorithm

V	0.003		
C	0.056		
#	Ω	⊖	⊠

$0.056 * P(C|C)P(\ominus|C)$
 $= 0.056 * 0.4 * 0.02 = 0.000448$
 $0.003 * P(C|V) * P(\ominus|C)$
 $= 0.003 * 0.9 * 0.02 = 0.000054$

Given these probabilities:

$P(C|#) = 0.7$

$P(C|C) = 0.4$

$P(V|V) = 0.1$

$P(\Omega|C) = 0.08$

$P(\Omega|V) = 0.01$

$P(\ominus|C) = 0.02$

$P(\ominus|V) = 0.14$

$P(\boxtimes|C) = 0.07$

$P(\boxtimes|V) = 0.20$

Viterbi Algorithm

V	0.003		
C	0.056		
#	Ω	ϑ	⊠

$$0.056 * P(V|C) * P(\vartheta|V)$$

$$= 0.056 * 0.6 * 0.14 = 0.004704$$

$$0.003 * P(V|V) * P(\vartheta|V)$$

$$= 0.003 * 0.1 * 0.14 = 0.000042$$

Given these probabilities:

$$P(C|\#) = 0.7$$

$$P(C|C) = 0.4$$

$$P(V|V) = 0.1$$

$$P(\Omega|C) = 0.08$$

$$P(\Omega|V) = 0.01$$

$$P(\vartheta|C) = 0.02$$

$$P(\vartheta|V) = 0.14$$

$$P(\boxtimes|C) = 0.07$$

$$P(\boxtimes|V) = 0.20$$

Viterbi Algorithm

V	0.003	0.004704	
C	0.056	0.000448	
#			

$$0.000448 * P(C|C)P(\uparrow | C)$$

$$= 0.000448 * 0.4 * 0.07 = 0.0000125$$

$$0.004704 * P(C|V) * P(\uparrow | C)$$

$$= 0.004704 * 0.9 * 0.07 = 0.000296$$

Given these probabilities:

$$P(C| \#) = 0.7$$

$$P(C|C) = 0.4$$

$$P(V|V) = 0.1$$

$$P(\ominus | C) = 0.08$$

$$P(\ominus | V) = 0.01$$

$$P(\omin� | C) = 0.02$$

$$P(\omin� | V) = 0.14$$

$$P(\boxtimes | C) = 0.07$$

$$P(\boxtimes | V) = 0.20$$

Viterbi Algorithm

V	0.003	0.004704	
C	0.056	0.000448	0.296
#			

$$\begin{aligned}
 & 0.000448 * P(V|C) * P(\uparrow|V) \\
 & = 0.000448 * 0.6 * 0.20 = 0.00005376 \\
 & 0.004704 * P(V|V) * P(\uparrow|V) \\
 & = 0.004704 * 0.1 * 0.20 = 0.0000941
 \end{aligned}$$

Given these probabilities:

$$P(C|\#) = 0.7$$

$$P(C|C) = 0.4$$

$$P(V|V) = 0.1$$

$$P(\ominus|C) = 0.08$$

$$P(\ominus|V) = 0.01$$

$$P(\omin�|C) = 0.02$$

$$P(\omin�|V) = 0.14$$

$$P(\boxplus|C) = 0.07$$

$$P(\boxplus|V) = 0.20$$

Viterbi Algorithm

V	0.003	0.004704	0.000094
C	0.056	0.000448	0.000296
#	Ω	⊖	⊠

Given these probabilities:

$$P(C \mid \#) = 0.7$$

$$P(C \mid C) = 0.4$$

$$P(V \mid V) = 0.1$$

$$P(\Omega \mid C) = 0.08$$

$$P(\Omega \mid V) = 0.01$$

$$P(\ominus \mid C) = 0.02$$

$$P(\ominus \mid V) = 0.14$$

$$P(\boxtimes \mid C) = 0.07$$

$$P(\boxtimes \mid V) = 0.20$$

Viterbi Algorithm

V	0.003	0.004704	0.000094
C	0.056	0.000448	0.000296
#	Ω	Ω	⊠

Take the highest probability on the last column and trace back the pointers

Viterbi Algorithm

V	0.003	0.004704	0.000094
C	0.056	0.000448	0.000296
#	Ω	Ω	⊠

Take the highest probability on the last column and trace back the pointers

C
⊠

Viterbi Algorithm

V	0.003	0.004704	0.000094
C	0.056	0.000448	0.000296
#	Ω	Ω	⊠

Take the highest probability on the last column and trace back the pointers

V C
Ω ⊠

Viterbi Algorithm

V	0.003	0.004704	0.000094
C	0.056	0.000448	0.000296
#	Ω	Ω	⊠

Take the highest probability on the last column and trace back the pointers

C V C
Ω Ω ⊠

Viterbi Algorithm

- Complexity?

For sequence of length T and HMM with n states

- Runtime is $O(n^2T)$: $n \times T$ cells, n computations per cell
- Space is $O(nT)$

- Why does it work?

That is, how does it find

$$\underset{\text{state sequences } s}{\operatorname{argmax}} P(s|o) = \underset{\text{state sequences } s}{\operatorname{argmax}} P(s)P(o|s)$$

Convince
yourself
of this.

without enumerating all exponentially many state sequences?

- Best state sequence for first few time steps will be part of eventual best state sequence. Only change comes from transition probability from time t to $t+1$

HMMs: Three Canonical Problems

- Finding the best sequence of states for an observation
 - Done: Viterbi Algorithm.
- Finding the total probability of an observation
 - Coming up on Thursday: Forward Algorithm
- Estimating the transition and emission probabilities from a corpus
 - Coming up next week: Baum-Welch

Why do we need to find the best state sequence?

- Part of speech tagging: resolving ambiguity

Noun Verb Det Noun
I shot an elephant
or Noun?

- Predicting pronunciations of letters

K AA R S
or S? or AE? or sil? or Z?
c a r S

Questions about Viterbi

- Difference between Markov Models and HMMs?
- What is the generative story for HMMs
- Where does $P(s)P(o | s)$ come from?
- What is end result produced by Viterbi?
- Where do the probabilities come from?
- Is there an algorithm-model relationship analogous to Viterbi-HMM?
- What does the Viterbi table really mean?

Why do we need to find the total probability of observation?

- Can use HMMs as a **language model**, instead of n-grams
- Recall: Language models assign probabilities to strings
 - Text Classification
 - Random Generation
 - Machine Translation
 - Speech Recognition

HMM Observation Probability

- Compute total probability of observation as generated by the HMM

$$P(o) = \sum_{\text{state sequences } s} P(s, o) = \sum_{\text{state sequences } s} P(s)P(o | s)$$

State Transition Model

Channel Model

HMM Best State Sequence

- Compare previous slide to this formula for the best state sequence for an observation (last class)



State Transition Model

$$\operatorname{argmax}_{\text{state sequences } s} P(s | o) = \operatorname{argmax}_{\text{state sequences } s} P(s)P(o | s)$$



Channel Model

Recall: Viterbi Algorithm

V

C

#

Ω

\mathcal{C}

\square

Prob of state sequence that maximizes $P(s)P(o/s)$ where

$$o = \underline{\Omega} \mathcal{C}$$

and

s ends with C

What small modification can we make to Viterbi?

Change what each cell represents...

Total Probability of Observation

V			
C			
#	Ω	ω	$\square \wedge$

Marginal prob. of all state sequences where
 $\omega = \Omega \omega$
and
 s ends with C

Forward Algorithm

V			
C			
#	Ω	⊖	⊠

$$P(C \mid \#)P(\Omega \mid C) \\ = 0.7 * 0.08 = 0.056$$

Given these probabilities:

$$P(C \mid \#) = 0.7$$

$$P(C \mid C) = 0.4$$

$$P(V \mid V) = 0.1$$

$$P(\Omega \mid C) = 0.08$$

$$P(\Omega \mid V) = 0.01$$

$$P(\ominus \mid C) = 0.02$$

$$P(\ominus \mid V) = 0.14$$

$$P(\boxtimes \mid C) = 0.07$$

$$P(\boxtimes \mid V) = 0.20$$

Forward Algorithm

V			
C	0.056		
#	Ω	⊖	⊠

$$P(V \mid \#)P(\underline{\Omega} \mid V) = 0.3 * 0.01 = 0.003$$

Given these probabilities:

$$P(C \mid \#) = 0.7$$

$$P(C \mid C) = 0.4$$

$$P(V \mid V) = 0.1$$

$$P(\underline{\Omega} \mid C) = 0.08$$

$$P(\underline{\Omega} \mid V) = 0.01$$

$$P(\ominus \mid C) = 0.02$$

$$P(\ominus \mid V) = 0.14$$

$$P(\boxplus \mid C) = 0.07$$

$$P(\boxplus \mid V) = 0.20$$

Forward Algorithm

V	0.003		
C	0.056		
#	Ω	⊖	⊠

$$0.056 * P(C|C)P(\ominus|C)$$

$$= 0.056 * 0.4 * 0.02 = 0.000448$$

$$0.003 * P(C|V) * P(\ominus|C)$$

$$= 0.003 * 0.9 * 0.02 = 0.000054$$

Given these probabilities:

$$P(C|\#) = 0.7$$

$$P(C|C) = 0.4$$

$$P(V|V) = 0.1$$

$$P(\Omega|C) = 0.08$$

$$P(\Omega|V) = 0.01$$

$$P(\ominus|C) = 0.02$$

$$P(\ominus|V) = 0.14$$

$$P(\boxplus|C) = 0.07$$

$$P(\boxplus|V) = 0.20$$

Forw

$$0.056 * P(V|C)P(\text{☉} | V)$$

$$= 0.056 * 0.6 * 0.14 = 0.004704$$

$$0.003 * P(V|V) * P(\text{☉} | V)$$

$$= 0.003 * 0.1 * 0.14 = 0.000042$$

V	0.003		
C	0.056	0.000502	
#	Ω	☉	⊠

Given these probabilities:

- $P(C|#)=0.7$
- $P(C|C)=0.4$
- $P(V|V)=0.1$

- $P(\text{Ω} | C)=0.08$
- $P(\text{Ω} | V)=0.01$
- $P(\text{☉} | C)=0.02$
- $P(\text{☉} | V)=0.14$
- $P(\text{⊠} | C)=0.07$
- $P(\text{⊠} | V)=0.20$

For

$$0.000502 * P(C|C)P(\boxplus|C) \\ = 0.000502 * 0.4 * 0.07 = 0.000014$$

$$0.004746 * P(C|V) * P(\boxplus|C) \\ = 0.004746 * 0.9 * 0.07 = 0.000299$$

V	0.003	0.004746	
C	0.056	0.000502	
#	Ω	∞	⊔

in these probabilities:

$$P(C|\#) = 0.7$$

$$P(C|C) = 0.4$$

$$P(V|V) = 0.1$$

$$P(\Omega|C) = 0.08$$

$$P(\Omega|V) = 0.01$$

$$P(\infty|C) = 0.02$$

$$P(\infty|V) = 0.14$$

$$P(\boxplus|C) = 0.07$$

$$P(\boxplus|V) = 0.20$$

For

$$0.000502 * P(V|C)P(\boxplus|V)$$

$$= 0.000502 * 0.6 * 0.2 = 0.000060$$

$$0.004746 * P(V|V) * P(\boxplus|V)$$

$$= 0.004746 * 0.1 * 0.2 = 0.000095$$

V	0.003	0.004746	
C	0.056	0.000502	0.000313
#	Ω	⊗	⊠

in these probabilities:

$$P(C|\#) = 0.7$$

$$P(C|C) = 0.4$$

$$P(V|V) = 0.1$$

$$P(\underline{\Omega} | C) = 0.08$$

$$P(\underline{\Omega} | V) = 0.01$$

$$P(\otimes | C) = 0.02$$

$$P(\otimes | V) = 0.14$$

$$P(\boxplus | C) = 0.07$$

$$P(\boxplus | V) = 0.20$$

Forward Algorithm

V	0.003	0.004746	0.000016
C	0.056	0.000502	0.000313
#	Ω	⊖	⊠

Given these probabilities:

$$P(C \mid \#) = 0.7$$

$$P(C \mid C) = 0.4$$

$$P(V \mid V) = 0.1$$

$$P(\Omega \mid C) = 0.08$$

$$P(\Omega \mid V) = 0.01$$




$$P(\ominus \mid C) = 0.02$$

$$P(\ominus \mid V) = 0.14$$

$$P(\boxplus \mid C) = 0.07$$

$$P(\boxplus \mid V) = 0.20$$

Forward Algorithm

V	0.003	0.004746	0.000016
C	0.056	0.000502	0.000313
#			

Total probability
of $\underline{\Omega}\underline{\Omega}\underline{\Omega}$ under
this model:

$$\begin{aligned} &0.000016 \\ &+ \\ &0.000313 \\ &= 0.000329 \end{aligned}$$

Where do the parameters come from?

- Who gives these to us?
- Ideally, we want to estimate them from data
- Remember the maximum likelihood principle?
 - Find parameters that **maximize probability of data**, compared to all other parameter values

Maximum Likelihood Estimation

- HMM parameters:
transition & emission probabilities
- How to find MLE parameters?
 - Suppose the data includes annotations of best hidden states for each time step

CVCVC VC V CCVVC CVC
◆ □ ◡ ☉ ⊠ † ◆ ☉ † □ ℳ ☉ ◆ ◡ ☉ ⊠

Easy: Just take relative frequencies!
(# times $C \rightarrow C$, # times $C \rightarrow \text{◡}$, etc. and normalize.)

Maximum Likelihood Estimation

- HMM parameters:
transition & emission probabilities
- How to find MLE parameters?
 - What if the data does *not* contain annotations of best hidden states (more realistic)?
 - Can we get relative counts of state transitions and emissions then?

◆ □ ♪ ☉ ▴ ♯ ◆ ☉ ♯ □ ♫ ☉ ◆ ♪ ☉ ▴

Yes!

Unsupervised Learning

- Learn structure from data without any annotations about the structure
 - Eg: learn **part-of-speech** state transitions and emissions, given only the observed **words**
- But why?
 1. Isn't this how *we* learn?
 2. Getting annotations is difficult, expensive, sometimes impossible

Sometimes, I have some unlabeled data, and I want to put labels on it.



UNSUPERVISED LEARNING, AS SHE IS IMPLEMENTED.



So I write down a generative model, and then tell the data to find parameters that explain the data to me. And if I am not satisfied with the likelihood of this explanation, I tell the data to do it again until I am.



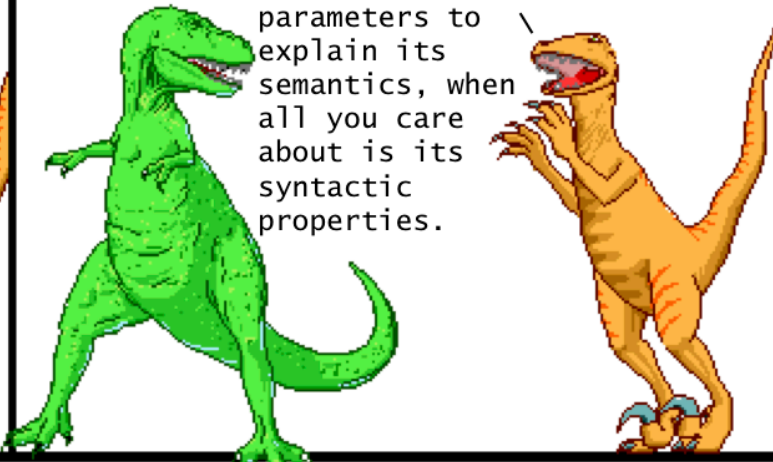
Wow, that sucks for the data.



I know, right? It's not the data's fault that I was too lazy to label it, right?



It seems like there are some deeper issues. Sometimes, most of the variation in the data comes from phenomena that are irrelevant to your desired labeling scheme. For example, the data might use its parameters to explain its semantics, when all you care about is its syntactic properties.



Why do we have to put labels on our data at all? Can't we just appreciate our data for who it is, and recognize that each datum is a unique and precious snowflake? Guys, all I'm saying, is maybe with a little supervision, our data can grow up to be whatever it wants to be!



What if...

- ... we had a corpus of observation annotated with the state sequences?
 - Can estimate the MLE parameters by counting
- ... we had the HMM parameters
 - Can find the best state sequences for the observations in the corpus with Viterbi

Expectation Maximization

1. Guess parameters of model
2. Tag observations with best state sequence
3. Re-estimate model parameters by counting
4. Repeat: go to step 2.

***Locking* parameters:**

If we like some of the original model parameters, can ignore the data and just not change those parameters.

Soft Expectation Maximization

1. Guess parameters of model
2. Tag observations with best state sequence
3. Re-estimate model parameters by counting
4. Repeat: go to step 2.

Bad guess can corner model
into bad local maximum

Instead of getting counts from the best state sequence, get **expected counts** (weighted average) from all possible state sequences with their probabilities

Building block: Backward Algorithm

- Forward algorithm computes total probability of observation going **left to right**
- Backward algorithm computes total probability of observation going **right to left**

Backward Algorithm

Given these probabilities:

$$P(C \mid \#) = 0.7$$

$$P(C \mid C) = 0.4$$

$$P(V \mid V) = 0.1$$

$$P(\Omega \mid C) = 0.08$$

$$P(\Omega \mid V) = 0.01$$

$$P(\mathcal{O} \mid C) = 0.02$$

$$P(\mathcal{O} \mid V) = 0.14$$

$$P(\boxtimes \mid C) = 0.07$$

$$P(\boxtimes \mid V) = 0.20$$

V			
C			
	Ω	\mathcal{O}	\boxtimes

Probability of starting at C at this position and generating the remainder of the word

Backward Algorithm

Given these probabilities:

$$P(C | \#) = 0.7$$

$$P(C | C) = 0.4$$

$$P(V | V) = 0.1$$

$$P(\underline{\Omega} | C) = 0.08$$

$$P(\underline{\Omega} | V) = 0.01$$

$$P(\mathcal{G} | C) = 0.02$$

$$P(\mathcal{G} | V) = 0.14$$

$$P(\boxtimes | C) = 0.07$$

$$P(\boxtimes | V) = 0.20$$

V			1
C			1
	$\underline{\Omega}$	\mathcal{G}	\boxtimes

Backward Algorithm

Given these probabilities:

$$P(C|\#)=0.7$$

$$P(C|C)=0.4$$

$$P(V|V)=0.1$$

$$P(\underline{\Omega} | C)=0.08$$

$$P(\underline{\Omega} | V)=0.01$$

$$P(\mathcal{G} | C)=0.02$$

$$P(\mathcal{G} | V)=0.14$$

$$P(\boxtimes | C)=0.07$$

$$P(\boxtimes | V)=0.20$$

V			1
C			1
	$\underline{\Omega}$	\mathcal{G}	\boxtimes

$P(C|C)*P(\boxtimes|C)*1$
 $= 0.4*0.07*1 = 0.028$
 $P(V|C)*P(\boxtimes|V)*1$
 $= 0.6*0.20*1 = 0.12$

Backward Algorithm

Given these probabilities:

$$P(C|\#)=0.7$$

$$P(C|C)=0.4$$

$$P(V|V)=0.1$$

$$P(\underline{\Omega} | C)=0.08$$

$$P(\underline{\Omega} | V)=0.01$$

$$P(\mathcal{O} | C)=0.02$$

$$P(\mathcal{O} | V)=0.14$$

$$P(\boxtimes | C)=0.07$$

$$P(\boxtimes | V)=0.20$$

V			1
			1
	$\underline{\Omega}$	\mathcal{O}	\boxtimes

$P(C|V)*P(\boxtimes | C)*1$
 $= 0.9*0.07*1 = 0.063$
 $P(V|V)*P(\boxtimes | V)*1$
 $= 0.1*0.20*1 = 0.02$

.148

Backward Algorithm

Given these probabilities:

$$P(C|\#)=0.7$$

$$P(C|C)=0.4$$

$$P(V|V)=0.1$$

$$P(\ominus|C)=0.08$$

$$P(\ominus|V)=0.01$$

$$P(\odot|C)=0.02$$

$$P(\odot|V)=0.14$$

$$P(\boxtimes|C)=0.07$$

$$P(\boxtimes|V)=0.20$$

V		0.083	1
C		0.148	1
	\odot		\boxtimes

$$P(C|C) * P(\odot|C) * 0.148$$

$$= 0.4 * 0.02 * 0.148 = 0.001184$$

$$P(V|C) * P(\odot|V) * 0.083$$

$$= 0.6 * 0.14 * 0.083 = 0.006972$$

Backward

Given these probabilities:

$$P(C|\#)=0.7$$

$$P(C|C)=0.4$$

$$P(V|V)=0.1$$

$$P(\underline{\Omega} | C)=0.08$$

$$P(\underline{\Omega} | V)=0.01$$

$$P(\textcircled{\Omega} | C)=0.02$$

$$P(\textcircled{\Omega} | V)=0.14$$

$$P(\square \wedge | C)=0.07$$

$$P(\square \wedge | V)=0.20$$

$$P(C|V) * P(\textcircled{\Omega} | C) * 0.148$$

$$= 0.9 * 0.02 * 0.148 = 0.002664$$

$$P(V|V) * P(\textcircled{\Omega} | V) * 0.083$$

$$= 0.1 * 0.14 * 0.083 = 0.001162$$

V		0.083	1
C	0.00816	0.148	1
	$\underline{\Omega}$	$\textcircled{\Omega}$	$\square \wedge$

Backward Algorithm

Total
probability
of $\underline{\Omega} \textcircled{9} \square \wedge$:

$$\begin{aligned}
 &P(C|\#) \\
 &*P(\underline{\Omega} | C) \\
 &*0.00816 \\
 &+ \\
 &P(V|\#) \\
 &*P(\textcircled{9} | V) \\
 &*0.00383 \\
 &= 0.00045696 \\
 &+ 0.00001149 \\
 &= 0.000468
 \end{aligned}$$

V	0.00383	0.083	1
C	0.00816	0.148	1
	$\underline{\Omega}$	$\textcircled{9}$	$\square \wedge$

Forward Prob = Backward Prob

V	0.003	0.004746	1.55×10^{-4}
C	0.056	0.000502	3.13×10^{-4}
#	d	a	y

Total
Probability
of $\Omega \circlearrowleft \boxtimes =$
 4.68×10^{-4}

Soft Expectation Maximization

1. Guess parameters of model
2. Tag observations with best state sequence
3. Re-estimate model parameters by counting
4. Repeat: go to step 2.

Bad guess can corner model
into bad local maximum

Instead of getting counts from the best state sequence, get **expected counts** (weighted average) from all possible state sequences with their probabilities

Forward Backward Algorithm

- How to avoid enumerating all exponentially many state sequences in the E-Step?
 - Dynamic Programming magic!

Emission Counts

- What is the total probability of *all* state sequences where state_2 is **C** given $\Omega \subseteq \Sigma \square$?
- Gives expected number of times $C \rightarrow \subseteq$ at state_2

C	C	C
C	C	V
V	C	V
V	C	C

$$\begin{aligned}
 & P(\text{day}, \text{CCC}) \\
 & + P(\text{day}, \text{CCV}) \\
 & + P(\text{day}, \text{VCV}) \\
 & + P(\text{day}, \text{VCC})
 \end{aligned}$$

Ω	\subseteq	\square
----------	-------------	-----------

Emission Counts

V	0.003	0.004746	1.55×10^{-4}
C	0.056	0.000502	3.13×10^{-4}
	$\underline{\Omega}$	\mathcal{O}	\triangleleft

V	0.00383	0.083	1
C	0.00816	0.148	1
	$\underline{\Omega}$	\mathcal{O}	\triangleleft

$$\begin{aligned}
 & P(\underline{\Omega} \mathcal{O} \triangleleft, \text{state}_2 = C) \\
 &= \text{Forward}(\text{state}_2 = C) \\
 & \quad * \text{Backward}(\text{state}_2 = C) \\
 &= 0.000502 \\
 & \quad * 0.148 \\
 &= 7.4296 \times 10^{-5}
 \end{aligned}$$

$$\begin{aligned}
 & P(\text{state}_2 = C \mid \underline{\Omega} \mathcal{O} \triangleleft) \\
 &= \frac{P(\underline{\Omega} \mathcal{O} \triangleleft, \text{state}_2 = C)}{P(\underline{\Omega} \mathcal{O} \triangleleft)} \\
 &= 7.4296 \times 10^{-5} / 4.68 \times 10^{-4} \\
 &= \mathbf{0.1588}
 \end{aligned}$$

Emission Counts

- What is the total probability of *all* state sequences where state₂ is **V** given $\Omega \subseteq \Omega \boxtimes$?
- Gives expected number of times $V \rightarrow \subseteq$ at state₂

C	V	C
C	V	V
V	V	V
V	V	C

P(day, CVC)
+ P(day, CVV)
+ P(day, VVV)
+ P(day, VVC)

Ω	\subseteq	\boxtimes
----------	-------------	-------------

Emission Counts

V	0.003	0.004746	1.55×10^{-4}
C	0.056	0.000502	3.13×10^{-4}
	$\underline{\Omega}$	\mathcal{O}	\triangleleft

V	0.00383	0.083	1
C	0.00816	0.148	1
	$\underline{\Omega}$	\mathcal{O}	\triangleleft

$$\begin{aligned}
 &P(\underline{\Omega} \mathcal{O} \triangleleft, \text{state}_2 = V) \\
 &= \text{Forward}(\text{state}_2 = V) \\
 &\quad * \text{Backward}(\text{state}_2 = V) \\
 &= 0.004746 \\
 &\quad * 0.083 \\
 &= 3.939 \times 10^{-4}
 \end{aligned}$$

$$\begin{aligned}
 &P(\text{state}_2 = V \mid \underline{\Omega} \mathcal{O} \triangleleft) \\
 &= \frac{P(\underline{\Omega} \mathcal{O} \triangleleft, \text{state}_2 = V)}{P(\underline{\Omega} \mathcal{O} \triangleleft)} \\
 &= 3.939 \times 10^{-4} / 4.68 \times 10^{-4} \\
 &= \mathbf{0.8412}
 \end{aligned}$$

Emission Counts

V	0.003	0.004746	1.55×10^{-4}
C	0.056	0.000502	3.13×10^{-4}
	$\underline{\Omega}$	$\underline{\Theta}$	$\underline{\Lambda}$

V	0.00383	0.083	1
Sanity Check #2: $P(\text{state}_i=C \text{obs}) + P(\text{state}_i=V \text{obs}) = 1$			

$$\begin{aligned}
 &P(\underline{\Omega} \underline{\Theta} \underline{\Lambda}, \text{state}_2 = V) \\
 &= \text{Forward}(\text{state}_2 = V) \\
 &\quad * \text{Backward}(\text{state}_2 = V) \\
 &= 0.004746 \\
 &\quad * 0.083 \\
 &= 3.939 \times 10^{-4}
 \end{aligned}$$

$$\begin{aligned}
 &P(\text{state}_2 = V | \underline{\Omega} \underline{\Theta} \underline{\Lambda}) \\
 &= \frac{P(\underline{\Omega} \underline{\Theta} \underline{\Lambda}, \text{state}_2 = V)}{P(\underline{\Omega} \underline{\Theta} \underline{\Lambda})} \\
 &= 3.939 \times 10^{-4} / 4.68 \times 10^{-4} \\
 &= \mathbf{0.8412}
 \end{aligned}$$

Expected Emission Counts

- To get expected number of times $V \rightarrow \mathfrak{S}$
 - For every string in corpus, at all positions i where \mathfrak{S} occurs, compute total conditional probability of V at position i
 - Add up these probabilities

How about transition counts?

- What is the total probability of all state sequences where state₁ is C and state₂ is V given $\underline{\Omega} \circlearrowleft \boxdot$?
- Gives expected number of times $C \rightarrow V$ from state₁ to state₂

C	V	C
C	V	V

P(day, CVC)
+ P(day, CVV)

$\underline{\Omega}$	\circlearrowleft	\boxdot
----------------------	--------------------	-----------

Transition Counts

V	0.003	0.004746	1.55×10^{-4}
C	0.056	0.000502	3.13×10^{-4}
	$\underline{\Omega}$	\mathcal{O}	\triangleleft

V	0.00383	0.083	1
C	0.00816	0.148	1
	$\underline{\Omega}$	\mathcal{O}	\triangleleft

$$\begin{aligned}
 &P(\underline{\Omega} \mathcal{O} \triangleleft, \\
 &\text{state}_1 = C, \text{state}_2 = V) \\
 &= \text{Forward}(\text{state}_1 = C) \\
 &* P(V|C) * P(\mathcal{O}|V) \\
 &* \text{Backward}(\text{state}_2 = V) \\
 &= 0.056 \\
 &* 0.6 * 0.14 \\
 &* 0.083 = 3.904 \times 10^{-4}
 \end{aligned}$$

$$\begin{aligned}
 &P(\text{state}_1 = C, \text{state}_2 = V | \\
 &\underline{\Omega} \mathcal{O} \triangleleft) \\
 &= P(\underline{\Omega} \mathcal{O} \triangleleft, \text{state}_1 = C, \\
 &\text{state}_2 = V) / P(\underline{\Omega} \mathcal{O} \triangleleft) \\
 &= 3.904 \times 10^{-4} / 4.68 \times 10^{-4} \\
 &= 0.8341
 \end{aligned}$$

Transition Counts

V	0.003	0.004746	1.55×10^{-4}
C	0.056	0.000502	3.13×10^{-4}
	$\underline{\Omega}$	\mathcal{O}	\triangleleft

V	0.00383	0.083	1
C	0.00816	0.148	1
	$\underline{\Omega}$	\mathcal{O}	\triangleleft

$$\begin{aligned}
 &P(\underline{\Omega} \mathcal{O} \triangleleft, \\
 &\text{state}_1 = C, \text{state}_2 = C) \\
 &= \text{Forward}(\text{state}_1 = C) \\
 &* P(C|C) * P(\mathcal{O}|C) \\
 &* \text{Backward}(\text{state}_2 = C) \\
 &= 0.056 \\
 &* 0.4 * 0.02 \\
 &* 0.148 = 6.6304 \times 10^{-5}
 \end{aligned}$$

$$\begin{aligned}
 &P(\text{state}_1 = C, \text{state}_2 = C | \\
 &\underline{\Omega} \mathcal{O} \triangleleft) \\
 &= P(\underline{\Omega} \mathcal{O} \triangleleft, \text{state}_1 = C, \\
 &\text{state}_2 = C) / P(\underline{\Omega} \mathcal{O} \triangleleft) \\
 &= 6.6304 \times 10^{-5} / 4.68 \times 10^{-4} \\
 &= 0.1417
 \end{aligned}$$

Transition Counts

V	0.003	0.004746	1.55×10^{-4}
C	0.056	0.000502	3.13×10^{-4}
	$\underline{\Omega}$	\mathcal{C}	\triangleleft

V	0.00383	0.083	1
C	0.00816	0.148	1
	$\underline{\Omega}$	\mathcal{C}	\triangleleft

$$\begin{aligned}
 &P(\underline{\Omega} \mathcal{C} \triangleleft, \\
 &\text{state}_1 = V, \text{state}_2 = C) \\
 &= \text{Forward}(\text{state}_1 = V) \\
 &* P(C|V) * P(\mathcal{C}|C) \\
 &* \text{Backward}(\text{state}_2 = C) \\
 &= 0.003 \\
 &* 0.9 * 0.02 \\
 &* 0.148 = 7.992 \times 10^{-6}
 \end{aligned}$$

$$\begin{aligned}
 &P(\text{state}_1 = V, \text{state}_2 = C | \\
 &\underline{\Omega} \mathcal{C} \triangleleft) \\
 &= P(\underline{\Omega} \mathcal{C} \triangleleft, \text{state}_1 = V, \\
 &\text{state}_2 = C) / P(\underline{\Omega} \mathcal{C} \triangleleft) \\
 &= 7.992 \times 10^{-6} / 4.68 \times 10^{-4} \\
 &= 0.0171
 \end{aligned}$$

Transition Counts

V	0.003	0.004746	1.55×10^{-4}
C	0.056	0.000502	3.13×10^{-4}
	$\underline{\Omega}$	\mathcal{O}	\triangleleft

V	0.00383	0.083	1
C	0.00816	0.148	1
	$\underline{\Omega}$	\mathcal{O}	\triangleleft

$$\begin{aligned}
 &P(\underline{\Omega} \mathcal{O} \triangleleft, \\
 &\text{state}_1 = V, \text{state}_2 = V) \\
 &= \text{Forward}(\text{state}_1 = V) \\
 &* P(V|V) * P(\mathcal{O}|V) \\
 &* \text{Backward}(\text{state}_2 = V) \\
 &= 0.003 \\
 &* 0.1 * 0.14 \\
 &* 0.083 = 3.486 \times 10^{-6}
 \end{aligned}$$

$$\begin{aligned}
 &P(\text{state}_1 = V, \text{state}_2 = V | \\
 &\underline{\Omega} \mathcal{O} \triangleleft) \\
 &= P(\underline{\Omega} \mathcal{O} \triangleleft, \text{state}_1 = V, \\
 &\text{state}_2 = V) / P(\underline{\Omega} \mathcal{O} \triangleleft) \\
 &= 3.486 \times 10^{-6} / 4.68 \times 10^{-4} \\
 &= 0.0074
 \end{aligned}$$

Transition Counts Sanity Check

- $P(\text{state}_1 = C, \text{state}_2 = C \mid \underline{\Omega} \mathcal{C} \boxtimes) = 0.1417$
- $P(\text{state}_1 = C, \text{state}_2 = V \mid \underline{\Omega} \mathcal{C} \boxtimes) = 0.8341$
- $P(\text{state}_1 = V, \text{state}_2 = V \mid \underline{\Omega} \mathcal{C} \boxtimes) = 0.0074$
- $P(\text{state}_1 = V, \text{state}_2 = C \mid \underline{\Omega} \mathcal{C} \boxtimes) = 0.0171$

Sanity Check #3:
These should add to 1

Expected Transition Counts

- To get expected number of times $C \rightarrow V$
 - For every string in corpus, at all positions i compute total conditional probability of C at position i and V at position $i+1$
 - Add up these probabilities

What are HMMs used for?

- *Any* sequence modeling problem!
 - Part of speech tagging
 - Speech recognition
 - Machine translation
 - Gene prediction in computational biology
 - Modeling human gait or facial expressions
 - Financial analysis

Techniques are even more general

- Dynamic Programming
 - Shortest path in graphs
(map route-finding, network routing)
 - Audio alignment
 - AI game-playing
 - Fast matrix multiplication

Techniques are even more general

- Expectation Maximization:
any problem to find unknown structure of data
 - Clustering
 - Computer vision
 - Signal processing
 - Financial modeling

Variants

	Discrete States	Continuous States
Discrete Observations	HMMs with Discrete Emission Probabilities (Text, Biology)	Kalman Filters, Particle Filters (Finance, GPS, Astronomy)
Continuous Observations	HMMs with Continuous Emission Probabilities (Speech, Vision)	

Dynamic Bayesian Networks

- HMMs assume each observation symbol is generated from a single state
- What if there's more to the story?
 - Word generated from part-of-speech, topic, sentiment, previous word
 - Sound generated from phoneme, voice quality, gender

Machine Learning Fundamentals

- **Inference**: answering questions about data under an **existing model**
 - The cross entropy of a text under an n-gram language model
 - The best HMM state sequence for an observation
 - Computing the perplexity of a text under a language model
 - Edit distance two strings under a set of edit costs
 - Best parse of a sentence under a grammar

Machine Learning Fundamentals

- **Learning**: computing the **parameters of a model** from a collection of data
 - Estimating an n-gram language model from a corpus
 - Estimating the channel model spelling error costs from a corpus of spelling mistakes
 - Estimating the HMM transition and emission probabilities
- Our favorite strategy: maximum likelihood