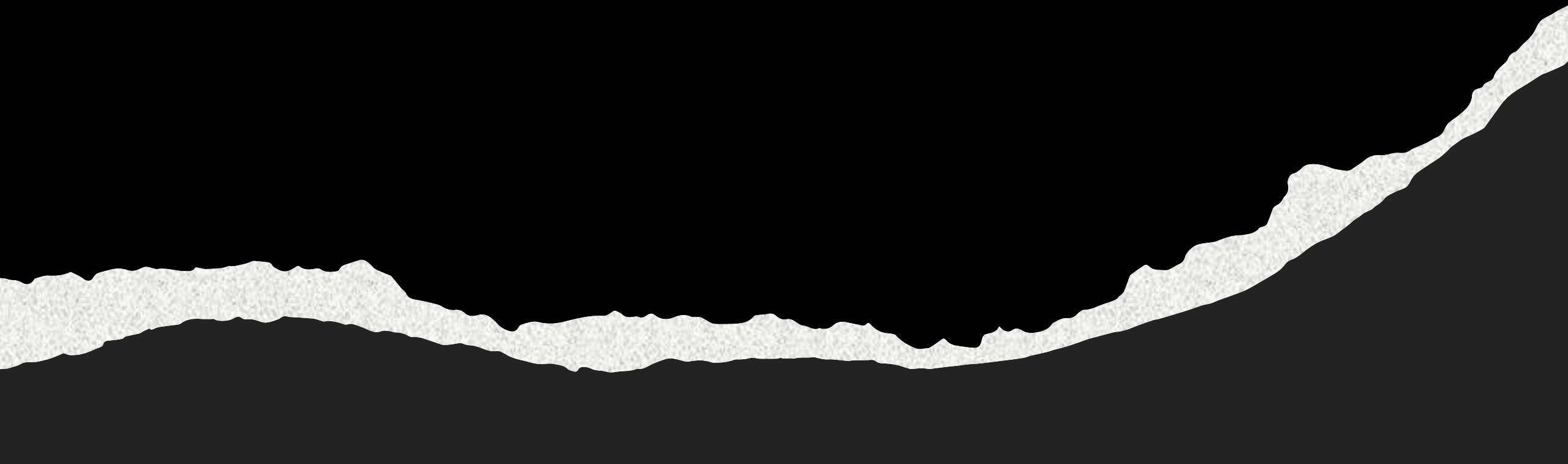


Hidden Markov Models



HMM Applications

- Malware detection
- Stock prediction
- Bioinformatics
- Speech recognition (ex, Siri)
- Cryptanalysis
- Machine translation
- Handwriting recognition
- Time series analysis

and many more...

- A **hidden Markov model (HMM)** includes a **Markov process** that is “hidden”

We cannot directly observe the state of the process

- But we do have access to **a series of observations** that are **probabilistically related to the underlying Markov model**

Probabilities

- The notation “|” denotes “*given*” information, so that $P(B|A)$ is read as “the probability of B, given A”

For any two events A and B, we have:

$$P(A \text{ and } B) = P(A)P(B|A)$$

For example, suppose that we draw two cards without replacement from a standard 52-card deck

- Let $A = \{\text{1st card is ace}\}$ and $B = \{\text{2nd card is ace}\}$. Then:

$$P(A \text{ and } B) = P(A)P(B|A) = 4/52 * 3/51 = 12/2652 = 1/221$$

- In this example, $P(B)$ depends on what happens in the first event, so we say that A and B are **dependent events**

Probabilities

On the other hand, suppose we flip a fair coin twice

- The probability that the second flip comes up heads is $\frac{1}{2}$
Regardless of the outcome of the first coin flip, so these events are independent
- For dependent events, the “*given*” information is relevant when determining the sample space.
Consequently, in such cases we can view the information to the right of the “*given*” sign as defining the space over which probabilities will be compute

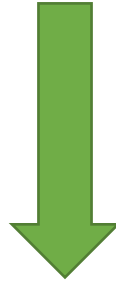
Probabilities

- Note that $P(B|A)$ would define a **First-order Markov Model**
- We can also have **Second-order Markov Model**: $P(C|A,B)$
- And **Third-order Markov Model**: $P(D|A,B,C)$

And so on...

- We rewrite:

$$P(A \text{ and } B) = P(A) P(B | A)$$



- As:

$$P(B | A) = \frac{P(A \text{ and } B)}{P(A)}$$

- This expression can be viewed as the definition of **conditional probability**

We can see it as an **intersection**

This will be useful later...

$$P(A \text{ and } B) = P(A, B) = P(A \cap B)$$

Markov Process example

We have three sentences:

1. "I like cats"
2. "I like you "
3. "I love SJSU"

And 6 states: 'I', 'like', 'cats', 'you', 'love', 'SJSU'

Then:

$\pi('I') = 1$, $\pi(w) = 0$, where $w = \{'like', 'cats', 'you', 'love', 'SJSU'\}$

so: $\pi = [1, 0, 0, 0, 0, 0]$

$$P('like' \mid 'I') = 2/3$$

$$P('cats' \mid 'like') = 1/2$$

$$P('love' \mid 'I') = 1/3$$

$$P('you' \mid 'like') = 1/2$$

HMM Example

- Suppose we want to determine the average annual temperature at a particular location on Earth over a series of years in the distant past
To simplify the problem, we only consider “hot” and “cold” for the average annual temperature (binary)

Modern evidence indicates that:

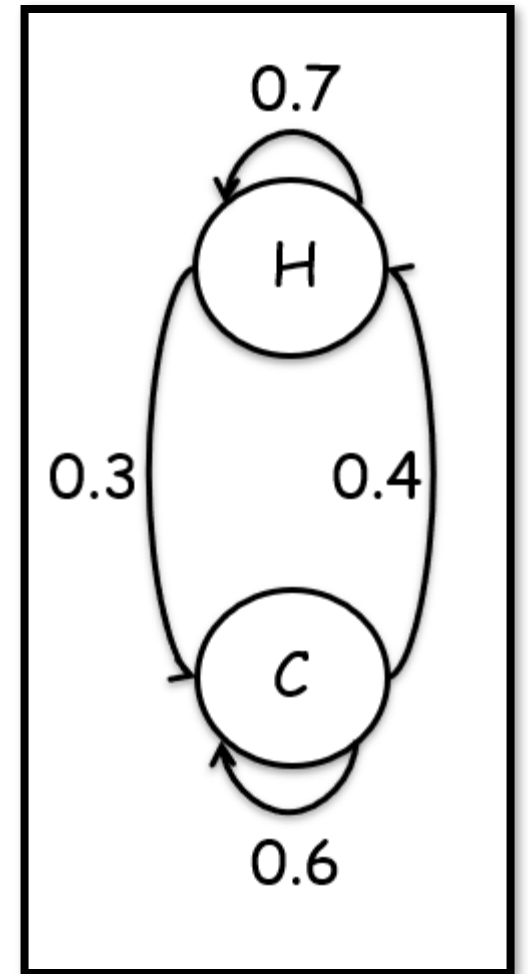
- The probability of a hot year followed by another hot year is 0.7
- The probability that a cold year is followed by another cold year is 0.6

We'll assume that these probabilities also held in the distant past

HMM Example

- This information can be summarized as:

$$\begin{array}{c} H \\ C \end{array} \begin{array}{cc} H & C \\ \left(\begin{array}{cc} 0.7 & 0.3 \\ 0.4 & 0.6 \end{array} \right) \end{array}$$



- The transition from one state to the next is a **Markov process**
Since the next state depends only on the current state and the fixed probabilities given in the matrix

HMM Example

- Next, suppose that current research indicates a correlation between the size of tree growth rings and temperature.
- For simplicity, we only consider three different tree ring sizes, *small*, *medium*, and *large*, denoted **S**, **M** and **L**



HMM Example

- Based on currently available evidence, the probabilistic relationship between annual temperature and tree ring sizes is given by:

$$\begin{array}{c} H \\ C \end{array} \begin{array}{ccc} S & M & L \\ \left(\begin{array}{ccc} 0.1 & 0.4 & 0.5 \\ 0.7 & 0.2 & 0.1 \end{array} \right) \end{array}$$

HMM Example

$$\begin{array}{c} H \\ C \end{array} \begin{array}{ccc} S & M & L \\ \left(\begin{array}{ccc} 0.1 & 0.4 & 0.5 \\ 0.7 & 0.2 & 0.1 \end{array} \right) \end{array}$$

- Although we can't observe the state (temperature) in the past, we can observe the size of tree rings
- Since the underlying states are hidden, this type of system is known as a **hidden Markov model (HMM)**

HMM Model

- **Transition matrix:**

$$A = \begin{pmatrix} 0.7 & 0.3 \\ 0.4 & 0.6 \end{pmatrix}$$

$$\begin{matrix} & H & C \\ \begin{matrix} H \\ C \end{matrix} & \begin{bmatrix} 0.7 & 0.3 \\ 0.4 & 0.6 \end{bmatrix} \end{matrix}$$

- **Observation matrix:**

$$B = \begin{pmatrix} 0.1 & 0.4 & 0.5 \\ 0.7 & 0.2 & 0.1 \end{pmatrix}$$

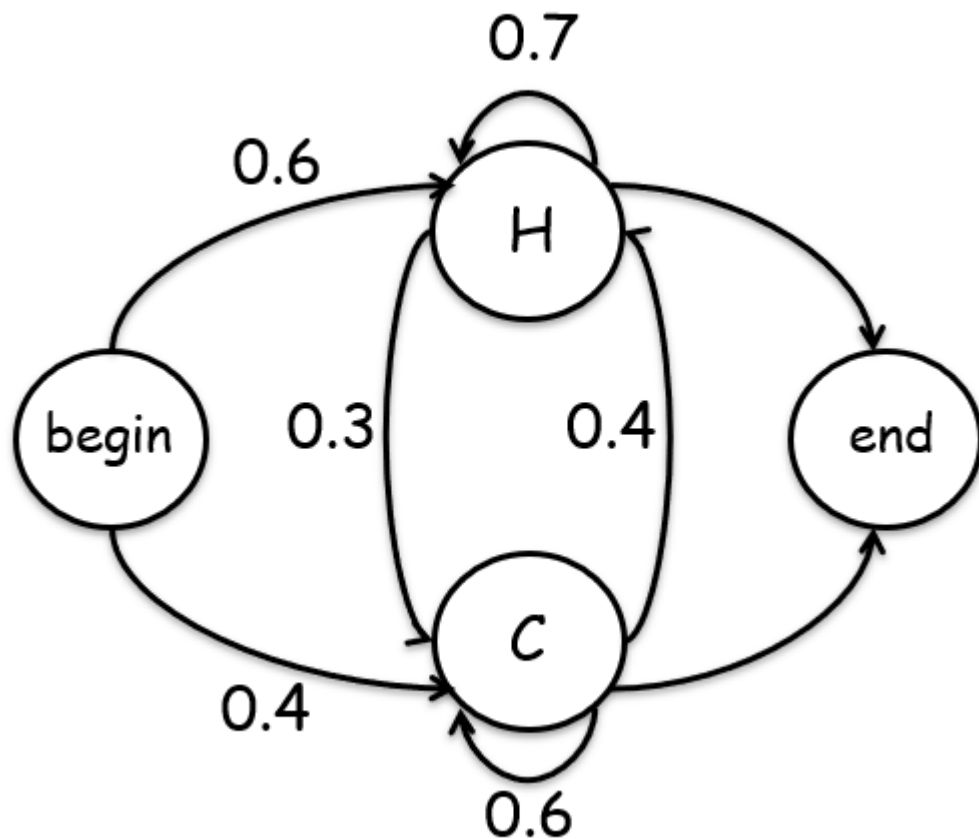
$$\begin{matrix} & S & M & L \\ \begin{matrix} H \\ C \end{matrix} & \begin{bmatrix} 0.1 & 0.4 & 0.5 \\ 0.7 & 0.2 & 0.1 \end{bmatrix} \end{matrix}$$

- **Initial state distribution:**

$$\pi = (0.6 \quad 0.4)$$

$$\pi = \begin{bmatrix} H & C \\ 0.6 & 0.4 \end{bmatrix}$$

HMM Model



$$\pi = (0.6 \quad 0.4) \qquad \pi = \begin{bmatrix} H & C \\ 0.6 & 0.4 \end{bmatrix}$$

- The initial distribution matrix indicates the chance that we start in the H state is 0.6 and the chance that we start in the C state is 0.4

HMM Example

- Now, suppose that we consider a particular four-year period of interest from the distant past
- For this particular four-year period, we observe the series of tree ring sizes **S, M, S, L**. Letting:
 - 0 represent S
 - 1 represent M
 - 2 represent L

$$\mathcal{O} = (0, 1, 0, 2)$$

HMM Example

- We might want to determine **the most likely state sequence** of the Markov process **given these observations**

That is, we might want to know the **most likely** average annual temperatures over this four-year period of interest

But what “*most likely*” means?

1. We could define “*most likely*” as the state sequence with the highest probability from among all possible state sequences of length 4

Dynamic programming (DP) can be used to efficiently solve this problem

2. Or, “*most likely*” could be the state sequence that maximizes the expected number of correct states

➤ An **HMM** can be used to find the most likely hidden state sequence in this latter sense

HMM Example

- It's important to realize that the **DP** and **HMM** solutions to this problem are **not necessarily the same**
- For example, the DP solution must, by definition, **include valid state transitions**, while this is **not the case for the HMM**

And **even if all state transitions are valid**, the **HMM** solution can still differ from the **DP solution**, as we'll illustrate in an example

The Notation

Notation	Explanation
T	Length of the observation sequence
N	Number of states in the model
M	Number of observation symbols
Q	Distinct states of the Markov process, q_0, q_1, \dots, q_{N-1}
V	Possible observations, assumed to be $0, 1, \dots, M - 1$
A	State transition probabilities
B	Observation probability matrix
π	Initial state distribution
\mathcal{O}	Observation sequence, $\mathcal{O}_0, \mathcal{O}_1, \dots, \mathcal{O}_{T-1}$

HMM Notation

Note that for simplicity, observations taken from $\mathbf{V} = \{0, 1, \dots, M-1\}$

That is:

$$\mathcal{O}_i \in V \text{ for } i = 0, 1, \dots, T-1$$

- The transition matrix $\mathbf{A} = \{a_{ij}\}$ is $\mathbf{N} \times \mathbf{N}$, where

$$a_{ij} = P(\text{state } q_j \text{ at } t+1 \mid \text{state } q_i \text{ at } t)$$

- The observation matrix $\mathbf{B} = \{b_j(k)\}$ is $\mathbf{N} \times \mathbf{M}$ where

$$b_j(k) = P(\text{observation } k \text{ at } t \mid \text{state } q_j \text{ at } t).$$

\mathbf{V} = possible observations

\mathbf{N} = # of states

\mathbf{M} = # of observation symbols

\mathbf{q} = state in the Markov process

\mathcal{O} = observation sequence

\mathbf{T} = length of \mathcal{O}

Consider our temperature example...

- What are the observations?

$V = \{0, 1, 2\}$, corresponding to S, M, L

- What are states of Markov process?

$Q = \{H, C\}$

- What are A, B, π , and T ?

A, B, π on previous slides
 T is number of tree rings measured

- What are N and M ?

$N = 2$ and $M = 3$

V = possible observations

M = # observation symbols

N = # of states

M = # of observation symbols

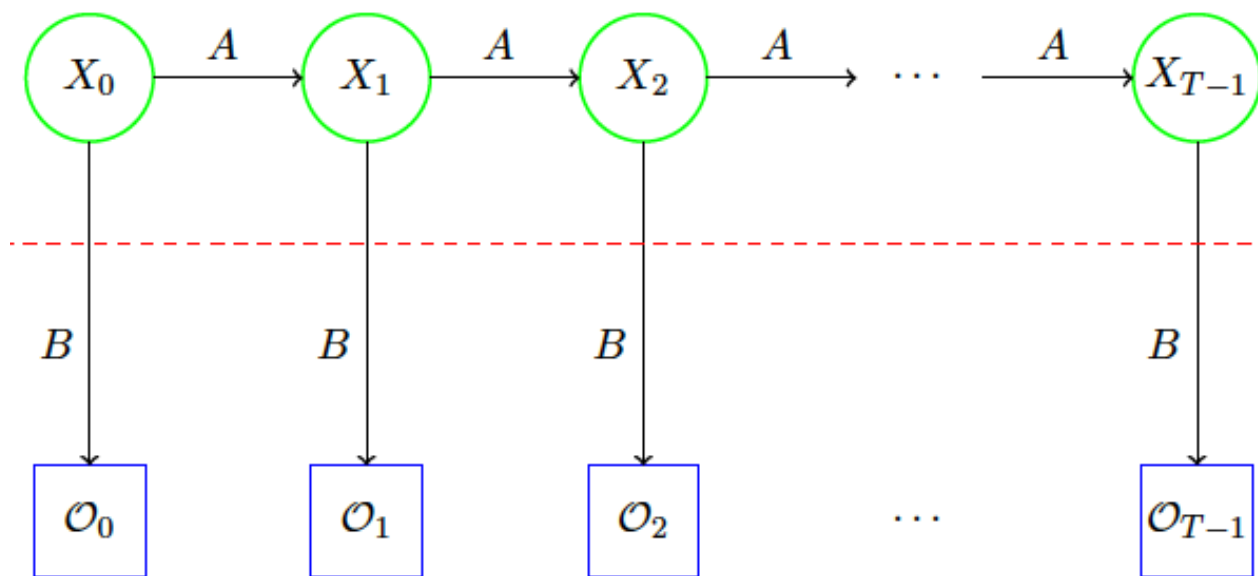
q = state in the Markov process

\mathcal{O} = observation sequence

T = length of \mathcal{O}

Generic HMM

Generic view of HMM



HMM defined by \mathbf{A} , \mathbf{B} , and $\boldsymbol{\pi}$

- We denote HMM “model” as $\boldsymbol{\lambda} = (\mathbf{A}, \mathbf{B}, \boldsymbol{\pi})$

HMM Example

- Imagine that you have a sequence of states $X = (X_0, X_1, X_2, X_3)$
- We want to compute $P(X)$

That is, how likely is to observe this exact sequence of states

$$P(X) = \pi_{x_0} b_{x_0}(\mathcal{O}_0) a_{x_0, x_1} b_{x_1}(\mathcal{O}_1) a_{x_1, x_2} b_{x_2}(\mathcal{O}_2) a_{x_2, x_3} b_{x_3}(\mathcal{O}_3)$$

Where:

- π_{x_0} is the probability of starting in state X_0
- $b_{x_0}(\mathcal{O}_0)$ is the probability of initial observation
- a_{x_0, x_1} is the probability of transition X_0 to X_1

HMM Example

$$P(X) = \pi_{x_0} b_{x_0}(\mathcal{O}_0) a_{x_0, x_1} b_{x_1}(\mathcal{O}_1) a_{x_1, x_2} b_{x_2}(\mathcal{O}_2) a_{x_2, x_3} b_{x_3}(\mathcal{O}_3)$$

Let's Assume that we observe: (0, 1, 0, 2) or (S, M, S, L)

- What is probability of: H H C C ?

$$P(HHCC) = 0.6(0.1)(0.7)(0.4)(0.3)(0.7)(0.6)(0.1) = 0.000212.$$

$$\pi = (0.6 \quad 0.4)$$

$$A = \begin{pmatrix} 0.7 & 0.3 \\ 0.4 & 0.6 \end{pmatrix}$$

$$B = \begin{pmatrix} 0.1 & 0.4 & 0.5 \\ 0.7 & 0.2 & 0.1 \end{pmatrix}$$

Initial state distribution

Transition matrix

Observation matrix

HMM Example

$$P(X) = \pi_{x_0} b_{x_0}(\mathcal{O}_0) a_{x_0, x_1} b_{x_1}(\mathcal{O}_1) a_{x_1, x_2} b_{x_2}(\mathcal{O}_2) a_{x_2, x_3} b_{x_3}(\mathcal{O}_3)$$

Let's Assume that we observe: (0, 1, 0, 2) or (S, M, S, L)

- What is probability of: H H C C ?

$$P(HHCC) = 0.6(0.1)(0.7)(0.4)(0.3)(0.7)(0.6)(0.1) = 0.000212.$$

$$\pi = \begin{bmatrix} H & C \\ 0.6 & 0.4 \end{bmatrix}$$

Initial state distribution

$$\begin{matrix} & H & C \\ H & \begin{bmatrix} 0.7 & 0.3 \end{bmatrix} \\ C & \begin{bmatrix} 0.4 & 0.6 \end{bmatrix} \end{matrix}$$

Transition matrix A

$$\begin{matrix} & S & M & L \\ H & \begin{bmatrix} 0.1 & 0.4 & 0.5 \end{bmatrix} \\ C & \begin{bmatrix} 0.7 & 0.2 & 0.1 \end{bmatrix} \end{matrix}$$

Observation matrix B

HMM Example

- We repeat it for all the $2^4 = 16$ combinations of 4-state sequences
- We find that the winner is: **CCCH**
- This is the result applying **Dynamic Programming (DP)**
Basically, just picking the maximum
- But HMM works differently...

state	probability	normalized probability
<i>HHHH</i>	.000412	.042787
<i>HHHC</i>	.000035	.003635
<i>HHCH</i>	.000706	.073320
<i>HHCC</i>	.000212	.022017
<i>HCHH</i>	.000050	.005193
<i>HCHC</i>	.000004	.000415
<i>HCCH</i>	.000302	.031364
<i>HCCC</i>	.000091	.009451
<i>CHHH</i>	.001098	.114031
<i>CHHC</i>	.000094	.009762
<i>CHCH</i>	.001882	.195451
<i>CHCC</i>	.000564	.058573
<i>CCHH</i>	.000470	.048811
<i>CCHC</i>	.000040	.004154
<i>CCCH</i>	.002822	.293073
<i>CCCC</i>	.000847	.087963

HMM Example

- To find the optimal state sequence in the **HMM** sense, we choose **the most probable symbol at each position**
- So, we **sum the normalized probabilities that have an H in the first position**
- The sum is equal to:
0.188182

state	probability	normalized probability
<i>HHHH</i>	.000412	.042787
<i>HHHC</i>	.000035	.003635
<i>HHCH</i>	.000706	.073320
<i>HHCC</i>	.000212	.022017
<i>HCHH</i>	.000050	.005193
<i>HCHC</i>	.000004	.000415
<i>HCCH</i>	.000302	.031364
<i>HCCC</i>	.000091	.009451
<i>CHHH</i>	.001098	.114031
<i>CHHC</i>	.000094	.009762
<i>CHCH</i>	.001882	.195451
<i>CHCC</i>	.000564	.058573
<i>CCHH</i>	.000470	.048811
<i>CCHC</i>	.000040	.004154
<i>CCCH</i>	.002822	.293073
<i>CCCC</i>	.000847	.087963

SUM
0.188181

HMM Example

- We repeat the same for C in the first position
- The normalized sum in this case is:
0.811818

state	probability	normalized probability	
<i>HHHH</i>	.000412	.042787	<div>SUM ↓ 0.188181</div>
<i>HHHC</i>	.000035	.003635	
<i>HHCH</i>	.000706	.073320	
<i>HHCC</i>	.000212	.022017	
<i>HCHH</i>	.000050	.005193	
<i>HCHC</i>	.000004	.000415	
<i>HCCH</i>	.000302	.031364	
<i>HCCC</i>	.000091	.009451	<div>SUM ↓ 0.811818</div>
<i>CHHH</i>	.001098	.114031	
<i>CHHC</i>	.000094	.009762	
<i>CHCH</i>	.001882	.195451	
<i>CHCC</i>	.000564	.058573	
<i>CCHH</i>	.000470	.048811	
<i>CCHC</i>	.000040	.004154	
<i>CCCH</i>	.002822	.293073	
<i>CCCC</i>	.000847	.087963	

HMM Example





- We continue, now with all the cases where H is in second position:

0.519576

- And C in second position:

0.480424

And so on...

state	probability	normalized probability	
<i>HHHH</i>	.000412	.042787	 SUM 0.519576
<i>HHHC</i>	.000035	.003635	
<i>HHCH</i>	.000706	.073320	
<i>HHCC</i>	.000212	.022017	
<i>HCHH</i>	.000050	.005193	 SUM 0.480424
<i>HCHC</i>	.000004	.000415	
<i>HCCH</i>	.000302	.031364	
<i>HCCC</i>	.000091	.009451	
<i>CHHH</i>	.001098	.114031	 SUM 0.519576
<i>CHHC</i>	.000094	.009762	
<i>CHCH</i>	.001882	.195451	
<i>CHCC</i>	.000564	.058573	
<i>CCHH</i>	.000470	.048811	 SUM 0.480424
<i>CCHC</i>	.000040	.004154	
<i>CCCH</i>	.002822	.293073	
<i>CCCC</i>	.000847	.087963	

HMM Example

- We continue with all the other positions, until we built a table like this one:

	Position in state sequence			
	0	1	2	3
$P(H)$	0.188182	0.519576	0.228788	0.804029
$P(C)$	0.811818	0.480424	0.771212	0.195971

- Finally, for each position, we select the most likely symbol. That is, our sequence (in HMM terms) would be:

CHCH

HMM Example

- **HMM** solution gives us **CHCH**
- While **DP** solution is **CCCH**

Which solution is better?

- Neither solution is better
Just using different definitions of “best”

HMM Paradox?

- **HMM** maximizes *expected number of correct states*
Whereas DP chooses “best” overall path
- Possible for **HMM** to choose a “path” that is impossible
Could be a transition probability of 0
- We cannot get impossible path with DP
Is this a flaw with HMM?
No, it’s a feature

The Three Problems



The Three Problems

HMMs are used to solve 3 problems:

- **Problem 1:** Given a model $\lambda = (\mathbf{A}, \mathbf{B}, \boldsymbol{\pi})$ and observation sequence \mathbf{O} , find $P(\mathbf{O}|\lambda)$
 - We score an observation sequence to see how well it fits a given model
- **Problem 2:** Given $\lambda = (\mathbf{A}, \mathbf{B}, \boldsymbol{\pi})$ and \mathbf{O} , find an **optimal state sequence** (in HMM sense)
 - Uncover hidden part (like previous example)
- **Problem 3:** Given \mathbf{O} , \mathbf{N} (# states), and \mathbf{M} (# of symbols), find the model λ that maximizes probability of \mathbf{O}
 - That is, **train** a model to fit observations

HMMs in Practice

Often, we use HMMs as follows:

1. We compute an observation sequence
 - We also assume that (hidden) Markov process exists
2. We train a model based on observations
 - That is, we solve Problem 3
 - “Best” N (# states) can be found by trial and error
3. Then given a (usually different, but ‘*compatible*’) sequence of observations, we score it versus the model we trained
 - This is Problem 1 — high score implies similar to training data, low score says it’s not similar

HMMs in Practice

Previous slide gives sense in which HMM is a “machine learning” technique

- To train model, we do not need to specify anything except the parameter N
- “Best” N often found by trial and error (like a “hyperparameter”)

So, we don't need to think too much

- Just train HMM and then use it
- Practical, since we have efficient algorithms for training HMM and scoring

The Three Solutions

- We give detailed solutions to 3 problems

However, we must find **efficient** solutions

The 3 problems:

- **Problem 1**: Score an observation sequence versus a given model
- **Problem 2**: Given a model, “uncover” *most likely* hidden part
- **Problem 3**: Given an observation sequence, train a model

Recall that we considered example for 2 and 1 in the previous slides

But direct solutions are soooooo inefficient

Solution 1

- Score observations versus a given model

➤ Given model $\lambda = (\mathbf{A}, \mathbf{B}, \boldsymbol{\pi})$ and observation sequence $\mathbf{O} = (\mathbf{O}_0, \mathbf{O}_1, \dots, \mathbf{O}_{T-1})$, find $P(\mathbf{O}|\lambda)$

Denote hidden states as: $X = (X_0, X_1, \dots, X_{T-1})$

And from definition of \mathbf{A} and $\boldsymbol{\pi}$: $P(X|\lambda) = \pi_{x_0} a_{x_0, x_1} a_{x_1, x_2} \dots a_{x_{T-2}, x_{T-1}}$

Then from definition of \mathbf{B} : $P(\mathbf{O}|X, \lambda) = b_{x_0}(\mathbf{O}_0) b_{x_1}(\mathbf{O}_1) \dots b_{x_{T-1}}(\mathbf{O}_{T-1})$

$$\boldsymbol{\pi} = \begin{array}{cc} & \begin{matrix} H & C \end{matrix} \\ \begin{matrix} H \\ C \end{matrix} & \begin{bmatrix} 0.6 & 0.4 \end{bmatrix} \end{array}$$

$$\mathbf{A} = \begin{array}{cc} & \begin{matrix} H & C \end{matrix} \\ \begin{matrix} H \\ C \end{matrix} & \begin{bmatrix} 0.7 & 0.3 \\ 0.4 & 0.6 \end{bmatrix} \end{array}$$

$$\mathbf{B} = \begin{array}{ccc} & \begin{matrix} S & M & L \end{matrix} \\ \begin{matrix} H \\ C \end{matrix} & \begin{bmatrix} 0.1 & 0.4 & 0.5 \\ 0.7 & 0.2 & 0.1 \end{bmatrix} \end{array}$$

Remember:

$$P(A \text{ and } B) = P(A, B) = P(A \cap B)$$

$$P(B | A) = \frac{P(A \text{ and } B)}{P(A)}$$

- Since:

$$P(\mathcal{O}, X | \lambda) = \frac{P(\mathcal{O} \cap X \cap \lambda)}{P(\lambda)}$$

- And:

$$P(\mathcal{O} | X, \lambda)P(X | \lambda) = \frac{P(\mathcal{O} \cap X \cap \lambda)}{P(X \cap \lambda)} \cdot \frac{P(X \cap \lambda)}{P(\lambda)} = \frac{P(\mathcal{O} \cap X \cap \lambda)}{P(\lambda)}$$

- We have:

$$P(\mathcal{O}, X | \lambda) = P(\mathcal{O} | X, \lambda)P(X | \lambda)$$

The problem that we want to solve
over all the possible X combinations

Definition of B times
definition of A and π

Solution 1

- The probability $P(\mathcal{O} | \lambda)$ is given summing together all possible state sequence probabilities:

$$\begin{aligned} P(\mathcal{O} | \lambda) &= \sum_X P(\mathcal{O}, X | \lambda) \\ &= \sum_X P(\mathcal{O} | X, \lambda) P(X | \lambda) \\ &= \sum_X \pi_{X_0} b_{X_0}(\mathcal{O}_0) a_{X_0, X_1} b_{X_1}(\mathcal{O}_1) \cdots a_{X_{T-2}, X_{T-1}} b_{X_{T-1}}(\mathcal{O}_{T-1}) \end{aligned}$$

Definition of B times
definition of A and π

$$P(X) = \pi_{x_0} b_{x_0}(\mathcal{O}_0) a_{x_0, x_1} b_{x_1}(\mathcal{O}_1) a_{x_1, x_2} b_{x_2}(\mathcal{O}_2) a_{x_2, x_3} b_{x_3}(\mathcal{O}_3)$$

- The direct computation here is generally infeasible

Since the **number of multiplications is about $2TN^T$**

Where T is typically large and $N \geq 2$

Forward Algorithm

- Instead, we can use ***forward algorithm***
Also known as: “**alpha pass**”

- For $\mathbf{t} = \mathbf{0}, \mathbf{1}, \dots, \mathbf{T}-\mathbf{1}$ and $\mathbf{i} = \mathbf{0}, \mathbf{1}, \dots, \mathbf{N}-\mathbf{1}$, let:

$$\alpha_t(i) = P(\mathcal{O}_0, \mathcal{O}_1, \dots, \mathcal{O}_t, X_t = q_i \mid \lambda)$$

Instead of $P(\mathcal{O}, X \mid \lambda)$,
that considers all the
states, we now stop at
the states up to time t

- This is **probability of partial observation sequence up to time t** , where the underlying Markov process is in state \mathbf{q}_i at time t
Can be computed recursively and efficiently

Forward Algorithm

1: **Given:**

Model $\lambda = (A, B, \pi)$

Observations $\mathcal{O} = (O_0, \mathcal{O}_1, \dots, \mathcal{O}_{T-1})$

2: **for** $i = 0, 1, \dots, N - 1$ **do**

3: $\alpha_0(i) = \pi_i b_i(\mathcal{O}_0)$

4: **end for**

5: **for** $t = 1, 2, \dots, T - 1$ **do**

6: **for** $i = 0, 1, \dots, N - 1$ **do**

7: $\alpha_t(i) = \left(\sum_{j=0}^{N-1} \alpha_{t-1}(j) a_{ji} \right) b_i(\mathcal{O}_t)$

8: **end for**

9: **end for**

	H	C
H	0.7	0.3
C	0.4	0.6

	S	M	L
H	0.1	0.4	0.5
C	0.7	0.2	0.1

	H	C
π	0.6	0.4

This requires only N^2T multiplications

Forward Algorithm

1: **Given:**

Model $\lambda = (A, B, \pi)$

Observations $\mathcal{O} = (O_0, \mathcal{O}_1, \dots, \mathcal{O}_{T-1})$

2: **for** $i = 0, 1, \dots, N - 1$ **do**

3: $\alpha_0(i) = \pi_i b_i(\mathcal{O}_0)$

4: **end for**

5: **for** $t = 1, 2, \dots, T - 1$ **do**

6: **for** $i = 0, 1, \dots, N - 1$ **do**

7: $\alpha_t(i) = \left(\sum_{j=0}^{N-1} \alpha_{t-1}(j) a_{ji} \right) b_i(\mathcal{O}_t)$

8: **end for**

9: **end for**

This is also called:
Initialization Step

This instead is called:
Induction Step

	H	C
H	0.7	0.3
C	0.4	0.6

	S	M	L
H	0.1	0.4	0.5
C	0.7	0.2	0.1

	H	C
$\pi =$	0.6	0.4

This requires only N^2T multiplications

Forward Algorithm

- From the formula:

$$\alpha_t(i) = P(\mathcal{O}_0, \mathcal{O}_1, \dots, \mathcal{O}_t, X_t = q_i \mid \lambda)$$

- It follows that:

$$P(\mathcal{O} \mid \lambda) = \sum_{i=0}^{N-1} \alpha_{T-1}(i)$$

Hence, the forward algorithm gives us an efficient way to compute a score for a given sequence \mathcal{O} , relative to a given model λ

Solution 2

- Given a model, find hidden states

Given $\lambda = (\mathbf{A}, \mathbf{B}, \boldsymbol{\pi})$ and \mathbf{O} , find an optimal state sequence

- Recall that optimal means “maximize expected number of correct states”
- In contrast, DP finds best scoring path

- For temp/tree ring example, we solved this

But hopelessly inefficient approach

- A better way: ***backward algorithm***

Also known as: “**beta pass**”

Backward Algorithm

- For $\mathbf{t} = \mathbf{0}, \mathbf{1}, \dots, \mathbf{T}-\mathbf{1}$ and $\mathbf{i} = \mathbf{0}, \mathbf{1}, \dots, \mathbf{N}-\mathbf{1}$, let:

$$\beta_t(i) = P(\mathcal{O}_{t+1}, \mathcal{O}_{t+2}, \dots, \mathcal{O}_{T-1} \mid X_t = q_i, \lambda)$$

- This is probability of partial observation sequence starting from time \mathbf{t} , where the underlying Markov process is in state \mathbf{q}_i at time \mathbf{t}
- Analogous to the forward algorithm
As with forward algorithm, this can be computed recursively and efficiently

Backward Algorithm

1: **Given:**

Model $\lambda = (A, B, \pi)$

Observations $\mathcal{O} = (O_0, \mathcal{O}_1, \dots, \mathcal{O}_{T-1})$

2: **for** $i = 0, 1, \dots, N - 1$ **do**

3: $\beta_{T-1}(i) = 1$

4: **end for**

5: **for** $t = T - 2, T - 3, \dots, 0$ **do**

6: **for** $i = 0, 1, \dots, N - 1$ **do**

7: $\beta_t(i) = \sum_{j=0}^{N-1} a_{ij} b_j(\mathcal{O}_{t+1}) \beta_{t+1}(j)$

8: **end for**

9: **end for**

	H	C
H	0.7	0.3
C	0.4	0.6

	S	M	L
H	0.1	0.4	0.5
C	0.7	0.2	0.1

	H	C
$\pi =$	0.6	0.4

Note: it goes backward now

Solution 2

- Now, for $\mathbf{t} = \mathbf{0}, \mathbf{1}, \dots, \mathbf{T}-\mathbf{1}$ and $\mathbf{i} = \mathbf{0}, \mathbf{1}, \dots, \mathbf{N}-\mathbf{1}$ define:

$$\gamma_t(i) = P(X_t = q_i \mid \mathcal{O}, \lambda)$$

Most likely state at \mathbf{t} is \mathbf{q}_i that maximizes $\gamma_t(\mathbf{i})$

- Since $\alpha_t(i)$ measures the relevant probability up to time t and $\beta_t(i)$ measures the relevant probability after time t , we have:

$$\gamma_t(i) = \frac{\alpha_t(i)\beta_t(i)}{P(\mathcal{O} \mid \lambda)}$$

And recall:

$$P(\mathcal{O} \mid \lambda) = \sum_{i=0}^{N-1} \alpha_{T-1}(i)$$

Solution 2

$$\gamma_t(i) = \frac{\alpha_t(i)\beta_t(i)}{P(\mathcal{O} | \lambda)}$$

Why is it necessary to normalize gamma by dividing by $P(\mathcal{O} | \lambda)$?

- Because these probabilities are computed assuming the observation sequence is known (given \mathcal{O}), as opposed to being computed relative to the larger probability space

Solution 2

- From the definition of $\gamma_t(i)$ it follows that **the most likely state at time t is the state for which $\gamma_t(i)$ is maximum**, where the maximum is taken over the index i

Then **the most likely state at time t** is given by:

$$\tilde{X}_t = \max_i \gamma_t(i)$$

- The bottom line?
 - Forward algorithm solves Problem 1
 - Forward/backward algorithms solve Problem 2

Solution 3

- Train a model: Given \mathbf{O} , \mathbf{N} , and \mathbf{M} , find $\boldsymbol{\lambda}$ that maximizes probability of \mathbf{O}
- We'll iteratively adjust $\boldsymbol{\lambda} = (\mathbf{A}, \mathbf{B}, \boldsymbol{\pi})$ to better fit the given observations \mathbf{O}
 - The size of matrices are fixed (N and M)
 - But elements of matrices can change
- It is nice that this works...
...and amazing that it's efficient!

Solution 3

- For $\mathbf{t=0,1,...,T-2}$ and $\mathbf{i,j}$ in $\{\mathbf{0,1,...,N-1}\}$, define “**di-gammas**” as:

$$\gamma_t(i, j) = P(X_t = q_i, X_{t+1} = q_j \mid \mathcal{O}, \lambda)$$

➤ Note: $\gamma_t(i, j)$ is probability of being in state $\mathbf{q_i}$ at time \mathbf{t} and transiting to state $\mathbf{q_j}$ at $\mathbf{t+1}$

Then we can rewr

$$\gamma_t(i, j) = \frac{\alpha_t(i) a_{ij} b_j(\mathcal{O}_{t+1}) \beta_{t+1}(j)}{P(\mathcal{O} \mid \lambda)}$$

And:

$$\gamma_t(i) = \sum_{j=0}^{N-1} \gamma_t(i, j)$$

Solution 3

1: **Given:**

$\gamma_t(i)$, for $t = 0, 1, \dots, T - 1$ and $i = 0, 1, \dots, N - 1$

$\gamma_t(i, j)$, for $t = 0, 1, \dots, T - 2$ and $i, j \in \{0, 1, \dots, N - 1\}$

2: **for** $i = 0, 1, \dots, N - 1$ **do**

3: $\pi_i = \gamma_0(i)$

4: **end for**

5: **for** $i = 0, 1, \dots, N - 1$ **do**

6: **for** $j = 0, 1, \dots, N - 1$ **do**

7: $a_{ij} = \sum_{t=0}^{T-2} \gamma_t(i, j) / \sum_{t=0}^{T-2} \gamma_t(i)$

8: **end for**

9: **end for**

10: **for** $j = 0, 1, \dots, N - 1$ **do**

11: **for** $k = 0, 1, \dots, M - 1$ **do**

12: $b_j(k) = \sum_{\substack{t \in \{0, 1, \dots, T-1\} \\ O_t = k}} \gamma_t(j) / \sum_{t=0}^{T-1} \gamma_t(j)$

13: **end for**

14: **end for**

$$\pi = \begin{bmatrix} H & C \\ & \end{bmatrix}$$

$$\begin{matrix} & H & C \\ H & & \\ C & & \end{matrix} \begin{bmatrix} \\ \\ \end{bmatrix}$$

$$\begin{matrix} & S & M & L \\ H & & & \\ C & & & \end{matrix} \begin{bmatrix} \\ \\ \end{bmatrix}$$

These two sums are very similar, but in the nominator we consider **only t for which $O_t = k$** and only these values are counted in the numerator

That is, per each individual symbol $[0, M-1]$ that appears in O

Solution 3

1: Given:

$\gamma_t(i)$, for $t = 0, 1, \dots, T - 1$ and $i = 0, 1, \dots, N - 1$

$\gamma_t(i, j)$, for $t = 0, 1, \dots, T - 2$ and $i, j \in \{0, 1, \dots, N - 1\}$

2: for $i = 0, 1, \dots, N - 1$ do

3: $\pi_i = \gamma_0(i)$

4: end for

5: for $i = 0, 1, \dots, N - 1$ do

6: for $j = 0, 1, \dots, N - 1$ do

7: $a_{ij} = \sum_{t=0}^{T-2} \gamma_t(i, j) / \sum_{t=0}^{T-2} \gamma_t(i)$

8: end for

9: end for

10: for $j = 0, 1, \dots, N - 1$ do

11: for $k = 0, 1, \dots, M - 1$ do

12: $b_j(k) = \sum_{\substack{t \in \{0, 1, \dots, T-1\} \\ O_t = k}} \gamma_t(j) / \sum_{t=0}^{T-1} \gamma_t(j)$

13: end for

14: end for

- The numerator of the re-estimated $b_j(k)$ is the expected number of times the model is in state q_j with observation k
- While the denominator is the expected number of times the model is in state q_j
- Therefore, the ratio is the probability of observing symbol k , given that the model is in state q_j , and this is the desired value for $b_j(k)$

Model Re-estimation

Re-estimation is an iterative process

- First, we **initialize** $\lambda = (\mathbf{A}, \mathbf{B}, \pi)$ with a reasonable guess, or, if no reasonable guess is available, we choose **random values** such that

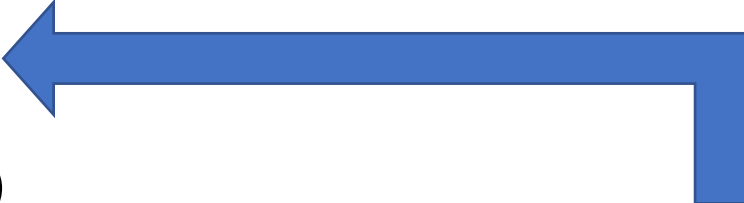
- $\pi \approx \frac{1}{N}$
- $a_{ij} \approx \frac{1}{N}$
- $b_j(k) \approx \frac{1}{M}$

It's critical that A, B and π are randomized

- Since **exactly uniform values** will result in a **local maximum** from which the **model cannot climb**

Solution 3

To summarize:

1. Initialize $\lambda = (\mathbf{A}, \mathbf{B}, \pi)$
2. Compute $\alpha_t(i)$, $\beta_t(i)$, $\gamma_t(i, j)$, $\gamma_t(i)$ 
3. Re-estimate the model $\lambda = (\mathbf{A}, \mathbf{B}, \pi)$
4. If $P(\mathcal{O} | \lambda)$ increases by more than ϵ (where ϵ is small), goto 2

Solution 3

Model initialization

- If we have a good guess for $\lambda = (\mathbf{A}, \mathbf{B}, \boldsymbol{\pi})$ then we can use it for initialization
- If not, let: $\pi_i \approx \frac{1}{N}$ $a_{ij} \approx \frac{1}{N}$ $b_j(k) \approx \frac{1}{M}$
- Subject to row stochastic conditions
- But do **not** initialize to exactly uniform values

Stopping conditions

- Stop after some number of iterations and/or...
- Stop if increase in $P(\mathcal{O} | \lambda)$ is too small (less than ϵ)

HMM as Discrete Hill Climb

- Algorithm on previous slides shows that **HMM** is a “**discrete hill climb**”
- HMM consists of **discrete states** \mathbf{X}_t
“Climb” on the elements of the matrices
- And re-estimation process improves model by modifying parameters
So, “climbs” toward improved model
This happens in a high-dimensional space

Dynamic Programming

Brief detour to show the close relationship between dynamic programming (DP) and HMMs

- For $\lambda = (\mathbf{A}, \mathbf{B}, \boldsymbol{\pi})$ it's easy to define a dynamic program
DP is like the forward algorithm (alpha-pass) but with “sum” replaced by “max”

Dynamic Programming

- Here we see the **dynamic programming algorithm**, which is also known as the **Viterbi algorithm**

1: **Given:**

Model $\lambda = (A, B, \pi)$

Observations $\mathcal{O} = (O_0, O_1, \dots, O_{T-1})$

2: **for** $i = 0, 1, \dots, N - 1$ **do**

3: $\delta_0(i) = \pi_i b_i(O_0)$

4: **end for**

5: **for** $t = 1, 2, \dots, T - 1$ **do**

6: **for** $i = 0, 1, \dots, N - 1$ **do**

7: $\delta_t(i) = \max_{j \in \{0, 1, \dots, N-1\}} (\delta_{t-1}(j) a_{ji} b_i(O_t))$

8: **end for**

9: **end for**

- Note that at each t , the DP **computes best path for each state**, up to that point
So, **probability of best path** is:

$$\max_{j \in \{0, 1, \dots, N-1\}} \delta_{T-1}(j)$$

- This max gives the highest probability, but not the best path

Dynamic Programming

$$\max_{j \in \{0,1,\dots,N-1\}} \delta_{T-1}(j)$$

- It is important to realize that this formula **only gives the optimal probability**, **not the corresponding path**
- To determine optimal path:
 1. While computing deltas, keep track of pointers to previous state
 2. When finished, construct optimal path by tracing back points

Dynamic Programming

- For example, let's consider again the temperature example

Recall that we observe $(0, 1, 0, 2)$

- The initial probabilities (path of length 1) are:

$$P(H) = \pi_0 b_0(0) = 0.6(0.1) = 0.06$$

$$P(C) = \pi_1 b_1(0) = 0.4(0.7) = 0.28$$

- The probabilities of the paths of length two are given by:

$$P(HH) = 0.06(0.7)(0.4) = 0.0168$$

$$\begin{array}{c} H \\ C \end{array} \begin{array}{cc} H & C \end{array} \begin{bmatrix} 0.7 & 0.3 \\ 0.4 & 0.6 \end{bmatrix}$$

$$\begin{array}{c} H \\ C \end{array} \begin{array}{ccc} S & M & L \end{array} \begin{bmatrix} 0.1 & 0.4 & 0.5 \\ 0.7 & 0.2 & 0.1 \end{bmatrix}$$

$$\pi = \begin{array}{cc} H & C \end{array} \begin{bmatrix} 0.6 & 0.4 \end{bmatrix}$$

Dynamic Programming

- Probabilities for each path of length 2

$$P(HH) = 0.06(0.7)(0.4) = 0.0168$$

$$P(HC) = 0.06(0.3)(0.2) = 0.0036$$

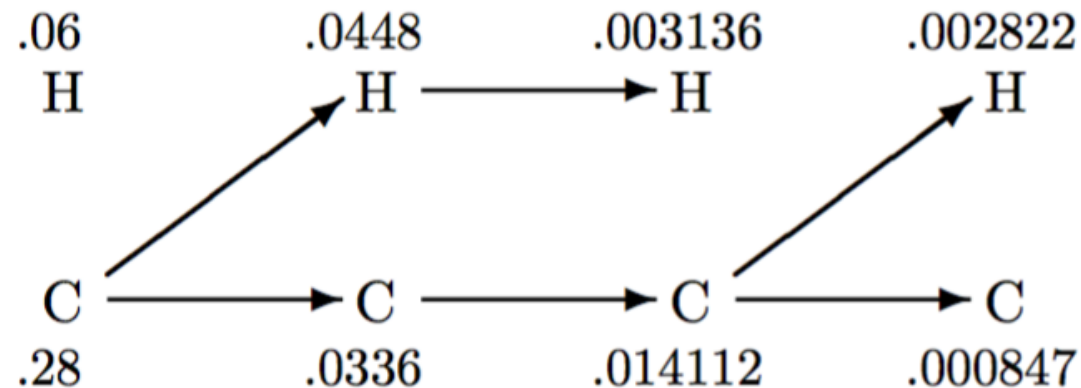
$$P(CH) = 0.28(0.4)(0.4) = 0.0448$$

$$P(CC) = 0.28(0.6)(0.2) = 0.0336$$

- Best path of length 2 ending with H is CH
- Best path of length 2 ending with C is CC

Dynamic Program

1. Continuing, we compute best path ending at H and C at each step
2. And then we save pointers
 - Note that at each stage, the dynamic programming algorithm only needs to maintain the highest-scoring path ending at each state—not a list of all possible paths (this is the key to the efficiency of the algorithm)



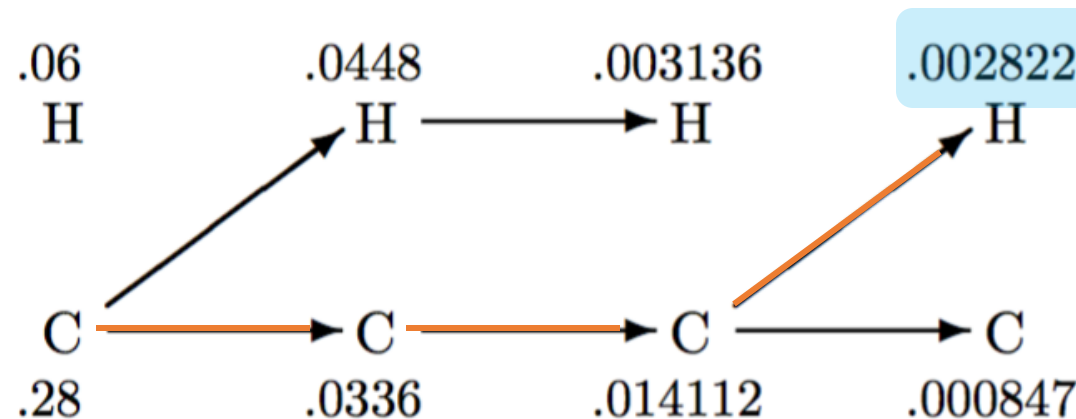
Dynamic Program

- Best final score is .002822

And thanks to pointers, best path is **CCCH**

- But what about **underflow**?

A serious problem in bigger cases



state	probability
<i>HHHH</i>	.000412
<i>HHHC</i>	.000035
<i>HHCH</i>	.000706
<i>HHCC</i>	.000212
<i>HCHH</i>	.000050
<i>HCHC</i>	.000004
<i>HCCH</i>	.000302
<i>HCCC</i>	.000091
<i>CHHH</i>	.001098
<i>CHHC</i>	.000094
<i>CHCH</i>	.001882
<i>CHCC</i>	.000564
<i>CCHH</i>	.000470
<i>CCHC</i>	.000040
<i>CCCH</i>	.002822
<i>CCCC</i>	.000847

Underflow Resistant DP

- Common trick to prevent underflow:

Instead of multiplying probabilities...

- We add logarithms of probabilities

- Why does this work?

Because $\log(xy) = \log x + \log y$

➤ Adding logs does not tend to 0

- Note that these logs are negative... and we must avoid 0 probabilities

Underflow Resistant DP

1: **Given:**

Model $\lambda = (A, B, \pi)$

Observations $\mathcal{O} = (O_0, \mathcal{O}_1, \dots, \mathcal{O}_{T-1})$

2: **for** $i = 0, 1, \dots, N - 1$ **do**

3: $\hat{\delta}_0(i) = \log(\pi_i b_i(\mathcal{O}_0))$

4: **end for**

5: **for** $t = 1, 2, \dots, T - 1$ **do**

6: **for** $i = 0, 1, \dots, N - 1$ **do**

7: $\hat{\delta}_t(i) = \max_{j \in \{0, 1, \dots, N-1\}} \left(\hat{\delta}_{t-1}(j) + \log(a_{ji}) + \log(b_i(\mathcal{O}_t)) \right)$

8: **end for**

9: **end for**

- Again, the optimal score is given by: $\max_{j \in \{0, 1, \dots, N-1\}} \hat{\delta}_{T-1}(j)$

HMM Scaling

- The three HMM solutions all involve products of probabilities
It's very easy to see, for example, that $\alpha_t(i)$ tends to 0 exponentially as T increases

```
for  $t = 1, 2, \dots, T - 1$  do
  for  $i = 0, 1, \dots, N - 1$  do
    
$$\alpha_t(i) = \left( \sum_{j=0}^{N-1} \alpha_{t-1}(j) a_{ji} \right) b_i(\mathcal{O}_t)$$

  end for
end for
```

- Therefore, any attempt to implement the HMM algorithms will inevitably result in **underflow**

HMM Scaling

- The solution to this underflow problem is to scale the numbers
However, care must be taken to ensure that the algorithms remain valid

First, consider the computation of $\alpha_t(i)$

- The basic recurrence is:

$$\alpha_t(i) = \sum_{j=0}^{N-1} \alpha_{t-1}(j) a_{ji} b_i(\mathcal{O}_t)$$

- It seems sensible to normalize each $\alpha_t(i)$ by dividing by: $\sum_{j=0}^{N-1} \alpha_t(j)$

HMM Scaling

- Following this approach, we compute **scaling factors** c_t and the **scaled** $\alpha_t(i)$, which we denote as $\hat{\alpha}_t(i)$
- To verify, we first note that $\hat{\alpha}_0(i) = c_0 \alpha_0(i)$

Next slide...

- Now suppose that for some t we have:

$$\hat{\alpha}_t(i) = c_0 c_1 \cdots c_t \alpha_t(i)$$

- Then:

$$\begin{aligned} \hat{\alpha}_{t+1}(i) &= c_{t+1} \tilde{\alpha}_{t+1}(i) \\ &= c_{t+1} \sum_{j=0}^{N-1} \hat{\alpha}_t(j) a_{ji} b_i(\mathcal{O}_{t+1}) \\ &= c_0 c_1 \cdots c_t c_{t+1} \sum_{j=0}^{N-1} \alpha_t(j) a_{ji} b_i(\mathcal{O}_{t+1}) \\ &= c_0 c_1 \cdots c_{t+1} \alpha_{t+1}(i) \end{aligned}$$

- And hence $\hat{\alpha}_t(i) = c_0 c_1 \cdots c_t \alpha_t(i)$ holds, **by induction**, for all t


```

1: Given:
     $\alpha_t(i)$ , for  $t = 0, 1, \dots, T - 1$  and  $i = 0, 1, \dots, N - 1$ 
2: for  $i = 0, 1, \dots, N - 1$  do
3:      $\tilde{\alpha}_0(i) = \alpha_0(i)$ 
4: end for
5:  $c_0 = 1 / \sum_{j=0}^{N-1} \tilde{\alpha}_0(j)$ 
6: for  $i = 0, 1, \dots, N - 1$  do
7:      $\hat{\alpha}_0(i) = c_0 \tilde{\alpha}_0(i)$ 
8: end for
9: for  $t = 1, 2, \dots, T - 1$  do
10:    for  $i = 0, 1, \dots, N - 1$  do
11:         $\tilde{\alpha}_t(i) = \sum_{j=0}^{N-1} \hat{\alpha}_{t-1}(j) a_{ji} b_i(\mathcal{O}_t)$ 
12:    end for
13:     $c_t = 1 / \sum_{j=0}^{N-1} \tilde{\alpha}_t(j)$ 
14:    for  $i = 0, 1, \dots, N - 1$  do
15:         $\hat{\alpha}_t(i) = c_t \tilde{\alpha}_t(i)$ 
16:    end for
17: end for

```

```

1: Given:
    $\alpha_t(i)$ , for  $t = 0, 1, \dots, T - 1$  and  $i = 0, 1, \dots, N - 1$ 
2: for  $i = 0, 1, \dots, N - 1$  do
3:    $\tilde{\alpha}_0(i) = \alpha_0(i)$ 
4: end for
5:  $c_0 = 1 / \sum_{j=0}^{N-1} \tilde{\alpha}_0(j)$ 
6: for  $i = 0, 1, \dots, N - 1$  do
7:    $\hat{\alpha}_0(i) = c_0 \tilde{\alpha}_0(i)$ 
8: end for
9: for  $t = 1, 2, \dots, T - 1$  do
10:  for  $i = 0, 1, \dots, N - 1$  do
11:     $\tilde{\alpha}_t(i) = \sum_{j=0}^{N-1} \hat{\alpha}_{t-1}(j) a_{ji} b_i(\mathcal{O}_t)$ 
12:  end for
13:   $c_t = 1 / \sum_{j=0}^{N-1} \tilde{\alpha}_t(j)$ 
14:  for  $i = 0, 1, \dots, N - 1$  do
15:     $\hat{\alpha}_t(i) = c_t \tilde{\alpha}_t(i)$ 
16:  end for
17: end for

```

Note that (obviously):

$$\sum_{j=0}^{N-1} \hat{\alpha}_{T-1}(j) = 1$$

Also, remembering that:

$$\hat{\alpha}_t(i) = c_0 c_1 \cdots c_t \alpha_t(i).$$

We have:

$$\begin{aligned}
\sum_{j=0}^{N-1} \hat{\alpha}_{T-1}(j) &= c_0 c_1 \cdots c_{T-1} \sum_{j=0}^{N-1} \alpha_{T-1}(j) \\
&= c_0 c_1 \cdots c_{T-1} P(\mathcal{O} \mid \lambda)
\end{aligned}$$

$$\sum_{j=0}^{N-1} \hat{\alpha}_{T-1}(j) = 1$$

$$\begin{aligned} \sum_{j=0}^{N-1} \hat{\alpha}_{T-1}(j) &= c_0 c_1 \cdots c_{T-1} \sum_{j=0}^{N-1} \alpha_{T-1}(j) \\ &= c_0 c_1 \cdots c_{T-1} P(\mathcal{O} \mid \lambda) \end{aligned}$$

Combining these results, we obtain:

$$P(\mathcal{O} \mid \lambda) = 1 / \prod_{j=0}^{T-1} c_j$$

It follows that we can compute the log of $P(\mathcal{O} \mid \lambda)$ directly from the scaling factors c_t as:

Just “logging” the two sides of the equation

$$\log(P(\mathcal{O} \mid \lambda)) = - \sum_{j=0}^{T-1} \log c_j$$

$$\sum_{j=0}^{N-1} \hat{\alpha}_{T-1}(j) = 1$$

$$\sum_{j=0}^{N-1} \hat{\alpha}_{T-1}(j) = c_0 c_1 \cdots c_{T-1} \sum_{j=0}^{N-1} \alpha_{T-1}(j)$$

$$\Rightarrow c_0 c_1 \cdots c_{T-1} P(\mathcal{O} | \lambda)$$

Combining these results, we obtain:

$$P(\mathcal{O} | \lambda) = \frac{1}{\prod_{j=0}^{T-1} c_j}$$

It follows that we can compute the log of $P(\mathcal{O} | \lambda)$ directly from the scaling factors c_t as:

Just “logging” the two sides of the equation

$$\log(P(\mathcal{O} | \lambda)) = - \sum_{j=0}^{T-1} \log c_j$$

HMM Scaling

- Similarly, scale betas as $ct\beta_t(i)$
- For re-estimation:
Compute $\gamma_t(i, j)$ and $\gamma_t(i)$ using original formulas, but with scaled alphas, betas
- This gives us new values for $\lambda = (A, B, \pi)$
“Easy exercise” to show re-estimate is exact when scaled alphas and betas used
- Also, $P(O|\lambda)$ cancels from formula
Use $\log P(O|\lambda) = -\sum \log(c_j)$ to decide if iterate improves

All Together Now

Here we see the complete pseudo code for Solution 3

Given: $(O_0, O_1, \dots, O_{T-1})$ and N and M

Initialize: $\lambda = (A, B, \pi)$

- A is $N \times N$, B is $N \times M$ and π is $1 \times N$
- $\pi_i \approx 1/N$, $a_{ij} \approx 1/N$, $b_j(k) \approx 1/M$, each matrix row stochastic, but not uniform

Initialize:

- `maxIters` = max number of re-estimation steps
- `iters` = 0
- `oldLogProb` = $-\infty$

Solution to problem 1



Solution to problem 2






Solution to problem 3



Forward Algorithm

- Forward algorithm
With scaling

Solution to problem 1 
Solution to problem 2 
Solution to problem 3 

```
// compute  $\alpha_0(i)$   
 $c_0 = 0$   
for  $i = 0$  to  $N - 1$   
     $\alpha_0(i) = \pi(i)b_i(\mathcal{O}_0)$   
     $c_0 = c_0 + \alpha_0(i)$   
next  $i$ 
```


```
// scale the  $\alpha_0(i)$   
 $c_0 = 1/c_0$   
for  $i = 0$  to  $N - 1$   
     $\alpha_0(i) = c_0\alpha_0(i)$   
next  $i$ 
```

```
// compute  $\alpha_t(i)$   
for  $t = 1$  to  $T - 1$   
     $c_t = 0$   
    for  $i = 0$  to  $N - 1$   
         $\alpha_t(i) = 0$   
        for  $j = 0$  to  $N - 1$   
             $\alpha_t(i) = \alpha_t(i) + \alpha_{t-1}(j)a_{ji}$   
        next  $j$   
         $\alpha_t(i) = \alpha_t(i)b_i(\mathcal{O}_t)$   
         $c_t = c_t + \alpha_t(i)$   
    next  $i$ 
```

```
// scale  $\alpha_t(i)$   
 $c_t = 1/c_t$   
for  $i = 0$  to  $N - 1$   
     $\alpha_t(i) = c_t\alpha_t(i)$   
next  $i$   
next  $t$ 
```




Backward Algorithm

- Backward algorithm or “beta pass”
With scaling
- Note: same scaling factor as alphas



```
// Let  $\beta_{T-1}(i) = 1$ , scaled by  $c_{T-1}$ 
for  $i = 0$  to  $N - 1$ 
     $\beta_{T-1}(i) = c_{T-1}$ 
next  $i$ 




//  $\beta$ -pass
for  $t = T - 2$  to  $0$  by  $-1$ 
    for  $i = 0$  to  $N - 1$ 
         $\beta_t(i) = 0$ 
        for  $j = 0$  to  $N - 1$ 
             $\beta_t(i) = \beta_t(i) + a_{ij}b_j(\mathcal{O}_{t+1})\beta_{t+1}(j)$ 
        next  $j$ 
        // scale  $\beta_t(i)$  with same scale factor as  $\alpha_t(i)$ 
         $\beta_t(i) = c_t\beta_t(i)$ 
    next  $i$ 
next  $t$ 
```

Solution to problem 1 
Solution to problem 2 
Solution to problem 3 

NOTE: Solution to problem 2 is NOT shown completely in this pseudo-code

Gammas

- Using scaled alphas and betas




Solution to problem 1 
Solution to problem 2 
Solution to problem 3 

```
for t = 0 to T - 2
  denom = 0
  for i = 0 to N - 1
    for j = 0 to N - 1
      denom = denom +  $\alpha_t(i)a_{ij}b_j(\mathcal{O}_{t+1})\beta_{t+1}(j)$ 
    next j
  next i
  for i = 0 to N - 1
     $\gamma_t(i) = 0$ 
    for j = 0 to N - 1
       $\gamma_t(i, j) = (\alpha_t(i)a_{ij}b_j(\mathcal{O}_{t+1})\beta_{t+1}(j))/\text{denom}$ 
       $\gamma_t(i) = \gamma_t(i) + \gamma_t(i, j)$ 
    next j
  next i
next t

// Special case for  $\gamma_{T-1}(i)$ 
denom = 0
for i = 0 to N - 1
  denom = denom +  $\alpha_{T-1}(i)$ 
next i
for i = 0 to N - 1
   $\gamma_{T-1}(i) = \alpha_{T-1}(i)/\text{denom}$ 
next i
```

Re-Estimation

- Again, using scaled gammas
- So formulas unchanged

Solution to problem 1 
Solution to problem 2 
Solution to problem 3 

```
// re-estimate  $\pi$ 
for  $i = 0$  to  $N - 1$ 
     $\pi_i = \gamma_0(i)$ 
next  $i$ 

// re-estimate  $A$ 
for  $i = 0$  to  $N - 1$ 
    denom = 0
    for  $t = 0$  to  $T - 2$ 
        denom = denom +  $\gamma_t(i)$ 
    next  $t$ 
    for  $j = 0$  to  $N - 1$ 
        numer = 0
        for  $t = 0$  to  $T - 2$ 
            numer = numer +  $\gamma_t(i, j)$ 
        next  $t$ 
         $a_{ij} = \text{numer} / \text{denom}$ 
    next  $j$ 
next  $i$ 

// re-estimate  $B$ 
for  $i = 0$  to  $N - 1$ 
    denom = 0
    for  $t = 0$  to  $T - 1$ 
        denom = denom +  $\gamma_t(i)$ 
    next  $t$ 
    for  $j = 0$  to  $M - 1$ 
        numer = 0
        for  $t = 0$  to  $T - 1$ 
            if ( $\mathcal{O}_t == j$ ) then
                numer = numer +  $\gamma_t(i)$ 
            end if
        next  $t$ 
         $b_i(j) = \text{numer} / \text{denom}$ 
    next  $j$ 
next  $i$ 
```

Stopping Criteria

- Check that probability increases

In practice, we want:

$$\text{logProb} > \text{oldLogProb} + \varepsilon$$

- And don't exceed max iterations

```
// Compute  $\log[P(\mathcal{O} | \lambda)]$ 
```

```
logProb = 0
```

```
for  $i = 0$  to  $T - 1$ 
```

```
    logProb = logProb +  $\log(c_i)$ 
```

```
next  $i$ 
```

```
logProb =  $-\text{logProb}$ 
```

```
// To iterate or not to iterate, that is the question...
```

```
iters = iters + 1
```

```
if (iters < maxIters and logProb > oldLogProb) then
```

```
    oldLogProb = logProb
```

```
    goto  $\alpha$ -pass
```

```
else
```

```
    output  $\lambda = (\pi, A, B)$ 
```

```
end if
```

Solution to problem 1



Solution to problem 2



Solution to problem 3



Recap solutions

- To problem 1:

$$\log(P(\mathcal{O} \mid \lambda)) = - \sum_{j=0}^{T-1} \log c_j$$

- To problem 2:

$$\gamma_t(i) = \frac{\alpha_t(i)\beta_t(i)}{P(\mathcal{O} \mid \lambda)} \qquad \tilde{X}_t = \max_i \gamma_t(i)$$

- To problem 3:

The full code seen before...

HMM Applications

- Applications ***not*** chosen because they are the most ingenious or clever, but because they:
 - Illustrate the basics
 - Show off the strengths of technique
 - Easy to understand and appreciate
- Academic publications usually favor novel, different, clever, (weird?), ...



HMM for English Text Analysis

- Marvin the Martian arrives on Earth
 - Marvin sees written English text
 - Wants to learn something about it
 - Martians know HMMs, but not English
- So, strip out all non-letters, make all letters lower-case
 - 27 symbols (26 letters, word-space)
 - Train HMM on long sequence of “symbols”

Training

- For first training case, initialize
 - $N = 2$ and $M = 27$
 - Elements of A and π are all $\sim 1/2$
 - Elements of B are each $\sim 1/27$
- We'll use 50,000 symbols for training
- After 1st iteration: $\log P(O|\lambda) \approx -165097$
- After 100th iteration: $\log P(O|\lambda) \approx -137305$

The A and π Matrices

- Matrices A and π converge to:

$$\pi = \begin{bmatrix} 0.00000 & 1.00000 \end{bmatrix} \quad A = \begin{bmatrix} 0.25596 & 0.74404 \\ 0.71571 & 0.28429 \end{bmatrix}$$

- What does this tell us?
 - Started in hidden state 1 (not state 0)
 - And we know transition probabilities between hidden states
- Nothing too interesting here
 - We don't (yet) know about hidden states

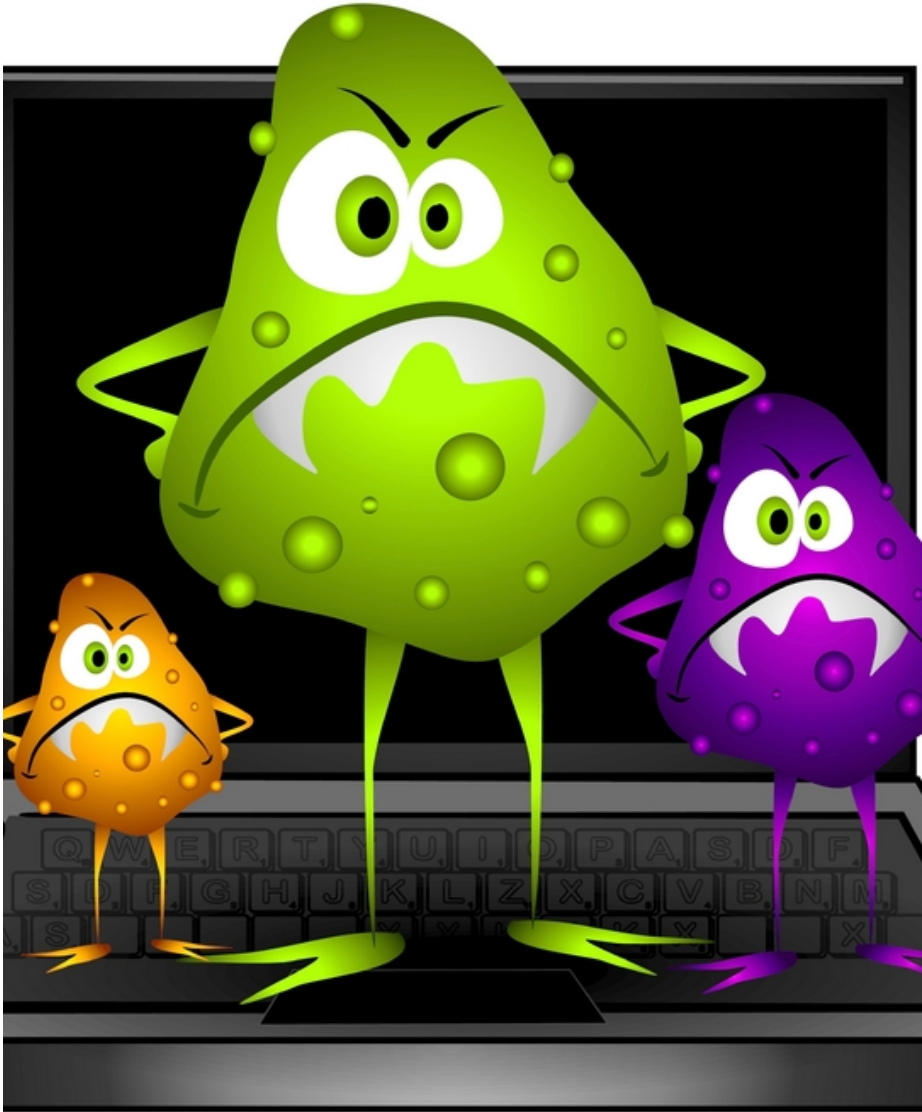
The B Matrix (Transpose)

- What does B matrix tell us?
- This is **very interesting!**
Why???

	Initial		Final	
a	0.03735	0.03909	0.13845	0.00075
b	0.03408	0.03537	0.00000	0.02311
c	0.03455	0.03537	0.00062	0.05614
d	0.03828	0.03909	0.00000	0.06937
e	0.03782	0.03583	0.21404	0.00000
f	0.03922	0.03630	0.00000	0.03559
g	0.03688	0.04048	0.00081	0.02724
h	0.03408	0.03537	0.00066	0.07278
i	0.03875	0.03816	0.12275	0.00000
j	0.04062	0.03909	0.00000	0.00365
k	0.03735	0.03490	0.00182	0.00703
l	0.03968	0.03723	0.00049	0.07231
m	0.03548	0.03537	0.00000	0.03889
n	0.03735	0.03909	0.00000	0.11461
o	0.04062	0.03397	0.13156	0.00000
p	0.03595	0.03397	0.00040	0.03674
q	0.03641	0.03816	0.00000	0.00153
r	0.03408	0.03676	0.00000	0.10225
s	0.04062	0.04048	0.00000	0.11042
t	0.03548	0.03443	0.01102	0.14392
u	0.03922	0.03537	0.04508	0.00000
v	0.04062	0.03955	0.00000	0.01621
w	0.03455	0.03816	0.00000	0.02303
x	0.03595	0.03723	0.00000	0.00447
y	0.03408	0.03769	0.00019	0.02587
z	0.03408	0.03955	0.00000	0.00110
space	0.03688	0.03397	0.33211	0.01298

Conclusion

- The B matrix tells Marvin about consonants and vowels
- He made no assumption *a priori*
- The training itself “learned” this info
- This is the essence of machine learning!
 - We don’t have to think too much
 - Machine does all the hard work
 - Machine “learns” (and so do we)



Detecting “Undetectable” Malware

- Malware is malicious software
Designed to do bad things
- Most common form of detection is based on signatures
What is a signature?
- Malware writers try to avoid detection
 - How to avoid signature detection?
One option is to morph/modify code

Metamorphic Malware

- Malware that morphs with each new infection is called “**metamorphic**”
- How to write **metamorphic** malware?
Standalone generator vs malware that “carries its own morphing engine”
- How to detect metamorphic malware?
Also not easy, it’s a research topic
 - Machine learning

Real-World Metamorphism

Examples of metamorphic malware

	Evol	Zmist	Zperm	Regswap	MetaPHOR
Instruction substitution	—	—	—	✓	—
Instruction permutation	✓	✓	—	—	✓
Garbage code	✓	✓	—	—	✓
Variable substitution	✓	✓	—	✓	✓
Alter control flow	—	✓	✓	—	✓

- Why “garbage code” (or “dead code”)?
- Why is substitution not used more?

Metamorphic Example

From malware known as
NGVCK

3 morphed versions

All do the same thing

Each has different opcode
sequence

What can you say about
signature(s)?

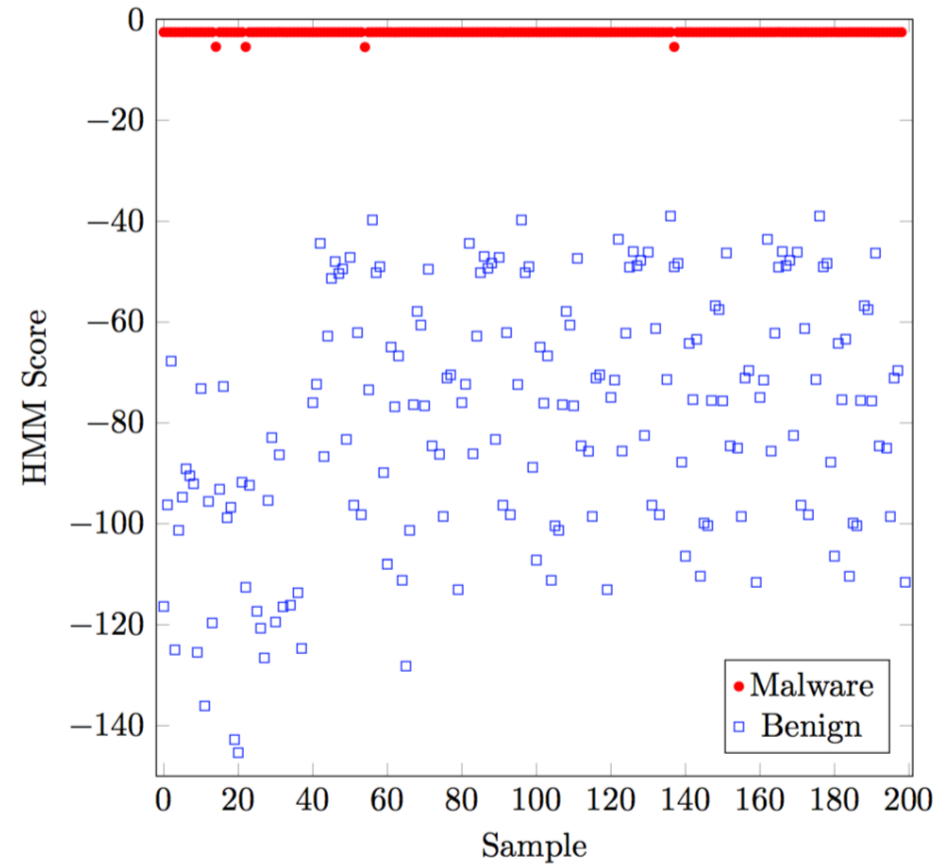
Base Version	
	call delta
delta:	pop ebp sub ebp, offset delta
Morphed Version 1	
	call delta
delta:	sub dword ptr[esp], offset delta pop eax mov ebp, eax
Morphed Version 2	
	add ecx, 0031751B ; junk call delta
delta:	sub dword ptr[esp], offset delta sub ebx, 00000909 ; junk mov edx, [esp] xchg ecx, eax ; junk add esp, 00000004 and ecx, 00005E44 ; junk xchg edx, ebp

Experiments

- Seed the generator with NGVCK virus
- Generated 200 morphed copies
 1. Assemble each morphed *asm* into *exe*
 2. Verify that seed virus detected by AV...
 3. ...and morphed copies **not detected** by AV
- Disassemble *exes*, extract opcodes
 1. Train HMMs, using 5-fold cross validation
 2. Score each model vs 40 benign samples

Scatterplot of Results

- Ideal separation!
Best-case scenario
- Using HMM...
Easy to separate benign
from virus
Easy to detect!
- Undetectable using
signatures...



What's Going On Here?

- Why can we detect this morphed malware using HMMs
But not using signatures?
- Signatures are easily disrupted by simple transposition strategy
- HMM not affected by transposition
 - HMM “sees” differences between viruses and benign, in spite of this morphing

Statistical properties are not significantly affected by morphing

Conclusion

- Transposition is highly effective anti-signature strategy
- But ineffective for machine learning (or statistical-based) analysis
- Can virus writer defeat both signatures and machine learning?
 - Yes, but something more is required.
 - There is much more to say about this...