

# Sequence Models - 2

Spring 2021

COS 484/584

- **Topics:** Lectures up to March 8 (RNNs, Neural LMs)
- Logistics of the exam will be announced on Canvas tomorrow
- Last year's midterm will be made available
  - Not all topics are relevant for this year, but you will get a sense for format and types of questions
- Midterm review: COS 484 precept this week (March 5)

## Midterm





- 1. Set of states  $S = \{1, 2, ..., K\}$  and set of observations O
- 2. Initial state probability distribution  $\pi(s_1)$
- 3. Transition probabilities  $P(s_{t+1} | s_t)$  (OR  $\theta_{s_t \rightarrow s_{t+1}}$ )
- 4. Emission probabilities  $P(o_t | s_t)$  (OR  $\phi_{s_t \rightarrow o_t}$ )

## Recap: Hidden Markov Models

Strong assumptions



## Maximum Entropy Markov Models

## Generative vs Discriminative

• HMM is a generative model

• Can we model  $P(s_1, \ldots, s_n | o_1, \ldots, o_n)$  directly?

Generative

Text classification

Naive Bayes: P(c)P(d | c)

Sequence prediction

HMM:

 $P(s_1,\ldots,s_n)P(o_1,\ldots,o_n \mid s_1,\ldots,s_n)$ 

Discriminative

Logistic Regression:  $P(c \mid d)$ 

MEMM:  $P(s_1, ..., s_n | o_1, ..., o_n)$ 

(No factorization)



• Compute the posterior directly:

• 
$$\hat{S} = \arg \max_{S} P(S \mid O) \approx \arg \max_{S} \hat{S}$$

• Use features:  $P(s_i | o_i, s_{i-1}) \propto \exp(w \cdot f(s_i, o_i, s_{i-1}))$ 

Maximum Entropy Markov Model



 $O = \langle o_1, o_2, \dots, o_n \rangle$ 

 $\prod_{i} P(s_i | o_i, s_{i-1}) \longleftarrow$ Features ~ weights

No factorization into transition, emission

(Bigram MEMM)







$$\hat{S} = \arg\max_{S} P(S \mid O) = \arg\max_{S} \prod_{i} P(s_i \mid o_n, o_{i-1}, \dots, o_1, s_{i-1}, \dots, s_1)$$

$$P(s_i | s_{i-1}, \dots, s_1, O) \propto \exp(w \cdot f(s_i, s_{i-1}, \dots, s_1, O))$$

## MEMM



• In general, we can use all observations and all previous states:

Why couldn't we do this with HMMs?



$$\langle t_i, w_{i-2} \rangle, \langle t_i, w_{i-1} \rangle, \langle t_i, w_i \rangle, \langle t_i, w_{i+1} \rangle, \langle t_i, w_{i+2} \rangle \langle t_i, t_{i-1} \rangle, \langle t_i, t_{i-2}, t_{i-1} \rangle, \langle t_i, t_{i-1}, w_i \rangle, \langle t_i, w_{i-1}, w_i \rangle \langle t_i, w_i, w_{i+1} \rangle,$$

Feature templates

## Features in an MEMM

 $t_i = VB$  and  $w_{i-2} = Janet$  $t_i = VB$  and  $w_{i-1} = will$  $t_i = VB$  and  $w_i = back$  $t_i = VB$  and  $w_{i+1} = the$  $t_i$  = VB and  $w_{i+2}$  = bill  $t_i = VB$  and  $t_{i-1} = MD$  $t_i$  = VB and  $t_{i-1}$  = MD and  $t_{i-2}$  = NNP  $t_i = VB$  and  $w_i = back$  and  $w_{i+1} = the$ 

### t = tags (states) w = words (observations)

### Features



DT JJ Incorrect

DT NN Correct

The old

Which of these feature templates would help most to tag 'old' correctly? A)  $\langle t_i, w_i \rangle$ B)  $\langle t_i, w_i, w_{i-1} \rangle$ C)  $\langle t_i, w_i, w_{i-1}, w_{i+1} \rangle$ D)  $\langle t_i, w_i, w_{i-1}, w_{i+1}, w_{i+2} \rangle$ 

## Features in an MEMM

- NN DT NN
- VB DT NN
- the man boat

t = tagsw = words





$$\hat{S} = \arg\max_{S} P(S \mid O) = \arg\max_{S} P(S \mid O)$$

(assume features only on previous time step and current obs)

• Greedy decoding:



 $\operatorname{rgmax}_{S} \Pi_{i} P(s_{i} \mid o_{i}, s_{i-1})$ 

s = argmax P(s|The)

$$\hat{S} = \arg\max_{S} P(S \mid O) = \arg$$

• Greedy decoding:



 $g \max_{G} \prod_{i} P(s_i \mid o_i, s_{i-1})$ S

 $S_2 = \alpha Agmax P(S|cat, DT)$ 

= NN

$$\hat{S} = \arg\max_{S} P(S \mid O) = \arg$$

• Greedy decoding:



 $g\max_{S} \prod_{i} P(s_i \mid o_i, s_{i-1})$ S

$$\hat{S} = \arg\max_{S} P(S \mid O)$$

Greedy decoding

• Viterbi decoding:

k  $\bigwedge$ DP Lattice (Best sequence ending in  $s_i$ )

 $= \arg\max_{S} \prod_{i} P(s_i \mid o_i, s_{i-1})$ 

 $M[i,j] = \max_{i} M[i-1,k] P(s_{j} | o_{i}, s_{k}) \quad 1 \le k \le K \quad 1 \le i \le n$ # states # timesteps

(or equivalent log form)

# What do you think of the computational complexity of Viterbi decoding for bigram MEMMs compared to decoding for bigram HMMs?

- A) More operations in MEMM
- B) More operations in HMM
- C) Equal

D) Depends on number of features in MEMM

$$M[i,j] = \max_{k} M[i-1,k] F$$

$$M[i,j] = \max_{k} M[i-1,k]$$





# MEMM: Learning

- **Gradient descent:** similar to logistic regression!
  - $P(s_i | s_1, \ldots, s_{i-1}, O)$
- Given: pairs of (S, O) wher
  - Loss for one sequence, L

$$= \frac{\exp(w \cdot f(s_1, \dots, s_{i-1}, s_i, O))}{\sum_{s'} \exp(w \cdot f(s_1, \dots, s_{i-1}, s', O))}$$

re each 
$$S = \langle s_1, s_2, \ldots, s_n \rangle$$

$$= -\sum_{i} \log P(s_{i} | s_{1}, \dots, s_{i-1}, O)$$

• Compute gradients with respect to weights w and update

## Label bias



HMM

$$P(JJ \mid DT) P(\text{old} \mid JJ) P(N)$$

$$P(NN \mid DT) P(\text{old} \mid NN) P(N)$$

Low entropy transitions between labels may override the effect of observations



The/? old/? man/? the/? boat/?



### **Stanford Parser**

Please enter a sentence to be parsed:



### Your query

The old man the boat

Tagging

The/DT old/JJ man/NN the/DT boat/NN

	/
Contonco	Deree
	Parse
	i uioc

### Solution? **Conditional Random Fields** (advanced)



## Expectation Maximization

## Expectation Maximization

- Unsupervised learning method
- Can train a model without any labeled data
- Treat unknowns as "hidden" or "latent" variables
- Maximize (expected) likelihood of observed data

# EM: Some intuition

- Let's say I have 3 coins in my pocket,
  - Coin 0 has probability  $\lambda$  of heads Coin 1 has probability  $p_1$  of heads Coin 2 has probability  $p_2$  of heads
- For each trial:
  - First I toss Coin 0 If coin 0 turns up **heads**, I toss coin 1 three times If coin 0 turns up **tails**, I toss coin 2 three times
    - I don't tell you the results of the coin 0 toss, or whether coin **1 or coin 2 was tossed**, but I tell you how many heads/tails are seen after each trial
- You see the following sequence:  $\langle H, H, H \rangle, \langle T, T, T \rangle, \langle H, H, H \rangle, \langle T, T, T \rangle, \langle H, H, H \rangle$



What would you estimate as the values for  $\lambda, p_1, p_2$  ? A) 1/2, 1, 0 B) 3/5, 1, 0 C) 1/2, 1/2, 1/2



## Maximum Likelihood Estimate

- Data points  $x_1, x_2, \ldots, x_n$  from (finite or countable) set  $\mathcal{X}$
- Parameter vector  $\theta$
- Parameter space  $\Omega$ , i.e.  $\theta \in \Omega$
- We have a distribution  $P(x | \theta)$  for any  $\theta \in \Omega$ , such that

 $x \in \mathcal{X}$ 

- $P(x \mid \theta^*)$  for some  $\theta^* \in \Omega$ 
  - This  $\theta^*$  is the MLE

 $\sum P(x \mid \theta) = 1 \text{ and } P(x \mid \theta) \ge 0 \quad \forall x$ 

• Assume data points are drawn independently and identically distributed from a distribution

# Log Likelihood

- Data points  $x_1, x_2, \ldots, x_n$  from (finite or countable) set  $\mathcal{X}$
- Parameter vector heta and a parameter space  $\Omega$
- Probability distribution  $P(x | \theta)$  for any  $\theta \in \Omega$
- Data Likelihood( $\theta$ ) =  $P(x_1,$

• Log-likelihood,  $L(\theta) = \sum_{i=1}^{n} \log P(x_i | \theta)$ i=1

$$x_2, \dots, x_n | \theta) = \prod_{i=1}^n P(x_i | \theta)$$
 (Each  $x_i$  is a data |



## Example I: Coin Tossing

- heads and tails, e.g.
  - НТНТННННТТТ
- coming up heads
- Parameter space  $\Omega = [0,1]$
- Distribution  $P(x | \theta) = \begin{cases} \theta \text{ if } x = H \\ 1 \theta \text{ if } x = T \end{cases}$



•  $\mathscr{X} = \{H, T\}$ . Our data points  $x_1, x_2, \ldots, x_n$  are a sequence of

• Parameter vector  $\theta$  is a single parameter, i.e probability of coin

What distribution is this? A) Binomial

B) Bernoulli

C) Multinomial

D) Gaussian



## Example 2: Markov chains



- state  $\phi$  and initial transition  $\phi \rightarrow s_1$  (how many parameters?)
- Let  $T(\alpha) \subset T$  be all the transitions of the form  $\alpha \to \beta$  (i.e. all transitions from state  $\alpha$ )
- for all  $\alpha \in S$ ,  $\sum \theta_t = 1$  $t \in T(\alpha)$

•  ${\mathscr X}$  is the set of all possible state (e.g tag) sequences created by the underlying generative process. Our sample is n sequences  $X_1, \ldots, X_n$  such that each  $X_i \in \mathcal{X}$ , consists of a sequence of states  $s_1, s_2, s_3, \ldots$ 

•  $\theta_T$  is the vector of all transition  $(s_i \rightarrow s_j)$  parameters. Without loss of generality, assume a dummy start

• Then, parameter space  $\Omega$  is the set of  $\theta \in [0,1]^{|S+1||S|}$  where S is set of all states (tags), such that: (why?)

• Now, if  $\theta_T$  is the vector of all transition parameters

• Then, we have for a sequence X:  $P(X \mid \theta) = \prod_{t \in \mathcal{O}} \theta_t^{Count(X,t)}$  $t \in T$ 

sequence X

$$\implies \log P(X \mid \theta) = \sum_{t \in T} Cout$$

Example 2: Markov chains

- where Count(X, t) is the number of times transition t is seen in

 $nt(X, t) \log \theta_t$ 

## MLE for Markov chains

### • We have $\log P(X|\theta) = \sum Count(X,t) \log \theta_t$ $t \in T$

in sequence X

 And,  $L(\theta) = \sum \log P(X_i | \theta) = \sum \sum Count(X_i, t) \log \theta_t$  $i \qquad i \quad t \in T$ 

where Count(X, t) is the number of times transition t is seen



## MLE for Markov chains

• 
$$L(\theta) = \sum_{i} \log P(X_i | \theta) = \sum_{i} \sum_{t \in T} Cour$$

• Solve  $\theta_{MLE} = \underset{\theta \in \Omega}{\arg \max L(\theta)}$ 

$$\implies \text{find } \theta \quad \text{s. t. } \frac{\partial L(\theta)}{\partial \theta} = 0 \text{ with ap}$$

• This gives: 
$$\theta_t = \frac{\sum_i Count(X_i, t)}{\sum_i \sum_{t' \in T(\alpha)} Count(X_i, t)}$$

where *t* is of the form  $\alpha \rightarrow \beta$  for some  $\beta$ 

a transition.

 $nt(X_i, t) \log \theta_t$ 

opropriate probability constraints

t')

• Intuitively, the denominator is simply counting all occurrences of state  $\alpha$  at the start of

- Now say we have two sets  $\mathcal{X}$  and  $\mathcal{Y}$ , and a joint distribution  $P(x, y | \theta)$
- If we had **fully observable data**,  $(x_i, y_i)$  pairs, then log likelihood can be estimated as:  $L(\theta) = \sum \log P(x_i, y_i | \theta)$
- If we have **partially observable data**, *x<sub>i</sub>* examples only, then  $L(\theta) = \sum \log P(x_i | \theta)$  $= \sum \log \sum P(x_i, y | \theta)$

y∈¥

## Models with hidden variables

(y is hidden)

**Unsupervised Learning** 

## **Expectation Maximization**

then

 $L(\theta) =$ 

for finding

Maximization

• If we have **partially observable data**, x<sub>i</sub> examples only,

$$\sum_{i} \log \sum_{y \in \mathcal{Y}} P(x_i, y \mid \theta)$$

The EM (Expectation Maximization) algorithm is a method



- In the three coins example,  $\mathcal{Y} = \{H, T\}$  (possible outcomes of coin 0)  $\mathcal{X} = \{HHH, TTT, HTT, THH, HHT, TTH, HTH, THT\}$  $\theta = \{\lambda, p_1, p_2\}$
- and  $P(x, y | \theta) = P(y | \theta) P(x | y, \theta)$ where

and

$$P(x | y, \theta) = \begin{cases} p_1^h \\ p_2^h \end{cases}$$

### The three coins example

(all possible observations of length 3)

 $P(y | \theta) = \begin{cases} \lambda \text{ if } y = H \\ 1 - \lambda \text{ if } y = T \end{cases}$ 

 $(1 - p_1)^t$  if y = H $(1 - p_2)^t$  if y = T

• Calculating various probabilities:  $P(x = THT, y = H | \theta) = \lambda p_1 (1 - p_1)^2$  $P(x = THT, y = T | \theta) = (1 - \lambda)p_2(1 - p_2)^2$ 

$$P(x = THT | \theta) = P(x = THT, y = \lambda p_1 (1 - p_1)^2 + \lambda p_1 (1 - p$$

$$P(y = H | x = THT, \theta) = \frac{P(x = T)}{P(x)}$$

### The three coins example

 $= H | \theta) + P(x = THT, y = T | \theta)$  $(1 - \lambda)p_2(1 - p_2)^2$ 

 $THT, y = H[\theta]$  $= THT | \theta$ )  $\lambda p_1 (1 - p_1)^2$  $\lambda p_1 (1-p_1)^2 + (1-\lambda) p_2 (1-p_2)^2$ 

- Fully observed data might look like:
- In this case, maximum likelihood estimates are:

 $(\langle HHH \rangle, H), (\langle TTT \rangle, T), (\langle HHH \rangle, H), (\langle TTT \rangle, T), (\langle HHH \rangle, H)$ 

 $\lambda = \frac{3}{5}$  $p_1 = \frac{9}{9}$  $p_{\gamma}$ 0

• Partially observed data might look like:

 $\langle HHH \rangle, \langle TTT \rangle, \langle HHH \rangle, \langle TTT \rangle, \langle HHH \rangle$ 

• How do we find the MLE parameters?

## The three coins example

EM!

• Partially observed data might look like:

 $\langle HHH \rangle, \langle TTT \rangle, \langle HHH \rangle, \langle TTT \rangle, \langle HHH \rangle$ 

• Assume we guess the current parameters to be some  $\lambda, p_1, p_2$ . Then:

$$P(y = H | x = \langle HHH \rangle) = \frac{P(\langle HHH \rangle, H)}{P(\langle HHH \rangle, H)}$$
$$= \frac{\lambda p_1^3}{\lambda p_1^3 + (1 - \lambda)p}$$
$$P(y = H | x = \langle TTT \rangle) = \frac{P(\langle TTT \rangle, H) + \lambda(1 - \lambda)p}{P(\langle TTT \rangle, H) + \lambda(1 - \lambda)p}$$
$$= \frac{\lambda(1 - \mu)^3 + \mu(1 - \lambda)p}{\lambda(1 - \mu)^3 + \mu(1 - \lambda)p}$$

### The three coins example

 $HH\rangle, H)$  $+ P(\langle HHH \rangle, T)$ 

 $p_2^3$  $\langle T \rangle, H \rangle$  $(p_1)^3 - \lambda)(1 - p_2)^3$ 

• If the current parameters are  $\lambda, p_1, p_2$  $P(y = H | x = \langle HHH \rangle) = \frac{P(\langle HHH \rangle, H)}{P(\langle HHH \rangle, H) + P(\langle HHH \rangle, T)}$  $=\frac{\lambda p_1^3}{\lambda p_1^3 + (1-\lambda)p_2^3}$  $P(y = H | x = \langle TTT \rangle) = \frac{P(\langle HHH \rangle, H)}{P(\langle TTT \rangle, H) + P(\langle TTT \rangle, T)}$  $= \frac{\lambda(1-p_1)^3}{\lambda(1-p_1)^3 + (1-\lambda)(1-p_2)^3}$ • If  $\lambda = 0.3$ ,  $p_1 = 0.3$ ,  $p_2 = 0.6$ 

P If 
$$\lambda = 0.3$$
,  $p_1 = 0.3$ ,  $p_2 = 0.6$ :  
 $P(y = H | x = 0.0)$   
 $P(y = H | x = 0.0)$ 

### The three coins example

 $\langle HHH \rangle = 0.0508$  $\langle TTT \rangle$ ) = 0.6967

- observed data might look like:
  - $(\langle \text{HHH} \rangle, H) = P(y = \text{H} | \text{HHH}) = 0.0508$
  - ((HHH), T) = P(y = T | HHH) = 0.9492

  - $((\text{HHH}), H) \quad P(y = \text{H} | \text{HHH}) = 0.0508$
  - $((\text{HHH}), T) \quad P(y = T | \text{HHH}) = 0.9492$
  - $((TTT), H) \quad P(y = H | TTT) = 0.6967$
  - ((TTT), T) = P(y = T | TTT) = 0.3033

  - $((\text{HHH}), T) \qquad P(y = T | \text{HHH}) = 0.9492$

Treat this as a pseudo-annotated dataset (with appropriate weights) and use MLE!

## The three coins example

After filling in hidden variables for each example, partially

- $((TTT), H) \quad P(y = H | TTT) = 0.6967$
- ((TTT), T) = P(y = T | TTT) = 0.3033

((HHH), H) = P(y = H | HHH) = 0.0508

each pair sums to 1

$(\langle \text{HHH} \rangle, H)$	P(y = H   HHH) = 0.0508
$(\langle \text{HHH} \rangle, T)$	$P(y = T \mid HHH) = 0.9492$
$(\langle TTT \rangle, H)$	P(y = H   TTT) = 0.6967
$(\langle TTT \rangle, T)$	P(y = T   TTT) = 0.3033
$(\langle HHH \rangle, H)$	P(y = H   HHH) = 0.0508
$(\langle \text{HHH} \rangle, T)$	$P(y = T \mid HHH) = 0.9492$
$(\langle TTT \rangle, H)$	P(y = H   TTT) = 0.6967
$(\langle TTT \rangle, T)$	P(y = T   TTT) = 0.3033
$(\langle \text{HHH} \rangle, H)$	P(y = H   HHH) = 0.0508
$(\langle HHH \rangle, T)$	$P(y = T \mid HHH) = 0.9492$

• New estimates:





$(\langle \text{HHH} \rangle, H)$	P(y = H   HHH) = 0.0508
$(\langle \text{HHH} \rangle, T)$	$P(y = T \mid HHH) = 0.9492$
$(\langle TTT \rangle, H)$	P(y = H   TTT) = 0.6967
$(\langle TTT \rangle, T)$	$P(y = T \mid TTT) = 0.3033$
$(\langle HHH \rangle, H)$	P(y = H   HHH) = 0.0508
$(\langle \text{HHH} \rangle, T)$	$P(y = T \mid HHH) = 0.9492$
$(\langle TTT \rangle, H)$	P(y = H   TTT) = 0.6967
$(\langle TTT \rangle, T)$	P(y = T   TTT) = 0.3033
$(\langle \text{HHH} \rangle, H)$	$P(y = H \mid HHH) = 0.0508$
$(\langle \text{HHH} \rangle, T)$	$P(y = T \mid HHH) = 0.9492$

• New estimates:







$(\langle \text{HHH} \rangle, H)$	P(y = H   HHH) = 0.0508
$(\langle \text{HHH} \rangle, T)$	$P(y = T \mid HHH) = 0.9492$
$(\langle TTT \rangle, H)$	P(y = H   TTT) = 0.6967
$(\langle TTT \rangle, T)$	$P(y = T \mid TTT) = 0.3033$
$(\langle \text{HHH} \rangle, H)$	P(y = H   HHH) = 0.0508
$(\langle \text{HHH} \rangle, T)$	$P(y = T \mid HHH) = 0.9492$
$(\langle TTT \rangle, H)$	P(y = H   TTT) = 0.6967
$(\langle \mathrm{TTT} \rangle, T)$	P(y = T   TTT) = 0.3033
$(\langle \text{HHH} \rangle, H)$	P(y = H   HHH) = 0.0508
$(\langle \text{HHH} \rangle, T)$	$P(y = T \mid HHH) = 0.9492$

• New estimates:





- Begin with parameters:  $\lambda = 0.3, p_1 = 0.3, p_2 = 0.6$
- Fill in hidden variables, using  $P(y = H | x = \langle HHH \rangle) = 0.0508$  $P(y = H | x = \langle TTT \rangle) = 0.6967$
- This gives us a pseudo-annotated dataset with **fractional** counts
- Re-estimate parameters to be  $\lambda = 0.3092, p_1 = 0.0987, p_2 = 0.8244$

### Summary

## EM iterations

Iteration	λ	$p_1$	$p_2$	$\tilde{p}_1$	$\tilde{p}_2$	$\bar{p}_3$	$\tilde{p}_4$
0	0.3000	0.3000	0.6000	0.0508	0.6967	0.0508	0.6967
1	0.3738	0.0680	0.7578	0.0004	0.9714	0.0004	0.9714
2	0.4859	0.0004	0.9722	0.0000	1.0000	0.0000	1.0000
3	0.5000	0.0000	1.0000	0.0000	1.0000	0.0000	1.0000

which always shows tails, and is picking between them with equal probability ( $\lambda = 0.5$ ).

 $x_2$  and  $x_4$ , whereas coin 2 generated  $x_1$  and  $x_3$ 

The coin example for  $x = \{\langle HHH \rangle, \langle TTT \rangle, \langle HHH \rangle, \langle TTT \rangle\}$ . The solution that EM reaches is intuitively correct: the coin tosser has two coins, one which always shows heads, and another

Posterior probabilities  $\bar{p}_i$  show that we are certain that coin 1 (tail-biased) generate

### EM iterations

Iteration	$\lambda$	$p_1$	$p_2$	$\tilde{p}_1$	$\tilde{p}_2$	$\tilde{p}_3$	$\tilde{p}_4$	$\tilde{p}_5$
0	0.3000	0.3000	0.6000	0.0508	0.6967	0.0508	0.6967	0.0508
1	0.3092	0.0987	0.8244	0.0008	0.9837	0.0008	0.9837	0.0008
2	0.3940	0.0012	0.9893	0.0000	1.0000	0.0000	1.0000	0.0000
3	0.4000	0.0000	1.0000	0.0000	1.0000	0.0000	1.0000	0.0000

Coin example for  $\{\langle HHH \rangle, \langle TTT \rangle, \langle HHH \rangle, \langle TTT \rangle, \langle HHH \rangle\}$ 

 $\lambda$  is now 0.4, indicating that coin 0 has a probability 0.4 of selecting the tailbiased coin

## EM iterations

Iteration	λ	$p_1$	$p_2$	$\tilde{p}_1$	$\tilde{p}_2$	$\tilde{p}_3$	$\tilde{p}_4$
0	0.3000	0.3000	0.6000	0.1579	0.6967	0.0508	0.6967
1	0.4005	0.0974	0.6300	0.0375	0.9065	0.0025	0.9065
2	0.4632	0.0148	0.7635	0.0014	0.9842	0.0000	0.9842
3	0.4924	0.0005	0.8205	0.0000	0.9941	0.0000	0.9941
4	0.4970	0.0000	0.8284	0.0000	0.9949	0.0000	0.9949

EM selects a tails-only coin, and a coin which is heavily heads-biased (tail-biased) is far more likely.

### Coin example for $x = \{\langle HHT \rangle, \langle TTT \rangle, \langle HHH \rangle, \langle TTT \rangle \}$ .

 $(p_2 = 0.8284)$ . It's certain that  $x_1$  and  $x_3$  were generated by coin 2 since they contain heads.  $x_2$  and  $x_4$  could have been generated by either coin but coin 1

## Initialization matters

Iteration	$\lambda$	$p_1$	$p_2$	$ ilde p_1$	$\tilde{p}_2$	$\tilde{p}_3$	$\tilde{p}_4$
0	0.3000	0.7000	0.7000	0.3000	0.3000	0.3000	0.3000
1	0.3000	0.5000	0.5000	0.3000	0.3000	0.3000	0.3000
2	0.3000	0.5000	0.5000	0.3000	0.3000	0.3000	0.3000
3	0.3000	0.5000	0.5000	0.3000	0.3000	0.3000	0.3000
4	0.3000	0.5000	0.5000	0.3000	0.3000	0.3000	0.3000
5	0.3000	0.5000	0.5000	0.3000	0.3000	0.3000	0.3000
6	0.3000	0.5000	0.5000	0.3000	0.3000	0.3000	0.3000

Coin example for  $x = \{\langle HHH \rangle, \langle TTT \rangle, \langle HHH \rangle, \langle TTT \rangle \}$ .

In this case, EM is stuck at a **saddle point**.

Iteration	$\lambda$	$p_1$	$p_2$	$\tilde{p}_1$	$\tilde{p}_2$	$\tilde{p}_3$	$\tilde{p}_4$
0	0.3000	0.7001	0.7000	0.3001	0.2998	0.3001	0.2998
1	0.2999	0.5003	0.4999	0.3004	0.2995	0.3004	0.2995
2	0.2999	0.5008	0.4997	0.3013	0.2986	0.3013	0.2986
3	0.2999	0.5023	0.4990	0.3040	0.2959	0.3040	0.2959
4	0.3000	0.5068	0.4971	0.3122	0.2879	0.3122	0.2879
5	0.3000	0.5202	0.4913	0.3373	0.2645	0.3373	0.2645
6	0.3009	0.5605	0.4740	0.4157	0.2007	0.4157	0.2007
7	0.3082	0.6744	0.4223	0.6447	0.0739	0.6447	0.0739
8	0.3593	0.8972	0.2773	0.9500	0.0016	0.9500	0.0016
9	0.4758	0.9983	0.0477	0.9999	0.0000	0.9999	0.0000
10	0.4999	1.0000	0.0001	1.0000	0.0000	1.0000	0.0000
11	0.5000	1.0000	0.0000	1.0000	0.0000	1.0000	0.0000

Coin example for  $x = \{\langle HHH \rangle, \langle TTT \rangle, \langle HHH \rangle, \langle TTT \rangle \}$ .

If we initialize  $p_1$  and  $p_2$  even a small amount away from the saddle point  $p_1 = p_2$ , EM diverges and eventually reaches the global maximum