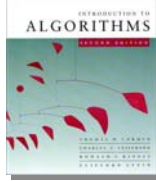


LECTURE 17
Hidden Markov Models II

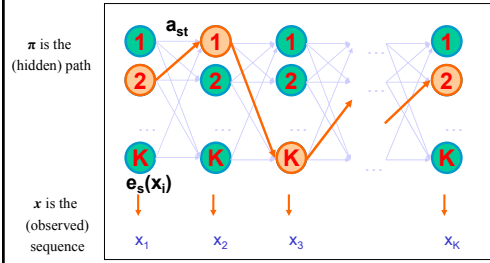


- Review:
 - Path $P(x, \pi) = a_{0\pi_1} \prod_i e_{\pi_i}(x_i) a_{\pi_i \pi_{i+1}}$.
 - Viterbi $\pi^* = \text{argmax}_{\pi} P(x, \pi)$.
 - Forward $P(x) = \sum_{\pi} P(x, \pi)$.
- Posterior decoding $P(\pi_i = k | x)$.
- Supervised learning $\max_{\Lambda} P(x, \pi | \Lambda)$.
- Unsupervised learning $\max_{\Lambda} P(x | \Lambda)$.

Prof. Manolis Kellis

April 10, 2008

Notation: path π , emissions x , state k , time i .

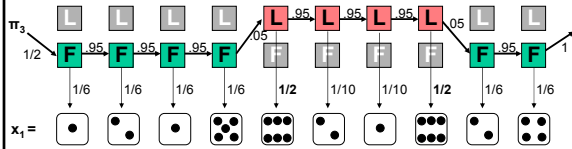


Probability of observing emissions x with known path π .

$$P(x, \pi) = a_{0\pi_1} * \prod_i e_{\pi_i}(x_i) * a_{\pi_i \pi_{i+1}}$$

start emission transition

Example: the dishonest casino



$$P(x_1, \pi_3) = P(\text{transition}) * P(\text{emission}) * P(\text{transition}) * P(\text{emission}) * \dots$$

$$= \frac{1}{2} \times P(1 | \text{Fair}) P(\text{Fair}_{t+1} | \text{Fair}_t) P(2 | \text{Fair}) P(\text{Fair} | \text{Fair}) \dots P(4 | \text{Fair})$$

$$= \frac{1}{2} \times (1/10)^2 \times (1/2)^2 \times (1/6)^6 \times (0.95)^7 \times (0.05)^2$$

$$= 4.6 \times 10^{-11} \quad (\text{least likely})$$

$$P(x_1, \text{all-loaded}) = 7.9 \times 10^{-10} \quad (\text{more likely})$$

$$P(x_1, \text{all-fair}) = 0.5 \times 10^{-9} \quad (\text{most likely})$$

One path

1. Scoring x , one path

$$P(x, \pi)$$

Prob of a path, emissions

Scoring

All paths

2. Scoring x , all paths

$$P(x) = \sum_{\pi} P(x, \pi)$$

Prob of emissions, over all paths

Decoding

3. Viterbi decoding

$$\pi^* = \text{argmax}_{\pi} P(x, \pi)$$

Most likely path

4. Posterior decoding

$$\pi^{\wedge} = \{\pi_i | \pi_i = \text{argmax}_k \sum_{\pi} P(\pi_i = k | x)\}$$

Path containing the most likely state at any time point.

Learning

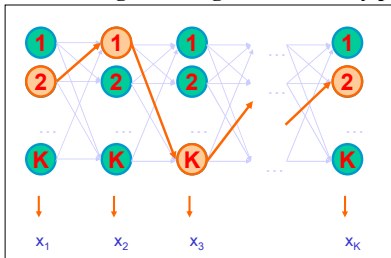
5. Supervised learning, given π
 $\Lambda^* = \text{argmax}_{\Lambda} P(x, \pi | \Lambda)$
6. Unsupervised learning.
 $\Lambda^* = \text{argmax}_{\Lambda} \max_{\pi} P(x, \pi | \Lambda)$
Viterbi training, best path

6. Unsupervised learning

$$\Lambda^* = \text{argmax}_{\Lambda} \sum_{\pi} P(x, \pi | \Lambda)$$

Baum-Welch training, over all paths

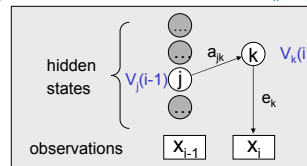
Viterbi decoding - Finding the most likely path



- Find path π^* that maximizes total joint probability $P(x, \pi)$
 - $\pi^* = \text{argmax}_{\pi} P(x, \pi) = a_{0\pi_1} * \prod_i e_{\pi_i}(x_i) * a_{\pi_i \pi_{i+1}}$
- start emission transition

Viterbi algorithm: Calculate $\max P(x, \pi)$ recursively

Define $V_k(i)$ = Probability of the most likely path through state $\pi_i = k$.
Compute $V_k(i+1)$ recursively, as a function of $\max_j \{ V_k(i) \}$.



- Assume we know all V_j values for previous time step $i-1$.

$$V_k(i) = e_k(x_i) * \max_j (V_j(i-1) * a_{jk})$$

current max this emission max ending in state j at step i-1 Transition from state j

The Viterbi Algorithm

Input: $x = x_1 \dots x_N$

Traceback: Follow max pointers back

Initialization: $V_0(0)=1, V_k(0) = 0$, for all $k > 0$

Iteration: $V_k(i) = e_k(x_i) \times \max_j a_{jk} V_j(i-1)$

Termination: $P(x, \pi^*) = \max_k V_k(N)$

In practice: Use log scores for computation

Running time and space: Time: $O(K^2N)$ Space: $O(KN)$

Derivation of Viterbi algorithm

Let $V_k(i) = \max_{\{\pi_1, \dots, \pi_{i-1}\}} P[x_1 \dots x_{i-1}, \pi_1 \dots \pi_{i-1}, x_i, \pi_i = k]$
 = Probability of most likely sequence of states ending at state $\pi_i = k$

Calculate $V_k(i+1)$ recursively as a function of $V_j(i)$

$$V_k(i+1) = \max_{\{\pi_1, \dots, \pi_i\}} P[x_1 \dots x_i, \pi_1, \dots, \pi_i, x_{i+1}, \pi_{i+1} = 1]$$

from definition

$$= \max_{\{\pi_1, \dots, \pi_i\}} P(x_{i+1}, \pi_{i+1} = 1 | x_1 \dots x_i, \pi_1, \dots, \pi_i) P[x_1 \dots x_i, \pi_1, \dots, \pi_i]$$

from Bayes's Law

$$= \max_{\{\pi_1, \dots, \pi_i\}} P(x_{i+1}, \pi_{i+1} = 1 | \pi_i) P[x_1 \dots x_i, \pi_1, \dots, \pi_i]$$

from Markov property (no memory)

$$= \max_k P(x_{i+1}, \pi_{i+1} = 1 | \pi_i = k) \max_{\{\pi_1, \dots, \pi_i\}} P[x_1 \dots x_i, \pi_1, \dots, \pi_i, x_{i+1}, \pi_{i+1} = k]$$

from commutativity of multiplication, max

$$= e_k(x_{i+1}) \max_k a_{ki} V_k(i)$$

from recursive definition of V_k variable

	One path	All paths
Scoring	Scoring x , one path $P(x, \pi)$ Prob of a path, emissions	2. Scoring x , all paths $P(x) = \sum_{\pi} P(x, \pi)$ Prob of emissions, over all paths
Decoding	3. Viterbi decoding $\pi^* = \operatorname{argmax}_{\pi} P(x, \pi)$ Most likely path	4. Posterior decoding $\pi^\wedge = \{\pi_i \pi_i = \operatorname{argmax}_{\pi} P(\pi_i = k x)\}$ Path containing the most likely state at any time point.
Learning	5. Supervised learning, given π $\Lambda^* = \operatorname{argmax}_{\Lambda} P(x, \pi \Lambda)$ 6. Unsupervised learning. $\Lambda^* = \operatorname{argmax}_{\Lambda} \max_{\pi} P(x, \pi \Lambda)$ Viterbi training, best path	6. Unsupervised learning $\Lambda^* = \operatorname{argmax}_{\Lambda} \sum_{\pi} P(x, \pi \Lambda)$ Baum-Welch training, over all paths

2. Model evaluation:

Total $P(x|M)$, summed over all paths

Forward algorithm

Simple: Given the model, generate some sequence x

Given a HMM, we can generate a sequence of length n as follows:

1. Start at state π_1 according to prob $a_{0\pi_1}$
2. Emit letter x_1 according to prob $e_{\pi_1}(x_1)$
3. Go to state π_2 according to prob $a_{\pi_1 \pi_2}$
4. ... until emitting x_n

We have some sequence x that can be emitted by p . Can calculate its likelihood. However, in general, many different paths may emit this same sequence x . How do we find the **total probability** of generating a given x , over any path?

Complex: Given x , was it generated by the model?

Given a sequence x ,
 What is the probability that x was generated by the model (using any path)?

$$P(x) = \sum_{\pi} P(x, \pi)$$

- Challenge: exponential number of paths
- (cheap) alternative:
 - Calculate probability over maximum (Viterbi) path π^*
- (real) solution
 - Calculate sum iteratively using principles of dynamic programming

Calculate total probability $\sum_{\pi} P(x, \pi)$ recursively

- Assume we know f_j for the previous time step (i-1)
- Calculate $f_k(i) = e_k(x_i) * \sum_j (f_j(i-1) * a_{jk})$
 - current max
 - this emission
 - sum ending in state j at step i
 - transition from state j
 - every possible previous state j

The Forward Algorithm – derivation

Define the forward probability:

$$f_i(i) = P(x_1 \dots x_i, \pi_i = i)$$

$$= \sum_{\pi_1 \dots \pi_{i-1}} P(x_1 \dots x_{i-1}, \pi_1 \dots \pi_{i-2}, \pi_{i-1}, \pi_i = i) e_i(x_i)$$

$$= \sum_k \sum_{\pi_1 \dots \pi_{i-2}} P(x_1 \dots x_{i-1}, \pi_1, \dots, \pi_{i-2}, \pi_{i-1} = k) a_{ki} e_i(x_i)$$

$$= \sum_k \boxed{f_k(i-1)} a_{ki} e_i(x_i)$$

$$= e_i(x_i) \sum_k f_k(i-1) a_{ki}$$

The Forward Algorithm

Input: $x = x_1 \dots x_N$

In practice: Sum of log scores is difficult
 → approximate $\exp(1+p+q)$
 → scaling of probabilities

Initialization:
 $f_0(0)=1, f_k(0) = 0$, for all $k > 0$

Iteration:
 $f_k(i) = e_k(x_i) * \sum_j a_{jk} f_j(i-1)$

Termination:
 $P(x, \pi^*) = \sum_k f_k(N)$

Running time and space:
 Time: $O(K^2N)$
 Space: $O(KN)$

The Backward Algorithm

Input: $x = x_1 \dots x_N$

In practice: Sum of log scores is difficult
 → approximate $\exp(1+p+q)$
 → scaling of probabilities

Initialization:
 $b_k(N) = a_{k0}$, for all k

Iteration:
 $b_k(i) = \sum_l e_l(x_{i+1}) a_{kl} b_l(i+1)$

Termination:
 $P(x) = \sum_i a_{0i} e_i(x_i) b_i(1)$

Running time and space:
 Time: $O(K^2N)$
 Space: $O(KN)$

Calculate total end probability recursively

- Assume we know b_l for the next time step (i+1)
- Calculate $b_k(i) = \sum_l (e_l(x_{i+1}) * a_{kl} * b_l(i+1))$
 - current max
 - next emission
 - transition to next state
 - prob sum from state l to end
 - sum over all possible next states

The Backward Algorithm – derivation

Define the backward probability:

$$b_k(i) = P(x_{i+1} \dots x_N | \pi_i = k)$$

$$= \sum_{\pi_{i+1} \dots \pi_N} P(x_{i+1}, x_{i+2}, \dots, x_N, \pi_{i+1}, \dots, \pi_N | \pi_i = k)$$

$$= \sum_l \sum_{\pi_{i+2} \dots \pi_N} P(x_{i+1}, x_{i+2}, \dots, x_N, \pi_{i+1} = l, \pi_{i+2}, \dots, \pi_N | \pi_i = k)$$

$$= \sum_l e_l(x_{i+1}) a_{kl} \sum_{\pi_{i+2} \dots \pi_N} P(x_{i+2}, \dots, x_N, \pi_{i+2}, \dots, \pi_N | \pi_{i+1} = l)$$

$$= \sum_l e_l(x_{i+1}) a_{kl} \boxed{b_l(i+1)}$$

4. Decoding, all paths

Find the likelihood an emission x_i is generated by a state

Calculate $P(\pi_7 = C_p G + \mid x_7 = G)$

- With no knowledge (no characters)
 - $P(\pi_i = k)$ = most likely state (**prior**)
 - Time spent in markov chain states
- With very little knowledge (just that character)
 - $P(\pi_i = k \mid x_i = G)$ = (prior) * (most likely emission)
 - Emission probabilities adjusted for time spent
- With knowledge of entire sequence (all characters)
 - $P(\pi_i = k \mid x = AGCGCG...GATTATCGTCGTA)$
 - Sum over all paths that emit 'G' at position 7
 - **Posterior decoding**

Combining the Forward and Backward algorithms

We want to compute:

$P(\pi_i = k \mid x)$, the probability distribution on the i^{th} position, given x

We can expand it into its forward and backward components:

$$P(\pi_i = k, x) = P(x_1 \dots x_i, \pi_i = k, x_{i+1} \dots x_N)$$

$$= P(x_1 \dots x_i, \pi_i = k) P(x_{i+1} \dots x_N \mid x_1 \dots x_i, \pi_i = k)$$

$$= P(x_1 \dots x_i, \pi_i = k) \mathbf{P}(x_{i+1} \dots x_N \mid \pi_i = k)$$

Forward, $f_k(i)$ Backward, $b_k(i)$

Putting it all together: Posterior decoding

- $P(k) = P(\pi_i = k \mid x) = f_k(i) * b_k(i) / P(x)$
 - Probability that i^{th} state is k , given all emissions x
- Posterior decoding
 - Define most likely state for every of sequence x
 - $\pi_i^* = \text{argmax}_k P(\pi_i = k \mid x)$
- Posterior decoding 'path' π_i^*
 - For classification, more informative than Viterbi path π^*
 - More refined measure of "which hidden states" generated x
 - However, it may give an invalid sequence of states
 - Not all $j \rightarrow k$ transitions may be possible

Summary this far

- Generative model. Hidden states, observed emissions.
 - Generate a random sequence
 - Choose random transition, choose random emission (#0)
- Scoring the likelihood of a sequence
 - Calculate likelihood of annotated path and sequence
 - Multiply emission and transition probabilities (#1)
 - Without specifying a path, total probability of generating x
 - Sum probabilities over all paths
 - Forward algorithm (#3)
- Decoding: Finding the most likely path, given a sequence
 - What is the most likely path generating entire sequence?
 - Viterbi algorithm (#2)
 - What is the most probable state at each time step?
 - Forward + backward algorithms, posterior decoding (#4)
- Next: Learning (#5 and #6)

	One path	All paths
Scoring	1. Scoring x , one path $P(x, \pi)$ Prob of a path, emissions	2. Scoring x , all paths $P(x) = \sum_{\pi} P(x, \pi)$ Prob of emissions, over all paths
	3. Viterbi decoding $\pi^* = \text{argmax}_{\pi} P(x, \pi)$ Most likely path	4. Posterior decoding $\pi^{\wedge} = \{\pi_i \mid \pi_i = \text{argmax}_k \sum_{\pi} P(\pi_i = k \mid x)\}$ Path containing the most likely state at any time point.
Learning	5. Supervised learning, given π $\Lambda^* = \text{argmax}_{\Lambda} P(x, \pi \mid \Lambda)$	6. Unsupervised learning $\Lambda^* = \text{argmax}_{\Lambda} \sum_{\pi} P(x, \pi \mid \Lambda)$
	6. Unsupervised learning. $\Lambda^* = \text{argmax}_{\Lambda} \max_{\pi} P(x, \pi \mid \Lambda)$ Viterbi training, best path	Baum-Welch training, over all paths

Two types of learning: Supervised / Unsupervised

5. Supervised learning
infer model parameters given **labeled** training data

- GIVEN:
 - a HMM M , with unspecified transition/emission probs.
 - labeled sequence x ,
- FIND:
 - parameters $\theta = (E_i, A_{ij})$ that maximize $P[x, \pi | \theta]$

→ Simply count frequency of each emission and transition, as observed in the training data

6. Unsupervised learning
infer model parameters given **unlabeled** training data

- GIVEN:
 - a HMM M , with unspecified transition/emission probs.
 - unlabeled sequence x ,
- FIND:
 - parameters $\theta = (E_i, A_{ij})$ that maximize $P[x | \theta]$

→ Viterbi training:
guess parameters, find optimal Viterbi path (#2), update parameters (#5), iterate

→ Baum-Welch training:
guess parameters, sum over all paths (#4), update parameters (#5), iterate

5: Supervised learning

Estimate model parameters based on **labeled** training data

Two learning scenarios

Case 1. Estimation when the "right answer" is known

Examples:

GIVEN: a genomic region $x = x_1 \dots x_{1,000,000}$, where we have good (experimental) annotations of the CpG islands

GIVEN: the casino player allows us to observe him one evening, as he changes dice and produces 10,000 rolls

Case 2. Estimation when the "right answer" is unknown

Examples:

GIVEN: the porcupine genome; we don't know how frequent are the CpG islands there, neither do we know their composition

GIVEN: 10,000 rolls of the casino player, but we don't see when he changes dice

QUESTION: Update the parameters θ of the model to maximize $P(x|\theta)$

Case 1. When the right answer is known

Given $x = x_1 \dots x_N$
for which the true $\pi = \pi_1 \dots \pi_N$ is known,

Define:

A_{kl} = # times $k \rightarrow l$ transition occurs in π
 $E_k(b)$ = # times state k in π emits b in x

We can show that the maximum likelihood parameters θ are:

$$a_{kl} = \frac{A_{kl}}{\sum_i A_{ki}} \quad e_k(b) = \frac{E_k(b)}{\sum_c E_k(c)}$$

Case 1. When the right answer is known

Intuition: When we know the underlying states, Best estimate is the average frequency of transitions & emissions that occur in the training data

Drawback:
Given little data, there may be **overfitting**:
 $P(x|\theta)$ is maximized, but θ is unreasonable
0 probabilities - VERY BAD

Example:
Given 10 casino rolls, we observe
 $x = 2, 1, 5, 6, 1, 2, 3, 6, 2, 3$
 $\pi = F, F, F, F, F, F, F, F, F, F$

Then:
 $a_{FF} = 1; \quad a_{F1} = 0$
 $e_1(1) = e_1(3) = .2,$
 $e_1(2) = .3; e_1(4) = 0; e_1(5) = e_1(6) = .1$

Pseudocounts

Solution for small training sets:

Add pseudocounts


$$A_{kl} = \text{# times } k \rightarrow l \text{ transition occurs in } \pi + r_{kl}$$

$$E_k(b) = \text{# times state } k \text{ in } \pi \text{ emits } b \text{ in } x + r_k(b)$$

$r_{kl}, r_k(b)$ are pseudocounts representing our prior belief

Larger pseudocounts \Rightarrow Strong prior belief

Small pseudocounts ($\epsilon < 1$): just to avoid 0 probabilities



Pseudocounts

Example: dishonest casino

We will observe player for one day, 500 rolls

Reasonable pseudocounts:

$$r_{0F} = r_{0L} = r_{F0} = r_{L0} = 1;$$

$$r_{FL} = r_{LF} = r_{FF} = r_{LL} = 1;$$


$$r_F(1) = r_F(2) = \dots = r_F(6) = 20 \quad (\text{strong belief fair is fair})$$

$$r_L(1) = r_L(2) = \dots = r_L(6) = 5 \quad (\text{wait and see for loaded})$$

Above #s pretty arbitrary – assigning priors is an art

6: Unsupervised learning

Estimate model parameters based on **unlabeled** training data




Learning case 2. When the right answer is unknown

We don't know the true $A_{kl}, E_k(b)$

Idea:

- We estimate our “best guess” on what $A_{kl}, E_k(b)$ are
- We update the parameters of the model, based on our guess
- We repeat



Case 2. When the right answer is unknown


Starting with our best guess of a model M, parameters θ :

Given $x = x_1 \dots x_N$
for which the true $\pi = \pi_1 \dots \pi_N$ is unknown,

We can get to a provably more likely parameter set θ

Principle: **EXPECTATION MAXIMIZATION**

1. Estimate $A_{kl}, E_k(b)$ in the training data
2. Update θ according to $A_{kl}, E_k(b)$
3. Repeat 1 & 2, until convergence



Estimating new parameters

To estimate A_{kl} :

At each position i of sequence x ,


Find probability transition $k \rightarrow l$ is used:

$$P(\pi_i = k, \pi_{i+1} = l | x) = [1/P(x)] \times P(\pi_i = k, \pi_{i+1} = l, x_1 \dots x_N) = Q/P(x)$$

where $Q = P(x_1 \dots x_i, \pi_i = k, \pi_{i+1} = l, x_{i+1} \dots x_N) =$
 $= P(\pi_{i+1} = l, x_{i+1} \dots x_N | \pi_i = k) P(x_1 \dots x_i, \pi_i = k) =$
 $= P(\pi_{i+1} = l, x_{i+1} x_{i+2} \dots x_N | \pi_i = k) f_k(i) =$
 $= P(x_{i+2} \dots x_N | \pi_{i+1} = l) P(x_{i+1} | \pi_{i+1} = l) P(\pi_{i+1} = l | \pi_i = k) f_k(i) =$
 $= b_l(i+1) e_l(x_{i+1}) a_{kl} f_k(i)$

So: $P(\pi_i = k, \pi_{i+1} = l | x, \theta) = \frac{f_k(i) a_{kl} e_l(x_{i+1}) b_l(i+1)}{P(x | \theta)}$

(For one such transition, at time step $i \rightarrow i+1$)



Estimating new parameters

(Sum over all $k \rightarrow l$ transitions, at any time step i)

So,

$$A_{kl} = \sum_i P(\pi_i = k, \pi_{i+1} = l | x, \theta) = \sum_i \frac{f_k(i) a_{kl} e_l(x_{i+1}) b_l(i+1)}{P(x | \theta)}$$

Similarly,

$$E_k(b) = [1/P(x)] \sum_{\{i | x_i = b\}} f_k(i) b_k(i)$$

Estimating new parameters

(Sum over all training seqs, all $k \rightarrow l$ transitions, all time steps i)

If we have several training sequences, x^1, \dots, x^M , each of length N ,

$$A_{kl} = \sum_x \sum_i P(\pi_i = k, \pi_{i+1} = l | x, \theta) = \sum_x \sum_i \frac{f_k(i) a_{kl} e_i(x_{i+1}) b_l(i+1)}{P(x | \theta)}$$

Similarly,

$$E_k(b) = \sum_x (1/P(x)) \sum_{\{i | x_i = b\}} f_k(i) b_k(i)$$

The Baum-Welch Algorithm

Initialization:
Pick the best-guess for model parameters (or arbitrary)

Iteration:

- Forward
- Backward
- Calculate $A_{kl}, E_k(b)$
- Calculate new model parameters $a_{kl}, e_k(b)$
- Calculate new log-likelihood $P(x | \theta)$

GUARANTEED TO BE HIGHER BY EXPECTATION-MAXIMIZATION

Until $P(x | \theta)$ does not change much

The Baum-Welch Algorithm – comments

Time Complexity:
iterations $\times O(K^2N)$

- Guaranteed to increase the log likelihood of the model

$$P(\theta | x) = P(x, \theta) / P(x) = P(x | \theta) / (P(x) P(\theta))$$

- Not guaranteed to find globally best parameters
- Converges to local optimum, depending on initial conditions
- Too many parameters / too large model: Overtraining

Alternative: Viterbi Training

Initialization: Same

Iteration:

- Perform Viterbi, to find π^*
- Calculate $A_{kl}, E_k(b)$ according to π^* + pseudocounts
- Calculate the new parameters $a_{kl}, e_k(b)$

Until convergence

Notes:

- Convergence is guaranteed – Why?
- Does not maximize $P(x | \theta)$
- In general, worse performance than Baum-Welch

	One path	All paths
Scoring	1. Scoring x , one path $P(x, \pi)$ Prob of a path, emissions	2. Scoring x , all paths $P(x) = \sum_{\pi} P(x, \pi)$ Prob of emissions, over all paths
	3. Viterbi decoding $\pi^* = \operatorname{argmax}_{\pi} P(x, \pi)$ Most likely path	4. Posterior decoding $\pi^* = \{\pi_i \pi_i = \operatorname{argmax}_k \sum_{\pi} P(\pi_i = k x)\}$ Path containing the most likely state at any time point.
Learning	5. Supervised learning, given π $\Lambda^* = \operatorname{argmax}_{\Lambda} P(x, \pi \Lambda)$	6. Unsupervised learning $\Lambda^* = \operatorname{argmax}_{\Lambda} \sum_{\pi} P(x, \pi \Lambda)$
	6. Unsupervised learning. $\Lambda^* = \operatorname{argmax}_{\Lambda} \max_{\pi} P(x, \pi \Lambda)$ Viterbi training, best path	Baum-Welch training, over all paths

What have we learned ?

- Generative model. Hidden states, observed emissions.
 - Generate a random sequence
 - Choose random transition, choose random emission (#0)
- Scoring: Finding the likelihood of a given sequence
 - Calculate likelihood of annotated path and sequence
 - Multiply emission and transition probabilities (#1)
 - Without specifying a path, total probability of generating x
 - Sum probabilities over all paths
 - Forward algorithm (#3)
- Decoding: Finding the most likely path, given a sequence
 - What is the most likely path generating entire sequence?
 - Viterbi algorithm (#2)
 - What is the most probable state at each time step?
 - Forward + backward algorithms, posterior decoding (#4)
- Learning: Estimating HMM parameters from training data
 - When state sequence is known
 - Simply compute maximum likelihood A and E (#5)
 - When state sequence is not known
 - Baum-Welch: Iterative estimation of all paths / frequencies (#6a)
 - Viterbi training: Iterative estimation of best path / frequencies (#6b)

The main questions on HMMs	
<p>1. Scoring x, one path = Joint probability of a sequence and a path, given the model</p> <ul style="list-style-type: none"> - GIVEN a HMM M, a path π, and a sequence x - FIND $\text{Prob}[x, \pi M]$ <p>→ "Running the model", simply multiply emission and transition probabilities → Application: "all fair" vs. "all loaded" comparisons</p>	SCORING
<p>2. Scoring x, all paths = total probability of a sequence, summed across all paths</p> <ul style="list-style-type: none"> - GIVEN a HMM M, a sequence x - FIND the total probability $P(x M)$ summed across all paths <p>→ Forward algorithm, sum score over all paths (same result as backward)</p>	PARSING
<p>3. Viterbi decoding = parsing a sequence into the optimal series of hidden states</p> <ul style="list-style-type: none"> - GIVEN a HMM M, and a sequence x, - FIND the sequence π^* of states that maximizes $P[x, \pi M]$ <p>→ Viterbi algorithm, dynamic programming, max score over all paths, trace pointers find path</p> <p>4. Posterior decoding = total prob that emission x_t came from state k, across all paths</p> <ul style="list-style-type: none"> - GIVEN a HMM M, a sequence x - FIND the total probability $P[\pi_t = k x, M]$ <p>→ Posterior decoding: run forward & backward algorithms to & from state $\pi_t = k$</p>	
<p>5. Supervised learning = optimize parameters of a model given training data</p> <ul style="list-style-type: none"> - GIVEN a HMM M, with unspecified transition/emission probs., labeled sequence x, - FIND parameters $\theta = (e, a_t)$ that maximize $P[x \theta]$ <p>→ Simply count frequency of each emission and transition observed in the training data</p> <p>6. Unsupervised learning = optimize parameters of a model given training data</p> <ul style="list-style-type: none"> - GIVEN a HMM M, with unspecified transition/emission probs., unlabeled sequence x, - FIND parameters $\theta = (e, a_t)$ that maximize $P[x \theta]$ <p>→ Viterbi training: guess parameters, find optimal Viterbi path (#2), update parameters (#5), iterate → Baum-Welch training: guess, sum over all emissions/transitions (#4), update (#5), iterate</p>	LEARNING

	One path	All paths
Scoring	<p>1. Scoring x, one path</p> $P(x, \pi)$ <p>Prob of a path, emissions</p>	<p>2. Scoring x, all paths</p> $P(x) = \sum_{\pi} P(x, \pi)$ <p>Prob of emissions, over all paths</p>
Decoding	<p>3. Viterbi decoding</p> $\pi^* = \text{argmax}_{\pi} P(x, \pi)$ <p>Most likely path</p>	<p>4. Posterior decoding</p> $\pi^{\wedge} = \{\pi_i \pi_i = \text{argmax}_k \sum_{\pi} P(\pi_i = k x)\}$ <p>Path containing the most likely state at any time point.</p>