

# Training HMMs

Goal: Given: • Long sequence of observations  $X = x_1 \dots x_L$   
• Structure of an HMM: Set of states  $S$   
Set of possible trans.  $T$

Find: { Emission probabilities :  $E$   
Transition prob. :  $T$   
Initial state prob :  $I$  }

$$Pr[A] = \sum_b Pr[A, B=b]$$

→ Such that  $(E, T, I)$  best represent the observed seq.  $X$

Maximum Likelihood estimation

←  $Pr[X = x_1 \dots x_L | E, T, I]$  is maximized

$$Pr[X = x_1 \dots x_L | E, T, I] = \sum_{\text{path } P = p_1 \dots p_L} Pr[X, P | E, T, I]$$

↓ We know how to calculate

Can be calculated in  $O(L \cdot n^3)$  using Forward algorithm

~~Class~~

Simplified version:

Given:  $X = x_1 \dots x_L$

$P = p_1 \dots p_L \leftarrow$  annotation of  $X$

Find:  $E, T, I$  s.t.  $P_r[X, P | E, T, I]$  is max

Let  $N_e(s, s') =$  # of times transition happened  
btw  $s$  and  $s'$

$=$  # of positions  $i$  s.t.  $p_i = s \wedge p_{i+1} = s'$

Then, choose  $T(s, s') = \frac{N_e(s, s')}{\sum_{x \in S} N_e(s, x)}$

$\rightarrow$  total # of  
times in states

Example:  $X = \text{ACA} \text{ TAC} \text{ TGA} \text{ CTG}$   
 $P = \text{GGG} \text{ IIII} \text{ IGGG} \text{ II}$

$$T(G, I) = \frac{2}{7}$$

$$T(G, G) = \frac{5}{7}$$

For emissions:  $N_e(s, a) =$  # times  $a$  was emitted from  $s$   
 $=$  # positions  $i$  s.t.  $p_i = s \wedge x_i = a$

$$\text{Choose } E(s, a) = \frac{N_e(s, a)}{\sum_b N_e(s, b)}$$

$$\text{Example: } E(G, A) = \frac{4}{7}$$



Back to problem where  $P$  is not given

## ① Viterbi training

① Choose "reasonable"  $E, T, I$  (either from prior knowledge or randomly)

Repeat until convergence

② Use Viterbi algo to find best path for  $X = x_1 \dots x_n$  assuming  $E, T, I \Rightarrow$  Produces path  $P$

③ Re-estimate  $E, T, I$  from  $P, X$  as done previously

Problem: Algo. frequently gets stuck in local optima

## ② Baum-Welsh algorithm

Bases re-estimation of  $E, T, I$  on all paths (weighted by their probability)

$\Rightarrow$  Reduces the prob of getting stuck in local optima