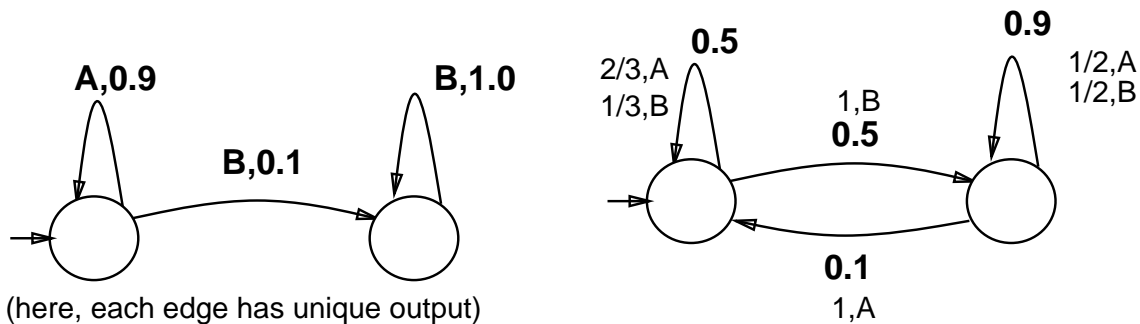


# Hidden Markov Models

Process for generating string of characters.



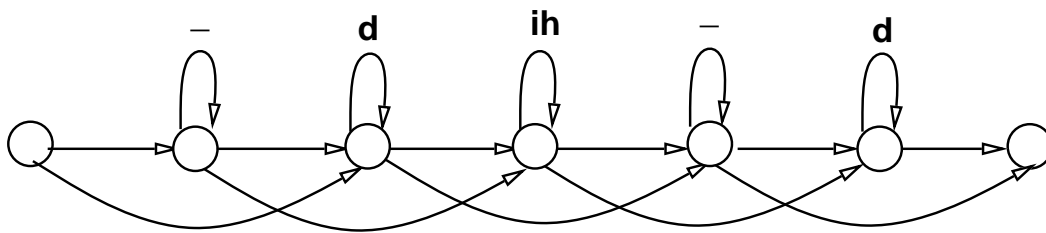
- Set  $S$  of states, one marked as initial. (Sometimes also one marked as final)
- Transition probabilities:  $a_{ij}$  = prob go to state  $j$  given that in state  $i$ .  
(For each  $i$ ,  $\sum_j a_{ij} = 1$ .)
- Output probabilities:  $b_{ij}(X)$  = prob of outputting  $X$  given that on arc from  $i$  to  $j$ .

Typical uses:

- model the result of taking an action in an environment.
- model a phoneme, word, or sentence.

# Examples

---



HMM for 'did'

## Typical HMM questions

---

- “Evaluation Problem” : Find  $P(\text{string} \mid \text{model})$ .

E.g., say you have HMMs for each of 50 words. You hear some sounds and want the most likely word.

$$P(\text{word} \mid \text{sounds}) \propto P(\text{sounds} \mid \text{word}) \times P(\text{word}).$$

- “Decoding problem” : Find most likely path given string, model.
- “Learning Problem” : Given a model without the probabilities and an observation, what is the best setting of the probabilities to maximize the likelihood of the observation:  $P(\text{observation} \mid \text{model})$ .

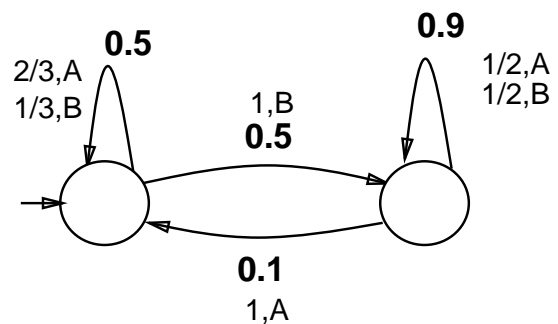
## Evaluation problem: how to calculate $P(\text{observation} \mid \text{model})$ ?

---

One way: for a given model,

$$P(\text{obs}) = \sum_{\text{paths } S \text{ of length } |\text{obs}|} P(S) \cdot P(\text{obs} \mid S)$$

E.g., “BBA” for model:



One path is  $S = \text{stay, stay, stay}$ .

$$P(S) = (0.5)^3, P(\text{obs} \mid S) = (1/3)(1/3)(2/3).$$

Another?

**Problem:** Could be exponentially many paths.

# Evaluation problem: how to calculate $P(\text{observation} \mid \text{model})$ ?

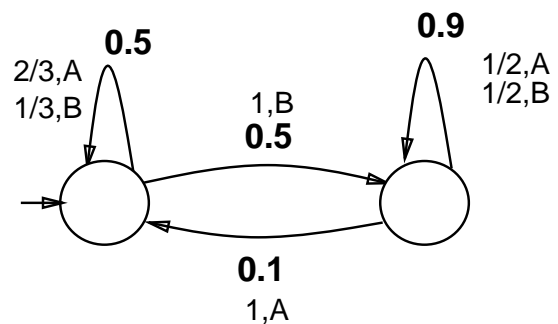
---

Idea: Use dynamic programming approach.

- Calculate  $P(\text{at time } t \text{ we're in state } i \text{ and have produced first } t \text{ characters of observation} \mid \text{model})$ .  
- called  $\alpha_i(t)$ .
- $\alpha_i(0)$  is initial distribution.
- $\alpha_j(t) = \sum_i [P(\text{in state } i \text{ at time } t - 1 \text{ and have produced first } t - 1 \text{ chars}) \times P(\text{then go from } i \text{ to } j \text{ and output } t\text{th char})]$

$$= \sum_i \alpha_i(t - 1) a_{ij} b_{ij}(\text{obs}_t).$$

- Example: "BBA" for



- Final prob =  $\sum_i \alpha_i(|\text{obs}|)$ .

## Learning problem (parameter estimation)

Goal: Try to find probabilities that maximize likelihood of observation.

(Could be lots of independent runs, or one big run, or combination)

**Method 1:** Gradient descent. (Not used much for HMMs)

- Start with some guess and compute  $P(\text{obs} \mid \text{model})$ .
- Try tweaking values to see if increases or decreases, and hill-climb.
- Problem: slow convergence. Hard to get exact derivatives.

## Learning problem contd.

Goal: Try to find probabilities that maximize likelihood of observation.

### **Method 2:** Baum-Welsh/Forward-backward/EM

- Start with some guess of model.
- E:** Use model & observation to calculate expected number of traversals across each arc.  
(and expected # traversals across each arc while producing a given character)
- M:** Use calculation to update model to model most likely to produce these ratios.

## Baum-Welsh, more detail

---

**E:** Want to calculate for given model:

- \*  $E_{ij} = E[\# \text{ times cross } i \rightarrow j \mid \text{obs}]$ .

- \*  $E_{ij}(X) = E[\# \text{ times cross } i \rightarrow j \text{ while producing letter } X \mid \text{obs}]$ .

**Idea:** calculate  $P[\text{cross } i \rightarrow j \text{ at time } t \mid \text{obs}]$  and add.

(Add over all  $t$  for  $E_{ij}$ . Add over  $t$  such that  $obs_t = X$  to get  $E_{ij}(X)$ .)

Details on next slide...

**M:** Find what  $a_{ij}$ ,  $b_{ij}(X)$  are most likely to produce these expectations.

- \*  $\bar{a}_{ij} = E_{ij} / \sum_k E_{ik}$ .

- \*  $\bar{b}_{ij}(X) = E_{ij}(X) / E_{ij}$ .



## Baum-Welsh, contd.

---

Details of E step: (high level idea more important than details)

To calculate  $P[\text{cross } i \rightarrow j \text{ at time } t \mid \text{obs}]$ ,

- Rewrite as:  $\frac{P[\text{cross } i \rightarrow j \text{ at time } t \ \& \ \text{produce obs}]}{P[\text{produce obs}]}$
- Numerator =  $P[\text{at state } i \text{ at time } t \text{ and output first } t \text{ chars}] \times a_{ij} \times b_{ij}(\text{obs}_t) \times P[\text{output remaining chars} \mid \text{at state } j \text{ at time } t + 1]$ .

**Theorem:** This algorithm will converge to local optimum.