

CISC 636 Computational Biology & Bioinformatics

(Fall 2016)

Hidden Markov Models (III)

- Viterbi training
- Baum-Welch algorithm
- Maximum Likelihood
- Expectation Maximization

Model building

- Topology
 - Requires domain knowledge
- Parameters
 - When states are labeled for sequences of observables
 - Simple counting (Maximum Likelihood):
$$a_{kl} = A_{kl} / \sum_l A_{kl}, \text{ and } e_k(b) = E_k(b) / \sum_b E_k(b')$$

- When states are not labeled

Method 1 (Viterbi training)

1. Assign random parameters
2. Use Viterbi algorithm for labeling/decoding
2. Do counting to collect new a_{kl} and $e_k(b)$;
3. Repeat steps 2 and 3 until stopping criterion is met.

Method 2 (Baum-Welch algorithm)

Baum-Welch algorithm (Expectation-Maximization)

- An iterative procedure similar to Viterbi training
- Probability that a_{kl} is used at position i in sequence j .

$$P(\pi_i = k, \pi_{i+1} = l \mid \mathbf{x}, \theta) = f_k(i) a_{kl} e_l(\mathbf{x}_{i+1}) b_l(i+1) / P(\mathbf{x}^j)$$

Calculate the expected number of times that is used by summing over all position and over all training sequences.

$$A_{kl} = \sum_j \{ (1/P(\mathbf{x}^j)) [\sum_i f_k^j(i) a_{kl} e_l(\mathbf{x}_{i+1}^j) b_l^j(i+1)] \}$$

Similarly, calculate the expected number of times that symbol b is emitted in state k .

$$E_k(b) = \sum_j \{ (1/P(\mathbf{x}^j)) [\sum_{\{i \mid \mathbf{x}_i^j = b\}} f_k^j(i) b_k^j(i)] \}$$

Maximum Likelihood

Define $L(\theta) = P(x | \theta)$

Estimate θ such that the distribution with the estimated θ best agrees with or support the data observed so far.

$$\theta^{\text{ML}} = \underset{\theta}{\operatorname{argmax}} L(\theta)$$

E.g. There are red and black balls in a box. What is the probability P of picking up a black ball?

Do sampling (with replacement).

Maximum Likelihood

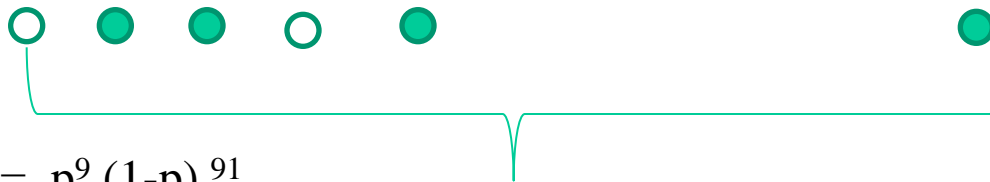
Define $L(\theta) = P(x | \theta)$

Estimate such that the distribution with the estimated best agrees with or supports the data observed so far.

$$\theta^{ML} = \operatorname{argmax}_{\theta} L(\theta)$$

When $L(\theta)$ is differentiable,
$$\frac{\partial L(\theta)}{\partial \theta} \Big|_{\theta^{ML}} = 0$$

For example, want to know the ratio: # of blackball/# of whiteball, in other words, the probability P of picking up a black ball. Sampling (with replacement):



$$\text{Prob (iid) } = p^9 (1-p)^{91}$$

$$\text{Likelihood } L(p) = p^9 (1-p)^{91}.$$

100
times

Counts: whiteball 91,
blackball 9

$$\frac{\partial L(p)}{\partial p} = 9p^8(1-p)^{91} - 91p^9(1-p)^{90} = 0$$

$\Rightarrow P^{ML} = 9/100 = 9\%$. The ML estimate of P is just the frequency.

A proof that the observed frequency \rightarrow ML estimate of probabilities for polynomial distribution

Let Counts n_i for outcome i

The observed frequencies $\theta_i = n_i / N$, where $N = \sum_i n_i$

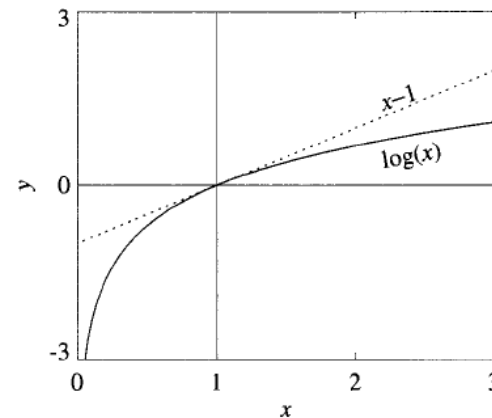
If $\theta_i^{ML} = n_i / N$, then $P(n | \theta^{ML}) > p(n | \theta)$ for any $\theta \neq \theta^{ML}$

Proof:

$$\log \frac{P(n | \theta^{ML})}{P(n | \theta)} = \log \frac{\prod_i (\theta_i^{ML})^{n_i}}{\prod_i (\theta_i)^{n_i}} = \log \prod_i \left(\frac{\theta_i^{ML}}{\theta_i} \right)^{n_i}$$

$$= \sum_i n_i \log \left(\frac{\theta_i^{ML}}{\theta_i} \right) = N \sum_i \frac{n_i}{N} \log \left(\frac{\theta_i^{ML}}{\theta_i} \right) = \sum_i \theta_i^{ML} \log \left(\frac{\theta_i^{ML}}{\theta_i} \right)$$

$$= H(\theta^{ML} \| \theta) \geq 0$$



Maximum Likelihood: pros and cons

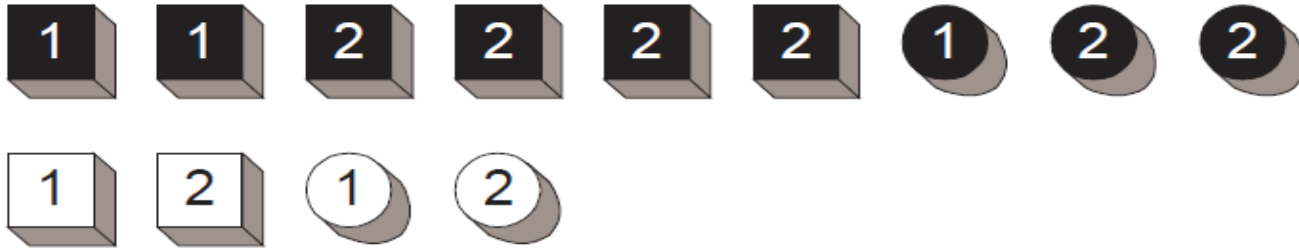
- Consistent, i.e., in the limit of a large amount of data, ML estimate converges to the true parameters by which the data are created.
- Simple
- Poor estimate when data are insufficient.
e.g., if you roll a die for less than 6 times, the ML estimate for some numbers would be zero.

Pseudo counts:

$$\theta_i = \frac{n_i + \alpha_i}{N + A},$$

where $A = \sum_i \alpha_i$

Conditional Probability and Joint Probability



$$P(\text{one}) = 5/13$$

$$P(\text{square}) = 8/13$$

$$P(\text{one, square}) = 3/13$$

$$P(\text{one} \mid \text{square}) = 3/8 = P(\text{one, square}) / P(\text{square})$$

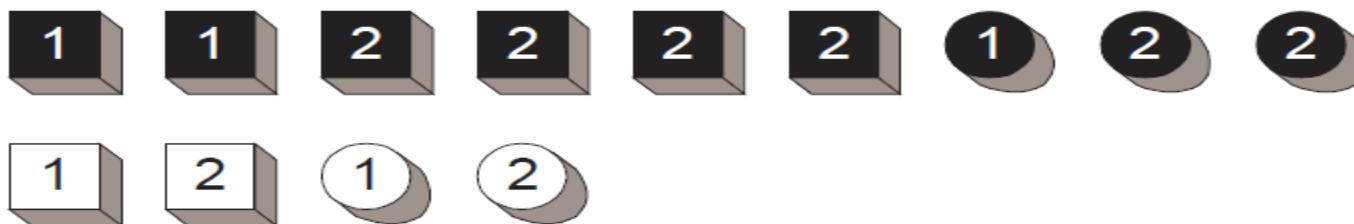
In general, $P(D, M) = P(D \mid M)P(M) = P(M \mid D)P(D)$

=> **Baye's Rule:**

$$P(M \mid D) = \frac{P(D \mid M)P(M)}{P(D)}$$

$$\begin{aligned} P(\text{One} \mid \text{Black}) &= \frac{P(\text{Black} \mid \text{One})P(\text{One})}{P(\text{Black} \mid \text{One})P(\text{One}) + P(\text{Black} \mid \text{Two})P(\text{Two})} \\ &= \frac{\left(\frac{3}{5}\right)\left(\frac{5}{13}\right)}{\left(\frac{3}{5}\right)\left(\frac{5}{13}\right) + \left(\frac{6}{8}\right)\left(\frac{8}{13}\right)} = \frac{1}{3}, \end{aligned}$$

Conditional Probability and Conditional Independence



$$P(\text{One}) = \frac{5}{13}$$

$$P(\text{One}|\text{Square}) = \frac{3}{8}$$

$$P(\text{One}|\text{Black}) = \frac{3}{9} = \frac{1}{3}$$

$$P(\text{One}|\text{Square} \cap \text{Black}) = \frac{2}{6} = \frac{1}{3}$$

$$P(\text{One}|\text{White}) = \frac{2}{4} = \frac{1}{2}$$

$$P(\text{One}|\text{Square} \cap \text{White}) = \frac{1}{2}$$

So One and Square are not independent, but they are conditionally independent given Black and given White.

Baye's Rule:

$$P(M | D) = \frac{P(D | M)P(M)}{P(D)}$$

Example: disease diagnosis/inference

$P(\text{Leukemia} | \text{Fever}) = ?$

$P(\text{Fever} | \text{Leukemia}) = 0.85$

$P(\text{Fever}) = 0.9$

$P(\text{Leukemia}) = 0.005$

$P(\text{Leukemia} | \text{Fever}) = P(F | L)P(L)/P(F) = 0.85 * 0.01 / 0.9 = 0.0047$

Bayesian Inference

Maximum a posterior estimate

$$\theta^{MAP} = \arg \max_{\theta} P(\theta | \mathbf{x})$$

Expectation Maximization

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

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

$$\log P(x | \theta) = \log P(x, y | \theta) - \log P(y | x, \theta)$$

$$\sum_y P(y | x, \theta^t) \left(\log P(x | \theta) = \sum_y P(y | x, \theta^t) \log P(x, y | \theta) - \sum_y P(y | x, \theta^t) \log P(y | x, \theta) \right) \quad \textbf{Expectation}$$

$$\log P(x | \theta) = \sum_y P(y | x, \theta^t) \log P(x, y | \theta) - \sum_y P(y | x, \theta^t) \log P(y | x, \theta)$$

$$Q(\theta | \theta^t) = \sum_y P(y | x, \theta^t) \log P(x, y | \theta)$$

$$\log P(x | \theta) - \log P(x | \theta^t)$$

$$= Q(\theta | \theta^t) - Q(\theta^t | \theta^t) + \sum_y P(y | x, \theta^t) \log \frac{P(y | x, \theta^t)}{P(y | x, \theta)}$$

$$\geq Q(\theta | \theta^t) - Q(\theta^t | \theta^t)$$

$$\theta^{t+1} = \arg \max_{\theta} Q(\theta | \theta^t)$$

Maximization

EM explanation of the Baum-Welch algorithm

We like to
maximize by
choosing θ

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

But state path π is
hidden variable. Thus,
EM.

$$Q(\theta | \theta^t) = \sum_{\pi} P(\pi | x, \theta^t) \log P(x, \pi | \theta)$$

$$P(x, \pi | \theta) = \prod_{k=1}^M \prod_b [e_k(b)]^{E_k(b, \pi)} \prod_{k=0}^M \prod_{l=1}^M a_{kl}^{A_{kl}(\pi)},$$

$$Q(\theta | \theta^t) = \sum_{\pi} P(\pi | x, \theta^t) \times \left[\sum_{k=1}^M \sum_b E_k(b, \pi) \log e_k(b) + \sum_{k=0}^M \sum_{l=1}^M A_{kl}(\pi) \log a_{kl} \right].$$

EM Explanation of the Baum-Welch algorithm

$$E_k(b) = \sum_{\pi} P(\pi | x, \theta^t) E_k(b, \pi) \quad \text{and} \quad A_{kl} = \sum_{\pi} P(\pi | x, \theta^t) A_{kl}(\pi).$$

$$Q(\theta | \theta^t) = \sum_{k=1}^M \sum_b E_k(b) \log e_k(b) + \sum_{k=0}^M \sum_{l=1}^M A_{kl} \log a_{kl}.$$

E-term

A-term

A-term is maximized if

$$a_{kl}^{EM} = \frac{A_{kl}}{\sum_{l'} A_{kl'}}$$

E-term is maximized if

$$e_k^{EM}(b) = \frac{\sum_{l'} E_k(b)}{\sum_{b'} E_k(b')}$$