







- **Evaluation**: given *x*, what is the probability *p*(*x*) that it was produced by the model?
- *Decoding*: given *x*, what is the most probable path that produces *x* in the model?
- **Learning**: given *x*, what are the parameters (transitional probabilities and emission probabilities) of the model that maximize p(x).





Computing p(x)

- Before, $p(x) = a_{sx1} \prod_{i=2...n} a_{x^{i}x^{i}}$
- Now, $p(x) = \Sigma_{\pi} p(x, \pi)$
- Enumerating all π is exponential!
- Use Viterbi, same as before, but change \mbox{max} to $\pmb{\Sigma}$







Forward evaluation algorithm

- Initialization $f_0(0) = 1, f_k(0) = 0$ for k > 0
- Main iteration for i = 1...n $f_i(i) = e_i(x_i).\Sigma_k (f_k(i-1).a_{kl})$
- Termination $p(x) = \sum_k f_k(n) a_{k0}$



Problem of small numbers

- In Viterbi and Forward algorithm we multiply probabilities → numbers will soon be very small and we loose precision.
- Use log space → addition instead of multiplication.



Log space Viterbi

 $v_{i}(i) = e_{i}(x_{i}).max_{k}(v_{k}(i-1).a_{k})$

Let $V_i(i) = \log v_i(i)$

```
 \begin{array}{ll} V_k(i) & = \log \left[ e_k(x_i) . \max_k \left( v_k(i-1) . a_{ki} \right) \right] \\ & = \log \left[ e_i(x_i) + \log \left[ \max_k \left( v_k(i-1) . a_{ki} \right) \right] \\ & = \log \left[ e_i(x_i) + \max_k \log \left[ \left( v_k(i-1) . a_{ki} \right) \right] \\ & = \log \left[ e_i(x_i) + \max_k \left( V_k(i-1) + \log a_{ki} \right) \right] \end{array}
```

Log space Forward

 $f_i(i) = e_i(x_i) \cdot \Sigma_k \left(f_k(i-1) \cdot a_{kl} \right)$

Let $F_i(i) = \log f_i(i)$

$$\begin{array}{l} F_{k}(i) & = \log \left[e_{k}(x_{i}).\Sigma_{k} \left(f_{k}(i-1).a_{kl} \right) \right] \\ & = \log \left[e_{k}(x_{i}) + \log \Sigma_{k} \left(f_{k}(i-1).a_{kl} \right) \right] \\ & \neq \log \left[e_{k}(x_{i}) + \Sigma_{k} \log \left[\left(f_{k}(i-1).a_{kl} \right) \right] \end{array}$$

$$= \log e_i(x_i) + \log \sum_k e^{(Fk(i-1) + \log akl)}$$

Back to the most probable path

- The Viterbi algorithm finds it!
- The most probable path might not be the most appropriate basis for judgment.
- We might want, for instance, the most probable state for an observation *x*_{*r*}
- More generally, we are interested in $p(\pi_i = k \mid x)$



Computing $p(\pi_i = k \mid x)$

- $p(\pi_i = k \mid x) = p(x, \pi_i = k)/p(x)$
- I know how to compute *p*(*x*): forward alg.
- $p(x, \pi_i = k)$
 - $= p(x_1...x_i, \pi_i = k).p(x_{i+1}...x_n \mid x_1...x_i, \pi_i = k)$ $= p(x_1...x_i, \pi_i = k).p(x_{i+1}...x_n \mid \pi_i = k)$
 - $= p(\mathbf{x}_1 \dots \mathbf{x}_i, \mathbf{x}_i \mathbf{x}) \cdot p(\mathbf{x}_{i+1} \dots \mathbf{x}_n + \mathbf{x}_i)$ $= f_k(\mathbf{i}) \cdot b_k(\mathbf{i})$

Backward evaluation algorithm • Initialization $b_k(n) = a_{k0}$ for all k• Main iteration for $i = n - 1 \dots 1$ $b_k(i) = \Sigma_l a_{kl} e_l(x_{i+1}) b_l(i+1)$ • Termination $p(x) = \Sigma_l a_{0l} e_l(x_1) b_l(1)$

Learning (training the HMM)

- Let θ be the parameters of the HMM (transition probabilities and emission probabilities, the a's and e's)
- Given independent sequences x¹, ..., xⁿ, we would like to find θ that will maximize:

 $\log p(x^1 \dots x^n \mid \theta) = \sum_{j=1\dots n} \log p(x^j \mid \theta)$

This is called the maximum likelihood parameters.



State sequence is known

- Assume the path for each xi is known
 For instance, we have sequences in which CpG islands are already labeled
- Paths are known, let
 A_{kl} = number of transitions from k to l
 - $E_k(b) =$ number of times b emitted in state k
- The maximum likelihood parameters are given by:
 a_{kl} = A_{kl} / Σ_tA_{kl},
 e_k(b) = E_k(b) / Σ_tE_k(b)

Maximum likelihood from counts

- Assume we have a sequence of independent observations x₁...x_n and that we count n_i occurrences of outcome i, i=1...k.
- Let θ_i = probability of *i*.
- Then $\theta^{ML} = \{\theta_i = n/n, i=1...k\}$ is the maximum likelihood solution for θ .
- Consider any other θ . We want to show that $p(x \mid \theta^{ML}) > p(x \mid \theta)$



Proof

$$\log \frac{p(x \mid \theta^{ML})}{p(x \mid \theta)} = \log \frac{\prod_{i} (\theta_{i}^{ML})^{n_{i}}}{\prod_{i} \theta_{i}^{n_{i}}}$$
$$= \sum_{i} n_{i} \log \frac{\theta_{i}^{ML}}{\theta_{i}}$$
$$= n \sum_{i} \theta_{i}^{ML} \log \frac{\theta_{i}^{ML}}{\theta_{i}} > 0$$

The last summation is the relative entropy of θ^{ML} and θ which is always positive and 0 iff $\theta^{\text{ML}} = \theta \text{ (from information theory)}$

Some problems

- Maximum likelihood are vulnerable to overfitting if insufficient data.
- For instance, if a state *k* was never used in the set of training sequences, then $-a_{kl} = 0$ for all *l* $-e_k(b) = 0$ for all *b* .
- To avoid such problem, start with pseudocounts of r_{kl} for A_{kl} and $r_k(b)$ for $E_k(b)$. •
- Large pseudocount indicates strong prior belief about the probabilities (will require more data to modify) ٠
- Small pseudocount just to avoid zero probability

Example

Dishonest Casino HMM

 $r_{0F} = r_{0L} = r_{F0} = r_{L0} = 1;$ [avoid zero probability] $r_{\rm FL} = r_{\rm LF} = r_{\rm FF} = r_{\rm LL} = 1; \quad [{\rm avoid \ zero \ probability}]$

 $r_{\rm F}(1) = r_{\rm F}(2) = \ldots = r_{\rm F}(6) = 20$ [strong belief that fair is fair]

 $r_{\rm L}(1) = r_{\rm L}(2) = \dots = r_{\rm L}(6) = 5$ [wait and see for loaded]

New species comes in...

- New species with different distribution of CpG islands.
- We do not have labeled genomic sequences for the new species.
- Need to find maximum likelihood θ of HMM without knowing the paths!



Baum–Welsh algorithm

 $\begin{array}{l} \text{start at iteration 0 with some } \theta, \text{ call it } \theta^0 \\ L^0 \leftarrow \log p(x^1...x^n \mid \theta^0) \\ i \leftarrow 0 \\ \hline \textbf{repeat} \\ i \leftarrow i + 1 \\ A_{ki}^{-i} \leftarrow \mathbb{E} \left[A_{ki} \mid x^1...x^n, \theta^{i\cdot1} \right] \quad (\text{expected value}) \\ E_k(b)^i \leftarrow \mathbb{E} \left[E_k(b) \mid x^1...x^n, \theta^{i\cdot1} \right] \quad (\text{expected value}) \\ \text{calculate } \theta^i \text{ using maximum likelihood estimators} \\ \text{from counts } A_{ki}^{-i} \text{ and } E_k(b)^i \text{ as before.} \\ L^i \leftarrow \log p(x^1...x^n \mid \theta) \quad (\text{new likelihood}) \\ \end{array}$

until $L^i - L^{i-1} < \text{threshold}$



Saad Mns

What is the guarantee?

- Baum–Welsh algorithm is a special case of a general algorithm known as Expectation Maximization (EM)
- EM guarantees that $p(X | \theta^{i+1}) \ge p(X | \theta^{i})$
- It will therefore converge to a local maximum (not necessarily the maximum)

We need to...

Compute:

 $- \mathsf{E} \left[A_{kl} \mid x^1 \dots x^n, \theta \right]$

 $- \mathsf{E} \left[E_k(b) \mid x^1 \dots x^n, \theta \right]$



$\mathsf{E}\left[A_{kl} \mid x^{1} \dots x^{n}, \theta\right]$

By linearity of expectation: E $[A_{kl} | x^1...x^n, \theta] = \sum_j E [A_{kl} / | x^j, \theta]$

By linearity of expectation, again: $E [A_{kl}^{i} | x^{i}, \theta] = \sum_{i} E[\# \text{ of } k \rightarrow l \text{ at } x_{i}^{i} | x^{i}, \theta]$ $= \sum_{i} p(k \rightarrow l \text{ at } x_{i}^{i} | x^{i}, \theta)$

 $p(k \rightarrow l \text{ at } x_i^j \mid x^j, \theta) = p(\pi_i = k, \pi_{i+1} = l \mid x^j, \theta)$

 $p(\pi_{i} = k, \pi_{i+1} = l \mid x^{i}, \theta) = p(\pi_{i} = k, \pi_{i+1} = l, x^{i} \mid \theta) / p(x^{i} \mid \theta)$ = $f_{k}^{i}(l)a_{k}e^{l}(x_{i+1})b^{i}(i+1) / p(x^{i} \mid \theta)$



Viterbi training

start at iteration 0 with some $\theta,$ call it θ^0 i \leftarrow 0

repeat *i* ← *i* + 1

 $A_{k^{j}} \leftarrow \text{number of transitions } k \rightarrow l \text{ on the most probable paths } \pi^{1\star}, \, ..., \, \pi^{j\star}$

 $E_k(b)^i \leftarrow$ number of times k emits b on the most probable paths $\pi^{1*}, ..., \pi^{i*}$

calculate θ^i using maximum likelihood estimators from counts $A_{kj}^{\ i}$ and $E_k(b)^i$ as before.

until none of the optimal paths change



Saad Mne

What is the guarantee

- It will converge
- It will not necessarily maximize the true likelihood $p(x_1...x^n \mid \theta)$, but $p(x_1...x^n \mid \theta, \pi^{1*}, ..., \pi^{i*})$
- Usually performs less well than Baum–Welsh
- Practical, don't have to perform Forward and Backward algorithms, only Viterbi!
- Makes sense if we are using only Viterbi decoding