

## Sum-product for E step in the HMM learning problem (review)

Given all  $\alpha(z_n)$  and  $\beta(z_n)$

$$\gamma(z_n) = \frac{\alpha(z_n)\beta(z_n)}{p(X)}$$

$$p(X) = \sum_{z_n} \alpha(z_n)\beta(z_n)$$

$$\xi(z_{n-1}, z_n) = \frac{\alpha(z_{n-1})p(x_n|z_n)p(z_n|z_{n-1})\beta(z_n)}{p(X)}$$

## Rescaled alpha beta (Bishop, 13.2.4)

The alpha-beta algorithm has similar precision problems to the ones for EM where we discussed the fix of scaling log quantities by the max, before exponentiation for normalizing.

One way to handle this is to reformulate the alpha-beta algorithm in terms of:

$$\hat{\alpha}(z_n) = p(z_n | x_1, \dots, x_n) = \frac{\alpha(z_n)}{p(x_1, \dots, x_n)}$$

$$\hat{\beta}(z_n) = \frac{p(x_{n+1}, \dots, x_N | z_n)}{p(x_{n+1}, \dots, x_N | x_1, \dots, x_n)} = \frac{\beta(z_n)}{p(x_{n+1}, \dots, x_N | x_1, \dots, x_n)}$$

## Rescaled alpha beta (Bishop, 13.2.4)

The alpha-beta algorithm has similar precision problems to the ones for EM where we discussed the fix of scaling log quantities by the max, before exponentiation for normalizing.

One way to handle this is to reformulate the alpha-beta algorithm in terms of:

$$\hat{\alpha}(z_n) = p(z_n | x_1, \dots, x_n) = \frac{\alpha(z_n)}{p(x_1, \dots, x_n)}$$

Pieces of  $p(X)$

$$\hat{\beta}(z_n) = \frac{p(x_{n+1}, \dots, x_N | z_n)}{p(x_{n+1}, \dots, x_N | x_1, \dots, x_n)} = \frac{\beta(z_n)}{p(x_{n+1}, \dots, x_N | x_1, \dots, x_n)}$$

## Rescaled alpha beta (Bishop, 13.2.4)

$$\hat{\alpha}(z_n) = p(z_n | x_1, \dots, x_n) = \frac{\alpha(z_n)}{p(x_1, \dots, x_n)}$$

Let  $c_n = p(x_n | x_1, \dots, x_{n-1})$

and note that  $p(x_1, \dots, x_n) = \prod_{m=1}^n c_m$ . Then

$$c_n \hat{\alpha}(z_n) = p(x_n | z_n) \sum_{z_{n-1}} \hat{\alpha}(z_{n-1}) p(z_n | z_{n-1})$$

and we get  $c_n$  as the normalizer of the RHS.

(See Bishop for the betas).

## Classic HMM computational problems

Given data, what is the HMM (**learning**). ✓

Given an HMM, what is the **distribution over the state** variables. Also, **how likely** are the observations, given the model. ✓

Given an HMM, what is the most likely **state sequence** for some data?

## Viterbi algorithm (special case of max-sum)

### Recall max-sum

Forward direction is like sum-product, except

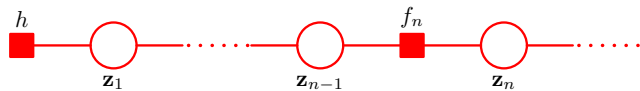
We take the max instead of sum

We use sum of logs instead of product

We remember incoming variable values that give max (\*)

Backwards direction is simply backtracking on (\*).

## Recall simplified factor graph

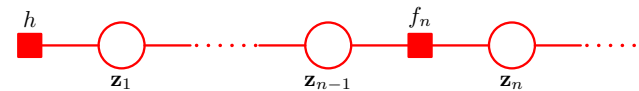


$$h = p(z_1)p(x_1|z_1) \quad f_n = p(z_n|z_{n-1})p(x_n|z_n)$$

Left to right messages

$$\omega(z_n) \equiv \mu_{f_n \rightarrow z_n}(z_n) = ?$$

## Recall simplified factor graph



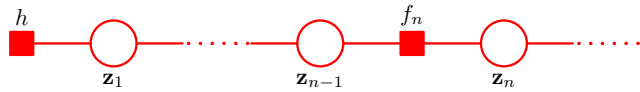
$$h = p(z_1)p(x_1|z_1) \quad f_n = p(z_n|z_{n-1})p(x_n|z_n)$$

Left to right messages

$$\omega(z_n) = \log(x_n|z_n) + \max_{z_{n-1}} \left\{ \log(p(z_n|z_{n-1}) + \omega(z_{n-1})) \right\}$$

$$\omega(z_1) = \log(p(z_1)) + \log(p(x_1|z_1))$$

## Intuitive understanding



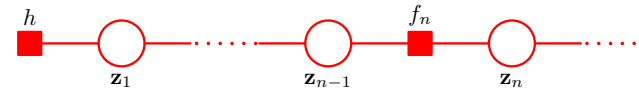
$$\omega(z_n) = \log(x_n | z_n) + \max_{z_{n-1}} \left\{ \log(p(z_n | z_{n-1}) + \omega(z_{n-1})) \right\}$$

Consider all possible paths to each of the  $k$  states for time  $n$ .

The message encodes the probabilities for the maximum probability path for each of the  $K$  states.

EG, if you are in state  $k$ , this records is the probability of being there by via the maximal probably sequence.

## Intuitive understanding



The message is the vector of probabilities for the maximum probability path for each of the  $K$  states.

$$\omega(z_n) = \log(x_n | z_n) + \max_{z_{n-1}} \left\{ \log(p(z_n | z_{n-1}) + \omega(z_{n-1})) \right\}$$

For each state  $k$

Consider getting there from each previous state  $k'$

The message is the vector of probabilities for the maximum probability path for each of the  $K$  states.

$$\omega(z_n) = \log(x_n | z_n) + \max_{z_{n-1}} \left\{ \log(p(z_n | z_{n-1}) + \omega(z_{n-1})) \right\}$$

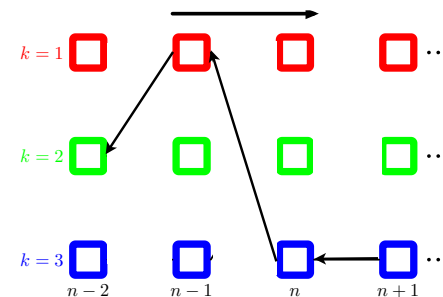
For each state  $k$

Consider getting there from each previous state  $k'$

We can see that this is the new maximum

For Viterbi, we need to remember the previous state,  $k'$ , for each  $k$ .

## Intuitive understanding



The max path is shown (but we only know it when we get to the end).

To find the path, we need to chase the back pointers.

## Classic HMM computational problems

- Given data, what is the HMM (**learning**). ✓
- Given an HMM, what is the **distribution over the state** variables. Also, **how likely** are the observations, given the model. ✓
- Given an HMM, what is the most likely **state sequence** for some data? ✓

## Final comments on learning

In many applications, the states have specified meaning, and are available in training data, so EM is not needed.

(Most authors still call this an HMM because states are hidden when the model is used).

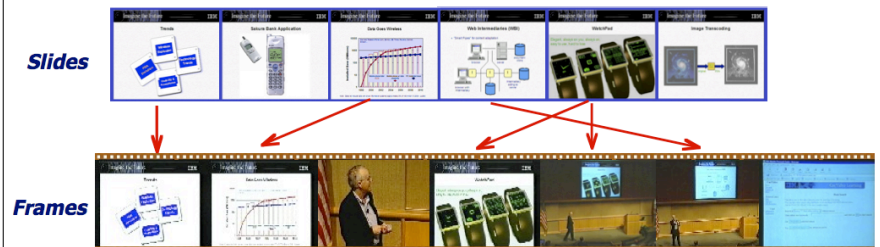
We described training the HMM based on a single data sequence, but often multiple sequences that come from the same HMM are used (modifying the algorithm is very straightforward).

## Two HMM examples (specified states)

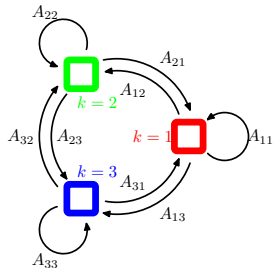
Domain is SLIC (Semantically Linked Instructional Content).

- 1) Temporal information for matching video frames to slides.
- 2) Aligning noisy speech transcripts with slides.

## Matching slides to video frames



## Matching slides to video frames



Our state sequence corresponds to what slide is being shown.

$$p(X, Z | \theta) = p(z_1 | \pi) \left[ \prod_{n=2}^N p(z_n | z_{n-1}, A) \right] \prod_{m=1}^N p(x_m | z_m, \phi)$$

## Matching slides to video frames

$$p(X, Z | \theta) = p(z_1 | \pi) \left[ \prod_{n=2}^N p(z_n | z_{n-1}, A) \right] \prod_{m=1}^N p(x_m | z_m, \phi)$$

From image matching

## Matching slides to video frames

$$p(X, Z | \theta) = p(z_1 | \pi) \left[ \prod_{n=2}^N p(z_n | z_{n-1}, A) \right] \prod_{m=1}^N p(x_m | z_m, \phi)$$

$$p(z_n | z_{n-1}, A) = f(z_n - z_{n-1}) \quad \text{encodes slide jump statistics.}$$

We assume that only the jump matters. IE, going from slide 6 to 8 has the same chance of going from 10 to 12.

$$p(z_1 | \pi) \left[ \prod_{n=2}^N p(z_n | z_{n-1}, A) \right] \quad \text{says how likely a sequence is, without looking at the images.}$$

## Aligning speech to slides

Why bother?

Mistakes in speech transcripts can be corrected.  
Speech transcripts are noisy and too poorly on jargon  
But jargon words often appear on slides.

We can highlight or auto-laser-point what the speaker is pointing to

We can improve close-captioning.

## Aligning speech to slides

A reasonable model for some speakers is that they say some approximation of their bullet points, with some extra stuff before and after.

Automated speech recognizers try to produce results that are plausible on a phoneme level.

If a slide word is used, its phoneme sequence will likely be approximated in the phonemes in the speech transcript.

We can calibrate the phoneme “confusion matrix.”

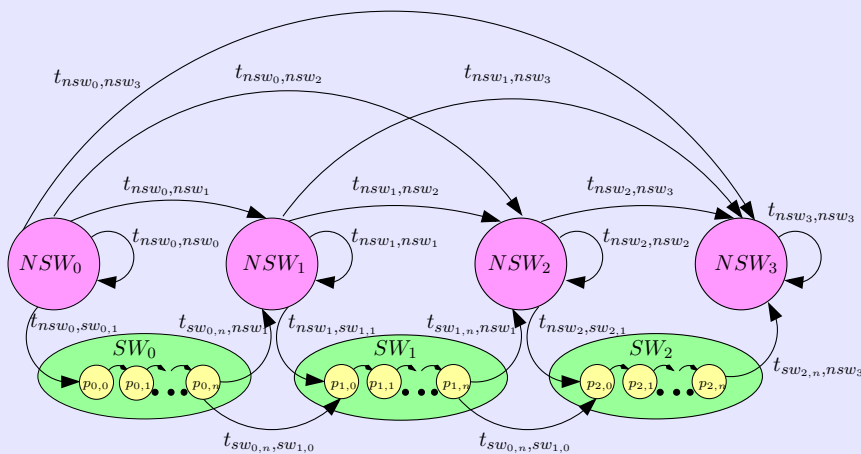
## Aligning speech to slides

We assume that going backwards does not happen.

We have an HMM state for each slide word

We also have an HMM state for emitting phonemes between slide words.

## SEQUENTIAL ALIGNMENT MODEL



## Aligned speech for correction

Speaker says : maliciousness

ASR produces: my dishes nests

*m ay d ih sh ah z n eh s t*

Slide word : maliciousness

*m ah l ih sh ah s n ah s*

with slide word phoneme sequence

If the same mistake is made later, where the word is not on the slide, we can propagate the correction.