

Statistical Methods in Natural Language Processing

9. Hidden Markov Models and Baum Welch.

Pavel Pecina, Jindřich Helcl

16 December, 2025

Course Segments

1. Introduction, probability, essential information theory
2. Statistical language modelling (n-gram)
3. Statistical properties of words
4. Word representations
5. Hidden Markov models, Tagging

Recap from Last Week

Markov Properties

- Markov Chain can generalize to any process (not just words):
 - Sequence of random variables: $X = (X_1, X_2, \dots, X_T)$
 - Sample space S (*states*), size N : $S = \{s_0, s_1, s_2, \dots, s_N\}$
- Two properties
 1. Limited history (context, horizon):

$$\forall i \in 1..T; P(X_i | X_1, \dots, X_{i-1}) = P(X_i | X_{i-1})$$

1 7 3 7 9 0 6 7 3 4 5... 1 7 3 7 9 0 6 7 3 4 5...

2. Time invariance (Markov Chain is stationary, homogeneous)

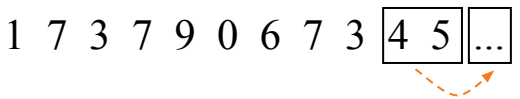
$$\forall i \in 1..T, \forall y, x \in S; P(X_i=y | X_{i-1}=x) = p(y|x)$$

1 7 3 7 9 0 6 7 3 4 5...

ok ... same distribution

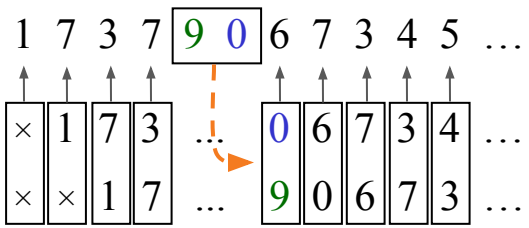
Long History Possible

- What if we want trigrams:



- Formally, use transformation:
 - Define new variables Q_i , such that $X_i = \{Q_{i-1}, Q_i\}$
 - And then $P(X_i|X_{i-1}) = P(Q_{i-1}, Q_i|Q_{i-2}, Q_{i-1}) = P(Q_i|Q_{i-2}, Q_{i-1})$

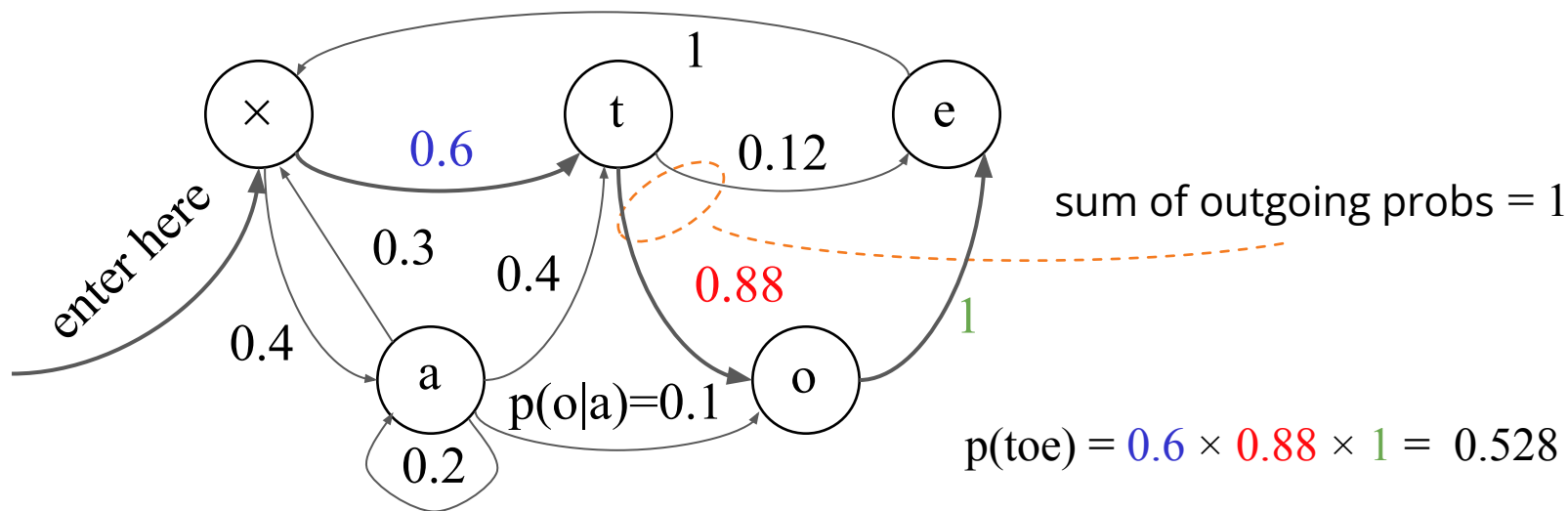
Predicting (X_i):



History ($X_{i-1} = \{Q_{i-2}, Q_{i-1}\}$):

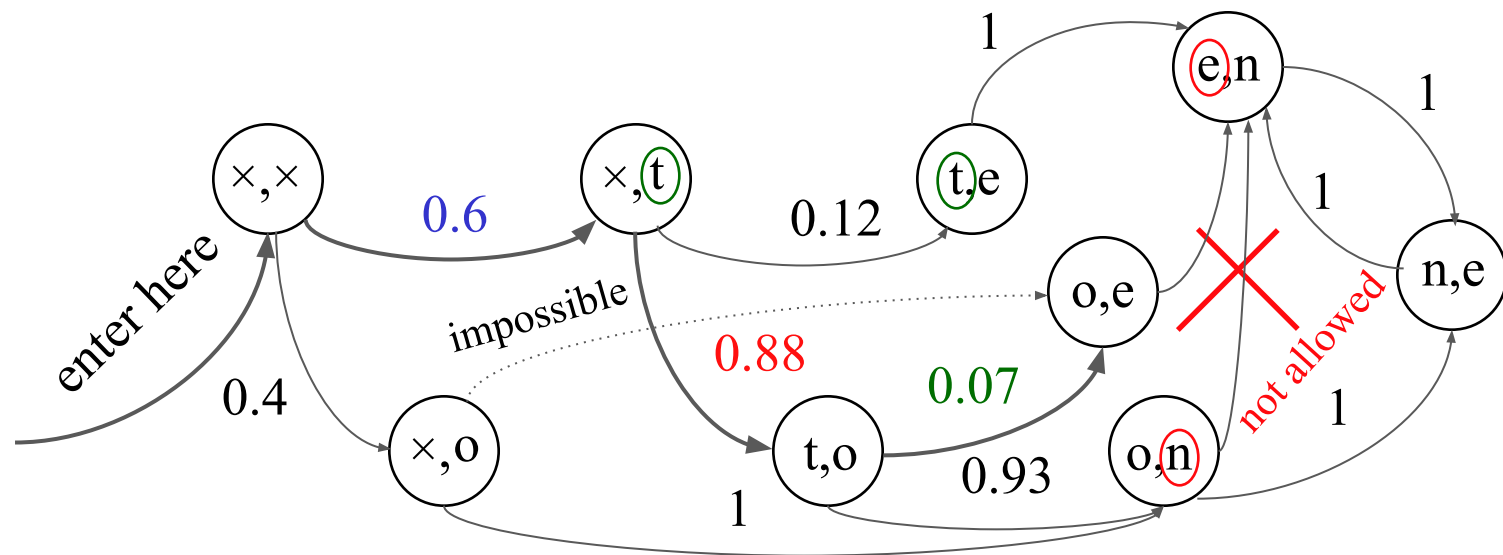
Markov Models: Bigram Case

- Nodes: States, $S = \{s_0, s_1, s_2, \dots, s_N\}$
- Arcs: Transitions with probabilities, $P(X_i|X_{i-1})$, X_i generates s_i



Markov Models: Trigram Case

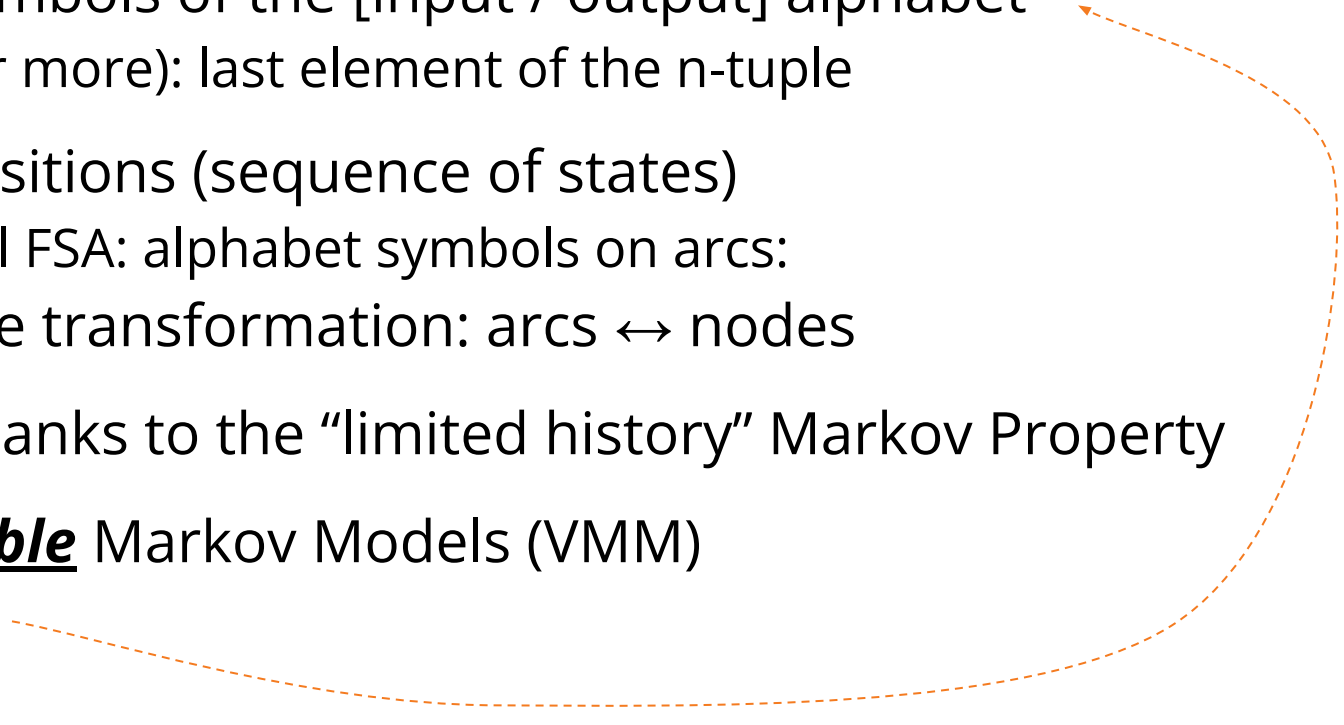
- Nodes: Pairs of states (s_k, s_l) , $S = \{s_0, s_1, s_2, \dots, s_N\}$
- Arcs: Transitions with probabilities, $P(X_i | X_{i-1})$, X_i generates (s_k, s_l)



$$p(\text{toe}) = 0.6 \times 0.88 \times 0.07 = 0.037$$

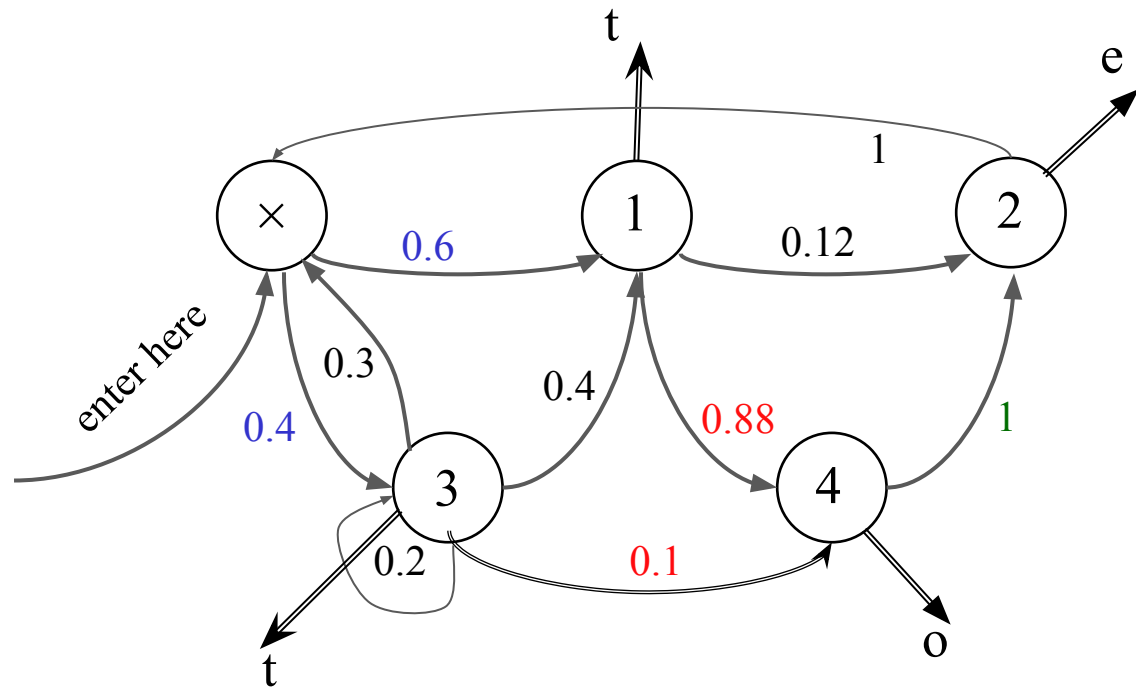
$$p(\text{one}) = ?$$

Finite State Automaton

- States ~ symbols of the [input / output] alphabet
 - pairs (or more): last element of the n-tuple
 - Arcs ~ transitions (sequence of states)
 - Classical FSA: alphabet symbols on arcs:
 - possible transformation: arcs \leftrightarrow nodes
 - Possible thanks to the “limited history” Markov Property
 - So far: **Visible** Markov Models (VMM)
- 

Hidden Markov Models

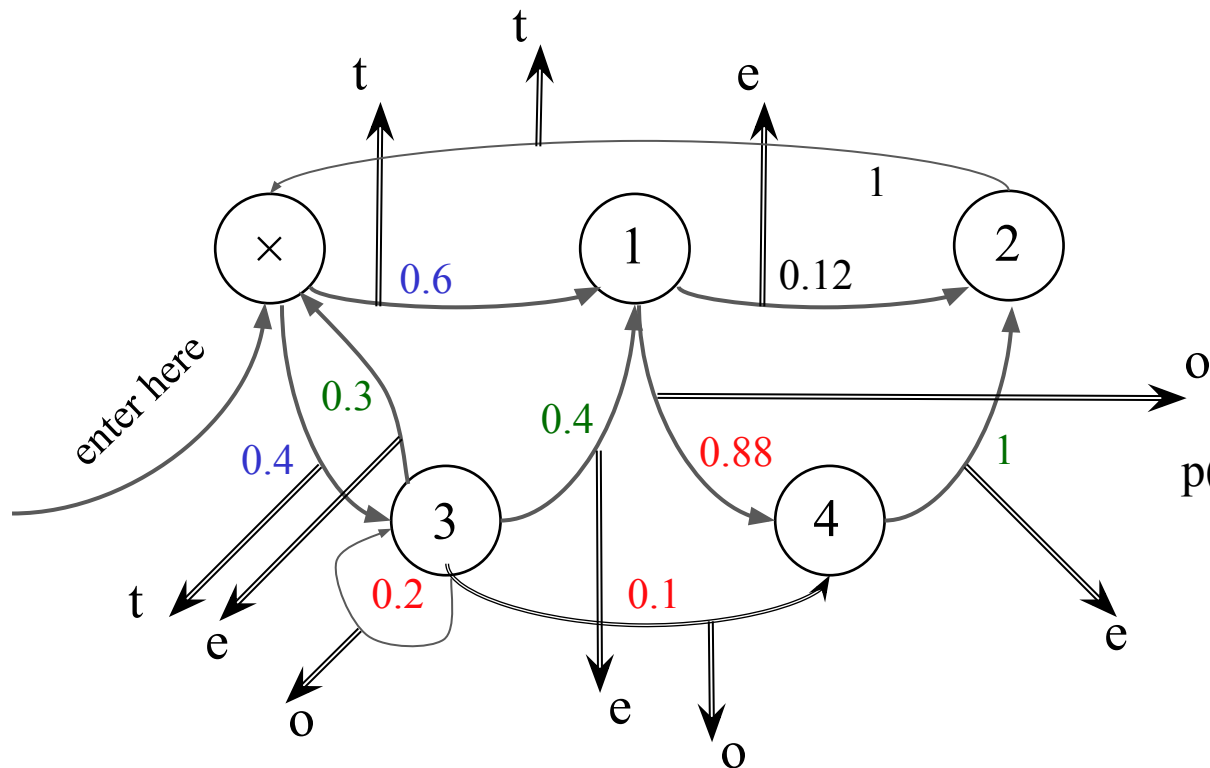
- The simplest HMM: states generate [*observable*] output (using the “data” alphabet) but remain “invisible”



$$p(\text{toe}) = 0.6 \times 0.88 \times 1 + 0.4 \times 0.1 \times 1 \cong 0.568$$

Output from Arcs...

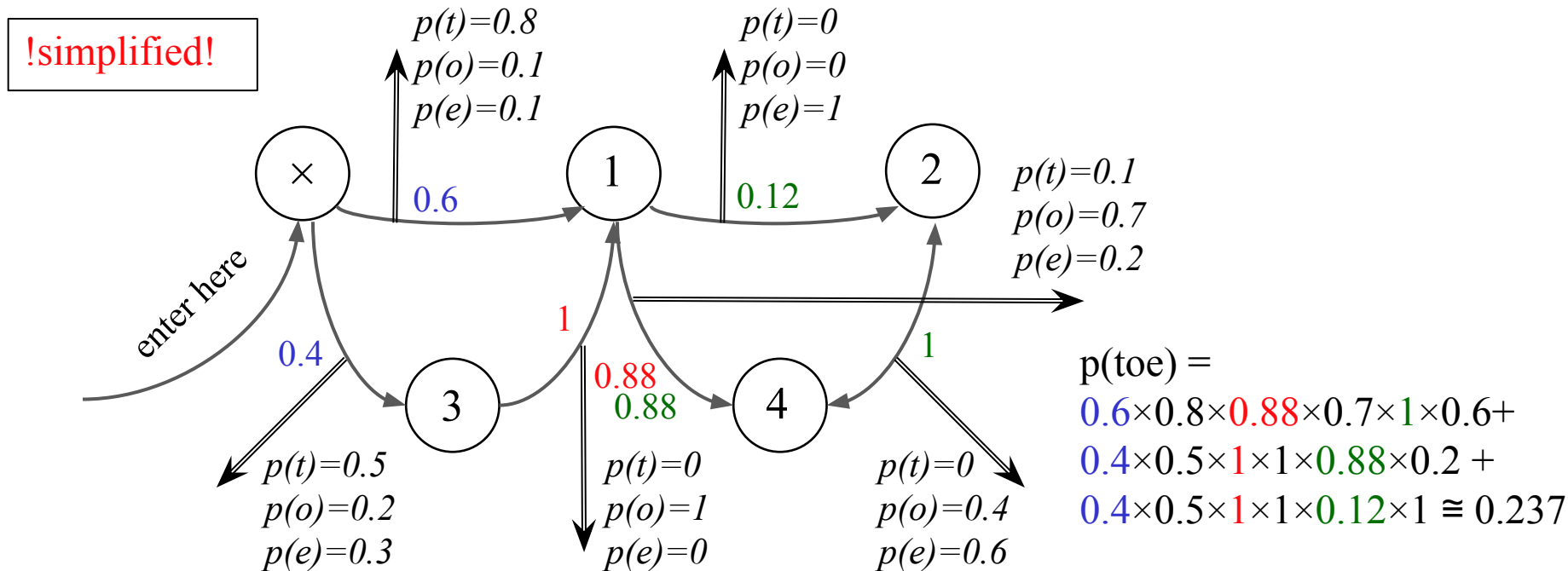
- Added flexibility: Generate output from arcs, not states:



$$\begin{aligned} p(\text{toe}) = & 0.6 \times 0.88 \times 1 + \\ & 0.4 \times 0.1 \times 1 + \\ & 0.4 \times 0.2 \times 0.3 + \\ & 0.4 \times 0.2 \times 0.4 \approx 0.624 \end{aligned}$$

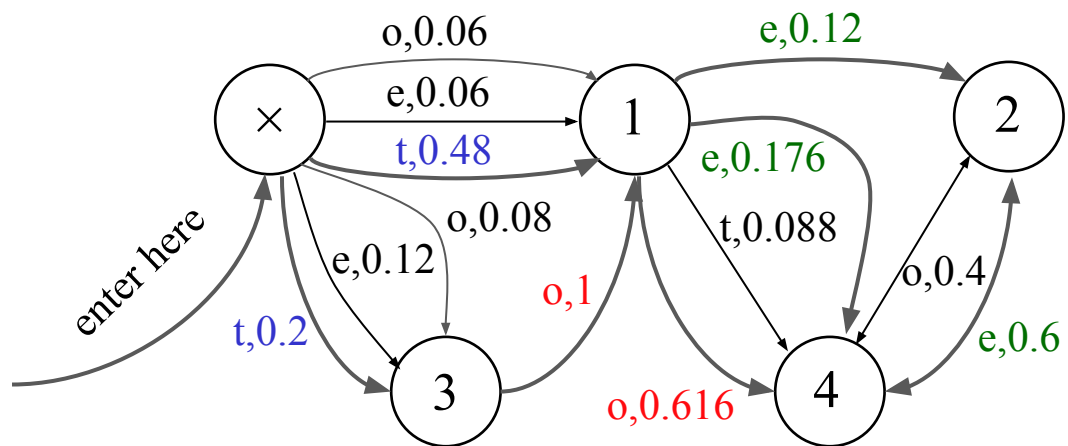
... and Finally, Add Output Probabilities

- Maximum flexibility: [Unigram] distribution (sample space: output alphabet) at each output arc:



Slightly Different View

- Allow for multiple arcs from $s_i \rightarrow s_j$, mark them by output symbols, get rid of output distributions:



$$\begin{aligned} p(\text{toe}) = & 0.48 \times 0.616 \times 0.6 + \\ & 0.2 \times 1 \times 0.176 + \\ & 0.2 \times 1 \times 0.12 \cong 0.237 \end{aligned}$$

Formalization

- HMM (the general case): five-tuple (S, s_0, Y, P_S, P_Y) , where:
 - $S = \{s_1, s_2, \dots, s_T\}$ is the set of states, s_0 is the initial state,
 - $Y = \{y_1, y_2, \dots, y_V\}$ is the output alphabet,
 - $P_S(s_j | s_i)$ is the set of prob. distributions of transitions,
 - size of P_S : $|S|^2$
 - $P_Y(y_k | s_i, s_j)$ is the set of output (emission) probability distributions
 - size of P_Y : $|S|^2 \times |Y|$
- Example:
 - $S = \{x, 1, 2, 3, 4\}, s_0 = x$
 - $Y = \{t, o, e\}$

Formalization - Example

- Example (for graph on slide 36):
 - $S = \{x, 1, 2, 3, 4\}, s_0 = x$
 - $Y = \{t, o, e\}$
 - $P_S:$

	x	1	2	3	4
x	0	0.6	0	0.4	0
1	0	0	0.12	0	0.88
2	0	0	0	0	1
3	0	1	0	0	0
4	0	0	1	0	0

⇒ $\Sigma = 1$

$P_Y:$

	e	x	1	2	3	4
o	x	1	2	3	4	
t	x	1	2	3	4	0.2
x		0.8		0.5		0.7
1					0.1	
2					0	
3		0				
4			0			

$\Sigma = 1$

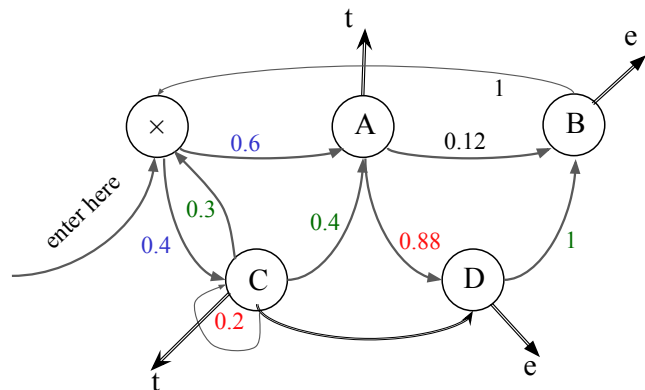
HMM: The Two Tasks

- HMM (the general case): five-tuple (S, S_0, Y, P_S, P_Y) , where:
 - $S = \{s_1, s_2, \dots, s_T\}$ is the set of states, S_0 is the initial state,
 - $Y = \{y_1, y_2, \dots, y_V\}$ is the output alphabet,
 - $P_S(s_j | s_i)$ is the set of prob. distributions of transitions,
 - $P_Y(y_k | s_i, s_j)$ is the set of output (emission) probability distributions
- Given an HMM & an output sequence $Y = \{y_1, y_2, \dots, y_k\}$:
 - (Task 1) compute the probability of Y ;
 - (Task 2) compute the most likely sequence of states which has generated Y .

Trellis

Trellis: HMM “roll-out”

HMM:



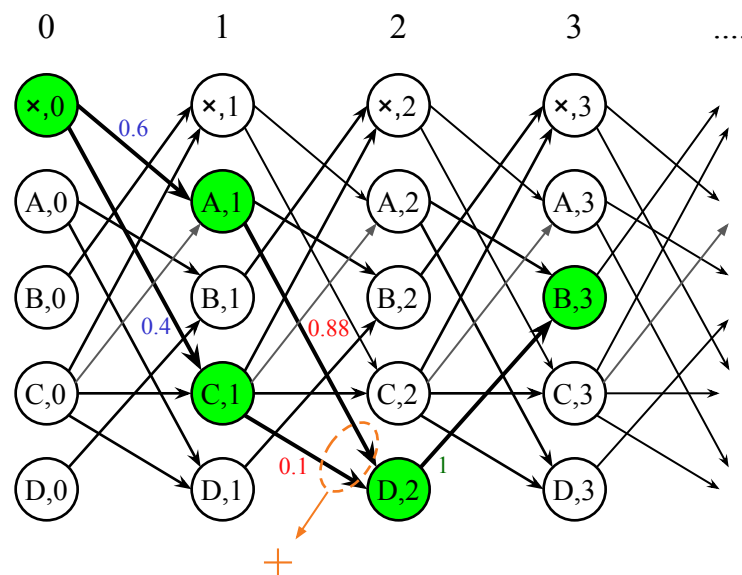
$$p(\text{toe}) = 0.6 \times 0.88 \times 1 + 0.4 \times 1 \times 1 \cong 0.568$$

- Trellis state: (HMM state, position)
- each state: holds **one** number (prob): α
- probability of Y: $\sum \alpha$ in the last state

Trellis:

time/position:

“roll-out”



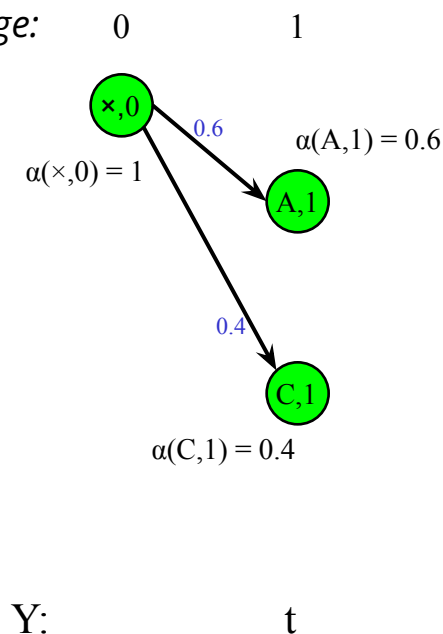
Y:

$$\alpha(\times,0) = 1 \quad \alpha(A,1) = 0.6 \quad \alpha(D,2) = 0.568 \quad \alpha(B,3) = 0.568$$

$$\alpha(C,1) = 0.4$$

Creating the Trellis: The Start

- Start in the start state (\times),
 - set its $\alpha(\times, 0)$ to 1
- Create the first stage:
 - get the first “output” symbol y_1
 - create the first stage (column)
 - but only those Trellis states which generate y_1
 - set their $\alpha(state, 1)$ to the $P_s(state|\times) \alpha(\times, 0)$
- and forget about the 0 -th stage



Trellis: The Next Step

- Suppose we are in stage i
- Creating the next stage:
 - create all trellis states in the next stage which generate y_{i+1} , but only those reachable from any of the stage- i states
 - set their $\alpha(state, i+1)$ to:
 $P_S(state|prev.state) \times \alpha(prev.state, i)$
(add up all such numbers on arcs going to a common Trellis state)
...and forget about stage i

position/stage:

1

2

$$\alpha(A,1) = 0.6$$

A,1

C,1

$$\alpha(C,1) = 0.4$$

Y:

t

o

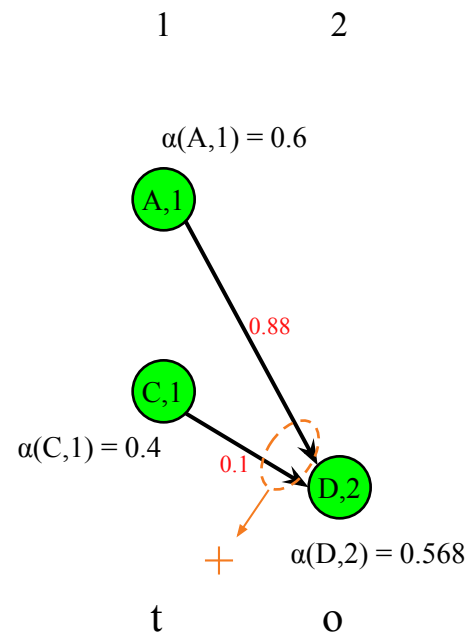
Trellis: The Next Step

- Suppose we are in stage i
- Creating the next stage: *position/stage:*

1	2
---	---

 - create all trellis states in the next stage which generate y_{i+1} , but only those reachable from any of the stage- i states
 - set their $\alpha(state, i+1)$ to:
 $P_S(state|prev.state) \times \alpha(prev.state, i)$
(add up all such numbers on arcs going to a common Trellis state)
...and forget about stage i

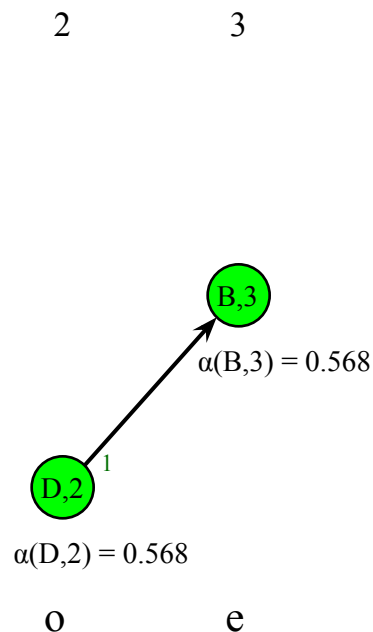
Y:



Trellis: The Last Step

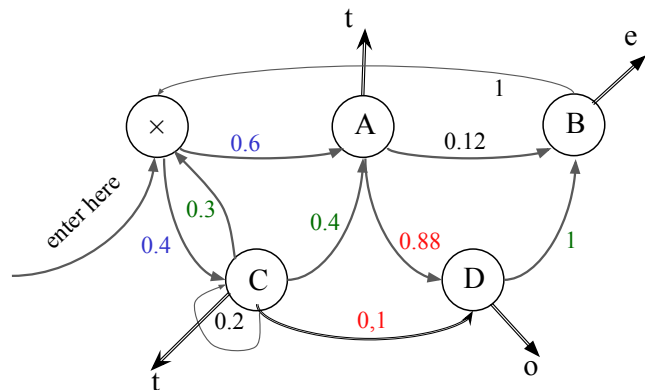
- Continue until output exhausted
 - $|Y| = 3$: until stage 3 *position/stage:*
- Add together all the $\alpha(state, |Y|)$
- That's the $P(Y)$
- Observation (pleasant):
 - memory usage max: $2|S|$
 - multiplications max: $|S|^2|Y|$

Y:



Trellis (again)

HMM:



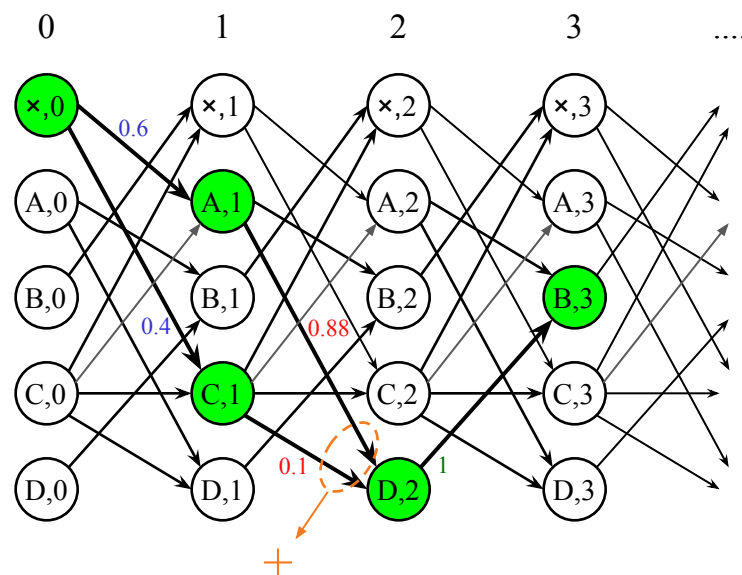
$$p(\text{toe}) = 0.6 \times 0.88 \times 1 + 0.4 \times 1 \times 1 \cong 0.568$$

- Trellis state: (HMM state, position)
- each state: holds **one** number (prob): α
- probability of Y: $\sum \alpha$ in the last state

Trellis:

time/position:

"roll-out"



Y:

$$\alpha(\times,0) = 1 \quad \alpha(A,1) = 0.6 \quad \alpha(D,2) = 0.568 \quad \alpha(B,3) = 0.568$$

$$\alpha(C,1) = 0.4$$

HMM: The Two Tasks

- HMM (the general case): five-tuple (S, S_0, Y, P_S, P_Y) , where:
 - $S = \{s_1, s_2, \dots, s_T\}$ is the set of states, S_0 is the initial state,
 - $Y = \{y_1, y_2, \dots, y_V\}$ is the output alphabet,
 - $P_S(s_j | s_i)$ is the set of prob. distributions of transitions,
 - $P_Y(y_k | s_i, s_j)$ is the set of output (emission) probability distributions
- Given an HMM & an output sequence $Y = \{y_1, y_2, \dots, y_k\}$:
 - (Task 1) compute the probability of Y ;
 - (Task 2) compute the most likely sequence of states which has generated Y .

Trellis: The General Case (still, bigrams)

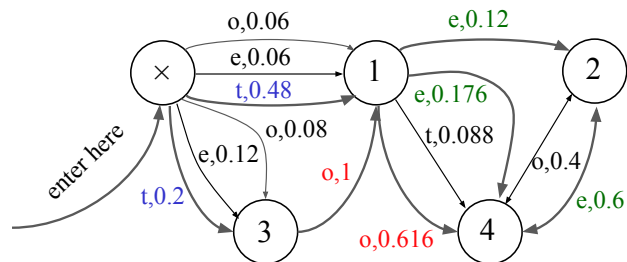
- Start as usual:

- start state (\times), set its $\alpha(\times, 0)$ to 1.

position/stage: 0



$$\alpha(\times, 0) = 1$$



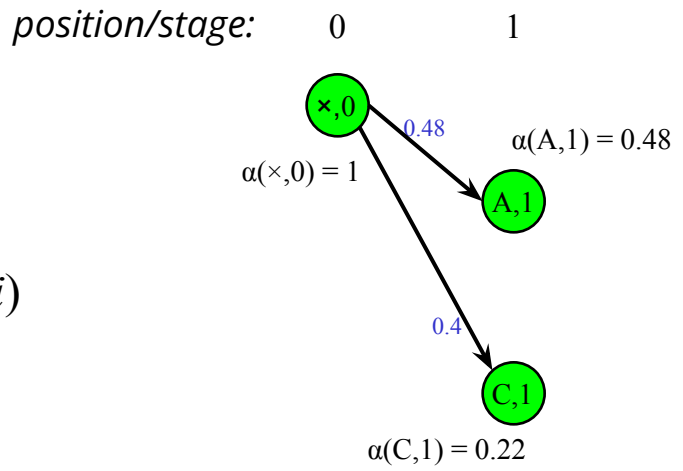
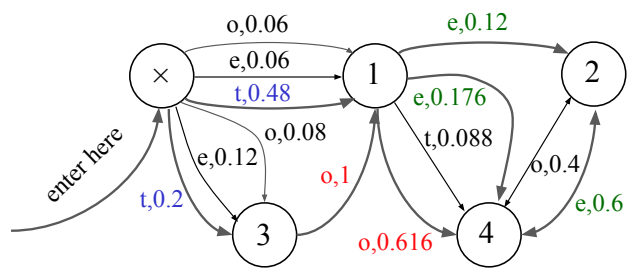
$$\begin{aligned} p(\text{toe}) &= 0.48 \times 0.616 \times 0.6 + \\ &\quad 0.2 \times 1 \times 0.176 + \\ &\quad 0.2 \times 1 \times 0.12 \cong 0.237 \end{aligned}$$

Y:

General Trellis: The Next Step

We are in stage i :

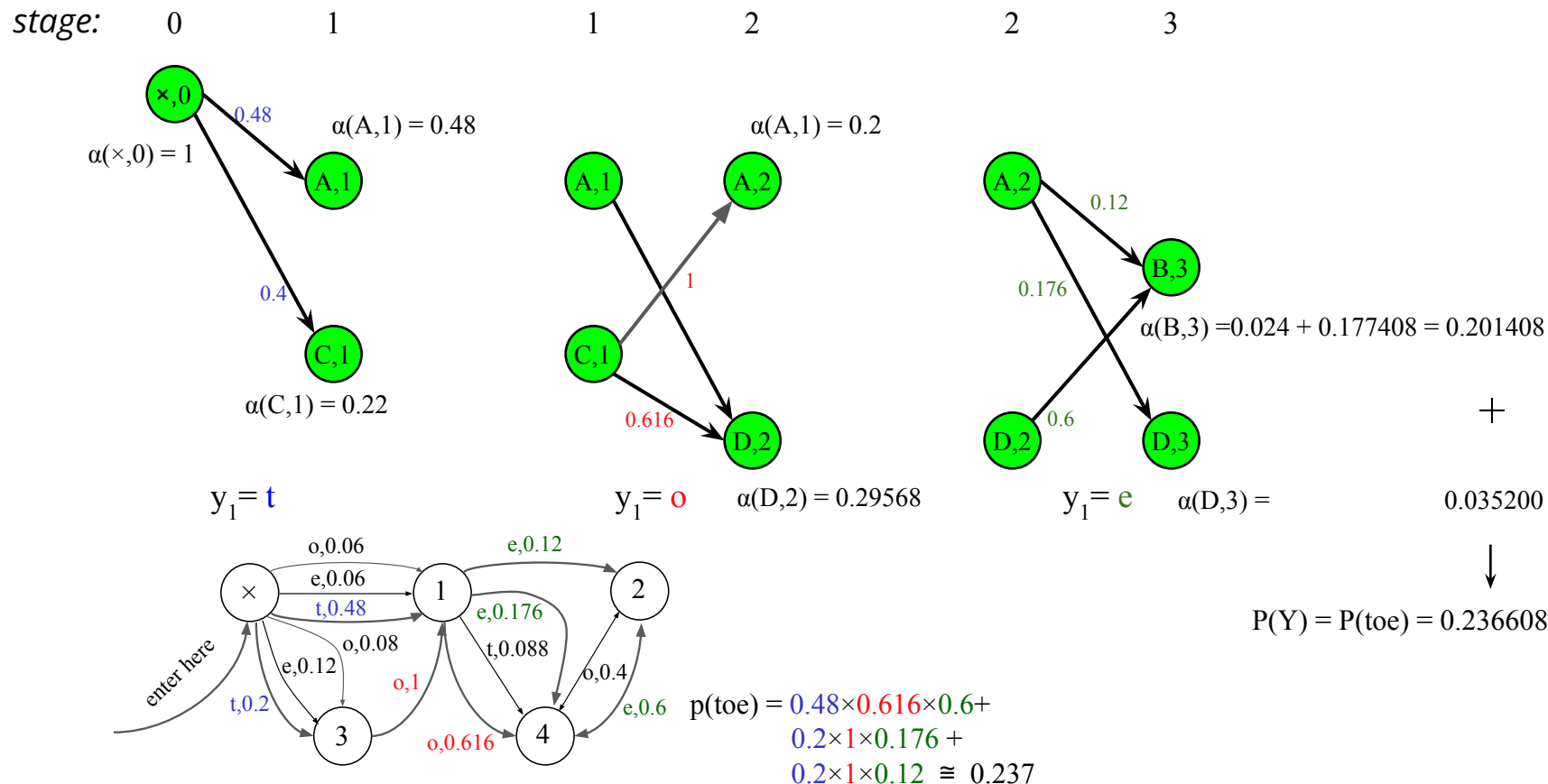
- Generate the next stage $i+1$ as before (except now arcs generate output, thus use only those arcs marked by the output symbol y_{i+1})
- For each generated *state*, compute $\alpha(state, i+1)$
 $= \sum_{inc. \text{ arcs}} P_Y(y_{i+1} | state, prev.state) \times \alpha(prev.state, i)$



Y: t

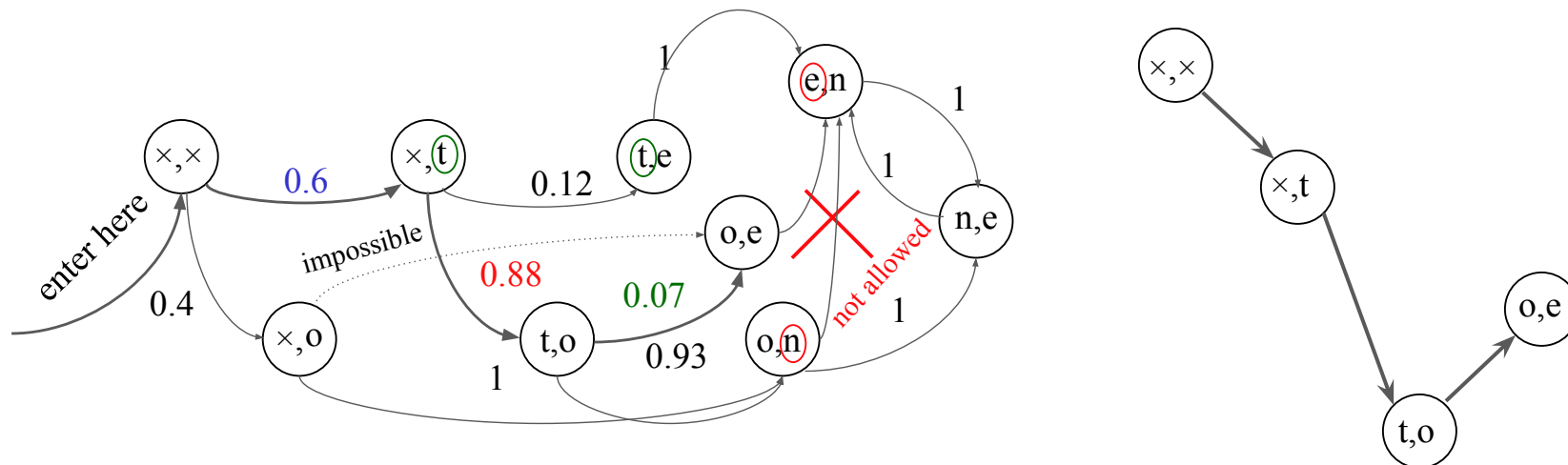
And forget about the previous state ...

Trellis: The Complete Example



The Case of Trigrams

- Like before, but:
 - states correspond to bigrams
 - output function always emits the second output symbol of the pair (state) to which the arc goes:



- Multiple paths not possible \rightarrow trellis not really needed

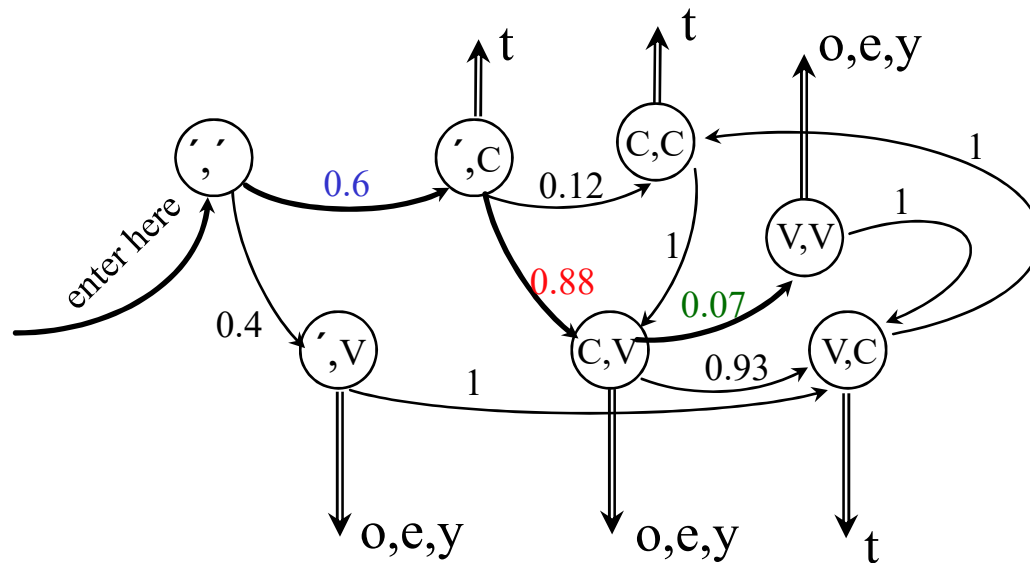
Trigrams with Classes

- More interesting:

- n-gram class LM: $p(w_i|w_{i-2},w_{i-1}) = p(w_i|c_i) p(c_i|c_{i-2},c_{i-1})$

- states are pairs of classes (c_{i-1},c_i) , and emit “words”:

(letters in our example)



$p(t|C) = 1$ usual,
 $p(o|V) = .3$ non-
 $p(e|V) = .6$ overlapping
 $p(y|V) = .1$ classes

$$p(\text{toe}) = .6 \cdot 1 \cdot .88 \cdot .3 \cdot .07 \cdot .6 \cong .00665$$

$$p(\text{teo}) = .6 \cdot 1 \cdot .88 \cdot .6 \cdot .07 \cdot .3 \cong .00665$$

$$p(\text{toy}) = .6 \cdot 1 \cdot .88 \cdot .3 \cdot .07 \cdot .1 \cong .00111$$

$$p(\text{tty}) = .6 \cdot 1 \cdot .12 \cdot 1 \cdot 1 \cdot .1 \cong .0072$$

Class Trigrams: the Trellis

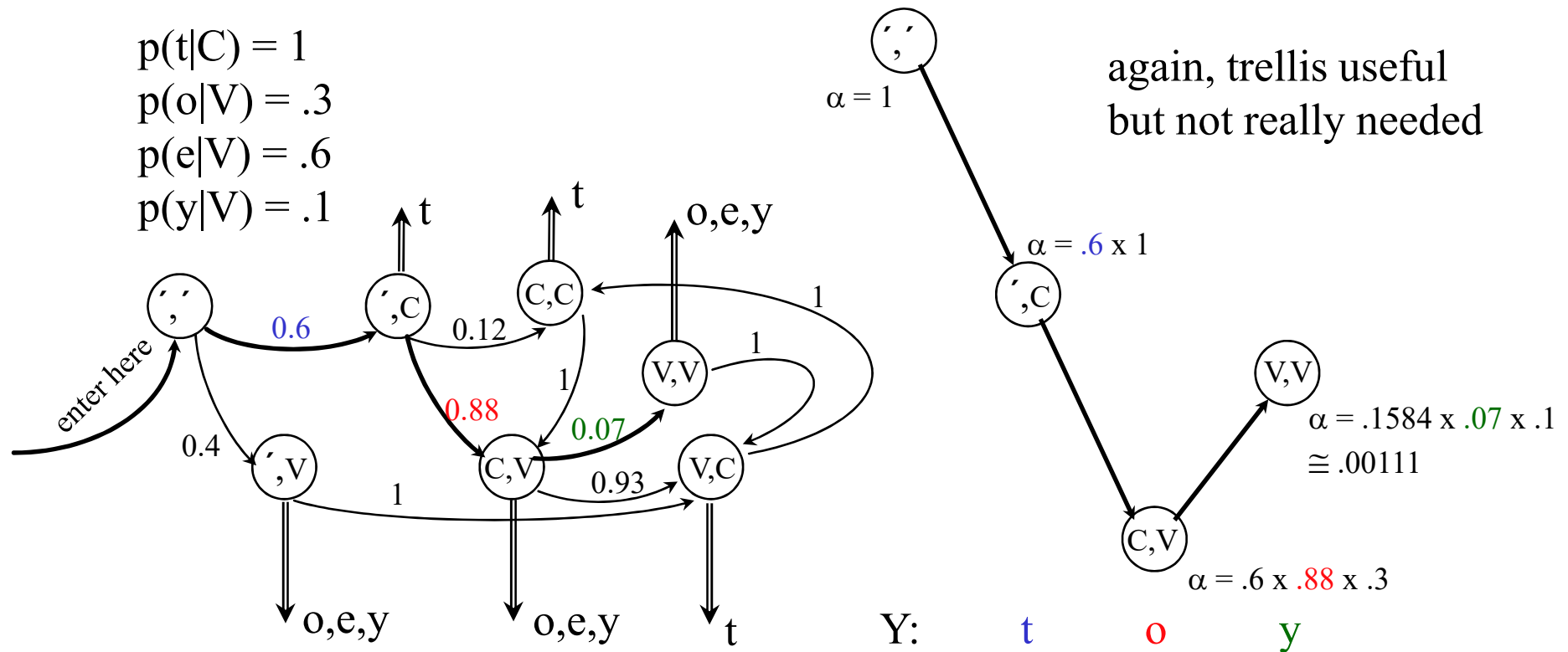
- Trellis generation (Y = “toy”):

$$p(t|C) = 1$$

$$p(o|V) = .3$$

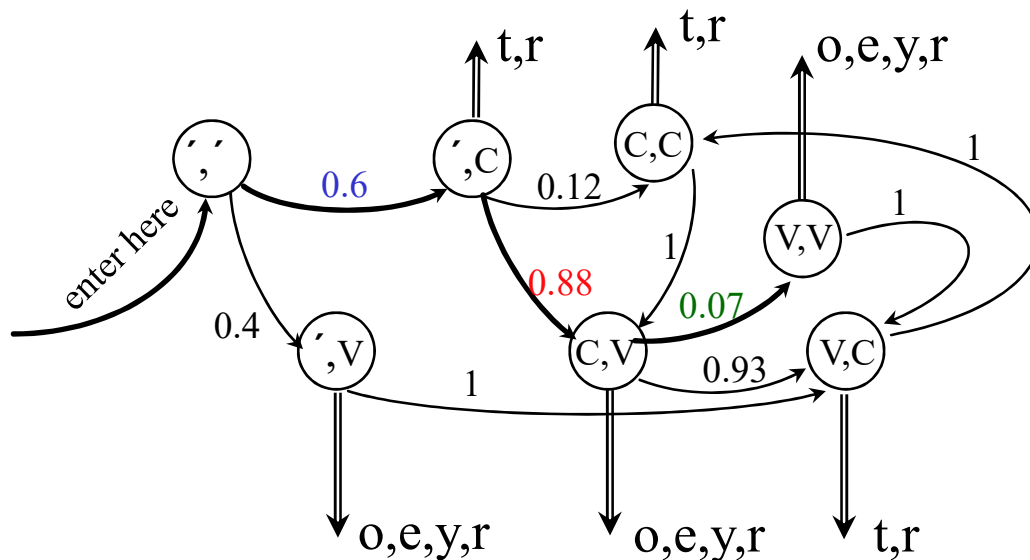
$$p(e|V) = .6$$

$$p(y|V) = .1$$



Overlapping Classes

- Imagine that classes may overlap
 - e.g. 'r' is sometimes vowel sometimes consonant, belongs to V as well as C:

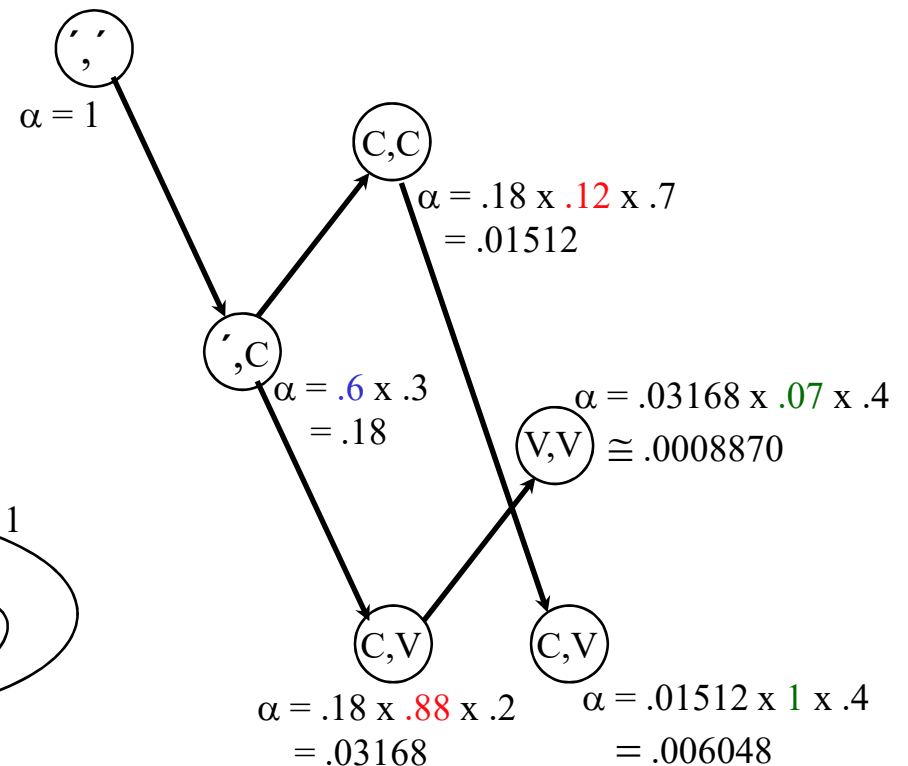
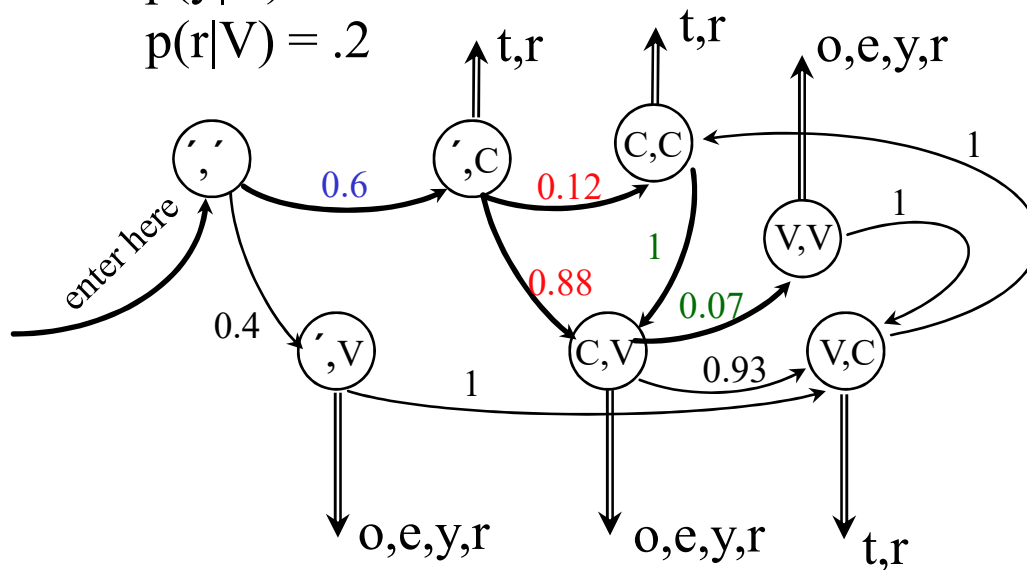


$$\begin{aligned}
 p(t|C) &= .3 \\
 p(r|C) &= .7 \\
 p(o|V) &= .1 \\
 p(e|V) &= .3 \\
 p(y|V) &= .4 \\
 p(r|V) &= .2
 \end{aligned}$$

$$p(\text{try}) = ?$$

Overlapping Classes: Trellis Example

$$\begin{aligned} p(t|C) &= .3 \\ p(r|C) &= .7 \\ p(o|V) &= .1 \\ p(e|V) &= .3 \\ p(y|V) &= .4 \\ p(r|V) &= .2 \end{aligned}$$



Y: t r y $p(Y) = \underline{.006935}$

Trellis: Remarks

- So far, we went left to right (computing α)
- Same result: going right to left (computing β)
 - supposed we know where to start (finite data)
- In fact, we might start in the middle going left and right
- Important for parameter estimation
(Forward-Backward Algorithm alias Baum-Welch)
- Implementation issues:
 - scaling/normalizing probabilities, to avoid too small numbers
& addition problems with many transitions

The Viterbi Algorithm

- Solving the task of finding the most likely sequence of states which generated the observed data
- i.e., finding

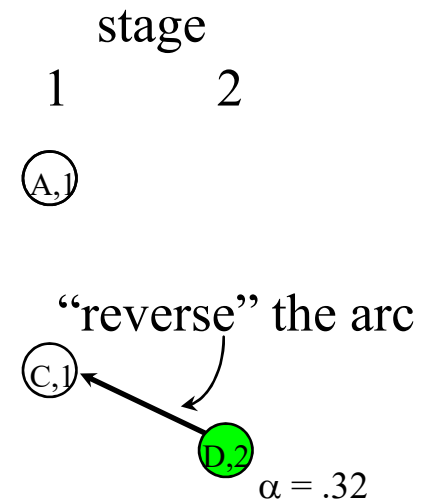
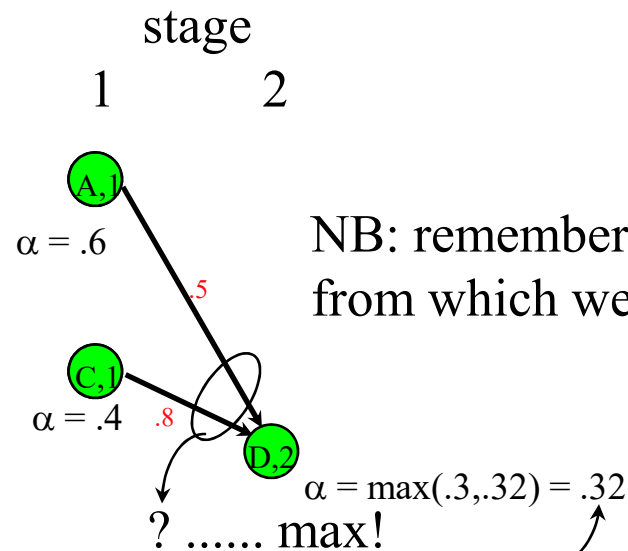
$$S_{\text{best}} = \operatorname{argmax}_S P(S|Y)$$

which is equal to (Y is constant and thus $P(Y)$ is fixed):

$$\begin{aligned} S_{\text{best}} &= \operatorname{argmax}_S P(S, Y) = \\ &= \operatorname{argmax}_S P(s_0, s_1, s_2, \dots, s_k, y_1, y_2, \dots, y_k) = \\ &= \operatorname{argmax}_S \prod_{i=1..k} p(y_i | s_i, s_{i-1}) p(s_i | s_{i-1}) \end{aligned}$$

The Crucial Observation

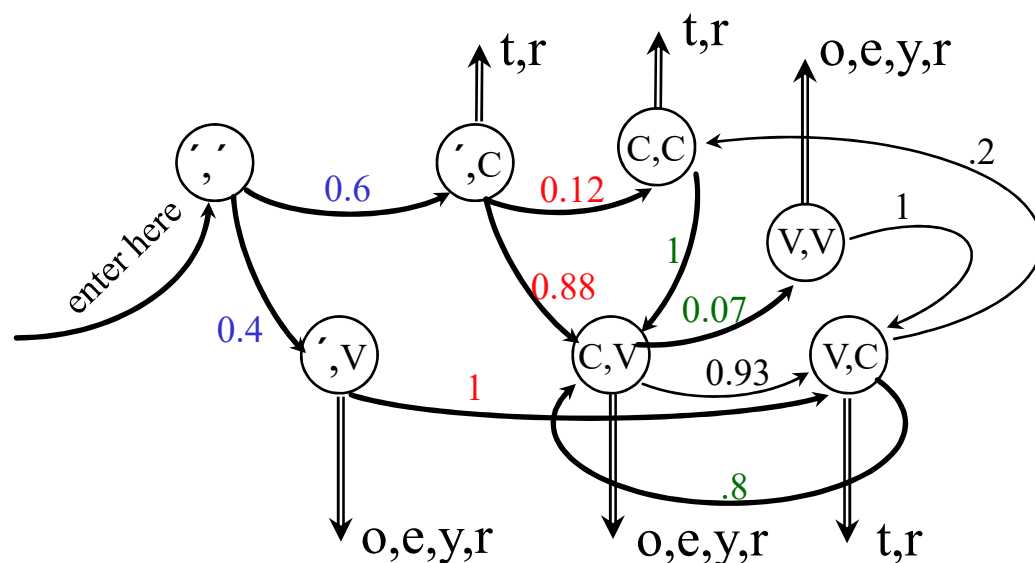
- Imagine the trellis build as before (but do not compute the α s yet; assume they are o.k.); stage i :



this is certainly the “backwards” maximum to (D,2)... but
it cannot change even whenever we go forward (M. Property: Limited History)

Viterbi Example

- ‘r’ classification (C or V?, sequence?):



$$p(t|C) = .3$$

$$p(r|C) = .7$$

$$p(o|V) = .1$$

$$p(e|V) = .3$$

$$p(y|V) = .4$$

$$p(r|V) = .2$$

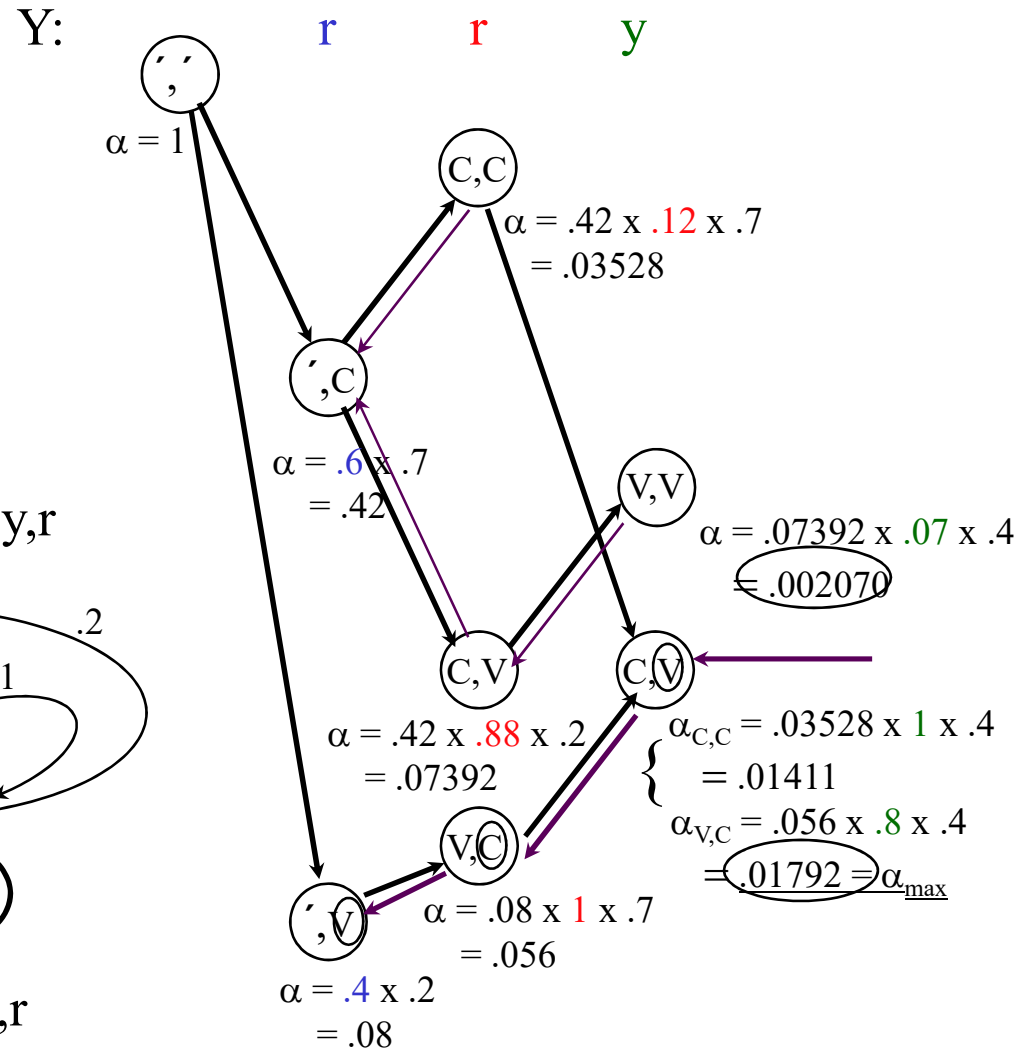
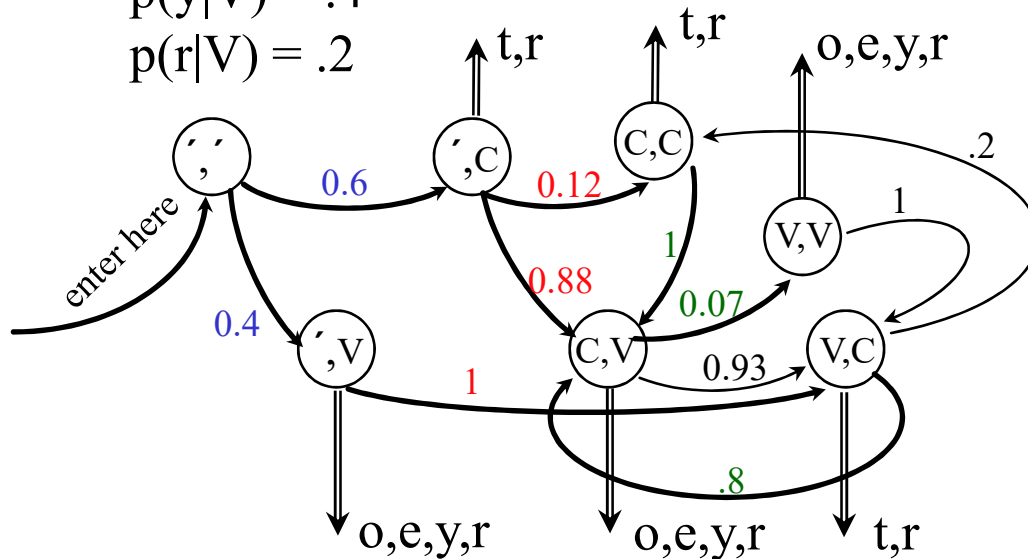
$$\operatorname{argmax}_{XYZ} p(rry|XYZ) = ?$$

Possible state seq.: (.,v)(v,c)(c,v)[VCV], (.,c)(c,c)(c,v)[CCV], (.,c)(c,v)(v,v) [CVV]

Viterbi Computation

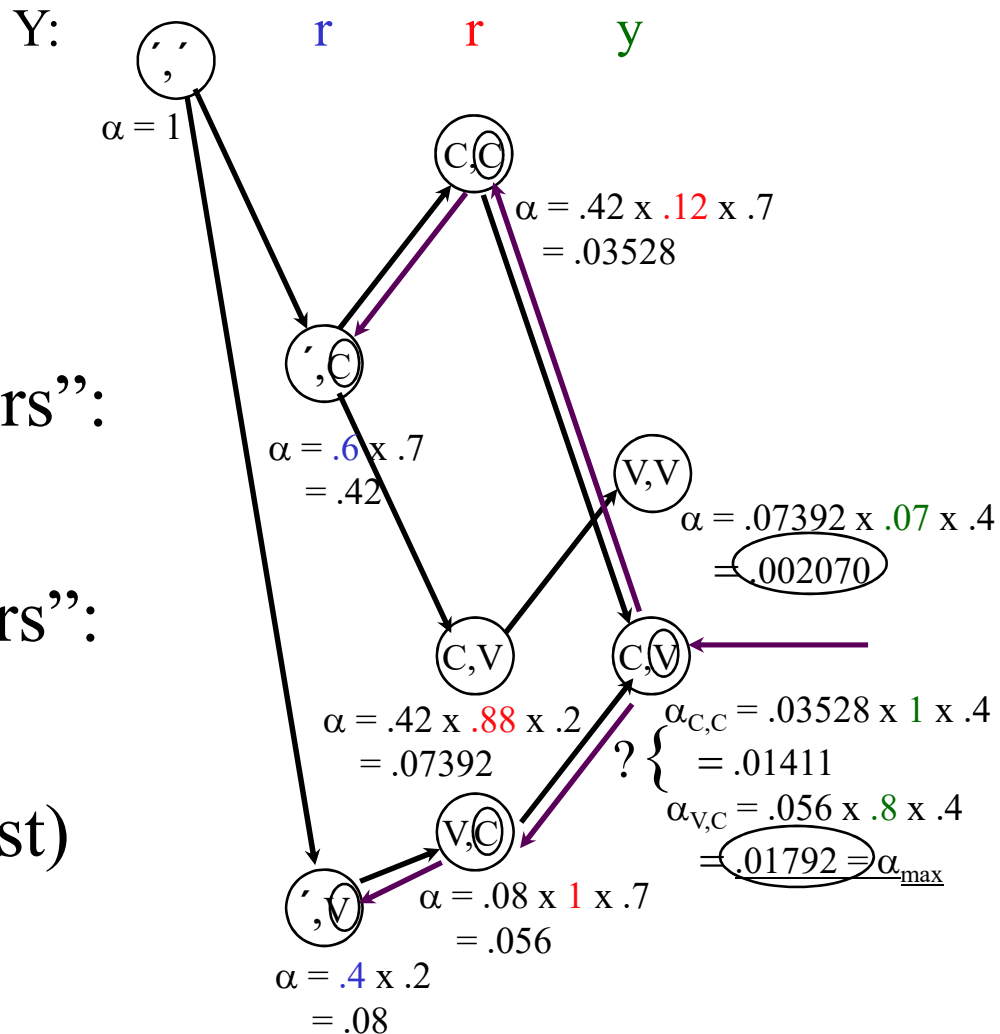
$$\begin{aligned} p(t|C) &= .3 \\ p(r|C) &= .7 \\ p(o|V) &= .1 \\ p(e|V) &= .3 \\ p(y|V) &= .4 \\ p(r|V) &= .2 \end{aligned}$$

α in trellis
state:
best prob
from start
to here



n-best State Sequences

- Keep track of n best “back pointers”:
- Ex.: $n=2$:
Two “winners”:
VCV (best)
CCV (2nd best)

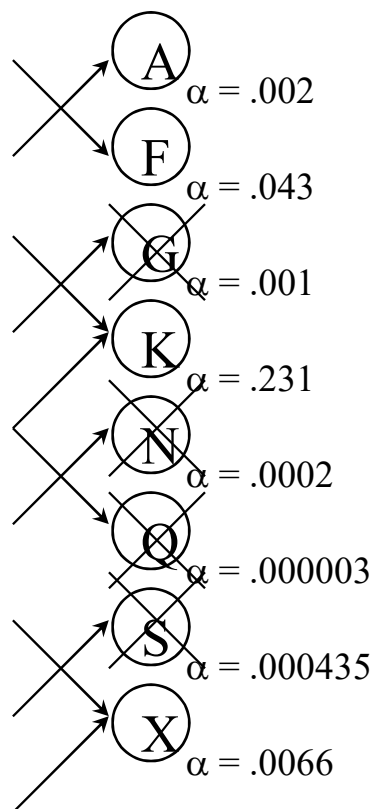


Tracking Back the n-best paths

- Backtracking-style algorithm:
 - **Start at the end, in the best of the n states (s_{best})**
 - **Put the other n-1 best nodes/back pointer pairs on stack, except those leading from s_{best} to the same best-back state.**
- Follow the back “beam” towards the start of the data, spitting out nodes on the way (backwards of course) using always only the best back pointer.
- At every beam split, push the diverging node/back pointer pairs onto the stack (node/beam width is sufficient!).
- When you reach the start of data, close the path, and pop the top-most node/back pointer(width) pair from the stack.
- Repeat until the stack is empty; expand the result tree if necessary.

Pruning

- Sometimes, too many trellis states in a stage:



criteria: (a) $\alpha < \text{threshold}$
(b) $\Sigma\pi < \text{threshold}$
(c) # of states $> \text{threshold}$
(get rid of smallest α)

HMM Parameter Estimation: the Baum-Welch Algorithm

HMM: The Tasks

- HMM (the general case):
 - five-tuple (S, S_0, Y, P_S, P_Y) , where:
 - $S = \{s_1, s_2, \dots, s_T\}$ is the set of states, S_0 is the initial state,
 - $Y = \{y_1, y_2, \dots, y_V\}$ is the output alphabet,
 - $P_S(s_j | s_i)$ is the set of prob. distributions of transitions,
 - $P_Y(y_k | s_i, s_j)$ is the set of output (emission) probability distributions.
- Given an HMM & an output sequence $Y = \{y_1, y_2, \dots, y_k\}$:
 - ✓ (Task 1) compute the probability of Y ;
 - ✓ (Task 2) compute the most likely sequence of states which has generated Y .
 - (Task 3) Estimating the parameters (transition/output distributions)

A Variant of EM

- Idea (\sim EM, for another variant see LM smoothing):
 - Start with (possibly random) estimates of P_S and P_Y .
 - Compute (fractional) “counts” of state transitions/emissions taken, from P_S and P_Y , given data Y .
 - Adjust the estimates of P_S and P_Y from these “counts” (using the MLE, i.e. relative frequency as the estimate).
- Remarks:
 - many more parameters than the simple four-way smoothing
 - no proofs here; see Jelinek, Chapter 9

Setting

- HMM (without P_S, P_Y) (S, S_0, Y), and data $T = \{y^i \in Y\}_{i=1..|T|}$
 - **will use $T \sim |T|$**
- HMM structure is given: (S, S_0)
- P_S : Typically, one wants to allow “fully connected” graph
 - **(i.e. no transitions forbidden ~ no transitions set to hard 0)**
 - **why? → we better leave it on the learning phase, based on the data!**
 - **sometimes possible to remove some transitions ahead of time**
- P_Y : should be restricted (if not, we will not get anywhere!)
 - **restricted ~ hard 0 probabilities of $p(y|s,s')$**
 - **“Dictionary”: states \leftrightarrow words, “m:n” mapping on $S \times Y$ (in general)**

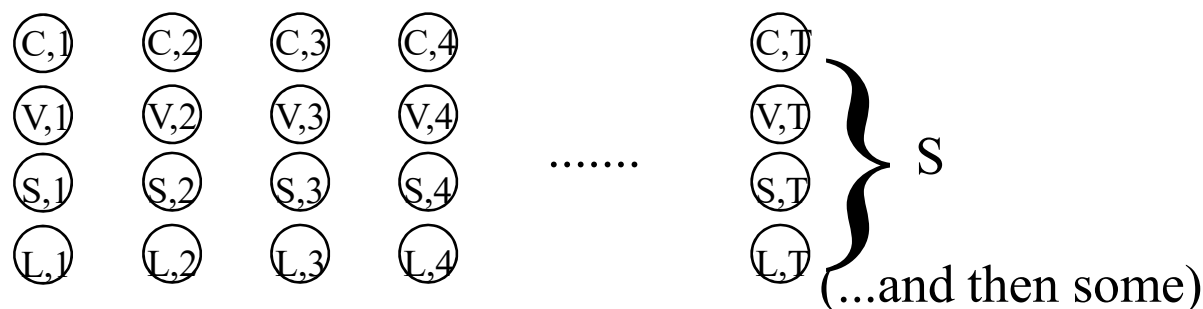
Initialization

- For computing the initial expected “counts”
- Important part
 - EM guaranteed to find a *local* maximum only (albeit a good one in most cases)
- P_Y initialization more important
 - fortunately, often easy to determine
 - **together with dictionary \leftrightarrow vocabulary mapping, get counts, then MLE**
- P_S initialization less important
 - e.g. uniform distribution for each $p(.|s)$

Data Structures

- Will need storage for:
 - The predetermined structure of the HMM
(unless fully connected \rightarrow need not to keep it!)
 - The parameters to be estimated (P_S, P_Y)
 - The expected counts (same size as P_S, P_Y)
 - The training data $T = \{y^i \in Y\}_{i=1..T}$
 - The trellis (if f.c.): $\uparrow T$ Size: $T' \times S$ (Precisely, $|T|' \times |S|$)

Each trellis state:
two [float] numbers
 (forward/backward)



The Algorithm Part I

1. Initialize P_S, P_Y

2. Compute “forward” probabilities:

- follow the procedure for trellis (summing), compute $\alpha(s,i)$
- use the current values of P_S, P_Y ($p(s'|s), p(y|s,s')$):

$$\alpha(s',i) = \sum_{s \rightarrow s'} \alpha(s,i-1) \times p(s'|s) \times p(y_i|s,s')$$

- **NB: do not throw away the previous stage!**

3. Compute “backward” probabilities

- start at all nodes of the last stage, proceed backwards, $\beta(s,i)$
- i.e., probability of the “tail” of data from stage i to the end of data

$$\beta(s',i) = \sum_{s \leftarrow s'} \beta(s,i+1) \times p(s|s') \times p(y_{i+1}|s',s)$$

- also, keep the $\beta(s,i)$ at all trellis states

The Algorithm Part II

4. Collect counts:

- for each output/transition pair compute

$$c(y, s, s') = \sum_{i=0..k-1, y=y_{i+1}} \alpha(s, i) \underbrace{p(s'|s) p(y_{i+1}|s, s')}_{\text{this transition prob}} \beta(s', i+1)$$

one pass through data, only stop at (output) y prefix prob. ' output prob tail prob

$$c(s, s') = \sum_{y \in Y} c(y, s, s') \text{ (assuming all observed } y_i \text{ in } Y)$$

$$c(s) = \sum_{s' \in S} c(s, s')$$

5. Reestimate: $p'(s'|s) = c(s, s')/c(s)$ $p'(y|s, s') = c(y, s, s')/c(s, s')$

6. Repeat 2-5 until desired convergence limit is reached.

Baum-Welch: Tips & Tricks

- Normalization badly needed
 - long training data → extremely small probabilities
- Normalize α, β using the same norm. factor:

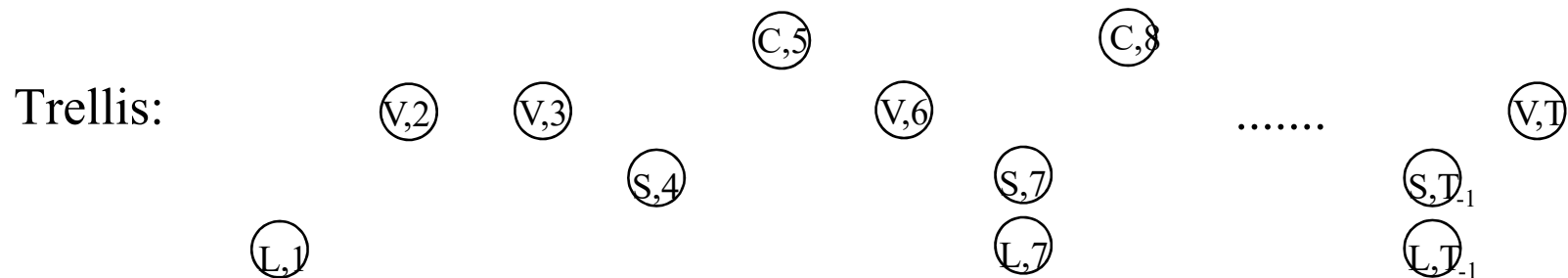
$$N(i) = \sum_{s \in S} \alpha(s, i)$$

as follows:

- **compute $\alpha(s, i)$ as usual (Step 2 of the algorithm), computing the sum $N(i)$ at the given stage i as you go.**
- **at the end of each stage, recompute all α s (for each state s):**
 - $\alpha^*(s, i) = \alpha(s, i) / N(i)$
- **use the same $N(i)$ for β s at the end of each backward (Step 3) stage:**
 - $\beta^*(s, i) = \beta(s, i) / N(i)$

Example

- Task: pronunciation of “the”
- Solution: build HMM, fully connected, 4 states:
 - S - short article, L - long article, C,V - starting w/consonant, vowel
 - thus, only “the” is ambiguous (a, an, the - not members of C,V)
- Output from states only ($p(w|s,s') = p(w|s')$)
- Data Y: an egg and a piece of the big the end



Example: Initialization

- Output probabilities:
 $p_{\text{init}}(w|c) = c(c,w) / c(c)$; where $c(S,\text{the}) = c(L,\text{the}) = c(\text{the})/2$
(other than that, everything is deterministic)
- Transition probabilities:
 - $p_{\text{init}}(c'|c) = 1/4$ (uniform)
- Don't forget:
 - about the space needed
 - initialize $\alpha(X,0) = 1$ (X : the never-occurring front buffer st.)
 - initialize $\beta(s,T) = 1$ for all s (except for $s = X$)

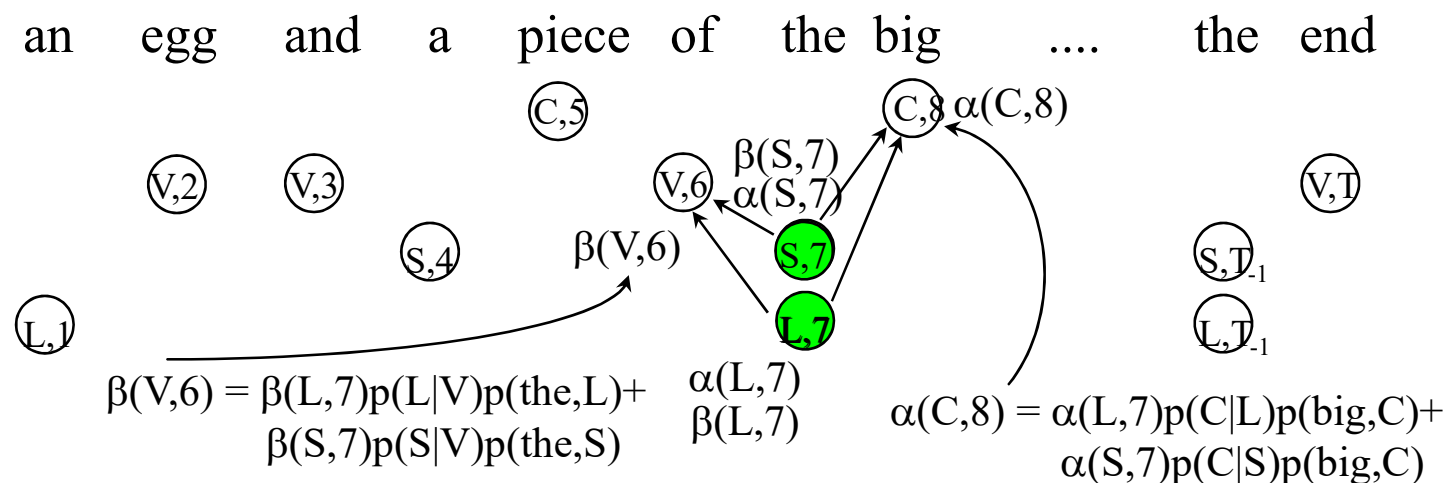
Fill in alpha, beta

- Left to right, alpha:

$$\alpha(s', i) = \sum_{s \rightarrow s'} \alpha(s, i-1) \times p(s'|s) \times p(w_i|s')$$

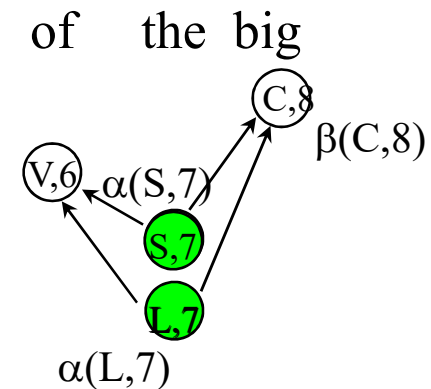
output from states

- Remember normalization (N(i)).
- Similarly, beta (on the way back from the end).



Counts & Reestimation

- One pass through data
 - At each position i , go through all pairs (s_i, s_{i+1})
 - Increment appropriate counters by frac. counts (Step 4):
 - $\text{inc}(y_{i+1}, s_i, s_{i+1}) = a(s_i, i) p(s_{i+1}|s_i) p(y_{i+1}|s_{i+1}) b(s_{i+1}, i+1)$
 - $c(y, s_i, s_{i+1}) += \text{inc}$ (for y at pos $i+1$)
 - $c(s_i, s_{i+1}) += \text{inc}$ (always)
 - $c(s_i) += \text{inc}$ (always)
- $\text{inc}(\text{big}, L, C) = \alpha(L, 7) p(C|L) p(\text{big}, C) \beta(C, 8)$
 $\text{inc}(\text{big}, S, C) = \alpha(S, 7) p(C|S) p(\text{big}, C) \beta(C, 8)$
- Reestimate $p(s'|s)$, $p(y|s)$
 - and hope for increase in $p(C|S)$ and $p(V|L)$...!!



HMM: Final Remarks

- Parameter “tying”:
 - keep certain parameters same (\sim just one “counter” for all of them)
 - any combination in principle possible
 - ex.: smoothing (just one set of lambdas)
- Real Numbers Output
 - Y of infinite size (\mathbb{R} , \mathbb{R}^n):
 - **parametric (typically: few) distribution needed (e.g., “Gaussian”)**
- “Empty” transitions: do not generate output
 - **\sim vertical arcs in trellis; do not use in “counting”**