# Introduction to
# Natural Language Processing (600.465)

# Markov Models

Dr. Jan Hajič

CS Dept., Johns Hopkins Univ.

`hajic@cs.jhu.edu`

`www.cs.jhu.edu/~hajic`

# Review: Markov Process

- Bayes formula (chain rule):

$$P(W) = P(w_1, w_2, \ldots, w_T) = \prod_{i=1..T} p(w_i | \cancel{w_1, w_2, \ldots,} w_{i-n+1}, \ldots, w_{i-1})$$

*approximation*

- n-gram language models:
  - Markov process (chain) of the order n-1:

$$P(W) = P(w_1, w_2, \ldots, w_T) = \prod_{i=1..T} p(w_i | w_{i-n+1}, w_{i-n+2}, \ldots, w_{i-1})$$

Using just <u>one</u> distribution (Ex.: trigram model: $p(w_i | w_{i-2}, w_{i-1})$):

| | | | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **Positions:** | 1 | 2 | <u>3</u> | <u>4</u> | <u>5</u> | 6 | 7 | 8 | 9 | 10 | 11 | <u>12</u> | <u>13</u> | <u>14</u> | 15 | 16 |
| **Words:** | My | car | broke down | and | within | hours | Bob | 's | car | broke down | too | . |

$$p(, | \textbf{broke down}) = p(w_5 | w_3, w_4)) = p(w_{14} | w_{12}, w_{13})$$

# Markov Properties

- Generalize to any process (not just words/LM):
  - Sequence of random variables: $X = (X_1, X_2, ..., X_T)$
  - Sample space S (*states*), size N: $S = \{s_0, s_1, s_2, ..., s_N\}$

1. Limited History (Context, Horizon):

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

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 (M.C. 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:

  1  7  3  7  9  0  | 6  7 | 3 |  4  5...
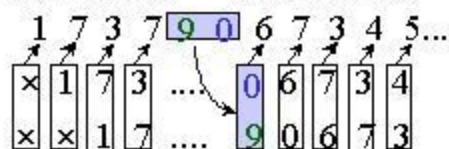
- Formally, use transformation:

  **Define new variables $Q_i$, such that $X_i = \{Q_{i-1}, Q_i\}$:**

  **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 ($X_i$):**  1  7  3  7  9  0  6  7  3  4  5...

  | × | 1 | 7 | 3 | .... | 0 | 6 | 7 | 3 | 4 |
  |---|---|---|---|------|---|---|---|---|---|
  | × | × | 1 | 7 | .... | 9 | 0 | 6 | 7 | 3 |

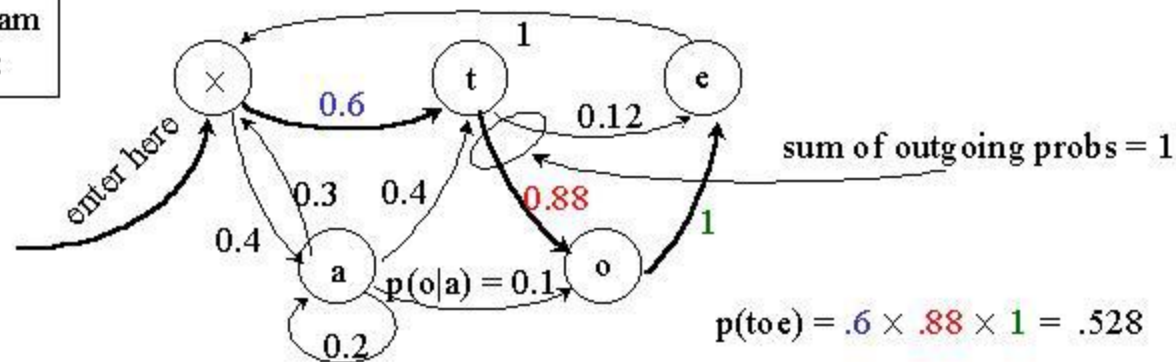  **History ($X_{i-1} = \{Q_{i-2}, Q_{i-1}\}$):**

# Graph Representation: State Diagram

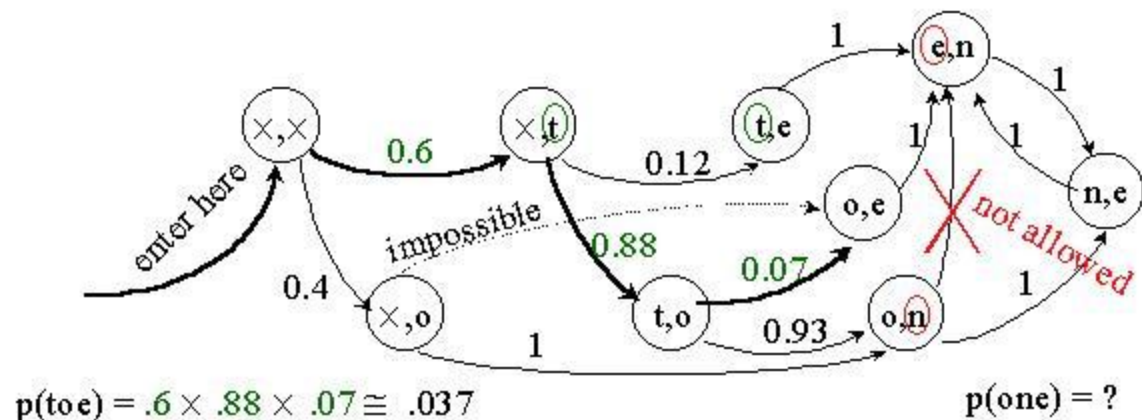- $S = \{s_0, s_1, s_2, ..., s_N\}$: states
- Distribution $P(X_i|X_{i-1})$:
  - transitions (as arcs) with probabilities attached to them:

Bigram case:



sum of outgoing probs = 1
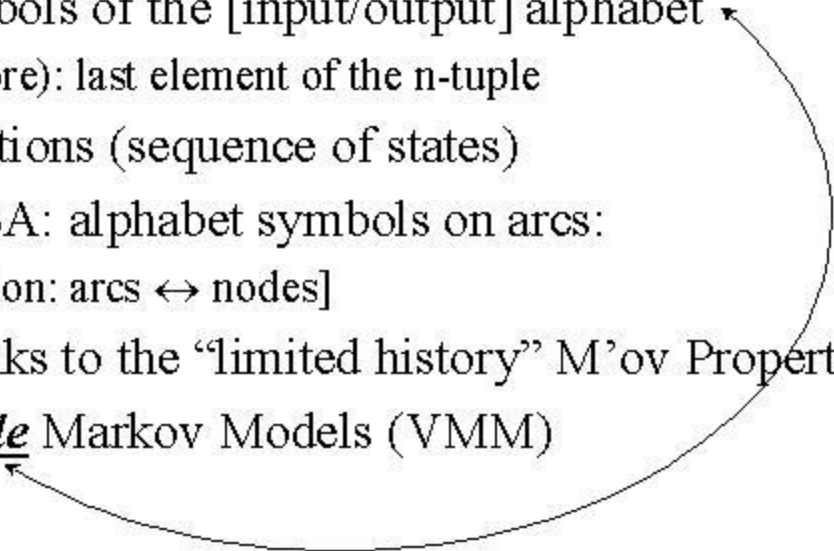
$p(o|a) = 0.1$

$p(to\,e) = .6 \times .88 \times 1 = .528$

# The Trigram Case

- $S = \{s_0, s_1, s_2, \ldots, s_N\}$: states: pairs $s_i = (x, y)$
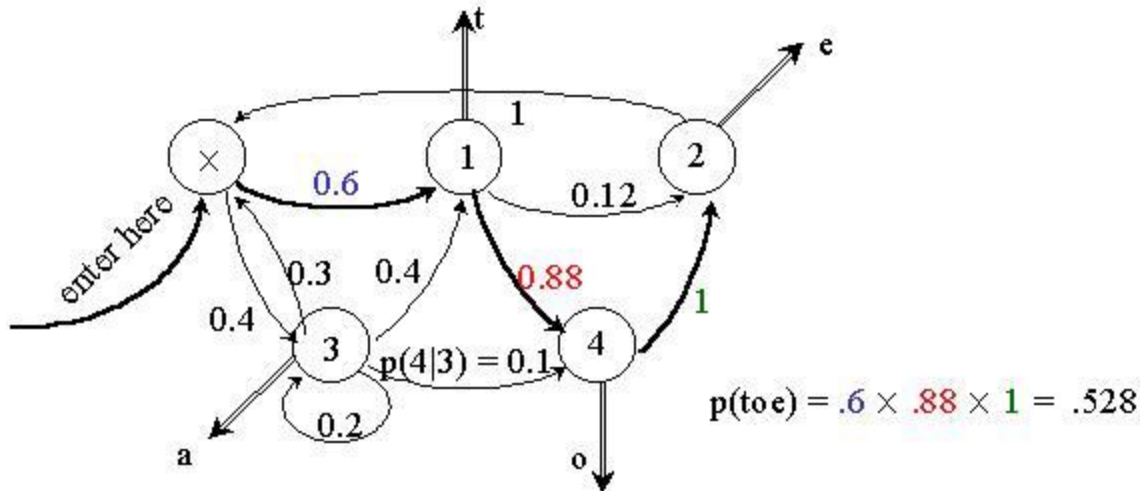- Distribution $P(X_i | X_{i-1})$: (r.v. X: generates pairs $s_i$)



$p(toe) = .6 \times .88 \times .07 \cong .037$

$p(one) = ?$

# 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:
  - transformation: arcs ↔ nodes]
- Possible thanks to the "limited history" M'ov Property
- So far: *__Visible__* Markov Models (VMM)

# Hidden Markov Models

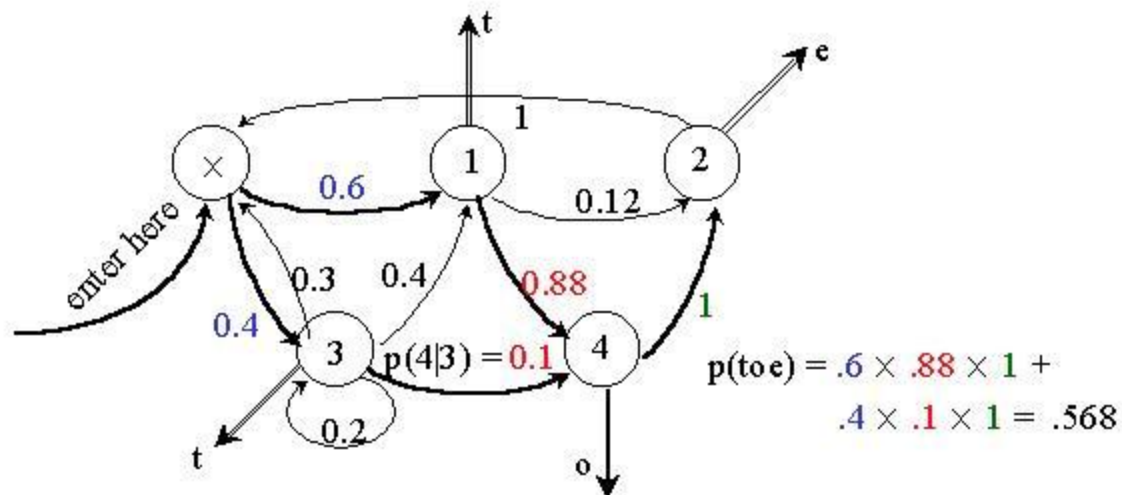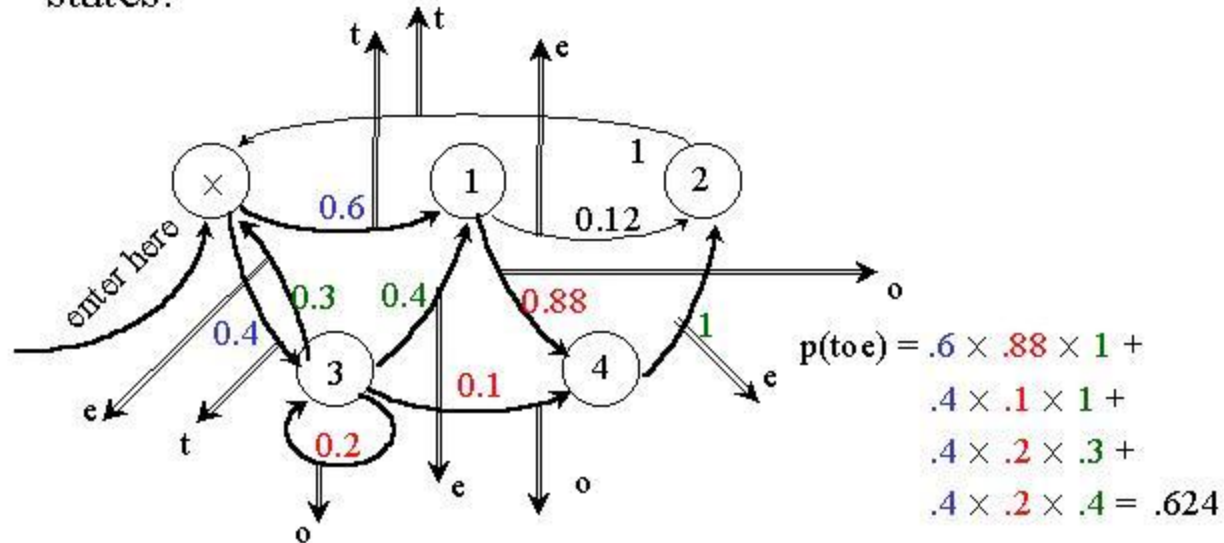- The simplest HMM: states generate [observable] output (using the "data" alphabet) but remain "invisible":



$p(to\,e) = .6 \times .88 \times 1 = .528$

# Added Flexibility

- So far, no change; but different states may generate the same output (why not?):



$p(t o e) = .6 \times .88 \times 1 +$
$.4 \times .1 \times 1 = .568$

# Output from Arcs...

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



$p(toe) = .6 \times .88 \times 1 +$
$.4 \times .1 \times 1 +$
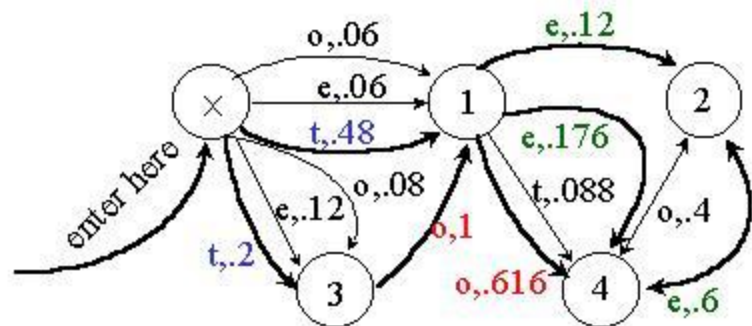$.4 \times .2 \times .3 +$
$.4 \times .2 \times .4 = .624$

# ... and Finally, Add Output Probabilities

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



!simplified!

$p(t)=.8$
$p(o)=.1$
$p(e)=.1$

$p(t)=0$
$p(o)=0$
$p(e)=1$

$p(t)=.1$
$p(o)=.7$
$p(e)=.2$

enter here

0.6

0.12

1

0.4

0.88
0.88

1

$p(t)=.5$
$p(o)=.2$
$p(e)=.3$

$p(t)=0$
$p(o)=1$
$p(e)=0$

$p(t)=0$
$p(o)=.4$
$p(e)=.6$

$\mathbf{p(toe)} = .6 \times .8 \times .88 \times .7 \times 1 \times .6 +$
$.4 \times .5 \times 1 \times 1 \times .88 \times .2 +$
$.4 \times .5 \times 1 \times 1 \times .12 \times 1$
$\cong .237$

# Slightly Different View

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



$p(toe) = .48 \times .616 \times .6 +$
$.2 \times 1 \times .176 +$
$.2 \times 1 \times .12 \cong .237$

**In the future, we will use the view more convenient for the problem at hand.**

# Formalization

- HMM (the most general case):
  - five-tuple $(S, s_0, Y, P_S, P_Y)$, where:
    - $S = \{s_0, s_1, s_2, ..., s_T\}$ is the set of states, $s_0$ is the initial state,
    - $Y = \{y_1, y_2, ..., y_V\}$ 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 $P_Y$: $|S|^2$ x $|Y|$
- Example:
  - $S = \{x, 1, 2, 3, 4\}$, $s_0 = x$
  - $Y = \{t, o, e\}$

# Formalization - Example

- Example (for graph, see foils 11,12):
  - $S = \{x, 1, 2, 3, 4\}$, $s_0 = x$
  - $Y = \{ e, o, t \}$
  - $P_S$:

| | x | 1 | 2 | 3 | 4 |
|---|---|---|---|---|---|
| x | 0 | .6 | 0 | .4 | 0 |
| 1 | 0 | 0 | .12 | 0 | .88 |
| 2 | 0 | 0 | 0 | 0 | 1 |
| 3 | 0 | 1 | 0 | 0 | 0 |
| 4 | 0 | 0 | 1 | 0 | 0 |

$\longrightarrow \Sigma = 1$

$P_Y$:

| e | x | 1 | 2 | 3 | 4 |
|---|---|---|---|---|---|

| o | x | 1 | 2 | 3 | 4 |
|---|---|---|---|---|---|

| t | x | 1 | 2 | 3 | 4 |
|---|---|---|---|---|---|
| x | | .8 | | .5 | |
| 1 | | | 0 | | .1 |
| 2 | | | | 0 | |
| 3 | | 0 | | | |
| 4 | | | 0 | | |

$\Sigma = 1$, .2, .7

# Using the HMM

- The generation algorithm (of limited value :-)):

    1. Start in $s = s_0$.

    2. Move from s to s' with probability $P_S(s'|s)$.

    3. Output (emit) symbol $y_k$ with probability $P_S(y_k|s,s')$.

    4. Repeat from step 2 (until somebody says enough).

- More interesting usage:

    – Given an output sequence $Y = \{y_1,y_2,...,y_k\}$, compute its probability.

    – Given an output sequence $Y = \{y_1,y_2,...,y_k\}$, compute the most likely sequence of states which has generated it.

    – ...plus variations: e.g., n best state sequences

# Introduction to
# Natural Language Processing (600.465)

# HMM Algorithms: Trellis and Viterbi

Dr. Jan Hajič

CS Dept., Johns Hopkins Univ.
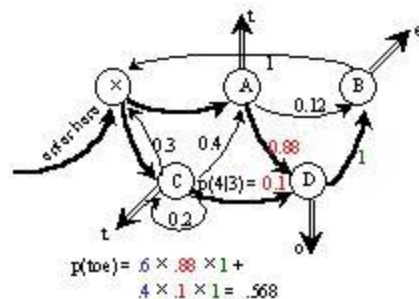
`hajic@cs.jhu.edu`

`www.cs.jhu.edu/~hajic`

# HMM: The Two Tasks

- HMM (the general case):
    - five-tuple $(S, S_0, Y, P_S, P_Y)$, where:
        - $S = \{s_1, s_2, ..., s_T\}$ is the set of states, $S_0$ is the initial state,
        - $Y = \{y_1, y_2, ..., y_V\}$ 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 $Y = \{y_1, y_2, ..., y_k\}$:

    (Task 1) compute the probability of Y;

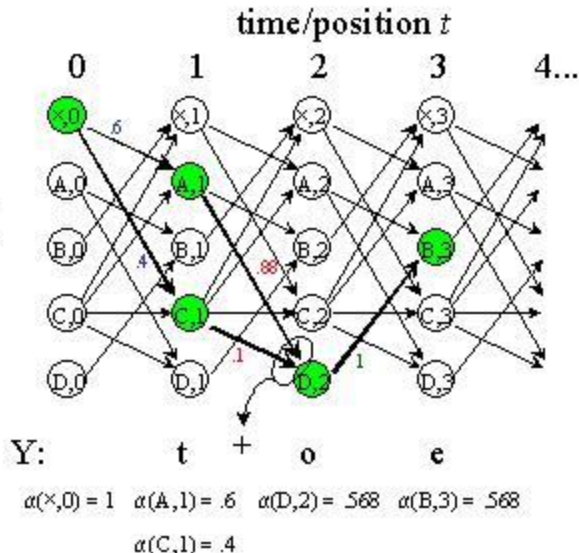    (Task 2) compute the most likely sequence of states which has generated Y.

# Trellis - Deterministic Output

HMM:

Trellis:
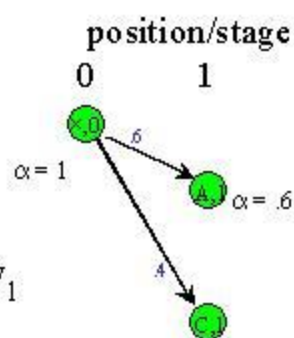
time/position $t$

"rollout"



- trellis state: (HMM state, position)
- each state: holds **_one_** number (prob): $\alpha$
- probability or Y: $\Sigma\alpha$ in the last state

Y:       t  +  o       e

$\alpha(\times,0) = 1$   $\alpha(A,1) = .6$   $\alpha(D,2) = .568$   $\alpha(B,3) = .568$

$\alpha(C,1) = .4$

# Creating the Trellis: The Start

position/stage

0          1

- 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 $P_S(state|\times)\ \underbrace{\alpha(\times,0)}_{1}$
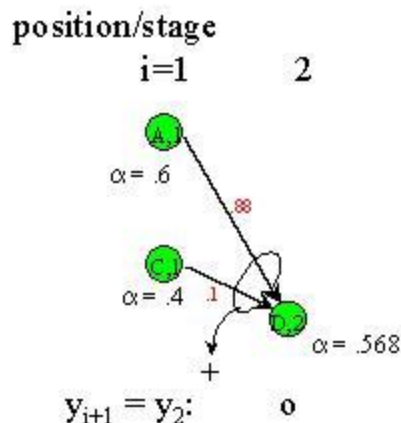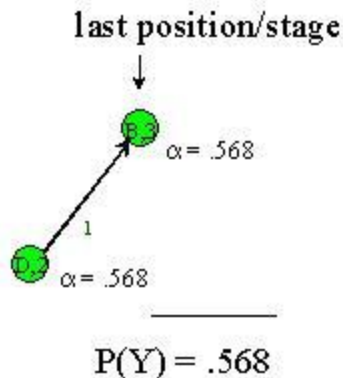- ...and forget about the $0$-th  stage

$y_1$:          t

# 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
i=1        2



$\alpha = .6$

$\alpha = .4$    .1

$\alpha = .568$

$y_{i+1} = y_2$:    o

# Trellis: The Last Step

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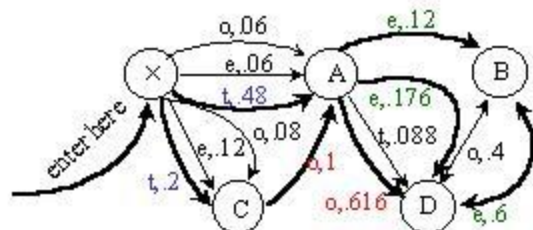  - memory usage max: 2|S|
  - multiplications max: $|S|^2|Y|$

last position/stage

↓

$\alpha = .568$

1

$\alpha = .568$

_____

P(Y) = .568

# Trellis: The General Case (still, bigrams)

- Start as usual:
  - start state (×), set its $\alpha(×,0)$ to 1.

$\alpha = 1$



$$p(toe) = .48×.616×.6+$$
$$.2×1×.176 +$$
$$.2×1×.12 \cong .237$$

# General Trellis: The Next Step

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

position/stage



$y_1$: t

...and forget about stage $i$ as usual.

# Trellis: The Complete Example

Stage:



$y_1$: **t**     $y_2$: **o**     $y_3$: **e**

$\alpha = 1$     $\alpha = .48$     $\alpha = .2$     $\alpha = .2$     $\alpha \cong .29568$

$\alpha = .024 + .177408 = .201408$

$+$

$= .035200$
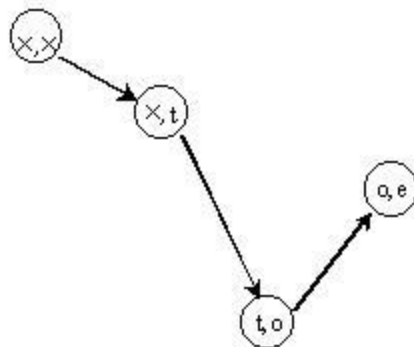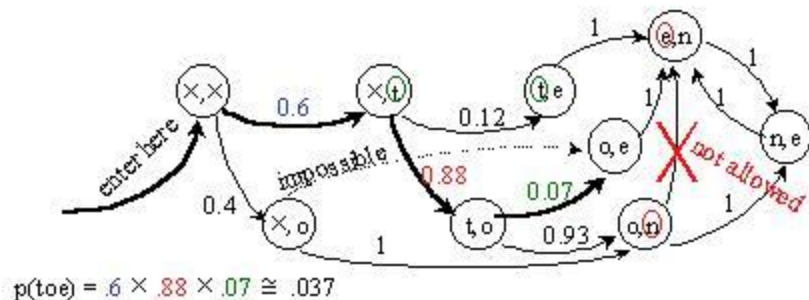
$P(Y) = P(toe) = .236608$
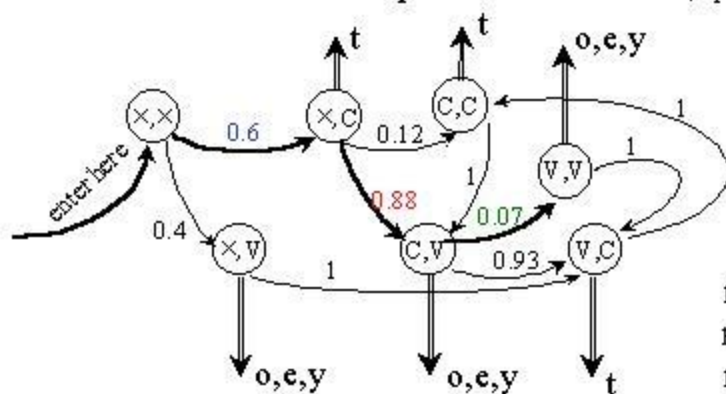
# The Case of Trigrams

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



$p(toe) = .6 \times .88 \times .07 \cong .037$

Multiple paths not possible $\rightarrow$ trellis not really needed

# Trigrams with Classes

- More interesting:
  - n-gram class LM: $p(w_i|w_{i-2},w_{i-1}) = p(w_i|c_i) \, p(c_i|c_{i-2},c_{i-1})$
    - $\rightarrow$ states are pairs of classes $(c_{i-1},c_i)$, and emit "words":
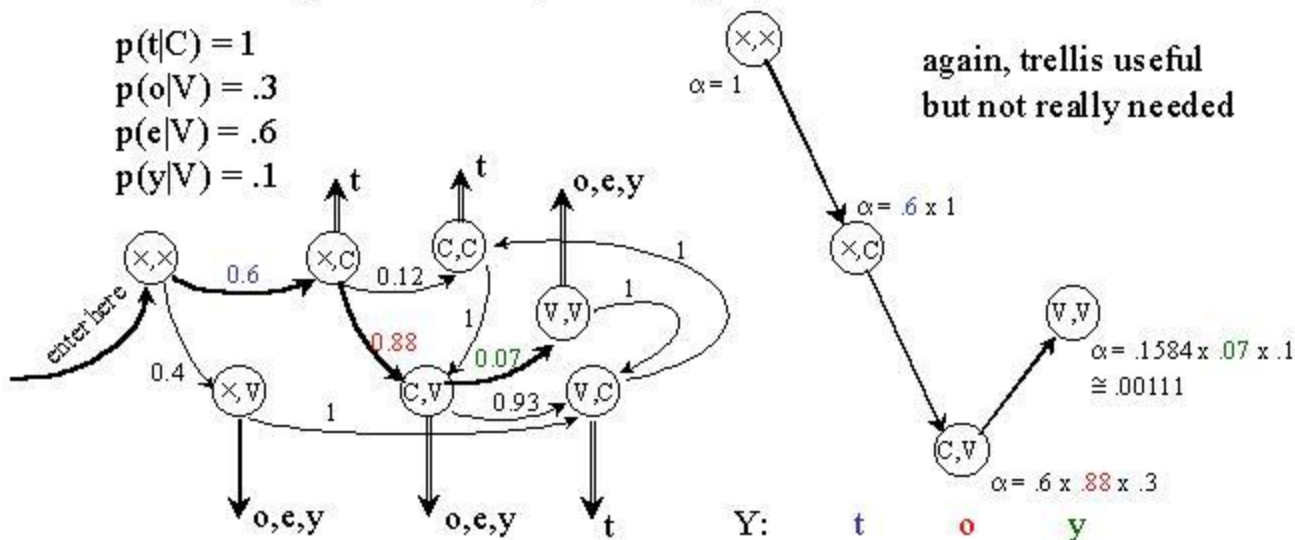
(letters in our example)



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

$p(toe) = .6 \times 1 \times .88 \times .3 \times .07 \times .6 \cong .00665$
$p(teo) = .6 \times 1 \times .88 \times .6 \times .07 \times .3 \cong .00665$
$p(toy) = .6 \times 1 \times .88 \times .3 \times .07 \times .1 \cong .00111$
$p(tty) = .6 \times 1 \times .12 \times 1 \times 1 \times .1 \cong .0072$

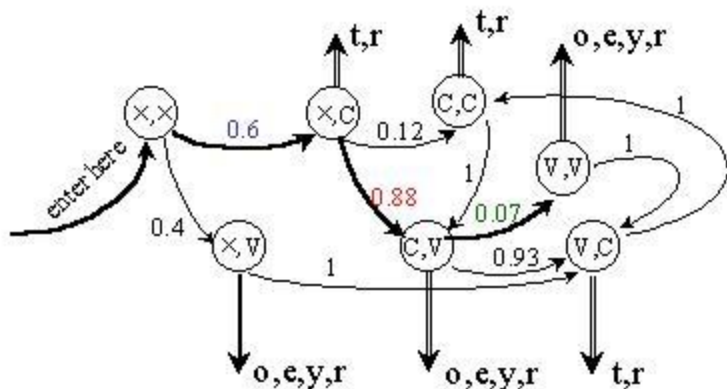# Class Trigrams: the Trellis

- Trellis generation (Y = "toy"):



p(t|C) = 1
p(o|V) = .3
p(e|V) = .6
p(y|V) = .1

again, trellis useful
but not really needed

# Overlapping Classes

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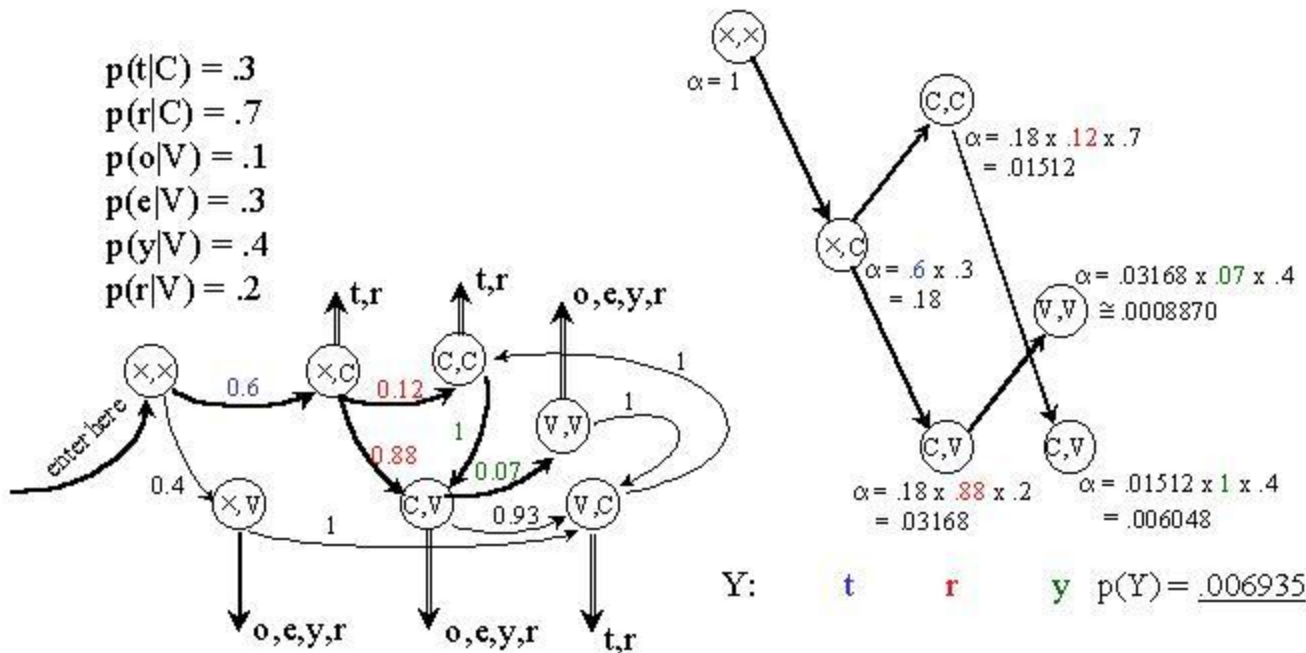  - e.g. 'r' is sometimes vowel sometimes consonant, belongs to V as well as C:



$p(t|C) = .3$
$p(r|C) = .7$
$p(o|V) = .1$
$p(e|V) = .3$
$p(y|V) = .4$
$p(r|V) = .2$

$p(try) = ?$

# Overlapping Classes: Trellis Example



$p(t|C) = .3$
$p(r|C) = .7$
$p(o|V) = .1$
$p(e|V) = .3$
$p(y|V) = .4$
$p(r|V) = .2$

$\alpha = 1$

$\alpha = .18 \times .12 \times .7$
$= .01512$

$\alpha = .6 \times .3$
$= .18$

$\alpha = .03168 \times .07 \times .4$
$\cong .0008870$

$\alpha = .18 \times .88 \times .2$
$= .03168$

$\alpha = .01512 \times 1 \times .4$
$= .006048$

enter here

0.6    0.12    1    0.88    0.07    0.4    1    0.93    1    1

t,r    t,r    o,e,y,r    o,e,y,r    o,e,y,r    t,r

Y:      t          r          y   p(Y) = .006935

# Trellis: Remarks

- So far, we went left to right (computing $\alpha$)
- Same result: going right to left (computing $\beta$)
  - supposed we know where to start (finite data)
- In fact, we might start in the middle going left <u>and</u> right
- Important for parameter estimation

  (Forward-Backward Algortihm 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
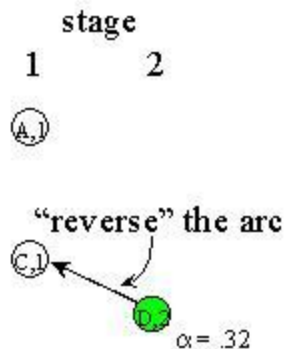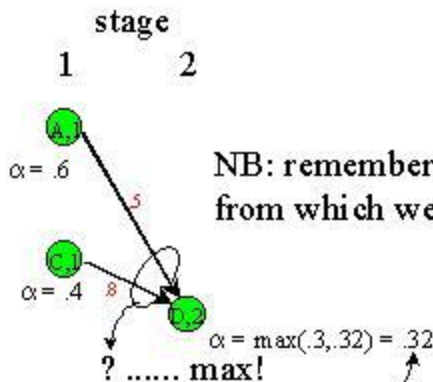
- i.e., finding

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

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

$$S_{best} = \text{argmax}_S P(S,Y) =$$
$$= \text{argmax}_S P(s_0, s_1, s_2, \ldots, s_k, y_1, y_2, \ldots, y_k) =$$
$$= \text{argmax}_S \Pi_{i=1..k}\ p(y_i|s_i, s_{i-1})p(s_i|s_{i-1})$$

# The Crucial Observation

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



NB: remember previous state from which we got the maximum:

α = .6

.5

α = .4    .8

α = max(.3,.32) = .32
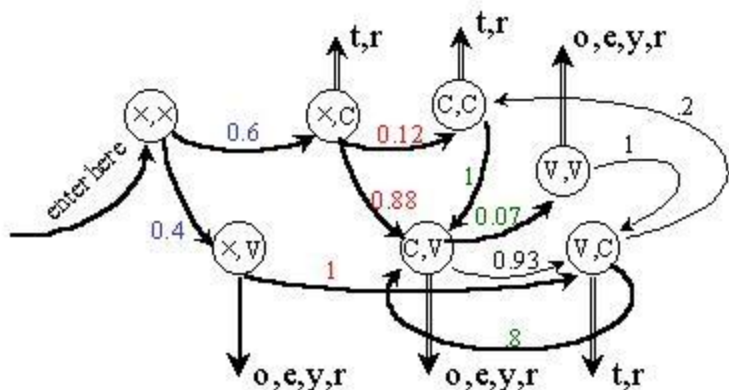
? ...... max!

"reverse" the arc

α = .32

this is certainly the "backwards" maximum to (D,2)... but
it <u>cannot change even whenever we go forward</u> (M. Property: Limited History)

# Viterbi Example

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



$p(t|C) = .3$
$p(r|C) = .7$
$p(o|V) = .1$
$p(e|V) = .3$
$p(y|V) = .4$
$p(r|V) = .2$

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

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

# Viterbi Computation



Y:    **r**    **r**    **y**

α in trellis state:
best prob from start to here

p(t|C) = .3
p(r|C) = .7
p(o|V) = .1
p(e|V) = .3
p(y|V) = .4
p(r|V) = .2

α = 1

α = .42 x .12 x .7
= .03528

α = .6 x .7
= .42

α = .07392 x .07 x .4
= .002070

α = .42 x .88 x .2
= .07392

α_c,c = .03528 x 1 x .4
= .01411

α_v,c = .056 x .8 x .4
= .01792 = α_max

α = .08 x 1 x .7
= .056

α = .4 x .2
= .08

t,r    t,r    o,e,y,r

×,×
enter here
0.6    ×,C    0.12    C,C
0.88    1
0.4    0.07
×,V    C,V    0.93    V,C
1
.2
1
.8

o,e,y,r    o,e,y,r    t,r

# n-best State Sequences



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

Y:  ×,×   r   r   y

α = 1

c,c
α = .42 x .12 x .7 = .03528

×,c
α = .6 x .7 = .42

v,v
α = .07392 x .07 x .4 = .002070

C,V
α = .42 x .88 x .2 = .07392

c,v
$\alpha_{c,c}$ = .03528 x 1 x .4 = .01411

? {

$\alpha_{V,c}$ = .056 x .8 x .4 = .01792 = $\alpha_{max}$

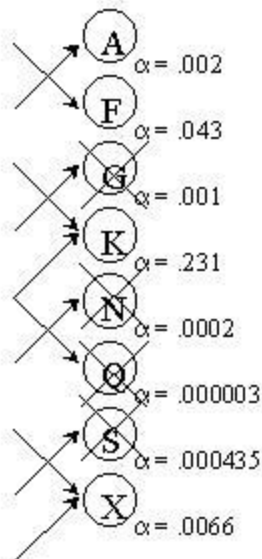v,c
α = .08 x 1 x .7 = .056

×,v
α = .4 x .2 = .08

# 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 <u>best</u> 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:



A $\alpha = .002$
F $\alpha = .043$
G $\alpha = .001$
K $\alpha = .231$
N $\alpha = .0002$
Q $\alpha = .000003$
S $\alpha = .000435$
X $\alpha = .0066$

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