

finna bless us
y = [S] V V V

Tags: "V"erb and pr"O"noun (and [S]tart)
Let's use three feature templates:

Transition features: for example $f_{VV}(x,y)$ = number of V-V transitions in y	Word-tag observation features: for example $f_{V,dog}(x,y)$ = number of tokens that are word "dog" under a Verb tag	"ends with s"-tag features: $f_{V,-s}(x,y)$ = number of tokens that end with -s and are tagged as Verb
--	---	---

(Global features have to be COUNTS: the reason why is further below.)
For 3 word vocabulary and 2 tag types, that's J=14 total features.
Assume we have fixed model weights θ and would like to score the goodness of
the above tag sequence.

Global feature vector $f(x,y) =$

f_{SV}	f_{SO}	f_{VV}	f_{VO}	f_{OV}	f_{OO}	$f_{V,finna}$	$f_{V,bless}$	$f_{V,us}$	$f_{O,finna}$	$f_{O,bless}$	$f_{O,us}$	$f_{V,-s}$	$f_{O,-s}$
1	0	2	0	0	0	1	0	1	1	0	0	2	0

Model parameters $\theta =$

θ_{SV}	θ_{SO}	θ_{VV}	θ_{VO}	θ_{OV}	θ_{OO}	$\theta_{V,finna}$	$\theta_{V,bless}$	$\theta_{V,us}$	$\theta_{O,finna}$	$\theta_{O,bless}$	$\theta_{O,us}$	$\theta_{V,-s}$	$\theta_{O,-s}$
-0.2	-0.8	+0.1	+0.5	+4.3	-0.3	-1.2	-0.1	+0.1	+5.3	-4.1	-0.3	+1.1	+2.2

Goodness score $G(y) = \theta' f(x,y) = \sum_{j=1}^J \theta_j f_j(x,y)$
 $= -0.2 + 0 + 0.2 + 0 + 0 + 0 - 1.2 + 0 + 0.1 + 5.3 + 0 + 0 + 2.2 + 0$

Global feature vector is from the sum of local feature vectors

$f(x,y) = \sum_t f_t(y_{t-1}, y_t, x_t)$
 $f_t(y_{t-1}, y_t, x_t)$ = local feature vector including the transition between these two tags,
and the observation of word at position t.

The local features are, for example:
 $f_{VV}(y_{prev}, y_{cur}, curword) = \{1 \text{ if } y_{prev}=V \text{ and } y_{cur}=V, \text{ else } 0\}$
 $f_{V,dog}(y_{prev}, y_{cur}, curword) = \{1 \text{ if } y_{cur}=V \text{ and } curword="dog", \text{ else } 0\}$
 $f_{V,-s}(y_{prev}, y_{cur}, curword) = \{1 \text{ if } y_{cur}=V \text{ and } curword \text{ ends in "s", else } 0\}$

And so on, repeated for different tags and words.

Example

.....trans. feats..... obs. feats.....

	f_{SV}	f_{SO}	f_{VV}	f_{VO}	f_{OV}	f_{OO}	$f_{V,finna}$	$f_{V,bless}$	$f_{V,us}$	$f_{O,finna}$	$f_{O,bless}$	$f_{O,us}$	$f_{V,-s}$	$f_{O,-s}$
$f(\text{START}, V, \text{finna})$	1	0	0	0	0	0	1	0	0	0	0	0	0	0
+ $f(V, V, \text{bless})$	0	0	1	0	0	0	0	1	0	0	0	0	1	0
+ $f(V, V, \text{us})$	0	0	1	0	0	0	0	0	1	0	0	0	1	0
= $f(x=\text{finna bless us}, y=V V V) =$	1	0	2	0	0	0	1	0	1	1	0	0	2	0

Local feature decomposition implies that the scoring function decomposes, too.

$G(y) = \theta' f(x,y) = \theta' \sum_t f_t(y_{t-1}, y_t, x_t) = \sum_t \theta' f_t(y_{t-1}, y_t, x_t)$
 $= \theta' f(\text{START}, V, \text{finna}) + \theta' f(V, V, \text{bless}) + \theta' f(V, V, \text{us})$

= dotprod (

-0.2	-0.8	+0.1	+0.5	+4.3	-0.3	-1.2	-0.1	+0.1	+5.3	-4.1	-0.3	+1.1	+2.2
1	0	0	0	0	0	1	0	0	0	0	0	0	0

+ dotprod (

-0.2	-0.8	+0.1	+0.5	+4.3	-0.3	-1.2	-0.1	+0.1	+5.3	-4.1	-0.3	+1.1	+2.2
0	0	1	0	0	0	0	1	0	0	0	0	1	0

+ dotprod (

-0.2	-0.8	+0.1	+0.5	+4.3	-0.3	-1.2	-0.1	+0.1	+5.3	-4.1	-0.3	+1.1	+2.2
0	0	1	0	0	0	0	0	1	0	0	0	1	0

)

In the assignment we're compiling these local feature scoring functions into A and B_t functions.

A(y_{prev}, y_{cur}) is the matrix of transition logprobs.
(If this was an HMM, A would be the able of log condprobs for possible transitions.)

$A(y_{t-1}, y_t) = \sum_{j \in \text{trans feats}} \theta_j f_j(y_{t-1}, y_t, x_t)$

B_t(y_{cur}) is the log-prob weight for tag y_{cur}, according to observation information at t.
(If this was an HMM, B_t would be the column of log condprobs for the word.)

$B_t(y_t) = \sum_{j \in \text{obs feats}} \theta_j f_j(\text{null}, y_t, x_t)$

I'm writing "null" for the prev tag because the obs features don't depend on the prev tag: only the current one.

In HW4, note we separate the local feature function into two feature functions: one for the transition feature, and one for observation features. In lecture I tried to talk a single local feature function combining transition and observation features at the position just because that's how the Chen blogpost does it.

Why'd we remove the x's from the math notation? We could've kept them if we wanted to. We just want to set up the graph weights for the Viterbi algorithm. Viterbi doesn't need to know the words. It only has to know the log-prob weights for different tags at each timestep, and the log-prob weights for different tag transitions.

Sparse vectors: The feature vectors are big long vectors where most elements are 0. In practice, you should never actually allocate an array of length J in your computer's memory. Instead, use a "sparse representation" where you use a Python dict where you have key-value pairs only for features with a nonzero value.