

# CS542 (Fall 2021) Structured Perceptron Worked Toy Example

Prof. Nikhil Krishnaswamy

## Structured Perceptrons

(You may find the discussion in Chapter 7.5.1 of the Eisenstein book helpful.)

Suppose we are given the following weight matrix  $\Theta$ :

$\Theta$	$y_i = \dots$		
	NN	VB	DT
$y_{i-1} = \langle S \rangle$	-0.3	-0.7	0.3
$y_{i-1} = \text{NN}$	-0.7	0.3	-0.3
$y_{i-1} = \text{VB}$	-0.3	-0.7	0.3
$y_{i-1} = \text{DT}$	0.3	-0.3	-0.7
$x_i = \text{Alice}$	-0.3	-0.7	0.3
$x_i = \text{admired}$	0.3	-0.3	-0.7
$x_i = \text{Dorothy}$	-0.3	0.3	-0.7
$x_i = \text{every}$	-0.7	-0.3	0.3
$x_i = \text{dwarf}$	0.3	-0.7	-0.3
$x_i = \text{cheered}$	-0.7	0.3	-0.3

(If you think this matrix looks like  $\pi$ ,  $\mathbf{A}$ , and  $\mathbf{B}$  stacked on top of each other, you are right! Note that we don't need to account for the unknown word  $\langle \text{UNK} \rangle$ , and for simplicity, we will ignore the bias term.)

## 1 Training: Viterbi

Suppose we are given the following training sentence:

Alice/NN admired/VB Dorothy/NN

Use the Viterbi algorithm to compute the best tag sequence. As part of your answer, you should fill in the Viterbi trellis below. You should also keep track of backpointers, either using arrows or in a separate table.

	Alice	admired	Dorothy
NN	•	•	•
VB	•	•	•
DT	•	•	•

## 1.1 Answer

(a) We fill in the first column in our Viterbi trellis, this time using values from the weight matrix  $\Theta$ , and adding rather than multiplying:

Viterbi( $y_i, x_i$ )	Alice	admired	Dorothy
NN	(a)	(d)	(g)
VB	(b)	(e)	(h)
DT	(c)	(f)	(i)

1.  $\Theta(<S>, NN) + \Theta(\text{Alice}, NN) = -0.3 - 0.3 = -0.6$
2.  $\Theta(<S>, VB) + \Theta(\text{Alice}, VB) = -0.7 - 0.7 = -1.4$
3.  $\Theta(<S>, DT) + \Theta(\text{Alice}, DT) = 0.3 + 0.3 = 0.6$

Then we can fill in the second column in our trellis:

1.  $\max((\text{NN}, \text{Viterbi}(\text{NN}, \text{Alice}) + \Theta(\text{NN}, \text{NN}) + \Theta(\text{admired}, \text{NN})),$   
 $(\text{VB}, \text{Viterbi}(\text{VB}, \text{Alice}) + \Theta(\text{VB}, \text{NN}) + \Theta(\text{admired}, \text{NN})),$   
 $(\text{DT}, \text{Viterbi}(\text{DT}, \text{Alice}) + \Theta(\text{DT}, \text{NN}) + \Theta(\text{admired}, \text{NN})))$   
 $= \max((\text{NN}, -0.6 - 0.7 + 0.3), (\text{VB}, -1.4 - 0.3 + 0.3), (\text{DT}, 0.6 + 0.3 + 0.3))$   
 $= \max((\text{NN}, -1.0), (\text{VB}, -1.4), (\text{DT}, 1.2))$   
 $= (\text{DT}, 1.2)$
2.  $\max((\text{NN}, \text{Viterbi}(\text{NN}, \text{Alice}) + \Theta(\text{NN}, \text{VB}) + \Theta(\text{admired}, \text{VB})),$   
 $(\text{VB}, \text{Viterbi}(\text{VB}, \text{Alice}) + \Theta(\text{VB}, \text{VB}) + \Theta(\text{admired}, \text{VB})),$   
 $(\text{DT}, \text{Viterbi}(\text{DT}, \text{Alice}) + \Theta(\text{DT}, \text{VB}) + \Theta(\text{admired}, \text{VB})))$   
 $= \max((\text{NN}, -0.6 + 0.3 - 0.3), (\text{VB}, -1.4 - 0.7 - 0.3), (\text{DT}, 0.6 - 0.3 - 0.3))$   
 $= \max((\text{NN}, -0.6), (\text{VB}, -2.4), (\text{DT}, 0.0))$   
 $= (\text{DT}, 0.0)$
3.  $\max((\text{NN}, \text{Viterbi}(\text{NN}, \text{Alice}) + \Theta(\text{NN}, \text{DT}) + \Theta(\text{admired}, \text{DT})),$   
 $(\text{VB}, \text{Viterbi}(\text{VB}, \text{Alice}) + \Theta(\text{VB}, \text{DT}) + \Theta(\text{admired}, \text{DT})),$   
 $(\text{DT}, \text{Viterbi}(\text{DT}, \text{Alice}) + \Theta(\text{DT}, \text{DT}) + \Theta(\text{admired}, \text{DT})))$   
 $= \max((\text{NN}, -0.6 - 0.3 - 0.7), (\text{VB}, -1.4 + 0.3 - 0.7), (\text{DT}, 0.6 - 0.7 - 0.7))$   
 $= \max((\text{NN}, -1.6), (\text{VB}, -1.8), (\text{DT}, -0.8))$   
 $= (\text{DT}, -0.8)$

Finally we can fill in the third column in our trellis:

1.  $\max((\text{NN}, \text{Viterbi}(\text{NN}, \text{admired}) + \Theta(\text{NN}, \text{NN}) + \Theta(\text{Dorothy}, \text{NN})),$   
 $(\text{VB}, \text{Viterbi}(\text{VB}, \text{admired}) + \Theta(\text{VB}, \text{NN}) + \Theta(\text{Dorothy}, \text{NN})),$

- $$\begin{aligned}
& (\text{DT}, \text{Viterbi}(\text{DT}, \text{admired}) + \Theta(\text{DT}, \text{NN}) + \Theta(\text{Dorothy}, \text{NN})) \\
&= \max((\text{NN}, 1.2 - 0.7 - 0.3), (\text{VB}, 0.0 - 0.3 - 0.3), (\text{DT}, -0.8 + 0.3 - 0.3)) \\
&= \max((\text{NN}, 0.2), (\text{VB}, -0.6), (\text{DT}, -0.8)) \\
&= (\text{NN}, 0.2)
\end{aligned}$$
2.  $\max((\text{NN}, \text{Viterbi}(\text{NN}, \text{admired}) + \Theta(\text{NN}, \text{VB}) + \Theta(\text{Dorothy}, \text{VB})),$   
 $(\text{VB}, \text{Viterbi}(\text{VB}, \text{admired}) + \Theta(\text{VB}, \text{VB}) + \Theta(\text{Dorothy}, \text{VB})),$   
 $(\text{DT}, \text{Viterbi}(\text{DT}, \text{admired}) + \Theta(\text{DT}, \text{VB}) + \Theta(\text{Dorothy}, \text{VB})))$   
 $= \max((\text{NN}, 1.2 + 0.3 + 0.3), (\text{VB}, 0.0 - 0.7 + 0.3), (\text{DT}, -0.8 - 0.3 + 0.3))$   
 $= \max((\text{NN}, 1.8), (\text{VB}, -0.4), (\text{DT}, -0.8))$   
 $= (\text{NN}, 1.8)$
  3.  $\max((\text{NN}, \text{Viterbi}(\text{NN}, \text{admired}) + \Theta(\text{NN}, \text{DT}) + \Theta(\text{Dorothy}, \text{DT})),$   
 $(\text{VB}, \text{Viterbi}(\text{VB}, \text{admired}) + \Theta(\text{VB}, \text{DT}) + \Theta(\text{Dorothy}, \text{DT})),$   
 $(\text{DT}, \text{Viterbi}(\text{DT}, \text{admired}) + \Theta(\text{DT}, \text{DT}) + \Theta(\text{Dorothy}, \text{DT})))$   
 $= \max((\text{NN}, 1.2 - 0.3 - 0.7), (\text{VB}, 0.0 + 0.3 - 0.7), (\text{DT}, -0.8 - 0.7 - 0.7))$   
 $= \max((\text{NN}, 0.2), (\text{VB}, -0.4), (\text{DT}, -2.2))$   
 $= (\text{NN}, 0.2)$

We put the max in the Viterbi trellis and the argmax in the backpointer trellis:

$\text{Viterbi}(y_i, x_i)$	Alice	admired	Dorothy
NN	-0.6	1.2	0.2
VB	-1.4	0.0	1.8
DT	0.6	-0.8	0.2

$\text{backpointer}(y_i, x_i)$	Alice	admired	Dorothy
NN	0	DT	NN
VB	0	DT	NN
DT	0	DT	NN

The best last tag is the argmax of the last column of the Viterbi trellis, i.e. VB. The best previous tag is then  $\text{backpointer}(\text{VB}, \text{Dorothy}) = \text{NN}$ . The best previous tag before that is then  $\text{backpointer}(\text{NN}, \text{admired}) = \text{DT}$ . Since we have reached the beginning of the sentence, the best tag sequence is:

Alice/DT admired/NN Dorothy/VB

## 2 Training: Weight Update

Update the weight matrix. Use a constant learning rate  $\eta = 1$ .

### 2.1 Answer

(b) The correct tag sequence is:

Alice/NN admired/VB Dorothy/NN

But the predicted tag sequence is:

Alice/DT admired/NN Dorothy/VB

Therefore, we increment the weights for features in the correct sequence:

$$\Theta(\langle S \rangle, \text{NN}) = -0.3 + 1 = 0.7$$

$$\Theta(\text{Alice}, \text{NN}) = -0.3 + 1 = 0.7$$

$$\Theta(\text{NN}, \text{VB}) = 0.3 + 1 = 1.3$$

$$\Theta(\text{admired}, \text{VB}) = -0.3 + 1 = 0.7$$

$$\Theta(\text{VB}, \text{NN}) = -0.3 + 1 = 0.7$$

$$\Theta(\text{Dorothy}, \text{NN}) = -0.3 + 1 = 0.7$$

And decrement the weights for features in the predicted sequence:

$$\Theta(\langle S \rangle, \text{DT}) = 0.3 - 1 = -0.7$$

$$\Theta(\text{Alice}, \text{DT}) = 0.3 - 1 = -0.7$$

$$\Theta(\text{DT}, \text{NN}) = 0.3 - 1 = -0.7$$

$$\Theta(\text{admired}, \text{NN}) = 0.3 - 1 = -0.7$$

$$\Theta(\text{NN}, \text{VB}) = 1.3 - 1 = 0.3$$

$$\Theta(\text{Dorothy}, \text{VB}) = 0.3 - 1 = -0.7$$

Note that  $\Theta(\text{NN}, \text{VB})$  is a feature in both the correct sequence and the predicted sequence, and therefore was both incremented and decremented. The overall effect is that there is no change in the weight.

### 3 Testing

Suppose we are given the following testing sentence:

Alice cheered

Use the Viterbi algorithm to compute the best tag sequence. Again, you should fill in the Viterbi trellis below, and keep track of backpointers.

	Alice	cheered
NN	•	•
VB	•	•
DT	•	•

#### 3.1 Answer

We fill in our Viterbi trellis one more time, this time with our updated weights:

Viterbi( $y_i, x_i$ )	Alice	cheered
NN	(a)	(d)
VB	(b)	(e)
DT	(c)	(f)

1.  $\Theta(<S>, NN) + \Theta(\text{Alice}, NN) = 0.7 + 0.7 = 1.4$
2.  $\Theta(<S>, VB) + \Theta(\text{Alice}, VB) = -0.7 - 0.7 = -1.4$
3.  $\Theta(<S>, DT) + \Theta(\text{Alice}, DT) = -0.7 - 0.7 = -1.4$
4.  $\max((NN, \text{Viterbi}(NN, \text{Alice}) + \Theta(NN, NN) + \Theta(\text{cheered}, NN)),$   
 $(VB, \text{Viterbi}(VB, \text{Alice}) + \Theta(VB, NN) + \Theta(\text{cheered}, NN)),$   
 $(DT, \text{Viterbi}(DT, \text{Alice}) + \Theta(DT, NN) + \Theta(\text{cheered}, NN)))$   
 $= \max((NN, 1.4 - 0.7 - 0.7), (VB, -1.4 + 0.7 - 0.7), (DT, -1.4 - 0.7 - 0.7))$   
 $= \max((NN, 0.0), (VB, -1.4), (DT, -2.8))$   
 $= (NN, 0.0)$
5.  $\max((NN, \text{Viterbi}(NN, \text{Alice}) + \Theta(NN, VB) + \Theta(\text{cheered}, VB)),$   
 $(VB, \text{Viterbi}(VB, \text{Alice}) + \Theta(VB, VB) + \Theta(\text{cheered}, VB)),$   
 $(DT, \text{Viterbi}(DT, \text{Alice}) + \Theta(DT, VB) + \Theta(\text{cheered}, VB)))$   
 $= \max((NN, 1.4 + 0.3 + 0.3), (VB, -1.4 - 0.7 + 0.3), (DT, -1.4 - 0.3 + 0.3))$   
 $= \max((NN, 2.0), (VB, -1.8), (DT, -1.4))$   
 $= (NN, 2.0)$
6.  $\max((NN, \text{Viterbi}(NN, \text{Alice}) + \Theta(NN, DT) + \Theta(\text{cheered}, DT)),$   
 $(VB, \text{Viterbi}(VB, \text{Alice}) + \Theta(VB, DT) + \Theta(\text{cheered}, DT)),$   
 $(DT, \text{Viterbi}(DT, \text{Alice}) + \Theta(DT, DT) + \Theta(\text{cheered}, DT)))$   
 $= \max((NN, 1.4 - 0.3 - 0.3), (VB, -1.4 + 0.3 - 0.3), (DT, -1.4 - 0.7 - 0.3))$   
 $= \max((NN, 0.8), (VB, -1.4), (DT, -2.4))$   
 $= (NN, 0.8)$

We put the max in the Viterbi trellis and the argmax in the backpointer trellis:

$\text{Viterbi}(y_i, x_i)$	Alice	cheered
NN	1.4	0.0
VB	-1.4	2.0
DT	-1.4	0.8

$\text{backpointer}(y_i, x_i)$	Alice	cheered
NN	0	NN
VB	0	NN
DT	0	NN

The best last tag is the argmax of the last column of the Viterbi trellis, i.e. VB. The best previous tag is then  $\text{backpointer}(VB, \text{cheered}) = NN$ . Since we have reached the beginning of the sentence, the best tag sequence is:

Alice/NN cheered/VB