

On “Structured Perceptron with Inexact Search”, *NAACL 2012*

John Hewitt

CIS 700-006 : Structured Prediction for NLP

2017-09-23

Setup: tagging and perceptrons

x: The pizza runs

Setup: tagging and perceptrons

x: The pizza runs

y: *Determiner, Noun, Verb* : [D,V,N] (part of speech (POS) tagging)

Setup: tagging and perceptrons

x: The pizza runs

y: *Determiner, Noun, Verb* : [D,V,N] (part of speech (POS) tagging)

$\Phi(\mathbf{x}, [N,V,N])$: a featurization of the input sequence and output tags
Ex. previous word, previous POS, 2-character suffix

Setup: tagging and perceptrons

\mathbf{x} : The pizza runs

\mathbf{y} : *Determiner, Noun, Verb* : [D,V,N] (part of speech (POS) tagging)

$\Phi(\mathbf{x}, [N,V,N])$: a featurization of the input sequence and output tags
Ex. previous word, previous POS, 2-character suffix

$\mathbf{w}^* \Phi(\mathbf{x}, [N,V,N])$: score given to the output sequence [N,V,N] by the model for input \mathbf{x} .

Setup: tagging and perceptrons

\mathbf{x} : The pizza runs

\mathbf{y} : *Determiner, Noun, Verb* : [D,V,N] (part of speech (POS) tagging)

$\Phi(\mathbf{x}, [N,V,N])$: a featurization of the input sequence and output tags
Ex. previous word, previous POS, 2-character suffix

$\mathbf{w}^* \Phi(\mathbf{x}, [N,V,N])$: score given to the output sequence [N,V,N] by the model for input \mathbf{x} .

Perceptron update rule (*informal*): when you detect that the model makes a mistake, update the weight vector, \mathbf{w} .

Perceptron learning algorithm

1. Initialize weight vector $\mathbf{w} = \mathbf{0}$.
2. For each of i iterations (i in 1,4, 8, 16, etc.):
 - a. For each $(\mathbf{x}, \mathbf{y}) \in$ observation set \mathbf{D}

Perceptron learning algorithm

1. Initialize weight vector $\mathbf{w} = \mathbf{0}$.
2. For each of i iterations (maximum may be 1, 4, 8, 16, etc.):
 - a. For each $(\mathbf{x}, \mathbf{y}) \in$ observation set \mathbf{D}
 - i. Compute the most likely label sequence according to the model:

$$\mathbf{y}^* = \arg \max_{\mathbf{y}} \mathbf{w}^T \Phi(\mathbf{x}, \mathbf{y})$$

Perceptron learning algorithm

1. Initialize weight vector $\mathbf{w} = \mathbf{0}$.
2. For each of i iterations (maximum may be 1, 4, 8, 16, etc.):
 - a. For each $(\mathbf{x}, \mathbf{y}) \in$ observation set \mathbf{D}
 - i. Compute the most likely label sequence according to the model:

$$\mathbf{y}^* = \arg \max_{\mathbf{y}} \mathbf{w}^T \Phi(\mathbf{x}, \mathbf{y})$$

- ii. If $\mathbf{y}^* \neq \mathbf{y}$, this is considered a *violation*. Update the weight vector (to be described below.)

Perceptron learning algorithm

1. Initialize weight vector $\mathbf{w} = \mathbf{0}$.
2. For each of i iterations (maximum may be 1, 4, 8, 16, etc.):
 - a. For each $(\mathbf{x}, \mathbf{y}) \in$ observation set \mathbf{D}
 - i. Compute the most likely label sequence according to the model:

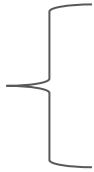
$$\mathbf{y}^* = \arg \max_{\mathbf{y}} \mathbf{w}^T \Phi(\mathbf{x}, \mathbf{y})$$

- ii. If $\mathbf{y}^* \neq \mathbf{y}$, this is considered a **violation**. Update the weight vector (to be described below.)

Foreshadowing: the *argmax* isn't always tractable. Do we actually need an *argmax*, or do we just need to find **some** incorrect label sequence that's ranked more highly than the correct sequence (a violation)?

A correct violation detection

Scoring
according
to the
model



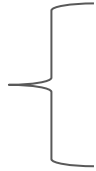
low

x: The pizza runs

high

A correct violation detection

Scoring
according
to the
model



x: The pizza runs

low $w^*\Phi(x, [N,V,N])$

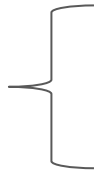
$w^*\Phi(x, [D,N,V])$

high

↑
Correct
sequence

A correct violation detection

Scoring
according
to the
model

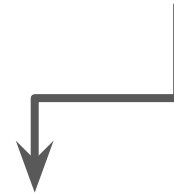


x: The pizza runs



↑
Correct
sequence

↑
Violation



Update weight vector: $\mathbf{w} \leftarrow \mathbf{w} + \Phi(x, [D,N,V]) - \Phi(x, [N,N,N])$

The problem is in inference

$$\arg \max_y w^T \Phi(\mathbf{x}, \mathbf{y}) \longrightarrow \text{low } w^* \Phi(\mathbf{x}, [N, V, N]) \quad w^* \Phi(\mathbf{x}, [D, N, V]) \quad w^* \Phi(\mathbf{x}, [N, N, N]) \quad \text{high}$$

The difficulty of solving this argmax (or *argtop_k*) depends heavily on the problem.

One consideration is the size of the space $\mathbf{Y}(\mathbf{x})$.

For some problems with only local features, $\mathbf{Y}(\mathbf{x})$ can be enumerated in its entirety.

The problem is in inference

$$\arg \max_y w^T \Phi(x, y) \longrightarrow \text{low } w^* \Phi(x, [N, V, N]) \quad w^* \Phi(x, [D, N, V]) \quad w^* \Phi(x, [N, N, N]) \quad \text{high}$$

The difficulty of solving this argmax (or argtop_k) depends heavily on the problem.

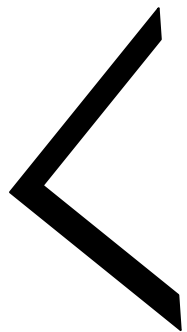
One consideration is the size of the space $\mathbf{Y}(\mathbf{x})$.

For some problems with only local features, $\mathbf{Y}(\mathbf{x})$ can be enumerated in its entirety.

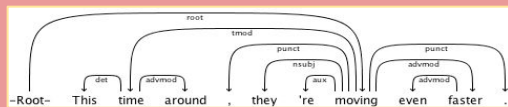
Trigram POS Tagging

45 tags
Bigram condition \rightarrow search
space is $45^3 = 91125$

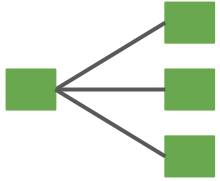
Exact Inference Possible


(difficulty)

Incremental Parsing



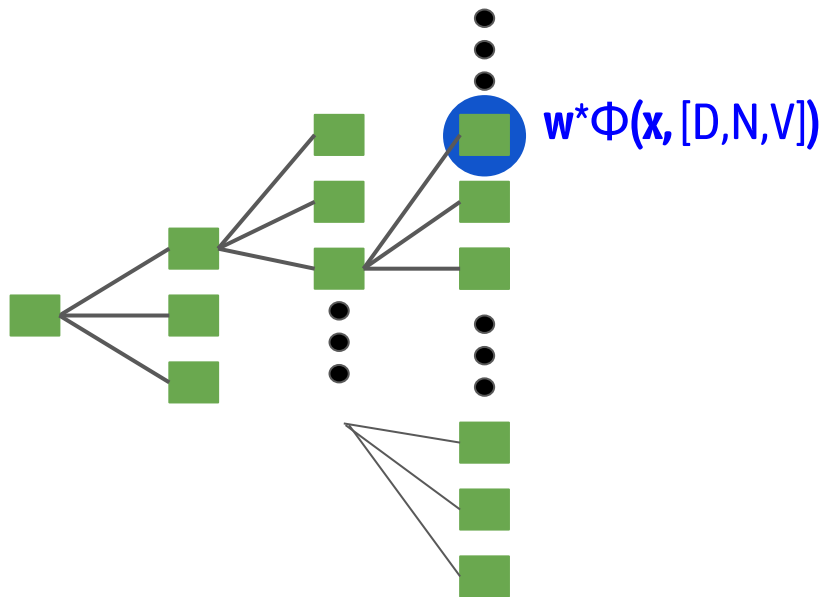
Not so much



Exact Search

Finds true argmax

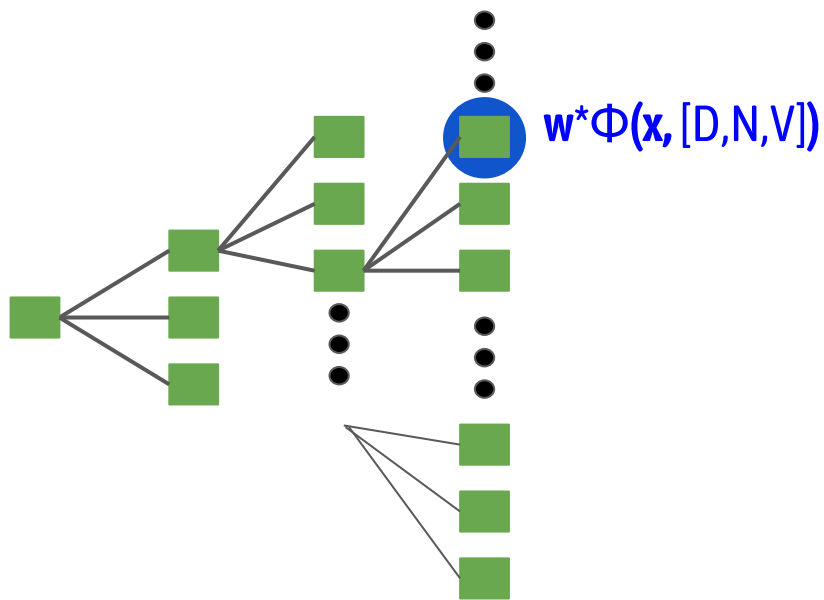
■: subsequence up to this point will be scored by the model



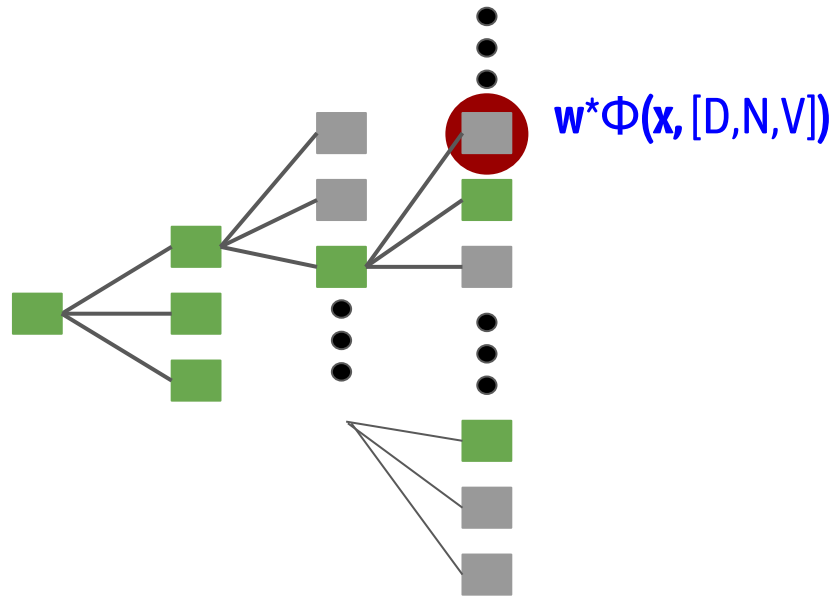
Exact Search

Finds true argmax

■: subsequence up to this point will be scored by the model



Exact Search
Finds true argmax

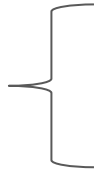


Beam Search (beam = 3)
May prune correct hypothesis due to low scores early in search

■ : subsequence up to this point will be scored by the model

An *incorrect* violation detection

Scoring
according
to the
model



x: The pizza runs

low $w^*\Phi(x, [N,V,N])$

↑
*Not a
Violation*

$w^*\Phi(x, [D,N,V])$

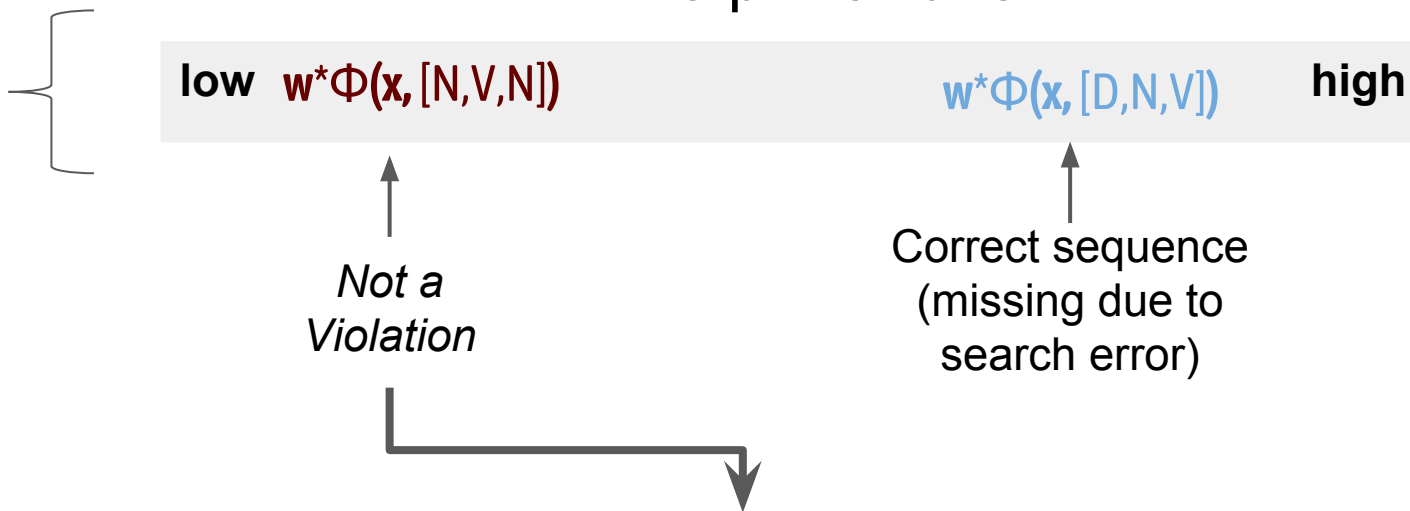
↑
Correct sequence
(missing due to
search error)

high

An *incorrect* violation detection

Scoring
according
to the
model

x: The pizza runs



We still update the weight vector, even though it was correct:

$$\mathbf{w} \leftarrow \mathbf{w} + \Phi(x, [D, N, V]) - \Phi(x, [N, N, N])$$

Perceptron Convergence

Perceptrons are proven to converge in a finite number of updates as long as they are given valid violations with which to update weight vectors. If one cannot guarantee a proposed violation is valid, convergence *may be lost**.

Perceptron Convergence

Perceptrons are proven to converge in a finite number of updates as long as they are given valid violations with which to update weight vectors. If one cannot guarantee a proposed violation is valid, convergence *may be lost**

*In this paper, it is assumed that convergence is lost without valid violations found during inference, but there are relaxations of this requirement not discussed in the paper.

Notation for structured perceptron proofs

Notation 1:

Let \mathbf{x} be an input sequence, \mathbf{y} be the correct output sequence, and \mathbf{z} be some incorrect output sequence. Then we denote the difference between the featurization of the correct hypothesis and that of \mathbf{z} to be the difference:

$$\Delta\Phi(\mathbf{x}, \mathbf{y}, \mathbf{z}) = \Phi(\mathbf{x}, \mathbf{y}) - \Phi(\mathbf{x}, \mathbf{z})$$

Notation for structured perceptron proofs

Notation 1:

Let \mathbf{x} be an input sequence, \mathbf{y} be the correct output sequence, and \mathbf{z} be some incorrect output sequence. Then we denote the difference between the featurization of the correct hypothesis and that of \mathbf{z} to be the difference:

$$\Delta\Phi(\mathbf{x}, \mathbf{y}, \mathbf{z}) = \Phi(\mathbf{x}, \mathbf{y}) - \Phi(\mathbf{x}, \mathbf{z})$$

Definition 1: Confusion Set

Let \mathbf{D} be a dataset. Let \mathbf{x} be an input sequence and $Y(\mathbf{x})$ the set of possible output sequences for \mathbf{x} . Let \mathbf{y} be the correct output sequence, and \mathbf{z} be some incorrect output sequence. Then the confusion set of \mathbf{D} is the set of triples of \mathbf{x} , \mathbf{y} , and any sequence \mathbf{z} .

$$C(\mathbf{D}) := \{(x, y, z) \mid (x, y) \in \mathbf{D}, z \in Y(x) - \{y\}\}.$$

Notation for structured perceptron proofs

Definition 2:

A dataset \mathbf{D} is said to be **linearly separable** with margin δ if there exists an oracle vector \mathbf{u} , $\|\mathbf{u}\| = 1$, such that for all \mathbf{x} , \mathbf{y} , and \mathbf{z} in \mathbf{D} , the weight vector \mathbf{u} scores the sequence \mathbf{y} at least δ better than \mathbf{z} .

$$\mathbf{u} \cdot \Delta\Phi(\mathbf{x}, \mathbf{y}, \mathbf{z}) \geq \delta$$

Notation for structured perceptron proofs

Definition 2:

A dataset \mathbf{D} is said to be **linearly separable** with margin δ if there exists an oracle vector \mathbf{u} , $\|\mathbf{u}\| = 1$, such that for all \mathbf{x} , \mathbf{y} , and \mathbf{z} in \mathbf{D} , the weight vector \mathbf{u} scores the sequence \mathbf{y} at least δ better than \mathbf{z} .

$$\mathbf{u} \cdot \Delta\Phi(\mathbf{x}, \mathbf{y}, \mathbf{z}) \geq \delta$$

Definition 3:

The **diameter**, denoted \mathbf{R} , of dataset \mathbf{D} is the largest norm of the vector difference between the featurization of any pair (\mathbf{x}, \mathbf{y}) and (\mathbf{x}, \mathbf{z}) .

$$R = \max_{(x,y,z)} \|\Delta\Phi(x, y, z)\|$$

Theorem 1: Structured Perceptron Convergence

For a dataset \mathbf{D} separable under Φ with margin δ and diameter R , the perceptron with exact search is guaranteed to converge after k updates, where

$$k \leq R^2/\delta^2$$

Theorem 1: Structured Perceptron Convergence

For a dataset \mathbf{D} separable under Φ with margin δ and diameter R , the perceptron with exact search is guaranteed to converge after k updates, where

$$k \leq R^2/\delta^2$$

Proof:

We bound the norm of w , $|w_k|$, from two directions. First:

Recall that the weight update when a violation \mathbf{z} is found is

$$\mathbf{w}_{i+1} \leftarrow \mathbf{w}_i + \Delta\Phi(\mathbf{x}, \mathbf{y}, \mathbf{z})$$

Theorem 1: Structured Perceptron Convergence

For a dataset \mathbf{D} separable under Φ with margin δ and diameter R , the perceptron with exact search is guaranteed to converge after k updates, where

$$k \leq R^2/\delta^2$$

Proof:

We bound the norm of w , $|w_k|$, from two directions. First:

Recall that the weight update when a violation \mathbf{z} is found is

$$\mathbf{w}_{i+1} \leftarrow \mathbf{w}_i + \Delta\Phi(\mathbf{x}, \mathbf{y}, \mathbf{z})$$

Let \mathbf{u} be an oracle weight vector that achieves δ separation. Then dot both sides of this equation by \mathbf{u} .

$$\begin{aligned} \mathbf{u} \cdot \mathbf{w}_{i+1} &= \mathbf{u} \cdot \mathbf{w}_i + \mathbf{u} \cdot \Delta\Phi(\mathbf{x}, \mathbf{y}, \mathbf{z}) \\ &\geq \mathbf{u} \cdot \mathbf{w}_i + \delta \end{aligned}$$

Theorem 1: Structured Perceptron Convergence

For a dataset \mathbf{D} separable under Φ with margin δ and diameter R , the perceptron with exact search is guaranteed to converge after k updates, where

$$k \leq R^2/\delta^2$$

Proof:

We bound the norm of w , $|w_k|$, from two directions. First:

Recall that the weight update when a violation \mathbf{z} is found is

$$\mathbf{w}_{i+1} \leftarrow \mathbf{w}_i + \Delta\Phi(\mathbf{x}, \mathbf{y}, \mathbf{z})$$

Let \mathbf{u} be an oracle weight vector that achieves δ separation. Then dot both sides of this equation by \mathbf{u} .

$$\begin{aligned} \mathbf{u} \cdot \mathbf{w}_{i+1} &= \mathbf{u} \cdot \mathbf{w}_i + \mathbf{u} \cdot \Delta\Phi(\mathbf{x}, \mathbf{y}, \mathbf{z}) \\ &\geq \mathbf{u} \cdot \mathbf{w}_i + \delta \end{aligned}$$

By induction, we have $\mathbf{u} \cdot \mathbf{w}_{i+1} \geq k\delta$. Further, we recall that $|\mathbf{a}||\mathbf{b}| \geq \mathbf{a} \cdot \mathbf{b}$, for any \mathbf{a}, \mathbf{b} .

$$|\mathbf{u}||\mathbf{w}_{i+1}| \geq \mathbf{u} \cdot \mathbf{w}_{i+1} \geq k\delta$$

$$|\mathbf{w}_{i+1}| \geq k\delta \text{ (since } \mathbf{u} \text{ is a unit vector.)}$$

Theorem 1: *Structured Perceptron Convergence*

For a dataset \mathbf{D} separable under Φ with margin δ and diameter R , the perceptron with exact search is guaranteed to converge after k updates, where

$$k \leq R^2/\delta^2$$

Proof:

Now we bound $|w|$ from below.

Recall that $|\mathbf{a}+\mathbf{b}|^2 = |\mathbf{a}|^2 + |\mathbf{b}|^2 + 2\mathbf{a}\cdot\mathbf{b}$

Theorem 1: Structured Perceptron Convergence

For a dataset \mathbf{D} separable under Φ with margin δ and diameter R , the perceptron with exact search is guaranteed to converge after k updates, where

$$k \leq R^2/\delta^2$$

Proof:

Now we bound $|w|$ from below.

Recall that $|\mathbf{a}+\mathbf{b}|^2 = |\mathbf{a}|^2 + |\mathbf{b}|^2 + 2\mathbf{a}\cdot\mathbf{b}$

Thus,

$$\begin{aligned} |w_{i+1}|^2 &= |w_i + \Delta\Phi(\mathbf{x}, \mathbf{y}, \mathbf{z})|^2 \\ &= |w_i|^2 + |\Delta\Phi(\mathbf{x}, \mathbf{y}, \mathbf{z})|^2 + 2w_i \cdot \Delta\Phi(\mathbf{x}, \mathbf{y}, \mathbf{z}) \end{aligned}$$

Theorem 1: Structured Perceptron Convergence

For a dataset \mathcal{D} separable under Φ with margin δ and diameter R , the perceptron with exact search is guaranteed to converge after k updates, where

$$k \leq R^2/\delta^2$$

Proof:

Now we bound $|w|$ from below.

Recall that $|a+b|^2 = |a|^2 + |b|^2 + 2a \cdot b$

Thus,

$$\begin{aligned} |w_{i+1}|^2 &= |w_i + \Delta\Phi(x, y, z)|^2 \\ &= |w_i|^2 + |\Delta\Phi(x, y, z)|^2 + 2w_i \cdot \Delta\Phi(x, y, z) \\ &\leq |w_i|^2 + \underbrace{R^2}_{\text{Replacing each term with the maximum (diameter.)}} + \underbrace{0}_{\text{By the definition of a valid violation, } w^* \Phi(x, y) < w^* \Phi(x, z).} \rightarrow \text{induction gives } |w_i| \leq kR^2 \end{aligned}$$

Replacing each term with the maximum (diameter.)

By the definition of a **valid violation**, $w^* \Phi(x, y) < w^* \Phi(x, z)$.

Theorem 1: Structured Perceptron Convergence

For a dataset \mathbf{D} separable under Φ with margin δ and diameter R , the perceptron with exact search is guaranteed to converge after k updates, where

$$k \leq R^2/\delta^2$$

Proof:

Now we bound $|w|$ from below.

Recall that $|a+b|^2 = |a|^2 + |b|^2 + 2a \cdot b$

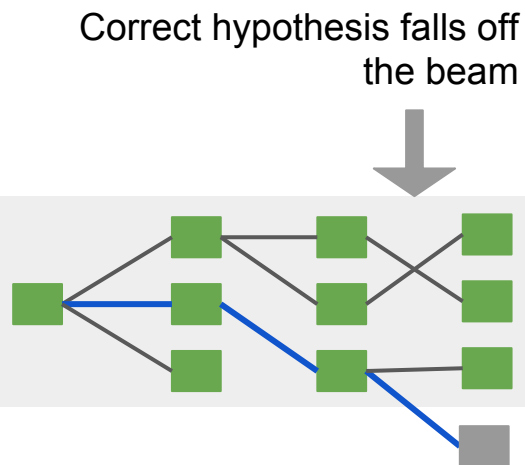
Thus,

$$\begin{aligned} |w_{i+1}|^2 &= |w_i + \Delta\Phi(x, y, z)|^2 \\ &= |w_i|^2 + |\Delta\Phi(x, y, z)|^2 + 2w_i \cdot \Delta\Phi(x, y, z) \\ &\leq |w_i|^2 + \underbrace{R^2}_{\text{Replacing each term with the maximum (diameter.)}} + \underbrace{0}_{\text{By the definition of a violation, } w^* \Phi(x, y) < w^* \Phi(x, z)} \rightarrow \text{induction gives } |w_i| \leq kR^2 \end{aligned}$$

Now, we combine our two bounds to get $k^2\delta^2 \leq |w_{i+1}| \leq kR^2$, which gives us $k \leq R^2/\delta^2$

Eliminating Invalid Updates

We want ways to ensure the validity of violations without requiring exact search.



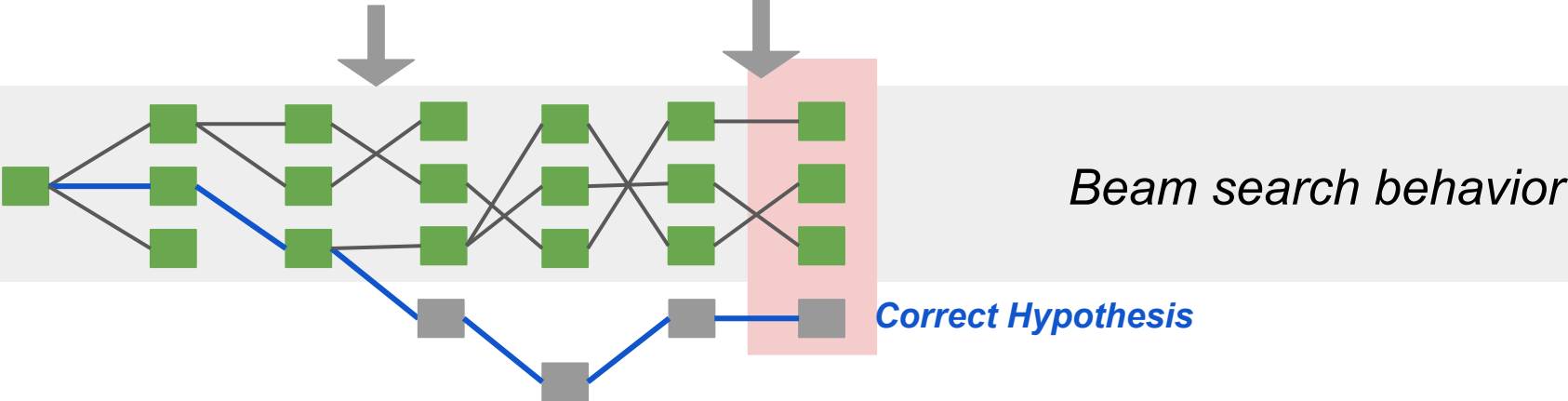
Beam search behavior

Eliminating Invalid Updates

We want ways to ensure the validity of violations without requiring exact search.

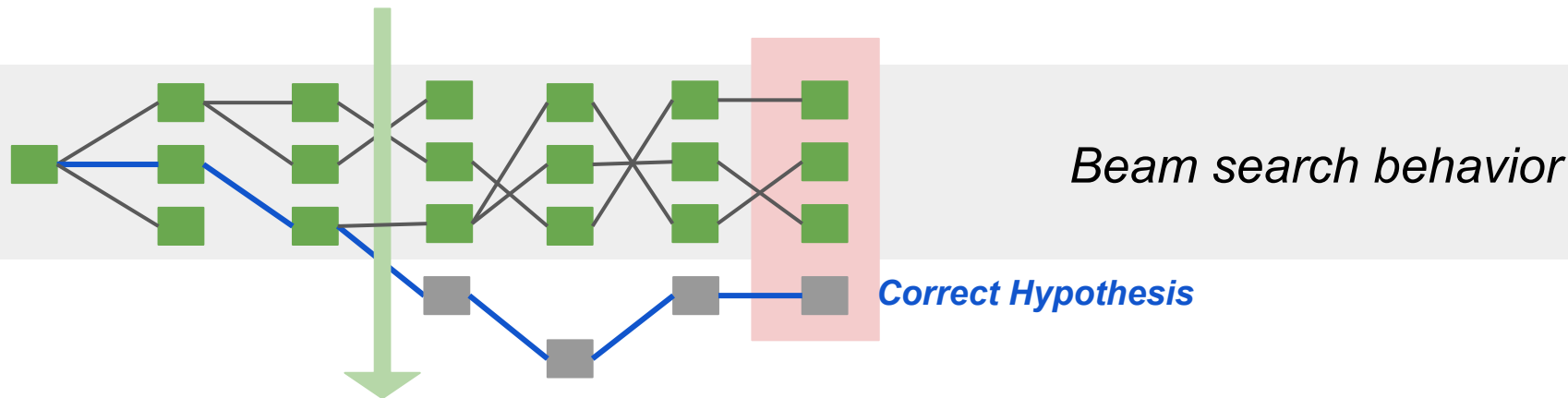
Let that from this point on, the model would have scored the correct sequence above all alternatives, making any found violation invalid

Correct hypothesis falls off the beam



Eliminating Invalid Updates: Early Update

Correct hypothesis falls off the beam. **Early Update Here!**

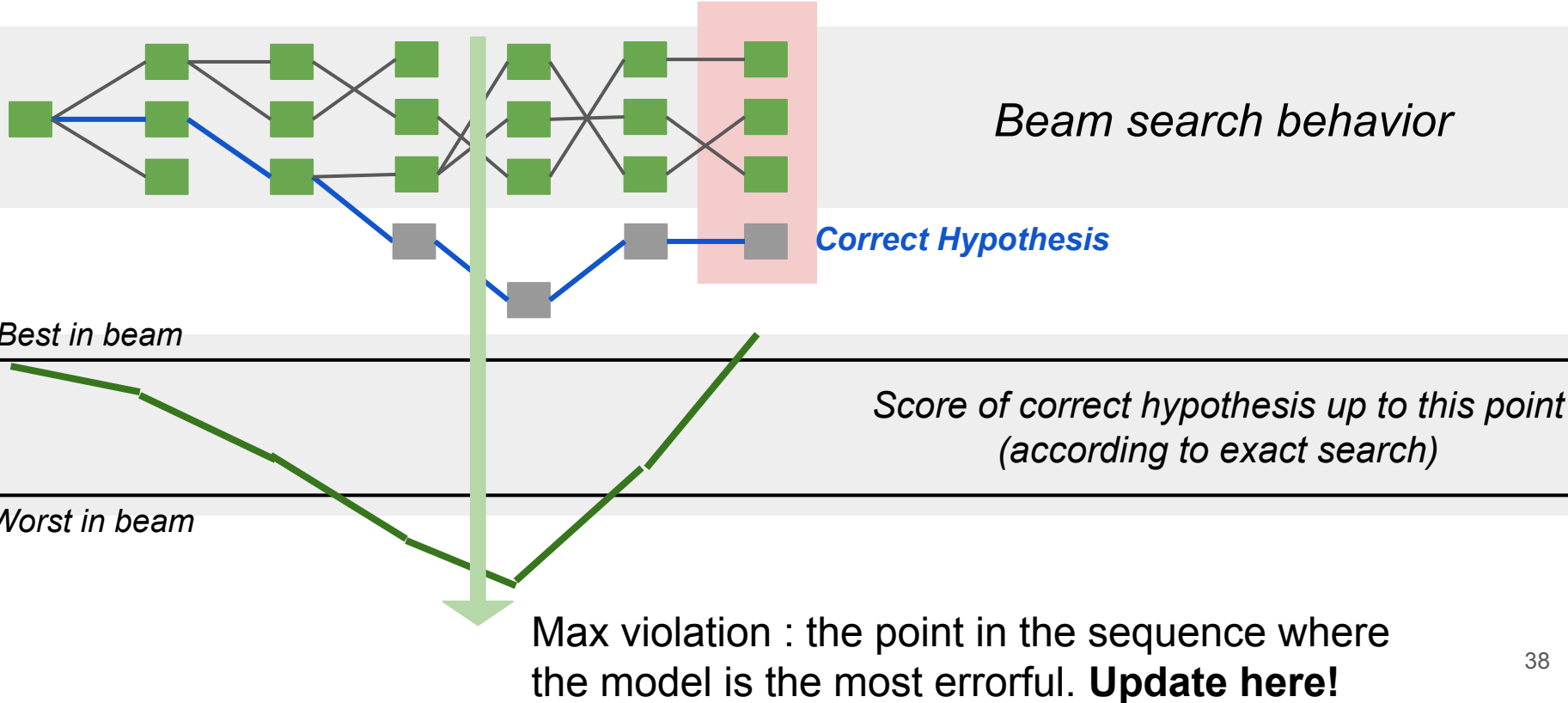


When the correct hypothesis falls off the beam, we have $\mathbf{x}_{[1:i]}$, $\mathbf{y}_{[1:i]}$, $\mathbf{z}_{[1:i]}$ such that

$$\mathbf{w}^* \Delta \Phi(\mathbf{x}_{[1:i]}, \mathbf{y}_{[1:i]}, \mathbf{z}_{[1:i]}) < 0.$$

That is, the model has clearly made an error, as there are enough subsequences $\mathbf{z}_{[1:i]}$ scored higher than $\mathbf{y}_{[1:i]}$ as to push it off the beam.

Eliminating Invalid Updates: Max Violation



Recap

Under exact search, a valid violation, if it exists, is always guaranteed to be found and used to update \mathbf{w} .

Recap

Under exact search, a valid violation, if it exists, is always guaranteed to be found and used to update \mathbf{w} .

However, *any violation prefix sequence will do*, because all satisfy $\mathbf{w}^* \Delta \Phi$
 $(\mathbf{x}_{[1:i]} \mathbf{y}_{[1:i]} \mathbf{z}_{[1:i]}) < 0$, which is the crucial condition in perceptron convergence.

Recap

Under exact search, a valid violation, if it exists, is always guaranteed to be found and used to update \mathbf{w} .

However, *any violation prefix sequence will do*, because all satisfy $\mathbf{w}^* \Delta \Phi(\mathbf{x}_{[1:i]}, \mathbf{y}_{[1:i]}, \mathbf{z}_{[1:i]}) < 0$, which is the crucial condition in perceptron convergence.

One way to find such a prefix is to take the prefix of some hypothesis on the beam as soon as the correct hypothesis falls off. This is called **early update**.

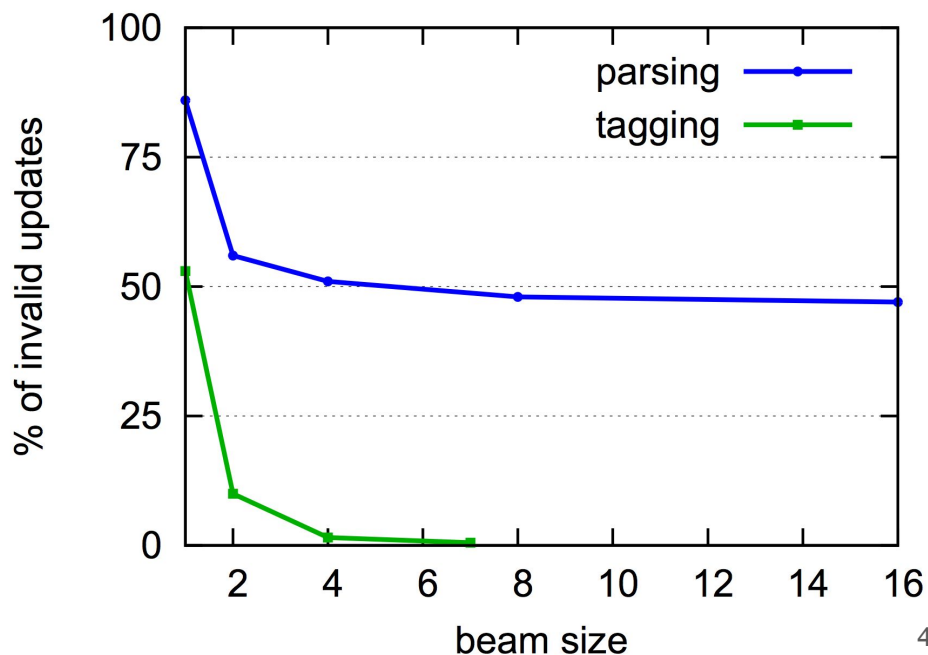
Another way to find such a prefix is to keep scoring the correct hypothesis after it falls off the beam, and take some candidate prefix on the beam at the prefix length at which the correct hypothesis is scored worst. This is called **max violation**.

Results

“I don’t believe any of this for a second.” I want to train my perceptron using full sequences \mathbf{x} , \mathbf{y} , and \mathbf{z} , and ignore potentially invalid updates.

Turns out, this isn’t a huge problem for trigram POS tagging.

For parsing, however, the story is much worse.



Results

method	<i>b</i>	it	time	dev	test
early*		38	15.4 h	92.24	92.09
standard		1	0.4 h	78.99	79.14
hybrid	8	11	5.6 h	92.26	91.95
latest		9	4.5 h	92.06	91.96
max-viol.		12	5.5 h	92.32	92.18
early	8	Huang & Sagae 2010			92.1

Parsing Results:
Ensuring valid
violations is crucial

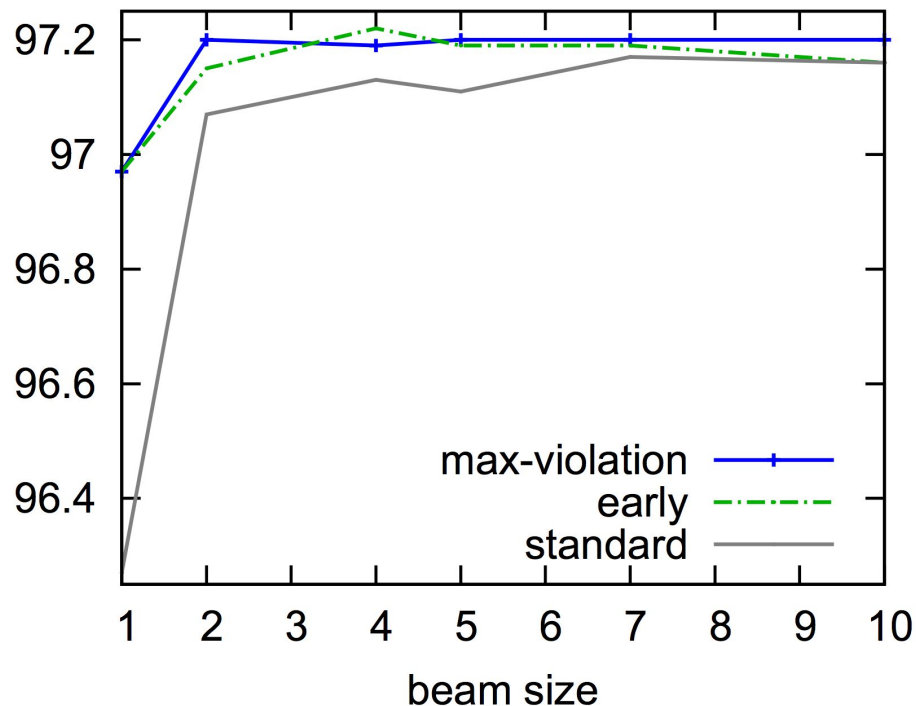
Results

method	b	it	time	dev	test
early*		38	15.4 h	92.24	92.09
standard		1	0.4 h	78.99	79.14
hybrid	8	11	5.6 h	92.26	91.95
latest		9	4.5 h	92.06	91.96
max-viol.		12	5.5 h	92.32	92.18
early	8	Huang & Sagae 2010			92.1

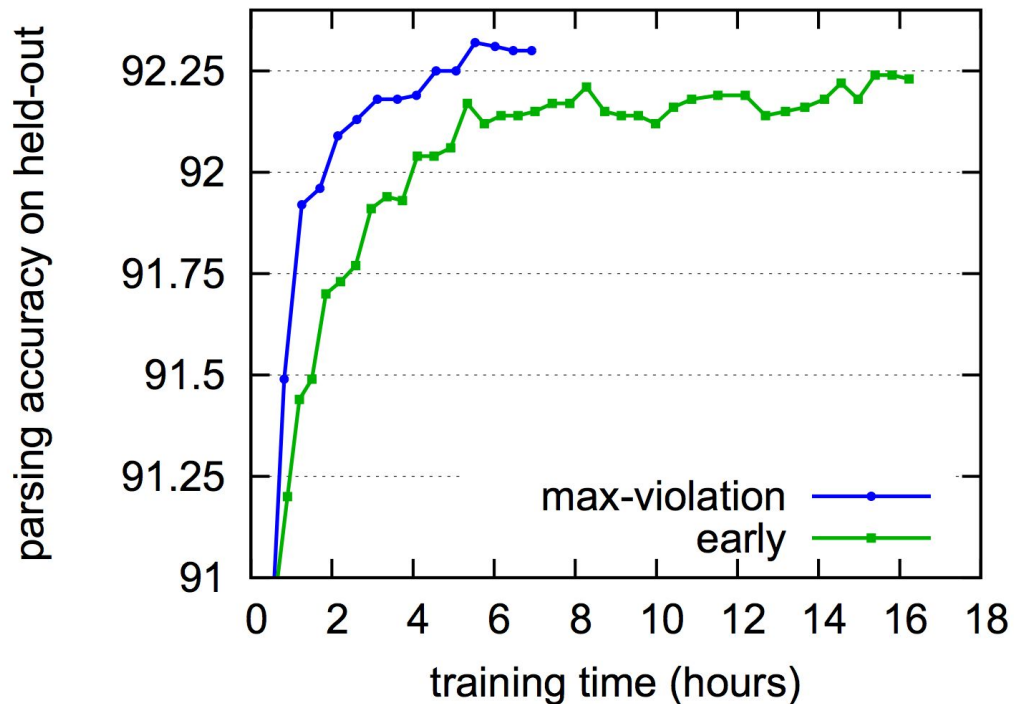
Parsing Results:
Ensuring valid
violations is crucial

Tagging Results:
Ensuring valid violations is “helpful”

best tagging accuracy on held-out



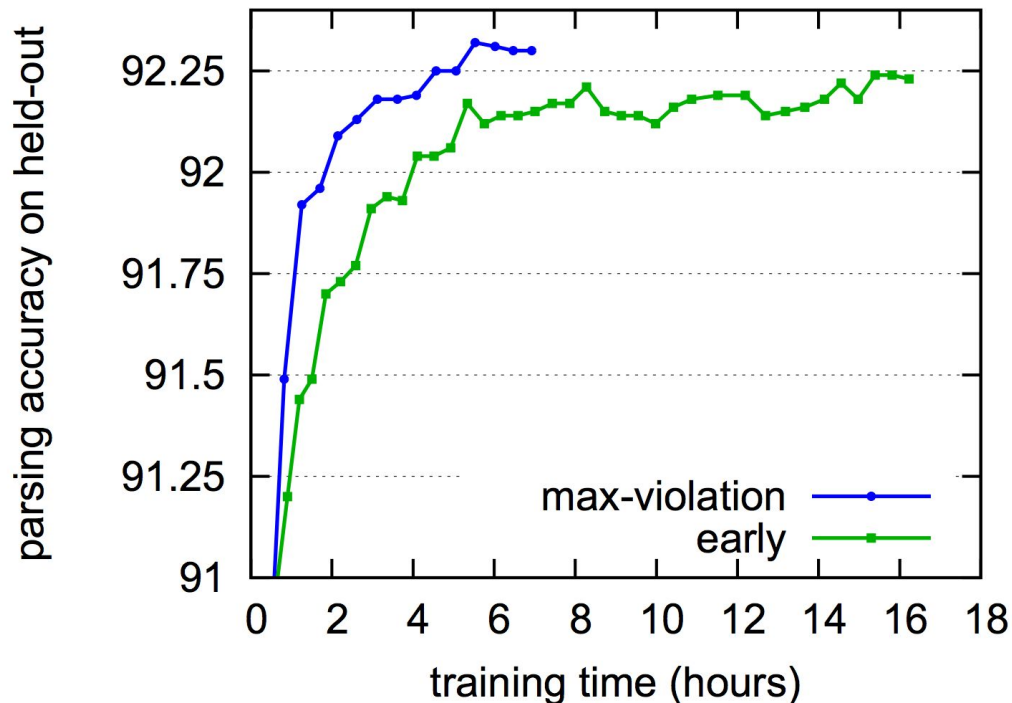
Results



Training time much faster for max-violation than early.

Why?

Results



Training time much faster for max-violation than early.

Why? Probably has to do with the margin of the dataset and informativeness of each update! Early update takes minimally informative updates.