

Machine Learning

CUNY Graduate Center, Spring 2013

Lectures 9-10: Structured Learning

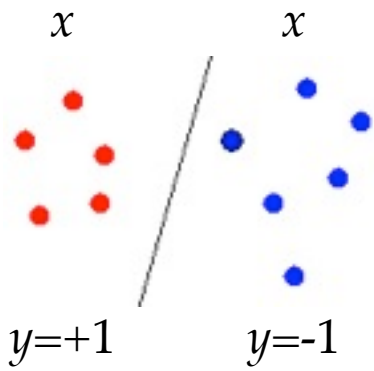
(structured perceptron, HMM, learning w/ inexact search)

Professor Liang Huang

huang@cs.qc.cuny.edu

<http://acl.cs.qc.edu/~lhuang/teaching/machine-learning>

Structured Prediction



0 1 2 3 4 5 6 7 8 9

the man bit the dog



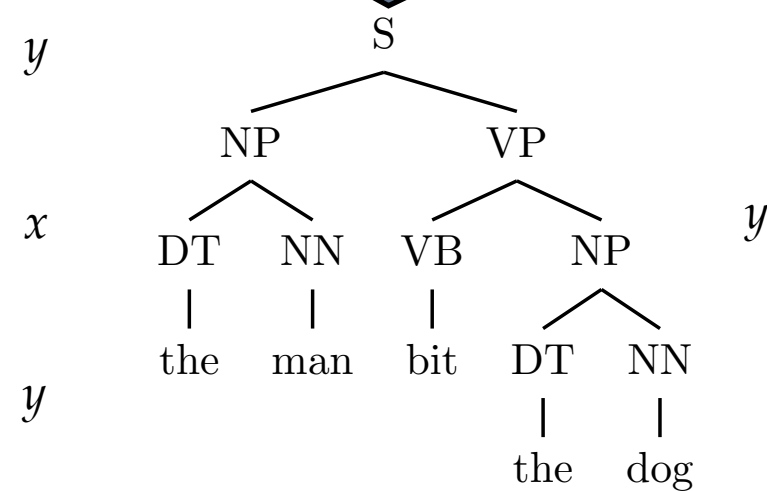
DT NN VBD DT NN

the man bit the dog



那人咬了狗

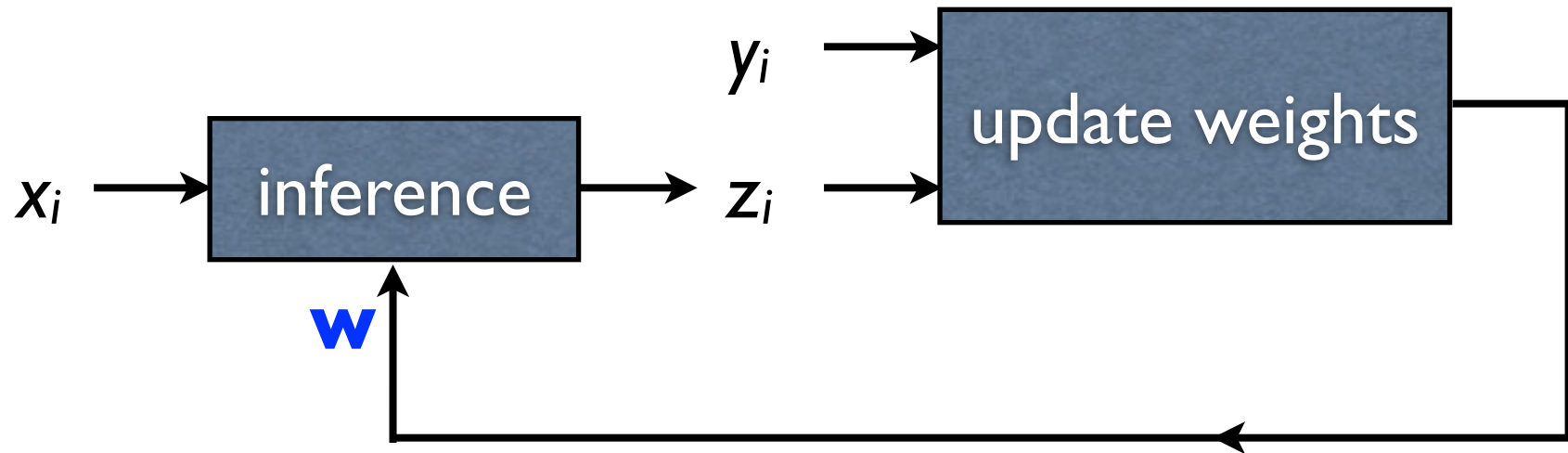
the man bit the dog



- binary classification: output is binary
- multiclass classification: output is a number (small # of classes)
- structured classification: output is a structure (seq., tree, graph)
 - part-of-speech tagging, parsing, summarization, translation
 - exponentially many classes: *search* (inference) efficiency is crucial!₂

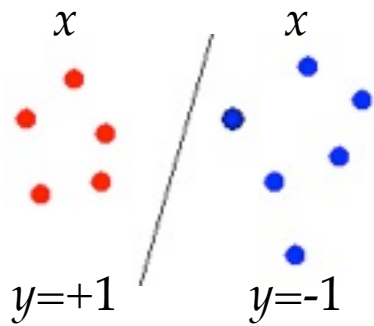
Generic Perceptron

- online-learning: one example at a time
- learning by doing
 - find the best output under the current weights
 - update weights at mistakes



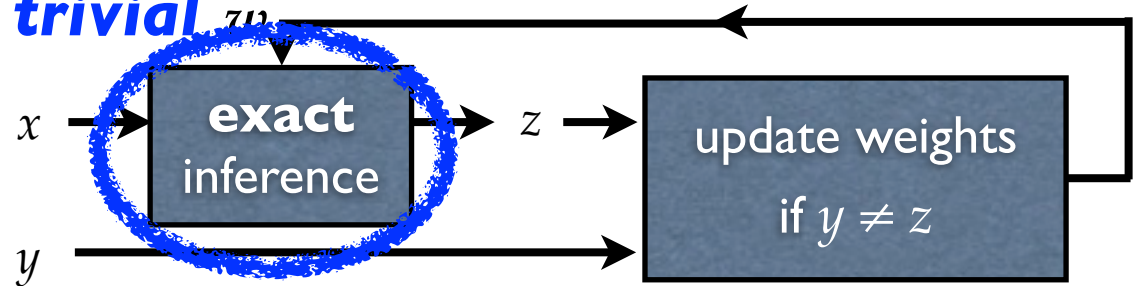
Perceptron: from binary to structured

binary classification



2 classes

trivial 2^2

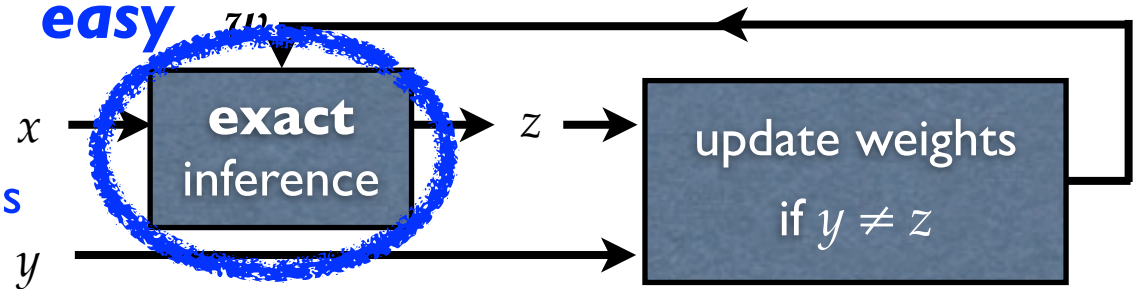


multiclass classification

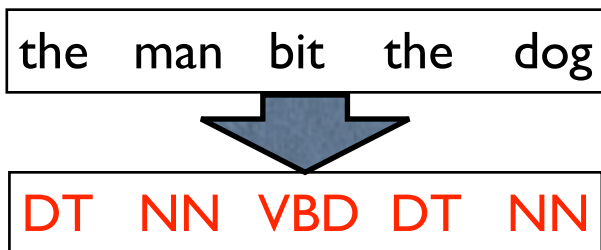


constant
of classes

easy 2^2

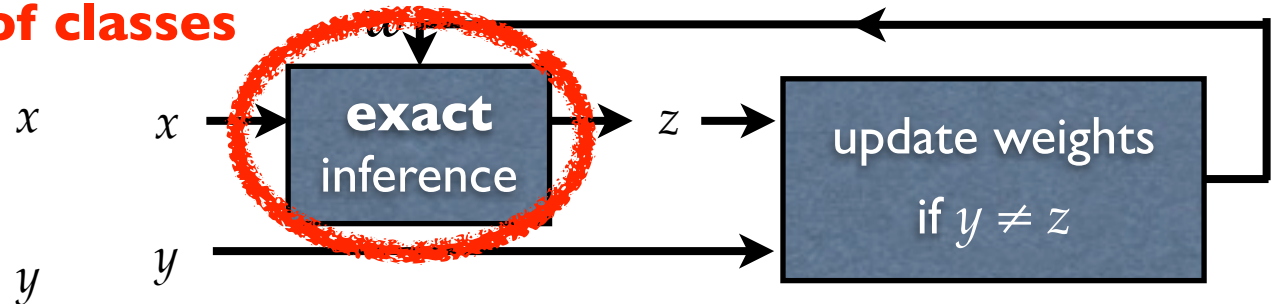


structured classification

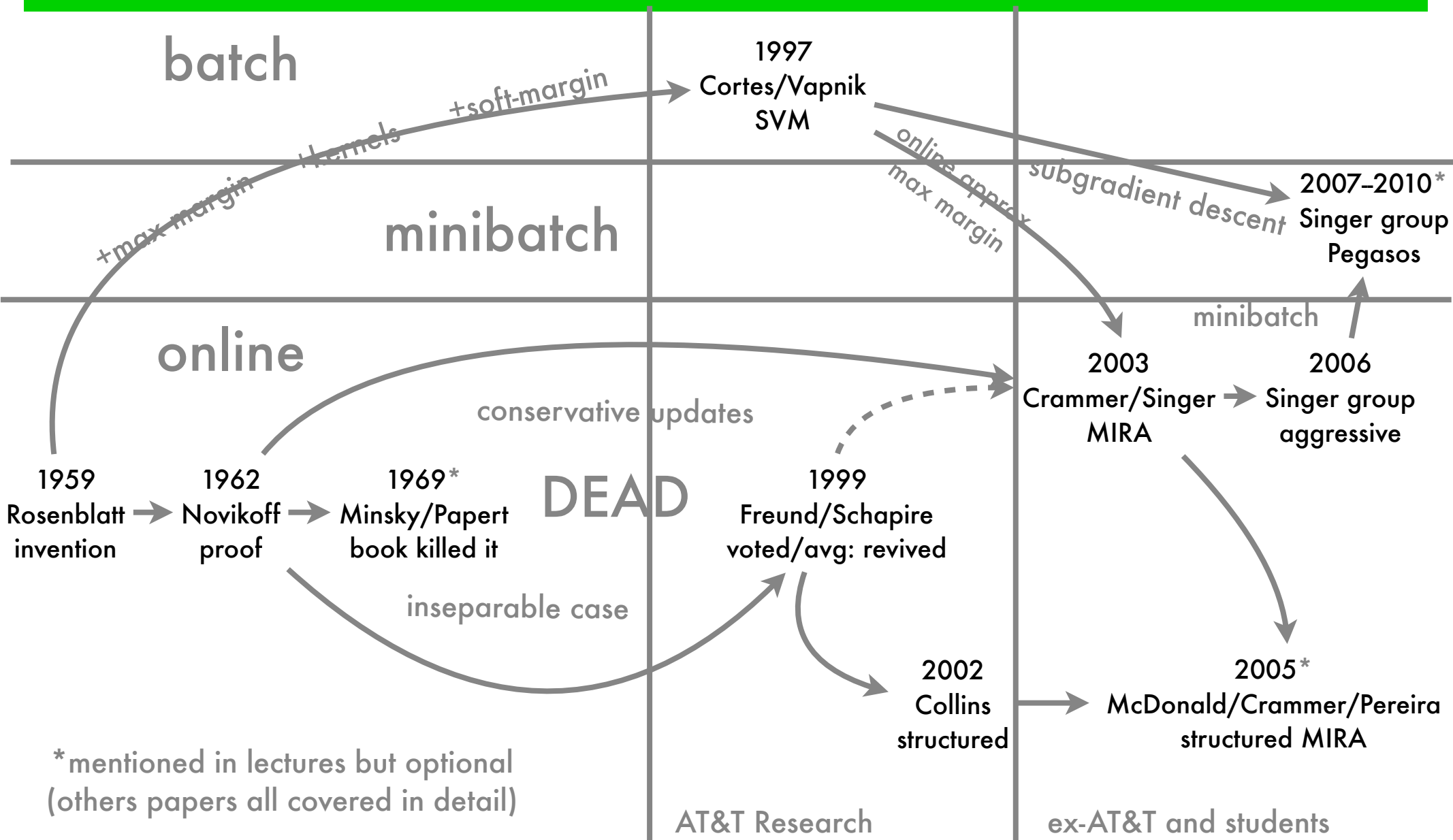


exponential
of classes

hard 2^2



Brief History of Perceptron

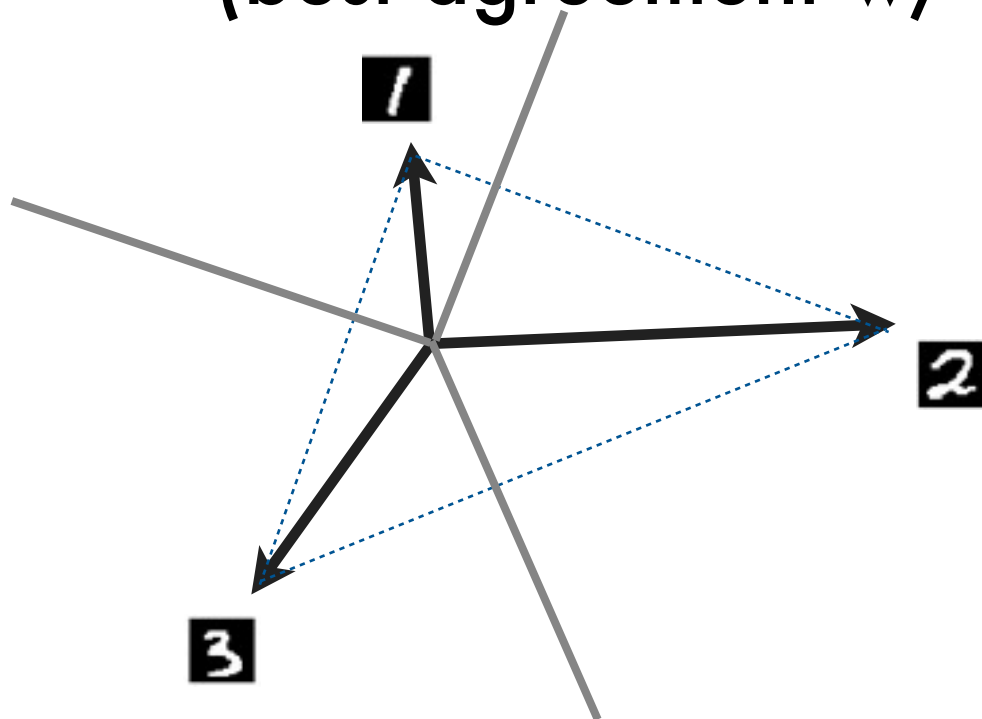


Multiclass Classification: Review

- one weight vector (“prototype”) for each class:

$$\mathbf{w} = (\mathbf{w}^{(1)}, \mathbf{w}^{(2)}, \dots, \mathbf{w}^{(M)}),$$

- multiclass decision rule: $\hat{y} = \operatorname{argmax}_{z \in 1 \dots M} w^{(z)} \cdot x$
(best agreement w/ prototype)



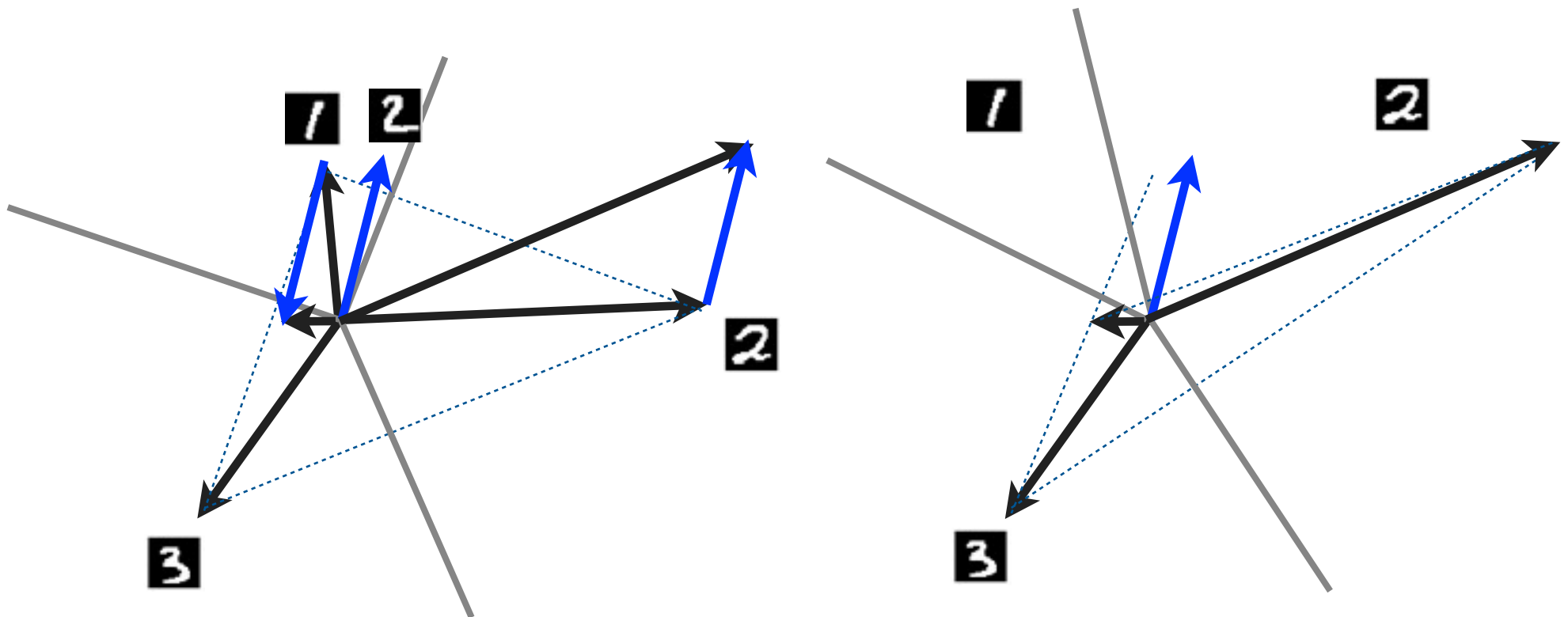
Q1: what about 2-class?

Q2: do we still need augmented space?

0 1 2 3 4 5 6 7 8 9

Multiclass Perceptron: Review

- on an error, penalize the weight for the wrong class, and reward the weight for the true class



Convergence of Multiclass

0 1 2 3 4 5 6 7 8 9

$$\mathbf{w} = (\mathbf{w}^{(1)}, \mathbf{w}^{(2)}, \dots, \mathbf{w}^{(M)}),$$

where $\mathbf{w}^{(i)}$ is used to calculate the functional margin for training example with label i ;

for a given training example \mathbf{x} and a label y , we define feature map function Φ as

$$\Phi(\mathbf{x}, y) = (\mathbf{0}^{(1)}, \dots, \mathbf{0}^{(y-1)}, \mathbf{x}, \mathbf{0}^{(y+1)}, \dots, \mathbf{0}^{(M)}).$$

such that $\mathbf{w} \cdot \Phi(\mathbf{x}, y) = \mathbf{w}^{(y)} \cdot \mathbf{x}$.

We also define that, with a given training example \mathbf{x} , the difference between two feature vectors for labels y and z as $\Delta\Phi$:

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

update rule:

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

separability:

$$\exists \mathbf{u}, \text{ s.t. } \forall (\mathbf{x}, y) \in D, z \neq y$$

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

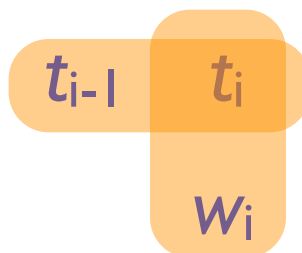
Example: POS Tagging

• gold-standard: DT NN VBD DT NN y
• the man bit the dog x $\Phi(x, y)$

• current output: DT NN NN DT NN z
• the man bit the dog x $\Phi(x, z)$

• assume only two feature classes

• tag bigrams

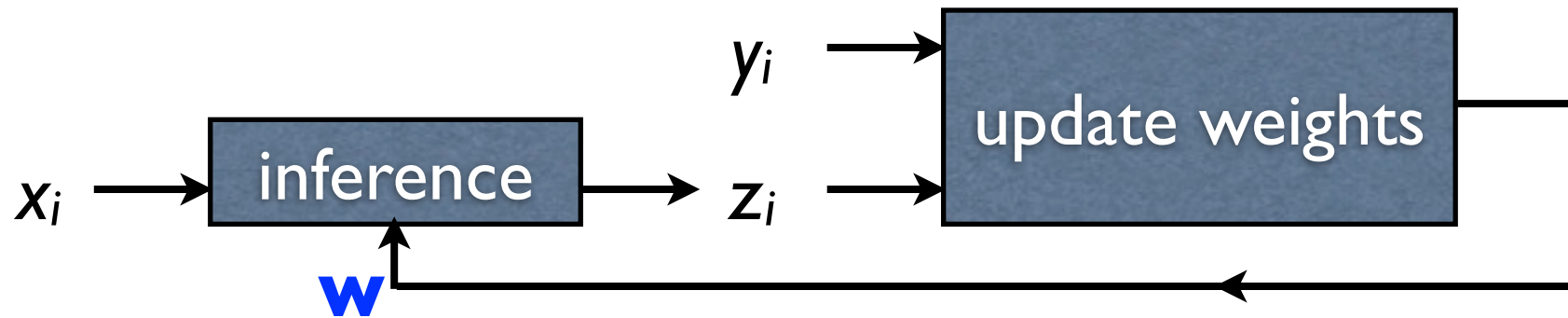


• word/tag pairs

• weights ++: (NN, VBD) (VBD, DT) (VBD \rightarrow bit)

• weights --: (NN, NN) (NN, DT) (NN \rightarrow bit)

Structured Perceptron



Inputs: Training set (x_i, y_i) for $i = 1 \dots n$

Initialization: $\mathbf{W} = 0$

Define: $F(x) = \operatorname{argmax}_{y \in \mathbf{GEN}(x)} \Phi(x, y) \cdot \mathbf{W}$

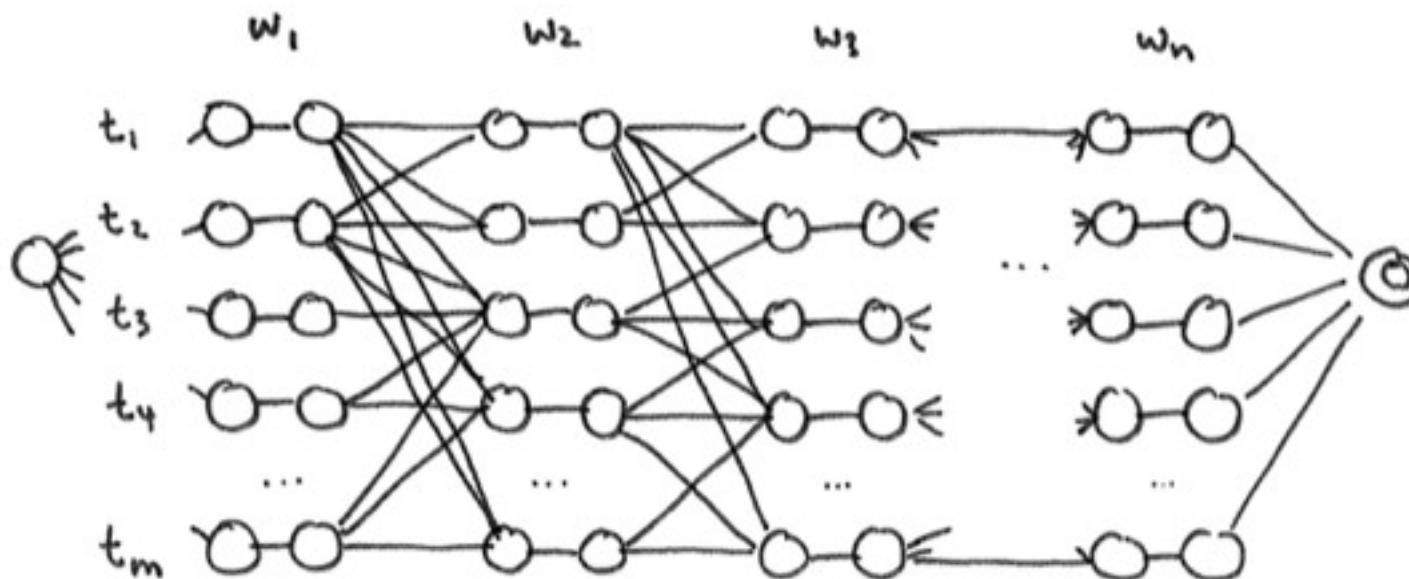
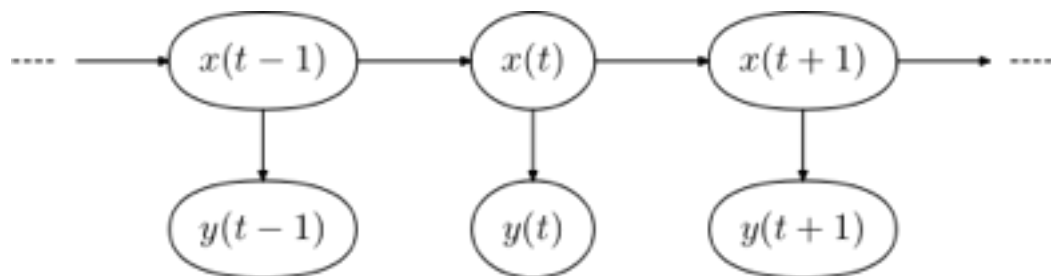
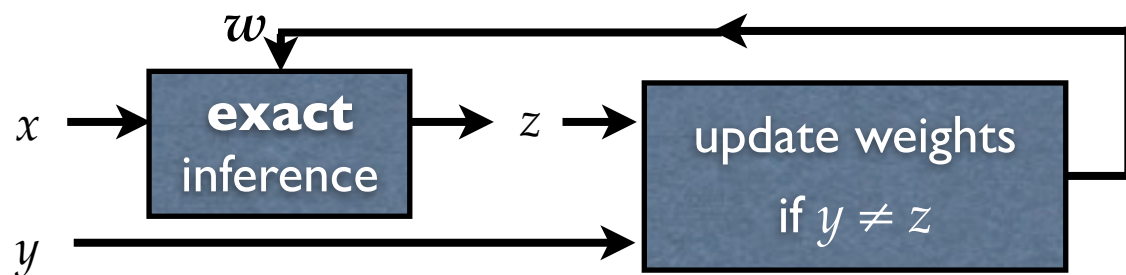
Algorithm: For $t = 1 \dots T, i = 1 \dots n$

$$z_i = F(x_i)$$

$$\text{If } (z_i \neq y_i) \quad \mathbf{W} \leftarrow \mathbf{W} + \Phi(x_i, y_i) - \Phi(x_i, z_i)$$

Output: Parameters \mathbf{W}

Inference: Dynamic Programming



Viterbi for argmax

Viterbi search for $\operatorname{argmax}_{t \dots t} P(t \dots t) \cdot P(w \dots w | t \dots t) :$

```
for j = 1 to m
  Q[1,j] = P(tj) · P(w1 | tj)
```

```
for i = 2 to n
  for j = 1 to m
    Q[i,j] = 0
    best-prev[i,j] = 0
    best-score = -∞
    for k = 1 to m
      r = P(tj | tk) · P(wi | tj) · Q[i-1,k]
      if r > best-score
        best-score = r
        best-prev[i,j] = k
        Q[i,j] = r
```

$Q[i,j]$ = cost of shortest path ending with word i getting assigned tag j .

sets back pointers

```
final-best = 0
final-score = -∞
for j = 1 to m
  if Q[n,j] > final-score
    final-score = Q[n,j]
    final-best = j
```

```
print tfinal-best}
current = final-best
for i = n-1 down to 1
  current = best-prev[i+1, current]
  print tcurrent
```

prints best tags in reverse order

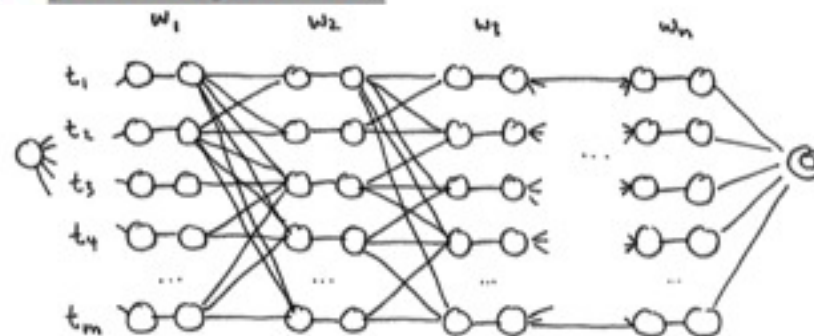
how about unigram?

Python implementation

Complete this Python code implementing the Viterbi algorithm for part-of-speech tagging. It should print a list of word/tag pairs, e.g. [('a', 'D'), ('can', 'N'), ('can', 'A'), ('can', 'V'), ('a', 'D'), ('can', 'N')].

```
1 from collections import defaultdict
2
3 best = defaultdict(lambda : defaultdict(float))
4 best[0]["<s>"] = 1
5 back = defaultdict(dict)
6
7 words = "<s> a can can can a can </s>".split()
8
9 tags = {"a": ["D"], "can": ["N", "A", "V"], "</s>": ["</s>"]} # possible tags for each word
10 ptag = {"D": {"N": 1}, "V": {"</s>": 0.5, "D": 0.5}, ... } # ptag[x][y] = p(y | x)
11 pword = {"D": {"a": 0.5}, "N": {"can": 0.1}, ... } # pword[x][w] = p(w | x)
12
13 for i, word in enumerate(words[1:], 1): # i=1,2,...; word=a,can,...
14     for tag in tags[word]:
15         for prev in best[i-1]:
16             if tag in ptag[prev]:
17                 score = best[i-1][prev] * ptag[prev][tag] * pword[tag][word]
18                 if score > best[i][tag]:
19                     best[i][tag] = score
20                     back[i][tag] = prev
21
22 def backtrack(i, tag):
23     if i == 0:
24         return []
25     return backtrack(i-1, back[i][tag]) + [(words[i], tag)]
26
27 print backtrack(len(words)-1, "</s>")[:-1]
```

Q: what about top-down recursive + memoization?



Trigram HMM

for $j = 1$ to m
 $Q_1[1, j] = \dots$

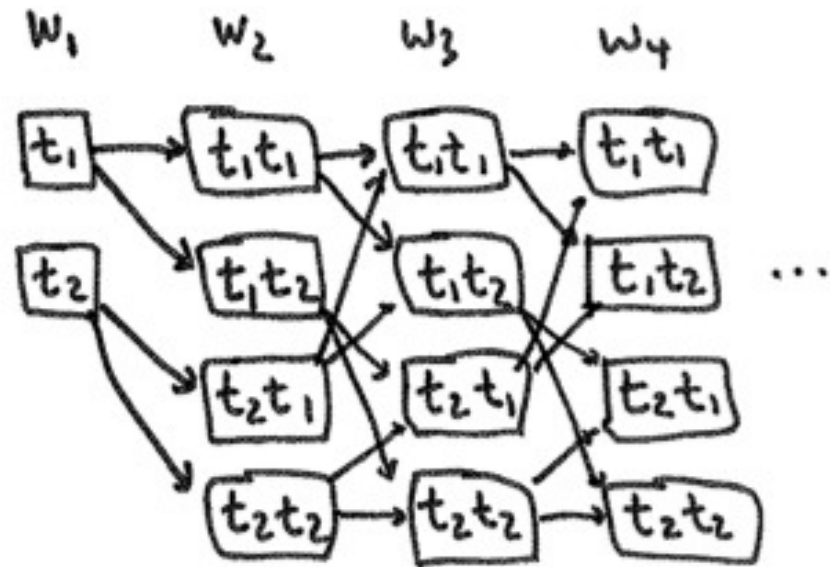
for $j = 1$ to m
 for $j_2 = 1$ to m
 $Q[2, j, j_2] = \dots$

for $i = 3$ to n
 for $j = 1$ to m
 for $j_2 = 1$ to m
 $Q[i, j, j_2] = 0$
 $\text{best-pred}[i, j, j_2] = 0$
 $\text{best-score} = -\infty$

for $k = 1$ to m

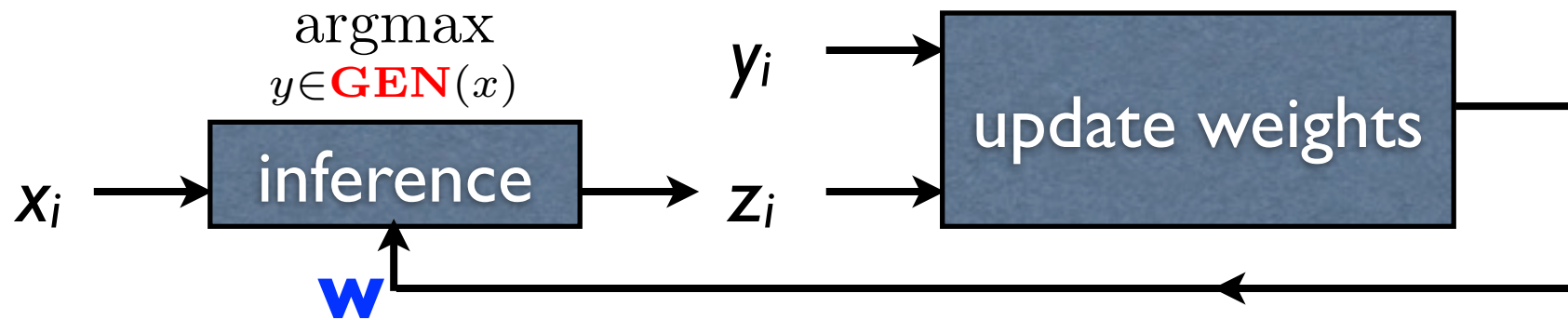
$$r = P(t_{j_2} | t_j) \cdot P(w_i | t_{j_2}) \cdot Q[i-1, k, j]$$

if $r > \text{best-score} \dots$

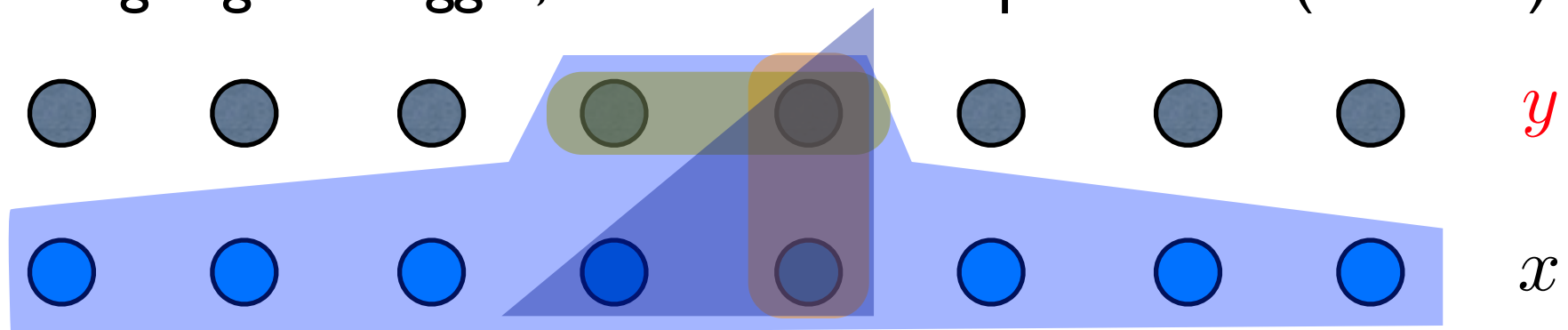


time complexity: $O(nT^3)$
 in general: $O(nT^g)$ for g -gram

Efficiency vs. Expressiveness



- the **inference** (argmax) must be efficient
 - either the search space **GEN**(x) is small, or **factored**
 - features must be local to y (but can be global to x)
 - e.g. bigram tagger, but look at all input words (cf. CRFs)



Averaged Perceptron

Inputs: Training set (x_i, y_i) for $i = 1 \dots n$

Initialization: $\mathbf{W}_0 = 0$

Define: $F(x) = \operatorname{argmax}_{y \in \mathbf{GEN}(x)} \Phi(x, y) \cdot \mathbf{W}$

Algorithm: For $t = 1 \dots T, i = 1 \dots n$
 $z_i = F(x_i)$
If $(z_i \neq y_i)$ $\mathbf{W}_{j+1} \leftarrow \mathbf{W}_j + \Phi(x_i, y_i) - \Phi(x_i, z_i)$

Output: Parameters $\mathbf{W} = \sum_j \mathbf{W}_j$

- more stable and accurate results
- approximation of voted perceptron (Freund & Schapire, 1999)

Averaging Tricks

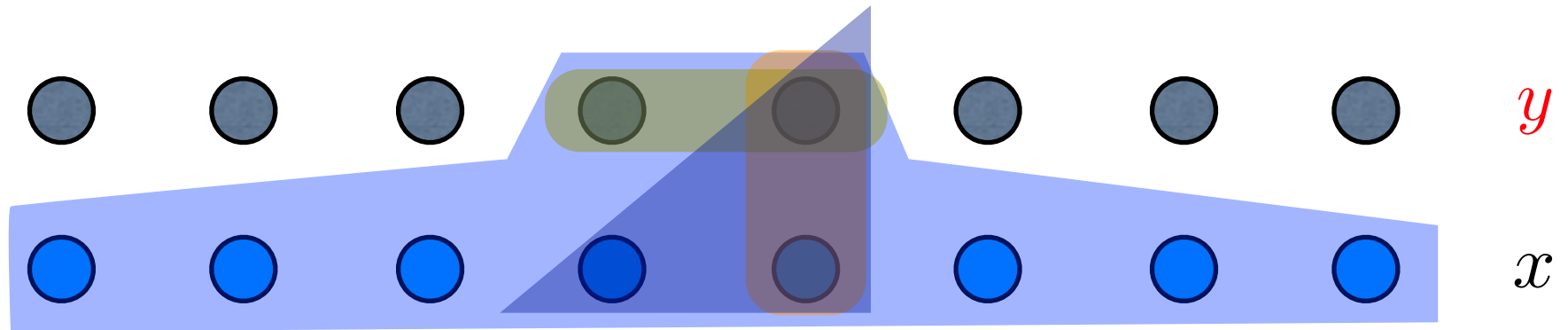
- Daume (2006, PhD thesis)

Algorithm AVERAGEDSTRUCTUREDPERCEPTRON($x_{1:N}, y_{1:N}, I$)

```
1:  $w_0 \leftarrow \langle 0, \dots, 0 \rangle$ 
2:  $w_a \leftarrow \langle 0, \dots, 0 \rangle$ 
3:  $c \leftarrow 1$ 
4: for  $i = 1 \dots I$  do
5:   for  $n = 1 \dots N$  do
6:      $\hat{y}_n \leftarrow \arg \max_{y \in \mathcal{Y}} w_0^\top \Phi(x_n, y)$ 
7:     if  $y_n \neq \hat{y}_n$  then
8:        $w_0 \leftarrow w_0 + \Phi(x_n, y_n) - \Phi(x_n, \hat{y}_n)$ 
9:        $w_a \leftarrow w_a + c\Phi(x_n, y_n) - c\Phi(x_n, \hat{y}_n)$ 
10:    end if
11:     $c \leftarrow c + 1$ 
12:  end for
13: end for
14: return  $w_0 - w_a/c$ 
```

Figure 2.3: The averaged structured perceptron learning algorithm.

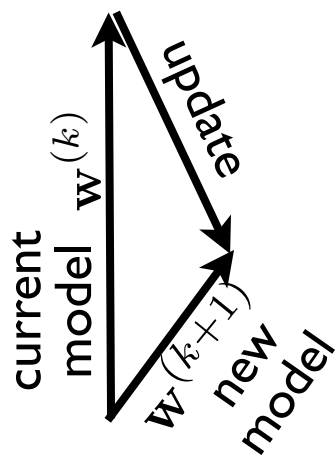
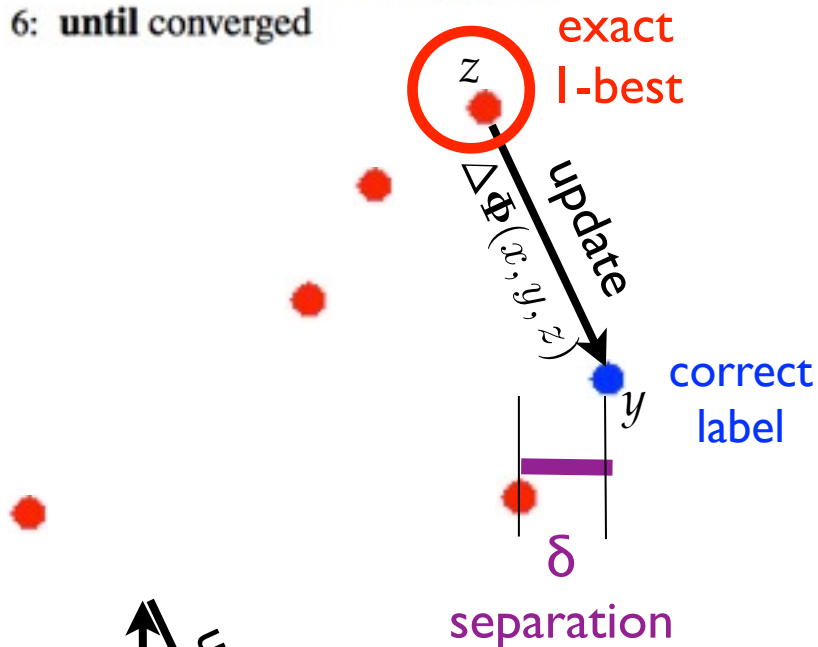
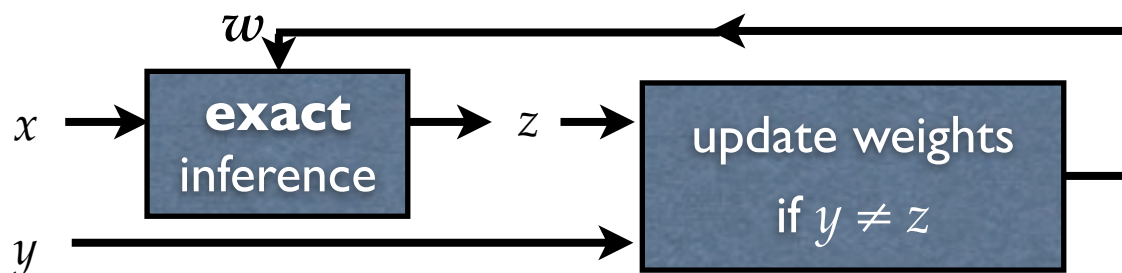
Do we need smoothing?



- smoothing is much easier in discriminative models
- just make sure for each feature template, its subset templates are also included
 - e.g., to include $(t_0 w_0 w_{-1})$ you must also include
 - $(t_0 w_0)$ $(t_0 w_{-1})$ $(w_0 w_{-1})$
 - and maybe also $(t_0 t_{-1})$ because t is less sparse than w

Geometry of Convergence Proof pt I

- 1: repeat
- 2: for each example (x, y) in D do
- 3: $z \leftarrow \text{EXACT}(x, \mathbf{w})$
- 4: if $z \neq y$ then
- 5: $\mathbf{w} \leftarrow \mathbf{w} + \Delta\Phi(x, y, z)$
- 6: until converged



perceptron update:

$$\mathbf{w}^{(k+1)} = \mathbf{w}^{(k)} + \Delta\Phi(x, y, z)$$

$$\mathbf{u} \cdot \mathbf{w}^{(k+1)} = \mathbf{u} \cdot \mathbf{w}^{(k)} + \boxed{\mathbf{u} \cdot \Delta\Phi(x, y, z) \geq \delta \text{ margin}}$$

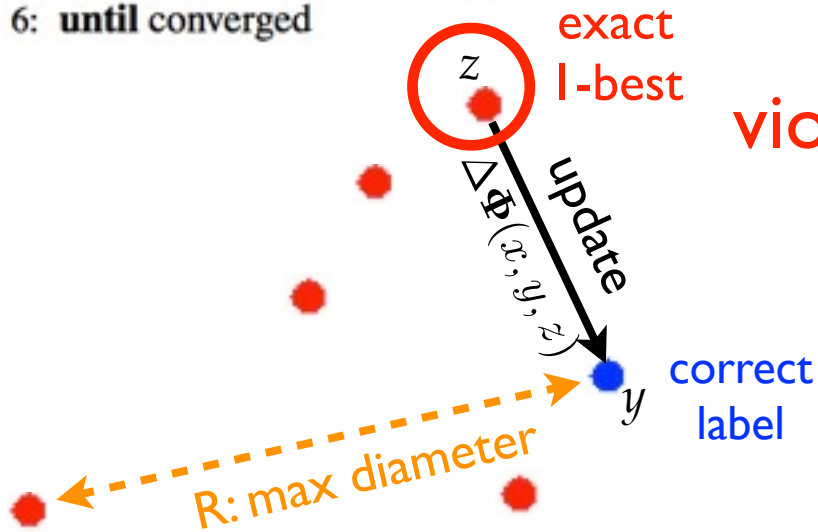
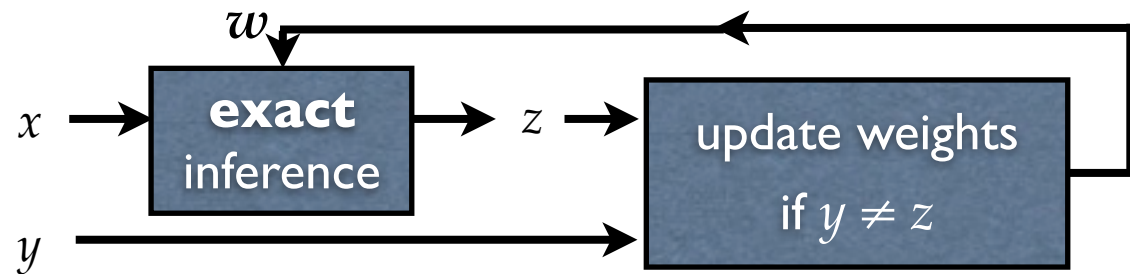
$$\mathbf{u} \cdot \mathbf{w}^{(k+1)} \geq k\delta \quad (\text{by induction})$$

$$\|\mathbf{u}\| \|\mathbf{w}^{(k+1)}\| \geq \mathbf{u} \cdot \mathbf{w}^{(k+1)} \geq k\delta$$

$$\|\mathbf{w}^{(k+1)}\| \geq k\delta \quad (\text{part I: upperbound})$$

Geometry of Convergence Proof pt 2

- 1: repeat
- 2: for each example (x, y) in D do
- 3: $z \leftarrow \text{EXACT}(x, \mathbf{w})$
- 4: if $z \neq y$ then
- 5: $\mathbf{w} \leftarrow \mathbf{w} + \Delta\Phi(x, y, z)$
- 6: until converged



violation: incorrect label scored higher

perceptron update:

$$\mathbf{w}^{(k+1)} = \mathbf{w}^{(k)} + \Delta\Phi(x, y, z)$$

$$\|\mathbf{w}^{(k+1)}\|^2 = \|\mathbf{w}^{(k)} + \Delta\Phi(x, y, z)\|^2$$

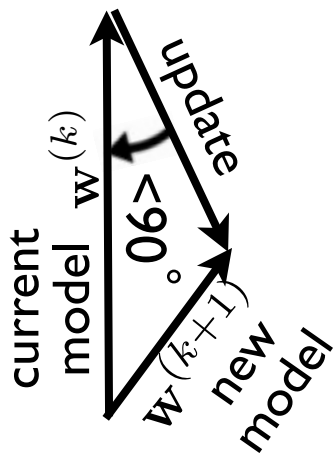
$$= \|\mathbf{w}^{(k)}\|^2 + \|\Delta\Phi(x, y, z)\|^2 + 2 \mathbf{w}^{(k)} \cdot \Delta\Phi(x, y, z)$$

$$\leq R^2$$

diameter

$$\leq 0$$

violation



by induction: $\|\mathbf{w}^{(k+1)}\|^2 \leq kR^2$ (part 2: upperbound)

parts 1+2 => update bounds:

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

Experiments

Experiments: Tagging

- (almost) identical features from (Ratnaparkhi, 1996)
 - trigram tagger: current tag t_i , previous tags t_{i-1} , t_{i-2}
 - current word w_i and its spelling features
 - surrounding words w_{i-1} w_{i+1} w_{i-2} w_{i+2} ..

Method	Error rate/%	Numits
Perc, avg, cc=0	2.93	10
Perc, noavg, cc=0	3.68	20
Perc, avg, cc=5	3.03	6
Perc, noavg, cc=5	4.04	17
ME, cc=0	3.4	100
ME, cc=5	3.28	200

Experiments: NP Chunking

- B-I-O scheme

Rockwell International Corp.

B I I

's Tulsa unit said it signed

B I I O B O

a tentative agreement ...

B I I

- features:
 - unigram model
 - surrounding words and POS tags

Current word	w_i	& t_i
Previous word	w_{i-1}	& t_i
Word two back	w_{i-2}	& t_i
Next word	w_{i+1}	& t_i
Word two ahead	w_{i+2}	& t_i
Bigram features	w_{i-2}, w_{i-1}	& t_i
	w_{i-1}, w_i	& t_i
	w_i, w_{i+1}	& t_i
	w_{i+1}, w_{i+2}	& t_i
Current tag	p_i	& t_i
Previous tag	p_{i-1}	& t_i
Tag two back	p_{i-2}	& t_i
Next tag	p_{i+1}	& t_i
Tag two ahead	p_{i+2}	& t_i
Bigram tag features	p_{i-2}, p_{i-1}	& t_i
	p_{i-1}, p_i	& t_i
	p_i, p_{i+1}	& t_i
	p_{i+1}, p_{i+2}	& t_i
Trigram tag features	p_{i-2}, p_{i-1}, p_i	& t_i
	p_{i-1}, p_i, p_{i+1}	& t_i
	p_i, p_{i+1}, p_{i+2}	& t_i

Experiments: NP Chunking

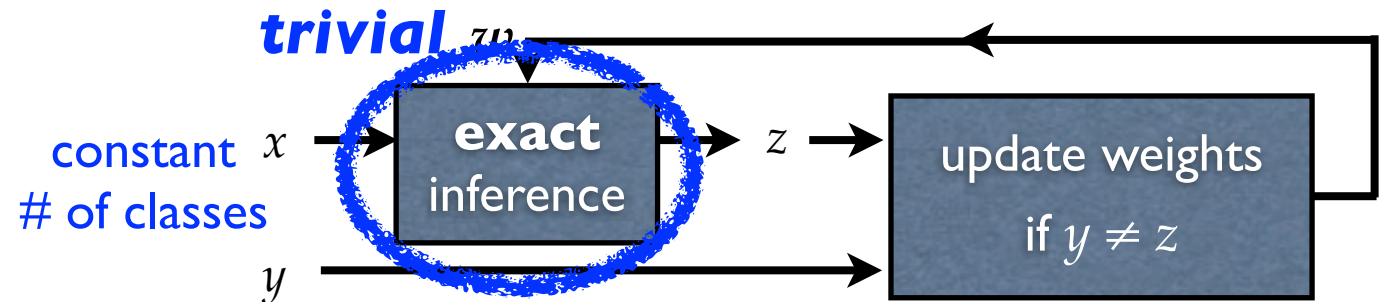
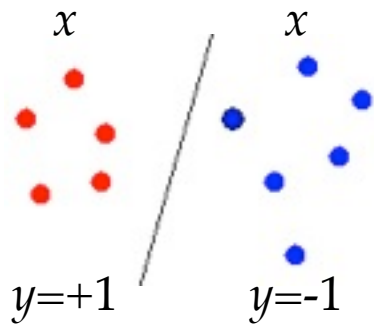
- results

Method	F-Measure	Numits
Perceptron, avg, cc=0	93.53	13
Perceptron, noavg, cc=0	93.04	35
Perceptron, avg, cc=5	93.33	9
Perceptron, noavg, cc=5	91.88	39
Max-ent, cc=0	92.34	900
Max-ent, cc=5	92.65	200

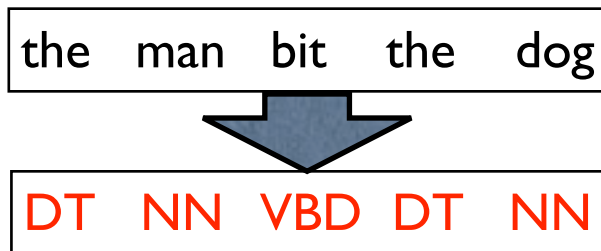
- (Sha and Pereira, 2003) **trigram** tagger
 - voted perceptron: 94.09% vs. CRF: 94.38%

Structured Perceptron (Collins 02)

binary classification

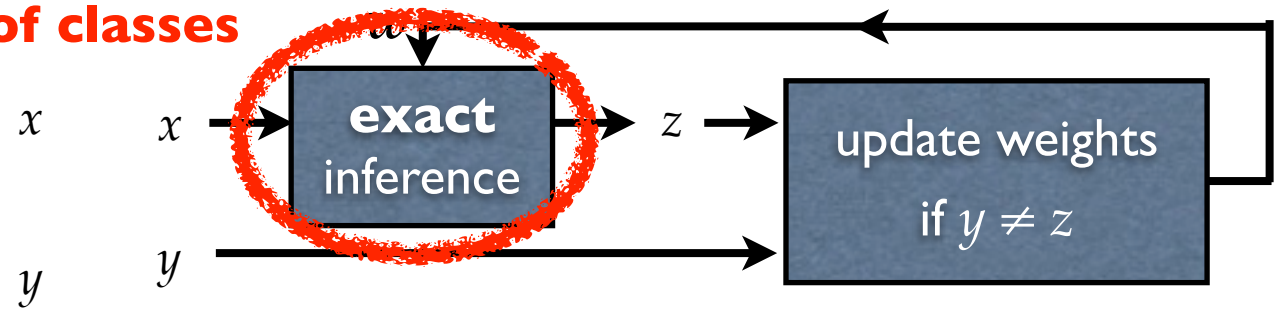


structured classification



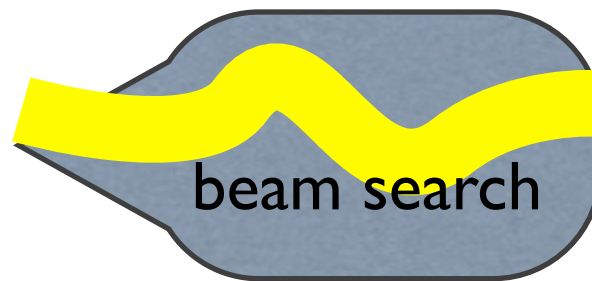
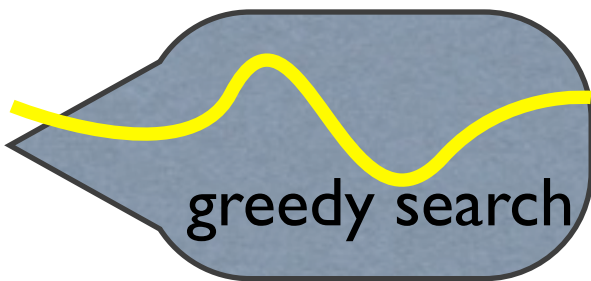
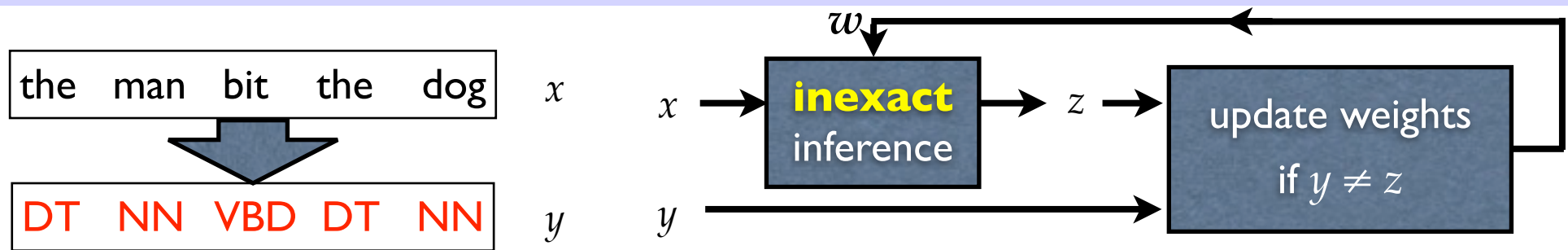
exponential
of classes

hard



- challenge: search efficiency (exponentially many classes)
 - often use dynamic programming (DP)
 - but still too slow for repeated use, e.g. parsing is $O(n^3)$
 - and can't use non-local features in DP

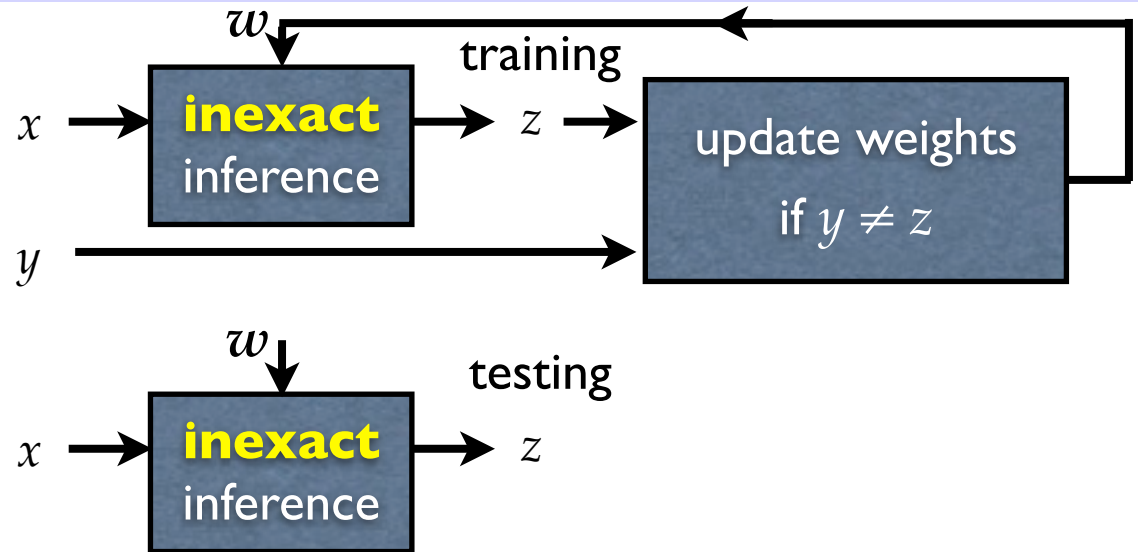
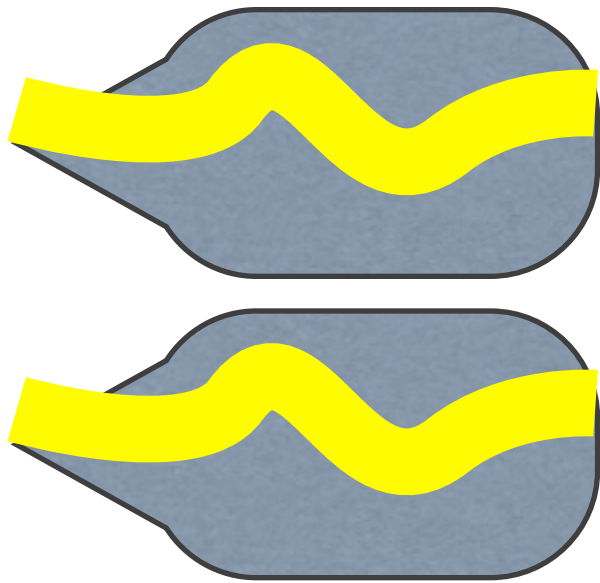
Learning w/ Inexact Inference (Huang et al 2012)



does it still work???

- routine use of inexact inference in NLP (e.g. beam search)
- how does structured perceptron work with inexact search?
 - so far most structured learning theory assume exact search
 - would search errors break these learning properties?
 - if so how to modify learning to accommodate inexact search?

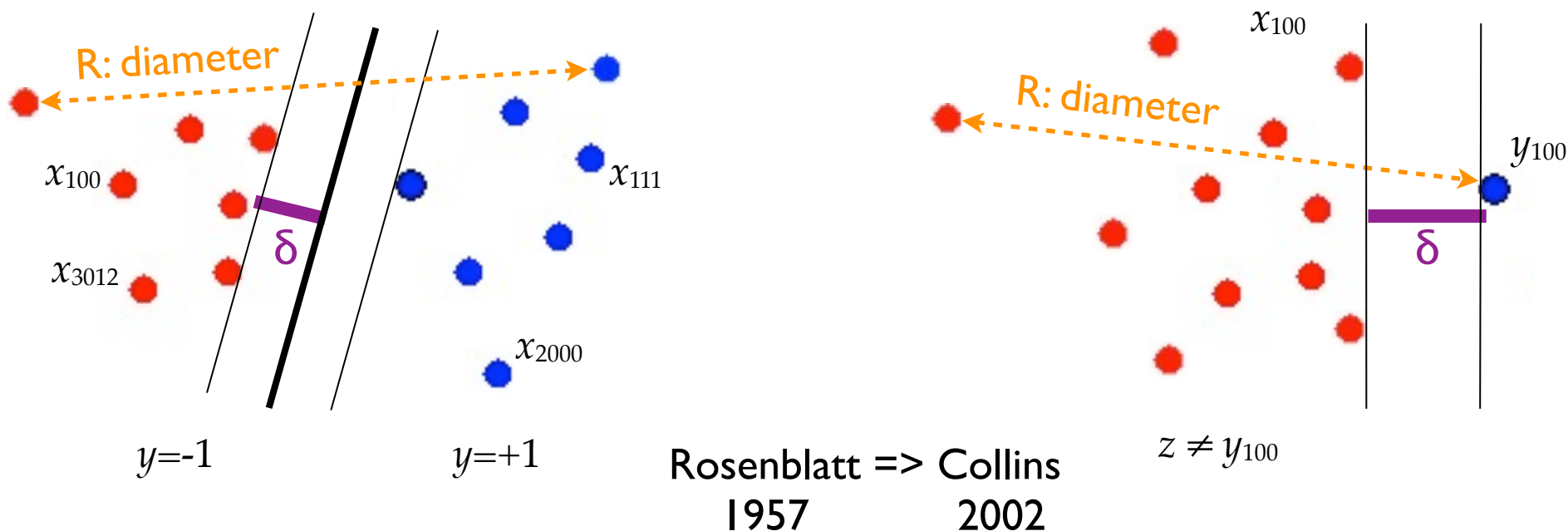
Idea: Search-Error-Robust Model



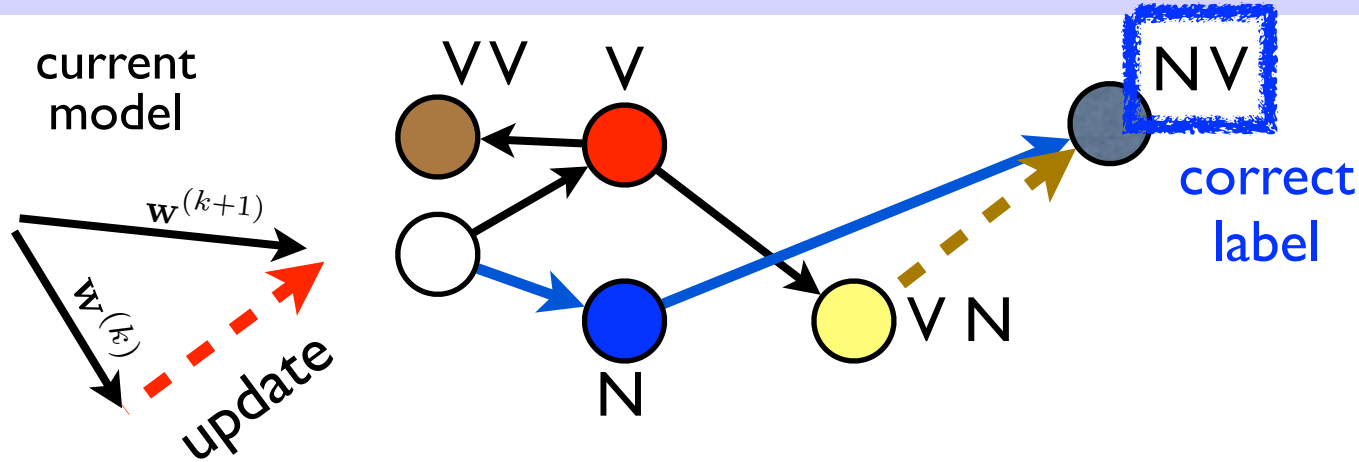
- train a “search-specific” or “search-error-robust” model
 - we assume the same “search box” in training and testing
 - model should “live with” search errors from search box
- exact search => convergence; greedy => no convergence
 - how can we make perceptron converge w/ greedy search?

Convergence with Exact Search

- linear classification: converges iff. data is separable
- structured: converges iff. data separable & search exact
 - there is an oracle vector that correctly labels all examples
 - one vs the rest (correct label better than all incorrect labels)
- **theorem**: if separable, then **# of updates** $\leq R^2 / \delta^2$ R: diameter



Convergence with Exact Search



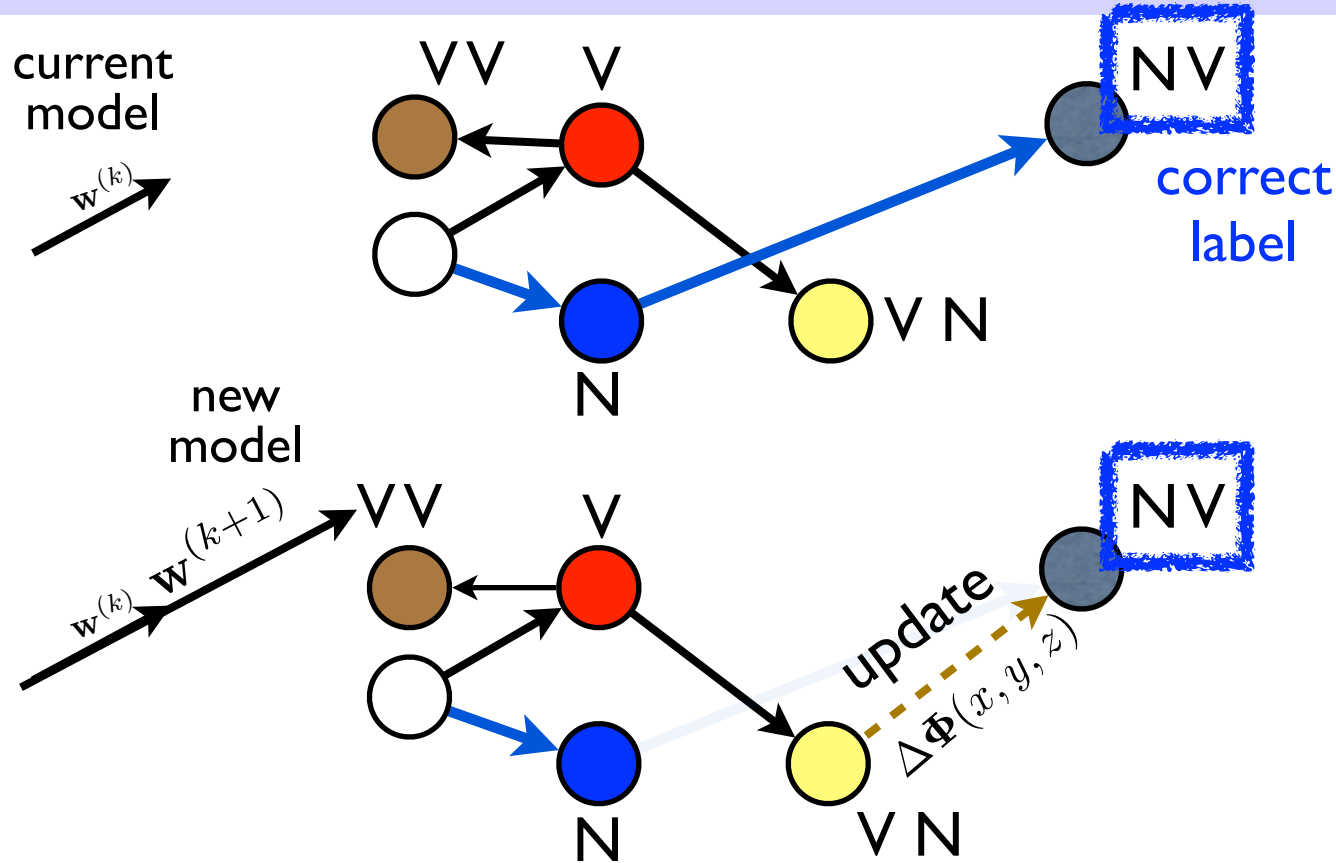
training example

time	flies
N	V

output space
 $\{N, V\} \times \{N, V\}$

standard perceptron
converges with exact
search

No Convergence w/ Greedy Search



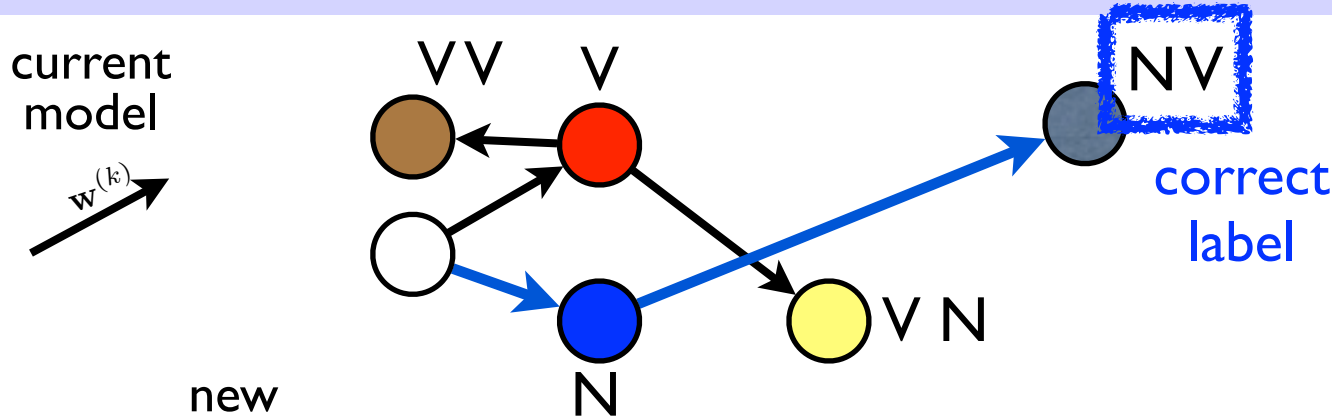
training example

time flies
N V

output space
 $\{N, V\} \times \{N, V\}$

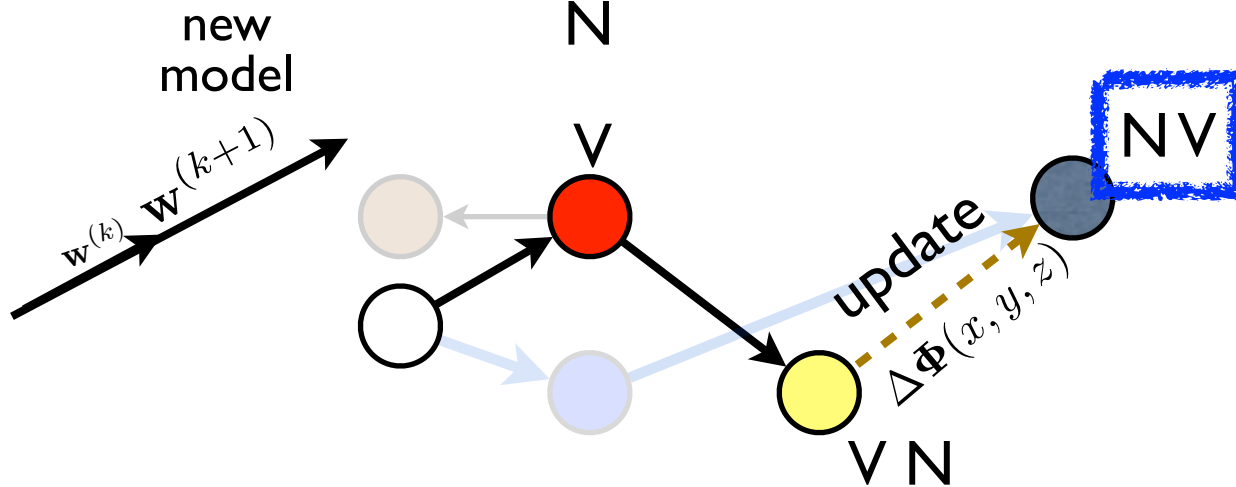
standard perceptron
does not converge
with greedy search

Early update (Collins/Roark 2004) to rescue

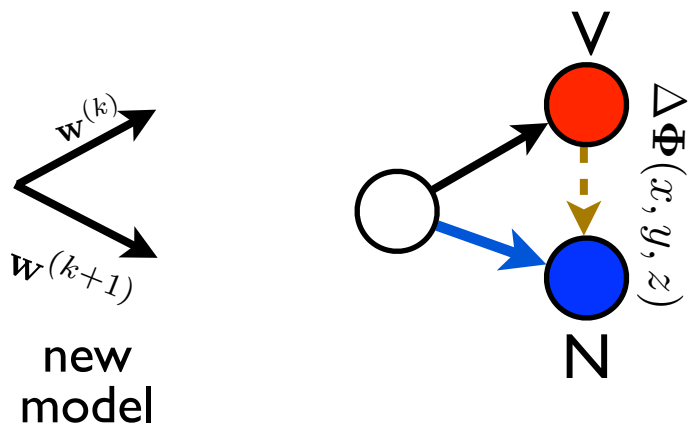


training example
time flies
N V

output space
 $\{N, V\} \times \{N, V\}$



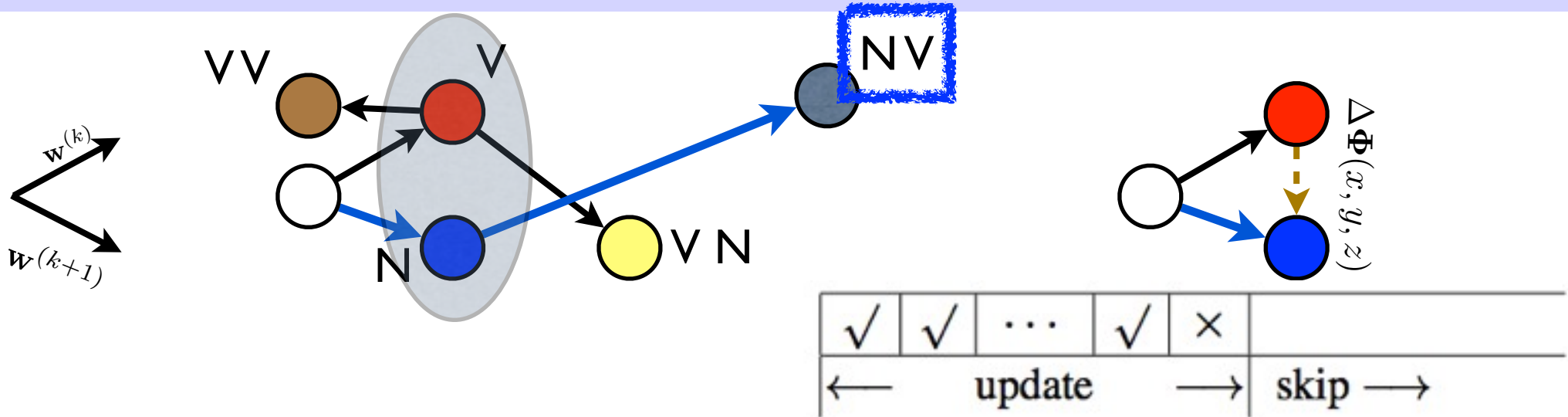
standard perceptron
does not converge
with greedy search



✓	✓	...	✓	×	
← update →				← skip →	

stop and update at the first mistake

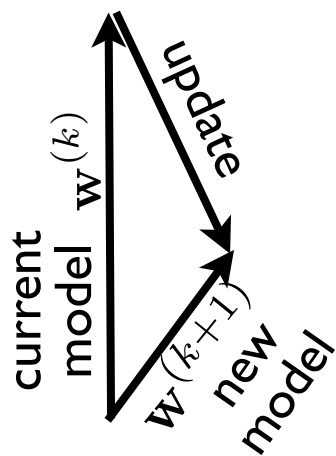
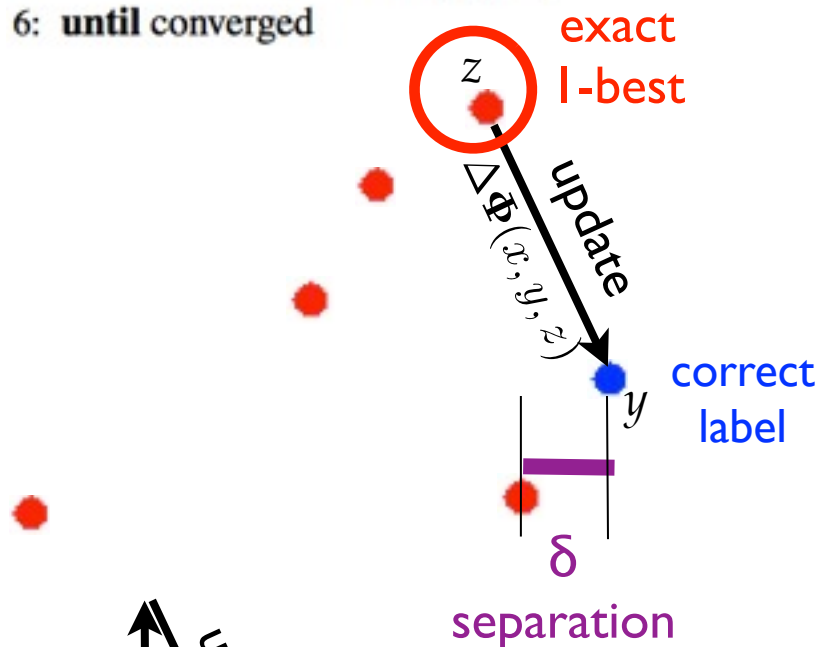
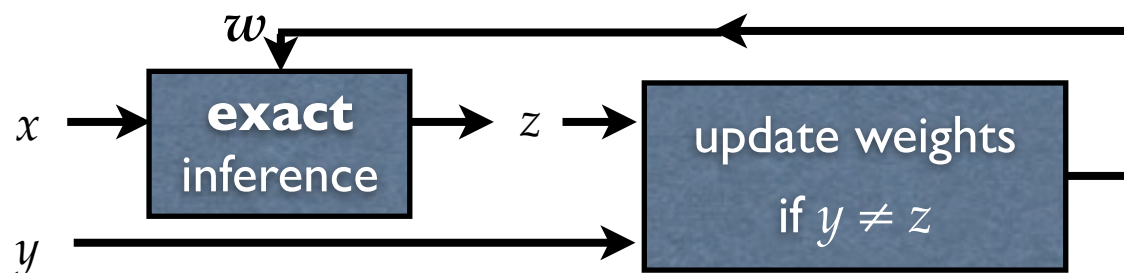
Why?



- why does inexact search break convergence property?
 - what is required for convergence? exactness?
- why does early update (Collins/Roark 04) work?
 - it works well in practice and is now a standard method
 - but there has been no theoretical justification
- we answer these Qs by inspecting the convergence proof

Geometry of Convergence Proof pt I

- 1: repeat
- 2: for each example (x, y) in D do
- 3: $z \leftarrow \text{EXACT}(x, \mathbf{w})$
- 4: if $z \neq y$ then
- 5: $\mathbf{w} \leftarrow \mathbf{w} + \Delta\Phi(x, y, z)$
- 6: until converged



perceptron update:

$$\mathbf{w}^{(k+1)} = \mathbf{w}^{(k)} + \Delta\Phi(x, y, z)$$

$$\mathbf{u} \cdot \mathbf{w}^{(k+1)} = \mathbf{u} \cdot \mathbf{w}^{(k)} + \boxed{\mathbf{u} \cdot \Delta\Phi(x, y, z) \geq \delta \text{ margin}}$$

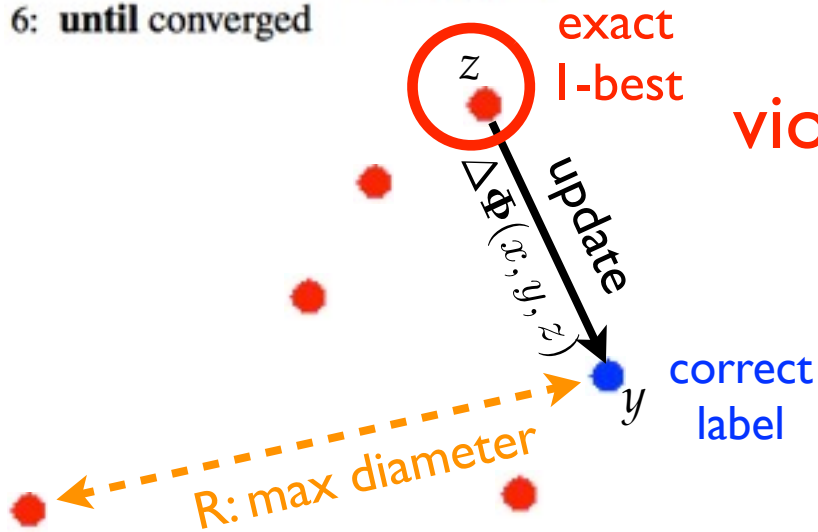
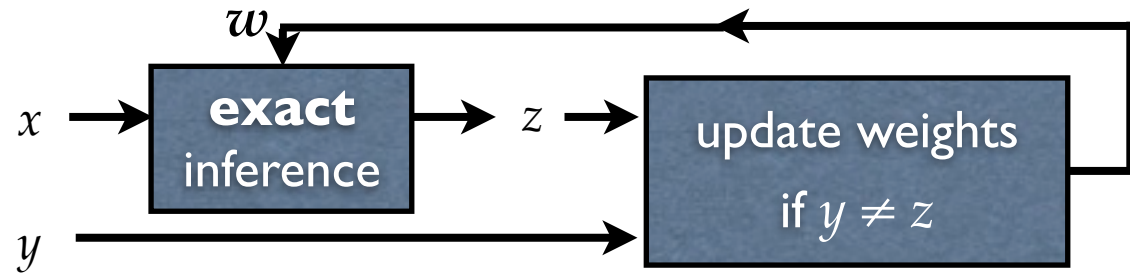
$$\mathbf{u} \cdot \mathbf{w}^{(k+1)} \geq k\delta \quad (\text{by induction})$$

$$\|\mathbf{u}\| \|\mathbf{w}^{(k+1)}\| \geq \mathbf{u} \cdot \mathbf{w}^{(k+1)} \geq k\delta$$

$$\|\mathbf{w}^{(k+1)}\| \geq k\delta \quad (\text{part I: upperbound})$$

Geometry of Convergence Proof pt 2

- 1: repeat
- 2: for each example (x, y) in D do
- 3: $z \leftarrow \text{EXACT}(x, \mathbf{w})$
- 4: if $z \neq y$ then
- 5: $\mathbf{w} \leftarrow \mathbf{w} + \Delta\Phi(x, y, z)$
- 6: until converged



violation: incorrect label scored higher

perceptron update:

$$\mathbf{w}^{(k+1)} = \mathbf{w}^{(k)} + \Delta\Phi(x, y, z)$$

$$\|\mathbf{w}^{(k+1)}\|^2 = \|\mathbf{w}^{(k)} + \Delta\Phi(x, y, z)\|^2$$

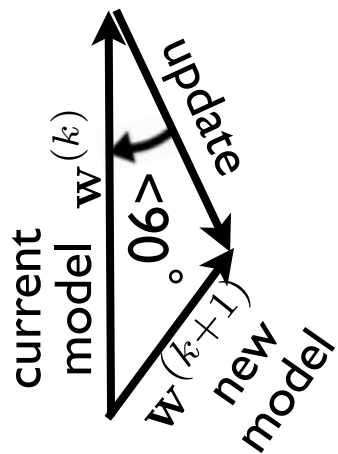
$$= \|\mathbf{w}^{(k)}\|^2 + \|\Delta\Phi(x, y, z)\|^2 + 2 \mathbf{w}^{(k)} \cdot \Delta\Phi(x, y, z)$$

$$\leq R^2$$

diameter

$$\leq 0$$

violation



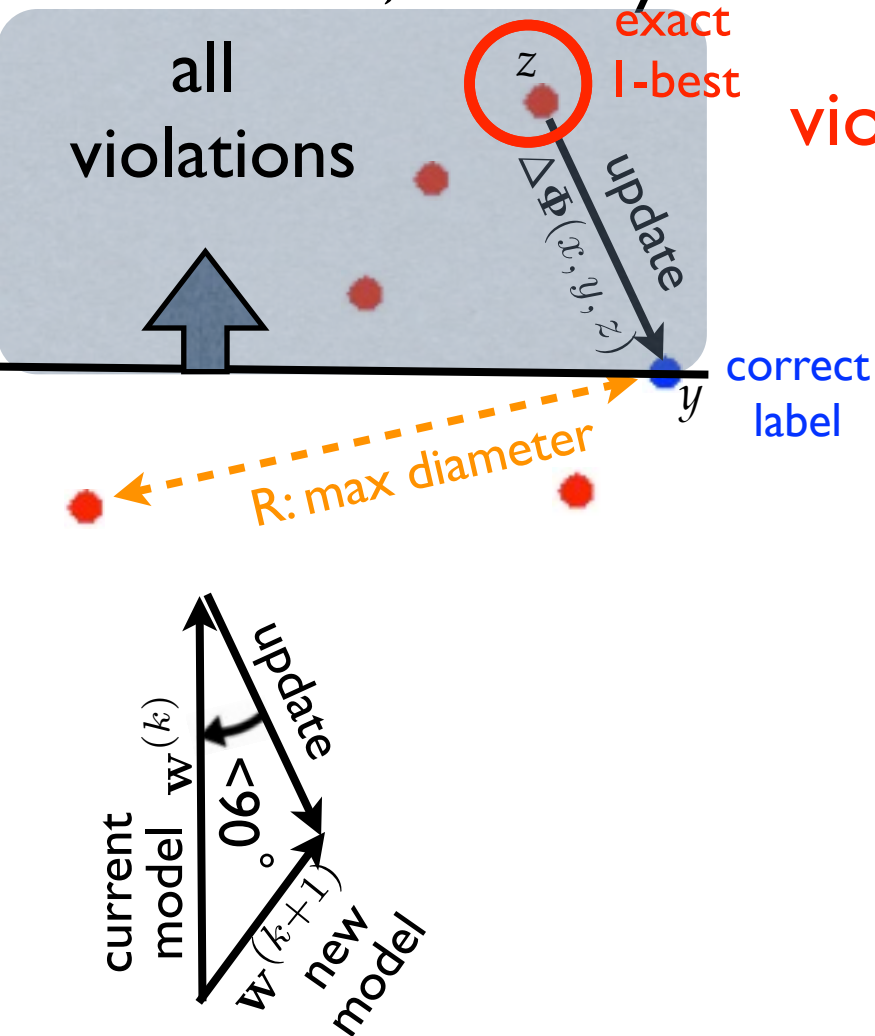
by induction: $\|\mathbf{w}^{(k+1)}\|^2 \leq kR^2$ (part 2: upperbound)

parts 1+2 => update bounds:

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

Violation is All we need!

- exact search is **not** really required by the proof
- rather, it is only used to ensure violation!



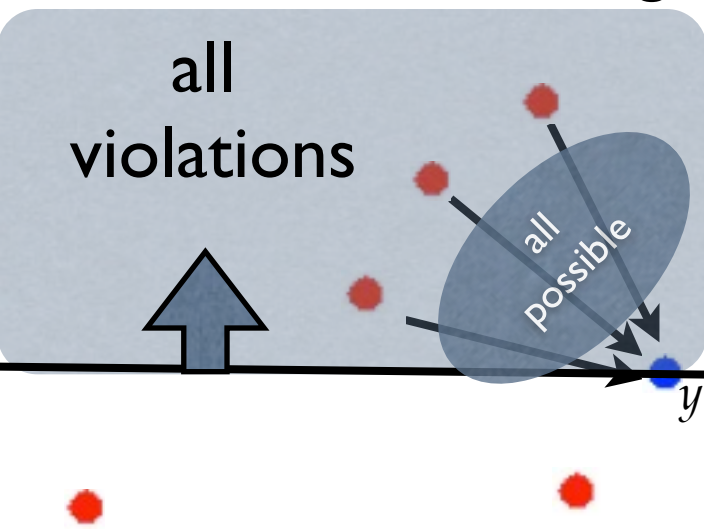
violation: incorrect label scored higher

the proof only uses 3 facts:

1. separation (margin)
2. diameter (always finite)
3. violation (but no need for exact)

Violation-Fixing Perceptron

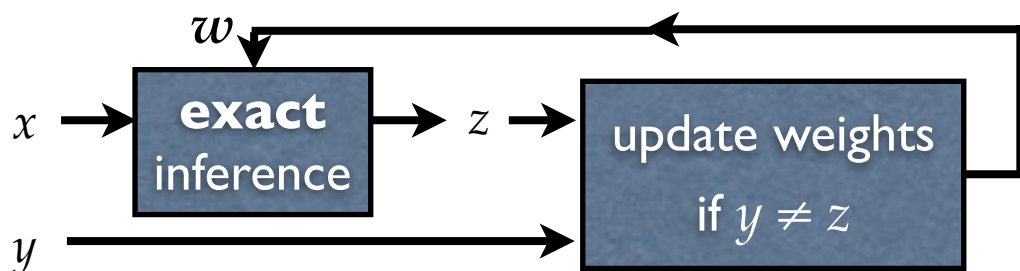
- if we guarantee violation, we don't care about exactness!
 - violation is good b/c we can at least fix a mistake



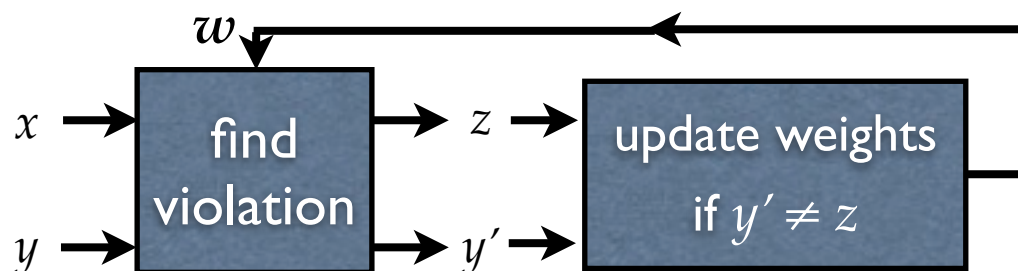
same mistake bound as before!

- 1: **repeat**
- 2: **for each example** (x, y) **in** D **do**
- 3: $(x, y', z) = \text{FINDVIOLATION}(x, y, \mathbf{w})$
- 4: **if** $z \neq y$ **then** $\triangleright (x, y', z)$ is a viol
- 5: $\mathbf{w} \leftarrow \mathbf{w} + \Delta \Phi(x, y', z)$
- 6: **until** converged

standard perceptron

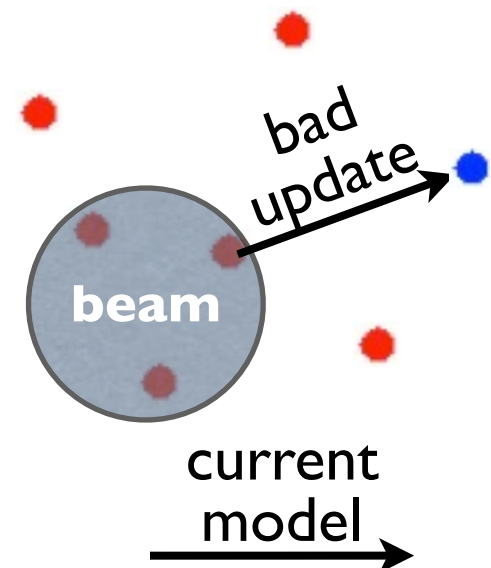


violation-fixing perceptron

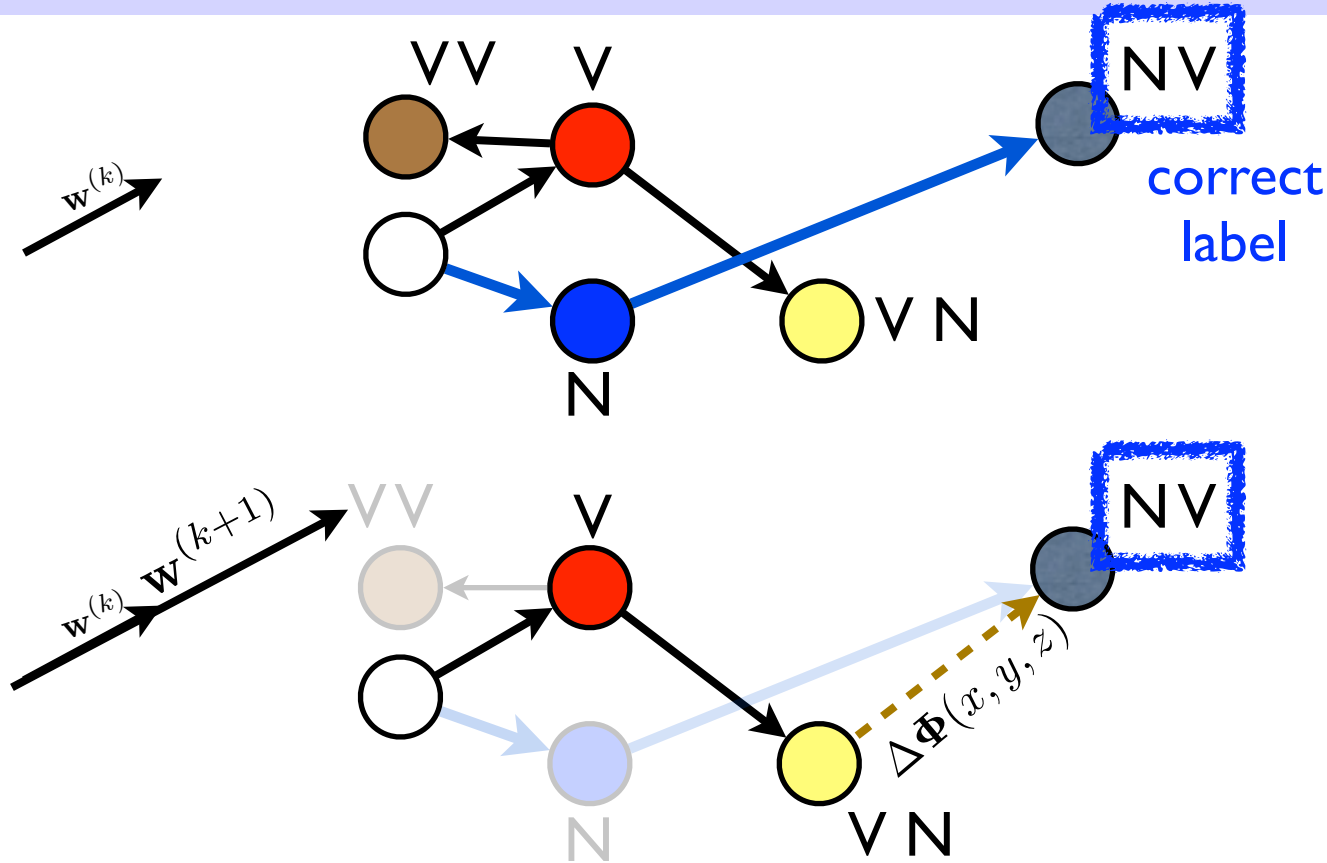


What if can't guarantee violation

- this is why perceptron doesn't work well w/ inexact search
 - because not every update is guaranteed to be a violation
 - thus the proof breaks; no convergence guarantee
- example: beam or greedy search
 - the model might prefer the correct label (if exact search)
 - but the search prunes it away
 - such a **non-violation update is "bad"** because it doesn't fix any mistake
 - the new model still misguides the search



Standard Update: No Guarantee



training example

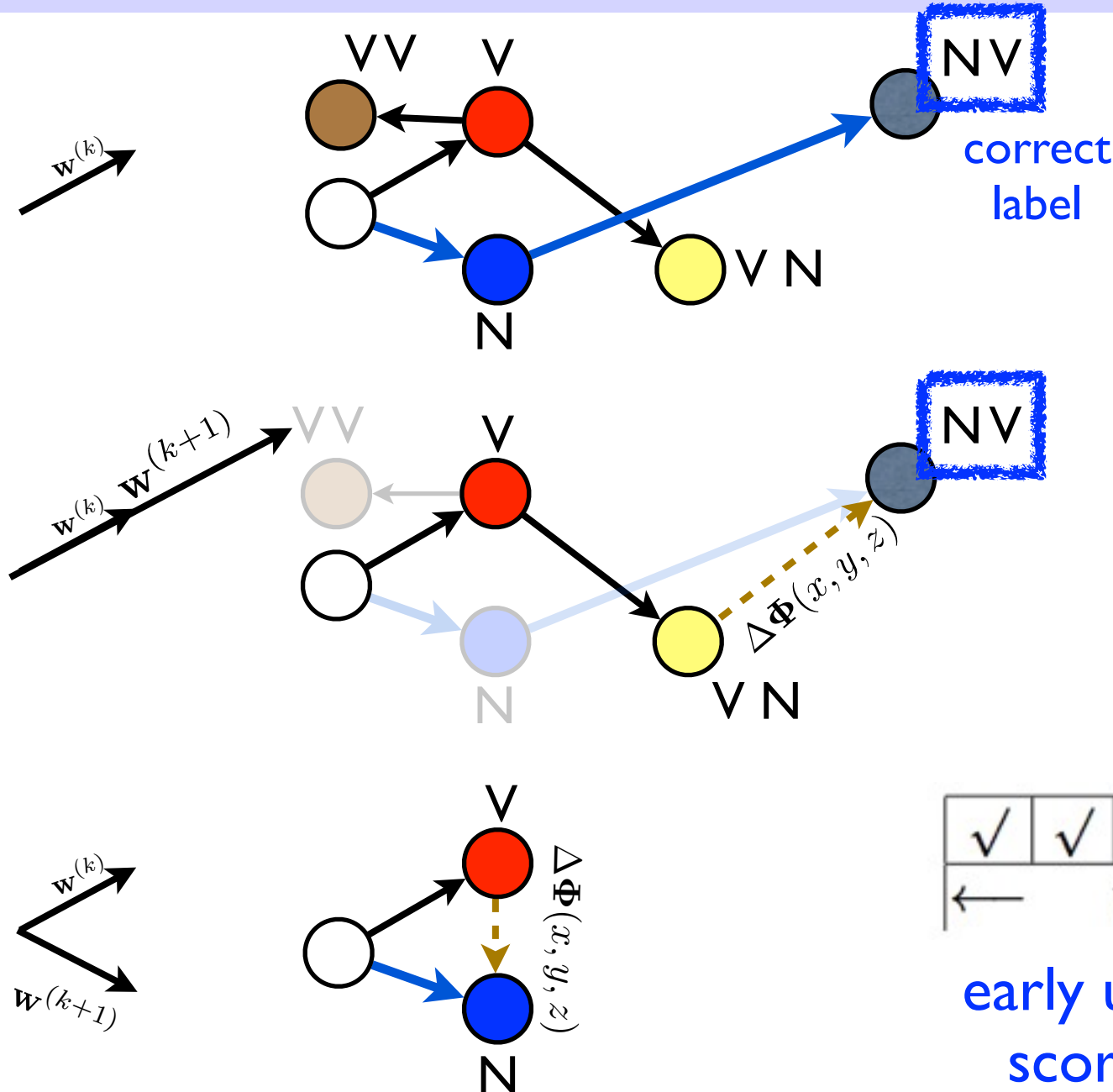
time flies
N V

output space
 $\{N, V\} \times \{N, V\}$

standard update
doesn't converge
b/c it doesn't
guarantee violation

correct label scores higher.
non-violation: bad update!

Early Update: Guarantees Violation



training example

time flies
N V

output space
 $\{N, V\} \times \{N, V\}$

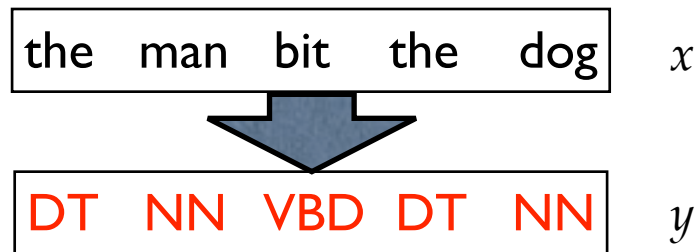
standard update
doesn't converge
b/c it doesn't
guarantee violation

✓	✓	...	✓	×	
← update				skip	→

early update: incorrect prefix
scores higher: a violation!

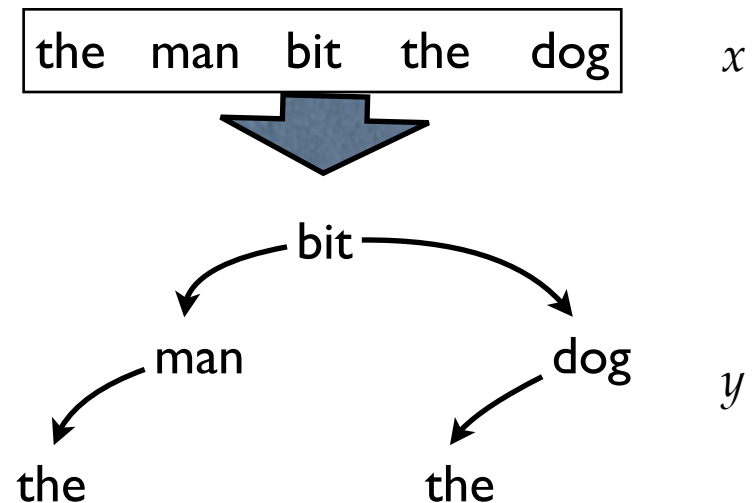
Experiments

trigram part-of-speech tagging



local features only,
exact search tractable
(proof of concept)

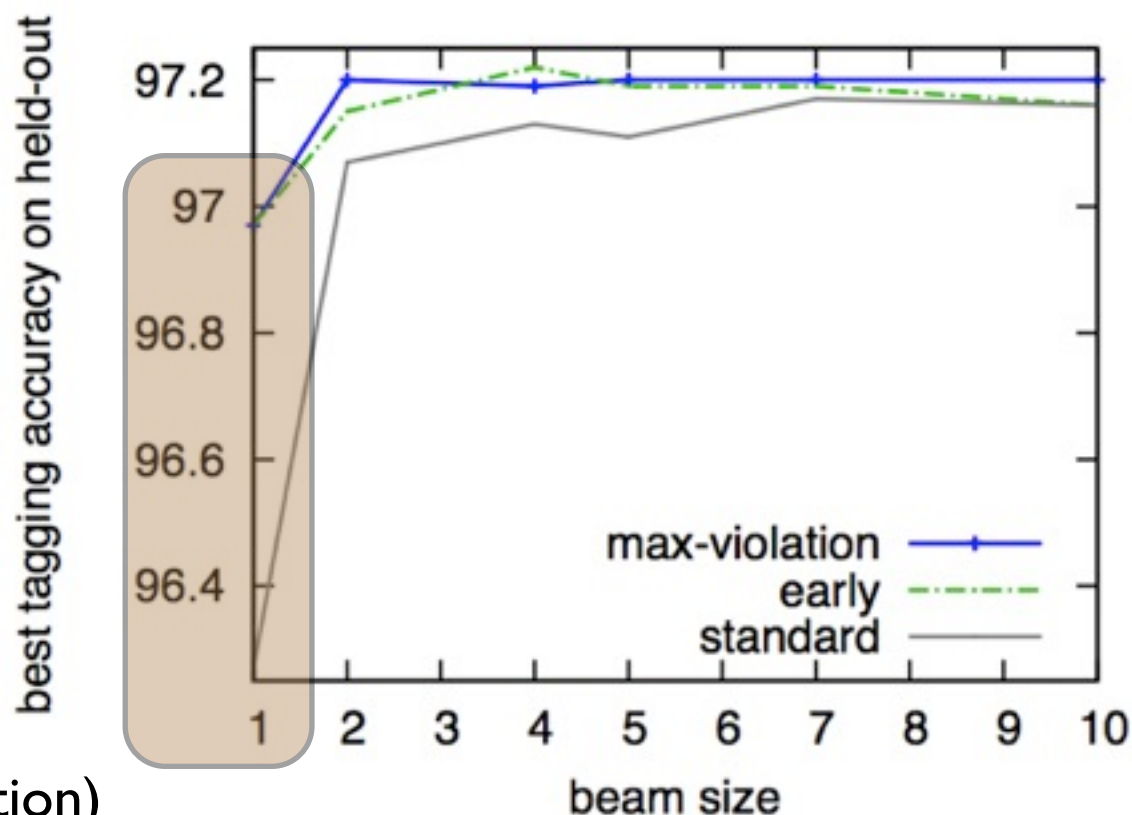
incremental dependency parsing



non-local features,
exact search intractable
(real impact)

I) Trigram Part of Speech Tagging

- standard update performs terribly with greedy search ($b=1$)
- because search error is severe at $b=1$: half updates are bad!
- no real difference beyond $b=2$: search error becomes rare

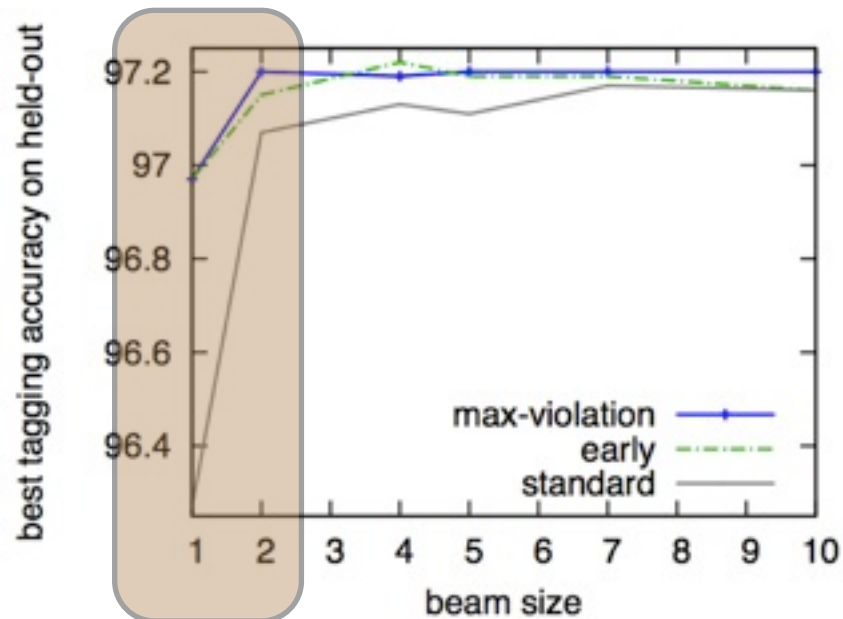


% of bad (non-violation)
standard updates

53% 10% 1.5% 0.5%

Max-Violation Reduces Training Time

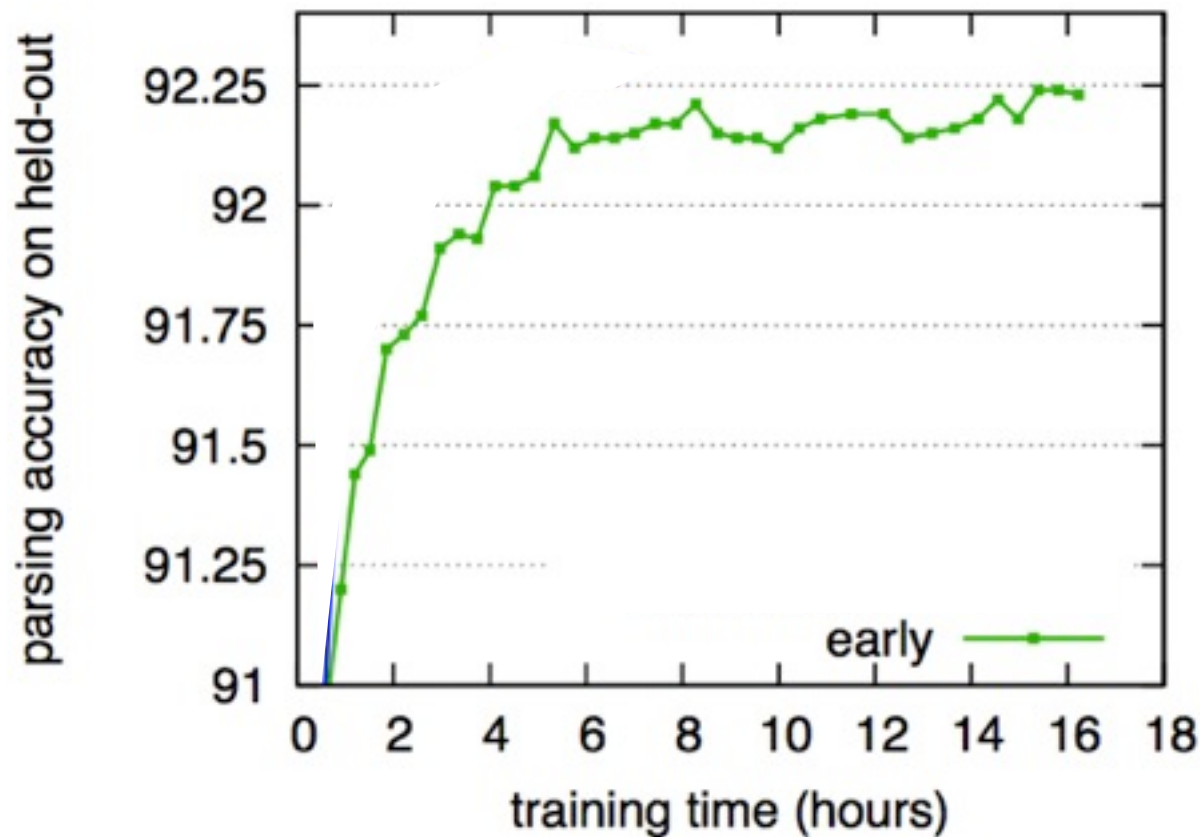
- max-violation peaks at $b=2$, greatly reduced training time
- early update achieves the highest dev/test accuracy
 - comparable to best published accuracy (Shen et al '07)
- future work: add non-local features to tagging



	<i>beam</i>	<i>iter</i>	<i>time</i>	<i>test</i>
standard	-	6	162m	97.28
early	4	6	37m	97.27
max-violation	2	3	26m	97.27
Shen et al (2007)				97.33

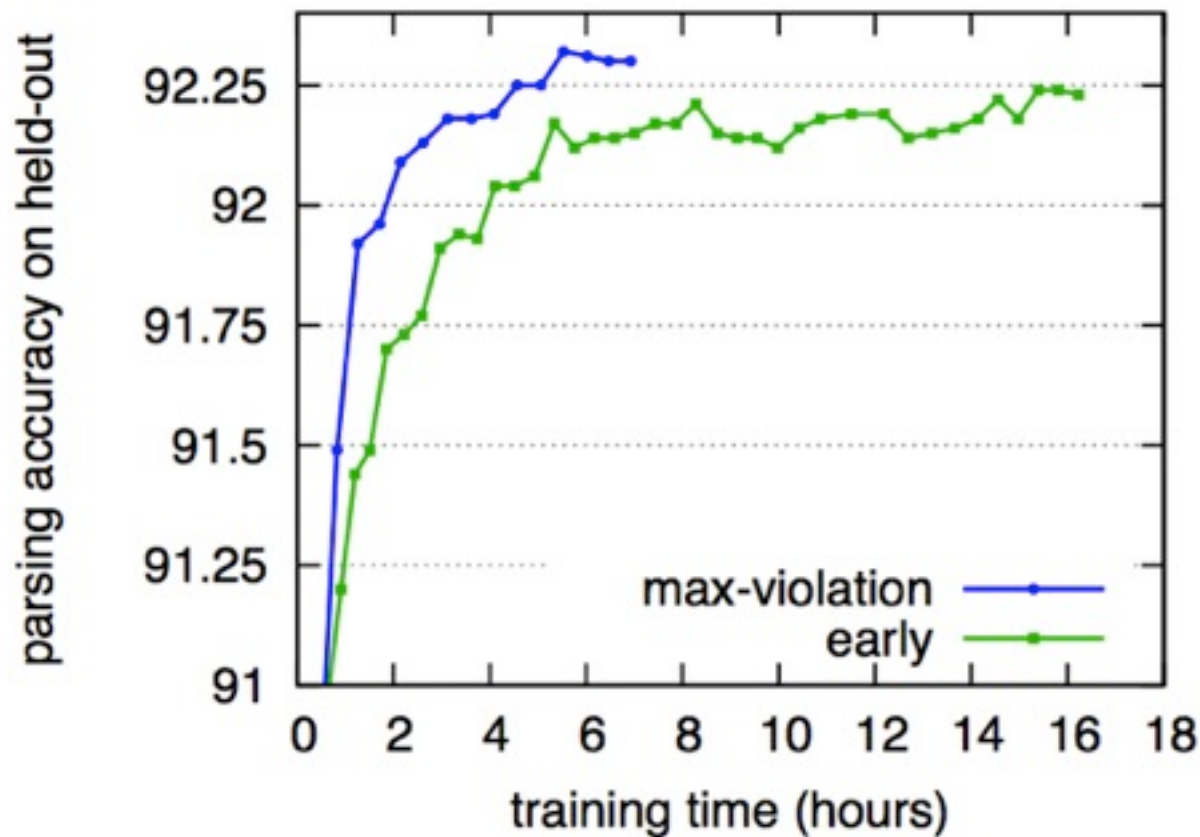
2) Incremental Dependency Parsing

- DP incremental dependency parser (Huang and Sagae 2010)
- non-local history-based features rule out exact DP
 - we use beam search, and search error is severe
 - baseline: early update. extremely slow: 38 iterations



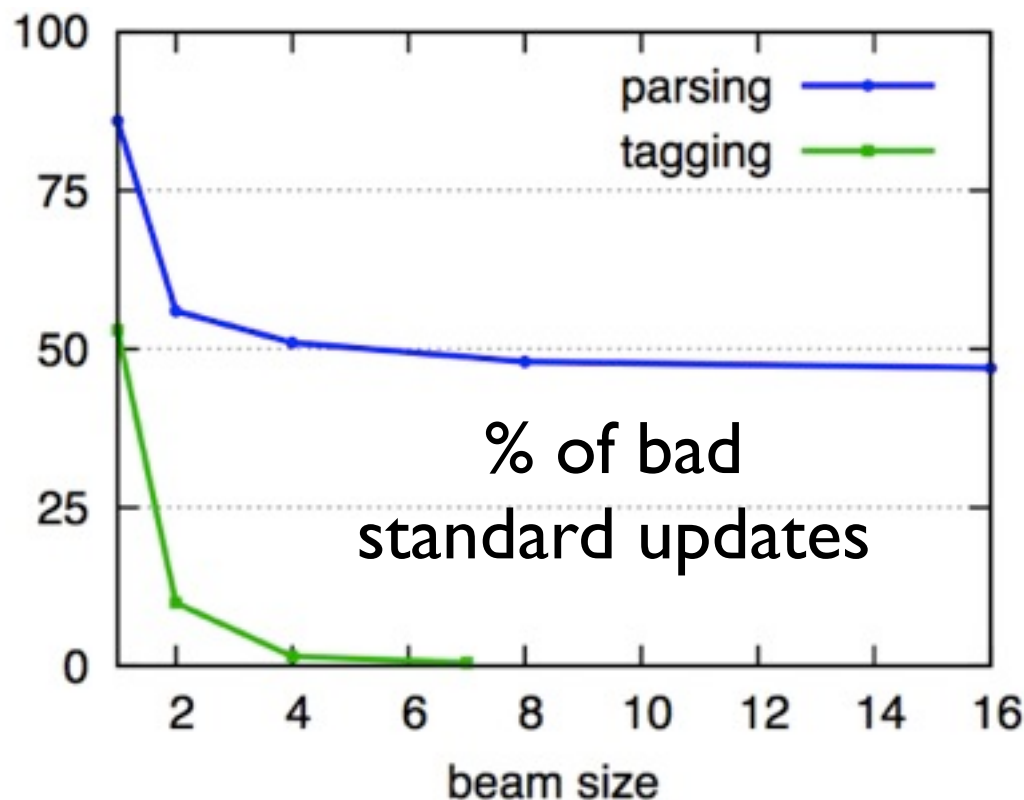
Max-violation converges much faster

- early update: 38 iterations, 15.4 hours (92.24)
- **max-violation**: 10 iterations, 4.6 hours (92.25)
12 iterations, 5.5 hours (92.32)

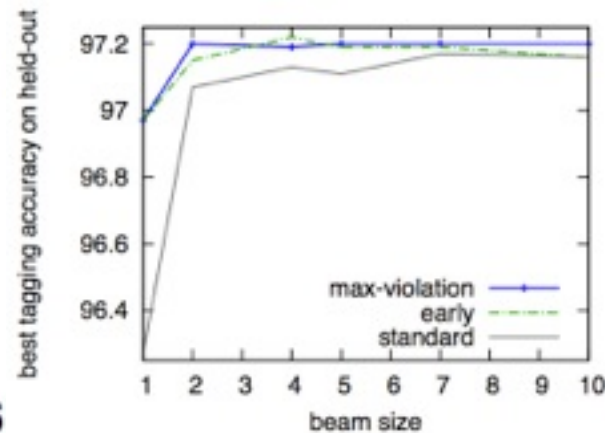


Comparison b/w tagging & parsing

- search error is much more severe in parsing than in tagging
- standard update is OK in tagging except greedy search ($b=1$)
- but performs **horribly** in parsing even at large beam ($b=8$)
 - because $\sim 50\%$ of standard updates are bad (non-violation)!



take-home message:
our methods are more helpful
for harder search problems!



	<i>test</i>
standard	79.1
early	92.1
max-violation	92.2

Annotated Bibliography

- Collins, 2002. Discriminative Training for Hidden Markov Models. In Proceedings of EMNLP. (“structured perceptron”)
- Collins and Roark, 2004. Perceptron Algorithm for Incremental Parsing. In Proceedings of ACL (“early-update”)
- Daume and Marcu, 2005. Learning as Search Optimization. In Proceedings of ICML (“LaSO”).
- Daume, 2006. PhD Thesis. (fast “averaging” trick)
- Huang, Phayong, and Guo, 2012. Structured Perceptron with Inexact Search. In Proceedings of NAACL. (“violation-fixing” framework and proofs, “max-violation”)