

# Machine Learning

Fall 2017

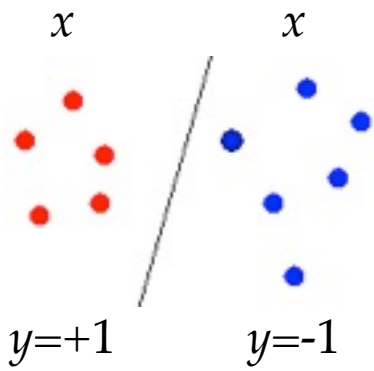
## **Structured Prediction**

(structured perceptron, HMM, structured SVM)

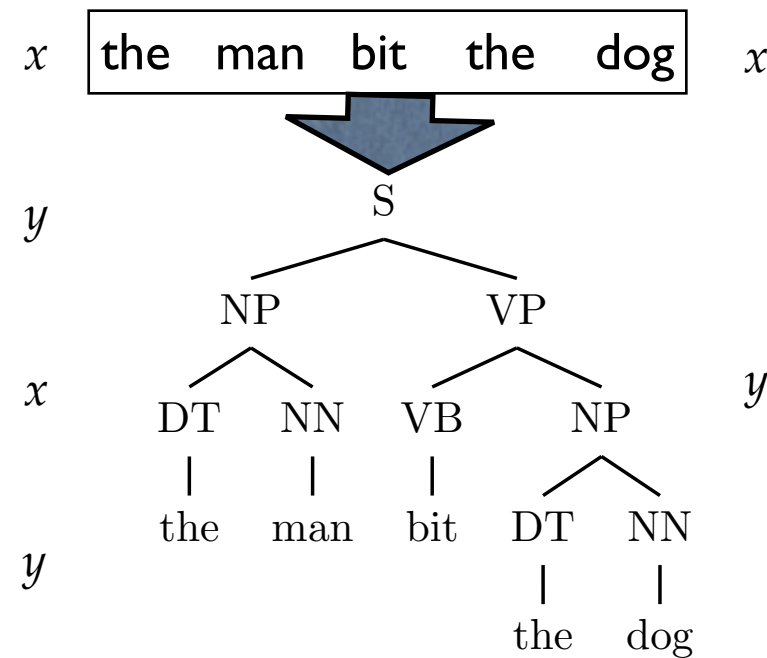
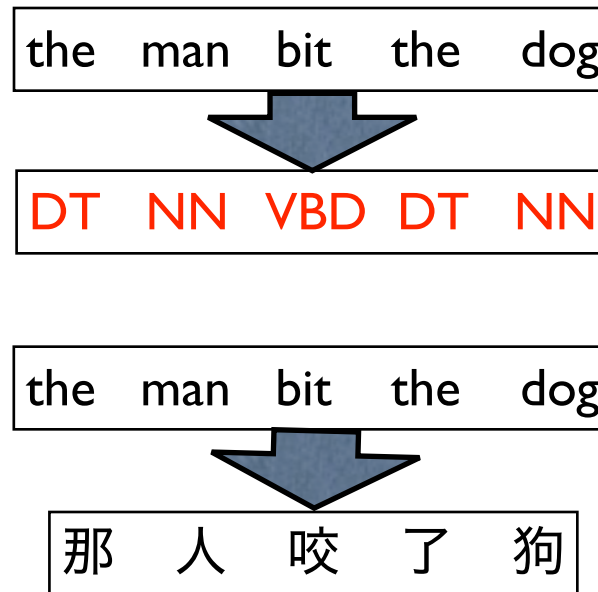
Professor Liang Huang

(Chap. 17 of CIML)

# Structured Prediction



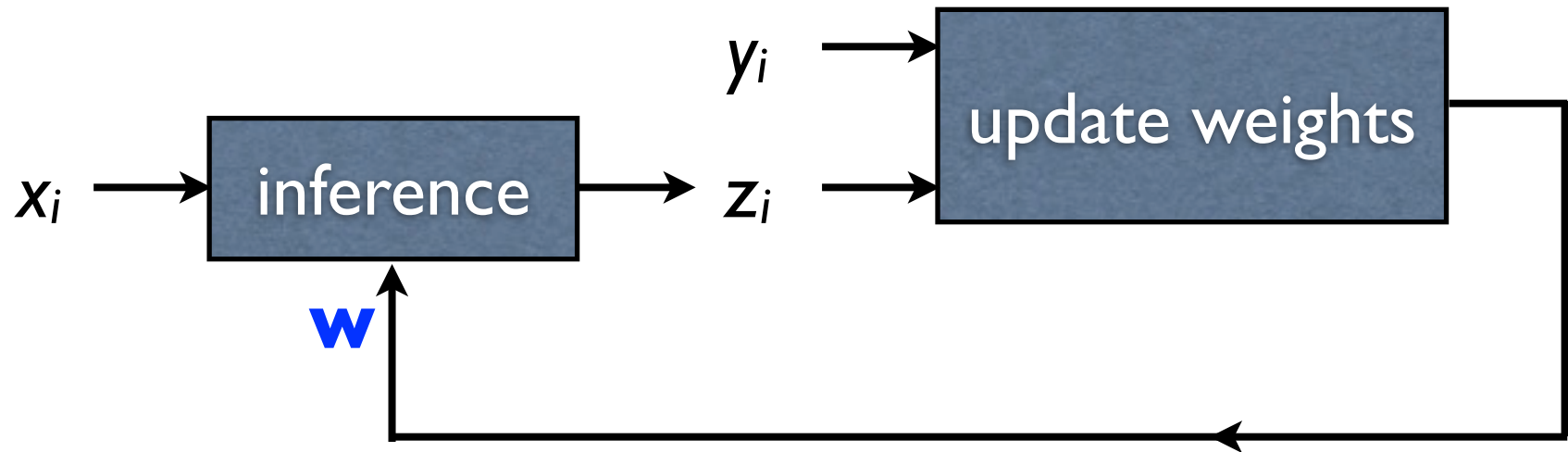
0 1 2 3 4 5 6 7 8 9



- 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!<sub>2</sub>

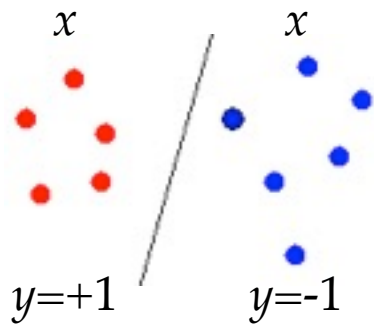
# Generic Perceptron

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



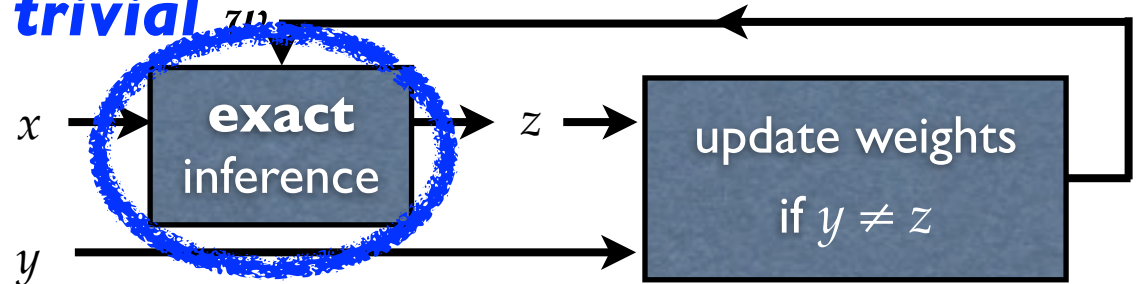
# Perceptron: from binary to structured

## binary classification



2 classes

**trivial**

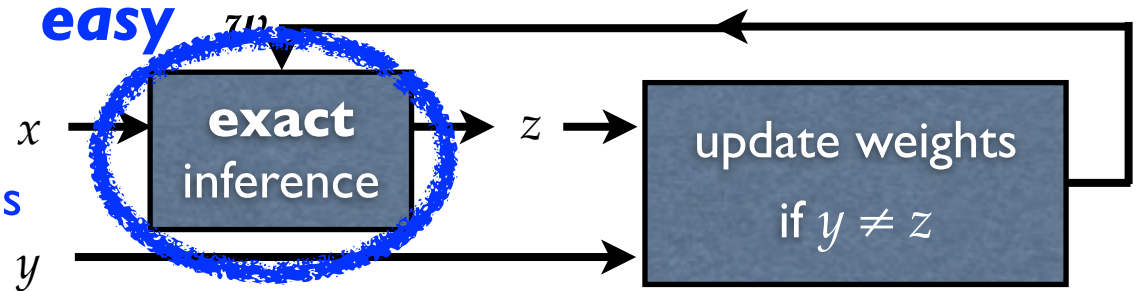


## multiclass classification



constant  
# of classes

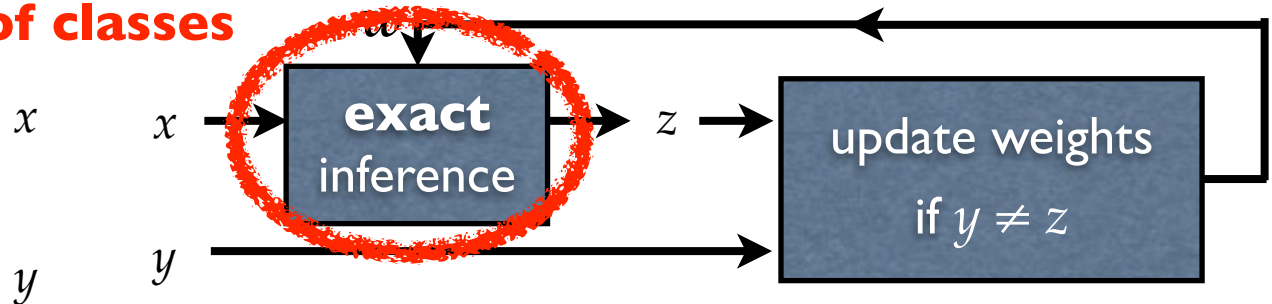
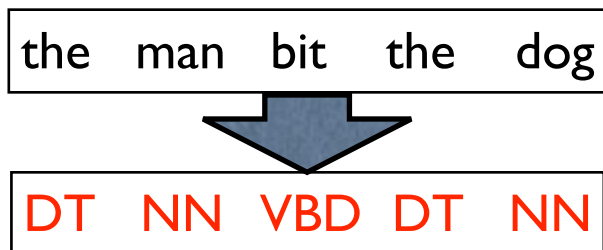
**easy**



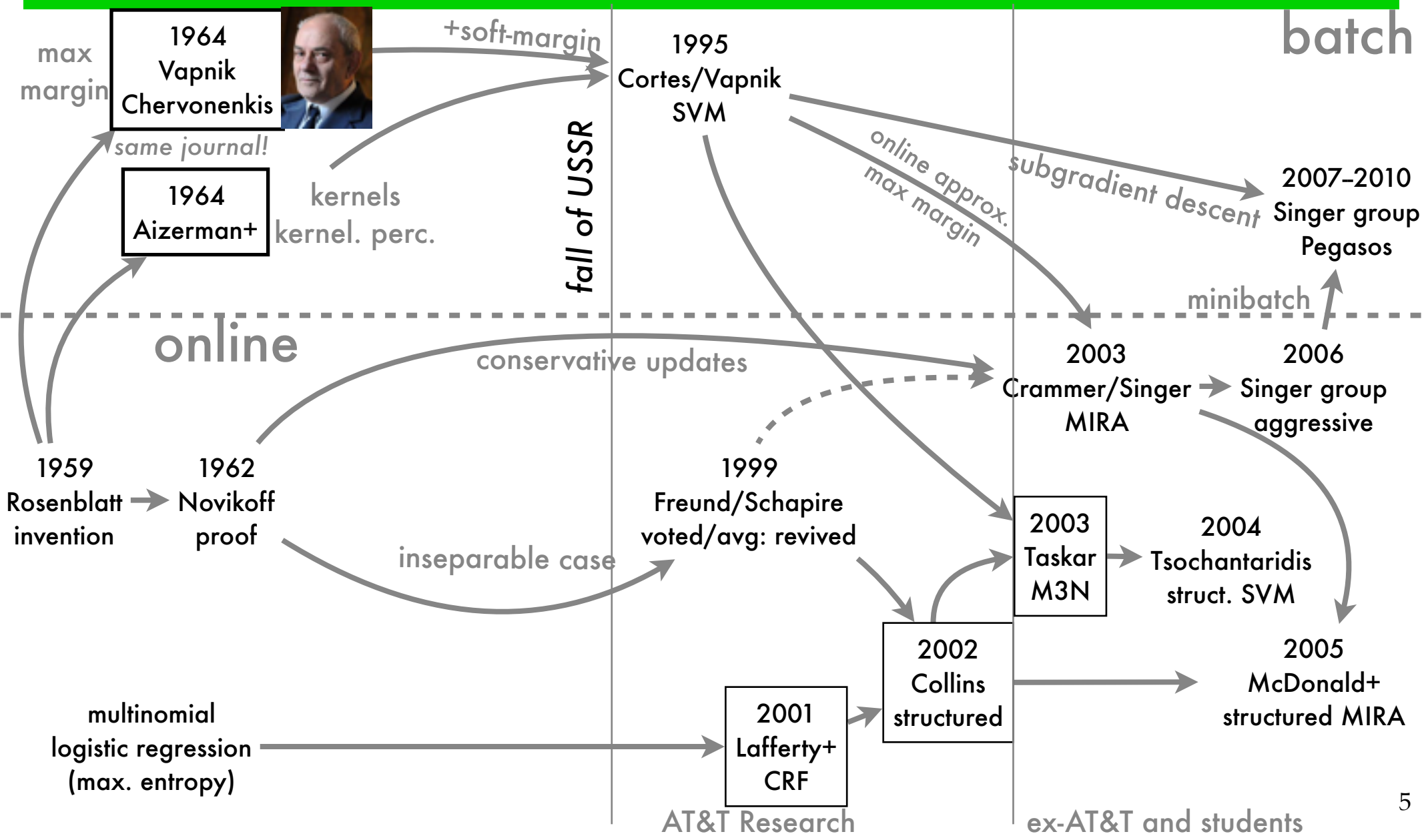
## structured classification

**exponential  
# of classes**

**hard**



# From Perceptron to SVM

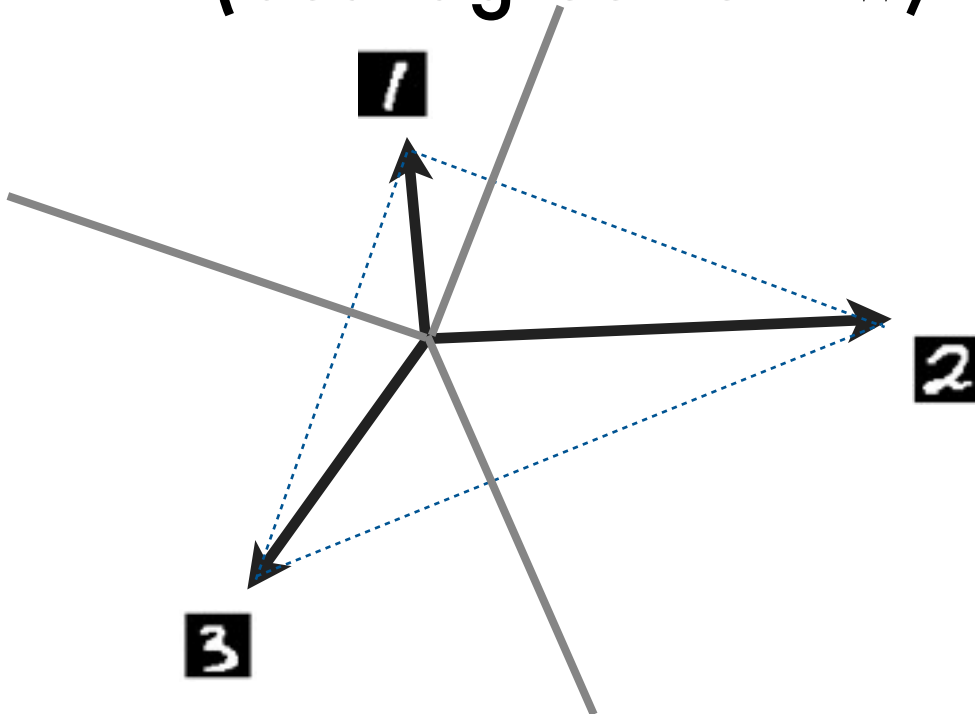


# Multiclass Classification

- 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)



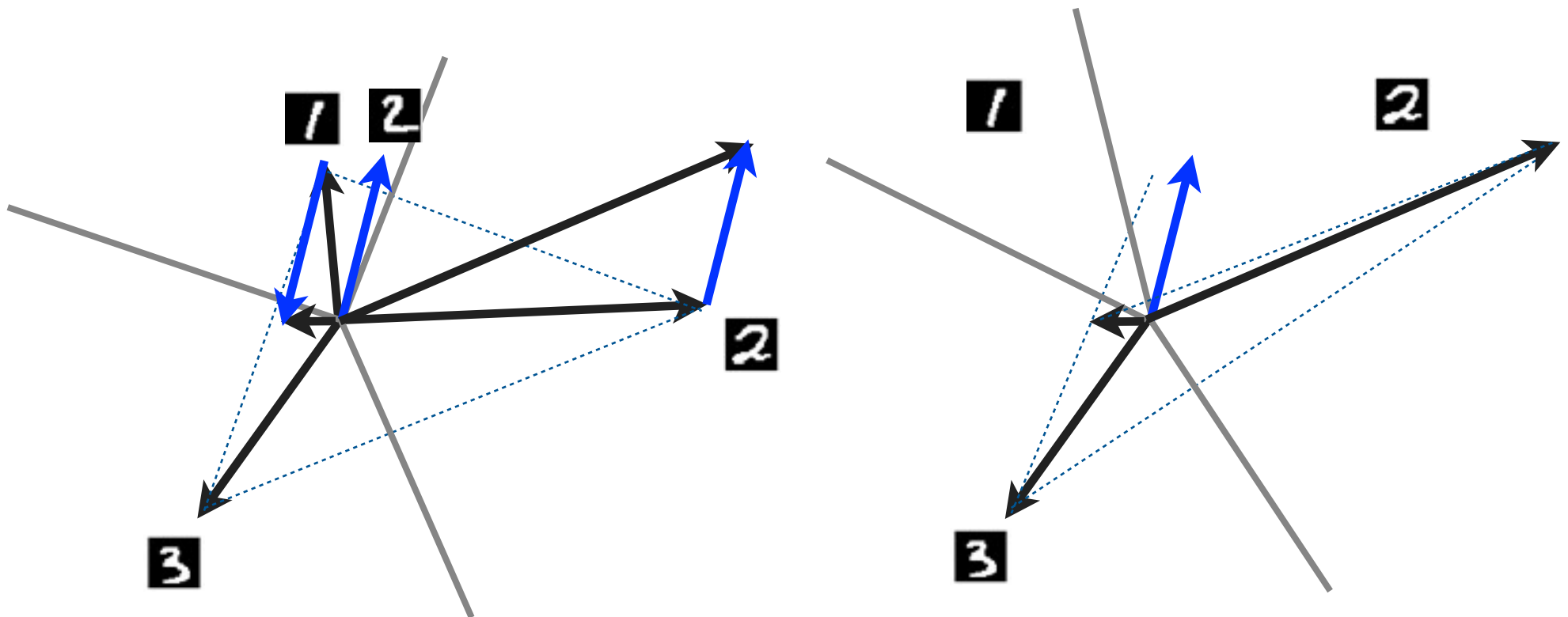
Q1: what about 2-class?

Q2: do we still need  
augmented space?

0 1 2 3 4 5 6 7 8 9

# Multiclass Perceptron

- 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  $\Phi$  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:

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

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



# Example: POS Tagging

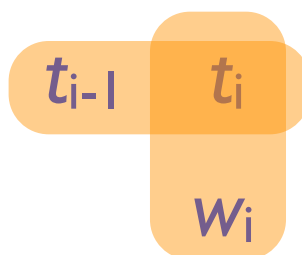
- gold-standard: DT NN VBD DT NN  $y$   
 the man bit the dog  $x$   $\Phi(x, y)$

- 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



$$\phi(x, y) = \{(\langle s \rangle, DT): 1, (DT, NN): 2, \dots, (NN, \langle /s \rangle): 1, (DT, the): 1, \dots, (VBD, bit): 1, \dots, (NN, dog): 1\}$$

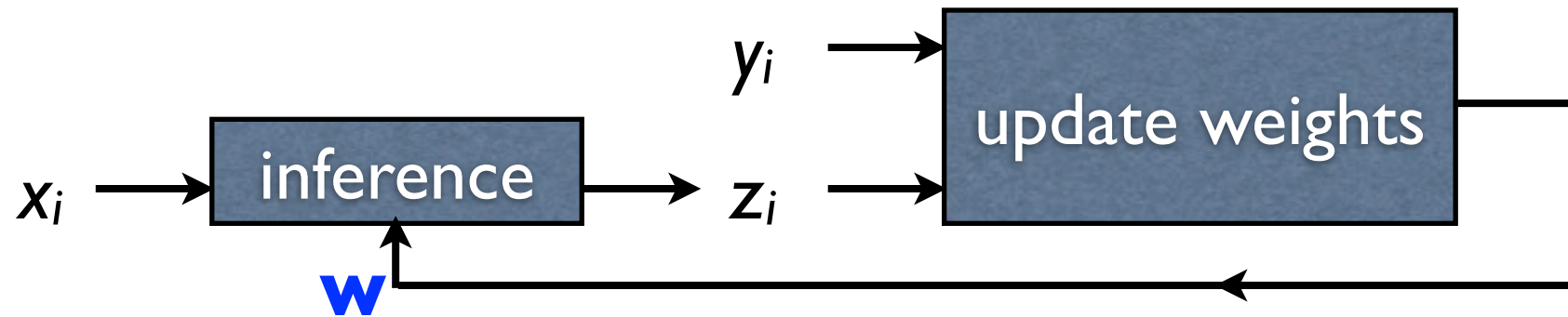
$$\phi(x, z) = \{(\langle s \rangle, DT): 1, (DT, NN): 2, \dots, (NN, \langle /s \rangle): 1, (DT, the): 1, \dots, (NN, bit): 1, \dots, (NN, dog): 1\}$$

- weights ++: (NN, VBD) (VBD, DT) (VBD, bit)

- weights --: (NN, NN) (NN, DT) (NN, bit)

$$\phi(x, y) - \phi(x, z)$$

# 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  $t = 1 \dots T, i = 1 \dots n$

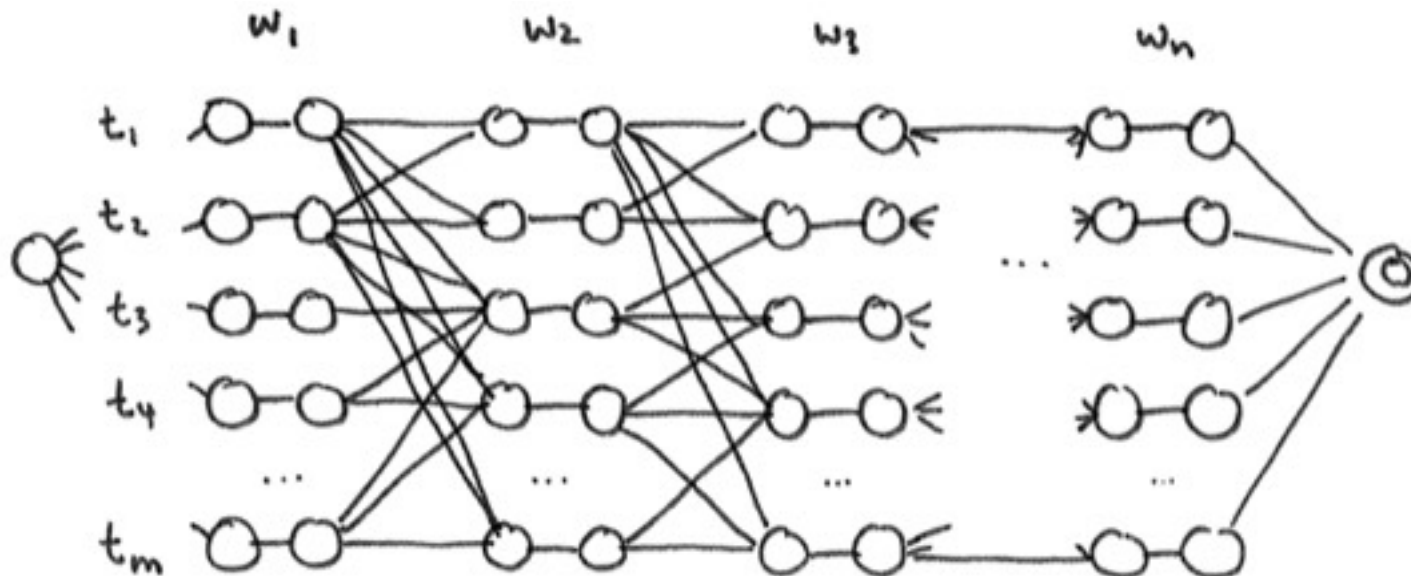
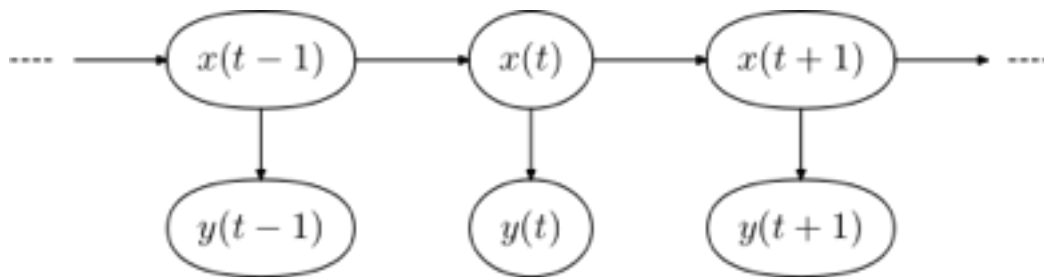
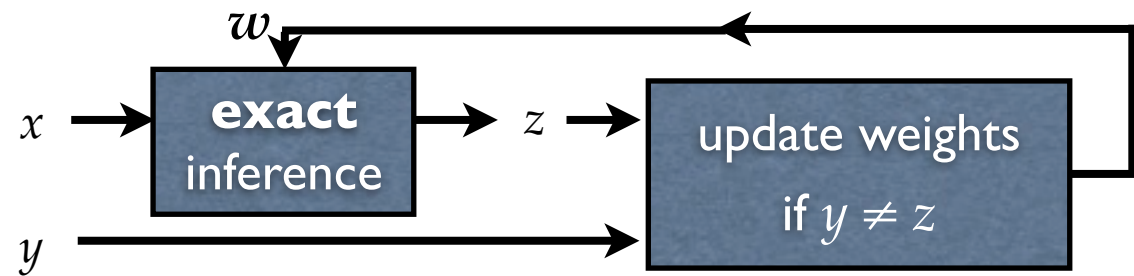
$$z_i = F(x_i)$$

If  $(z_i \neq y_i)$

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

**Output:** Parameters  $\mathbf{W}$

# Inference: Dynamic Programming

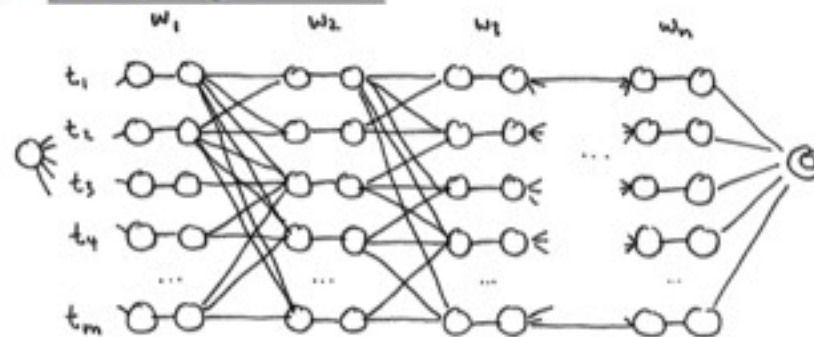


# Python implementation

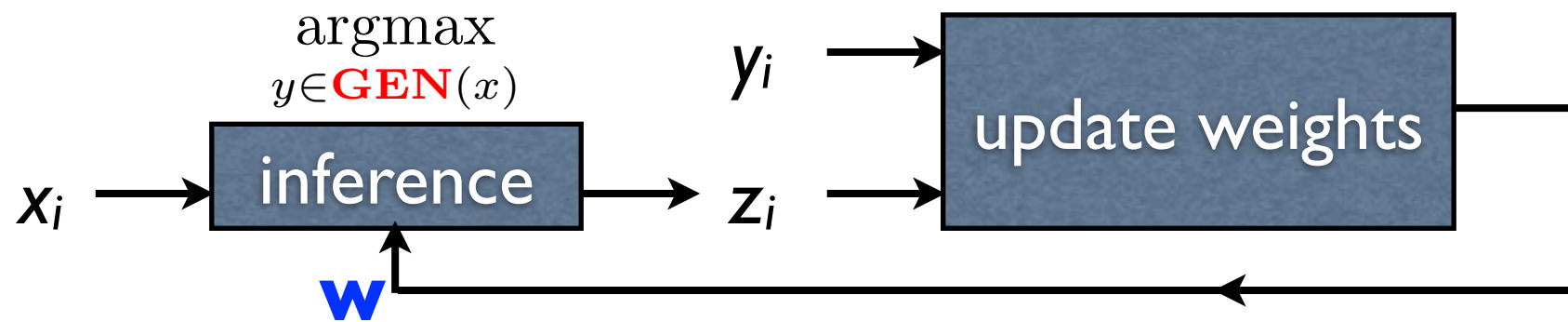
Complete this Python code implementing the Viterbi algorithm for part-of-speech tagging. It should print a list of word/tag pairs, e.g. [('a', 'D'), ('can', 'N'), ('can', 'A'), ('can', 'V'), ('a', 'D'), ('can', 'N')].

```
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]
```

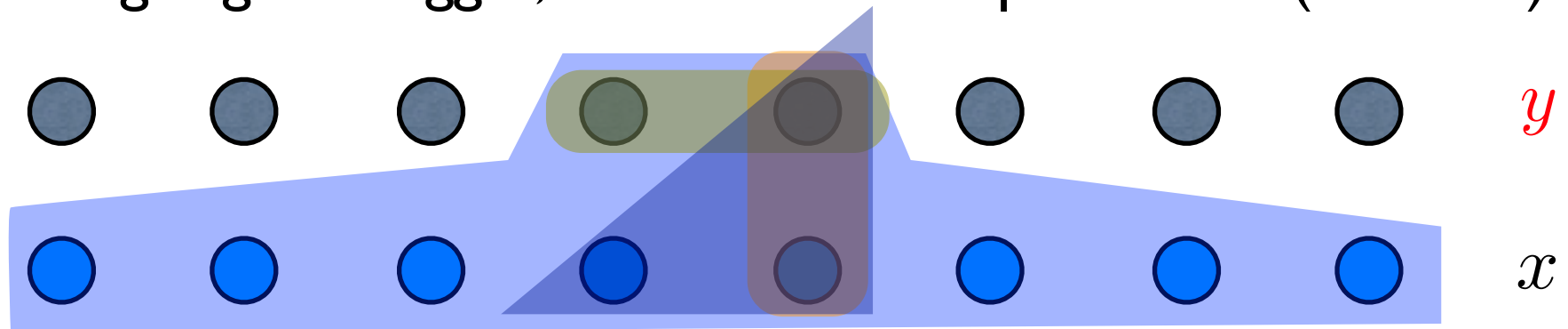
Q: what about top-down recursive + memoization?



# Efficiency vs. Expressiveness



- the **inference** (argmax) must be efficient
  - either the search space  $\text{GEN}(x)$  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) sparse vector: defaultdict

**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  $w_0 - w_a/c$ 

```

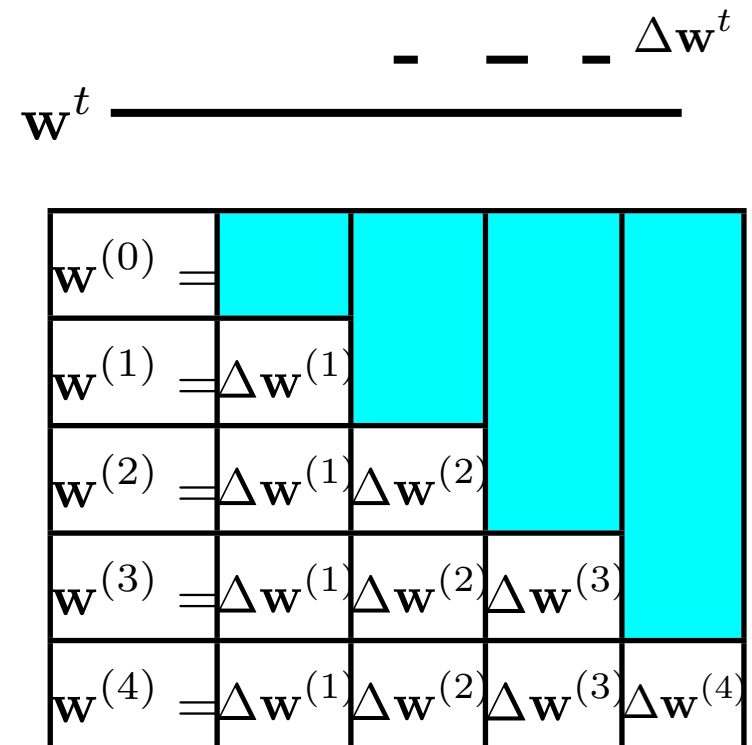
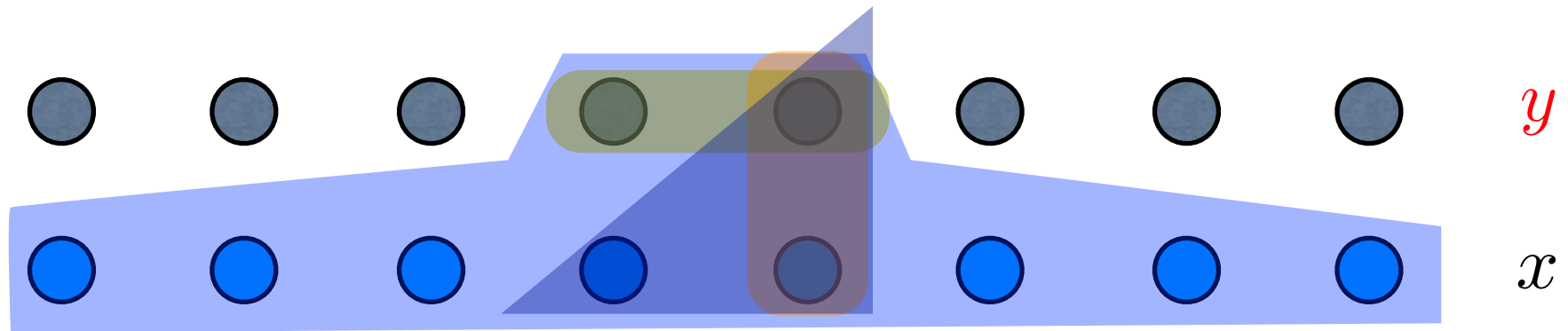


Figure 2.3: The averaged structured perceptron learning algorithm.

# Do we need smoothing?

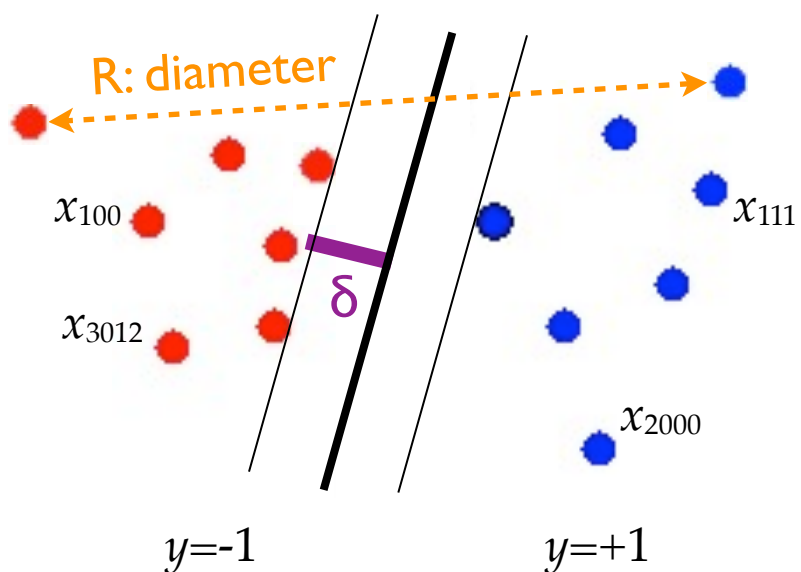


- smoothing is much easier in discriminative models
- just make sure for each feature template, its subset templates are also included
  - e.g., to include  $(t_0 w_0 w_{-1})$  you must also include
    - $(t_0 w_0)$   $(t_0 w_{-1})$   $(w_0 w_{-1})$
    - and maybe also  $(t_0 t_{-1})$  because  $t$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

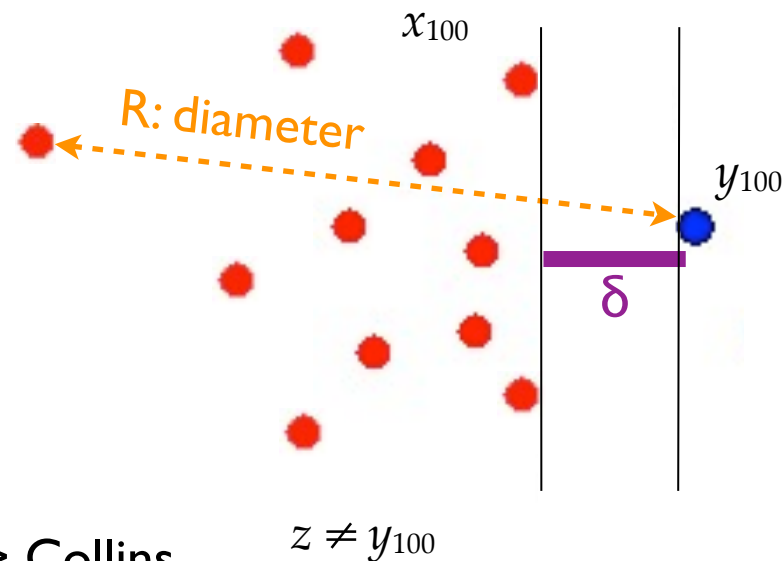


# Convergence with Exact Search

- 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



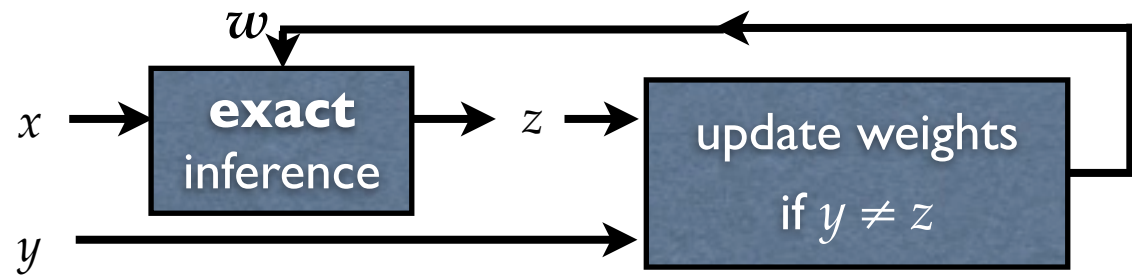
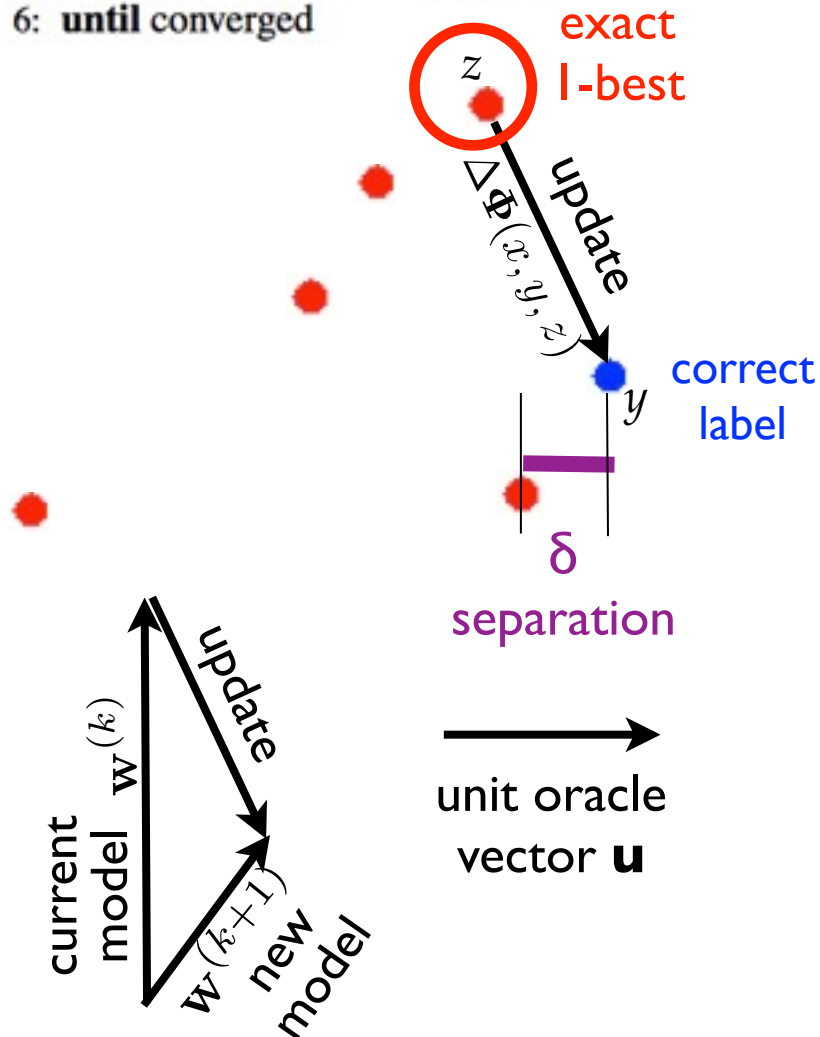
Rosenblatt  $\Rightarrow$  Collins  
1957 2002



# 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:

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

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

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

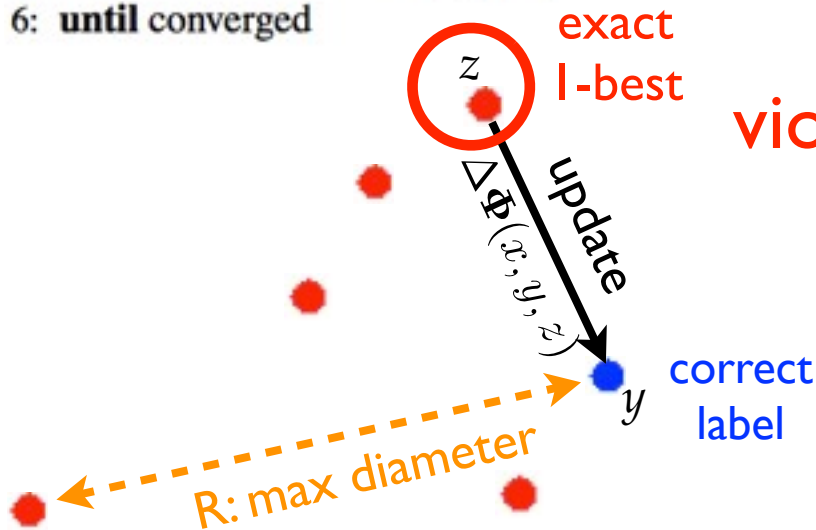
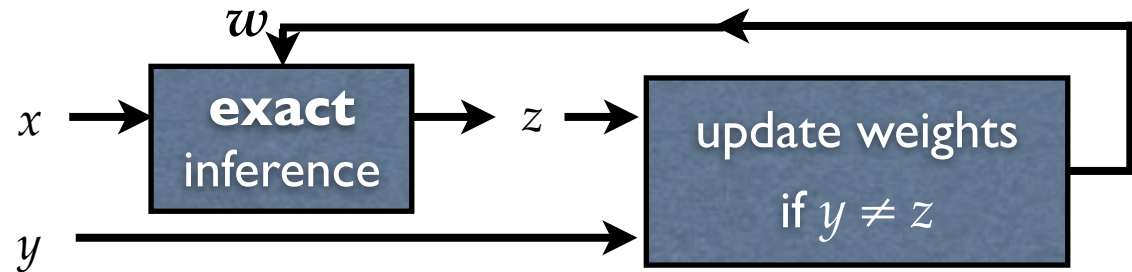
$$\|\mathbf{u}\| \|\mathbf{w}^{(k+1)}\| \geq \mathbf{u} \cdot \mathbf{w}^{(k+1)} \geq k\delta$$

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

# 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
    
```



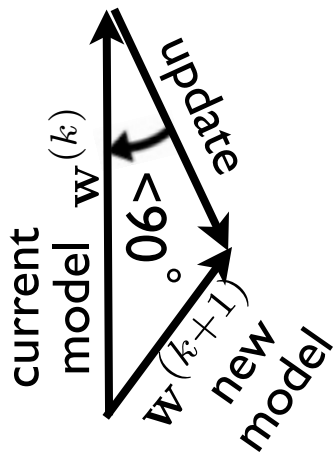
violation: incorrect label scored higher

perceptron update:

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

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

$$= \|\mathbf{w}^{(k)}\|^2 + \underbrace{\|\Delta\Phi(x, y, z)\|^2}_{\leq R^2 \text{ diameter}} + 2 \underbrace{\mathbf{w}^{(k)} \cdot \Delta\Phi(x, y, z)}_{\leq 0 \text{ violation}}$$



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	<b>2.93</b>	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	<b>3.28</b>	200

# Experiments: NP Chunking

- B-I-O scheme

Rockwell International Corp.

B I I O B O

's Tulsa unit said it signed

a tentative agreement ...

- 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	<b>93.53</b>	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	<b>92.65</b>	200

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

# Structured SVM

- structured perceptron:  $w \cdot \Delta\phi(x,y,z) > 0$
- SVM: for all  $(x,y)$ , functional margin  $y(w \cdot x) \geq 1$
- structured SVM version 1: simple loss
  - for all  $(x,y)$ , for all  $z \neq y$ , margin  $w \cdot \Delta\phi(x,y,z) \geq 1$
  - correct  $y$  has to score higher than any wrong  $z$  by 1
- structured SVM version 2: structured loss
  - for all  $(x,y)$ , for all  $z \neq y$ , margin  $w \cdot \Delta\phi(x,y,z) \geq \ell(y,z)$
  - correct  $y$  has to score higher than any wrong  $z$  by  $\ell(y,z)$ , a distance metric such as hamming loss



# Loss-Augmented Decoding

- want for all  $z$ :  $w \cdot \phi(x,y) \geq w \cdot \phi(x,z) + \ell(y,z)$
- same as:  $w \cdot \phi(x,y) \geq \max_z w \cdot \phi(x,z) + \ell(y,z)$
- loss-augmented decoding:  $\operatorname{argmax}_z w \cdot \phi(x,z) + \ell(y,z)$
- if  $\ell(y,z)$  factors in  $z$  (e.g. hamming), just modify DP

---

**Algorithm 41** `STOCHSUBGRADSTRUCTSVM(D, MaxIter,  $\lambda$ ,  $\ell$ )`

---

CIML version

```

1:  $w \leftarrow 0$  // initialize weights
2: for  $iter = 1 \dots \text{MaxIter}$  do
3:   for all  $(x,y) \in D$  do
4:      $\hat{y} \leftarrow \operatorname{argmax}_{\hat{y} \in \mathcal{Y}(x)} w \cdot \phi(x, \hat{y}) + \ell(y, \hat{y})$  // loss-augmented prediction ← modified DP
5:     if  $\hat{y} \neq y$  then
6:        $w \leftarrow w + \phi(x, y) - \phi(x, \hat{y})$  // update weights ← should have learning rate!
7:     end if
8:      $w \leftarrow w - \frac{\lambda}{N} w$  // shrink weights due to regularizer  $\lambda = 1/(2C)$ 
9:   end for
10: end for
11: return  $w$  // return learned weights

```

*very similar to Pegasos; but should use Pegasos framework instead*

# Correct Version following Pegasos

- want for all  $z$ :  $w \cdot \phi(x,y) \geq w \cdot \phi(x,z) + \ell(y,z)$
- same as:  $w \cdot \phi(x,y) \geq \max_z w \cdot \phi(x,z) + \ell(y,z)$
- loss-augmented decoding:  $\operatorname{argmax}_z w \cdot \phi(x,z) + \ell(y,z)$
- if  $\ell(y,z)$  factors in  $z$  (e.g. hamming), just modify DP

---

**Algorithm 41** `STOCHSUBGRADSTRUCTSVM(D, MaxIter,  $\lambda$ ,  $\ell$ )`

---

```
1:  $w \leftarrow 0$  // initialize weights
2: for  $iter = 1 \dots \text{MaxIter}$  do
3:   for all  $(x,y) \in D$  do
4:      $w \leftarrow w - \lambda/t \cdot w$  // shrink weights due to regularizer
5:      $\hat{y} \leftarrow \operatorname{argmax}_{\hat{y} \in \mathcal{Y}(x)} w \cdot \phi(x, \hat{y}) + \ell(y, \hat{y})$  // loss-augmented prediction
6:     if  $\hat{y} \neq y$  then
7:        $w \leftarrow w + \lambda C/2t (\phi(x, y) - \phi(x, \hat{y}))$  // update weights
8:     end if
9:   end for
10: end for
11: return  $w$  // return learned weights
```

---

$N=|D|$ ,  $C$  is from SVM  
 $t += 1$  for each example

# Struct. Perceptron vs Struct. SVM

- tagging, ATIS (train: 488 sent); SVM < avg perc << perc

## perceptron

epoch 1	updates 102,	W =291,	train_err 3.90%,	dev_err 9.36%	avg_err 6.14%
epoch 2	updates 91,	W =334,	train_err 3.33%,	dev_err 8.19%	avg_err <b>4.97%</b>
epoch 3	updates 78,	W =347,	train_err 2.92%,	dev_err <b>5.85%</b>	avg_err 4.97%
epoch 4	updates 81,	W =368,	train_err 3.11%,	dev_err 6.73%	avg_err 5.85%
epoch 5	updates 78,	W =378,	train_err 2.70%,	dev_err 6.14%	avg_err 5.56%
epoch 6	updates 63,	W =385,	train_err 2.26%,	dev_err 6.14%	avg_err 5.56%
epoch 7	updates 69,	W =385,	train_err 2.43%,	dev_err 7.02%	avg_err 5.56%
epoch 8	updates 60,	W =388,	train_err 2.15%,	dev_err 6.73%	avg_err 5.56%
epoch 9	updates 59,	W =390,	train_err 2.04%,	dev_err 6.14%	avg_err 5.56%
epoch 10	updates 64,	W =394,	train_err 2.15%,	dev_err 5.85%	avg_err 5.26%

## SVM C=1

epoch 1	updates 116,	W =311,	train_err 4.55%,	dev_err 5.85%
epoch 2	updates 82,	W =328,	train_err 3.05%,	dev_err 4.97%
epoch 3	updates 78,	W =334,	train_err 2.92%,	dev_err 5.56%
epoch 4	updates 77,	W =339,	train_err 2.92%,	dev_err 5.26%
epoch 5	updates 80,	W =344,	train_err 2.94%,	dev_err 5.56%
epoch 6	updates 73,	W =345,	train_err 2.75%,	dev_err <b>4.68%</b>
epoch 7	updates 72,	W =347,	train_err 2.75%,	dev_err 4.97%
epoch 8	updates 75,	W =352,	train_err 2.86%,	dev_err 4.97%
epoch 9	updates 74,	W =353,	train_err 2.78%,	dev_err 4.97%
epoch 10	updates 72,	W =354,	train_err 2.78%,	dev_err 4.97%