



Department of Computer Science

UNIVERSITY OF COLORADO **BOULDER**



Structure Learning

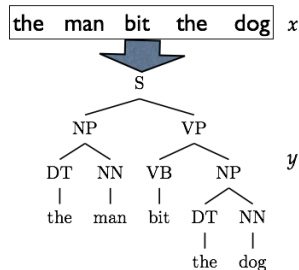
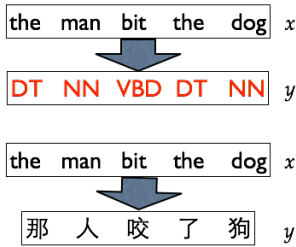
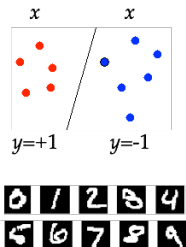
Jordan Boyd-Graber
University of Colorado Boulder

8. DECEMBER 2014

Slides adapted from Liang Huang

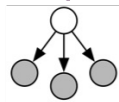
Roadmap

- Structured learning
- Alternative to generative models





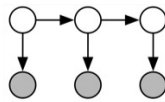
binary/multiclass



naive bayes



structured learning



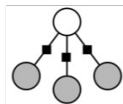
HMMs

generative

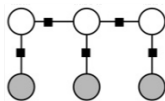
(count & divide)

Conditional

Conditional



logistic regression (maxent)



CRFs

discriminative

(expectations)

Online+
Viterbi

Online+
Viterbi

perceptron



structured perceptron

(argmax)

max
margin

max
margin

SVM

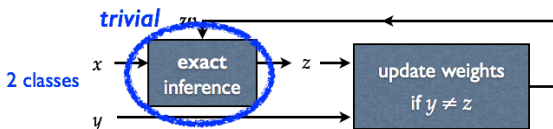
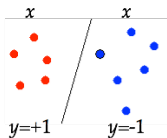


structured SVM

(loss-augmented
argmax)

Binary to Structure

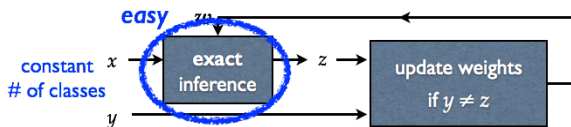
binary perceptron
(Rosenblatt, 1959)



Binary to Structure

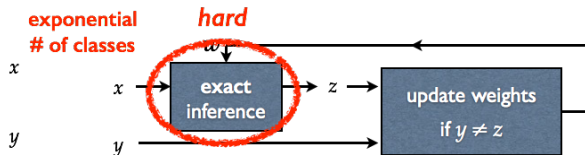
multiclass perceptron
(Freund/Schapire, 1999)

0 1 2 3 4 5 6 7 8 9



Binary to Structure

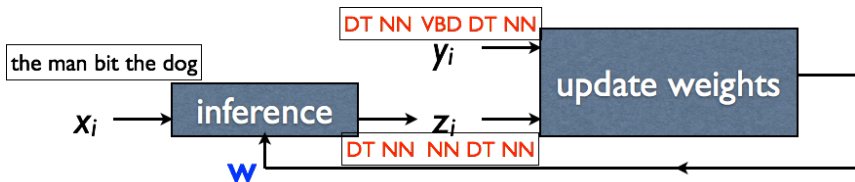
structured perceptron
(Collins, 2002)



Generic Perceptron

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

Structured Perceptron



Perceptron Algorithm

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}

POS Example

• 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

t_{i-1}

t_i

• word/tag pairs

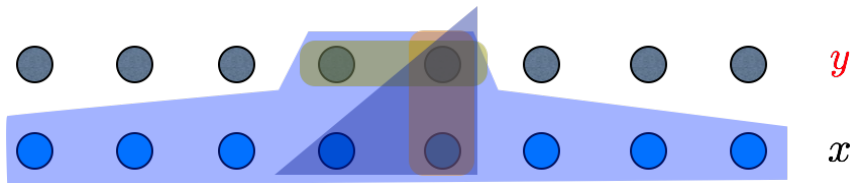
w_i

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

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

What must be true?

- Finding highest scoring structure must be really fast (you'll do it often)
- Requires some sort of dynamic programming algorithm
- For tagging: features must be local to y (but can be global to x)



Averaging is Good

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

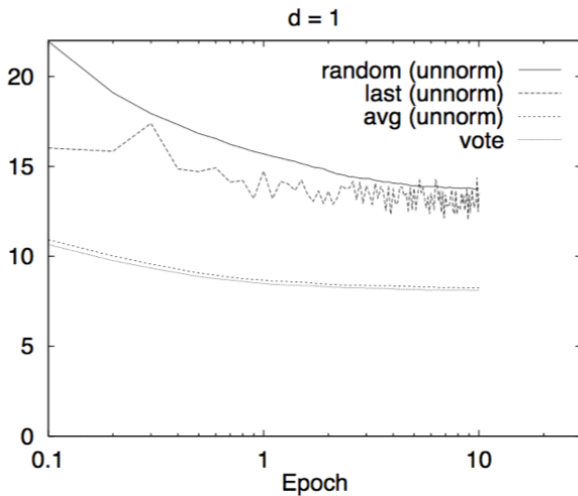
Initialization: $\mathbf{W}_0 = 0$

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

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 $\mathbf{W} = \sum_j \mathbf{W}_j$

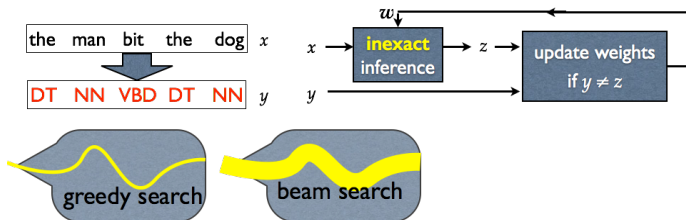
Averaging is Good



Smoothing

- Must include subset templates for features
- For example, if you have feature (t_0, w_0, w_{-1}) , you must also have
 - (t_0, w_0) ; (t_0, w_{-1}) ; (w_0, w_{-1})

Inexact Search?



- Sometimes search is too hard
- So we use beam search instead
- How to create algorithms that respect this relaxation: track when right answer falls off the beam

Structured Learning

- Sometimes you want discriminative method for complex y
- Learning those models are difficult
- Need to be scalable