



10-708 Probabilistic Graphical Models

Machine Learning Department
School of Computer Science
Carnegie Mellon University



Structured Perceptron

+

Structured SVM

Matt Gormley
Lecture 10
Mar. 3, 2021

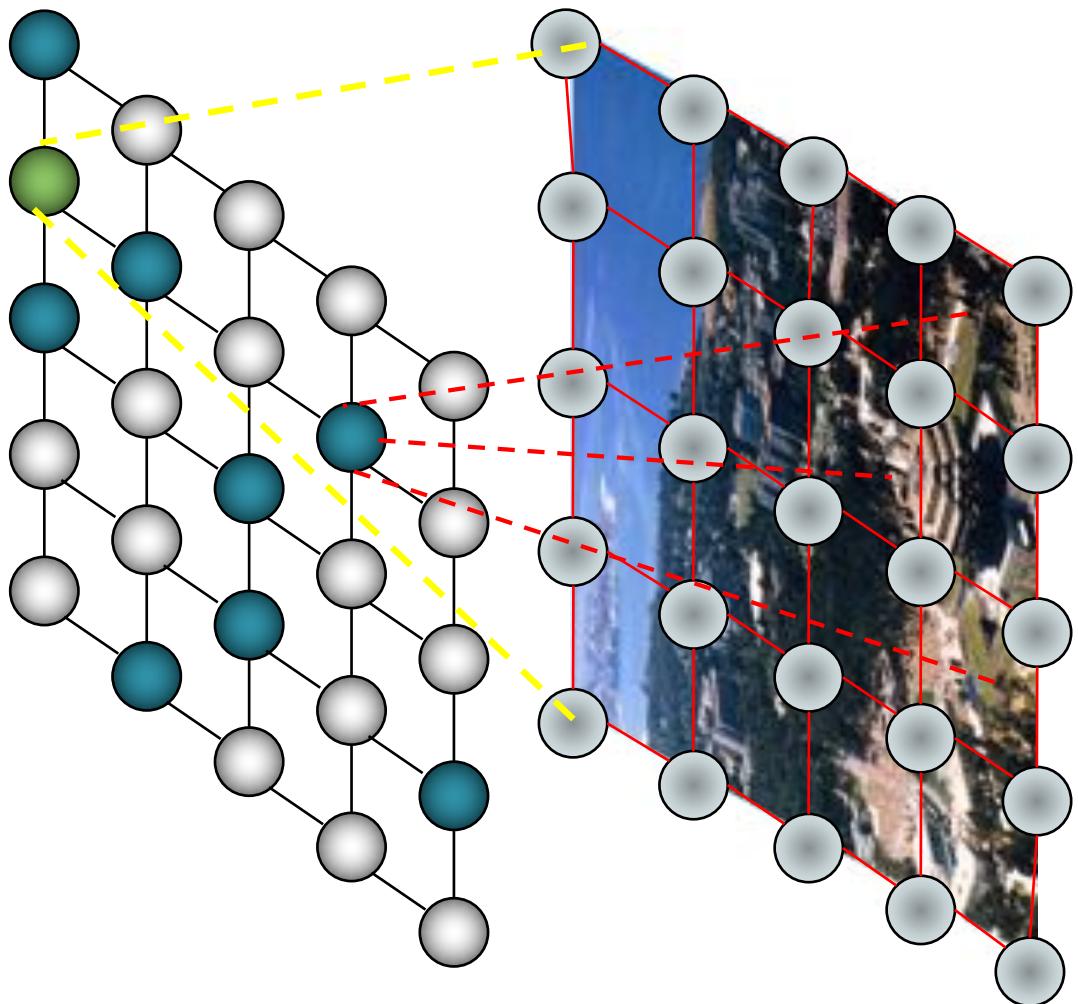
Reminders

- **Homework 2: Exact inference and supervised learning (CRF+RNN)**
 - Out: Wed, Feb. 24
 - Due: Wed, Mar. 10 at 11:59pm
- **Schedule Changes**
 - Quiz 1: Mon, Mar. 15
 - Quiz 2: Mon, Apr. 12
 - Friday Lectures on Apr. 02 and Apr. 23
 - More rational ordering for HWs and quizzes

Case #2: Multiclass Variables

MAP INFERENCE AS MATHEMATICAL PROGRAMMING

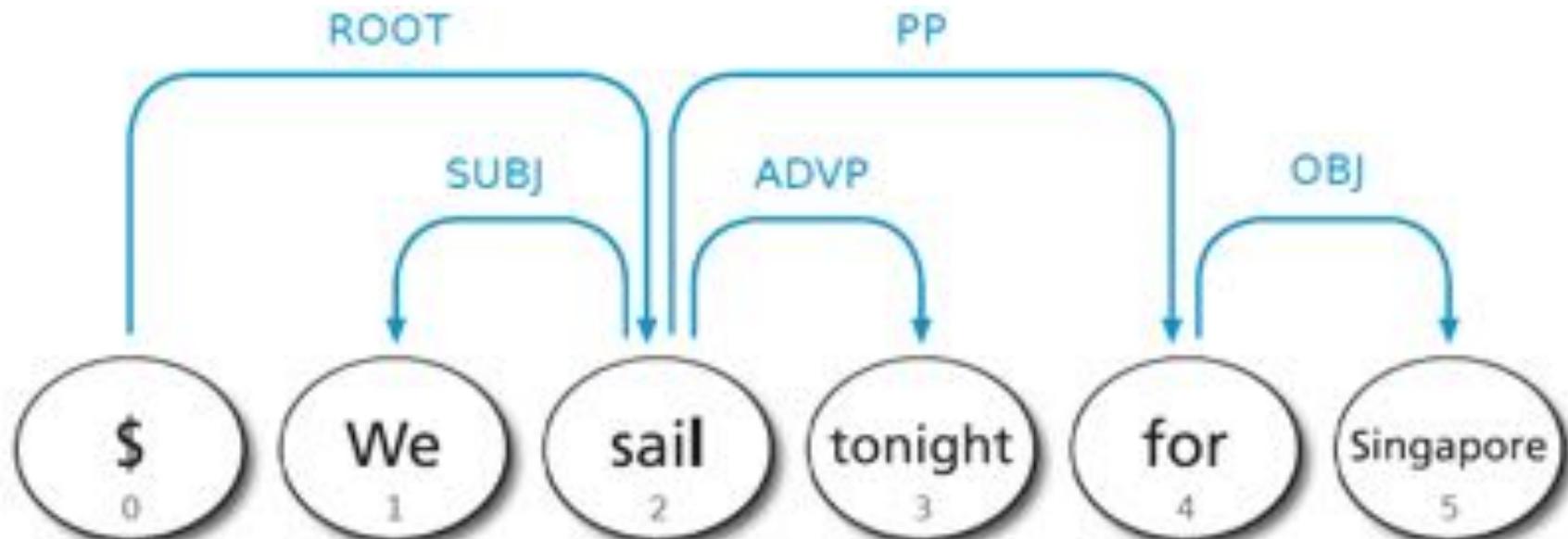
Image Segmentation



$$p_{\theta}(y|x) = \frac{1}{Z(\theta, x)} \exp \left\{ \sum_c \theta_c f_c(x, y_c) \right\}$$

- Jointly segmenting/annotating images
- Image-image matching, image-text matching
- Problem:
 - Given structure (feature), learning $\vec{\theta}$
 - Learning sparse, interpretable, **predictive** structures/features

Dependency parsing of Sentences



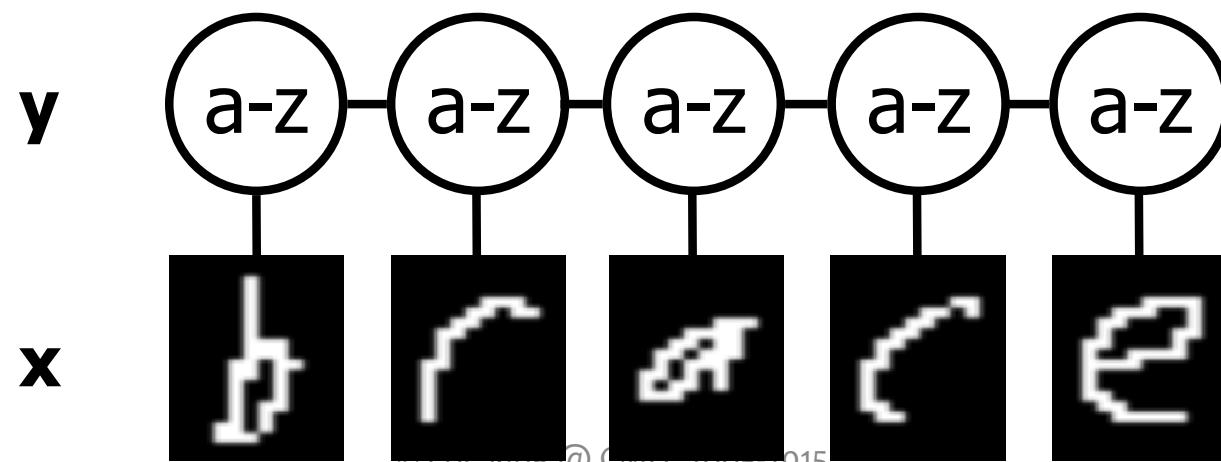
Challenge:

Structured outputs, and globally constrained to be a valid tree

OCR example



Sequential structure

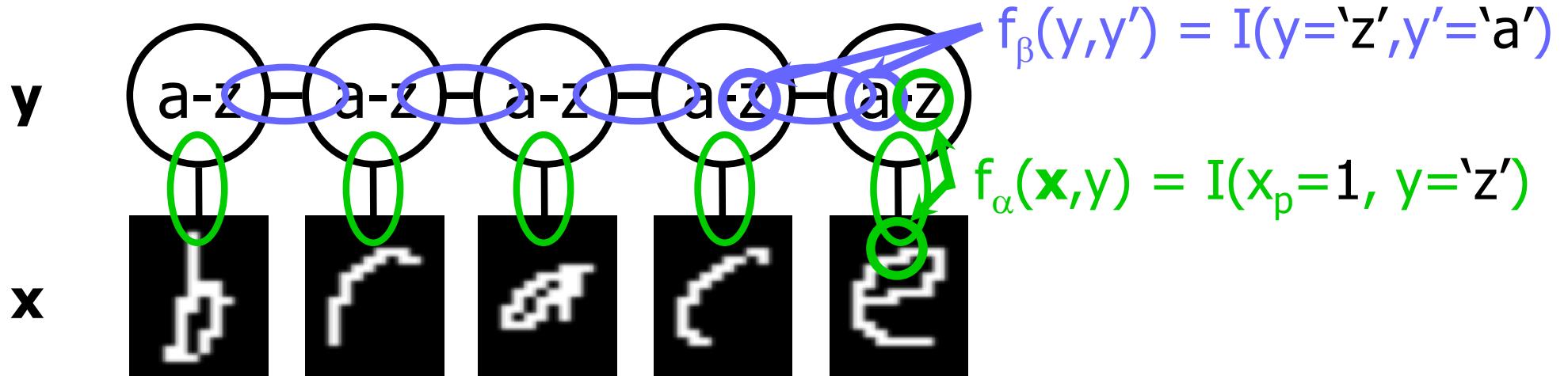


Linear-chain CRF for OCR

$$P(\mathbf{y} | \mathbf{x}) = \frac{1}{Z(\mathbf{x})} \underbrace{\prod_i \phi(\mathbf{x}_i, y_i)}_{\text{green}} \underbrace{\prod_i \phi(y_i, y_{i+1})}_{\text{blue}}$$

$$\phi(\mathbf{x}_i, y_i) = \exp\{\sum_\alpha w_\alpha f_\alpha(\mathbf{x}_i, y_i)\}$$

$$\phi(y_i, y_{i+1}) = \exp\{\sum_\beta w_\beta f_\beta(y_i, y_{i+1})\}$$



$y \Rightarrow z$ map for linear chain structures

OCR example: $y = 'ABABB'$;

z 's are the indicator variables for the corresponding classes (alphabet)

	$z_1(m)$	$z_2(m)$	$z_3(m)$	$z_4(m)$	$z_5(m)$
A	1	0	1	0	0
B	0	1	0	1	1
:	:	:	:	:	:
Z	0	0	0	0	0

	$z_{12}(m, n)$	$z_{23}(m, n)$	$z_{34}(m, n)$	$z_{45}(m, n)$
A	0 1 . 0	0 0 . 0	0 1 . 0	0 0 . 0
B	0 0 . 0	1 0 . 0	0 0 . 0	0 1 . 0
:	. . . 0	. . . 0	. . . 0	. . . 0
Z	0 0 0 0	0 0 0 0	0 0 0 0	0 0 0 0

A	B	.	Z
A	B	.	Z

$y \Rightarrow z$ map for linear chain structures

$$\max_{\mathbf{y}} \sum_j \mathbf{w}^T f_{\text{node}}(x_j, y_j) + \sum_{j,k} \mathbf{w}^T f_{\text{edge}}(\mathbf{x}_{jk}, y_j, y_k)$$

Rewriting the maximization function in terms of indicator variables:

$$\begin{aligned} \max_{\mathbf{z}} \quad & \sum_{j,m} z_j(m) [\mathbf{w}^\top \mathbf{f}_{\text{node}}(\mathbf{x}_j, m)] \\ & + \sum_{jk,m,n} z_{jk}(m, n) [\mathbf{w}^\top \mathbf{f}_{\text{edge}}(\mathbf{x}_{jk}, m, n)] \end{aligned}$$

$$z_k(n)$$

$$z_j(m) \geq 0; z_{jk}(m, n) \geq 0;$$

0	1	0	0
---	---	---	---

normalization

$$\sum_m z_j(m) = 1$$

$$z_j(m)$$

0
0
1
0

0	0	0	0
0	0	0	0
0	1	0	0
0	0	0	0

agreement

$$\sum_n z_{jk}(m, n) = z_j(m)$$

integer

$$z_j(m) \in \mathcal{Z}, z_{jk}(m, n) \in \mathcal{Z}$$

$$z_{jk}(m, n)$$

$y \Rightarrow z$ map for linear chain structures

$$\max_{\mathbf{y}} \sum_j \mathbf{w}^T f_{\text{node}}(x_j, y_j) + \sum_{j,k} \mathbf{w}^T f_{\text{edge}}(\mathbf{x}_{jk}, y_j, y_k)$$

Rewriting the maximization function in terms of indicator variables:

$$\begin{aligned} \max_{\mathbf{z}} \quad & \sum_{j,m} z_j(m) [\mathbf{w}^\top \mathbf{f}_{\text{node}}(\mathbf{x}_j, m)] \\ & + \sum_{jk,m,n} z_{jk}(m, n) [\mathbf{w}^\top \mathbf{f}_{\text{edge}}(\mathbf{x}_{jk}, m, n)] \end{aligned} \quad \left. \right\} (\mathbf{F}^\top \mathbf{w})^\top \mathbf{z}$$

$$\begin{array}{c} z_k(n) \\ \hline \begin{matrix} 0 & 1 & 0 & 0 \end{matrix} \end{array} \quad \begin{array}{c} z_j(m) \\ \hline \begin{matrix} 0 \\ 0 \\ 1 \\ 0 \end{matrix} \end{array} \quad \begin{array}{c} z_j(m) \geq 0; z_{jk}(m, n) \geq 0; \\ \text{normalization} \quad \sum_m z_j(m) = 1 \end{array} \quad \left. \right\} \mathbf{A}\mathbf{z} = \mathbf{b}$$

$$\begin{array}{c} z_{jk}(m, n) \\ \hline \begin{matrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{matrix} \end{array}$$

$$\text{agreement} \quad \sum_n z_{jk}(m, n) = z_j(m)$$

$$\max_{A\mathbf{z}=\mathbf{b}} (\mathbf{F}^\top \mathbf{w})^\top \mathbf{z}$$

MAP Inference

Suppose we want to predict the highest likelihood structure y , given observations x and parameters w .

$$\begin{aligned}\hat{\mathbf{y}} &= \operatorname{argmax}_{\mathbf{y}} \log p_w(\mathbf{y}|\mathbf{x}) \\ &= \operatorname{argmax}_{\mathbf{y}} \sum_j \mathbf{w}^T f_{\text{node}}(x_j, y_j) + \sum_{j,k} \mathbf{w}^T f_{\text{edge}}(\mathbf{x}_{jk}, y_j, y_k)\end{aligned}$$

Idea:

1. Reformulate the problem as an integer linear program (ILP) – note that this is just going to be a new way of writing down the problem: $\mathbf{y} \rightarrow \mathbf{z}$
2. Then remove the integer constraints (i.e. solve the linear program (LP) relaxation)

Lemma: (Wainwright et al., 2002) If there is a unique MAP assignment, the LP relaxation of the ILP above is guaranteed to have an integer solution, which is exactly the MAP solution!

STRUCTURED PERCEPTRON

Linear Models

Setting: training examples are (\vec{x}, y)
where $\vec{x} \in \mathbb{R}^P$ $y \in \{1, \dots, K\}$

Model: parameters $\Theta \in \mathbb{R}^M$
feature function $f(\vec{x}, y) \in \mathbb{R}^M$

Predict: $\hat{y} = h_\Theta(\vec{x}) = \underset{y \in \{1, \dots, K\}}{\operatorname{argmax}} \Theta^T f(\vec{x}, y)$

Ex#1: $f(\vec{x}, y) = \text{vectorize} \begin{pmatrix} 0 & 0 & \dots & 0 \\ \vdots & & & \\ x_1 & x_2 & \dots & x_P \\ \vdots & & & \\ 0 & 0 & \dots & 0 \end{pmatrix}$

KxP matrix

only the y^{th} row is non-zero

$$= [0 0 \dots 0 \ x_1 \ x_2 \ \dots \ x_P \ 0 0 \dots 0]^T$$

$\Rightarrow M = K \times P$

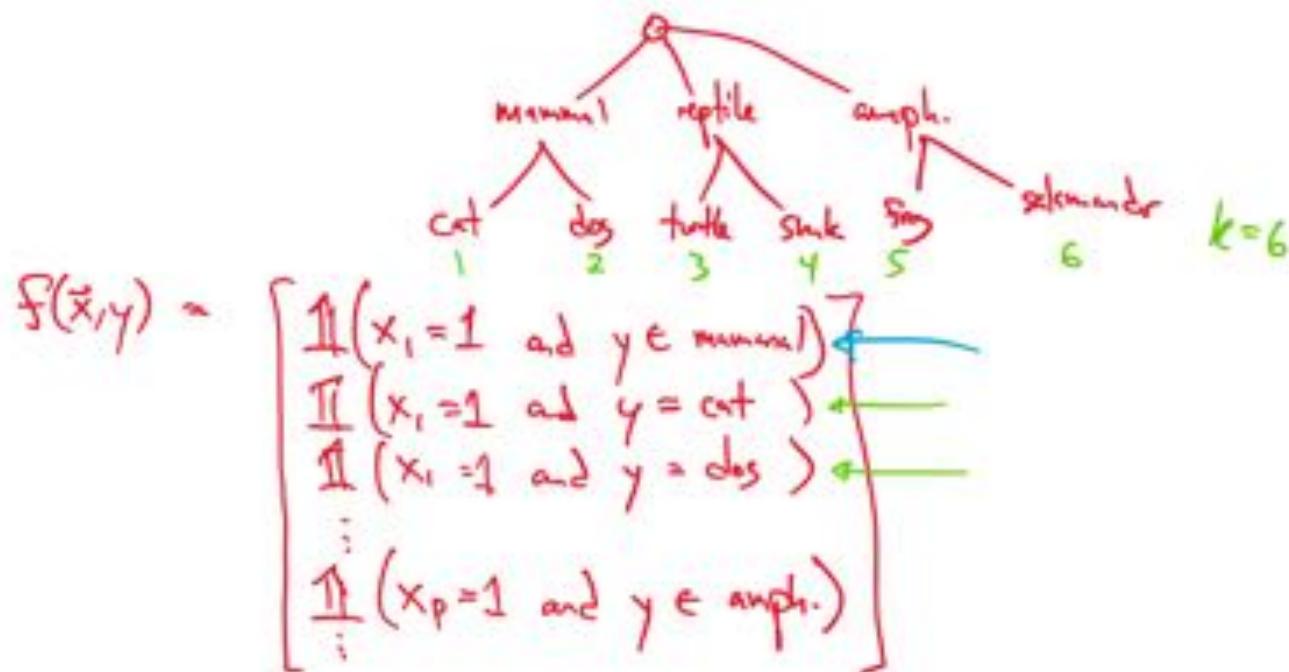
Linear Models

Setting: training examples are (\vec{x}, y)
where $\vec{x} \in \mathbb{R}^P$ $y \in \{1, \dots, k\}$

Model: parameters $\Theta \in \mathbb{R}^M$
feature function $f(\vec{x}, y) \in \mathbb{R}^M$

Predict: $\hat{y} = h_\Theta(\vec{x}) = \underset{y \in \{1, \dots, k\}}{\operatorname{argmax}} \Theta^T f(\vec{x}, y)$

Ex 4c: Suppose $\{1, \dots, k\}$ exist in some hierarchy and $\vec{x} \in \{0, 1\}^P$



Structured Perceptron

Whiteboard

- Multiclass Perceptron
- Structured Perceptron

Structured Perceptron

Mistake Bound:

Definition 1 Let $\overline{\text{GEN}}(x_i) = \text{GEN}(x_i) - \{y_i\}$. In other words $\overline{\text{GEN}}(x_i)$ is the set of incorrect candidates for an example x_i . We will say that a training sequence (x_i, y_i) for $i = 1 \dots n$ is **separable with margin $\delta > 0$** if there exists some vector \mathbf{U} with $\|\mathbf{U}\| = 1$ such that

$$\forall i, \forall z \in \overline{\text{GEN}}(x_i), \quad \mathbf{U} \cdot \Phi(x_i, y_i) - \mathbf{U} \cdot \Phi(x_i, z) \geq \delta \quad (3)$$

($\|\mathbf{U}\|$ is the 2-norm of \mathbf{U} , i.e., $\|\mathbf{U}\| = \sqrt{\sum_s \mathbf{U}_s^2}$.)

Theorem 1 For any training sequence (x_i, y_i) which is separable with margin δ , then for the perceptron algorithm in figure 2

$$\text{Number of mistakes} \leq \frac{R^2}{\delta^2}$$

where R is a constant such that $\forall i, \forall z \in \overline{\text{GEN}}(x_i) \quad \|\Phi(x_i, y_i) - \Phi(x_i, z)\| \leq R$.

Structured Perceptron

- Results from Collins (2002) on two **sequence tagging** problems
- Metrics:
 - **F-measure:** higher is better
 - **Error:** lower is better
- Comparison of...
 - Structured Perceptron **with** and **without** averaging
 - Maximum entropy Markov model (**MEMM**)
- Takeaways:
 - incredibly **easy to implement**
 - typically **blazing fast**

Table from Collins (2002)

NP Chunking Results

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

POS Tagging Results

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

Figure 4: Results for various methods on the part-of-speech tagging and chunking tasks on development data. All scores are error percentages. Numits is the number of training iterations at which the best score is achieved. Perc is the perceptron algorithm, ME is the maximum entropy method. Avg/noavg is the perceptron with or without averaged parameter vectors. cc=5 means only features occurring 5 times or more in training are included, cc=0 means all features in training are included.

aka. Max-Margin Markov Networks (M^3Ns)

STRUCTURED SVM

Support Vector Machines

Binary SVM

Data: $D = \{(x^{(i)}, y^{(i)})\}_{i=1}^N$ where $\vec{x} \in \mathbb{R}^M$ $y^{(i)} \in \{+1, -1\}$

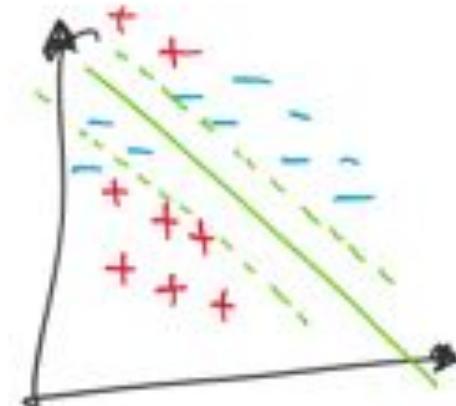
Model: $\hat{y} = h_{w,b}(x) = \text{sign}(w^T x^{(i)} + b)$

Quadratic Program (QP):

$$\min_{w,e} \frac{1}{2} (\|w\|_2)^2 + C \sum_{i=1}^N e_i$$

$$\text{s.t. } y^{(i)}(w^T x^{(i)} + b) \geq 1 - e_i \quad \forall i$$

$$e_i \geq 0 \quad \forall i$$



Binary SVM: Hinge Loss

Hinge Loss: $l^{\text{hinge}}(y, \hat{y}) = \max(0, 1 - y\hat{y})$

Unconstrained Opt. Problem:

Observe: $e_i \geq 1 - y^{(i)}(\underbrace{w^T x^{(i)} + b}_{\hat{y}})$] $\Rightarrow e_i \geq \max(0, 1 - y^{(i)}(w^T x^{(i)} + b))$
 and $e_i \geq 0$

$$\min_w \underbrace{\frac{1}{2} (\|w\|_2)^2}_{\text{large margin}} + C \underbrace{\sum_{i=1}^N \max(0, 1 - y^{(i)}(w^T x^{(i)} + b))}_{\text{few/small errors}}$$

Structured SVM

Whiteboard

- Structured Large Margin
- Structured Hinge Loss
- Gradient of Structured Hinge Loss
- SGD for Structured SVM
- Loss Augmented MAP Inference