



Transductive Support Vector Machines for Structured Variables

Alexander Zien



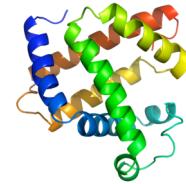
Ulf Brefeld



Tobias Scheffer



Label Sequence Learning



- Protein secondary structure prediction:

$x = \text{"XSITKTELGD ILPLVARGKV..."} \rightarrow y = \text{" ss TT SS EEEE SS... "}$

- Named entity recognition (NER):

$x = \text{"Tom comes from London.} \rightarrow y = \text{"Person, -, -, Location"}$

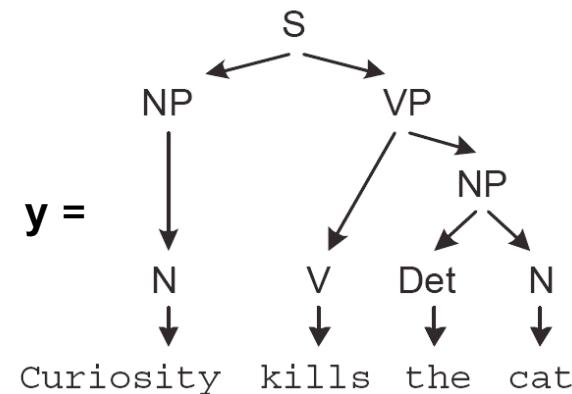
$x = \text{"S of PTH and CT..."} \rightarrow y = \text{"-, -, Gene, -, Gene, ..."}$

- Part-of-speech (POS) tagging:

$x = \text{"Curiosity kills the cat.} \rightarrow y = \text{"noun, verb, det, noun"}$

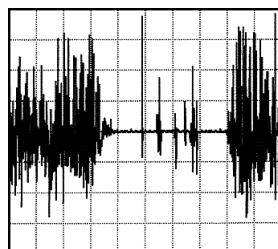
Natural Language Parsing

$x = \text{"Curiosity kills the cat"}$ → $y =$



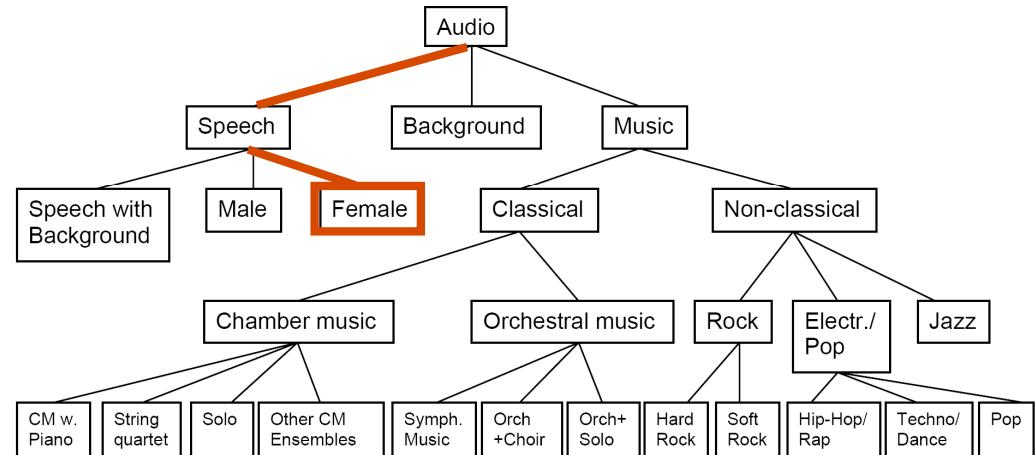
Classification with Taxonomies

$x =$



→

$y =$



Semi-supervised Discriminative Learning

- Labeled training data is scarce and expensive.
 - ◆ E.g., experiments in computational biology.
 - ◆ Need for expert knowledge.
 - ◆ Tedious and time consuming.
- Unclassified instances are abundant and cheap.
 - ◆ Extract sentences from www (POS-tagging, NER).
 - ◆ Assess primary structure of proteins from DNA/RNA.
 - ◆ ...

**There is a need for semi-supervised
techniques in structural learning!**

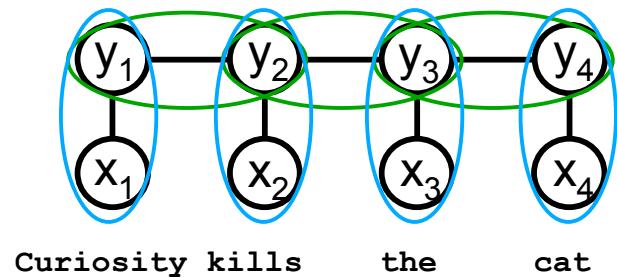
Overview

- Recap: SVM with structured output variables.
 - ◆ Convex optimization problem.
- Transductive SVM with structured variables.
 - ◆ Uses information from unlabeled data.
 - ◆ Combinatorial, non-convex optimization problem.
- Efficient optimization:
 - ◆ Transform, remove discrete parameters.
 - ◆ Differentiable, continuous optimization problem.
- Empirical results.

Structural Learning

- Given:
 - ◆ Labeled pairs $(\mathbf{x}_1, \mathbf{y}_1), \dots, (\mathbf{x}_n, \mathbf{y}_n)$, drawn from $P_{\mathcal{X}\mathcal{Y}}$.
- Learn decision function: $f : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$ with low risk.
- Learn generalized linear model:
 - ◆ $f(\mathbf{x}, \mathbf{y}) = \langle \mathbf{w}, \Phi(\mathbf{x}, \mathbf{y}) \rangle$
 - ◆ with low regularized empirical risk:
 - ◆ $\hat{R}(f) = \sum_{i=1}^n \ell(\mathbf{x}_i, \mathbf{y}_i, f) + \|f\|^2$
- Inference, decoding:
 - ◆ Compute output $\hat{\mathbf{y}} = \operatorname{argmax}_{\mathbf{y} \in \mathcal{Y}} f(\mathbf{x}, \mathbf{y})$

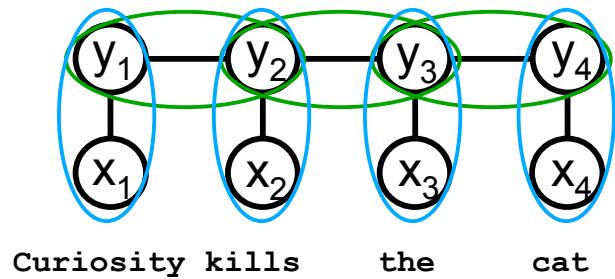
Label Sequence Learning



- Attribute for each pair of values of adjacent labels y_t und y_{t+1} .
 - ◆ $\varphi_{123}(y_t, y_{t+1}) = [[y_t = \text{"Noun"} \wedge y_{t+1} = \text{"Verb"}]]$
- Attribute for each pair of values for input x_t and output y_t .

$$\bar{\varphi}_{234}(x_t, y_t) = [[y_t = \text{"Noun"} \wedge x_t = \text{"cat"}]]$$

Label Sequence Learning



- Attribute for each pair of values of adjacent labels y_t und y_{t+1} .
 - ◆ $\varphi_{123}(y_t, y_{t+1}) = [[y_t = \text{Noun} \wedge y_{t+1} = \text{Verb}]]$
- Attribute for each pair of values for input x_t and output y_t .

$$\bar{\varphi}_{234}(x_t, y_t) = [[y_t = \text{Noun} \wedge x_t = \text{cat}]]$$
- Label-label: $\sum_t \varphi_i(y_t, y_{t+1})$.
- Label-observation $\sum_t \bar{\varphi}_i(x_t, y_t)$.
- Joint feature representation $\Phi(x, y) = \sum_t (\dots, \varphi_{123}(y_t, y_{t+1}), \dots, \bar{\varphi}_{234}(x_t, y_t), \dots)^T$
- Weight vector $w = (\dots, w_{123}, \dots, w_{234}, \dots)^T$
- Use Viterbi decoder to find argmax.

Structural Learning

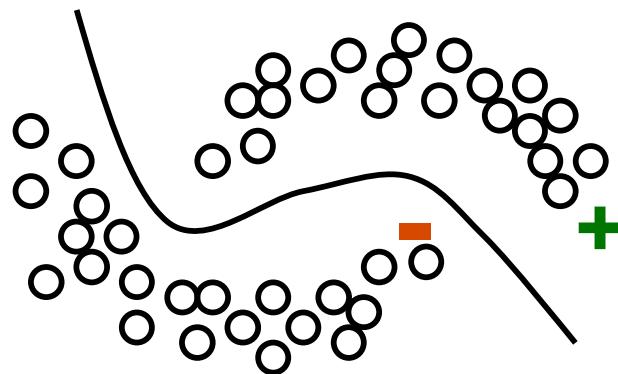
- Given:
 - ◆ Labeled pairs $(\mathbf{x}_1, \mathbf{y}_1), \dots, (\mathbf{x}_n, \mathbf{y}_n)$, drawn from $P_{\mathcal{X}\mathcal{Y}}$.
- Learn decision function: $f : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$ with low risk.
- Learn generalized linear model:
 - ◆ $f(\mathbf{x}, \mathbf{y}) = \langle \mathbf{w}, \Phi(\mathbf{x}, \mathbf{y}) \rangle$
 - ◆ with low regularized empirical risk:
 - ◆ $\hat{R}(f) = \sum_{i=1}^n \ell(\mathbf{x}_i, \mathbf{y}_i, f) + \|f\|^2$
- Inference, decoding:
 - ◆ Compute output $\hat{\mathbf{y}} = \operatorname{argmax}_{\mathbf{y} \in \mathcal{Y}} f(\mathbf{x}, \mathbf{y})$

Structural Learning

- Given:
 - ◆ Labeled pairs $(\mathbf{x}_1, \mathbf{y}_1), \dots, (\mathbf{x}_n, \mathbf{y}_n)$, drawn from $P_{\mathcal{X}\mathcal{Y}}$.
 - ◆ Unlabeled instances $\mathbf{x}_{n+1}, \dots, \mathbf{x}_{n+m}$.
- Learn decision function: $f : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$ with low risk.
- Learn generalized linear model:
 - ◆ $f(\mathbf{x}, \mathbf{y}) = \langle \mathbf{w}, \Phi(\mathbf{x}, \mathbf{y}) \rangle$
 - ◆ with low regularized empirical risk:
 - ◆ $\hat{R}(f) = \sum_{i=1}^n \ell(\mathbf{x}_i, \mathbf{y}_i, f) + \|f\|^2 + r(\mathbf{x}_{n+1}, \dots, \mathbf{x}_{n+m})$
 - ◆ Add data-dependent regularizer.
 - ◆ Dependent on unlabeled pairs, linked to low risk.

Data-Dependent Regularizer

- Separate unlabeled instances by large margin.
- Decision boundary not to cross high-density regions.



- Add slack terms for unlabeled instances to optimization problem.
- Treat class labels as additional discrete parameters.

Inductive Optimization Criterion

- Over all \mathbf{w} :

$$\begin{aligned} \min_{\mathbf{w}, \xi_k} \quad & \frac{1}{2} \mathbf{w}^\top \mathbf{w} + C \sum_i \xi_i \\ \text{s.t.} \quad & \forall \bar{\mathbf{y}}_i \neq \mathbf{y}_i : \mathbf{w}^\top [\Phi(\mathbf{x}_i, \mathbf{y}_i) - \Phi(\mathbf{x}_i, \bar{\mathbf{y}}_i)] \geq 1 - \xi_i, \quad \xi_i \geq 0 \end{aligned}$$

- Labeled examples $(\mathbf{x}_i, \mathbf{y}_i)$.

Inductive Optimization Criterion

- Over all \mathbf{w} :

$$\begin{aligned} \min_{\mathbf{w}, \xi_k} \quad & \frac{1}{2} \mathbf{w}^\top \mathbf{w} + C \sum_i \xi_i \\ \text{s.t.} \quad & \forall \bar{\mathbf{y}}_i \neq \mathbf{y}_i : \mathbf{w}^\top [\Phi(\mathbf{x}_i, \mathbf{y}_i) - \Phi(\mathbf{x}_i, \bar{\mathbf{y}}_i)] \geq 1 - \xi_i, \quad \xi_i \geq 0 \end{aligned}$$

- Iterative working-set algorithm.
 - ◆ For each \mathbf{x}_i , find highest-scoring offending $\bar{\mathbf{y}} \neq \mathbf{y}_i$.
 - ◆ If margin constraint is violated, add to working set.
 - ◆ Repeat until no more margin violators.

Inductive Optimization Criterion

- Over all \mathbf{w} :

$$\begin{aligned} \min_{\mathbf{w}, \xi_k} \quad & \frac{1}{2} \mathbf{w}^\top \mathbf{w} + C \sum_i \xi_i \\ \text{s.t.} \quad & \forall \bar{\mathbf{y}}_i \neq \mathbf{y}_i : \mathbf{w}^\top [\Phi(\mathbf{x}_i, \mathbf{y}_i) - \Phi(\mathbf{x}_i, \bar{\mathbf{y}}_i)] \geq 1 - \xi_i, \quad \xi_i \geq 0 \end{aligned}$$

- Labeled examples $(\mathbf{x}_i, \mathbf{y}_i)$.

Transductive Optimization Criterion

- Over all \mathbf{w} , \mathbf{y}_j :

$$\begin{aligned} \min_{\mathbf{w}, \mathbf{y}_j, \xi_k} \quad & \frac{1}{2} \mathbf{w}^\top \mathbf{w} + C \sum_i \xi_i + C^* \sum_j \xi_j \\ \text{s.t.} \quad & \forall \bar{\mathbf{y}}_i \neq \mathbf{y}_i : \mathbf{w}^\top [\Phi(\mathbf{x}_i, \mathbf{y}_i) - \Phi(\mathbf{x}_i, \bar{\mathbf{y}}_i)] \geq 1 - \xi_i, \quad \xi_i \geq 0 \\ & \forall \bar{\mathbf{y}}_j \neq \mathbf{y}_j : \mathbf{w}^\top [\Phi(\mathbf{x}_j, \mathbf{y}_j) - \Phi(\mathbf{x}_j, \bar{\mathbf{y}}_j)] \geq 1 - \xi_j, \quad \xi_j \geq 0 \end{aligned}$$

- Labeled examples $(\mathbf{x}_i, \mathbf{y}_i)$.
- Unlabeled examples \mathbf{x}_j .

Transductive Optimization Criterion

- Over all \mathbf{w} , \mathbf{y}_j :

$$\begin{aligned}
 \min_{\mathbf{w}, \mathbf{y}_j, \xi_k} \quad & \frac{1}{2} \mathbf{w}^\top \mathbf{w} + C \sum_i \xi_i + C^* \sum_j \xi_j \\
 \text{s.t.} \quad & \forall \bar{\mathbf{y}}_i \neq \mathbf{y}_i : \mathbf{w}^\top [\Phi(\mathbf{x}_i, \mathbf{y}_i) - \Phi(\mathbf{x}_i, \bar{\mathbf{y}}_i)] \geq 1 - \xi_i, \quad \xi_i \geq 0 \\
 & \forall \bar{\mathbf{y}}_j \neq \mathbf{y}_j : \mathbf{w}^\top [\Phi(\mathbf{x}_j, \mathbf{y}_j) - \Phi(\mathbf{x}_j, \bar{\mathbf{y}}_j)] \geq 1 - \xi_j, \quad \xi_j \geq 0
 \end{aligned}$$

- Non-convex problem.
- Variables \mathbf{y}_j are discrete.
- Integer program; NP-complete!

Transductive Optimization Criterion

- Over all \mathbf{w} , \mathbf{y}_j :

$$\begin{aligned} \min_{\mathbf{w}, \mathbf{y}_j, \xi_k} \quad & \frac{1}{2} \mathbf{w}^\top \mathbf{w} + C \sum_i \xi_i + C^* \sum_j \xi_j \\ \text{s.t.} \quad & \forall \bar{\mathbf{y}}_i \neq \mathbf{y}_i : \mathbf{w}^\top [\Phi(\mathbf{x}_i, \mathbf{y}_i) - \Phi(\mathbf{x}_i, \bar{\mathbf{y}}_i)] \geq 1 - \xi_i, \quad \xi_i \geq 0 \\ & \forall \bar{\mathbf{y}}_j \neq \mathbf{y}_j : \mathbf{w}^\top [\Phi(\mathbf{x}_j, \mathbf{y}_j) - \Phi(\mathbf{x}_j, \bar{\mathbf{y}}_j)] \geq 1 - \xi_j, \quad \xi_j \geq 0 \end{aligned}$$

- Non-convex problem.
- Variables \mathbf{y}_j are discrete
- Integer program; NP-complete!
- → Unconstrained, differentiable optimization problem.

Continuous Optimization

Solve for slack terms, closed expression of loss.

$$\begin{aligned} \min_{\mathbf{w}, \mathbf{y}_j, \xi_k} \quad & \frac{1}{2} \mathbf{w}^\top \mathbf{w} + C \sum_i \xi_i + C^* \sum_j \xi_j \\ \text{s.t.} \quad & \begin{cases} \forall \bar{\mathbf{y}}_i \neq \mathbf{y}_i : \mathbf{w}^\top [\Phi(\mathbf{x}_i, \mathbf{y}_i) - \Phi(\mathbf{x}_i, \bar{\mathbf{y}}_i)] \geq 1 - \xi_i, & \xi_i \geq 0 \\ \forall \bar{\mathbf{y}}_j \neq \mathbf{y}_j : \mathbf{w}^\top [\Phi(\mathbf{x}_j, \mathbf{y}_j) - \Phi(\mathbf{x}_j, \bar{\mathbf{y}}_j)] \geq 1 - \xi_j, & \xi_j \geq 0 \end{cases} \\ & \Leftrightarrow \\ & \xi_i = \max_{\bar{\mathbf{y}}_i \neq \mathbf{y}_i} \max \left\{ 1 - \mathbf{w}^\top [\Phi(\mathbf{x}_i, \mathbf{y}_i) - \Phi(\mathbf{x}_i, \bar{\mathbf{y}}_i)], 0 \right\} \\ & \xi_j = \min_{\mathbf{y}_j} \max_{\bar{\mathbf{y}}_j \neq \mathbf{y}_j} \max \left\{ 1 - \mathbf{w}^\top [\Phi(\mathbf{x}_j, \mathbf{y}_j) - \Phi(\mathbf{x}_j, \bar{\mathbf{y}}_j)], 0 \right\} \end{aligned}$$

Continuous Optimization

- Solve for slack terms, closed expression of loss.

$$\begin{aligned} \min_{\mathbf{w}, \mathbf{y}_j, \xi_k} \quad & \frac{1}{2} \mathbf{w}^\top \mathbf{w} + C \sum_i \xi_i + C^* \sum_j \xi_j \\ \text{s.t.} \quad & \begin{cases} \forall \bar{\mathbf{y}}_i \neq \mathbf{y}_i : \mathbf{w}^\top [\Phi(\mathbf{x}_i, \mathbf{y}_i) - \Phi(\mathbf{x}_i, \bar{\mathbf{y}}_i)] \geq 1 - \xi_i, & \xi_i \geq 0 \\ \forall \bar{\mathbf{y}}_j \neq \mathbf{y}_j : \mathbf{w}^\top [\Phi(\mathbf{x}_j, \mathbf{y}_j) - \Phi(\mathbf{x}_j, \bar{\mathbf{y}}_j)] \geq 1 - \xi_j, & \xi_j \geq 0 \end{cases} \\ & \Leftrightarrow \\ & \xi_i = \max_{\bar{\mathbf{y}}_i \neq \mathbf{y}_i} \max \left\{ 1 - \mathbf{w}^\top [\Phi(\mathbf{x}_i, \mathbf{y}_i) - \Phi(\mathbf{x}_i, \bar{\mathbf{y}}_i)], 0 \right\} \\ & \xi_j = \min_{\mathbf{y}_j} \max_{\bar{\mathbf{y}}_j \neq \mathbf{y}_j} \max \left\{ 1 - \mathbf{w}^\top [\Phi(\mathbf{x}_j, \mathbf{y}_j) - \Phi(\mathbf{x}_j, \bar{\mathbf{y}}_j)], 0 \right\} \end{aligned}$$

Now written as closed, unconstrained optimization problem.

Continuous Optimization

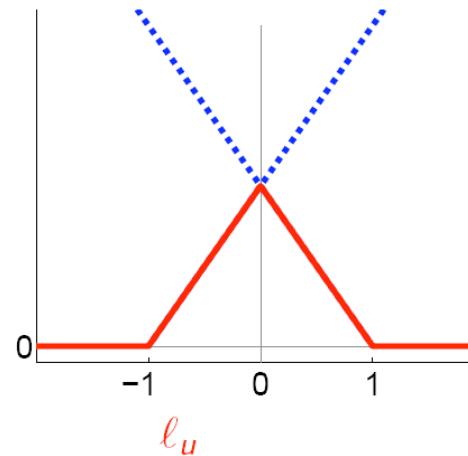
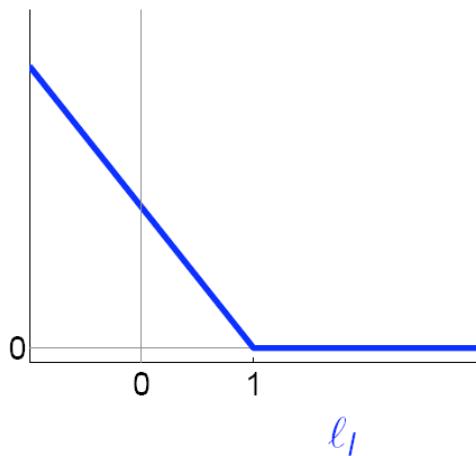
- Make differentiable.

$$\min_{\mathbf{w}, \mathbf{y}_j, \xi_k} \quad \frac{1}{2} \mathbf{w}^\top \mathbf{w} + C \sum_i \xi_i + C^* \sum_j \xi_j$$

s.t.

$$\xi_i = \max_{\bar{\mathbf{y}}_i \neq \mathbf{y}_i} \max \left\{ 1 - \mathbf{w}^\top [\Phi(\mathbf{x}_i, \mathbf{y}_i) - \Phi(\mathbf{x}_i, \bar{\mathbf{y}}_i)], 0 \right\}$$

$$\xi_j = \min_{\mathbf{y}_j} \max_{\bar{\mathbf{y}}_j \neq \mathbf{y}_j} \max \left\{ 1 - \mathbf{w}^\top [\Phi(\mathbf{x}_j, \mathbf{y}_j) - \Phi(\mathbf{x}_j, \bar{\mathbf{y}}_j)], 0 \right\}$$



Continuous Optimization

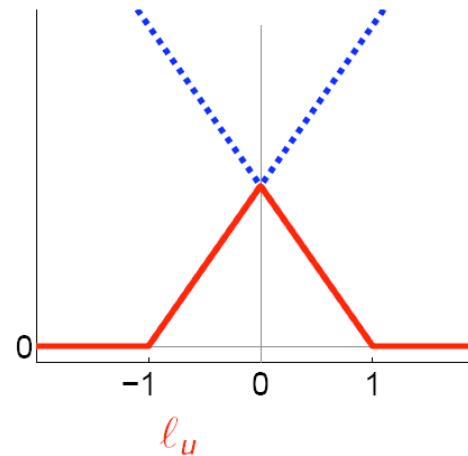
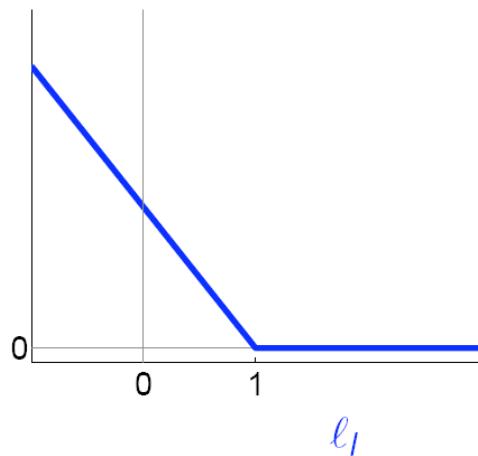
- Make differentiable.

$$\min_{\mathbf{w}, \mathbf{y}_j, \xi_k} \quad \frac{1}{2} \mathbf{w}^\top \mathbf{w} + C \sum_i \xi_i + C^* \sum_j \xi_j$$

s.t.

$$\xi_i = \max_{\bar{\mathbf{y}}_i \neq \mathbf{y}_i} \ell_I \left(\mathbf{w}^\top \Phi(\mathbf{x}_i, \mathbf{y}_i) - \mathbf{w}^\top \Phi(\mathbf{x}_i, \bar{\mathbf{y}}_i) \right)$$

$$\xi_j = \min_{\mathbf{y}_j} \max_{\bar{\mathbf{y}}_j \neq \mathbf{y}_j} \ell_u \left(\mathbf{w}^\top \Phi(\mathbf{x}_j, \mathbf{y}_j) - \mathbf{w}^\top \Phi(\mathbf{x}_j, \bar{\mathbf{y}}_j) \right)$$



Continuous Optimization

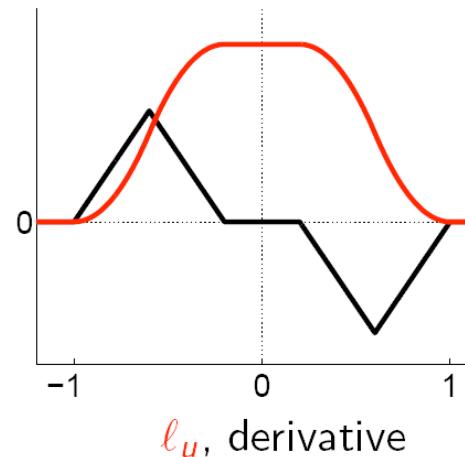
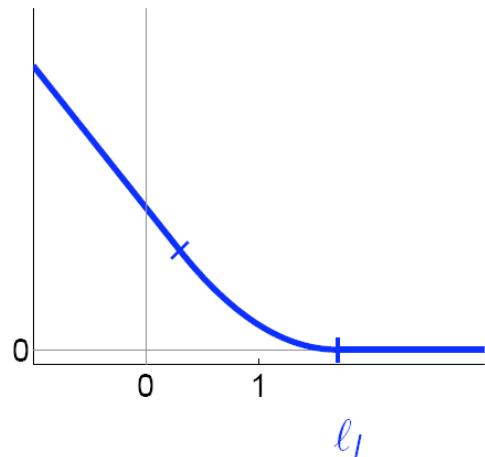
- Make differentiable.

$$\min_{\mathbf{w}, \mathbf{y}_j, \xi_k} \quad \frac{1}{2} \mathbf{w}^\top \mathbf{w} + C \sum_i \xi_i + C^* \sum_j \xi_j$$

s.t.

$$\xi_i = \max_{\bar{\mathbf{y}}_i \neq \mathbf{y}_i} \ell_I \left(\mathbf{w}^\top \Phi(\mathbf{x}_i, \mathbf{y}_i) - \mathbf{w}^\top \Phi(\mathbf{x}_i, \bar{\mathbf{y}}_i) \right)$$

$$\xi_j = \min_{\mathbf{y}_j} \max_{\bar{\mathbf{y}}_j \neq \mathbf{y}_j} \ell_u \left(\mathbf{w}^\top \Phi(\mathbf{x}_j, \mathbf{y}_j) - \mathbf{w}^\top \Phi(\mathbf{x}_j, \bar{\mathbf{y}}_j) \right)$$



Continuous Optimization

- Make differentiable.

$$\min_{\mathbf{w}, \mathbf{y}_j, \xi_k} \quad \frac{1}{2} \mathbf{w}^\top \mathbf{w} + C \sum_i \xi_i + C^* \sum_j \xi_j$$

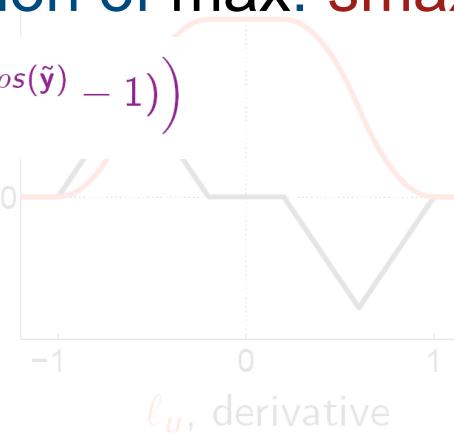
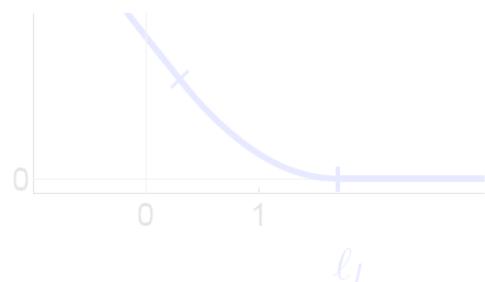
s.t.

$$\xi_i = \underset{\bar{\mathbf{y}}_i \neq \mathbf{y}_i}{\text{smax}} \ell_I \left(\mathbf{w}^\top \Phi(\mathbf{x}_i, \mathbf{y}_i) - \mathbf{w}^\top \Phi(\mathbf{x}_i, \bar{\mathbf{y}}_i) \right)$$

$$\xi_j = \min_{\mathbf{y}_j} \underset{\bar{\mathbf{y}}_j \neq \mathbf{y}_j}{\text{smax}} \ell_U \left(\mathbf{w}^\top \Phi(\mathbf{x}_j, \mathbf{y}_j) - \mathbf{w}^\top \Phi(\mathbf{x}_j, \bar{\mathbf{y}}_j) \right)$$

- Differentiable approximation of max: **smax**.

$$\text{smax}(s(\tilde{\mathbf{y}})) = \frac{1}{\rho} \log \left(1 + \sum_{\tilde{\mathbf{y}} \neq \mathbf{y}_k} (e^{\rho s(\tilde{\mathbf{y}})} - 1) \right)$$



Continuous Optimization

- Unconstraint, differentiable optimization criterion.

$$\begin{aligned} \min_{\mathbf{w}, \xi_k} \quad & \frac{1}{2} \mathbf{w}^\top \mathbf{w} \\ & + C \sum_i \underset{\bar{\mathbf{y}}_i \neq \mathbf{y}_i}{\text{smax}} \ell_l \left(\mathbf{w}^\top \Phi(\mathbf{x}_i, \mathbf{y}_i) - \mathbf{w}^\top \Phi(\mathbf{x}_i, \bar{\mathbf{y}}_i) \right) \\ & + C^* \sum_j \underset{\mathbf{y}_j}{\min} \underset{\bar{\mathbf{y}}_j \neq \mathbf{y}_j}{\text{smax}} \ell_u \left(\mathbf{w}^\top \Phi(\mathbf{x}_j, \mathbf{y}_j) - \mathbf{w}^\top \Phi(\mathbf{x}_j, \bar{\mathbf{y}}_j) \right) \end{aligned}$$

- Iterative working-set algorithm for labeled and unlabeled data, until no more margin violators.

Continuous Optimization

- Unconstraint, differentiable optimization criterion.

$$\begin{aligned} \min_{\mathbf{w}, \xi_k} \quad & \frac{1}{2} \mathbf{w}^\top \mathbf{w} \\ & + C \sum_i \operatorname{smax}_{\bar{\mathbf{y}}_i \neq \mathbf{y}_i} \ell_I \left(\mathbf{w}^\top \Phi(\mathbf{x}_i, \mathbf{y}_i) - \mathbf{w}^\top \Phi(\mathbf{x}_i, \bar{\mathbf{y}}_i) \right) \\ & + C^* \sum_j \min_{\mathbf{y}_j} \operatorname{smax}_{\bar{\mathbf{y}}_j \neq \mathbf{y}_j} \ell_U \left(\mathbf{w}^\top \Phi(\mathbf{x}_j, \mathbf{y}_j) - \mathbf{w}^\top \Phi(\mathbf{x}_j, \bar{\mathbf{y}}_j) \right) \end{aligned}$$

- Gradient descent over primal parameters:

- ◆ Follow $\frac{\partial obj}{\partial \mathbf{w}}$.

Continuous Optimization

- Unconstraint, differentiable optimization criterion.

$$\begin{aligned} \min_{\mathbf{w}, \xi_k} \quad & \frac{1}{2} \mathbf{w}^\top \mathbf{w} \\ & + C \sum_i \underset{\bar{\mathbf{y}}_i \neq \mathbf{y}_i}{\text{smax}} \ell_I \left(\mathbf{w}^\top \Phi(\mathbf{x}_i, \mathbf{y}_i) - \mathbf{w}^\top \Phi(\mathbf{x}_i, \bar{\mathbf{y}}_i) \right) \\ & + C^* \sum_j \underset{\mathbf{y}_j}{\min} \underset{\bar{\mathbf{y}}_j \neq \mathbf{y}_j}{\text{smax}} \ell_U \left(\mathbf{w}^\top \Phi(\mathbf{x}_j, \mathbf{y}_j) - \mathbf{w}^\top \Phi(\mathbf{x}_j, \bar{\mathbf{y}}_j) \right) \end{aligned}$$

- Gradient descent over parameters.
- Invoke Representer Theorem:

$$\mathbf{w} = \sum_{k=1}^{n+m} \sum_{\mathbf{y} \in \mathcal{Y}(\mathbf{x}_k)} \alpha_{k,y} \Phi(\mathbf{x}_k, \mathbf{y})$$

Continuous Optimization

- Unconstraint, differentiable optimization criterion.

$$\begin{aligned} \min_{\mathbf{w}, \xi_k} \quad & \frac{1}{2} \mathbf{w}^\top \mathbf{w} \\ & + C \sum_i \underset{\bar{\mathbf{y}}_i \neq \mathbf{y}_i}{\text{smax}} \ell_I \left(\mathbf{w}^\top \Phi(\mathbf{x}_i, \mathbf{y}_i) - \mathbf{w}^\top \Phi(\mathbf{x}_i, \bar{\mathbf{y}}_i) \right) \\ & + C^* \sum_j \underset{\mathbf{y}_j}{\min} \underset{\bar{\mathbf{y}}_j \neq \mathbf{y}_j}{\text{smax}} \ell_U \left(\mathbf{w}^\top \Phi(\mathbf{x}_j, \mathbf{y}_j) - \mathbf{w}^\top \Phi(\mathbf{x}_j, \bar{\mathbf{y}}_j) \right) \end{aligned}$$

- Chain rule gives dual gradient.

◆ Calculate $\frac{\partial \text{obj}}{\partial \alpha_{k,y}} = \frac{\partial \text{obj}}{\partial \mathbf{w}} \cdot \frac{\partial \mathbf{w}}{\partial \alpha_{k,y}}$

◆ With

$$\mathbf{w}^\top \Phi(\mathbf{x}_i, \mathbf{y}_i) = \sum_k \sum_y \alpha_{k,y} \underbrace{\Phi(\mathbf{x}_k, \mathbf{y})^\top \Phi(\mathbf{x}_i, \mathbf{y}_i)}_{k((\mathbf{x}_k, \mathbf{y}), (\mathbf{x}_i, \mathbf{y}_i))}$$

Continuous Optimization

■ Continuous TSVM for structured variables.

Input: labeled points $\{(\mathbf{x}_i, \mathbf{y}_i)\}$, unlabeled points $\{\mathbf{x}_j\}$.

Output: working set \mathcal{W} and associated $\alpha_{k,y}$.

Initialize $\mathcal{W} \leftarrow \{(\mathbf{x}_i, \mathbf{y}_i)\}$.

Alternate until convergence:

① Augment working set \mathcal{W}

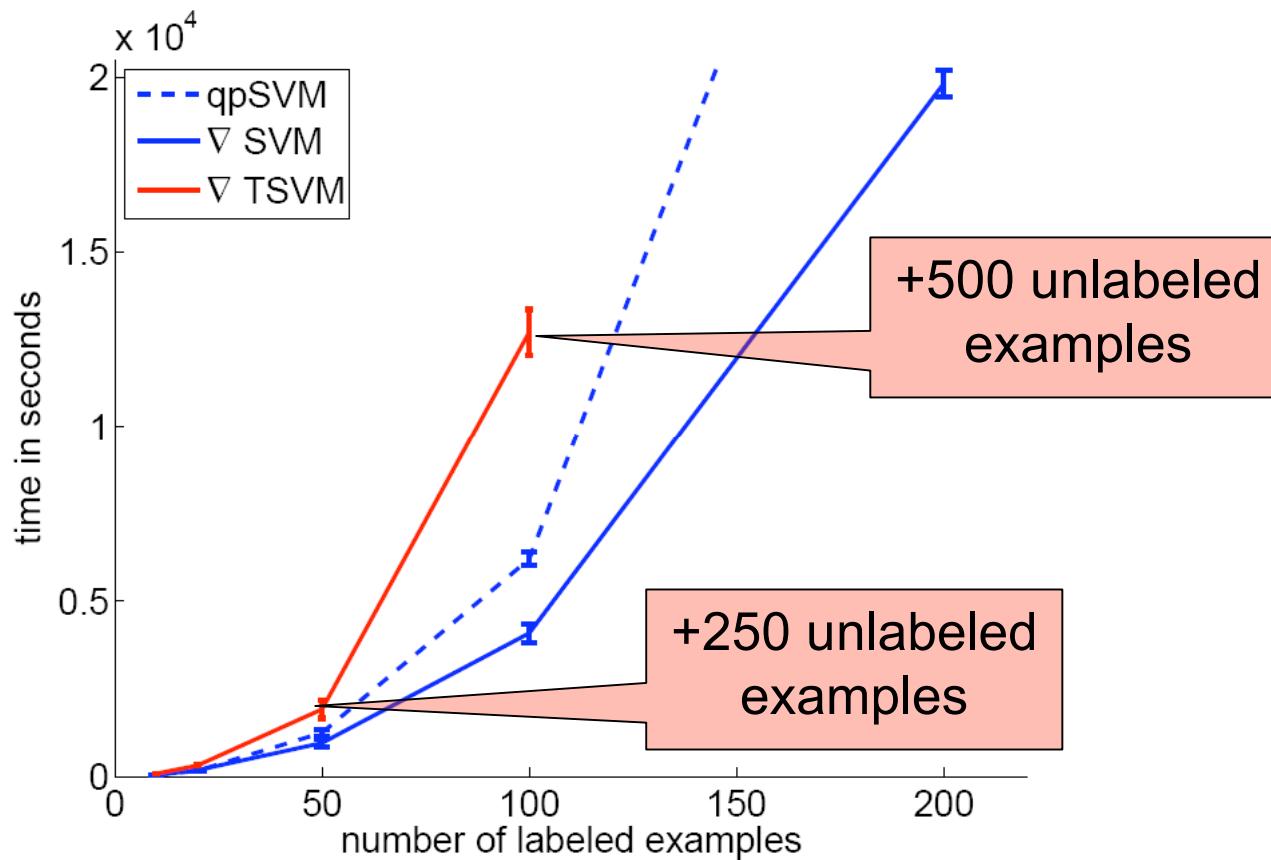
- add $\{(\mathbf{x}_i, \bar{\mathbf{y}}_i^*)\}$ to \mathcal{W} (worst margin violators)
- find $\{\mathbf{y}_j^*\}$ (highest scoring labels)
- add $\{(\mathbf{x}_j, \bar{\mathbf{y}}_j^*)\}$ to \mathcal{W} (2nd highest scoring labels)

② Optimize α by preconditioned Conjugate Gradient.

Experiments

- Execution time QP versus continuous optimization.
- Accuracy supervised vs. semi-supervised learning.
 - ◆ Multi-class classification: text classification.
 - ◆ Label-sequence learning: NER.
- Comparison / combination with Laplacian kernel.

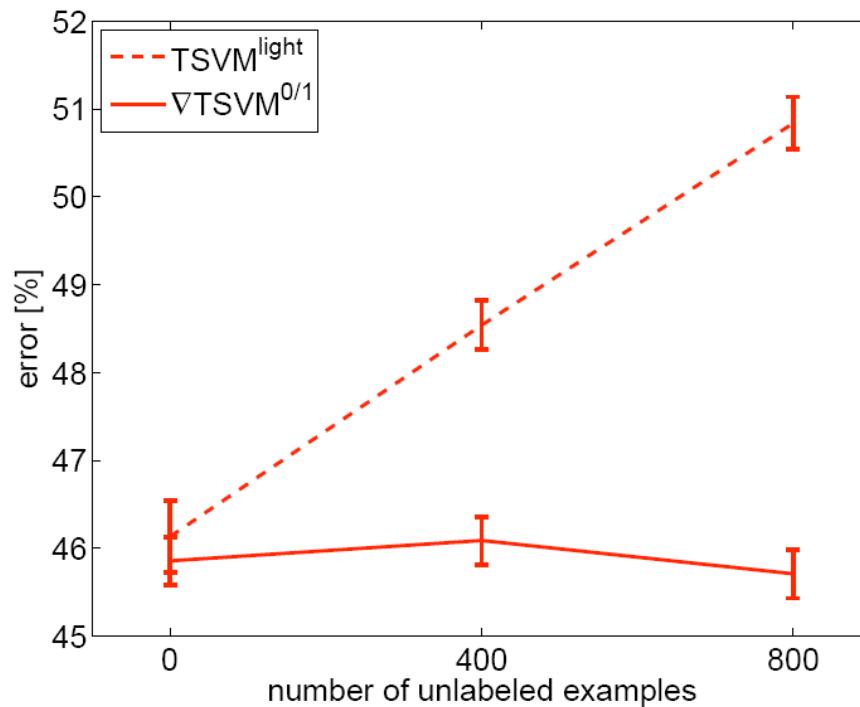
Execution Time



- Continuous optimization many times faster!

Accuracy: Cora Text Classification

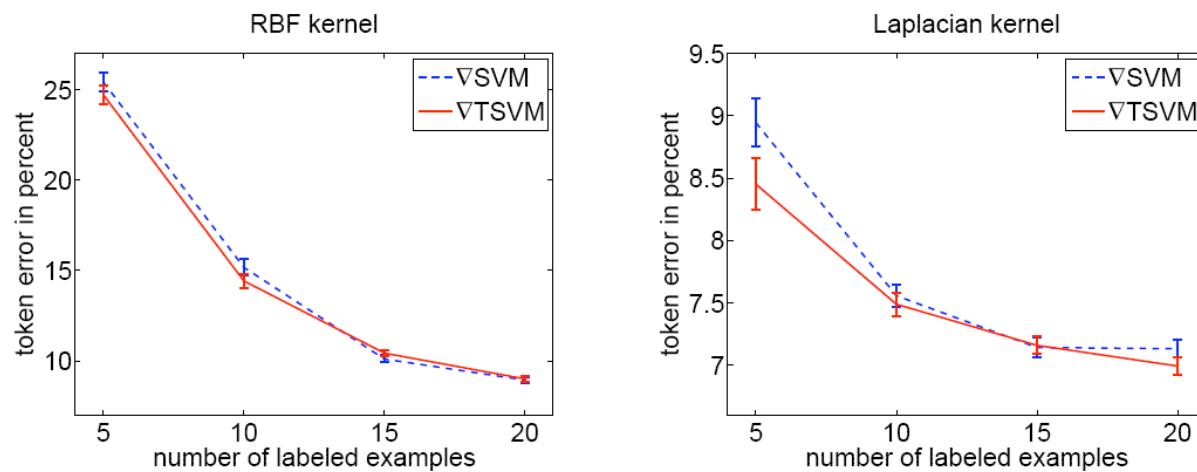
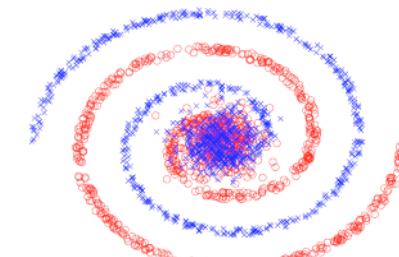
- Multiclass, 8 classes, 200 unlabeled examples.



- Combinatorial optimization: much higher error!
- Continuous optimization: higher or same error!

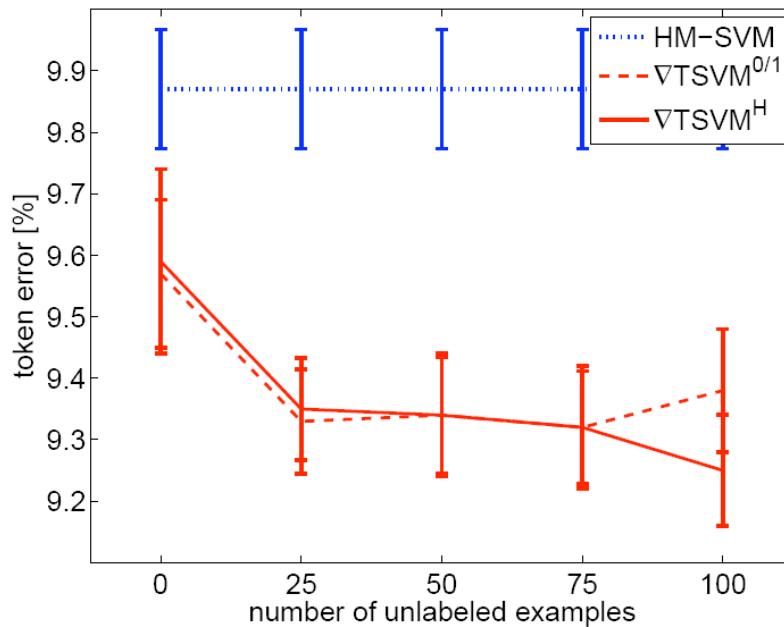
Accuracy: Galaxies Problem

- Artificial data set [Lafferty et al., IMCL 2004].
- Laplacian kernel: huge benefit.
- TSVM: marginal benefit.
- Laplacian+TSVM: slight additional benefit.



Accuracy: Spanish News Wire NER

- Sequence learning, 9 labels.

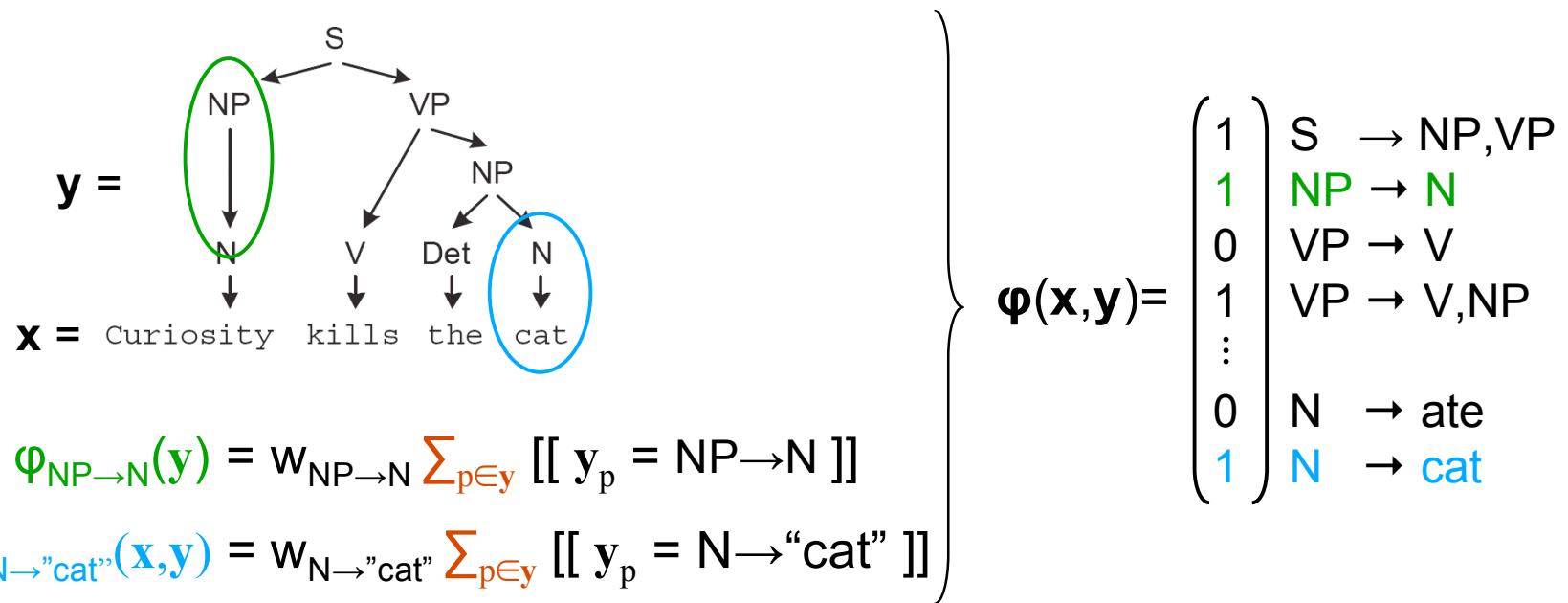


- Gradient TSVM clearly outperforms SVM.
- Adding unlabeled data further decreases error.

Conclusions

- TSVM for structured outputs.
 - ◆ Use information from unlabeled (test) examples.
 - ◆ Unconstrained, differentiable optimization criterion.
 - ◆ Efficient conjugate gradient descent algorithm.
- Empirically:
 - ◆ Sometimes no improvement.
 - ◆ Sometimes, unlabeled data increase accuracy significantly.
- Because:
 - ◆ SVM criterion is convex;
 - ◆ TSVM criterion has many local minima.

Natural Language Parsing



- Joint feature representation: $\Phi(x, y) = (\dots, \Phi_{NP \rightarrow N}(y_p), \dots, \Phi_{N \rightarrow "cat"}(x, y_p), \dots)^T$
- Weight vector $w = (\dots, w_{NP \rightarrow N}, \dots, w_{N \rightarrow "cat"}, \dots)^T$.
- Use chart parser as decoder to find argmax.