

Large Margin Methods for Structured and Interdependent Output Variables

I. Tsochantaris, T. Joachims, T. Hofmann, Y. Altun
Journal of Machine Learning Research (JMLR), 6(Sep):1453-1484, 2005.

Presented by
Thorsten Joachims

Cornell University
Department of Computer Science

Goal of Paper

- Learn function

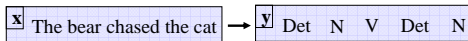
$$h : X \rightarrow Y$$

where X is the input space and Y is some structured output space (e.g. sequences, trees, equivalence relations).

- Paper proposes Support Vector Machine (SVM) method
 - that applies to a large class of structured outputs Y
 - Sequences (i.e. Hidden Markov Models)
 - Trees (i.e. Weighted Context-Free Grammars)
 - Hierarchical classification
 - Sequence alignment (i.e. Edit-Distance cost function)
 - allows the use of fairly general loss functions
 - is a generalization of multi-class SVMs
 - has polynomial time training algorithm.

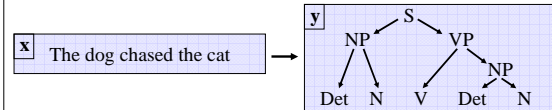
Examples of Complex Output Spaces

- Part-of-Speech Tagging
 - Given a sequence of words x , predict sequence of tags y .
 - Dependencies from tag-tag transitions in Markov model.



Examples of Complex Output Spaces

- Natural Language Parsing
 - Given a sequence of words x , predict the parse tree y .
 - Dependencies from structural constraints, since y has to be a tree.



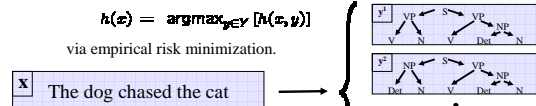
Structured Output Prediction as Multi-Class Classification

- Learning Task: $P(X, Y) = P(X) P(Y|X)$
 - Input Space: X (i.e. feature vectors, word sequence, etc.)
 - Output Space: Y (i.e. class, tag sequence, parse tree, etc.)
 - Training Data: $S = ((x_1, y_1), \dots, (x_n, y_n)) \sim_{iid} P(X, Y)$
- Approach: view as multi-class classification task
 - Every complex output $y \in Y$ is one class
- Goal: Find $h : X \rightarrow Y$ with low expected loss
 - Loss function: $\Delta(y, y')$ (penalty for predicting y' if y correct)
 - Expected loss (i.e. risk, prediction error):

$$Err_P(h) = \sum_{x, y} \Delta(y, h(x)) P(X=x, Y=y)$$

Natural Language Parsing as Multi-Class Classifications

- Approach: view as multi-class classification task
 - Every complex output $y \in Y$ is one class
 - Learn discriminant function

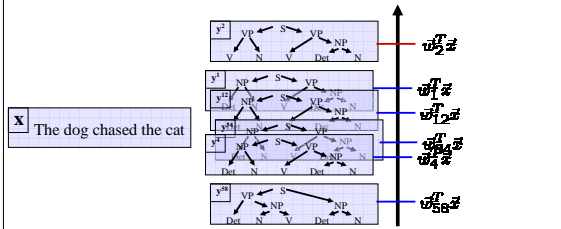


- Challenges: exponentially many classes
 - How to compactly represent model?
 - How to do efficient prediction with model (i.e. $\arg\max_{y \in Y} [h(x, y)]$)?
 - How to effectively estimate model from data? (e.g. compute $h_w = \arg\min_{h \in \mathcal{H}} \sum_{(x, y) \in S} \Delta(y, h_w(x))$)

Multi-Class Linear Discriminant

- Linear discriminant function of the form: $h_w(x, y) = w_y^T x$
- Learn one weight vector w_y for each class $y \in Y$

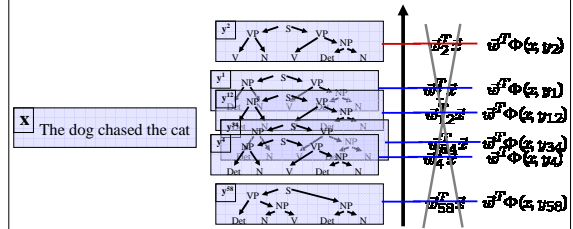
$$h(\vec{x}) = \operatorname{argmax}_{y \in Y} [w_y^T x]$$



Joint Feature/Kernel Map

- Feature vector $\Phi(x, y)$ that describes match between x and y
- Linear discriminant function of the form: $h_w(x, y) = w^T \Phi(x, y)$

$$h(\vec{x}) = \operatorname{argmax}_{y \in Y} [w^T \Phi(x, y)]$$

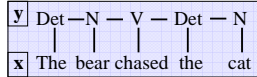


Joint Feature Map for Sequences

- Linear Chain Model (HMM)
 - Only local dependencies
 - Score for each adjacent label/label and word/label pair
 - Find highest scoring sequence

$$h(\vec{x}) = \operatorname{argmax}_{y \in Y} [w^T \Phi(x, y)]$$

Viterbi

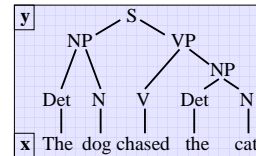


$$\Phi(x, y) = \begin{pmatrix} 1 & N \rightarrow V \\ 0 & Det \rightarrow V \\ 2 & Det \rightarrow N \\ 1 & V \rightarrow Det \\ \vdots & \\ 0 & Det \rightarrow bear \\ 2 & Det \rightarrow the \\ 1 & N \rightarrow bear \\ 1 & V \rightarrow chased \\ 1 & N \rightarrow cat \end{pmatrix}$$

Joint Feature Map for Trees

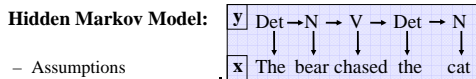
- Weighted Context Free Grammar
 - Each rule r_i (e.g. $S \rightarrow NP VP$) has a weight w_i
 - Score of a tree is the sum of its weights
 - Find highest scoring tree $h(\vec{x}) = \operatorname{argmax}_{y \in Y} [w^T \Phi(x, y)]$

CKY Parser



$$\Phi(x, y) = \begin{pmatrix} 1 & S \rightarrow NP VP \\ 0 & S \rightarrow NP \\ 2 & NP \rightarrow Det N \\ 1 & VP \rightarrow V NP \\ \vdots & \\ 0 & Det \rightarrow dog \\ 2 & Det \rightarrow the \\ 1 & N \rightarrow dog \\ 1 & V \rightarrow chased \\ 1 & N \rightarrow cat \end{pmatrix}$$

Connection to Graphical Models



Assumptions

$$P(y = (y^{(1)}, \dots, y^{(l)})) = \prod_{t=1}^l P(y_t = y^{(t)} | y_{t-1} = y^{(t-1)})$$

$$P(x = (x^{(1)}, \dots, x^{(l)}) | y = (y^{(1)}, \dots, y^{(l)})) = \prod_{t=1}^l P(x_t = x^{(t)} | y_t = y^{(t)})$$

$$\rightarrow \text{Rule: } h(\vec{x}) = \operatorname{argmax}_{y \in Y} [P(X=x | Y=y) P(Y=y)]$$

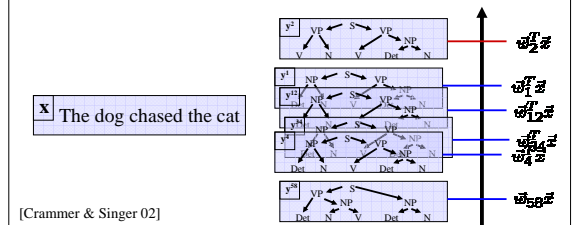
$$= \operatorname{argmax}_{(y^{(1)}, \dots, y^{(l)}) \in Y} \left[\prod_{t=1}^l P(y_t = y^{(t)} | y_{t-1} = y^{(t-1)}) P(x_t = x^{(t)} | y_t = y^{(t)}) \right]$$

$$= \operatorname{argmax}_{(y^{(1)}, \dots, y^{(l)}) \in Y} [w^T \Phi(x, y)]$$

with $w_{ab} = -\log[P(Y_c = a | Y_c = b)]$ and $w_{cd} = -\log[P(X_c = c | Y_c = d)]$ and $\Phi(x, y)$ histogram

Multi-Class SVM

- Training Examples: $(\vec{x}_1, y_1), \dots, (\vec{x}_n, y_n)$ $\vec{x} \in \mathbb{R}^N$ $y \in \{1, \dots, k\}$
- Hypothesis Space: $h(\vec{x}) = \operatorname{argmax}_{i \in \{1, \dots, k\}} [w_i^T \vec{x}]$



[Crammer & Singer 02]

Training: Find $\langle \bar{w}_1, \dots, \bar{w}_k \rangle$ that solve

$$\min_{\bar{w}_1, \dots, \bar{w}_k, \xi} \sum_{i=1}^k \bar{w}_i^T \bar{w}_i + C \sum_{i=1}^n \xi_i$$

s.t. $\forall j \neq y_1 : \bar{w}_{y_1}^T \bar{x}_1 \geq \bar{w}_j^T \bar{x}_1 + 1 - \xi_1$
 \dots
 $\forall j \neq y_n : \bar{w}_{y_n}^T \bar{x}_n \geq \bar{w}_j^T \bar{x}_n + 1 - \xi_n$

X The dog chased the cat

[Crammer & Singer 02]

Structural Support Vector Machine

- Joint features $\Phi(x, y)$ describe match between x and y
- Learn weights \bar{w} so that $\bar{w}^T \Phi(x, y)$ is max for correct y

Structural Support Vector Machine

Hard-margin optimization problem:

- $\min_{\bar{w}} \frac{1}{2} \bar{w}^T \bar{w}$
- s.t. $\forall y \in Y \setminus y_1 : \bar{w}^T \Phi(x_1, y_1) \geq \bar{w}^T \Phi(x_1, y) + 1$
- \dots
- $\forall y \in Y \setminus y_n : \bar{w}^T \Phi(x_n, y_n) \geq \bar{w}^T \Phi(x_n, y) + 1$

Soft-Margin Struct SVM (Margin Rescaling)

- Loss function $\Delta(y_i, y)$ measures match between target and prediction.

Soft-Margin Struct SVM (Margin Rescaling)

Soft-margin optimization problem:

- $\min_{\bar{w}, \xi} \frac{1}{2} \bar{w}^T \bar{w} + C \sum_{i=1}^n \xi_i$
- s.t. $\forall y \in Y \setminus y_1 : \bar{w}^T \Phi(x_1, y_1) \geq \bar{w}^T \Phi(x_1, y) + \Delta(y_1, y) - \xi_1$
- \dots
- $\forall y \in Y \setminus y_n : \bar{w}^T \Phi(x_n, y_n) \geq \bar{w}^T \Phi(x_n, y) + \Delta(y_n, y) - \xi_n$

Lemma: The training loss is upper bounded by

$$Err_S(h) = \frac{1}{n} \sum_{i=1}^n \Delta(y_i, h(\bar{x}_i)) \leq \frac{1}{n} \sum_{i=1}^n \xi_i$$

Soft-Margin Struct SVM (Slack Rescaling)

- Loss function $\Delta(y_i, y)$ measures match between target and prediction.

Soft-Margin Struct SVM (Slack Rescaling)

Soft-margin optimization problem:

$$\min_{\tilde{w}, \xi} \frac{1}{2} \tilde{w}^T \tilde{w} + C \sum_{i=1}^n \xi_i$$

s.t.

$$\forall y \in Y_{y_1} : \tilde{w}^T \Phi(x_1, y_1) \geq \tilde{w}^T \Phi(x_1, y) + 1 - \frac{\xi_1}{\Delta(y_1, y)}$$

$$\dots$$

$$\forall y \in Y_{y_n} : \tilde{w}^T \Phi(x_n, y_n) \geq \tilde{w}^T \Phi(x_n, y) + 1 - \frac{\xi_n}{\Delta(y_n, y)}$$

Lemma: The training loss is upper bounded by

$$Err_S(h) = \frac{1}{n} \sum_{i=1}^n \Delta(y_i, h(\tilde{x}_i)) \leq \frac{1}{n} \sum_{i=1}^n \xi_i$$

(x₁, y₁) (x₂, y₂) (x₃, y₃) (x_n, y_n)

Cutting-Plane Algorithm for Structural SVM

- Input: (x₁, y₁), ..., (x_n, y_n), C, ε
- S ← ∅, w̃ ← 0, ξ̃ ← 0
- REPEAT
 - FOR i = 1, ..., n
 - Find most violated constraint
 - Violated by more than ε?
 - compute $\hat{y} = \operatorname{argmax}_{y \in Y} \{ \Delta(y_i, y) + \tilde{w}^T \Phi(x_i, y) \}$
 - IF $(\Delta(y_i, \hat{y}) - \tilde{w}^T [\Phi(x_i, y_i) - \Phi(x_i, \hat{y})]) > \xi_i + \epsilon$
 - S ← S ∪ { $\tilde{w}^T [\Phi(x_i, y_i) - \Phi(x_i, \hat{y})] \geq \Delta(y_i, \hat{y}) - \xi_i$ }
 - [w̃, ξ̃] ← optimize StructSVM over S
 - ADD constraint to working set
 - ENDIF
- UNTIL S has not changed during iteration

[AltHo03] [Jo03] [TsoJoHoAlt05]

Polynomial Sparsity Bound

- Theorem:** The sparse-approximation algorithm finds a solution to the soft-margin optimization problem after adding at most

$$\max \left\{ \frac{2nA}{\epsilon}, \frac{8nCA R^2}{\epsilon^2} \right\}$$
 constraints to the working set S, so that the Kuhn-Tucker conditions are fulfilled up to a precision ε. The loss has to be bounded 0 ≤ Δ(y_i, ŷ) ≤ A, and ||Φ(x, ŷ)|| ≤ R.

[Jo 03] [Tsochantaridis et al. 04] [Tsochantaridis et al. 05]

Experiment: Natural Language Parsing

- Implementation**
 - Implemented Sparse-Approximation Algorithm in SVM^{light}
 - Incorporated modified version of Mark Johnson's CKY parser
 - Learned weighted CFG with ε = 0.01, C = 1
- Data**
 - Penn Treebank sentences of length at most 10 (start with POS)
 - Train on Sections 2-22: 4098 sentences
 - Test on Section 23: 163 sentences

Method	Test Accuracy		Training Efficiency		
	Acc	F ₁	CPU-h	Iter	Const
PCFG with MLE	55.2	86.0	0	N/A	N/A
SVM with (1-F ₁)-Loss	58.9	88.5	3.4	12	8043

[Tsochantaridis et al. 05]

More Expressive Features

- Linear composition:**

$$\Phi(x, y) = \sum_{i=1}^l \phi(x, y_i)$$
- General form:**

$$\phi(x, y_i) = \phi_{\text{kernel}}(\phi(x, [\text{rule}, \text{start}, \text{end}]))$$

$$K(a, b) = \phi_{\text{kernel}}(a)^T \phi_{\text{kernel}}(b)$$
- So far:**

$$\phi(x, y_i) = \begin{pmatrix} 0 \\ 1 \\ 0 \\ \dots \\ 0 \end{pmatrix} \text{ if rule}(y_i) = 'S \leftarrow NP VP'$$
- Example:**

$$\phi(x, y_i) = \begin{pmatrix} 1 & \text{if } y_i = NP \wedge \text{start} = ' \\ (\text{start} - \text{end})^2 & \text{if } y_i = NP \\ 1 & \text{if } y_i = S \wedge \text{span contains } x_c = \text{'and'}$$

see [Taskar et al. 05]

Applying Structural SVM to New Problem

- Application specific**
 - Loss function Δ(y_i, ŷ)
 - Representation Φ(x, y)
 - Algorithms to compute

$$\hat{y} = \operatorname{argmax}_{y \in Y} \{ \tilde{w}^T \Phi(x_i, y) \}$$

$$\hat{y} = \operatorname{argmax}_{y \in Y} \{ \Delta(y_i, y) + \tilde{w}^T \Phi(x_i, y) \}$$
- Implementation SVM-struct:** <http://svmlight.joachims.org>
 - Context-free grammars
 - Sequence alignment
 - Classification with multivariate loss (e.g. F1, ROC Area)
 - General API for other problems

Summary

- **Support Vector Machine approach to training**
 - Hidden Markov Models
 - Weighted Context-Free Grammars
 - Sequence Alignment cost functions
 - Etc.
- **Incorporate loss functions via**
 - Margin rescaling
 - Slack rescaling
- **General training algorithm based on cutting-plane method**
 - Efficient for all linear discriminant models where argmax efficient