

Learning to Predict Trees, Sequences, and other Structured Outputs

Tutorial at NESCAI 2006

Thorsten Joachims

Cornell University

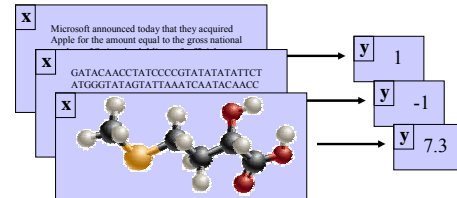
Department of Computer Science

Supervised Learning

- Find function from input space X to output space Y

$$h : X \rightarrow Y$$

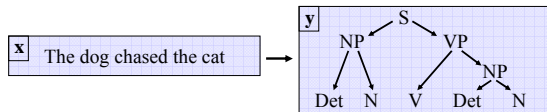
such that the prediction error is low.



Examples of Complex Output Spaces

Natural Language Parsing

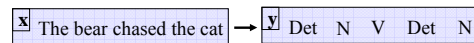
- Given a sequence of words x , predict the parse tree y .
- Dependencies from structural constraints, since y has to be a tree.



Examples of Complex Output Spaces

Part-of-Speech Tagging

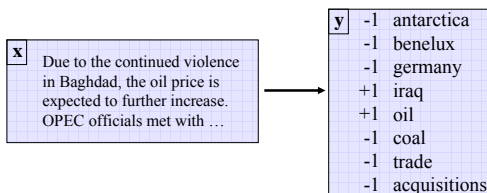
- Given a sequence of words x , predict sequence of tags y .
- Dependencies from tag-tag transitions in Markov model.



Examples of Complex Output Spaces

Multi-Label Classification

- Given a (bag-of-words) document x , predict a set of labels y .
- Dependencies between labels from correlations between labels ("iraq" and "oil" in newswire corpus)



Examples of Complex Output Spaces

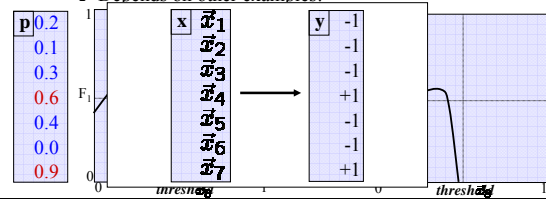
Non-Standard Performance Measures (e.g. F_1 -score, Lift)

- F_1 -score: harmonic average of precision and recall

$$F_1 = \frac{2 \text{Prec} \text{Rec}}{\text{Prec} + \text{Rec}}$$

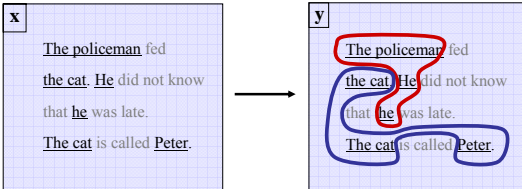
- New example vector \vec{x}_8 . Predict $y_8=1$, if $P(y_8=1/\vec{x}_8)=0.4$?

→ Depends on other examples!



Examples of Complex Output Spaces

- **Noun-Phrase Co-reference**
 - Given a set of noun phrases x , predict a clustering y .
 - Structural dependencies, since prediction has to be an equivalence relation.
 - Correlation dependencies from interactions.



Overview

- **Task: Learning to predict complex outputs**
- **Formalizing the problem**
 - Multi-class classification: Generative vs. Discriminative
 - Compact models for problems with structured outputs
- **Training models with structured outputs**
 - Generative training
 - SVMs and large-margin training
 - Conditional likelihood training
- **Example 1: Learning to parse natural language**
 - Learning weighted context free grammar
- **Example 2: Optimizing F1-score in text classification**
 - Predict vector of class labels
- **Example 3: Learning to cluster**
 - Learning a clustering function that produces desired clusterings
- **Summary & Reading**

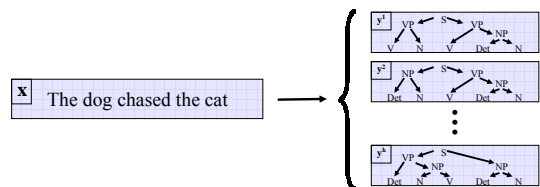
Why do we Need Research on Complex Outputs?

- **Important applications for which conventional methods don't fit!**
 - Noun-phrase co-reference: two step approaches of pair-wise classification and clustering as postprocessing, e.g [Ng & Cardie, 2002]
 - Directly optimize complex loss functions (e.g. F1, AvgPrec)
- **Improve upon existing methods!**
 - Natural language parsing: generative models like probabilistic context-free grammars
 - SVM outperforms naïve Bayes for text classification [Joachims, 1998] [Dumais et al., 1998]

	Precision/Recall Break-Even Point	Naïve Bayes	Linear SVM
• Mor			
- Reuters		72.1	87.5
• Tra			
- WebKB		82.0	90.3
- Ohsumed		62.4	71.6
- Support Vector Machines			

Natural Language Parsing as Multi-Class Classifications

- **Input Space X:** Sequences of words
- **Output Space Y:** Each tree is a class



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):

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

Generative Model: Model $P(X, Y)$

- **Bayes' Decision Rule:** Optimal Decision is

$$h(x) = \underset{y \in Y}{\operatorname{argmin}} \left[\sum_{y' \in Y} \Delta(y, y') P(Y = y' | X = x) \right]$$
- **Equivalent Reformulations:** For 0/1-Loss $\Delta(y, y') = 1$, if $y \neq y'$, 0 else

$$h(x) = \underset{y \in Y}{\operatorname{argmax}} [P(Y = y | X = x)]$$

$$= \underset{y \in Y}{\operatorname{argmax}} \left[\frac{P(X = x | Y = y) P(Y = y)}{P(X = x)} \right]$$

$$= \underset{y \in Y}{\operatorname{argmax}} [P(X = x | Y = y) P(Y = y)]$$

$$= \underset{y \in Y}{\operatorname{argmax}} [P(X = x, Y = y)]$$

$$= \underset{y \in Y}{\operatorname{argmax}} [P(Y = y | X = x) P(X = x)]$$
- **Learning: maximum likelihood (or MAP, or Bayesian)**
 - Assume model class $P(X, Y | \omega)$ with parameters $\omega \in \Omega$
 - Find $\omega' = \underset{\omega \in \Omega}{\operatorname{argmax}} \prod_{i=1}^n [P(Y = y_i, X = x_i | \omega)]$

Naïve Bayes' Classifier (Multivariate)

- **Input Space X: Feature Vector**

- **Output Space Y: {1,-1}**

- **Model:**

- Prior class probabilities

$$P(Y = +1) \quad P(Y = -1)$$

- Class conditional model (one for each class)

$$P(X=x|Y=+1) = \prod_{j=1}^N P(X^{(j)}=x^{(j)}|Y=+1)$$

$$P(X=x|Y=-1) = \prod_{j=1}^N P(X^{(j)}=x^{(j)}|Y=-1)$$

- **Classification rule:**

$$h_{naive}(x) = \underset{y \in \{+1, -1\}}{\operatorname{argmax}} \left\{ P(Y=y) \prod_{j=1}^N P(X^{(j)}=x^{(j)}|Y=y) \right\}$$

fever	cough	pukes	flu?
{3}	{2}	{2}	
high	yes	no	1
high	no	yes	1
low	yes	no	-1
low	yes	yes	1
high	no	yes	???

Estimating the Parameters of Naïve Bayes

- **Count frequencies in training data**

- n : number of training examples
- n_+ / n_- : number of pos/neg examples
- $\#(X^{(j)}=x^{(j)}, y)$: number of times feature $X^{(j)}$ takes value $x^{(j)}$ for examples in class y
- $|X^{(j)}|$: number of values attribute of $X^{(j)}$

fever	cough	pukes	flu?
{3}	{2}	{2}	
high	yes	no	1
high	no	yes	1
low	yes	no	-1
low	yes	yes	1
high	no	yes	???

- **Estimating: $\omega^j = \underset{\omega \in \Omega}{\operatorname{argmax}} \prod_{i=1}^n [P(Y=y_i)] \prod_{j=1}^N [P(X_i^{(j)}=x_i^{(j)}|Y=y_i)]$**

- P(Y): Maximum Likelihood Estimate

$$P(Y=1) = \frac{n_+}{n} \quad P(Y=-1) = \frac{n_-}{n}$$

- P(X|Y): Maximum Likelihood Estimate

$$P(X^{(j)}=x^{(j)}|Y=y) = \frac{\#(X^{(j)}=x^{(j)}, y)}{n_y}$$

- P(X|Y): Smoothing with Laplace estimate

$$P(X^{(j)}=x^{(j)}|Y=y) = \frac{\#(X^{(j)}=x^{(j)}, y) + 1}{n_y + |X^{(j)}|}$$

Discriminative Model: Model P(Y|X)

- **Bayes' Decision Rule:**

- Assume 0/1 Loss $\Delta(y,y') = 1$, if $y \neq y'$, 0 else

- Optimal Decision: $h(x) = \underset{y \in Y}{\operatorname{argmax}} [P(Y=y|X=x)]$

- **Learning: maximum likelihood (or MAP, or Bayesian)**

- Assume model class $P(Y|X, \omega)$ with parameters $\omega \in \Omega$

- Find

$$\omega^j = \underset{\omega \in \Omega}{\operatorname{argmax}} \prod_{i=1}^n [P(Y=y_i|X=x_i, \omega)]$$

- **Example: Logistic regression classifier**

- Assume $P(Y=y|X=x, \omega) = \frac{e^{\omega^T x}}{\sum_{y \in Y} e^{\omega^T x}}$

Discriminative Model:

Model Discriminant Function h Directly

- **Discriminant Function: $h_\omega: X \times Y \rightarrow \mathcal{R}$**

$$h(x) = \underset{y \in Y}{\operatorname{argmax}} \left[\sum_{y' \in Y} \Delta(y, y') P(Y=y'|X=x) \right]$$

$$= \underset{y \in Y}{\operatorname{argmax}} [h(x, y)]$$

- **Consistency of Empirical Risk:**

- Empirical Risk (i.e. training error): $Err_S(h) = \sum_{i=1}^n \Delta(y_i, h(x_i))$

- For sufficiently "small" H_Ω and "large" S : Rule $h' \in H$ with best $Err_S(h')$ has $Err_P(h')$ close to $\min_{h \in H} [Err_P(h)]$

- **Learning: Empirical Risk Minimization (ERM)**

- Assume class H_Ω of discriminant functions $h_\omega: X \rightarrow Y$

- Find

$$h_\omega^j = \underset{h_\omega \in H_\Omega}{\operatorname{argmin}} \sum_{i=1}^n \Delta(y_i, h_\omega(x_i))$$

Generative vs. Discriminative Models

Learning Task:

- Generator: Generate descriptions according to distribution $P(X)$.
- Teacher: Assigns a value to each description based on $P(Y|X)$.

Training Examples $(x_1, y_1), \dots, (x_n, y_n) \sim P(X, Y)$

Discriminative Model

- Model $P(Y|X)$ with $P(Y|X, \omega)$
 - Find ω e.g. via MLE
 - Examples: Log. Reg., CRF
- Model discriminant functions
 - Find $h_\omega \in H$ with low train loss (e.g. Emp. Risk Min.)
 - Examples: SVM, Dec. Tree

Generative Model

- Model $P(X, Y)$ with distributions $P(Y, X| \omega)$
 - Find ω that best matches $P(X, Y)$ on training data (e.g. MLE)
 - Examples: naive Bayes, HMM

Prediction: $h(x) = \underset{y \in Y}{\operatorname{argmax}} [h(x, y)]$

Challenges in Learning with Complex Outputs

- **Approach: view as multi-class classification task**

- Every complex output $y \in Y$ is one class

X The bear chased the cat $\rightarrow Y$ Det \rightarrow N \rightarrow V \rightarrow Det \rightarrow N

- **Problem: Exponentially many classes!**

- Generative Model: $P(X, Y| \omega)$
- Discriminative Model: $P(Y|X, \omega)$
- Discriminant Functions: $h_\omega: X \times Y \rightarrow \mathcal{R}$

- **Challenges**

- How to compactly represent model?

- How to do efficient inference with model (i.e. $\underset{y \in Y}{\operatorname{argmax}} [h(x, y)]$)?

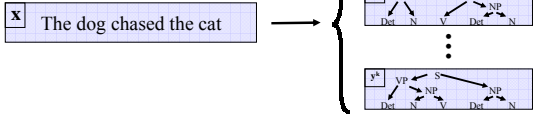
- How to effectively estimate model from data? (e.g. compute $h_\omega^j = \underset{h_\omega \in H_\Omega}{\operatorname{argmin}} \sum_{i=1}^n \Delta(y_i, h_\omega(x_i))$)

Natural Language Parsing

- **Input Space X:** Sequences of words
- **Output Space Y:** Trees over sequences
- **Problem:** How to compute predictions

$$h(x) = \operatorname{argmax}_{y \in Y} [h(x, y)]$$

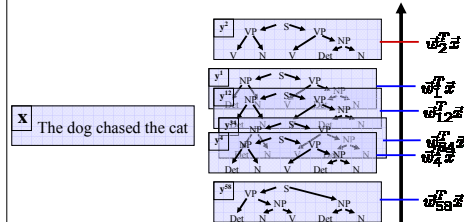
efficiently?



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 $y \in Y$

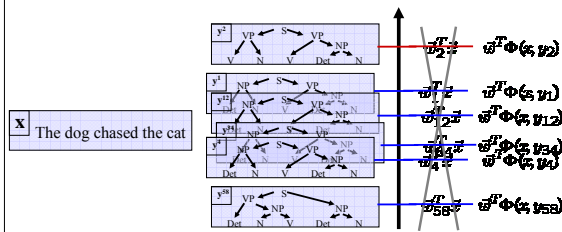
$$h(\vec{x}) = \operatorname{argmax}_{y \in Y} [w_y^T \vec{x}]$$



Joint Feature Map

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

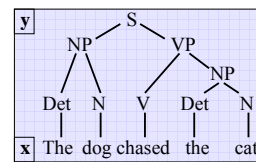
$$h(\vec{x}) = \operatorname{argmax}_{y \in Y} [w^T \Phi(x, y)]$$



Joint Feature Map for 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



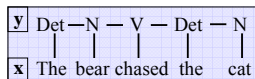
$$\Phi(x, y) = \begin{pmatrix} 1 & S \rightarrow NP VP \\ 0 & S \rightarrow NP \\ 2 & NP \rightarrow Det N \\ 1 & VP \rightarrow V NP \\ \vdots & \\ 0 & Det \rightarrow dog \\ 2 & Det \rightarrow the \\ 1 & N \rightarrow dog \\ 1 & V \rightarrow chased \\ 1 & N \rightarrow cat \end{pmatrix}$$

Joint Feature Map for Sequences

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

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

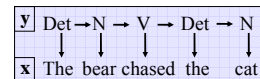
Viterbi



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

Connection to Graphical Models

Hidden Markov Model:



- Assumptions

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

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

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

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

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

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

Overview

- **Task: Learning to predict complex outputs**
- **Formalizing the problem**
 - Multi-class classification: Generative vs. Discriminative
 - Compact models for problems with structured outputs
- ➔ • **Training models with structured outputs**
 - Generative training
 - SVMs and large-margin training
 - Conditional likelihood training
- **Example 1: Learning to parse natural language**
 - Learning weighted context free grammar
- **Example 2: Optimizing F₁-score in text classification**
 - Predict vector of class labels
- **Example 3: Learning to cluster**
 - Learning a clustering function that produces desired clusterings
- **Summary & Reading**

Training Generative Model: HMM

- **Assume:**
 - model class $P(X, Y | \omega)$ with parameters $\omega \in \Omega$
- **Maximum Likelihood (or alternative estimator)**
 - Find $\omega' = \underset{\omega \in \Omega}{\operatorname{argmax}} \prod_{i=1}^n [P(Y = y_i, X = x_i | \omega)]$
 $= \underset{\omega \in \Omega}{\operatorname{argmax}} \left[\prod_{i=1}^n \prod_{j=1}^k P(Y_i = y_i^{(j)} | Y_{i-1} = y_{i-1}^{(j-1)}) P(X_i = x_i^{(j)} | Y_i = y_i^{(j)}) \right]$
- **Example: Hidden Markov Model**
 - Closed-form solutions

$$P(Y_i = y_a | Y_r = y_b) = \frac{\# \text{ of Times State A Follows State B}}{\# \text{ of Times State B Occurs}}$$

$$P(X_i = x_a | Y_i = y_b) = \frac{\# \text{ of Times Output Observed In State B}}{\# \text{ of Times State B Occurs}}$$
 - Need for smoothing the estimates (e.g. max a posteriori)

Support Vector Machine

- Training Examples: $(\vec{x}_1, y_1), \dots, (\vec{x}_n, y_n)$ $\vec{x} \in \mathbb{R}^N$ $y \in \{+1, -1\}$
- Hypothesis Space: $h(\vec{x}) = \operatorname{sgn}[\vec{w}^T \vec{x} + b]$ with $\vec{w} = \sum_{i=1}^n \alpha_i y_i \vec{x}_i$
- Training: Find hyperplane $\langle \vec{w}, b \rangle$ with minimal $\frac{1}{\rho^2} + C \sum_{i=1}^n \xi_i$

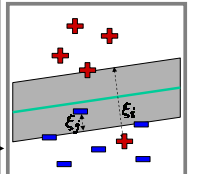
Optimization Problem:

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

$$\text{s.t. } y_1(\vec{w}^T \vec{x}_1 + b) \geq 1 - \xi_1$$

$$\dots$$

$$y_n(\vec{w}^T \vec{x}_n + b) \geq 1 - \xi_n$$

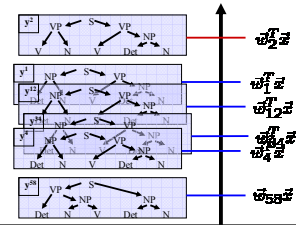


[Vapnik et al.]

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: $h(\vec{x}) = \operatorname{argmax}_{i \in \{1, \dots, k\}} [\vec{w}_i^T \vec{x}]$

X The dog chased the cat



[Crammer & Singer 02]

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

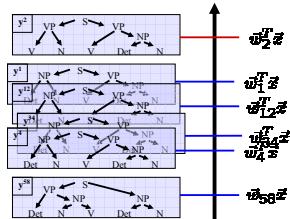
$$\min_{\vec{w}_1, \dots, \vec{w}_k, \xi} \sum_{i=1}^k \vec{w}_i^T \vec{w}_i + C \sum_{i=1}^n \xi_i$$

$$\text{s.t. } \forall j \neq y_1 : \vec{w}_{y_1}^T \vec{x}_1 \geq \vec{w}_j^T \vec{x}_1 + 1 - \xi_1$$

$$\dots$$

$$\forall j \neq y_n : \vec{w}_{y_n}^T \vec{x}_n \geq \vec{w}_j^T \vec{x}_n + 1 - \xi_n$$

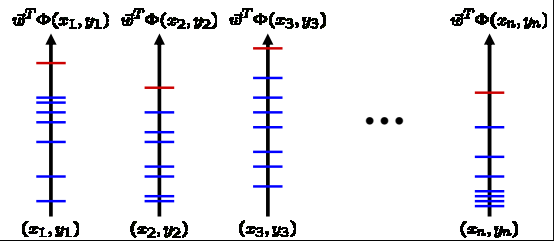
X The dog chased the cat



[Crammer & Singer 02]

Structural Support Vector Machine

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



Structural Support Vector Machine

Hard-margin optimization problem:

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

Loss Functions: Soft-Margin Struct SVM

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

Loss Functions: Soft-Margin Struct SVM

Soft-margin optimization problem:

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

Lemma: The training loss is upper bounded by

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

Training Approach 1: Factored QP

- Assume:
 - Linearly decomposable loss function: $\Delta(y, \hat{y}) = \sum_{i=1}^I \delta(y^{(i)}, \hat{y}^{(i)})$
 - Linear program solution is integral:

$$\hat{y} = \underset{y \in Y}{\operatorname{argmax}} \{ \Delta(y, \hat{y}) + \vec{w}^T \Phi(x_i, \hat{y}) \}$$
- Algorithm:
 - Min-Max Formulation:

$$\min_{\vec{w}, \xi} \frac{1}{2} \vec{w}^T \vec{w} - C \left(\sum_{i=1}^n \underbrace{\vec{w}^T \Phi(x_i, y_i) - \max_{y \in Y} [\Delta(y, y_i) + \vec{w}^T \Phi(x_i, y)]}_{= \xi_i} \right)$$
 - Plug in linear program for "max"
 - Quadratic program can be rewritten so that it has
 - Polynomially many variables
 - Polynomially many constraints

[Taskar et al. 03/04]

Training Approach 2: Cutting-Plane Algorithm for Structural SVM

- Input: $(x_1, y_1), \dots, (x_n, y_n), C, \epsilon$
- $S \leftarrow \emptyset, \vec{w} \leftarrow \mathbf{0}, \vec{\xi} \leftarrow \mathbf{0}$
- REPEAT
 - FOR $i = 1, \dots, n$
 - Find most violated constraint
 - Violated by more than ϵ ?
 - compute $\hat{y} = \underset{y \in Y}{\operatorname{argmax}} \{ \Delta(y, y_i) + \vec{w}^T \Phi(x_i, y) \}$
 - IF $(\Delta(y_i, \hat{y}) - \vec{w}^T [\Phi(x_i, y_i) - \Phi(x_i, \hat{y})]) > \xi_i + \epsilon$
 - $S \leftarrow S \cup \{ \vec{w}^T [\Phi(x_i, y_i) - \Phi(x_i, \hat{y})] \geq \Delta(y_i, \hat{y}) - \xi_i \}$
 - $[\vec{w}, \xi] \leftarrow$ optimize StructSVM over S
 - ENDFOR
 - ENDFOR
- UNTIL S has not changed during iteration

[Altho03] [Jo03] [TsoJoHoAl05]

Training Approach 2: Polynomial Sparsity Bound

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

$$\frac{4CA^2R^2}{\epsilon^2}$$
 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 $\mathbf{0} \leq \Delta(y_i, \hat{y}) \leq A$, and $\|\Phi(x_i, y)\| \leq R$.

[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 $\epsilon = 0.01, C = 1$
- **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

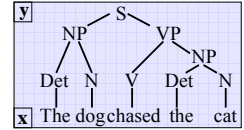
Method	Test Accuracy		Training Efficiency		
	Acc	F_1	CPU-h	Iter	Const
PCFG with MLE	55.2	86.0	0	N/A	N/A
SVM with $(1-F_1)$ -Loss	58.9	88.5	3.4	12	8043

[Tsochantaridis et al. 05]

More Expressive Features

- **Linear composition:**

$$\Phi(x, y) = \sum_{i=1}^I \phi(x, y_i)$$



- **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:**

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

- **Example:**

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

see [Taskar et al. 05]

Applying Structural SVM to New Problem

- **Application specific**
 - Loss function $\Delta(y_i, y)$
 - Representation $\Phi(x, y)$
 - Algorithms to compute

$$\hat{y} = \underset{y \in Y}{\text{argmax}} \{ \tilde{w}^T \Phi(x, y) \}$$

$$\hat{y} = \underset{y \in Y}{\text{argmax}} \{ \Delta(y_i, y) + \tilde{w}^T \Phi(x, y) \}$$
- **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

[Lafferty et al. 01]

Conditional Random Field (CRF)

- **Assume:**

- model class $P(Y|X, \omega)$ with parameters $\omega \in \Omega$
- In particular,
$$P(Y = y | X = x, \omega) = \frac{e^{\tilde{w}^T \Phi(x, y)}}{\sum_{y' \in Y} e^{\tilde{w}^T \Phi(x, y')}}$$

- **Training:** Maximum A Posteriori

- Objective

$$\tilde{w}' = \underset{\tilde{w}}{\text{argmax}} e^{-\frac{1}{2\sigma^2} \tilde{w}'^T \tilde{w}'} \prod_{i=1}^n \left[\frac{e^{\tilde{w}'^T \Phi(x_i, y_i)}}{\sum_{y' \in Y} e^{\tilde{w}'^T \Phi(x_i, y')}} \right]$$

$$= \underset{\tilde{w}}{\text{argmin}} \frac{1}{2\sigma^2} \tilde{w}'^T \tilde{w}' - \sum_{i=1}^n \left[\tilde{w}'^T \Phi(x_i, y_i) - \log \left[\sum_{y' \in Y} e^{\tilde{w}'^T \Phi(x_i, y')} \right] \right]$$
- Gradient

$$\frac{\delta \mathcal{L}}{\delta w_j} = \frac{\tilde{w}'_j}{\sigma^2} - \sum_{i=1}^n \left[\Phi_j(x_i, y_i) - \frac{\delta}{\delta w_j} \left[\log \left[\sum_{y' \in Y} e^{\tilde{w}'^T \Phi(x_i, y')} \right] \right] \right]$$

- Methods: iterative scaling, quasi Newton, conjugate gradient [Sha & Pereira 03] [Wallach 02] [Malouf 02] [Minka 03]

Applying CRF to New Problem

- **Application specific**
 - Representation $\Phi(x, y)$
 - Algorithms to compute

$$\hat{y} = \underset{y \in Y}{\text{argmax}} \{ \tilde{w}^T \Phi(x, y) \}$$

$$\frac{\delta Z(x_i)}{\delta w_j} = \frac{\delta}{\delta w_j} \left[\log \left[\sum_{y' \in Y} e^{\tilde{w}^T \Phi(x_i, y')} \right] \right]$$
- **Implementation of CRF:** <http://mallet.cs.umass.edu/>

Relationship CRF and Structural SVM

- **Objective Functions:**

- CRF:

$$\tilde{w}' = \underset{\tilde{w}}{\text{argmin}} \frac{1}{2\sigma^2} \tilde{w}'^T \tilde{w}' - \sum_{i=1}^n \left[\tilde{w}'^T \Phi(x_i, y_i) - \log \left[\sum_{y' \in Y} e^{\tilde{w}'^T \Phi(x_i, y')} \right] \right]$$
- Structural SVM:

$$\tilde{w}' = \underset{\tilde{w}}{\text{argmin}} \frac{1}{2} \tilde{w}'^T \tilde{w}' - C \left(\sum_{i=1}^n \tilde{w}'^T \Phi(x_i, y_i) - \max_{y' \in Y} [\Delta(y_i, y') + \tilde{w}'^T \Phi(x_i, y')] \right)$$

- **Basic Cases for two classes**

- CRF: Regularized logistic regression
 - Structural SVM: binary classification SVM

Overview

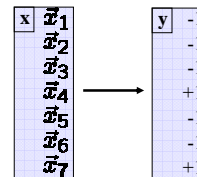
- **Task: Learning to predict complex outputs**
- **Formalizing the problem**
 - Multi-class classification: Generative vs. Discriminative
 - Compact models for problems with structured outputs
- **Training models with structured outputs**
 - Generative training
 - SVMs and large-margin training
 - Conditional likelihood training
- **Example 1: Learning to parse natural language**
 - Learning weighted context free grammar
- **Example 2: Optimizing F_1 -score in text classification**
 - Predict vector of class labels
- **Example 3: Learning to cluster**
 - Learning a clustering function that produces desired clusterings
- **Summary & Reading**

Examples of Complex Output Spaces

- **Non-Standard Performance Measures (e.g. F_1 -score, Lift)**
 - F_1 -score: harmonic average of precision and recall

$$F_1 = \frac{2 \text{Prec Rec}}{\text{Prec} + \text{Rec}}$$

- New example vector \vec{x}_8 . Predict $y_8=1$, if $P(y_8=1|\vec{x}_8)=0.4?$
 - Depends on other examples!



Struct SVM for Optimizing F_1 -Score

- **Loss Function**
 - $\Delta(\vec{y}, \hat{y}) = (1 - F_1) = \left(1 - \frac{2 \text{Prec Rec}}{\text{Prec} + \text{Rec}}\right)$
- **Representation**
 - $\vec{x} = (\vec{x}_1, \vec{x}_2, \vec{x}_3, \dots, \vec{x}_n)$
 - $\vec{y} = (y_1, y_2, y_3, \dots, y_n)$
 - Joint feature map $\Phi(\vec{x}, \vec{y}) = \sum_{i=1}^n y_i \vec{x}_i$
- **Prediction**
 - $\hat{y} = \underset{y \in \mathcal{Y}}{\text{argmax}} \{w^T \Phi(\vec{x}, y)\} \Leftrightarrow \hat{y} = \begin{cases} \text{sign}(w^T \vec{x}_1) \\ \text{sign}(w^T \vec{x}_2) \\ \vdots \\ \text{sign}(w^T \vec{x}_n) \end{cases}$
- **Find most violated constraint**
 - $\hat{y} = \underset{y \in \mathcal{Y}}{\text{argmax}} \{ \Delta(\vec{y}, \hat{y}) + w^T \Phi(\vec{x}, \hat{y}) \}$
 - Only n^2 different contingency tables → search brute force

[Jo 05]

Struct SVM for Optimizing F_1 -Score

- **Loss F**
 - $\Delta(\vec{y}, \hat{y}) = (1 - F_1) = \left(1 - \frac{2 \text{Prec Rec}}{\text{Prec} + \text{Rec}}\right)$
- **Repre**
 - $\vec{x} = (\vec{x}_1, \vec{x}_2, \vec{x}_3, \dots, \vec{x}_n)$
 - $\vec{y} = (y_1, y_2, y_3, \dots, y_n)$
 - Joint feature map $\Phi(\vec{x}, \vec{y}) = \sum_{i=1}^n y_i \vec{x}_i$
- **Predic**
 - $\hat{y} = \underset{y \in \mathcal{Y}}{\text{argmax}} \{w^T \Phi(\vec{x}, y)\} \Leftrightarrow \hat{y} = \begin{cases} \text{sign}(w^T \vec{x}_1) \\ \text{sign}(w^T \vec{x}_2) \\ \vdots \\ \text{sign}(w^T \vec{x}_n) \end{cases}$
- **Find m**
 - $\hat{y} = \underset{y \in \mathcal{Y}}{\text{argmax}} \{ \Delta(\vec{y}, \hat{y}) + w^T \Phi(\vec{x}, \hat{y}) \}$
 - Only n^2 different contingency tables → search brute force

[Jo 05]

Experiment: Text Classification

- **Dataset: Reuters-21578 (ModApte)**
 - 9603 training / 3299 test examples, 90 categories
 - TFIDF unit vectors (no stemming, no stopword removal)
- **Experiment Setup**
 - Classification SVM with optimal C in hindsight
 - Linear Cost SVM [Morik et al., 1999] with C^*/C via 2-CV
 - F_1 -loss SVM with C via 2-CV
- **Results**

Method	Test F_1	Training Efficiency		
		CPU-min	Const	SV
Classification SVM	52.8	0.1	N/A	264
Linear Cost SVM	56.1	0.1	N/A	371
SVM with $(1-F_1)$ -Loss	62.0	32.2	173	86

[Jo 05]

Multivariate SVM Generalizes Classification SVM

Theorem: The solutions of the multivariate SVM with number of errors as the loss function and an (unbiased) classification SVM are equal.

Multivariate SVM optimizing Error Rate:

$$\min_{\vec{w}, \xi \geq 0} \frac{1}{2} \vec{w}^T \vec{w} + C \xi$$

$$s.t. \forall \vec{y} \in \mathcal{Y} : \vec{w}^T \Phi(\vec{x}, \vec{y}) \geq \vec{w}^T \Phi(\vec{x}, \hat{y}) + 2E_{\tau}(\vec{y}, \hat{y}) - \xi$$

\Leftrightarrow

Classification SVM (unbiased):

$$\min_{\vec{w}, \xi \geq 0} \frac{1}{2} \vec{w}^T \vec{w} + 2C \sum_{i=1}^n \xi_i$$

$$s.t. \forall_1 (w^T x_1) \geq 1 - \xi_1, \dots, \forall_n (w^T x_n) \geq 1 - \xi_n$$

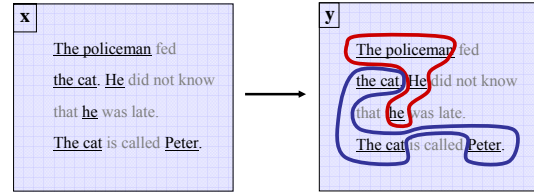
[Jo 05]

Overview

- **Task: Learning to predict complex outputs**
- **Formalizing the problem**
 - Multi-class classification: Generative vs. Discriminative
 - Compact models for problems with structured outputs
- **Training models with structured outputs**
 - Generative training
 - SVMs and large-margin training
 - Conditional likelihood training
- **Example 1: Learning to parse natural language**
 - Learning weighted context free grammar
- **Example 2: Optimizing F₁-score in text classification**
 - Predict vector of class labels
- **Example 3: Learning to cluster**
 - Learning a clustering function that produces desired clusterings
- **Summary & Reading**

Learning to Cluster

- **Noun-Phrase Co-reference**
 - Given a set of noun phrases x , predict a clustering y .
 - Structural dependencies, since prediction has to be an equivalence relation.
 - Correlation dependencies from interactions.



[Finley & Jo 05], see also [McCallum & Wellner 04]

Struct SVM for Supervised Clustering

- **Representation**
 - $x = \begin{pmatrix} \bar{x}_{11} & \bar{x}_{12} & \dots & \bar{x}_{1n} \\ \bar{x}_{21} & \bar{x}_{22} & \dots & \bar{x}_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ \bar{x}_{n1} & \bar{x}_{n2} & \dots & \bar{x}_{nn} \end{pmatrix} \quad y = \begin{pmatrix} y_{11} & y_{12} & \dots & y_{1n} \\ y_{21} & y_{22} & \dots & y_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ y_{n1} & y_{n2} & \dots & y_{nn} \end{pmatrix} \quad y_{ij} \in \{0, 1\}$
 - y is reflexive ($y_{ii}=1$), symmetric ($y_{ij}=y_{ji}$), and transitive (if $y_{ij}=1$ and $y_{jk}=1$, then $y_{ik}=1$)
 - Joint feature map $\Phi(x, y) = \sum_{i=1}^n \sum_{j=1}^n y_{ij} \bar{x}_{ij}$
- **Loss Function**
 - $\Delta(x, y) = \|y - \bar{y}\|_1$
- **Prediction**
 - $\hat{y} = \underset{y \in \mathcal{Y}}{\operatorname{argmax}} \{ \bar{y}^T \Phi(x, \bar{y}) \} = \underset{y \in \mathcal{Y}}{\operatorname{argmax}} \{ \sum_{i,j} y_{ij} (\bar{y}^T \bar{x}_{ij}) \}$
 - NP hard, use linear relaxation instead [Demaine & Immerlica, 2003]
- **Find most violated constraint**
 - $\hat{y} = \underset{y \in \mathcal{Y}}{\operatorname{argmax}} \{ \Delta(x, \bar{y}) - \bar{y}^T \Phi(x, \bar{y}) \}$
 - NP hard, use linear relaxation instead [Demaine & Immerlica, 2003]

[Finley & Jo 05], see also [McCallum & Wellner 04]

Struct SVM for Supervised Clustering

- **Representation**
 - $\bar{y}^T \bar{x}_{ij} = \begin{pmatrix} 1.0 & 2.5 & 0.0 & 0.0 & -1.1 & -0.1 & 0.0 \\ 0.0 & 1.0 & 3.0 & 0.1 & 0.1 & 0.3 & 0.0 \\ 0.0 & 0.0 & 1.0 & 0.3 & 2.2 & -0.9 & 0.0 \\ 0.0 & -0.4 & -0.1 & 1.0 & -1.1 & -0.1 & -2.0 \\ 0.0 & 0.0 & 0.0 & -0.3 & 1.0 & 2.3 & 0.0 \\ 0.1 & 0.1 & 0.0 & 0.1 & 0.1 & 1.0 & 1.0 \\ -1.0 & -0.1 & 0.0 & -0.2 & -1.1 & 0.1 & 1.0 \end{pmatrix} \rightarrow \begin{pmatrix} y_1 & 1 & 1 & 1 & 0 & 0 & 0 \\ 1 & 1 & 1 & 1 & 0 & 0 & 0 \\ 1 & 1 & 1 & 1 & 0 & 0 & 0 \\ 1 & 1 & 1 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 1 & 1 \\ 0 & 0 & 0 & 0 & 1 & 1 & 1 \\ 0 & 0 & 0 & 0 & 1 & 1 & 1 \end{pmatrix} \in \{0, 1\}$
- **Loss**
 - $\Delta(x, y) = \|y - \bar{y}\|_1$
- **Prediction**
 - $\hat{y} = \begin{pmatrix} y_1 & 1 & 1 & 1 & 0 & 0 & 0 \\ 1 & 1 & 1 & 1 & 0 & 0 & 0 \\ 1 & 1 & 1 & 1 & 0 & 0 & 0 \\ 1 & 1 & 1 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 1 & 1 \\ 0 & 0 & 0 & 0 & 1 & 1 & 1 \\ 0 & 0 & 0 & 0 & 1 & 1 & 1 \end{pmatrix}$
 - NI
- **Find r**
 - $\hat{y} = \begin{pmatrix} y_1 & 1 & 1 & 0 & 0 & 0 & 0 \\ 1 & 1 & 1 & 0 & 0 & 0 & 0 \\ 1 & 1 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 1 & 1 & 1 \\ 0 & 0 & 0 & 1 & 1 & 1 & 1 \\ 0 & 0 & 0 & 1 & 1 & 1 & 1 \end{pmatrix}$
 - NI

Summary

- **Learning to predict structured and interdependent output**
 - Discriminant function: $h(x) = \underset{y \in \mathcal{Y}}{\operatorname{argmax}} [h(x, y)]$
- **Training:**
 - Generative
 - Structural SVM
 - Conditional Random Field
- **Examples**
 - Learning to predict trees (natural language parsing)
 - Optimize to non-standard performance measures (imbalanced classes)
 - Learning to cluster (noun-phrase coreference resolution)
- **Software:**
 - SVM^{struct}: <http://svmlight.joachims.org/>
 - Mallet: <http://mallet.cs.umass.edu/>

Reading

- **Generative training**
 - Hidden-Markov models [Manning & Schuetze, 1999]
 - Probabilistic context-free grammars [Manning & Schuetze, 1999]
 - Markov random fields [Geman & Geman, 1984]
 - Etc.
- **Discriminative training**
 - Multivariate output regression [Izeman, 1975] [Breiman & Friedman, 1997]
 - Kernel Dependency Estimation [Weston et al. 2003]
 - Conditional HMM [Krogh, 1994]
 - Transformer networks [LeCun et al., 1998]
 - Conditional random fields [Lafferty et al., 2001] [Sutton & McCallum, 2005]
 - Perceptron training of HMM [Collins, 2002]
 - Structural SVMs / Maximum-margin Markov networks [Taskar et al., 2003] [Tsochantaridis et al., 2004, 2005] [Taskar 2004]