

# Discriminative Learning for Structured Output Prediction

CS4780/5780 – Machine Learning  
Fall 2014

Thorsten Joachims  
Cornell University

Reading:

T. Joachims, T. Hofmann, Yisong Yue, Chun-Nam Yu,  
*Predicting Structured Objects with Support Vector Machines*,  
Communications of the ACM, Research Highlight, 52(11):97-104, 2009.

[http://www.cs.cornell.edu/people/tj/publications/joachims\\_etal\\_09b.pdf](http://www.cs.cornell.edu/people/tj/publications/joachims_etal_09b.pdf)

# Discriminative vs. Generative

## Bayes Decision Rule

$$\begin{aligned} - h_{bayes}(x) &= \operatorname{argmax}_{y \in Y} [P(Y = y|X = x)] \\ &= \operatorname{argmax}_{y \in Y} [P(X = x|Y = y)P(Y = y)] \end{aligned}$$

## Generative:

- Idea: Make assumptions about  $P(X = x|Y = y), P(Y = y)$
- Method: Estimate parameters of the two distributions, then apply Bayes decision rule.

## Discriminative:

- Idea: Define set of prediction rules (i.e. hypotheses)  $H$ , then search for  $h \in H$  that best approximates

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

- Method: find  $h \in H$  that minimizes training error.

Question: Can we train HMM's discriminately?

# Idea for Discriminative Training of HMM

Idea:

- $h_{bayes}(x) = \operatorname{argmax}_{y \in Y} [P(Y = y|X = x)]$   
 $= \operatorname{argmax}_{y \in Y} [P(X = x|Y = y)P(Y = y)]$
- Model  $P(Y = y|X = x)$  with  $\vec{w} \cdot \phi(x, y)$  so that  
 $(\operatorname{argmax}_{y \in Y} [P(Y = y|X = x)]) = (\operatorname{argmax}_{y \in Y} [\vec{w} \cdot \phi(x, y)])$

Hypothesis Space:

$$h(x) = \operatorname{argmax}_{y \in Y} [\vec{w} \cdot \phi(x, y)] \text{ with } \vec{w} \in \mathfrak{R}^N$$

Intuition:

- Tune  $\vec{w}$  so that correct  $y$  has the highest value of  $\vec{w} \cdot \phi(x, y)$
- $\phi(x, y)$  is a feature vector that describes the match between  $x$  and  $y$

# Training HMMs with Structural SVM

- HMM

$$P(x, y) = P(y_1)P(x_1|y_1) \prod_{i=2}^l P(x_i|y_i)P(y_i|y_{i-1})$$

$$\log P(x, y) = \log P(y_1) + \log P(x_1|y_1) \sum_{i=2}^l \log P(x_i|y_i) + \log P(y_i|y_{i-1})$$

- Define  $\phi(x, y)$  so that model is isomorphic to HMM
  - One feature for each possible start state
  - One feature for each possible transition
  - One feature for each possible output in each possible state
  - Feature values are counts

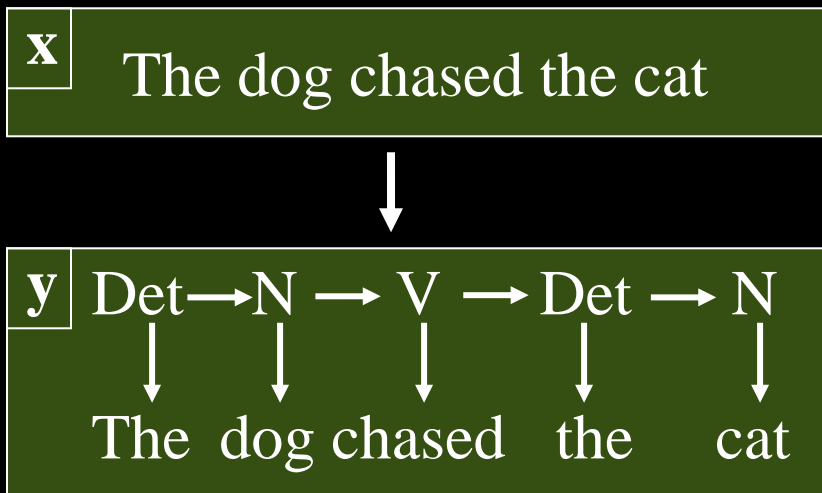
# Joint Feature Map for Sequences

- Linear Chain HMM

- Each transition and emission has a weight
- Score of a sequence is the sum of its weights

- Find highest scoring sequence  $h(x) = \operatorname{argmax}_{y \in Y} [\vec{w} \cdot \phi(x, y)]$

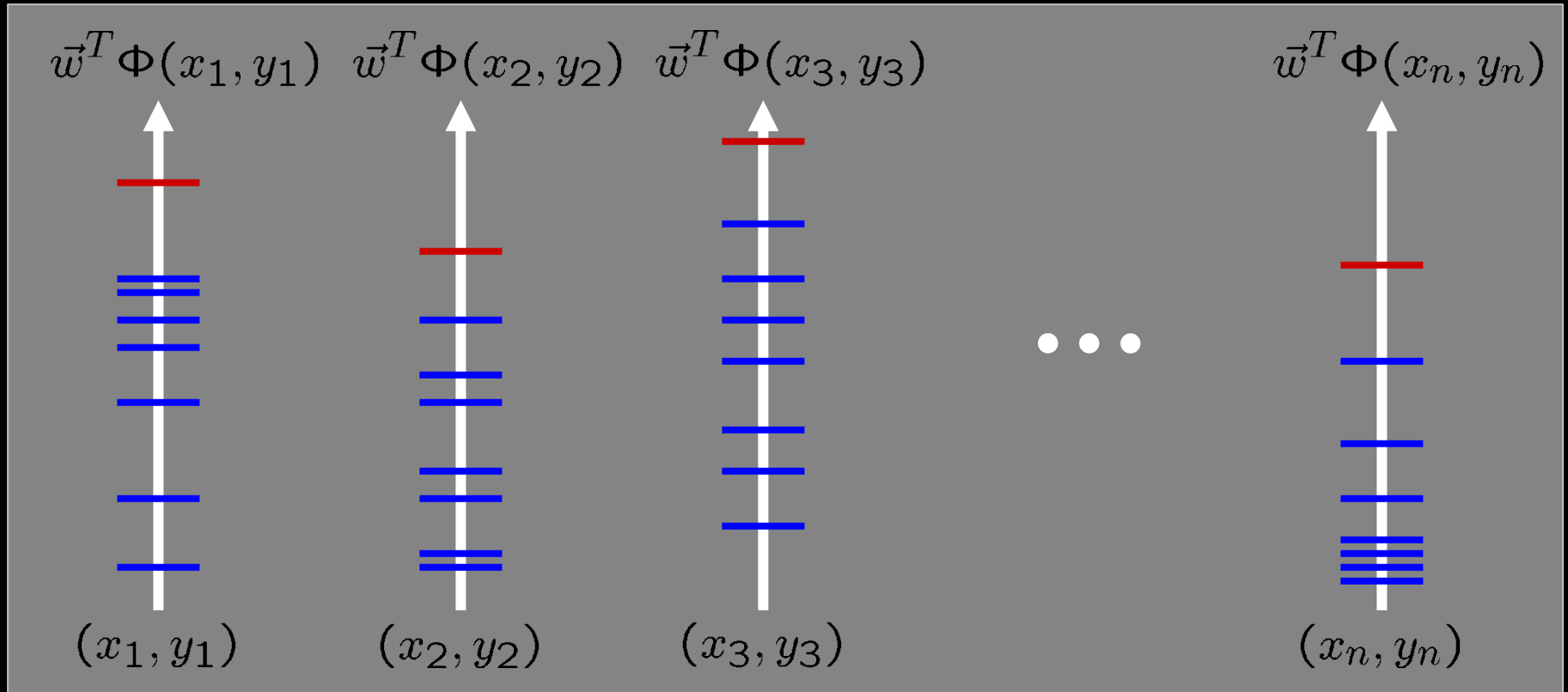
Viterbi



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

# Structural Support Vector Machine

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



# Structural SVM Training Problem

Hard-margin optimization problem:

$$\min_{\vec{w}} \quad \frac{1}{2} \vec{w}^T \vec{w}$$

$$s.t. \quad \forall y \in Y \setminus y_1 : \vec{w}^T \Phi(x_1, y_1) \geq \vec{w}^T \Phi(x_1, y) + 1$$

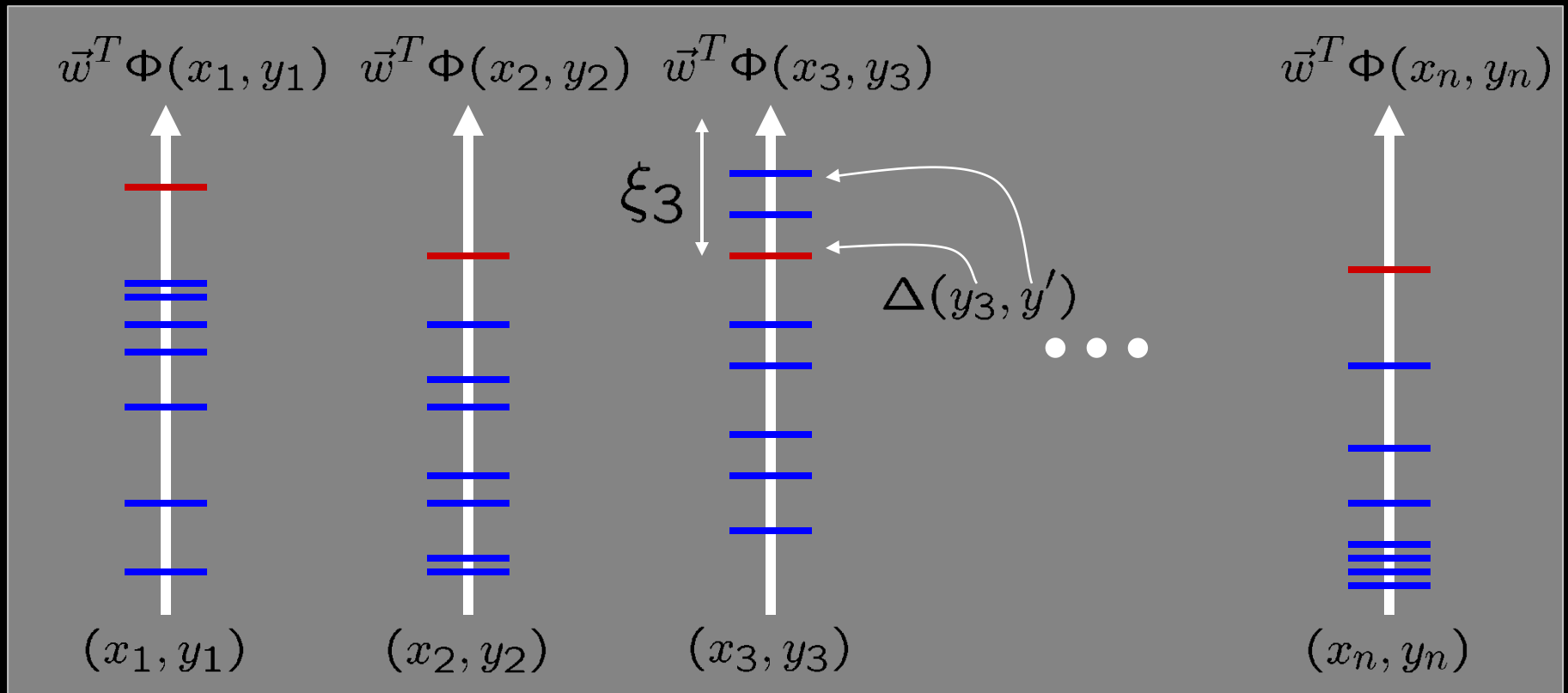
...

$$\forall y \in Y \setminus y_n : \vec{w}^T \Phi(x_n, y_n) \geq \vec{w}^T \Phi(x_n, y) + 1$$

- Training Set:  $(x_1, y_1), \dots, (x_n, y_n)$
- Prediction Rule:  $h_{svm}(x) = \operatorname{argmax}_{y \in Y} [\vec{w} \cdot \phi(x, y)]$
- Optimization:
  - Correct label  $y_i$  must have higher value of  $\vec{w} \cdot \phi(x, y)$  than any incorrect label  $y$
  - Find weight vector with smallest norm

# Soft-Margin Structural SVM

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





# Soft-Margin Structural SVM

**Soft-margin optimization problem:**

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

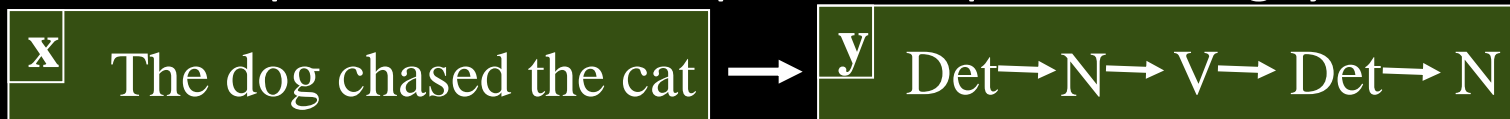
**Lemma: The training loss is upper bounded by**

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

# Experiment: Part-of-Speech Tagging

- **Task**

- Given a sequence of words  $x$ , predict sequence of tags  $y$ .



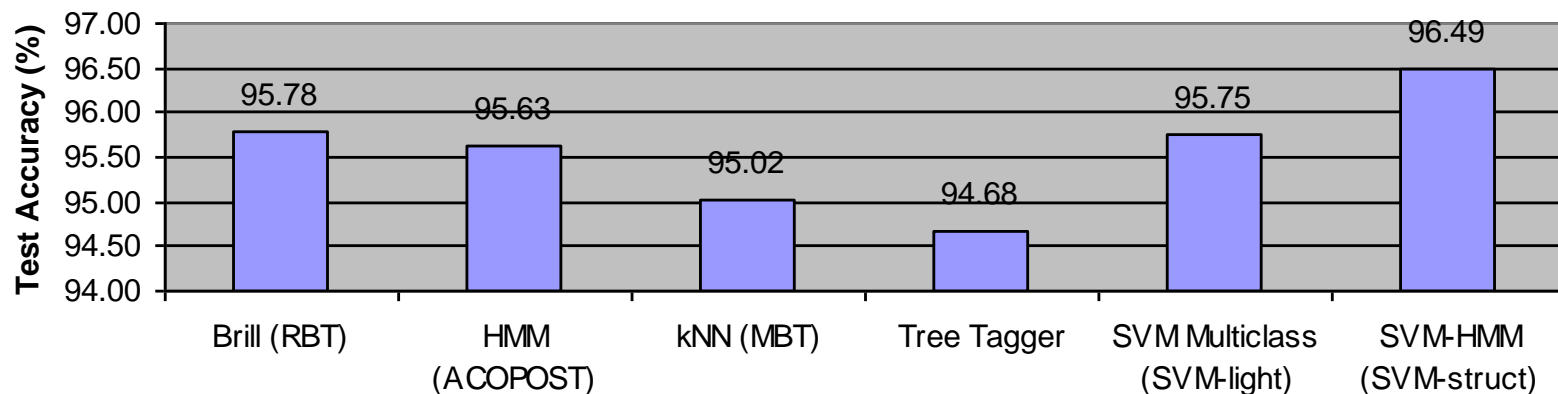
- Dependencies from tag-tag transitions in Markov model.

- **Model**

- Markov model with one state per tag and words as emissions
- Each word described by  $\sim 250,000$  dimensional feature vector (all word suffixes/prefixes, word length, capitalization ...)

- **Experiment (by Dan Fleisher)**

- Train/test on 7966/1700 sentences from Penn Treebank



# NE Identification

- Identify all named locations, named persons, named organizations, dates, times, monetary

The delegation, which included the commander of the U.N. troops in Bosnia, Lt. Gen. Sir Michael Rose, went to the Serb stronghold of Pale, near Sarajevo, for talks with Bosnian Serb leader Radovan Karadzic.

Este ha sido el primer comentario publico del presidente Clinton respecto a la crisis de Oriente Medio desde que el secretario de Estado, Warren Christopher, decidiera regresar precipitadamente a Washington para impedir la ruptura del proceso de paz tras la violencia desatada en el sur de Libano.

1. Locations
2. Persons
3. Organizations

**Figure 1.1 Examples.** Examples of correct labels for English text and for Spanish text.

# Experiment: Named Entity Recognition

- Data
  - Spanish Newswire articles
  - 300 training sentences
  - 9 tags
    - no-name,
    - beginning and continuation of person name, organization, location, misc name
  - Output words are described by features (e.g. starts with capital letter, contains number, etc.)
- Error on test set (% mislabeled tags):
  - Generative HMM: 9.36%
  - Support Vector Machine HMM: 5.08%

# Cutting-Plane Algorithm for Structural SVM

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

Find most violated constraint

Violated by more than  $\epsilon$  ?

Add constraint to working set

→ Polynomial Time Algorithm (SVM-struct)

# General Problem: Predict Complex Outputs

- Supervised Learning from Examples
  - Find function from input space  $X$  to output space  $Y$

$$h: X \rightarrow Y$$

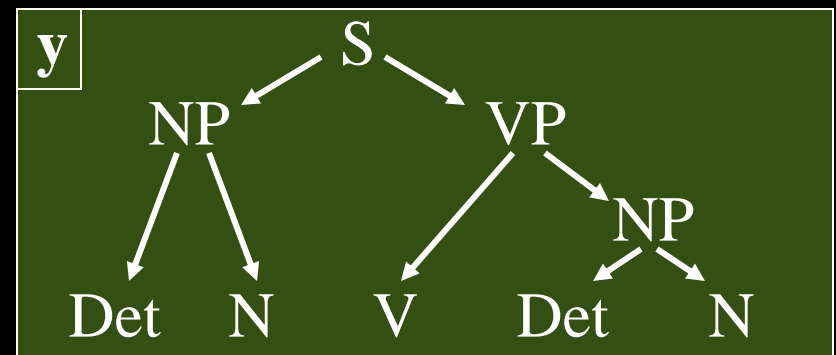
such that the prediction error is low.

- Typical
  - Output space is just a single number
    - Classification: -1,+1
    - Regression: some real number
- General
  - Predict outputs that are complex objects

# 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.

$x$  The dog chased the cat



# 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)

**x** Due to the continued violence in Baghdad, the oil price is expected to further increase. OPEC officials met with ...



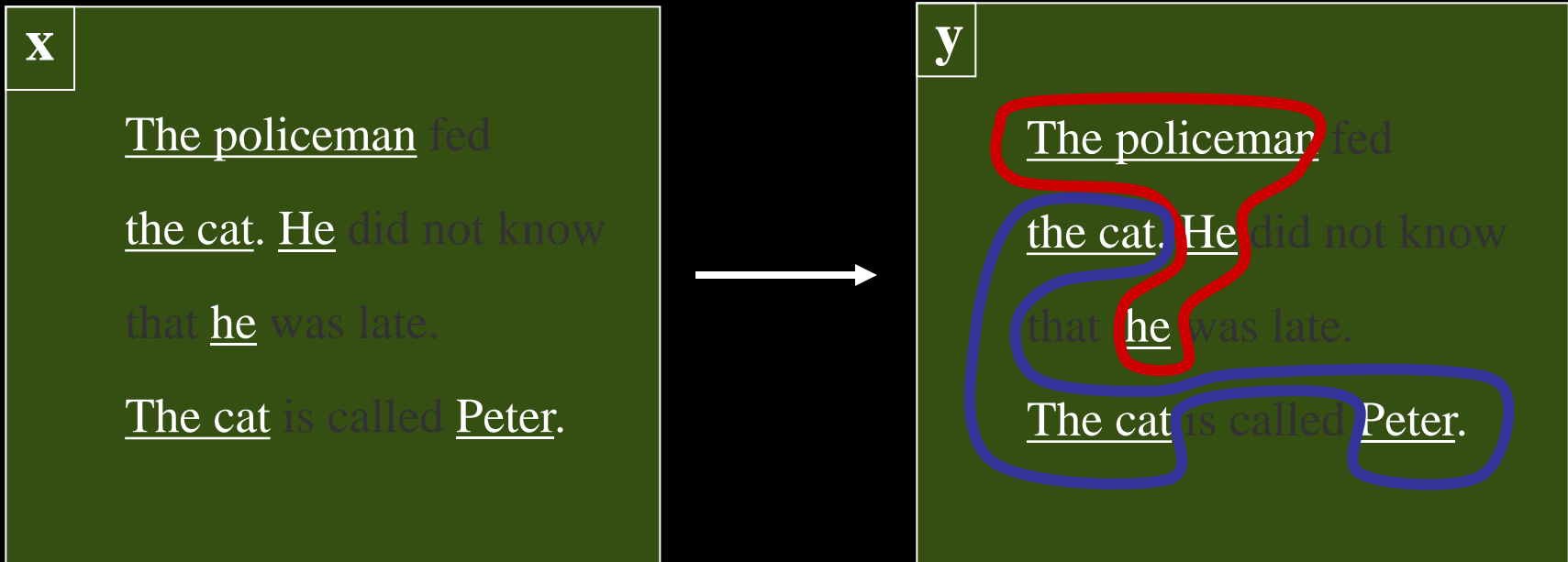
**y**

|    |              |
|----|--------------|
| -1 | antarctica   |
| -1 | benelux      |
| -1 | germany      |
| +1 | iraq         |
| +1 | oil          |
| -1 | coal         |
| -1 | trade        |
| -1 | acquisitions |



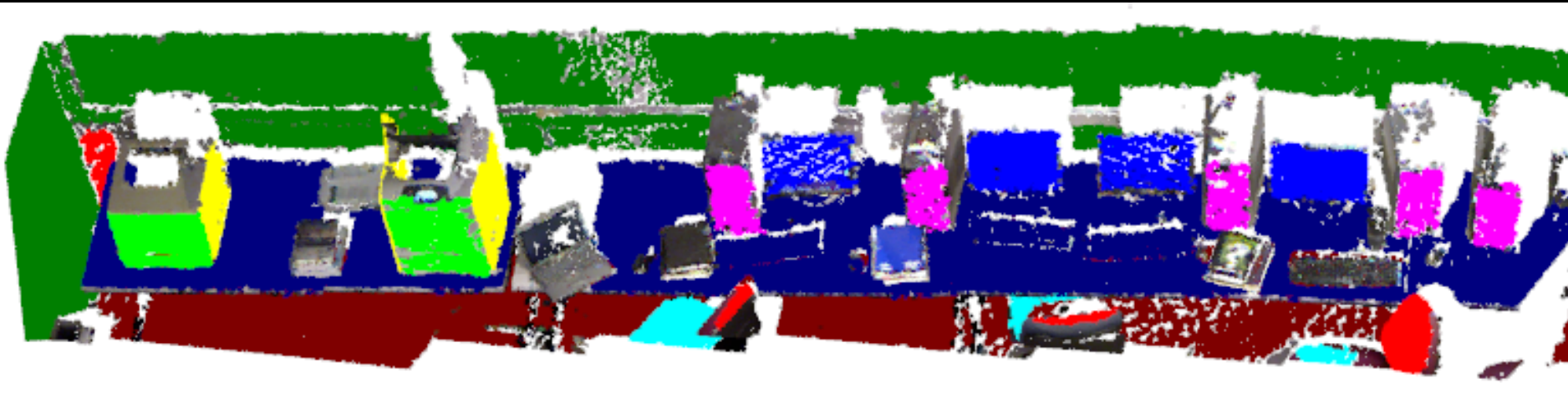
# Examples of Complex Output Spaces

- 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.



# Examples of Complex Output Spaces

- Scene Recognition
  - Given a 3D point cloud with RGB from Kinect camera
  - Segment into volumes
  - Geometric dependencies between segments (e.g. monitor usually close to keyboard)



Wrap-Up

# Batch Learning for Classification

- Discriminative
    - Decision Trees →
    - Perceptron →
    - Linear SVMs →
    - Kernel SVMs →
  - Other Methods
    - Logical rule learning
    - Stochastic Gradient
    - Logistic Regression
    - Neural Networks
    - RBF Networks
    - Boosting
    - Bagging
    - Parametric (Graphical) Models
    - Non-Parametric Models
  - Generative
    - Multinomial Naïve Bayes →
    - Multivariate Naïve Bayes →
    - Less Naïve Bayes →
    - Linear Discriminant →
    - Nearest Neighbor →
- Methods + Theory + Practice
- \*-Regression
  - \*-Multiclass
-

# Online Learning

- Expert Setting
    - Halving Algorithm
    - Weighted Majority
    - Exponentiated Gradient
  - Bandit Setting
    - EXP3
  - Other Settings
    - Dueling Bandits
    - Coactive Learning
  - Other Methods
    - Hedge
    - Follow the Leader
    - UCB
    - Zooming
    - Partial Monitoring
- CS4786 Machine Learning for Data Science
- CS6783 Machine Learning Theory

# Structured Prediction

- Discriminative
    - Structural SVMs
  - Generative
    - Hidden Markov Model
  - Other Methods
    - Maximum Margin Markov Networks
    - Conditional Random Fields
    - Markov Random Fields
    - Bayesian Networks
    - Statistical Relational Learning
- CS4786 Machine Learning for Data Science
- CS4782 Probabilistic Graphical Models

# Unsupervised Learning

- Clustering
    - Hierarchical Agglomerative Clustering
    - K-Means
    - Mixture of Gaussians and EM-Algorithm
  - Other Methods
    - Spectral Clustering
    - Latent Dirichlet Allocation
    - Latent Semantic Analysis
    - Multi-Dimensional Scaling
  - Other Tasks
    - Outlier Detection
    - Novelty Detection
    - Dimensionality Reduction
    - Non-Linear Manifold Detection
- CS4786 Machine Learning for Data Science
- CS4850 Math Found for the Information Age

# Other Learning Problems and Applications

- Recommender Systems, Search Ranking, etc.
  - CS4300 Information Retrieval
- Reinforcement Learning and Markov Decision Processes
  - CS4758 Robot Learning
- Computer Vision
  - CS4670 Intro Computer Vision
- Natural Language Processing
  - CS4740 Intro Natural Language Processing



# Other Machine Learning Courses at Cornell

- INFO 3300 - Data-Driven Web Pages
- CS 4700 - Introduction to Artificial Intelligence
- CS 4780/5780 - Machine Learning (for Intelligent Systems)
- **CS 4786/5786 - Machine Learning for Data Science**
- CS 4758 - Robot Learning
- CS 4782 - Probabilistic Graphical Models
- OR 4740 - Statistical Data Mining
- CS 6780 - Advanced Machine Learning
- CS 6783 - Machine Learning Theory
- CS 6784 - Advanced Topics in Machine Learning
- CS 6756 - Advanced Topics in Robot Learning
- ORIE 6740 - Statistical Learning Theory for Data Mining
- ORIE 6750 - Optimal learning
- ORIE 6780 - Bayesian Statistics and Data Analysis
- MATH 7740 - Statistical Learning Theory