Statistical Learning for Text Data Analytics Sequence Labeling and Structured Output Learning: Conditional Models and Global Classifiers

#### Yangqiu Song

#### Hong Kong University of Science and Technology

yqsong@cse.ust.hk

#### Spring 2018

 $\ast \textsc{Contents}$  are based on materials created by Vivek Srikumar, Dan Roth, Andrew Ng

- Vivek Srikumar. CS 6355 Structured Prediction. https: //svivek.com/teaching/structured-prediction/spring2018/
- Dan Roth. CS546: Machine Learning and Natural Language . http://l2r.cs.uiuc.edu/~danr/Teaching/CS546-16/
- Andrew Ng. CS229: Machine Learning. http://cs229.stanford.edu/



- Representation: language models, word embeddings, topic models
- Learning: supervised learning, semi-supervised learning, sequence models, deep learning, optimization techniques
- Inference: constraint modeling, joint inference, search algorithms

NLP applications: tasks introduced in Lecture 1

Yangqiu Song (HKUST)

Learning for Text Analytics

## Overview

#### 1 Hidden Markov Models

- Representation
- Learning
- Inference

#### Conditional Models and Local Classifiers

- Conditional Models for Predicting Sequences
- Log-linear Models for Multiclass Classification

#### • Maximum Entropy Markov Models

• The Label Bias Problem

#### Global Models

- Conditional Random Fields
- Structured Perceptron for Sequences
- Structural Support Vector Machine

$$P(y_n|y_{n-1}, y_{n-2}, \ldots, \mathbf{x}_n, \mathbf{x}_{n-1}, \mathbf{x}_{n-2}, \ldots) = P(y_n|y_{n-1}, \mathbf{x}_n)$$



• This assumption lets us write the conditional probability of the output as \_\_\_\_

$$P(y_{1:N}|\mathbf{x}_{1:N}) = \prod_{n} P(y_n|y_{n-1},\mathbf{x}_n)$$

Goal: Compute  $P(y_{1:N}|\mathbf{x}_{1:N}, \mathbf{w}) = \prod_n P(y_n|y_{n-1}, \mathbf{x}_{1:N})$  where  $P(y_n|y_{n-1}, \mathbf{x}_{1:N}) \propto \exp(\mathbf{w}^\top \phi(\mathbf{x}, n, y_n, y_{n-1}))$ 



The prediction task: Using the entire input and the current label, predict the next label

Yangqiu Song (HKUST)

Learning for Text Analytics

Spring 2018 6 / 52

Goal: Compute 
$$P(y_{1:N}|\mathbf{x}_{1:N}, \mathbf{w}) = \prod_n P(y_n|y_{n-1}, \mathbf{x}_{1:N})$$
 where  
 $P(y_n|y_{n-1}, \mathbf{x}_{1:N}) \propto \exp(\mathbf{w}^\top \phi(\mathbf{x}, n, y_n, y_{n-1}))$ 



Goal: Compute 
$$P(y_{1:N}|\mathbf{x}_{1:N}, \mathbf{w}) = \prod_n P(y_n|y_{n-1}, \mathbf{x}_{1:N})$$
 where  
 $P(y_n|y_{n-1}, \mathbf{x}_{1:N}) \propto \exp(\mathbf{w}^\top \phi(\mathbf{x}, n, y_n, y_{n-1}))$ 



Goal: Compute 
$$P(y_{1:N}|\mathbf{x}_{1:N}, \mathbf{w}) = \prod_n P(y_n|y_{n-1}, \mathbf{x}_{1:N})$$
 where  
 $P(y_n|y_{n-1}, \mathbf{x}_{1:N}) \propto \exp(\mathbf{w}^\top \phi(\mathbf{x}, n, y_n, y_{n-1}))$ 



Goal: Compute 
$$P(y_{1:N}|\mathbf{x}_{1:N}, \mathbf{w}) = \prod_n P(y_n|y_{n-1}, \mathbf{x}_{1:N})$$
 where  
 $P(y_n|y_{n-1}, \mathbf{x}_{1:N}) \propto \exp(\mathbf{w}^\top \phi(\mathbf{x}, n, y_n, y_{n-1}))$ 



Goal: Compute 
$$P(y_{1:N}|\mathbf{x}_{1:N}, \mathbf{w}) = \prod_n P(y_n|y_{n-1}, \mathbf{x}_{1:N})$$
 where  
 $P(y_n|y_{n-1}, \mathbf{x}_{1:N}) \propto \exp(\mathbf{w}^\top \phi(\mathbf{x}, n, y_n, y_{n-1}))$ 



Goal: Compute 
$$P(y_{1:N}|\mathbf{x}_{1:N}, \mathbf{w}) = \prod_n P(y_n|y_{n-1}, \mathbf{x}_{1:N})$$
 where  
 $P(y_n|y_{n-1}, \mathbf{x}_{1:N}) \propto \exp(\mathbf{w}^\top \phi(\mathbf{x}, n, y_n, y_{n-1}))$ 



Goal: Compute 
$$P(y_{1:N}|\mathbf{x}_{1:N}, \mathbf{w}) = \prod_n P(y_n|y_{n-1}, \mathbf{x}_{1:N})$$
 where  
 $P(y_n|y_{n-1}, \mathbf{x}_{1:N}) \propto \exp(\mathbf{w}^\top \phi(\mathbf{x}, n, y_n, y_{n-1}))$ 



## Overview

#### 1 Hidden Markov Models

- Representation
- Learning
- Inference

#### Conditional Models and Local Classifiers

- Conditional Models for Predicting Sequences
- Log-linear Models for Multiclass Classification
- Maximum Entropy Markov Models
   The Label Bias Problem

#### **B** Global Models

- Conditional Random Fields
- Structured Perceptron for Sequences
- Structural Support Vector Machine

### But ... Local Classifiers $\rightarrow$ Label Bias Problem

• Recall: the independence assumption ("Next-state" classifiers are locally normalized)

$$P(y_n|y_{n-1}, y_{n-2}, ..., \mathbf{x}_n, \mathbf{x}_{n-1}, \mathbf{x}_{n-2}, ...) = P(y_n|y_{n-1}, \mathbf{x}_n)$$



### But ... Local Classifiers $\rightarrow$ Label Bias Problem

• The robot wheels Fred round



- The path scores are the same
- Even if the word Fred is never observed as a verb in the data, it will be predicted as one
- The input Fred does not influence the output at all

- States with a single outgoing transition effectively ignore their input
  - States with lower-entropy next states are less influenced by observations
- Why?
  - Because each the next-state classifiers are locally normalized
  - If a state has fewer next states, each of those will get a higher probability mass
    - and hence preferred
- Surprisingly doesn't affect some tasks
  - E.g., POS tagging

- Conditional models
- Use rich features in the mode
- Possibly suffer from label bias problem

## Overview

#### Hidden Markov Models

- Representation
- Learning
- Inference

#### Conditional Models and Local Classifiers

- Conditional Models for Predicting Sequences
- Log-linear Models for Multiclass Classification
- Maximum Entropy Markov Models
  - The Label Bias Problem

#### **3** Global Models

#### • Conditional Random Fields

- Structured Perceptron for Sequences
- Structural Support Vector Machine

- Hidden Markov models
  - Pros: Decomposition of total probability with tractable
  - Cons: Doesn't allow use of features for representing inputs
    - Also, generative model is not really a downside, but we may get better performance with conditional models if we care only about predictions
- Local, Conditional Markov Models
  - Pros: Conditional model, allows features to be used
  - Cons: Label bias problem

- Train the predictor globally
  - Instead of training local decisions independently
- Normalize globally
  - Make each edge in the model undirected
  - Not associated with a probability, but just a "score"
- Recall the difference between local vs. global for multiclass

### HMM vs. A Local Model vs. A Global Model



- Each node is a random variable
- We observe some nodes and the rest are unobserved
- The goal: To characterize a probability distribution over the unobserved variables, given the observed



- Each node is a random variable
- We observe some nodes and the rest are unobserved
- The goal: To characterize a probability distribution over the unobserved variables, given the observed



• Each clique is associated with a score

Yangqiu Song (HKUST)

### Multi-class Classification as a Special Case

- Rewrite input features and weight vector
  - Define a feature vector for label *i* being associated to input **x**
  - Stack all weight vectors into an *nK*-dimensional vector

$$\phi(\mathbf{x},i) = \begin{bmatrix} \mathbf{0}_n \\ \vdots \\ \mathbf{x} \\ \vdots \\ \mathbf{0}_n \end{bmatrix}_{nK \times 1} \mathbf{w} = \begin{bmatrix} \mathbf{w}_1 \\ \mathbf{w}_2 \\ \vdots \\ \mathbf{w}_K \end{bmatrix}_{nK \times 1}$$

This is called the Kesler construction

# Conditional Random Field

- Each node is a random variable
- We observe some nodes and the rest are unobserved
- The goal: To characterize a probability distribution over the unobserved variables, given the observed



- Each clique is associated with a score
- We can put arbitrary features, as with local conditional models

Yangqiu Song (HKUST)

# Conditional Random Field

- Each node is a random variable
- We observe some nodes and the rest are unobserved
- The goal: To characterize a probability distribution over the unobserved variables, given the observed



Each clique factor is associated with a score
We can put arbitrary features, as with local conditional models

Yangqiu Song (HKUST)

Spring 2018 20 / 52

# Conditional Random Field

- Each node is a random variable
- We observe some nodes and the rest are unobserved
- The goal: To characterize a probability distribution over the unobserved variables, given the observed



• A different factorization: Recall decomposition of structures into parts. Same idea

Yangqiu Song (HKUST)

Spring 2018 21 / 52

### Conditional Random Field for Sequences



• The conditional probability is

$$P(y_{1:N}|\mathbf{x}_{1:N}) = \frac{1}{Z} \prod_{n} \exp(\mathbf{w}^{\top} \phi(\mathbf{x}_{1:N}, y_n, y_{n-1}))$$

where Z is a normalization constant

$$Z = \sum_{y_{1:N}} \prod_{n} \exp(\mathbf{w}^{\top} \phi(\mathbf{x}_{1:N}, y_n, y_{n-1}))$$

Yangqiu Song (HKUST)

- Input:  $\mathbf{x}_{1:N}$ , Output:  $y_{1:N}$ , both sequences (for now)
- Define a feature vector for the entire input and output sequence:  $\phi(\mathbf{x}_{1:N}, y_{1:N})$
- Define a giant log-linear model,  $P(y_{1:N}|\mathbf{x}_{1:N})$  parameterized by  $\mathbf{w}$

$$P(y_{1:N}|\mathbf{x}_{1:N}) = \frac{1}{Z} \prod_{n} \exp(\mathbf{w}^{\top} \phi(\mathbf{x}_{1:N}, y_n, y_{n-1})) \\ \propto \exp\left(\mathbf{w}^{\top} \sum_{n} \phi(\mathbf{x}_{1:N}, y_n, y_{n-1})\right)$$

- Just like any other log-linear model, except
  - Space of y is the set of all possible sequences of the correct length
  - Normalization constant sums over all sequences
  - In an MEMM, probabilities were locally normalized

• The feature function decomposes over the sequence

$$\phi(\mathbf{x}_{1:N}, y_{1:N}) = \sum_{n} \phi(\mathbf{x}_{1:N}, y_{n}, y_{n-1})$$



## **CRF** Prediction

Given

$$P(y_{1:N}|\mathbf{x}_{1:N}) = \frac{1}{Z} \exp\left(\mathbf{w}^{\top} \phi(\mathbf{x}_{1:N}, y_{1:N})\right)$$

• Goal: To predict most probable sequence  $y_{1:N}$  given an input  $\mathbf{x}_{1:N}$ 

$$\begin{array}{ll} \arg\max_{y_{1:N}} P(y_{1:N}|\mathbf{x}_{1:N}) &= \arg\max_{y_{1:N}} \exp\mathbf{w}^{\top}\phi(\mathbf{x}_{1:N},y_{1:N}) \\ &= \arg\max_{y_{1:N}}\mathbf{w}^{\top}\phi(\mathbf{x}_{1:N},y_{1:N}) \end{array}$$

• The score decomposes as  $\mathbf{w}^{\top} \phi(\mathbf{x}_{1:N}, y_{1:N}) = \mathbf{w}^{\top} \sum_{n} \phi(\mathbf{x}_{1:N}, y_{n}, y_{n-1})$ 

• So we do prediction via Viterbi (with sum instead of product)

$$\operatorname{score}_{1}(y_{1}) = \mathbf{w}^{\top} \sum_{n} \phi(\mathbf{x}_{1:N}, y_{1})$$

$$\operatorname{score}_{n}(y_{n}) = \max_{y_{n-1}}(\mathbf{w}^{\top}\phi(\mathbf{x}_{1:N}, y_{n}, y_{n-1}) + \operatorname{score}_{n-1}(y_{n-1}))$$

#### Input:

- Dataset with labeled sequences:  $\mathcal{D} = \{\mathbf{x}_{1:N_i}^{(i)}, \mathbf{y}_{1:N_i}^{(i)}\}$
- A definition of the feature function
- How do we train?
  - Maximize the (regularized) log-likelihood

$$\max_{\mathbf{w}} - \frac{\lambda}{2} \mathbf{w}^{\top} \mathbf{w} + \sum_{i} \log P(y_{1:N_i}^{(i)} | \mathbf{x}_{1:N_i}^{(i)}, \mathbf{w})$$

# Training with Inference

Given

$$\max_{\mathbf{w}} - \frac{\lambda}{2} \mathbf{w}^{\top} \mathbf{w} + \sum_{i} \log P(y_{1:N_i}^{(i)} | \mathbf{x}_{1:N_i}^{(i)}, \mathbf{w})$$

and

$$\mathsf{P}(y_{1:N}|\mathbf{x}_{1:N}) = \frac{1}{Z} \exp\left(\mathbf{w}^{\top} \phi(\mathbf{x}_{1:N}, y_{1:N})\right)$$

- Many methods for training
  - Numerical optimization
  - Can use a gradient or hessian based method
- Simple gradient ascent

$$\mathbf{w} \leftarrow \mathbf{w} + \sum_{i} \left( \phi(\mathbf{x}_{1:N_{i}}^{(i)}, y_{1:N_{i}}^{(i)}) - \sum_{\hat{y}_{1:N}} P(\hat{y}_{1:N} | \mathbf{x}_{1:N_{i}}^{(i)}, \mathbf{w}) \phi(\mathbf{x}_{1:N_{i}}^{(i)}, \hat{y}_{1:N}) \right)$$

- Red part: Training involves inference!
  - A different kind than what we have seen so far
  - Summing over all sequences is just like Viterbi (with summation instead of maximization)

Yangqiu Song (HKUST)

Learning for Text Analytics

Spring 2018 27 / 52

#### • An undirected graphical model

- Decompose the score over the structure into a collection of factors
- Each factor assigns a score to assignment of the random variables it is connected to
- Training and prediction
  - Final prediction via arg max<sub> $y_{1:N}$ </sub>  $\mathbf{w}^{\top} \phi(\mathbf{x}_{1:N}, y_{1:N})$
  - Train by maximum (regularized) likelihood
- Relation to other models
  - Effectively a linear classifier
  - A generalization of logistic regression to structures
  - An instance of Markov Random Field, with some random variables observed

## Overview

#### Hidden Markov Models

- Representation
- Learning
- Inference

#### Conditional Models and Local Classifiers

- Conditional Models for Predicting Sequences
- Log-linear Models for Multiclass Classification
- Maximum Entropy Markov Models
  - The Label Bias Problem

#### **3** Global Models

- Conditional Random Fields
- Structured Perceptron for Sequences
- Structural Support Vector Machine

# HMM is also a Linear Classifier



• Consider log  $P(\mathbf{x}_{1:N}, y_{1:N}) = \sum_{n} \log P(y_n | y_{n-1}) + P(\mathbf{x}_n | y_n)$ 



- log  $P(\mathbf{x}_{1:N}, y_{1:N})$  is linear scoring function  $\mathbf{w}^{\top} \phi(\mathbf{x}_{1:N}, y_{1:N})$ 
  - w: parameters of the model
  - $\phi(\mathbf{x}_{1:N}, y_{1:N})$ : properties of this output and the input

Yangqiu Song (HKUST)

Learning for Text Analytics

Spring 2018 30 / 52

- HMM is a linear classifier
  - Can we treat it as any linear classifier for training?
  - If so, we could add additional features that are global properties
    - As long as the output can be decomposed for easy inference
- The Viterbi algorithm calculates max  $\mathbf{w}^{\top} \phi(\mathbf{x}_{1:N}, y_{1:N})$ 
  - Viterbi only cares about scores to structures (not necessarily normalized)
- We could push the learning algorithm to train for un-normalized scores
  - $\bullet\,$  If we need normalization, we could always normalize by exponentiating and dividing by Z
  - That is, the learning algorithm can effectively just focus on the score of  $y_{1:N}$  for a particular  $\mathbf{x}_{1:N}$
  - Train a discriminative model instead of the generative one!

# Structured Perceptron Algorithm

- Given a training set  $\mathcal{D} = \{\mathbf{x}_{1:N}^{(i)}, y_{1:N}^{(i)}\}$
- (In practice, good to shuffle data before running)
- Initialize w = 0
- For epoch  $= 1, \ldots, T$ :
- (T is a hyperparameter to the algorithm)
  - For each training example  $(\mathbf{x}_{1:N}, y_{1:N}) \in \mathcal{D}$ 
    - Predict  $y'_{1:N} = \arg \max_{y'} \mathbf{w}^\top \phi(\mathbf{x}_{1:N}, y'_{1:N})$
    - (Inference in the training loop)
    - If  $y'_{1:N} \neq y_{1:N}$ , update  $\mathbf{w} \leftarrow \mathbf{w} + \eta(\phi(\mathbf{x}_{1:N}, y_{1:N}) \phi(\mathbf{x}_{1:N}, y'_{1:N}))$
- Return w
- Prediction:  $\arg \max_{y} = \mathbf{w}^{\top} \phi(\mathbf{x}_{1:N}, y_{1:N})$

- Mistake bound for separable data, just like perceptron
- In practice, use averaging for better generalization
  - Initialize **a** = 0
  - $\bullet\,$  After each step, whether there is an update or not,  $\mathbf{a} \leftarrow \mathbf{a} + \mathbf{w}\,$ 
    - $\bullet\,$  Note, we still check for mistake using w not a
  - Return  ${\bf a}$  at the end instead of  ${\bf w}$
- Global update
  - One weight vector for entire sequence (not for each position)

• Stochastic gradient descent update for CRF: expectation vs max

$$\mathbf{w} \leftarrow \mathbf{w} + \eta(\phi(\mathbf{x}_{1:N_i}, y_{1:N_i}) - \sum_{\hat{y}_{1:N}} P(\hat{y}_{1:N} | \mathbf{x}_{1:N_i}, \mathbf{w}) \phi(\mathbf{x}_{1:N_i}, \hat{y}_{1:N}))$$

$$\mathbf{w} \leftarrow \mathbf{w} + \eta(\phi(\mathbf{x}_{1:N_i}, y_{1:N_i}) - \mathbb{E}_{P(\hat{y}_{1:N} | \mathbf{x}_{1:N_i}, \mathbf{w})}[\phi(\mathbf{x}_{1:N_i}, \hat{y}_{1:N})])$$

Structured perceptron

$$\mathbf{w} \leftarrow \mathbf{w} + \eta(\phi(\mathbf{x}_{1:N}, y_{1:N}) - \phi(\mathbf{x}_{1:N}, y'_{1:N}))$$

where  $y'_{1:N} = \arg \max_{y'} \mathbf{w}^\top \phi(\mathbf{x}_{1:N}, y'_{1:N})$ 

• Caveat: Adding regularization will change the CRF update, averaging changes the perceptron update

# The lay of the land

HMM: A generative model, assigns probabilities to sequences

Structured Perceptron/Structured SVM	Conditional Random field
Hidden Markov Models are actually just linear classifiers	Model probabilities via exponential functions. Gives us the log-linear representation
Dont really care whether we are predicting probabilities. We are assigning scores to a full output for a given input (like multiclass)	Log-probabilities for sequences for a given input
Generalize algorithms for linear classifiers. Sophisticated models that can use arbitrary features	Learn by maximizing likelihood. Sophisticated models that can use arbitrary features

Applicable beyond sequences. Eventually, similar objective minimized with different loss functions

Yangqiu Song (HKUST)

# Overview

#### Hidden Markov Models

- Representation
- Learning
- Inference

#### Conditional Models and Local Classifiers

- Conditional Models for Predicting Sequences
- Log-linear Models for Multiclass Classification
- Maximum Entropy Markov Models
  - The Label Bias Problem

#### **3** Global Models

- Conditional Random Fields
- Structured Perceptron for Sequences
- Structural Support Vector Machine

# Recall: Binary SVM

# • Functional margin $\gamma^{(i)} = y^{(i)} (\mathbf{w}^{\top} \mathbf{x}^{(i)} + b)$

- If  $y^{(i)} = 1$ , we want  $\mathbf{w}^{\top} \mathbf{x}^{(i)} + b$  to be a large positive number
- If  $y^{(i)} = -1$ , we want  $\mathbf{w}^{\top} \mathbf{x}^{(i)} + b$  to be a large negative number
- If  $y^{(i)}(\mathbf{w}^{\top}\mathbf{x}^{(i)} + b) > 0$ , then our prediction on this example is correct
- If we replace  $(\mathbf{w}, b)$  with  $(2\mathbf{w}, 2b)$ , it does not change the decision
  - So it makes more sense to impose some sort of normalization condition such as that  $||{\bf w}||_2=1$
- The functional margin of a dataset is defined as the smallest of the function margins of individual exmaples

$$\gamma = \min_{i} \gamma^{(i)}$$

# Recall: Binary SVM (Cont'd)



- Geometric margin
  - The points lie on the decision boundary satisfy the equation  $\mathbf{w}^{\top}\mathbf{x} + b = 0$
  - The points lie on the margin satisfy  $\mathbf{w}^{\top}(\mathbf{x} \gamma \frac{\mathbf{w}}{||\mathbf{w}||}) + b = 0$
  - $\bullet\,$  Solving for  $\gamma$  yields

$$\gamma = \frac{\mathbf{w}^{\top} \mathbf{x} + b}{||\mathbf{w}||} = \left(\frac{\mathbf{w}}{||\mathbf{w}||}\right)^{\top} \mathbf{x} + \frac{b}{||\mathbf{w}||}$$

• If  $||\mathbf{w}|| = 1$ , then the functional margin equals the geometric margin. The geometric margin is invariant to margin of the parameters

- The geometric margin is invariant to rescaling of the parameters
  - We replace w with 2w and b with 2b, then the geometric margin does not change

# Recall: Binary SVM (Cont'd)

• To maximize geometric margin, we can pose the following optimization problem

$$\begin{array}{ll} \max_{\gamma, \mathbf{w}, b} & \gamma \\ s.t. & y^{(i)} (\mathbf{w}^\top \mathbf{x}^{(i)} + b) \geq \gamma \\ & ||\mathbf{w}|| = 1 \end{array}$$

- The constraint  $||\mathbf{w}|| = 1$  is non-convex. So we convert it as following nicer one  $\max_{\gamma,\mathbf{w},b} \quad \frac{\gamma}{||\mathbf{w}||}$ s.t.  $\gamma^{(i)}(\mathbf{w}^{\top}\mathbf{x}^{(i)} + b) \geq \gamma$
- Note that we can add an arbitrary scaling constraint on **w** and *b* without changing anything
- We introduce the scaling constraint that the functional margin of **w** and *b* with respect to the training set must be 1:

$$\begin{array}{ll} \min_{\gamma, \mathbf{w}, b} & \frac{1}{2} ||\mathbf{w}||^2 \\ s.t. & y^{(i)}(\mathbf{w}^\top \mathbf{x}^{(i)} + b) \geq 1 \end{array}$$

where  $||\mathbf{w}||^2$  makes it easier (convex) for optimization and remain the same optimization problem

Yangqiu Song (HKUST)

Spring 2018 39 / 52

#### SVM optimizes

$$\begin{array}{l} \min_{\gamma, \mathbf{w}, b} \quad \frac{1}{2} ||\mathbf{w}||^2 \\ s.t. \qquad y^{(i)} (\mathbf{w}^\top \mathbf{x}^{(i)} + b) \geq 1 \end{array}$$

• With Lagrange multiplier, a equivalent formulation is

$$\frac{1}{2}||\mathbf{w}||^2 + C\sum_i \max(0, 1 - y^{(i)}(\mathbf{w}^\top \mathbf{x}^{(i)} + b))$$

where  $\max(0, 1 - y(\mathbf{w}^{\top}\mathbf{x} + b))$  is called hinge loss • As a comparison,  $\max(0, -y(\mathbf{w}^{\top}\mathbf{x} + b))$  is called perceptron loss

### Compare losses



### Multiclass SVM in the Separable Case

• For dataset with K classes, our optimization problem is

$$\begin{array}{ll} \min_{\gamma, \mathbf{w}} & \frac{1}{2} \sum_{k}^{K} ||\mathbf{w}_{k}||^{2} \\ s.t. & \mathbf{w}_{y^{(i)}}^{\top} \mathbf{x}^{(i)} - \mathbf{w}_{k}^{\top} \mathbf{x}^{(i)} \geq 1 \quad \forall i, k \neq y^{(i)} \end{array}$$

The score for the true label is higher than the score for any other label by  $1 \ensuremath{$ 

• With Kesler construction, we have

$$\begin{array}{ll} \min_{\gamma, \mathbf{w}} & \frac{1}{2} ||\mathbf{w}||^2 \\ s.t. & \mathbf{w}^\top (\phi(\mathbf{x}^{(i)}, y^{(i)}) - \phi(\mathbf{x}^{(i)}, k)) \geq 1 \ \forall i, k \neq y^{(i)} \end{array}$$

$$\phi(\mathbf{x},i) = \begin{bmatrix} \mathbf{0}_n \\ \vdots \\ \mathbf{x} \\ \vdots \\ \mathbf{0}_n \end{bmatrix}_{nK \times 1} \mathbf{w} = \begin{bmatrix} \mathbf{w}_1 \\ \mathbf{w}_2 \\ \vdots \\ \mathbf{w}_K \end{bmatrix}_{nK \times 1}$$

Yangqiu Song (HKUST)

## Structural SVM: A First Attempt

- Suppose we have some definition of a structure (a factor graph)
  - And feature definitions for each "part" p as  $\phi_p(\mathbf{x}, \mathbf{y}_p)$
  - Remember: we can talk about the feature vector for the entire structure
    - Features sum over the parts

$$\phi(\mathbf{x},\mathbf{y}) = \sum_{p \in parts(\mathbf{x})} \phi_p(\mathbf{x},\mathbf{y}_p)$$



• We also have a dataset  $\{\mathbf{x}^{(i)}, \mathbf{y}^{(i)}\}$ 

## Structural SVM: A First Attempt



What do we want from training (following the multiclass idea)?

- For each example,
  - The annotated structure  $\mathbf{y}$  gets the highest score among all structures
  - The structure y gets a score of at least one more than any other

$$\mathbf{w}^{ op}\phi(\mathbf{x},\mathbf{y})\geq\mathbf{w}^{ op}\phi(\mathbf{x},\mathbf{y}')+1 \;\; orall \mathbf{y}'
eq \mathbf{y}$$

Goal

maximize margin

s.t. score for gold structure  $\geq$  score for other structure + 1 for every training example

Corresponding to

$$\begin{array}{ll} \min_{\mathbf{w}} & \frac{1}{2} \mathbf{w}^\top \mathbf{w} \\ s.t. & \mathbf{w}^\top \phi(\mathbf{x}^{(i)}, \mathbf{y}^{(i)}) \geq \mathbf{w}^\top \phi(\mathbf{x}^{(i)}, \mathbf{y}') + 1 \ \forall \mathbf{y}' \neq \mathbf{y}^{(i)} \end{array}$$

Problem:



• Structure B has is more wrong, but this formulation will be happy if both A ans B are scored one less than gold!

Yangqiu Song (HKUST)

Learning for Text Analytics

Spring 2018 45 / 52

# Structural SVM: Second Attempt

First attempt

$$\begin{array}{ll} \min_{\mathbf{w}} & \frac{1}{2} \mathbf{w}^{\top} \mathbf{w} \\ s.t. & \mathbf{w}^{\top} \phi(\mathbf{x}^{(i)}, \mathbf{y}^{(i)}) \geq \mathbf{w}^{\top} \phi(\mathbf{x}^{(i)}, \mathbf{y}') + 1 \ \forall \mathbf{y}' \neq \mathbf{y}^{(i)} \end{array}$$

• Introduce Hamming distance between structures

$$\begin{array}{ll} \min_{\mathbf{w}} & \frac{1}{2} \mathbf{w}^\top \mathbf{w} \\ s.t. & \mathbf{w}^\top \phi(\mathbf{x}^{(i)}, \mathbf{y}^{(i)}) \geq \mathbf{w}^\top \phi(\mathbf{x}^{(i)}, \mathbf{y}') + \Delta(\mathbf{y}', \mathbf{y}^{(i)}) \ \forall \mathbf{y}' \neq \mathbf{y}^{(i)} \end{array}$$

where  $\Delta(\mathbf{y}', \mathbf{y}^{(i)})$  is defined as Hamming distance between structures • The Hamming distance between:

- "karolin" and "kathrin" is 3.
- "karolin" and "kerstin" is 3.
- 1011101 and 1001001 is 2.
- 2173896 and 2233796 is 3.
- Intuition
  - Structures that are more different from the true structure should be scored much lower

Yangqiu Song (HKUST)

Spring 2018 46 / 52

# Structural SVM: Third Attempt

- Problem? What if the data is not separable?
  - What if these constraints are not satisfied for any **w** for a given dataset?
- Slack variable for each example  $\xi_i$ , must be positive

$$\min_{\mathbf{w},\xi_i} \quad \frac{1}{2} \mathbf{w}^\top \mathbf{w} + C \sum_i \xi_i \\ s.t. \quad \mathbf{w}^\top \phi(\mathbf{x}^{(i)}, \mathbf{y}^{(i)}) \ge \mathbf{w}^\top \phi(\mathbf{x}^{(i)}, \mathbf{y}') + \Delta(\mathbf{y}', \mathbf{y}^{(i)}) - \xi_i \quad \forall \mathbf{y}' \neq \mathbf{y}^{(i)} \\ \xi_i \ge 0$$

- Slack variables allow some examples to be misclassified
- Minimizing the slack forces this to happen as few times as possible
- An equivalent formulation

$$\min_{\mathbf{w}} \frac{1}{2} \mathbf{w}^\top \mathbf{w} + C \sum_i \max_{\mathbf{y}'} \left( \mathbf{w}^\top \phi(\mathbf{x}^{(i)}, \mathbf{y}') + \Delta(\mathbf{y}', \mathbf{y}^{(i)}) - \mathbf{w}^\top \phi(\mathbf{x}^{(i)}, \mathbf{y}^{(i)}) \right)$$

- Other slightly different formulations exist
  - Generally same principle
- Multiclass is a special case of structure
  - Structural SVM strictly generalizes multiclass SVM
- Can be seen as minimizing structured version of hinge loss
- Learning as optimization

- Collect some annotated data. More is generally better
- Pick a hypothesis class (also called model)
  - Decide how the score decomposes over the parts of the output
- Choose a loss function
  - Decide on how to penalize incorrect decisions
- Learning = minimize empirical risk + regularizer
  - Typically an optimization procedure needed here

### Structured Classifiers: Different Learning Objectives

Structural SVM

$$\min_{\mathbf{w}} \frac{1}{2} \mathbf{w}^\top \mathbf{w} + C \sum_{i} \max_{\mathbf{y}'} \left( \mathbf{w}^\top \phi(\mathbf{x}^{(i)}, \mathbf{y}') + \Delta(\mathbf{y}', \mathbf{y}^{(i)}) - \mathbf{w}^\top \phi(\mathbf{x}^{(i)}, \mathbf{y}^{(i)}) \right)$$

• CRF:  $\max_{\mathbf{w}} -\frac{\lambda}{2}\mathbf{w}^{\top}\mathbf{w} + \sum_{i} \log P(y_{1:N_{i}}^{(i)}|\mathbf{x}_{1:N_{i}}^{(i)},\mathbf{w})$ 

$$\max_{\mathbf{w}} -\frac{\lambda}{2} \mathbf{w}^{\top} \mathbf{w} + \sum_{i} \left( \mathbf{w}^{\top} \phi(\mathbf{x}_{1:N}^{(i)}, y_{1:N}^{(i)}) - \log \sum_{y_{1:N}^{\prime}} \left( \mathbf{w}^{\top} \phi(\mathbf{x}_{1:N}^{(i)}, y_{1:N}^{\prime}) \right) \right)$$

$$\min_{\mathbf{w}} \frac{\lambda}{2} \mathbf{w}^{\top} \mathbf{w} + \sum_{i} \left( \log \sum_{y'_{1:N}} \left( \mathbf{w}^{\top} \phi(\mathbf{x}_{1:N}^{(i)}, y'_{1:N}) \right) - \mathbf{w}^{\top} \phi(\mathbf{x}_{1:N}^{(i)}, y_{1:N}^{(i)}) \right)$$

• Structured Perceptron

$$\min_{\mathbf{w}} C \sum_{i} \max_{\mathbf{y}'} \left( \mathbf{w}^{\top} \phi(\mathbf{x}^{(i)}, \mathbf{y}') - \mathbf{w}^{\top} \phi(\mathbf{x}^{(i)}, \mathbf{y}^{(i)}) \right)$$

Yangqiu Song (HKUST)

# Further Reading

- Ng (2017). CS229 Lecture notes: Support Vector Machines. http://cs229.stanford.edu/notes/cs229-notes3.pdf
- Punyakanok and Roth (2000). The Use of Classifiers in Sequential Inference.
- McCallum et al. (2000). Maximum Entropy Markov Models for Information Extraction and Segmentation.
- Lafferty et al. (2001). Conditional Random Fields: Probabilistic Models for Segmenting and Labeling Sequence Data.
- Taskar et al. (2003). Max-Margin Markov Networks.
- Collins (2011). Course notes for COMS w4705: Log-linear models, MEMMs, and CRFs, 2011. http://www.cs.columbia.edu/~mcollins/crf.pdf
- Smith (2004). Log-linear models, 2004. https://homes.cs. washington.edu/~nasmith/papers/smith.tut04.pdf

(日) (同) (三) (三)

- Collins, M. (2011). Log-linear models, MEMMs, and CRFs. Technical report, Columbia University.
- Lafferty, J. D., McCallum, A., and Pereira, F. C. N. (2001). Conditional random fields: Probabilistic models for segmenting and labeling sequence data. In *ICML*, pages 282–289.
- McCallum, A., Freitag, D., and Pereira, F. C. N. (2000). Maximum entropy markov models for information extraction and segmentation. In *ICML*, pages 591–598.
- Ng, A. (2017). Cs229 lecture notes: Support vector machines. Technical report, Stanford University.
- Punyakanok, V. and Roth, D. (2000). The use of classifiers in sequential inference. In NIPS.
- Smith, N. A. (2004). Log-linear models. Technical report, University of Washington.
- Taskar, B., Guestrin, C., and Koller, D. (2003). Max-margin markov networks. In *NIPS*, pages 25–32.

(日) (周) (三) (三)