

CS 446/ECE 449: Machine Learning

A. G. Schwing

University of Illinois at Urbana-Champaign, 2020

L24: Structured Latent Variable Models

Goals of this lecture

- Getting to know Structured Latent Variable Models
- Learning about Hidden Markov Models (HMMs)

Reading material:

- K. Murphy; Machine Learning: A Probabilistic Perspective; Chapter 17

Recap:

Generative:

$$\ln p_{\theta}(x^{(i)}) = \ln \sum_{\mathbf{z}} p_{\theta}(x^{(i)}, \mathbf{z})$$

Discriminative:

$$\ln p_{\mathbf{w}}(\mathbf{y}|x^{(i)})$$

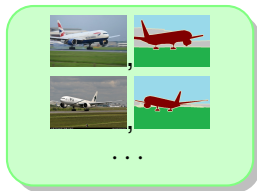
Observations:

- **z never observed** (assignment to cluster center in GMM)
- **z very simple categorical variable**
- **y** always completely observed

Questions: structure in **z** or parts of **y** unobserved

Recap: Learning with full observations

- Training set of data pairs $(x^{(i)}, \mathbf{y}^{(i)})$



- Inference

$$\arg \max_{\hat{\mathbf{y}}} F(\mathbf{w}, x^{(i)}, \hat{\mathbf{y}}) = \arg \max_{\hat{\mathbf{y}}} \sum_r f_r(\mathbf{w}, x^{(i)}, \hat{\mathbf{y}}_r)$$

- Learning target

$$\forall \hat{\mathbf{y}} \quad F(\mathbf{w}, x^{(i)}, \hat{\mathbf{y}}) \leq F(\mathbf{w}, x^{(i)}, \mathbf{y}^{(i)})$$

$$\max_{\hat{\mathbf{y}}} F(\mathbf{w}, x^{(i)}, \hat{\mathbf{y}}) \leq F(\mathbf{w}, x^{(i)}, \mathbf{y}^{(i)})$$

Recap: Learning with full observations

Hinge loss: penalize whenever maximum is within a margin $L(\hat{\mathbf{y}}, \mathbf{y}^{(i)})$
of the data $(\mathbf{x}^{(i)}, \mathbf{y}^{(i)})$ score:

$$\max_{\hat{\mathbf{y}}} \left(F(\mathbf{w}, \mathbf{x}^{(i)}, \hat{\mathbf{y}}) + L(\hat{\mathbf{y}}, \mathbf{y}^{(i)}) \right) \geq F(\mathbf{w}, \mathbf{x}^{(i)}, \mathbf{y}^{(i)})$$

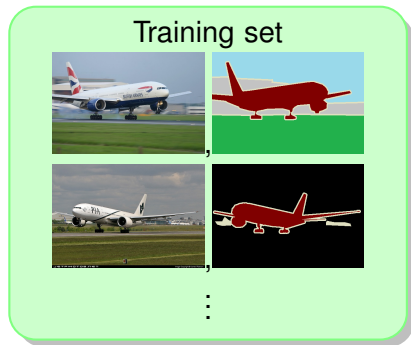
$$\min_{\mathbf{w}} \frac{C}{2} \|\mathbf{w}\|_2^2 + \sum_{i \in \mathcal{D}} \left(\underbrace{\max_{\hat{\mathbf{y}}} \left(F(\mathbf{w}, \mathbf{x}^{(i)}, \hat{\mathbf{y}}) + L(\hat{\mathbf{y}}, \mathbf{y}^{(i)}) \right)}_{\text{Loss-augmented inference}} - F(\mathbf{w}, \mathbf{x}^{(i)}, \mathbf{y}^{(i)}) \right)$$

Parts of \mathbf{y} unobserved/not annotated:

Fully labeled:



Weakly labeled:



Complete Data:

$$\mathbf{y} = (\mathbf{s}, \mathbf{z})$$

Parts of \mathbf{y} unobserved/not annotated:

- Complete Data:

$$\mathbf{y} = (\mathbf{s}, \mathbf{z})$$

- Weakly labeled hinge loss: penalize whenever best overall prediction exceeds best prediction with annotation being clamped

$$\max_{\hat{\mathbf{s}}, \hat{\mathbf{z}}} F(\mathbf{w}, x^{(i)}, \hat{\mathbf{s}}, \hat{\mathbf{z}}) \geq \max_{\hat{\mathbf{z}}} F(\mathbf{w}, x^{(i)}, \mathbf{s}^{(i)}, \hat{\mathbf{z}})$$

- Latent SSVM (LSSVM) [Yu and Joachims 2009]:

$$\frac{C}{2} \|\mathbf{w}\|_2^2 + \sum_{i \in \mathcal{D}} \max_{\hat{\mathbf{s}}, \hat{\mathbf{z}}} F(\mathbf{w}, x^{(i)}, \hat{\mathbf{s}}, \hat{\mathbf{z}}) - \sum_{i \in \mathcal{D}} \max_{\hat{\mathbf{z}}} F(\mathbf{w}, x^{(i)}, \mathbf{s}^{(i)}, \hat{\mathbf{z}})$$

Structured Prediction with Latent Variables

$$\min_{\mathbf{w}} \frac{C}{2} \|\mathbf{w}\|_2^2 + \sum_i \epsilon \ln \sum_{\hat{\mathbf{s}}, \hat{\mathbf{z}}} \exp \frac{F(\mathbf{w}, x^{(i)}, \hat{\mathbf{s}}, \hat{\mathbf{z}}) + L_i(\hat{\mathbf{s}}, \hat{\mathbf{z}})}{\epsilon} - \\ - \sum_i \epsilon \ln \sum_{\hat{\mathbf{z}}} \exp \frac{F(\mathbf{w}, x^{(i)}, \mathbf{s}^{(i)}, \hat{\mathbf{z}}) + L_i^c(\mathbf{s}^{(i)}, \hat{\mathbf{z}})}{\epsilon}$$

- Soft-max function

$$\epsilon \ln \sum_{\hat{\mathbf{s}}, \hat{\mathbf{z}}} \exp \frac{F(\mathbf{w}, x^{(i)}, \hat{\mathbf{s}}, \hat{\mathbf{z}})}{\epsilon} \xrightarrow{\epsilon \rightarrow 0} \max_{\hat{\mathbf{s}}, \hat{\mathbf{z}}} F(\mathbf{w}, x^{(i)}, \hat{\mathbf{s}}, \hat{\mathbf{z}})$$

- LSSVM & hidden CRFs (HCRFs) [Lafferty et al. 2001]

$\epsilon = 0$	hinge-loss	max-margin	LSSVM
$\epsilon = 1$	log-loss	max-likelihood	HCRF

- Margin L

Recall: variational expression for partition function

$$\epsilon \ln \sum_{\mathbf{z}} \exp \frac{F(\mathbf{w}, \mathbf{x}^{(i)}, \mathbf{s}^{(i)}, \mathbf{z})}{\epsilon} = \max_{q(\mathbf{z}) \in \Delta} \left(\sum_{\mathbf{z}} q(\mathbf{z}) F(\mathbf{w}, \mathbf{x}^{(i)}, \mathbf{s}^{(i)}, \mathbf{z}) + \epsilon H(q(\mathbf{z})) \right)$$

- Similarity to structured inference
- Similar algorithms can be employed

$$\min_{\mathbf{w}} \quad \frac{C}{2} \|\mathbf{w}\|_2^2 + \sum_i \left(\epsilon \ln \sum_{\hat{\mathbf{s}}, \hat{\mathbf{z}}} \exp \frac{F(\mathbf{w}, \mathbf{x}^{(i)}, \hat{\mathbf{s}}, \hat{\mathbf{z}})}{\epsilon} + L_i(\hat{\mathbf{s}}, \hat{\mathbf{z}}) \right. \\ \left. - \max_{q_i(\mathbf{z}) \in \Delta} \left(\sum_{\mathbf{z}} q_i(\mathbf{z}) \left(F(\mathbf{w}, \mathbf{x}^{(i)}, \mathbf{s}^{(i)}, \mathbf{z}) + L_i^c(\mathbf{s}^{(i)}, \mathbf{z}) \right) + \epsilon H(q_i(\mathbf{z})) \right) \right)$$

Alternating optimization between q and \mathbf{w}

Algorithmic structure

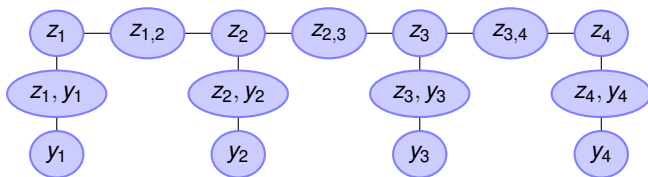
HCRF, LSSVM:

```
repeat
  repeat
    latent variable prediction to obtain  $q(\mathbf{z})$ 
  until convergence
  repeat
    update  $w$ 
  until convergence
until convergence
```

Irrespective of whether \mathbf{z} sometimes observed or never observed

Example of model with latent variables:

Hidden Markov Model



Algorithm:

- Inferring $q(\mathbf{z})$: forward/backward pass on a chain graph
- Updating \mathbf{w} : loss-augmented inference on a tree graph

Applications of HMMs:

- Time-series data in general
- Video data
- Language models

Directly applicable to generative modeling with structured \mathbf{z}

Quiz:

- Variational expression of partition function?
- Structured prediction with latent variables?
- Algorithmic structure of general framework?

Important topics of this lecture

- General framework for structured prediction with latent variables
- Similarity between generative and discriminative techniques

What's next

- Variational Auto-Encoders