

CS 446/ECE 449: Machine Learning

A. G. Schwing

University of Illinois at Urbana-Champaign, 2020

L14: Conditional Random Fields, Structured SVMs, Deep Structured Nets

Goals of this lecture

- Learning of structured distributions

Reading material:

- D. Koller and N. Friedman; Probabilistic Graphical Models: Principles and Techniques;

Recap: General framework for learning:

$$\min_{\mathbf{w}} \frac{C}{2} \|\mathbf{w}\|_2^2 + \sum_{i \in \mathcal{D}} \epsilon \ln \sum_{\hat{y}} \exp \frac{L(y^{(i)}, \hat{y}) + F(\mathbf{w}, x^{(i)}, \hat{y})}{\epsilon} - F(\mathbf{w}, x^{(i)}, y^{(i)})$$

Attention:

- Scoring function $F(\mathbf{w}, x, y)$
- Loss function (log-loss, hinge-loss)
- Taskloss L

How to get to

- Logistic regression
- Binary SVM
- Multiclass regression
- Multiclass SVM
- Deep Learning

Recap: Inference (how to find the highest scoring configuration):

$$y^* = \arg \max_{\hat{y}} F(\mathbf{w}, x, y)$$

Structured Prediction

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

Efficiency and accuracy of inference algorithms is problem dependent:

- Exhaustive search
- Dynamic programming
- Integer linear program
- Linear programming relaxation
- Message passing
- Graph-cut

Question for today:

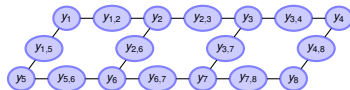
How to put more complex objects $\mathbf{y}^{(i)}$ (as opposed to $y^{(i)}$) into our learning formulation?

$$\min_{\mathbf{w}} \frac{C}{2} \|\mathbf{w}\|_2^2 + \sum_{i \in \mathcal{D}} \epsilon \ln \sum_{\hat{y}} \exp \frac{L(y^{(i)}, \hat{y}) + F(\mathbf{w}, x^{(i)}, \hat{y})}{\epsilon} - F(\mathbf{w}, x^{(i)}, y^{(i)})$$

Let's start with our linear multi-class formulation:

$$\min_{\mathbf{w}} \frac{C}{2} \|\mathbf{w}\|_2^2 + \sum_{i \in \mathcal{D}} \epsilon \ln \sum_{\hat{y}} \exp \frac{L(y^{(i)}, \hat{y}) + \mathbf{w}^T \psi(x^{(i)}, \hat{y})}{\epsilon} - \mathbf{w}^T \psi(x^{(i)}, y^{(i)})$$

Example: Semantic segmentation



- A classifier provides evidence for every pixel (unary potentials)
- Smoothness is encouraged via correlations (pairwise potentials)

How much weight should we give to each of the two cues?

Feature vector for structured objects:

$$f(x^{(i)}, \mathbf{y}) = \begin{bmatrix} \frac{f_1(x^{(i)}, \mathbf{y})}{f_M(x^{(i)}, \mathbf{y})} \\ \vdots \\ \frac{f_M(x^{(i)}, \mathbf{y})}{f_M(x^{(i)}, \mathbf{y})} \end{bmatrix} = \begin{bmatrix} \frac{\sum_r f_{1,r}(x^{(i)}, \mathbf{y}_r)}{\sum_r f_{M,r}(x^{(i)}, \mathbf{y}_r)} \\ \vdots \\ \frac{\sum_r f_{M,r}(x^{(i)}, \mathbf{y}_r)}{\sum_r f_{M,r}(x^{(i)}, \mathbf{y}_r)} \end{bmatrix} = \begin{bmatrix} y_1 & y_2 & y_3 & y_4 \\ \vdots & \vdots & \vdots & \vdots \\ y_{1,3} & y_{1,4} & y_{2,3} & y_{2,4} \end{bmatrix}$$

- Every vector element f_m is a graphical model
- Learn to weight parts of the combined graphical model

$$F(\mathbf{w}, x^{(i)}, \mathbf{y}) = \mathbf{w}^\top f(x^{(i)}, \mathbf{y}) = \sum_{m=1}^M w_m \underbrace{\left(\sum_r f_{m,r}(x^{(i)}, \mathbf{y}_r) \right)}_{\text{individual models}} = \sum_r \underbrace{\hat{f}_r(\mathbf{w}, x^{(i)}, \mathbf{y}_r)}_{\text{combined model}}$$

Recall: How did we derive multi-class SVM objective?

Learning target with groundtruth $\mathbf{y}^{(i)}$:

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

How many constraints are there?

Alternatively:

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

Hinge loss:

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

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

- Max-Margin Markov Network [Taskar et al. 2003]
- Structured Support Vector Machine [Tsochantaridis et al. 2004]

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

How to optimize this?

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

E.g., with (sub-)gradient descent:

Iterate:

- 1 Loss-augmented inference:

$$\arg \max_{\hat{\mathbf{y}}} \left(\mathbf{w}^\top f(x^{(i)}, \hat{\mathbf{y}}) + L(\hat{\mathbf{y}}, \mathbf{y}^{(i)}) \right)$$

- 2 Perform gradient step:

$$\mathbf{w} \leftarrow \mathbf{w} - \alpha \nabla_{\mathbf{w}} L(\mathbf{w})$$

How complicated is this?

Solve one structured prediction task per sample per iteration.

Structured Learning

$$\min_{\mathbf{w}} \frac{C}{2} \|\mathbf{w}\|_2^2 + \sum_{i \in \mathcal{D}} \epsilon \ln \sum_{\hat{\mathbf{y}}} \exp \frac{L(\mathbf{y}^{(i)}, \hat{\mathbf{y}}) + \mathbf{w}^\top f(x, \hat{\mathbf{y}})}{\epsilon} - \mathbf{w}^\top f(x^{(i)}, \mathbf{y}^{(i)})$$

Subsumes:

- linear structured SVMs
- linear conditional random fields (we condition on data x)
- binary/multi-class SVM
- logistic/multi-class regression

What changed compared to our previous formulation?

Structured Learning

$$\min_{\mathbf{w}} \frac{C}{2} \|\mathbf{w}\|_2^2 + \sum_{i \in \mathcal{D}} \epsilon \ln \sum_{\hat{\mathbf{y}}} \exp \frac{L(\mathbf{y}^{(i)}, \hat{\mathbf{y}}) + \mathbf{w}^\top f(x, \hat{\mathbf{y}})}{\epsilon} - \mathbf{w}^\top f(x^{(i)}, \mathbf{y}^{(i)})$$

What's the main limitation?

What happens if we make our score function non-linear in its parameters \mathbf{w} ?

$$\min_{\mathbf{w}} \frac{C}{2} \|\mathbf{w}\|_2^2 + \sum_{i \in \mathcal{D}} \epsilon \ln \sum_{\hat{\mathbf{y}}} \exp \frac{L(\mathbf{y}^{(i)}, \hat{\mathbf{y}}) + F(\mathbf{w}, x^{(i)}, \hat{\mathbf{y}})}{\epsilon} - F(\mathbf{w}, x^{(i)}, \mathbf{y}^{(i)})$$

$F(\mathbf{w}, x^{(i)}, \hat{\mathbf{y}})$ is a score function represented by a computation graph,
e.g., a deep net

How to optimize this program:

$$\min_{\mathbf{w}} \frac{C}{2} \|\mathbf{w}\|_2^2 + \sum_{i \in \mathcal{D}} \epsilon \ln \sum_{\hat{\mathbf{y}}} \exp \frac{\cancel{L(\mathbf{y}^{(i)}, \hat{\mathbf{y}})} + F(\mathbf{w}, x, \hat{\mathbf{y}})}{\epsilon} - F(\mathbf{w}, x^{(i)}, \mathbf{y}^{(i)})$$

Let's try gradient descent: (for simplicity $\epsilon = 1$)

$$p(\hat{\mathbf{y}}; x, \mathbf{w}) = \frac{\exp F(\mathbf{w}, x, \hat{\mathbf{y}})}{\sum_{\tilde{\mathbf{y}}} \exp F(\mathbf{w}, x, \tilde{\mathbf{y}})}$$

$$\begin{aligned} & C\mathbf{w} + \sum_{i \in \mathcal{D}} \left(\sum_{\hat{\mathbf{y}}} p(\hat{\mathbf{y}}; x^{(i)}, \mathbf{w}) \nabla_{\mathbf{w}} F(\mathbf{w}, x^{(i)}, \hat{\mathbf{y}}) - \nabla_{\mathbf{w}} F(\mathbf{w}, x^{(i)}, \mathbf{y}^{(i)}) \right) \\ = & C\mathbf{w} + \sum_{i \in \mathcal{D}, \hat{\mathbf{y}}} \left(p(\hat{\mathbf{y}}; x^{(i)}, \mathbf{w}) - \underbrace{\delta(\hat{\mathbf{y}} = \mathbf{y}^{(i)})}_{p_{\mathbf{y}^{(i)}}^{\text{GT}}(\hat{\mathbf{y}})} \right) \nabla_{\mathbf{w}} F(\mathbf{w}, x^{(i)}, \hat{\mathbf{y}}) \end{aligned}$$

Algorithm

Repeat until stopping criteria

- ➊ Forward pass to compute $F(\mathbf{w}, x^{(i)}, \hat{\mathbf{y}})$
- ➋ Compute $p(\hat{\mathbf{y}}; x, \mathbf{w})$ via soft-max
- ➌ Backward pass via chain rule to obtain gradient
- ➍ Update parameters \mathbf{w}

Where are the challenges?

- How do we represent $F(\mathbf{w}, x, \hat{\mathbf{y}})$ if output space \mathcal{Y} is large?

Output space size of some applications:

Tag prediction Segmentation



$$|\mathcal{Y}| = 2^{\#tags}$$

$$|\mathcal{Y}| = C^{\#pixels}$$

Observation:

- Interest in jointly predicting multiple variables $\mathbf{y} = (y_1, \dots, y_D)$

Solution:

- assume scoring function decomposes additively

$$F(\mathbf{w}, x, \mathbf{y}) = F(\mathbf{w}, x, y_1, \dots, y_D) = \sum_r f_r(\mathbf{w}, x, \mathbf{y}_r)$$

Every $f_r(\mathbf{w}, x, \mathbf{y}_r)$ is an arbitrary composite function, e.g., a CNN

$$F(\mathbf{w}, x, \mathbf{y}) = F(\mathbf{w}, x, y_1, \dots, y_D) = \sum_r f_r(\mathbf{w}, x, \mathbf{y}_r)$$

How to compute gradient:

$$\begin{aligned} \frac{\partial}{\partial \mathbf{w}} & \sum_{i \in \mathcal{D}} \left(\log \sum_{\hat{\mathbf{y}}} \exp F(\mathbf{w}, x^{(i)}, \hat{\mathbf{y}}) - F(\mathbf{w}, x^{(i)}, \mathbf{y}^{(i)}) \right) \\ &= \sum_{i \in \mathcal{D}, \hat{\mathbf{y}}} \left(p(\hat{\mathbf{y}}; x^{(i)}, \mathbf{w}) - \delta(\hat{\mathbf{y}} = \mathbf{y}^{(i)}) \right) \nabla_{\mathbf{w}} F(\mathbf{w}, x^{(i)}, \hat{\mathbf{y}}) \\ &= \sum_{i \in \mathcal{D}, r, \hat{\mathbf{y}}_r} \left(p_r(\hat{\mathbf{y}}_r; x^{(i)}, \mathbf{w}) - \delta_r(\hat{\mathbf{y}}_r = \mathbf{y}_r^{(i)}) \right) \nabla_{\mathbf{w}} f_r(\mathbf{w}, x^{(i)}, \hat{\mathbf{y}}_r) \end{aligned}$$

How to obtain marginals $p_r(\hat{\mathbf{y}}_r; x, \mathbf{w})$?

Approximate marginals $b_r(\hat{\mathbf{y}}_r; x, \mathbf{w})$ via:

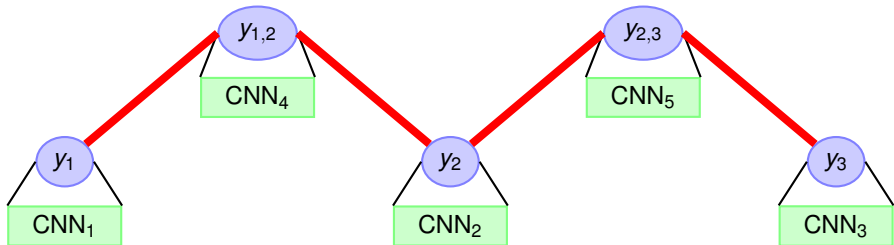
- LP relaxation
- Message passing

General algorithm for learning:

```
repeat
  repeat
    update marginals  $b_r$ 
  until convergence
  update parameters  $\mathbf{w}$ 
until convergence
```

What operation is used for computation of marginals in multi-class prediction? Soft-max

Intuition: Deep Structured Learning



Algorithm summary:

Repeat until stopping criteria

- 1 Forward passes to compute $f_r(\mathbf{w}, x, \hat{\mathbf{y}}_r)$
- 2 **Estimate beliefs** $b_r(\hat{\mathbf{y}}_r, x, \mathbf{w})$ **exactly/approximately**
- 3 Compute difference between estimated and groundtruth beliefs
- 4 Backpropagation of difference to obtain gradient
- 5 Update parameters \mathbf{w}

Example: Find five letters within distorted images



banal

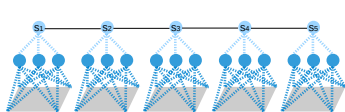


julep

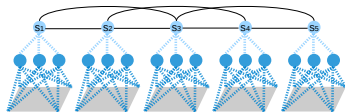


resty

- Graphical model: 1st or 2nd order Markov

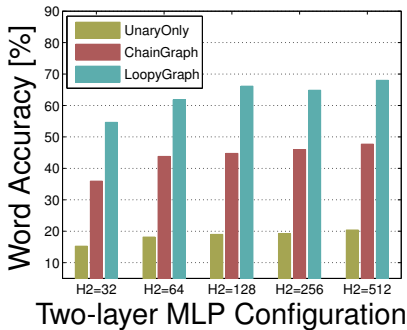
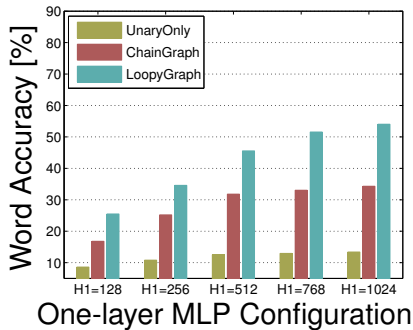


1st order Markov



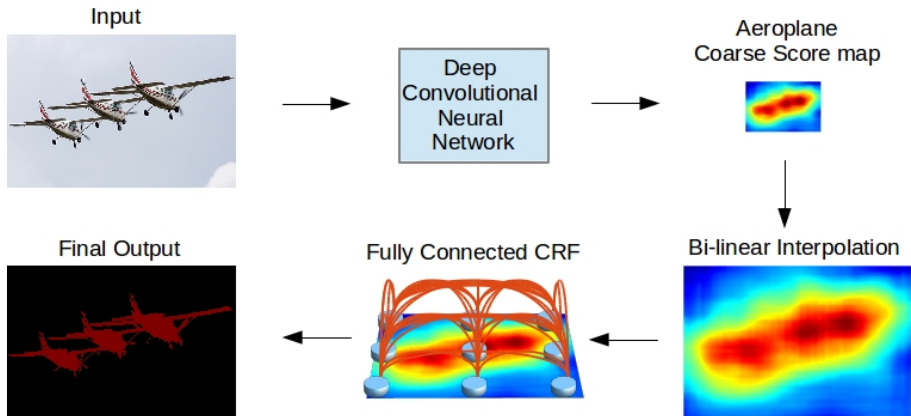
2nd order Markov

- Unary function: multi-layer perceptron (MLP)
- Pairwise function: linear or non-linear MLP
- $|\mathcal{Y}| = 26^5$

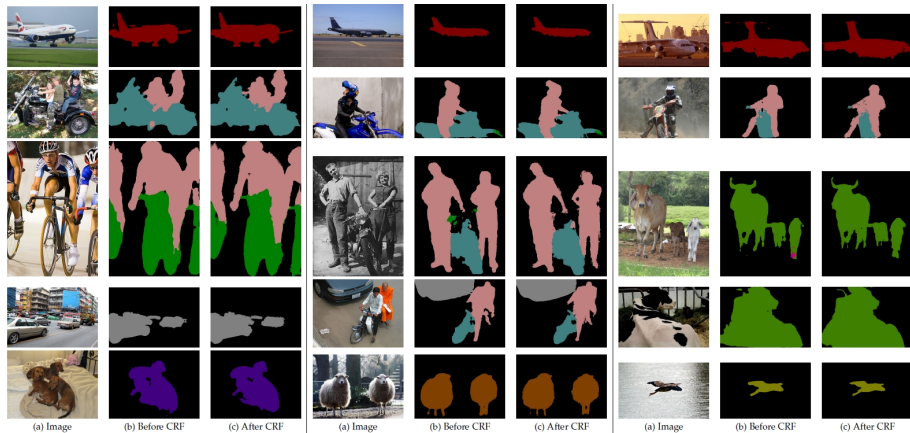


Deeper and more structured \Rightarrow better performance

Example: Semantic segmentation with DeepLab



Results:



Recap:

Inference:

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

Learning:

$$\min_{\mathbf{w}} \frac{C}{2} \|\mathbf{w}\|_2^2 + \sum_{i \in \mathcal{D}} \epsilon \ln \sum_{\hat{\mathbf{y}}} \exp \frac{L(\mathbf{y}^{(i)}, \hat{\mathbf{y}}) + F(\mathbf{w}, x^{(i)}, \hat{\mathbf{y}})}{\epsilon} - F(\mathbf{w}, x^{(i)}, \mathbf{y}^{(i)})$$

repeat

repeat

update marginals b_r

until convergence

update parameters \mathbf{w}

until convergence

Quiz:

- What's our feature vector for structured output space data?
- Describe algorithms for learning with structured output space data?
- What makes learning with structured output spaces hard?
- How can deep nets and structured output spaces be combined?

Important topics of this lecture

- Feature vectors for structured output space data
- Learning with structured output space data

Up next:

- Learning Theory