

# Max-Margin Markov Networks

Ben Taskar  
 Carlos Guestrin  
 Daphne Koller

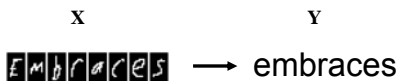
Presenter: Nam Nguyen

# Overview

- *Max-margin Markov networks captures the strengths of the two frameworks:*
  - *Kernel-based approaches (i.e. SVMs): the ability to use high-dimensional feature spaces; strong theoretical guarantees.*
  - *Graphical models (i.e. Markov networks): the ability to represent correlations between labels.*

# Problem Formulation

- Multi-label supervised learning:
  - Input:  $\mathbf{X} = (x_1, x_2, \dots, x_l)$
  - Output:  $\mathbf{Y} = (Y_1, Y_2, \dots, Y_l)$
- Example problem:



# Discriminative vs. Generative Models

- Linear Discriminative Model:

$$h_w(\mathbf{x}) = \arg \max_y \sum_{i=1}^n w_i f_i(\mathbf{x}, \mathbf{y}) = \arg \max_y \mathbf{w}^T \mathbf{f}(\mathbf{x}, \mathbf{y})$$

- Generative Model:

$$\bar{\mathbf{y}} = \arg \max_y P(\mathbf{x}, \mathbf{y}) = \arg \max_y P(\mathbf{y} | \mathbf{x}) P(\mathbf{x}) = \arg \max_y P(\mathbf{y} | \mathbf{x})$$

# Markov Network

- A pairwise Markov network
  - A graph  $G=(Y,E)$
  - A potential function associated with each edge
 
$$\psi_{ij}(\mathbf{x}, y_i, y_j) = \exp[\sum_{k=1}^n w_k f_k(\mathbf{x}, y_i, y_j)] = \exp[\mathbf{w}^T \mathbf{f}(\mathbf{x}, y_i, y_j)]$$
- The joint conditional probability distribution
 
$$P(\mathbf{y} | \mathbf{x}) \propto \prod_{(i,j) \in E} \psi_{ij}(\mathbf{x}, y_i, y_j) = \prod_{(i,j) \in E} \exp[\sum_{k=1}^n w_k f_k(\mathbf{x}, y_i, y_j)]$$

$$= \exp[\sum_{(i,j) \in E} \sum_{k=1}^n w_k f_k(\mathbf{x}, y_i, y_j)] = \exp[\sum_{k=1}^n w_k \underbrace{\sum_{(i,j) \in E} f_k(\mathbf{x}, y_i, y_j)}_{\mathbf{f}_k(\mathbf{x}, \mathbf{y})}]$$

$$= \exp[\sum_{k=1}^n w_k f_k(\mathbf{x}, \mathbf{y})] = \exp(\mathbf{w}^T \mathbf{f}(\mathbf{x}, \mathbf{y}))$$

# Marginal-based Markov networks

- Goal:
 
$$\arg \max_y P(\mathbf{y} | \mathbf{x}) = \arg \max_y [\exp(\mathbf{w}^T \mathbf{f}(\mathbf{x}, \mathbf{y}))] = \arg \max_y \mathbf{w}^T \mathbf{f}(\mathbf{x}, \mathbf{y})$$
- Start with formulation of SVMs for a *single-label* binary classification:
 
$$\text{maximize } \gamma$$

$$\text{s.t. } \|\mathbf{w}\| \leq 1; \quad \mathbf{w}^T \Delta \mathbf{f}_x(\mathbf{y}) \geq \gamma, \quad \forall \mathbf{x} \in S, \quad \forall \mathbf{y} \neq \mathbf{t}(\mathbf{x})$$

where  $\Delta \mathbf{f}_x(\mathbf{y}) = \mathbf{f}(\mathbf{x}, \mathbf{t}(\mathbf{x})) - \mathbf{f}(\mathbf{x}, \mathbf{y})$

## Margin-based Markov networks

- In the multi-label setting, the margin between  $\mathbf{t}(\mathbf{x})$  and  $\mathbf{y}$  scales linearly with the number of wrong labels in  $\mathbf{y}$ ,  $\Delta \mathbf{t}_{\mathbf{x}}(\mathbf{y})$ :  
 maximize  $\gamma$   
 s.t.  $\|\mathbf{w}\| \leq 1$ ;  $\mathbf{w}^\top \Delta \mathbf{f}_{\mathbf{x}}(\mathbf{y}) \geq \gamma \Delta \mathbf{t}_{\mathbf{x}}(\mathbf{y})$ ,  $\forall \mathbf{x} \in S$ ,  $\forall \mathbf{y}$   
 where  $\Delta \mathbf{t}_{\mathbf{x}}(\mathbf{y}) = \sum_{i=1}^l \Delta t_{\mathbf{x}}(y_i)$  and  $\Delta t_{\mathbf{x}}(y_i) \equiv I(y_i \neq (\mathbf{t}(\mathbf{x}))_i)$
- The QP form:  
 minimize  $\frac{1}{2} \|\mathbf{w}\|^2$   
 s.t.  $\mathbf{w}^\top \Delta \mathbf{f}_{\mathbf{x}}(\mathbf{y}) \geq \Delta \mathbf{t}_{\mathbf{x}}(\mathbf{y})$ ,  $\forall \mathbf{x} \in S$ ,  $\forall \mathbf{y}$

## Margin-based Markov networks

- Introduction of slack variables, the primal formulation:

$$\min \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{\mathbf{x}} \xi_{\mathbf{x}};$$

$$\text{s.t. } \mathbf{w}^\top \Delta \mathbf{f}_{\mathbf{x}}(\mathbf{y}) \geq \Delta \mathbf{t}_{\mathbf{x}}(\mathbf{y}) - \xi_{\mathbf{x}}, \forall \mathbf{x}, \mathbf{y}$$

- The dual formulation:

$$\max \sum_{\mathbf{x}, \mathbf{y}} \alpha_{\mathbf{x}}(\mathbf{y}) \Delta \mathbf{t}_{\mathbf{x}}(\mathbf{y}) - \frac{1}{2} \left\| \sum_{\mathbf{x}, \mathbf{y}} \alpha_{\mathbf{x}}(\mathbf{y}) \Delta \mathbf{f}_{\mathbf{x}}(\mathbf{y}) \right\|^2;$$

$$\text{s.t. } \sum_{\mathbf{y}} \alpha_{\mathbf{x}}(\mathbf{y}) = C, \forall \mathbf{x}; \alpha_{\mathbf{x}}(\mathbf{y}) \geq 0, \forall \mathbf{x}, \mathbf{y}.$$

## Exploiting Structure in M<sup>3</sup> networks

- Note: the number of constraints in the primal QP and the number of variables in the dual QP are exponential in the number of labels.
- Main Idea: the variables  $\alpha_{\mathbf{x}}(\mathbf{y})$  in the dual formulation can be interpreted as a density function over  $\mathbf{y}$  conditional on  $\mathbf{x}$ .

## Exploiting Structure in M3 networks

- Define the marginal dual variables:

$$\mu_{\mathbf{x}}(y_i, y_j) = \sum_{\mathbf{y} \sim [y_i, y_j]} \alpha_{\mathbf{x}}(\mathbf{y}), \quad \forall (i, j) \in E, \forall y_i, y_j, \forall \mathbf{x};$$

$$\mu_{\mathbf{x}}(y_i) = \sum_{\mathbf{y} \sim [y_i]} \alpha_{\mathbf{x}}(\mathbf{y}), \quad \forall i, \forall y_i, \forall \mathbf{x};$$

- Require conditions:

– Consistency:  $\sum_{y_i} \mu_{\mathbf{x}}(y_i, y_j) = \mu_{\mathbf{x}}(y_j)$ ,  $\forall y_j, \forall (i, j) \in E, \forall \mathbf{x}$

– The Markov network is a forest.

- The equivalent dual QP:

$$\max \sum_{\mathbf{x}} \sum_{\mathbf{y}} \mu_{\mathbf{x}}(\mathbf{y}) \Delta \mathbf{t}_{\mathbf{x}}(\mathbf{y}) - \frac{1}{2} \sum_{\mathbf{x}, \mathbf{x}'} \sum_{(i, j) \in E} \sum_{y_i, y_j, y_i', y_j'} \mu_{\mathbf{x}}(y_i, y_j) \mu_{\mathbf{x}'}(y_i', y_j') \mathbf{f}_{\mathbf{x}}(y_i, y_j)^\top \mathbf{f}_{\mathbf{x}'}(y_i', y_j');$$

$$\text{s.t. } \sum_{y_i} \mu_{\mathbf{x}}(y_i, y_j) = \mu_{\mathbf{x}}(y_j); \quad \sum_{y_i} \mu_{\mathbf{x}}(y_i) = C; \quad \mu_{\mathbf{x}}(y_i, y_j) \geq 0. \quad (10)$$

## Generalization Error Bound

- A  $\gamma$ -margin per-label loss measures the worst per-label loss on  $\mathbf{x}$  made by any classifier  $\mathbf{z}$  which is perturbed from  $\mathbf{w}^\top \mathbf{f}_{\mathbf{x}}$  by at most a  $\gamma$ -margin per-label.

$$\mathcal{L}^\gamma(\mathbf{w}, \mathbf{x}) = \sup_{\mathbf{z}} \sum_{i=1}^l \max_{y_i} \{z(y_i) - \mathbf{w}^\top \mathbf{f}_{\mathbf{x}}(y_i)\}_{\leq \gamma} \Delta t_{\mathbf{x}}(\mathbf{y}); \quad \forall \mathbf{y} \frac{1}{\gamma} \Delta \mathbf{t}_{\mathbf{x}}(\arg \max_{\mathbf{y}} \mathbf{z}(\mathbf{y}))$$

**Theorem 6.1** If the edge features have bounded 2-norm,  $\max_{(i, j), y_i, y_j} \|\mathbf{f}_{\mathbf{x}}(y_i, y_j)\|_2 \leq R_{\text{edge}}$ , then for a family of hyperplanes parameterized by  $\mathbf{w}$ , and any  $\delta > 0$ , there exists a constant  $K$  such that for any  $\gamma > 0$  per-label margin, and  $m > 1$  samples, the per-label loss is bounded by:

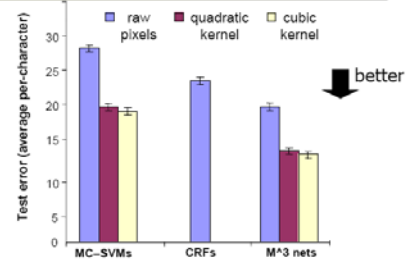
$$E_{\mathbf{x}} \mathcal{L}(\mathbf{w}, \mathbf{x}) \leq E_S \mathcal{L}^\gamma(\mathbf{w}, \mathbf{x}) + \sqrt{\frac{K}{m} \left[ \frac{R_{\text{edge}}^2 \|\mathbf{w}\|_2^2 q^2}{\gamma^2} [\ln m + \ln l + \ln q + \ln k] + \ln \frac{1}{\delta} \right]};$$

with probability at least  $1 - \delta$ , where  $q = \max_i |\{(i, j) \in E\}|$  is the maximum edge degree in the network,  $k$  is the number of classes in a label, and  $l$  is the number of labels. ■

## Handwriting Recognition

Length: ~8 chars  
 Letter: 16x8 pixels  
 10-fold Train/Test  
 5000/50000 letters  
 600/6000 words

Models:  
 Multiclass-SVMs  
 CRFs  
 M<sup>3</sup> nets



\*Taskar et al 03

## Conclusion

- M3 networks integrate kernel methods with graphical models.
- M3 networks effectively train by exploiting the network structure, i.e. a forest.
- Authors provide theoretical guarantees on the average *per-label* generalization error of the model in term of training set margin.