Partially Observed Maximum Entropy Discrimination Markov Networks

Jun Zhu

Dept. of Comp. Sci. & Tech., Tsinghua University Machine Learning Dept., Carnegie Mellon University

Joint work with Eric P. Xing and Bo Zhang

Outline

Introduction

- Structured Prediction
- Max-margin Markov Networks

• Max-Entropy Discrimination Markov Networks (MaxEnDNet)

- Basic Theorems
- Partially Observed MaxEnDNet
- Experimental Results
- Summary

Classification

- Inputs:
 - a set of training samples $\mathcal{D} = \{(\mathbf{x}^i, y^i)\}_{i=1}^N$, where $\mathbf{x}^i = [x_1^i, x_2^i, \cdots, x_d^i]^\top$ and $y^i \in C \triangleq \{c_1, c_2, \cdots, c_L\}$
- Outputs:
 - a predictive function $h(\mathbf{x})$: $y^* = h(\mathbf{x}) \triangleq \arg \max_{\mathbf{x}} F(\mathbf{x}, y; \mathbf{w})$
- Examples $(F(\mathbf{x}, y; \mathbf{w}) = \mathbf{w}^{\top} \mathbf{f}(\mathbf{x}, y))$:



Support Vector Machine (SVM) • Max-margin learning $\min_{\mathbf{w},\xi} \frac{1}{2} \mathbf{w}^{\top} \mathbf{w} + C \sum_{i=1}^{N} \xi_i;$ s.t. $\mathbf{w}^{\top} \Delta \mathbf{f}_i(y) \ge 1 - \xi_i, \ \forall i, \forall y \ne y^i.$

Structured Prediction

• Inputs:

- a set of training samples: $\mathcal{D} = \{(\mathbf{x}^i, \mathbf{y}^i)\}_{i=1}^N$, where

 $\mathbf{x}^{i} = (\mathbf{x}_{1}^{i}, \mathbf{x}_{2}^{i}, \cdots, \mathbf{x}_{\ell_{i}}^{i}) \text{ and } \mathbf{x}_{j}^{i} = [x_{j1}^{i}, \cdots, x_{jd}^{i}]^{\top}$ $\mathbf{y}^{i} = (y_{1}^{i}, y_{2}^{i}, \cdots, y_{\ell_{i}}^{i}) \text{ and } y_{j}^{i} \in C \triangleq \{c_{1}, c_{2}, \cdots, c_{L}\}$

- Examples:
 - Part-of-speech (POS) Tagging:

 $\mathbf{X}=$ "Do you want fries with that?" $\ _{->}\mathbf{y}=$ <verb pron verb noun prep pron>

– Image segmentation

$$\mathbf{x} = \begin{pmatrix} x_{11} & x_{12} & \dots \\ x_{21} & x_{22} & \dots \\ \vdots & \vdots & \dots \end{pmatrix} \qquad \mathbf{y} = \begin{pmatrix} y_{11} & y_{12} & \dots \\ y_{21} & y_{22} & \dots \\ \vdots & \vdots & \dots \end{pmatrix}$$

• Outputs:

- a predictive function $h(\mathbf{x}) : \mathbf{y}^{\star} = h(\mathbf{x}) \triangleq \arg \max_{\mathbf{y}} F(\mathbf{x}, \mathbf{y}; \mathbf{w})$

Structured Prediction Models

- Conditional Random Fields (CRFs) (Lafferty et al., 2001)
 - Based on Logistic Regression
 - Max-likelihood estimation (point-estimate)

$$\max_{\mathbf{w}} \mathcal{L}(\mathcal{D}; \mathbf{w}) \triangleq \sum_{i=1}^{N} \log p(\mathbf{y}^{i} | \mathbf{x}^{i})$$

$$p(\mathbf{y}|\mathbf{x}) = \frac{\exp\{\mathbf{w}^{\top}\mathbf{f}(\mathbf{x},\mathbf{y})\}}{\sum_{\mathbf{y}'}\exp\{\mathbf{w}^{\top}\mathbf{f}(\mathbf{x},\mathbf{y}')\}}$$

- Max-margin Markov Networks (M³Ns) (Taskar et al., 2003)
 - Based on SVM
 - Max-margin learning (point-estimate)

P0 (M³N) : $\min_{\mathbf{w},\xi} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^N \xi_i$ s.t. $\forall i, \forall \mathbf{y} \neq \mathbf{y}^i : \mathbf{w}^\top \Delta \mathbf{f}_i(\mathbf{y}) \ge \Delta \ell_i(\mathbf{y}) - \xi_i, \ \xi_i \ge 0$, where $\mathbf{w}^\top \Delta \mathbf{f}_i(\mathbf{y})$ denotes the margin and $\Delta \ell_i(\mathbf{y})$ is a loss function.

Markov properties are encoded in the feature functions f(x, y)



Between MLE and max-margin learning

• Likelihood-based estimation

- Probabilistic (joint/conditional likelihood model)
- Easy to perform Bayesian learning, and consider prior knowledge, missing data

Max-margin learning

- Non-probabilistic (concentrate on inputoutput mapping)
- Not obvious how to perform Bayesian learning or consider prior, and missing data
- Sound theoretical guarantee with limited samples
- Maximum Entropy Discrimination (MED) (Jaakkola, et al., 1999)
 - A Bayesian learning approach $\hat{y} = \operatorname{sign} \int p(\mathbf{w}) F(x; \mathbf{w}) \, \mathrm{d}\mathbf{w}$ $(y \in \{+1, -1\})$
 - The optimization problem (binary classification)

$\overset{\min KL(p(\Theta)||p_0(\Theta))}{\text{MED subsumes SVM}}.$

s.t. $p(\mathfrak{G})[g_ir(x, \mathbf{w}) - \varsigma_i] \mathfrak{u} \mathfrak{G} \geq 0, \forall i,$

where Θ is the parameter \mathbf{w} when ξ are kept fixed or the pair (\mathbf{w}, ξ) when we want to optimize over ξ

Machine Learning Lunch @ CMU

MaxEnt Discrimination Markov networks

 MaxEnt Discrimination Markov Networks (MaxEnDNet): P1(MaxEnDNet): min KL(p(w)||p₀(w)) + U(ξ) s.t. p(w) ∈ F₁, ξ_i ≥ 0, ∀i.

$$P(1) = 1, S_t = 3, \dots$$

- Generalized maximum entropy or regularized KL-divergence
- Subspace of distributions defined with *expected* margin constraints

 $\mathcal{F}_1 = \{ p(\mathbf{w}) : \int p(\mathbf{w}) [\Delta F_i(\mathbf{y}; \mathbf{w}) - \Delta \ell_i(\mathbf{y})] \, \mathrm{d}\mathbf{w} \ge -\xi_i, \, \forall i, \forall \mathbf{y} \neq \mathbf{y}^i \},\$

• Bayesian-style Prediction

 $h_1(\mathbf{x}) = \arg \max_{\mathbf{y} \in \mathcal{Y}(\mathbf{x})} \int p(\mathbf{w}) F(\mathbf{x}, \mathbf{y}; \mathbf{w}) \, \mathrm{d}\mathbf{w}$

Machine Learning Lunch @ CMU

 $D(p, p_0) = KL(p || p_0)$

 p_0

p

Solution to MaxEnDNet

- Theorem 1 (Solution to MaxEnDNet):
 - Posterior Distribution:

$$p(\mathbf{w}) = \frac{1}{Z(\alpha)} p_0(\mathbf{w}) \exp\left\{\sum_{i,\mathbf{y}} \alpha_i(\mathbf{y}) [\Delta F_i(\mathbf{y}; \mathbf{w}) - \Delta \ell_i(\mathbf{y})]\right\}$$

- Dual Optimization Problem:

D1: $\max_{\alpha} -\log Z(\alpha) - U^{\star}(\alpha)$ s.t. $\alpha_i(\mathbf{y}) \ge 0, \ \forall i, \ \forall \mathbf{y},$

U^{*}(·) is the conjugate of the U(·), i.e., U^{*}(α) = sup_ξ (Σ_{i,y} α_i(y)ξ_i - U(ξ))
Convex conjugate (closed proper convex φ(μ))

- Def: $\phi^{\star}(\nu) = \sup_{\mu} [\nu^{\top} \mu - \phi(\mu)]$ - Ex: $\phi(p(\mu)) = KL(p(\mu)||p_0(\mu)) \qquad \phi^{\star}(\nu) = \log Z(\nu)$ $\phi(\xi) = C \sum_i |\xi_i| \qquad \phi^{\star}(\nu) = \mathbb{I}_{\infty}(|\nu_i| \le C, \forall i)$

Machine Learning Lunch @ CMU

Reduction to M³Ns

Theorem 2 (Reduction of MaxEnDNet to M³Ns):

- Assume $F(\mathbf{x}, \mathbf{y}; \mathbf{w}) = \mathbf{w}^{\top} \mathbf{f}(\mathbf{x}, \mathbf{y}), U(\xi) = C \sum_{i} \xi_{i}, and p_{0}(\mathbf{w}) = \mathcal{N}(\mathbf{w}|0, I)$

Posterior distribution: $p(\mathbf{w}) = \mathcal{N}(\mathbf{w}|\mu_{\mathbf{w}}, I)$, where $\mu_{\mathbf{w}} = \sum_{i, \mathbf{y}} \alpha_i(\mathbf{y}) \Delta \mathbf{f}_i(\mathbf{y})$ Dual optimization:

$$\max_{\alpha} \sum_{i,\mathbf{y}} \alpha_i(\mathbf{y}) \Delta \ell_i(\mathbf{y}) - \frac{1}{2} \| \sum_{i,\mathbf{y}} \alpha_i(\mathbf{y}) \Delta \mathbf{f}_i(\mathbf{y}) \|^2$$

s.t.
$$\sum_{\mathbf{y}} \alpha_i(\mathbf{y}) = C; \ \alpha_i(\mathbf{y}) \ge 0, \ \forall i, \ \forall \mathbf{y},$$

Predictive rule:

$$h_1(\mathbf{x}) = \arg \max_{\mathbf{y} \in \mathcal{Y}(\mathbf{x})} \int p(\mathbf{w}) F(\mathbf{x}, \mathbf{y}; \mathbf{w}) \, \mathrm{d}\mathbf{w} = \arg \max_{\mathbf{y} \in \mathcal{Y}(\mathbf{x})} \mu_{\mathbf{w}}^{\top} \mathbf{f}(\mathbf{x}, \mathbf{y})$$

- Thus, MaxEnDNet subsumes M³Ns and admits all the merits of max-margin learning
- Furthermore, MaxEnDNet has at least three advantages ...

Three Advantages

- PAC-Bayesian prediction error guarantee $\Pr_Q(M(h, \mathbf{x}, \mathbf{y}) \le 0) \le \Pr_D(M(h, \mathbf{x}, \mathbf{y}) \le \gamma) + O\left(\sqrt{\frac{\gamma^{-2}KL(p||p_0)\ln(N|\mathcal{Y}|) + \ln N + \ln \delta^{-1}}{N}}\right)$
- Introduce regularization effects, such as sparsity bias
 - Laplace prior => Posterior shrinkage effects (Laplace M³N, ICML'08)



- An elegant approach to incorporate latent variables and structures
 - Partially observed MaxEnDNet

10

Motivating Example

🖉 shop. scholastic. com - Scholastic Store - Micro. . . 📃 • Web data extraction <u>File Edit View Favorites Tools Help</u> Results 1 - 3 of 3 • Goal: Name, Image, Price, Description, etc. Barney Safety Songs Hooray For Barney Driver hree Barney board Learn safety tips by books will delight singing and playing Y_0 vour child with happ[,] along with a purple photos. Given Data Record dinosaur Our Price: \$12.95 Our Price: \$12.99 You Save: 28% You Save: 35% Ages Birth - 3 Electronic Toy Y_1 Ages 1 - 5 Data Record 1 Data Record 2 Web Page Tail Tail [<mark>Head]</mark> {Info Block} Y_5 Hierarchical model {Repeat block} {Note} [{Note} Y_6 Advantages: Data Record Data Record Computational efficiency 0 o Long-range dependency {name, price} {image} {name, price} {image} o Joint extraction name}][Image) Desc {price} {desc} {name} Image Desc Name Name) Desc) (Price) Note Note Desc)(Desc Price Note 1/15/2009

Note

Learn hierarchical model with latent variables

- Can we learn a hierarchical model from partially labeled data?
 - Yes!
 - Partially observed CRFs for object recognition (NIPS'04)
 - Dynamic Hierarchical MRFs for web data extraction (ICML'07)
- How about max-margin learning?
 - Yes!
 - Easy with MaxEnDNet

Partially observed MaxEnDNet

• $\mathcal{D} = \{\langle \mathbf{x}^i, \mathbf{y}^i \rangle\}_{i=1}^N$ augmented with hidden variables $\{\mathbf{z}\} = (\mathbf{z}^1, \cdots, \mathbf{z}^N)$

P1(MaxEnDNet) : $\min_{p(\mathbf{w}),\xi} KL(p(\mathbf{w})||p_0(\mathbf{w})) + U(\xi)$ s.t. $p(\mathbf{w}) \in \mathcal{F}_1, \ \xi_i \ge 0, \forall i.$

• Def of PoMEN:

P2(PoMEN): $\min_{p(\mathbf{w}, \{\mathbf{z}\}), \xi} KL(p(\mathbf{w}, \{\mathbf{z}\}) || p_0(\mathbf{w}, \{\mathbf{z}\})) + U(\xi)$ s.t. $p(\mathbf{w}, \{\mathbf{z}\}) \in \mathcal{F}_2, \ \xi_i \ge 0, \forall i.$

$$\mathcal{F}_{2} = \left\{ p(\mathbf{w}, \{\mathbf{z}\}) : \sum_{\mathbf{z}} \int p(\mathbf{w}, \mathbf{z}) [\Delta F_{i}(\mathbf{y}, \mathbf{z}; \mathbf{w}) - \Delta \ell_{i}(\mathbf{y})] \, \mathrm{d}\mathbf{w} \geq -\xi_{i}, \, \forall i, \forall \mathbf{y} \neq \mathbf{y}^{i} \right\},\$$

Prediction:

$$h_2(\mathbf{x}) = \arg \max_{\mathbf{y} \in \mathcal{Y}(\mathbf{x})} \sum_{\mathbf{z}} \int p(\mathbf{w}, \mathbf{z}) F(\mathbf{x}, \mathbf{y}, \mathbf{z}; \mathbf{w}) \, \mathrm{d}\mathbf{w}$$

Alternating Minimization Alg.

• Factorization assumption: $p_0(\mathbf{w}, \{\mathbf{z}\}) = p_0(\mathbf{w}) \prod_{i=1}^{i=1} p_0(\mathbf{z}_i)$

$$p(\mathbf{w}, \{\mathbf{z}\}) = p(\mathbf{w}) \prod_{i=1}^{N} p(\mathbf{z}_i)$$

• Alternating minimization:

• Step 1: keep
$$p(\mathbf{z})$$
 fixed, optimize over $p(\mathbf{w})$

$$\min_{p(\mathbf{w}),\xi} KL(p(\mathbf{w})||p_0(\mathbf{w})) + C\sum_{i} \xi_{i}$$
o Normal prior
s.t. $p(\mathbf{w}) \in \mathcal{F}'_{1}, \xi_{i} \ge 0, \forall i.$
 $\mathcal{F}'_{1} = \{p(\mathbf{w}): \int p(\mathbf{w})E_{p(\mathbf{z})}[\Delta F_{i}(\mathbf{y},\mathbf{z};\mathbf{w}) - \Delta \ell_{i}(\mathbf{y})] \, d\mathbf{w} \ge -\xi_{i}, \forall i, \forall \mathbf{y}\}$
o Laplace prior
• Laplace M^3N problem (VB)
• Step 2: keep $p(\mathbf{w})$ fixed, optimize over $p(\mathbf{z})$

$$\min_{p(\mathbf{w}),\xi} KL(p(\mathbf{z})||p_0(\mathbf{z})) + C\xi_{i}$$

$$p(\mathbf{z}) = \frac{1}{Z(\beta)}p_{0}(\mathbf{z}) \exp\{\sum_{\mathbf{y}} \beta(\mathbf{y})[\mu_{\mathbf{w}}^{\top}\Delta f_{i}(\mathbf{y},\mathbf{z}) - \Delta \ell_{i}(\mathbf{y})]\}$$
s.t. $p(\mathbf{z}) \in \mathcal{F}_{1}^{\star}, \xi_{i} \ge 0.$

$$max - \log(\sum_{\mathbf{z}} p_{0}(\mathbf{z}) \exp\{\sum_{\mathbf{y}} \beta(\mathbf{y})[\mu_{\mathbf{w}}^{\top}\Delta f_{i}(\mathbf{y},\mathbf{z}) - \Delta \ell_{i}(\mathbf{y})]\}$$

$$\mathcal{F}_{1}^{\star} = \{p(\mathbf{z}): \sum_{\mathbf{z}} p(\mathbf{z}) \int p(\mathbf{w})[\Delta F_{i}(\mathbf{y},\mathbf{z};\mathbf{w}) - \Delta \ell_{i}(\mathbf{y})] \, d\mathbf{w} \ge -\xi_{i}, \forall i, \forall \mathbf{y}\}$$
s.t. $\sum_{\mathbf{y}} \beta(\mathbf{y}) = C, \beta(\mathbf{y}) \ge 0, \forall \mathbf{y}.$
Machine Learning Lunch @ CMU
equivalently reduced to a NLP with a polynomial number of constraints

Experimental Results

- Web data extraction:
 - Name, Image, Price, Description
 - Methods:
 - Hierarchical CRFs, Hierarchical M[^]3N
 - PoMEN, Partially observed HCRFs
 - Pages from 37 templates
 - Training: 185 (5/per template) pages, or 1585 data records
 - Testing: 370 (10/per template) pages, or 3391 data records
 - Record-level Evaluation
 - Leaf nodes are labeled
 - Page-level Evaluation
 - Supervision Level 1:
 - Leaf nodes and data record nodes are labeled
 - Supervision Level 2:
 - Level 1 + the nodes above data record nodes



Record-Level Evaluations

- Overall performance:
 - Avg F1:

0.9

0.85

0.8

0.7

0.65

0.6

0

Machin

缸 0.75

- avg F1 over all attributes
- Block instance accuracy:
 - % of records whose *Name*, *Image*, and *Price* are correct
- Attribute performance:

Name

HCRF

НМЗN

PoM3N

40

20

Training Ratio

PoHCRF

Page-Level Evaluations

- Supervision Level 1:
 - Leaf nodes and data record nodes are labeled

- Supervision Level 2:
 - Level 1 + the nodes above data record nodes

Summary

- MaxEnt Discrimination Markov Networks (MaxEnDNet)

 PAC-Bayesian performance guarantee
 Sparsity regularization effects
 Incorporating latent variables and structures
- Experimental results show the advantages of max-margin learning over likelihood methods with latent variables

Margin-based Learning Paradigms

Machine Learning Lunch (a) CMU