



Recent Advances in Learning **SPARSE** Structured I/O Models: models, algorithms, and applications

Eric Xing

epxing@cs.cmu.edu

Machine Learning Dept./Language Technology Inst./Computer Science Dept.
Carnegie Mellon University

VLPR 2009 @ Beijing, China

8/6/2009

1



Structured Prediction Problem

- Unstructured prediction



$$\mathbf{x} = (x_{11} \ x_{12} \ \dots)$$

$$\mathbf{y} = y_1$$

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

$$\mathbf{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \end{pmatrix}$$

- Structured prediction

- Part of speech tagging

$$\mathbf{x} = \text{"Do you want sugar in it?"} \Rightarrow \mathbf{y} = \langle \text{verb pron verb noun prep pron} \rangle$$

- Image segmentation



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

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

VLPR 2009 @ Beijing, China

8/6/2009

2



Classical Predictive Models

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

- Logistic Regression, Bayes classifiers

- Max-likelihood estimation

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

$$p(y|x) = \frac{\exp\{\mathbf{w}^T \mathbf{f}(x, y)\}}{\sum_{y'} \exp\{\mathbf{w}^T \mathbf{f}(x, y')\}}$$

- Support Vector Machines (SVM)

- Max-margin learning

$$\min_{\mathbf{w}, \xi} \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_{i=1}^N \xi_i;$$

$$\text{s.t. } \mathbf{w}^T \Delta \mathbf{f}_i(y) \geq 1 - \xi_i, \forall i, \forall y \neq y^i.$$

Advantages:

1. Full probabilistic semantics
2. Straightforward Bayesian or direct regularization VLP
3. Hidden structures or generative hierarchy

Advantages:

1. Dual sparsity: few support vectors
2. Kernel tricks
3. Strong empirical results



Structured Prediction Models

• Conditional Random Fields (CRFs) (Lafferty et al 2001)

- Based on Logistic Regression
- Max-likelihood estimation (point-estimate)

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

$$p(\mathbf{y} | \mathbf{x}) = \frac{\exp\{\mathbf{w}^T \mathbf{f}(\mathbf{x}, \mathbf{y})\}}{\sum_{\mathbf{y}'} \exp\{\mathbf{w}^T \mathbf{f}(\mathbf{x}, \mathbf{y}')\}}$$

• Max-margin Markov Networks (M³Ns) (Taskar et al 2003)

- Based on SVM
- Max-margin learning (point-estimate)

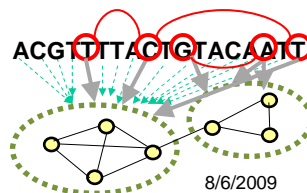
$$P_0(\text{M}^3\text{N}) : \min_{\mathbf{w}, \xi} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^N \xi_i$$

$$\text{s.t. } \forall i, \forall \mathbf{y} \neq \mathbf{y}^i : \mathbf{w}^T \Delta \mathbf{f}_i(\mathbf{y}) \geq \Delta \ell_i(\mathbf{y}) - \xi_i, \xi_i \geq 0,$$

where $\mathbf{w}^T \Delta \mathbf{f}_i(\mathbf{y} | \mathbf{x}_i)$ denotes the margin and $\Delta \ell_i(\mathbf{y})$ is a loss function.

Challenges:

- **SPARSE** prediction model
- Prior information of structures
- Scalable to large-scale problems (e.g., 10⁴ input/output dimension)



SAILING LAB
State Key Laboratory of Intelligent Information Processing

Outline

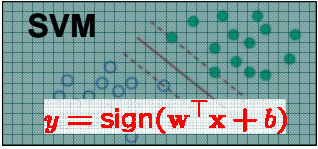
- Structured sparse regression
 - Graph-guided fused lasso: unlinked SNPs to trait networks (Kim and Xing, PLoS Genetics)
 - Temporally-smoothed graph regression: learning time-varying graphs (Ahmed and Xing, PNAS 2009, Kolar and Xing, under review, Annals of Statistics)
- Maximum entropy discrimination Markov networks
 - General Theorems (Zhu and Xing, JMLR submitted)
 - Gaussian MEDN: reduction to M³N (Zhu, Xing and Zhang, ICML 08)
 - Laplace MEDN: a sparse M³N (Zhu, Xing and Zhang, ICML 08)
 - Partially observed MEDN: (Zhu, Xing and Zhang, NIPS 08)
 - Max-margin/Max entropy topic model: (Zhu, Ahmed, and Xing, ICML 09)

VLPR 2009 @ Beijing, China 8/6/2009 5

SAILING LAB
State Key Laboratory of Intelligent Information Processing

Max-Margin Learning Paradigms

SVM



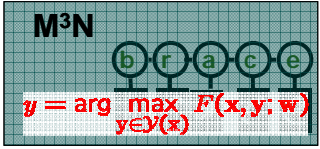
$y = \text{sign}(w^T x + b)$

$$\min_{w, \xi} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^m \xi_i$$

$y^i (w^T x^i + b) \geq 1 - \xi_i, \forall i$

→

M³N

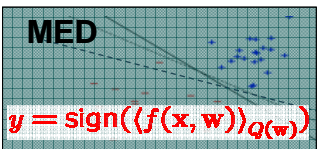


$y = \arg \max_{y \in \mathcal{Y}(x)} F(x, y; w)$

$$\min_{w, \xi} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^m \xi_i$$

$w^T [f(x^i, \cdot)] - f(x^i, y^i) \geq \ell(y^i, y) - \xi_i, \forall i, \forall y \neq y^i$

MED



$y = \text{sign}((f(x, w))_Q(w))$

$$\min_Q \text{KL}(Q \| Q_0)$$

$y^i (f(x^i))_Q \geq \xi_i, \forall i$


→

MED-MN

= SMED + “Bayesian” M³N

Primal and Dual Sparse!

VLPR 2009 @ Beijing, China 8/6/2009 6



Primal and Dual Problems of M³Ns

• Primal problem:

$$P0 (M^3N) : \min_{\mathbf{w}, \xi} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^N \xi_i$$

s.t. $\forall i, \forall \mathbf{y} \neq \mathbf{y}^i : \mathbf{w}^\top \Delta \mathbf{f}_i(\mathbf{y}) \geq \Delta \ell_i(\mathbf{y}) - \xi_i,$
 $\xi_i \geq 0,$

• Algorithms

- Cutting plane
- Sub-gradient
- ...

• Dual problem:

$$D0 (M^3N) : \max_{\alpha} \sum_{i, \mathbf{y}} \alpha_i(\mathbf{y}) \Delta \ell_i(\mathbf{y}) - \frac{1}{2} \eta^\top \eta$$

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

where $\eta = \sum_{i, \mathbf{y}} \alpha_i(\mathbf{y}) \Delta \mathbf{f}_i(\mathbf{y}).$


• Algorithms:

- SMO
- Exponentiated gradient
- ...

$$\mathbf{w}^* = \eta^* = \sum_{i, \mathbf{y}} \alpha_i^*(\mathbf{y}) \Delta \mathbf{f}_i(\mathbf{y}).$$

• So, M³N is dual sparse!

VLPR 2009 @ Beijing, China
8/6/2009
7



MLE versus max-margin learning

• Likelihood-based estimation

- Probabilistic (joint/conditional likelihood model)
- Easy to perform Bayesian learning, and incorporate prior knowledge, latent structures, missing data
- Bayesian regularization!!

• Max-margin learning

- Non-probabilistic (concentrate on input-output 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)

- Model averaging $\hat{y} = \text{sign} \int p(\mathbf{w}) F(x; \mathbf{w}) d\mathbf{w} \quad (y \in \{+1, -1\})$
- The optimization problem (binary classification)


$$\min_{\Theta} KL(p(\Theta) || p_0(\Theta))$$

MED subsumes SVM.

$$\text{s.t. } \int p(\Theta) |y_i F(x; \mathbf{w}) - \xi_i| d\Theta \geq 0, \forall i.$$

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

VLPR 2009 @ Beijing, China
8/6/2009
8



MaxEnt Discrimination Markov Network

- Structured MaxEnt Discrimination (SMED):

$$P1: \min_{p(\mathbf{w}), \xi} KL(p(\mathbf{w}) || p_0(\mathbf{w})) + U(\xi)$$

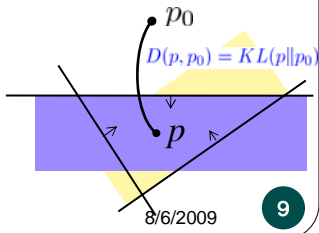
s.t. $p(\mathbf{w}) \in \mathcal{F}_1, \xi_i \geq 0, \forall i.$

generalized maximum entropy or regularized KL-divergence
- Feasible subspace of weight distribution:


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

expected margin constraints.
- Average from distribution of M³Ns

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



VLPR 2009 @ Beijing, China 8/6/2009 9



Solution to MaxEnDNet

- Theorem 1:
 - 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:

$$D1: \max_{\alpha} -\log Z(\alpha) - U^*(\alpha)$$

s.t. $\alpha_i(\mathbf{y}) \geq 0, \forall i, \forall \mathbf{y},$

$U^*(\cdot)$ is the conjugate of the $U(\cdot)$, i.e., $U^*(\alpha) = \sup_{\xi} (\sum_{i, \mathbf{y}} \alpha_i(\mathbf{y}) \xi_i - U(\xi))$

VLPR 2009 @ Beijing, China 8/6/2009 10

SAILING LAB
State Key Laboratory of Intelligent Information Processing

Gaussian MaxEnDNet (reduction to M³N)

- Theorem 2
 - Assume $F(\mathbf{x}, \mathbf{y}; \mathbf{w}) = \mathbf{w}^T \mathbf{f}(\mathbf{x}, \mathbf{y})$, $U(\xi) = C \sum_i \xi_i$, 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:

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

$$\text{s.t. } \sum_{\mathbf{y}} \alpha_i(\mathbf{y}) = C; \alpha_i(\mathbf{y}) \geq 0, \forall i, \forall \mathbf{y},$$
 - Predictive rule:

$$h_1(\mathbf{x}) = \arg \max_{\mathbf{y} \in \mathcal{Y}(\mathbf{x})} \int p(\mathbf{w}) F(\mathbf{x}, \mathbf{y}; \mathbf{w}) d\mathbf{w} = \arg \max_{\mathbf{y} \in \mathcal{Y}(\mathbf{x})} \mu_{\mathbf{w}}^T \mathbf{f}(\mathbf{x}, \mathbf{y})$$
- Thus, MaxEnDNet subsumes M³Ns and admits all the merits of max-margin learning
- Furthermore, MaxEnDNet has at least **three advantages** ...

VLPR 2009 @ Beijing, China 8/6/2009 **11**

SAILING LAB
State Key Laboratory of Intelligent Information Processing

Three Advantages

- An averaging Model: PAC-Bayesian prediction error guarantee

$$\Pr_Q(M(h, \mathbf{x}, \mathbf{y}) \leq 0) \leq \Pr_P(M(h, \mathbf{x}, \mathbf{y}) \leq \gamma) + O\left(\sqrt{\frac{\gamma^{-2} KL(p||p_0) \ln(N|\mathcal{Y}^I) + \ln N + \ln \delta^{-1}}{N}}\right).$$
- Entropy regularization: Introducing useful biases
 - Standard Normal prior => reduction to standard M³N (we've seen it)
 - Laplace prior => Posterior shrinkage effects (sparse M³N)

$$\forall k, \langle w_k \rangle_p = \frac{2\eta_k}{\lambda - \eta_k^2}$$
- Integrating Generative and Discriminative principles
 - Incorporate latent variables and structures (PoMEN)
 - Semisupervised learning (with partially labeled data)

VLPR 2009 @ Beijing, China 8/6/2009 **12**



I: Generalization Guarantee

- MaxEntNet is an averaging model
- Theorem 3 (PAC-Bayes Bound)

Let p_0 be any continuous probability distribution over \mathcal{H} and $\delta \in (0, 1)$

If $\forall F \in \mathcal{H} : \mathcal{X} \times \mathcal{Y} \mapsto [-c, c]$

Then with probability at least $1 - \delta$ over random samples \mathcal{D} of Q , for very distribution p over \mathcal{H} and for all margin thresholds $\gamma > 0$

$$\Pr_Q(M(h, \mathbf{x}, \mathbf{y}) \leq 0) \leq \Pr_{\mathcal{D}}(M(h, \mathbf{x}, \mathbf{y}) \leq \gamma) + O\left(\sqrt{\frac{\gamma^{-2} KL(p||p_0) \ln(N|\mathcal{Y}|) + \ln N + \ln \delta^{-1}}{N}}\right).$$

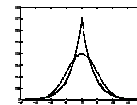
where $\Pr_Q(\cdot)$ and $\Pr_{\mathcal{D}}(\cdot)$ represent the prob. of event under dist. Q and D , respectively.



II: Laplace MaxEnDNet (primal sparse M^3N)

- Laplace Prior:

$$p_0(\mathbf{w}) = \prod_{k=1}^K \frac{\sqrt{\lambda}}{2} e^{-\sqrt{\lambda}|w_k|} = \left(\frac{\sqrt{\lambda}}{2}\right)^K e^{-\sqrt{\lambda}\|\mathbf{w}\|}$$



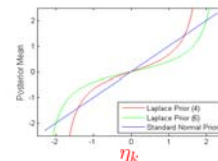
- Corollary 4:

- Under a Laplace MaxEnDNet, the posterior mean of parameter vector \mathbf{w} is:

$$\forall k, \langle w_k \rangle_p = \frac{2\eta_k}{\lambda - \eta_k^2}$$


where the vector η is a linear combination of "support vectors":

$$\eta = \sum_{\alpha} \alpha_i(\mathbf{y}) \Delta f_i(\mathbf{y})$$



- The Gaussian MaxEnDNet and the regular M^3N has no such shrinkage
 - there, we have

$$\langle \mathbf{w} \rangle_p = \eta \iff \forall k, \langle w_k \rangle_p = \eta_k$$



$$\min_{\mu, \xi} |\mu| + C \sum_{i=1}^N \xi_i$$

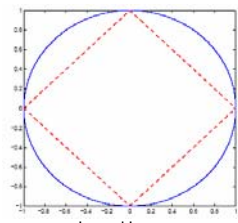
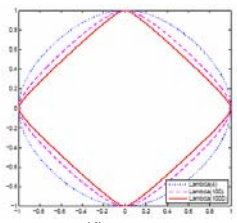
$$\text{s.t. } \mu^\top \Delta f_i(\mathbf{y}) \geq \Delta \ell_i(\mathbf{y}) - \xi_i; \xi_i \geq 0, \forall i, \forall \mathbf{y} \neq \mathbf{y}^i.$$

LapMEDN vs. L_2 and L_1 regularization

- Corollary 5: LapMEDN corresponding to solving the following primal optimization problem:

$$\min_{\mu, \xi} \sqrt{\lambda} \sum_{k=1}^K \left(\sqrt{\mu_k^2 + \frac{1}{\lambda}} - \frac{1}{\sqrt{\lambda}} \log \frac{\sqrt{\lambda \mu_k^2 + 1} + 1}{2} \right) + C \sum_{i=1}^N \xi_i$$


$$\text{s.t. } \mu^\top \Delta f_i(\mathbf{y}) \geq \Delta \ell_i(\mathbf{y}) - \xi_i; \xi_i \geq 0, \forall i, \forall \mathbf{y} \neq \mathbf{y}^i.$$
- KL norm: $\|\mu\|_{KL} \triangleq \sum_{k=1}^K \left(\sqrt{\mu_k^2 + \frac{1}{\lambda}} - \frac{1}{\sqrt{\lambda}} \log \frac{\sqrt{\lambda \mu_k^2 + 1} + 1}{2} \right)$

VLPR 2009 @ Beijing, China

8/6/2009

15



Variational Learning of LapMEDN

- Exact dual function is hard to optimize

$$\max_{\alpha} L - \sum_{k=1}^K \log \frac{\lambda}{\lambda - \eta_k^2}$$
- Use the hierarchical representation, we get:

$$KL(p||p_0) = -H(p) - \langle \log \int p(\mathbf{w}|\tau)p(\tau|\lambda) d\tau \rangle_p$$

$$\leq -H(p) - \langle \int q(\tau) \log \frac{p(\mathbf{w}|\tau)p(\tau|\lambda)}{q(\tau)} d\tau \rangle_p \triangleq \mathcal{L}(p(\mathbf{w}), q(\tau))$$
- We optimize an upper bound:

$$\min_{p(\mathbf{w}) \in \mathcal{F}_1; q(\tau); \xi} \mathcal{L}(p(\mathbf{w}), q(\tau)) + U(\xi)$$
- Why is it easier?
 - Alternating minimization leads to nicer optimization problems

Keep $q(\tau)$ fixed	Keep $p(\mathbf{w})$ fixed
The effective prior is normal $\forall k: p_0(w_k \tau_k) = \mathcal{N}(w_k 0, (\frac{1}{\tau_k})^{-1})$	Closed form solution of $q(\tau)$ and its expectation $\langle \frac{1}{\tau_k} \rangle_q = \sqrt{\frac{1}{\eta_k^2}}$

VLPR 2009 @ Beijing, China

8/6/2009

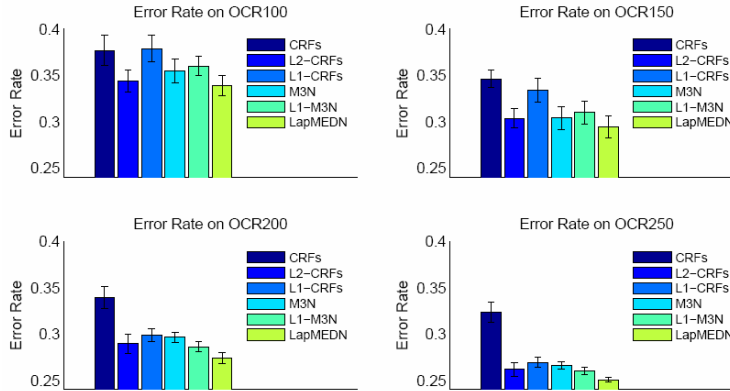
16



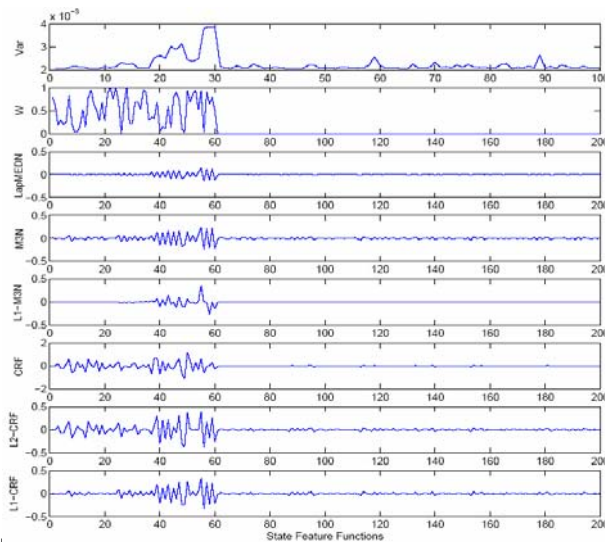
Experimental results on OCR datasets


(CRFs, L_1 -CRFs, L_2 -CRFs, M^3 Ns, L_1 - M^3 Ns, and LapMEDN)

- We randomly construct OCR100, OCR150, OCR200, and OCR250 for 10 fold CV.

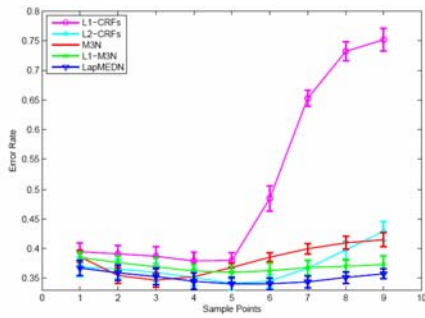


Feature Selection






Sensitivity to Regularization Constants



- L_1 -CRF and L_2 -CRF:
 - 0.001, 0.01, 0.1, 1, 4, 9, 16
- M^3N and Lap M^3N :
 - 1, 4, 9, 16, 25, 36, 49, 64, 81

- L_1 -CRFs are much sensitive to regularization constants; the others are more stable
- Lap M^3N is the most stable one

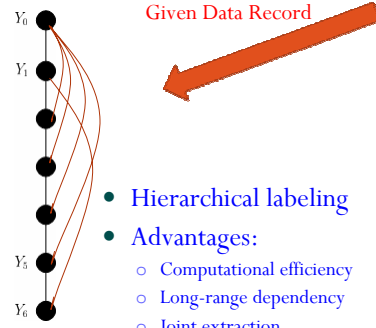
VLPR 2009 @ Beijing, China
8/6/2009
19




III: Latent Hierarchical MaxEnDNet

- Web data extraction
 - Goal: *Name, Image, Price, Description, etc.*

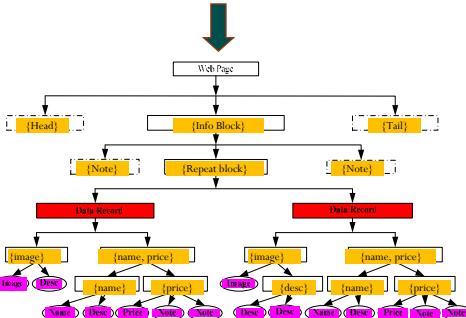
Given Data Record



- Hierarchical labeling
- Advantages:
 - Computational efficiency
 - Long-range dependency
 - Joint extraction



Web Page



VLPR 2009 @ Beijing, China



Partially Observed MaxEnDNet (PoMEN)

- Now we are given partially labeled data: $\mathcal{D} = \{ \langle \mathbf{x}^i, \mathbf{y}^i, \mathbf{z}^i \rangle \}_{i=1}^N$

- PoMEN: learning $p(\mathbf{w}, \mathbf{z})$

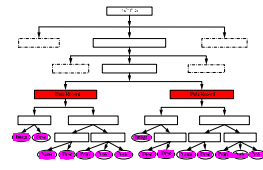
$$P2(\text{PoMEN}) : \min_{p(\mathbf{w}, \{\mathbf{z}\}), \xi} KL(p(\mathbf{w}, \{\mathbf{z}\}) || p_0(\mathbf{w}, \{\mathbf{z}\})) + U(\xi)$$

$$\text{s.t. } p(\mathbf{w}, \{\mathbf{z}\}) \in \mathcal{F}_2, \xi_i \geq 0, \forall i.$$

$$\mathcal{F}_2 = \{ p(\mathbf{w}, \{\mathbf{z}\}) : \sum_{\mathbf{z}} \int p(\mathbf{w}, \mathbf{z}) [\Delta F_i(\mathbf{y}, \mathbf{z}; \mathbf{w}) - \Delta \ell_i(\mathbf{y})] d\mathbf{w} \geq -\xi_i, \forall i, \forall \mathbf{y} \neq \mathbf{y}^i \},$$

- Prediction:

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



Alternating Minimization Alg.

- Factorization assumption:

$$p_0(\mathbf{w}, \{\mathbf{z}\}) = p_0(\mathbf{w}) \prod_{i=1}^N p_0(\mathbf{z}_i) \quad p(\mathbf{w}, \{\mathbf{z}\}) = p(\mathbf{w}) \prod_{i=1}^N p(\mathbf{z}_i)$$

- Alternating minimization:

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

$$\min_{p(\mathbf{w}), \xi} KL(p(\mathbf{w}) || p_0(\mathbf{w})) + C \sum_i \xi_i$$

$$\text{s.t. } p(\mathbf{w}) \in \mathcal{F}_1^*, \xi_i \geq 0, \forall i.$$

- o Normal prior
 - M³N problem (QP)
- o Laplace prior
 - Laplace M³N problem (VB)

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

- Step 2: keep $p(\mathbf{w})$ fixed, optimize over $p(\mathbf{z})$

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

$$\text{s.t. } p(\mathbf{z}) \in \mathcal{F}_1^*, \xi_i \geq 0.$$

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

Equivalently reduced to an LP with a polynomial number of constraints

SAILING LAB
Advanced Computational Methods for Complex Systems Research

Record-Level Evaluations

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

The top two graphs show overall performance. The left graph plots Average F1 (0.86 to 0.92) and the right graph plots Block Instance Accuracy (0.65 to 0.85) against Training Ratio (0 to 40). Both graphs compare HCRF (blue), PoHCRF (green), HM3N (magenta), and PoM3N (red). PoM3N consistently shows the highest performance, followed by HM3N, HCRF, and PoHCRF.

The bottom four graphs show attribute performance for Name, Image, Price, and Description. Each graph plots F1 (0.6 to 0.98) against Training Ratio (0 to 40). The same four models are compared. PoM3N generally performs best across all attributes, with particularly high F1 scores for Name and Image.

VLPR 2009 @ 3

SAILING LAB
Advanced Computational Methods for Complex Systems Research

VI: Max-Margin/Max Entropy Topic Model – MED-LDA

(from images.google.cn)

"Arts"	"Budgets"	"Children"	"Education"
NEW	MILLION	CHILDREN	SCHOOL
FILM	TAX	WOMEN	STUDENTS
SHOW	PROGRAM	PEOPLE	SCHOOLS
MUSIC	BUDGET	CHILD	EDUCATION
MOVIE	BILLION	YEARS	TEACHERS
PLAY	FEDERAL	FAMILIES	HIGH
MUSICAL	YEAR	WORK	PUBLIC
REST	SPENDING	PARENTS	TEACHER
ACTOR	NEW	SAYS	BENNETT
FIRST	STATE	FAMILY	MANGAT
YORK	PLAN	WELFARE	SAMPHY
OPERA	MONEY	MEN	STATE
THEATER	PROGRAMS	PERCENT	PRESIDENT
ADDRESS	GOVERNMENT	CARE	ELEMENTARY
LOVE	CONGRESS	LIFE	HAITI

The **William Randolph Hearst Foundation** will give \$1.25 million to Lincoln Center, Metropolitan Opera Co, New York Philharmonic and Juilliard School. "Our board felt that we had a real opportunity to make a mark on the future of the performing arts with these grants as art every bit as important as our traditional areas of support in health, medical research, education and the social services," Hearst Foundation President Randolph A. Hearst said Monday in announcing the grants. Lincoln Center's share will be \$200,000 for its new building, which will house young artists and provide new public facilities. The Metropolitan Opera Co. and New York Philharmonic will receive \$400,000 each. The Juilliard School, where music and the performing arts are taught, will get \$250,000. The Hearst Foundation, a leading supporter of the Lincoln Center Consolidated Corporate Fund, will make its usual annual \$100,000 donation, too.

VLPR 2009 @ Beijing, China 8/6/2009 24

SAILING LAB

LDA: a generative story for documents

- Bag-of-word representation of documents
- Each word is generated by ONE topic
- Each document is a random mixture over topics

Document #1: gif jpg image current
file color images ground power file
current format file formats circuit gif
images

Document #2: wire currents file format
ground power image format wire circuit
current wiring ground circuit images
files...

Mixture Components Mixture Weights Bayesian Approach
Dirichlet Prior **LDA**

VLPR 2009 @ Beijing, China 8/6/2009 25

SAILING LAB

LDA: Latent Dirichlet Allocation

(Blei et al., 2003)

Dirichlet parameter α Per-document topic proportions θ_d Per-word topic assignment $Z_{d,n}$ observed word $W_{d,n}$ Topics β_k

θ

θ_1	...	θ_D
0.8	...	0.3
0.2	...	0.7

β : $\begin{pmatrix} 0.70 & 0.05 & 0.08 & \dots \\ 0.12 & 0.52 & 0.05 & \dots \end{pmatrix}$

- Generative Procedure:
 - For each document d :
 - Sample a topic proportion $\theta_d \sim \text{Dir}(\alpha)$
 - For each word:
 - Sample a topic $Z_{d,n} \sim \text{Mult}(\theta_d)$
 - Sample a word $W_{d,n} \sim \text{Mult}(\beta_{z_{d,n}})$

• Joint Distribution: $p(\theta, \mathbf{z}, \mathbf{W} | \alpha, \beta) = \prod_{d=1}^D p(\theta_d | \alpha) \left(\prod_{n=1}^N p(z_{dn} | \theta_d) p(w_{dn} | z_{dn}, \beta) \right)$
exact inference intractable!

• Variational Inference with $q(\mathbf{z}, \theta) \sim p(\mathbf{z}, \theta | \mathbf{W}, \alpha, \beta)$
 $\mathcal{L}(q) \triangleq -E_q[\log p(\theta, \mathbf{z}, \mathbf{W} | \alpha, \beta)] - \mathcal{H}(q(\mathbf{z}, \theta)) \geq -\log p(\mathbf{W} | \alpha, \beta)$

• Minimize the variational bound to estimate parameters and infer the posterior distribution

VLPR 2009 @ Beijing, China 8/6/2009 26

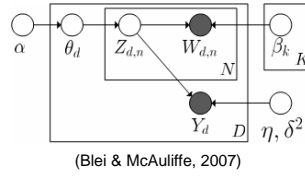


Supervised Topic Model (sLDA)

- LDA ignores documents' side information (e.g., categories or rating score), thus lead to suboptimal topic representation for supervised tasks
- Supervised Topic Models handle such problems, e.g., **sLDA** (Blei & McAuliffe, 2007) and DiscLDA (Simon et al., 2008)

- Generative Procedure (sLDA):

- For each document d :
 - Sample a topic proportion $\theta_d \sim \text{Dir}(\alpha)$
 - For each word:
 - Sample a topic $Z_{d,n} \sim \text{Mult}(\theta_d)$
 - Sample a word $W_{d,n} \sim \text{Mult}(\beta_{z_{d,n}})$
 - Sample y_d



(Blei & McAuliffe, 2007)

$$y_d \sim \begin{cases} \mathcal{N}(\eta^\top Z_d, \delta^2), & \text{if } y_d \text{ is continuous} \\ \text{GLM}(Z_d, \eta, \delta^2), & \text{otherwise} \end{cases}$$

- Joint distribution:

$$p(\theta, \mathbf{z}, \mathbf{y}, \mathbf{W} | \alpha, \beta, \eta, \delta^2) = \prod_{d=1}^D p(\theta_d | \alpha) \left(\prod_{n=1}^N p(z_{dn} | \theta_d) p(w_{dn} | z_{dn}, \beta) \right) p(y_d | \eta^\top z_d, \delta^2)$$

- Variational inference:


$$\mathcal{L}(q) \triangleq -E_q[\log p(\theta, \mathbf{z}, \mathbf{y}, \mathbf{W} | \alpha, \beta, \eta, \delta^2)] - \mathcal{H}(q(\mathbf{z}, \theta)) \geq -\log p(\mathbf{y}, \mathbf{W} | \alpha, \beta, \eta, \delta^2)$$



The big picture

Max-Likelihood Estimation	Max-Margin and Max-Likelihood
sLDA	MedLDA

- How to integrate the max-margin principle into a probabilistic latent variable model?



MedLDA Regression Model

- Generative Procedure (Bayesian sLDA):
 - Sample a parameter $\eta \sim p_0(\eta)$
 - For each document d :
 - Sample a topic proportion $\theta_d \sim \text{Dir}(\alpha)$
 - For each word:
 - Sample a topic $Z_{d,n} \sim \text{Mult}(\theta_d)$
 - Sample a word $W_{d,n} \sim \text{Mult}(\beta_{z_{d,n}})$
 - Sample y_d : $y_d \sim \mathcal{N}(\eta^\top \bar{Z}_d, \delta^2)$
- Def:

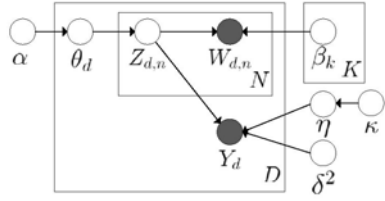
$$P1(\text{MedLDA}^r) : \min_{\alpha, \beta, \delta^2, \xi, \xi^*} \underbrace{-\log p(\mathbf{y}, \mathbf{W} | \alpha, \beta, \delta^2)}_{\text{intractable}} + C \sum_{d=1}^D (\xi_d + \xi_d^*)$$

$$\text{s.t. } \forall d : \begin{cases} y_d - E[\eta^\top \bar{Z}_d] \leq \epsilon + \xi_d, \mu_d \\ -y_d + E[\eta^\top \bar{Z}_d] \leq \epsilon + \xi_d^*, \mu_d^* \\ \xi_d \geq 0, v_d \\ \xi_d^* \geq 0, v_d^* \end{cases}$$

model fitting

predictive accuracy
- Predictive Rule:


$$\bar{y} = E[Y | w_{1:N}, \alpha, \beta, \delta^2] = E_{q(Z, \eta)}[\eta^\top \bar{Z} | w_{1:N}, \alpha, \beta, \delta^2]$$



VLPR 2009 @ Beijing, China

8/6/2009

29



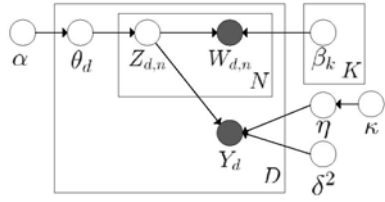
MedLDA Regression Model

- Generative Procedure (Bayesian sLDA):
 - Sample a parameter $\eta \sim p_0(\eta)$
 - For each document d :
 - Sample a topic proportion $\theta_d \sim \text{Dir}(\alpha)$
 - For each word:
 - Sample a topic $Z_{d,n} \sim \text{Mult}(\theta_d)$
 - Sample a word $W_{d,n} \sim \text{Mult}(\beta_{z_{d,n}})$
 - Sample y_d : $y_d \sim \mathcal{N}(\eta^\top \bar{Z}_d, \delta^2)$
 - Def:

$$P1(\text{MedLDA}^r) : \min_{q, \alpha, \beta, \delta^2, \xi, \xi^*} \underbrace{\mathcal{L}(q)}_{\text{intractable}} + C \sum_{d=1}^D (\xi_d + \xi_d^*)$$

$$\text{s.t. } \forall d : \begin{cases} y_d - E[\eta^\top \bar{Z}_d] \leq \epsilon + \xi_d, \mu_d \\ -y_d + E[\eta^\top \bar{Z}_d] \leq \epsilon + \xi_d^*, \mu_d^* \\ \xi_d \geq 0, v_d \\ \xi_d^* \geq 0, v_d^* \end{cases}$$
 - η is a hidden random variable
 - Variational bound $q(\theta, \mathbf{z}, \eta | \gamma, \phi) \sim p(\theta, \mathbf{z}, \eta | \alpha, \beta, \delta^2, \mathbf{y}, \mathbf{W})$
$$\mathcal{L}(q) \triangleq -E[\log p(\theta, \mathbf{z}, \eta, \mathbf{y}, \mathbf{W} | \alpha, \beta, \delta^2)] - \mathcal{H}(q(\mathbf{z}, \theta, \eta)) \geq -\log p(\mathbf{y}, \mathbf{W} | \alpha, \beta, \delta^2)$$
 - Predictive Rule:

$$\bar{y} = E[Y | w_{1:N}, \alpha, \beta, \delta^2] = E_{q(Z, \eta)}[\eta^\top \bar{Z} | w_{1:N}, \alpha, \beta, \delta^2]$$



VLPR 2009 @ Beijing, China

8/6/2009

30



Variational EM Alg.

- **E-step**: infer the posterior distribution of hidden r.v. $(\theta, \mathbf{z}, \eta)$
- **M-step**: estimate unknown parameters $(\alpha, \beta, \delta^2)$

- Independence assumption: $q(\theta, \mathbf{z}, \eta|\gamma, \phi) = q(\eta) \prod_{d=1}^D q(\theta_d|\gamma_d) \prod_{n=1}^N q(z_{dn}|\phi_{dn})$

$$L(\gamma, \phi, q(\eta), \alpha, \beta, \delta^2, \xi, \xi^*, \mu, \mu^*, v, v^*) = \mathcal{L}(q) + C \sum_{d=1}^D (\xi_d + \xi_d^*) - \sum_{d=1}^D \sum_{i=1}^N c_{di} (\sum_{j=1}^K \phi_{dij} - 1) - \sum_{d=1}^D \mu_d (\epsilon + \xi_d - y_d + E[\eta^\top \bar{Z}_d]) - \sum_{d=1}^D (\mu_d^* (\epsilon + \xi_d^* + y_d - E[\eta^\top \bar{Z}_d]) + v_d \xi_d + v_d^* \xi_d^*)$$

- Optimize L over ϕ :

$$\phi_{di} \propto \exp \left(\underbrace{E[\log \theta|\gamma]} + \underbrace{E[\log p(w_{di}|\beta)]}_{\frac{y_d}{N\delta^2} E[\eta]} - \underbrace{\frac{2E[\eta^\top \phi_{d,-i}] + E[\eta \circ \eta]}{2N^2\delta^2}} + \underbrace{\frac{E[\eta]}{N} (\mu_d - \mu_d^*)}_{\text{regularizer}} \right)$$

- The first two terms are the same as in LDA
- The third and fourth terms are similar to those of sLDA, but in **expected** version. The **variance** matters!
- The last term is a **regularizer**. Only **support vectors** affect the topic proportions
- Optimize L over other variables. See our paper for details!



MedLDA Classification Model

- Normalization factor in GLM makes inference harder
- We use LDA as the underlying topic model

- Multiclass MedLDA Classification Model:

$$P2(\text{MedLDA}^c) : \min_{q, q(\eta), \alpha, \beta, \xi} \mathcal{L}(q) + KL(q(\eta)||p_0(\eta)) + C \sum_{d=1}^D \xi_d$$

$$\text{s.t. } \forall d, y \neq y_d : E[\eta^\top \Delta \mathbf{f}_d(y)] \geq 1 - \xi_d; \xi_d \geq 0,$$

model fitting

predictive accuracy

- Variational upper bound $(q(\theta, \mathbf{z}|\gamma, \phi) \sim p(\theta, \mathbf{z}|\mathbf{W}, \alpha, \beta))$

$$\mathcal{L}(q) \triangleq -E[\log p(\theta, \mathbf{z}, \mathbf{W}|\alpha, \beta)] - \mathcal{H}(q(\theta, \mathbf{z})) \geq -\log p(\mathbf{W}|\alpha, \beta)$$

- **Expected** margin constraints. $\Delta \mathbf{f}_d(y) = \mathbf{f}(y_d, \bar{Z}_d) - \mathbf{f}(y, \bar{Z}_d)$

- Predictive Rule: $y^* = \arg \max_y E[\eta^\top \mathbf{f}(y, \bar{Z})|\alpha, \beta]$



Variational EM Alg.

- Independence assumption:

$$q(\theta, \mathbf{z}|\gamma, \phi) = \prod_{d=1}^D q(\theta_d|\gamma_d) \prod_{n=1}^N q(z_{dn}|\phi_{dn})$$

- Lagrangian function:

$$L(q, q(\eta), \xi, \mu_d(y), c_{di}) = \mathcal{L}(q) + KL(q(\eta)||p_0(\eta)) + C \sum_{d=1}^D \xi_d - \sum_{d=1}^D v_d \xi_d - \sum_{d=1}^D \sum_{y \neq y_d} \mu_d(y) (E[\eta^\top \Delta \mathbf{f}_d(y)] + \xi_d - 1) - \sum_{d=1}^D \sum_{i=1}^N c_{di} (\sum_{j=1}^K \phi_{dij} - 1)$$

- Optimize L over ϕ :

$$\phi_{di} \propto \exp \left(\underbrace{E[\log \theta|\gamma]}_{\text{LDA}} + E[\log p(w_{di}|\beta)] + \underbrace{\frac{1}{N} \sum_{y \neq y_d} \mu_d(y) E[\eta_{y_d} - \eta_y]}_{\text{MedLDA}} \right).$$

Only support vectors matter!

- Optimize L over other variables. See the paper for details



MedTM: a general framework

- MedLDA can be generalized to arbitrary topic models:
 - Unsupervised or supervised
 - Generative or undirected random fields (e.g., Harmoniums)

- MED Topic Model (MedTM):

$$P(\text{MedTM}) : \min_{q(H), q(\Upsilon), \Psi, \xi} \underbrace{\mathcal{L}(q(H))}_{\text{model fitting}} + \underbrace{KL(q(\Upsilon)||p_0(\Upsilon))}_{\text{predictive accuracy}} + U(\xi)$$

s.t. *expected margin constraints*

- H : hidden r.v.s in the underlying topic model, e.g., (θ, \mathbf{z}) in LDA
- Υ : parameters in predictive model, e.g., η in sLDA
- Ψ : parameters of the topic model, e.g., α in LDA
- \mathcal{L} : an variational upper bound of the log-likelihood
- U : a convex function over slack variables



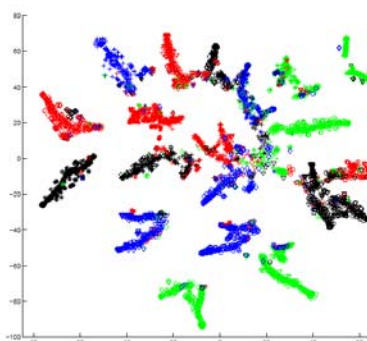
Experiments

- Goal:
 - To **qualitatively** and **quantitatively** evaluate how the max-margin estimates of MedLDA affect its topic discovering procedure
- Data Sets:
 - 20 Newsgroups (**classification**)
 - Documents from 20 categories
 - ~ 20,000 documents in each group
 - Remove stop word as listed in UMASS Mallet
 - Movie Review (**regression**)
 - 5006 documents, and 1.6M words
 - Dictionary: 5000 terms selected by tf-idf
 - Preprocessing to make the response approximately normal (Blei & McAuliffe, 2007)

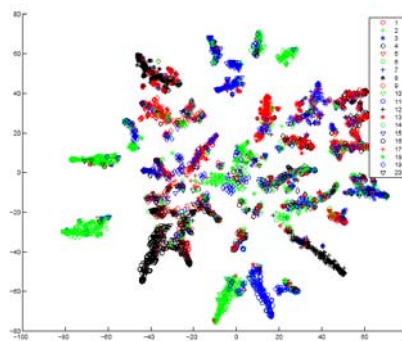


Document Modeling

- Data Set: 20 Newsgroups
- 110 topics + 2D embedding with t-SNE (van der Maaten & Hinton, 2008)



MedLDA



LDA

Document Modeling (cont')

comp.graphics

politics.mideast

MedLDA			LDA		
T 69	T 11	T 80	T 59	T 104	T 31
image	graphics	db	image	ftp	card
jpeg	image	key	jpeg	pub	monitor
gif	data	chip	color	graphics	dos
file	ftp	encryption	file	mail	video
color	software	clipper	gif	version	apple
files	pub	system	images	tar	windows
bit	mail	government	format	file	drivers
images	package	keys	bit	information	vga
format	fax	law	files	send	cards
program	images	escrow	display	server	graphics

T 30	T 40	T 51	T 42	T 78	T 47
israel	turkish	israel	israel	jews	armenian
israeli	armenian	lebanese	israeli	jewish	turkish
jews	armenians	israeli	peace	israel	armenians
arab	armenia	lebanon	writes	israeli	armenia
writes	people	people	article	arab	turks
people	turks	attacks	arab	people	genocide
article	greek	soldiers	war	arabs	russian
jewish	turkey	villages	lebanese	center	soviet
state	government	peace	lebanon	jew	people
rights	soviet	writes	people	nazi	muslim


VLPR 2009 @ Beijing, China 8/6/2009 37

Classification

- Data Set:** 20Newsgroups
 - Binary classification: "alt.atheism" and "talk.religion.misc" (Simon et al., 2008)
 - Multiclass Classification: all the 20 categories
- Models:** DiscLDA, sLDA (Binary ONLY! Classification sLDA (Wang et al., 2009)), LDA+SVM (baseline), MedLDA, MedLDA+SVM
- Measure:** Relative Improvement Ratio

$$RR(M) = \frac{precision(M)}{precision(LDA + SVM)} - 1$$

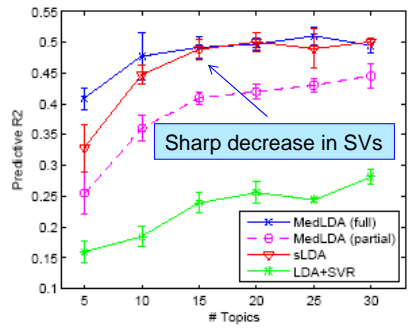
VLPR 2009 @ Beijing, China 8/6/2009 38



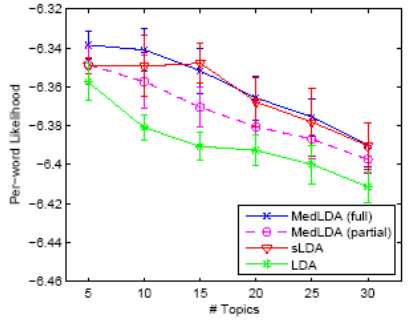
Regression

- **Data Set:** Movie Review (Blei & McAuliffe, 2007)
- **Models:** MedLDA(*partial*), MedLDA(*full*), sLDA, LDA+SVR
- **Measure:** predictive R^2 and per-word log-likelihood


$$pR^2 = 1 - \frac{\sum_d (y_d - \hat{y}_d)^2}{\sum_d (y_d - \bar{y}_d)^2}$$



Sharp decrease in SVs



VLPR 2009 @ Beijing, China
8/6/2009
39



Summary

- A general framework of MaxEnDNet for learning structured input/output models
 - Subsumes the standard M^3 Ns
 - Model averaging: PAC-Bayes theoretical error bound
 - Entropic regularization: sparse M^3 Ns
 - Generative + discriminative: latent variables, semi-supervised learning on partially labeled data
- Laplace MaxEnDNet: simultaneously primal and dual sparse
 - Can perform as well as sparse models on synthetic data
 - Perform better on real data sets
 - More stable to regularization constants
- PoMEN
 - Provides an elegant approach to incorporate latent variables and structures under max-margin framework
 - Experimental results show the advantages of max-margin learning over likelihood methods with latent variables

VLPR 2009 @ Beijing, China
8/6/2009
40

SAILING LAB
Laboratory for Statistical Artificial Intelligence & Integrative Genomics

Margin-based Learning Paradigms

SVM

$y = \text{sign}(\mathbf{w}^T \mathbf{x} + b)$

Structured prediction

Bayes learning

$\min_{\mathbf{w}, \xi} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^N \xi_i$
s.t. $y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1 - \xi_i, \forall i.$

M³N

$\mathbf{y}^* = \arg \max_{\mathbf{y}} \mathbf{w}^T \mathbf{f}(\mathbf{x}, \mathbf{y}; \mathbf{w})$

Structured prediction

Bayes learning

$\min_{\mathbf{w}, \xi} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^N \xi_i$
s.t. $\mathbf{w}^T \Delta \mathcal{F}_i(\mathbf{y}) \geq \Delta \mathcal{L}_i(\mathbf{y}) - \xi_i, \forall i, \forall \mathbf{y} \neq \mathbf{y}'$

MED

$y = \text{sign}(\langle \mathbf{w}^T \mathbf{f}(\mathbf{x}) \rangle_{p(\mathbf{w})})$

Structured prediction

Bayes learning

$\min_{p, \xi} KL(p||p_0) + C \sum_{i=1}^N \xi_i$
s.t. $y_i(\mathbf{f}(\mathbf{x}_i))_{p(\mathbf{w})} \geq 1 - \xi_i, \forall i.$

MaxEnDNet

$\mathbf{y}^* = \arg \max_{\mathbf{y}} \langle \mathbf{w}^T \mathbf{f}(\mathbf{x}, \mathbf{y}; \mathbf{w}) \rangle_{p(\mathbf{w})}$


Structured prediction

Bayes learning

$\min_{p(\mathbf{w}), \xi} KL(p||p_0) + U(\xi)$
s.t. $\int p(\mathbf{w})[\Delta \mathcal{F}_i(\mathbf{y}; \mathbf{w}) - \Delta \mathcal{L}_i(\mathbf{y})] d\mathbf{w} \geq -\xi_i, \xi_i \geq 0, \forall i, \forall \mathbf{y} \neq \mathbf{y}'.$


VLPR 2009 @ Beijing, China
8/6/2009
41

Acknowledgement




Laboratory for Statistical Artificial Intelligence & Integrative Genomics

Postdocs




Seyoung Kim Postdoc (PHD, Postdoc, UCSD)
Le Song Postdoc (PHD, USydney)

Visitors



Jun Zhu Tsinghua University

PhD students



Amr Ahmed (CI), Wang Fu (CSD), Fan Guo (CSD), Steve Hannink (MLD), Jude Hoernig (Comp Bio), Hutunandan K. (CSD), Maden Kolar (CI), Andrej Martins (CI / UT Austin, UT), Koki Purohit (CSD), Pradipta Ray (Comp Bio), Sayan Choudhury (CSD), Hyung Ah Kim (CSD)

<http://www.sailing.cs.cmu.edu/>

VLPR 2009 @ Beijing, China
8/6/2009
42




Thanks!

Reference:



Markov Chain Prior

$$P(\mathbf{c}) = P(c_1) \prod_{j=2}^J P(c_j | c_{j-1})$$




Markov Chain Prior

$$P(c) = P(c_1) \prod_{j=2}^J P(c_j | c_{j-1})$$

- $c_j = c_{j-1}$ if
 - 1) the **distance** between the two SNPs is small, or
 - 2) the **recombination rate** between the two SNPs is small

VLPR 2009 @ Beijing, China
8/6/2009
45



Markov Chain Prior

$$P(c) = P(c_1) \prod_{j=2}^J P(c_j | c_{j-1})$$


Poisson process

$$P(c_j | c_{j-1}) = \exp(-d_j \rho_j) \delta(c_j, c_{j-1}) + (1 - \exp(-d_j \rho_j)) \Pi_{c_{j-1}, c_j}$$

- ρ_j : Recombination rate at j th SNP
- d_j : Distance between j th and $(j-1)$ th SNP
- Π : Transition probability matrix

$$\begin{pmatrix} \pi_0 & 1 - \pi_0 \\ 1 - \pi_1 & \pi_1 \end{pmatrix}$$

VLPR 2009 @ Beijing, China
8/6/2009
46



Variational Bayesian Learning (Cont')

$$\min_{p(\mathbf{w}) \in \mathcal{F}_1; q(\tau); \xi} \mathcal{L}(p(\mathbf{w}), q(\tau)) + U(\xi)$$

Initialize $\langle \mathbf{w} \rangle_p^1 \leftarrow 0, \Sigma_w^1 \leftarrow I$

Solve an M³N Problem (Σ_w^t)

Update Σ_w^{t+1}

$t \leftarrow t + 1$

$$p(\mathbf{w}) \propto \exp\left\{ \int q(\tau) \log p(\mathbf{w}|\tau) d\tau - b \right\} \cdot \exp\{\mathbf{w}^\top \eta - L\}$$

$$\propto \exp\left\{ -\frac{1}{2} \mathbf{w}^\top (A^{-1})_q \mathbf{w} - b + \mathbf{w}^\top \eta - L \right\}$$

$$= \mathcal{N}(\mathbf{w} | \mu_w, \Sigma_w), \text{ where } L = \sum \alpha_i(y) \Delta \ell_i(y), A = \text{diag}(\tau_k),$$

$$\mu_w = \Sigma_w \eta \text{ and } \Sigma_w = \left(\langle A^{-1} \rangle_q \right)^{-1} = \langle \mathbf{w} \mathbf{w}^\top \rangle_p - \langle \mathbf{w} \rangle_p \langle \mathbf{w} \rangle_p^\top$$

$$\max_{\alpha} \sum_{i,y} \alpha_i(y) \Delta \ell_i(y) - \frac{1}{2} \eta^\top \Sigma_w \eta$$

$$\text{s.t. } \sum_y \alpha_i(y) = C; \alpha_i(y) \geq 0, \forall i, \forall y.$$

$$\min_{w, \xi} \frac{1}{2} \mathbf{w}^\top \Sigma_w^{-1} \mathbf{w} + C \sum_{i=1}^N \xi_i$$

$$\text{s.t. } \mathbf{w}^\top \Delta \ell_i(y) \geq \Delta \ell_i(y) - \xi_i; \xi_i \geq 0, \forall i, \forall y \neq y^i.$$


$$q(\tau) = \prod_{k=1}^K q(\tau_k)$$

$$q(\tau_k) \propto p(\tau_k | \lambda) \exp\{(\log p(w_k | \tau_k))_p\}$$

$$\propto \mathcal{N}(\sqrt{\langle w_k^2 \rangle_p} | 0, \tau_k) \exp(-\frac{1}{2} \lambda \tau_k).$$

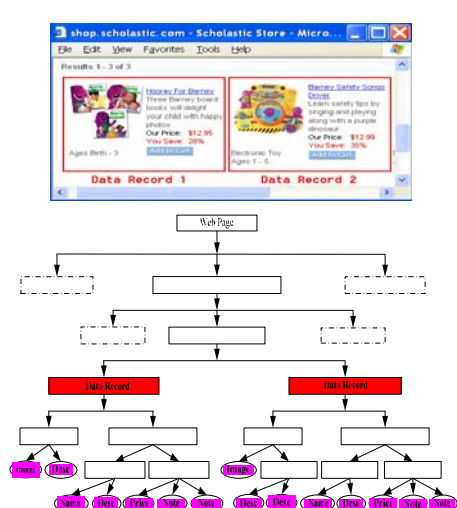
$$\langle \frac{1}{\tau_k} \rangle_q = \int \frac{1}{\tau_k} q(\tau_k) d\tau_k = \sqrt{\langle w_k^2 \rangle_p}$$

VLPR 2009 @ Beijing, China
8/6/2009
47



Experimental Results

- Web data extraction:
 - Name, Image, Price, Description
- Methods:
 - Hierarchical CRFs, Hierarchical M³N
 - 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



VLPR 2009 @ Beijing, China
8/6/2009
48