

# Multi-Task Sparsity via Maximum Entropy Discrimination

Tony Jebara  
Columbia University

July 26, 2011

- 1 Task Constraints in Generative Learning
- 2 Maximum Entropy Discrimination
- 3 Feature & Kernel Sparsity
- 4 Multi-Task Margin Constraints
- 5 Sequential Quadratic Programming
- 6 Graphical Model Structure Estimation
- 7 Relative Margin Constraints
- 8 Conclusions

# Generative Learning

- Given  $\mathcal{D} = (\mathbf{x}_t, y_t)_{t=1}^T$  sampled *iid* from unknown  $P(\mathbf{x}, y)$
- Find rule producing  $\hat{y}$  from  $\mathbf{x}$  with low error

# Generative Learning

- Given  $\mathcal{D} = (\mathbf{x}_t, y_t)_{t=1}^T$  sampled *iid* from unknown  $P(\mathbf{x}, y)$
- Find rule producing  $\hat{y}$  from  $\mathbf{x}$  with low error
- Generative Bayesian approach:
  - Assume  $p(\mathbf{x}, y|\Theta)$  and  $p(\Theta)$
  - Get  $p(\Theta|\mathcal{D}) \propto \prod_{t=1}^T p(\mathbf{x}_t, y_t|\Theta)p(\Theta)$
  - Predict  $\hat{y} = \arg \max_y \int_{\Theta} p(\mathbf{x}, y|\Theta)p(\Theta|\mathcal{D})$

# Generative Learning

- Given  $\mathcal{D} = (\mathbf{x}_t, y_t)_{t=1}^T$  sampled *iid* from unknown  $P(\mathbf{x}, y)$
- Find rule producing  $\hat{y}$  from  $\mathbf{x}$  with low error
- Generative Bayesian approach:
  - Assume  $p(\mathbf{x}, y|\Theta)$  and  $p(\Theta)$
  - Get  $p(\Theta|\mathcal{D}) \propto \prod_{t=1}^T p(\mathbf{x}_t, y_t|\Theta)p(\Theta)$
  - Predict  $\hat{y} = \arg \max_y \int_{\Theta} p(\mathbf{x}, y|\Theta)p(\Theta|\mathcal{D})$
- Conditional Bayesian approach:
  - Assume  $p(y|\mathbf{x}, \Theta)$  and  $p(\Theta)$
  - Get  $p(\Theta|\mathcal{D}) \propto \prod_{t=1}^T p(y_t|\mathbf{x}_t, \Theta)p(\Theta)$
  - Predict  $\hat{y} = \arg \max_y \int_{\Theta} p(y|\mathbf{x}, \Theta)p(\Theta|\mathcal{D})$

# Generative Learning

- Given  $\mathcal{D} = (\mathbf{x}_t, y_t)_{t=1}^T$  sampled *iid* from unknown  $P(\mathbf{x}, y)$
- Find rule producing  $\hat{y}$  from  $\mathbf{x}$  with low error
- Generative Bayesian approach:
  - Assume  $p(\mathbf{x}, y|\Theta)$  and  $p(\Theta)$
  - Get  $p(\Theta|\mathcal{D}) \propto \prod_{t=1}^T p(\mathbf{x}_t, y_t|\Theta)p(\Theta)$
  - Predict  $\hat{y} = \arg \max_y \int_{\Theta} p(\mathbf{x}, y|\Theta)p(\Theta|\mathcal{D})$
- Conditional Bayesian approach:
  - Assume  $p(y|\mathbf{x}, \Theta)$  and  $p(\Theta)$
  - Get  $p(\Theta|\mathcal{D}) \propto \prod_{t=1}^T p(y_t|\mathbf{x}_t, \Theta)p(\Theta)$
  - Predict  $\hat{y} = \arg \max_y \int_{\Theta} p(y|\mathbf{x}, \Theta)p(\Theta|\mathcal{D})$
- Problem: high test error when assumptions are wrong!
- Solution: add task constraint that correct  $y_t$  wins by margin  $\gamma$

# Discriminating with Task Constraints

Conditional Bayes,  $\hat{y} = \arg \max_y \int_{\Theta} q(\Theta) p(y|\mathbf{x}, \Theta)$  via

$$\min_{q(\Theta)} \mathcal{KL} \left( q(\Theta) \parallel \frac{1}{Z} \prod_{t=1}^T p(y_t|\mathbf{x}_t, \Theta) p(\Theta) \right)$$

# Discriminating with Task Constraints

Conditional Bayes,  $\hat{y} = \arg \max_y \int_{\Theta} q(\Theta) p(y|\mathbf{x}, \Theta)$  via

$$\min_{q(\Theta)} \mathcal{KL} \left( q(\Theta) \parallel \frac{1}{Z} \prod_{t=1}^T p(y_t|\mathbf{x}_t, \Theta) p(\Theta) \right)$$

$$s.t. \int_{\Theta} q(\Theta) p(y_t|\mathbf{x}_t, \Theta) \geq \max_{y \neq y_t} \int_{\Theta} q(\Theta) p(y|\mathbf{x}_t, \Theta) + \gamma \quad \forall t$$



# Discriminating with Task Constraints

Log Conditional Bayes,  $\hat{y} = \arg \max_y \int_{\Theta} q(\Theta) \ln p(y|\mathbf{x}, \Theta)$  via

$$\min_{q(\Theta)} \mathcal{KL} \left( q(\Theta) \parallel \frac{1}{Z} \prod_{t=1}^T p(y_t|\mathbf{x}_t, \Theta) p(\Theta) \right)$$

$$s.t. \int_{\Theta} q(\Theta) \ln p(y_t|\mathbf{x}_t, \Theta) \geq \max_{y \neq y_t} \int_{\Theta} q(\Theta) \ln p(y|\mathbf{x}_t, \Theta) + \gamma \quad \forall t$$

# Discriminating with Slackened Task Constraints

Log Conditional Bayes,  $\hat{y} = \arg \max_y \int_{\Theta} q(\Theta) \ln p(y|\mathbf{x}, \Theta)$  via

$$\min_{q(\Theta)} \mathcal{KL} \left( q(\Theta) \parallel \frac{1}{Z} \prod_{t=1}^T p(y_t|\mathbf{x}_t, \Theta) p(\Theta) \right) + C \sum_{t=1}^T \xi_t$$

$$\text{s.t. } \int_{\Theta} q(\Theta) \ln p(y_t|\mathbf{x}_t, \Theta) \geq \max_{y \neq y_t} \int_{\Theta} q(\Theta) \ln p(y|\mathbf{x}_t, \Theta) + \gamma - \xi_t \quad \forall t$$

# Discriminating with Slackened Task Constraints

Log Conditional Bayes,  $\hat{y} = \arg \max_y \int_{\Theta} q(\Theta) \ln p(y|\mathbf{x}, \Theta)$  via

$$\min_{q(\Theta)} \mathcal{KL} \left( q(\Theta) \parallel \frac{1}{Z} \prod_{t=1}^T p(y_t|\mathbf{x}_t, \Theta) p(\Theta) \right) + C \sum_{t=1}^T \xi_t$$

$$\text{s.t. } \int_{\Theta} q(\Theta) \ln p(y_t|\mathbf{x}_t, \Theta) \geq \max_{y \neq y_t} \int_{\Theta} q(\Theta) \ln p(y|\mathbf{x}_t, \Theta) + \gamma - \xi_t \quad \forall t$$

Easy, Maximum Entropy Discrimination (Jaakkola Meila Jebara 99)  
Slack allows some misclassification of training data

# Primal and Dual MED

## Primal

$$\min_{q(\Theta)} \mathcal{KL}(q(\Theta) \parallel \hat{p}(\Theta)) + C \sum_{t=1}^T \xi_t$$

$$s.t. \int_{\Theta} q(\Theta) \ln p(y_t | \mathbf{x}_t, \Theta) \geq \max_{y \neq y_t} \int_{\Theta} q(\Theta) \ln p(y | \mathbf{x}_t, \Theta) + \gamma - \xi_t \quad \forall t$$

## Dual

$$\max_{\lambda \in [0, C]} -\ln \int_{\Theta} \hat{p}(\Theta) \prod_{t=1}^T \left( \frac{p(y_t | \mathbf{x}_t, \Theta)}{\max_{y \neq y_t} p(y | \mathbf{x}_t, \Theta)} \right)^{\lambda_t} \exp(-\gamma \lambda_t)$$

$$q(\Theta) = \frac{1}{Z(\lambda)} \hat{p}(\Theta) \prod_{t=1}^T \left( \frac{p(y_t | \mathbf{x}_t, \Theta)}{\max_{y \neq y_t} p(y | \mathbf{x}_t, \Theta)} \right)^{\lambda_t} \exp(-\gamma \lambda_t)$$

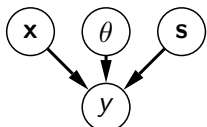
# MED for Support Vector Machines

- Set  $p(y|\mathbf{x}, \Theta) \propto \exp(y/2(\mathbf{x}^\top \theta + b))$  and  $\Theta = \{\theta, b\}$
- Set  $\hat{p}(\Theta) = \mathcal{N}(b|0, \infty)\mathcal{N}(\theta|\mathbf{0}, \mathbf{I})$
- MED dual produces support vector machine optimization

$$\max_{\substack{\lambda \in [0, c] \\ \sum_i y_i \lambda_i = 0}} \gamma \sum_i \lambda_i - \frac{1}{2} \sum_{ij} \lambda_i \lambda_j y_i y_j \mathbf{x}_i^\top \mathbf{x}_j$$

- MED prediction becomes the same
 
$$\hat{y} = \text{sign} \left( \sum_t y_t \lambda_t \mathbf{x}_t^\top \mathbf{x} + b \right)$$
- Solvable in  $\mathcal{O}(T^3)$  with quadratic programming
- Faster solution to  $\epsilon$  accuracy with e.g. Pegasos

## MED for Feature Selection

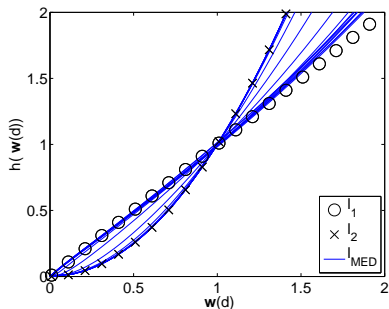


- Model  $\Theta = \{\theta, b, \mathbf{s}\}$  where  $\mathbf{s} \in \mathbb{B}^D$  sparsifies  $\theta \in \mathbb{R}^D$
- Set  $p(y|\mathbf{x}, \Theta, \mathbf{s}) \propto \exp(y/2(\sum_d \mathbf{s}(d)\mathbf{x}(d)\theta(d) + b))$
- Set  $\hat{p}(\Theta) = \mathcal{N}(b|0, \infty)\mathcal{N}(\theta|\mathbf{0}, \mathbf{I}) \prod_d \rho^{\mathbf{s}(d)}(1 - \rho)^{1-\mathbf{s}(d)}$
- Parameter  $\rho$  (or  $\alpha = \frac{1-\rho}{\rho}$ ) is prior % of non-sparse features
- MED dual is

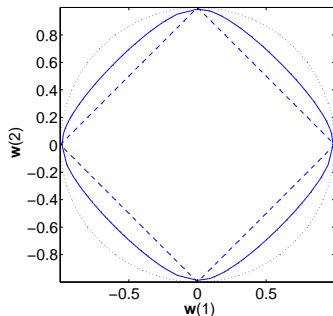
$$\max_{\substack{\lambda \in [0, C] \\ \sum_i y_i \lambda_i = 0}} \gamma \sum_i \lambda_i - \sum_{d=1}^D \ln \left( \alpha + e^{\frac{1}{2} \sum_{ij} y_i y_j \lambda_i \lambda_j \mathbf{x}_i(d) \mathbf{x}_j(d)} \right)$$

- MED prediction is  $\hat{y} = \text{sign} \left( \sum_{di} y_i \lambda_i \mathbf{x}_i(d) \mathbf{x}(d) \hat{\mathbf{s}}(d) + b \right)$   
 where  $\hat{\mathbf{s}}(d) = \left( 1 + \alpha \exp \left( -\frac{1}{2} \sum_{ij} y_i y_j \lambda_i \lambda_j \mathbf{x}_i(d) \mathbf{x}_j(d) \right) \right)^{-1}$

## MED for Feature Selection



(a) One dimensional plot



(b) Two dimensional contour plot

**Figure:** (a) Various norms on the weight vector shown as  $\alpha$  varies from  $\alpha = 0$  which mimics an  $l_2$  norm to  $\alpha$  large which mimics an  $l_1$  norm. (b) a two-dimensional contour plot of the  $l_1$  penalty (dot-dash line), the  $l_2$  penalty (dotted line) and the  $l_{MED}$  penalty with  $\alpha = 2$  (solid line).

## MED for Feature Selection

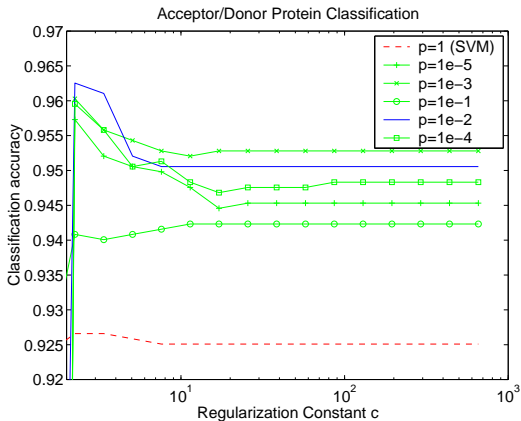
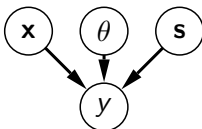


Figure: Acceptor/donor sequence classification accuracy.



# MED for Kernel Selection

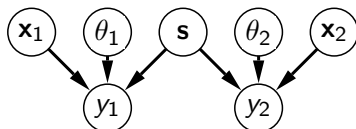


- Select from  $\phi_1(\mathbf{x}), \dots, \phi_D(\mathbf{x})$  mappings to Hilbert space
- Model  $\Theta = \{\theta_1, \dots, \theta_D, b, \mathbf{s}\}$  where  $\mathbf{s} \in \mathbb{B}^D$  selects  $\theta_d \in \mathcal{H}$
- Set  $p(y|\mathbf{x}, \Theta, \mathbf{s}) \propto \exp(y/2(\sum_d \mathbf{s}(d)\theta_d^\top \phi_d(\mathbf{x}) + b))$
- MED dual is

$$\max_{\substack{\lambda \in [0, c] \\ \sum_i y_i \lambda_i = 0}} \gamma \sum_i \lambda_i - \sum_{d=1}^D \ln \left( \alpha + e^{\frac{1}{2} \sum_{ij} y_i y_j \lambda_i \lambda_j k_d(\mathbf{x}_i, \mathbf{x}_j)} \right)$$

- MED prediction is  $\hat{y} = \text{sign} \left( \sum_{di} y_i \lambda_i \hat{\mathbf{s}}(d) k_d(\mathbf{x}, \mathbf{x}_i) + b \right)$  where  $\hat{\mathbf{s}}(d) = \left( 1 + \alpha \exp \left( -\frac{1}{2} \sum_{ij} y_i y_j \lambda_i \lambda_j k_d(\mathbf{x}, \mathbf{x}_i) \right) \right)^{-1}$

## Multi-Task Margin Constraints for Kernel Selection

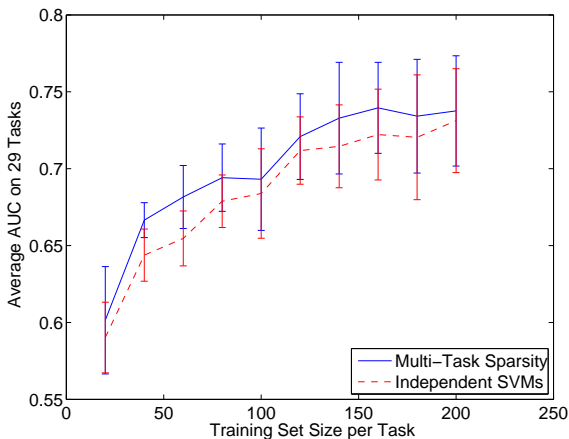


- Have  $M$  models  $\Theta = \{\Theta_1, \dots, \Theta_M, \mathbf{s}\}$  and sparsifier  $\mathbf{s} \in \mathbb{B}^D$
- Each model is  $\Theta_m = \{\theta_{m1}, \dots, \theta_{mD}, b_m\}$
- Set  $p(y|\mathbf{x}, \Theta_m, \mathbf{s}) \propto \exp(y/2(\sum_d \mathbf{s}(d)\theta_{md}^\top \phi_d(\mathbf{x}) + b_m))$
- MED dual is

$$\max_{\substack{\lambda \in [0, C] \\ \sum_i y_{mi} \lambda_{mi} = 0}} \gamma \sum_{mi} \lambda_{mi} - \sum_{d=1}^D \ln \left( \alpha + e^{\frac{1}{2} \sum_{mij} y_{mi} y_{mj} \lambda_{mi} \lambda_{mj} k_d(\mathbf{x}_{mi}, \mathbf{x}_{mj})} \right)$$

- Task  $m$  predicts  $\hat{y} = \text{sign} \left( \sum_{di} y_{mi} \lambda_{mi} \hat{\mathbf{s}}(d) k_d(\mathbf{x}, \mathbf{x}_{mi}) + b_m \right)$   
 $\hat{\mathbf{s}}(d) = \left( 1 + \alpha \exp \left( -\frac{1}{2} \sum_m \sum_{ij} y_{mi} y_{mj} \lambda_{mi} \lambda_{mj} k_d(\mathbf{x}_{mi}, \mathbf{x}_{mj}) \right) \right)^{-1}$

# Multi-Task Margin Constraints for Kernel Selection



**Figure:** Feature and RBF kernel selection on the Landmine dataset (Xue Liao Carin Krishnapauram 07). Values for  $C$  and  $\alpha$  obtained by cross-validation on held out data.

# Sequential Quadratic Programming

- How to optimize MED when it's not a QP? For example,  $\max_{\lambda \in [0, C], \sum_i y_i \lambda_i = 0} \gamma \sum_i \lambda_i - \sum_{d=1}^D \ln(\alpha + e^{\lambda^\top H \lambda})$
- Lower bound – In terms with a quadratic, get sequential QP

## Theorem (Jebara 11)

For any  $\mathbf{v}$ , the term  $-\ln(\alpha + e^{\mathbf{u}^\top \mathbf{u}})$  is lower bounded by

$$-\ln(\alpha + e^{\mathbf{v}^\top \mathbf{v}}) - 2 \frac{\mathbf{v}^\top (\mathbf{u} - \mathbf{v})}{\alpha e^{-\mathbf{v}^\top \mathbf{v}} + 1} - (\mathbf{u} - \mathbf{v})^\top (\mathcal{G} \mathbf{v} \mathbf{v}^\top + I) (\mathbf{u} - \mathbf{v})$$

when  $\mathcal{G} \geq \frac{\tanh(\frac{1}{2} \ln(\alpha \exp(-\mathbf{v}^\top \mathbf{v})))}{\ln(\alpha \exp(-\mathbf{v}^\top \mathbf{v}))}$ .

# Sequential Quadratic Programming

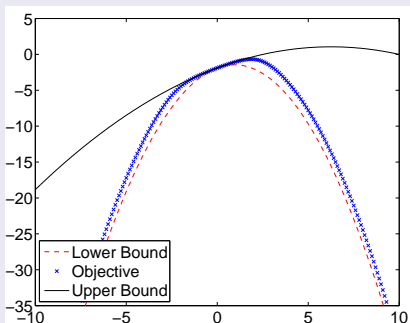
0	Input: dataset $\mathcal{D}$ , $C > 0$ , $\alpha \geq 0$ , $0 < \epsilon < 1$
1	Initialize Lagrange multipliers to zero, $\lambda = \mathbf{0}$ .
2	Store $\tilde{\lambda} = \lambda$ .
3	Apply bound on all $-\ln\left(\alpha + e^{\lambda^\top H \lambda}\right)$ terms at $\tilde{\lambda}$ .
4	Solve resulting (fast) SVM problem to get $\lambda$ .
5	If $\ \lambda - \tilde{\lambda}\  > \epsilon \ \lambda\ $ go to 2.
6	Output $\lambda$ .

## Theorem (Jebara 11)

*SQP achieves  $(1 - \epsilon)J(\lambda^*)$  within  $\left\lceil \frac{\log(1/\epsilon)}{\log(\min(1 + \frac{1}{\alpha}, 2))} \right\rceil$  iterations.  
 Sparse multi-task is a constant factor more work than SVMs.*

# Sequential Quadratic Programming

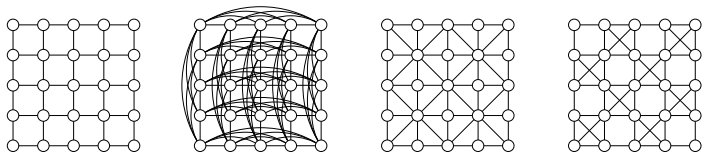
Proof.



**Figure:** Quadratic bounding sandwich. Compare upper and lower bound curvatures to bound maximum # of iterations.



# Graphical Model Structure Estimation



- Consider an Ising model with distribution  $p(\mathbf{x})$  proportional to

$$\exp \left( \sum_{m=1}^D \eta(m) \mathbf{x}(m) + \sum_{m=1}^D \sum_{n=1}^D \mathbf{s}(m, n) \theta(m, n) \mathbf{x}(m) \mathbf{x}(n) \right)$$

- $\mathbf{s} \in \mathbb{B}^{D \times D}$  a symmetric matrix with zero on its diagonal
- $\theta \in \mathbb{R}^{D \times D}$  is a symmetric matrix with zero on its diagonal
- $\eta \in \mathbb{R}^D$  a vector
- Goal: given  $T$  binary vectors  $\mathbf{x}_1, \dots, \mathbf{x}_T$  where  $\mathbf{x}_t \in \mathbb{B}^D$  sampled *iid* from an unknown  $p(\mathbf{x})$  find  $\hat{\mathbf{s}}$

# Graphical Model Structure Estimation

- Recall  $\ell_1$  method of Wainwright Ravikumar Lafferty 07
- Recover  $\hat{\mathbf{s}}$  by solving  $D$   $\ell_1$  sparse regressions

$$\min_{\theta \in \mathbb{R}^D} \nu \sum_{d \neq m} |\theta(d)| + \sum_{t=1}^T \log(1 + e^{\sum_{d \neq m} \theta(d) \mathbf{x}_t(d) + \theta(m)} - \mathbf{x}_t(m) (\sum_{d \neq m} \theta(d) \mathbf{x}_t(d) + \theta(m)))$$

- Learn to regress each dimension from all others
- Symmetrize by  $\hat{\mathbf{s}}(m, n) = [\hat{\theta}_m(n) \hat{\theta}_n(m)]$  (AND operation)
- Symmetrize by  $\hat{\mathbf{s}}(m, n) = [\hat{\theta}_m(n) + \hat{\theta}_n(m)]$  (OR operation)
- Both have nice asymptotic consistency properties



# Graphical Model Structure Estimation

- How to symmetrize *during* learning in Wainwright et al. ?
- MED approach: assume the following predictive distribution

$$p(y|m, \mathbf{x}, \theta, b, \mathbf{s}) \propto \exp\left(\frac{y}{2} \sum_{d \neq m} \mathbf{s}(m, d) \mathbf{x}(d) \theta(m, d) + b(m)\right)$$

- Use sparse symmetric Bernoulli prior on  $\mathbf{s}$
- Use Gaussian prior on  $\theta$  allowing it to be asymmetric
- MED's dual is a single joint constrained convex optimization

$$\max_{\lambda \in \Lambda} \lambda^\top \mathbf{1} - \sum_{m=1}^D \sum_{d=m+1}^D \log \left( \alpha + e^{\frac{1}{2} \lambda^\top H_{m,d} \lambda} \right)$$

- Solve with bound and sequential quadratic programming

# Graphical Model Structure Estimation

- Symmetrizing the graph structure via MED helps!

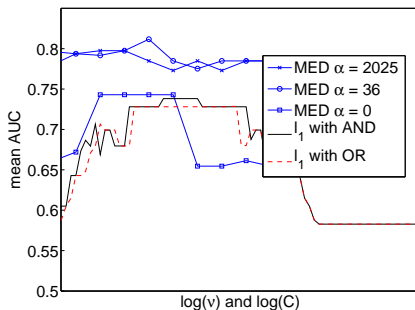


Figure: Graphical model structure recovery from synthetic experiments

# Relative Margin Constraints

- Why stop with just  $\gamma$  classification constraints?
- Limit spread between correct to weakest by  $\beta$  constraint
- Relative margin machines (Shivaswamy & Jebara 10)

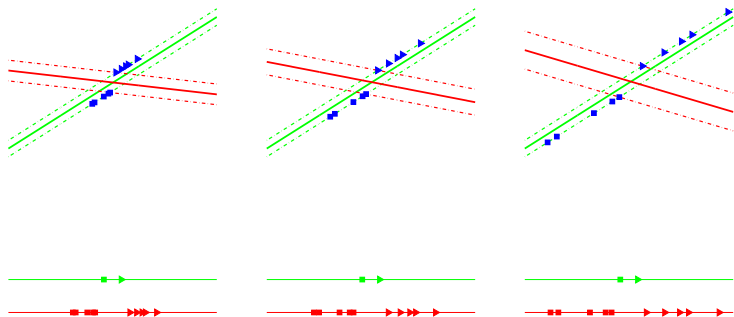
$$\min_{q(\Theta)} \mathcal{KL}(q(\Theta) \parallel \hat{p}(\Theta))$$

$$s.t. \int_{\Theta} q(\Theta) \ln p(y_t | \mathbf{x}_t, \Theta) \geq \max_{y \neq y_t} \int_{\Theta} q(\Theta) \ln p(y | \mathbf{x}_t, \Theta) + \gamma \quad \forall t$$

$$s.t. \int_{\Theta} q(\Theta) \ln p(y_t | \mathbf{x}_t, \Theta) \leq \min_{y \neq y_t} \int_{\Theta} q(\Theta) \ln p(y | \mathbf{x}_t, \Theta) + \beta \quad \forall t$$

- The SVM optimization (with slack omitted) then becomes  $\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2$  subject to  $\beta \geq y_i(\mathbf{w}^\top \mathbf{x}_i + b) \geq \gamma$

# Relative Margin Constraints



**Figure:** Top: As the data is scaled, the SVM (red or dark shade) deviates from the RMM (green or light shade). Bottom: The projections of the examples (that is  $\mathbf{w}^\top \mathbf{x} + b$ ) on the real line for the SVM and the RMM.

## Relative Margin Constraints in Binary Classification

Dataset	SVM	KFDA	$\Sigma$ -SVM	RMM (C=D)	RMM
banana	10.5±0.4	10.8±0.5	10.5±0.4	<b>10.4±0.4</b>	<b>10.4±0.4*</b>
b.cancer	<b>25.3±4.6*</b>	26.6±4.8	28.8±4.6	25.9±4.5	27.2±4.8
diabetes	<b>23.1±1.7</b>	<b>23.2±1.8</b>	24.2±1.9	<b>23.1±1.7</b>	<b>23.0±1.7*</b>
f.solar	<b>32.3±1.8</b>	33.1±1.6	34.6±2.0	<b>32.3±1.8*</b>	33.1±2.5
german	<b>23.4±2.2</b>	24.1±2.4	25.9±2.4	<b>23.4±2.1</b>	<b>23.2±2.2*</b>
heart	15.5±3.3	15.7±3.2	19.9±3.6	15.4±3.3	<b>15.2±3.1*</b>
image	<b>3.0±0.6</b>	3.1±0.6	3.3±0.7	3.0±0.6	<b>2.9±0.7</b>
ringnorm	1.5±0.1	<b>1.5±0.1</b>	1.5±0.1	<b>1.5±0.1</b>	<b>1.5±0.1*</b>
splice	10.9±0.7	10.6±0.7	10.8±0.6	10.8±0.6	10.8±0.6
thyroid	4.7±2.1	<b>4.2±2.1</b>	4.5±2.1	<b>4.2±1.8*</b>	<b>4.2±2.2</b>
titanic	22.3±1.1	<b>22.0±1.3*</b>	23.1±2.2	22.3±1.1	<b>22.2±1.3</b>
twonorm	2.4±0.1	2.4±0.2	2.5±0.2	2.4±0.1	<b>2.3±0.1*</b>
waveform	9.9±0.4	9.9±0.4	10.5±0.5	10.0±0.4	<b>9.7±0.4*</b>

# Relative Margin Constraints in Structured Prediction

MED extended to structured prediction (Zhu & Xing 09)

Add relative margin to structured prediction (Shivaswamy & J 09)

Multi-class classification error

Kernel	StructSVM	StructRMM	p-value
Poly 1	3.78 $\pm$ 0.54	3.85 $\pm$ 0.62	0.55
Poly 2	2.11 $\pm$ 0.43	<b>1.46 <math>\pm</math> 0.34</b>	0.00
Ploy 3	1.73 $\pm$ 0.37	<b>1.24 <math>\pm</math> 0.43</b>	0.00
Poly 4	1.55 $\pm$ 0.45	<b>1.18 <math>\pm</math> 0.43</b>	0.00

Sequence label error (Named Entity Rec. & Part of Speech)

	CRF	StructSVM	StructRMM	p-value
NER	5.13 $\pm$ 0.28	5.09 $\pm$ 0.32	<b>5.05 <math>\pm</math> 0.28</b>	0.07
POS	11.34 $\pm$ 0.64	11.14 $\pm$ 0.60	<b>10.42 <math>\pm</math> 0.47</b>	0.00

# Conclusions

- Added task-relevant margin constraints in generative learning
- We call this Maximum Entropy Discrimination
- Straightforward multi-task feature and kernel selection
- Produces sparse symmetric graphical model structure
- A blend between  $\ell_1$  and  $\ell_2$  norms
- Optimizations via sequential quadratic programming
- Only constant time more work than SVM
- Relative margin constraints yield further improvements