

MULTI-VIEW LEARNING OVER STRUCTURED AND NON-IDENTICAL OUTPUTS

Kuzman Ganchev ¹ João V. Graça ² John Blitzer ¹ Ben Taskar ¹

¹Computer & Information Science
University of Pennsylvania

²INESC-ID
Lisboa, Portugal

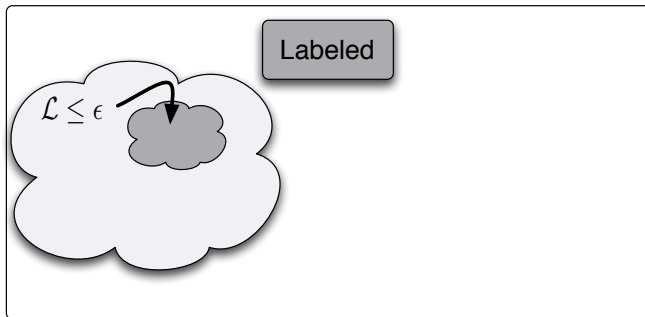
July 11, 2008

SUPERVISED LEARNING



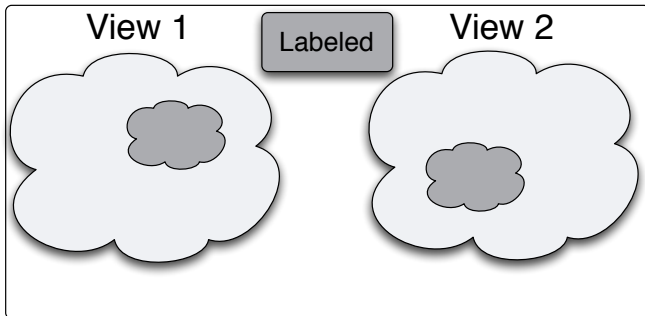
- We have a hypothesis class

SUPERVISED LEARNING

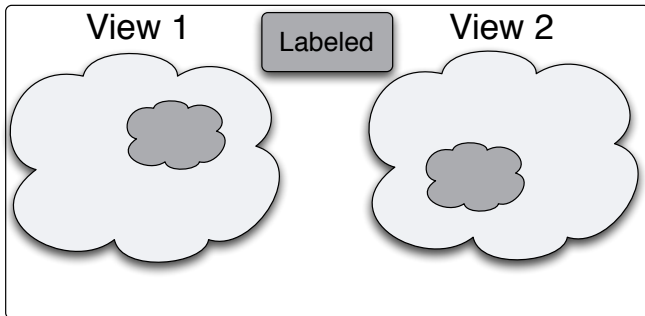


- We have a hypothesis class
labeled data to choose hypothesis

TWO VIEW LEARNING

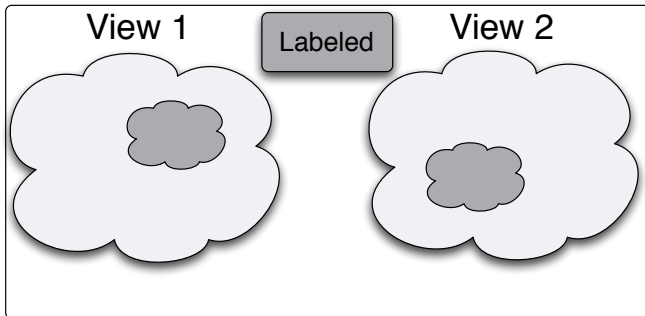


TWO VIEW LEARNING



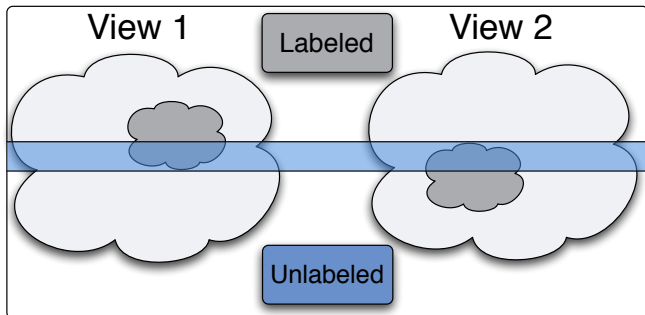
- each view performs well alone

TWO VIEW LEARNING



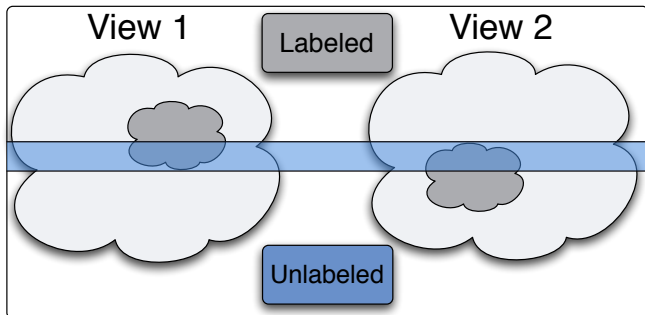
- each view performs well alone
 ⇒ correct models should agree on unlabeled data

TWO VIEW LEARNING



- each view performs well alone
 - ⇒ correct models should agree on unlabeled data
- views don't share too much extra information
 - ⇒ can further reduce hypothesis space

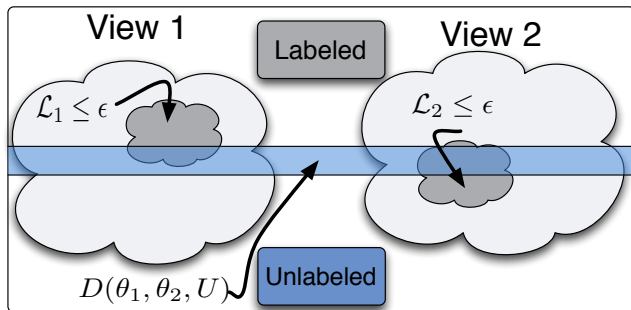
TWO VIEW LEARNING



- each view performs well alone
 \implies correct models should agree on unlabeled data
- views don't share too much extra information
 \implies can further reduce hypothesis space

Assumptions: (Blum & Mitchell, 1998; Balkan & Blum, 2006; Kakade and Foster, 2007)

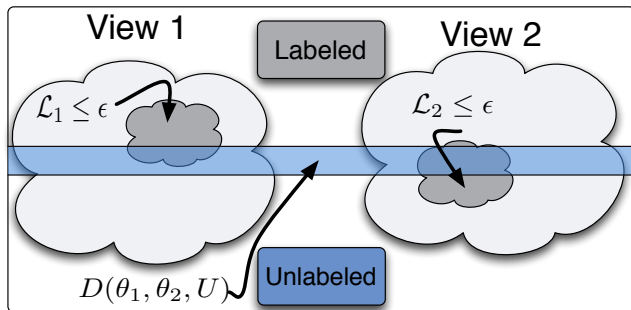
HOW TO LEARN MODELS THAT AGREE



$$\min_{\theta_1, \theta_2} \mathcal{L}(\theta_1, L) + \mathcal{L}(\theta_2, L) + D(\theta_1, \theta_2, U)$$

- Learning probabilistic classifiers

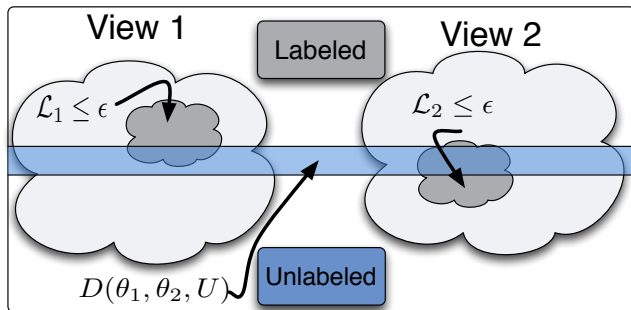
HOW TO LEARN MODELS THAT AGREE



$$\min_{\theta_1, \theta_2} \mathcal{L}(\theta_1, L) + \mathcal{L}(\theta_2, L) + D(\theta_1, \theta_2, U)$$

- Learning probabilistic classifiers
- \mathcal{L} : log-loss on labeled data L

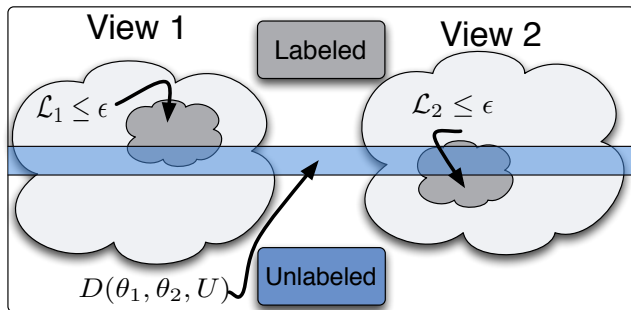
HOW TO LEARN MODELS THAT AGREE



$$\min_{\theta_1, \theta_2} \mathcal{L}(\theta_1, L) + \mathcal{L}(\theta_2, L) + D(\theta_1, \theta_2, U)$$

- Learning probabilistic classifiers
- \mathcal{L} : log-loss on labeled data L
- $\theta = \theta_1, \theta_2$: model parameters

HOW TO LEARN MODELS THAT AGREE



$$\min_{\theta_1, \theta_2} \mathcal{L}(\theta_1, L) + \mathcal{L}(\theta_2, L) + D(\theta_1, \theta_2, U)$$

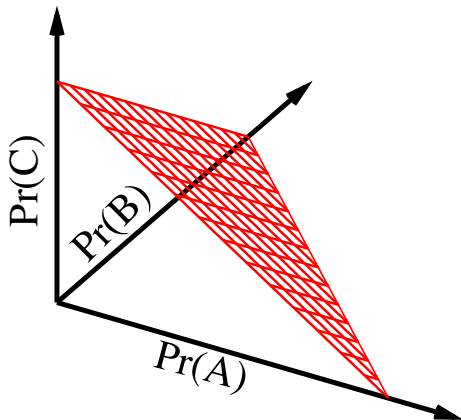
- Learning probabilistic classifiers
- \mathcal{L} : log-loss on labeled data L
- $\theta = \theta_1, \theta_2$: model parameters
- D : co-regularizer (encouraging agreement on unlabeled data U)

CO-REGULARIZER

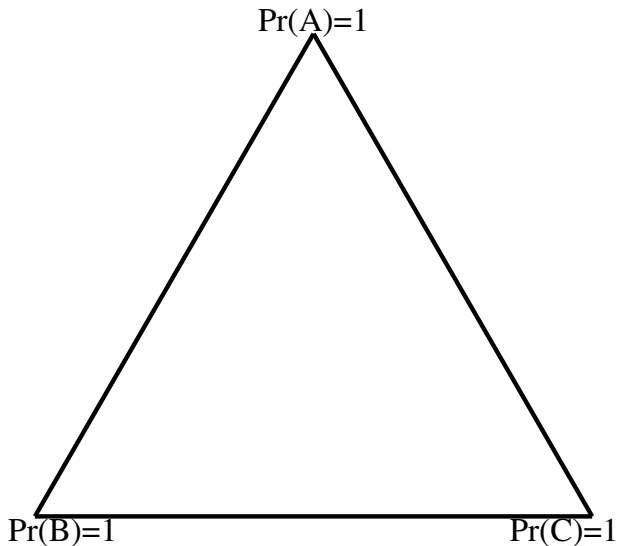
The coregularizer $D...$

- Based on KL distance to a consensus $q = \mathbf{agree}(p_1, p_2)$
- p_i is distribution given by model i
- Illustrative to think in terms of consensus q

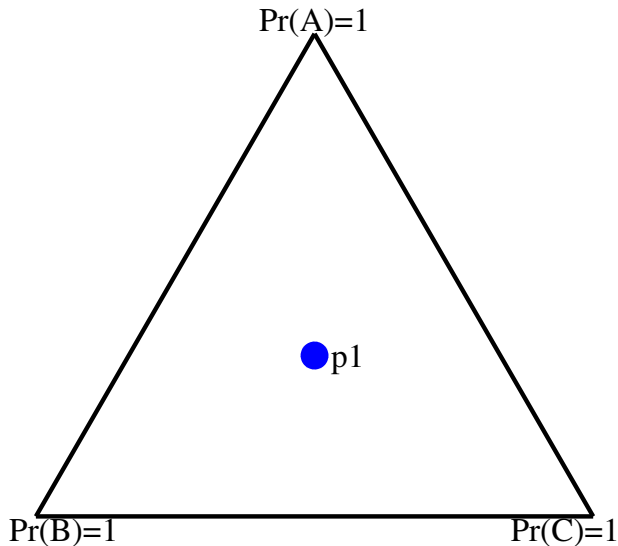
PROBABILISTIC COREGULARIZATION



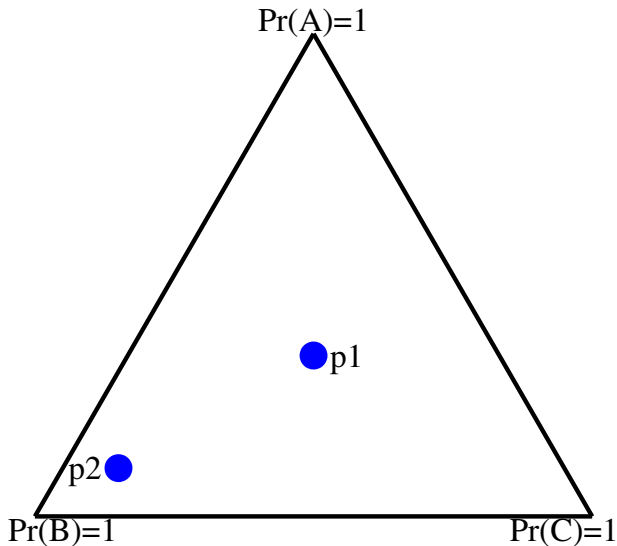
PROBABILISTIC COREGULARIZATION



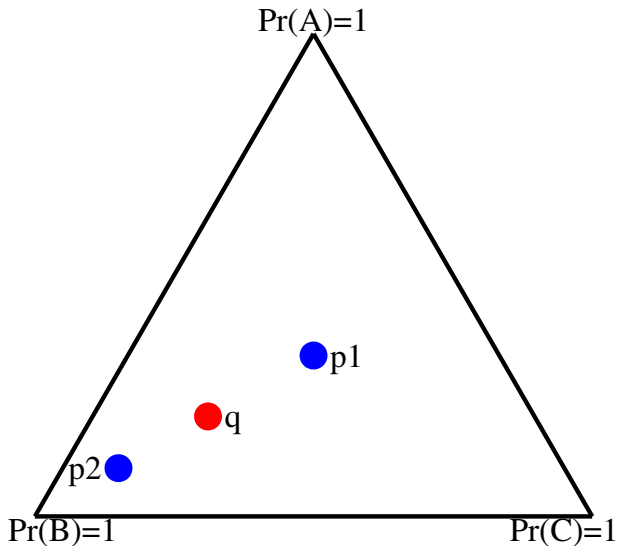
PROBABILISTIC COREGULARIZATION



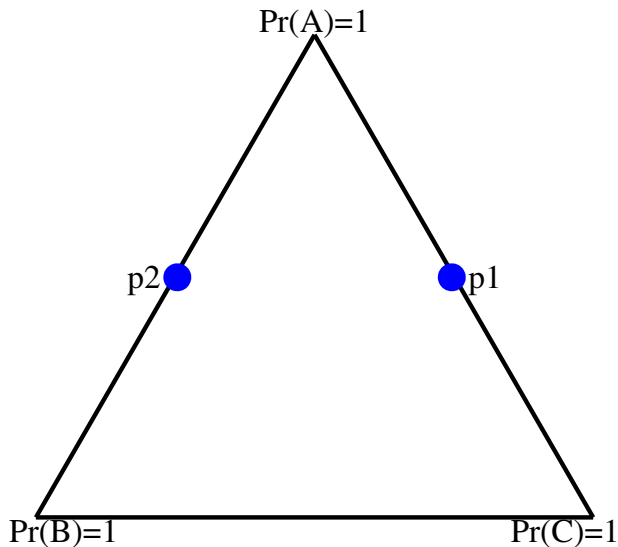
PROBABILISTIC COREGULARIZATION



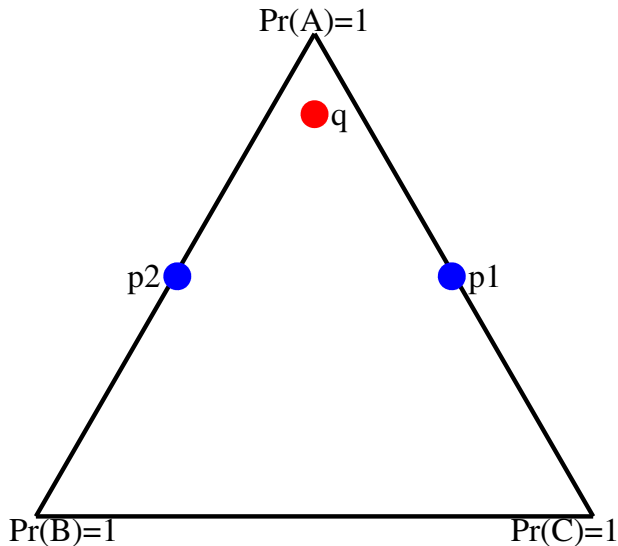
PROBABILISTIC COREGULARIZATION



PROBABILISTIC COREGULARIZATION



PROBABILISTIC COREGULARIZATION



OUR AGREE FUNCTION

$$\mathbf{agree}(p_1, p_2) = \arg \min_q \text{KL}(q \parallel p_1) + \text{KL}(q \parallel p_2)$$

OUR AGREE FUNCTION

$$\mathbf{agree}(p_1, p_2) = \arg \min_q \text{KL}(q \parallel p_1) + \text{KL}(q \parallel p_2)$$

THEOREM: $\mathbf{agree}(p_1, p_2) \propto \sqrt{p_1 \times p_2}$

ALGORITHM

- 1: $\theta_i \leftarrow \min_{\theta} \mathcal{L}_i(\theta)$
- 2: **for** n iterations **do**
- 3: $q(y_1|\mathbf{x}) \leftarrow \mathbf{agree}(p_1(y_1|x), p_2(y_2|x)) \quad \forall x \in U$
- 4: $\theta_i \leftarrow \min_{\theta} \mathcal{L}_i(\theta) - c \mathbf{E}_{x,y \sim U,q} [\log p_i(y_i|x; \theta)]$
- 5: **end for**

ALGORITHM

- 1: $\theta_i \leftarrow \min_{\theta} \mathcal{L}_i(\theta)$
- 2: **for** n iterations **do**
- 3: $q(y_1|\mathbf{x}) \leftarrow \mathbf{agree}(p_1(y_1|x), p_2(y_2|x)) \quad \forall x \in U$
- 4: $\theta_i \leftarrow \min_{\theta} \mathcal{L}_i(\theta) - c \mathbf{E}_{x,y \sim U,q} [\log p_i(y_i|x; \theta)]$
- 5: **end for**

ALGORITHM

- 1: $\theta_i \leftarrow \min_{\theta} \mathcal{L}_i(\theta)$
- 2: **for** n iterations **do**
- 3: $q(y_1|\mathbf{x}) \leftarrow \mathbf{agree}(p_1(y_1|x), p_2(y_2|x)) \quad \forall x \in U$
- 4: $\theta_i \leftarrow \min_{\theta} \mathcal{L}_i(\theta) - c \mathbf{E}_{x, y \sim U, q} [\log p_i(y_i|x; \theta)]$
- 5: **end for**

ALGORITHM

- 1: $\theta_i \leftarrow \min_{\theta} \mathcal{L}_i(\theta)$
- 2: **for** n iterations **do**
- 3: $q(y_1|\mathbf{x}) \leftarrow \mathbf{agree}(p_1(y_1|x), p_2(y_2|x)) \quad \forall x \in U$
- 4: $\theta_i \leftarrow \min_{\theta} \mathcal{L}_i(\theta) - c \mathbf{E}_{x,y \sim U,q} [\log p_i(y_i|x; \theta)]$
- 5: **end for**

ALGORITHM

- 1: $\theta_i \leftarrow \min_{\theta} \mathcal{L}_i(\theta)$
- 2: **for** n iterations **do**
- 3: $q(y_1|\mathbf{x}) \leftarrow \mathbf{agree}(p_1(y_1|x), p_2(y_2|x)) \quad \forall x \in U$
- 4: $\theta_i \leftarrow \min_{\theta} \mathcal{L}_i(\theta) - c \mathbf{E}_{x, y \sim U, q} [\log p_i(y_i|x; \theta)]$
- 5: **end for**

THEOREM: this minimizes co-regularized loss:

$$\mathcal{L}_1(\theta) + \mathcal{L}_2(\theta) + c \mathbf{E}_U \left[\min_q \text{KL}(q \| p_1) + \text{KL}(q \| p_1) \right] .$$

ALGORITHM

THEOREM: this minimizes co-regularized loss:

$$\mathcal{L}_1(\theta) + \mathcal{L}_2(\theta) + c \mathbf{E}_U \left[\min_q \text{KL}(q \| p_1) + \text{KL}(q \| p_1) \right] .$$

ALGORITHM

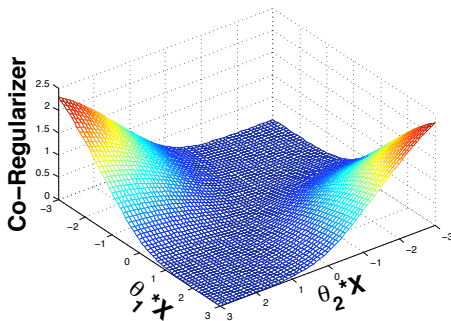
THEOREM: this minimizes co-regularized loss:

$$\begin{aligned} & \mathcal{L}_1(\theta) + \mathcal{L}_2(\theta) + c \mathbf{E}_U \left[\min_q \text{KL}(q \| p_1) + \text{KL}(q \| p_2) \right] . \\ &= \mathcal{L}_1(\theta) + \mathcal{L}_2(\theta) + c \mathbf{E}_U \left[-\log \sum_y \sqrt{p(y; \theta_1)p(y; \theta_2)} \right] . \end{aligned}$$

Bhattacharyya distance

LINEAR MODEL COREGULARIZER

Stochastic Agreement Regularizer



- log-linear models:

$$p_i(1) \propto \exp(\theta_i \cdot x)$$

$$p_i(-1) \propto \exp(-\theta_i \cdot x)$$

OTHER APPROACHES

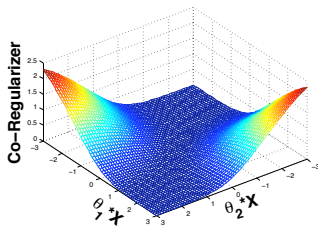
- CoBoosting (Collins and Singer, 1999),
CoPerceptron (Brefeld et al., 2005)

OTHER APPROACHES

- CoBoosting (Collins and Singer, 1999),
CoPerceptron (Brefeld et al., 2005)
 - Different regularized loss functions

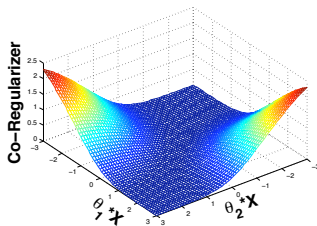
DIFFERENT LOSS FUNCTIONS

Stochastic Agreement Regularizer

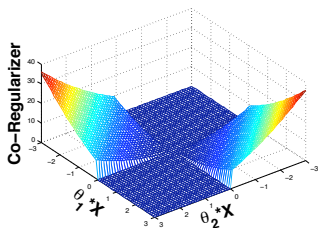


DIFFERENT LOSS FUNCTIONS

Stochastic Agreement Regularizer

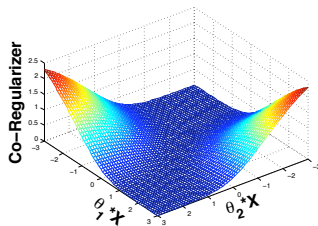


CoPerceptron

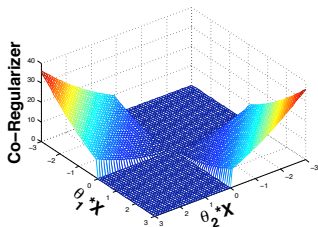


DIFFERENT LOSS FUNCTIONS

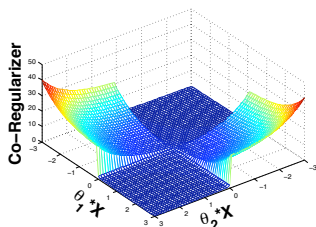
Stochastic Agreement Regularizer



CoPerceptron



CoBoosting



OTHER APPROACHES

- CoBoosting (Collins and Singer, 1999),
CoPerceptron (Brefeld et al., 2005)
 - Different regularized loss functions

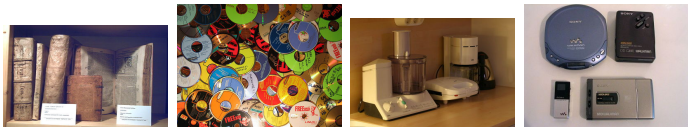
OTHER APPROACHES

- CoBoosting (Collins and Singer, 1999),
CoPerceptron (Brefeld et al., 2005)
 - Different regularized loss functions
 - Hard assignment on unlabeled data

OTHER APPROACHES

- CoBoosting (Collins and Singer, 1999),
CoPerceptron (Brefeld et al., 2005)
 - Different regularized loss functions
 - Hard assignment on unlabeled data
- Many others (Blum & Mitchell, 1998; Sindhwani et al., 2005;
Kakade & Foster, 2007; Suzuki et al., 2007)

SENTIMENT CLASSIFICATION - DOMAIN ADAPTATION (BLITZER ET AL, 2007)



- Product reviews from Amazon.com
 - Books, DVDs, Kitchen Appliances, Electronics
 - 2000 labeled, 3000 - 6000 unlabeled reviews per domain
- Binary classification problem
 - Positive if 4 stars or more, negative if 2 or less
- Transfer learning task
- Views: random split of features



SENTIMENT CLASSIFICATION

Domains	MIRA	SCL	CoBoost	CoPerc	SAR
books→dvds	77.2				
dvds→books	72.8				
books→electr	70.8				
electr→books	70.7				
books→kitchn	74.5				
kitchn→books	70.9				
dvds→electr	73.0				
electr→dvds	70.6				
dvds→kitchn	74.0				
kitchn→dvds	72.7				
electr→kitchn	84.0				
kitchn→electr	82.7				
Total					



SENTIMENT CLASSIFICATION

Domains	MIRA	SCL	CoBoost	CoPerc	SAR
books→dvds	77.2	-1.4			
dvds→books	72.8	6.9			
books→electr	70.8	5.1			
electr→books	70.7	4.7			
books→kitchn	74.5	4.4			
kitchn→books	70.9	-2.3			
dvds→electr	73.0	1.1			
electr→dvds	70.6	5.6			
dvds→kitchn	74.0	7.4			
kitchn→dvds	72.7	4.2			
electr→kitchn	84.0	1.9			
kitchn→electr	82.7	4.1			
Total		4			



SENTIMENT CLASSIFICATION

Domains	MIRA	SCL	CoBoost	CoPerc	SAR
books→dvds	77.2	-1.4	1.6		
dvds→books	72.8	6.9	7.0		
books→electr	70.8	5.1	6.2		
electr→books	70.7	4.7	0.3		
books→kitchn	74.5	4.4	3.5		
kitchn→books	70.9	-2.3	-1.1		
dvds→electr	73.0	1.1	2.3		
electr→dvds	70.6	5.6	2.9		
dvds→kitchn	74.0	7.4	5.0		
kitchn→dvds	72.7	4.2	-2.6		
electr→kitchn	84.0	1.9	1.0		
kitchn→electr	82.7	4.1	0.3		
Total		4	1		



SENTIMENT CLASSIFICATION

Domains	MIRA	SCL	CoBoost	CoPerc	SAR
books→dvds	77.2	-1.4	1.6	-1.7	
dvds→books	72.8	6.9	7.0	1.7	
books→electr	70.8	5.1	6.2	-1.5	
electr→books	70.7	4.7	0.3	-3.2	
books→kitchn	74.5	4.4	3.5	2.0	
kitchn→books	70.9	-2.3	-1.1	-4.3	
dvds→electr	73.0	1.1	2.3	-1.8	
electr→dvds	70.6	5.6	2.9	-7.3	
dvds→kitchn	74.0	7.4	5.0	4.3	
kitchn→dvds	72.7	4.2	-2.6	-12.2	
electr→kitchn	84.0	1.9	1.0	-0.7	
kitchn→electr	82.7	4.1	0.3	-2.2	
Total		4	1	0	



SENTIMENT CLASSIFICATION

Domains	MIRA	SCL	CoBoost	CoPerc	SAR
books→dvds	77.2	-1.4	1.6	-1.7	2.6
dvds→books	72.8	6.9	7.0	1.7	8.5
books→electr	70.8	5.1	6.2	-1.5	4.7
electr→books	70.7	4.7	0.3	-3.2	3.6
books→kitchn	74.5	4.4	3.5	2.0	6.5
kitchn→books	70.9	-2.3	-1.1	-4.3	1.9
dvds→electr	73.0	1.1	2.3	-1.8	3.5
electr→dvds	70.6	5.6	2.9	-7.3	2.4
dvds→kitchn	74.0	7.4	5.0	4.3	8.8
kitchn→dvds	72.7	4.2	-2.6	-12.2	0.1
electr→kitchn	84.0	1.9	1.0	-0.7	1.8
kitchn→electr	82.7	4.1	0.3	-2.2	2.8
Total		4	1	0	6



SENTIMENT CLASSIFICATION

Domains	MIRA	SCL	CoBoost	CoPerc	SAR
books→dvds	77.2	-1.4	1.6	-1.7	2.6
dvds→books	72.8	6.9	7.0	1.7	8.5
books→electr	70.8	5.1	6.2	-1.5	4.7
electr→books	70.7	4.7	0.3	-3.2	3.6
books→kitchn	74.5	4.4	3.5	2.0	6.5
kitchn→books	70.9	-2.3	-1.1	-4.3	1.9
dvds→electr	73.0	1.1	2.3	-1.8	3.5
electr→dvds	70.6	5.6	2.9	-7.3	2.4
dvds→kitchn	74.0	7.4	5.0	4.3	8.8
kitchn→dvds	72.7	4.2	-2.6	-12.2	0.1
electr→kitchn	84.0	1.9	1.0	-0.7	1.8
kitchn→electr	82.7	4.1	0.3	-2.2	2.8
Total		4	1	0	6

NAMED ENTITY DISAMBIGUATION



- Classification of CoNLL 2003 named entites:
 - Person, location, organization, miscellaneous

NAMED ENTITY DISAMBIGUATION



- Classification of CoNLL 2003 named entites:
 - Person, location, organization, miscellaneous
- View 1 - Content
 - Features look only inside named entity
- View 2 - Context
 - Features look only outside named entity

NAMED ENTITY DISAMBIGUATION

Data size	mx-ent	SAR (RRE)
500	74.4	
1000	80.0	
2000	83.4	

- Prior variance and unlabeled weigh choose by cross-validation

NAMED ENTITY DISAMBIGUATION

Data size	mx-ent	SAR (RRE)
500	74.4	76.4 (9.2%)
1000	80.0	81.7 (8.5%)
2000	83.4	84.8 (8.4%)

- Prior variance and unlabeled weigh choose by cross-validation

HOW TO GENERALIZE TWO VIEW IDEA

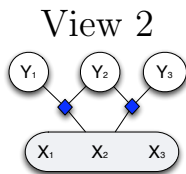
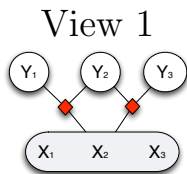
- Structured Output
- Partial Agreement Scenarios
- Both

STRUCTURED OUTPUT

$$\mathbf{agree}(p_1, p_2) = \arg \min_q \text{KL}(q \parallel p_1) + \text{KL}(q \parallel p_2)$$

STRUCTURED OUTPUT

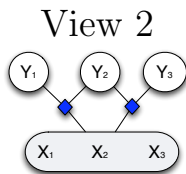
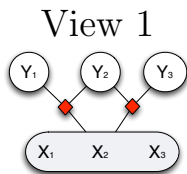
$$\mathbf{agree}(p_1, p_2) = \arg \min_q \text{KL}(q \parallel p_1) + \text{KL}(q \parallel p_2)$$



$$p_1(y \mid x) \propto \prod_c \phi_1(y_c, x) \quad p_2(y \mid x) \propto \prod_c \phi_2(y_c, x)$$

STRUCTURED OUTPUT

$$\mathbf{agree}(p_1, p_2) = \arg \min_q \text{KL}(q \parallel p_1) + \text{KL}(q \parallel p_2)$$



$$p_1(y \mid x) \propto \prod_c \phi_1(y_c, x) \quad p_2(y \mid x) \propto \prod_c \phi_2(y_c, x)$$

THEOREM: $q_i(y) \propto \prod_c \sqrt{\phi_1(y_c, x) \phi_2(y_c, x)}$

STRUCTURED PREDICTION

Small experiments on structured task.

- English NP-chunking from CoNLL 2000
- 500 sentences test data
- views are:
 - current word and POS tag
 - previous/next word and POS tag

STRUCTURED PREDICTION

size	Perc	coPerc	CRF	SAR(RRE)
10	69.4			
20	74.4			
50	80.1			
100	86.1			
200	89.3			
500	90.8			
1000	91.5			

STRUCTURED PREDICTION

size	Perc	coPerc	CRF	SAR(RRE)
10	69.4	71.2		
20	74.4	76.8		
50	80.1	84.1		
100	86.1	88.1		
200	89.3	89.7		
500	90.8	90.9		
1000	91.5	91.8		

STRUCTURED PREDICTION

size	Perc	coPerc	CRF	SAR(RRE)
10	69.4	71.2	73.2	
20	74.4	76.8	79.4	
50	80.1	84.1	86.3	
100	86.1	88.1	88.5	
200	89.3	89.7	89.6	
500	90.8	90.9	91.3	
1000	91.5	91.8	91.6	

STRUCTURED PREDICTION

size	Perc	coPerc	CRF	SAR(RRE)
10	69.4	71.2	73.2	78.2 (19%)
20	74.4	76.8	79.4	84.2 (23%)
50	80.1	84.1	86.3	86.9 (4%)
100	86.1	88.1	88.5	88.9 (3%)
200	89.3	89.7	89.6	89.6 (0%)
500	90.8	90.9	91.3	90.6 (-8%)
1000	91.5	91.8	91.6	91.1 (-6%)

DIFFERENT OUTPUT SPACES

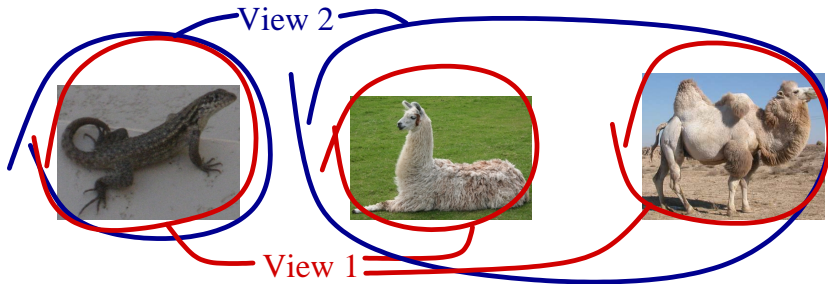


DIFFERENT OUTPUT SPACES



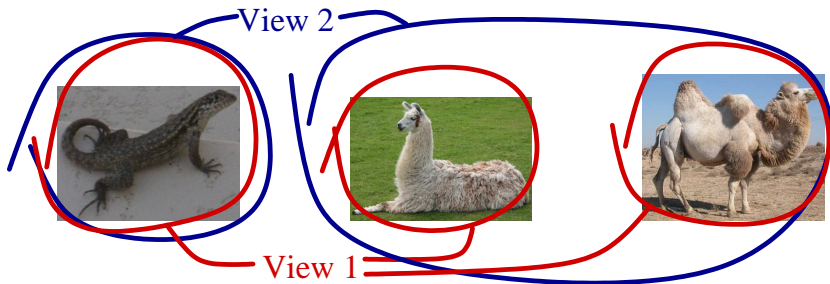
- Different Tag Sets between views

DIFFERENT OUTPUT SPACES



- Different Tag Sets between views

DIFFERENT OUTPUT SPACES



- Different Tag Sets between views
- Partial mapping between labels

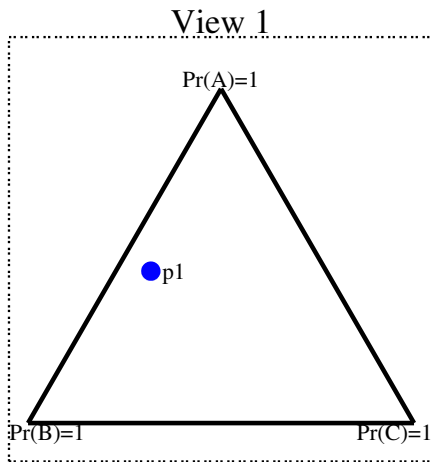
ALGORITHM

- Algorithm stays the same
- Consensus changes:

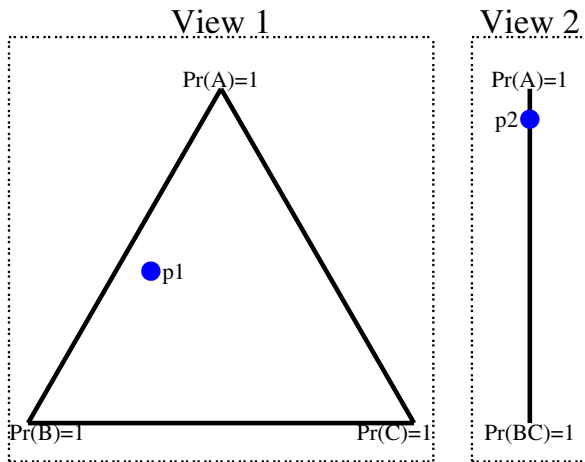
$$\begin{aligned}
 \mathbf{agree}(p_1, p_2) &= \arg \min_q \text{KL}(q(y_1, y_2) || p_1(y_1)p_2(y_2)) \\
 \mathbf{s.t.} \quad &q(z_1, z_2) = q(z_1)q(z_2) \\
 \mathbf{where} \quad &z_i = g(y_i)
 \end{aligned}$$

$g_i(y_i)$: mapping from output of model i to common space

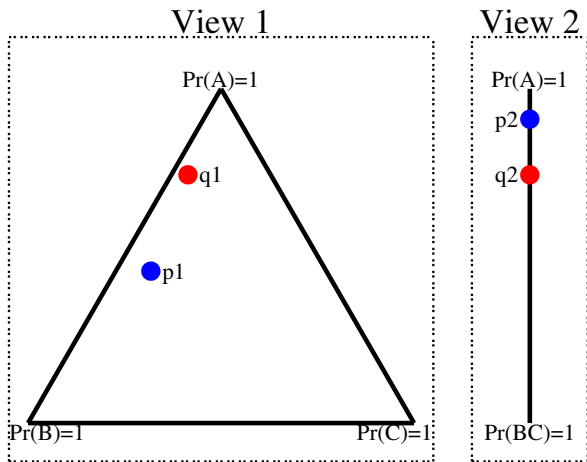
DIFFERENT OUTPUT SPACES



DIFFERENT OUTPUT SPACES



DIFFERENT OUTPUT SPACES



SUMMARY

- New two-view learning algorithm

SUMMARY

- New two-view learning algorithm
- Probabilistic interpretation

SUMMARY

- New two-view learning algorithm
- Probabilistic interpretation
- Generalizes naturally to structured data

SUMMARY

- New two-view learning algorithm
- Probabilistic interpretation
- Generalizes naturally to structured data, partial agreement

SUMMARY

- New two-view learning algorithm
- Probabilistic interpretation
- Generalizes naturally to structured data, partial agreement
- Empirically better than non-smooth alternatives

THANKS!

