

An RKHS for Multi-View Learning and Manifold Co-Regularization

¹Vikas Sindhwani and ²*David S. Rosenberg*

July 7, 2008

¹IBM T. J. Watson Research Center, Yorktown Heights, NY

²Department of Statistics, University of California, Berkeley, CA

Outline

Multi-View Semi-Supervised Learning

An RKHS for L^2 Co-Regularization

Manifold Co-Regularization

Generalization

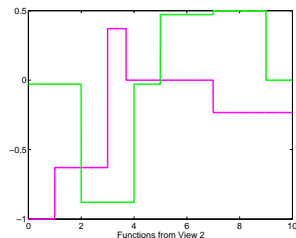
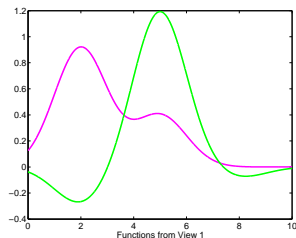
Semi-Supervised Learning

- ▶ Input space \mathcal{X}
- ▶ Output space \mathcal{Y}
- ▶ Distribution $P_{\mathcal{X} \times \mathcal{Y}}$ on $\mathcal{X} \times \mathcal{Y}$
- ▶ Training Data:
 - labeled data : $(x_1, y_1), (x_2, y_2), \dots, (x_\ell, y_\ell)$ i.i.d. $P_{\mathcal{X} \times \mathcal{Y}}$
 - unlabeled data : $x_{\ell+1}, x_{\ell+2}, \dots, x_{\ell+u}$ i.i.d. $P_{\mathcal{X}}$
- ▶ Goal: Predict y given x

Two-View Learning

Basic Framework

- ▶ Two sets of prediction functions (the views): \mathcal{F} and \mathcal{G}



- ▶ Want to find good functions $f^* \in \mathcal{F}$ and $g^* \in \mathcal{G}$
- ▶ Take final prediction to be

$$\varphi^* = \frac{1}{2} (f^* + g^*)$$

Two-View Learning

Motivating Idea

- ▶ Write $\eta(x)$ for the “target” prediction function.

Motivating Idea of Multi-View Learning

- ▶ Given “good” functions $f^* \in \mathcal{F}$ and $g^* \in \mathcal{G}$ s.t.

$$f^*(x) \approx \eta(x) \quad g^*(x) \approx \eta(x),$$

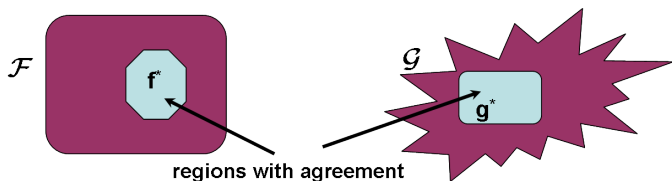
- ▶ then

$$f^*(x) \approx g^*(x)$$

⇒ “All good prediction functions are about equal.”

Two-View Learning

- ▶ Restrict search to “agreement regions”



- ▶ How to find agreement regions?

Two-View Learning with L^2 Co-Regularization

- ▶ RKHS's $(\mathcal{F}, k_{\mathcal{F}})$ and $(\mathcal{G}, k_{\mathcal{G}})$
- ▶ Objective Function

$$\begin{aligned}(f^*, g^*) := \arg \min_{f \in \mathcal{F}, g \in \mathcal{G}} & \hat{L} \left(\frac{f + g}{2} \right) \\ & + \gamma_{\mathcal{F}} \|f\|_{\mathcal{F}}^2 + \gamma_{\mathcal{G}} \|g\|_{\mathcal{G}}^2 \\ & + \lambda \sum_{i=\ell+1}^{\ell+u} |f(x_i) - g(x_i)|^2\end{aligned}$$

- ▶ By Representer theorem,
 - ▶ If $\hat{L}(h) = \sum_{i=1}^{\ell} (h(x_i) - y_i)^2$, solve by minimizing a quadratic
 - ▶ If $\hat{L}(h) = \sum_{i=1}^{\ell} (1 - h(x_i)y_i)_+$, get a QP

Co-Regularization Reformulated

- ▶ Final prediction is

$$h^* = \frac{f^* + g^*}{2}$$

- ▶ Define

$$\tilde{\mathcal{H}} = \left\{ \tilde{h} = \frac{f + g}{2} : f \in \mathcal{F}, g \in \mathcal{G} \right\}$$

- ▶ Can show that

$$h^* = \arg \min_{\tilde{h} \in \tilde{\mathcal{H}}} \left[\hat{L}(\tilde{h}) + \|\tilde{h}\|_{\tilde{\mathcal{H}}}^2 \right]$$

where

$$\|\tilde{h}\|_{\tilde{\mathcal{H}}}^2 := \min_{\substack{f \in \mathcal{F}, g \in \mathcal{G} \\ \text{s.t. } \tilde{h} = \frac{1}{2}(f+g)}} \gamma_{\mathcal{F}} \|f\|_{\mathcal{F}}^2 + \gamma_{\mathcal{G}} \|g\|_{\mathcal{G}}^2 + \lambda \sum_{i=l+1}^{\ell+u} |f(x_i) - g(x_i)|^2$$

The Co-Regularized RKHS

Theorem

- ▶ *There exists an inner product on $\tilde{\mathcal{H}}$ for which $\tilde{\mathcal{H}}$ is an RKHS with norm $\|\cdot\|_{\tilde{\mathcal{H}}}$ (defined earlier).*
- ▶ *The reproducing kernel $\tilde{k} : \mathcal{X} \times \mathcal{X} \rightarrow \mathbf{R}$ is given by*

$$\tilde{k}(x, z) = \gamma_{\mathcal{F}}^{-1} k_{\mathcal{F}}(x, z) + \gamma_{\mathcal{G}}^{-1} k_{\mathcal{G}}(x, z) - \lambda d'_x H d_z$$

where

$$d_x = \gamma_{\mathcal{F}}^{-1} k_{\mathcal{F}}(x, \text{unlab}) - \gamma_{\mathcal{G}}^{-1} k_{\mathcal{G}}(x, \text{unlab})$$

and

$$H = \left(I + \frac{\lambda}{\gamma_{\mathcal{F}}} K_{\mathcal{F}}(\text{unlab}, \text{unlab}) + \frac{\lambda}{\gamma_{\mathcal{G}}} K_{\mathcal{G}}(\text{unlab}, \text{unlab}) \right)^{-1}$$

What have we gained?

- ▶ Can plug-in this “multi-view” kernel to any kernel algorithm
 - ▶ multiview semi-supervised KPCA
 - ▶ multiview semi-supervised ***
- ▶ Reproduce and improve results in [R. & Bartlett, AISTATS'07]
 - ▶ Rademacher complexity bounds follow from standard theory RKHS theory
 - ▶ Localized Rademacher complexity theory improves the generalization bounds

Corollary

Rademacher Complexity of the L^2 Co-Regularized Function Class \mathcal{J}

- ▶ For simplicity, take $\gamma_{\mathcal{F}} = \gamma_{\mathcal{G}} = 1$.

Theorem

For the L^2 co-regularized function class \mathcal{J} ,

$$\frac{1}{\sqrt[4]{2}} \frac{U}{\ell} \leq \hat{R}_{\ell}(\mathcal{J}) \leq \frac{U}{\ell},$$

where

$$U = \sqrt{\text{tr}[K_{\mathcal{F}}(lab, lab)] + \text{tr}[K_{\mathcal{G}}(lab, lab)] - \Delta(\lambda)},$$

and $\Delta(\lambda) \geq 0$ (explicit form given on next slide).

Interpretation of Complexity Reduction Term

Theorem (Continued)

For $\lambda = 0$,

$$\Delta(\lambda) = 0.$$

For $\lambda > 0$,

$$\Delta(\lambda) = \sum_{i=1}^{\ell} d^2 [k_{\mathcal{F}}(x_i, \text{unlab}), k_{\mathcal{G}}(x_i, \text{unlab})],$$

where

- ▶ $d(\cdot, \cdot)$ is a metric on \mathbf{R}^u defined by

$$d^2(c, f) = (c - f)' (\lambda^{-1}I + M)^{-1} (c - f),$$

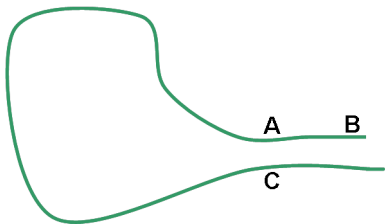
and

$$M := K_{\mathcal{F}}(\text{unlab}, \text{unlab}) + K_{\mathcal{G}}(\text{unlab}, \text{unlab})$$

Manifold Regularization

Geodesic Distance

- ▶ Say data lie on a manifold \mathcal{M} in \mathbf{R}^2



- ▶ In the ambient space \mathbf{R}^2 : $d_{\mathcal{A}}(A, C) < d_{\mathcal{A}}(A, B)$
- ▶ On the manifold, $d_{\mathcal{M}}(A, B) \ll d_{\mathcal{M}}(A, C)$
- ▶ Hypothesis: Good prediction function is smooth w.r.t. “manifold distance”
- ▶ Lots of work with this idea
 - ▶ “Manifold learning”: ISOMAP (Tenebaum et al., *Science* 2000), LLE (Roweis & Saul, *Science* 2000), Laplacian Eigenmaps (Belkin & Niyogi, NIPS 2001)
 - ▶ “Manifold transduction”: Joachims et al. (ICML 2003), Belkin et al. (JLMR 2006), Zhou et al. (NIPS 2004), Sindhwani et al. (IMCL 2005)

Manifold Co-Regularization

The Two Views

- ▶ *Manifold regularization* (Belkin et al., 2006)

$$Q(f) = \hat{L}(f) + \gamma_{\mathcal{A}} \|f\|_{\mathcal{A}}^2 + \gamma_{\mathcal{I}} \mathbf{f}^T M \mathbf{f}$$

- ▶ Ambient View $\mathcal{H}_{\mathcal{A}}$: norm controls smoothness on \mathcal{X}
- ▶ Intrinsic View: controls smoothness w.r.t.
 $\mathcal{G} = \{x_1, \dots, x_{\ell}, x_{\ell+1}, \dots, x_{\ell+u}\}$

- ▶ *Manifold co-regularization*
Find $(f^*, \mathbf{g}^*) \in \mathcal{H}_{\mathcal{A}} \times \mathbf{R}^{\ell+u}$ minimizing

$$Q(f, \mathbf{g}) = \hat{L}(f, \mathbf{g}) + \gamma_{\mathcal{A}} \|f\|_{\mathcal{A}}^2 + \gamma_{\mathcal{I}} \mathbf{g}^T M \mathbf{g} + \lambda \sum_{i=\ell+1}^{\ell+u} [f(x_i) - \mathbf{g}_i]^2$$

Experiment: Manifold Regularization vs CoMR

- ▶ 5 Datasets
 - ▶ LINES
 - ▶ G50C: 50-dim Gaussians s.t. Bayes error is 5%
 - ▶ USPS: 10-class digit recognition
 - ▶ COIL20: 32 x 32 gray scale images of 20 objects, varying angles
 - ▶ PCMAC: from 20 newsgroup dataset

Table: Datasets with d features and c classes. 10 random data splits were created with l labeled, u unlabeled, t test, and v validation examples.

Dataset	d	c	l	u	t	v
LINES	2	2	2	500	250	250
G50C	50	2	50	338	112	50
USPST	256	10	50	1430	477	50
COIL20	241	20	40	1320	40	40
PCMAC	7511	2	50	1385	461	50

Manifold Regularization vs CoMR

Table: Error Rates on Test Data

Dataset	MR	CoMR
LINES	7.7 (1.2)	1.0 (1.5)
G50C	5.8 (2.8)	5.5 (2.3)
USPST	18.2 (1.5)	14.1 (1.6)
COIL20	23.8 (11.1)	14.8 (8.8)
PCMAC	11.9 (3.4)	8.9 (2.6)

Table: $\lambda = 1$ for CoMR. Linear kernel for LINES, RBF for others.

Dataset	nn	σ	p	MR	CoMR
				γ_1, γ_2	γ_1, γ_2
LINES	10	—	1	0.01, 10^{-6}	10^{-4} , 100
G50C	50	17.5	5	1, 100	10, 10
USPST	10	9.4	2	0.01, 0.01	10^{-6} , 10^{-4}
COIL20	2	0.6	1	10^{-4} , 10^{-6}	10^{-6} , 10^{-6}
PCMAC	50	2.7	5	10, 100	1, 10

Generalization of RKHS Theorem

- ▶ RKHS views $\mathcal{H}^1, \dots, \mathcal{H}^m$
- ▶ Arbitrary (fixed) linear combination of views: For $a_1, \dots, a_m \in \mathbf{R}$,

$$\tilde{\mathcal{H}} := \left\{ a_1 f^1 + \dots + a_m f^m \mid f^i \in \mathcal{H}^i, i = 1, \dots, m \right\}$$

- ▶ Arbitrary quadratic “coupling” term.

$$\underset{\varphi \in \tilde{\mathcal{H}}}{\operatorname{argmin}} \quad \underset{\substack{f^1, \dots, f^m \text{ s.t.} \\ \varphi = a_1 f^1 + \dots + a_m f^m}}{\min} \quad L(\varphi) + \sum_{i=1}^m \gamma_i \|f^i\|_{\mathcal{H}^i}^2 + \lambda \underline{\mathbf{f}}^T M \underline{\mathbf{f}} \quad (1)$$

where

$$\underline{\mathbf{f}} := (f^1(x_1), \dots, f^1(x_n), \dots, f^m(x_1), \dots, f^m(x_n))^T$$

Generalization of Theorem

Theorem (Rosenberg et al. 2008)

- ▶ There exists an inner product on $\tilde{\mathcal{H}}$ for which $\tilde{\mathcal{H}}$ is an RKHS with norm

$$\|\varphi\|_{\tilde{\mathcal{H}}}^2 = \min_{\substack{f^1, \dots, f^m \\ \varphi = a_1 f^1 + \dots + a_m f^m}} \sum_{i=1}^m \gamma_i \|f^i\|_{\mathcal{H}^i}^2 + \lambda \underline{\mathbf{f}}^T \underline{\mathbf{M}} \underline{\mathbf{f}} \quad (2)$$

- ▶ The reproducing kernel $\tilde{k} : \mathcal{X} \times \mathcal{X} \rightarrow \mathbf{R}$ is given by

$$\tilde{k}(z, x) = \sum_{j=1}^m \frac{a_j^2}{\gamma_j} k^j(z, x) - \lambda \underline{\mathbf{k}}_x^T \underline{\mathbf{A}} (\mathbf{I} + \lambda \underline{\mathbf{M}} \underline{\mathbf{K}})^{-1} \underline{\mathbf{M}} \underline{\mathbf{A}} \underline{\mathbf{G}}^{-1} \underline{\mathbf{k}}_z,$$

where

$$\begin{aligned} \underline{\mathbf{K}} &= \text{diag}(K^1, \dots, K^m) \in \mathbf{R}^{nm \times nm} & \underline{\mathbf{k}}_x &= (k^1(x_1, x), \dots, k^1(x_n, x), \dots, k^m(x_1, x), \dots) \\ \underline{\mathbf{A}} &= \text{diag}(\underbrace{a_1, \dots, a_1}_{\gamma_1}, \dots, \underbrace{a_m, \dots, a_m}_{\gamma_m}) & \underline{\mathbf{G}} &= \text{diag}(\underbrace{\gamma_1, \dots, \gamma_1}_{\gamma_1}, \dots, \underbrace{\gamma_m, \dots, \gamma_m}_{\gamma_m}) \end{aligned}$$

Summary and Conclusion

- ▶ Presented the co-regularized RKHS with closed form kernel
 - ▶ Easily convert any kernel method into a “multi-view semi-supervised” method
 - ▶ Can apply standard RKHS theory to get improved generalization bounds
- ▶ Presented Manifold Co-Regularization
 - ▶ reinterpretation of manifold regularization
 - ▶ improved empirical results
- ▶ Generalization to
 - ▶ multiple views
 - ▶ arbitrary linear combinations
 - ▶ quadratic coupling term