

Nonparametric Bayesian Kernel Models

Feng Liang

Joint work with Ming Liao, Sayan Mukherjee, Mike West

Institute of Statistics and Decision Sciences
Duke University

NISS Affiliate Annual Meeting, Athens GA, 2006.

Outline

- 1 Introduction
- 2 A Nonparametric Bayesian Approach
- 3 Semi-supervised learning
- 4 Conclusion

Outline

- 1 Introduction
- 2 A Nonparametric Bayesian Approach
- 3 Semi-supervised learning
- 4 Conclusion

- Assumptions: $(\mathbf{x}_i, y_i)_{i=1}^n \sim P(\mathbf{x}, y)$.
- Find a function $f(\mathbf{x}) \rightarrow y$.
- Risk : $R[f] = \mathbb{E}L(Y, f(\mathbf{X}))$

$$f^*(\mathbf{x}) = \operatorname{argmin}_f R[f]$$

e.g., $f^*(\mathbf{x}) = \mathbb{E}(Y \mid \mathbf{X} = \mathbf{x})$ under squared error loss.

- Empirical risk minimization

$$\hat{f}_n(\mathbf{x}) = \operatorname{argmin}_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n L(y_i, f(\mathbf{x}_i))$$

Regularization

- Regularization:

$$\min_{f \in \mathcal{F}} \hat{R}_n[f] + \lambda \Omega[f]$$

- Consider a linear regression model: $f(\mathbf{x}) = \mathbf{x}^t \beta$,

$$\hat{\beta}_{OLS} = \operatorname{argmin}_{\beta \in \mathbb{R}^p} \sum_i (y_i - \mathbf{x}_i^t \beta)^2.$$

- Large p small n : $p \gg n$.
- $\Omega(\beta) = \|\beta\|^2$ (ridge); $|\beta|$ (LASSO).

Reproducing Kernel Hilbert Space (RKHS)

- $\mathcal{H}_K = \{f(\mathbf{x}), \mathbf{x} \in \mathcal{X}\}$.
- $K(\cdot, \cdot)$ is a **semi-positive definite** bivariate symmetric function defined on $\mathcal{X} \times \mathcal{X}$, i.e.

$$\sum_{i=1}^m \sum_{j=1}^m a_i a_j K(\mathbf{x}_i, \mathbf{x}_j) \geq 0, \quad \forall a_i \in \mathbb{R}, \forall \mathbf{x}_i \in \mathcal{X}, \forall m \in \mathbb{N}.$$

- Denote $k_{\mathbf{u}}(\cdot) = K(\mathbf{u}, \cdot)$.

$$\mathcal{H}_K = \overline{\text{span}\{k_{\mathbf{u}}(\cdot), \mathbf{u} \in \mathcal{X}\}}.$$

- **Reproducing property:** $\forall f \in \mathcal{H}_K, \forall \mathbf{x}_0 \in \mathcal{X}$

$$f(\mathbf{x}_0) = \langle f, k_{\mathbf{x}_0} \rangle = \langle f, K(\mathbf{x}_0, \cdot) \rangle.$$

A Simple Examples of RKHS

- Consider all lineal functions in \mathbb{R}^2 passing the origin, i.e.

$$\mathcal{H}_K = \{f_\theta(\mathbf{x}) = \theta^t \mathbf{x} = \theta_1 x_1 + \theta_2 x_2, \theta \in \mathbb{R}^2\}$$

with

$$\langle f_\theta, f_\lambda \rangle = \theta^t \lambda, \quad K(\mathbf{x}, \mathbf{x}') = \mathbf{x}^t \mathbf{x}' = x_1 x_1' + x_2 x_2'.$$

- Reproducing property : $\forall f_\theta$,

$$f_\theta(\mathbf{x}_0) = \theta^t \mathbf{x}_0 = \langle f_\theta, f_{\mathbf{x}_0} \rangle.$$

- Semi-positive definite :

$$\sum_i \sum_j a_i a_j K(\mathbf{x}_i, \mathbf{x}_j) = \left(\sum_i a_i \mathbf{x}_i \right)^t \cdot \left(\sum_j a_j \mathbf{x}_j \right) = \left\| \sum_i a_i \mathbf{x}_i \right\|^2$$

More Examples of RKHS

- Any RKHS is a space of *linear* functions and $K(\cdot, \cdot)$ represents a dot product

More Examples of RKHS

- Any RKHS is a space of *linear* functions and $K(\cdot, \cdot)$ represents a dot product in the **feature space**.

More Examples of RKHS

- Any RKHS is a space of *linear* functions and $K(\cdot, \cdot)$ represents a dot product in the **feature space**.
- Feature mapping $\phi: \mathbf{x} \in \mathbb{R}^p \longrightarrow \phi(\mathbf{x}) \in \mathbb{R}^N$.

$$(x_1, x_2) \in \mathbb{R}^2 \longrightarrow (1, x_1, x_2, x_1^2, x_2^2, x_1x_2) \in \mathbb{R}^6$$

- $\mathcal{H}_K = \{f_\theta(\mathbf{x}) = \theta^t \cdot \phi(\mathbf{x})\}$ and $K(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x})^t \cdot \phi(\mathbf{x}')$.
- Examples
 - **Polynomial** kernels: $K(\mathbf{x}, \mathbf{x}') = (1 + \mathbf{x}^t \mathbf{x}')^d$.
 - **Gaussian radial basis** kernels:

$$K(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{2\sigma^2}\right).$$

Representer Theorem

(Kimeldorf and Wahba, 1971)

$$f^*(\mathbf{x}) = \operatorname{argmin}_{f \in \mathcal{H}_K} \frac{1}{n} \sum_{i=1}^n L(y_i, f(\mathbf{x}_i)) + \lambda \|f\|_{\mathcal{H}_K}^2$$

Representer Theorem

(Kimeldorf and Wahba, 1971)

$$\begin{aligned} f^*(\mathbf{x}) &= \operatorname{argmin}_{f \in \mathcal{H}_K} \frac{1}{n} \sum_{i=1}^n L(y_i, f(\mathbf{x}_i)) + \lambda \|f\|_{\mathcal{H}_K}^2 \\ &= \sum_{i=1}^n \alpha_i k_{\mathbf{x}_i}(\mathbf{x}) = \sum_{i=1}^n \alpha_i K(\mathbf{x}, \mathbf{x}_i) \end{aligned}$$

Representer Theorem

(Kimeldorf and Wahba, 1971)

$$\begin{aligned} f^*(\mathbf{x}) &= \operatorname{argmin}_{f \in \mathcal{H}_K} \frac{1}{n} \sum_{i=1}^n L(y_i, f(\mathbf{x}_i)) + \lambda \|f\|_{\mathcal{H}_K}^2 \\ &= \sum_{i=1}^n \alpha_i k_{\mathbf{x}_i}(\mathbf{x}) = \sum_{i=1}^n \alpha_i K(\mathbf{x}, \mathbf{x}_i) \end{aligned}$$

Proof : Let $\mathcal{H}_1 = \operatorname{span}\{k_{\mathbf{x}_1}, \dots, k_{\mathbf{x}_n}\}$ and $\mathcal{H}_2 = \mathcal{H}_1^\perp$,

$$f = f_1 + f_2, \quad \mathcal{H}_K = \mathcal{H}_1 \oplus \mathcal{H}_2$$

- $\|f\|^2 \geq \|f_1\|^2$
- $f(\mathbf{x}_i) = f_1(\mathbf{x}_i)$ because

$$\langle f, k_{\mathbf{x}_i} \rangle = \langle f_1 + f_2, k_{\mathbf{x}_i} \rangle = \langle f_1, k_{\mathbf{x}_i} \rangle$$

Outline

- 1 Introduction
- 2 A Nonparametric Bayesian Approach**
- 3 Semi-supervised learning
- 4 Conclusion

Primer on Bayesian Analysis

- Data and a parametric family: $\mathbf{y} \sim p(\cdot | \theta)$

Primer on Bayesian Analysis

- Data and a parametric family: $\mathbf{y} \sim p(\cdot | \theta)$
- Prior on θ : $\pi(\theta)$

Primer on Bayesian Analysis

- Data and a parametric family: $\mathbf{y} \sim p(\cdot | \theta)$
- Prior on θ : $\pi(\theta)$
- Joint distribution on (\mathbf{y}, θ) : $p(\mathbf{y} | \theta)\pi(\theta)$

Primer on Bayesian Analysis

- Data and a parametric family: $\mathbf{y} \sim p(\cdot | \theta)$
- Prior on θ : $\pi(\theta)$
- Joint distribution on (\mathbf{y}, θ) : $p(\mathbf{y} | \theta)\pi(\theta)$
- Posterior inference:

$$\pi(\theta | \mathbf{y}) = \frac{p(\mathbf{y} | \theta)\pi(\theta)}{\int p(\mathbf{y} | \theta)\pi(\theta)d\theta} \propto p(\mathbf{y} | \theta)\pi(\theta)$$

Connection to Regularization

- Log posterior = $\log p(\mathbf{y} \mid \theta) + \log \pi(\theta) + \dots$

$$\text{Posterior Mode} = \operatorname{argmin}_{\theta} \sum_i (y_i - f_{\theta}(\mathbf{x}_i))^2 - c \log \pi(\theta)$$

- Regularization

$$\operatorname{argmin}_{f} \sum_{i=1} (y_i - f(\mathbf{x}_i))^2 + \lambda \Omega[f]$$

Connection to Regularization

- Log posterior = $\log p(\mathbf{y} \mid \theta) + \log \pi(\theta) + \dots$

$$\text{Posterior Mode} = \operatorname{argmin}_{\theta} \sum_i (y_i - f_{\theta}(\mathbf{x}_i))^2 - c \log \pi(\theta)$$

- Regularization

$$\operatorname{argmin}_{f} \sum_{i=1} (y_i - f(\mathbf{x}_i))^2 + \lambda \Omega[f]$$

- $\pi(\theta) \propto \exp\{-\lambda \Omega[f]\}$. For example, when $f_{\theta} = \theta^t \mathbf{x}$,
 - Ridge ($\|\theta\|^2$) \iff normal
 - LASSO ($|\theta|$) \iff double exponential

Previous Work

For example, Tipping 2001, Chakraborty et al. 2005, and others.

- Start with the **finite** representation from the representer

Theorem:

$$\sum_{i=1}^n \alpha_i K(\mathbf{x}, \mathbf{x}_i) \quad (*)$$

- Specify priors on the coefficients α_i 's.
- Their models change when sample size changes **without** a coherent argument.
- Can we justify **(*)** using the connection between regularization and posterior mode?

An Orthonormal Representation of \mathcal{H}_K

- For Mercer kernels,

$$K(\mathbf{x}, \mathbf{u}) = \sum_{j=1}^{\infty} \lambda_j \phi_j(\mathbf{x}) \phi_j(\mathbf{u}),$$

where ϕ_j is a sequence of orthonormal functions and $\lambda_1 \geq \lambda_2 \geq \dots \geq 0$.

- (Cucker and Smale, 2001) $\forall f \in \mathcal{H}_K$,

$$f(\mathbf{x}) = \sum_j \alpha_j \phi_j(\mathbf{x}), \quad \sum_j \frac{\alpha_j^2}{\lambda_j} < \infty.$$

That is, \mathcal{H}_K can be parameterized by

$$\mathcal{A} = \{(\alpha_j)_{j=1}^{\infty} : \sum_j \alpha_j^2 / \lambda_j < \infty\}.$$

An Overcomplete Representation of \mathcal{H}_K

- Recall

$$\mathcal{H}_K = \overline{\text{span}\{K(\cdot, \mathbf{u}), \mathbf{u} \in \mathcal{X}\}}$$

- Start with a larger space

$$f(\mathbf{x}) = \int_{\mathcal{X}} K(\mathbf{x}, \mathbf{u}) d\Gamma(\mathbf{u}),$$

where $\Gamma(\mathbf{u})$ is a sign measure on \mathcal{X} .

- In this talk, we focus on the following representation

$$f(\mathbf{x}) = \int_{\mathcal{X}} w(\mathbf{u}) K(\mathbf{x}, \mathbf{u}) dF(\mathbf{u}),$$

where $w(\mathbf{u})$ denotes the coefficient at location \mathbf{u} and F denotes the distribution function of the location \mathbf{u} .

An Overcomplete Representation of \mathcal{H}_K

- Recall

$$\mathcal{H}_K = \overline{\text{span}\{K(\cdot, \mathbf{u}), \mathbf{u} \in \mathcal{X}\}}$$

- Start with a larger space

$$f(\mathbf{x}) = \int_{\mathcal{X}} K(\mathbf{x}, \mathbf{u}) d\Gamma(\mathbf{u}),$$

where $\Gamma(\mathbf{u})$ is a sign measure on \mathcal{X} .

- In this talk, we focus on the following representation

$$f(\mathbf{x}) = \int_{\mathcal{X}} w(\mathbf{u}) K(\mathbf{x}, \mathbf{u}) dF_{\mathbf{X}}(\mathbf{u}),$$

where $w(\mathbf{u})$ denotes the coefficient at location \mathbf{u} and $F_{\mathbf{X}}$ denotes the distribution function of explanatory variable \mathbf{X} .

Dirichlet Process Priors

- **Beta** (α, β) on $x \in [0, 1]$

$$f(x) \propto x^{\alpha-1}(1-x)^{\beta-1}$$

- **Dir** $(\alpha_1, \dots, \alpha_k)$ on (x_1, \dots, x_k) where $x_i \in [0, 1]$ and $\sum_i x_i = 1$

$$f(x) \propto x_1^{\alpha_1-1} \dots x_k^{\alpha_k-1}$$

- **DP** (α_0, F_0) on \mathbf{F} (note \mathbf{F} is a random distribution on \mathcal{X}):
for any measurable partition of \mathcal{X} ,

$$\mathbf{F}(B_1), \mathbf{F}(B_2), \dots, \mathbf{F}(B_k) \sim \text{Dir}(\alpha_1, \dots, \alpha_k),$$

where $\alpha_i = \alpha_0 F_0(B_i)$.

A Bayesian Representer Theorem

Given $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n \sim \mathbf{F}$, with $\text{DP}(\alpha_0, F_0)$ prior on \mathbf{F} , the posterior distribution of \mathbf{F} is $\text{DP}(\alpha_0 + n, \alpha_0 F_0 + \sum \delta_{\mathbf{x}_i})$. So

$$\mathbb{E}[\mathbf{F} \mid \mathbf{x}_1, \dots, \mathbf{x}_n] = \frac{\alpha_0}{\alpha_0 + n} F_0 + \frac{1}{\alpha_0 + n} \sum_{i=1}^n \delta_{\mathbf{x}_i}.$$

Bayesian Representer Theorem

For $f(x) = \int w(\mathbf{u}) K(\mathbf{x}, \mathbf{u}) d\mathbf{F}_X(\mathbf{u})$, under a Dirichlet prior on \mathbf{F}_X ,

$$\mathbb{E}[f(\mathbf{x})] \approx \sum_{i=1}^n w_i K(\mathbf{x}, \mathbf{x}_i).$$

Prior Specification

- Likelihood: $\mathbf{Y} \mid \mathcal{N}(\mathbf{w}_0 + \mathbf{K}\mathbf{w}, \sigma^2 I_n)$ where $\mathbf{K}_{n \times n}$ be the centered kernel matrix at the n data points.
- Non-informative prior on (\mathbf{w}_0, σ^2) ,

$$\pi(\mathbf{w}_0, \sigma^2) \propto 1/\sigma^2$$

- Generalized g-prior on \mathbf{w} (West 2002)

$$\begin{aligned} \mathbf{w} \mid \mathbf{T} &\sim \mathcal{N}(0, U\Delta^{-1}T^{-1}\Delta^{-1}U^t) \\ \tau_1, \dots, \tau_n &\sim \text{Gamma}\left(\frac{s_0}{2}, \frac{s_0}{2}v\right), \quad v \sim \text{Exp}(\alpha_0) \end{aligned}$$

where U and Δ come from $\mathbf{K} = U\Delta U^t$.

- $\mathbf{K}\mathbf{w} = U\beta$, then the prior on \mathbf{w} corresponds to a student t distribution on β .

Model Fitting via MCMC

- Gibbs sampling
 - $w_0 = \dots$
 - Draw $\beta \sim \mathcal{N}(\cdot, \cdot)$
 - Draw $T = \text{diag}(\tau_1, \dots, \tau_n) \sim Ga(\cdot, \cdot)$
 - Draw $v \sim Ga(\cdot, \cdot)$
- Probit model: $(\mathbf{x}, y) \longrightarrow (\mathbf{x}, y, z)$ and

$$P(y = 1) = \Phi(z).$$

Outline

- 1 Introduction
- 2 A Nonparametric Bayesian Approach
- 3 Semi-supervised learning**
- 4 Conclusion

Semi-supervised Learning

- Supervised learning (**labelled data**):

$$(\mathbf{x}_i, y_i) \sim P(\mathbf{x}, y)$$

- Unsupervised learning (**unlabelled data**):

$$(\mathbf{x}_i) \sim P_{\mathbf{X}}(\mathbf{x})$$

- Semi-supervised learning:

$$(\mathbf{x}_i, y_i)_{i=1}^n \sim P, \quad (\mathbf{x}_i)_{i=n+1}^{n+m} \sim P_{\mathbf{X}}$$

- How's it different from missing data?

The Role of Unlabelled Data

- Data $D = (\mathbf{X}, Y, \mathbf{X}^U)$ where

$$\begin{aligned}(\mathbf{X}, Y) &\sim P(\mathbf{x}, y) = p(\mathbf{x} \mid \phi) p(y \mid \mathbf{x}, \theta) \\ \mathbf{X}^U &\sim P_{\mathbf{X}}(\mathbf{x}) = p(\mathbf{x} \mid \phi)\end{aligned}$$

- Prediction of \mathbf{y}^* at a new location \mathbf{x}^* ,

$$\begin{aligned}\mathbb{E}[Y^* \mid \mathbf{x}^*, D] &= \int y^* p(y^* \mid \mathbf{x}^*, D) dy^* \\ &= \int y^* p(y^* \mid \mathbf{x}^*, \theta) \pi(\theta \mid D) d\theta dy^*.\end{aligned}$$

- The key to understand the role of \mathbf{X}^m is

$$\begin{aligned}\pi(\theta \mid D) &= \int \pi(\theta, \phi \mid D) d\phi \\ &\propto \int p(\mathbf{X}, \mathbf{X}^U \mid \phi) p(Y \mid \mathbf{X}, \theta) \pi(\theta, \phi) d\phi\end{aligned}$$

An Intimate Relationship

- $P(\mathbf{x}, y) = p(\mathbf{x} \mid \phi)p(y \mid \mathbf{x}, \theta)$
- Recall that $y \mid \mathbf{x} \sim \mathcal{N}(f(\mathbf{x}), \sigma^2)$ where

$$f(\mathbf{x}) = \int w(\mathbf{u})K(\mathbf{x}, \mathbf{u})dF_{\mathbf{X}}(\mathbf{u}).$$

So in our model, $\theta = (\phi, \dots)$. Therefore unlabelled data will be relevant and should be incorporated into prediction.

- By our Bayesian representer Theorem, given $(\mathbf{x}_i, y_i)_{i=1}^n$ and $(\mathbf{x}_j)_{j=n+1}^{n+m}$,

$$f(\mathbf{x}) \approx \sum_{i=1}^n w_i K(\mathbf{x}, \mathbf{x}_i) + \sum_{j=n+1}^{n+m} w_j K(\mathbf{x}, \mathbf{x}_j).$$

Connection to Regularization

- Transductive SVM (TSVM) (Joachims, 1999)

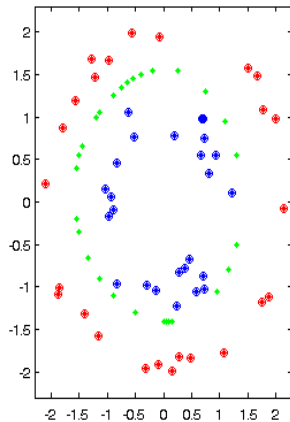
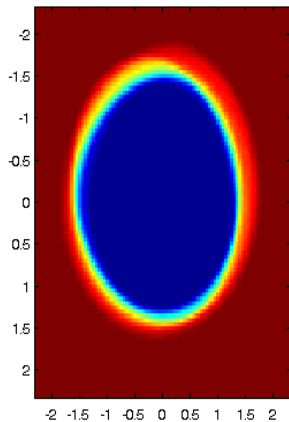
$$\operatorname{argmin} C \sum_{i=1}^n L(y_i, f(\mathbf{x}_i)) + C^* \sum_{i=n+1}^{n+m} L(y_i, f(\mathbf{x}_i)) + \|f\|_{\mathcal{H}_K}^2$$

- Manifold regularization (Belkin et al, 2005)

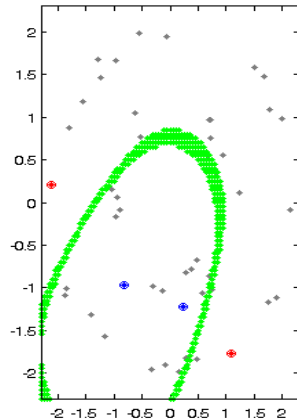
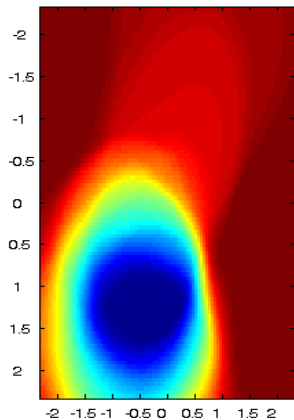
$$\operatorname{argmin} \sum_{i=1}^n L(y_i, f(\mathbf{x}_i)) + \lambda_1 \|f\|_{\mathcal{H}_K}^2 + \lambda_2 \|f\|_I^2,$$

where $\|f\|_I^2$ measures the intrinsic structure of $F_{\mathbf{X}}$ and is approximated on all the data (include the unlabelled ones) using graph Laplacian.

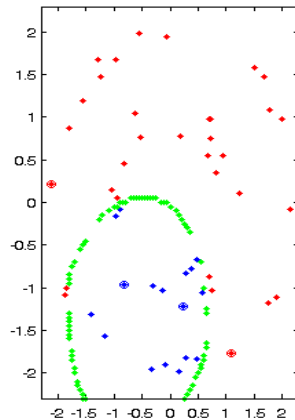
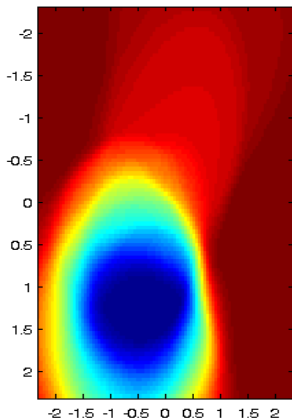
A Toydata



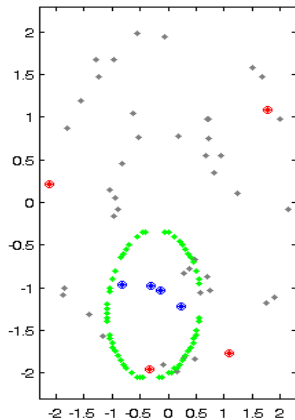
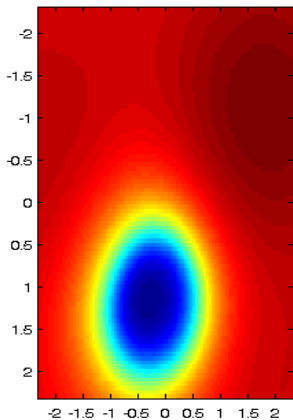
A Toydata



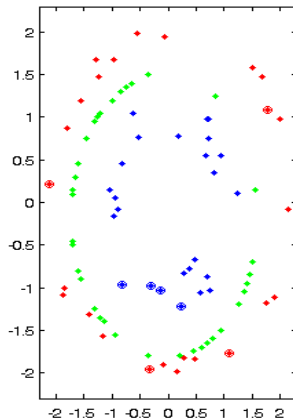
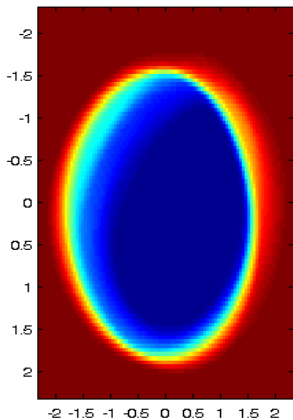
A Toydata



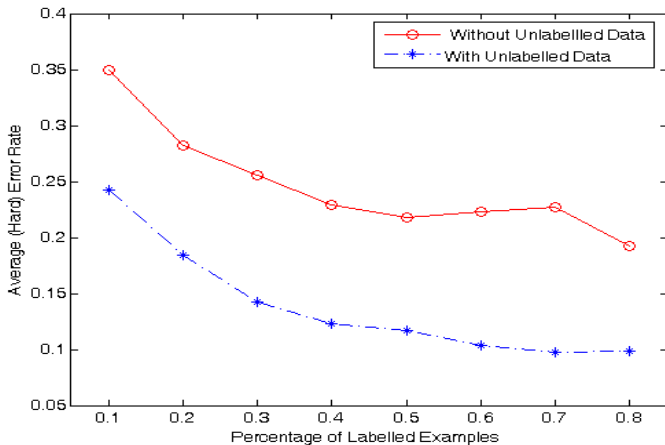
A Toydata



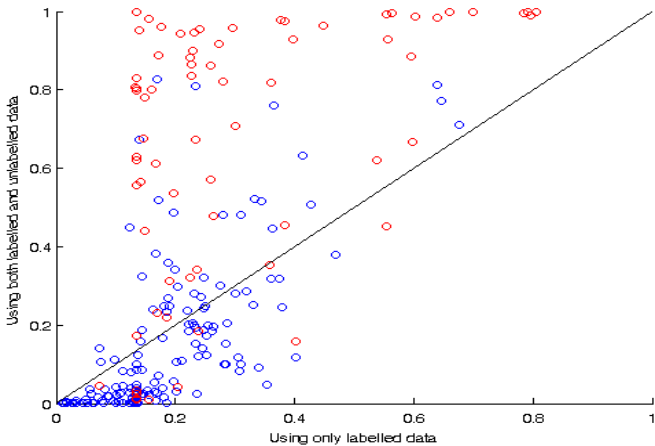
A Toydata



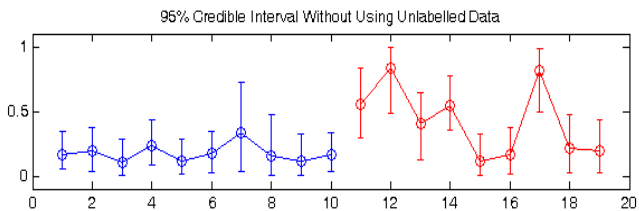
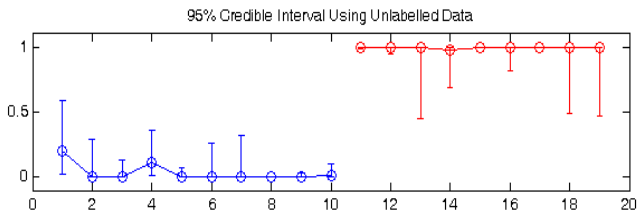
Cancer Data



Cancer Data



Cancer Data



Outline

- 1 Introduction
- 2 A Nonparametric Bayesian Approach
- 3 Semi-supervised learning
- 4 Conclusion**

Conclusion

- A coherent Bayesian approach on RKHS:
 - characterize RKHS using an overcomplete representation;
 - specify priors on the whole RKHS;
 - incorporate the relevant information from the unlabelled data.

Conclusion

- A coherent Bayesian approach on RKHS:
 - characterize RKHS using an overcomplete representation;
 - specify priors on the whole RKHS;
 - incorporate the relevant information from the unlabelled data.
- Future work:
 - Other choice of priors and sensitivity study
 - Feature selection