

# Active Kernel Learning

---

**Rong Jin**

Department of Computer Science and Engineering  
Michigan State University, USA

Joint work with Steven H. C. Hoi from Nanyang Technological University

# Kernel Learning

---

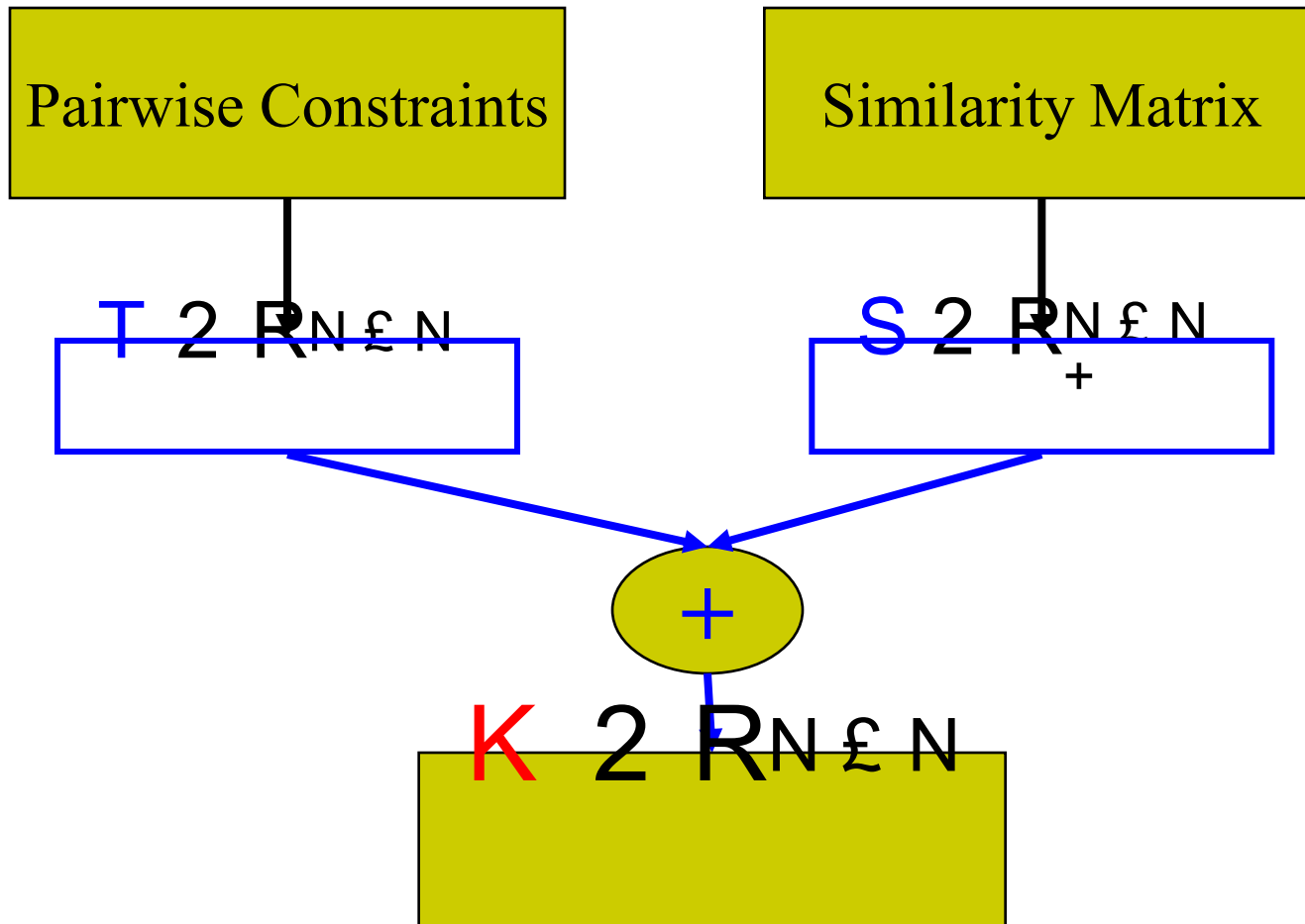
- Kernel is the key to many learning schemes
- Kernel learning
  - Learning parametric kernel functions  
e.g., RBF kernels and diffusion kernels
  - Non-parametric kernel learning
    - No parametric assumption about kernel function
    - Learning kernel matrices

# Kernel Learning Methods

---

- Kernel alignment (Cristianini et al., 2001)
  - Maximize the alignment with the assigned class labels
- Learning kernel with SDP ( Lanckriet et al., 2004)
  - Maximize the classification margin
- Nonparametric graph kernel (Zhu et al., 2005)
  - Kernel alignment + order constraint
- Learning low rank kernel matrices (Kulis et al., 2006)
  - Initial kernel + consistent with the assigned class labels
- Non-parametric kernel learning (Hoi et. al., 2007)
  - Initial similarity + pairwise constraints

# Non-parametric Kernel Learning



# Active Kernel Learning

---

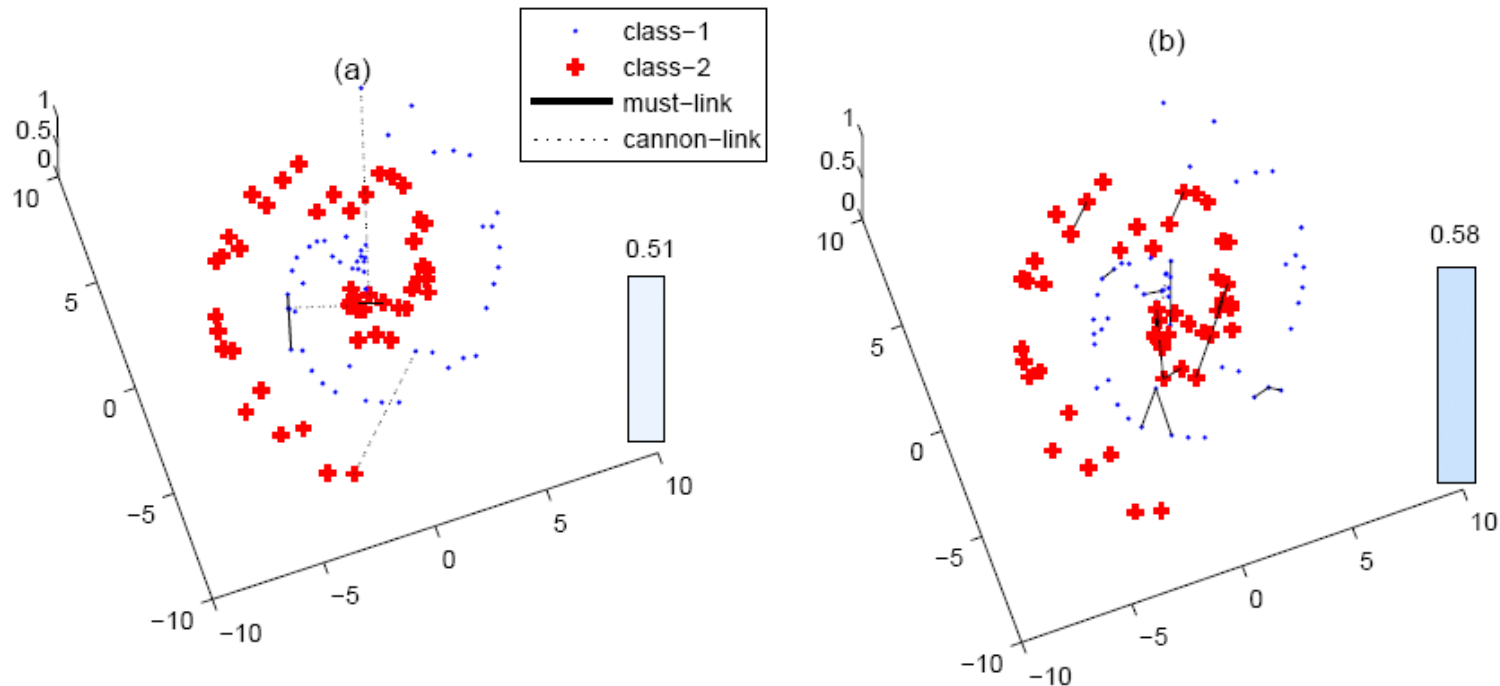
- Existing kernel learning schemes are passive
  - Assume labeling information is given beforehand
  - Less efficient, may require more labeling efforts
- Active kernel learning
  - Select the most informative example pairs from a pool of unlabeled examples
  - Extension of non-parametric kernel learning (Hoi et. al., 2007)

# Active Kernel Learning

---

- Key challenge
  - How to measure the informativeness for an example pair ?
- Simple strategy: uncertainty principle
  - Select an example pair with the least confidence
    - Example pairs with least  $\sum_{i,j} K_{i,j}$
  - Main drawback
    - Lacks of principles
    - Could result in the preference of must-link constraints

# Active Kernel Learning



Uncertainty principle leads to the preference of must-link constraints



# Contribution

---

- Propose a min-max framework for active kernel learning
- Present a convex relaxation for the proposed framework of active kernel learning



# Non-parametric Kernel Learning: Review

$$U = (x_1; x_2; \dots; x_N); \quad x_i \in \mathbb{R}^d$$

□ Data

$$T \in \mathbb{R}^{N \times N}$$

□ **Pairwise constraints:**

$$T_{i;j} = \begin{cases} < 1 & x_i \text{ and } x_j \text{ in the same class} \\ i > 1 & x_i \text{ and } x_j \text{ in different classes} \\ 0 & \text{otherwise} \end{cases} \quad \begin{array}{l} \text{must-link} \\ \text{cannot-link} \end{array}$$

$$S \in \mathbb{R}_+^{N \times N}$$

□ **Similarity measurements:**

■  $S_{i;j}$ : similarity between  $x_i$  and  $x_j$

# Non-parametric Kernel Learning: Review

$$U = (x_1; x_2; \dots; x_N); \quad x_i \in \mathbb{R}^d$$

□ Data

□ **Pairwise constraints:**

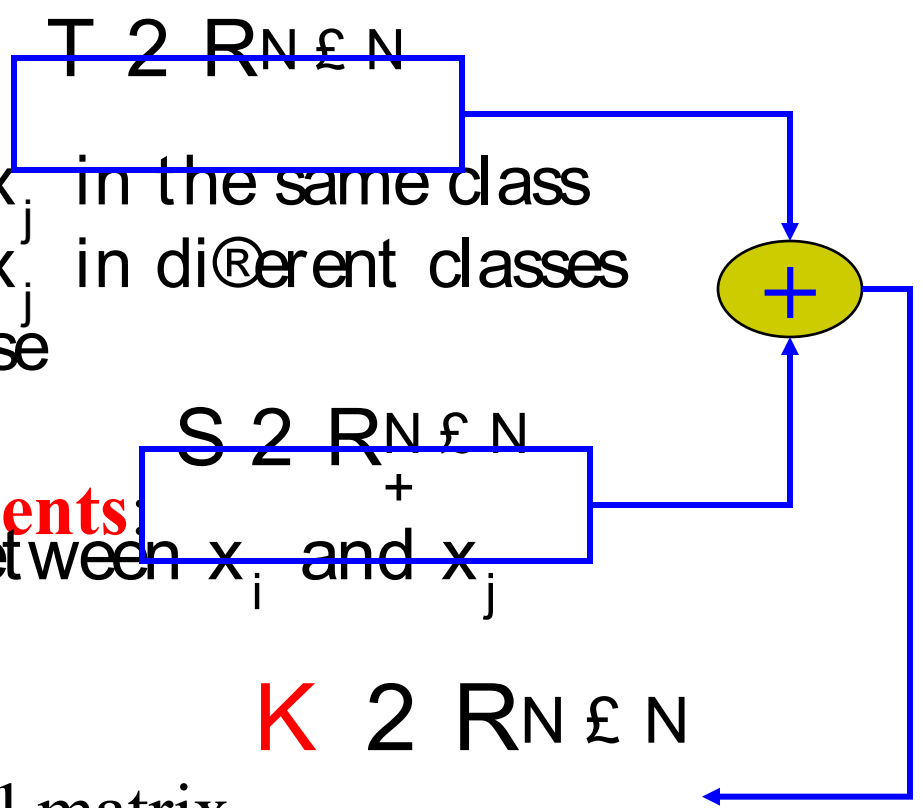
$$T_{i;j} = \begin{cases} < 1 & x_i \text{ and } x_j \text{ in the same class} \\ i & x_i \text{ and } x_j \text{ in different classes} \\ 0 & \text{otherwise} \end{cases}$$

□ **Similarity measurements:**

■  $S_{i;j}$ : similarity between  $x_i$  and  $x_j$

$$K \in \mathbb{R}^{N \times N}$$

Learning a kernel matrix



# Non-parametric Kernel Learning

$$\begin{aligned}
 \arg \min_{Z; \gamma} & \quad \text{tr}(LZ) + \frac{c}{2} \sum_{(i;j) \in T} \|Z_{i;j} - \gamma_{i;j}\|_2^2 \\
 \text{s. t.} & \quad Z \succeq 0
 \end{aligned}$$

Consistency with similarity matrix  $S$

Consistency with labeled example pairs

- $Z$ : target kernel matrix
- $L$ : graph Laplacian for similarity matrix  $S$

# Min-Max Framework

- Informativeness of an unlabeled pair  $(x_k; x_l)$

$$! ((k;l); y) = \min_{Z; \alpha} \text{tr}(LZ) + \frac{c}{2} \sum_{i;j} \alpha_{i;j}^2 + \frac{c}{2} \alpha_{k;l}^2$$

$$\text{s. t. } \sum_{i;j} \alpha_{i;j} Z_{i;j} = \mathbf{1}; \sum_{i;j} \alpha_{i;j}^2 \leq 8; Z \succeq 0$$

$$yZ_{k;l} = \mathbf{1}; \alpha_{k;l} = 1$$

Measures how example pair  $(k,l)$  with label  $y$  will affect the overall objective function.

# Min-Max Framework

- Informativeness of an unlabeled pair  $(x_k; x_l)$

$$\begin{aligned}
 \min_{Z; \gamma} \quad & \text{tr}(LZ) + \frac{c}{2} \sum_{(i;j) \in T} \|z_{i;j}\|^2 + \frac{c}{2} \sum_{(k;l) \in T} \|z_{k;l}\|^2 \\
 \text{s. t.} \quad & \sum_{(i;j) \in T} z_{i;j} = 1; \quad \sum_{(i;j) \in T} \|z_{i;j}\|^2 = 1 \\
 & \gamma z_{k;l} = 1; \quad \|z_{k;l}\|^2 = 0
 \end{aligned}$$

$(k,l)$  is highly consistent with the current kernel  $Z$

↓

$\gamma((k;l); y)$  will be small with appropriate choice of  $y$ , and will be large with incorrect  $y$

# Min-Max Framework

- Informativeness of an unlabeled pair  $(x_k; x_l)$

$$!((k;l); y) = \min_{Z} \text{tr}(LZ) + \frac{c}{2} \sum_{(i;j) \in T} \|z_i - z_j\|^2 + \frac{c}{2} \|z_k - z_l\|^2$$

s. t.  $\sum_{(i;j) \in T} \|z_i - z_j\|^2 \leq 1$

Measure the uninformative-ness of an example pair

example pairs highly consistent with  $Z \rightarrow$  large  $!((k;l); y)$

$$!((k;l)) = \max_{y \in \{-1, +1\}} !((k;l); y)$$

$$y \in \{-1, +1\}$$

# Min-Max Framework

$$\begin{aligned}
 (k; l)^* &= \arg \min_{(k; l) \in T} \cdot (k; l) \\
 &= \arg \min_{(k; l) \in T} \max_{t \in T} \{ (k; l); t \} \quad (P)
 \end{aligned}$$

- A challenging optimization problem
  - is not a closed-form function
  - Min-max optimization

# Min-Max Framework: Simplification

**Theorem 1:**  $(P)$  is equivalent to

$$\min_{Z; (k;l) \in T} \text{tr}(LZ) + \frac{1}{2} \sum_{(i;j) \in T} \|z_{i;j}\|_2^2 + \frac{c}{2} \sum_{(k;l)} \|z_{k;l}\|_2^2$$

s. t.

$$\begin{aligned} & T_{i;j} z_{i;j} \leq 1 \quad \forall (i;j) \in T \\ & \|z_{k;l}\|_2 \leq 1 + \sum_j z_{k;l} \quad \forall (k;l) \in T \\ & Z \succeq 0 \end{aligned} \quad (P1)$$

**Corollary 2:** When  $Z$  is fixed,  $(P1)$  is equivalent to

$$(k;l)^\alpha = \arg \min_{(k;l) \in T} \sum_j z_{k;l}$$



# Min-Max Framework: Simplification

**Theorem 1:**  $(P)$  is equivalent to

$$\min_{Z; (k;l) \in T} \text{tr}(LZ) + \frac{1}{2} \sum_{(i;j) \in T} \|z_{i;j}\|_2^2 + \frac{c}{2} \sum_{(k;l)} \|z_{k;l}\|_2^2$$

s. t.

$$\begin{aligned} T_{i;j} z_{i;j} &\leq 1 \quad \forall (i;j) \in T \\ \|z_{k;l}\|_2 &\leq 1 + \rho_{k;l} \quad \forall (k;l) \in T \\ Z &\succeq 0 \end{aligned} \quad (P1)$$

1. Introduce  $\rho_{k;l}$  to weight each example pair  $(k,l)$
2. Aggregate all the optimization problems together

# Min-Max Framework: Simplification

**Theorem 3:** (P1) is equivalent to

$$\begin{aligned}
 \min_{Z \succeq 0; p} \quad & \text{tr}(LZ) + \frac{c}{2} \sum_{(i;j) \in T} \|Z_{i;j}\|_2^2 + \frac{c}{2} \sum_{(k;l) \in T} p_{k;l}^2 \\
 \text{s. t.} \quad & T_{i;j} Z_{i;j} \succeq 1_i \otimes I_{i;j}; \quad \forall (i;j) \in T \\
 & \otimes_{k;l} I_{k;l} \succeq Z_{k;l} \succeq 1_i \otimes I_{k;l}; \quad \forall (k;l) \in T \\
 & p_{k;l} \succeq 1; \quad p_{k;l} \succeq 0; \quad \forall (k;l) \in T
 \end{aligned} \tag{P2}$$

**But, this is non-convex optimization !**

# Min-Max Framework: Simplification

$$\min_{\substack{X \\ (k;l) \in T}} \sum_{(k;l) \in T} p_{k;l} \frac{1}{p_{k;l}}$$

$$h_{k;l} = p_{k;l}^{-1}$$

$$\begin{aligned} \min_{Z \succeq 0; p;} \quad & \text{tr}(LZ) + \sum_{(i;j) \in T} \frac{c}{2} \|z_{i;j}\|^2 + \sum_{(k;l) \in T} \frac{c}{2} p_{k;l}^2 \\ \text{s. t.} \quad & T_{i;j} Z_{i;j} \succeq 1 I_{i;j}; \quad \forall (i;j) \in T \\ & \|z_{k;l}\| \leq 1; \quad \forall (k;l) \in T \\ & p_{k;l} \leq 1; \quad p_{k;l} \geq 0; \quad \forall (k;l) \in T \end{aligned}$$

# Min-Max Framework: Simplification

$$\begin{aligned}
 \min_{Z \succeq 0; h; \gamma} \quad & \text{tr}(LZ) + \frac{c}{2} \sum_{(i;j) \in \mathcal{T}} \|Z_{i;j}\|_2^2 + \frac{c}{2} \sum_{(k;l) \in \mathcal{T}} \|h_{k;l}\|_2^2 \\
 \text{s. t.} \quad & T_{i;j} Z_{i;j} \succeq 1_i \gamma_{i;j}; \gamma_{i;j} \succeq 0; \gamma_{i;j} \succeq \delta_{i;j} \mathbf{I} \\
 & \gamma_{k;l} \succeq 1_i \gamma_{k;l}; \gamma_{k;l} \succeq \delta_{k;l} \mathbf{I} \\
 & h_{k;l} \succeq m^2; h_{k;l} \succeq 1; \gamma_{k;l} \succeq \delta_{k;l} \mathbf{I}
 \end{aligned}$$

$$h_{k;l} = \rho_{k;l}^{-1}$$

## Theorem 4:

- $(P3)$  is convex (semi-definite programming)
- The optimal value of  $(P2)$  is upper bounded by that of  $(P3)$

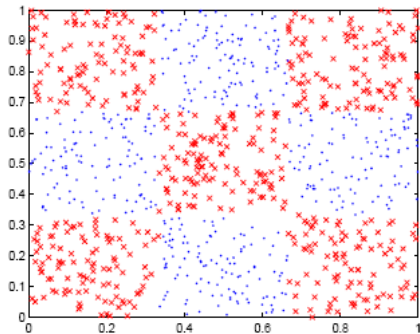
# Min-Max Framework: Simplification

$$\begin{aligned}
 \min_{Z \succeq 0; h} \quad & \text{tr}(LZ) + \frac{c}{2} \sum_{(i;j) \in T} \|z_{i;j}\|^2 + \frac{c}{2} \sum_{(k;l) \in T} \|h_{k;l}\|^2 \\
 \text{s. t.} \quad & T_{i;j} Z_{i;j} \succeq 1_i - \|z_{i;j}\|^2; \quad \|z_{i;j}\|^2 \leq 0; \quad \forall (i;j) \in T \\
 & \|h_{k;l}\|^2 \leq 1; \quad Z_{k;l} \succeq 1_i - \|h_{k;l}\|^2; \quad \forall (k;l) \in T \\
 & h_{k;l} \cdot m^2; \quad h_{k;l} \succeq 1; \quad \forall (k;l) \in T
 \end{aligned} \tag{P3}$$

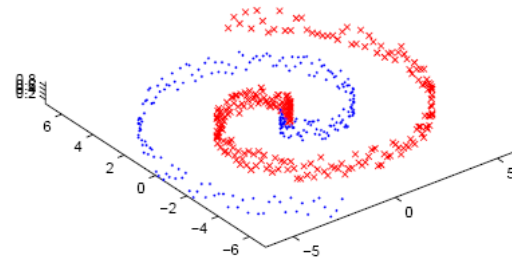
(P3) can be optimized more efficiently using the dual formulation (details in paper)

# Experiments: Datasets

| Dataset       | #Classes | #Instances | #Features |
|---------------|----------|------------|-----------|
| Chessboard    | 2        | 100        | 2         |
| Double-Spiral | 2        | 100        | 3         |
| Glass         | 6        | 214        | 9         |
| Heart         | 2        | 270        | 13        |
| Iris          | 3        | 150        | 4         |
| Protein       | 6        | 116        | 20        |
| Sonar         | 2        | 208        | 60        |
| Soybean       | 4        | 47         | 35        |
| Wine          | 3        | 178        | 12        |



Chessboard



Double Spiral

# Experiments: Setup

---

## □ Baselines

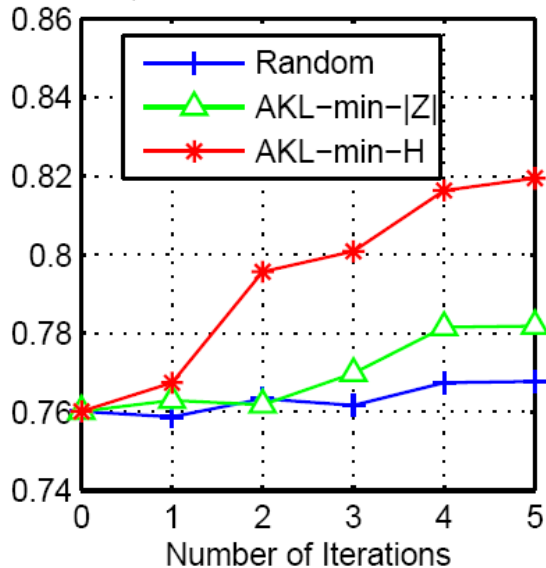
- **Random**: randomly samples example pairs
- **AKL-min-|Z|**: chooses pair examples with the least  $|Z_{k,l}|$
- **AKL-min-H**: the proposed AKL algorithm. selects example pairs with least  $h_{k,l}$  that corresponds to maximal  $p_{k,l}$ .

## □ Evaluation

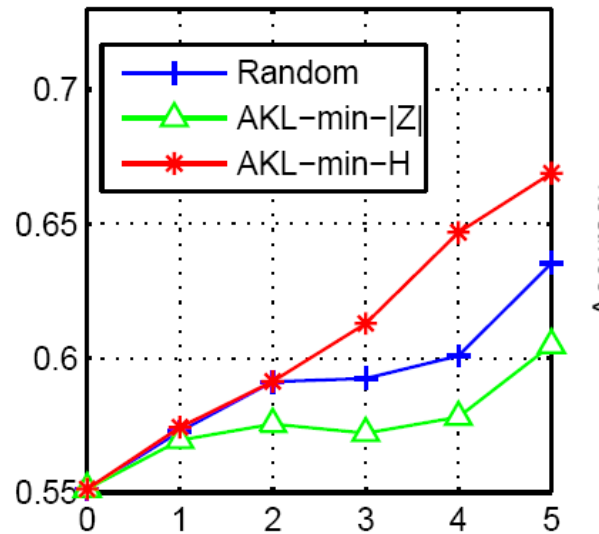
- Evaluate kernel matrices by data clustering using kernel K-means
- Clustering accuracy = 
$$\frac{\sum_{i>j} \mathbb{1}\{c_i = c_j\}}{0.5n(n-1)}$$

# Experimental Results

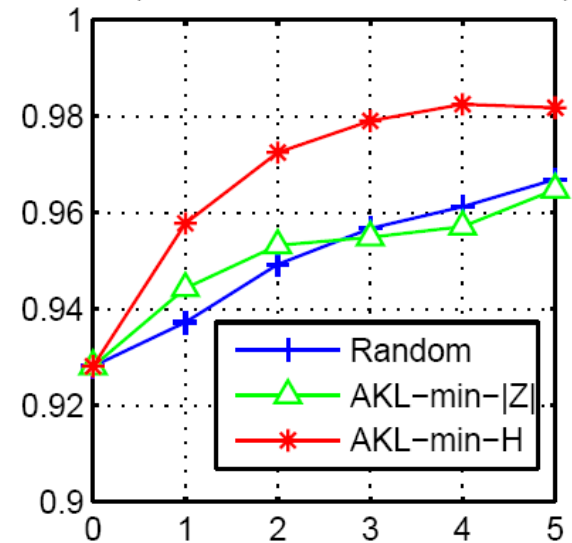
Wine (N=178, C=3, D=12, Nc=100)



Sonar (N=208, C=2, D=60, Nc=100)



Iris (N=150, C=3, D=4, Nc=100)



1. Nc: number of initially given pairwise constraints
2. 20 additional pairwise constraints are provided each iteration

$$\text{AKL-min-H} > \text{Random} \approx \text{AKL-min-|Z|}$$





# Conclusion

---

- Propose a min-max framework for active learning, and apply it to active kernel learning
- Present a convex relaxation for the proposed min-max framework for active kernel learning
- Show encouraging experimental results using data clustering