



# Outline

1. Introduction
2. Multiple Kernel Learning
3. Localized Multiple Kernel Learning
4. Local Projection Kernels
5. Experiments
6. Conclusions





# Basics

- Observed data with outputs  $\{(\mathbf{x}_i, y_i)\}_{i=1}^N$
- Estimate output for an unseen test instance  $\mathbf{x}$

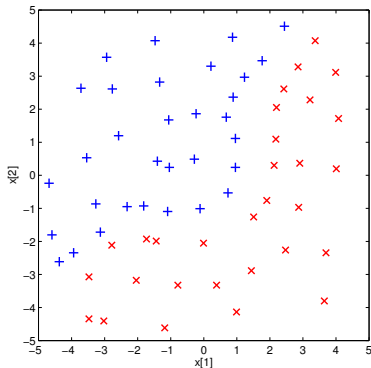
# Basics

- Observed data with outputs  $\{(\mathbf{x}_i, y_i)\}_{i=1}^N$
- Estimate output for an unseen test instance  $\mathbf{x}$
- Sample problems

Problem	Input	Features	Output
Digit recognition		Pixel values	$\{0, 1, \dots, 9\}$
Disease diagnosis	ATCGGT . . . TTA	Nucleotide bases	$\{\text{Healthy, Not healthy}\}$
Face recognition		Pixel values	$\{\text{Person1, Person2, } \dots \}$

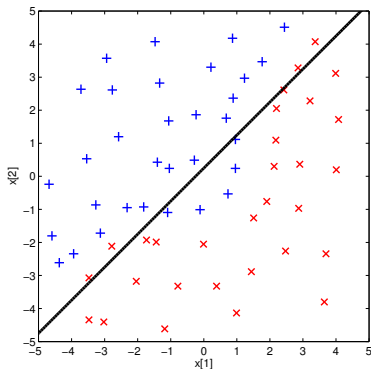
# Support Vector Machines

- $\{(\mathbf{x}_i, y_i)\}_{i=1}^N$  where  $\mathbf{x}_i \in \mathbb{R}^D$  and  $y_i \in \{-1, +1\}$



# Support Vector Machines

- $\{(\mathbf{x}_i, y_i)\}_{i=1}^N$  where  $\mathbf{x}_i \in \mathbb{R}^D$  and  $y_i \in \{-1, +1\}$
- $f(\mathbf{x}) = \langle \mathbf{w}, \mathbf{x} \rangle + b$



# Support Vector Machines

- $\{(\mathbf{x}_i, y_i)\}_{i=1}^N$  where  $\mathbf{x}_i \in \mathbb{R}^D$  and  $y_i \in \{-1, +1\}$
- $f(\mathbf{x}) = \langle \mathbf{w}, \mathbf{x} \rangle + b$
- Primal optimization problem (Vapnik, 1998)

$$\begin{aligned} \min. \quad & \underbrace{\frac{1}{2} \|\mathbf{w}\|_2^2}_{\text{regularization}} + C \underbrace{\sum_{i=1}^N \xi_i}_{\text{soft error}} \\ \text{w.r.t.} \quad & \mathbf{w} \in \mathbb{R}^D, \boldsymbol{\xi} \in \mathbb{R}_+^N, b \in \mathbb{R} \\ \text{s.t.} \quad & y_i(\langle \mathbf{w}, \mathbf{x}_i \rangle + b) \geq 1 - \xi_i \quad \forall i \end{aligned}$$

# Support Vector Machines

## ■ Dual optimization problem

$$\max. \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle$$

$$\text{w.r.t. } \boldsymbol{\alpha} \in \mathbb{R}_+^N$$

$$\text{s.t. } \sum_{i=1}^N \alpha_i y_i = 0$$

$$C \geq \alpha_i \geq 0 \quad \forall i$$

# Support Vector Machines

- Dual optimization problem

$$\max. \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle$$

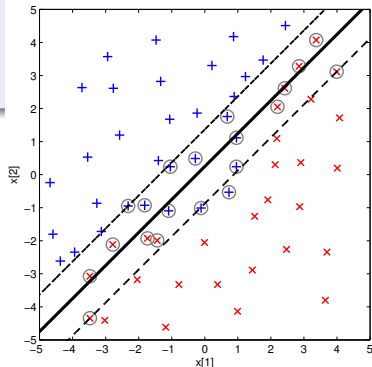
$$\text{w.r.t. } \boldsymbol{\alpha} \in \mathbb{R}_+^N$$

$$\text{s.t. } \sum_{i=1}^N \alpha_i y_i = 0$$

$$C \geq \alpha_i \geq 0 \quad \forall i$$

- $\mathbf{w} = \sum_{i=1}^N \alpha_i y_i \mathbf{x}_i$

- $f(\mathbf{x}) = \sum_{i=1}^N \alpha_i y_i \langle \mathbf{x}_i, \mathbf{x} \rangle + b$



# Kernel Functions

- Define a mapping function  $\Phi: \mathbb{R}^D \rightarrow \mathbb{R}^S$
- Replace  $\langle \mathbf{x}_i, \mathbf{x}_j \rangle$  with  $\underbrace{\langle \Phi(\mathbf{x}_i), \Phi(\mathbf{x}_j) \rangle}_{k(\mathbf{x}_i, \mathbf{x}_j)}$
- Choosing kernel function  $k(\cdot, \cdot)$  and regularization parameter  $C$

# Kernel Functions

■ Define a mapping function  $\Phi: \mathbb{R}^D \rightarrow \mathbb{R}^S$

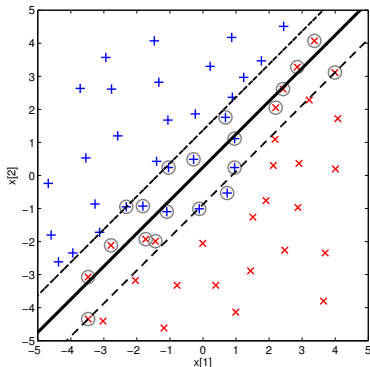
■ Replace  $\langle \mathbf{x}_i, \mathbf{x}_j \rangle$  with  $\underbrace{\langle \Phi(\mathbf{x}_i), \Phi(\mathbf{x}_j) \rangle}_{k(\mathbf{x}_i, \mathbf{x}_j)}$

■ Choosing kernel function  $k(\cdot, \cdot)$  and regularization parameter  $C$

■ Common kernel functions

▶ Linear kernel

$$k_L(\mathbf{x}_i, \mathbf{x}_j) = \langle \mathbf{x}_i, \mathbf{x}_j \rangle$$



# Kernel Functions

■ Define a mapping function  $\Phi: \mathbb{R}^D \rightarrow \mathbb{R}^S$

■ Replace  $\langle \mathbf{x}_i, \mathbf{x}_j \rangle$  with  $\underbrace{\langle \Phi(\mathbf{x}_i), \Phi(\mathbf{x}_j) \rangle}_{k(\mathbf{x}_i, \mathbf{x}_j)}$

■ Choosing kernel function  $k(\cdot, \cdot)$  and regularization parameter  $C$

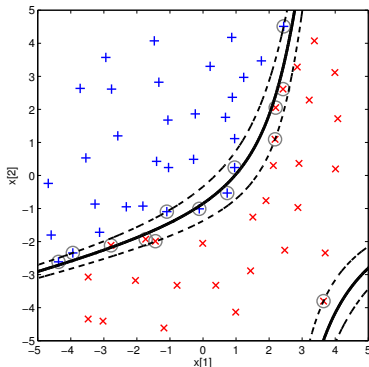
■ Common kernel functions

▶ Linear kernel

$$k_L(\mathbf{x}_i, \mathbf{x}_j) = \langle \mathbf{x}_i, \mathbf{x}_j \rangle$$

▶ Polynomial kernel

$$k_P(\mathbf{x}_i, \mathbf{x}_j) = (\langle \mathbf{x}_i, \mathbf{x}_j \rangle + 1)^q$$



# Kernel Functions

■ Define a mapping function  $\Phi: \mathbb{R}^D \rightarrow \mathbb{R}^S$

■ Replace  $\langle \mathbf{x}_i, \mathbf{x}_j \rangle$  with  $\underbrace{\langle \Phi(\mathbf{x}_i), \Phi(\mathbf{x}_j) \rangle}_{k(\mathbf{x}_i, \mathbf{x}_j)}$

■ Choosing kernel function  $k(\cdot, \cdot)$  and regularization parameter  $C$

■ Common kernel functions

▶ Linear kernel

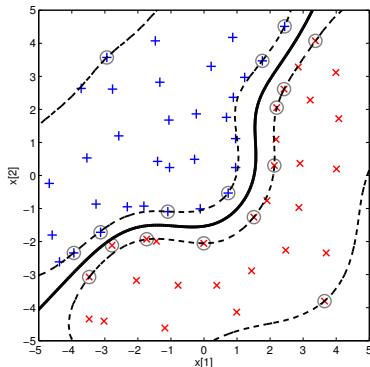
$$k_L(\mathbf{x}_i, \mathbf{x}_j) = \langle \mathbf{x}_i, \mathbf{x}_j \rangle$$

▶ Polynomial kernel

$$k_P(\mathbf{x}_i, \mathbf{x}_j) = (\langle \mathbf{x}_i, \mathbf{x}_j \rangle + 1)^q$$

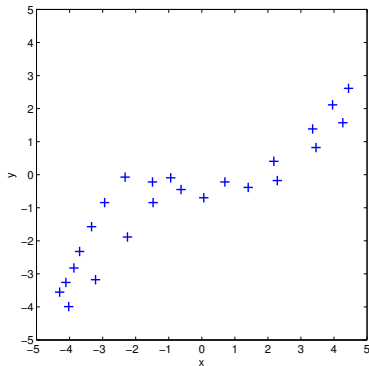
▶ Gaussian kernel

$$k_G(\mathbf{x}_i, \mathbf{x}_j) = \exp(-\|\mathbf{x}_i - \mathbf{x}_j\|_2^2/s^2)$$



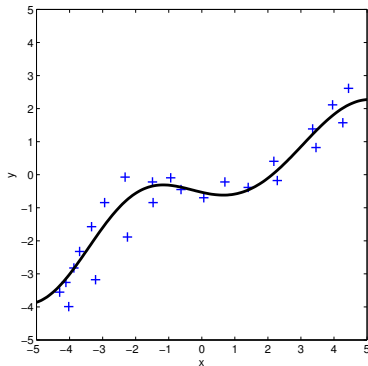
# Regression Support Vector Machines

- $\{(\mathbf{x}_i, y_i)\}_{i=1}^N$  where  $\mathbf{x}_i \in \mathbb{R}^D$  and  $y_i \in \mathbb{R}$



# Regression Support Vector Machines

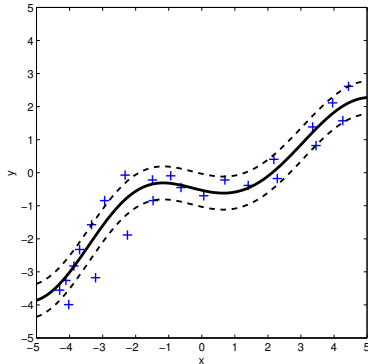
- $\{(\mathbf{x}_i, y_i)\}_{i=1}^N$  where  $\mathbf{x}_i \in \mathbb{R}^D$  and  $y_i \in \mathbb{R}$
- $f(\mathbf{x}) = \langle \mathbf{w}, \Phi(\mathbf{x}) \rangle + b$  where  $\Phi: \mathbb{R}^D \rightarrow \mathbb{R}^S$



# Regression Support Vector Machines

- $\{(\mathbf{x}_i, y_i)\}_{i=1}^N$  where  $\mathbf{x}_i \in \mathbb{R}^D$  and  $y_i \in \mathbb{R}$
- $f(\mathbf{x}) = \langle \mathbf{w}, \Phi(\mathbf{x}) \rangle + b$  where  $\Phi: \mathbb{R}^D \rightarrow \mathbb{R}^S$
- $\epsilon$ -insensitive error function (Vapnik, 1998)

$$e(y, f(\mathbf{x})) = \begin{cases} 0 & \text{if } |y - f(\mathbf{x})| \leq \epsilon \\ |y - f(\mathbf{x})| & \text{otherwise} \end{cases}$$



# Regression Support Vector Machines

- $\{(\mathbf{x}_i, y_i)\}_{i=1}^N$  where  $\mathbf{x}_i \in \mathbb{R}^D$  and  $y_i \in \mathbb{R}$
- $f(\mathbf{x}) = \langle \mathbf{w}, \Phi(\mathbf{x}) \rangle + b$  where  $\Phi: \mathbb{R}^D \rightarrow \mathbb{R}^S$
- $\epsilon$ -insensitive error function (Vapnik, 1998)

$$e(y, f(\mathbf{x})) = \begin{cases} 0 & \text{if } |y - f(\mathbf{x})| \leq \epsilon \\ |y - f(\mathbf{x})| & \text{otherwise} \end{cases}$$

- Primal optimization problem

$$\min. \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_{i=1}^N (\xi_i^+ + \xi_i^-)$$

$$\text{w.r.t. } \mathbf{w} \in \mathbb{R}^S, \boldsymbol{\xi}^+ \in \mathbb{R}_+^N, \boldsymbol{\xi}^- \in \mathbb{R}_+^N, b \in \mathbb{R}$$

$$\text{s.t. } \epsilon + \xi_i^+ \geq y_i - \langle \mathbf{w}, \Phi(\mathbf{x}_i) \rangle - b \quad \forall i$$

$$\epsilon + \xi_i^- \geq \langle \mathbf{w}, \Phi(\mathbf{x}_i) \rangle + b - y_i \quad \forall i$$

# Regression Support Vector Machines

## ■ Dual optimization problem

$$\begin{aligned} \max. \quad & \sum_{i=1}^N y_i(\alpha_i^+ - \alpha_i^-) - \epsilon \sum_{i=1}^N (\alpha_i^+ + \alpha_i^-) \\ & - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N (\alpha_i^+ - \alpha_i^-)(\alpha_j^+ - \alpha_j^-) k(\mathbf{x}_i, \mathbf{x}_j) \end{aligned}$$

$$\text{w.r.t. } \boldsymbol{\alpha}^+ \in \mathbb{R}_+^N, \boldsymbol{\alpha}^- \in \mathbb{R}_+^N$$

$$\text{s.t. } \sum_{i=1}^N (\alpha_i^+ - \alpha_i^-) = 0$$

$$C \geq \alpha_i^+, \alpha_i^- \geq 0 \quad \forall i$$

# Regression Support Vector Machines

- Dual optimization problem

$$\begin{aligned} \max. \quad & \sum_{i=1}^N y_i (\alpha_i^+ - \alpha_i^-) - \epsilon \sum_{i=1}^N (\alpha_i^+ + \alpha_i^-) \\ & - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N (\alpha_i^+ - \alpha_i^-) (\alpha_j^+ - \alpha_j^-) k(\mathbf{x}_i, \mathbf{x}_j) \end{aligned}$$

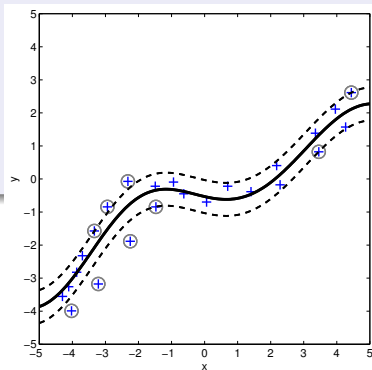
$$\text{w.r.t. } \boldsymbol{\alpha}^+ \in \mathbb{R}_+^N, \boldsymbol{\alpha}^- \in \mathbb{R}_+^N$$

$$\text{s.t. } \sum_{i=1}^N (\alpha_i^+ - \alpha_i^-) = 0$$

$$C \geq \alpha_i^+, \alpha_i^- \geq 0 \quad \forall i$$

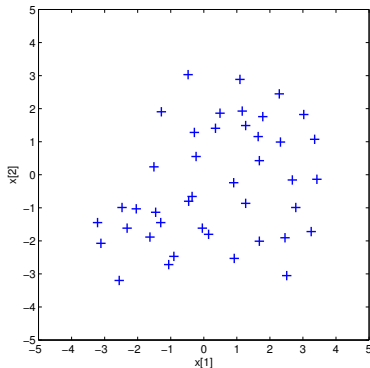
- $\mathbf{w} = \sum_{i=1}^N (\alpha_i^+ - \alpha_i^-) \Phi(\mathbf{x}_i)$

- $f(\mathbf{x}) = \sum_{i=1}^N (\alpha_i^+ - \alpha_i^-) k(\mathbf{x}_i, \mathbf{x}) + b$



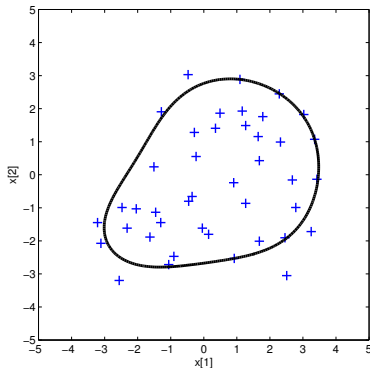
# One-Class Classification Support Vector Machines

- $\{\mathbf{x}_i\}_{i=1}^N$  where  $\mathbf{x}_i \in \mathbb{R}^D$



# One-Class Classification Support Vector Machines

- $\{\mathbf{x}_i\}_{i=1}^N$  where  $\mathbf{x}_i \in \mathbb{R}^D$
- $f(\mathbf{x}) = \langle \mathbf{w}, \Phi(\mathbf{x}) \rangle + b$  where  $\Phi: \mathbb{R}^D \rightarrow \mathbb{R}^S$



# One-Class Classification Support Vector Machines

- $\{\mathbf{x}_i\}_{i=1}^N$  where  $\mathbf{x}_i \in \mathbb{R}^D$
- $f(\mathbf{x}) = \langle \mathbf{w}, \Phi(\mathbf{x}) \rangle + b$  where  $\Phi: \mathbb{R}^D \rightarrow \mathbb{R}^S$
- Primal optimization problem (Schölkopf and Smola, 2002)

$$\begin{aligned} \min. \quad & \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_{i=1}^N \xi_i + b \\ \text{w.r.t.} \quad & \mathbf{w} \in \mathbb{R}^S, \boldsymbol{\xi} \in \mathbb{R}_+^N, b \in \mathbb{R} \\ \text{s.t.} \quad & \langle \mathbf{w}, \Phi(\mathbf{x}_i) \rangle + b + \xi_i \geq 0 \quad \forall i \end{aligned}$$

# One-Class Classification Support Vector Machines

- Dual optimization problem

$$\max. \quad -\frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j k(\mathbf{x}_i, \mathbf{x}_j)$$

$$\text{w.r.t. } \boldsymbol{\alpha} \in \mathbb{R}_+^N$$

$$\text{s.t. } \sum_{i=1}^N \alpha_i = 1$$

$$C \geq \alpha_i \geq 0 \quad \forall i$$

# One-Class Classification Support Vector Machines

- Dual optimization problem

$$\max. \quad -\frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j k(\mathbf{x}_i, \mathbf{x}_j)$$

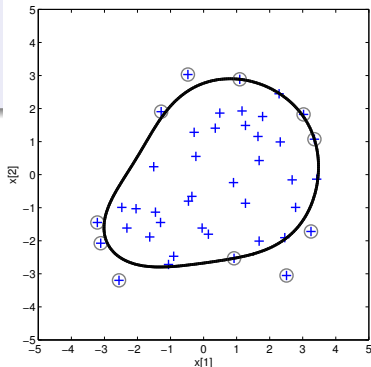
$$\text{w.r.t. } \boldsymbol{\alpha} \in \mathbb{R}_+^N$$

$$\text{s.t. } \sum_{i=1}^N \alpha_i = 1$$

$$C \geq \alpha_i \geq 0 \quad \forall i$$

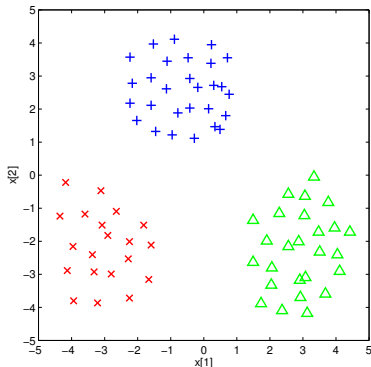
- $\mathbf{w} = \sum_{i=1}^N \alpha_i \Phi(\mathbf{x}_i)$

- $f(\mathbf{x}) = \sum_{i=1}^N \alpha_i k(\mathbf{x}_i, \mathbf{x}) + b$



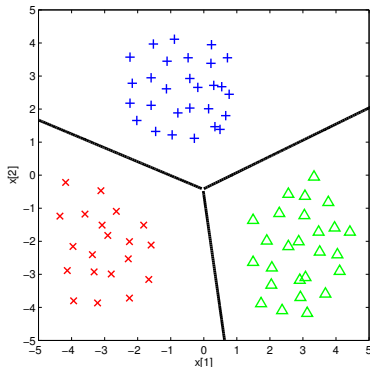
# Multiclass Classification Support Vector Machines

- $\{(\mathbf{x}_i, y_i)\}_{i=1}^N$  where  $\mathbf{x}_i \in \mathbb{R}^D$  and  $y_i \in \{1, \dots, K\}$



# Multiclass Classification Support Vector Machines

- $\{(\mathbf{x}_i, y_i)\}_{i=1}^N$  where  $\mathbf{x}_i \in \mathbb{R}^D$  and  $y_i \in \{1, \dots, K\}$
- $f^l(\mathbf{x}) = \langle \mathbf{w}^l, \Phi(\mathbf{x}) \rangle + b^l \quad \forall l$  where  $\Phi: \mathbb{R}^D \rightarrow \mathbb{R}^S$



# Multiclass Classification Support Vector Machines

- $\{(\mathbf{x}_i, y_i)\}_{i=1}^N$  where  $\mathbf{x}_i \in \mathbb{R}^D$  and  $y_i \in \{1, \dots, K\}$
- $f^l(\mathbf{x}) = \langle \mathbf{w}^l, \Phi(\mathbf{x}) \rangle + b^l \quad \forall l$  where  $\Phi: \mathbb{R}^D \rightarrow \mathbb{R}^S$
- Primal optimization problem (Vapnik, 1998; Weston and Watkins, 1998; Bredensteiner and Bennett, 1999)

$$\begin{aligned} \min. \quad & \frac{1}{2} \sum_{l=1}^K \|\mathbf{w}^l\|_2^2 + C \sum_{i=1}^N \sum_{l=1}^K \xi_i^l \\ \text{w.r.t.} \quad & \mathbf{w}^l \in \mathbb{R}^S, \boldsymbol{\xi}^l \in \mathbb{R}_+^N, b^l \in \mathbb{R} \\ \text{s.t.} \quad & f^{y_i}(\mathbf{x}_i) - f^l(\mathbf{x}_i) \geq 2 - \xi_i^l \quad \forall (i, l \neq y_i) \\ & \xi_i^{y_i} = 0 \quad \forall i \end{aligned}$$

# Multiclass Classification Support Vector Machines

## ■ Dual optimization problem

$$\max. \quad 2 \sum_{i=1}^N \sum_{l=1}^K \alpha_i^l - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \left( \delta_{y_i}^{y_j} A_i A_j - \sum_{l=1}^K \alpha_i^l (2\alpha_j^{y_i} - \alpha_j^l) \right) k(\mathbf{x}_i, \mathbf{x}_j)$$

$$\text{w.r.t. } \boldsymbol{\alpha}^l \in \mathbb{R}_+^N$$

$$\text{s.t. } \sum_{i=1}^N \alpha_i^l - \sum_{i=1}^N \delta_{y_i}^l A_i = 0 \quad \forall l \quad \text{where } A_i = \sum_{l=1}^K \alpha_i^l$$

$$(1 - \delta_{y_i}^l) C \geq \alpha_i^l \geq 0 \quad \forall (i, l)$$

# Multiclass Classification Support Vector Machines

- Dual optimization problem

$$\max. \quad 2 \sum_{i=1}^N \sum_{l=1}^K \alpha_i^l - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \left( \delta_{y_i y_j}^l A_i A_j - \sum_{l=1}^K \alpha_i^l (2\alpha_j^{y_i} - \alpha_j^l) \right) k(\mathbf{x}_i, \mathbf{x}_j)$$

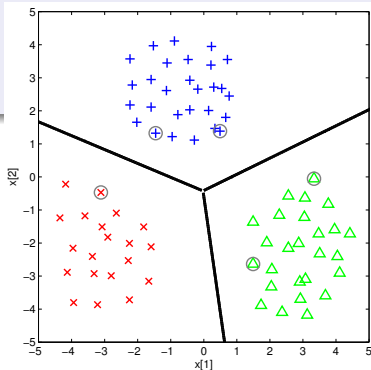
$$\text{w.r.t. } \alpha^l \in \mathbb{R}_+^N$$

$$\text{s.t. } \sum_{i=1}^N \alpha_i^l - \sum_{i=1}^N \delta_{y_i}^l A_i = 0 \quad \forall l$$

$$(1 - \delta_{y_i}^l) C \geq \alpha_i^l \geq 0 \quad \forall (i, l)$$

- $\mathbf{w}^l = \sum_{i=1}^N (\delta_{y_i}^l A_i - \alpha_i^l) \Phi(\mathbf{x}_i)$

- $f^l(\mathbf{x}) = \sum_{i=1}^N (\delta_{y_i}^l A_i - \alpha_i^l) k(\mathbf{x}_i, \mathbf{x}) + b^l$



# Posterior Probability Support Vector Machines

- Soft labels instead of hard labels (Tao et al., 2005)

$$\hat{y}_i = 2 \Pr(+1|\mathbf{x}_i) - 1 \quad \forall i$$

- $\hat{y}_i(\langle \mathbf{w}, \Phi(\mathbf{x}_i) \rangle + b) \geq \hat{y}_i^2 - \hat{y}_i^2 \xi_i \quad \forall i$

# Posterior Probability Support Vector Machines

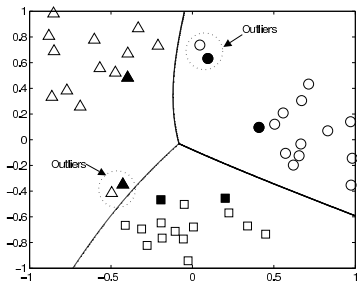
- Soft labels instead of hard labels (Tao et al., 2005)

$$\hat{y}_i = 2 \Pr(+1|\mathbf{x}_i) - 1 \quad \forall i$$

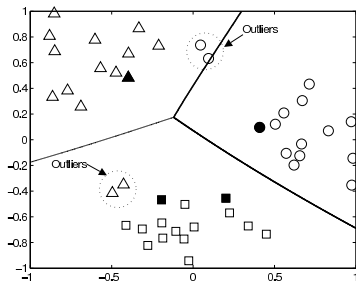
- $\hat{y}_i(\langle \mathbf{w}, \Phi(\mathbf{x}_i) \rangle + b) \geq \hat{y}_i^2 - \hat{y}_i^2 \xi_i \quad \forall i$

- Extended to multiclass classification (Gönen et al., 2008)

$$f^{y_i}(\mathbf{x}_i) - f^l(\mathbf{x}_i) \geq 2(\Pr(y_i|\mathbf{x}_i) - \Pr(l|\mathbf{x}_i)) - \xi_i^l \quad \forall (i, l \neq y_i)$$



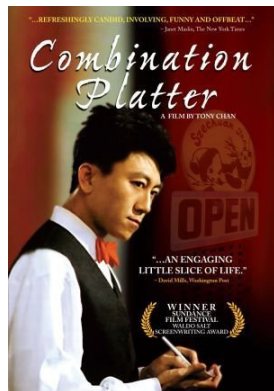
MCSVM



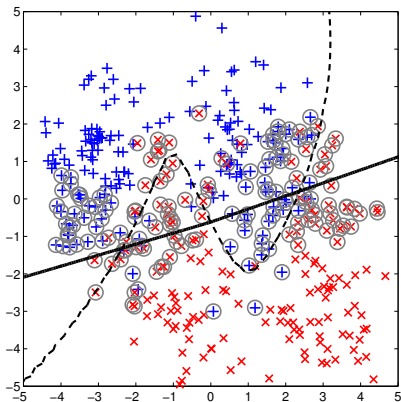
MCPSPVM

# Outline

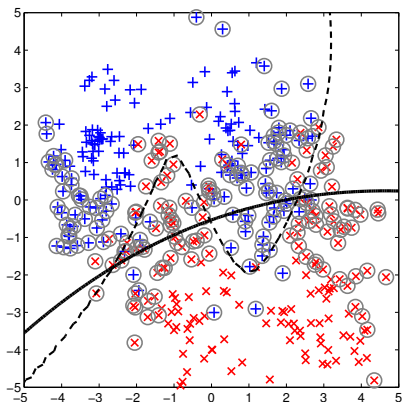
1. Introduction
2. Multiple Kernel Learning
3. Localized Multiple Kernel Learning
4. Local Projection Kernels
5. Experiments
6. Conclusions



# Multiple Kernel Learning

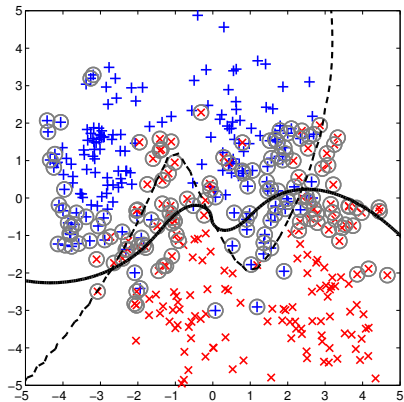


SVM ( $k_L$ )



SVM ( $k_P$ )

# Multiple Kernel Learning



MKL ( $k_L - k_P$ )

- Bayes optimal  
Accuracy: 91.25
- SVM ( $k_L$ )  
Accuracy:  $90.05 \pm 0.35$   
Support Vector:  $39.78 \pm 2.63$
- SVM ( $k_P$ )  
Accuracy:  $90.10 \pm 0.46$   
Support Vector:  $53.65 \pm 1.74$
- MKL ( $k_L - k_P$ )  
Accuracy:  $90.95 \pm 0.61$   
Support Vector:  $38.23 \pm 2.34$   
 $(\eta_L - \eta_P) = (0.31 - 0.69)$

# Multiple Kernel Learning

- $P$  feature representations (not necessarily different)

$$\mathbf{x}_i = \{\mathbf{x}_i^m\}_{m=1}^P \quad \text{where } \mathbf{x}_i^m \in \mathbb{R}^{D_m}$$

- MULTIPLE KERNEL LEARNING (MKL)

$$k_\eta(\mathbf{x}_i, \mathbf{x}_j) = f_\eta(\{k_m(\mathbf{x}_i^m, \mathbf{x}_j^m)\}_{m=1}^P | \Theta) \quad \text{where } f_\eta: \mathbb{R}^P \rightarrow \mathbb{R}$$

# Multiple Kernel Learning

- $P$  feature representations (not necessarily different)

$$\mathbf{x}_i = \{\mathbf{x}_i^m\}_{m=1}^P \quad \text{where } \mathbf{x}_i^m \in \mathbb{R}^{D_m}$$

- MULTIPLE KERNEL LEARNING (MKL)

$$k_\eta(\mathbf{x}_i, \mathbf{x}_j) = f_\eta(\{k_m(\mathbf{x}_i^m, \mathbf{x}_j^m)\}_{m=1}^P | \Theta) \quad \text{where } f_\eta: \mathbb{R}^P \rightarrow \mathbb{R}$$

- Why MKL?

- ▶ Different notions of similarity

$$\mathbf{x}_i^1 = \mathbf{x}_i^2 = \dots = \mathbf{x}_i^P$$
$$k_1 \neq k_2 \neq \dots \neq k_P$$

- ▶ Different representations from different sources or modalities

$$\mathbf{x}_i^1 \neq \mathbf{x}_i^2 \neq \dots \neq \mathbf{x}_i^P$$
$$k_1 = k_2 = \dots = k_P$$

# Multiple Kernel Learning

- Unweighted sum

$$k_{\eta}(\mathbf{x}_i, \mathbf{x}_j) = \sum_{m=1}^P k_m(\mathbf{x}_i^m, \mathbf{x}_j^m)$$

(Pavlidis et al., 2001)

# Multiple Kernel Learning

- Unweighted sum

$$k_{\eta}(\mathbf{x}_i, \mathbf{x}_j) = \sum_{m=1}^P k_m(\mathbf{x}_i^m, \mathbf{x}_j^m)$$

(Pavlidis et al., 2001)

- Weighted sum

$$k_{\eta}(\mathbf{x}_i, \mathbf{x}_j | \Theta = \eta) = \sum_{m=1}^P \eta_m k_m(\mathbf{x}_i^m, \mathbf{x}_j^m)$$

(Lanckriet et al., 2004; Bach et al., 2004)

# Multiple Kernel Learning

- Unweighted sum

$$k_{\eta}(\mathbf{x}_i, \mathbf{x}_j) = \sum_{m=1}^P k_m(\mathbf{x}_i^m, \mathbf{x}_j^m)$$

(Pavlidis et al., 2001)

- Weighted sum

$$k_{\eta}(\mathbf{x}_i, \mathbf{x}_j | \Theta = \eta) = \sum_{m=1}^P \eta_m k_m(\mathbf{x}_i^m, \mathbf{x}_j^m)$$

(Lanckriet et al., 2004; Bach et al., 2004)

- Locally weighted sum

$$k_{\eta}(\mathbf{x}_i, \mathbf{x}_j | \Theta = \mathbf{V}) = \sum_{m=1}^P \eta_m(\mathbf{x}_i | \mathbf{V}) k_m(\mathbf{x}_i^m, \mathbf{x}_j^m) \eta_m(\mathbf{x}_j | \mathbf{V})$$

(Gönen and Alpaydın, 2008)

# Multiple Kernel Learning

- Weighted MKL formulation of Bach et al. (2004)

- $$f(\mathbf{x}) = \sum_{m=1}^P \langle \mathbf{w}_m, \Phi_m(\mathbf{x}_i^m) \rangle + b$$

# Multiple Kernel Learning

- Weighted MKL formulation of Bach et al. (2004)

$$\blacksquare f(\mathbf{x}) = \sum_{m=1}^P \langle \mathbf{w}_m, \Phi_m(\mathbf{x}_i^m) \rangle + b$$

- Primal optimization problem

$$\min. \frac{1}{2} \left( \sum_{m=1}^P d_m \|\mathbf{w}_m\|_2 \right)^2 + C \sum_{i=1}^N \xi_i$$

$$\text{w.r.t. } \mathbf{w}_m \in \mathbb{R}^{S_m}, \boldsymbol{\xi} \in \mathbb{R}_+^N, b \in \mathbb{R}$$

$$\text{s.t. } y_i \left( \sum_{m=1}^P \langle \mathbf{w}_m, \Phi_m(\mathbf{x}_i^m) \rangle + b \right) \geq 1 - \xi_i \quad \forall i$$

- weighted  $l_1$ -norm on feature spaces
- $l_2$ -norm within each feature space

# Multiple Kernel Learning

## ■ Dual optimization problem

$$\min. \frac{1}{2}\gamma^2 - \sum_{i=1}^N \alpha_i$$

$$\text{w.r.t. } \gamma \in \mathbb{R}, \boldsymbol{\alpha} \in \mathbb{R}_+^N$$

$$\text{s.t. } \gamma^2 d_m^2 \geq \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j k_m(\mathbf{x}_i^m, \mathbf{x}_j^m) \quad \forall m$$

$$\sum_{i=1}^N \alpha_i y_i = 0$$

$$C \geq \alpha_i \geq 0 \quad \forall i$$

# Multiple Kernel Learning

## ■ Dual optimization problem

$$\min. \quad \frac{1}{2}\gamma^2 - \sum_{i=1}^N \alpha_i$$

$$\text{w.r.t. } \gamma \in \mathbb{R}, \boldsymbol{\alpha} \in \mathbb{R}_+^N$$

$$\text{s.t. } \gamma^2 d_m^2 \geq \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j k_m(\mathbf{x}_i^m, \mathbf{x}_j^m) \quad \forall m \quad \eta_m$$

$$\sum_{i=1}^N \alpha_i y_i = 0$$

$$C \geq \alpha_i \geq 0 \quad \forall i$$

$$\blacksquare \quad \mathbf{w}_m = \sum_{i=1}^N \alpha_i y_i \eta_m \Phi_m(\mathbf{x}_i^m)$$

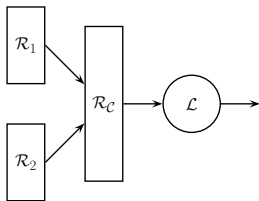
$$\blacksquare \quad f(\mathbf{x}) = \sum_{i=1}^N \alpha_i y_i \sum_{m=1}^P \eta_m k_m(\mathbf{x}_i^m, \mathbf{x}^m) + b \quad \text{where } \sum_{m=1}^P d_m^2 \eta_m = 1$$

# Multiple Kernel Learning as Intermediate Integration

- Three methods for combining multiple representations (Noble, 2004)

# Multiple Kernel Learning as Intermediate Integration

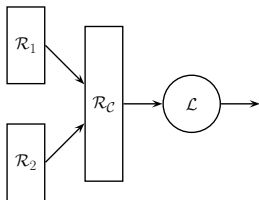
- Three methods for combining multiple representations (Noble, 2004)
  - ▶ Early integration



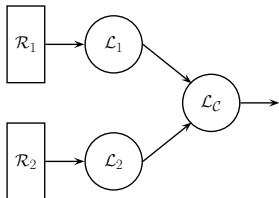
# Multiple Kernel Learning as Intermediate Integration

- Three methods for combining multiple representations (Noble, 2004)

- ▶ Early integration



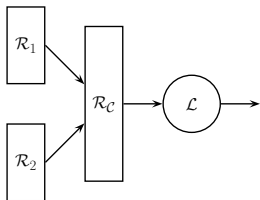
- ▶ Late integration



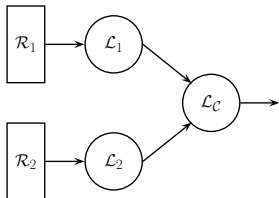
# Multiple Kernel Learning as Intermediate Integration

- Three methods for combining multiple representations (Noble, 2004)

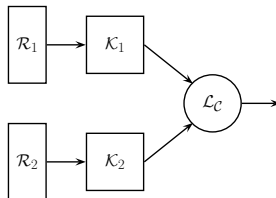
- ▶ Early integration



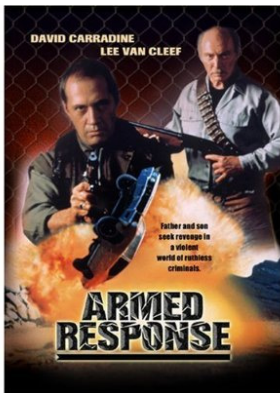
- ▶ Late integration



- ▶ Intermediate integration

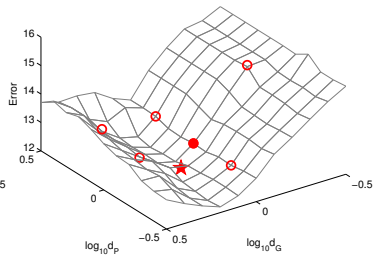
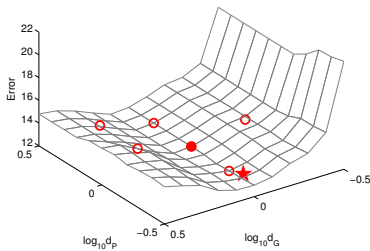


# Regularizing Multiple Kernel Learning Using Response Surface Methodology



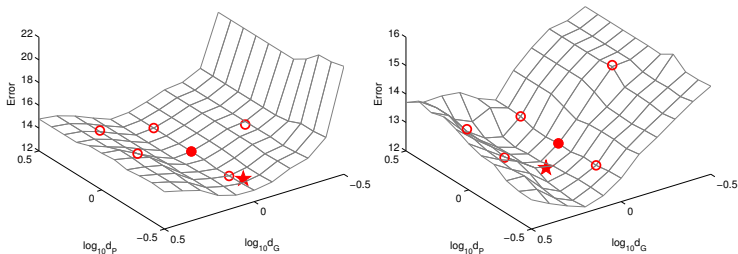
# Regularized Multiple Kernel Learning

- $\{d_m\}_{m=1}^P$  parameters in the formulation of Bach et al. (2004)



# Regularized Multiple Kernel Learning

- $\{d_m\}_{m=1}^P$  parameters in the formulation of Bach et al. (2004)



- Second-order model in the log scale

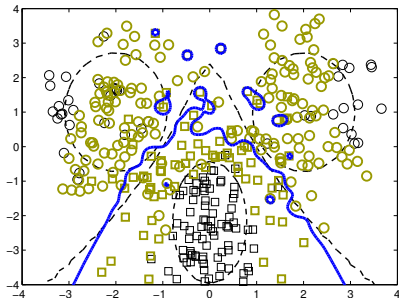
$$\text{error} = \beta_0 + \sum_{m=2}^P \beta_m \log_{10} d_m + \sum_{m=2}^P \sum_{h=m}^P \beta_{mh} \log_{10} d_m \log_{10} d_h$$

- Grid search versus RSM

- ▶  $L^{(P-1)}$  points versus  $\mathcal{O}(P^2)$  points

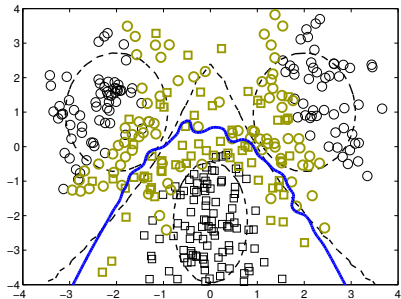
- ▶ predefined factor levels versus arbitrary factor values

# Regularized Multiple Kernel Learning



MKL

$$(\eta_L - \eta_P - \eta_G) = (0.00 - 0.25 - 0.75)$$

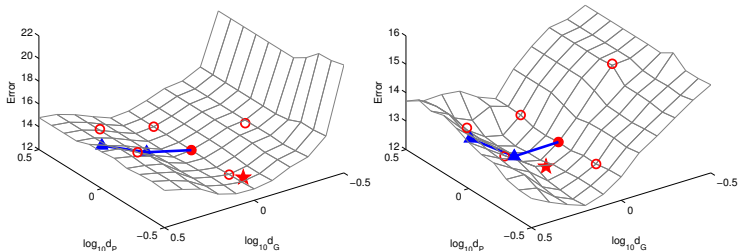


RMKL

$$(\eta_L - \eta_P - \eta_G) = (0.00 - 0.48 - 0.31)$$

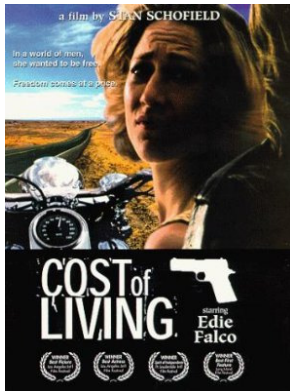
# Regularized Multiple Kernel Learning

- Bach et al. (2005) propose to select  $\{d_m\}_{m=1}^P$  by looking at eigenvalues of kernel matrices (EMKL)



- Perform RSM on  $\{\eta_m\}_{m=1}^P$  directly (RWKL)

# Cost-Conscious Multiple Kernel Learning



## Cost-Conscious Multiple Kernel Learning

- $\{d_m\}_{m=1}^P$  parameters in the formulation of Bach et al. (2004)
- Can be interpreted as the costs of kernels (Gönen and Alpaydın, 2010a)

# Cost-Conscious Multiple Kernel Learning

- $\{d_m\}_{m=1}^P$  parameters in the formulation of Bach et al. (2004)
- Can be interpreted as the costs of kernels (Gönen and Alpaydın, 2010a)
- The costs of evaluating kernels
  - ▶ preprocessing data
  - ▶ computation time

# Cost-Conscious Multiple Kernel Learning

- $\{d_m\}_{m=1}^P$  parameters in the formulation of Bach et al. (2004)
- Can be interpreted as the costs of kernels (Gönen and Alpaydın, 2010a)
- The costs of evaluating kernels
  - ▶ preprocessing data
  - ▶ computation time
- The costs of extracting/sensing the corresponding representations/signals
  - ▶ bioinformatics: protein sequence, protein structure, etc.
  - ▶ biometrics: face, fingerprint, iris, signature, etc.
  - ▶ speech recognition: acoustic input, visual input, etc.

# Cost-Conscious Multiple Kernel Learning

## ■ Representation selection



DYN



STA4



STA8



STA16

# Cost-Conscious Multiple Kernel Learning

## ■ Representation selection



DYN



STA4



STA8



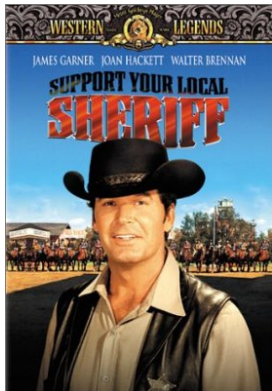
STA16

## ■ Kernel/representation selection

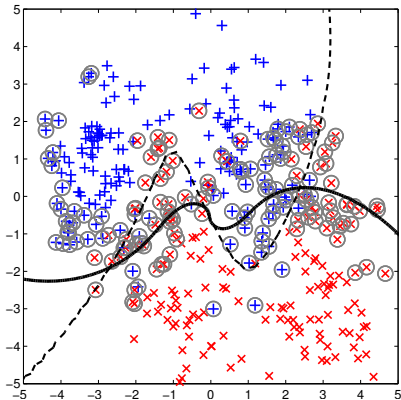
Kernel	Explanation	Data Source
$k_{SW}$	Smith-Waterman kernel	Protein sequences
$k_B$	BLAST kernel	Protein sequences
$k_{Pfam}$	Pfam kernel	Protein sequences
$k_{FFT}$	FFT kernel	Hydropathy profiles
$k_{LI}$	Linear kernel	Protein interactions
$k_D$	Diffusion kernel	Protein interactions
$k_E$	Gaussian kernel	Gene expression profiles

# Outline

1. Introduction
2. Multiple Kernel Learning
- 3. Localized Multiple Kernel Learning**
4. Local Projection Kernels
5. Experiments
6. Conclusions



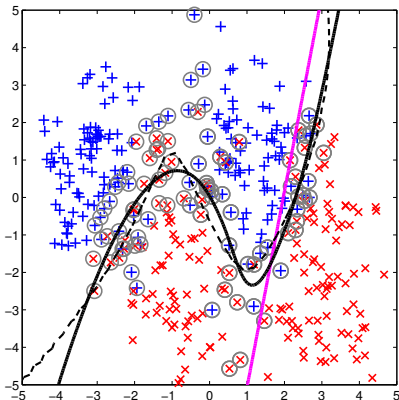
# Localized Multiple Kernel Learning



MKL ( $k_L - k_P$ )

Accuracy:  $90.95 \pm 0.61$

Support Vector:  $38.23 \pm 2.34$



LMKL ( $k_L - k_P$ )

Accuracy:  $91.83 \pm 0.24$

Support Vector:  $25.12 \pm 0.91$

## Localized Multiple Kernel Learning

- $f(\mathbf{x}) = \sum_{m=1}^P \eta_m(\mathbf{x}|\mathbf{V}) \langle \mathbf{w}_m, \Phi_m(\mathbf{x}^m) \rangle + b$     where  $\eta_m: \mathbb{R}^{D_G} \rightarrow \mathbb{R}$
- Similar to (but also different from) “mixture of experts” (Jacobs et al., 1991)

## Localized Multiple Kernel Learning

- $f(\mathbf{x}) = \sum_{m=1}^P \eta_m(\mathbf{x}|\mathbf{V}) \langle \mathbf{w}_m, \Phi_m(\mathbf{x}^m) \rangle + b$  where  $\eta_m: \mathbb{R}^{D_G} \rightarrow \mathbb{R}$

- Similar to (but also different from) “mixture of experts” (Jacobs et al., 1991)

- Softmax gating

$$\eta_m(\mathbf{x}|\mathbf{V}) = \frac{\exp(\langle \mathbf{v}_m, \mathbf{x}^G \rangle + v_{m0})}{\sum_{h=1}^P \exp(\langle \mathbf{v}_h, \mathbf{x}^G \rangle + v_{h0})} \quad \forall m$$

# Localized Multiple Kernel Learning

- $f(\mathbf{x}) = \sum_{m=1}^P \eta_m(\mathbf{x}|\mathbf{V}) \langle \mathbf{w}_m, \Phi_m(\mathbf{x}^m) \rangle + b$  where  $\eta_m: \mathbb{R}^{D_G} \rightarrow \mathbb{R}$

- Similar to (but also different from) “mixture of experts” (Jacobs et al., 1991)

- Softmax gating

$$\eta_m(\mathbf{x}|\mathbf{V}) = \frac{\exp(\langle \mathbf{v}_m, \mathbf{x}^G \rangle + v_{m0})}{\sum_{h=1}^P \exp(\langle \mathbf{v}_h, \mathbf{x}^G \rangle + v_{h0})} \quad \forall m$$

- Sigmoid gating

$$\eta_m(\mathbf{x}|\mathbf{V}) = \frac{1}{1 + \exp(-\langle \mathbf{v}_m, \mathbf{x}^G \rangle - v_{m0})} \quad \forall m$$

# Localized Multiple Kernel Learning

- Primal optimization problem (Gönen and Alpaydm, 2008)

$$\min. \frac{1}{2} \sum_{m=1}^P \|\mathbf{w}_m\|_2^2 + C \sum_{i=1}^N \xi_i$$

$$\text{w.r.t. } \mathbf{w}_m \in \mathbb{R}^{S_m}, \boldsymbol{\xi} \in \mathbb{R}_+^N, \mathbf{V} \in \mathbb{R}^{P \times (D_G + 1)}, b \in \mathbb{R}$$

$$\text{s.t. } y_i \left( \sum_{m=1}^P \eta_m(\mathbf{x}_i | \mathbf{V}) \langle \mathbf{w}_m, \Phi_m(\mathbf{x}_i^m) \rangle + b \right) \geq 1 - \xi_i \quad \forall i$$

# Localized Multiple Kernel Learning

- Primal optimization problem (Gönen and Alpaydm, 2008)

$$\min. \frac{1}{2} \sum_{m=1}^P \|\mathbf{w}_m\|_2^2 + C \sum_{i=1}^N \xi_i$$

$$\text{w.r.t. } \mathbf{w}_m \in \mathbb{R}^{S_m}, \boldsymbol{\xi} \in \mathbb{R}_+^N, \mathbf{V} \in \mathbb{R}^{P \times (D_G+1)}, b \in \mathbb{R}$$

$$\text{s.t. } y_i \left( \sum_{m=1}^P \eta_m(\mathbf{x}_i | \mathbf{V}) \langle \mathbf{w}_m, \Phi_m(\mathbf{x}_i^m) \rangle + b \right) \geq 1 - \xi_i \quad \forall i$$

- Nonconvexity due to gating model
- For a given  $\mathbf{V}$ , problem becomes convex

# Localized Multiple Kernel Learning

- Dual optimization problem for a given  $\mathbf{V}$

$$\max. J(\mathbf{V}) = \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j k_{\eta}(\mathbf{x}_i, \mathbf{x}_j)$$

$$\text{w.r.t. } \boldsymbol{\alpha} \in \mathbb{R}_+^N$$

$$\text{s.t. } \sum_{i=1}^N \alpha_i y_i = 0$$

$$C \geq \alpha_i \geq 0 \quad \forall i$$

# Localized Multiple Kernel Learning

- Dual optimization problem for a given  $\mathbf{V}$

$$\max. J(\mathbf{V}) = \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j k_{\eta}(\mathbf{x}_i, \mathbf{x}_j)$$

$$\text{w.r.t. } \boldsymbol{\alpha} \in \mathbb{R}_+^N$$

$$\text{s.t. } \sum_{i=1}^N \alpha_i y_i = 0$$

$$C \geq \alpha_i \geq 0 \quad \forall i$$

- *locally combined kernel function*

$$k_{\eta}(\mathbf{x}_i, \mathbf{x}_j) = \sum_{m=1}^P \eta_m(\mathbf{x}_i | \mathbf{V}) \underbrace{\langle \Phi_m(\mathbf{x}_i^m), \Phi_m(\mathbf{x}_j^m) \rangle}_{k_m(\mathbf{x}_i^m, \mathbf{x}_j^m)} \eta_m(\mathbf{x}_j | \mathbf{V})$$

# Localized Multiple Kernel Learning

- 1: Initialize  $\mathbf{V}^{(0)}$  randomly
- 2: **repeat**
- 3: Calculate  $\mathbf{K}_\eta^{(t)} = \{k_\eta(\mathbf{x}_i, \mathbf{x}_j)\}_{i,j=1}^N$  using  $\mathbf{V}^{(t)}$
- 4: Solve kernel machine with  $\mathbf{K}_\eta^{(t)}$
- 5: Determine step size,  $\Delta^{(t)}$ , using Armijo's rule
- 6:  $\mathbf{V}^{(t+1)} \leftarrow \mathbf{V}^{(t)} - \Delta^{(t)} \frac{\partial J(\mathbf{V})}{\partial \mathbf{V}}$
- 7: **until** convergence

$$\blacksquare f(\mathbf{x}) = \sum_{i=1}^N \alpha_i y_i \underbrace{\sum_{m=1}^P \eta_m(\mathbf{x}_i | \mathbf{V}) k_m(\mathbf{x}_i^m, \mathbf{x}^m) \eta_m(\mathbf{x} | \mathbf{V})}_{k_\eta(\mathbf{x}_i, \mathbf{x})} + b$$

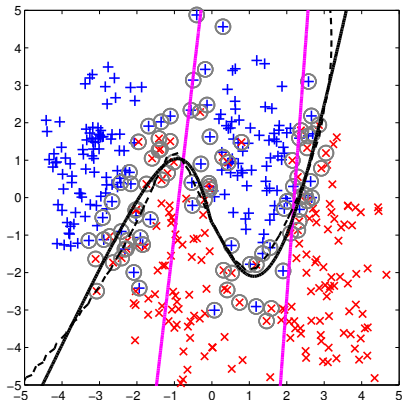
# Localized Multiple Kernel Learning

- 1: Initialize  $\mathbf{V}^{(0)}$  randomly
- 2: **repeat**
- 3: Calculate  $\mathbf{K}_\eta^{(t)} = \{k_\eta(\mathbf{x}_i, \mathbf{x}_j)\}_{i,j=1}^N$  using  $\mathbf{V}^{(t)}$
- 4: Solve kernel machine with  $\mathbf{K}_\eta^{(t)}$
- 5: Determine step size,  $\Delta^{(t)}$ , using Armijo's rule
- 6:  $\mathbf{V}^{(t+1)} \leftarrow \mathbf{V}^{(t)} - \Delta^{(t)} \frac{\partial J(\mathbf{V})}{\partial \mathbf{V}}$
- 7: **until** convergence

$$\blacksquare f(\mathbf{x}) = \sum_{i=1}^N \alpha_i y_i \underbrace{\sum_{m=1}^P \eta_m(\mathbf{x}_i | \mathbf{V}) k_m(\mathbf{x}_i^m, \mathbf{x}^m) \eta_m(\mathbf{x} | \mathbf{V})}_{k_\eta(\mathbf{x}_i, \mathbf{x})} + b$$

- Can be generalized to other kernel-based methods

# Localized Multiple Kernel Learning

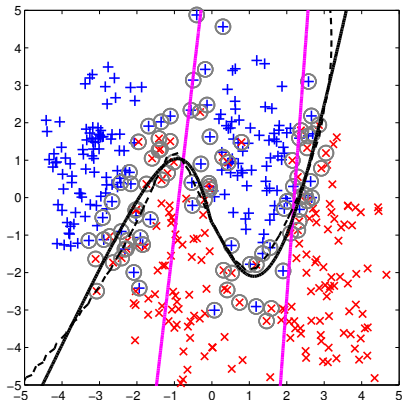


LMKL ( $k_L-k_L-k_L$ )

Accuracy:  $91.78 \pm 0.55$

Support Vector:  $23.82 \pm 1.20$

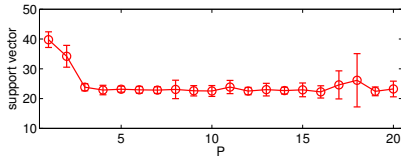
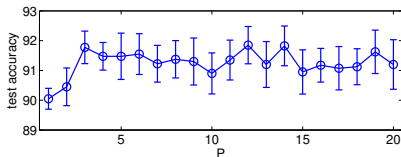
# Localized Multiple Kernel Learning



LMKL ( $k_L - k_L - k_L$ )

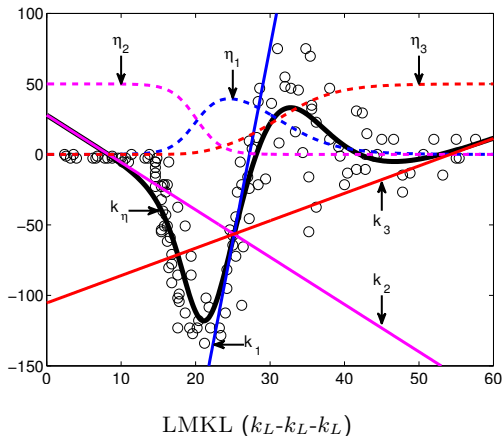
Accuracy:  $91.78 \pm 0.55$

Support Vector:  $23.82 \pm 1.20$



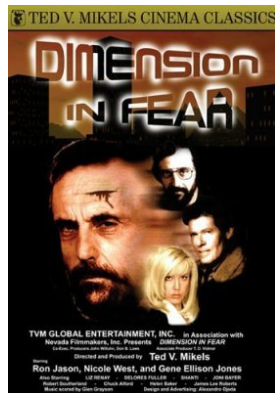
# Localized Multiple Kernel Regression

- LMKL on regression problems (Gönen and Alpaydm, 2010b)
- MOTORCYCLE data set



# Outline

1. Introduction
2. Multiple Kernel Learning
3. Localized Multiple Kernel Learning
4. Local Projection Kernels
5. Experiments
6. Conclusions



# Global Projection Kernels

- $z = \mathbf{W}^\top \mathbf{x}$     where  $\mathbf{W} \in \mathbb{R}^{D \times R}$

- $f(\mathbf{x}) = \langle \mathbf{w}, \Phi(\mathbf{W}^\top \mathbf{x}) \rangle + b$     where  $\Phi: \mathbb{R}^R \rightarrow \mathbb{R}^S$

# Global Projection Kernels

- $z = \mathbf{W}^\top \mathbf{x}$  where  $\mathbf{W} \in \mathbb{R}^{D \times R}$
- $f(\mathbf{x}) = \langle \mathbf{w}, \Phi(\mathbf{W}^\top \mathbf{x}) \rangle + b$  where  $\Phi: \mathbb{R}^R \rightarrow \mathbb{R}^S$
- Primal optimization problem (Chapelle et al., 2002)

$$\begin{aligned} \min. \quad & \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_{i=1}^N \xi_i \\ \text{w.r.t. } & \mathbf{w} \in \mathbb{R}^S, \boldsymbol{\xi} \in \mathbb{R}_+^N, \mathbf{W} \in \mathbb{R}^{D \times R}, b \in \mathbb{R} \\ \text{s.t. } & y_i (\langle \mathbf{w}, \Phi(\mathbf{W}^\top \mathbf{x}_i) \rangle + b) \geq 1 - \xi_i \quad \forall i \end{aligned}$$

# Global Projection Kernels

- $\mathbf{z} = \mathbf{W}^\top \mathbf{x}$  where  $\mathbf{W} \in \mathbb{R}^{D \times R}$
- $f(\mathbf{x}) = \langle \mathbf{w}, \Phi(\mathbf{W}^\top \mathbf{x}) \rangle + b$  where  $\Phi: \mathbb{R}^R \rightarrow \mathbb{R}^S$
- Primal optimization problem (Chapelle et al., 2002)

$$\begin{aligned} \min. \quad & \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_{i=1}^N \xi_i \\ \text{w.r.t. } & \mathbf{w} \in \mathbb{R}^S, \boldsymbol{\xi} \in \mathbb{R}_+^N, \mathbf{W} \in \mathbb{R}^{D \times R}, b \in \mathbb{R} \\ \text{s.t. } & y_i (\langle \mathbf{w}, \Phi(\mathbf{W}^\top \mathbf{x}_i) \rangle + b) \geq 1 - \xi_i \quad \forall i \end{aligned}$$

- For a given  $\mathbf{W}$ , problem becomes convex

## Global Projection Kernels

- Dual optimization problem for a given  $\mathbf{W}$

$$\max. J(\mathbf{W}) = \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j k(\mathbf{z}_i, \mathbf{z}_j)$$

$$\text{w.r.t. } \boldsymbol{\alpha} \in \mathbb{R}_+^N$$

$$\text{s.t. } \sum_{i=1}^N \alpha_i y_i = 0$$

$$C \geq \alpha_i \geq 0 \quad \forall i$$

# Global Projection Kernels

- Dual optimization problem for a given  $\mathbf{W}$

$$\max. J(\mathbf{W}) = \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j k(\mathbf{z}_i, \mathbf{z}_j)$$

$$\text{w.r.t. } \boldsymbol{\alpha} \in \mathbb{R}_+^N$$

$$\text{s.t. } \sum_{i=1}^N \alpha_i y_i = 0$$

$$C \geq \alpha_i \geq 0 \quad \forall i$$

- *global projection kernel function*

$$k(\mathbf{z}_i, \mathbf{z}_j) = \underbrace{\langle \Phi(\mathbf{W}^\top \mathbf{x}_i), \Phi(\mathbf{W}^\top \mathbf{x}_j) \rangle}_{k(\mathbf{W}^\top \mathbf{x}_i, \mathbf{W}^\top \mathbf{x}_j)}$$

# Global Projection Kernels

- 1: Initialize  $\mathbf{W}^{(0)}$  randomly
- 2: **repeat**
- 3: Calculate  $\mathbf{K}^{(t)} = \{k(\mathbf{z}_i, \mathbf{z}_j)\}_{i,j=1}^N$  using  $\mathbf{W}^{(t)}$
- 4: Solve kernel machine with  $\mathbf{K}^{(t)}$
- 5: Determine step size,  $\mu^{(t)}$ , using Armijo's rule
- 6:  $\mathbf{W}^{(t+1)} \leftarrow \mathbf{W}^{(t)} - \mu^{(t)} \frac{\partial J(\mathbf{W})}{\partial \mathbf{W}}$
- 7: **until** convergence

$$\blacksquare f(\mathbf{x}) = \sum_{i=1}^N \alpha_i y_i k(\mathbf{W}^\top \mathbf{x}_i, \mathbf{W}^\top \mathbf{x}) + b$$

# Global Projection Kernels

- 1: Initialize  $\mathbf{W}^{(0)}$  randomly
- 2: **repeat**
- 3: Calculate  $\mathbf{K}^{(t)} = \{k(\mathbf{z}_i, \mathbf{z}_j)\}_{i,j=1}^N$  using  $\mathbf{W}^{(t)}$
- 4: Solve kernel machine with  $\mathbf{K}^{(t)}$
- 5: Determine step size,  $\mu^{(t)}$ , using Armijo's rule
- 6:  $\mathbf{W}^{(t+1)} \leftarrow \mathbf{W}^{(t)} - \mu^{(t)} \frac{\partial J(\mathbf{W})}{\partial \mathbf{W}}$
- 7: **until** convergence

$$\blacksquare f(\mathbf{x}) = \sum_{i=1}^N \alpha_i y_i k(\mathbf{W}^\top \mathbf{x}_i, \mathbf{W}^\top \mathbf{x}) + b$$

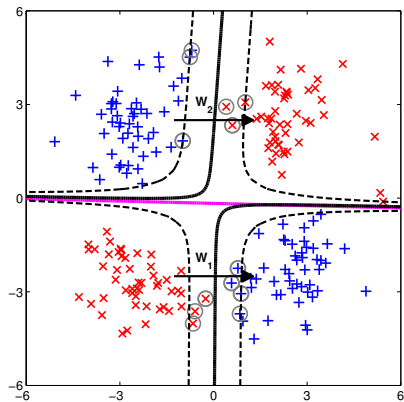
- Common kernel functions

$$k_L(\mathbf{z}_i, \mathbf{z}_j) = \langle \mathbf{z}_i, \mathbf{z}_j \rangle = \mathbf{x}_i^\top \mathbf{W} \mathbf{W}^\top \mathbf{x}_j$$

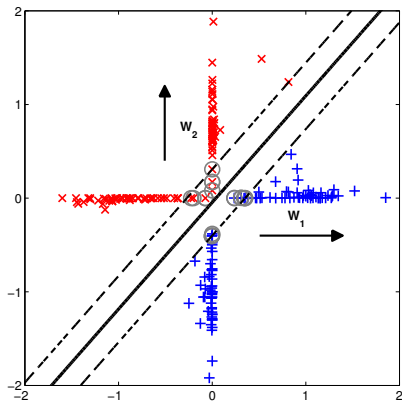
$$k_P(\mathbf{z}_i, \mathbf{z}_j) = (\langle \mathbf{z}_i, \mathbf{z}_j \rangle + 1)^q = (\mathbf{x}_i^\top \mathbf{W} \mathbf{W}^\top \mathbf{x}_j + 1)^q$$

$$\begin{aligned} k_G(\mathbf{z}_i, \mathbf{z}_j) &= \exp(-\|\mathbf{z}_i - \mathbf{z}_j\|_2^2 / s^2) \\ &= \exp(-(\mathbf{x}_i - \mathbf{x}_j)^\top \mathbf{W} \mathbf{W}^\top (\mathbf{x}_i - \mathbf{x}_j) / s^2) \end{aligned}$$

# Local Projection Kernels



Original Data



Projected Data

## Local Projection Kernels

■  $\mathbf{z}^m = \mathbf{W}_m^\top \mathbf{x}^m$     where  $\mathbf{W}_m \in \mathbb{R}^{D_m \times R_m}$

■  $f(\mathbf{x}) = \sum_{m=1}^P \eta_m(\mathbf{x}|\mathbf{V}) \langle \mathbf{w}_m, \Phi_m(\mathbf{W}_m^\top \mathbf{x}^m) \rangle + b$     where  $\eta_m: \mathbb{R}^{D_g} \rightarrow \mathbb{R}$

## Local Projection Kernels

■  $\mathbf{z}^m = \mathbf{W}_m^\top \mathbf{x}^m$  where  $\mathbf{W}_m \in \mathbb{R}^{D_m \times R_m}$

■  $f(\mathbf{x}) = \sum_{m=1}^P \eta_m(\mathbf{x}|\mathbf{V}) \langle \mathbf{w}_m, \Phi_m(\mathbf{W}_m^\top \mathbf{x}^m) \rangle + b$  where  $\eta_m: \mathbb{R}^{D_g} \rightarrow \mathbb{R}$

■ Primal optimization problem (Gönen and Alpaydm, 2010c)

$$\min. \frac{1}{2} \sum_{m=1}^P \|\mathbf{w}_m\|_2^2 + C \sum_{i=1}^N \xi_i$$

$$\text{w.r.t. } \mathbf{w}_m \in \mathbb{R}^{S_m}, \boldsymbol{\xi} \in \mathbb{R}_+^N, \mathbf{V} \in \mathbb{R}^{P \times (D_g+1)}, \mathbf{W}_m \in \mathbb{R}^{D_m \times R_m}, b \in \mathbb{R}$$

$$\text{s.t. } y_i \left( \sum_{m=1}^P \eta_m(\mathbf{x}_i|\mathbf{V}) \langle \mathbf{w}_m, \Phi_m(\mathbf{W}_m^\top \mathbf{x}_i^m) \rangle + b \right) \geq 1 - \xi_i \quad \forall i$$

## Local Projection Kernels

- $\mathbf{z}^m = \mathbf{W}_m^\top \mathbf{x}^m$  where  $\mathbf{W}_m \in \mathbb{R}^{D_m \times R_m}$

- $f(\mathbf{x}) = \sum_{m=1}^P \eta_m(\mathbf{x}|\mathbf{V}) \langle \mathbf{w}_m, \Phi_m(\mathbf{W}_m^\top \mathbf{x}^m) \rangle + b$  where  $\eta_m: \mathbb{R}^{D_g} \rightarrow \mathbb{R}$

- Primal optimization problem (Gönen and Alpaydmn, 2010c)

$$\min. \frac{1}{2} \sum_{m=1}^P \|\mathbf{w}_m\|_2^2 + C \sum_{i=1}^N \xi_i$$

$$\text{w.r.t. } \mathbf{w}_m \in \mathbb{R}^{S_m}, \boldsymbol{\xi} \in \mathbb{R}_+^N, \mathbf{V} \in \mathbb{R}^{P \times (D_g+1)}, \mathbf{W}_m \in \mathbb{R}^{D_m \times R_m}, b \in \mathbb{R}$$

$$\text{s.t. } y_i \left( \sum_{m=1}^P \eta_m(\mathbf{x}_i|\mathbf{V}) \langle \mathbf{w}_m, \Phi_m(\mathbf{W}_m^\top \mathbf{x}_i^m) \rangle + b \right) \geq 1 - \xi_i \quad \forall i$$

- For given  $\mathbf{V}$  and  $\{\mathbf{W}_m\}_{m=1}^P$ , problem becomes convex

## Local Projection Kernels

- Dual optimization problem for given  $\mathbf{V}$  and  $\{\mathbf{W}_m\}_{m=1}^P$

$$\max. J(\mathbf{V}, \{\mathbf{W}_m\}_{m=1}^P) = \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j k_{\eta}(\mathbf{x}_i, \mathbf{x}_j)$$

$$\text{w.r.t. } \boldsymbol{\alpha} \in \mathbb{R}_+^N$$

$$\text{s.t. } \sum_{i=1}^N \alpha_i y_i = 0$$

$$C \geq \alpha_i \geq 0 \quad \forall i$$

## Local Projection Kernels

- Dual optimization problem for given  $\mathbf{V}$  and  $\{\mathbf{W}_m\}_{m=1}^P$

$$\max. J(\mathbf{V}, \{\mathbf{W}_m\}_{m=1}^P) = \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j k_\eta(\mathbf{x}_i, \mathbf{x}_j)$$

$$\text{w.r.t. } \boldsymbol{\alpha} \in \mathbb{R}_+^N$$

$$\text{s.t. } \sum_{i=1}^N \alpha_i y_i = 0$$

$$C \geq \alpha_i \geq 0 \quad \forall i$$

- *local projection kernel function*

$$k_\eta(\mathbf{x}_i, \mathbf{x}_j) = \sum_{m=1}^P \eta_m(\mathbf{x}_i | \mathbf{V}) \underbrace{\langle \Phi_m(\mathbf{W}_m^\top \mathbf{x}_i^m), \Phi_m(\mathbf{W}_m^\top \mathbf{x}_j^m) \rangle}_{k_m(\mathbf{W}_m^\top \mathbf{x}_i^m, \mathbf{W}_m^\top \mathbf{x}_j^m)} \eta_m(\mathbf{x}_j | \mathbf{V})$$

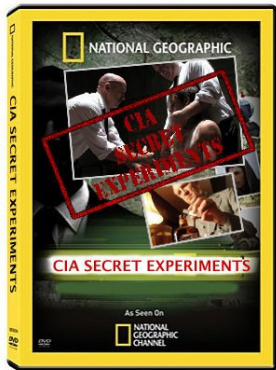
# Local Projection Kernels

- 1: Initialize  $\mathbf{V}^{(0)}$  and  $\{\mathbf{W}_m^{(0)}\}_{m=1}^P$  randomly
- 2: **repeat**
- 3: Calculate  $\mathbf{K}_\eta^{(t)} = \{k_\eta(\mathbf{x}_i, \mathbf{x}_j)\}_{i,j=1}^N$  using  $\mathbf{V}^{(t)}$  and  $\{\mathbf{W}_m^{(t)}\}_{m=1}^P$
- 4: Solve kernel machine with  $\mathbf{K}_\eta^{(t)}$
- 5: Determine step size,  $\mu^{(t)}$ , using Armijo's rule
- 6:  $\mathbf{W}_m^{(t+1)} \leftarrow \mathbf{W}_m^{(t)} - \mu^{(t)} \frac{\partial J(\mathbf{V}, \{\mathbf{W}_m\}_{m=1}^P)}{\partial \mathbf{W}_m} \quad \forall m$
- 7: Calculate  $\mathbf{K}_\eta^{(t)} = \{k_\eta(\mathbf{x}_i, \mathbf{x}_j)\}_{i,j=1}^N$  using  $\mathbf{V}^{(t)}$  and  $\{\mathbf{W}_m^{(t+1)}\}_{m=1}^P$
- 8: Solve kernel machine with  $\mathbf{K}_\eta^{(t)}$
- 9: Determine step size,  $\Delta^{(t)}$ , using Armijo's rule
- 10:  $\mathbf{V}^{(t+1)} \leftarrow \mathbf{V}^{(t)} - \Delta^{(t)} \frac{\partial J(\mathbf{V}, \{\mathbf{W}_m\}_{m=1}^P)}{\partial \mathbf{V}}$
- 11: **until** convergence

$$\blacksquare f(\mathbf{x}) = \sum_{i=1}^N \alpha_i y_i \underbrace{\sum_{m=1}^P \eta_m(\mathbf{x}_i | \mathbf{V}) k_m(\mathbf{W}_m^\top \mathbf{x}_i^m, \mathbf{W}_m^\top \mathbf{x}^m) \eta_m(\mathbf{x} | \mathbf{V})}_{k_\eta(\mathbf{x}_i, \mathbf{x})} + b$$

# Outline

1. Introduction
2. Multiple Kernel Learning
3. Localized Multiple Kernel Learning
4. Local Projection Kernels
- 5. Experiments**
6. Conclusions



## Intermediate Integration Experiments

- Protein stability change due to amino acid substitutions (Özen et al., 2009)
- We extract a data set from the ProTherm database (Gromiha et al., 2000)

Repr.	Original Feat.	New Repr.	New Feat.
SO	$\pm 3$ Neighbors ( $\pm 3NE$ ) Mutation (MUT) T/PH	SO*	PAM250 (PAM)
TO	Mutation (MUT) $C_\alpha$ Contacts (CA) SA/T/PH	TO*	PAM250 (PAM) $C_\alpha$ B-FACTORS (BFA) $C_\beta$ B-FACTORS (BFB) $C_\alpha$ and $C_\beta$ Contacts (CB)
ST	$\pm 3$ neighbors ( $\pm 3 NE$ ) Mutation (MUT) $C_\alpha$ Contacts (CA) SA/T/PH	ST*	PAM250 (PAM) $C_\alpha$ B-FACTORS (BFA) $C_\beta$ B-FACTORS (BFB) $C_\alpha$ and $C_\beta$ Contacts (CB)

# Intermediate Integration Experiments

## ■ Classification results

SO	SO*	TO	TO*	ST	ST*
0.800	0.799	0.805	0.802	0.806	0.804

## ■ Knowledge extraction results

SO	(0.19)1 <sub>NE</sub> + (0.20)2 <sub>NE</sub> + (0.23)3 <sub>NE</sub> + (0.27)M <sub>UT</sub> + (0.09)T + (0.03)P <sub>H</sub>
SO*	(0.19)1 <sub>NE</sub> + (0.20)2 <sub>NE</sub> + (0.22)3 <sub>NE</sub> + (0.27)M <sub>UT</sub> + (0.09)T + (0.03)P <sub>H</sub> + (0.00)P <sub>AM</sub>
TO	(0.19)M <sub>UT</sub> + (0.56)C <sub>A</sub> + (0.17)S <sub>A</sub> + (0.05)T + (0.02)P <sub>H</sub>
TO*	(0.21)M <sub>UT</sub> + (0.23)C <sub>A</sub> + (0.12)S <sub>A</sub> + (0.06)T + (0.02)P <sub>H</sub> + (0.00)P <sub>AM</sub> + (0.23)C <sub>B</sub> + (0.07)B <sub>FA</sub> + (0.06)B <sub>FB</sub>
ST	(0.04)1 <sub>NE</sub> + (0.03)2 <sub>NE</sub> + (0.04)3 <sub>NE</sub> + (0.21)M <sub>UT</sub> + (0.45)C <sub>A</sub> + (0.15)S <sub>A</sub> + (0.06)T + (0.02)P <sub>H</sub>
ST*	(0.02)1 <sub>NE</sub> + (0.02)2 <sub>NE</sub> + (0.03)3 <sub>NE</sub> + (0.21)M <sub>UT</sub> + (0.21)C <sub>A</sub> + (0.11)S <sub>A</sub> + (0.06)T + (0.02)P <sub>H</sub> + (0.00)P <sub>AM</sub> + (0.19)C <sub>B</sub> + (0.06)B <sub>FA</sub> + (0.06)B <sub>FB</sub>

# Combining General Purpose Kernels

- ACCEPTORS and DONORS  $\Rightarrow$  human splice site detection
- ARABIDOPSIS and VERTEBRATES  $\Rightarrow$  translation initiation site detection
- POLYADENYLATION  $\Rightarrow$  polyadenylation signal prediction

## Combining General Purpose Kernels

- ACCEPTORS and DONORS  $\Rightarrow$  human splice site detection
- ARABIDOPSIS and VERTEBRATES  $\Rightarrow$  translation initiation site detection
- POLYADENYLATION  $\Rightarrow$  polyadenylation signal prediction
- MKL versus RMKL with  $(k_L - k_P - k_G)$

	MKL		RMKL		$d_L - d_P - d_G$
	Test Acc.	SV	Test Acc.	SV	
	$\eta_L - \eta_P - \eta_G$		$\eta_L - \eta_P - \eta_G$		
ACCEPTORS	90.45 $\pm$ 0.42	50.12 $\pm$ 1.11	<b>91.19<math>\pm</math>0.38</b>	<b>44.90<math>\pm</math>0.82</b>	1.00-1.55-0.47
	0.30-0.70-0.00		0.00-0.00-4.59		
DONORS	94.77 $\pm$ 0.28	27.78 $\pm$ 0.47	94.78 $\pm$ 0.21	<b>26.25<math>\pm</math>0.45</b>	1.00-1.52-1.06
	0.11-0.03-0.86		0.17-0.00-0.73		
ARABIDOPSIS	85.29 $\pm$ 0.84	62.32 $\pm$ 1.10	<b>85.92<math>\pm</math>0.85</b>	62.39 $\pm$ 0.98	1.00-1.00-0.44
	0.25-0.75-0.00		0.00-0.00-5.16		
VERTEBRATES	86.22 $\pm$ 0.22	53.84 $\pm$ 0.72	86.15 $\pm$ 0.30	<b>40.86<math>\pm</math>0.44</b>	1.00-1.42-1.16
	0.20-0.80-0.00		0.70-0.15-0.00		
POLYADENYLATION	68.25 $\pm$ 1.24	72.14 $\pm$ 1.02	68.70 $\pm$ 1.34	72.81 $\pm$ 1.06	1.00-1.54-0.96
	0.01-0.16-0.84		0.00-0.00-1.09		

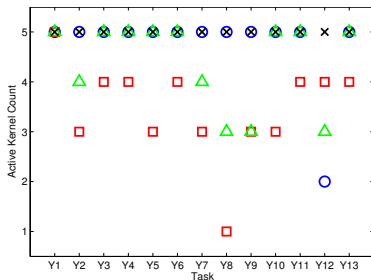
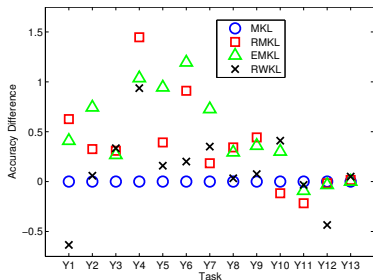
# Combining Domain Specific Kernels

- Protein function prediction
- Different kernel functions

Kernel	Explanation	Data Source
$k_{Pfam}$	Pfam kernel	Protein sequences
$k_{PfamE}$	Enriched Pfam kernel	Protein sequences
$k_{TAP}$	Diffusion kernel	Protein interactions
$k_{Phys}$	Diffusion kernel	Protein interactions
$k_{Gen}$	Diffusion kernel	Protein interactions
$k_{Exp}$	Correlation kernel	Gene expression profiles
$k_{ExpG}$	Gaussian kernel	Gene expression profiles
$k_{SW}$	Smith-Waterman kernel	Protein sequences

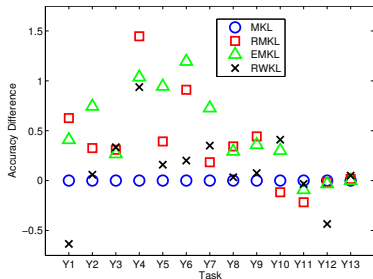
- We compare MKL, RMKL, EMKL, and RWKL

# Combining Domain Specific Kernels

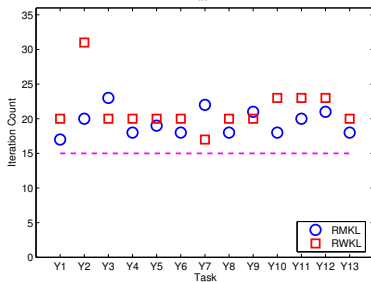
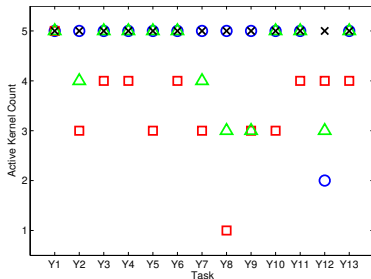


- RMKL produces smaller kernel ensembles without sacrificing accuracy

# Combining Domain Specific Kernels



- RMKL produces smaller kernel ensembles without sacrificing accuracy
- Grid search requires 81 ( $3^4$ ) iterations



## Representation Selection

$$\blacksquare \text{ cost} = \underbrace{\sum_{i=1}^N \mathbf{1}(\alpha_i > 0)}_{\text{the number of support vectors}} \underbrace{\sum_{m=1}^P \mathbf{1}(\eta_m > 0) \frac{d_m}{\sum_{h=1}^P d_h}}_{\text{the total normalized cost for active kernels}}$$

# Representation Selection

$$\blacksquare \text{ cost} = \underbrace{\sum_{i=1}^N \mathbf{1}(\alpha_i > 0)}_{\text{the number of support vectors}} \underbrace{\sum_{m=1}^P \mathbf{1}(\eta_m > 0) \frac{d_m}{\sum_{h=1}^P d_h}}_{\text{the total normalized cost for active kernels}}$$



$$d_{\text{DYN}} = 1.00$$



$$d_{\text{STA4}} = 1.00$$



$$d_{\text{STA8}} = 2.00$$



$$d_{\text{STA16}} = 4.00$$

# Representation Selection

$$\blacksquare \text{ cost} = \underbrace{\sum_{i=1}^N \mathbf{1}(\alpha_i > 0)}_{\text{the number of support vectors}} \underbrace{\sum_{m=1}^P \mathbf{1}(\eta_m > 0)}_{\text{the total normalized cost for active kernels}} \frac{d_m}{\sum_{h=1}^P d_h}$$



$d_{\text{DYN}} = 1.00$



$d_{\text{STA4}} = 1.00$



$d_{\text{STA8}} = 2.00$



$d_{\text{STA16}} = 4.00$

Data Set	Test Acc.	SV	Total Cost	$\eta_{\text{DYN}} - \eta_{\text{STA4}} - \eta_{\text{STA8}} - \eta_{\text{STA16}}$
PENDIGITS (0vs8)	99.50±0.11	10.70±0.11	10.70±0.33	0.24-0.16-0.33-0.26
	99.30±0.32	<b>5.80±0.42</b>	<b>2.90±0.21</b>	0.48-0.38-0.14-0.00
PENDIGITS (1vs7)	99.90±0.11	9.70±0.53	9.70±1.33	0.32-0.31-0.31-0.06
	99.70±0.32	<b>3.00±0.42</b>	<b>1.50±0.21</b>	0.45-0.51-0.04-0.00
PENDIGITS (5vs9)	99.30±0.11	12.90±0.95	12.90±0.04	0.35-0.12-0.10-0.44
	99.50±0.11	<b>10.30±0.11</b>	<b>2.58±0.03</b>	0.64-0.36-0.00-0.00

## Kernel/Representation Selection

- Protein location prediction with  $(k_{SW}-k_B-k_{Pfam}-k_{FFT}-k_{LI}-k_D-k_E)$

Kernel	Explanation	Data Source
$k_{SW}$	Smith-Waterman kernel	Protein sequences
$k_B$	BLAST kernel	Protein sequences
$k_{Pfam}$	Pfam kernel	Protein sequences
$k_{FFT}$	FFT kernel	Hydropathy profiles
$k_{LI}$	Linear kernel	Protein interactions
$k_D$	Diffusion kernel	Protein interactions
$k_E$	Gaussian kernel	Gene expression profiles

- $(d_{SW}-d_B-d_{Pfam}-d_{FFT}-d_{LI}-d_D-d_E) = (1.41-1.41-1.41-1.41-1.00-2.00-1.00)$

## Kernel/Representation Selection

- Protein location prediction with  $(k_{SW}-k_B-k_{Pfam}-k_{FFT}-k_{LI}-k_D-k_E)$

Kernel	Explanation	Data Source
$k_{SW}$	Smith-Waterman kernel	Protein sequences
$k_B$	BLAST kernel	Protein sequences
$k_{Pfam}$	Pfam kernel	Protein sequences
$k_{FFT}$	FFT kernel	Hydropathy profiles
$k_{LI}$	Linear kernel	Protein interactions
$k_D$	Diffusion kernel	Protein interactions
$k_E$	Gaussian kernel	Gene expression profiles

- $(d_{SW}-d_B-d_{Pfam}-d_{FFT}-d_{LI}-d_D-d_E) = (1.41-1.41-1.41-1.41-1.00-2.00-1.00)$

Task	Test Acc.	SV	Total Cost	$\eta_{SW}-\eta_B-\eta_{Pfam}-\eta_{FFT}-\eta_{LI}-\eta_D-\eta_E$
MEMBRANE	86.30±0.61	84.42±1.03	75.34±5.09	0.28-0.31-0.11-0.05-0.00-0.15-0.10
	85.40±0.55	<b>71.33±1.67</b>	<b>37.81±5.09</b>	0.00-0.26-0.06-0.00-0.07-0.00-0.60
RIBOSOMAL	98.99±0.26	18.56±1.56	6.86±1.96	0.01-0.00-0.00-0.16-0.03-0.00-0.80
	99.07±0.18	18.61±1.46	6.08±1.96	0.00-0.00-0.00-0.01-0.03-0.00-0.96

# Combining Multiple Feature Representations

- Multiple Features (MULTIFEAT)  
digit recognition data set
- Separate small ('0' - '4') digits from  
large ('5' - '9') digits

Name	Dim.	Data Source
FAC	216	Profile correlations
FOU	76	Fourier coefficients of the shapes
KAR	64	Karhunen-Loève coefficients
MOR	6	Morphological features
PIX	240	Pixel averages in $2 \times 3$ windows
ZER	47	Zernike moments

# Combining Multiple Feature Representations

- Multiple Features (MULTIFEAT) digit recognition data set
- Separate small ('0' - '4') digits from large ('5' - '9') digits

Name	Dim.	Data Source
FAC	216	Profile correlations
FOU	76	Fourier coefficients of the shapes
KAR	64	Karhunen-Loève coefficients
MOR	6	Morphological features
PIX	240	Pixel averages in $2 \times 3$ windows
ZER	47	Zernike moments

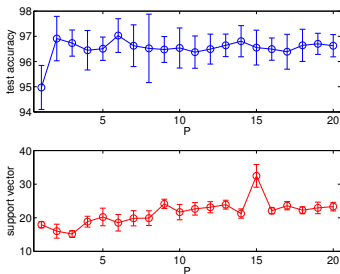
Method	Test Accuracy	Support Vector
SVM (FAC)	94.97±0.87	17.93±0.91
SVM (FOU)	90.54±1.11	28.90±1.69
SVM (KAR)	88.13±0.73	33.62±1.31
SVM (MOR)	69.61±0.14	61.90±0.49
SVM (PIX)	89.42±0.65	46.35±1.64
SVM (ZER)	89.12±0.63	26.27±1.67
SVM (ALL)	97.69±0.44	23.36±1.15
MKL	97.40±0.37	32.59±0.82
LMKL (softmax)	97.69±0.44	15.06±1.03
LMKL (sigmoid)	98.58±0.41	15.27±0.92

# Combining Multiple Feature Representations

- Multiple Features (MULTIFEAT) digit recognition data set
- Separate small ('0' - '4') digits from large ('5' - '9') digits

Name	Dim.	Data Source
FAC	216	Profile correlations
FOU	76	Fourier coefficients of the shapes
KAR	64	Karhunen-Loève coefficients
MOR	6	Morphological features
PIX	240	Pixel averages in $2 \times 3$ windows
ZER	47	Zernike moments

Method	Test Accuracy	Support Vector
SVM (FAC)	$94.97 \pm 0.87$	$17.93 \pm 0.91$
SVM (FOU)	$90.54 \pm 1.11$	$28.90 \pm 1.69$
SVM (KAR)	$88.13 \pm 0.73$	$33.62 \pm 1.31$
SVM (MOR)	$69.61 \pm 0.14$	$61.90 \pm 0.49$
SVM (PIX)	$89.42 \pm 0.65$	$46.35 \pm 1.64$
SVM (ZER)	$89.12 \pm 0.63$	$26.27 \pm 1.67$
SVM (ALL)	$97.69 \pm 0.44$	$23.36 \pm 1.15$
MKL	$97.40 \pm 0.37$	$32.59 \pm 0.82$
LMKL (softmax)	$97.69 \pm 0.44$	$15.06 \pm 1.03$
LMKL (sigmoid)	$98.58 \pm 0.41$	$15.27 \pm 0.92$



# Combining Multiple Feature Representations

- Internet Advertisements (ADVERT) data set
- Predict whether an image is an advertisement or not

Name	Dim.	Data Source
URL	457	Phrases in the URL
ORIGURL	495	Phrases in the URL of the image
ANCURL	472	Phrases in the anchor text
ALT	111	Phrases in the alternative text
CAPTION	19	Phrases in the caption terms

# Combining Multiple Feature Representations

- Internet Advertisements (ADVERT) data set
- Predict whether an image is an advertisement or not

Name	Dim.	Data Source
URL	457	Phrases in the URL
ORIGURL	495	Phrases in the URL of the image
ANCURL	472	Phrases in the anchor text
ALT	111	Phrases in the alternative text
CAPTION	19	Phrases in the caption terms

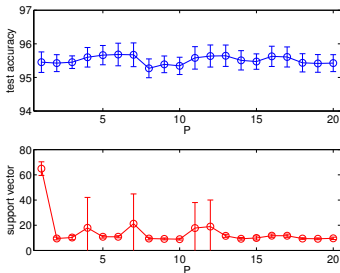
Method	Test Accuracy	Support Vector
SVM (URL)	94.67±0.24	83.32± 1.89
SVM (ORIGURL)	92.04±0.26	96.16± 0.51
SVM (ANCURL)	95.45±0.31	64.90± 5.41
SVM (ALT)	89.64±0.38	87.73± 1.17
SVM (CAPTION)	86.60±0.09	96.65± 0.42
SVM (ALL)	96.43±0.24	41.99± 1.76
MKL	96.32±0.50	35.82± 4.35
LMKL (softmax)	95.78±0.46	41.72±11.59
LMKL (sigmoid)	96.72±0.46	34.40± 1.51

# Combining Multiple Feature Representations

- Internet Advertisements (ADVERT) data set
- Predict whether an image is an advertisement or not

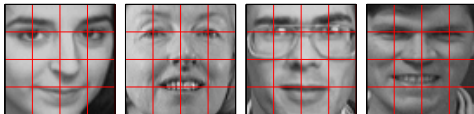
Name	Dim.	Data Source
URL	457	Phrases in the URL
ORIGURL	495	Phrases in the URL of the image
ANCURL	472	Phrases in the anchor text
ALT	111	Phrases in the alternative text
CAPTION	19	Phrases in the caption terms

Method	Test Accuracy	Support Vector
SVM (URL)	94.67±0.24	83.32± 1.89
SVM (ORIGURL)	92.04±0.26	96.16± 0.51
SVM (ANCURL)	95.45±0.31	64.90± 5.41
SVM (ALT)	89.64±0.38	87.73± 1.17
SVM (CAPTION)	86.60±0.09	96.65± 0.42
SVM (ALL)	96.43±0.24	41.99± 1.76
MKL	96.32±0.50	35.82± 4.35
LMKL (softmax)	95.78±0.46	41.72±11.59
LMKL (sigmoid)	96.72±0.46	34.40± 1.51



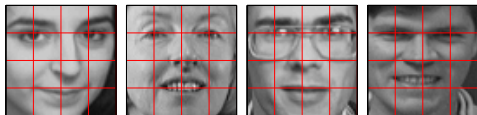
## Combining Multiple Input Patches

- Gender recognition experiments (Gönen and Alpaydın, 2009)



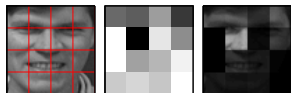
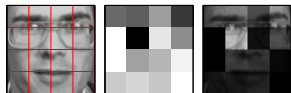
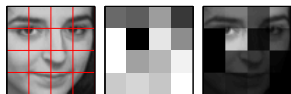
# Combining Multiple Input Patches

- Gender recognition experiments (Gönen and Alpaydın, 2009)



Method	Test Accuracy	Support Vector
SVM ( $\mathbf{x} = 4 \times 4$ )	93.28±0.65	21.70±0.93
SVM ( $\mathbf{x} = 8 \times 8$ )	97.50±1.16	20.13±1.04
SVM ( $\mathbf{x} = 16 \times 16$ )	97.03±0.93	19.91±1.01
SVM ( $\mathbf{x} = 32 \times 32$ )	97.97±1.48	23.71±1.39
SVM ( $\mathbf{x} = 64 \times 64$ )	97.66±1.41	25.94±1.01
<b>MKL</b>	<b>99.06±0.88</b>	<b>22.19±1.00</b>
LMKL (softmax and $\mathbf{x}^G = 4 \times 4$ )	97.19±2.81	16.65±3.34
LMKL (softmax and $\mathbf{x}^G = 8 \times 8$ )	97.19±2.48	22.54±4.56
<b>LMKL (softmax and <math>\mathbf{x}^G = 16 \times 16</math>)</b>	<b>99.22±1.33</b>	<b>16.38±1.50</b>
LMKL (sigmoid and $\mathbf{x}^G = 4 \times 4$ )	99.22±0.93	22.72±1.83
<b>LMKL (sigmoid and <math>\mathbf{x}^G = 8 \times 8</math>)</b>	<b>99.84±0.44</b>	<b>26.88±2.24</b>
LMKL (sigmoid and $\mathbf{x}^G = 16 \times 16$ )	99.38±1.34	21.65±1.44

## Combining Multiple Input Patches



(a)

(b)

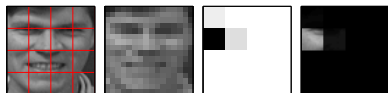
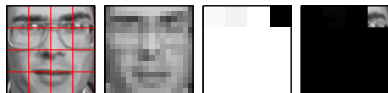
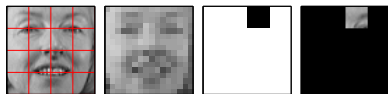
(c)

- (a)  $\Phi_m(\mathbf{x}^m)$ : features fed into kernels
- (b)  $\eta_m$ : MKL combination weights
- (c)  $\eta_m \Phi_m(\mathbf{x}^m)$ : features weighted with MKL combination weights

- MKL uses the same weights over the whole input space

# Combining Multiple Input Patches

## Softmax Gating



(a)

(b)

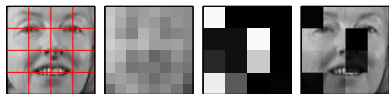
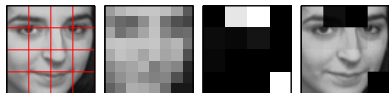
(c)

(d)

- (a)  $\Phi_m(\mathbf{x}^m)$ : features fed into kernels
  - (b)  $\mathbf{x}^g$ : features fed into softmax gating model
  - (c)  $\eta_m(\mathbf{x}|\mathbf{V})$ : softmax gating model outputs
  - (d)  $\eta_m(\mathbf{x}|\mathbf{V})\Phi_m(\mathbf{x}^m)$ : features weighted with softmax gating model outputs
- The softmax gating model generally activates a single patch
  - Eyes or eyebrows depending on the subject

# Combining Multiple Input Patches

## Sigmoid Gating



(a)

(b)

(c)

(d)

- (a)  $\Phi_m(\mathbf{x}^m)$ : features fed into kernels
- (b)  $\mathbf{x}^G$ : features fed into sigmoid gating model
- (c)  $\eta_m(\mathbf{x}|\mathbf{V})$ : sigmoid gating model outputs
- (d)  $\eta_m(\mathbf{x}|\mathbf{V})\Phi_m(\mathbf{x}^m)$ : features weighted with sigmoid gating model outputs

- Multiple patches are given nonzero weights
- Drives a high-resolution “eye” to regions of high saliency

# Combining Multiple Kernels

- Wine Quality (WHITEWINE) data set

Method	MSE	Support Vector
SVR ( $k_L$ )	0.59±0.00	66.83± 0.57
SVR ( $k_P$ and $q = 2$ )	0.54±0.01	66.22± 0.67
SVR ( $k_P$ and $q = 3$ )	0.54±0.00	66.14± 1.13
SVR ( $k_P$ and $q = 4$ )	0.52±0.01	66.55± 1.03
SVR ( $k_P$ and $q = 5$ )	0.52±0.01	66.27± 1.24
LMKL (softmax)	0.52±0.01	18.66±13.41
LMKL (sigmoid)	0.51±0.00	38.29± 2.34

# Combining Multiple Kernels

- Wine Quality (WHITEWINE) data set

- Concrete Compressive Strength (CONCRETE) data set

Method	MSE	Support Vector
SVR ( $k_L$ )	0.59±0.00	66.83± 0.57
SVR ( $k_P$ and $q = 2$ )	0.54±0.01	66.22± 0.67
SVR ( $k_P$ and $q = 3$ )	0.54±0.00	66.14± 1.13
SVR ( $k_P$ and $q = 4$ )	0.52±0.01	66.55± 1.03
SVR ( $k_P$ and $q = 5$ )	0.52±0.01	66.27± 1.24
LMKL (softmax)	0.52±0.01	18.66±13.41
LMKL (sigmoid)	0.51±0.00	38.29± 2.34

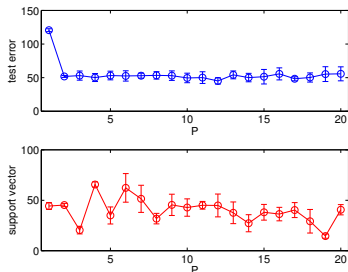
Method	MSE	Support Vector
SVR ( $k_L$ )	120.61±2.15	44.31±3.46
SVR ( $k_P$ and $q = 2$ )	92.57±4.19	36.22±1.24
SVR ( $k_P$ and $q = 3$ )	58.32±3.66	73.16±1.21
SVR ( $k_P$ and $q = 4$ )	63.83±9.58	52.52±2.40
SVR ( $k_P$ and $q = 5$ )	61.26±5.31	52.17±2.23
LMKL (softmax)	44.80±6.33	64.28±4.02
LMKL (sigmoid)	48.18±5.22	49.20±2.05

# Combining Multiple Kernels

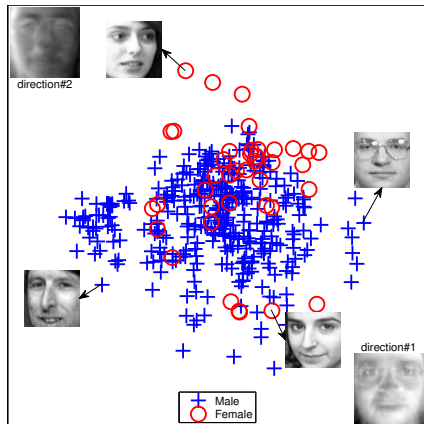
- Wine Quality (WHITEWINE) data set
- Concrete Compressive Strength (CONCRETE) data set

Method	MSE	Support Vector
SVR ( $k_L$ )	0.59±0.00	66.83± 0.57
SVR ( $k_P$ and $q = 2$ )	0.54±0.01	66.22± 0.67
SVR ( $k_P$ and $q = 3$ )	0.54±0.00	66.14± 1.13
SVR ( $k_P$ and $q = 4$ )	0.52±0.01	66.55± 1.03
SVR ( $k_P$ and $q = 5$ )	0.52±0.01	66.27± 1.24
LMKL (softmax)	0.52±0.01	18.66±13.41
LMKL (sigmoid)	0.51±0.00	38.29± 2.34

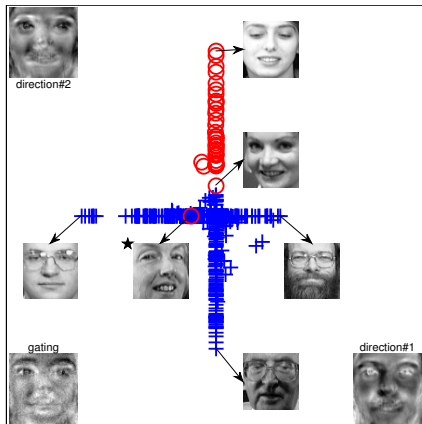
Method	MSE	Support Vector
SVR ( $k_L$ )	120.61±2.15	44.31±3.46
SVR ( $k_P$ and $q = 2$ )	92.57±4.19	36.22±1.24
SVR ( $k_P$ and $q = 3$ )	58.32±3.66	73.16±1.21
SVR ( $k_P$ and $q = 4$ )	63.83±9.58	52.52±2.40
SVR ( $k_P$ and $q = 5$ )	61.26±5.31	52.17±2.23
LMKL (softmax)	44.80±6.33	64.28±4.02
LMKL (sigmoid)	48.18±5.22	49.20±2.05



# Visualization Experiments



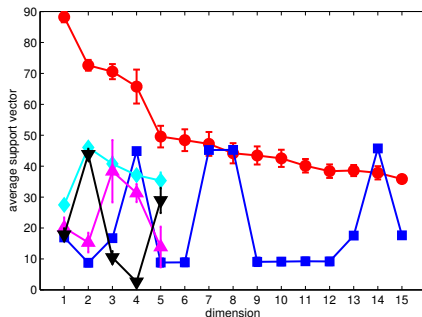
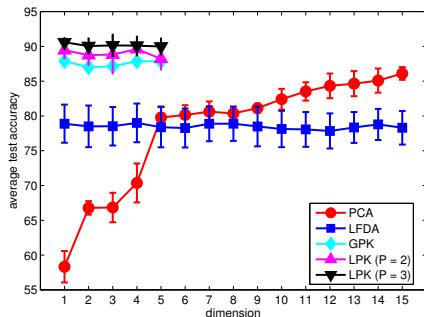
PCA



LPK

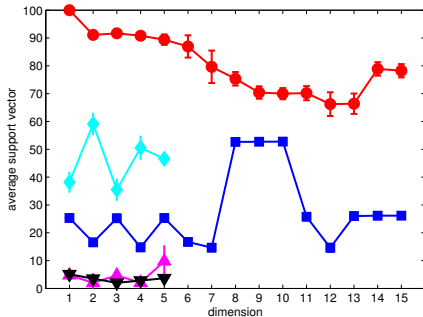
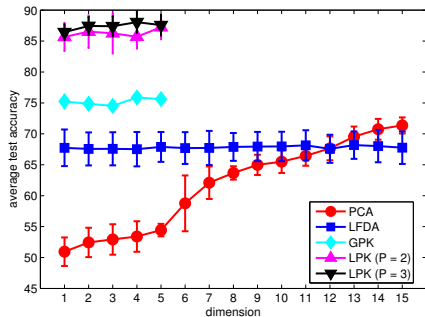
# Classification Experiments

- We compare PCA, LFDA, GPK, and LPK
- USPS digit recognition data set
- Separate even digits from odd digits



# Classification Experiments

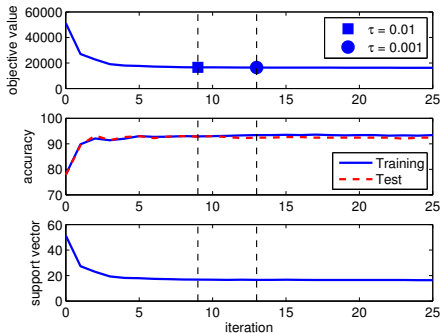
- USPS digit recognition data set
- Separate small ('0' - '4') digits from large ('5' - '9') digits



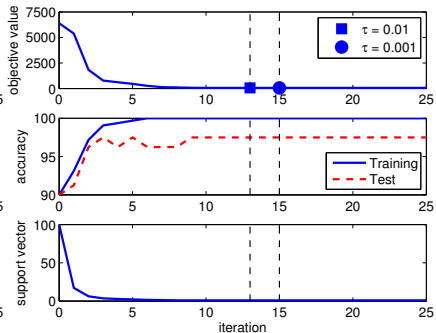
# Convergence Analysis

■ Convergence parameter  $\tau$

■ stop if  $(\text{objective value})^{(t)} > (1 - \tau)(\text{objective value})^{(t-1)}$



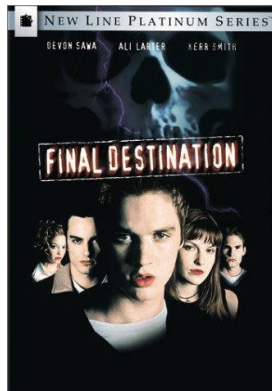
WAVEFORM



OLIVETTI

# Outline

1. Introduction
2. Multiple Kernel Learning
3. Localized Multiple Kernel Learning
4. Local Projection Kernels
5. Experiments
6. Conclusions



# Conclusions

- Regularized Multiple Kernel Learning
  - ▶ optimizes regularization parameters using RSM
  - ▶ obtains smoother decision functions
  - ▶ uses fewer kernels and support vectors
  - ▶ has advantage of RSM over grid search

# Conclusions

- Regularized Multiple Kernel Learning
  - ▶ optimizes regularization parameters using RSM
  - ▶ obtains smoother decision functions
  - ▶ uses fewer kernels and support vectors
  - ▶ has advantage of RSM over grid search
- Cost-Conscious Multiple Kernel Learning
  - ▶ includes the costs of kernels
  - ▶ eliminates costlier data representations

# Conclusions

- Localized Multiple Kernel Learning
  - ▶ input-dependent selection among multiple kernels or feature representations
  - ▶ two main ingredients
    - gating model for data-dependent kernel weights
    - kernel-based learning algorithm
  - ▶ two-step alternating optimization algorithm
  - ▶ applied to different gating models and learning problems

# Conclusions

- Localized Multiple Kernel Learning
  - ▶ input-dependent selection among multiple kernels or feature representations
  - ▶ two main ingredients
    - gating model for data-dependent kernel weights
    - kernel-based learning algorithm
  - ▶ two-step alternating optimization algorithm
  - ▶ applied to different gating models and learning problems
  - ▶ obtains better learners compared to MKL
  - ▶ improves accuracy and/or stores fewer support vectors
  - ▶ can combine multiple copies of the same kernel
  - ▶ acts as a saliency detector on image recognition problems

# Conclusions

## ■ Localized Multiple Kernel Learning

### ▶ computational complexity

- complexity of canonical solver & number of iterations

- Armijo's rule

- improvement in test time

$$f(\mathbf{x}) = \sum_{i=1}^N \alpha_i y_i \underbrace{\sum_{m=1}^P \eta_m(\mathbf{x}_i | \mathbf{V}) k_m(\mathbf{x}_i^m, \mathbf{x}^m) \eta_m(\mathbf{x} | \mathbf{V})}_{k_\eta(\mathbf{x}_i, \mathbf{x})} + b$$

# Conclusions

## ■ Localized Multiple Kernel Learning

- ▶ computational complexity
  - complexity of canonical solver & number of iterations
- Armijo's rule
- improvement in test time

$$f(\mathbf{x}) = \sum_{i=1}^N \alpha_i y_i \underbrace{\sum_{m=1}^P \eta_m(\mathbf{x}_i | \mathbf{V}) k_m(\mathbf{x}_i^m, \mathbf{x}^m) \eta_m(\mathbf{x} | \mathbf{V})}_{k_\eta(\mathbf{x}_i, \mathbf{x})} + b$$

- ▶ knowledge extraction
  - data-dependent relative importances

# Conclusions

## ■ Localized Multiple Kernel Learning

- ▶ computational complexity
  - complexity of canonical solver & number of iterations
- Armijo's rule
- improvement in test time

$$f(\mathbf{x}) = \sum_{i=1}^N \alpha_i y_i \underbrace{\sum_{m=1}^P \eta_m(\mathbf{x}_i | \mathbf{V}) k_m(\mathbf{x}_i^m, \mathbf{x}^m) \eta_m(\mathbf{x} | \mathbf{V})}_{k_\eta(\mathbf{x}_i, \mathbf{x})} + b$$

- ▶ knowledge extraction
  - data-dependent relative importances
- ▶ regularization
  - regularization on locally combined kernel
  - no overfitting with increasing number of kernels

# Conclusions

- Local Projection Kernels
  - ▶ localized dimensionality reduction on multimodal data sets
  - ▶ improves accuracy and/or stores fewer support vectors
  - ▶ computational complexity

# Conclusions

- Local Projection Kernels
  - ▶ localized dimensionality reduction on multimodal data sets
  - ▶ improves accuracy and/or stores fewer support vectors
  - ▶ computational complexity
- Future work
  - ▶ using RSM on validation performance for model selection problems
  - ▶ Bayesian reformulations of LMKL and LPK

# Conclusions

- Local Projection Kernels
  - ▶ localized dimensionality reduction on multimodal data sets
  - ▶ improves accuracy and/or stores fewer support vectors
  - ▶ computational complexity
- Future work
  - ▶ using RSM on validation performance for model selection problems
  - ▶ Bayesian reformulations of LMKL and LPK
- MATLAB source codes

<http://www.cmpe.boun.edu.tr/~gonen>

# References I

- Bach, F. R., Lanckriet, G. R. G. and Jordan, M. I. (2004). Multiple kernel learning, conic duality, and the SMO algorithm, *Proceedings of the 21st International Conference on Machine Learning*.
- Bach, F. R., Thibaux, R. and Jordan, M. I. (2005). Computing regularization paths for learning multiple kernels, *Advances in Neural Information Processing Systems 17*.
- Bredensteiner, E. J. and Bennett, K. P. (1999). Multicategory classification by support vector machines, *Computational Optimization and Applications* **12**(1-3): 53–79.
- Chapelle, O., Vapnik, V., Bousquet, O. and Mukherjee, S. (2002). Choosing multiple parameters for support vector machines, *Machine Learning* **46**(1–3): 131–159.
- Gönen, M. and Alpaydın, E. (2008). Localized multiple kernel learning, *Proceedings of the 25th International Conference on Machine Learning*.
- Gönen, M. and Alpaydın, E. (2009). Localized multiple kernel learning for image recognition, *NIPS Workshop on Understanding Multiple Kernel Learning Methods*.
- Gönen, M. and Alpaydın, E. (2010a). Cost-conscious multiple kernel learning, *Pattern Recognition Letters* **31**(9): 959–965.
- Gönen, M. and Alpaydın, E. (2010b). Localized multiple kernel regression, *Proceedings of the 20th International Conference on Pattern Recognition*.
- Gönen, M. and Alpaydın, E. (2010c). Supervised learning of local projection kernels, *Neurocomputing* **73**(10–12): 1694–1703.
- Gönen, M., Tanuğur, A. G. and Alpaydın, E. (2008). Multiclass posterior probability support vector machines, *IEEE Transactions on Neural Networks* **19**(1): 130–139.
- Gromiha, M. M., An, J., Kono, H., Oobatake, M., Uedaira, H., Prabakaran, P. and Sarai, A. (2000). Protherm, version 2.0: Thermodynamic database for proteins and mutants, *Nucleic Acids Research* **28**: 283–285.

## References II

- Jacobs, R. A., Jordan, M. I., Nowlan, S. J. and Hinton, G. E. (1991). Adaptive mixtures of local experts, *Neural Computation* **3**: 79–87.
- Lanckriet, G. R. G., Cristianini, N., Bartlett, P., Ghaoui, L. E. and Jordan, M. I. (2004). Learning the kernel matrix with semidefinite programming, *Journal of Machine Learning Research* **5**: 27–72.
- Noble, W. S. (2004). Support vector machine applications in computational biology, in B. Schölkopf, K. Tsuda and J.-P. Vert (eds), *Kernel Methods in Computational Biology*, The MIT Press, chapter 3.
- Özen, A., Gönen, M., Alpaydın, E. and Haliloğlu, T. (2009). Machine learning integration for predicting the effect of single amino acid substitutions on protein stability, *BMC Structural Biology* **9**: 66.
- Pavlidis, P., Weston, J., Cai, J. and Grundy, W. N. (2001). Gene functional classification from heterogeneous data, *Proceedings of the 5th Annual International Conference on Computational Molecular Biology*.
- Schölkopf, B. and Smola, A. J. (2002). *Learning with Kernels: Support Vector Machines, Regularization, Optimization, and Beyond*, The MIT Press, Cambridge, MA.
- Tao, Q., Wu, G.-W., Wang, F.-Y. and Wang, J. (2005). Posterior probability support vector machines for unbalanced data, *IEEE Transactions on Neural Networks* **16**(6): 1561–1573.
- Vapnik, V. (1998). *The Nature of Statistical Learning Theory*, John Wiley & Sons.
- Weston, J. and Watkins, C. (1998). Multi-class support vector machines, *Technical Report CSD-TR-98-04*, Royal Holloway, University of London, Department of Computer Science.