

# Efficient Multiple Kernel Learning

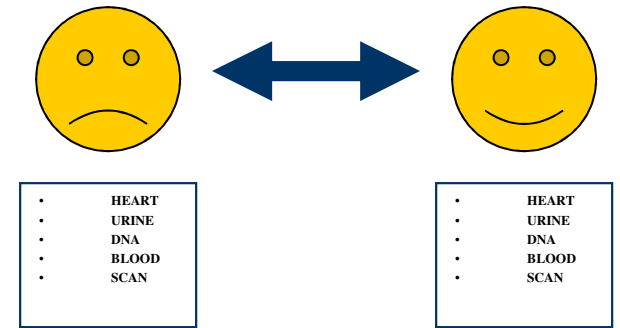
Lei Tang

# Outline

- What is Kernel Learning?
- What's the problem with existing formulation?
- Two new formulations for large scale kernel selection
  - SIL formulation (Cutting Planes)
  - More efficient MKL (Steepest Decent)

# Linear algorithm: binary classification

- **Data:**  $\{(x_i, y_i)\}_{i=1\dots n}$ 
  - $x \in \mathbb{R}^d$  = feature vector
  - $y \in \{-1, +1\}$  = label



- **Question:** design a classification rule

$$y = f(x)$$

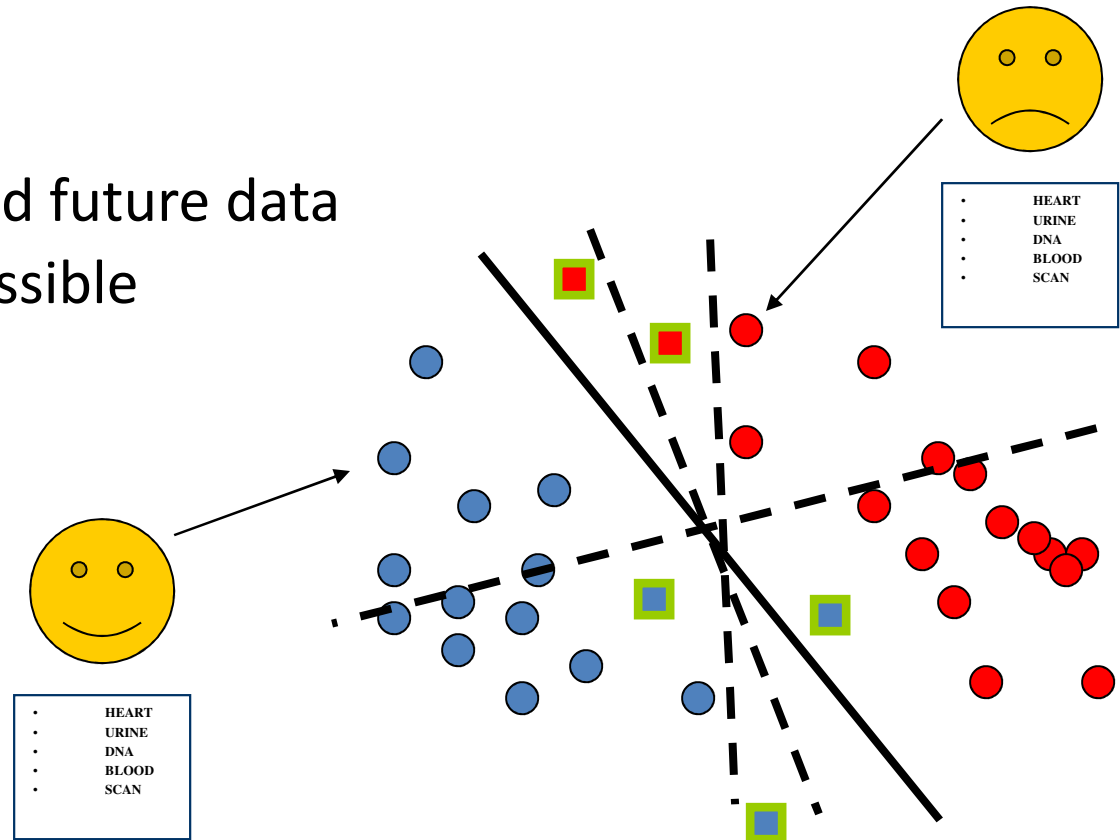
such that, given a new  $x$ , this predicts  $y$  with minimal probability of error

# Linear algorithm: binary classification

- Find **good hyperplane**

$(w, b) \in \mathbb{R}^{d+1}$

that classifies this and future data points as good as possible



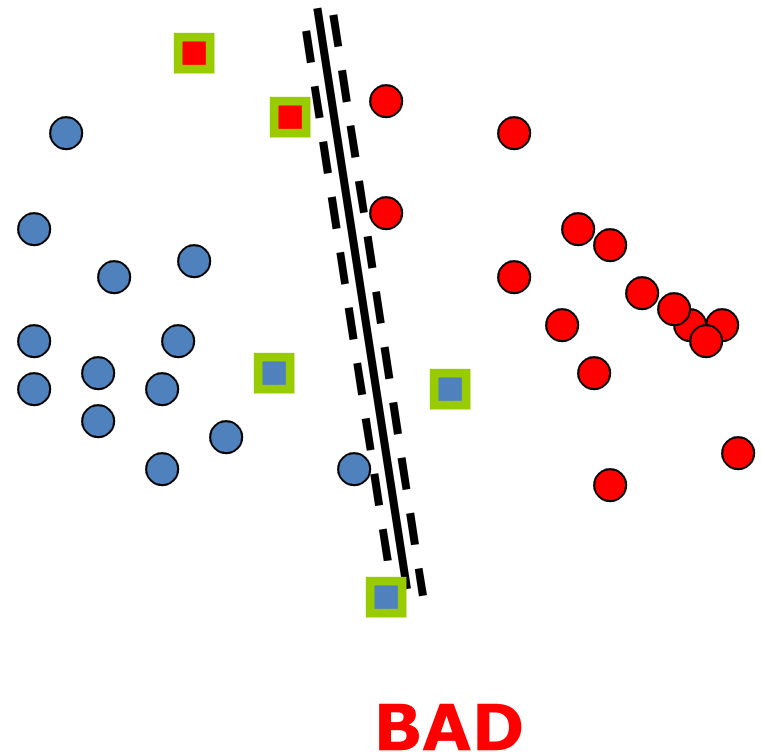
$$w^T x + b = 0$$

**Classification Rule:**

$$y = f(x) = \text{sign}(w^T x + b)$$

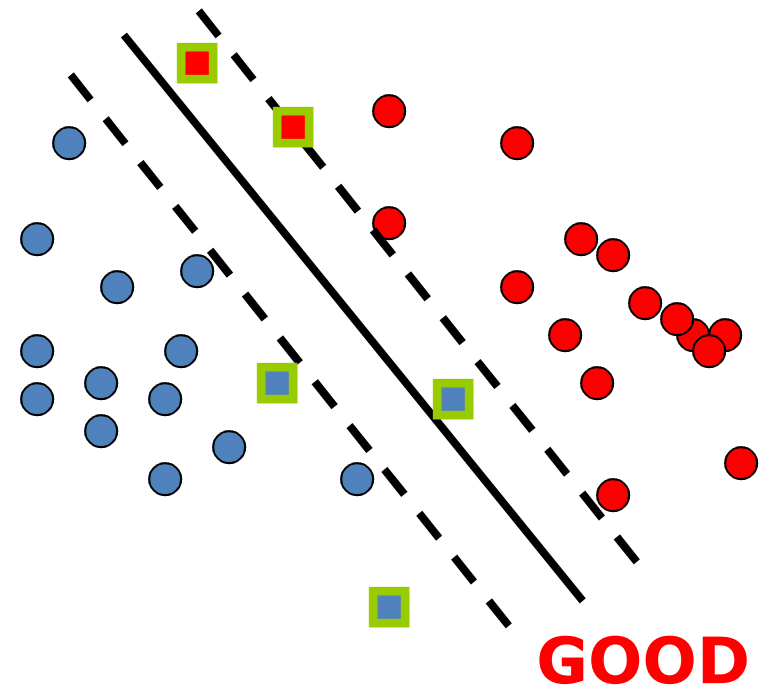
# Linear algorithm: binary classification

- **Intuition** (Vapnik, 1965) if linearly separable:
  - Separate the data
  - Place hyperplane “far” from the data: **large margin**



# Linear algorithm: binary classification

- **Intuition** (Vapnik, 1965) if linearly separable:
  - Separate the data
  - Place hyperplane “far” from the data: **large margin**

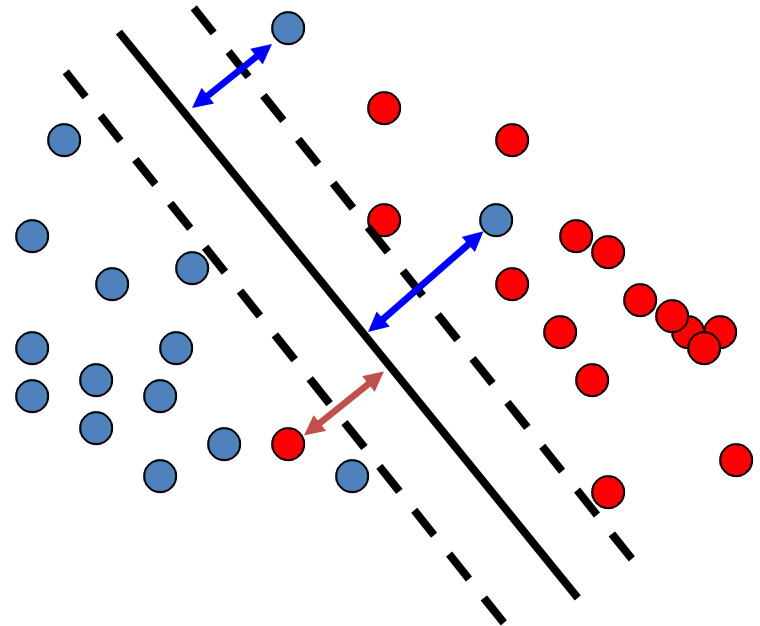


→ **Maximal Margin Classifier**

# Linear algorithm: binary classification

If **not linearly separable**:

- **Allow** some **errors**
- Still, try to place hyperplane “far” from each class



# SVM: Primal & Dual

$$f = w \cdot x + b$$

Primal:

$$\begin{aligned} \min_{w, b} \quad & \frac{1}{2} \|w\|_2^2 + C \sum_{i=1}^N \xi_i \\ \text{subject to} \quad & y_i (w \cdot x_i + b) \geq 1 - \xi_i \quad \forall i \\ & \xi_i \geq 0 \end{aligned}$$

Dual:

$$\begin{aligned} \max_{\alpha} \quad & \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j x_i^T x_j \\ \text{subject to} \quad & \sum_i \alpha_i y_i = 0, \quad C \geq \alpha_i \geq 0 \end{aligned}$$



# Linear algorithm: binary classification

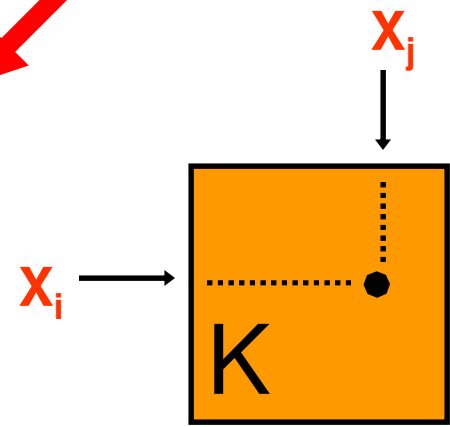
- **Training = convex optimization problem (QP):**

$$\begin{aligned} \max_{\alpha} \quad & \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \boxed{x_i^T x_j} \\ \text{subject to} \quad & \sum_i \alpha_i y_i = 0, \quad \alpha_i \geq 0 \end{aligned}$$

implicit  
embedding

$$K_{ij} = k(x_i, x_j)$$

$$\begin{aligned} \max_{\alpha} \quad & e^T \alpha - \frac{1}{2} \alpha^T D_y \boxed{K} D_y \alpha \\ \text{subject to} \quad & y^T \alpha = 0, \quad \alpha \geq 0 \end{aligned}$$



# Kernel algorithm: Support Vector Machine (SVM)

- **Training = convex optimization problem (QP):**

$$\max_{\alpha} \quad e^T \alpha - \frac{1}{2} \alpha^T D_y K D_y \alpha$$

subject to  $y^T \alpha = 0, \quad C \geq \alpha \geq 0$

- **Classification rule: classify new data point  $x$ :**

$$f(x) = \text{sign}(w^T x + b) = \text{sign}\left(\sum_{i=1}^{n_{SV}} \alpha_i y_i \boxed{x_i^T x} + b\right)$$

**Kernel algorithm !**

# Support Vector Machines (SVM)

- Hand-writing recognition (e.g., USPS)
- Computational biology (e.g., micro-array data)
- Text classification
- Face detection
- Face expression recognition
- Time series prediction (regression)
- Drug discovery (novelty detection)

# Different Kernels

- Various kinds of Kernel

- Linear kernel

$$K(X, Y) = \frac{\sum_i X_i Y_i}{\sqrt{\sum_i X_i X_i} \sqrt{\sum_i Y_i Y_i}}$$

- Gaussian kernel

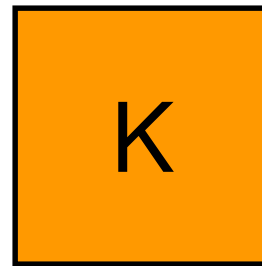
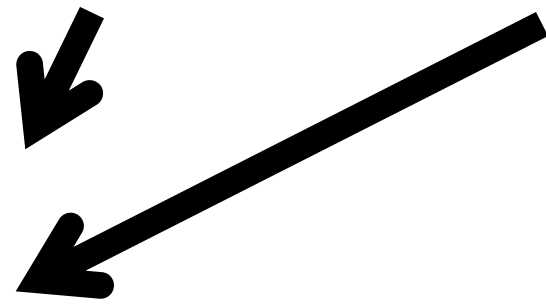
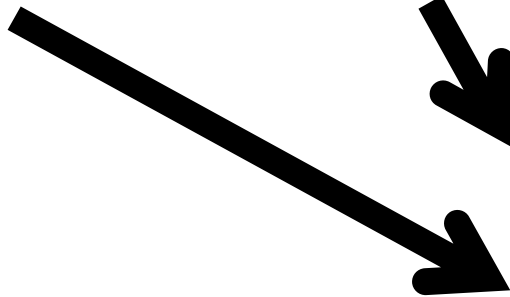
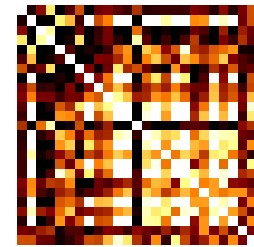
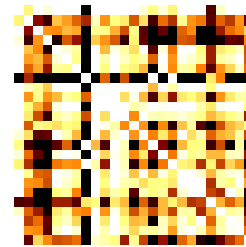
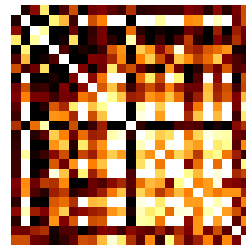
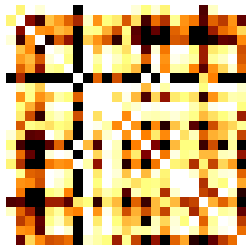
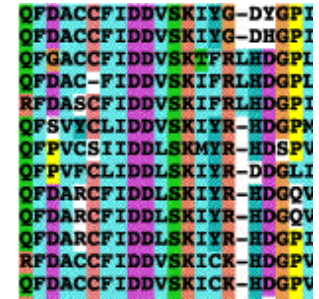
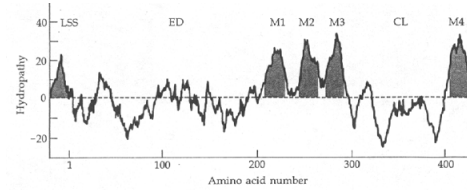
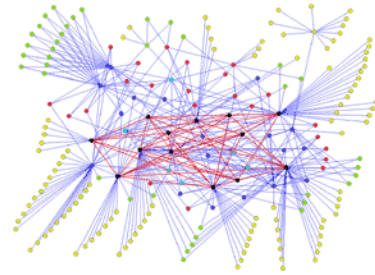
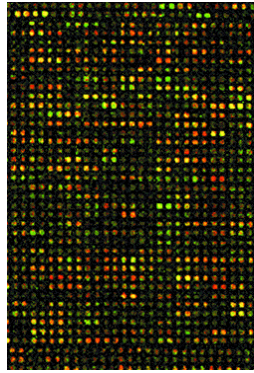
$$K(X, Y) = \exp\left(\frac{-\|X - Y\|^2}{2\sigma^2}\right)$$

- Diffusion kernel

- String Kernel

- .....

# Learning with Multiple Kernels



# Learning the optimal Kernel

## Overview of SVM with single kernel :

Binary classification of data  $x_i \in \mathbb{R}^d$ , labels  $y_i \in \{-1, 1\}$  ( $i = 1, \dots, n$ ) using  $f(x) = w^\top \phi(x) + b$ .

Primal problem

$$\begin{aligned} \min_{w,b} \quad & \frac{1}{2} \|w\|_2^2 + C \sum_i \xi_i \\ \text{s.t.} \quad & y_i (w^\top \phi(x_i) + b) \geq 1 - \xi_i, \quad i = 1, \dots, n \\ & \xi_i \geq 0, \quad i = 1, \dots, n \end{aligned}$$

Dual problem

$$\max_{\alpha} \quad \alpha^\top \mathbf{1} - \frac{1}{2} \alpha^\top K \alpha \quad \text{s.t.} \quad \alpha^\top y = 0, \quad \mathbf{0} \leq \alpha \leq C$$

G(K)



# Learning the optimal Kernel

$$K = \sum_i \mu_i K_i$$

**Learn a linear mix**

Upper bound:  
the smaller, the  
better the guaranteed  
performance

SVM, one kernel, dual formulation

$$\max_{\alpha, \alpha^\top y=0} \alpha^\top 1 - \frac{1}{2} \alpha^\top K \alpha \quad \text{s.t.} \quad 0 \leq \alpha \leq C$$

G(K)

SVM, multiple kernels, dual formulation

$$\min_{\eta, \eta^\top e=1} \left( \max_{\alpha, \alpha^\top y=0} \alpha^\top 1 - \frac{1}{2} \alpha^\top \left( \sum_j \eta_j K_j \right) \alpha \quad \text{s.t.} \quad 0 \leq \alpha \leq C \right)$$

To be Continued