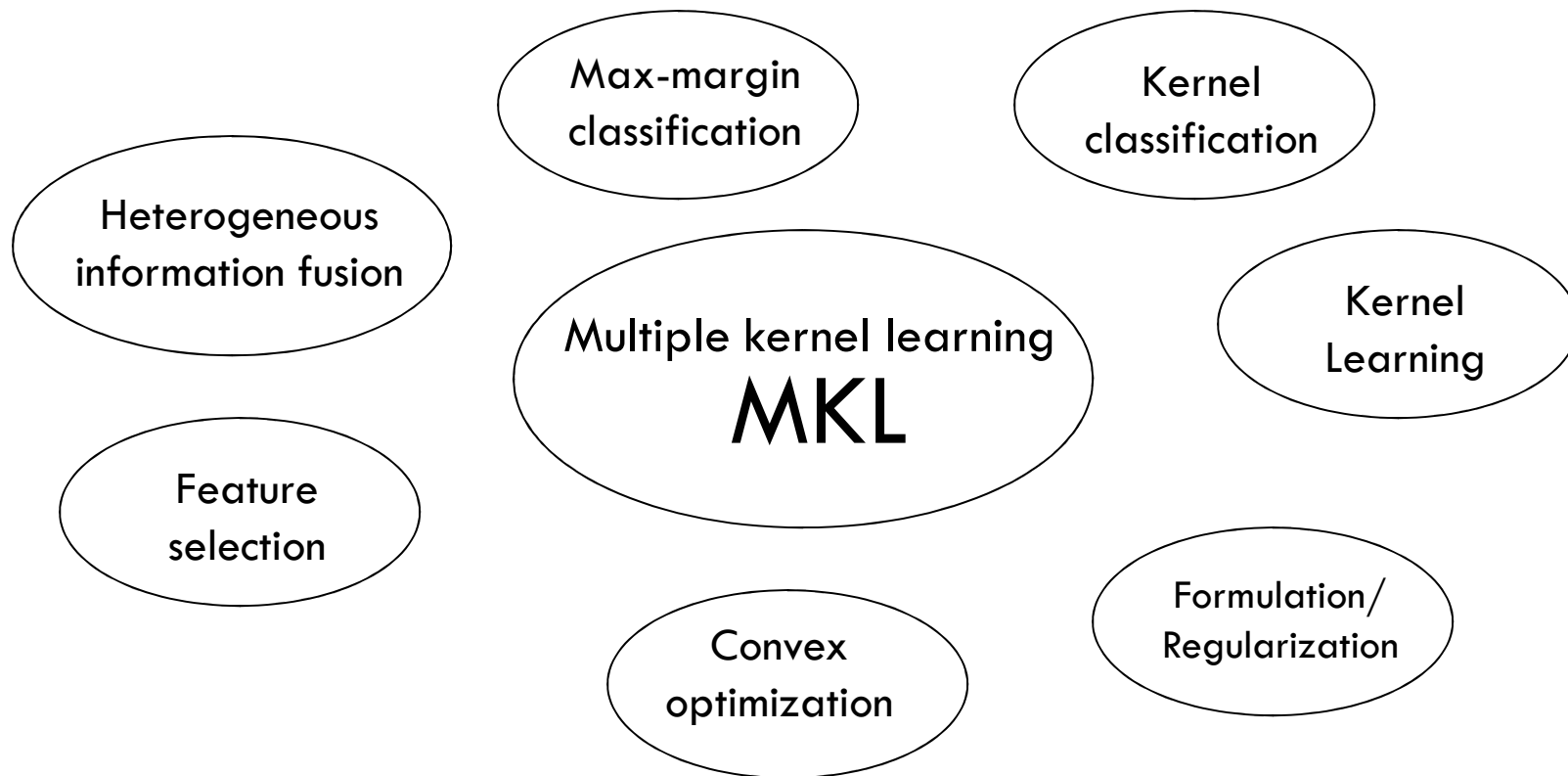


# MULTIPLE KERNEL LEARNING

CSE902

# Multiple Kernel Learning - keywords



MKL is used when there are heterogeneous sources (representations) of data for the task at hand (we consider classification)

# Outline

1. **Kernel Combination**
  1. *Heterogeneous Information Fusion*
  2. *Feature Selection*
2. Linear margin classifiers (SVM)
3. Kernel Classifiers
4. Kernel Learning / Multiple Kernel Learning
  1. MKL Formulation
  2. Sparse vs. Smooth MKL
  3. MKL Optimization
5. Experimental Results from Literature
6. Our Experiments
7. Conclusions

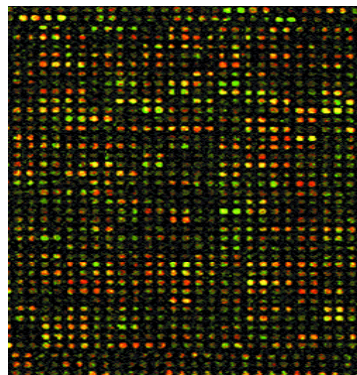
# Heterogeneous Information Fusion

## Web page categorization

- Data point = web page
- Sources of information about the webpage:
  - ▣ Content:
    - Text
    - Images
    - Structure
    - Sound
  - ▣ Relation to other web pages: links → network
  - ▣ Users (log data):
    - click behavior
    - origin

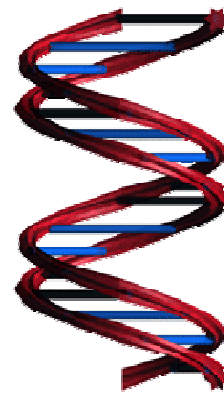
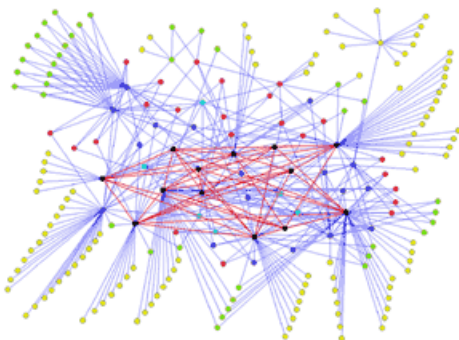
# Heterogeneous Information Fusion

## □ Bioinformatics

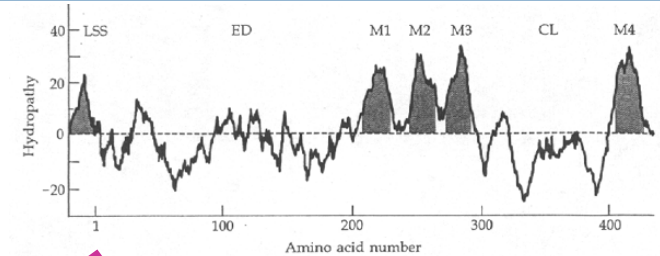


mRNA  
expression data

protein-protein  
interaction data



upstream region data  
(TF binding sites)



hydropobicity data

sequence data  
(gene, protein)

```
QFDACCFIDDVSKIYG-DYGPI
QFDACCFIDDVSKIYG-DHGPI
QFGACCFIDDVSKTFRLEDGPI
QFDAC-FIDDVSKIYRLEDGPI
RFDASCFIDDVSKIYRLEDGPI
QFSVYCLIDDVSKIYR-HDGPV
QFPVCSIIDDLISKIYR-HDSPV
QFPVFCLIDDLISKIYR-DDGLI
QFDARCFIDDLISKIYR-HDGPV
QFDARCFIDDLISKIYR-HDGPV
QFDARCFIDDLISKIYR-HDGPV
RFDACCFIDDVSKICK-HDGPV
QFDACCFIDDVSKICK-HDGPV
```

Lanckriet et al.

# Feature Selection



## Text mining

- Newsgroups data
  - ▣ Unstructured text
  - ▣ 1.3M Bag-of-Words features
- Log data for suspicious links
  - ▣ Structured text
  - ▣ 3.2M constant and variable features
- Real Time processing
- Using small number of features for prediction speed

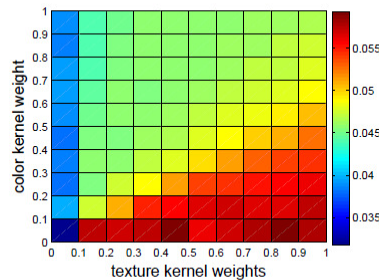
# Feature Selection

## Image categorization and retrieval

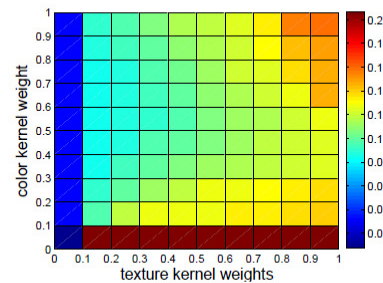


- Hundreds of feature types (SIFT, HOG, GIST)
- Select and combine features for improved prediction accuracy and speed

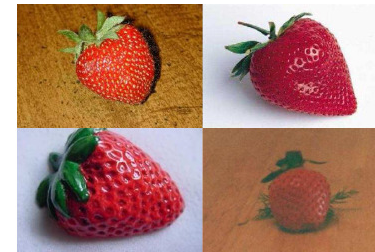
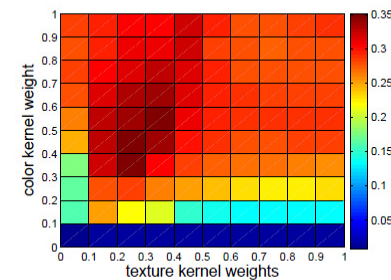
# Image Categorization: A Small Experiment



crocodile



snoopy



strawberry

- First row: Surface graphs (influence of different feature combination weights on MAP)
- Second row: 4 examples from each class



# Questions



- How to merge different information sources (images, sounds, click data, user feedback ...)
  - ▣ Same format (numerical), different scale
  - ▣ Different formats (ordinal, categorical, non-vectorial)
- How to eliminate irrelevant or redundant data
- How to combine the feature vectors (feature weighting)
- How to minimize the number of sources to improve computational efficiency (sparsity)

# MKL vs Alternatives



- MKL is proposed as an alternative to:
  - Cross validation
  - Feature Selection
  - Metric learning
  - Ensemble methods

# MKL vs Alternatives



- MKL is proposed as an alternative to:
  - ▣ Cross validation
  - ▣ Feature Selection
  - ▣ Metric learning
  - ▣ Ensemble methods
- Advantages of MKL
  - ▣ A technically sound way of combining features
  - ▣ Feature combination and classifier training is done simultaneously
  - ▣ Good learning bounds
  - ▣ Different data formats can be used in the same formulation
  - ▣ Non-linear learning

# Outline



1. Kernel Combination
  1. Heterogeneous Information Fusion
  2. Feature Selection
2. *Linear margin classifiers (SVM)*
3. Kernel Classifiers
4. Kernel Learning / Multiple Kernel Learning
  1. MKL Formulation
  2. Sparse vs. Smooth MKL
  3. MKL Optimization
5. Experimental Results from Literature
6. Our Experiments
7. Conclusions

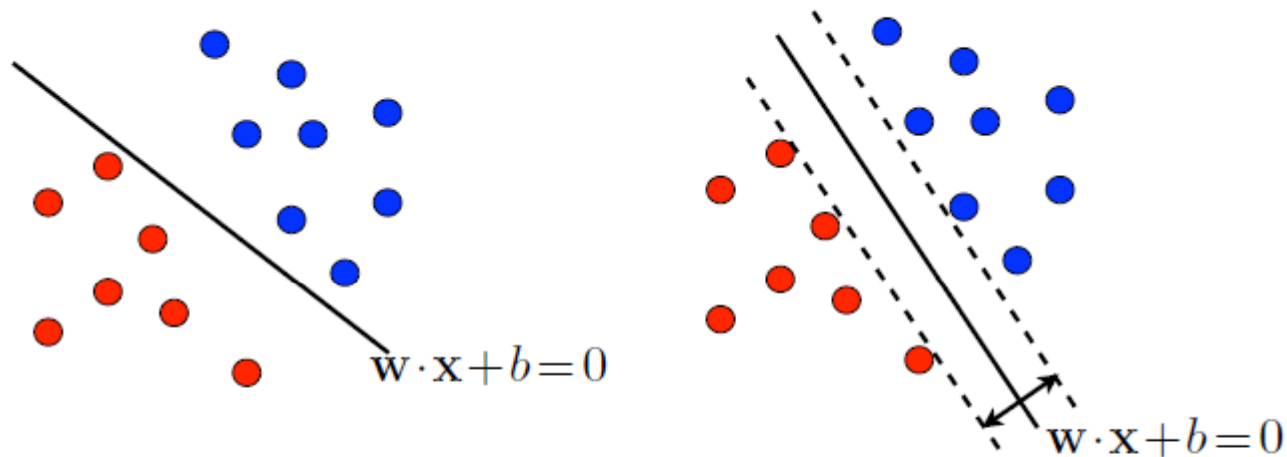
# Binary Classification

- **Training data:** sample drawn i.i.d. from set  $X \subseteq \mathbb{R}^N$  according to some distribution  $D$ ,

$$S = ((x_1, y_1), \dots, (x_m, y_m)) \in X \times \{-1, +1\}.$$

- **Problem:** find hypothesis  $h: X \mapsto \{-1, +1\}$  in  $H$  (classifier) with small generalization error  $R_D(h)$ .
- **Linear classification:**
  - Hypotheses based on hyperplanes.
  - Linear separation in high-dimensional space.

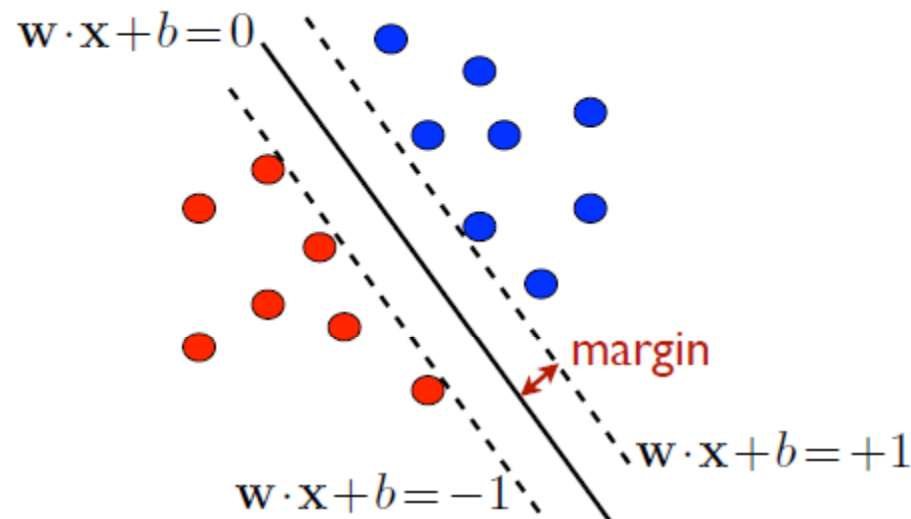
# Linear Separation



■ Classifiers:  $H = \{x \mapsto \text{sgn}(w \cdot x + b) : w \in \mathbb{R}^N, b \in \mathbb{R}\}$ .

Why is it important to maximize the margin?

# Margin Classifier



- **Canonical hyperplane:**  $w$  and  $b$  chosen such that for closest points  $|w \cdot x + b| = 1$ .
- **Margin:**  $\rho = \min_{x \in S} \frac{|w \cdot x + b|}{\|w\|} = \frac{1}{\|w\|}$ .

# Computing the Classifier

## ■ Constrained optimization:

$$\begin{aligned} \min_{\mathbf{w}, b} \quad & \frac{1}{2} \|\mathbf{w}\|^2 \\ \text{subject to} \quad & y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1, i \in [1, m]. \end{aligned}$$

## ■ Properties:

- Convex optimization (strictly convex).
- Unique solution for linearly separable sample.



# Computing the Classifier

## ■ Constrained optimization:

$$\begin{aligned} \min_{\mathbf{w}, b} \quad & \frac{1}{2} \|\mathbf{w}\|^2 \\ \text{subject to} \quad & y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1, i \in [1, m]. \end{aligned}$$

## ■ Properties:

- Convex optimization (strictly convex).
- Unique solution for linearly separable sample.

What to do when the data is not linearly separable?

# Computing the Classifier: for Linearly Non-separable Data

## ■ Constrained optimization:

$$\min_{\mathbf{w}, b, \xi} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m \underbrace{\xi_i}_{\text{Slack variable}}$$

subject to  $y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1 - \xi_i \wedge \xi_i \geq 0, i \in [1, m]$ .

## ■ Properties:

- $C \geq 0$  trade-off parameter.
- Convex optimization (strictly convex).
- Unique solution.

# Primal Optimization Problem

$$\min_{\mathbf{w}, b, \xi} \left[ \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m \xi_i \right]$$

subject to  $y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1 - \xi_i \wedge \xi_i \geq 0, i \in [1, m].$

Obj. Function

Regularizer + error term + constraints

- Regularizer: maximize the margin
- Error term: minimize the number of errors
- C: trade-of between overfitting and underfitting

# Dual Optimization Problem

□ Use Lagrangian to construct the dual problem

■ **Constrained optimization:**

$$\max_{\alpha} \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m \alpha_i \alpha_j y_i y_j (\mathbf{x}_i \cdot \mathbf{x}_j)$$

$$\text{subject to: } \alpha_i \geq 0 \wedge \sum_{i=1}^m \alpha_i y_i = 0, i \in [1, m].$$

■ **Solution:**

$$h(x) = \text{sgn}\left(\sum_{i=1}^m \alpha_i y_i (\mathbf{x}_i \cdot \mathbf{x}) + b\right),$$

$$\text{with } b = y_i - \sum_{j=1}^m \alpha_j y_j (\mathbf{x}_j \cdot \mathbf{x}_i) \text{ for any SV } \mathbf{x}_i.$$

# Dual Optimization Problem

□ Use Lagrangian to construct the dual problem

■ **Constrained optimization:**

$$\max_{\alpha} \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m \alpha_i \alpha_j y_i y_j (\mathbf{x}_i \cdot \mathbf{x}_j)$$

$$\text{subject to: } \alpha_i \geq 0 \wedge \sum_{i=1}^m \alpha_i y_i = 0, i \in [1, m].$$

■ **Solution:**

$$h(x) = \text{sgn}\left(\sum_{i=1}^m \alpha_i y_i (\mathbf{x}_i \cdot \mathbf{x}) + b\right),$$

$$\text{with } b = y_i - \sum_{j=1}^m \alpha_j y_j (\mathbf{x}_j \cdot \mathbf{x}_i) \text{ for any SV } \mathbf{x}_i.$$

**Which one to use:  
Dual or Primal?**

# Dual Optimization Problem

□ Use Lagrangian to construct the dual problem

■ **Constrained optimization:**

$$\max_{\alpha} \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m \alpha_i \alpha_j y_i y_j (\mathbf{x}_i \cdot \mathbf{x}_j)$$

$$\text{subject to: } \alpha_i \geq 0 \wedge \sum_{i=1}^m \alpha_i y_i = 0, i \in [1, m].$$

■ **Solution:**

$$h(x) = \text{sgn}\left(\sum_{i=1}^m \alpha_i y_i (\mathbf{x}_i \cdot \mathbf{x}) + b\right),$$

$$\text{with } b = y_i - \sum_{j=1}^m \alpha_j y_j (\mathbf{x}_j \cdot \mathbf{x}_i) \text{ for any SV } \mathbf{x}_i.$$

**What to do when  
linear classifiers  
do not work?**

Cortes et al.

# Outline



1. Kernel Combination
  1. Heterogeneous Information Fusion
  2. Feature Selection
2. Linear margin classifiers (SVM)
3. *Kernel Classifiers*
4. Kernel Learning / Multiple Kernel Learning
  1. MKL Formulation
  2. Sparse vs. Smooth MKL
  3. MKL Optimization
5. Experimental Results from Literature
6. Our Experiments
7. Conclusions

# Kernel Classifiers

## ■ Idea:

- Define  $K : X \times X \rightarrow \mathbb{R}$ , called **kernel**, such that:

$$\Phi(x) \cdot \Phi(y) = K(x, y).$$

- $K$  often interpreted as a similarity measure.

## ■ Benefits:

- **Efficiency:**  $K$  is often more efficient to compute than  $\Phi$  and the dot product.
- **Flexibility:**  $K$  can be chosen arbitrarily so long as the existence of  $\Phi$  is guaranteed (Mercer's condition).



# Kernel Trick: An Example

## □ Polynomial kernels

### ■ Definition:

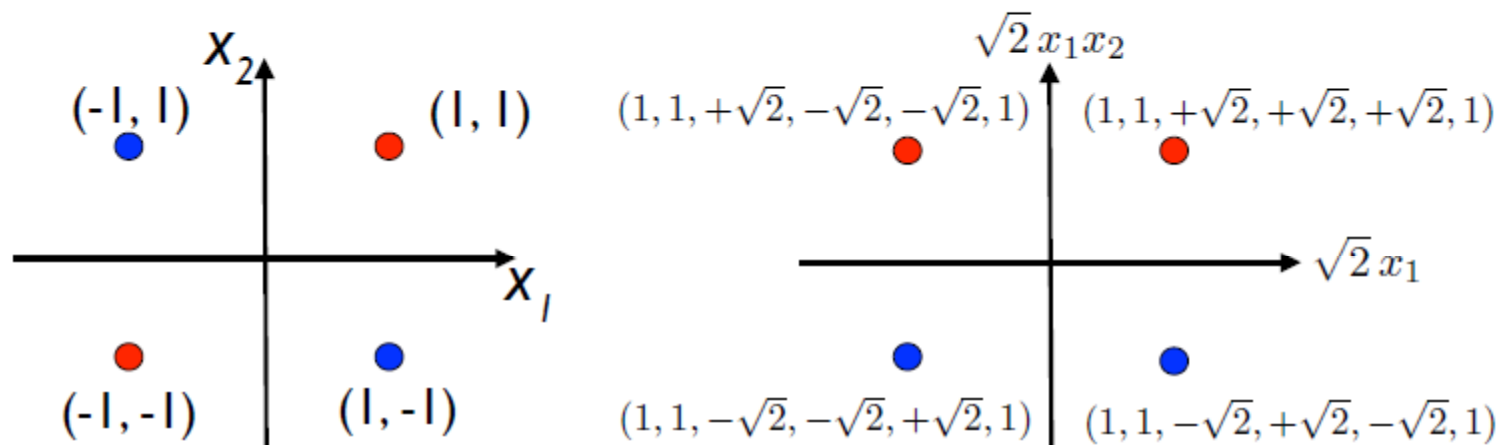
$$\forall x, y \in \mathbb{R}^N, K(x, y) = (x \cdot y + c)^d, \quad c > 0.$$

### ■ Example: for $N=2$ and $d=2$ ,

$$\begin{aligned} K(x, y) &= (x_1 y_1 + x_2 y_2 + c)^2 \\ &= \begin{bmatrix} x_1^2 \\ x_2^2 \\ \sqrt{2} x_1 x_2 \\ \sqrt{2c} x_1 \\ \sqrt{2c} x_2 \\ c \end{bmatrix} \cdot \begin{bmatrix} y_1^2 \\ y_2^2 \\ \sqrt{2} y_1 y_2 \\ \sqrt{2c} y_1 \\ \sqrt{2c} y_2 \\ c \end{bmatrix}. \end{aligned}$$

# XOR Problem: An example

- Use second-degree polynomial kernel with  $c = 1$ :



Linearly non-separable

Linearly separable by  
 $x_1x_2 = 0$ .

- Data is linearly separable in the transformed feature space (polynomial kernel)

# Dual Problem with Kernel Function

## ■ Constrained optimization:

$$\max_{\alpha} \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m \alpha_i \alpha_j y_i y_j K(x_i, x_j)$$

$$\text{subject to: } 0 \leq \alpha_i \leq C \wedge \sum_{i=1}^m \alpha_i y_i = 0, i \in [1, m].$$

## ■ Solution:

$$h(x) = \text{sgn}\left(\sum_{i=1}^m \alpha_i y_i K(x_i, x) + b\right),$$

$$\text{with } b = y_i - \sum_{j=1}^m \alpha_j y_j K(x_j, x_i) \text{ for any } x_i \text{ with } 0 < \alpha_i < C.$$

# Dual Problem with Kernel Function

## ■ Constrained optimization:

*Inner product replaced by the kernel function*

$$\max_{\alpha} \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m \alpha_i \alpha_j y_i y_j K(x_i, x_j)$$

$$\text{subject to: } 0 \leq \alpha_i \leq C \wedge \sum_{i=1}^m \alpha_i y_i = 0, i \in [1, m].$$

## ■ Solution:

*Decision function*

$$h(x) = \text{sgn}\left(\sum_{i=1}^m \alpha_i y_i K(x_i, x) + b\right),$$

$$\text{with } b = y_i - \sum_{j=1}^m \alpha_j y_j K(x_j, x_i) \text{ for any } x_i \text{ with } 0 < \alpha_i < C.$$

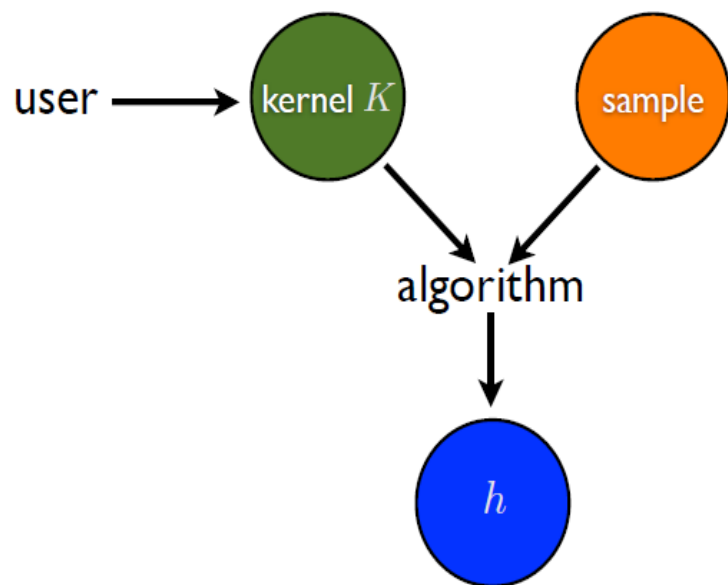
# Kernel-based Learning

- The inner product in the feature space (similarity score) is performed implicitly
- Any linear classification method can be extended to nonlinear feature space
- Non-vectorial data can be utilized (as long as kernel matrix is PSD)
- Which kernel to use? How to set the parameters?
- One kernel for each feature type or for all?

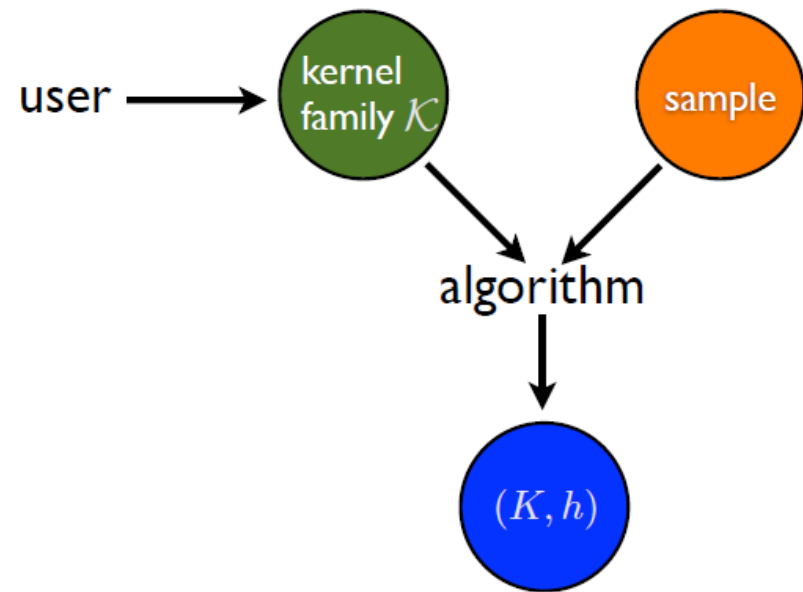
# Outline

1. Kernel Combination
  1. Heterogeneous Information Fusion
  2. Feature Selection
2. Linear margin classifiers (SVM)
3. Kernel Classifiers
4. *Kernel Learning / Multiple Kernel Learning*
  1. *MKL Formulation*
  2. *Sparse vs. Smooth MKL*
  3. *MKL Optimization*
5. Experimental Results from Literature
6. Our Experiments
7. Conclusions

# Kernel Learning



standard kernel  
classifier framework



kernel learning  
framework

# Multiple Kernel Learning (MKL)

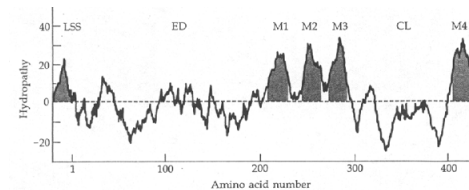
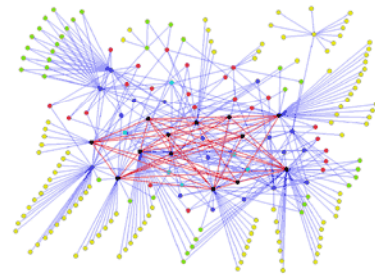
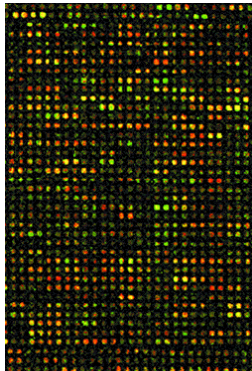


- Multiple kernel learning is a (**parametric**) kernel learning method.
- Some related approaches:
  - ▣ Feature selection
  - ▣ Metric Learning
  - ▣ Ensemble methods

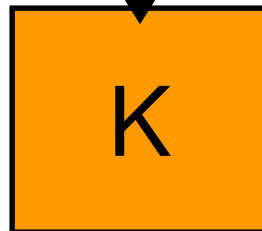


# Kernel Learning Framework

- Example: Protein classification with using heterogeneous information sources



QFDACCFIDVSKIYG-DYGPI  
QFDACCFIDVSKIYG-DHGPI  
QFGACCFIDVSKTFRLHDGPI  
QFDAC-FIDVSKIIFRLHDGPI  
RFDACCFIDVSKIFRLHDGPI  
QFSVYCLIDVSKIYR-HDGP  
QFPVCSIIDDL SKMYR-HDSPV  
QFPVFCLIDDL SKIYR-DDGLI  
QFDARCFIDDL SKIYR-HDGQV  
QFDARCFIDDL SKIYR-HDGQV  
QFDARCFIDDL SKIYR-HDGPI  
RFDACCFIDVSKICK-HDGPV  
QFDACCFIDVSKICK-HDGPV



Lankreit et al.

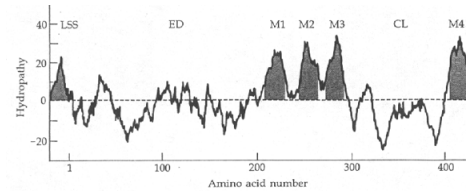
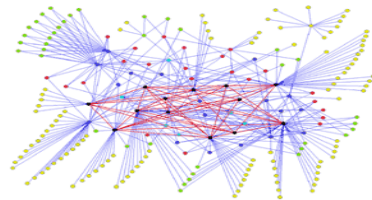
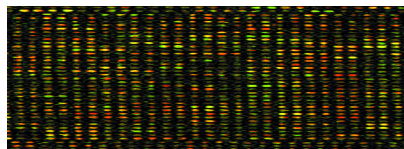
# Multiple Kernel Learning



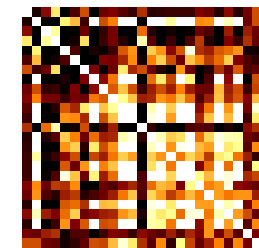
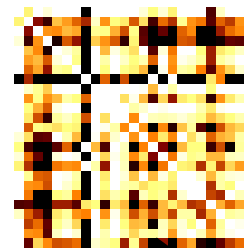
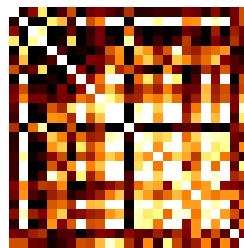
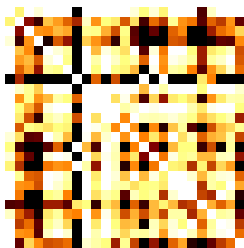
- The overall MKL framework:
  1. Extract features from all available sources
  2. Construct kernel matrices
    1. Different features
    2. Different kernel types
    3. Different kernel parameters
  3. Find the optimal kernel combination and the kernel classifier

# Kernel Learning Framework

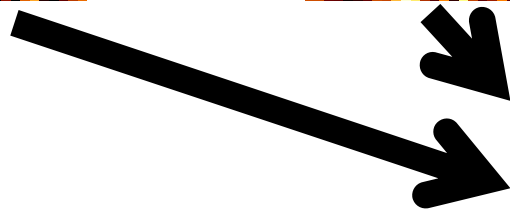
- Create individual kernels for each source (string kernel, diffusion kernel)



1



2



K



Lankreit et al.

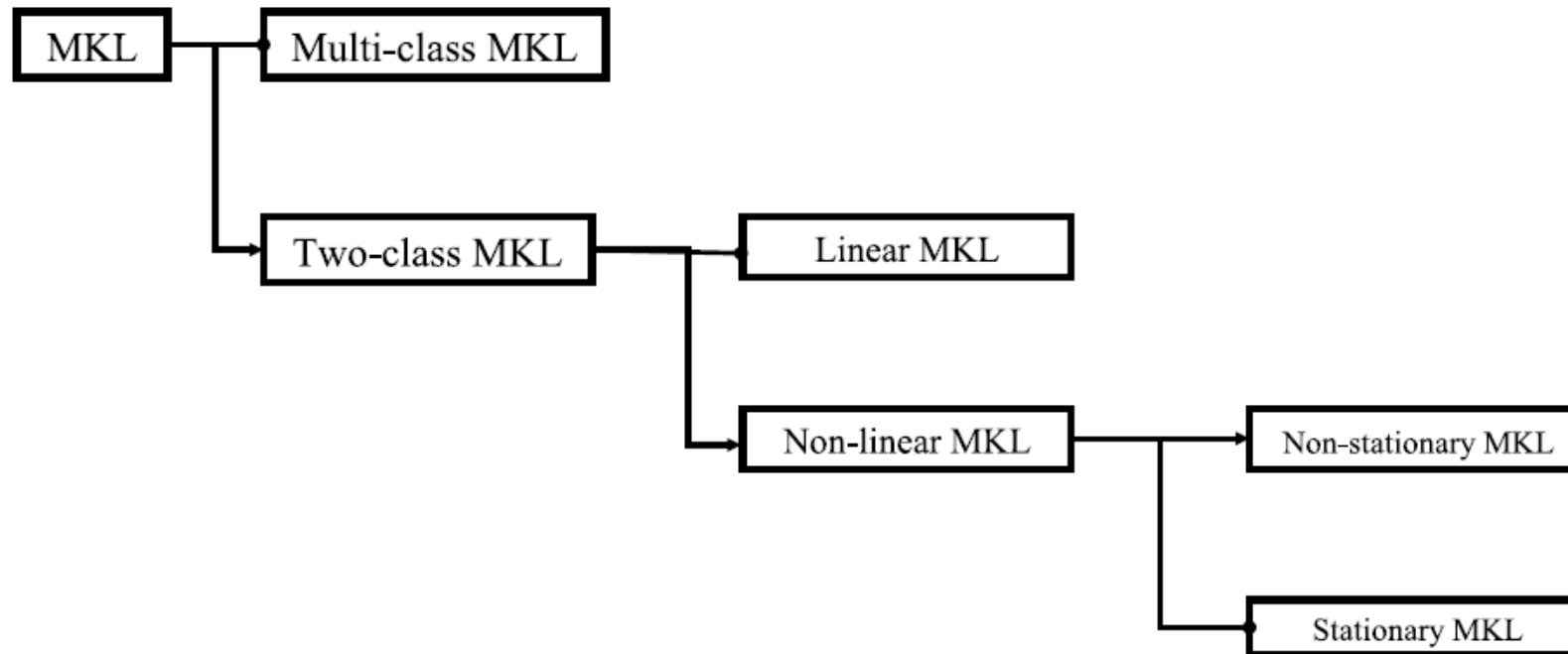
# Multiple Kernel Learning

- The overall MKL framework:
  1. Extract features from all available sources
  2. Construct kernel matrices
    1. Different features
    2. Different kernel types
    3. Different kernel parameters
  3. Find the optimal kernel combination and the kernel classifier

**How is the optimal  
combination computed?**

# Multiple Kernel Learning

- Multiple kernel learning is a **parametric** kernel learning method:



- Our focus is Two-class (binary) MKL

# Linear MKL

## □ Dual MKL formulation

$\circ$  : Hadamard product

$$\min_{\beta \in \Delta} \max_{\alpha \in \mathcal{Q}} \hat{\mathcal{L}}(\alpha, \beta) = \mathbf{1}^\top \alpha - \frac{1}{2} (\alpha \circ \mathbf{y})^\top \mathbf{K}(\beta) (\alpha \circ \mathbf{y}),$$

$$\mathbf{K}(\beta) = \sum_{j=1}^s \beta_j \mathbf{K}_j$$

## □ Similar to SVM dual problem, **BUT**:

- ▣ Minimization w.r.t. the kernel coefficient vector  $\beta$
- ▣ Constraint on the kernel coefficient  $\beta$
- ▣ (Linearly) Combined kernel instead of a single one.

# Regularization on Coefficient Vector

## □ Regularization:

- Leads to well-posed problem
- Affects the level of *sparsity*
- Can lead to *smooth* objective function (or not)

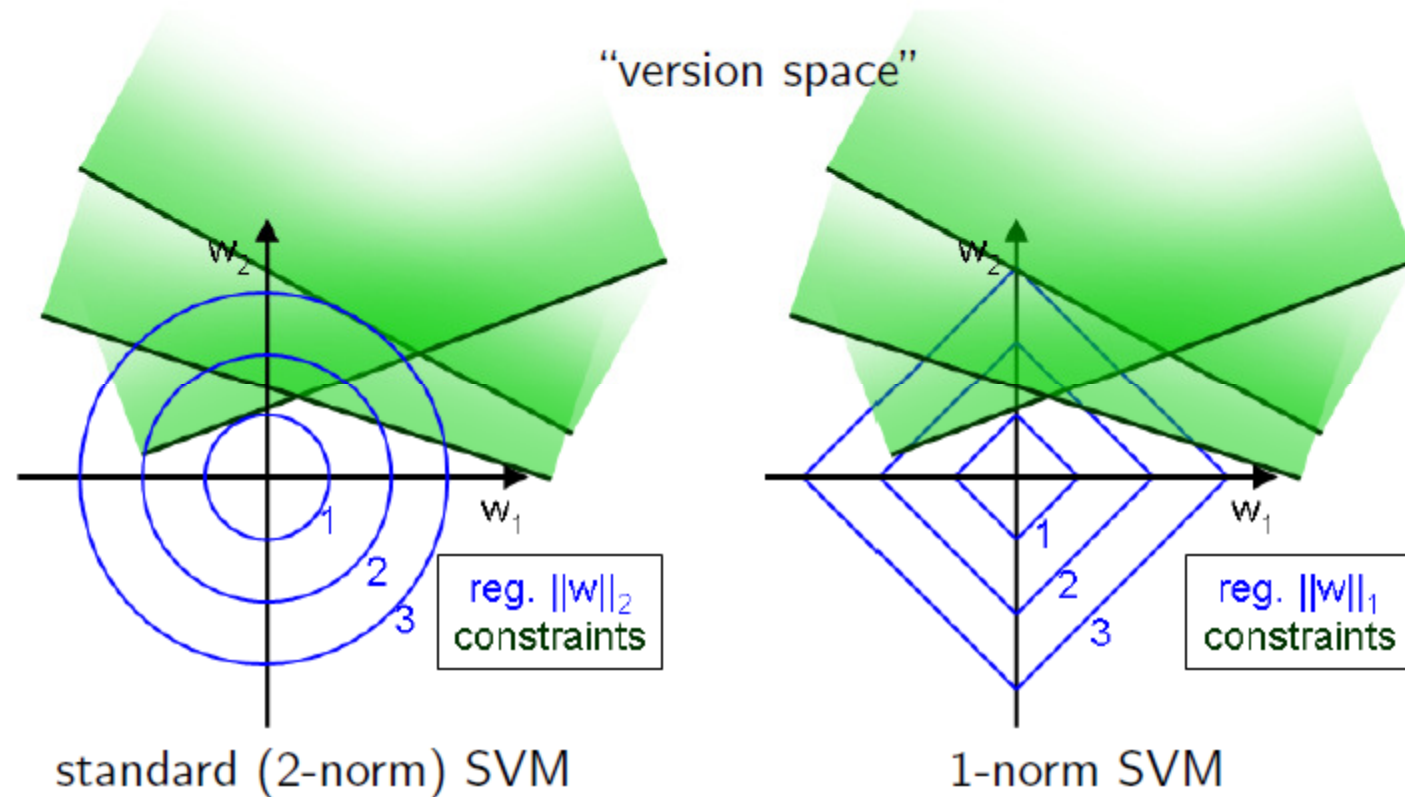
$$\Delta_1 = \left\{ \beta \in \mathbb{R}_+^s : \|\beta\|_1 = \sum_{j=1}^s |\beta^j| \leq 1 \right\}.$$

**L1-MKL**  
**SPARSE solutions**

$$\Delta_p = \left\{ \beta \in \mathbb{R}_+^s : \|\beta\|_p \leq 1 \right\}.$$

**Lp-MKL**  
**DENSE solutions**

# Why Does L1-norm Produces Sparse Solutions?

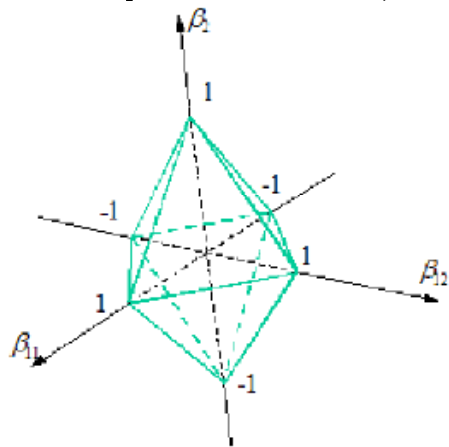


- Feasible solution meets the regularizers at corners

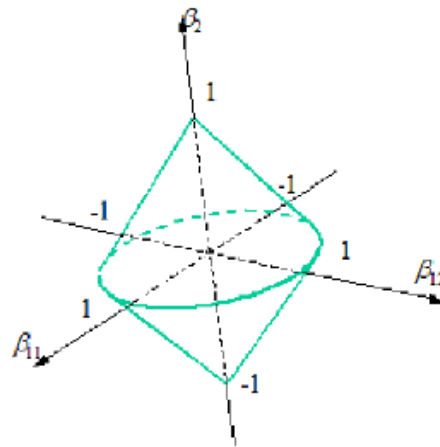


# MKL and Group Lasso

- L1-MKL is equivalent to feature selection with Group-Lasso (each kernel is a group of features)

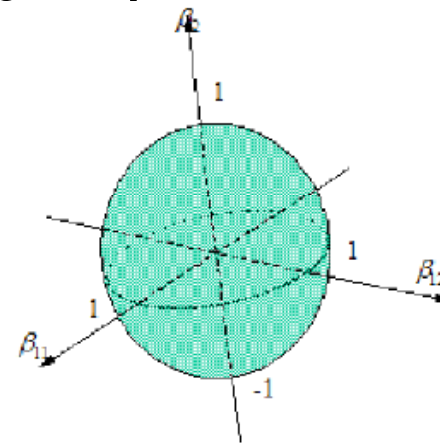


**1-norm SVM,  
lasso:**  
1-norm-constraints  
on all individual  
features



**standard MKL:**

- 1-norm-constraints *between* groups (ie kernels)
- 2-norm-constraints *within* feature groups



**standard SVM,  
ridge-regression:**  
2-norm-constraints  
on all features

# Sparse vs. Dense MKL



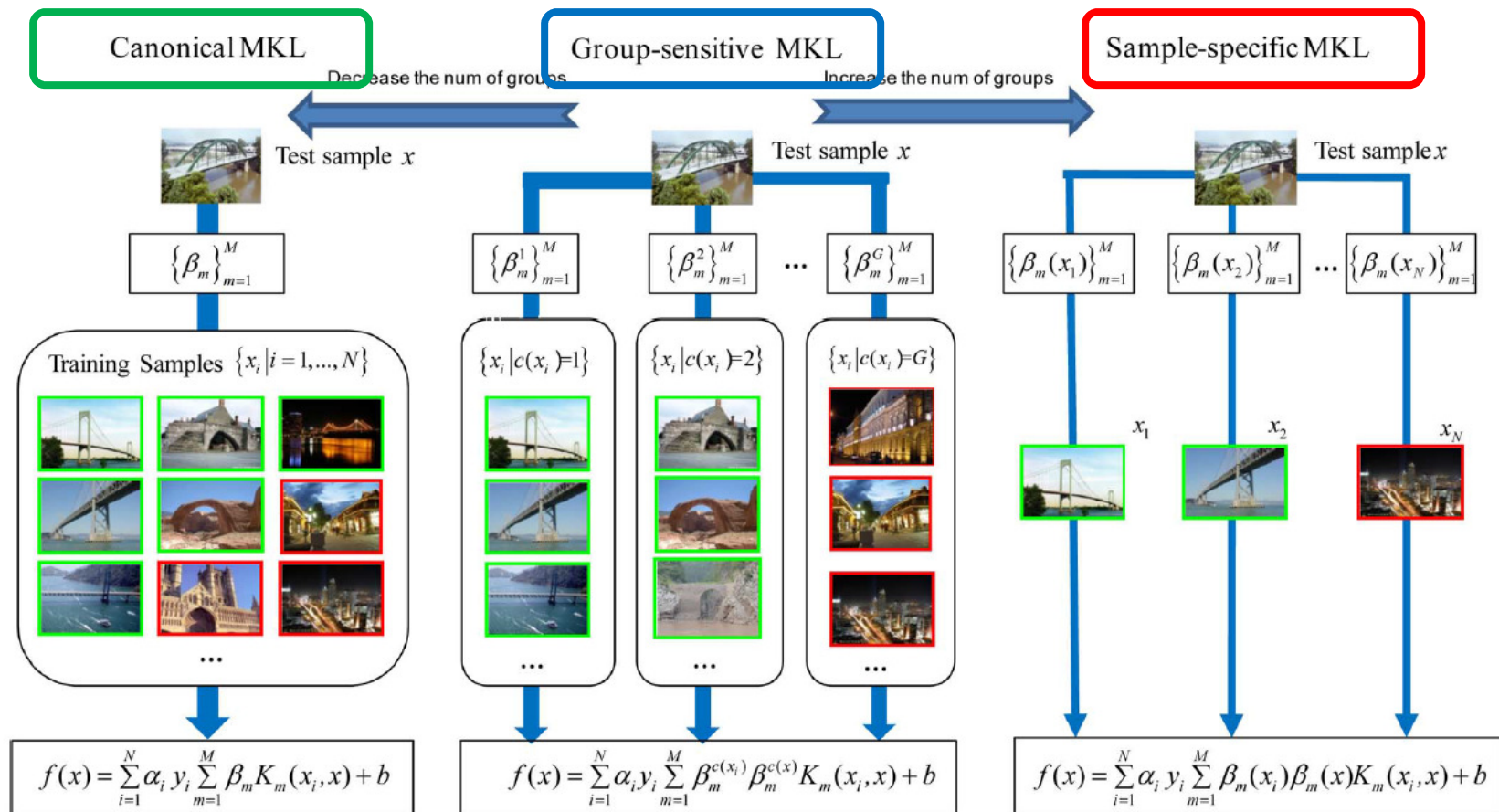
- Which one is better?
  - ▣ Prediction performance
  - ▣ Computational efficiency in training
  - ▣ Computational efficiency in prediction

# Non-linear MKL



- 2 categories:
  - ▣ Stationary: Learn one (and the same) kernel combination for all training samples
  - ▣ Non-stationary: Each sample (or group of samples) uses a different kernel combination

# Non-stationary MKL



(Left) **Canonical MKL**: one kernel combination for all; (Middle) **Group MKL**: one kernel combination for each group; (Right) **Sample-MKL**: one kernel combination for each sample

Yang et al. 2009

# Stationary Non-linear MKL

- Generalized MKL (Varma & Babu 2009 )

Gaussians:  $K_{\mu}(x_i, x_j) = \prod_{k=1}^p \exp(-\mu_k(x_{ik} - x_{jk})^2)$

- Polynomial kernels (Cortes et al. 2009)

Polynomials:  $K_{\mu,d}(x_i, x_j) = (1 + \sum_{k=1}^p \mu_k x_{ik} x_{jk})^d$

- Non-linear MKL methods:

- ▣ Have *high computational cost* due to **non-convexity**
- ▣ Do **NOT** have significant *accuracy improvement*

# Binary-MKL Optimization

- Dual formulation:

$$\min_{\beta \in \Delta} \max_{\alpha \in \mathcal{Q}} \hat{\mathcal{L}}(\alpha, \beta) = \mathbf{1}^\top \alpha - \frac{1}{2} (\alpha \circ \mathbf{y})^\top \mathbf{K}(\beta) (\alpha \circ \mathbf{y}),$$

- *Wrapper approaches* for MKL-L1 and MKL-Lp
- *Direct methods* for regularizers that make the objective function *smooth*

# Wrapper Methods



- 2 alternating steps:
  - ▣ Solve SVM for the fixed combined kernel
  - ▣ Update the kernel weights for the fixed dual variables
- **Pro:** Off-the-shelf SVM methods can be used
- **Con:** Requires high accuracy in intermediate SVM solutions

# Wrapper Methods

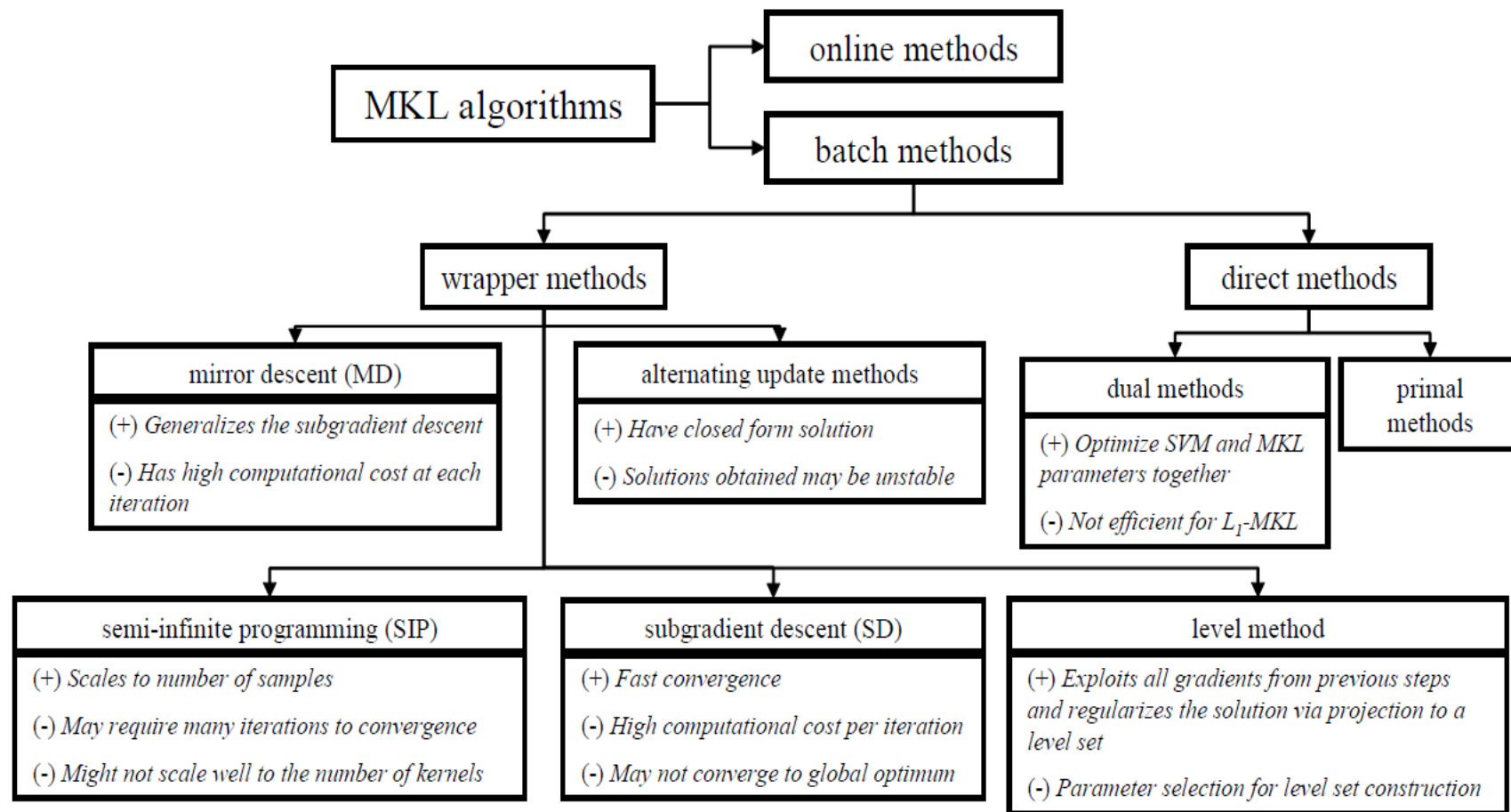


- Semi-infinite Programming (SIP) (Sonnenberg et al. 2005)
- Subgradient descend (Rakotomamonjy et al. 2007)
- Level-Method (Xu et al. 2009)
- Mirror-descend (Aflalo et al. 2011)
- Alternating update method (Group Lasso) (Xu et al. 2010)

They all solve the **same problem**: **NO** accuracy difference when the same regularizer is used!



# Binary-MKL Optimization



*They all solve the same problem (for the same regularizers!)*

# MKL-SMO

- Sequential minimization optimization (SMO) is a very efficient optimization tool that is used for SVM
- When  $\|\beta\|_p^2$  is used as a regularizer, the optimization problem becomes:

$$\max_{\alpha \in Q} \mathbf{1}^\top \alpha - \frac{1}{8\lambda} \left( \sum_{k=1}^s [(\alpha \circ \mathbf{y})^\top \mathbf{K}_k (\alpha \circ \mathbf{y})]^q \right)^{\frac{2}{q}}$$

- SMO can be used for this formulation because it is **smooth**

# Discussions



- The state-of-the-art is linear MKL
  - ▣ Non-linear and online MKL approaches are not mature
- Many alternative optimization methods that solve the same problem.
- Is MKL useful at all?
  - ▣ Compared to uniform (average) combination
- Which Regularizer should be used?
  - ▣ **Sparse** vs **Dense** MKL
- Which optimization method should be selected?

# Outline



1. Kernel Combination
  1. Heterogeneous Information Fusion
  2. Feature Selection
2. Linear margin classifiers (SVM)
3. Kernel Classifiers
4. Kernel Learning / Multiple Kernel Learning
  1. MKL Formulation
  2. Sparse vs. Smooth MKL
  3. MKL Optimization
5. ***Experimental Results from Literature***
6. Our Experiments
7. Conclusions

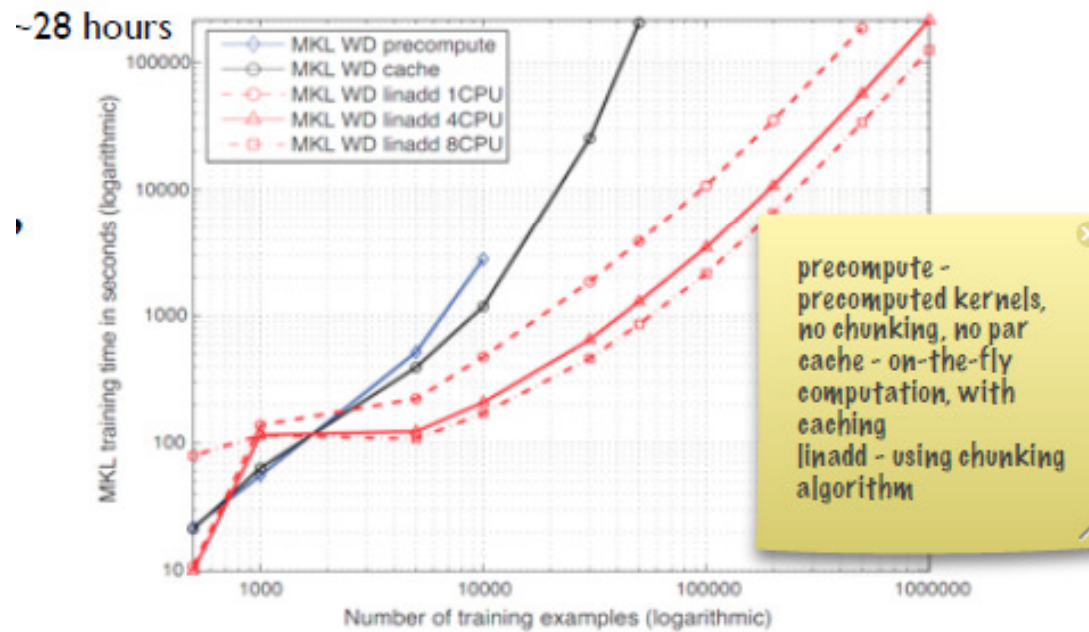
# Experimental Results



- Let's see the reported results from the literature.
- Since all methods solve the same problem, the evaluation is on:
  - ▣ Comparison of accuracy for different norms
  - ▣ Comparison of speed for different norms
  - ▣ Comparison of speed for different optimization methods

# SILP Algorithm- MKL-L1

- 1 M samples, 20 base kernels



- 1 M samples can be processed in less than 28 hours with caching, chunking, and parallelization.

Cortes et al. Sonnenburg et al. 2006

# SimpleMKL vs SILP for MKL-L1

## □ SimpleMKL Reduced gradient methods

Pima  $\ell = 538$   $M = 117$

Algorithm	# Kernel	Accuracy	Time (s)	# SVM eval	# Gradient eval
SILP	$11.6 \pm 1.0$	$76.5 \pm 2.3$	$224 \pm 37$	$95.6 \pm 13$	$95.6 \pm 13$
SimpleMKL	$14.7 \pm 1.4$	$76.5 \pm 2.6$	$79.0 \pm 13$	$314 \pm 44$	$24.3 \pm 4.8$
Grad. Desc.	$14.8 \pm 1.4$	$75.5 \pm 2.5$	$219 \pm 24$	$873 \pm 147$	$118 \pm 8.7$

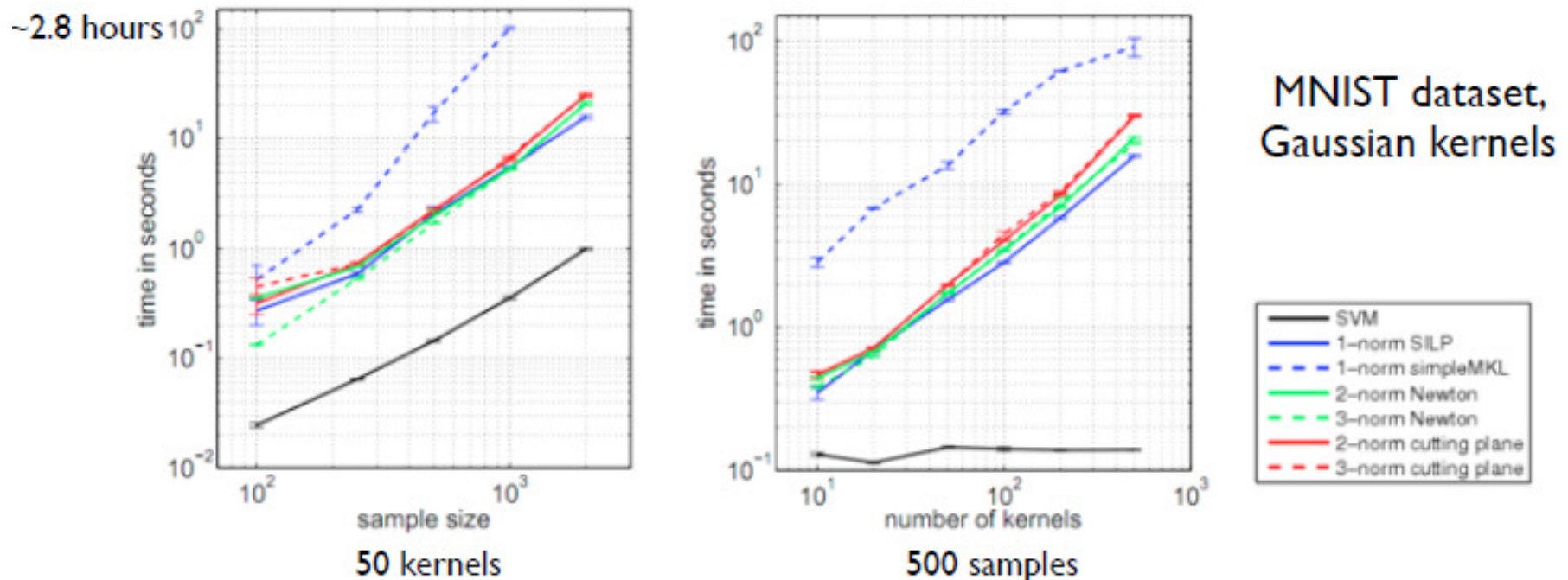
Sonar  $\ell = 146$   $M = 793$

Algorithm	# Kernel	Accuracy	Time (s)	# SVM eval	# Gradient eval
SILP	$33.5 \pm 3.8$	$80.5 \pm 5.1$	$2290 \pm 864$	$903 \pm 187$	$903 \pm 187$
SimpleMKL	$36.7 \pm 5.1$	$80.6 \pm 5.1$	$163 \pm 93$	$2770 \pm 1560$	$115 \pm 66$
Grad. Desc.	$35.7 \pm 3.9$	$80.2 \pm 4.7$	$469 \pm 90$	$7630 \pm 2600$	$836 \pm 99$

- SimpleMKL is faster for **small number** (hundreds) of **samples** and **large number** (hundreds) of **kernels**.

# MKL-Lp Optimization

- Wrapper method for MKL-Lp: Newton or cutting plane (Kloft 2009)

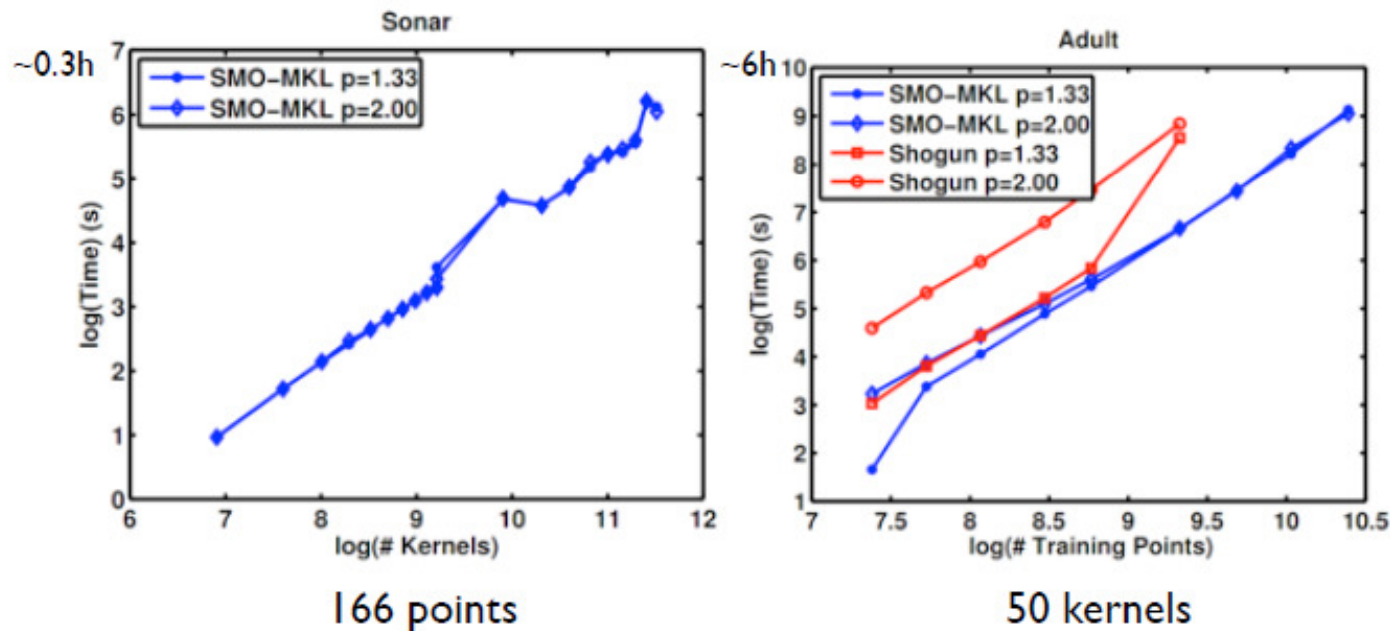


- Faster than MKL-L1 (Sonnenburg 2005)



# MKL-SMO vs. MKL-SIP (MKL-Lp)

- MKL-SMO is a direct approach
- MKL-SIP is a wrapper method

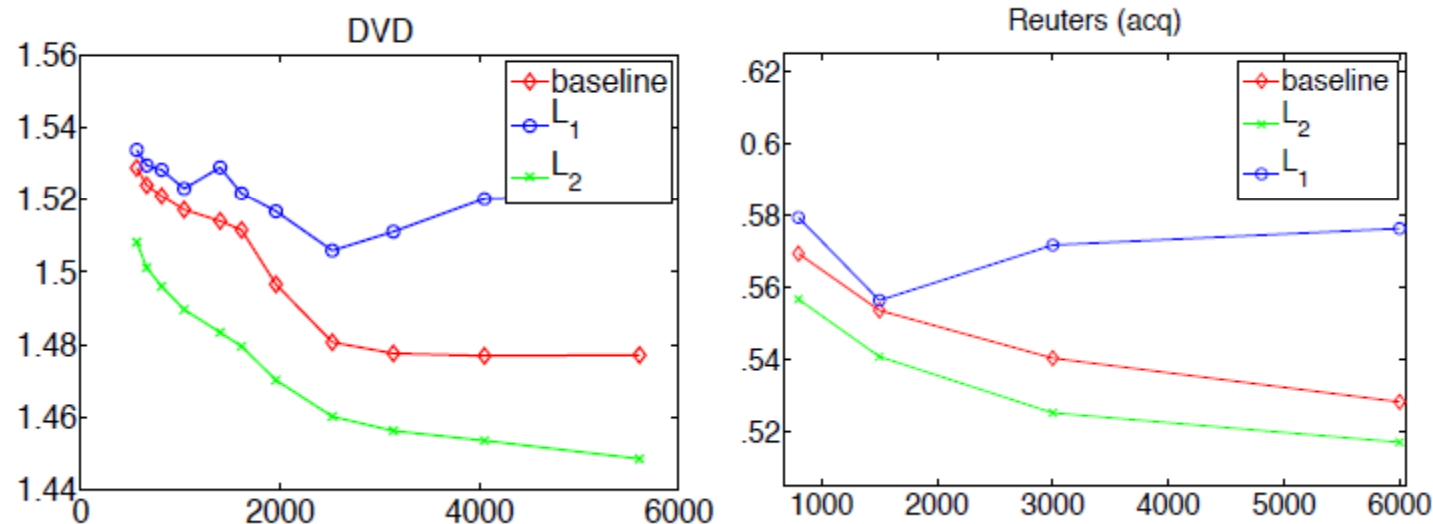


- MKL-SMO is more efficient
- The value of  $p$  does NOT have a big impact on efficiency

Cortes et al. Wishwanathan et al. 2010

# MKL-L1 vs. MKL-L2 (Accuracy)

- Combining kernels on single features – feature selection (Cortes et al. 2009)

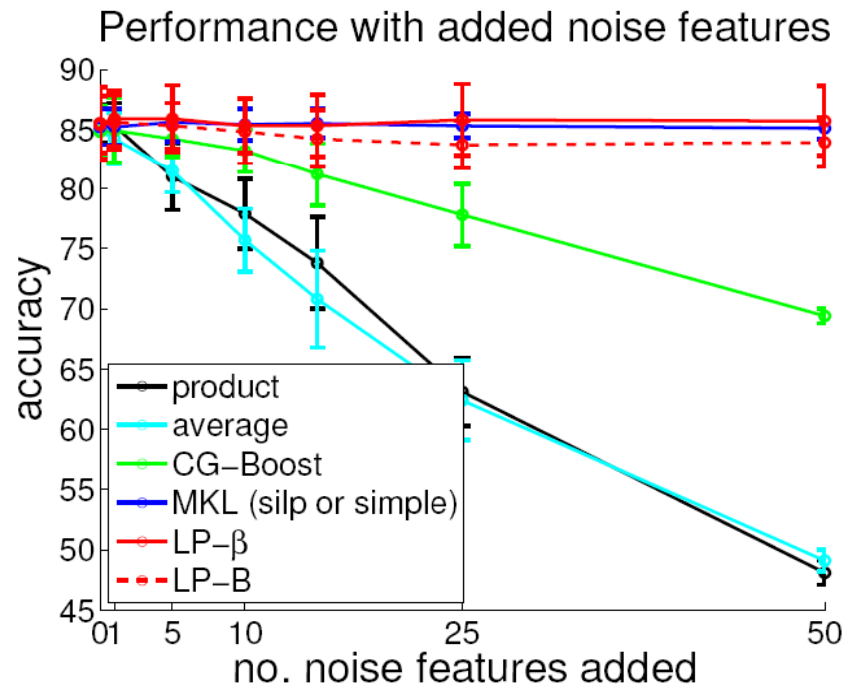


Number of kernels vs. classification error

- Dense combinations are better when using many kernels (MKL-L2 > Average kernel > MKL-L1)

# MKL-L1 vs. Dense Combination

- Oxford Flowers data set
- **Noisy** kernels are added



- MKL-L1 is more robust because it can **eliminate** weak kernels (**contradiction???**)

Cortes et al. Gehler et al.

# Learning Bounds for MKL-L1 & MKL-L2

□ We have the following theorem:

■ Assume that  $K_k(x, x) \leq R^2, \forall k$ . Then, for any  $\delta > 0$ , with probability at least  $1 - \delta$ , for any  $h \in H$

■  $L_2$  regularization :  $R(h) \leq \hat{R}_\rho(h) + O(M^{1/4} \sqrt{\frac{R^2/\rho^2}{n}})$

■  $L_1$  regularization :  $R(h) \leq \hat{R}_\rho(h) + O(\sqrt{(\log M) \frac{R^2/\rho^2}{n}})$

□ Meaning that, given a sufficiently large number of samples ( $n$ ), **MKL-L1 has a better bound** (more robust) w.r.t. the increased number of kernels.

Theorem from a presentation by AI Lab, University of Geneva et al.

# Conclusions So Far



- MKL-L2 is computationally more efficient than MKL-L1
- Dense (smooth) kernel combination performs better than sparse kernel combination
  - ▣ Unless there is a large number of noisy kernels
- MKL-L1 does not have advantage over average kernel baseline
- MKL-SMO is the fastest MKL method
  - ▣ But NOT available for MKL-L1

# Outline



1. Kernel Combination
  1. Heterogeneous Information Fusion
  2. Feature Selection
2. Linear margin classifiers (SVM)
3. Kernel Classifiers
4. Kernel Learning / Multiple Kernel Learning
  1. MKL Formulation
  2. Sparse vs. Smooth MKL
  3. MKL Optimization
5. Experimental Results from Literature
6. ***Our Experiments***
7. Conclusions

# Our Additional Experiments



## □ Data sets

- Caltech 101: 9,146 images, 101 classes, 48 base kernels
- VOC2007: 9,963 images, 20 classes, 15 base kernels
- Subset of ImageNet: 81,738 images, 18 classes, 10 base kernels

# Which Wrapper method is the Most Efficient

Caltech 101 dataset			
baseline	10 training instances per class		
	training	#iter	KerComb
GMKL- $L_1$	34.6 $\pm$ 8.6	38.4 $\pm$ 2.0	27.9 $\pm$ 7.7
SimpleMKL- $L_1$	55.7 $\pm$ 25.3	17.2 $\pm$ 6.8	46.1 $\pm$ 22.0
VSKL- $L_1$	14.1 $\pm$ 2.3	38.3 $\pm$ 4.3	11.1 $\pm$ 1.7
MKL-GL- $L_1$	21.9 $\pm$ 0.8	40.0 $\pm$ 0.0	19.5 $\pm$ 0.8
MKL-GL- $L_2$	5.3 $\pm$ 0.6	8.8 $\pm$ 1.0	4.8 $\pm$ 0.6
MKL-GL- $L_4$	<b>3.5 <math>\pm</math> 0.2</b>	<b>5.9 <math>\pm</math> 0.4</b>	3.2 $\pm$ 0.2
MKL-Level- $L_1$	8.0 $\pm$ 2.3	33.0 $\pm$ 9.5	5.5 $\pm$ 1.4
MKL-SIP- $L_1$	5.4 $\pm$ 0.9	39.4 $\pm$ 2.6	<b>2.1 <math>\pm</math> 0.3</b>
MKL-SIP- $L_2$	<b>3.8 <math>\pm</math> 1.2</b>	<b>5.6 <math>\pm</math> 0.9</b>	<b>2.4 <math>\pm</math> 1.1</b>
MKL-SIP- $L_4$	<b>3.3 <math>\pm</math> 0.6</b>	<b>4.4 <math>\pm</math> 0.5</b>	<b>1.8 <math>\pm</math> 0.6</b>

30 training instances per class			
baseline	30 training instances per class		
	training	#iter	KerComb
GMKL- $L_1$	256.7 $\pm$ 47.7	38.6 $\pm$ 1.8	212.5 $\pm$ 42.3
SimpleMKL- $L_1$	585.6 $\pm$ 204.7	19.0 $\pm$ 7.5	494.4 $\pm$ 174.7
VSKL- $L_1$	121.9 $\pm$ 22.4	36.6 $\pm$ 5.1	103.5 $\pm$ 17.7
MKL-GL- $L_1$	197.1 $\pm$ 9.1	39.8 $\pm$ 1.0	178.3 $\pm$ 8.5
MKL-GL- $L_2$	50.8 $\pm$ 5.6	<b>9.3 <math>\pm</math> 1.0</b>	46.3 $\pm$ 5.2
MKL-GL- $L_4$	32.5 $\pm$ 1.6	<b>5.9 <math>\pm</math> 0.3</b>	29.6 $\pm$ 1.5
MKL-Level- $L_1$	63.3 $\pm$ 22.1	27.5 $\pm$ 11.1	47.9 $\pm$ 14.9
MKL-SIP- $L_1$	44.3 $\pm$ 6.1	39.7 $\pm$ 2.9	23.2 $\pm$ 2.7
MKL-SIP- $L_2$	<b>30.4 <math>\pm</math> 4.2</b>	<b>6.3 <math>\pm</math> 1.0</b>	25.2 $\pm$ 3.9
MKL-SIP- $L_4$	<b>22.6 <math>\pm</math> 2.6</b>	<b>4.7 <math>\pm</math> 0.5</b>	<b>18.2 <math>\pm</math> 2.1</b>

- MKL-Lp is more efficient than MKL-L1
- Kernel combination consumes the most time
  - ▣ MKL-L1 might have advantage when the number of kernels is very large
- MKL-Lp converges in small number of iterations



# Wrapper vs. Direct MKL-Lp

- Performance comparison (training time in sec)
- MKL-SMO (direct MKL-Lp) is faster

		Number of training samples		
		<u><math>n = 10</math></u>	<u><math>n = 20</math></u>	<u><math>n = 30</math></u>
Caltech 101	MKL-SIP	$3.6 \pm 0.2$	$6.5 \pm 0.3$	$11.8 \pm 0.7$
	MKL-SMO	$0.2 \pm 0.1$	$2.3 \pm 0.2$	$3.8 \pm 0.5$
		<u>25%</u>	<u>50%</u>	<u>75%</u>
VOC 2007	MKL-SIP	$15.5 \pm 1.6$	$145.6 \pm 3.9$	$360.7 \pm 8.4$
	MKL-SMO	$3.5 \pm 0.7$	$14.2 \pm 1.8$	$33.1 \pm 3.0$
		Number of base kernels		
		<u><math>K = 48</math></u>	<u><math>K = 63</math></u>	<u><math>K = 108</math></u>
Caltech 101	MKL-SIP	$6.5 \pm 0.3$	$13.6 \pm 2.9$	$19.8 \pm 3.4$
	MKL-SMO	$2.3 \pm 0.2$	$3.2 \pm 0.8$	$6.3 \pm 1.0$
		<u><math>K = 15</math></u>	<u><math>K = 30</math></u>	<u><math>K = 75</math></u>
VOC 2007	MKL-SIP	$145.6 \pm 3.9$	$542.0 \pm 32.8$	$1412.1 \pm 63.4$
	MKL-SMO	$14.2 \pm 1.8$	$29.1 \pm 2.8$	$77.8 \pm 10.3$

# Classification Performance

## Caltech 101

Baseline	Norm	Number of training instances per class		
		10	20	30
Single		45.3 $\pm$ 0.9	55.2 $\pm$ 0.9	70.6 $\pm$ 0.9
Average		<b>59.0 <math>\pm</math> 0.7</b>	<b>69.7 <math>\pm</math> 0.6</b>	77.2 $\pm$ 0.5
MKL-SIP	$p = 1$	53.8 $\pm$ 0.6	63.8 $\pm$ 0.9	<b>83.9 <math>\pm</math> 0.7</b>
MKL-SIP	$p = 2$	<b>60.1 <math>\pm</math> 0.6</b>	<b>70.7 <math>\pm</math> 1.0</b>	79.1 $\pm$ 0.6

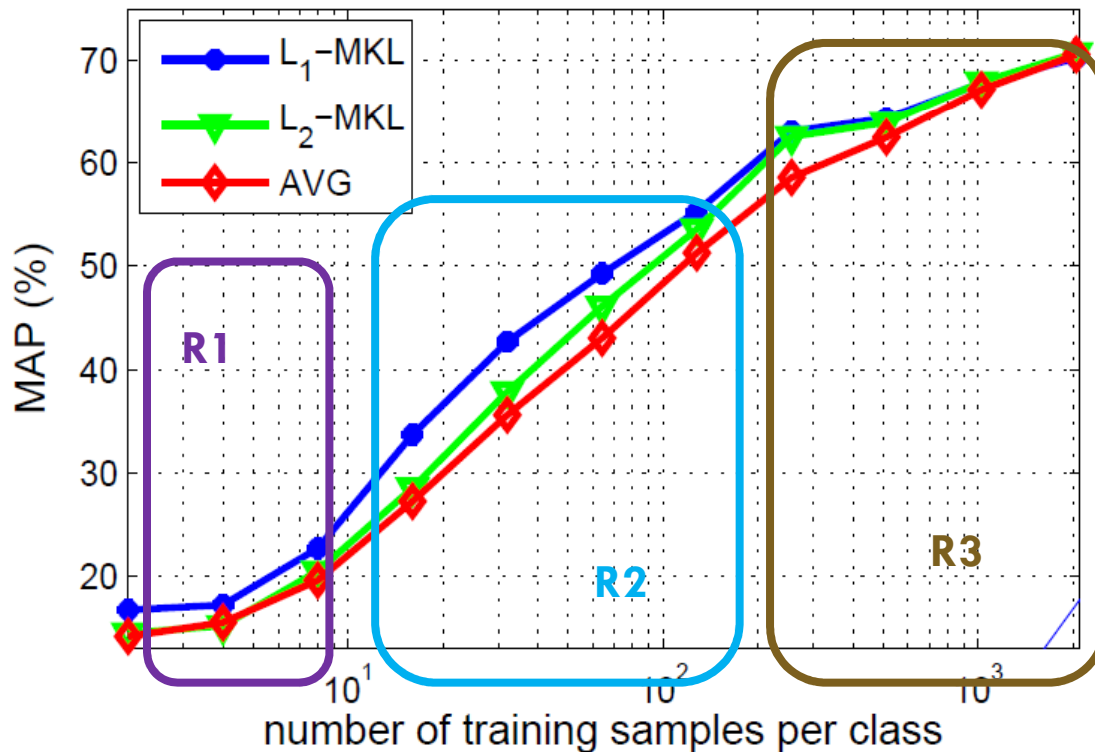
## VOC 2007

baseline	Percentage of the samples used for training			
	1%	25%	50%	75%
Single	<b>23.4 <math>\pm</math> 0.1</b>	44.7 $\pm$ 0.8	48.6 $\pm$ 0.8	50.0 $\pm$ 0.8
Average	21.9 $\pm$ 0.5	48.2 $\pm$ 0.8	54.5 $\pm$ 0.8	57.5 $\pm$ 0.8
$L_1$ -MKL	<b>23.5 <math>\pm</math> 0.7</b>	<b>51.9 <math>\pm</math> 0.4</b>	<b>57.4 <math>\pm</math> 0.4</b>	<b>59.9 <math>\pm</math> 0.9</b>
$L_2$ -MKL	22.7 $\pm$ 0.4	49.8 $\pm$ 0.2	<b>57.3 <math>\pm</math> 0.2</b>	<b>60.6 <math>\pm</math> 0.5</b>

MKL-L1 can perform better when there is a sufficient number of training samples

# ImageNet Results

- MKL-L1 is superior going from R1 to R2 (better kernel selection with increased number of samples)



- MKL-L1 has no advantage in R3:

- All base kernels are stronger with increased number of samples.

- Kernel selection has no edge.

# Adding Weaker Kernels

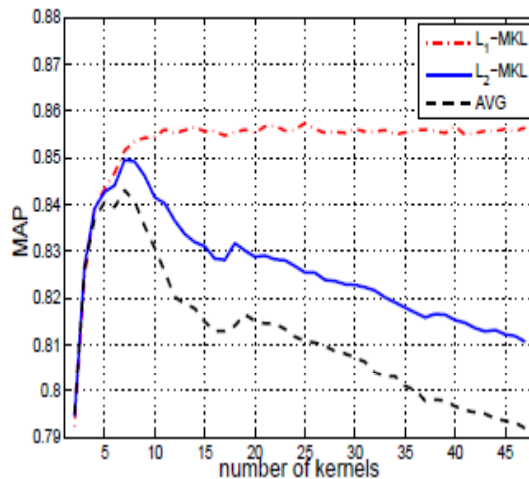


Fig. 4: The change in MAP score with respect to the number of base kernels for the Caltech 101 dataset.

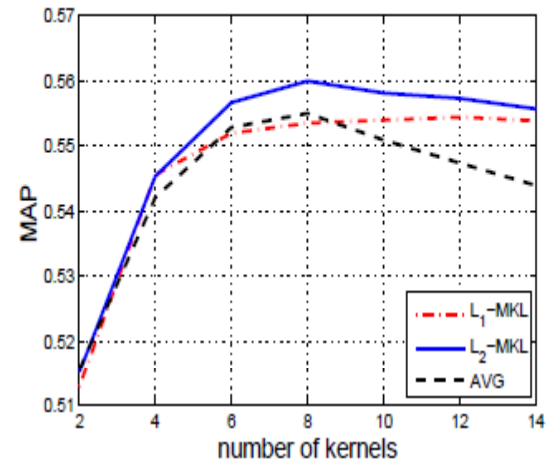


Fig. 5: The change in MAP score with respect to the number of base kernels for the VOC 2007 dataset.

- MKL-L1 is more robust against addition of weak (in terms of accuracy) base kernels

# Outline



1. Kernel Combination
  1. Heterogeneous Information Fusion
  2. Feature Selection
2. Linear margin classifiers (SVM)
3. Kernel Classifiers
4. Kernel Learning / Multiple Kernel Learning
  1. MKL Formulation
  2. Sparse vs. Smooth MKL
  3. MKL Optimization
5. Experimental Results from Literature
6. Our Experiments
7. **Conclusions**

# Conclusions I : Comparison of MKL Methods

- Linear MKL is preferred over non-linear
  - ▣ Unless you know the exact formulation you need
- MKL-L1 is preferred when:
  - ▣ The number of kernels is high (tens or more)
  - ▣ Kernel selection is needed
  - ▣ The number of training samples is sufficient (hundreds per class)
  - ▣ Prediction time is important (sparseness)
- MKL-SMO is the most efficient method (no MKL-L1)
- MKL-SIP is the fastest MKL-L1

# Conclusions II: Some Practical Advice

- Multiple MKL packages are available
  - ▣ MKL-SMO, Shogun etc.
- Kernel normalization is important
- MKL is scalable for thousands of training samples, hundreds of kernels
  - ▣ For more you might need to look into linear methods or kernel approximations
- Average kernel is claimed to perform comparable to sparse MKL
  - ▣ Even if it is true, sparseness gives advantage in prediction speed.

# Open Problems in MKL



- Improve the computational efficiency
  - ▣ Millions of training samples
  - ▣ Thousands of base kernels
- Prediction speed
  - ▣ Kernel methods are in general not very efficient in prediction (storage and use of support vectors)
- Going beyond linear MKL
  - ▣ Online MKL, non-linear MKL, multi-label MKL



# MKL Software

- Shogun (MKL-SIP):
  - ▣ <http://shogun-toolbox.org/>
- MKL-SMO:
  - ▣ <http://research.microsoft.com/en-us/um/people/manik/code/SMO-MKL/download.html>
- SimpleMKL (subgradient)
  - ▣ <http://asi.insa-rouen.fr/enseignants/~arakoto/code/mklindex.html>
- Multi-Label MKL:
  - <http://www.cse.msu.edu/~bucakser/software.html>
- Localized MKL (non-stationary MKL)
  - ▣ <http://users.ics.aalto.fi/gonen/>

# References

- S.S. Bucak, R. Jin, A.K. Jain, "Multiple Kernel Learning for Visual Object Recognition: A Review," *IEEE Transactions on Pattern Analysis and Machine Intelligence*
- M. Gnen and E. Alpaydin, "Multiple Kernel Learning Algorithms," *Journal of Machine Learning Research*, vol. 12, pp. 2211–2268, 2011.
- S. Vishwanathan, Z. Sun, N. Ampornpant, and M. Varma, "Multiple Kernel Learning and the SMO Algorithm," in *NIPS 23, 2010*, pp. 2361–2369.
- M. Kloft, U. Brefeld, S. Sonnenburg, and A. Zien, "Lp-norm Multiple Kernel Learning," *Journal of Machine Learning Research*, vol. 12, pp. 953–997, 2011.
- M. Varma and B. R. Babu, "More generality in efficient multiple kernel learning," in *ICML, June 2009*, pp. 1065–1072.
- J. Yang, Y. Li, Y. Tian, L. Duan, and W. Gao, "Group- Sensitive Multiple Kernel Learning for Object Categorization," in *ICCV, 2009*

# References

- P. V. Gehler and S. Nowozin, “On Feature Combination for Multiclass Object Classification,” in *ICML*, 2009
- C. Cortes, M. Mohri, and A. Rostamizadeh, “Learning Nonlinear Combinations of Kernels,” in *NIPS 22*, 2009, pp. 396–404.
- C. Cortes, M. Mohri, and A. Rostamizadeh, “Generalization Bounds for Learning Kernels,” in *ICML*, 2010.
- A. Rakotomamonjy, F. Bach, Y. Grandvalet, and S. Canu, “SimpleMKL,” *Journal of Machine Learning Research*, vol. 9, pp. 2491–2521, 2008.
- S. Sonnenburg, G. Rätsch, C. Schäfer, and B. Schölkopf, “Large Scale Multiple Kernel Learning,” *Journal of Machine Learning Research*, vol. 7, pp. 1531–1565, 2006.