# Kernel Learning for Novelty Detection

#### John Shawe-Taylor

University College London

NIPS Workshop Kernel Learning: Automatic Selection of Optimal Kernels, December 2008

Joint work with Zakria Hussain

# Introduction

- Motivating problem
- 1-class SVMs

#### 2 Multiple Kernel Learning

 Method 1: constraining the 1-norm of the weight vectors
Method 2: constraining a convex combination of the 1-norm and 2-norm of the weight vectors

#### 3 Experiments

- Assessing impact of µ
- Including negative examples

#### Conclusions

#### Introduction

- Motivating problem
- 1-class SVMs

#### Multiple Kernel Learning

- Method 1: constraining the 1-norm of the weight vectors
- Method 2: constraining a convex combination of the 1-norm and 2-norm of the weight vectors

#### 3 Experiments

- Assessing impact of  $\mu$
- Including negative examples

#### 4 Conclusions

#### Introduction

- Motivating problem
- 1-class SVMs

#### Multiple Kernel Learning

- Method 1: constraining the 1-norm of the weight vectors
- Method 2: constraining a convex combination of the 1-norm and 2-norm of the weight vectors

#### 3 Experiments

- Assessing impact of  $\mu$
- Including negative examples

#### Conclusions

#### Introduction

- Motivating problem
- 1-class SVMs

#### 2 Multiple Kernel Learning

- Method 1: constraining the 1-norm of the weight vectors
- Method 2: constraining a convex combination of the 1-norm and 2-norm of the weight vectors

#### 3 Experiments

- Assessing impact of  $\mu$
- Including negative examples

#### Conclusions

Motivating problem 1-class SVMs

# Outline



#### Introduction

- Motivating problem
- 1-class SVMs

- Method 1: constraining the 1-norm of the weight vectors
- Method 2: constraining a convex combination of the 1-norm and 2-norm of the weight vectors

- Assessing impact of  $\mu$
- Including negative examples

<<p>・

Motivating problem 1-class SVMs

# Content based image retrieval

- Consider problem of content based image retrieval (CBIR) using relevance feedback.
- There are many metrics under which we can compare images: colour, texture, objects included, etc.
- Learning to identify the target of the search is improved if we can identify the metric that best characterises the type of search: eg
  - sunset scene  $\Rightarrow$  colour,
  - Sunset over waterfall  $\Rightarrow$  colour & texture, etc.
- Baseline system is PicSOM uses 11 self-organising maps (SOMs) to represent database of images in 11 metrics – estimates a density of relevant vs irrelevant to weight the vertices of each SOM.
- Implicitly reweights the metrics via the density.

< 2 > < 2 >

Motivating problem 1-class SVMs

# Content based image retrieval

- Consider problem of content based image retrieval (CBIR) using relevance feedback.
- There are many metrics under which we can compare images: colour, texture, objects included, etc.
- Learning to identify the target of the search is improved if we can identify the metric that best characterises the type of search: eg
  - sunset scene  $\Rightarrow$  colour,
  - Sunset over waterfall  $\Rightarrow$  colour & texture, etc.
- Baseline system is PicSOM uses 11 self-organising maps (SOMs) to represent database of images in 11 metrics – estimates a density of relevant vs irrelevant to weight the vertices of each SOM.
- Implicitly reweights the metrics via the density.

Motivating problem 1-class SVMs

# Content based image retrieval

- Consider problem of content based image retrieval (CBIR) using relevance feedback.
- There are many metrics under which we can compare images: colour, texture, objects included, etc.
- Learning to identify the target of the search is improved if we can identify the metric that best characterises the type of search: eg
  - sunset scene  $\Rightarrow$  colour,
  - Sunset over waterfall  $\Rightarrow$  colour & texture, etc.
- Baseline system is PicSOM uses 11 self-organising maps (SOMs) to represent database of images in 11 metrics – estimates a density of relevant vs irrelevant to weight the vertices of each SOM.
- Implicitly reweights the metrics via the density.

Motivating problem 1-class SVMs

# Content based image retrieval

- Consider problem of content based image retrieval (CBIR) using relevance feedback.
- There are many metrics under which we can compare images: colour, texture, objects included, etc.
- Learning to identify the target of the search is improved if we can identify the metric that best characterises the type of search: eg
  - sunset scene  $\Rightarrow$  colour,
  - Sunset over waterfall  $\Rightarrow$  colour & texture, etc.
- Baseline system is PicSOM uses 11 self-organising maps (SOMs) to represent database of images in 11 metrics – estimates a density of relevant vs irrelevant to weight the vertices of each SOM.
- Implicitly reweights the metrics via the density.

Motivating problem 1-class SVMs

# Content based image retrieval

- Consider problem of content based image retrieval (CBIR) using relevance feedback.
- There are many metrics under which we can compare images: colour, texture, objects included, etc.
- Learning to identify the target of the search is improved if we can identify the metric that best characterises the type of search: eg
  - sunset scene  $\Rightarrow$  colour,
  - Sunset over waterfall  $\Rightarrow$  colour & texture, etc.
- Baseline system is PicSOM uses 11 self-organising maps (SOMs) to represent database of images in 11 metrics – estimates a density of relevant vs irrelevant to weight the vertices of each SOM.
- Implicitly reweights the metrics via the density.

Motivating problem 1-class SVMs

#### 1-class SVMs

- Negative data is sparse so consider 1-class learning initially.
- Metrics correspond to kernels: so task is about using combination of kernels to solve a retrieval task.
- If we include 'learning the kernel', we can automatically identify the relevant metrics for the particular search.
- Potential to scale to very large numbers of metrics/submetrics.

(日)

Motivating problem 1-class SVMs

#### 1-class SVMs

- Negative data is sparse so consider 1-class learning initially.
- Metrics correspond to kernels: so task is about using combination of kernels to solve a retrieval task.
- If we include 'learning the kernel', we can automatically identify the relevant metrics for the particular search.
- Potential to scale to very large numbers of metrics/submetrics.

・ロト ・ 日 ・ ・ 回 ・ ・ 日 ・

Motivating problem 1-class SVMs

#### 1-class SVMs

- Negative data is sparse so consider 1-class learning initially.
- Metrics correspond to kernels: so task is about using combination of kernels to solve a retrieval task.
- If we include 'learning the kernel', we can automatically identify the relevant metrics for the particular search.
- Potential to scale to very large numbers of metrics/submetrics.

・ロ・・ (日・・ (日・・ 日・)

Motivating problem 1-class SVMs

### 1-class SVMs

- Negative data is sparse so consider 1-class learning initially.
- Metrics correspond to kernels: so task is about using combination of kernels to solve a retrieval task.
- If we include 'learning the kernel', we can automatically identify the relevant metrics for the particular search.
- Potential to scale to very large numbers of metrics/submetrics.

Motivating problem 1-class SVMs

### Optimisation problem

$$\begin{array}{ll} \min_{\mathbf{w},\xi} & \frac{1}{2} \|\mathbf{w}\|_2^2 + C \, \|\xi\|_1 \\ \text{subject to} & \langle \mathbf{w}, \phi(\mathbf{x}_i) \rangle \geq 1 - \xi_i \\ & \xi_i \geq 0, \ i = 1, \dots, m \end{array}$$

John Shawe-Taylor Kernel Learning for Novelty Detection

臣

Method 1: constraining the 1-norm of the weight vectors Method 2: constraining a convex combination of the 1-norm a

(日)

# Outline

Introduction
Motivating problem
1-class SVMs

#### 2 Multiple Kernel Learning

- Method 1: constraining the 1-norm of the weight vectors
- Method 2: constraining a convex combination of the 1-norm and 2-norm of the weight vectors

#### 3 Experiments

- Assessing impact of µ
- Including negative examples

#### 4 Conclusions

Method 1: constraining the 1-norm of the weight vectors Method 2: constraining a convex combination of the 1-norm a

・ロト ・ 日 ・ ・ 回 ・ ・ 日 ・

### A linear combination of kernels

Let  $\kappa_k$  denote the *k*th kernel from a set  $\mathbb{K} = {\kappa_1, \ldots, \kappa_{|\mathbb{K}|}}$  of kernels. We define a weighted combination of kernels like so:

$$\kappa_{\mathbf{z}} = \sum_{k=1}^{|\mathbb{K}|} z_k \kappa_k$$

where 
$$\mathbf{z} = ig(z_1, \dots, z_{|\mathbb{K}|}ig)$$
 ,  $z_i \in \mathbb{R}^+.$ 

Method 1: constraining the 1-norm of the weight vectors Method 2: constraining a convex combination of the 1-norm a

・ロト ・ 日 ・ ・ 回 ・ ・ 日 ・

MKL Optimisation: Constraining the 1-norm (Recap)

Let  $K = |\mathbb{K}|$ , then we have the following 1-class SVM for MKL when regularising over the weight vector using the 1-norm (primal):

$$\begin{array}{ll} \min_{\mathbf{w}_{k},\xi} & \frac{1}{2} \left( \sum_{k=1}^{K} \|\mathbf{w}_{k}\|_{2} \right)^{2} + C \|\xi\|_{1} \\ \text{subject to} & \sum_{k=1}^{K} \left\langle \mathbf{w}_{k}, \phi_{k}(\mathbf{x}_{i}) \right\rangle \geq 1 - \xi_{i} \\ & \xi_{i} \geq 0, \ i = 1, \dots, m \end{array}$$

Q

Method 1: constraining the 1-norm of the weight vectors Method 2: constraining a convex combination of the 1-norm a

・ロ・ ・ 四・ ・ 回・ ・ 日・

臣

Constraining the 1-norm continued

The dual becomes:

 $\min_{\beta} \max_{\alpha}$  subject to

$$\sum_{\substack{i,j=1\\ j \neq 1}}^{m} \alpha_i \alpha_j \kappa_k(\mathbf{x}_i, \mathbf{x}_j) \le \beta,$$
  
$$\sum_{\substack{i=1\\ i=1}}^{m} \alpha_i = 1,$$
  
$$0 \le \alpha_i \le C, i = 1, \dots, m.$$

Method 1: constraining the 1-norm of the weight vectors Method 2: constraining a convex combination of the 1-norm a

(日)

- We ran experiments with a set of Gaussian kernels composed of different width parameters.
- Problem: only chose 1 kernel for learning, namely the Gaussian kernel with the largest width parameter. This is also true for experiments conducted with the VOC data sets (cat, cow, dog).
- Our conjecture: except in degenerate cases will only choose 1 kernel.
- A solution: Constrain a convex combination of the 1-norm and 2-norm in the optimisation problem.

Method 1: constraining the 1-norm of the weight vectors Method 2: constraining a convex combination of the 1-norm a

(日)

- We ran experiments with a set of Gaussian kernels composed of different width parameters.
- Problem: only chose 1 kernel for learning, namely the Gaussian kernel with the largest width parameter. This is also true for experiments conducted with the VOC data sets (cat, cow, dog).
- Our conjecture: except in degenerate cases will only choose 1 kernel.
- A solution: Constrain a convex combination of the 1-norm and 2-norm in the optimisation problem.

Method 1: constraining the 1-norm of the weight vectors Method 2: constraining a convex combination of the 1-norm a

・ロ・ ・ 四・ ・ 回・ ・ 回・

- We ran experiments with a set of Gaussian kernels composed of different width parameters.
- Problem: only chose 1 kernel for learning, namely the Gaussian kernel with the largest width parameter. This is also true for experiments conducted with the VOC data sets (cat, cow, dog).
- Our conjecture: except in degenerate cases will only choose 1 kernel.
- A solution: Constrain a convex combination of the 1-norm and 2-norm in the optimisation problem.

Method 1: constraining the 1-norm of the weight vectors Method 2: constraining a convex combination of the 1-norm a

・ロ・ ・ 四・ ・ 回・ ・ 回・

- We ran experiments with a set of Gaussian kernels composed of different width parameters.
- Problem: only chose 1 kernel for learning, namely the Gaussian kernel with the largest width parameter. This is also true for experiments conducted with the VOC data sets (cat, cow, dog).
- Our conjecture: except in degenerate cases will only choose 1 kernel.
- A solution: Constrain a convex combination of the 1-norm and 2-norm in the optimisation problem.

Method 1: constraining the 1-norm of the weight vectors Method 2: constraining a convex combination of the 1-norm a

# MKL Optimisation: Constraining a combination of the 1-norm and 2-norm

Constraining using both these norms gives us the following primal problem, with  $\mu$  controlling the sparsity trade-off,

$$\begin{array}{ll} \min_{\mathbf{w}_{k},\xi} & \frac{\mu}{2} \left( \sum_{k=1}^{K} \|\mathbf{w}_{k}\|_{2} \right)^{2} + \frac{1-\mu}{2} \sum_{k=1}^{K} \|\mathbf{w}_{k}\|_{2}^{2} + C \|\xi\|_{1} \\ \text{subject to} & \sum_{k=1}^{K} \langle \mathbf{w}_{k}, \phi_{k}(\mathbf{x}_{i}) \rangle \geq 1 - \xi_{i} \\ & \xi_{i} \geq 0, \ i = 1, \dots, m \end{array}$$

Method 1: constraining the 1-norm of the weight vectors Method 2: constraining a convex combination of the 1-norm a

・ロ・ ・ 四・ ・ 回・ ・ 回・

# Constraining a combination of the 1-norm and 2-norm continued

Let  $D = \sum_{k=1}^{K} \|\mathbf{w}_k\|_2$ , then the dual is:

$$\begin{array}{ll} \max_{\boldsymbol{\alpha}} & W(\boldsymbol{\alpha}) = \sum_{i=1}^{m} \alpha_i - \frac{A}{2} \sum_{k=1}^{K} \beta_k + \frac{B}{2} \left( \sum_{k \in J} \sqrt{\beta_k} \right)^2 \\ \text{subject to} & \beta_k = \sum_{i,j=1}^{m} \alpha_i \alpha_j \kappa_k(\mathbf{x}_i, \mathbf{x}_j) \\ & 0 \le \alpha_i \le C, \ i = 1, \dots, m \end{array}$$

where  $A = 1/(1 - \mu)$  and  $B = ((|J| - 1)\mu^2 + \mu)/((1 - \mu)(1 - \mu + \mu|J|)^2)$ , where  $J = \{k : z_k \neq 0\}$ , is the set of indices *k*, for which

$$\beta_k > \mu^2 D^2,$$

Method 1: constraining the 1-norm of the weight vectors Method 2: constraining a convex combination of the 1-norm a

・ロ・ ・ 四・ ・ 回・ ・ 日・

크

Combination of 1- and 2-norm continued

From the Lagrangian we get,

$$z_k = \max\left\{0, \frac{1}{1-\mu}\left(\frac{\sqrt{\beta_k}}{D} - \mu\right)\right\},$$

Also,

$$D = \frac{\sum_{k \in J} \sqrt{\beta_k}}{1 - \mu + \mu |J|}$$

Method 1: constraining the 1-norm of the weight vectors Method 2: constraining a convex combination of the 1-norm a

Algorithm for combination of 1- and 2-norm

#### We perform coordinate-wise descent in the $\alpha$ vector. Writing

$$g_i(\alpha_i) = \frac{\partial W(\alpha)}{\partial \alpha_i},$$

where  $\alpha_i$  is the *i*-th coordinate of  $\alpha$  in the argument of  $W(\cdot)$ , we seek the solution of  $g_i(\alpha_i) = 0$  as the new value for  $\alpha_i$ .

Method 1: constraining the 1-norm of the weight vectors Method 2: constraining a convex combination of the 1-norm a

Algorithm for combination of 1- and 2-norm

We expand  $g_i(\alpha_i)$  in a Taylor series around the current values  $\alpha^0$ :

$$g_i(\alpha_i) \approx \frac{\partial W(\alpha^0)}{\partial \alpha_i} + \frac{\partial W^2(\alpha^0)}{\partial \alpha_i^2} (\alpha_i - \alpha_i^0) = \mathbf{0}$$

and solve for  $\alpha_i$ . Hence our update rule for each  $\alpha_i$  becomes:

$$\alpha_i = \alpha_i^{\mathbf{0}} - \frac{\frac{\partial W(\alpha^0)}{\partial \alpha_i}}{\frac{\partial W^2(\alpha^0)}{\partial \alpha_i^2}}.$$

Method 1: constraining the 1-norm of the weight vectors Method 2: constraining a convex combination of the 1-norm a

# Algorithm for combination of 1- and 2-norm

Initialise  $\alpha^0$  vector to zero with one element, say  $\alpha_1^0 > 0$ .

- Repeat until KKT conditions satisfied or ||α<sup>n</sup> − α<sup>n−1</sup>||<sub>2</sub> < ε, where ε is a small positive real number
  - Compute update rule for each component of  $\alpha$  using:

$$\alpha_i = \alpha_i^{\mathbf{0}} - \frac{\frac{\partial W(\alpha^0)}{\partial \alpha_i}}{\frac{\partial W^2(\alpha^0)}{\partial \alpha_i^2}}.$$

• update z and D

Decision function:

$$f(\mathbf{x}) = \sum_{i=1}^{m} \alpha_i \sum_{k \in J} \frac{z_k}{\mu + (1-\mu)z_k} \kappa_k(\mathbf{x}_i, \mathbf{x})$$

Assessing impact of  $\mu$  Including negative examples

# Outline

Introduction
Motivating problem
1-class SVMs

#### 2 Multiple Kernel Learning

- Method 1: constraining the 1-norm of the weight vectors
- Method 2: constraining a convex combination of the 1-norm and 2-norm of the weight vectors

### 3 Experiments

- Assessing impact of  $\mu$
- Including negative examples

#### Conclusions

<<p>・

Assessing impact of  $\mu$  Including negative examples

#### Datasets considered

#### Considered the PASCAL VOC data: cat, cow, dog

11 feature sets extracted for PICSOM

・ロト ・ 日 ・ ・ 回 ・ ・ 日 ・

Assessing impact of  $\mu$  Including negative examples

#### Datasets considered

- Considered the PASCAL VOC data: cat, cow, dog
- 11 feature sets extracted for PICSOM

Feature	dimensions
DCT coefficients of average colour in rectangular grid	12
CIE L*a*b* colour of two dominant colour clusters	6
Histogram of local edge statistics	80
Haar transform of quantized HSV colour histogram	256
Histogram of interest point SIFT features	256
Average CIE L*a*b* colour	15
Three central moments of CIE L*a*b* colour distribution	45
Histogram of four Sobel edge directions	20
Co-occurrence matrix of four Sobel edge directions	80
Magnitude of the 16 $\times$ 16 FFT of Sobel edge image	128
Histogram of relative brightness of neighboring pixels	40

・ロト ・雪 ・ ・ ヨ ・

Assessing impact of  $\mu$  Including negative examples

### Effect of $\mu$ on sparsity



Figure: Sparsity as function of  $\mu$  for cats

John Shawe-Taylor Kernel Learning for Novelty Detection

・ロト ・ 日 ・ ・ 回 ・ ・ 日 ・

E

Assessing impact of  $\mu$  Including negative examples

#### Effect of $\mu$ on retrieval



#### Figure: Average precision 20 against $\mu$ for cats

#### Note that PicSOM uses negative examples

< 3 >

Assessing impact of  $\mu$  Including negative examples

#### Effect of $\mu$ on retrieval



Figure: Average precision 20 against  $\mu$  for cats

Note that PicSOM uses negative examples

Assessing impact of  $\mu$  Including negative examples

## Including negative examples

 Included negatives by negating features, i.e. negating kernel entries between differently labelled images



#### Figure: Average precision 20 against $\mu$ for cats

Assessing impact of  $\mu$  Including negative examples

#### Complete precision/recall curves



#### Figure: Precision/recall curve for cats

John Shawe-Taylor Kernel Learning for Novelty Detection

< 17 ▶

-

Assessing impact of  $\mu$  Including negative examples

#### Average Precision scores

Obj.	MKL 2-class (µ : 0.5)			MKL 1-class (µ : 0.5)			1-class SVM		PicSOM	
	AP20	AP50	#ker	AP20	AP50	#ker	AP20	AP50	AP20	AP50
Cat	0.52	0.46	3	0.34	0.24	2	0.14	0.13	0.25	0.25
Cow	0.29	0.20	5	0.17	0.14	3	0.14	0.12	0.25	0.20
Dog	0.37	0.36	11	0.11	0.13	2	0.11	0.12	0.28	0.28

John Shawe-Taylor Kernel Learning for Novelty Detection

・ロ・ ・ 四・ ・ 回・ ・ 日・

臣

# Outline

Introduction
Motivating problem
1-class SVMs

#### 2 Multiple Kernel Learning

- Method 1: constraining the 1-norm of the weight vectors
- Method 2: constraining a convex combination of the 1-norm and 2-norm of the weight vectors

#### 3 Experiments

- Assessing impact of µ
- Including negative examples

# 4 Conclusions

• (1) • (

# Conclusions

- Considered CBIR task in which learning the search metric corresponds to learning the kernel
- In 1-class MKL don't get variable sparsity by varying C
- Flexible mix of 1-norm and 2-norm regularisation gives natural control of sparsity with good performance against PicSOM and 1-class SVM on VOC cats
- Using negative (non-relevant) examples improves performance.
- SOM uses density learning to weight metrics should compare with same approach in kernel methods

(日)

# Conclusions

- Considered CBIR task in which learning the search metric corresponds to learning the kernel
- In 1-class MKL don't get variable sparsity by varying C
- Flexible mix of 1-norm and 2-norm regularisation gives natural control of sparsity with good performance against PicSOM and 1-class SVM on VOC cats
- Using negative (non-relevant) examples improves performance.
- SOM uses density learning to weight metrics should compare with same approach in kernel methods

・ロト ・ 日 ・ ・ 回 ・ ・ 日 ・

# Conclusions

- Considered CBIR task in which learning the search metric corresponds to learning the kernel
- In 1-class MKL don't get variable sparsity by varying C
- Flexible mix of 1-norm and 2-norm regularisation gives natural control of sparsity with good performance against PicSOM and 1-class SVM on VOC cats
- Using negative (non-relevant) examples improves performance.
- SOM uses density learning to weight metrics should compare with same approach in kernel methods

・ロ・ ・ 四・ ・ 回・ ・ 日・

# Conclusions

- Considered CBIR task in which learning the search metric corresponds to learning the kernel
- In 1-class MKL don't get variable sparsity by varying C
- Flexible mix of 1-norm and 2-norm regularisation gives natural control of sparsity with good performance against PicSOM and 1-class SVM on VOC cats
- Using negative (non-relevant) examples improves performance.
- SOM uses density learning to weight metrics should compare with same approach in kernel methods

・ロ・ ・ 四・ ・ 回・ ・ 日・

# Conclusions

- Considered CBIR task in which learning the search metric corresponds to learning the kernel
- In 1-class MKL don't get variable sparsity by varying C
- Flexible mix of 1-norm and 2-norm regularisation gives natural control of sparsity with good performance against PicSOM and 1-class SVM on VOC cats
- Using negative (non-relevant) examples improves performance.
- SOM uses density learning to weight metrics should compare with same approach in kernel methods