

# Online-Batch Strongly Convex Multi Kernel Learning

Francesco Orabona<sup>1</sup>   Luo Jie<sup>2,3</sup>   Barbara Caputo<sup>2</sup>

<sup>1</sup>DSI, Università degli Studi di Milano, Milano, Italy

<sup>2</sup>Idiap Research Institute, Martigny, Switzerland

<sup>3</sup>EPF Lausanne

23<sup>rd</sup> IEEE Conference on Computer Vision and Pattern Recognition

# Outline

- 1 Multi Kernel Learning
  - Notation
  - Previous work
  - Sparsity?
  - Dual?
- 2 OBSCURE
  - A different MKL formulation
  - The algorithm
- 3 Experimental Results
  - Caltech-101

# Outline

- 1 Multi Kernel Learning
  - Notation
  - Previous work
  - Sparsity?
  - Dual?
- 2 OBSCURE
  - A different MKL formulation
  - The algorithm
- 3 Experimental Results
  - Caltech-101

# Problem definition

- We are given  $N$  training samples  $\{\mathbf{x}_i, y_i\}_{i=1}^N$ ,  $\mathbf{x}_i \in \mathbb{X}$  and  $y_i \in \mathbb{Y} = \{1, \dots, M\}$ , where  $M$  is the number of classes.
- We also have  $F$  kernels corresponding to different features, e.g. color, shape, etc.
- We want to learn a *score function*  $s(\mathbf{x}, y)$  that classifies a sample  $\mathbf{x}$  as

$$\arg \max_y s(\mathbf{x}, y).$$

# Linear classification

- We consider score functions of the form

$$s(\mathbf{x}, y) = \sum_{j=1}^F s^j(\mathbf{x}, y)$$

- Defining joint feature maps  $\phi^j(\mathbf{x}, y)$  on data  $\mathbb{X}$  and labels  $\mathbb{Y}$  [Tsochantaridis et al, 2004].

$$s^j(\mathbf{x}, y) = \mathbf{w}^j \cdot \phi^j(\mathbf{x}, y),$$

- Defining with  $\bar{\mathbf{w}} = [\mathbf{w}^1, \mathbf{w}^2, \dots, \mathbf{w}^F]$ , and  $\bar{\phi}(\mathbf{x}, y) = [\phi^1(\mathbf{x}, y), \dots, \phi^F(\mathbf{x}, y)]$ , we have

$$s(\mathbf{x}, y) = \bar{\mathbf{w}} \cdot \bar{\phi}(\mathbf{x}, y).$$

# Multi Kernel Learning

- In MKL we minimize

$$\lambda \|\bar{\mathbf{w}}\|_{2,1}^2 + \frac{1}{N} \sum_{i=1}^N \ell(\bar{\mathbf{w}}, \mathbf{x}_i, y_i).$$

where  $\|\bar{\mathbf{w}}\|_{2,1}^2 = \|\|\mathbf{w}^1\|_2, \|\mathbf{w}^2\|_2, \dots, \|\mathbf{w}^F\|_2\|_1$

- The regularization induces sparsity in the domain of the kernels.
- All the proposed algorithms use an alternating optimization strategy, through the dual formulation.

# Multi Kernel Learning

- In MKL we minimize

$$\lambda \|\bar{\mathbf{w}}\|_{2,1}^2 + \frac{1}{N} \sum_{i=1}^N \ell(\bar{\mathbf{w}}, \mathbf{x}_i, y_i).$$

where  $\|\bar{\mathbf{w}}\|_{2,1}^2 = \|\|\mathbf{w}^1\|_2, \|\mathbf{w}^2\|_2, \dots, \|\mathbf{w}^F\|_2\|_1$

- The regularization induces **sparsity** in the domain of the kernels.
- All the proposed algorithms use an alternating optimization strategy, through the **dual formulation**.

# Sparsity vs “PhD Student Kernels”

- But, are we sure this is the right regularization?
- Is sparsity really needed?
- Do we really want to use just a subset of the available kernels, given that each one is the result of *years* of research??





# Why using the dual?

- Historically, dual formulation for SVM has been introduced to have an *easier* optimization problem and to use *kernels*.

# Why using the dual?

- Historically, dual formulation for SVM has been introduced to have an *easier* optimization problem and to use *kernels*.
- However dual is not needed for neither of the two!
- Stochastic Sub-Gradient Descent algorithms for the primal have been proven to be better than optimizing the dual [Shalev-Shwartz and Srebro ICML08]!

# Use your favorite loss!

- With stochastic sub-gradient descent methods you can use easily *any* loss.
- Computational efficient for large dataset.
- If the objective function is strongly convex functions we can prove fast convergence rate bound to the optimal solution.
  - The algorithm will converge to the optimal solution with a rate  $\mathcal{O}(\frac{1}{T})$ .
  - For alternating optimization methods this is not possible.

# Use your favorite loss!

- With stochastic sub-gradient descent methods you can use easily *any* loss.
- Computational efficient for large dataset.
- If the objective function is strongly convex functions we can prove fast convergence rate bound to the optimal solution.
  - The algorithm will converge to the optimal solution with a rate  $\mathcal{O}(\frac{1}{T})$ .
  - For alternating optimization methods this is not possible.

However the group norm  $(2, 1)$  is not strongly convex...

# Outline

- 1 Multi Kernel Learning
  - Notation
  - Previous work
  - Sparsity?
  - Dual?
- 2 **OBSCURE**
  - A different MKL formulation
  - The algorithm
- 3 Experimental Results
  - Caltech-101

## $(2, p)$ group norm for MKL

- We propose to generalize the MKL formulation using the  $(2, p)$  group norm

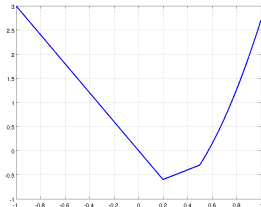
$$\frac{\lambda}{2} \|\bar{\mathbf{w}}\|_{2,p}^2 + \frac{1}{N} \sum_{i=1}^N \ell(\bar{\mathbf{w}}, \mathbf{x}_i, y_i),$$

where  $\|\bar{\mathbf{w}}\|_{2,p} = \left\| \|\mathbf{w}^1\|_2, \|\mathbf{w}^2\|_2, \dots, \|\mathbf{w}^F\|_2 \right\|_p$ .

- When  $p = 1$  we recover the sparse MKL formulation,  $p = 2$  corresponds to using the sum of the kernels.
- A similar formulation has been proposed in [Kloft et al. NIPS09].
- If  $p \in (1, 2]$  this new formulation is  $(1 - 1/p)$ -strongly convex.

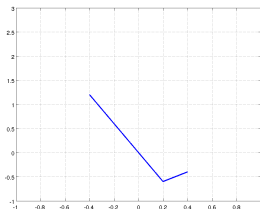
# A small ball is better than a big one...

- We want to minimize a convex function.
- If someone tells us that the solution is living in a small ball the problem is easier.
  - We can use this information with proximal regularization methods.



# A small ball is better than a big one...

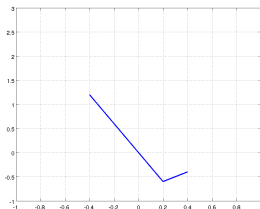
- We want to minimize a convex function.
- If someone tells us that the solution is living in a small ball the problem is easier.
  - We can use this information with proximal regularization methods.





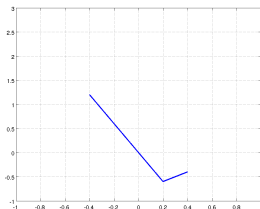
# A small ball is better than a big one...

- We want to minimize a convex function.
- If someone tells us that the solution is living in a small ball the problem is easier.
  - We can use this information with proximal regularization methods.
- But how to estimate this ball?



# A small ball is better than a big one...

- We want to minimize a convex function.
- If someone tells us that the solution is living in a small ball the problem is easier.
  - We can use this information with proximal regularization methods.
- But how to estimate this ball?
- **Solution:** use a fast online algorithm!



# Online-Batch Strongly Convex mUlti keRnel lEarning

- Start a quick online  $(2, \rho)$  MKL algorithm.
- Stop it at any time to obtain an estimate of the radius of the ball,  $R$ , where the optimal solution lives.
- Start a stochastic gradient descent algorithm for the  $(2, \rho)$  MKL problem, using the previous solution as starting point and the information on the radius of the ball.

# Convergence rate for OBSCURE

## Theorem

Let  $1 < p \leq 2$ , and  $q = \frac{p}{p-1}$ ,  $R$  the value returned by the online stage. Then in expectation after  $T$  iterations of the 2nd stage of the OBSCURE algorithm, the gap from the optimal solution is

$$\mathcal{O} \left( F^{1/q} \min \left( \frac{q}{\lambda T}, \frac{R\sqrt{q}}{\sqrt{T}} \right) \right)$$

# Convergence rate for OBSCURE

## Theorem

Let  $1 < p \leq 2$ , and  $q = \frac{p}{p-1}$ ,  $R$  the value returned by the online stage. Then in expectation after  $T$  iterations of the 2nd stage of the OBSCURE algorithm, the gap from the optimal solution is

$$\mathcal{O} \left( F^{1/q} \min \left( \frac{q}{\lambda T}, \frac{R\sqrt{q}}{\sqrt{T}} \right) \right)$$

Moreover, if the problem is linearly separable by a hyperplane  $\bar{\mathbf{u}}$ , the first stage will stop after  $4qF^{2/q} \|\bar{\mathbf{u}}\|_{2,p}^2$  updates,  $R$  will overestimate the radius of the ball at most by a factor of 4.

# A draft of the general algorithm

**Input:**  $q, \bar{\theta}_1, \bar{\mathbf{w}}_1, R, \lambda$

**for**  $t = 1, 2, \dots, T$  **do**

Sample at random  $(\mathbf{x}_t, y_t)$

Theory tells us how to set  $\eta_t$  and  $\alpha_t$

$$\bar{\theta}_{t+\frac{1}{2}} = \alpha_t \bar{\theta}_t + \eta_t \partial \ell(\bar{\mathbf{w}}_t, \mathbf{x}_t, y_t)$$

$$\mathbf{w}_{t+1}^j = \frac{1}{q} \left( \frac{\|\theta_{t+1}^j\|_2}{\|\bar{\theta}_{t+1}\|_{2,q}} \right)^{q-2} \theta_{t+1}^j, \quad \forall j = 1, \dots, F$$

**end for**

# $\alpha_t$ and $\eta_t$ are the core of the algorithm

- The choice of  $\alpha_t$  and  $\eta_t$  are critical to guarantee fast convergence to the optimal solution.
- Our particular choice is given by the theory: the details are in the paper.
- We just want to try it? Fine! Grab the source code at:  
<http://dogma.sourceforge.net>
  - Discriminative Online (Good?) Matlab Algorithms
  - The library is explicitly designed to have easy to modify algorithms.

# Outline

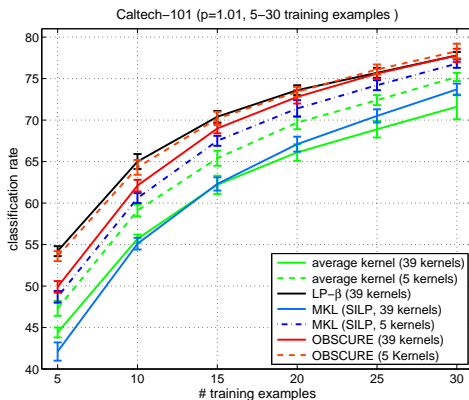
- 1 Multi Kernel Learning
  - Notation
  - Previous work
  - Sparsity?
  - Dual?
- 2 OBSCURE
  - A different MKL formulation
  - The algorithm
- 3 Experimental Results
  - Caltech-101



# Baseline

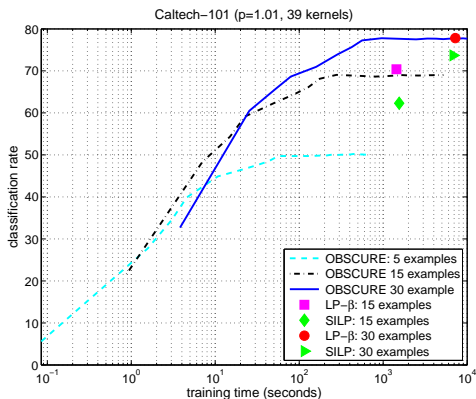
- We compared OBSCURE to SILP [Sonnenburg et al. JMLR06], LP- $\beta$  [Gehler and Nowozin ICCV09] and to SVM using average of all the kernels.
- We used the Caltech-101 with 39 kernels, as in [Gehler and Nowozin ICCV09].

# Caltech-101 Experiments: Performance



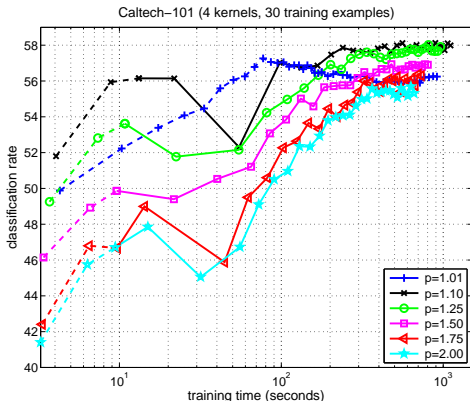
- OBSCURE is better than average kernel.
- Performance on par of LP- $\beta$ .

# Caltech-101 Experiments: Time



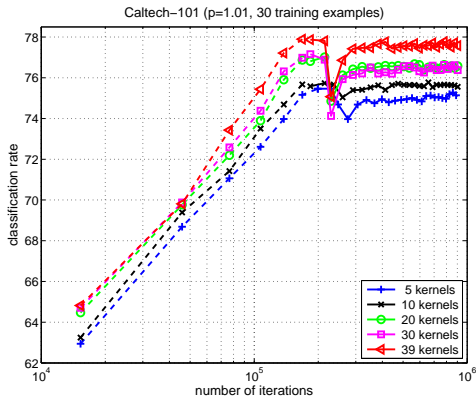
- With 15 samples, time similar to LP- $\beta$  and SILP.
- With 30 samples, OBSCURE is 7-10 times faster than LP- $\beta$  and SILP.

# Different settings of $p$



- When there are few good kernel, the sparse solution is worst.
- The optimal one corresponds to  $p = 1.1$ .

# More kernels = faster convergence



- We reach a given classification rate faster if we use more kernels.

# Summary

- We have introduced a new formulation for MKL problems and an algorithm to solve it.
- The online stage of OBSCURE quickly estimates the region where the solution lives.
- The second stage reaches the solution with a guaranteed convergence rate.

## Future work

- Extending OBSCURE to work with hierarchical losses.

# Thanks for your attention

Code: `http://dogma.sourceforge.net`  
My website: `http://francesco.orabona.com`