SpicyMKL Efficient multiple kernel learning method using dual augmented Lagrangian





The University of Tokyo

Graduate School of Information Science and Technology

Department of Mathematical Informatics

23rd Oct., 2009@ISM

Multiple Kernel Learning

Fast Algorithm SpicyMKL



- Gaussian, polynomial, chi-square,
 - Parameters : Gaussian width, polynomial degree

Features

 Computer Vision: color, gradient, sift (sift, hsvsift, huesift, scaling of sift), Geometric Blur, image regions, . . .

Multiple Kernel Learning (Lanckriet et al. 2004) Select important kernels and combine them

Multiple Kernel Learning (MKL) Single Kernel $f(x) = \sum_{i=1}^{N} \alpha_{i} k(x_{i}, x)$

$$f(x) = \sum_{i=1}^{n} \alpha_i k(x_i, x)$$

Multiple Kernel

$$f(x) = \sum_{i=1}^{N} \alpha_i \left(\sum_{m=1}^{M} \frac{d_m k_m(x_i, x)}{m} \right)$$

d_m:convex combination, Sparse (many 0 components)



Finite dimensional problem

$$(K_m)_{i,j} := k_m(x_i, x_j)$$
 : Gram matrix of *m*th kernel
 $\|\alpha_m\|_m := \sqrt{\alpha_m^T K_m \alpha_m}$



- Scales well against # of kernels.
- Formulated in a general framework.

Outline

- Introduction
- Details of Multiple Kernel Learning
- Our method
 - Proximal minimization
 - Skipping inactive kernels
- Numerical Experiments
 - Bench-mark datasets
 - Sparse or dense

Details of Multiple Kernel Learning (MKL)

Generalized Formulation





Rank 1 decomposition

If g is monotonically increasing and f_{ℓ} is diffible at the optimum, the solution of MKL is rank 1:

$$\exists d_m \in \mathbb{R}, \exists \beta \in \mathbb{R}^N \text{ such that}$$

$$\alpha_m^* = d_m \beta \quad (\forall m) \qquad \sum_{m=1}^M d_m = 1, \ d_m \ge 0$$

$$f^*(x) = \sum_{i=1}^N \beta_i \left(\sum_{m=1}^M d_m k_m(x_i, x)\right)$$
convex combination of kernels

Proof Derivative w.r.t. α_m

$$K_m \nabla f_\ell(K\alpha^*) + \partial g(\|\alpha_m^*\|_m) \frac{K_m}{\|\alpha_m^*\|_m} \alpha_m^* = 0$$

$$\implies \qquad \alpha_m^* = \frac{\|\alpha_m^*\|_m}{\partial g(\|\alpha_m^*\|_m)} \left(-\nabla f_\ell(K\alpha^*)\right)$$
13

Existing approach 1

Constraints based formulation:

[Lanckriet et al. 2004, JMLR] [Bach, Lanckriet & Jordan 2004, ICML]

primal

$$\begin{array}{l}
\min_{\alpha} f_{\ell}(K\alpha) + \frac{1}{2} \left(\sum_{m=1}^{M} \|\alpha_{m}\|_{m} \right)^{2} \\
\end{array}$$
Dual

$$\begin{array}{l}
\max_{\rho} - f_{\ell}^{*}(-\rho) + \frac{1}{2} \left(\max_{m} \|\rho\|_{m} \right)^{2} \\
\end{array}$$

$$\begin{array}{l}
\max_{\rho,r} - f_{\ell}^{*}(-\rho) + \frac{1}{2}r^{2} \\
\end{array}$$
s.t.
$$\|\rho\|_{m} \leq r \quad (\forall m)$$

- Lanckriet et al. 2004: SDP
- Bach, Lanckriet & Jordan 2004: SOCP

Existing approach 2

Upper bound based formulation:

[Sonnenburg et al. 2006, JMLR] [Rakotomamonjy et al. 2008, JMLR] [Chapelle & Rakotomamonjy 2008, NIPS workshop]

primal

$$\begin{array}{c}
\min_{\alpha} f_{\ell}(K\alpha) + \frac{1}{2} \left(\sum_{m=1}^{M} \|\alpha_{m}\|_{m} \right)^{2} \\
\end{array}$$

$$\begin{array}{c}
\min_{\alpha,d} f_{\ell}(K\alpha) + \frac{1}{2} \sum_{m=1}^{M} \frac{\alpha_{m}^{\top} K_{m} \alpha_{m}}{d_{m}} & \text{(Jensen's inequality)} \\
\text{s.t. } \sum_{m=1}^{M} d_{m} = 1, \ d_{m} \ge 0 \\
\end{array}$$

$$\begin{array}{c}
\min_{d} \min_{\alpha} \left\{ f_{\ell}(K(d)\beta) + \frac{1}{2} \beta^{\top} K(d)\beta \right\} & \theta(d) \\
\text{s.t. } \sum_{m=1}^{M} d_{m} = 1, \ d_{m} \ge 0, \ K(d) = \sum_{m=1}^{M} d_{m} K_{m}.
\end{array}$$

- Sonnenburg et al. 2006: Cutting plane (SILP)
- Rakotomamonjy et al. 2008: Gradient descent (SimpleMKL)
- Chapelle & Rakotomamonjy 2008: Newton method (HessianMKL)

Problem of existing methods

- Do not make use of sparsity during the optimization.
 - → Do not perform efficiently when
 # of kernels is large.

Outline

- Introduction
- Details of Multiple Kernel Learning
- Our method
 - Proximal minimization
 - Skipping inactive kernels
- Numerical Experiments
 - Bench-mark datasets
 - Sparse or dense

Our method SpicyMKL

- DAL (Tomioka & Sugiyama 09)
 - Dual Augmented Lagrangian
 - Lasso, group Lasso, trace norm regularization
- SpicyMKL = DAL + MKL
 - Kernelization of DAL



Relax by Proximal Minimization (Rockafellar, 1976)

Our algorithm

SpicyMKL: Proximal Minimization (Rockafellar, 1976)



Taking its dual, the update step can be solved efficiently

$$\begin{aligned} \overline{\mathsf{Fenchel's duality theorem}} \\ \min_{x} \{f(Ax) + h(x)\} &= \max_{y} \{-f^{*}(-y) - h^{*}(A^{\top}y)\} \end{aligned}$$

$$\begin{aligned} \min_{\alpha \in \mathbb{R}^{NM}} \quad f_{\ell}(K\alpha) + \sum_{m=1}^{M} g(\|\alpha_{m}\|_{m}) + \frac{1}{2\gamma_{t}} \sum_{m=1}^{M} \|\alpha_{m} - \alpha_{m}^{(t)}\|_{m}^{2} \end{aligned}$$

$$\begin{aligned} f_{\ell}(K\alpha) \quad \bigoplus_{f \in \mathbb{R}^{N}} \quad f_{\ell}(\alpha) \end{aligned}$$

$$\begin{aligned} \max_{\rho \in \mathbb{R}^{N}} \quad -f_{\ell}^{*}(-\rho) - \phi_{t}^{*}(K^{\top}\rho) \end{aligned}$$



Inner optimization

• Newton method is applicable.

$$\max_{\rho \in \mathbb{R}^{N}} -f_{\ell}^{*}(-\rho) - \phi_{t}^{*}(K^{\top}\rho)$$

Twice Differentiable Twice Differentiable (a.e.

Even if f_{ℓ}^* is not differentiable, we can apply almost the same algorithm.

Rigorous Derivation of ϕ_t^*

$$\begin{array}{c|c} \displaystyle \max_{\rho \in \mathbf{R}^N} & -f_\ell^* \left(-\rho \right) - \phi_t^* (K^\top \rho) \\ \hline \\ \hline \\ \textbf{Convolution and duality} & (h_1 \Box h_2)(z) \coloneqq \inf_{z=x+y} \{h_1(x) + h_2(y)\} \\ & (f+g)^* = f^* \Box g^* \end{array}$$

$$\begin{split} \phi_t^*(K^{\top}\rho) &= \sum_{m=1}^K (g(\cdot) + \frac{1}{2\gamma_t} \| \cdot -\alpha_m^{(t)} \|_m^2)^* (K_m \rho) \\ &= \sum_{m=1}^M \min_{u_m \in \mathbb{R}^N} \left(g^*(\|u_m\|_m) + \frac{\gamma_t}{2} \left\| u_m - \rho - \frac{\alpha_m^{(t)}}{\gamma_t} \right\|_m^2 - \frac{1}{2\gamma_t} \|\alpha_m^{(t)}\|_m^2 \right) \\ &= \sum_{m=1}^M \gamma_t \min_{u_m \in \mathbb{R}^N} \left(\frac{g^*(\|u_m\|_m)}{\gamma_t} + \frac{1}{2} \left\| u_m - \rho - \frac{\alpha_m^{(t)}}{\gamma_t} \right\|_m^2 \right) + (\text{const.}) \end{split}$$

Moreau envelope

$$\begin{split} \phi_t^* & \longrightarrow \min_{u_m \in \mathbb{R}^N} \left(\frac{g^*(\|u_m\|_m)}{\gamma_t} + \frac{1}{2} \left\| u_m - \rho - \frac{\alpha_m^{(t)}}{\gamma_t} \right\|_m^2 \right) \\ &= \min_{q_m \in \mathbb{R}} \left(\frac{g^*(q_m)}{\gamma_t} + \frac{1}{2} (q_m - \|\rho + \frac{\alpha_m^{(t)}}{\gamma_t}\|_m)^2 \right) \\ & u_m = q_m \frac{\rho + \frac{\alpha_m^{(t)}}{\gamma_t}}{\|\rho + \frac{\alpha_m^{(t)}}{\gamma_t}\|_m} \\ \hline \mathbf{Moreau envelope (Moreau 65)} \\ & \text{For a convex function } \phi, \text{ we define its Moreau envelope as} \\ & \mathcal{M}\phi(z) := \min_{x \in \mathbb{R}} \left(\frac{1}{2} (x - z)^2 + \phi(x) \right) \\ & = \mathcal{M} \frac{g^*}{\gamma_t} \left(\left\| \rho + \frac{\alpha_m^{(t)}}{\gamma_t} \right\|_m \right) \end{split}$$

26



What we have observed

We arrived at ...

Dual of update rule in proximal minimization



- Differentiable
- Only active kernels need to be calculated. (next slide)



If
$$\|\rho + \frac{\alpha_m^{(t)}}{\gamma_t}\|_m \leq C$$
 (inactive), $\mathcal{M}\frac{g^*}{\gamma_t}(\|\rho + \frac{\alpha_m^{(t)}}{\gamma_t}\|_m) = 0$,
and its gradient also vanishes.

Derivatives

• We can apply **Newton method** for the dual optimization.

M

The objective function:

$$\psi(\rho) = -f_{\ell}^*(-\rho) - \gamma_t \sum_{m=1}^{\infty} \mathcal{M} \frac{g^*}{\gamma_t} \left(\|\rho + \frac{\alpha_m^{(t)}}{\gamma_t}\|_m \right)$$

Its derivatives:

$$\begin{aligned} \nabla_{\rho}\psi(\rho) &= -\nabla_{\rho}^{2}f_{\ell}^{*}(-\rho) - \gamma_{t}\sum_{\substack{m:\text{active}}} q_{m}\frac{K_{m}p_{m}}{\|p_{m}\|_{m}} \\ \nabla_{\rho}^{2}\psi(\rho) &= -\nabla_{\rho}^{2}f_{\ell}^{*}(-\rho) - \gamma_{t}\sum_{\substack{m:\text{active}}} \left\{ (r_{m} - q_{m}\|p_{m}\|_{m})\frac{K_{m}p_{m}p_{m}^{\top}K_{m}}{\|p_{m}\|_{m}^{2}} + \gamma_{t}q_{m}\frac{K_{m}}{\|p_{m}\|_{m}} \right\} \end{aligned}$$

where
$$p_m := \|\rho + \frac{\alpha_m^{(t)}}{\gamma_t}\|_m$$
, $q_m := \operatorname{prox}_{\gamma_t g}(\|\alpha_m^{(t)} + \gamma_t \rho\|_m)$, $r_m := \frac{\operatorname{dprox}_{\gamma_t g}(x)}{\operatorname{d}x}\Big|_{x = \|\alpha_m^{(t)} + \gamma_t \rho\|_m}$

Only active kernels contributes their computations. → Efficient for large number of kernels.



Non-differentiable loss

• If loss f_{ℓ}^* is non-differentiable, almost the same algorithm is available by utilizing Moreau envelope of f_{ℓ}^* .



Dual of elastic net



Why 'Spicy' ?

- SpicyMKL = DAL + MKL.
- DAL also means a major Indian cuisine.
 - \rightarrow hot, spicy
 - → SpicyMKL



from Wikipedia

Outline

- Introduction
- Details of Multiple Kernel Learning
- Our method
 - Proximal minimization
 - Skipping inactive kernels
- Numerical Experiments
 - Bench-mark datasets
 - Sparse or dense

Numerical Experiments

IDA data set, L1 regularization.

CPU time against # of kernels



Scales well against # of kernels.

IDA data set, L1 regularization.

CPU time against # of samples



 Scaling against # of samples is almost the same as the existing methods.

UCI dataset



Comparison in elastic net (Sparse v.s. Dense)

Sparse or Dense ?

Elastic Net



Compare performances of sparse and dense solutions in a computer vision task.

caltech101 dataset (Fei-Fei et al., 2004)

- object recognition task
- anchor, ant, cannon, chair, cup



- 10 classification problems $(10=5 \times (5-1)/2)$
- # of kernels: <u>1760</u> = 4 (features) × 22 (regions) ×
 2 (kernel functions) × 10 (kernel parameters)
 - sift features [Lowe99]: colorDescriptor [van de Sande et al.10]
 hsvsift, sift (scale = 4px, 8px, automatic scale)
 - regions: whole region, 4 division, 16 division, spatial pyramid. (computed histogram of features on each region.)
 - kernel functions: Gaussian kernels, χ^2 kernels with 10 parameters.

[
[

Performances are averaged over all classification problems.



Conclusion and Future works

Conclusion

- We proposed a new MKL algorithm that is efficient when # of kernels is large.
 - proximal minimization
 - neglect 'in-active' kernels
- Medium density showed the best performance, but sparse solution also works well.

Future work

• Second order update

Technical report

T. Suzuki & R. Tomioka: SpicyMKL. arXiv: <u>http://arxiv.org/abs/0909.5026</u>

• DAL

R. Tomioka & M. Sugiyama: Dual Augmented Lagrangian Method for Efficient Sparse Reconstruction. *IEEE Signal Processing Letters*, 16 (12) pp. 1067--1070, 2009.

Thank you for your attention!

Some properties



•
$$\mathcal{M}g(z) + \mathcal{M}g^*(z) = \frac{||z||^2}{2}$$
 (Fenchel duality th.)
• $\operatorname{prox}_g(z) + \operatorname{prox}_{g^*}(z) = z$ (Fenchel duality th.)
• $\nabla_z \mathcal{M}g^*(z) = z - \operatorname{prox}_{g^*}(z) = \operatorname{prox}_g(z)$
Differentiable !
 $\nabla_z \mathcal{M}g^*(x) = \partial_z(\operatorname{obj.}) + \nabla_z \operatorname{prox}_{g^*}(z) \frac{\partial(\operatorname{obj.})}{\partial x}\Big|_{x = \operatorname{prox}_{g^*}(z)}$
 $= z - \operatorname{prox}_{g^*}(z)$

........