

# Dual Augmented Lagrangian, Proximal Minimization, and MKL

Ryota Tomioka<sup>1</sup>, Taiji Suzuki<sup>1</sup>, and Masashi Sugiyama<sup>2</sup>

<sup>1</sup>University of Tokyo

<sup>2</sup>Tokyo Institute of Technology

2009-09-15 @ TU Berlin

# Outline

## 1 Introduction

- Lasso, group lasso and MKL
- Objective

## 2 Method

- Proximal minimization algorithm
- Multiple Kernel Learning

## 3 Experiments

## 4 Summary

# Outline

## 1 Introduction

- Lasso, group lasso and MKL
- Objective

## 2 Method

- Proximal minimization algorithm
- Multiple Kernel Learning

## 3 Experiments

## 4 Summary

# Outline

- 1 Introduction
  - Lasso, group lasso and MKL
  - Objective

- 2 Method
  - Proximal minimization algorithm
  - Multiple Kernel Learning

- 3 Experiments

- 4 Summary

# Minimum norm reconstruction

- $\mathbf{w} \in \mathbb{R}^n$ : unknown.
- $y_1, y_2, \dots, y_m$ : observations,  
where  $y_i = \langle \mathbf{x}_i, \mathbf{w} \rangle + \epsilon_i$ ,  $\epsilon_i \sim \mathcal{N}(0, 1)$ .
- Recover  $\mathbf{w}$  from  $\mathbf{y}$ .

Underdetermined (when  $n > m$ )

$$\underset{\mathbf{w}}{\text{minimize}} \quad \|\mathbf{w}\|_0, \quad \text{s.t.} \quad L(\mathbf{w}) \leq C,$$

where

$$L(\mathbf{w}) = \frac{1}{2} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2,$$

and  $\|\mathbf{w}\|_0$ : the number of non-zero elements in  $\mathbf{w}$ .

Moreover, it is often good to know which features are useful.

However, this is NP hard!

# Minimum norm reconstruction

- $\mathbf{w} \in \mathbb{R}^n$ : unknown.
- $y_1, y_2, \dots, y_m$ : observations,  
where  $y_i = \langle \mathbf{x}_i, \mathbf{w} \rangle + \epsilon_i$ ,  $\epsilon_i \sim \mathcal{N}(0, 1)$ .
- Recover  $\mathbf{w}$  from  $\mathbf{y}$ .

Underdetermined (when  $n > m$ )

$$\underset{\mathbf{w}}{\text{minimize}} \quad \|\mathbf{w}\|_0, \quad \text{s.t.} \quad L(\mathbf{w}) \leq C,$$

where

$$L(\mathbf{w}) = \frac{1}{2} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2,$$

and  $\|\mathbf{w}\|_0$ : the number of non-zero elements in  $\mathbf{w}$ .

Moreover, it is often good to know which features are useful.

However, this is NP hard!

# Minimum norm reconstruction

- $\mathbf{w} \in \mathbb{R}^n$ : unknown.
- $y_1, y_2, \dots, y_m$ : observations,  
where  $y_i = \langle \mathbf{x}_i, \mathbf{w} \rangle + \epsilon_i$ ,  $\epsilon_i \sim \mathcal{N}(0, 1)$ .
- Recover  $\mathbf{w}$  from  $\mathbf{y}$ .

Underdetermined (when  $n > m$ )

$$\underset{\mathbf{w}}{\text{minimize}} \quad \|\mathbf{w}\|_0, \quad \text{s.t.} \quad L(\mathbf{w}) \leq C,$$

where

$$L(\mathbf{w}) = \frac{1}{2} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2,$$

and  $\|\mathbf{w}\|_0$ : the number of non-zero elements in  $\mathbf{w}$ .

Moreover, it is often good to know which features are useful.

However, this is NP hard!

# Minimum norm reconstruction

- $\mathbf{w} \in \mathbb{R}^n$ : unknown.
- $y_1, y_2, \dots, y_m$ : observations,  
where  $y_i = \langle \mathbf{x}_i, \mathbf{w} \rangle + \epsilon_i$ ,  $\epsilon_i \sim \mathcal{N}(0, 1)$ .
- Recover  $\mathbf{w}$  from  $\mathbf{y}$ .

Underdetermined (when  $n > m$ )

$$\underset{\mathbf{w}}{\text{minimize}} \quad \|\mathbf{w}\|_0, \quad \text{s.t.} \quad L(\mathbf{w}) \leq C,$$

where

$$L(\mathbf{w}) = \frac{1}{2} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2,$$

and  $\|\mathbf{w}\|_0$ : the number of non-zero elements in  $\mathbf{w}$ .

Moreover, it is often good to know which features are useful.

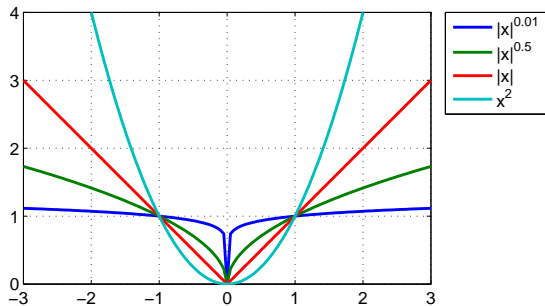
**However, this is NP hard!**



# Convex relaxation

- $p$ -norm like functions

$$\|\mathbf{w}\|_p^p = \sum_{j=1}^n |w_j|^p : \begin{cases} \text{If } p \geq 1 & \text{convex} \\ \text{If } p < 1 & \text{non-convex} \end{cases}$$

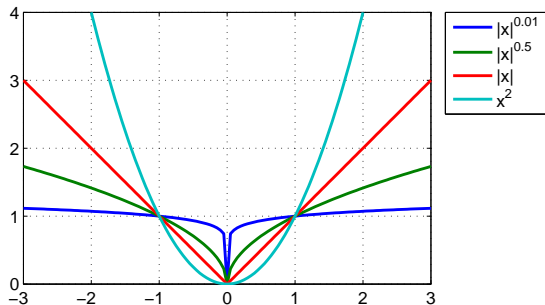


$\|\cdot\|_1$ -regularization is the closest to  $\|\cdot\|_0$  within convex norm-like functions

# Convex relaxation

- $p$ -norm like functions

$$\|\mathbf{w}\|_p^p = \sum_{j=1}^n |w_j|^p : \begin{cases} \text{If } p \geq 1 & \text{convex} \\ \text{If } p < 1 & \text{non-convex} \end{cases}$$



$\|\cdot\|_1$ -regularization is the closest to  $\|\cdot\|_0$  within convex norm-like functions

# Lasso regression

- Problem 1:

$$\text{minimize } \|\mathbf{w}\|_1, \quad \text{s.t. } L(\mathbf{w}) \leq C.$$

- Problem 2:

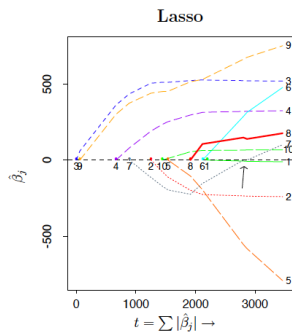
$$\text{minimize } L(\mathbf{w}), \quad \text{s.t. } \|\mathbf{w}\|_1 \leq C'.$$

- Problem 3:

$$\text{minimize } L(\mathbf{w}) + \lambda \|\mathbf{w}\|_1.$$

Note:

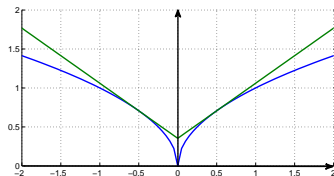
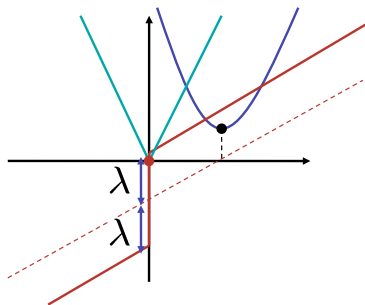
- Above three problems are equivalent to each other.
- Monotone operation preserves equivalence.
- We focus on the third problem.



[From Efron et al. (2003)]

# Why $\ell_1$ -regularization?

- The closest to  $\|\cdot\|_0$  within convex norm-like functions.
- Non-differentiable at the origin (truncation with finite  $\lambda$ ).
- Non-convex regularizers ( $p < 1$ )  
→ Iteratively solve (weighted)  $\ell_1$ -regularization.
- Bayesian sparse models (type-II ML)  
→ Iteratively solve (weighted)  $\ell_1$ -regularization (in special cases) .  
(Wipf&Nagarajan, 08)

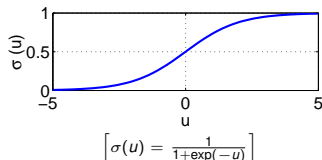


# Generalizations

- Generalize the loss term... e.g.,  $\ell_1$ -logistic regression

$$\underset{\mathbf{w} \in \mathbb{R}^n}{\text{minimize}} \quad \sum_{i=1}^m -\log P(y_i | \mathbf{x}_i; \mathbf{w}) + \lambda \|\mathbf{w}\|_1$$

$$\text{where } P(y | \mathbf{x}; \mathbf{w}) = \sigma(y \langle \mathbf{w}, \mathbf{x} \rangle) \\ (y \in \{-1, +1\})$$



- Generalize the reg. term... e.g., group lasso (Yuan&Lin,06)

$$\underset{\mathbf{w} \in \mathbb{R}^n}{\text{minimize}} \quad L(\mathbf{w}) + \lambda \sum_{g \in \mathcal{G}} \|\mathbf{w}_g\|_2$$

where,  $\mathcal{G}$  is a partition of  $\{1, \dots, n\}$ ,  $\mathbf{w} = \begin{pmatrix} (\mathbf{w}_{g_1}) \\ (\mathbf{w}_{g_2}) \\ \vdots \\ (\mathbf{w}_{g_q}) \end{pmatrix}$ ,  $q = |\mathcal{G}|$ .

# Introducing Kernels

Multiple Kernel Learning (MKL) (Lanckriet, Bach, et al., 04)

Let  $\mathcal{H}_1, \mathcal{H}_2, \dots, \mathcal{H}_n$  be RKHSs and  $K_1, K_2, \dots, K_n$  be the kernel

functions. Use functions  $f = \underbrace{f_1}_{\in \mathcal{H}_1} + \underbrace{f_2}_{\in \mathcal{H}_2} + \dots + \underbrace{f_n}_{\in \mathcal{H}_n}$

$$\underset{f_j \in \mathcal{H}_j, b \in \mathbb{R}}{\text{minimize}} \quad L(f_1 + f_2 + \dots + f_n + b) + \lambda \sum_{j=1}^n \|f_j\|_{\mathcal{H}_j}$$

↓ representer theorem

$$\underset{\alpha_j \in \mathbb{R}^m, b \in \mathbb{R}}{\text{minimize}} \quad \ell \left( \sum_{j=1}^n \mathbf{K}_j \alpha_j + b \mathbf{1} \right) + \lambda \sum_{j=1}^n \|\alpha_j\|_{\mathbf{K}_j}$$

where,  $\|\alpha_j\|_{\mathbf{K}_j} = \sqrt{\alpha_j^\top \mathbf{K}_j \alpha_j}$ .

... nothing but a kernel-weighted group lasso

# Introducing Kernels

Multiple Kernel Learning (MKL) (Lanckriet, Bach, et al., 04)

Let  $\mathcal{H}_1, \mathcal{H}_2, \dots, \mathcal{H}_n$  be RKHSs and  $K_1, K_2, \dots, K_n$  be the kernel

functions. Use functions  $f = \underbrace{f_1}_{\in \mathcal{H}_1} + \underbrace{f_2}_{\in \mathcal{H}_2} + \dots + \underbrace{f_n}_{\in \mathcal{H}_n}$

$$\underset{f_j \in \mathcal{H}_j, b \in \mathbb{R}}{\text{minimize}} \quad L(f_1 + f_2 + \dots + f_n + b) + \lambda \sum_{j=1}^n \|f_j\|_{\mathcal{H}_j}$$

↓ representer theorem

$$\underset{\alpha_j \in \mathbb{R}^m, b \in \mathbb{R}}{\text{minimize}} \quad f_\ell \left( \sum_{j=1}^n \mathbf{K}_j \alpha_j + b \mathbf{1} \right) + \lambda \sum_{j=1}^n \|\alpha_j\|_{\mathbf{K}_j}$$

where,  $\|\alpha_j\|_{\mathbf{K}_j} = \sqrt{\alpha_j^\top \mathbf{K}_j \alpha_j}$ .

... nothing but a kernel-weighted group lasso

# Modeling assumptions

In many cases the loss term  $L(\cdot)$  can be decomposed into a loss function  $f_\ell$  and a design matrix  $\mathbf{A}$ .

- Squared loss

$$f_\ell^Q(\mathbf{z}) = \frac{1}{2} \|\mathbf{y} - \mathbf{z}\|^2, \quad \mathbf{A} = \begin{pmatrix} \mathbf{x}_1^\top \\ \vdots \\ \mathbf{x}_m^\top \end{pmatrix}$$

$$\Rightarrow f_\ell^Q(\mathbf{A}\mathbf{w}) = \frac{1}{2} \sum_{i=1}^m (y_i - \langle \mathbf{x}_i, \mathbf{w} \rangle)^2$$

- Logistic loss

$$f_\ell^L(\mathbf{z}) = \sum_{i=1}^m \log(1 + \exp(-y_i z_i)), \quad \mathbf{A} = \begin{pmatrix} \mathbf{x}_1^\top \\ \vdots \\ \mathbf{x}_m^\top \end{pmatrix}$$

$$\Rightarrow f_\ell^L(\mathbf{A}\mathbf{w}) = \sum_{i=1}^m -\log \sigma(y_i \langle \mathbf{w}, \mathbf{x}_i \rangle)$$



# Outline

## 1 Introduction

- Lasso, group lasso and MKL
- **Objective**

## 2 Method

- Proximal minimization algorithm
- Multiple Kernel Learning

## 3 Experiments

## 4 Summary

# Objective

Develop an optimization algorithm for the problem:

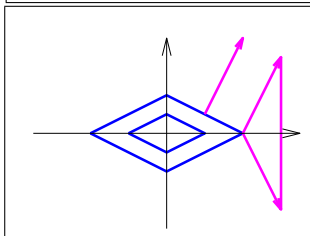
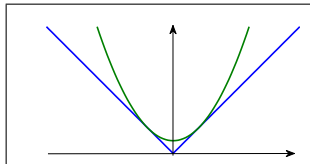
$$\underset{\mathbf{w} \in \mathbb{R}^n}{\text{minimize}} \quad f_\ell(\mathbf{A}\mathbf{w}) + \phi_\lambda(\mathbf{w}).$$

- $\mathbf{A} \in \mathbb{R}^{m \times n}$  ( $m$ : #observations,  $n$ : #unknowns) .
- $f_\ell$  is convex and twice differentiable.
- $\phi_\lambda(\mathbf{w})$  is convex but possibly non-differentiable, e.g.,  
 $\phi_\lambda(\mathbf{w}) = \lambda \|\mathbf{w}\|_1$ .
- $\eta\phi_\lambda = \phi_{\eta\lambda}$ .
- We are interested in algorithms for general  $f_\ell$  ( $\leftrightarrow$  LARS).

# Where does the difficulty come from?

*Conventional view:* the **non-differentiability** of  $\phi_\lambda(\mathbf{w})$

- Upper bound the regularizer from above with a differentiable function.
  - FOCUSS (Rao & Kreutz-Delgado, 99)
  - Majorization-Minimization (Figueiredo et al., 07)
- Explicitly handle the non-differentiability.
  - Sub-gradient L-BFGS (Andrew & Gao, 07; Yu et al., 08)



*Our view:* the **coupling between variables introduced by  $\mathbf{A}$** .

# Where does the difficulty come from?

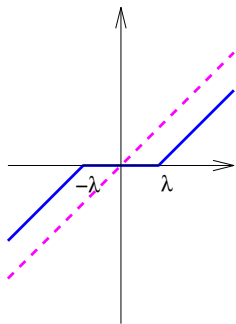
*Our view:* the coupling between variables introduced by  $\mathbf{A}$ .

In fact, when  $\mathbf{A} = \mathbf{I}_n$

$$\min_{\mathbf{w} \in \mathbb{R}^n} \left( \frac{1}{2} \|\mathbf{y} - \mathbf{w}\|_2^2 + \lambda \|\mathbf{w}\|_1 \right) = \sum_{j=1}^n \min_{w_j \in \mathbb{R}} \left( \frac{1}{2} (y_j - w_j)^2 + \lambda |w_j| \right).$$

$$\begin{aligned} \Rightarrow w_j^* &= \text{ST}_\lambda(y_j) \\ &= \begin{cases} y_j - \lambda & (\lambda \leq y_j), \\ 0 & (-\lambda \leq y_j \leq \lambda), \\ y_j + \lambda & (y_j \leq -\lambda). \end{cases} \end{aligned}$$

min is obtained analytically!



We focus on  $\phi_\lambda$  for which the above min can be obtained analytically

# Earlier study

Iterative Shrinkage/Thresholding (IST) (Figueiredo&Nowak, 03; Daubechies et al., 04;...):

## Algorithm

- 1 Choose an initial solution  $\mathbf{w}^0$ .
- 2 Repeat until some stopping criterion is satisfied:

$$\mathbf{w}^{t+1} \leftarrow \underset{\mathbf{w} \in \mathbb{R}^n}{\operatorname{argmin}} (Q_{\eta_t}(\mathbf{w}; \mathbf{w}^t) + \phi_{\lambda}(\mathbf{w}))$$

where

$$Q_{\eta}(\mathbf{w}; \mathbf{w}^t) = \underbrace{L(\mathbf{w}^t) + \nabla L^{\top}(\mathbf{w}^t)(\mathbf{w} - \mathbf{w}^t)}_{(1) \text{ Linearly approximate the loss term.}} + \frac{1}{2\eta} \underbrace{\|\mathbf{w} - \mathbf{w}^t\|_2^2}_{(2) \text{ penalize dist}^2 \text{ from the last iterate.}}.$$

Note: minimizing  $Q_{\eta}(\mathbf{w}; \mathbf{w}^t)$  gives the **ordinary gradient step**.

# Earlier study: IST

$$\operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^n} (Q_{\eta_t}(\mathbf{w}; \mathbf{w}^t) + \phi_\lambda(\mathbf{w}))$$

$$= \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^n} \left( \text{const.} + \nabla L^\top(\mathbf{w}^t)(\mathbf{w} - \mathbf{w}^t) + \frac{1}{2\eta_t} \|\mathbf{w} - \mathbf{w}^t\|_2^2 + \phi_\lambda(\mathbf{w}) \right)$$

$$= \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^n} \left( \frac{1}{2\eta_t} \|\mathbf{w} - \tilde{\mathbf{w}}^t\|_2^2 + \phi_\lambda(\mathbf{w}) \right) =: \operatorname{ST}_{\eta_t \lambda}(\tilde{\mathbf{w}}^t)$$

assume this min can be obtained analytically

where  $\tilde{\mathbf{w}}^t = \mathbf{w}^t - \eta_t \nabla L(\mathbf{w}^t)$  (gradient step)

Finally,

$$\mathbf{w}^{t+1} \leftarrow \underbrace{\operatorname{ST}_{\eta_t \lambda}}_{\text{shrink}} \left( \underbrace{\mathbf{w}^t - \eta_t \nabla L(\mathbf{w}^t)}_{\text{gradient step}} \right)$$

- Pro : easy to implement.
- Con : bad for poorly conditioned  $\mathbf{A}$ .

# Earlier study: IST

$$\begin{aligned}
 & \underset{\mathbf{w} \in \mathbb{R}^n}{\operatorname{argmin}} \left( Q_{\eta_t}(\mathbf{w}; \mathbf{w}^t) + \phi_\lambda(\mathbf{w}) \right) \\
 &= \underset{\mathbf{w} \in \mathbb{R}^n}{\operatorname{argmin}} \left( \text{const.} + \nabla L^\top(\mathbf{w}^t)(\mathbf{w} - \mathbf{w}^t) + \frac{1}{2\eta_t} \|\mathbf{w} - \mathbf{w}^t\|_2^2 + \phi_\lambda(\mathbf{w}) \right) \\
 &= \underset{\mathbf{w} \in \mathbb{R}^n}{\operatorname{argmin}} \left( \frac{1}{2\eta_t} \|\mathbf{w} - \tilde{\mathbf{w}}^t\|_2^2 + \phi_\lambda(\mathbf{w}) \right) =: \text{ST}_{\eta_t\lambda}(\tilde{\mathbf{w}}^t)
 \end{aligned}$$

assume this min can be obtained analytically

where  $\tilde{\mathbf{w}}^t = \mathbf{w}^t - \eta_t \nabla L(\mathbf{w}^t)$  (gradient step)

Finally,

$$\mathbf{w}^{t+1} \leftarrow \underbrace{\text{ST}_{\eta_t\lambda}}_{\text{shrink}} \left( \underbrace{\mathbf{w}^t - \eta_t \nabla L(\mathbf{w}^t)}_{\text{gradient step}} \right)$$

- Pro : easy to implement.
- Con : bad for poorly conditioned  $\mathbf{A}$ .

# Earlier study: IST

$$\operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^n} (Q_{\eta_t}(\mathbf{w}; \mathbf{w}^t) + \phi_\lambda(\mathbf{w}))$$

$$= \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^n} \left( \text{const.} + \nabla L^\top(\mathbf{w}^t)(\mathbf{w} - \mathbf{w}^t) + \frac{1}{2\eta_t} \|\mathbf{w} - \mathbf{w}^t\|_2^2 + \phi_\lambda(\mathbf{w}) \right)$$

$$= \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^n} \left( \frac{1}{2\eta_t} \|\mathbf{w} - \tilde{\mathbf{w}}^t\|_2^2 + \phi_\lambda(\mathbf{w}) \right) =: \operatorname{ST}_{\eta_t \lambda}(\tilde{\mathbf{w}}^t)$$

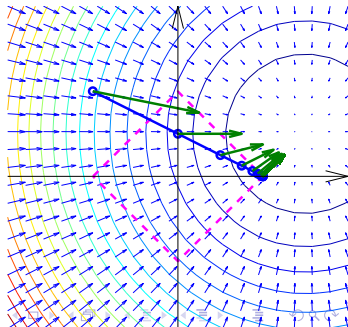
assume this min can be obtained analytically

where  $\tilde{\mathbf{w}}^t = \mathbf{w}^t - \eta_t \nabla L(\mathbf{w}^t)$  (gradient step)

Finally,

$$\mathbf{w}^{t+1} \leftarrow \underbrace{\operatorname{ST}_{\eta_t \lambda}}_{\text{shrink}} \left( \underbrace{\mathbf{w}^t - \eta_t \nabla L(\mathbf{w}^t)}_{\text{gradient step}} \right)$$

- **Pro** : easy to implement.
- **Con** : bad for poorly conditioned  $\mathbf{A}$ .





# Earlier study: IST

$$\operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^n} (Q_{\eta_t}(\mathbf{w}; \mathbf{w}^t) + \phi_\lambda(\mathbf{w}))$$

$$= \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^n} \left( \text{const.} + \nabla L^\top(\mathbf{w}^t)(\mathbf{w} - \mathbf{w}^t) + \frac{1}{2\eta_t} \|\mathbf{w} - \mathbf{w}^t\|_2^2 + \phi_\lambda(\mathbf{w}) \right)$$

$$= \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^n} \left( \frac{1}{2\eta_t} \|\mathbf{w} - \tilde{\mathbf{w}}^t\|_2^2 + \phi_\lambda(\mathbf{w}) \right) =: \operatorname{ST}_{\eta_t \lambda}(\tilde{\mathbf{w}}^t)$$

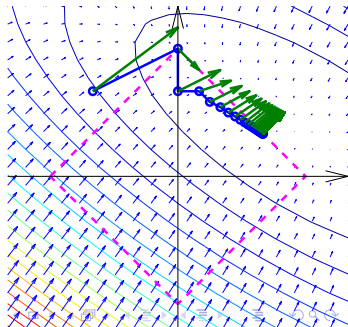
assume this min can be obtained analytically

where  $\tilde{\mathbf{w}}^t = \mathbf{w}^t - \eta_t \nabla L(\mathbf{w}^t)$  (gradient step)

Finally,

$$\mathbf{w}^{t+1} \leftarrow \underbrace{\operatorname{ST}_{\eta_t \lambda}}_{\text{shrink}} \left( \underbrace{\mathbf{w}^t - \eta_t \nabla L(\mathbf{w}^t)}_{\text{gradient step}} \right)$$

- **Pro** : easy to implement.
- **Con** : bad for poorly conditioned  $\mathbf{A}$ .



# Summary so far

We want to solve:

$$\underset{\mathbf{w} \in \mathbb{R}^n}{\text{minimize}} \quad f_\ell(\mathbf{A}\mathbf{w}) + \phi_\lambda(\mathbf{w}).$$

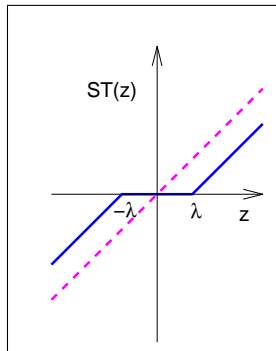
- $f_\ell$  is convex and twice differentiable.
- $\phi_\lambda(\mathbf{w})$  is a convex function for which the minimization:

$$\text{ST}_\lambda(\mathbf{z}) = \underset{\mathbf{w} \in \mathbb{R}^n}{\text{argmin}} \left( \frac{1}{2} \|\mathbf{w} - \mathbf{z}\|_2^2 + \phi_\lambda(\mathbf{w}) \right)$$

can be carried out analytically, e.g.,

$$\phi_\lambda(\mathbf{w}) = \lambda \|\mathbf{w}\|_1.$$

- **Exploit the non-differentiability of  $\phi_\lambda$ :**  
more sparsity  $\rightarrow$  more efficiency.
- Robustify against poor conditioning of  $\mathbf{A}$ .



# Outline

- 1 Introduction
  - Lasso, group lasso and MKL
  - Objective
- 2 Method**
  - Proximal minimization algorithm
  - Multiple Kernel Learning
- 3 Experiments
- 4 Summary

# Outline

- 1 Introduction
  - Lasso, group lasso and MKL
  - Objective
- 2 **Method**
  - **Proximal minimization algorithm**
  - Multiple Kernel Learning
- 3 Experiments
- 4 Summary

## Proximal Minimization (Rockafellar, 1976)

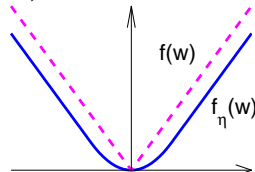
- 1 Choose an initial solution  $\mathbf{w}^0$ .
- 2 Repeat until some stopping criterion is satisfied:

$$\mathbf{w}^{t+1} \leftarrow \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^n} \left( \underbrace{f_\ell(\mathbf{A}\mathbf{w})}_{\text{No approximation}} + \phi_\lambda(\mathbf{w}) + \frac{1}{2\eta_t} \underbrace{\|\mathbf{w} - \mathbf{w}^t\|_2^2}_{\text{penalize dist}^2 \text{ from the last iterate.}} \right)$$

- Let

$$f_\eta(\mathbf{w}) = \min_{\tilde{\mathbf{w}} \in \mathbb{R}^n} \left( f_\ell(\mathbf{A}\tilde{\mathbf{w}}) + \phi_\lambda(\tilde{\mathbf{w}}) + \frac{1}{2\eta} \|\tilde{\mathbf{w}} - \mathbf{w}\|_2^2 \right).$$

- Fact 1:  $f_\eta(\mathbf{w}) \leq f(\mathbf{w}) = f_\ell(\mathbf{A}\mathbf{w}) + \phi_\lambda(\mathbf{w})$ .
- Fact 2:  $f_\eta(\mathbf{w}^*) = f(\mathbf{w}^*)$ .
- Linearly approximate the loss term  $\rightarrow$  IST



# The difference

- IST: linearly **approximates** the loss term:

$$f_\ell(\mathbf{A}\mathbf{w}) \simeq f_\ell(\mathbf{A}\mathbf{w}^t) + (\mathbf{w} - \mathbf{w}^t)^\top \mathbf{A}^\top \nabla f_\ell(\mathbf{A}\mathbf{w}^t)$$

→ tightest at the **current iterate**  $\mathbf{w}^t$

- DAL (proposed): linearly **lower bounds** the loss term:

$$f_\ell(\mathbf{A}\mathbf{w}) = \max_{\alpha \in \mathbb{R}^m} \left( -f_\ell^*(-\alpha) - \mathbf{w}^\top \mathbf{A}^\top \alpha \right)$$

→ tightest at the **next iterate**  $\mathbf{w}^{t+1}$

# The algorithm

## IST (Earlier study)

- 1 Choose an initial solution  $\mathbf{w}^0$ .
- 2 Repeat until some stopping criterion is satisfied:

$$\mathbf{w}^{t+1} \leftarrow \text{ST}_{\eta_t \lambda} \left( \mathbf{w}^t + \eta_t \mathbf{A}^\top (-\nabla f_\ell(\mathbf{A}\mathbf{w}^t)) \right)$$

## Dual Augmented Lagrangian (proposed)

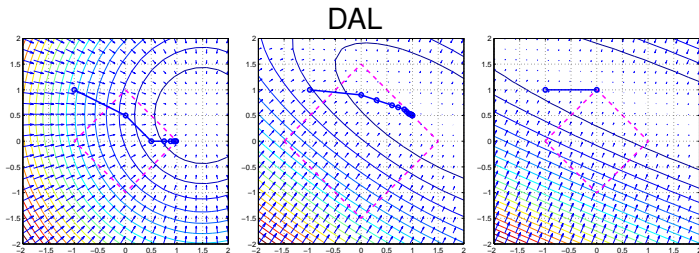
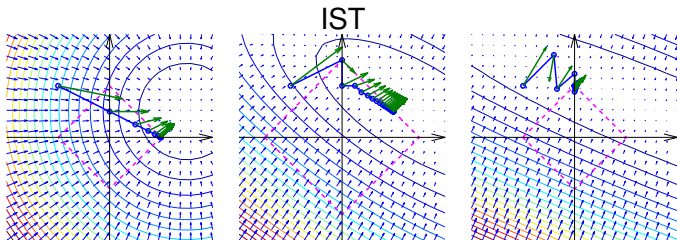
- 1 Choose an initial solution  $\mathbf{w}^0$  and a sequence  $\eta_0 \leq \eta_1 \leq \dots$ .
- 2 Repeat until some stopping criterion is satisfied:

$$\mathbf{w}^{t+1} \leftarrow \text{ST}_{\eta_t \lambda} \left( \mathbf{w}^t + \eta_t \mathbf{A}^\top \boldsymbol{\alpha}^t \right)$$

where

$$\boldsymbol{\alpha}^t = \underset{\boldsymbol{\alpha} \in \mathbb{R}^m}{\text{argmin}} \left( f_\ell^*(-\boldsymbol{\alpha}) + \frac{1}{2\eta_t} \|\text{ST}_{\eta_t \lambda}(\mathbf{w}^t + \eta_t \mathbf{A}^\top \boldsymbol{\alpha})\|_2^2 \right)$$

# Numerical examples





# Derivation

$$\begin{aligned}
 \mathbf{w}^{t+1} &\leftarrow \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^n} \left( f_{\ell}(\mathbf{A}\mathbf{w}) + \phi_{\lambda}(\mathbf{w}) + \frac{1}{2\eta_t} \|\mathbf{w} - \mathbf{w}^t\|_2^2 \right) \\
 &= \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^n} \left\{ \max_{\boldsymbol{\alpha} \in \mathbb{R}^m} \left( -f_{\ell}^*(-\boldsymbol{\alpha}) - \mathbf{w}^{\top} \mathbf{A}^{\top} \boldsymbol{\alpha} \right) + \phi_{\lambda}(\mathbf{w}) + \frac{1}{2\eta_t} \|\mathbf{w} - \mathbf{w}^t\|_2^2 \right\}
 \end{aligned}$$

Exchange the order of min and max, calculation, and calculation...

# Derivation

$$\begin{aligned}\mathbf{w}^{t+1} &\leftarrow \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^n} \left( f_\ell(\mathbf{A}\mathbf{w}) + \phi_\lambda(\mathbf{w}) + \frac{1}{2\eta_t} \|\mathbf{w} - \mathbf{w}^t\|_2^2 \right) \\ &= \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^n} \left\{ \max_{\boldsymbol{\alpha} \in \mathbb{R}^m} \left( -f_\ell^*(-\boldsymbol{\alpha}) - \mathbf{w}^\top \mathbf{A}^\top \boldsymbol{\alpha} \right) + \phi_\lambda(\mathbf{w}) + \frac{1}{2\eta_t} \|\mathbf{w} - \mathbf{w}^t\|_2^2 \right\}\end{aligned}$$

Exchange the order of min and max, calculation, and calculation...

# Augmented Lagrangian

Equality constrained problem:

$$\begin{aligned} \underset{\mathbf{x}}{\text{minimize}} \quad & f(\mathbf{x}), & \Leftrightarrow \quad & \underset{\mathbf{x}}{\text{minimize}} \quad L_{hard}(\mathbf{x}) = \begin{cases} f(\mathbf{x}) & (\text{if } \mathbf{c}(\mathbf{x}) = 0), \\ +\infty & (\text{otherwise}). \end{cases} \\ \text{s.t.} \quad & \mathbf{c}(\mathbf{x}) = 0. \end{aligned}$$

Ordinary Lagrangian:

$$L(\mathbf{x}, \mathbf{y}) = f(\mathbf{x}) + \mathbf{y}^\top \mathbf{c}(\mathbf{x})$$

Augmented Lagrangian:

$$L_\eta(\mathbf{x}, \mathbf{y}) = f(\mathbf{x}) + \mathbf{y}^\top \mathbf{c}(\mathbf{x}) + \frac{\eta}{2} \|\mathbf{c}(\mathbf{x})\|_2^2$$

$$L(\mathbf{x}, \mathbf{y}) \leq L_\eta(\mathbf{x}, \mathbf{y}) \leq L_{hard}(\mathbf{x})$$

$$\downarrow \min_{\mathbf{x}}$$

$$d(\mathbf{y}) \leq d_\eta(\mathbf{y}) \leq f(\mathbf{x}^*): \text{ primal optimum}$$

# Augmented Lagrangian

Equality constrained problem:

$$\begin{aligned} \underset{\mathbf{x}}{\text{minimize}} \quad & f(\mathbf{x}), & \Leftrightarrow \quad & \underset{\mathbf{x}}{\text{minimize}} \quad L_{hard}(\mathbf{x}) = \begin{cases} f(\mathbf{x}) & (\text{if } \mathbf{c}(\mathbf{x}) = 0), \\ +\infty & (\text{otherwise}). \end{cases} \\ \text{s.t.} \quad & \mathbf{c}(\mathbf{x}) = 0. \end{aligned}$$

Ordinary Lagrangian:

$$L(\mathbf{x}, \mathbf{y}) = f(\mathbf{x}) + \mathbf{y}^\top \mathbf{c}(\mathbf{x})$$

Augmented Lagrangian:

$$L_\eta(\mathbf{x}, \mathbf{y}) = f(\mathbf{x}) + \mathbf{y}^\top \mathbf{c}(\mathbf{x}) + \frac{\eta}{2} \|\mathbf{c}(\mathbf{x})\|_2^2$$

$$L(\mathbf{x}, \mathbf{y}) \leq L_\eta(\mathbf{x}, \mathbf{y}) \leq L_{hard}(\mathbf{x})$$

$$\downarrow \min_{\mathbf{x}}$$

$$d(\mathbf{y}) \leq d_\eta(\mathbf{y}) \leq f(\mathbf{x}^*): \text{ primal optimum}$$

# Augmented Lagrangian

Equality constrained problem:

$$\begin{aligned} \underset{\mathbf{x}}{\text{minimize}} \quad & f(\mathbf{x}), & \Leftrightarrow \quad & \underset{\mathbf{x}}{\text{minimize}} \quad L_{hard}(\mathbf{x}) = \begin{cases} f(\mathbf{x}) & (\text{if } \mathbf{c}(\mathbf{x}) = 0), \\ +\infty & (\text{otherwise}). \end{cases} \\ \text{s.t.} \quad & \mathbf{c}(\mathbf{x}) = 0. \end{aligned}$$

Ordinary Lagrangian:

$$L(\mathbf{x}, \mathbf{y}) = f(\mathbf{x}) + \mathbf{y}^\top \mathbf{c}(\mathbf{x})$$

Augmented Lagrangian:

$$L_\eta(\mathbf{x}, \mathbf{y}) = f(\mathbf{x}) + \mathbf{y}^\top \mathbf{c}(\mathbf{x}) + \frac{\eta}{2} \|\mathbf{c}(\mathbf{x})\|_2^2$$

$$L(\mathbf{x}, \mathbf{y}) \leq L_\eta(\mathbf{x}, \mathbf{y}) \leq L_{hard}(\mathbf{x})$$

$$\downarrow \min_{\mathbf{x}}$$

$$d(\mathbf{y}) \leq d_\eta(\mathbf{y}) \leq f(\mathbf{x}^*): \text{ primal optimum}$$

# Augmented Lagrangian

Equality constrained problem:

$$\begin{aligned} \underset{\mathbf{x}}{\text{minimize}} \quad & f(\mathbf{x}), & \Leftrightarrow \quad & \underset{\mathbf{x}}{\text{minimize}} \quad L_{hard}(\mathbf{x}) = \begin{cases} f(\mathbf{x}) & (\text{if } \mathbf{c}(\mathbf{x}) = 0), \\ +\infty & (\text{otherwise}). \end{cases} \\ \text{s.t.} \quad & \mathbf{c}(\mathbf{x}) = 0. \end{aligned}$$

Ordinary Lagrangian:

$$L(\mathbf{x}, \mathbf{y}) = f(\mathbf{x}) + \mathbf{y}^\top \mathbf{c}(\mathbf{x})$$

Augmented Lagrangian:

$$L_\eta(\mathbf{x}, \mathbf{y}) = f(\mathbf{x}) + \mathbf{y}^\top \mathbf{c}(\mathbf{x}) + \frac{\eta}{2} \|\mathbf{c}(\mathbf{x})\|_2^2$$

$$L(\mathbf{x}, \mathbf{y}) \leq L_\eta(\mathbf{x}, \mathbf{y}) \leq L_{hard}(\mathbf{x})$$

$$\downarrow \min_{\mathbf{x}}$$

$$d(\mathbf{y}) \leq d_\eta(\mathbf{y}) \leq f(\mathbf{x}^*): \text{ primal optimum}$$

# Augmented Lagrangian

Equality constrained problem:

$$\begin{aligned} \underset{\mathbf{x}}{\text{minimize}} \quad & f(\mathbf{x}), & \Leftrightarrow \quad & \underset{\mathbf{x}}{\text{minimize}} \quad L_{hard}(\mathbf{x}) = \begin{cases} f(\mathbf{x}) & (\text{if } \mathbf{c}(\mathbf{x}) = 0), \\ +\infty & (\text{otherwise}). \end{cases} \\ \text{s.t.} \quad & \mathbf{c}(\mathbf{x}) = 0. \end{aligned}$$

Ordinary Lagrangian:

$$L(\mathbf{x}, \mathbf{y}) = f(\mathbf{x}) + \mathbf{y}^\top \mathbf{c}(\mathbf{x})$$

Augmented Lagrangian:

$$L_\eta(\mathbf{x}, \mathbf{y}) = f(\mathbf{x}) + \mathbf{y}^\top \mathbf{c}(\mathbf{x}) + \frac{\eta}{2} \|\mathbf{c}(\mathbf{x})\|_2^2$$

$$L(\mathbf{x}, \mathbf{y}) \leq L_\eta(\mathbf{x}, \mathbf{y}) \leq L_{hard}(\mathbf{x})$$

$\downarrow \min_{\mathbf{x}}$

$$d(\mathbf{y}) \leq d_\eta(\mathbf{y}) \leq f(\mathbf{x}^*): \text{ primal optimum}$$

# Augmented Lagrangian

Equality constrained problem:

$$\begin{aligned} \underset{\mathbf{x}}{\text{minimize}} \quad & f(\mathbf{x}), & \Leftrightarrow \quad & \underset{\mathbf{x}}{\text{minimize}} \quad L_{hard}(\mathbf{x}) = \begin{cases} f(\mathbf{x}) & (\text{if } \mathbf{c}(\mathbf{x}) = 0), \\ +\infty & (\text{otherwise}). \end{cases} \\ \text{s.t.} \quad & \mathbf{c}(\mathbf{x}) = 0. \end{aligned}$$

Ordinary Lagrangian:

$$L(\mathbf{x}, \mathbf{y}) = f(\mathbf{x}) + \mathbf{y}^\top \mathbf{c}(\mathbf{x})$$

Augmented Lagrangian:

$$L_\eta(\mathbf{x}, \mathbf{y}) = f(\mathbf{x}) + \mathbf{y}^\top \mathbf{c}(\mathbf{x}) + \frac{\eta}{2} \|\mathbf{c}(\mathbf{x})\|_2^2$$

$$L(\mathbf{x}, \mathbf{y}) \leq L_\eta(\mathbf{x}, \mathbf{y}) \leq L_{hard}(\mathbf{x})$$

$$\downarrow \min_{\mathbf{x}}$$

$$d(\mathbf{y}) \leq d_\eta(\mathbf{y}) \leq f(\mathbf{x}^*): \text{ primal optimum}$$



## Augmented Lagrangian Algorithm (Hestenes, 69; Powell, 69)

- 1 Choose an initial multiplier  $\mathbf{y}^0$  and a sequence  $\eta_0 \leq \eta_1 \leq \dots$ .
- 2 Update the multiplier:

$$\mathbf{y}^{t+1} \leftarrow \mathbf{y}^t + \eta_t \mathbf{c}(\mathbf{x}^t)$$

where

$$\mathbf{x}^t = \operatorname{argmin}_{\mathbf{x}} L_{\eta_t}(\mathbf{x}, \mathbf{y}^t)$$

Note

- The multiplier  $\mathbf{y}^t$  is updated as long as the constraint is violated.
- AL method  $\Leftrightarrow$  proximal minimization in the dual (Rockafellar, 76).

$$\mathbf{y}^{t+1} \leftarrow \operatorname{argmax}_{\mathbf{y}} \left( \underbrace{d(\mathbf{y})}_{=\min_{\mathbf{x}}(f(\mathbf{x}) + \mathbf{y}^T \mathbf{c}(\mathbf{x}))} - \frac{1}{2\eta_t} \|\mathbf{y} - \mathbf{y}^t\|_2^2 \right)$$

# Outline

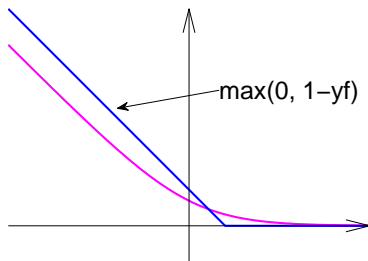
- 1 Introduction
  - Lasso, group lasso and MKL
  - Objective
- 2 **Method**
  - Proximal minimization algorithm
  - **Multiple Kernel Learning**
- 3 Experiments
- 4 Summary

# Objective

Learn a linear combination of kernels from data.

$$\underset{f \in \mathcal{H}, b \in \mathbb{R}, \mathbf{d} \in \mathbb{R}^n}{\text{minimize}} \quad \sum_{i=1}^m \ell_i(f(x_i) + b) + \frac{\lambda}{2} \|f\|_{\mathcal{H}(d)}^2$$

$$\text{s.t.} \quad \mathbf{K}(d) = \sum_{i=1}^n d_j \mathbf{K}_j, \quad d_j \geq 0, \quad \sum_j d_j \leq 1.$$



# Representer theorem

$$\begin{aligned} & \underset{\alpha \in \mathbb{R}^m, b \in \mathbb{R}, \mathbf{d} \in \mathbb{R}^n}{\text{minimize}} && L(\mathbf{K}(d)\alpha + b\mathbf{1}) + \frac{\lambda}{2} \alpha^\top \mathbf{K}(d)\alpha \\ & \text{s.t.} && \mathbf{K}(d) = \sum_{i=1}^n d_i \mathbf{K}_i, \quad d_j \geq 0, \quad \sum_j d_j \leq 1. \end{aligned}$$

# Relaxing the regularization term

$$\underset{\alpha \in \mathbb{R}^m, b \in \mathbb{R}, d \in \mathbb{R}^n}{\text{minimize}} \quad L(\mathbf{K}(d)\alpha + b\mathbf{1}) + \frac{\lambda}{2} \alpha^\top \mathbf{K}(d)\alpha$$

$$\text{s.t.} \quad \mathbf{K}(d) = \sum_{i=1}^n d_i \mathbf{K}_i, \quad d_i \geq 0, \quad \sum_i d_i \leq 1.$$

Introduce auxiliary variables  $\alpha_j$  ( $j = 1, \dots, n$ )

$$\alpha^\top \mathbf{K}(d)\alpha = \min_{\alpha_j \in \mathbb{R}^m} \left( \sum_{j=1}^n \frac{\alpha_j^\top \mathbf{K}_j \alpha_j}{d_j} \right) \quad \text{s.t.} \quad \sum_{j=1}^n \mathbf{K}_j \alpha_j = \mathbf{K}(d)\alpha$$

(Proof) Introduce Lagrangian multiplier  $\beta$  and minimize

$$\frac{1}{2} \sum_{j=1}^n \frac{\alpha_j^\top \mathbf{K}_j \alpha_j}{d_j} + \beta^\top \left( \mathbf{K}(d)\alpha - \sum_{j=1}^n \mathbf{K}_j \alpha_j \right).$$

$$\alpha_j = d_j \beta, \quad \beta = \alpha$$

# Relaxing the regularization term

$$\underset{\alpha \in \mathbb{R}^m, b \in \mathbb{R}, d \in \mathbb{R}^n}{\text{minimize}} \quad L(\mathbf{K}(d)\alpha + b\mathbf{1}) + \frac{\lambda}{2} \alpha^\top \mathbf{K}(d)\alpha$$

$$\text{s.t.} \quad \mathbf{K}(d) = \sum_{i=1}^n d_j \mathbf{K}_j, \quad d_j \geq 0, \quad \sum_j d_j \leq 1.$$

Introduce auxiliary variables  $\alpha_j$  ( $j = 1, \dots, n$ )

$$\alpha^\top \mathbf{K}(d)\alpha = \min_{\alpha_j \in \mathbb{R}^m} \left( \sum_{j=1}^n \frac{\alpha_j^\top \mathbf{K}_j \alpha_j}{d_j} \right) \quad \text{s.t.} \quad \sum_{j=1}^n \mathbf{K}_j \alpha_j = \mathbf{K}(d)\alpha$$

(Proof) Introduce Lagrangian multiplier  $\beta$  and minimize

$$\frac{1}{2} \sum_{j=1}^n \frac{\alpha_j^\top \mathbf{K}_j \alpha_j}{d_j} + \beta^\top \left( \mathbf{K}(d)\alpha - \sum_{j=1}^n \mathbf{K}_j \alpha_j \right).$$

$$\alpha_j = d_j \beta, \quad \beta = \alpha$$

# Relaxing the regularization term

$$\underset{\alpha \in \mathbb{R}^m, b \in \mathbb{R}, d \in \mathbb{R}^n}{\text{minimize}} \quad L(\mathbf{K}(d)\alpha + b\mathbf{1}) + \frac{\lambda}{2} \alpha^\top \mathbf{K}(d)\alpha$$

$$\text{s.t.} \quad \mathbf{K}(d) = \sum_{i=1}^n d_i \mathbf{K}_i, \quad d_i \geq 0, \quad \sum_i d_i \leq 1.$$

Introduce auxiliary variables  $\alpha_j$  ( $j = 1, \dots, n$ )

$$\alpha^\top \mathbf{K}(d)\alpha = \min_{\alpha_j \in \mathbb{R}^m} \left( \sum_{j=1}^n \frac{\alpha_j^\top \mathbf{K}_j \alpha_j}{d_j} \right) \quad \text{s.t.} \quad \sum_{j=1}^n \mathbf{K}_j \alpha_j = \mathbf{K}(d)\alpha$$

(Proof) Introduce Lagrangian multiplier  $\beta$  and minimize

$$\frac{1}{2} \sum_{j=1}^n \frac{\alpha_j^\top \mathbf{K}_j \alpha_j}{d_j} + \beta^\top \left( \mathbf{K}(d)\alpha - \sum_{j=1}^n \mathbf{K}_j \alpha_j \right).$$

$$\alpha_j = d_j \beta, \quad \beta = \alpha$$

# Minimization of the upper-bound

$$\begin{aligned} & \underset{\alpha_j \in \mathbb{R}^m, b \in \mathbb{R}, \mathbf{d} \in \mathbb{R}^n}{\text{minimize}} && L\left(\sum_{j=1}^n \mathbf{K}_j \alpha_j + b \mathbf{1}\right) + \frac{\lambda}{2} \sum_{j=1}^n \frac{\alpha_j^\top \mathbf{K}_j \alpha_j}{d_j} \\ & \text{s.t.} && d_j \geq 0, \quad \sum_j d_j \leq 1. \end{aligned}$$

$$\begin{aligned} \sum_{j=1}^n \frac{\alpha_j^\top \mathbf{K}_j \alpha_j}{d_j} &= \sum_{j=1}^n d_j \frac{\alpha_j^\top \mathbf{K}_j \alpha_j}{d_j^2} = \sum_{j=1}^n d_j \left( \frac{\|\alpha_j\|_{\kappa_j}}{d_j} \right)^2 \\ &\geq \left( \sum_{j=1}^n d_j \frac{\|\alpha_j\|_{\kappa_j}}{d_j} \right)^2 \quad \left( \sum_j d_j = 1 \right. \\ &\quad \left. \text{Jensen's inequality} \right) \\ &= \left( \sum_{j=1}^n \|\alpha_j\|_{\kappa_j} \right)^2 \end{aligned}$$

1 linear sum of RKHS norms



# Minimization of the upper-bound

$$\begin{aligned} & \underset{\alpha_j \in \mathbb{R}^m, b \in \mathbb{R}, \mathbf{d} \in \mathbb{R}^n}{\text{minimize}} && L\left(\sum_{j=1}^n \mathbf{K}_j \alpha_j + b \mathbf{1}\right) + \frac{\lambda}{2} \sum_{j=1}^n \frac{\alpha_j^\top \mathbf{K}_j \alpha_j}{d_j} \\ & \text{s.t.} && d_j \geq 0, \quad \sum_j d_j \leq 1. \end{aligned}$$

$$\begin{aligned} \sum_{j=1}^n \frac{\alpha_j^\top \mathbf{K}_j \alpha_j}{d_j} &= \sum_{j=1}^n d_j \frac{\alpha_j^\top \mathbf{K}_j \alpha_j}{d_j^2} = \sum_{j=1}^n d_j \left( \frac{\|\alpha_j\|_{\mathbf{K}_j}}{d_j} \right)^2 \\ &\geq \left( \sum_{j=1}^n d_j \frac{\|\alpha_j\|_{\mathbf{K}_j}}{d_j} \right)^2 \quad \left( \begin{array}{l} \sum_j d_j = 1 \\ \text{Jensen's inequality} \end{array} \right) \\ &= \left( \sum_{j=1}^n \|\alpha_j\|_{\mathbf{K}_j} \right)^2 \end{aligned}$$

↑ linear sum of RKHS norms

# Minimization of the upper-bound

$$\begin{aligned} & \underset{\alpha_j \in \mathbb{R}^m, b \in \mathbb{R}, \mathbf{d} \in \mathbb{R}^n}{\text{minimize}} && L\left(\sum_{j=1}^n \mathbf{K}_j \alpha_j + b \mathbf{1}\right) + \frac{\lambda}{2} \sum_{j=1}^n \frac{\alpha_j^\top \mathbf{K}_j \alpha_j}{d_j} \\ & \text{s.t.} && d_j \geq 0, \quad \sum_j d_j \leq 1. \end{aligned}$$

$$\begin{aligned} \sum_{j=1}^n \frac{\alpha_j^\top \mathbf{K}_j \alpha_j}{d_j} &= \sum_{j=1}^n d_j \frac{\alpha_j^\top \mathbf{K}_j \alpha_j}{d_j^2} = \sum_{j=1}^n d_j \left( \frac{\|\alpha_j\|_{\mathbf{K}_j}}{d_j} \right)^2 \\ &\geq \left( \sum_{j=1}^n d_j \frac{\|\alpha_j\|_{\mathbf{K}_j}}{d_j} \right)^2 \quad \left( \begin{array}{l} \sum_j d_j = 1 \\ \text{Jensen's inequality} \end{array} \right) \\ &= \left( \sum_{j=1}^n \|\alpha_j\|_{\mathbf{K}_j} \right)^2 \end{aligned}$$

↑ linear sum of RKHS norms

# Minimization of the upper-bound

$$\begin{aligned} & \underset{\alpha_j \in \mathbb{R}^m, b \in \mathbb{R}, \mathbf{d} \in \mathbb{R}^n}{\text{minimize}} && L\left(\sum_{j=1}^n \mathbf{K}_j \alpha_j + b \mathbf{1}\right) + \frac{\lambda}{2} \sum_{j=1}^n \frac{\alpha_j^\top \mathbf{K}_j \alpha_j}{d_j} \\ & \text{s.t.} && d_j \geq 0, \quad \sum_j d_j \leq 1. \end{aligned}$$

$$\begin{aligned} \sum_{j=1}^n \frac{\alpha_j^\top \mathbf{K}_j \alpha_j}{d_j} &= \sum_{j=1}^n d_j \frac{\alpha_j^\top \mathbf{K}_j \alpha_j}{d_j^2} = \sum_{j=1}^n d_j \left( \frac{\|\alpha_j\|_{\mathbf{K}_j}}{d_j} \right)^2 \\ &\geq \left( \sum_{j=1}^n d_j \frac{\|\alpha_j\|_{\mathbf{K}_j}}{d_j} \right)^2 && \left( \begin{array}{l} \sum_j d_j = 1 \\ \text{Jensen's inequality} \end{array} \right) \\ &= \left( \sum_{j=1}^n \|\alpha_j\|_{\mathbf{K}_j} \right)^2 \end{aligned}$$

↑ linear sum of RKHS norms

# Minimization of the upper-bound

$$\underset{\alpha_j \in \mathbb{R}^m, b \in \mathbb{R}, \mathbf{d} \in \mathbb{R}^n}{\text{minimize}} \quad L\left(\sum_{j=1}^n \mathbf{K}_j \alpha_j + b \mathbf{1}\right) + \frac{\lambda}{2} \sum_{j=1}^n \frac{\alpha_j^\top \mathbf{K}_j \alpha_j}{d_j}$$

$$\text{s.t.} \quad d_j \geq 0, \quad \sum_j d_j \leq 1.$$

$$\begin{aligned} \sum_{j=1}^n \frac{\alpha_j^\top \mathbf{K}_j \alpha_j}{d_j} &= \sum_{j=1}^n d_j \frac{\alpha_j^\top \mathbf{K}_j \alpha_j}{d_j^2} = \sum_{j=1}^n d_j \left( \frac{\|\alpha_j\|_{\mathbf{K}_j}}{d_j} \right)^2 \\ &\geq \left( \sum_{j=1}^n d_j \frac{\|\alpha_j\|_{\mathbf{K}_j}}{d_j} \right)^2 \quad \left( \begin{array}{l} \sum_j d_j = 1 \\ \text{Jensen's inequality} \end{array} \right) \\ &= \left( \sum_{j=1}^n \|\alpha_j\|_{\mathbf{K}_j} \right)^2 \end{aligned}$$

↑ linear sum of RKHS norms

# Minimization of the upper-bound

$$\begin{aligned} & \underset{\alpha_j \in \mathbb{R}^m, b \in \mathbb{R}, \mathbf{d} \in \mathbb{R}^n}{\text{minimize}} && L\left(\sum_{j=1}^n \mathbf{K}_j \alpha_j + b \mathbf{1}\right) + \frac{\lambda}{2} \sum_{j=1}^n \frac{\alpha_j^\top \mathbf{K}_j \alpha_j}{d_j} \\ & \text{s.t.} && d_j \geq 0, \quad \sum_j d_j \leq 1. \end{aligned}$$

$$\begin{aligned} \sum_{j=1}^n \frac{\alpha_j^\top \mathbf{K}_j \alpha_j}{d_j} &= \sum_{j=1}^n d_j \frac{\alpha_j^\top \mathbf{K}_j \alpha_j}{d_j^2} = \sum_{j=1}^n d_j \left( \frac{\|\alpha_j\|_{\mathbf{K}_j}}{d_j} \right)^2 \\ &\geq \left( \sum_{j=1}^n d_j \frac{\|\alpha_j\|_{\mathbf{K}_j}}{d_j} \right)^2 \quad \left( \begin{array}{l} \sum_j d_j = 1 \\ \text{Jensen's inequality} \end{array} \right) \\ &= \left( \sum_{j=1}^n \|\alpha_j\|_{\mathbf{K}_j} \right)^2 \end{aligned}$$

↑ linear sum of RKHS norms

# Equivalence of the two formulations

Penalizing the **square** of linear sum of RKHS norms (Bach et al.)

$$\underset{\alpha_j \in \mathbb{R}^m, b \in \mathbb{R}}{\text{minimize}} \quad L\left(\sum_{j=1}^n \mathbf{K}_j \alpha_j + b\mathbf{1}\right) + \frac{\lambda}{2} \left(\sum_{j=1}^n \|\alpha_j\|_{\mathbf{K}_j}\right)^2 \quad (\text{A})$$

Penalizing of the linear sum of RKHS norms (proposed)

$$\Leftrightarrow \underset{\alpha_j \in \mathbb{R}^m, b \in \mathbb{R}}{\text{minimize}} \quad L\left(\sum_{j=1}^n \mathbf{K}_j \alpha_j + b\mathbf{1}\right) + \lambda' \sum_{j=1}^n \|\alpha_j\|_{\mathbf{K}_j} \quad (\text{B})$$

Optimality of (A):

$$\nabla_{\alpha_j} L + \lambda \left(\sum_{j=1}^n \|\alpha_j\|_{\mathbf{K}_j}\right) \partial_{\alpha_j} \|\alpha_j\|_{\mathbf{K}_j} \ni 0$$

Optimality of (B):

$$\nabla_{\alpha_j} L + \lambda' \partial_{\alpha_j} \|\alpha_j\|_{\mathbf{K}_j} \ni 0$$

# SpicyMKL

DAL + MKL = SpicyMKL (Sparse Iterative MKL)

- The **bias term**  $b$ , and the **hinge-loss** need special care.
- Soft-thresholding per kernel ( $\leftrightarrow$  per variable)

$$ST_{\lambda}(\alpha_j) = \begin{cases} 0 & (\|\alpha_j\|_{\mathbf{K}_j} \leq \lambda) \\ \left(\|\alpha_j\|_{\mathbf{K}_j} - \lambda\right) \frac{\alpha_j}{\|\alpha_j\|_{\mathbf{K}_j}} & (\text{otherwise}) \end{cases}$$

# Outline

- 1 Introduction
  - Lasso, group lasso and MKL
  - Objective
- 2 Method
  - Proximal minimization algorithm
  - Multiple Kernel Learning
- 3 Experiments
- 4 Summary

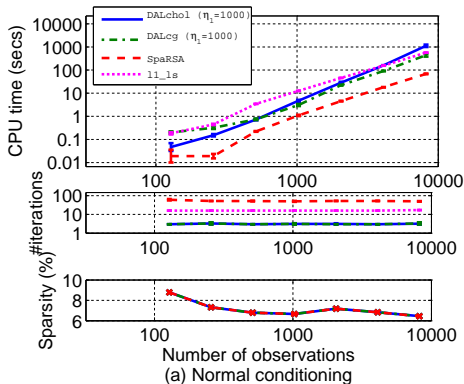


# Experimental setting

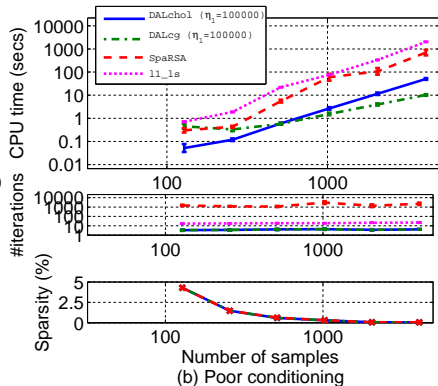
- Problem: lasso (square loss + L1 regularization)
- Comparison with:
  - `l1_ls` (interior-point method)
  - `SpaRSA` (step-size improved IST)
 (problem specific methods (e.g., LARS) are not considered.)
- Random design matrix  $A \in \mathbb{R}^{m \times n}$  ( $m$ : #observations,  $n$ : #unknowns) generated as:
  - $A = \text{randn}(m, n)$  ; (well conditioned)
  - $A = U * \text{diag}(1 ./ (1:m)) * V'$  ; (poorly conditioned)
- Two settings:
  - Medium Scale ( $n = 4m$ ,  $n < 10000$ )
  - Large Scale ( $m = 1024$ ,  $n < 1e+6$ )

## Results (medium scale)

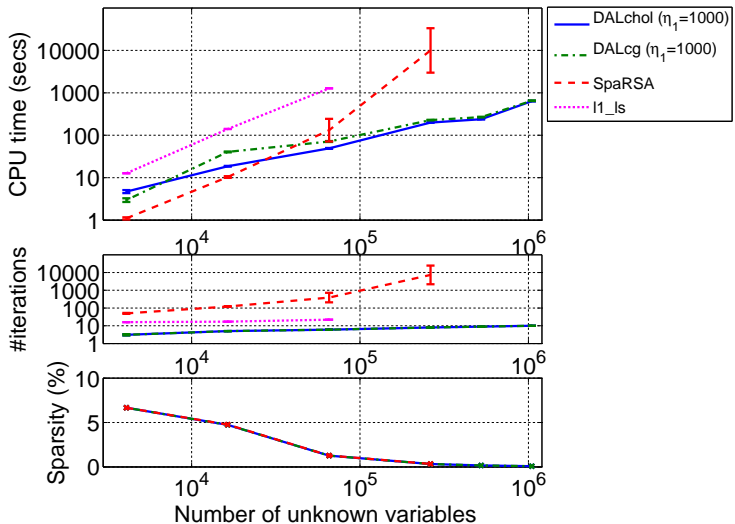
Well conditioned



Poorly conditioned

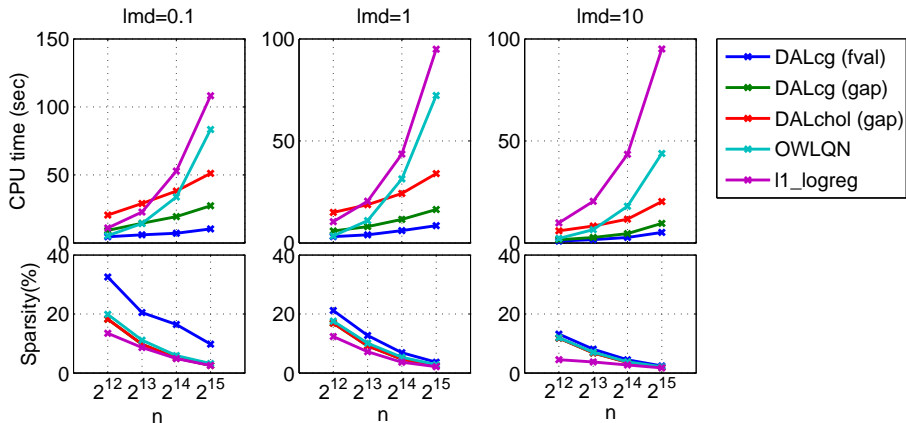


## Results (large scale)



# L1-logistic regression

- $m = 1,024$ .
- $n = 4,096 - 32,768$ .



# Image classification

- Picked five classes [anchor](#), [ant](#), [cannon](#), [chair](#), [cup](#) from Caltech 101 dataset (Fei-Fei et al., 2004).
- Ten 2-class classification problems.
- # kernels 1,760 = Feature extraction (4)  $\times$  Spatial subdivision (22)  $\times$  Kernel functions (20)
  - [Feature extraction](#): (a) hsvsift, (b) sift (scale auto), (c) sift (scale 4px fixed), (d) sift (scale 8px fixed) (used van de Sande's code)
  - [Spatial subdivision and integration](#): (a) whole image, (b) 2x2 grid, and (c) 4x4 grid + spatial pyramid kernel (Lazebnik et al., 06).
  - [Kernel functions](#): Gaussian RBF kernel and  $\chi^2$  kernels using 10 different scale parameters each.

(cf. Gehler & Nowozin, 09)

# Outline

- 1 Introduction
  - Lasso, group lasso and MKL
  - Objective
- 2 Method
  - Proximal minimization algorithm
  - Multiple Kernel Learning
- 3 Experiments
- 4 Summary

# Summary

## DAL (dual augmented Lagrangian)

- is a **dual method** in the **dual** (=primal proximal minimization)
- is efficient when  $m \ll n$ .
- tolerates poorly conditioned design matrix **A** better.
- exploits sparsity in the solution (not in the design).
- **Legendre transform**: linear lower bound instead of linear approximation.

## Tasks:

- theoretical analysis.
- cool application.