

Fast Projections onto $\ell_{1,q}$ -norm Balls (for grouped feature selection)

SUVRIT SRA

Max Planck Institute for Intelligent Systems, Tübingen

ECML 2011, Athens, Greece.

Regularized Optimization

$$\ell(\mathbf{x}) + \lambda r(\mathbf{x})$$

$$\ell(\mathbf{x}) \quad \text{s.t.} \quad r(\mathbf{x}) \leq \gamma$$

Regularized Optimization

$$\ell(x) + \lambda r(x)$$

$$\ell(x) \quad \text{s.t.} \quad r(x) \leq \gamma$$

Inverse
problems

**Comp.
Statistics**

Compr.
Sensing

**Machine
Learning**

Minimize

$$\ell(x) + \lambda r(x)$$

or

$$\ell(x) \quad \text{s.t.} \quad r(x) \leq \gamma$$

Minimize

$$\ell(x) + \lambda r(x)$$

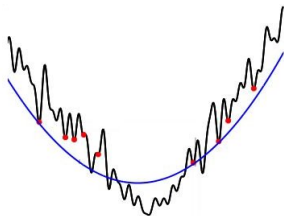
or

$$\ell(x) \quad \text{s.t.} \quad r(x) \leq \gamma$$

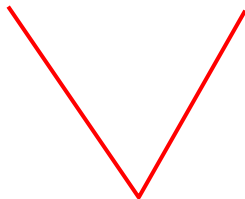
change composite compressed cooccurrence detection feature frequency graphical **group**
higher image ising joint **lasso learning** logistic loss machine
mixed mixed-norms model models **multitask** nonsmooth
norms order point potentials processing **regression** regularizer restoration
selection sensing **sparsity** statistics structured total variation

We solve

$$l(x) \quad \text{s.t.} \quad r(x) \leq \gamma$$



Nonconvex, Differentiable



Non-differentiable

Main iteration

$$\mathbf{x}^{t+1} = \Pi(\mathbf{x}^t - \alpha^t \nabla \ell(\mathbf{x}^t))$$

Main iteration

$$\mathbf{x}^{t+1} = \Pi(\mathbf{x}^t - \alpha^t \nabla \ell(\mathbf{x}^t))$$

α^t : Barzilai-Borwein spectral step-size

Main iteration

$$x^{t+1} = \Pi(x^t - \alpha^t \nabla \ell(x^t))$$

α^t : Barzilai-Borwein spectral step-size

$$\Pi(y) := \operatorname{argmin} \|x - y\|_2 \quad \text{s.t.} \quad r(x) \leq \gamma$$

Step-size is closed form

Main work: **projection**

Mixed norms

Restrict to: $r(\mathbf{X}) := \|\mathbf{X}\|_{1,q}$

Parameter (sub)vectors x_1, \dots, x_T

$$\|\mathbf{X}\|_{1,q} := \sum_i \|x_i\|_q$$

Mixed norms

Restrict to: $r(\mathbf{X}) := \|\mathbf{X}\|_{1,q}$

Parameter (sub)vectors x_1, \dots, x_T

$$\|\mathbf{X}\|_{1,q} := \sum_i \|x_i\|_q$$

- Enforce joint sparsity / feature selection
- $q = 1$: ordinary ℓ_1 ; $q = 2$: **group lasso**;
- $q = \infty$: most severe—tends to eliminate entire feature
- **multitask lasso**; **block compressed sensing**; etc.

We need to solve

$$\min \frac{1}{2} \|x - y\|_2^2 \quad \text{s.t.} \quad \|x\|_{1,q} \leq \gamma$$

We need to solve

$$\min \frac{1}{2} \|x - y\|_2^2 \quad \text{s.t.} \quad \|x\|_{1,q} \leq \gamma$$

Lagrangian duality

There is a θ^* for which we can instead solve

$$\min \frac{1}{2} \|x - y\|_2^2 + \theta^* (\|x\|_{1,q} - \gamma)$$

This is easier!

Projections for mixed norms

$$\|\mathbf{x}\|_{1,q} = \sum_i \|x_i\|_q; \text{ so problem } \textit{separable}$$

Separable ℓ_q -norm *proximity operations*

$$\min \frac{1}{2} \|\mathbf{x} - \mathbf{y}\|_2^2 + \theta^* \|\mathbf{x}\|_q$$

$q = 1$: soft-thresholding (closed-form)

$q = 2$: soft-thresholding (closed-form)

$q = \infty$: requires reformulation

What about θ^* ?

What about θ^* ?

Root finding

Let $g(\theta) := -\gamma + \|u(\theta)\|_{1,q}$
 $u(\theta)$ via *proximity operator*

Theorem: $\theta^* \in [0, \theta_{\max}]$

θ_{\max} via *dual-norm*

Root finding: bisection, inv. quadratic, secant

Solution to ϵ -accuracy in $O(\log(\theta_{\max}/\epsilon))$

$\ell_{1,\infty}$ projections

$$\min \|X - Y\|_F^2 \quad \text{s.t.} \quad \|X\|_{1,\infty} \leq \gamma$$

Projection methods

Quattoni et al., ICML 2009 – **QP**

Our method, ECML 2011 – **FP**

$l_{1,\infty}$ projections

$c\gamma$	QP	FP
.50	1982s	31s

$l_{1,\infty}$ projections

$c\gamma$	QP	FP
.50	1982s	31s
.20	11491s	25s
.10	17064s	23s
.05	20165s	24s

Multitask Lasso

- Let X_j be data matrix for task j
- Seek parameter matrix W
- Columns corr. to tasks, rows to features
- Regularize *shared* features *across* tasks

$$\sum_t \|y_t - X_t w_t\|_2^2 \text{ s.t. } \|W\|_{1,\infty} \leq \gamma$$

Small to medium scale data

Size	# projs	SPG-QP	SPG-FP
27G	16	1054s	320s

Small to medium scale data

Size	# projs	SPG-QP	SPG-FP
27G	16	1054s	320s
0.8G	29	2474s	84s

Small to medium scale data

Size	# projs	SPG-QP	SPG-FP
27G	16	1054s	320s
0.8G	29	2474s	84s
8M	171	826s	31s

MTL on news data

- Subset of CMU News
- Five simultaneous feature selection tasks
- 54K features, 5 sets of 3K documents

MTL on news data

Density	SPG-QP	SPG-FP
.01	6800s	507s
.10	9759s	650s
.20	4746s	554s

SPG-QP spends 96% time in projections

SPG-FP spends 20% time in projections

- Efficient projections onto $\ell_{1,q}$ -norm balls
- Application to mixed-norm regression
- Strong empirical results
- Several open issues remain

Summary

- Efficient projections onto $\ell_{1,q}$ -norm balls
- Application to mixed-norm regression
- Strong empirical results
- Several open issues remain



- Efficient projections onto $\ell_{1,q}$ -norm balls
- Application to mixed-norm regression
- Strong empirical results
- Several open issues remain

감사합니다 Natick
Grazie Danke Ευχαριστίες Dalu Obrigado
Thank You Köszönöm
Спасибо Dank Tack Gracias
谢谢 Merci Seé ありがとう