

Projection onto Minkowski Sums with Application to Constrained Learning

Joong-Ho (Johann) Won¹ Jason Xu² Kenneth Lange³

¹Department of Statistics, Seoul National University

²Department of Statistical Science, Duke University

³Departments of Biomathematics, Human Genetics, and Statistics, UCLA

June 11, 2019

International Conference on Machine Learning

Outline

- Minkowski sum and projection
- Why are Minkowski sums useful for constrained learning?
- Constrained learning via projection onto Minkowski sums
- Minkowski projection algorithm
- Applications to constrained learning
- Conclusion

Minkowski sum of sets

$$A + B \triangleq \{a + b : a \in A, b \in B\}, \quad A, B \subset \mathbb{R}^d$$

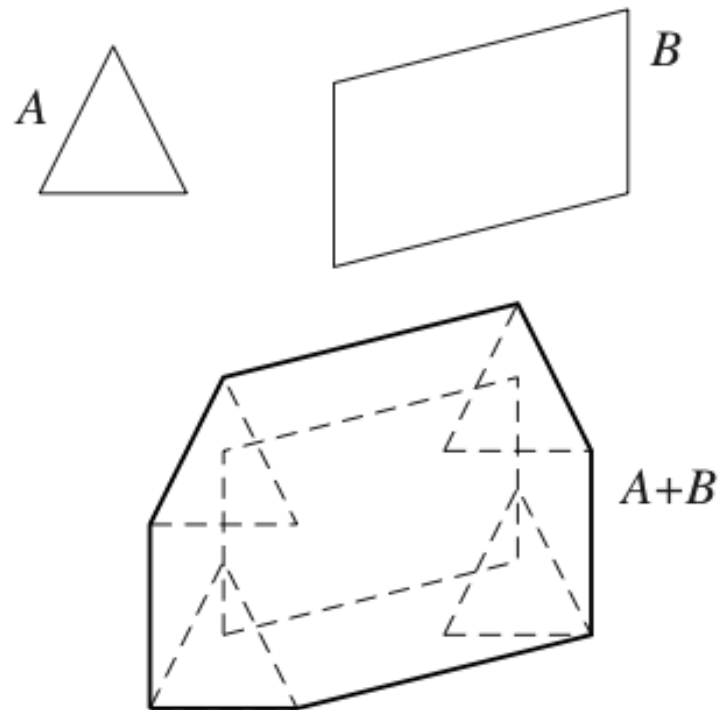


Image source: Christophe Weibel

<https://sites.google.com/site/christopheweibel/research/minkowski-sums>

Projection onto Minkowski sums

$$P_{A+B}(\mathbf{x}) = \operatorname{argmin}_{\mathbf{u} \in A+B} \frac{1}{2} \|\mathbf{u} - \mathbf{x}\|_2^2, \quad \mathbf{x} \notin A+B \quad (\text{P})$$

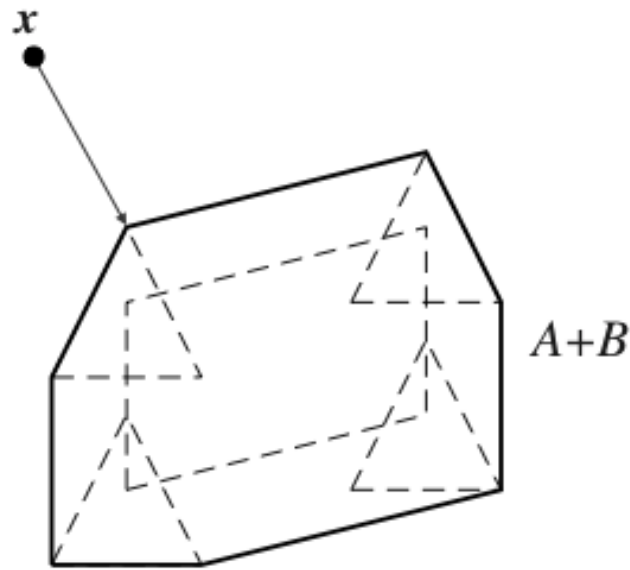


Image source: Christophe Weibel

<https://sites.google.com/site/christopheweibel/research/minkowski-sums>

Why are Minkowski sums useful for constrained learning?

Many penalized or constrained learning problems are of the form

$$\min_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x}) + \sum_{i=1}^k \sigma_{C_i}(\mathbf{x})$$

- $\sigma_C(\mathbf{x}) = \sup_{\mathbf{y} \in C} \langle \mathbf{x}, \mathbf{y} \rangle$ is the support function of convex set C .
- Example: elastic net $\min_{\mathbf{x}} f(\mathbf{x}) + \lambda_1 \|\mathbf{x}\|_1 + \lambda_2 \|\mathbf{x}\|_2$,
 $C_1 = \{\mathbf{x} : \|\mathbf{x}\|_\infty \leq \lambda_1\}$, $C_2 = \{\mathbf{x} : \|\mathbf{x}\|_2 \leq \lambda_2\}$ (dual norm balls)

Why are Minkowski sums useful for constrained learning?

Many penalized or constrained learning problems are of the form

$$\min_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x}) + \sum_{i=1}^k \sigma_{C_i}(\mathbf{x}) = \boxed{\min_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x}) + \sigma_{C_1 + \dots + C_k}(\mathbf{x})} \quad (1)$$

- Support functions are additive over Minkowski sums (Hiriart-Urruty and Lemaréchal 2012).
- New perspective on LHS: minimizing sum of *two* (convex) functions instead of $k + 1$ functions.

Multiple/overlapping norm penalties

$\ell_{1,p}$ group lasso/multitask learning (Yuan and Lin 2006) with overlaps allowed:

$$\min_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x}) + \lambda \sum_{i=1}^k \|\mathbf{x}_{i1}\|_p, \quad p \geq 1$$

where \mathbf{x}_{i1} = subvector of \mathbf{x} ; $i1 \subset \{1, \dots, d\}$ = group index.

- Involved sets: ℓ_q -norm disks.

$$C_i = \{\mathbf{y} = (\mathbf{y}_{i1}, \mathbf{y}_{i2}) : \|\mathbf{y}_{i1}\|_q \leq \lambda, \mathbf{y}_{i2} = \mathbf{0}\},$$
$$\frac{1}{p} + \frac{1}{q} = 1, \quad i2 = \{1, \dots, d\} \setminus i1. \quad (2)$$

- No distinction between overlapping vs. non-overlapping groups!

Conic constraints

$$\min_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x}) \text{ subject to } \mathbf{x} \in K_1^* \cap K_2^* \cap \dots \cap K_k^*$$

where $K_i^* = \{\mathbf{y} : \langle \mathbf{x}, \mathbf{y} \rangle \leq 0, \forall \mathbf{x} \in K_i\}$ is the polar cone of closed convex cone K_i .

- Use the fact $\iota_{K_i^*}(\mathbf{x}) = \sigma_{K_i}(\mathbf{x})$ to express it as

$$\min_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x}) + \sum_{i=1}^k \iota_{K_i^*}(\mathbf{x}) = \min_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x}) + \sum_{i=1}^k \sigma_{K_i}(\mathbf{x}).$$

- $\iota_S = 0/\infty$ indicator of set S

Constrained lasso: mix-and-match

$$\min_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x}) + \lambda \|\mathbf{x}\|_1 \text{ subject to } \mathbf{B}\mathbf{x} = \mathbf{0}, \mathbf{C}\mathbf{x} \leq \mathbf{0},$$

which subsumes the generalized lasso (Tibshirani and Taylor 2011) as a special case (James, Paulson, and Rusmevichientong 2013; Gaines, Kim, and Zhou 2018).

- Involved sets: cone, subspace, and ℓ_∞ -norm ball

$$\begin{aligned} C_1 &= \{\mathbf{x} : \mathbf{B}\mathbf{x} = \mathbf{0}\}^* = \{\mathbf{x} : \mathbf{B}\mathbf{x} = \mathbf{0}\}^\perp, \\ C_2 &= \{\mathbf{x} : \mathbf{C}\mathbf{x} \leq \mathbf{0}\}^*, \quad C_3 = \{\mathbf{x} : \|\mathbf{x}\|_\infty \leq \lambda\} \end{aligned} \tag{3}$$

Constrained learning via projection onto Minkowski sums

Contemporary methods for solving problem (1) (e.g., proximal gradient) requires computing the proximity operator of $\sigma_{C_1+\dots+C_k}$:

$$\text{prox}_{\gamma\sigma_{C_1+\dots+C_k}}(\mathbf{x}) = \underset{\mathbf{u} \in \mathbb{R}^d}{\text{argmin}} \sigma_{C_1+\dots+C_k}(\mathbf{u}) + \frac{1}{2\gamma} \|\mathbf{u} - \mathbf{x}\|_2^2$$

- Proximal gradient:

$$\mathbf{x}^{(t+1)} = \text{prox}_{\gamma_t\sigma_{C_1+\dots+C_k}} \left(\mathbf{x}^{(t)} - \gamma_t^{-1} \nabla f(\mathbf{x}^{(t)}) \right)$$

- Can be computed via Minkowski projection

- Duality:

$$\sigma_{C_1+\dots+C_k}^*(\mathbf{y}) = \iota_{C_1+\dots+C_k}(\mathbf{y}), \quad (\iota_S(\mathbf{u}) = 0 \text{ if } \mathbf{u} \in S, \infty \text{ otherwise})$$

if $C_1 + \dots + C_k$ is closed convex; $g^*(\mathbf{y}) = \sup_{\mathbf{x}} \langle \mathbf{x}, \mathbf{y} \rangle - g(\mathbf{x})$ is the Fenchel conjugate of g .

- Moreau's decomposition

$$\mathbf{x} = \text{prox}_{\gamma g}(\mathbf{x}) + \gamma \text{prox}_{\gamma^{-1} g^*}(\gamma^{-1} \mathbf{x})$$

In terms of Minkowski projection,

$$\begin{aligned} \text{prox}_{\gamma \sigma_{C_1+\dots+C_k}}(\mathbf{x}) &= \mathbf{x} - \gamma \text{prox}_{\gamma^{-1} \iota_{C_1+\dots+C_k}}(\gamma^{-1} \mathbf{x}) \\ &= \boxed{\mathbf{x} - \gamma P_{C_1+\dots+C_k}(\gamma^{-1} \mathbf{x})} \end{aligned}$$

Minkowski projection algorithm

Goal: to develop an efficient method for computing $P_{C_1+\dots+C_k}(\mathbf{x})$, in case projection onto each set $P_{C_i}(\mathbf{x})$ is simple.

MM algorithm:

- 1: Input: External point $\mathbf{x} \notin C_1 + \dots + C_k$;
Projection operator P_{C_i} onto set C_i , $i = 1, \dots, k$;
initial value \mathbf{a}_0^i , $i = 1, \dots, k$; viscosity parameter $\rho \geq 0$
- 2: Initialization: $n \leftarrow 0$
- 3: Repeat
- 4: For $i = 1, 2, \dots, k$
- 5: $\mathbf{a}_{n+1}^{(i)} \leftarrow P_{C_i} \left(\frac{1}{1+\rho} \left(\mathbf{x} - \sum_{j=1}^{i-1} \mathbf{a}_{n+1}^{(j)} - \sum_{j=i+1}^k \mathbf{a}_n^{(j)} \right) + \frac{\rho}{1+\rho} \mathbf{a}_n^{(i)} \right)$
- 6: End For
- 7: $n \leftarrow n + 1$
- 8: Until Convergence
- 9: Return $\sum_{i=1}^k \mathbf{a}_n^{(i)}$

Properties of the Algorithm

- Assume $k = 2$ for exposition purpose: $A = C_1, B = C_2$.

Proposition 1. If both A and B are closed and convex, and $A + B$ is closed, then the Algorithm with $\rho = 0$ generates a sequence converging to $P_{A+B}(\mathbf{x})$.

≫ Proof: paracontraction (Elsner, Koltracht, and Neumann 1992; Lange 2013).

Theorem 1. If in addition either A or B is *strongly convex*, then the sequence generated by Algorithm with $\rho = 0$ converges *linearly* to $P_{A+B}(\mathbf{x})$.

- ≫ Set $C \subset \mathbb{R}^d$ is α -strongly convex with respect to norm $\|\cdot\|$ if there is a constant $\alpha > 0$ such that for any \mathbf{a} and \mathbf{b} in C and any $\gamma \in [0, 1]$, C contains a ball of radius $r = \gamma(1 - \gamma)\frac{\alpha}{2}\|\mathbf{a} - \mathbf{b}\|^2$ centered at $\gamma\mathbf{a} + (1 - \gamma)\mathbf{b}$ (Garber and Hazan 2015).
- ≫ Ex) ℓ_q -norm ball for $q \in (1, 2]$

Theorem 2. If A and B are closed and subanalytic (possibly non-convex), and at least one of them is bounded, then the sequence generated by the Algorithm with $\rho > 0$ converges to a critical point of (P) regardless of the initial values.

≫ Proof: Kurdyka-Łojasiewicz inequality (Bolte, Daniilidis, and Lewis 2007).

Theorem 3. If $A + B$ is polyhedral, then the Algorithm with $\rho > 0$ generates a sequence converging *linearly* to $P_{A+B}(\mathbf{x})$.

≫ Proof: Luo-Tseng error bound (Karimi, Nutini, and Schmidt 2018).

≫ Ex) $\ell_{1,\infty}$ overlapping group penalty/multitask learning; polyhedra are not strongly convex

Applications to constrained learning

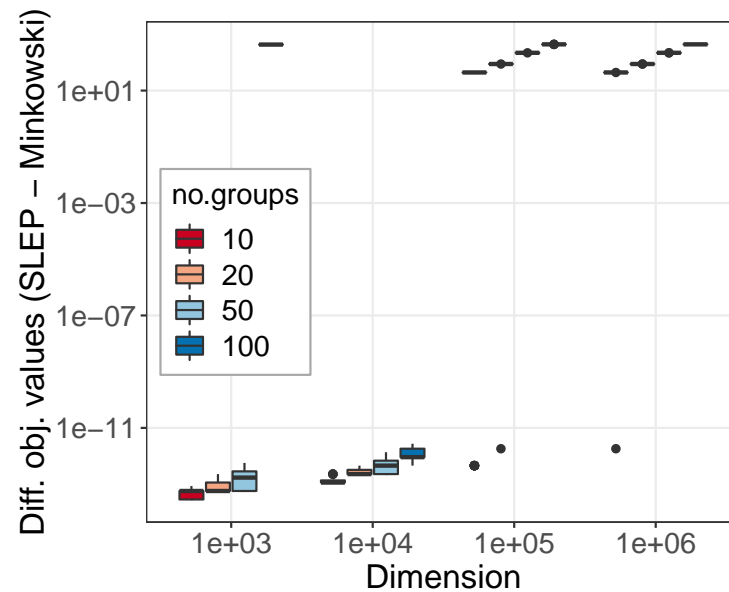
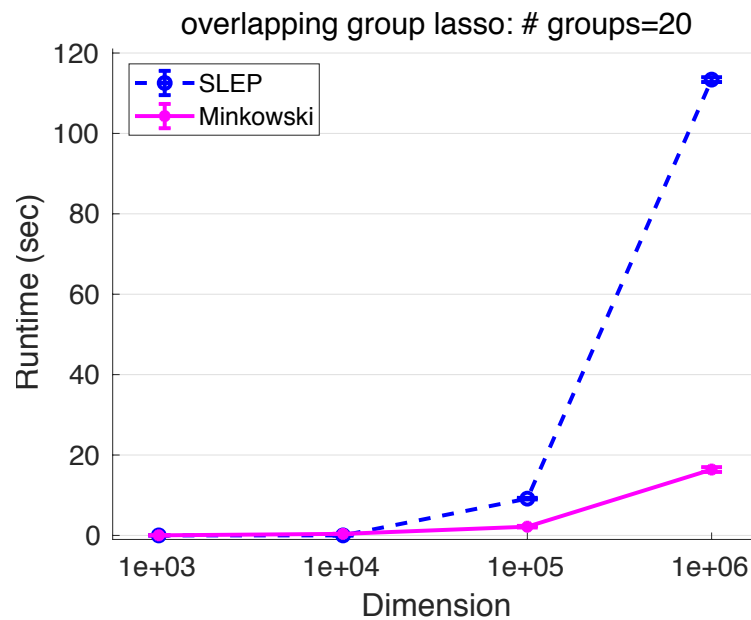
Overlapping group penalties/multitask learning

$$\min_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x}) + \lambda \sum_{i=1}^k \|\mathbf{x}_{i1}\|_p,$$

$$C_i = \{\mathbf{y} = (\mathbf{y}_{i1}, \mathbf{y}_{i2}) : \|\mathbf{y}_{i1}\|_q \leq \lambda, \mathbf{y}_{i2} = \mathbf{0}\}$$

- Overlaps automatically handled with Minkowski projection.
- If $p \in [2, \infty)$, dual ℓ_q -norm disks are strongly convex; if $p = \infty$, polyhedral (linear convergence)
- Fast and reliable algorithm for projection onto ℓ_q -norm disks available (Liu and Ye 2010).

- Comparison to the dual projected gradient method used in SLEP (Yuan, Liu, and Ye 2011; Liu, Ji, and Ye 2011; Zhou, Zhang, and So 2015):

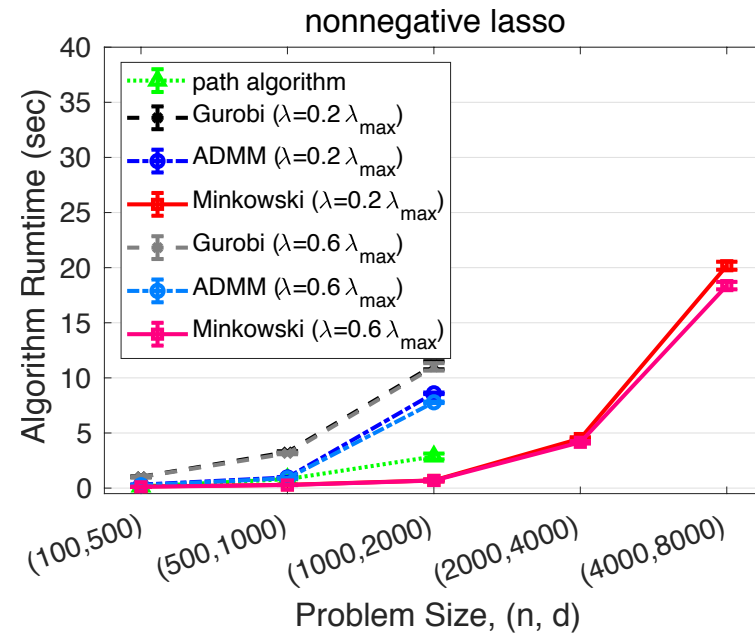
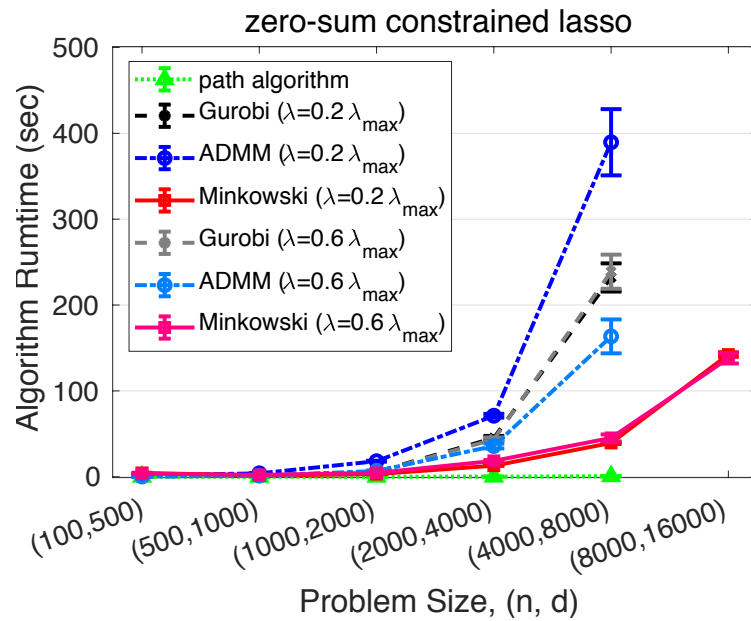


Constrained lasso

$$\min_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x}) + \lambda \|\mathbf{x}\|_1 \text{ subject to } \mathbf{B}\mathbf{x} = \mathbf{0}, \mathbf{C}\mathbf{x} \leq \mathbf{0},$$

- Zero-sum constrained lasso (Lin et al. 2014; Altenbuchinger et al. 2017): $C_1 = \{\mathbf{x} : \sum_{j=1}^d x_j = 0\}^\perp$, $C_2 = \{\mathbf{0}\}$, $C_3 = \{\mathbf{x} : \|\mathbf{x}\|_\infty \leq \lambda\}$ ($\mathbf{B} = \mathbf{1}^T$, $\mathbf{C} = \mathbf{0}$).
- Nonnegative lasso (Efron et al. 2004; El-Arini et al. 2013): $C_1 = \{\mathbf{0}\}$, $C_2 = \{\mathbf{x} : -\mathbf{x} \leq \mathbf{0}\}^*$, $C_3 = \{\mathbf{x} : \|\mathbf{x}\|_\infty \leq \lambda\}$ ($\mathbf{B} = \mathbf{0}$, $\mathbf{C} = -\mathbf{I}$).

- Comparison to generic methods by Gaines, Kim, and Zhou (2018), including path algorithm, ADMM, and commercial solver Gurobi:



Conclusion

- Reconsider constrained learning problems:
 - structural complexities such as non-separability can be handled gracefully via formulations involving Minkowski sums.
- Very simple and efficient algorithm for projecting points onto Minkowski sums of sets:
 - Linear rate of convergence whenever at least one summand is strongly convex or the Luo-Tseng error bound condition is satisfied.
- Our algorithm can serve as an inner loop in, e.g, proximal gradient:
 - Competitive performance
 - Fast (inner loop) convergence is crucial.

*References

- Altenbuchinger, Michael, Thorsten Rehberg, HU Zacharias, Frank Stämmler, Katja Dettmer, Daniela Weber, Andreas Hiergeist, Andre Gessner, Ernst Holler, Peter J Oefner, et al. 2017. “Reference point insensitive molecular data analysis.” *Bioinformatics* 33 (2): 219–226.
- El-Arini, Khalid, Min Xu, Emily B Fox, and Carlos Guestrin. 2013. “Representing documents through their readers.” In *Proceedings of the 19th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 14–22. ACM.
- Bolte, Jérôme, Aris Daniilidis, and Adrian Lewis. 2007. “The Łojasiewicz inequality for nonsmooth subanalytic functions with applications to subgradient dynamical systems.” *SIAM Journal on Optimization* 17 (4): 1205–1223.

- Boyd, Stephen, Neal Parikh, Eric Chu, Borja Peleato, and Jonathan Eckstein. 2010. “Distributed Optimization and Statistical Learning via the Alternating Direction Method of Multipliers.” *Foundations and Trends in Machine Learning* 3 (1): 1–122.
- Davis, Damek, and Wotao Yin. 2017. “A three-operator splitting scheme and its optimization applications.” *Set-valued and Variational Analysis* 25 (4): 829–858.
- Efron, Bradley, Trevor Hastie, Iain Johnstone, Robert Tibshirani, et al. 2004. “Least angle regression.” *Annals of Statistics* 32 (2): 407–499.
- Elsner, Ludwig, Israel Koltracht, and Michael Neumann. 1992. “Convergence of sequential and asynchronous nonlinear paracontractions.” *Numerische Mathematik* 62 (1): 305–319.
- Gaines, Brian R., Juhyun Kim, and Hua Zhou. 2018. “Algorithms for Fitting the Constrained Lasso.” *Journal of Computational and Graphical Statistics* 27 (4): 861–871.

- Garber, Dan, and Elad Hazan. 2015. “Faster rates for the Frank-Wolfe method over strongly-convex sets.” In *Proceedings of the 32nd International Conference on Machine Learning*, 37:541–549.
- Hiriart-Urruty, Jean-Baptiste, and Claude Lemaréchal. 2012. *Fundamentals of Convex Analysis*. Springer Science & Business Media.
- James, Gareth M, Courtney Paulson, and Paat Rusmevichientong. 2013. “Penalized and constrained regression.” *Unpublished Manuscript*, available at <http://www-bcf.usc.edu/~gareth/research/Research.html>.
- Karimi, Hamed, Julie Nutini, and Mark Schmidt. 2018. “Linear Convergence of Gradient and Proximal-Gradient Methods Under the Polyak-Łojasiewicz Condition.” *arXiv preprint arXiv:1608.04636 v3*.
- Lange, Kenneth. 2013. *Optimization*. 2nd. Springer.

- Lin, Wei, Pixu Shi, Rui Feng, and Hongzhe Li. 2014. “Variable selection in regression with compositional covariates.” *Biometrika* 101 (4): 785–797.
- Liu, Jun, Shuiwang Ji, and Jieping Ye. 2011. *SLEP: Sparse learning with efficient projections*. Technical report. Arizona State University. <https://github.com/jiayuzhou/SLEP>.
- Liu, Jun, and Jieping Ye. 2010. “Efficient ℓ_1/ℓ_q norm regularization.” *arXiv preprint arXiv:1009.4766*.
- Tibshirani, Ryan J., and Jonathan Taylor. 2011. “The solution path of the generalized lasso.” *Annals of Statistics* 39 (3): 1335–1371.
- Yuan, Lei, Jun Liu, and Jieping Ye. 2011. “Efficient methods for overlapping group lasso.” In *Advances in Neural Information Processing Systems*, 352–360.

- Yuan, Ming, and Yi Lin. 2006. “Model selection and estimation in regression with grouped variables.” *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 68 (1): 49–67.
- Zhou, Zirui, Qi Zhang, and Anthony Man-Cho So. 2015. “ $\ell_{1,p}$ -Norm Regularization: Error Bounds and Convergence Rate Analysis of First-Order Methods.” In *Proceedings of the 32nd International Conference on Machine Learning*, 37:1501–1510.

Comparison to other algorithms

- Splitting methods: ADMM (Boyd et al. 2010), Davis-Yin three-operator splitting (Davis and Yin 2017)
- Do not produce descent algorithms, and introduce additional variables as well as intermediate steps.
- We do not know whether these methods can achieve a linear convergence rate under, e.g., strong convexity of a summand set.
- Sublinear rates for non-strongly convex sets can be achieved with our algorithm with $\rho > 0$ as well.