

Lifted coordinate descent for learning with trace-norm regularization

Miro Dudik, Zaid Harchaoui, Jérôme Malick

► **To cite this version:**

Miro Dudik, Zaid Harchaoui, Jérôme Malick. Lifted coordinate descent for learning with trace-norm regularization. AISTATS - Proceedings of the Fifteenth International Conference on Artificial Intelligence and Statistics - 2012, Apr 2012, La Palma, Spain. 22, pp.327-336, 2012, JMLR Workshop and Conference Proceedings. <hal-00756802>

HAL Id: hal-00756802

<https://hal.inria.fr/hal-00756802>

Submitted on 23 Nov 2012

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Lifted coordinate descent for learning with trace-norm regularization

Miroslav Dudík
Yahoo! Research, NY

Zaid Harchaoui
INRIA and LJK, Grenoble

Jérôme Malick
CNRS and LJK, Grenoble

Abstract

We consider the minimization of a smooth loss with trace-norm regularization, which is a natural objective in multi-class and multi-task learning. Even though the problem is convex, existing approaches rely on optimizing a non-convex variational bound, which is not guaranteed to converge, or repeatedly perform singular-value decomposition, which prevents scaling beyond moderate matrix sizes. We lift the non-smooth convex problem into an infinitely dimensional smooth problem and apply coordinate descent to solve it. We prove that our approach converges to the optimum, and is competitive or outperforms state of the art.

1 Introduction

A large set of machine learning techniques including SVMs, logistic regression or boosting can be phrased as convex optimization among the set of linear predictors. The optimization objective has typically two parts: *empirical risk*, measuring a goodness of fit to the data, and *regularization*, measuring the complexity of the model and thus controlling its capacity to overfit. Here, we study variants where the model takes the form of a matrix, such as in multi-class or multi-task learning problems, and the regularization takes the form of the trace norm [36, 1, 2, 24, 30]. The trace norm (i.e., the sum of singular values) is the convex envelope of the rank of a matrix (on a special ball, see [15, 31]); thus it encourages matrices of low rank. Even though the resulting problems are convex, the existing techniques for trace-norm regularized optimization [36, 1, 24, 30] either do not scale beyond moderately sized matrices, or are not guaranteed to converge

to the optimum. Here we propose a simple approach based on coordinate descent which is guaranteed to converge to the minimum of the convex objective and also scales to large matrix sizes.

Previous approaches for trace-norm regularized problems fall into two categories. The first set of approaches adapts general-purpose convex optimization to trace norm. Here, the most scalable ones are composite optimization techniques [24, 30]. Composite optimization relies on the implementation of the *proximal operator*. For trace-norm regularization, the key operation in the proximal operator is singular-value decomposition (SVD), which is a bottleneck in scaling to large matrix sizes. The second set of approaches uses variational characterizations of trace norm [36, 1] and proceeds by alternating minimization. These techniques either rely on SVD calculations [1], which prevents their scaling, or lack global convergence guarantees [36], and thus they are sensitive to the starting point and may require extensive problem-specific tuning.

Our approach is based on coordinate descent, which has demonstrated extremely good performance in ℓ_1 -regularized optimization [17, 39]. Since the trace norm is a natural generalization of the ℓ_1 -norm to matrices with the coordinates replaced by rank-one matrices (see, e.g., [31]), we should expect that coordinate descent algorithms will generalize to matrix settings as long as we replace the coordinates by the suitable rank-one matrices. The challenge is determining *which* rank-one matrices. Ideally, we would like to use the rank-one matrices from the SVD of the solution. However, except for least-squares regression problems, it is not clear how to determine these matrices without first finding the solution. This conceptual obstacle has prevented the application of coordinate-descent techniques to matrix setting [35, Chapter 11].

In this paper, we overcome this conceptual obstacle by considering *all possible* (normalized) rank-one matrices as coordinates. This set of matrices forms an overcomplete and uncountable infinite basis of the space of matrices. We show that a simple strategy of performing a coordinate descent on this lifted space actually

Appearing in Proceedings of the 15th International Conference on Artificial Intelligence and Statistics (AISTATS) 2012, La Palma, Canary Islands. Volume XX of JMLR: W&CP XX. Copyright 2012 by the authors.

converges to the right solution. Picking the coordinate corresponding to the steepest descent direction amounts to calculating the top singular-vector pair. This operation is an order of magnitude faster than SVD. In our experiments we demonstrate that our approach is competitive and often outperforms existing approaches. Various tricks of the trade from ℓ_1 -based optimization, such as regularization path calculation [17], can be adapted to our setting.

The idea of coordinate descent in infinite dimensional (or functional) spaces is explored in the boosting literature [25, 11, 40]. In linear programming, this technique is known as column generation [12]. However, these techniques have not been previously applied and analyzed for trace-norm objectives. Several authors [19, 22, 33] have used rank-one updates in convex optimization settings, but their techniques were tailored to quadratic objectives and the focus was on the case of low-rank matrix approximation (e.g., in collaborative filtering) rather than the general purpose trace-norm optimization, which is the subject of this paper. Finally, there are some similarities between our algorithm and algorithms of [34] and [38] for ℓ_1 regularized least-squares regression. Both algorithms rely on coordinate descent techniques to identify a subspace of interesting “active” variables.

In Sec. 2, we describe our problem setting in more detail. In Sec. 3, we lift the non-smooth matrix objective to a smooth objective in infinite dimensions, describe our algorithm and prove its convergence. Finally, in Sec. 4 we evaluate our method on several synthetic and real-world data sets.

2 Supervised learning with trace-norm regularization

We begin by reviewing two supervised learning problems and the associated optimization objectives where the parameters naturally form a matrix. Then we introduce trace-norm regularized optimization.

2.1 Multi-class classification

Our first problem is multi-class classification by multinomial logistic regression. For example, in image classification classes correspond to different object categories, such as human faces, cars, or animals, and each training example is described by a vector of features derived from the pixel representation of the image.

Let $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)$ be a set of labeled training data, where $\mathbf{x}_i \in \mathbb{R}^d$ are feature vectors and $y_i \in \mathcal{Y} := \{1, \dots, k\}$ are the class labels. The linear classifier is specified by a separate weight vector $\mathbf{w}_y \in \mathbb{R}$ for each class. For a given test example

$\mathbf{x} \in \mathbb{R}^d$, the target class is predicted according to $\hat{y} = \text{Arg max}_y \mathbf{w}_y^\top \mathbf{x}$. Class-wise weight vectors form the weight matrix $\mathbf{W} = [\mathbf{w}_1, \dots, \mathbf{w}_k] \in \mathbb{R}^{d \times k}$. In regularized multinomial logistic regression, the classifier is obtained by solving the optimization problem:

$$\underset{\mathbf{W} \in \mathbb{R}^{d \times k}}{\text{Minimize}} \quad \lambda \Omega(\mathbf{W}) + \frac{1}{n} \sum_{i=1}^n L(\mathbf{W}; \mathbf{x}_i, y_i) \quad (1)$$

where $\Omega(\mathbf{W})$ is the regularization and

$$L(\mathbf{W}; \mathbf{x}, y) = \log \left(1 + \sum_{\ell \in \mathcal{Y} \setminus \{y\}} \exp \{ \mathbf{w}_\ell^\top \mathbf{x} - \mathbf{w}_y^\top \mathbf{x} \} \right)$$

is the multinomial logistic loss.

2.2 Multi-task classification

Our second example is multi-task learning, where the goal is to solve multiple classification problems simultaneously. Each individual problem is modeled by multinomial logistic regression. For example, in Sec. 4, we consider the problem of predicting user preferences for different products within some category (such as personal computers). Here, users correspond to tasks. For each user we collect relative preferences over various pairs of products within our category. Each product is described by a vector of features. Examples represent pairs of products, with the feature vector being the difference of the feature vectors of the two products, and the relative preference (for the first or the second product) being the target class. Thus, each individual task is a binary logistic regression problem.

Formally, we are given m data sets (tasks) indexed by $j = 1, \dots, m$, each consisting of n_j examples $(\mathbf{x}_{j,1}, y_{j,1}), \dots, (\mathbf{x}_{j,n_j}, y_{j,n_j})$ where $\mathbf{x}_{j,i} \in \mathbb{R}^d$ is the feature vector and $y_{j,i} \in \mathcal{Y}_j := \{1, \dots, k_j\}$ is a class label. For each task, we fit a matrix of linear predictors $\mathbf{W}_j = [\mathbf{w}_{j,1}, \dots, \mathbf{w}_{j,k_j}]$. We combine matrices across all tasks to obtain the joint matrix $\mathbf{W} = [\mathbf{W}_1, \dots, \mathbf{W}_m] \in \mathbb{R}^{d \times k}$, where $k = k_1 + \dots + k_m$, which is fitted by regularized multi-task logistic regression as follows:

$$\underset{\mathbf{W} \in \mathbb{R}^{d \times k}}{\text{Minimize}} \quad \lambda \Omega(\mathbf{W}) + \frac{1}{n} \sum_{j=1}^m \sum_{i=1}^{n_j} L_j(\mathbf{W}; \mathbf{x}_{j,i}, y_{j,i}) \quad (2)$$

where $\Omega(\mathbf{W})$ is the regularization, $n = n_1 + \dots + n_m$ is the total number of examples, and L_j is the multinomial logistic loss over the submatrix \mathbf{W}_j .

2.3 Matrix norms

Before we introduce trace-norm regularization and regularized optimization, some notation is in order. In a

vector space \mathbb{R}^p , we use notation $\|\cdot\|_2$, $\|\cdot\|_1$, $\|\cdot\|_\infty$, respectively, for ℓ_2 -norm (the Euclidean norm), ℓ_1 -norm (the sum of absolute values), and ℓ_∞ -norm (the maximum of absolute values).

For a matrix $\mathbf{W} \in \mathbb{R}^{d \times k}$, we write $\sigma(\mathbf{W})$ for the spectrum of the matrix, viewed as a vector of its singular values, and define the so-called Schatten p -norm as

$$\|\mathbf{W}\|_{\sigma,p} := \|\sigma(\mathbf{W})\|_p .$$

In this paper we only use $p = 1$ and $p = \infty$. For $p = \infty$, we obtain the maximum-singular value norm. For $p = 1$, we obtain the so-called trace norm of the matrix (also called nuclear norm). Note that for a positive-definite matrix \mathbf{W} , $\|\mathbf{W}\|_{\sigma,1}$ equals the trace of \mathbf{W} , hence the name trace norm.

2.4 Learning with trace-norm penalty

There are many choices of regularization functions in the two problems introduced in the previous sections. We focus on regularization by trace norm, i.e., $\Omega(\mathbf{W}) = \|\mathbf{W}\|_{\sigma,1}$. As we mention in the introduction, trace norm encourages low rank solutions. Hence, it corresponds to the assumption that the linear predictors lie in the same linear subspace of \mathbb{R}^d [1].

The topic of this paper is the optimization problem

$$\underset{\mathbf{W} \in \mathbb{R}^{d \times k}}{\text{Minimize}} \quad \phi_\lambda(\mathbf{W}) := \lambda \|\mathbf{W}\|_{\sigma,1} + \phi(\mathbf{W}) \quad (3)$$

where $\phi : \mathbb{R}^{d \times k} \rightarrow \mathbb{R}$ is the empirical risk (i.e., average loss across training examples). We assume it satisfies the following conditions:

- (A) **convexity**: ϕ is convex
- (B) **lower-boundedness**: ϕ is bounded below; we assume $\phi(\mathbf{W}) \geq 0$ (otherwise ϕ can be shifted)
- (C) **smoothness**: ϕ is differentiable and there exists a norm $\|\cdot\|$, and a constant $H > 0$ such that

$$\langle \mathbf{W}' - \mathbf{W}, \nabla \phi(\mathbf{W}') - \nabla \phi(\mathbf{W}) \rangle \leq H \|\mathbf{W}' - \mathbf{W}\|^2$$

for all $\mathbf{W}', \mathbf{W} \in \mathbb{R}^{d \times k}$.

For example, the empirical risks in Eqs. (1) and (2) satisfy these conditions. This result, based on properties of the multinomial logistic loss function L , is proved in Appendix B.

2.5 Matrix optimization

The learning problem (3) is a convex non-smooth matrix optimization problem. Let us briefly present some of its basic properties.

First, it is easy to see that a solution to the problem (3) always exists. This is because the minimization can

be restricted to the level-set $\{\mathbf{W} : \phi_\lambda(\mathbf{W}) \leq \phi_\lambda(\mathbf{0})\}$, which is compact (since ϕ is bounded below). By continuity, ϕ_λ attains a minimum over this set.

The necessary and sufficient condition for optimality of \mathbf{W} is that $\mathbf{0}$ lies in the subdifferential of ϕ_λ

$$\mathbf{0} \in \partial \phi_\lambda(\mathbf{W}) .$$

By subdifferential calculus [20], this is equivalent to

$$-\nabla \phi(\mathbf{W}) / \lambda \in \partial \|\mathbf{W}\|_{\sigma,1} .$$

Now, since $\partial \|\mathbf{W}\|_{\sigma,1}$ is exactly (see, e.g., [15])

$$\{\mathbf{M} \in \mathbb{R}^{d \times k} : \|\mathbf{M}\|_{\sigma,\infty} \leq 1, \langle \mathbf{M}, \mathbf{W} \rangle = \|\mathbf{W}\|_{\sigma,1}\},$$

we obtain the following proposition:

Proposition 2.1. *A matrix $\mathbf{W} \in \mathbb{R}^{d \times k}$ solves Eq. (3) if and only if*

- (i) $\|\nabla \phi(\mathbf{W})\|_{\sigma,\infty} \leq \lambda$, and
- (ii) $\langle \nabla \phi(\mathbf{W}), \mathbf{W} \rangle = -\lambda \|\mathbf{W}\|_{\sigma,1}$.

We say that \mathbf{W} is an ε -approximate solution, or simply an ε -solution, if the optimality conditions are approximately satisfied:

- (i') $\|\nabla \phi(\mathbf{W})\|_{\sigma,\infty} \leq \lambda + \varepsilon$, and
- (ii') $|\langle \nabla \phi(\mathbf{W}), \mathbf{W} \rangle + \lambda \|\mathbf{W}\|_{\sigma,1}| \leq \varepsilon \|\mathbf{W}\|_{\sigma,1}$.

3 Trace-norm optimization via lifted coordinate descent

This section presents our approach for solving the learning problem (3). The idea is to recast this non-smooth optimization problem in $\mathbb{R}^{d \times k}$ as a smooth optimization problem in an infinite dimensional space. We then exploit the simplicity of the new formulation to design a coordinate descent algorithm.

3.1 Lifting to an infinite dimensional space

We do not have a basis on which we could design a coordinate descent algorithm. So we construct, as follows, an *overcomplete and uncountable infinite "basis"* for the set of matrices living in $\mathbb{R}^{d \times k}$, by considering *all possible* (normalized) rank-one matrices.

Let \mathcal{M} denote the set of rank-one matrices

$$\mathcal{M} = \{\mathbf{u}\mathbf{v}^\top : \mathbf{u} \in \mathbb{R}^d, \mathbf{v} \in \mathbb{R}^k, \|\mathbf{u}\|_2 = \|\mathbf{v}\|_2 = 1\} .$$

Let \mathcal{I} be an index set for the elements of \mathcal{M} , i.e.,

$$\mathcal{M} = \{\mathbf{M}_i \in \mathbb{R}^{d \times k} : i \in \mathcal{I}\} = \{\mathbf{u}_i \mathbf{v}_i^\top : i \in \mathcal{I}\} .$$

A function from \mathcal{I} to \mathbb{R} can be written $\boldsymbol{\theta} = (\theta_i)_{i \in \mathcal{I}} \in \mathbb{R}^{\mathcal{I}}$; its support is the set of indices which are nonzero,

i.e., $\text{supp}(\boldsymbol{\theta}) = \{i \in \mathcal{I} : \theta_i \neq 0\}$. We consider the vector space of functions $\boldsymbol{\theta}$ with finite support, denoted as Θ ,

$$\Theta := \{\boldsymbol{\theta} \in \mathbb{R}^{\mathcal{I}} : \text{supp}(\boldsymbol{\theta}) \text{ is finite}\}$$

equipped with the natural ℓ_1 -norm defined by $\|\boldsymbol{\theta}\|_1 = \sum_{i \in \mathcal{I}} |\theta_i|$. Basic properties of the space and its elements are recalled in Appendix E. The connection between Θ and $\mathbb{R}^{d \times k}$ is the following: each $\boldsymbol{\theta} \in \Theta$ defines a unique matrix in $\mathbb{R}^{d \times k}$

$$\mathbf{W}_{\boldsymbol{\theta}} = \sum_{i \in \mathcal{I}} \theta_i \mathbf{M}_i . \quad (4)$$

The properties of the map $\boldsymbol{\theta} \mapsto \mathbf{W}_{\boldsymbol{\theta}}$ are simple, but important in our developments; we summarize them in the next proposition. It concerns the *non-negative orthant* of Θ :

$$\Theta^+ = \{\boldsymbol{\theta} \in \Theta : \theta_i \geq 0 \text{ for all } i \in \mathcal{I}\} .$$

Proposition 3.1. *The map $\boldsymbol{\theta} \mapsto \mathbf{W}_{\boldsymbol{\theta}}$ is a continuous linear map from Θ to $\mathbb{R}^{d \times k}$. Moreover, for all $\boldsymbol{\theta} \in \Theta^+$, we have*

$$\|\mathbf{W}_{\boldsymbol{\theta}}\|_{\sigma,1} \leq \sum_{i \in \mathcal{I}} \theta_i = \|\boldsymbol{\theta}\|_1$$

and for any $\mathbf{W} \in \mathbb{R}^{d \times k}$, the vector of its singular values corresponds to $\boldsymbol{\theta} \in \Theta^+$ such that $|\text{supp}(\boldsymbol{\theta})| = \text{rank}(\mathbf{W})$, $\mathbf{W}_{\boldsymbol{\theta}} = \mathbf{W}$ and $\|\boldsymbol{\theta}\|_1 = \|\mathbf{W}\|_{\sigma,1}$.

Thus the trace norm in $\mathbb{R}^{d \times k}$ and the ℓ_1 -norm in Θ^+ almost coincide [23]. It is tempting to replace the optimization in \mathbf{W} by optimization in $\boldsymbol{\theta}$ over Θ^+ . Consider $\psi(\boldsymbol{\theta}) := \phi(\mathbf{W}_{\boldsymbol{\theta}})$, the infinite dimensional version of ϕ , and the optimization problem

$$\underset{\boldsymbol{\theta} \in \Theta^+}{\text{Minimize}} \quad \psi_{\lambda}(\boldsymbol{\theta}) := \lambda \sum_{i \in \mathcal{I}} \theta_i + \phi(\mathbf{W}_{\boldsymbol{\theta}}) . \quad (5)$$

By Prop. 3.1, for all $\boldsymbol{\theta} \in \Theta^+$, we have an upper bound

$$\phi_{\lambda}(\mathbf{W}_{\boldsymbol{\theta}}) \leq \psi_{\lambda}(\boldsymbol{\theta}).$$

The next theorem shows that minimizing ϕ_{λ} and ψ_{λ} is actually equivalent.

Theorem 3.2. *The function $\psi_{\lambda}: \Theta \rightarrow \mathbb{R}$ is convex and differentiable. The following optimization problems are equivalent, i.e., they have the same optimal value and correspondence of optimal solutions as*

$$\hat{\boldsymbol{\theta}} \in \underset{\boldsymbol{\theta} \in \Theta^+}{\text{Arg min}} \psi_{\lambda}(\boldsymbol{\theta}) \quad \text{iff} \quad \mathbf{W}_{\hat{\boldsymbol{\theta}}} \in \underset{\mathbf{W} \in \mathbb{R}^{d \times k}}{\text{Arg min}} \phi_{\lambda}(\mathbf{W}) .$$

Note that the first-order optimality conditions for problem (5) are

- (a) $\forall i \in \mathcal{I} : \frac{\partial \psi}{\partial \theta_i}(\boldsymbol{\theta}) \geq -\lambda$
- (b) $\forall i \in \text{supp}(\boldsymbol{\theta}) : \frac{\partial \psi}{\partial \theta_i}(\boldsymbol{\theta}) = -\lambda$.

The ε -approximate optimality is defined by

- (a') $\forall i \in \mathcal{I} : \frac{\partial \psi}{\partial \theta_i}(\boldsymbol{\theta}) \geq -\lambda - \varepsilon$
- (b') $\forall i \in \text{supp}(\boldsymbol{\theta}) : \left| \frac{\partial \psi}{\partial \theta_i}(\boldsymbol{\theta}) + \lambda \right| \leq \varepsilon$

The correspondence of Theorem 3.2 between the optimal solutions of the two problems is in fact even stronger, as it extends to approximate solutions.

Theorem 3.3. *Let ε be such that $0 \leq \varepsilon \leq \lambda$. If $\boldsymbol{\theta}$ is an ε -solution of (5), then $\mathbf{W}_{\boldsymbol{\theta}}$ is an ε -solution of (3).*

We now use this infinite dimensional embedding to design our learning algorithm.

3.2 Coordinate descent algorithm

This section presents our coordinate descent algorithm for optimizing the function ψ_{λ} over Θ^+ , and thus solving the original learning problem (3).

At the current iterate $\boldsymbol{\theta}$, we pick the coordinate along which we can achieve the steepest descent while remaining in Θ^+ . For coordinates $i \notin \text{supp}(\boldsymbol{\theta})$, we can only move in the positive direction. So it suffices to pick $i \in \mathcal{I}$ with the smallest $\frac{\partial \psi_{\lambda}}{\partial \theta_i}(\boldsymbol{\theta})$. This problem is equivalent to calculating the top singular-vector pair of the matrix $-\nabla \phi(\mathbf{W})$ (where $\mathbf{W} = \mathbf{W}_{\boldsymbol{\theta}}$), since

$$\begin{aligned} \text{Arg min}_{i \in \mathcal{I}} \frac{\partial \psi_{\lambda}}{\partial \theta_i}(\boldsymbol{\theta}) &= \text{Arg min}_{i \in \mathcal{I}} (\lambda + \langle \mathbf{M}_i, \nabla \phi(\mathbf{W}) \rangle) \\ &= \text{Arg min}_{i \in \mathcal{I}} \langle \mathbf{u}_i \mathbf{v}_i^{\top}, \nabla \phi(\mathbf{W}) \rangle \\ &= \text{Arg max}_{i \in \mathcal{I}} \mathbf{u}_i^{\top} (-\nabla \phi(\mathbf{W})) \mathbf{v}_i . \end{aligned}$$

For coordinates $i \in \text{supp}(\boldsymbol{\theta})$, it is also possible to move in the negative direction. Here, we perform a traditional ℓ_1 -style coordinate descent. We can update coordinates either cyclically [17] (which is what we do here), or uniformly at random [27]. We optimize over $\text{supp}(\boldsymbol{\theta})$ until the optimality conditions are satisfied.

The final algorithm (Algorithm 1) is called **R1D**, which stands for *rank-one descent*. In our algorithm, we do not compute the steepest direction $\frac{\partial \psi_{\lambda}}{\partial \theta_i}(\boldsymbol{\theta})$, but only use a steep-enough direction (steepest up to $\varepsilon/2$). The following proposition shows that **R1D** is guaranteed to make fixed progress provided that the corresponding partial derivative is large enough.

Proposition 3.4. *There exist $\alpha, \delta > 0$ such that for all $\varepsilon > 0$, $\boldsymbol{\theta} \in \Theta^+$ and $i \in \mathcal{I}$ such that $\frac{\partial \psi_{\lambda}}{\partial \theta_i}(\boldsymbol{\theta}) \leq -\varepsilon$, we have*

$$\psi_{\lambda}(\boldsymbol{\theta} + \delta \mathbf{e}_i) \leq \psi_{\lambda}(\boldsymbol{\theta}) - \alpha \varepsilon^2 . \quad (6)$$

We deduce that the algorithm converges to an ε -optimal solution in a finite number of iterations.

Theorem 3.5. ***R1D** provides ε -optimal solutions $\boldsymbol{\theta}_{\varepsilon}$ and \mathbf{W}_{ε} after at most $8\psi_{\lambda}(\boldsymbol{\theta}_0)/\alpha\varepsilon^2$ iterations.*

Algorithm 1 $\mathbf{R1D}(\phi, \lambda, \theta_0, \varepsilon)$

Input: empirical risk ϕ , regularization λ
initial point \mathbf{W}_{θ_0} , convergence threshold ε

Output: ε -optimal \mathbf{W}_{θ}

Notation: $\mathbf{W}_t := \mathbf{W}_{\theta_t}$, $\mathbf{u}_t := \mathbf{u}_{i_t}$, $\mathbf{v}_t := \mathbf{v}_{i_t}$, $\mathbf{e}_t := \mathbf{e}_{i_t}$

Algorithm:

For $t = 0, 1, 2, \dots$:

1. Find an approximate top singular-vector pair of $(-\nabla\phi(\mathbf{W}_t))$, i.e., $i_t \in \mathcal{I}$ such that

$$\mathbf{u}_t^\top (-\nabla\phi(\mathbf{W}_t)) \mathbf{v}_t \geq \|\nabla\phi(\mathbf{W}_t)\|_{\sigma, \infty} - \varepsilon/2$$

2. Let $g_t := \frac{\partial\psi_\lambda}{\partial\theta_{i_t}}(\theta_t) = \lambda + \langle \nabla\phi(\mathbf{W}_t), \mathbf{u}_t \mathbf{v}_t^\top \rangle$

3. If $g_t \leq -\varepsilon/2$

$$\begin{aligned} \mathbf{W}_{t+1} &= \mathbf{W}_t + \delta \mathbf{u}_t \mathbf{v}_t^\top \text{ with } \delta \text{ given by Prop. 3.4} \\ \theta_{t+1} &= \theta_t + \delta \mathbf{e}_t \end{aligned}$$

4. Else (i.e., $g_t > -\varepsilon/2$)

If θ_t satisfies (b'), terminate and return θ_t

Otherwise, compute θ_{t+1} as an ε -solution of the restricted problem $\min_{\theta \in \mathbb{R}_+^{\text{supp}(\theta_t)}} \psi_\lambda(\theta)$

If the upper bound H is not known ahead of time, it is possible to use a search strategy as in [26]. For special loss functions, it is possible to derive a tighter upper bound than implied by Prop. 3.4. For instance, for multi-class loss, we use an upper bound along the lines of [14]. When a specialized bound is not available or too loose (as we found in the case of multi-task loss), we observed that it helps to augment the rule of Prop. 3.4 by line search. Details of these strategies will be provided in the extended version of this paper.

If instead of terminating $\mathbf{R1D}$ after reaching ε -optimality, we decrease ε according to a predefined sequence $(\varepsilon_\ell)_\ell$ converging to 0, we obtain an asymptotic convergence.

Theorem 3.6. *Let $\varepsilon_\ell \rightarrow 0$. Define \mathbf{W}_{θ_ℓ} as the solution generated by $\mathbf{R1D}$ with $\varepsilon = \varepsilon_\ell$. Then a subsequence of $(\mathbf{W}_{\theta_\ell})_\ell$ converges to a solution of (3).*

Running time We discuss the running time of our algorithm focusing on the specific losses associated with multi-class and multi-task learning. Recall that the number of training examples is n , the number of linear predictors (matrix columns) is k and the number of features is d . The final parameter is the size of the current $\text{supp}(\theta_t)$ which we denote r . The key operations are calculations of $\nabla\phi(\mathbf{W})$, approximate top singular-vector pair, and $\nabla_{\text{supp}(\theta)}\psi_\lambda(\theta)$ in Step 4. Their running times are as follows:

- $\nabla\phi(\mathbf{W})$: bottleneck of this operation is the

Algorithm 2 $\mathbf{ContR1D}(\phi, \lambda_0, \alpha, N)$

Input: empirical risk ϕ , initial regularization λ_0
multiplicative step $\alpha \in (0, 1)$
number of steps $N \geq 1$

Output: $(\mathbf{W}_{\theta_\ell})_{\ell=0}^N$ minimizing ϕ_{λ_ℓ} with $\lambda_\ell = \lambda_0 \alpha^\ell$

Algorithm:

$$\text{Let } \beta = \frac{1 - \alpha}{1 + \alpha}, \lambda_\ell = \lambda_0 \alpha^\ell, \varepsilon_\ell = \beta \lambda_\ell$$

$\theta_0 = \mathbf{R1D}(\phi, \lambda_0, \mathbf{0}, \varepsilon_0)$

For $\ell = 1, 2, \dots, N$:

$$\theta_\ell = \mathbf{R1D}(\phi, \lambda_\ell, \theta_{\ell-1}, \varepsilon_\ell)$$

matrix-vector product $\mathbf{W}^\top \mathbf{x}_i$ (in multi-class problem) or $\mathbf{W}_j^\top \mathbf{x}_i$ (in multi-task problem), respectively, for $i = 1, \dots, n$ and $j = 1, \dots, n_j$. Naively, the running time is $O(ndk)$ and $O(\sum_j n_j dk_j)$, but exploiting special properties (e.g., sparsity of feature vectors), this step can be much faster. The representation $\mathbf{W} = \mathbf{W}_\theta$ as a sum of rank-one matrices presents an additional opportunity which (even without special structure) immediately yields the running times $O(n(d+k)r)$ and $O(\sum_j n_j(d+k_j)r)$. If n is very large, the summation across examples can be parallelized.

- **approximate top singular-vector pair:** this can be calculated in time $O(dk)$ by a few steps of the power method or Lanczos iterations [9].
- $\nabla_{\text{supp}(\theta)}\psi_\lambda(\theta)$: without additional structure this can be done in time $O(nrk)$ and $O(\sum_j n_j rk_j)$, assuming that $\mathbf{u}_\ell^\top \mathbf{x}_i$ values are precomputed for all $\ell \in \text{supp}(\theta)$ and all $i = 1, \dots, n$ (this is amortized into the calculation of $\nabla\phi(\mathbf{W})$). Again, summation over n can be parallelized if needed.

Continuation It has been noted that solving ℓ_1 -regularized problems is faster when λ is large [17]. In fact, these instances give coordinate descent an edge over other techniques since the solutions for larger λ tend to be sparser. This can be extended to trace-norm problems, with the sparsity replaced by low rank. We can accelerate $\mathbf{R1D}$ by taking a sequence of problems with decreasing values of the λ , and using the intermediate solution as a warm start for the next problem. In addition to the benefit from the warm-starting, we obtain a sequence of models optimizing the same empirical risk with different values of regularization. Such a sequence is called a *regularization path* [18] and is in itself a useful output since in practice we typically choose among several values of λ by cross-validation.

ContR1D (Algorithm 2) is the continuation version of our $\mathbf{R1D}$. It returns a regularization path for

a geometrically spaced sequence of λ 's of the form $\lambda_\ell = \lambda_0 \alpha^\ell$ where $\alpha \in (0, 1)$. As a convergence criterion we use $\varepsilon_\ell = \beta \lambda_\ell$. We select the largest $\beta \in (0, 1)$ that guarantees that sets of ε_ℓ -approximate minimizers of ϕ_{λ_ℓ} *do not intersect* (with the exception of the case $\nabla \phi(\mathbf{0}) = \mathbf{0}$, for which the neighborhoods of $\mathbf{0}$ are approximate minimizers for any $\varepsilon_\ell > 0$). By investigating ε -optimality conditions for consecutive λ 's, we obtain the setting $\beta = (1 - \alpha)/(1 + \alpha)$.

Say that our goal is to calculate an $\bar{\varepsilon}$ -solution for a specific regularization coefficient $\bar{\lambda}$. Then we can use the algorithm **Contr1D** as follows. We set $\lambda_0 = \|\nabla \phi(\mathbf{0})\|_{\sigma, \infty}$ to guarantee that $\mathbf{0}$ is a minimizer of ϕ_{λ_0} . We set $\bar{\beta} = \bar{\varepsilon}/\bar{\lambda}$ and invert the formulas in **Contr1D** to obtain α and N such that $\lambda_N = \bar{\lambda}$ and $\varepsilon_N \leq \bar{\varepsilon}$. We run **Contr1D** with the calculated α and N , and obtain an $\bar{\varepsilon}$ -solution for $\bar{\lambda}$ as our last iterate \mathbf{W}_{θ_N} .

3.3 Previous algorithms

Existing approaches for trace-norm penalized learning fall into three categories: (1) proximal gradient methods [24, 30]; (2) alternating direction methods, based on variational characterizations of trace norm, such as the iterative rescaling [1], or low-norm factorization [36]; and (3) conditional gradient methods [22]. These methods are described in Appendix C. Here we stress the main differences between them and **R1D**.

Proximal gradient An iteration of a basic proximal gradient method [5, 3] consists of a gradient step on ϕ followed by a ‘‘correction’’ (proximal step) according to the trace norm. An accelerated version has good rate of convergence and has been shown to be effective for trace-norm problems [30]. Roughly speaking, $O(1/\sqrt{\varepsilon})$ iterations are required to achieve ε -accuracy, while **R1D** converges in $1/\varepsilon^2$ iterations. However, computing the proximal step for the trace norm requires solving an SVD, i.e., the running time $O(kd \text{rank}(\mathbf{W}))$, while **R1D** needs only an approximate top singular-vector pair with the running time $O(kd)$. As shown in Sec. 4, faster iterations are the key to scaling up to larger problems. Additionally, since our algorithm is incrementally adding rank-one matrices, it automatically maintains a matrix factorization of the approximate solution and has an explicit control over an upper bound on the rank.

Iterative rescaling The approach of [1] is based on reformulating the trace norm as an infimum of reweighted Frobenius norms, which bypasses the non-smoothness of the problem. While the resulting smooth convex optimization problems are easier to solve (by stabilized gradient-like methods for example), we lose (part of) the benefit of the trace-norm

regularization, which is that the trace-norm penalty is enforcing low rank. The numerical experiments will show that this method does not approximate well optimal solutions when they are of low rank.

Low-norm factorization The approach of [36] is to use the decomposition $\mathbf{W} = \mathbf{U}\mathbf{V}^\top$ to reformulate the trace norm as the minimum of the sum of squared Frobenius norms. A block-coordinate descent can then be applied since the minimizations with respect to each factor \mathbf{U} and \mathbf{V} are smooth optimization problems, tackled by stabilized gradient-like methods. This approach however breaks the convexity of the original problem (3), as the reformulated problem is not jointly convex with respect to (\mathbf{U}, \mathbf{V}) . As we will see in our experiments, the algorithm is very sensitive to starting points, and gets stuck in non-optimal critical points.

Conditional gradient Recently, conditional gradient algorithms were proposed for squared-loss problems with a trace-norm *constraint* [22], frequently used in collaborative filtering. We carried out preliminary experiments with the algorithm of [22], but we observed slow convergence. We attribute it to the fact that the algorithm was designed and evaluated on squared loss. While the recommended theoretical step-size works fine for squared loss, we believe it might be too conservative for multinomial logistic losses. Furthermore, it is difficult to conduct a fair experimental comparison since [22] works on the constrained formulation while our algorithm solves the penalized formulation. We defer a detailed comparison to future work.

4 Experimental results

We conducted experiments covering a wide range of problem scales, feature correlation amounts, and regularization amounts. We use the following acronyms for the four compared methods: **Contr1D** for our algorithm, **Prox++** for accelerated proximal gradient, **IR** for iterative rescaling, and **AM** for alternating minimization for the low-norm factorization objective.

4.1 Synthetic data

We first evaluate the methods on a synthetic multinomial logistic regression problem, following a similar protocol as in [35, Chapter 11].

We use $d = 250$ features and $k = 500$ classes. All our training data sets have 10 examples per class, yielding 5000 examples in total. Examples for each class are sampled from a separate multivariate normal distribution in \mathbb{R}^d with means determined by choosing the first $0.2d$ coordinates independently uniformly at random from $\{-1, 1\}$, and setting the remaining $0.8d$

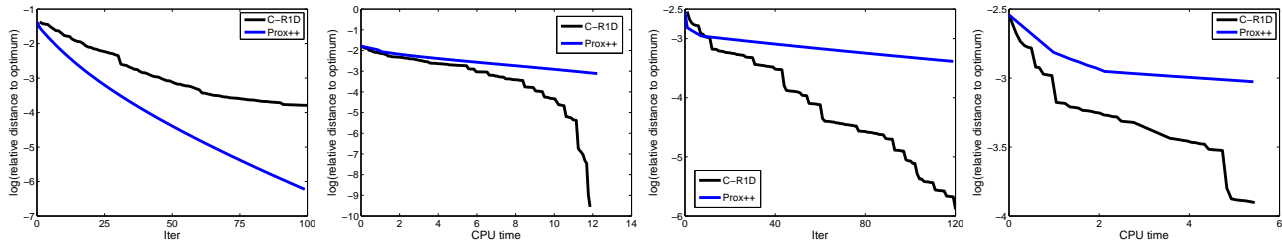


Figure 1: Optimization accuracy of **ContrR1D** and **Prox++** in high correlation settings. Left two plots: accuracy versus the number of iterations and CPU time for light regularization. Right two plots: heavy regularization.

coordinates to zero. The covariance matrix for each class is $(\Sigma)_{i,j} = \sigma^2 \rho^{|i-j|}$, where we use σ to control the separation of the classes and ρ to control correlation among features. We consider $\rho = 0.1$ and $\rho = 0.9$ to obtain *low correlation* and *high correlation* settings. We set σ so that the average distance among all pairs of class means is 3σ . We compared performance for $\lambda = 0.5$ (heavy regularization) and $\lambda = 0.001$ (light regularization). We report performance results averaged over 10 replications by plotting the relative accuracy $|(f - f^*)/f^*|$ against the computation time. We used the smallest value of the objective function attained by the best performing method after a large number of iterations as a proxy for f^* . As we observed similar or worse behavior for the other algorithms compared to **Prox++**, we only compare here our algorithm to **Prox++**.

In Fig. 1, we highlight the main strengths and weaknesses of our approach in high-correlation situations. We see that **ContrR1D** outperforms **Prox++** in terms of CPU time in both high and low regularization settings. As a function of iteration, **Prox++** converges faster for low regularization since its convergence rate is directly controlled by the Lipschitz constant of the gradient (fast convergence for gradient with low Lipschitz constant). However, since its iterations are more costly, **ContrR1D** achieves more accurate solutions faster. Note also the “staircase phenomenon” in the curves for **ContrR1D**, which can be attributed to alternation between Steps 3 and 4. Our approach achieved similar performance to **Prox++** in low correlation situations, so the results are omitted.

4.2 Real-world data

We considered two conjoint analysis datasets [1] and [8], which we refer to as *Conjoint (I)* and *Conjoint (II)*, and a subset of the *ImageNet* dataset 2010 [4]. In Fig. 2 (top), we compare our algorithm to the others in terms of optimization accuracy. In Fig. 2 (bottom), we plot the average test error as a function of training time (averaged over 10 cross-validation splits, using

the best performing regularization coefficient λ).

Conjoint analysis The goal of conjoint analysis is modeling people’s preferences among choices in some set (e.g., products in a certain category). We view it as a multi-task problem with individuals corresponding to tasks and *pairs of choices* corresponding to examples. Each individual choice is modeled by a feature vector, the pair is modeled as the difference of the two vectors, with the target class being the preferred choice between the two listed in the pair, i.e., each task is a binary logistic regression. Dataset *Conjoint (I)* was taken from a survey of 180 individuals, each providing on average preferences for 8 pairs of PCs (among 20 different PCs), parameterized by 13 features. Dataset *Conjoint (II)* is another survey regarding 1187 individuals, each providing on average preferences for 10 pairs of options (among the total of 5), parameterized by 22 features.

On *Conjoint (I)*, all algorithms show similar optimization performance. Since the alternating minimization algorithm **AM** is heavily dependent on the initialization because of the non-convexity of the objective, we only reported the best performance for this algorithm. In terms of test error, three algorithms **ContrR1D**, **Prox++**, and **IR** achieve similar performance, whereas the alternating minimization algorithm **AM** shows higher variance and worse accuracy. The algorithm gets easily trapped in local minima of the objective function which might not correspond to solutions yielding low test error.

On *Conjoint (II)*, the algorithms show different optimization performance. In particular, our algorithm clearly outperforms the other methods. The accelerated proximal gradient **Prox++** and the iterative rescaling **IR** show similar optimization performance. In terms of test error, our algorithm **ContrR1D** performs better than the other three and in fact manages to reach the lowest test error in the early stages of training. This phenomenon might be due to the fact that early iterates of our algorithm tend to have lower rank, which might yield better generalization than the

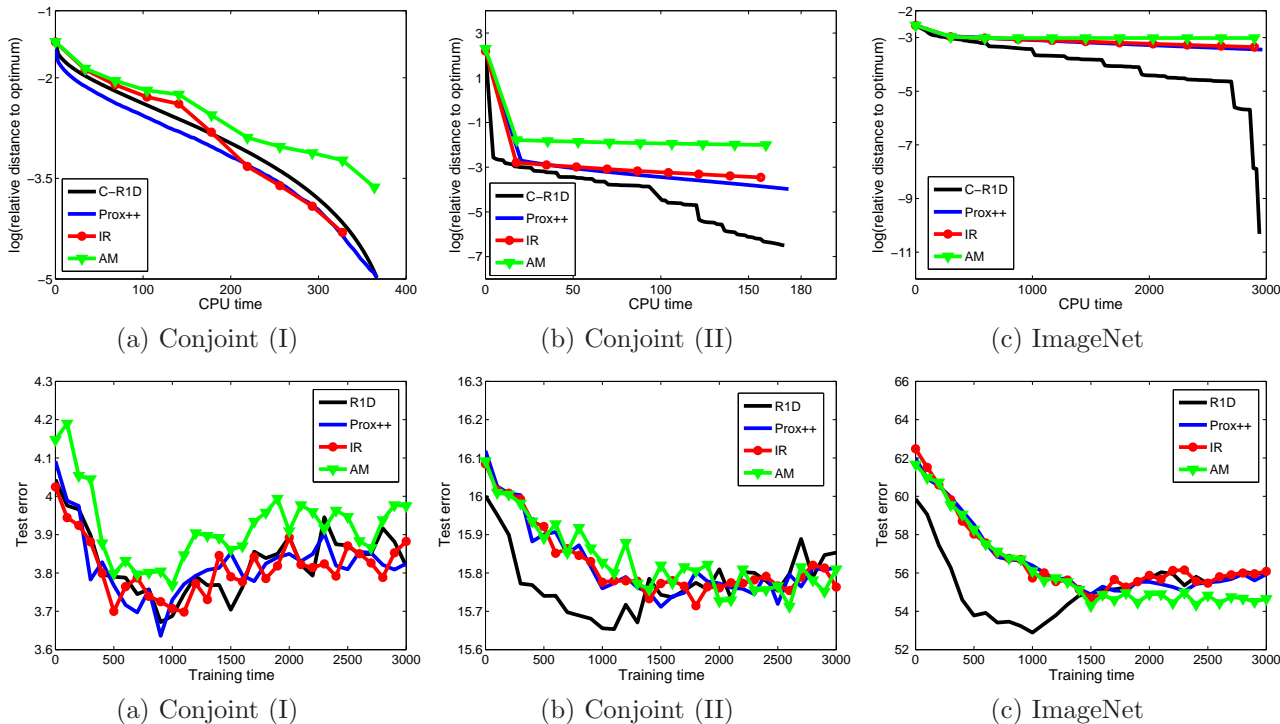


Figure 2: Comparison of optimization performance and test error (in percentage) on real-world data.

final solution of the trace-norm regularized optimization. This behavior is reminiscent of the behavior of coordinate descent algorithms for ℓ_1 -regularized problems (see, e.g., [39]), which tend to achieve better test error in early iterations as well.

ImageNet Here, we tackle a multi-class classification problem with the multinomial logistic loss. We chose a subset of 281 classes from the *ImageNet* dataset, corresponding to various kinds of birds, carnivores, and flowers. Each example is described by 4096 features. We use 10 training and 10 test examples per class.

Our algorithm shows superior optimization performance. This could be explained by the strong correlation of the visual features in this application. Such correlation harms the performance of the other algorithms whereas our algorithm remains robust.

As in the *Conjoint (II)* dataset, our coordinate descent algorithm reaches the lowest test error quickly, but the generalization performance then slightly deteriorates. It is worthwhile to note that, on this dataset, the alternating minimization algorithm **AM** reaches lower test error than the other algorithms after a long training time for a particular run. We interpret this by relating it to the low rank of the weight matrices at the end of the training phase. Since the optimization process of **AM** is not guaranteed to converge to the

global optimum and is strongly biased by the structure of the initial weight matrix, **AM** displays good test-error performance when this structure is particularly tailored to the data at hand. Here, we observe that for the best run of **AM** the initial matrix was of very low rank.

5 Conclusion

We have introduced a new fast coordinate descent algorithm for a wide range of trace-norm regularized learning problems. We have shown that in problems with large matrices, our approach is competitive or outperforms existing optimization algorithms. Our work paves the way for the design of efficient and scalable learning approaches for large-scale matrix problems such as the full *ImageNet* dataset.

Acknowledgements

This work was funded by a Math-STIC project from Grenoble University and the PASCAL 2 Network of Excellence.

References

- [1] A. Argyriou, T. Evgeniou, and M. Pontil. Convex multi-task feature learning. *Machine Learning*, 73(3):243–272, 2008.

- [2] F. Bach. Consistency of trace norm minimization. *JMLR*, 9:1019–1048, 2008.
- [3] A. Beck and M. Teboulle. A fast iterative shrinkage-thresholding algorithm for linear inverse problems. *SIAM Journal on Imaging Sciences*, 2(1):183–202, 2009.
- [4] A. Berg, J. Deng, and F.-F. Li. ImageNet large scale visual recognition challenge, 2010. <http://www.image-net.org/>.
- [5] D. Bertsekas. *Nonlinear Programming (2nd ed.)*. Athena Scientific, 2004.
- [6] S. Boyd and L. Vandenberghe. *Convex Optimization*. Cambridge UP, 2004.
- [7] V. Chandrasekaran. *Convex Optimization Methods for Graphs and Statistical Modeling*. PhD thesis, MIT, 2011.
- [8] O. Chapelle and Z. Harchaoui. A machine learning approach to conjoint analysis. In *Adv. NIPS*. 2005.
- [9] K. Chen. *Matrix Preconditioning Techniques and Applications*. Cambridge UP, 2005.
- [10] K. L. Clarkson. Coresets, sparse greedy approximation, and the Frank-wolfe algorithm. In *Proc. SODA*, 2008.
- [11] M. Collins, R. E. Schapire, and Y. Singer. Logistic regression, AdaBoost and Bregman distances. *Machine Learning*, 47, 2002.
- [12] A. Demiriz, K. P. Bennett, and J. Shawe-Taylor. Linear programming boosting via column generation. *Machine Learning*, 46, 2002.
- [13] V. Demyanov and A. Rubinov. *Approximate Methods in Optimization Problems*. American Elsevier, 1970.
- [14] M. Dudík, S. J. Phillips, and R. E. Schapire. Maximum entropy density estimation with generalized regularization and an application to species distribution modeling. *JMLR*, 8:1217–1260, 2007.
- [15] M. Fazel. *Matrix rank minimization with applications*. PhD thesis, Stanford, 2002.
- [16] M. Frank and P. Wolfe. An algorithm for quadratic programming. *Naval Research Logistics Quarterly*, 3:95–110, 1956.
- [17] J. Friedman, T. Hastie, and R. Tibshirani. Regularization paths for generalized linear models via coordinate descent. *Journal of Statistical Software*, 33(1), 2010.
- [18] T. Hastie, R. Tibshirani, and J. Friedman. *The Elements of Statistical Learning (2nd Ed.)*. Springer Series in Statistics. Springer, 2008.
- [19] E. Hazan. Sparse approximate solutions to semidefinite programs. In *Proc. 8th Latin American Conf. Theor. Informatics*, pages 306–316, 2008.
- [20] J.-B. Hiriart-Urruty and C. Lemaréchal. *Convex Analysis and Minimization Algorithms*. Springer Verlag, Heidelberg, 1993. Two volumes.
- [21] E. J. C. J-F Cai and Z. Shen. A singular value thresholding algorithm for matrix completion. *SIAM J. Optimization*, 20(4):1956–1982, 2008.
- [22] M. Jaggi and M. Sulovský. A simple algorithm for nuclear norm regularized problems. In *ICML*, 2010.
- [23] G. Jameson. *Summing and nuclear norms in Banach space theory*. London Mathematical Society Student Texts, 8. Cambridge University Press. XI, 1987.
- [24] S. Ji and J. Ye. An accelerated gradient method for trace norm minimization. In *ICML*, 2009.
- [25] L. Mason, P. Bartlett, J. Baxter, and M. Frean. Functional gradient techniques for combining hypotheses. In B. Schölkopf, A. Smola, P. Bartlett, and D. Schuurmans, editors, *Adv. Large Margin Classifiers*. MIT Press, 2000.
- [26] Y. Nesterov. Gradient methods for minimizing composite objective function. Technical report, CORE, 2007.
- [27] Y. Nesterov. Efficiency of coordinate descent methods on huge-scale optimization problems. Technical report, CORE, 2010.
- [28] G. Obozinski, B. Taskar, and M. I. Jordan. Joint covariate selection and joint subspace selection for multiple classification problems. *Statistics and Computing*, 2010.
- [29] R. Phelps. *Convex functions, monotone operators, and differentiability*. Lecture notes in mathematics. Springer-Verlag, 1993.
- [30] T. K. Pong, S. J. Paul Tseng, and J. Ye. Trace norm regularization: Reformulations, algorithms, and multi-task learning. *SIAM J. Optimization*, 20(6):3465–3489, 2010.
- [31] B. Recht, M. Fazel, and P. Parrilo. Guaranteed minimum-rank solutions of linear matrix equations via nuclear norm minimization. *SIAM Review*, 52(3):471–501, 2010.
- [32] R. Rockafellar. *Convex Analysis*. Princeton University Press, 1970.
- [33] S. Shalev-Shwartz, A. Gonen, and O. Shamir. Large-scale convex minimization with a low-rank constraint. In *ICML*, 2011.

- [34] W. Shi, G. Wahba, S. Wright, K. Lee, R. Klein, and B. Klein. Lasso-patternsearch algorithm with application to ophthalmology and genomic data. *ASA Proceedings of the Joint Statistical Meetings*, 2006.
- [35] S. Sra, S. Nowozin, and S. J. Wright. *Optimization for Machine Learning*. The MIT Press, 2010.
- [36] N. Srebro, J. D. M. Rennie, and T. S. Jaakola. Maximum-margin matrix factorization. In *Adv. NIPS*, 2005.
- [37] A. Tewari, P. K. Ravikumar, and I. S. Dhillon. Greedy algorithms for structurally constrained high dimensional problems. In *Adv. NIPS*, 2011.
- [38] S. Wright. Accelerated block-coordinate relaxation for regularized optimization. Technical report, 2010.
- [39] G.-X. Yuan, K.-W. Chang, C.-J. Hsieh, and C.-J. Lin. A comparison of optimization methods and software for large-scale l1-regularized linear classification. *JMLR*, 11, December 2010.
- [40] T. Zhang. Sequential greedy approximation for certain convex optimization problems. *IEEE Transaction on Information Theory*, 49:682–691, 2003.

Supplementary Material—Appendix

A Proofs

Proposition 3.1. *The map $\boldsymbol{\theta} \mapsto \mathbf{W}_\boldsymbol{\theta}$ is a continuous linear map from Θ to $\mathbb{R}^{d \times k}$. Moreover, for all $\boldsymbol{\theta} \in \Theta^+$, we have*

$$\|\mathbf{W}_\boldsymbol{\theta}\|_{\sigma,1} \leq \sum_{i \in \mathcal{I}} \theta_i = \|\boldsymbol{\theta}\|_1$$

and for any $\mathbf{W} \in \mathbb{R}^{d \times k}$, the vector of its singular values corresponds to $\boldsymbol{\theta} \in \Theta^+$ such that $|\text{supp}(\boldsymbol{\theta})| = \text{rank}(\mathbf{W})$, $\mathbf{W}_\boldsymbol{\theta} = \mathbf{W}$ and $\|\boldsymbol{\theta}\|_1 = \|\mathbf{W}\|_{\sigma,1}$.

Proof. The linearity is clear by definition of $\mathbf{W}_\boldsymbol{\theta}$. The continuity comes easily as follows. Since \mathcal{M} is compact, there exists a constant M such that $\|\mathbf{M}_i\|_{\sigma,1} \leq M$ for all $i \in \mathcal{I}$. So we write

$$\|\mathbf{W}_\boldsymbol{\theta}\|_{\sigma,1} \leq \sum_{i \in \mathcal{I}} |\theta_i| \|\mathbf{M}_i\|_{\sigma,1} \leq \sum_{i \in \mathcal{I}} |\theta_i| M = M \|\boldsymbol{\theta}\|_1,$$

which proves continuity. By definition of the trace norm, the SVD establishes that any $\mathbf{W} \in \mathbb{R}^{n \times k}$ has a non-negative representation in Θ . \square

Theorem 3.2. *The function $\psi_\lambda: \Theta \rightarrow \mathbb{R}$ is convex and differentiable. The following optimization problems are equivalent, i.e., they have the same optimal value and correspondence of optimal solutions as*

$$\hat{\boldsymbol{\theta}} \in \underset{\boldsymbol{\theta} \in \Theta^+}{\text{Arg min}} \psi_\lambda(\boldsymbol{\theta}) \quad \text{iff} \quad \mathbf{W}_{\hat{\boldsymbol{\theta}}} \in \underset{\mathbf{W} \in \mathbb{R}^{d \times k}}{\text{Arg min}} \phi_\lambda(\mathbf{W}).$$

Proof. The function ψ is the composition of $\boldsymbol{\theta} \mapsto \mathbf{W}_\boldsymbol{\theta}$ with ϕ . Since the first function is linear and the second convex, ψ is convex. Since the first function is continuous and linear, and the second differentiable, ψ is also differentiable. Moreover the function $\boldsymbol{\theta} \mapsto \sum_{i \in \mathcal{I}} \theta_i$ is obviously linear and continuous (so differentiable). So we can conclude that ψ_λ is convex and differentiable.

Let $\hat{\mathbf{W}}$ be a minimizer of (3) (the existence is proved at the end of Sec. 2.5). Let $\hat{\boldsymbol{\theta}}$ be the vector of Θ^+ made by the singular values of $\hat{\mathbf{W}}$, so that we have $\phi_\lambda(\hat{\mathbf{W}}) = \psi_\lambda(\hat{\boldsymbol{\theta}})$. Now write

$$\phi_\lambda(\hat{\mathbf{W}}) = \psi_\lambda(\hat{\boldsymbol{\theta}}) \geq \min \psi_\lambda(\boldsymbol{\theta}) \geq \min \phi_\lambda(\mathbf{W}) \geq \phi_\lambda(\hat{\mathbf{W}}).$$

All the above inequalities are in fact equalities, which proves that the optimal values coincide and in particular that ψ_λ has a minimizer, showing the “if” implication.

We now prove the converse implication. Let $\hat{\boldsymbol{\theta}}$ be a minimizer of (5). Since the optimal values coincide, we have

$$\psi_\lambda(\hat{\boldsymbol{\theta}}) = \min \psi_\lambda(\boldsymbol{\theta}) = \min_{\mathbf{W}} \phi_\lambda(\mathbf{W}) \leq \phi_\lambda(\mathbf{W}_{\hat{\boldsymbol{\theta}}}) \leq \psi_\lambda(\hat{\boldsymbol{\theta}}).$$

This proves the “only if” and finishes the proof. \square

Theorem 3.3. *Let ε be such that $0 \leq \varepsilon \leq \lambda$. If $\boldsymbol{\theta}$ is an ε -solution of (5), then $\mathbf{W}_\boldsymbol{\theta}$ is an ε -solution of (3).*

Proof. Corollary of Thm. D.3. \square

Proposition 3.4. *There exist $\alpha, \delta > 0$ such that for all $\varepsilon > 0$, $\boldsymbol{\theta} \in \Theta^+$ and $i \in \mathcal{I}$ such that $\frac{\partial \psi_\lambda}{\partial \theta_i}(\boldsymbol{\theta}) \leq -\varepsilon$, we have*

$$\psi_\lambda(\boldsymbol{\theta} + \delta \mathbf{e}_i) \leq \psi_\lambda(\boldsymbol{\theta}) - \alpha \varepsilon^2. \quad (7)$$

Proof. To simplify notation, we set $\mathbf{W} = \mathbf{W}_\boldsymbol{\theta}$ and $\mathbf{M} = \mathbf{M}_i$. We introduce also $M = \max_{\mathbf{M} \in \mathcal{M}} \|\mathbf{M}\|$ where $\|\cdot\|$ is the norm from assumption (C) (note that M exists and is finite by compactness of \mathcal{M}). We consider the function

$$f(t) = \phi(\mathbf{W} + t\mathbf{M}) = \psi(\boldsymbol{\theta} + t\mathbf{e}_i).$$

We have, for all t , $f'(t) = \langle \mathbf{M}, \nabla \phi(\mathbf{W} + t\mathbf{M}) \rangle$, and by assumption (C),

$$f'(t) - f'(0) \leq tH\|\mathbf{M}\|^2 \leq tHM^2.$$

Now we write for any $\delta > 0$

$$\begin{aligned} f(\delta) - f(0) &= \int_0^\delta f'(t) dt \\ &= \delta f'(0) + \int_0^\delta (f'(t) - f'(0)) dt \\ &\leq \delta f'(0) + HM^2 \delta^2 / 2. \end{aligned}$$

Observe now that the assumption $\frac{\partial \psi_\lambda}{\partial \theta_i}(\boldsymbol{\theta}) \leq -\varepsilon$ is equivalent to

$$f'(0) = \langle \mathbf{M}, \nabla \phi(\mathbf{W}) \rangle \leq -\varepsilon - \lambda.$$

The above two inequalities yield

$$\phi(\mathbf{W} + \delta \mathbf{M}) + \lambda \delta \leq \phi(\mathbf{W}) - \delta \varepsilon + HM^2 \delta^2 / 2.$$

Hence, for $\delta = \varepsilon / HM^2$,

$$\begin{aligned} \psi_\lambda(\boldsymbol{\theta} + \delta \mathbf{e}_j) &= \phi(\mathbf{W} + \delta \mathbf{M}) + \lambda \delta + \lambda \sum_{i \in \mathcal{I}} \theta_i \\ &\leq \phi(\mathbf{W}) - \frac{\varepsilon^2}{2HM^2} + \lambda \sum_{i \in \mathcal{I}} \theta_i \\ &= \psi_\lambda(\boldsymbol{\theta}) - \frac{\varepsilon^2}{2HM^2}. \end{aligned}$$

Therefore we have a guaranteed decrease with $\alpha = 1/2HM^2$. \square

Theorem 3.5. R1D *provides ε -optimal solutions $\boldsymbol{\theta}_\varepsilon$ and \mathbf{W}_ε after at most $8\psi_\lambda(\boldsymbol{\theta}_0)/\alpha\varepsilon^2$ iterations.*

Proof. The theorem follows easily from Prop. 3.4, as follows. The first observation is that it is not possible to have two iterations in a row where we enter Step 4. Suppose indeed that in iteration t we enter Step 4. Then:

- either $g_{t+1} \leq -\epsilon/2$ and the algorithm will not enter Step 4 in the next iteration;
- or $g_{t+1} > -\epsilon/2$, in which case the algorithm terminates, because $\boldsymbol{\theta}_{t+1}$ satisfies the condition (b'). This comes from the fact that the iterate $\boldsymbol{\theta}_{t+1}$ satisfies the optimality conditions of the restricted problem at the previous Step 4, namely

$$(a'') \forall i \in \text{supp}(\boldsymbol{\theta}_t) : \frac{\partial \psi}{\partial \theta_i}(\boldsymbol{\theta}) \geq -\lambda - \epsilon$$

$$(b'') \forall i \in \text{supp}(\boldsymbol{\theta}_{t+1}) : \left| \frac{\partial \psi}{\partial \theta_i}(\boldsymbol{\theta}) + \lambda \right| \leq \epsilon$$

We prove the bound on the number of iterations by contradiction. Assume that there have been more than $8\psi_\lambda(\boldsymbol{\theta}_0)/\alpha\epsilon^2$ iterations of the algorithm. By the first observation, this yields that there are more than $4\psi_\lambda(\boldsymbol{\theta}_0)/\alpha\epsilon^2$ iterations when we entered Step 3. Therefore, Prop. 3.4 implies that

$$\psi_\lambda(\boldsymbol{\theta}_t) \leq \psi_\lambda(\boldsymbol{\theta}_0) - \frac{4\psi_\lambda(\boldsymbol{\theta}_0)}{\alpha\epsilon^2} \cdot \frac{\alpha\epsilon^2}{4} = 0 .$$

This contradicts the fact that the function ψ is non-negative and completes the bound on the number of iterations.

To conclude the proof we need to argue that on termination, the algorithm returns an ϵ -solution. Above we have already shown that on termination the returned $\boldsymbol{\theta}_t$ satisfies the condition (b'). Condition (a') is implied by $-g_t < \epsilon/2$ and

$$\begin{aligned} \|\nabla\phi(\mathbf{W}_t)\|_{\sigma,\infty} &\leq \mathbf{u}_t^\top (-\nabla\phi(\mathbf{W}_t))\mathbf{v}_t + \epsilon/2 \\ &= -g_t + \lambda + \epsilon/2 \leq \lambda + \epsilon , \end{aligned}$$

which concludes the proof. \square

Theorem 3.6. *Let $\epsilon_\ell \rightarrow 0$. Define $\mathbf{W}_{\boldsymbol{\theta}_\ell}$ as the solution generated by **R1D** with $\epsilon = \epsilon_\ell$. Then a subsequence of $(\mathbf{W}_{\boldsymbol{\theta}_\ell})_\ell$ converges to a solution of (3).*

Proof. Since our algorithm is a descent method, we have $\psi_\lambda(\boldsymbol{\theta}_\ell) \leq \psi_\lambda(\boldsymbol{\theta}_0)$ for all ℓ . By the non-negativity of ϕ and by Thm. 3.2, this yields, for all ℓ ,

$$\|\mathbf{W}_{\boldsymbol{\theta}_\ell}\|_{\sigma,1} \leq \phi_\lambda(\boldsymbol{\theta}_\ell) \leq \psi_\lambda(\boldsymbol{\theta}_0).$$

Thus the sequence $(\mathbf{W}_{\boldsymbol{\theta}_\ell})_\ell$ is bounded (by $\psi_\lambda(\boldsymbol{\theta}_0)$). Let us extract a converging subsequence $(\mathbf{W}_\ell)_\ell$; let us show that its limit \mathbf{W}^* is a solution of (3).

With a slight abuse of notation, let us call again $(\epsilon_\ell)_\ell$ the subsequence associated to $(\mathbf{W}_\ell)_\ell$. By Thm. 3.3, we know that \mathbf{W}_ℓ is a ϵ_ℓ -solution to (3), that is, from (i') and (ii')

$$\begin{aligned} \|\nabla\phi(\mathbf{W}_\ell)\|_{\sigma,\infty} &\leq \lambda + \epsilon_\ell, \\ |\lambda\|\mathbf{W}_\ell\|_{\sigma,1} + \langle \nabla\phi(\mathbf{W}_\ell), \mathbf{W}_\ell \rangle| &\leq \epsilon_\ell\psi_\lambda(\boldsymbol{\theta}_0). \end{aligned}$$

Taking the limit $\epsilon_\ell \rightarrow 0$, we get by continuity of the functions

$$\begin{aligned} \|\nabla\phi(\mathbf{W}^*)\|_{\sigma,\infty} &\leq \lambda, \\ \lambda\|\mathbf{W}^*\|_{\sigma,1} + \langle \nabla\phi(\mathbf{W}^*), \mathbf{W}^* \rangle &= 0 \end{aligned}$$

which are the first order optimality conditions characterizing a solution of (3). \square

B Two loss functions

In this appendix, we establish that the two loss functions considered in this paper, namely the objective functions of the learning problems (1) and (2), satisfy the conditions (A–C). For these technical results, we need another matrix norm: the ℓ_1/ℓ_2 -operator norm defined as

$$\|\mathbf{D}\|_{1,2} = \sup_{\substack{\mathbf{v} \in \mathbb{R}^k \\ \mathbf{v} \neq \mathbf{0}}} \frac{\|\mathbf{D}\mathbf{v}\|_2}{\|\mathbf{v}\|_1} .$$

Proposition B.1. *Let*

$$\phi(\mathbf{W}) = \frac{1}{n} \sum_{i=1}^n L(\mathbf{W}; \mathbf{x}_i, y_i)$$

be the multi-class loss of Eq. (1) and $M = \sup_i \|\mathbf{x}_i\|_2$. Then ϕ satisfies conditions (A–C) for the norm $\|\mathbf{D}\| = \|\mathbf{D}\|_{1,2}$ and the Lipschitz constant $H = M^2$.

Proposition B.2. *Let*

$$\phi(\mathbf{W}) = \frac{1}{n} \sum_{j=1}^m \sum_{i=1}^{n_j} L_j(\mathbf{W}; \mathbf{x}_{ji}, y_{ji})$$

be the multitask loss of Eq. (2) and $M = \sup_{j,i} \|\mathbf{x}_{ji}\|_2$. Then ϕ satisfies conditions (A–C) for the norm $\|\mathbf{D}\| = \max_j \|\mathbf{D}_j\|_{1,2}$ and the Lipschitz constant $H = M^2$.

In fact, conditions (A) and (B) are clear from the convexity and non-negativeness of the multinomial logistic loss function $L(\cdot; \mathbf{x}, y)$. Condition (C) is a corollary of the following lemma.

Lemma B.3. *Let $\mathbf{W}, \mathbf{D} \in \mathbb{R}^{d \times k}$. For all $\mathbf{x} \in \mathbb{R}^d$, $y \in \mathcal{Y}$, we have $L(\mathbf{W}; \mathbf{x}, y) \geq 0$ and*

$$0 \leq \partial^2 L(\mathbf{W} + t\mathbf{D}; \mathbf{x}, y) / \partial t^2 \leq \|\mathbf{x}\|_2^2 \|\mathbf{D}\|_{1,2}^2 .$$

Proof. For a matrix $\mathbf{D} \in \mathbb{R}^{d \times k}$, let \mathbf{d}_ℓ denote its ℓ -th column and \mathbf{D}_j denote its j -th row. The first part follows by observation that $L(\mathbf{W}; \mathbf{x}, y) \geq \log 1 = 0$. For the second part, let $f(t) = L(\mathbf{W} + t\mathbf{D}; \mathbf{x}, y)$, and let

$$p_\ell(t) = \frac{\exp\{(\mathbf{w}_\ell + t\mathbf{d}_\ell)^\top \mathbf{x}\}}{\sum_{\ell' \in \mathcal{Y}} \exp\{(\mathbf{w}_{\ell'} + t\mathbf{d}_{\ell'})^\top \mathbf{x}\}} \quad \text{for } \ell \in \mathcal{Y} .$$

Note that $\sum_{\ell \in \mathcal{Y}} p_\ell(t) = 1$, i.e., $p_\ell(t)$ is a probability distribution over $\ell \in \mathcal{Y}$ for any fixed t . Furthermore,

$$f'(t) = \left(\sum_{\ell \in \mathcal{Y}} p_\ell(t) \mathbf{d}_\ell^\top \mathbf{x} \right) - \mathbf{d}_y^\top \mathbf{x} ,$$

$$f''(t) = \sum_{\ell \in \mathcal{Y}} p_\ell(t) \left(\mathbf{d}_\ell^\top \mathbf{x} - \sum_{\ell' \in \mathcal{Y}} p_{\ell'}(t) \mathbf{d}_{\ell'}^\top \mathbf{x} \right)^2 .$$

From the last identity, $f''(t) \geq 0$. Moreover, using that fact that the variance of a random variable in a range $[a, b]$ is at most $(b - a)^2/4$, we obtain

$$f''(t) \leq \frac{1}{4} \max_{\ell, \ell' \in \mathcal{Y}} (\mathbf{d}_\ell^\top \mathbf{x} - \mathbf{d}_{\ell'}^\top \mathbf{x})^2$$

$$\leq \|\mathbf{x}\|_2^2 (\max_{\ell \in \mathcal{Y}} \|\mathbf{d}_\ell\|_2)^2 = \|\mathbf{x}\|_2^2 \|\mathbf{D}\|_{1,2} .$$

where the last equality follows because

$$\max_{\ell \in \mathcal{Y}} \|\mathbf{d}_\ell\|_2 = \max_{\ell \in \mathcal{Y}} \sup_{\substack{\mathbf{u} \in \mathbb{R}^d: \\ \|\mathbf{u}\|_2 \leq 1}} \mathbf{u}^\top \mathbf{d}_\ell = \sup_{\substack{\mathbf{v} \in \mathbb{R}^k: \\ \|\mathbf{v}\|_1 \leq 1}} \sup_{\substack{\mathbf{u} \in \mathbb{R}^d: \\ \|\mathbf{u}\|_2 \leq 1}} \mathbf{u}^\top \mathbf{D} \mathbf{v}$$

$$= \sup_{\substack{\mathbf{v} \in \mathbb{R}^k: \\ \|\mathbf{v}\|_1 \leq 1}} \|\mathbf{D} \mathbf{v}\|_2 = \sup_{\substack{\mathbf{v} \in \mathbb{R}^k: \\ \mathbf{v} \neq \mathbf{0}}} \frac{\|\mathbf{D} \mathbf{v}\|_2}{\|\mathbf{v}\|_1} = \|\mathbf{D}\|_{1,2}$$

□

C Existing algorithms for trace norm

Here we give some details about the existing trace-norm algorithms mentioned in Sec. 3.3 and used in Sec. 4 in numerical experiments.

Proximal gradient algorithms Proximal gradient methods are specifically tailored to optimize an objective function which is the sum of a smooth function and a non-differentiable regularizer, such as trace norm. They have drawn increasing attention because of their (optimal) guaranteed convergence rate and their ability to deal with large non-smooth problems.

In our context, an iteration of the basic proximal gradient algorithm for solving (3) consists of

$$\mathbf{W}_{t+1} = \text{Prox}_{\sigma,1} \left(\mathbf{W}_t - \frac{1}{H} \nabla \phi(\mathbf{W}_t) \right) .$$

The proximal operator $\text{Prox}_{\sigma,1}$ associated with the trace norm (and parameter λ) is obtained by computing a SVD of the matrix and then replacing each singular value σ_i by $\max\{0, (1 - \lambda/|\sigma_i|)\sigma_i\}$ (its “soft-thresholding”, hence the name of the method of [21]). Accelerated versions of the algorithm [3] use a second variable and combine it with \mathbf{W}_t at marginal extra computational cost with information of previous step.

The basic proximal algorithm has a global convergence rate in $O(1/t)$ where t is the number of iterations of the

algorithm. The accelerated version has a convergence rate in $O(1/t^2)$. However, the computational burden of the SVD computed at each iteration is prohibitive for large-scale problems.

Variational formulation by iterative rescaling

The trace norm has the variational formulation as a reweighted Frobenius norm (see [1])

$$\|\mathbf{W}\|_{\sigma,1} = \min_{\mathbf{D} \succ \mathbf{0}} \text{trace}(\mathbf{W}^\top \mathbf{D}^{-1} \mathbf{W} + \mathbf{D})/2 .$$

The learning problem (3) can then be written as smooth optimization problem:

$$\min_{\mathbf{W} \in \mathbb{R}^{d \times k}} \min_{\mathbf{D} \succ \mathbf{0}} \lambda \text{trace}(\mathbf{W}^\top \mathbf{D}^{-1} \mathbf{W} + \mathbf{D}) + \phi(\mathbf{W}) . \quad (8)$$

A way to deal with the (open) constraint $\mathbf{D} \succ \mathbf{0}$ is to introduce the barrier function $\text{trace}(\mathbf{D}^{-1})$ controlled by a real parameter $\delta > 0$. So in practice, we consider the family of regularized smooth optimization problems parametrized by δ

$$\min_{\mathbf{W}, \mathbf{D}} \lambda \text{trace}(\mathbf{W}^\top \mathbf{D}^{-1} \mathbf{W} + \mathbf{D} + \delta \mathbf{D}^{-1})/2 + \phi(\mathbf{W}) ,$$

replacing the non-smooth learning problem (3).

The above formulation of the problem is particularly well-suited for an alternating direction approach, as follows. The minimization with respect to \mathbf{D}

$$\min_{\mathbf{D} \succ \mathbf{0}} \text{trace}(\mathbf{W}^\top \mathbf{D}^{-1} \mathbf{W} + \mathbf{D} + \delta \mathbf{D}^{-1})$$

has an explicit solution

$$\mathbf{D} = (\mathbf{W} \mathbf{W}^\top + \delta \mathbf{I}_k)^{1/2}$$

which is computed by SVD (of a $k \times k$ -matrix). The minimization over \mathbf{W} consists of minimizing a smooth and (strongly) convex function.

$$\min_{\mathbf{W} \in \mathbb{R}^{d \times k}} \lambda \text{trace}(\mathbf{W}^\top \mathbf{D}^{-1} \mathbf{W}) + \phi(\mathbf{W}) .$$

A wide range of algorithms can be applied to solve this problem, among them (accelerated and stabilized) gradient methods. In practice, rather than solving the problem in \mathbf{W} to optimality, we do only several iterations of such an algorithm.

While bypassing the non-smoothness of the problem, this algorithm loses (part of) the benefit of the trace-norm regularization. Numerical experiments show that this method produces worse solutions when the optimum is of low rank.

Variational formulation by factorization As observed in several works [15, 36, 31], the trace norm has a variational formulation by low-norm factorisation

$$\|\mathbf{W}\|_{\sigma,1} = \min_{\mathbf{W} = \mathbf{U} \mathbf{V}^\top} (\|\mathbf{U}\|^2 + \|\mathbf{V}\|^2)/2 .$$

The learning problem (3) can then be written as

$$\min_{\mathbf{U}, \mathbf{V}} \frac{\lambda}{2} (\|\mathbf{U}\|^2 + \|\mathbf{V}\|^2) + \phi(\mathbf{U}\mathbf{V}^\top). \quad (9)$$

Block-coordinate descent is then an appealing approach for solving the above problems: the minimization with respect to \mathbf{U} , and the one with respect to \mathbf{V} are mere smooth convex optimization problems. Again, accelerated and stabilized gradient methods are adapted for tackling them.

In contrast to the original problem (3) and its formulation (8), problem (9) is not jointly convex with respect to the couple (\mathbf{U}, \mathbf{V}) . As a result, the alternating algorithm can get stuck in local minima, or saddle-points. For example, $\mathbf{U} = \mathbf{V} = \mathbf{0}$ is a critical point but it is not the global minimum. In practice, we observed that the behavior of the algorithm is highly sensitive to the starting point. Problem-dependent tunings as in [36] might be necessary to overcome this weakness.

Conditional gradient approaches Our algorithm shares some similarities with a related family of algorithms, recently applied to learning problems with a *bounded trace-norm constraint* (or low-rank constraint): conditional gradient algorithms. Conditional gradient algorithms [16, 13], a.k.a. Frank-Wolfe algorithms, allow minimize a convex objective ϕ in a simple convex set S . At iteration $(t + 1)$, the conditional gradient algorithm first minimizes the linearized objective at the current iterate within the convex set

$$\tilde{\mathbf{W}}_{t+1} := \underset{\mathbf{W} \in S}{\text{Arg min}} \langle \mathbf{W} - \mathbf{W}_t, \nabla \phi(\mathbf{W}_t) \rangle,$$

then performs a line search over the line segment joining \mathbf{W}_t and $\tilde{\mathbf{W}}_{t+1}$ to obtain \mathbf{W}_{t+1} . Conditional gradient algorithms were applied to learning problems in [10, 19]. Recent works [22, 33] devised conditional gradient algorithms to learning problems with a trace-norm or low-rank constraint, with applications to collaborative filtering. However, these algorithms worked on *constrained formulations*, whereas we consider a *penalized formulation*. Penalized and constrained formulations are equivalent when the entire regularization path is calculated. Even though for each penalty coefficient there exists a constraint that yields the same solution, we are not aware of a method to obtain the matching constraint that does not involve solving the full optimization problem. In our experience, the regularization path calculation is more stable for penalized versions than for constrained version, and penalized versions are for example the state of the art in ℓ_1 -regularization literature.

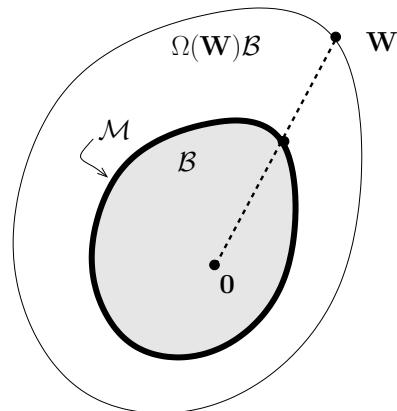


Figure 3: Illustration of the gauge function Ω . To evaluate $\Omega(\mathbf{W})$, we take a ray from the origin towards \mathbf{W} and compute the ratio between the distance to \mathbf{W} and the distance to the intersection of the ray (dotted) with the unit ball \mathcal{B} (in bold).

D Generalization to gauge regularization

In this appendix, we discuss how the optimization algorithm given in the paper generalizes to a broader class of regularization functions. As special cases, we recover coordinate descent for lasso [18], block-coordinate descent for group lasso [28], and rank-one descent discussed in this paper for trace norm. See also [7, 37] for independent, related work.

The specific regularization examples are:

$$\Omega_{\text{lasso}}(\mathbf{W}) = \sum_{j=1}^d \sum_{\ell=1}^k |\mathbf{W}_{j\ell}| \quad (10)$$

$$\Omega_{\text{gr-lasso}}(\mathbf{W}) = \sum_{j=1}^d \|\mathbf{W}_j\|_2 \quad (11)$$

$$\Omega_{\text{trace}}(\mathbf{W}) = \|\mathbf{W}\|_{\sigma,1} \quad (12)$$

where \mathbf{W}_j denotes the j -th row of the matrix. All of them can be naturally defined using the following construction.

Let $\mathcal{M} = \{\mathbf{M}_i \in \mathbb{R}^{d \times k} : i \in \mathcal{I}\}$ be a compact set of matrices, called *atoms*, and let $\mathcal{B} := \text{conv } \mathcal{M}$ be its convex hull. We assume that \mathcal{M} is chosen such that $\mathbf{0} \in \text{int } \mathcal{B}$. We think of \mathcal{M} as an “overcomplete basis” and \mathcal{B} as a “unit ball”. The *gauge function* Ω and *support function* Ω° associated with \mathcal{B} are convex functions defined as (see illustration in Fig. 3; for further details, see [32, 20, 6])

- $\Omega(\mathbf{W}) := \inf\{t \geq 0 : \mathbf{W} \in t\mathcal{B}\}$
- $\Omega^\circ(\mathbf{G}) := \sup_{\mathbf{M} \in \mathcal{B}} \langle \mathbf{M}, \mathbf{G} \rangle = \sup_{\mathbf{M} \in \mathcal{M}} \langle \mathbf{M}, \mathbf{G} \rangle.$

The key property of the gauge function is *sublinearity*:

- $\Omega(t\mathbf{W}) = t\Omega(\mathbf{W})$ for all \mathbf{W} and $t \geq 0$
- $\Omega(\mathbf{W} + \mathbf{W}') \leq \Omega(\mathbf{W}) + \Omega(\mathbf{W}')$ for all \mathbf{W} and \mathbf{W}' .

In addition, by assuming $\mathbf{0} \in \text{int } \mathcal{B}$, we also obtain:

- $\Omega(\mathbf{W}) \geq 0$, with equality if and only if $\mathbf{W} = \mathbf{0}$
- $\{\mathbf{W} : \Omega(\mathbf{W}) \leq t\} = t\mathcal{B}$ for $t \geq 0$, i.e., level sets are compact.

Unlike norms, gauges are not required to be symmetric. The support function plays the role of the dual norm in that $\langle \mathbf{W}, \mathbf{G} \rangle \leq \Omega(\mathbf{W})\Omega^\circ(\mathbf{G})$ for all $\mathbf{W}, \mathbf{G} \in \mathbb{R}^{d \times k}$.

The three examples Eqs. (10)–(12) are obtained by:

$$\begin{aligned} \mathcal{M}_{\text{lasso}} &= \{s\mathbf{e}_j\mathbf{e}_\ell^\top : s \in \{-1, 1\} \\ &\quad j \in \{1, \dots, d\}, \ell \in \{1, \dots, k\}\} \\ \mathcal{M}_{\text{gr-lasso}} &= \{\mathbf{e}_j\mathbf{v}^\top : j \in \{1, \dots, d\}, \mathbf{v} \in \mathbb{R}^k, \|\mathbf{v}\|_2 = 1\} \\ \mathcal{M}_{\text{trace}} &= \{\mathbf{u}\mathbf{v}^\top : \mathbf{u} \in \mathbb{R}^d, \mathbf{v} \in \mathbb{R}^k, \|\mathbf{u}\|_2 = \|\mathbf{v}\|_2 = 1\} \end{aligned}$$

where \mathbf{e}_j is the j -th vector of the Euclidean basis.

Positing the same assumptions on $\phi(\mathbf{W})$ as in Sec. 2.4, we consider minimization of the regularized objective

$$\underset{\mathbf{W} \in \mathbb{R}^{d \times k}}{\text{Minimize}} \quad \phi_\lambda(\mathbf{W}) := \lambda\Omega(\mathbf{W}) + \phi(\mathbf{W}) . \quad (13)$$

By compactness of level sets of Ω , lower-boundedness of ϕ , and continuity, the minimum is attained. Furthermore, the subdifferential of Ω is

$$\partial\Omega(\mathbf{W}) = \{\mathbf{M} \in \mathbb{R}^{d \times k} : \Omega^\circ(\mathbf{M}) \leq 1, \langle \mathbf{M}, \mathbf{W} \rangle = \Omega(\mathbf{W})\}$$

hence the ε -optimality is defined as:

- (i') $\Omega^\circ(-\nabla\phi(\mathbf{W})) \leq \lambda + \varepsilon$, and
- (ii') $|\langle \nabla\phi(\mathbf{W}), \mathbf{W} \rangle + \lambda\Omega(\mathbf{W})| \leq \varepsilon\Omega(\mathbf{W})$.

We define Θ , Θ^+ and \mathbf{W}_θ as before, and the lifted problem as

$$\underset{\theta \in \Theta^+}{\text{Minimize}} \quad \psi_\lambda(\theta) := \lambda \sum_{i \in \mathcal{I}} \theta_i + \phi(\mathbf{W}_\theta) . \quad (14)$$

The ε -optimality for Eq. (14) is defined as before:

- (a') $\forall i \in \mathcal{I} : (-\frac{\partial\psi}{\partial\theta_i}(\theta)) \leq \lambda + \varepsilon$
- (b') $\forall i \in \text{supp}(\theta) : \left| \frac{\partial\psi}{\partial\theta_i}(\theta) + \lambda \right| \leq \varepsilon$

The following is the generalization of Prop. 3.1.

Proposition D.1. *The map $\theta \mapsto \mathbf{W}_\theta$ is a continuous linear map from Θ to $\mathbb{R}^{d \times k}$. Moreover, for all $\theta \in \Theta^+$, we have*

$$\Omega(\mathbf{W}_\theta) \leq \sum_{i \in \mathcal{I}} \theta_i = \|\theta\|_1$$

and for any $\mathbf{W} \in \mathbb{R}^{d \times k}$ there exists $\theta \in \Theta^+$ such that $|\text{supp}(\theta)| \leq (dk + 1)$, $\mathbf{W}_\theta = \mathbf{W}$ and $\|\theta\|_1 = \Omega(\mathbf{W})$.

Proof. The linearity is clear by definition of \mathbf{W}_θ . The continuity comes easily as follows. Consider a norm $\|\cdot\|$ in $\mathbb{R}^{d \times k}$ (all norms are equivalent). Since \mathcal{M} is compact, there exists a constant M such that $\|\mathbf{M}_i\| \leq M$ for all $i \in \mathcal{I}$. So we write

$$\|\mathbf{W}_\theta\| \leq \sum_{i \in \mathcal{I}} |\theta_i| \|\mathbf{M}_i\| \leq \sum_{i \in \mathcal{I}} |\theta_i| M = M \|\theta\|_1 ,$$

which proves continuity.

We next show that any $\mathbf{W} \in \mathbb{R}^{d \times k}$ has a non-negative representation in Θ . The statement is true for $\mathbf{W} = \mathbf{0}$. Now, take $\mathbf{W} \neq \mathbf{0}$, we have $\Omega(\mathbf{W}) \neq 0$. So, we set $\mathbf{W}' = \mathbf{W}/\Omega(\mathbf{W})$. Since \mathbf{W}' lies in $\mathcal{B} = \text{conv } \mathcal{M}$, it can be written as a convex combination of matrices \mathbf{M}_i . By Carathéodory's theorem [20], there exists $\theta' \in \Theta^+$ such that $\sum_{i \in \mathcal{I}} \theta'_i = 1$, $\mathbf{W}' = \mathbf{W}_{\theta'}$, and $|\text{supp}(\theta')| \leq (dk + 1)$. Now, define $\theta = \Omega(\mathbf{W})\theta'$. Observe that $\theta \in \Theta^+$, $|\text{supp}(\theta)| \leq (dk + 1)$, $\mathbf{W}_\theta = \Omega(\mathbf{W})\mathbf{W}_{\theta'} = \mathbf{W}$, and $\|\theta\|_1 = \sum_{i \in \mathcal{I}} \theta_i = \Omega(\mathbf{W}) \sum_{i \in \mathcal{I}} \theta'_i = \Omega(\mathbf{W})$.

Finally, the inequality comes from the sublinearity of Ω and non-negativity of θ as follows:

$$\Omega(\mathbf{W}_\theta) = \Omega \left(\sum_{i \in \mathcal{I}} \theta_i \mathbf{M}_i \right) \leq \sum_{i \in \mathcal{I}} \theta_i \Omega(\mathbf{M}_i) \leq \sum_{i \in \mathcal{I}} \theta_i . \quad \square$$

From Prop. D.1, we obtain

$$\phi_\lambda(\mathbf{W}_\theta) \leq \psi_\lambda(\theta) .$$

We also obtain the equivalence similar to Thm. 3.2 and the sufficiency of lifted ε -optimality similar to Thm. 3.3.

Theorem D.2. *The function $\psi_\lambda : \Theta \rightarrow \mathbb{R}$ is convex and differentiable. The following optimization problems are equivalent, i.e., they have the same optimal value and correspondence of optimal solutions as*

$$\hat{\theta} \in \underset{\theta \in \Theta^+}{\text{Arg min}} \psi_\lambda(\theta) \quad \text{iff} \quad \mathbf{W}_{\hat{\theta}} \in \underset{\mathbf{W} \in \mathbb{R}^{d \times k}}{\text{Arg min}} \phi_\lambda(\mathbf{W}) .$$

Proof. The proof is identical to proof of Thm. 3.2. \square

Theorem D.3. *Let ε be such that $0 \leq \varepsilon \leq \lambda$. If θ is an ε -solution of (14), then \mathbf{W}_θ is an ε -solution of (13).*

Proof. Assume that θ satisfies conditions (a') and (b'). Note that $\frac{\partial\psi}{\partial\theta_i}(\theta) = \langle \mathbf{M}_i, \nabla\phi(\mathbf{W}_\theta) \rangle$. Hence, condition (a') implies

$$\forall \mathbf{M} \in \mathcal{M} : \langle \mathbf{M}, -\nabla\phi(\mathbf{W}_\theta) \rangle \leq \lambda + \varepsilon \quad (15)$$

which yields

$$\Omega^\circ(-\nabla\phi(\mathbf{W}_\theta)) = \sup_{\mathbf{M} \in \mathcal{M}} \langle \mathbf{M}, -\nabla\phi(\mathbf{W}_\theta) \rangle \leq \lambda + \varepsilon ,$$

i.e., condition (i') holds.

It remains to show condition (ii'). First note that

$$\begin{aligned} \langle \nabla\phi(\mathbf{W}_\theta), \mathbf{W}_\theta \rangle &= \sum_{i \in \mathcal{I}} \theta_i \langle \nabla\phi(\mathbf{W}_\theta), \mathbf{M}_i \rangle \\ &= \sum_{i \in \mathcal{I}} \theta_i \left(\frac{\partial\psi}{\partial\theta_i}(\theta) \right) \\ &\leq (-\lambda + \varepsilon) \sum_{i \in \mathcal{I}} \theta_i \leq (-\lambda + \varepsilon) \Omega(\mathbf{W}_\theta) \end{aligned} \quad (16)$$

where the last two inequalities follow by (b') and by Prop. D.1. We also know by Prop. D.1 that there exists $\theta^* \in \Theta^+$ such that $\mathbf{W}_{\theta^*} = \mathbf{W}_\theta$, and $\Omega(\mathbf{W}_\theta) = \sum_{i \in \mathcal{I}} \theta_i^*$. We can write

$$\begin{aligned} \langle \nabla\phi(\mathbf{W}_\theta), \mathbf{W}_\theta \rangle &= \sum_{i \in \mathcal{I}} \theta_i^* \langle \nabla\phi(\mathbf{W}_\theta), \mathbf{M}_i \rangle \\ &\geq (-\lambda - \varepsilon) \sum_{i \in \mathcal{I}} \theta_i^* = (-\lambda - \varepsilon) \Omega(\mathbf{W}_\theta) \end{aligned}$$

where the above inequality follows by Eq. (15). Combining with Eq. (16), we obtain

$$(-\lambda - \varepsilon) \Omega(\mathbf{W}_\theta) \leq \langle \nabla\phi(\mathbf{W}_\theta), \mathbf{W}_\theta \rangle \leq (-\lambda + \varepsilon) \Omega(\mathbf{W}_\theta)$$

i.e., condition (ii') holds as well. \square

Using Thm. D.3, we can derive **AtomDescent** (Algorithm 3), a gauge version of **R1D** (Algorithm 1). The only difference is that the computation of top singular-vector pair is now replaced by the *extremal point* evaluation. For $\mathbf{G} \in \mathbb{R}^{d \times k}$, we define

$$\text{Ext}(\mathbf{G}) = \text{Arg max}_{\mathbf{M} \in \mathcal{M}} \langle \mathbf{M}, \mathbf{G} \rangle .$$

Note that if $\mathbf{M}^* \in \text{Ext}(\mathbf{G})$, we have $\langle \mathbf{M}^*, \mathbf{G} \rangle = \Omega^\circ(\mathbf{G})$. For our three examples, we obtain:

- *lasso*: $\mathbf{M}^* = s^* \mathbf{e}_{j^*} \mathbf{e}_{\ell^*}^\top$
where $(j^*, \ell^*) = \text{Arg max}_{(j, \ell)} |G_{j\ell}|$,
 $s^* = \text{sign } G_{j^* \ell^*}$
- *group lasso*: $\mathbf{M}^* = \mathbf{e}_{j^*} \mathbf{v}^{*\top}$
where $j^* = \text{Arg max}_j \|\mathbf{G}_j\|_2$,
 $\mathbf{v}^{*\top} = \mathbf{G}_{j^*} / \|\mathbf{G}_{j^*}\|_2$
- *trace norm*: $\mathbf{M}^* = \mathbf{u}^* \mathbf{v}^{*\top}$
where $(\mathbf{u}^*, \mathbf{v}^*)$ is the top singular-vector pair of \mathbf{G}

Hence, for lasso we obtain coordinate descent; for group lasso, block-coordinate descent; and for trace norm, rank-one descent. As in **R1D**, also in

Algorithm 3 AtomDescent($\phi, \Omega, \lambda, \theta_0, \varepsilon$)

Input: empirical risk ϕ , gauge Ω , regularization λ
initial point \mathbf{W}_{θ_0} , convergence threshold ε

Output: ε -optimal \mathbf{W}_θ

Notation: $\mathbf{W}_t := \mathbf{W}_{\theta_t}$, $\mathbf{M}_t := \mathbf{M}_{i_t}$, $\mathbf{e}_t := \mathbf{e}_{i_t}$

Algorithm:

For $t = 0, 1, 2, \dots$:

1. Find $i_t \in \mathcal{I}$ corresponding to the coordinate of θ with approximately steepest descent in positive direction, i.e.,

$$\langle \mathbf{M}_t, -\nabla\phi(\mathbf{W}_t) \rangle \geq \Omega^\circ(-\nabla\phi(\mathbf{W}_t)) - \varepsilon/2$$

2. Let $g_t := \frac{\partial\psi}{\partial\theta_{i_t}}(\theta_t) = \lambda + \langle \mathbf{M}_t, \nabla\phi(\mathbf{W}_t) \rangle$

3. If $g_t \leq -\varepsilon/2$

$$\begin{aligned} \mathbf{W}_{t+1} &= \mathbf{W}_t + \delta \mathbf{M}_t \text{ with } \delta \text{ given by Prop. 3.4} \\ \theta_{t+1} &= \theta_t + \delta \mathbf{e}_t \end{aligned}$$

4. Else (i.e., $g_t > -\varepsilon/2$)

If θ_t satisfies (b'), terminate and return θ_t

Otherwise, compute θ_{t+1} as an ε -solution of the restricted problem $\min_{\theta \in \mathbb{R}_+^{\text{supp}(\theta_t)}} \psi_\lambda(\theta)$

AtomDescent, we only insist on approximate extremal points. All of the convergence results for **R1D** also apply to **AtomDescent**, because analysis of **R1D** only concerned the lifted problem Eq. (5) and did not depend on the particular linear map $\theta \mapsto \mathbf{W}_\theta$ (which is the only change between **R1D** and **AtomDescent**).

E Infinite dimensional space Θ

In this appendix, we briefly recall some additional material (especially about differentiability and optimality conditions) to demystify the infinite dimensional space Θ . We assume the gauge setting introduced in the previous section.

The completion of the normed space $(\Theta, \|\cdot\|_1)$ is the complete normed space $(\ell_1(\mathcal{I}), \|\cdot\|_1)$, the space of $(\theta_i)_{i \in \mathcal{I}}$ such that $\sum_{i \in \mathcal{I}} |\theta_i| < +\infty$. The two spaces are in duality with the space $(\ell_\infty(\mathcal{I}), \|\cdot\|_\infty)$ equipped with

$$\|\delta\|_\infty = \max_{i \in \mathcal{I}} |\delta_i|$$

through the bracket notation

$$\langle \delta, \theta \rangle = \sum_{i \in \mathcal{I}} \delta_i \theta_i \leq \|\delta\|_\infty \|\theta\|_1 .$$

Let $\psi : \Theta \rightarrow \mathbb{R}$ be a differentiable function. Its differential $d\psi(\theta) \in \ell_\infty(\mathcal{I})$ can be written with the help of partial derivatives as $d\psi(\theta) = (\frac{\partial\psi}{\partial\theta_i}(\theta))_{i \in \mathcal{I}}$. The

general optimality conditions in this context are the following; they are used in the proof of Prop. E.2.

Proposition E.1. *Let $\psi : \Theta \rightarrow \mathbb{R}$ be a convex differentiable function, and K a convex subset of Θ . Then θ^* is a minimum of ψ over K if and only if*

$$\sum_{i \in \mathcal{I}} \frac{\partial \psi}{\partial \theta_i}(\theta^*)(\theta_i - \theta_i^*) \geq 0 \quad \text{for all } \theta \in K.$$

Proof. The proof is based on the following basic property of convex functions (see [29]). Let $\psi : \Theta \rightarrow \mathbb{R}$ be a convex differentiable function; then

$$\psi(\eta) \geq \psi(\theta) + \langle d\psi(\theta), (\eta - \theta) \rangle \quad \text{for all } \eta, \theta. \quad (17)$$

With the help of the above inequality, the implication “if” is obvious. To prove the “only if” implication, take $t > 0$ and write for any $\theta \in K$, by definition of differentiability,

$$\frac{\psi(\theta^* - t(\theta - \theta^*))}{t} = \langle d\psi(\theta^*), (\theta - \theta^*) \rangle + \frac{o(t)}{t}.$$

Note that $t > 0$ so $\theta^* - t(\theta - \theta^*) \in K$ and then the left-hand-side is nonnegative. Taking the limit $t \rightarrow 0$, we obtain $\langle d\psi(\theta^*), (\theta - \theta^*) \rangle \geq 0$. \square

The key assumption is the differentiability of the mapping ψ , that we get, in our context, simply by construction.

In spite of the infinite dimension, the space Θ and the new optimization problem (14) have simple-looking structures, and they share many properties with finite-dimensional analogs. In particular, the optimality conditions are as expected.

Proposition E.2. *The three following properties are equivalent*

- (i) $\bar{\theta}$ is an optimal solution to problem (14)
- (ii) $\forall i \in \mathcal{I} : \frac{\partial \psi_\lambda}{\partial \theta_i}(\bar{\theta}) \geq 0$
and $\forall i \in \text{supp}(\bar{\theta}) : \frac{\partial \psi_\lambda}{\partial \theta_i}(\bar{\theta}) = 0$
- (iii) $\min_{i \in \mathcal{I}} \frac{\partial \psi_\lambda}{\partial \theta_i}(\bar{\theta}) \geq 0$
and $\bar{\theta} \in \text{Arg min}_{\theta \in \mathbb{R}_+^{\text{supp}(\bar{\theta})}} \psi_\lambda(\theta)$

Proof. To prove the equivalence between (i) and (ii), we apply both implications of Prop. E.1 with $\psi = \psi_\lambda$ and $K = \Theta^+$. We show first (i) \Leftarrow (ii). Let $\theta \in K$; we have

$$\sum_{i \in \mathcal{I}} \frac{\partial \psi}{\partial \theta_i}(\theta^*)(\theta_i - \theta_i^*) = \sum_{i \notin \text{supp}(\theta^*)} \frac{\partial \psi}{\partial \theta_i}(\theta^*)\theta_i \leq 0.$$

This is the optimality condition of (14), so we have (i).

We now prove the converse (i) \Rightarrow (ii). For all $i \in \mathcal{I}$, we write the optimality condition with $\eta \in \Theta$ defined by

$$\eta_\ell = \begin{cases} \theta_\ell^* & \text{if } \ell \neq i \\ \theta_i^* + 1 & \text{otherwise} \end{cases}$$

to get $\frac{\partial \psi}{\partial \theta_i}(\theta^*) \geq 0$. Similarly for all $i \in \mathcal{I}$ such that $\theta_i^* > 0$, we write the optimality condition with $\eta \in \Theta$ defined by

$$\eta_\ell = \begin{cases} \theta_\ell^* & \text{if } \ell \neq i \\ \theta_i^*/2 & \text{otherwise} \end{cases}$$

to get $\frac{\partial \psi}{\partial \theta_i}(\theta^*) \leq 0$, and we can conclude. \square