



**HAL**  
open science

## A response to “Fast OSCAR and OWL Regression via Safe Screening Rules” by Bao et al.

Clément Elvira, Cédric Herzet

► **To cite this version:**

Clément Elvira, Cédric Herzet. A response to “Fast OSCAR and OWL Regression via Safe Screening Rules” by Bao et al.. CentraleSupélec; Inria Rennes – Bretagne Atlantique. 2021. hal-03875689

**HAL Id: hal-03875689**

**<https://centralesupelec.hal.science/hal-03875689>**

Submitted on 29 Nov 2022

**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.

# A response to “Fast OSCAR and OWL Regression via Safe Screening Rules” by Bao *et al.*

Clément Elvira and Cédric Herzet

## I. INTRODUCTION

In this note, we discuss a recent contribution [2] of Bao *et al.* which addressed the so-called “Ordered Weighted L-One Linear Regression” (OWL) problem:

$$\min_{\mathbf{x} \in \mathbb{R}^n} (\mathbf{x}) \triangleq \frac{1}{2} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2 + r_{\text{OWL}}(\mathbf{x}) \quad (1\text{-OWL})$$

where

$$r_{\text{OWL}}(\mathbf{x}) \triangleq \sum_{k=1}^n \gamma_k |\mathbf{x}|_{[k]} \quad (2)$$

with

$$\gamma_1 > 0, \quad \gamma_1 \geq \dots \geq \gamma_n \geq 0, \quad (\text{H-3})$$

and  $|\mathbf{x}|_{[k]}$  refers to the  $k$ th largest element of  $\mathbf{x}$  in absolute value.

In their paper, the authors extend the so-called “safe screening” methodology introduced by El Ghaoui [5] to the OWL problem. The concept of safe screening has already been well-studied in the literature for “sparsity-promoting” convex problems (see [6] and references therein): it consists in devising simple “tests” that allow to detect zero entries in the solution of an optimization problem. The terminology “safe” refers to the fact that all the elements passing the test are *guaranteed* to correspond to zeros of the solution.

Nevertheless, we emphasize in this note that, contrarily to what is claimed by the authors, the test proposed in [2] does not satisfy the “safe” property. In particular, we provide two counterexamples in which the screening tests proposed in [2] leads to the screening of some *nonzero* coefficients. These counterexamples question the correctness of their result and suggest the presence of technical flaws in their derivations.

The note is organized as follows. In Section II, we introduce the notational conventions used throughout the document. Section III presents the rationale of the safe screening proposed by Bao *et al.* in [2] and recall their main result. Section IV discusses the correctness of this result by leveraging two counterexamples. All technical proofs are postponed to Appendices A and B.

C. Elvira is with SCEE/IETR UMR CNRS 6164, CentraleSupélec, Cesson Sévigné, France. [clement.elvira@centralesupelec.fr](mailto:clement.elvira@centralesupelec.fr), <https://c-elvira.github.io/>.

C. Herzet is with Inria centre Rennes - Bretagne Atlantique, Rennes, France. [cedric.herzet@inria.fr](mailto:cedric.herzet@inria.fr), <http://people.rennes.inria.fr/Cedric.Herzet/>

## II. NOTATIONAL CONVENTIONS

We will use the following notations throughout the paper. Vectors are denoted by lowercase bold letters (*e.g.*,  $\mathbf{x}$ ) and matrices by uppercase bold letters (*e.g.*,  $\mathbf{A}$ ). The “all-zero” (respectively “all-ones”) vector of dimension  $n$  is written  $\mathbf{0}_n$  (resp.  $\mathbf{1}_n$ ). Similarly,  $\mathbf{1}_{m \times n}$  denotes the “all-one” matrix of  $\mathbb{R}^{m \times n}$ . We use symbol  $^T$  to denote the transpose of a vector or a matrix.  $x_j$  refers to the  $j$ th component of  $\mathbf{x}$ . When referring to the sorted entries of a vector, we use bracket subscripts; more precisely, the notation  $\mathbf{x}_{[k]}$  refers to the  $k$ th largest value of  $\mathbf{x}$ . For matrices, we use  $\mathbf{a}_j$  to denote the  $j$ th column of  $\mathbf{A}$ . We use the notation  $|\mathbf{x}|$  to denote the vector made up of the absolute value of the components of  $\mathbf{x}$ . The sign function is defined for all scalars  $x$  as  $\text{sign}(x) = x/|x|$  with the convention  $\text{sign}(x) = 0$ .

Calligraphic letters are used to denote sets (*e.g.*,  $\mathcal{J}$ ). If  $a < b$  are two integers,  $[[a, b]]$  is used as a shorthand notation for the set  $\{a, a+1, \dots, b\}$ .

For all  $\mathbf{x} \in \mathbb{R}^n$  and  $j \in [[1, n]]$ ,  $r(\mathbf{x}, j)$  is such that  $\mathbf{x}_{(j)} = \mathbf{x}_{[r(\mathbf{x}, j)]}$ , that is it refers to the position of the  $j$ th entry of  $\mathbf{x}$  when  $\mathbf{x}$  is sorted. For instance, if  $\mathbf{x} = [2, 3, 1, 0]$   $r(\mathbf{x}, 1) = 2$ ,  $r(\mathbf{x}, 2) = 1$ ,  $r(\mathbf{x}, 3) = 3$  and  $r(\mathbf{x}, 4) = 4$ . If some of the coordinates of  $\mathbf{x}$  are equal, we break ties arbitrarily. In the sequel, we will refer to  $r$  as the “rank” function.

## III. RATIONALE OF SAFE SCREENING RULE FOR OWL

In this section, we summarize the main ingredients grounding the safe screening rules for OWL proposed in [2]. To distinguish between our derivations and the results of [2], all equation numbers referring to a result of the latter paper are prefixed by “B-” (*e.g.*, (B-7)).

The screening procedure proposed in [2] leverages the solution of the Fenchel dual problem of OWL. In particular, Bao *et al.* claim that [2, Eq. (6d)]

$$\mathbf{u}^* = \arg \max_{\mathbf{u} \in \mathcal{U}} D(\mathbf{u}) \triangleq -\frac{1}{2} \|\mathbf{u}\|_2^2 - \mathbf{u}^T \mathbf{y}, \quad (\text{B-4})$$

is the (Fenchel) dual problem to (1-OWL) where

$$\mathcal{U} = \{\mathbf{u} \in \mathbb{R}^m : \forall j \in [[1, n]], |\mathbf{a}_j^T \mathbf{u}| \leq \gamma_{r(\mathbf{x}, j)}\} \quad (\text{B-5})$$

is the dual feasible set and  $\mathbf{x}^* \in \mathbb{R}^n$  is a minimizer of (1-OWL). At this stage, let us note that  $\mathbf{u}^*$  exists and is unique (*i.e.*, the equality in (B-4) is well-defined) because  $D$  is a continuous strongly-concave function and  $\mathcal{U}$  a closed set.

Besides,  $\mathbf{u}^*$  is connected to any minimizer  $\mathbf{x}^*$  of (1-OWL) through the identity

$$\mathbf{u}^* = \mathbf{y} - \mathbf{A}\mathbf{x}^* \quad (6)$$

as a consequence of standard primal-dual optimality conditions [3, Th. 19.1].

We now describe the safe screening test for OWL proposed in [2]. One easily verifies that (1-OWL) admits at least one minimizer so that the screening problem is well-posed. Consider therefore a minimizer  $\mathbf{x}^*$  of (1-OWL). Bao *et al.* claim that<sup>1</sup> [2, Eq. (11)]

$$|\mathbf{a}_\ell^T \mathbf{u}^*| < \gamma_{r(\mathbf{x}^*, \ell)} \implies \mathbf{x}_\ell^* = 0. \quad (\text{B-7})$$

In other words, (B-7) describes the following sufficient condition for *safe screening*: if the inner product between the  $\ell$ th column of  $\mathbf{A}$  and  $\mathbf{u}^*$  – the maximizer of the dual problem (B-4) – is lower than some threshold, the  $\ell$ th entry of any minimizer of OWL can be set to zero. We have nevertheless identified some technical flaws in the derivation of (B-7) which prevent test (B-7) from being safe. This will be discussed in Section IV.

We may notice that (B-7) requires the knowledge of both  $\mathbf{u}^*$  and the ordering of a solution of OWL. On the one hand, computing  $\mathbf{u}^*$  is usually as difficult as solving (1-OWL). On the other hand, the ordering of  $\mathbf{x}^*$  is obviously not known beforehand. As a consequence, (B-7) is of poor practical interest and Bao *et al.* devised a relaxed version of their test.

The relaxed test leverages the following two ingredients. First, since  $\gamma_{r(\mathbf{x}^*, \ell)} \geq \gamma_n$  for all index  $\ell \in \llbracket 1, n \rrbracket$  by (H-3), one immediately deduces that

$$|\mathbf{a}_\ell^T \mathbf{u}^*| < \gamma_n \implies \mathbf{x}_\ell^* = 0. \quad (\text{B-8})$$

Second, the knowledge of  $\mathbf{u}^*$  can be circumvented if one has access to a couple of primal / dual feasible vectors. In particular if  $\mathbf{x} \in \mathbb{R}^n$  and  $\mathbf{u} \in \mathcal{U}$  are known, the test simplifies as

$$|\mathbf{a}_\ell^T \mathbf{c}| + \|\mathbf{a}_\ell\|_2 \sqrt{2(P(\mathbf{x}) - D(\mathbf{u}))} < \gamma_n \implies \mathbf{x}_\ell^* = 0, \quad (\text{B-9})$$

see [2, Eq. (19)]. Bao *et al.* then proposed to interleave the relaxed screening test (B-9) with the iterations of several iterative algorithms to assess its effectiveness.

#### IV. DISCUSSION

In this section, we discuss the arguments developed in [2] to show the safeness of test (B-8). In particular, Bao *et al.* ground their reasoning on the following optimality conditions

$$\mathbf{a}_j^T \mathbf{u}^* + \text{sign}(\mathbf{x}_j^*) \gamma_{r(\mathbf{x}^*, j)} = 0 \text{ if } \mathbf{x}_j^* \neq 0 \quad (\text{B-10a})$$

$$|\mathbf{a}_j^T \mathbf{u}^*| \leq \gamma_{r(\mathbf{x}^*, j)} \text{ otherwise,} \quad (\text{B-10b})$$

and used the series of implications

$$(\text{B-10b}) \implies (\text{B-7}) \implies (\text{B-8}) \quad (11)$$

<sup>1</sup>Note that the authors did not mention how to handle the design of the rank function in presence of multiple maximizers.

to prove the safeness of (B-8). We show below that neither (B-10b) nor (B-7) characterize the minimizers of generic OWL problems, thus invalidating their proof of safeness.

We first provide a counterexample where a minimizer of OWL violates (B-7):

**Counterexample 1.** Consider the OWL problem with dimensions  $m, n \geq 2$ , parameters  $\mathbf{y} = \mathbf{1}_m$ ,  $\mathbf{A} = \mathbf{1}_{m \times n}$  and weighting coefficients  $\{\gamma_k\}_{k=1}^n$  such that

$$1 \geq \gamma_1 > \gamma_2 > \dots > \gamma_n > 0. \quad (12)$$

In this setup, we show in Appendix A that the unique solution  $\mathbf{x}^*$  of the OWL problem writes  $\mathbf{x}^* = x^* \mathbf{1}_n$  where

$$x^* = \frac{mn - \sum_{k=1}^n \gamma_k}{mn^2} > 0. \quad (13)$$

As a consequence of (6), we have

$$\mathbf{A}^T \mathbf{u}^* = \frac{\sum_{k=1}^n \gamma_k}{n} \mathbf{1}_n. \quad (14)$$

Finally, let  $\ell \in \llbracket 1, n \rrbracket$  be an index such that  $r(\mathbf{x}^*, \ell) = 1$ . Using (12) we obtain

$$|\mathbf{a}_\ell^T \mathbf{u}^*| < \gamma_1. \quad (15)$$

Hence  $\mathbf{x}_\ell^* = 0$  by (B-7) which leads to a contradiction.

Counterexample 1 describes a scenario where both the minimizer  $\mathbf{x}^*$  of OWL and the maximizer  $\mathbf{u}^*$  of the dual problem can be written in closed-form so that the safeness of screening test (B-7) can be easily assessed. We note that, although all entries of  $\mathbf{x}^*$  are positive, (at least) one entry of  $\mathbf{x}^*$  is screened by test (B-7). This shows that implication (B-7) is violated in this particular setup.

We next provide a counter-example showing that (B-10a) and (B-10b) do not correctly describe the optimality condition of (1-OWL):

**Counterexample 2.** Consider the OWL problem with dimensions  $m = n \geq 2$ , parameters  $\mathbf{y} = \mathbf{1}_n$ ,  $\mathbf{A} = \mathbf{I}_n$  the identity matrix of  $\mathbb{R}^{n \times n}$  and weighting coefficients  $\{\gamma_k\}_{k=1}^n$  satisfying

$$\gamma_1 = 1 \text{ and } \gamma_k = 0 \forall k \in \llbracket 2, n \rrbracket. \quad (16)$$

Then, the OWL problem (1-OWL) reduces to

$$\min_{\mathbf{x} \in \mathbb{R}^n} \frac{1}{2} \|\mathbf{1}_n - \mathbf{x}\|_2^2 + \|\mathbf{x}\|_\infty. \quad (17)$$

We show in Appendix B that the solution of (17) is unique and writes

$$\mathbf{x}^* = \frac{n-1}{n} \mathbf{1}_n. \quad (18)$$

One also has (as a consequence of (6))

$$\mathbf{u}^* = \frac{1}{n} \mathbf{1}_n. \quad (19)$$

Hence, all entries of  $\mathbf{x}^*$  are positive. Yet, for all  $j \in \llbracket 1, n \rrbracket$ ,

$$\mathbf{a}_j^T \mathbf{u}^* + \text{sign}(\mathbf{x}_j^*) \gamma_{r(\mathbf{x}^*, j)} \geq \frac{1}{n} > 0 \quad (20)$$

which is in contradiction with (B-10a).

Similarly to Counterexample 1, Counterexample 2 describes a simple setup where both  $\mathbf{x}^*$  and  $\mathbf{u}^*$  are unique and can

be exhibited in closed-form. However, it appears that the latter couple of primal-dual solutions does not satisfies the optimality conditions derived in [2] for the OWL problem.

Prior to conclude our note, we provide some elements regarding the fallacy in the derivation of the optimality conditions (B-10a) and (B-10b). An attentive reading of [2, Section 3.1] suggests that Bao *et al.* implicitly assume that addressing the following optimization problem (see [2, Eq. (8)])

$$\min_{\mathbf{x} \in \mathbb{R}^n} \mathbf{x}^T \mathbf{A}^T \mathbf{u}^* + \sum_{k=1}^n \gamma_k |\mathbf{x}|_{[k]} \quad (\text{B-21})$$

is equivalent to solve

$$\min_{\mathbf{x} \in \mathbb{R}^n} \mathbf{x}^T \mathbf{A}^T \mathbf{u}^* + \sum_{j=1}^n \gamma_{\varepsilon(\mathbf{x}^*, j)} |\mathbf{x}_j|, \quad (\text{B-22})$$

where  $\mathbf{x}^*$  denotes a minimizer of (1-OWL). In other words, the authors assume that if the order of the entries (in absolute value) of the solution of (B-21) is known, then one can equivalently solve a simpler optimization problem that already takes into account the permutation of the entries. Besides, the authors did not mention how to handle the design of the rank function in presence of multiple minimizers.

To support our disagreement with this statement, consider again the setup described in Counterexample 2. Since all entries of  $\mathbf{x}^*$  are equal (and positive),  $[1, \dots, n]$  is a valid ordering. Yet, solving (B-21) is not equivalent to solving

$$\min_{\mathbf{x} \in \mathbb{R}^n} \mathbf{x}^T \mathbf{A}^T \mathbf{u}^* + \gamma_1 |\mathbf{x}_1| \quad (\text{B-23})$$

which admits  $-\infty$  as minimum.

## APPENDIX A PROOF RELATED TO COUNTEREXAMPLE 1

In this appendix, we exhibit the (unique) solution of the OWL problem from Counterexample 1. Our proof is organized as follows. We first state in Section A-A a technical lemma. We show in Section A-B that, in this example, all minimizers of OWL have nonnegative entries. Then, we demonstrate in Section A-C that all minimizers are proportional to the ‘‘all-one’’ vector. We finally solve in Section A-D the OWL problem and show that the minimizer is in fact unique.

In the following, we let  $\mathbf{y} \in \mathbb{R}^m$ ,  $\mathbf{A} \in \mathbb{R}^{m \times n}$  and  $\{\gamma_k\}_{k=1}^n$  be defined as in Counterexample 1.

### A. A technical lemma

The following results will be used in our subsequent derivations:

**Lemma 1.** *Let  $\mathbf{u}^* \in \mathbb{R}^m$  be the (unique) maximizer of the dual problem of OWL. Then, in the setup of Counterexample 1, we have:*

$$\mathbf{1}_m^T \mathbf{u}^* \geq 0. \quad (\text{24})$$

*Proof.* See first that the (primal) objective function  $P$  rewrites  $\forall \mathbf{x} \in \mathbb{R}^n$  as  $P(\mathbf{x}) = f(\mathbf{A}\mathbf{x}) + r_{\text{OWL}}(\mathbf{x})$  and recall that  $r_{\text{OWL}}$

defines a norm under (H-3), see [4, Proposition 1.1]. Hence, its Fenchel dual problem writes [1, Eq. (1.4)]

$$\begin{aligned} \arg \max_{\mathbf{u} \in \mathbb{R}^m} D(\mathbf{u}) = f^*(\mathbf{u}) &= \frac{1}{2} \|\mathbf{y}\|_2^2 - \frac{1}{2} \|\mathbf{y} - \mathbf{u}\|_2^2 \\ &\text{s.t. } r_{\text{OWL}}^*(\mathbf{A}^T \mathbf{u}) \leq 1 \end{aligned} \quad (\text{25})$$

where  $f^*$  denotes the Fenchel conjugate function of  $f$  and  $r_{\text{OWL}}^*$  the dual norm of  $r_{\text{OWL}}$ .

Since  $r_{\text{OWL}}^*$  is also a norm, we have  $r_{\text{OWL}}^*(\mathbf{A}^T \mathbf{0}_m) = 0 \leq 1$  so that the vector of zeros  $\mathbf{0}_m$  is admissible. Using the fact that  $f^*$  is concave and that  $\nabla f^*(\mathbf{0}_m) = \mathbf{y}$ , we have for all  $\mathbf{u} \in \mathbb{R}^m$ :

$$f^*(\mathbf{u}) \leq f^*(\mathbf{0}_m) + \mathbf{y}^T \mathbf{u}. \quad (\text{26})$$

Recalling that  $\mathbf{y} = \mathbf{1}_m$  and particularizing the latter inequality to  $\mathbf{u} = \mathbf{u}^*$ , we obtain

$$\mathbf{1}_m^T \mathbf{u}^* \geq f^*(\mathbf{u}^*) - f^*(\mathbf{0}_m) \geq 0 \quad (\text{27})$$

since  $\mathbf{u}^*$  is the maximizer of (25). This concludes the proof.  $\square$

### B. All minimizers have positive entries

Let  $\mathbf{x}^* \in \mathbb{R}^n$  be a minimizer of OWL. In the following, we show by contradiction that  $\mathbf{x}_j^* \geq 0$  for all  $j \in \llbracket 1, n \rrbracket$ .

Assume that there exists  $j_0 \in \llbracket 1, n \rrbracket$  such that  $\mathbf{x}_{j_0}^* < 0$ . Suppose moreover (without loss of generality and up to a permutation of the indices) that the entries of  $\mathbf{x}^*$  are ordered (in absolute value), *i.e.*, that

$$|\mathbf{x}_1^*| \geq |\mathbf{x}_2^*| \geq \dots \geq |\mathbf{x}_n^*| \quad (\text{28})$$

and that

$$|\mathbf{x}_{j_0}^*| > |\mathbf{x}_{j_0+1}^*| \quad (\text{29})$$

with the convention  $\mathbf{x}_{n+1}^* = 0$ . Define  $\mathbf{x}' \in \mathbb{R}^n$  as

$$\mathbf{x}' = \mathbf{x}^* + \varepsilon_1 \mathbf{e}_{j_0} \quad (\text{30})$$

where  $\mathbf{e}_{j_0}$  denote the  $j_0$ th vector of the canonical basis of  $\mathbb{R}^n$  and  $\varepsilon_1$  satisfies

$$0 < \varepsilon_1 < \min \left( |\mathbf{x}_{j_0}^*| - |\mathbf{x}_{j_0+1}^*|, \frac{2\gamma_{j_0}}{n} \right). \quad (\text{31})$$

Note that such  $\varepsilon_1$  exists since  $\gamma_{j_0} > 0$  by hypothesis and that  $|\mathbf{x}_{j_0}^*| > |\mathbf{x}_{j_0+1}^*|$ . Moreover, the ordering of  $\mathbf{x}'$  is also conserved, that is

$$|\mathbf{x}'_1| \geq |\mathbf{x}'_2| \geq \dots \geq |\mathbf{x}'_n| \quad (\text{32})$$

and  $|\mathbf{x}'_{j_0}| = |\mathbf{x}_{j_0}^*| - \varepsilon_1$  by design of  $\varepsilon_1$ . Hence,

$$r_{\text{OWL}}(\mathbf{x}') = r_{\text{OWL}}(\mathbf{x}^*) - \gamma_{j_0} \varepsilon_1. \quad (\text{33})$$

Using the fact that  $\mathbf{A} \mathbf{e}_{j_0} = \mathbf{a}_{j_0} = \mathbf{1}_m$  by definition of  $\mathbf{A}$ , we have

$$\begin{aligned} P(\mathbf{x}') &= P(\mathbf{x}^*) + \frac{\varepsilon_1^2}{2} \|\mathbf{1}_m\|_2^2 - \varepsilon_1 \mathbf{1}_m^T (\mathbf{y} - \mathbf{A} \mathbf{x}^*) - \gamma_{j_0} \varepsilon_1 \\ &\stackrel{(6)}{=} P(\mathbf{x}^*) + \frac{\varepsilon_1^2}{2} \|\mathbf{1}_m\|_2^2 - \varepsilon_1 \mathbf{1}_m^T \mathbf{u}^* - \gamma_{j_0} \varepsilon_1 \\ &\stackrel{\text{Lemma 1}}{\leq} P(\mathbf{x}^*) + \varepsilon_1 \left( \frac{n}{2} \varepsilon_1 - \gamma_{j_0} \right). \end{aligned} \quad (\text{34})$$

By choice of  $\varepsilon_1$  (see (31)) we have

$$\varepsilon_1 \left( \frac{n}{2} \varepsilon_1 - \gamma_{j_0} \right) < 0 \quad (35)$$

so that  $P(\mathbf{x}') < P(\mathbf{x}^*)$  which is the desired contradiction. Therefore, the entries of all minimizers of OWL have non-negative entries.

### C. All minimizers are proportional to $\mathbf{1}_n$

We now show that all the solutions of OWL are proportional to vector  $\mathbf{1}_n$ . We again proceed by contradiction: let  $\mathbf{x}^*$  be a minimizer of (1-OWL) and assume that there exists  $j_0 \in \llbracket 1, n-1 \rrbracket$  such that  $\mathbf{x}_{j_0}^* > \mathbf{x}_{j_0+1}^*$  (recall that all entries of  $\mathbf{x}^*$  are nonnegative). Assume moreover and without loss of generality (up to a permutations of the columns of  $\mathbf{A}$ ) that

$$\mathbf{x}_1^* \geq \mathbf{x}_2^* \geq \dots \geq \mathbf{x}_n^*. \quad (36)$$

We follow a similar rationale as in Section A-B: define the vector  $\mathbf{x}' \in \mathbb{R}^n$  such that

$$\mathbf{x}' = \mathbf{x}^* - \varepsilon_2 \mathbf{e}_{j_0} + \varepsilon_2 \mathbf{e}_{j_0+1} \quad (37)$$

where  $\mathbf{e}_{j_0}, \mathbf{e}_{j_0+1}$  denote the  $j_0, j_0+1$ th vectors of the canonical basis of  $\mathbb{R}^n$  and  $\varepsilon_2$  is a scalar such that

$$0 < \varepsilon_2 < \frac{\mathbf{x}_{j_0}^* - \mathbf{x}_{j_0+1}^*}{2}. \quad (38)$$

Note that such a scalar exists since  $\mathbf{x}_{j_0}^* > \mathbf{x}_{j_0+1}^*$  by assumption. See also that  $\mathbf{x}'$  is nonnegative by definition, and that

$$\mathbf{x}'_1 \geq \mathbf{x}'_2 \geq \dots \geq \mathbf{x}'_n \geq 0 \quad (39)$$

Hence,

$$r_{\text{OWL}}(\mathbf{x}') = r_{\text{OWL}}(\mathbf{x}^*) + \varepsilon_2 (\gamma_{j_0+1} - \gamma_{j_0}) < r_{\text{OWL}}(\mathbf{x}^*) \quad (40)$$

by definition of  $\varepsilon_2$  and the sequence of weights  $\{\gamma_k\}_{k=1}^n$ .

Recalling that  $\mathbf{A}\mathbf{e}_j = \mathbf{1}_m$  for all  $j \in \llbracket 1, n \rrbracket$ , we have

$$\begin{aligned} P(\mathbf{x}') &= \frac{1}{2} \|\mathbf{y} - \mathbf{A}\mathbf{x}'\|_2^2 + r_{\text{OWL}}(\mathbf{x}') \\ &\stackrel{(40)}{<} \frac{1}{2} \|\mathbf{y} - \mathbf{A}\mathbf{x}^*\|_2^2 + r_{\text{OWL}}(\mathbf{x}^*) \\ &= P(\mathbf{x}^*) \end{aligned} \quad (41)$$

which leads to the desired contradiction. Hence all solutions of OWL are proportional to the vector of ones.

### D. Solving the OWL problem

Combining the results of Sections A-B and A-C, we have that solving (1-OWL) is equivalent<sup>2</sup> to solve

$$\arg \min_{x \in \mathbb{R}_+} \tilde{P}(x) \triangleq \frac{1}{2} \|\mathbf{y} - x\mathbf{A}\mathbf{1}_n\|_2^2 + x \sum_{k=1}^n \gamma_k. \quad (42)$$

We let the reader check that the scalar  $x^*$  defined in (13) satisfies  $\tilde{P}'(x^*) = 0$ . Moreover, one sees that  $x^* > 0$  using the fact that  $\gamma_1 \leq 1$  and (H-3). Therefore,  $x^*$  satisfies the optimality condition of (the strictly-convex) optimization problem (42) and is therefore its unique minimizer.

<sup>2</sup>In the sense that there exists a bijection between the set of minimizers.

## APPENDIX B

### PROOF RELATED TO COUNTEREXAMPLE 2

Let  $\mathbf{y} \in \mathbb{R}^m$ ,  $\mathbf{A} \in \mathbb{R}^{m \times n}$  and  $\{\gamma_k\}_{k=1}^n$  be defined as in Counterexample 2. We first note that, for such choices of the parameters, OWL reduces to (17) since  $r_{\text{OWL}}$  simplifies as

$$\forall \mathbf{x} \in \mathbb{R}^n, \quad r_{\text{OWL}}(\mathbf{x}) = \|\mathbf{x}\|_{[1]} = \max_j \{|\mathbf{x}_j|\} = \|\mathbf{x}\|_\infty. \quad (43)$$

Note also that the solution of (17) exists and is unique because the function to minimize is continuous, strongly-convex and  $\mathbb{R}^n$  is a nonempty, convex and closed set.

Let  $\mathbf{x}^* = (n-1)/n \mathbf{1}_n$ . Since the function to minimize is lower-semicontinuous, a direct application of Fermat's rule [3, Theorem 16.3] to problem (45) shows that  $\mathbf{x}^*$  is a minimizer of (17) if and only if

$$\mathbf{y} - \mathbf{x}^* = \frac{1}{n} \mathbf{1}_n \in \partial \|\mathbf{x}^*\|_\infty. \quad (44)$$

Since the function  $\mathbf{x} \rightarrow \|\mathbf{x}\|_\infty$  defines a norm, its subdifferential is well defined for all  $\mathbf{x} \in \mathbb{R}^n$  and writes [1, Eq. (1.4)]:

$$\partial \|\mathbf{x}\|_\infty = \{\mathbf{g} \in \mathbb{R}^n : \|\mathbf{g}\|_1 \leq 1 \text{ and } \mathbf{x}^T \mathbf{g} = \|\mathbf{x}\|_\infty\}. \quad (45)$$

Let us show now that  $\mathbf{y} - \mathbf{x}^* \in \partial \|\mathbf{x}^*\|_\infty$ . Since  $\mathbf{y} - \mathbf{x}^*$  has nonnegative entries, we first have

$$\|\mathbf{y} - \mathbf{x}^*\|_1 = \frac{1}{n} \|\mathbf{1}_n\|_1 = 1. \quad (46)$$

See then that

$$\mathbf{x}^{*T} (\mathbf{y} - \mathbf{x}^*) = \frac{n-1}{n} = \|\mathbf{x}^*\|_\infty. \quad (47)$$

Hence,  $\mathbf{y} - \mathbf{x}^* \in \partial \|\mathbf{x}^*\|_\infty$  which concludes the proof.

## REFERENCES

- [1] F. BACH, R. JENATTON, J. MAIRAL, AND G. OBOZINSKI, *Convex optimization with sparsity-inducing norms*, in Optimization for Machine Learning, Neural information processing series, MIT Press, 2011, pp. 19–49, <https://doi.org/10.7551/mitpress/8996.003.0004>.
- [2] R. BAO, B. GU, AND H. HUANG, *Fast OSCAR and OWL with safe screening rules*, in Proceedings of the 37th International Conference on Machine Learning, 2020.
- [3] H. H. BAUSCHKE AND P. L. COMBETTES, *Convex Analysis and Monotone Operator Theory in Hilbert Spaces*, Springer International Publishing, 2017, <https://doi.org/10.1007/978-3-319-48311-5>.
- [4] M. BOGDAN, E. VAN DEN BERG, W. SU, AND E. CANDÈS, *Statistical estimation and testing via the sorted  $l_1$  norm*, 2013, <https://arxiv.org/abs/1310.1969>.
- [5] L. GHAOUI, V. VIALON, AND T. RABBANI, *Safe feature elimination in sparse supervised learning*, Pacific Journal of Optimization, 8 (2010).
- [6] E. NDIAYE, O. FERCOQ, ALEX, RE GRAMFORT, AND J. SALMON, *Gap safe screening rules for sparsity enforcing penalties*, Journal of Machine Learning Research, 18 (2017), pp. 1–33, <http://jmlr.org/papers/v18/16-577.html>.