

## A. Preliminaries on Convex Envelopes and Subdifferentials

In this section we review some basic facts about convex envelopes and their sub-differentials that we will use in our theory. Throughout the section we will assume that any infimum is attained. This is true for example if the function is lower semi continuous with bounded level sets, which is the case for our objective function  $f(\mathbf{x}) = g(\text{card}(\mathbf{x})) + \|\mathbf{x}\|^2$ . In addition the quadratic term grows faster than any linear term of the type  $\langle \mathbf{x}, \mathbf{y} \rangle$  and therefore this is also true when we add linear terms.

The convex envelope  $f^{**}$  of a function  $f$  is the largest convex function that fulfills  $f^{**} \leq f$ . For a convex function we should have  $f^{**}(\sum_i \lambda_i \mathbf{x}^j) \leq \sum_j \lambda_j f^{**}(\mathbf{x}^j)$ ,  $0 \leq \lambda_j$ ,  $\sum_j \lambda_j = 1$ . At a point  $\mathbf{x} = \sum_j \lambda_j \mathbf{x}^j$  where  $f(\mathbf{x}) > \sum_j \lambda_j f(\mathbf{x}^j)$  we compute the value  $f^{**}(\mathbf{x})$  by minimizing over convex combinations of points using  $f^{**}(\mathbf{x}) =$

$$\min \left\{ \sum_{j=1}^{d+1} \lambda_j f(\mathbf{x}^j); \quad \sum_{j=1}^{d+1} \lambda_j \mathbf{x}^j = \mathbf{x}, \quad \sum_{j=1}^{d+1} \lambda_j = 1, \quad \lambda_j > 0 \right\}. \quad (20)$$

It can be shown (using Caratheodory's Theorem) that it is enough to consider combinations of  $d + 1$  points if  $\mathbf{x} \in \mathbb{R}^d$ . Figure 4 shows one example of convex envelope. Here the two functions  $f^{**}$  and  $f$  coincide at  $x = 0$  and  $|x| > 1$ . If  $x \in (0, 1)$  where the functions differ and the value of  $f^{**}$  is computed using the convex combination  $(1 - x)f(0) + xf(1)$ . Note that when  $f$  and  $f^{**}$  differs the function  $f^{**}$  will be affine in some direction.

An alternative way of computing  $f^{**}$  is using supporting hyperplanes and the conjugate function

$$f^*(\mathbf{y}) = \max_{\mathbf{x}} \langle \mathbf{x}, \mathbf{y} \rangle - f(\mathbf{x}). \quad (21)$$

From the definition it is clear that

$$f^*(\mathbf{y}) \geq \langle \mathbf{x}, \mathbf{y} \rangle - f(\mathbf{x}), \quad (22)$$

for all  $\mathbf{x}$ . Rearranging terms we get

$$f(\mathbf{x}) \geq \langle \mathbf{x}, \mathbf{y} \rangle - f^*(\mathbf{y}), \quad (23)$$

which is an affine function in  $\mathbf{x}$  and therefore a supporting hyperplane to  $f$ . Figure 4 shows three supporting hyperplanes for  $f$ . Note that these touch  $f^{**}$  in (at least) one point. For each  $\mathbf{x}$  we can find a hyperplane that touches  $f^{**}(\mathbf{x})$  which means that

$$f^{**}(\mathbf{x}) = \max_{\mathbf{y}} \langle \mathbf{x}, \mathbf{y} \rangle - f^*(\mathbf{y}), \quad (24)$$

that is, the convex envelope is the conjugate of the conjugate function.

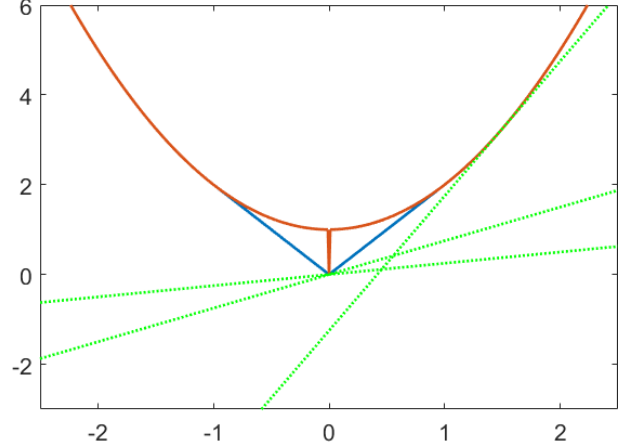


Figure 4: An example of convex envelope. Here  $f^{**}(x) = \mu - \max(\sqrt{\mu} - |x|, 0)^2 + x^2$  (blue curve) and  $f(x) = \mu \text{card}(x) + x^2$  (orange curve). Green dotted lines are supporting hyperplanes of the form  $xy - f^*(y)$  for three different  $y$ .

For a convex function  $f^{**}$  the set of sub-gradients  $\partial f^{**}(\mathbf{x})$  at a point  $\mathbf{x}$  is defined as all vectors  $\mathbf{y}$  such that

$$f^{**}(\mathbf{x}') - f^{**}(\mathbf{x}) \geq \langle \mathbf{y}, \mathbf{x}' - \mathbf{x} \rangle, \quad \forall \mathbf{x}'. \quad (25)$$

or equivalently

$$\langle \mathbf{y}, \mathbf{x} \rangle - f^{**}(\mathbf{x}) \geq \langle \mathbf{y}, \mathbf{x}' \rangle - f^{**}(\mathbf{x}'), \quad \forall \mathbf{x}'. \quad (26)$$

Since we clearly have equality when  $\mathbf{x}' = \mathbf{x}$  this means that  $\langle \mathbf{y}, \mathbf{x} \rangle - f^{**}(\mathbf{x}) = \max_{\mathbf{x}'} \langle \mathbf{y}, \mathbf{x}' \rangle - f^{**}(\mathbf{x}') = f^{***}(\mathbf{x}) = f^*(\mathbf{x})$ . Rearranging terms shows that

$$\langle \mathbf{y}, \mathbf{x} \rangle - f^*(\mathbf{y}) = f^{**}(\mathbf{x}) = \max_{\mathbf{y}'} \langle \mathbf{y}', \mathbf{x} \rangle - f^*(\mathbf{y}'). \quad (27)$$

Thus the set of sub-gradients at a point  $\mathbf{x}$  are all the vectors  $\mathbf{y}$  that achieves the maximal value in the second conjugation. In points where  $f^{**}$  is non-differentiable the function has several sub-gradients. In a differentiable point the only sub-gradient is the standard gradient.

The following result does not appear to be standard but is crucial for our main theorem. Therefore we state it somewhat more formally below.

**Lemma A.1.** *Suppose that for a point  $\mathbf{x}$  we have  $f(\mathbf{x}) > f^{**}(\mathbf{x})$ . Then there is a set of points  $\{\mathbf{x}^j\}$  such that*

$$\mathbf{x} = \sum_j \lambda_j \mathbf{x}^j, \quad 0 \leq \lambda_j \leq 1, \quad \sum_j \lambda_j = 1, \quad (28)$$

$f^{**}(\mathbf{x}^j) = f(\mathbf{x}^j)$  and

$$f^{**}(\mathbf{x}) = \sum_j \lambda_j f(\mathbf{x}^j). \quad (29)$$

In addition  $\partial f^{**}(\mathbf{x}) \subset \bigcap_j \partial f^{**}(\mathbf{x}^j)$ .

*Proof.* Consider the convex combination  $\mathbf{x} = \sum_j \lambda_j \mathbf{x}^j$  that solves the minimization in (20).

We have that  $f(\mathbf{x}^j) \geq f^{**}(\mathbf{x}^j)$ . Assume further that  $f(\mathbf{x}^j) > f^{**}(\mathbf{x}^j)$  for some  $j$ . Then we have

$$f^{**}(\sum_j \lambda_j \mathbf{x}^j) = f^{**}(\mathbf{x}) = \sum_j \lambda_j f(\mathbf{x}^j) > \sum_j \lambda_j f^{**}(\mathbf{x}^j), \quad (30)$$

which contradicts the convexity of  $f^{**}$ . Therefore  $f(\mathbf{x}^j) = f^{**}(\mathbf{x}^j)$  for all  $j$ .

Now consider a subgradient  $\mathbf{y} \in \partial f^{**}(\mathbf{x})$ . By definition we have that

$$f^{**}(\mathbf{x}') \geq f^{**}(\mathbf{x}) + \langle \mathbf{y}, \mathbf{x}' - \mathbf{x} \rangle. \quad (31)$$

Now assume that

$$f^{**}(\mathbf{x}^j) > f^{**}(\mathbf{x}) + \langle \mathbf{y}, \mathbf{x}^j - \mathbf{x} \rangle, \quad (32)$$

for some  $j$ . Then we have

$$\underbrace{\sum_j \lambda_j f^{**}(\mathbf{x}^j)}_{=f^{**}(\mathbf{x})} > \underbrace{\sum_j \lambda_j f^{**}(\mathbf{x})}_{=f^{**}(\mathbf{x})} + \underbrace{\sum_j \lambda_j \langle \mathbf{y}, \mathbf{x}^j - \mathbf{x} \rangle}_{=0}, \quad (33)$$

which shows that we must have

$$f^{**}(\mathbf{x}^j) = f^{**}(\mathbf{x}) + \langle \mathbf{y}, \mathbf{x}^j - \mathbf{x} \rangle. \quad (34)$$

This gives us

$$f^{**}(\mathbf{x}^j) + \langle \mathbf{y}, \mathbf{x}' - \mathbf{x}^j \rangle = f^{**}(\mathbf{x}) + \langle \mathbf{y}, \mathbf{x}' - \mathbf{x} \rangle \leq f^{**}(\mathbf{x}'), \quad (35)$$

which shows that  $\mathbf{y} \in \partial f^{**}(\mathbf{x}^j)$  for all  $j$ .  $\square$

### A.1. The Conjugate of $f$

We now consider our class of functions  $f(\mathbf{x}) = g(\text{card}(\tilde{\mathbf{x}})) + \|\tilde{\mathbf{x}}\|^2$ . Consider the conjugate (21). Since  $f(\mathbf{x})$  only depends on  $\tilde{\mathbf{x}}$  and not the signs or ordering of the elements it is clear that the elements of  $\mathbf{x}$  should have the same sign and ordering as those in  $\mathbf{y}$  to maximize the term  $\langle \mathbf{x}, \mathbf{y} \rangle$ . Therefore we get

$$f^*(\mathbf{y}) = \max_{\tilde{\mathbf{x}}} \langle \tilde{\mathbf{x}}, \tilde{\mathbf{y}} \rangle - g(\text{card}(\tilde{\mathbf{x}})) - \|\tilde{\mathbf{x}}\|^2. \quad (36)$$

This can equivalently be written

$$\max_k \max_{\|\tilde{\mathbf{x}}\|_0=k} \langle \tilde{\mathbf{x}}, \tilde{\mathbf{y}} \rangle - \sum_{i=1}^k (g_i + \tilde{x}_i^2). \quad (37)$$

Completing squares gives

$$\max_k \max_{\|\tilde{\mathbf{x}}\|_0=k} -\|\tilde{\mathbf{x}} - \frac{1}{2}\tilde{\mathbf{y}}\|^2 + \frac{1}{4}\|\tilde{\mathbf{y}}\|^2 - \sum_{i=1}^k g_i. \quad (38)$$

It is clear that the inner maximization is solved by letting  $\tilde{x}_i = \frac{\tilde{y}_i}{2}$  if  $i \leq k$  and  $\tilde{x}_i = 0$  otherwise. After some simple manipulations this gives the conjugate function

$$f^*(\mathbf{y}) = \sum_{i=1}^n \max\left(\frac{1}{4}\tilde{y}_i^2 - g_i, 0\right). \quad (39)$$

Not that the computations for the matrix case are close to identical. In this case we maximize the scalar product  $\langle X, Y \rangle$  when  $X$  and  $Y$  SVDs with the same  $U$  and  $V$  matrices (von Neumann's trace theorem), in this case  $\langle X, Y \rangle = \langle \tilde{\mathbf{x}}, \tilde{\mathbf{y}} \rangle$ .

### A.2. The biconjugate of $f$

Taking the conjugate once more gives

$$f^{**}(\mathbf{x}) = \max_{\mathbf{y}} \langle \mathbf{x}, \mathbf{y} \rangle - \sum_{i=1}^n \max\left(\frac{1}{4}\tilde{y}_i^2 - g_i, 0\right). \quad (40)$$

Again the second term only depends on the elements of  $\tilde{\mathbf{y}}$  and therefore

$$f^{**}(\mathbf{x}) = \max_{\tilde{\mathbf{y}}} \langle \tilde{\mathbf{x}}, \tilde{\mathbf{y}} \rangle - \sum_{i=1}^n \max\left(\frac{1}{4}\tilde{y}_i^2 - g_i, 0\right). \quad (41)$$

For ease of notation we let  $\tilde{\mathbf{y}} = 2\tilde{\mathbf{z}}$  which gives

$$f^{**}(\mathbf{x}) = \max_{\tilde{\mathbf{z}}} 2\langle \tilde{\mathbf{x}}, \tilde{\mathbf{z}} \rangle - \sum_{i=1}^n \max(\tilde{z}_i^2 - g_i, 0). \quad (42)$$

The maximization over  $\tilde{\mathbf{z}}$  does in general not have any closed form solution but has to be evaluated numerically. Note however that it is a concave maximization problem that we can solve efficiently. One exception where we can find  $\tilde{\mathbf{z}}$  is for points  $\mathbf{x}$  where  $f^{**}(\mathbf{x}) = f(\mathbf{x})$ . In what follows we will derive some properties of the maximizing  $\mathbf{z}$  that simplifies the optimization.

We first consider the elements of  $\tilde{\mathbf{z}}$  independently without regard for their ordering. Each one has an objective function of the form

$$c_i(\tilde{z}_i) := 2\tilde{x}_i\tilde{z}_i - \max(\tilde{z}_i^2 - g_i, 0). \quad (43)$$

Figure 5 shows  $c_i(\tilde{z}_i)$  for different values of  $\tilde{x}_i$ . When  $\tilde{x}_i \geq \sqrt{g_i} \geq 0$  there is a unique maximizing point in  $\tilde{z}_i = \tilde{x}_i$ . If  $0 < \tilde{x}_i \leq \sqrt{g_i}$  the maximizing point is  $\tilde{z}_i = \sqrt{g_i}$ . In the last case where  $0 = \tilde{x}_i \leq \sqrt{g_i}$  any

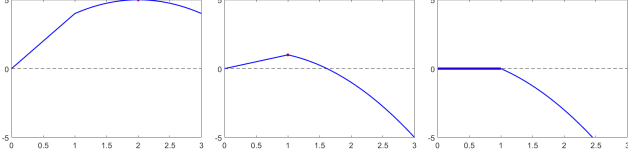


Figure 5: The objective function (43) for  $g_i = 1$  and  $x_i = 0, \frac{1}{2}$  and 2. The maximizing points are shown in red in each case.

$\tilde{z}_i \in [0, \sqrt{g_i}]$  is a maximizer. Suppose that we select  $k$  such that

$$\tilde{x}_i \geq \sqrt{g_i} \quad i \leq k \quad (44)$$

$$\tilde{x}_i < \sqrt{g_i} \quad i > k. \quad (45)$$

Then the unconstrained minimizers  $u_i^*$  of (43) can be written

$$u_i^* \in \begin{cases} \tilde{x}_i & i \leq k \\ \sqrt{g_i} & i > k, \tilde{x}_i \neq 0 \\ [0, \sqrt{g_i}] & i > k, \tilde{x}_i = 0 \end{cases} \quad (47)$$

Before we proceed any further we note that if  $\tilde{x}_i \notin (0, \sqrt{g_i})$  the second case will not occur. We can then select the elements of  $\tilde{z}$  so that each  $\tilde{z}_i$  maximizes  $c_i(\tilde{z}_i)$  without violating the ordering constraint.

**Lemma A.2.** *If  $\tilde{x}_i \notin (0, \sqrt{g_i})$  then the vectors maximizing (42) are given by*

$$\tilde{z}_i^* = \begin{cases} \tilde{x}_i & i \leq k \\ s_i & i > k, \end{cases} \quad (48)$$

where  $s_i, i = k+1, \dots, n$  is non-increasing and  $s_i \in [0, \min(\tilde{x}_k, \sqrt{g_{k+1}})]$ . In this case we also have

$$\begin{aligned} f^{**}(\mathbf{x}) &= \sum_{i=1}^k c_i(\tilde{x}_i) = \sum_{i=1}^k 2\tilde{x}_i \tilde{z}_i - \max(\tilde{x}_i^2 - g_i, 0) \\ &= \sum_{i=1}^k g_i + \tilde{x}_i^2 = f(\mathbf{x}). \end{aligned} \quad (49)$$

Before we proceed to the general case we note that if  $\tilde{x}_i \in (0, \sqrt{g_i})$  for some  $i$  then

$$c_i(u_i^*) = \begin{cases} \tilde{x}_i^2 + g_i & i \leq k \\ 2\tilde{x}_i \sqrt{g_i} & i > k, \tilde{x}_i \neq 0 \\ 0 & i > k, \tilde{x}_i = 0 \end{cases} \quad (50)$$

Since  $2\tilde{x}_i \sqrt{g_i} < \tilde{x}_i^2 + g_i$  if  $x_i < \sqrt{g_i}$  it is clear that this implies that

$$f^{**}(\mathbf{x}) < \sum_{i=1}^{\text{card}(\mathbf{x})} g_i + \tilde{x}_i^2 = f(\mathbf{x}). \quad (51)$$

For the general case the unconstrained minimizers are not non-increasing. To handle this we consider the best value  $\tilde{z}_i^*$  given values for its neighbors  $\tilde{z}_{i-1}^*$  and  $\tilde{z}_{i+1}^*$ . It is clear from the figures above that if the unconstrained minimizer  $u_i^*$  is unique then the best choice is of  $\tilde{z}_i^*$  is

$$\tilde{z}_i^* = \begin{cases} \tilde{z}_{i+1}^* & u_i^* \leq \tilde{z}_{i+1}^* \\ u_i^* & u_i^* \in [\tilde{z}_{i+1}^*, \tilde{z}_{i-1}^*] \\ \tilde{z}_{i-1}^* & u_i^* \geq \tilde{z}_{i-1}^* \end{cases} \quad (52)$$

Here we have adopted the convention that  $\tilde{z}_0 = \infty$  and  $\tilde{z}_{n+1} = 0$ . In the non-unique case we similarly have that  $\tilde{z}_i^* \in [0, \sqrt{g_i}] \cap [\tilde{z}_{i+1}^*, \tilde{z}_{i-1}^*]$  if this intersection is non-empty or  $\tilde{z}_i^* = \tilde{z}_{i+1}^*$ .

**Lemma A.3.** *Suppose that  $\{u_i^*\}$  is not monotone for all  $i$  such that  $\tilde{x}_i = 0$ . Let  $p$  be defined so that the sequence  $\{u_i^*\}$  is non-increasing for  $i \leq p$  and non-decreasing for  $i > p$ , whenever  $\tilde{x}_i \neq 0$ . The constrained maximizers  $\tilde{z}_i^*$  will then fulfill*

$$\tilde{z}_i^* \in \begin{cases} \max(u_i^*, \tilde{z}_{i+1}^*) & i \leq p \\ \tilde{z}_{i+1}^* & i > p, \tilde{x}_i \neq 0 \\ [0, \min(u_i^*, \tilde{z}_{i-1}^*)] & i > p, \tilde{x}_i = 0 \end{cases} \quad (53)$$

*Proof.* We first consider  $i > k$  with  $\tilde{x}_i = 0$ . Since  $\tilde{x}_i$  is non-increasing it is clear we can make  $c_i(\tilde{z}_i) = \dots = c_n(\tilde{z}_n) = 0$  by letting  $\tilde{z}_i = \dots = \tilde{z}_n = 0$ , regardless of what  $\tilde{z}_{i-1}$  is. Any optimal solution therefore has to have  $c_i(\tilde{z}_i) = \dots = c_n(\tilde{z}_n) = 0$ , which is achieved when  $\tilde{z}_i^* \in [0, \min(u_i^*, \tilde{z}_{i-1}^*)]$ .

Since  $\tilde{z}_0^* = \infty$  we have that  $u_1^* \notin [\tilde{z}_0^*, \tilde{z}_2^*]$  if and only if  $\tilde{z}_2^* > u_1^*$ . Therefore it is clear that  $\tilde{z}_1^* = \max(u_1^*, \tilde{z}_2^*)$ . Now suppose  $\tilde{z}_{i-1}^* = \max(u_{i-1}^*, \tilde{z}_i^*)$  for  $i \leq p$  then  $\tilde{z}_{i-1}^* \geq u_{i-1}^* \geq u_i^*$ , which means that either  $u_i^* \in [\tilde{z}_{i+1}^*, \tilde{z}_{i-1}^*]$ , in which case  $\tilde{z}_i^* = u_i^*$ , or  $u_i^* \leq \tilde{z}_{i+1}^*$  which gives  $\tilde{z}_i^* = \tilde{z}_{i+1}^*$ . This proves the first case in (53).

Now suppose that  $i \geq p+1$  has  $\tilde{x}_i \neq 0$ . If  $\tilde{z}_i^* \leq u_i^*$  then  $\tilde{z}_i^* \leq u_{i+1}^*$  since  $\{u_i^*\}$  is not decreasing and therefore  $\tilde{z}_{i+1}^* = \tilde{z}_i^*$  according to (53). If  $\tilde{z}_i^* > u_i^*$  we again have  $\tilde{z}_{i+1}^* = \tilde{z}_i^*$  according to (53).  $\square$

**Corollary A.4.** *The constrained minimizers  $\tilde{z}_i^*$  can be written*

$$\tilde{z}_i^* = \begin{cases} \max(u_i^*, s) & i \leq p \\ s & i > p, \tilde{x}_i \neq 0 \\ s_i & i > p, \tilde{x}_i = 0 \end{cases} \quad (54)$$

Here  $u_p^* \leq s$ ,  $\{s_i\}$ ,  $i = k+1, \dots, n$  is non-increasing and  $s_i \in [0, \min(\sqrt{g_i}, s)]$ .

*Proof.* The first two cases in (54) are fairly obvious. First it is clear that  $s := \tilde{z}_p^* = \tilde{z}_i^*$ , for all  $i > p$  with  $\tilde{x}_i \neq 0$ , which is the middle case in (54). Next we see that  $\tilde{z}_{p-1}^* = \max(u_{p-1}^*, \tilde{z}_k^*) = \max(u_{p-1}^*, \max(u_p^*, s)) = \max(u_{k-1}^*, s)$  since  $u_{p-1}^* \geq u_p^*$ . Repeating the same argument again shows the first case in (54).

Finally we note that  $s_i$  non increasing and  $s_i \in [0, \min(\sqrt{g_i}, s)]$  implies that  $s_i \leq s$  and therefore also that  $s_i \in [0, \min(\sqrt{g_i}, s)] \cap [0, s_{i-1}] = [0, \min(\sqrt{g_i}, s_{i-1})]$ , which shows the third case of (54).

To see that  $u_p^* \leq s$  we note that all residuals  $c_i$  are non-decreasing with  $s$  when  $s < u^*$ .  $\square$

## B. Proof of Theorem 3.1

In this section we give the proof of Theorem 3.1 which shows that "fixed cardinality/rank" solutions are stationary in our relaxation (7). The proofs for vector and matrix cases are somewhat different and therefore we treat them separately.

### B.1. The vector case

*Proof of Theorem 3.1.* The objective function of (7) can be written  $f^{**}(\mathbf{x}) + h(\mathbf{x})$  where  $h(\mathbf{x}) = -\|\mathbf{x}\|^2 + \|\mathbf{A}\mathbf{x} - \mathbf{b}\|^2$ . A stationary point therefore fulfills  $-\nabla h(\mathbf{x}) \in \partial f^{**}(\mathbf{x})$ . We have  $\nabla h(\mathbf{x}) = -2\mathbf{x} + 2A^T(\mathbf{A}\mathbf{x} - \mathbf{b})$  which yields

$$-A^T(\mathbf{A}\mathbf{x} - \mathbf{b}) = \mathbf{z} - \mathbf{x}, \quad (55)$$

where  $2\mathbf{z} \in \partial f^{**}(\mathbf{x})$ . Now suppose that  $\mathbf{o}$  fulfills the requirements of the theorem and let  $S$  be the set of nonzero elements of  $\mathbf{o}$ . The sub-differential  $\partial f^{**}(\mathbf{o})$  consists of the maximizing  $\mathbf{z}$ -vectors given in Lemma A.2. The the vector  $\mathbf{z} - \mathbf{o}$  is zero for every element in  $S$ . To see that the same is true for  $A^T(\mathbf{A}\mathbf{o} - \mathbf{b})$  we note that

$$\mathbf{o} = \arg \min_{\text{supp}(\mathbf{x})=S} \|\mathbf{A}\mathbf{x} - \mathbf{b}\|^2 = \arg \min_{\text{supp}(\mathbf{x})=S} \|A_S \mathbf{x} - \mathbf{b}\|^2, \quad (56)$$

where  $A_S$  is constructed by taking  $A$  and setting the columns no in  $S$  to zero. Therefore the normal equations  $A_S^T(A_S \mathbf{o} - \mathbf{b}) = A_S^T(\mathbf{A}\mathbf{o} - \mathbf{b}) = 0$  hold which shows that the elements of  $A^T(\mathbf{A}\mathbf{o} - \mathbf{b})$  that are in  $S$  all vanish.

It now remains to show that the elements in the complement of  $S$  are smaller than  $\min\{\tilde{x}_k, \sqrt{g_{k+1}}\}$ . This is however clear since by assumption

$$\|A^T(A^T \mathbf{o} - \mathbf{b})\| \leq \|A\| \|\mathbf{o}\| \leq \min\{\tilde{x}_k, \sqrt{g_{k+1}}\}. \quad (57)$$

We remark that estimating the size of the elements by the vector norm is a simple but very crude estimation and the result is therefore likely to hold under much more generous conditions.  $\square$

### B.2. The matrix case

*Proof of Theorem 3.1.* Similar to the vector case we need to show that

$$-\mathcal{A}^*(\mathbf{A}\mathbf{O} - \mathbf{b}) = \mathbf{Z} - \mathbf{O}, \quad (58)$$

where  $2\mathbf{Z} \in \partial f^{**}(\mathbf{O})$ . The matrix  $\mathbf{Z}$  is in the sub differential of  $f^{**}(\mathbf{O})$  if we can find orthogonal matrices  $U$  and  $V$  such that  $\mathbf{O} = UD_{\tilde{\mathbf{o}}}V^T$  and  $\mathbf{Z} = UD_{\tilde{\mathbf{z}}}V^T$ . Here  $\tilde{\mathbf{o}}$  and  $\tilde{\mathbf{z}}$  are the singular values of  $\mathbf{O}$  and  $\mathbf{Z}$  respectively. The matrices  $D_{\tilde{\mathbf{o}}}$  and  $D_{\tilde{\mathbf{z}}}$  are diagonal matrices with elements  $\tilde{\mathbf{o}}$  and  $\tilde{\mathbf{z}}$ . Note that  $\mathbf{O}$  is typically of low rank  $D_{\tilde{\mathbf{o}}}$  and  $D_{\tilde{\mathbf{z}}}$  can be partitioned into block matrices

$$D_{\tilde{\mathbf{o}}} = \begin{bmatrix} \Sigma & 0 \\ 0 & 0 \end{bmatrix} \text{ and } D_{\tilde{\mathbf{z}}} = \begin{bmatrix} \Sigma & 0 \\ 0 & \Delta \end{bmatrix}. \quad (59)$$

Here  $\Sigma$  contains the  $k$  non-zero singular values of  $\mathbf{O}$ . Due to Lemma A.2 the  $D_{\tilde{\mathbf{z}}}$  also contains this block. The matrix  $\Delta$  contains the singular values of  $\tilde{\mathbf{z}}$  that correspond to zeros in  $\tilde{\mathbf{o}}$ . We can make a corresponding partition of the  $U$  and  $V$  matrices into

$$U = [\bar{U} \quad \bar{U}_\perp] \text{ and } V = [\bar{V} \quad \bar{V}_\perp], \quad (60)$$

where  $\bar{U}$  and  $\bar{V}$  are the first  $k$  columns of  $U$  and  $V$  respectively. Note that only  $\bar{U}$  and  $\bar{V}$  are uniquely determined by  $\mathbf{O}$ . The matrices  $\bar{U}_\perp$  and  $\bar{V}_\perp$  can be selected arbitrarily as long as they are orthogonal to  $\bar{U}$  and  $\bar{V}$  respectively. Any choice of  $\bar{U}$ ,  $\bar{V}$  and  $\Delta$ , where the elements of  $\Delta$  are less than  $\min\{\tilde{o}_k, \sqrt{g_{k+1}}\}$  gives us a  $\mathbf{Z}$  that is in the sub differential. Consequently we have

$$\mathbf{Z} - \mathbf{O} = [\bar{U} \quad \bar{U}_\perp] \begin{bmatrix} 0 & 0 \\ 0 & \Delta \end{bmatrix} \begin{bmatrix} \bar{V}^T \\ \bar{V}_\perp^T \end{bmatrix} = \bar{U}_\perp \Delta \bar{V}_\perp^T. \quad (61)$$

$\square$

We now consider term  $\mathcal{A}^*(\mathbf{A}\mathbf{O} - \mathbf{b})$ . We have

$$\begin{aligned} \|\mathcal{A}(\mathbf{O} + t\mathbf{H}) - \mathbf{b}\|^2 &= t^2 \|\mathbf{A}\mathbf{H}\|^2 + 2t \langle \mathbf{H}, \mathcal{A}^*(\mathbf{A}\mathbf{O} - \mathbf{b}) \rangle \\ &\quad + \|\mathbf{A}\mathbf{O} - \mathbf{b}\|^2. \end{aligned} \quad (62)$$

Recall that  $\mathbf{O}$  minimizes the left hand side over all matrices with rank at most  $k$ . Since the linear term dominates the quadratic one for small  $t$  we must have

$$\langle \mathbf{H}, \mathcal{A}^*(\mathbf{A}\mathbf{O} - \mathbf{b}) \rangle \geq 0 \quad (63)$$

for all  $\mathbf{H}$  such that  $\text{rank}(\mathbf{O} + t\mathbf{H}) \leq k$ . Since  $\Sigma$  has full rank it is clear that any matrix of the form

$$\mathbf{H} = [\bar{U} \quad \bar{U}_\perp] \begin{bmatrix} H_{11} & H_{12} \\ H_{21} & 0 \end{bmatrix} \begin{bmatrix} \bar{V}^T \\ \bar{V}_\perp^T \end{bmatrix} \quad (64)$$

fulfills this requirement. It is now easy to see that

$$-\mathcal{A}^*(\mathcal{A}O - \mathbf{b}) = U_{\perp} M V_{\perp}^T, \quad (65)$$

where  $M$  is some matrix. Furthermore since  $U_{\perp}$  and  $V_{\perp}$  can be selected freely (as long as they are perpendicular to  $U$  and  $V$  respectively) we can assume that  $M$  is diagonal. What remains is therefore to estimate its singular values, which similarly to the vector case is done by

$$\|\mathcal{A}^*(\mathcal{A}O - \mathbf{b})\|_2 \leq \|\mathcal{A}\| \|\epsilon\| \leq \min\{\bar{o}_k, \sqrt{g_{k+1}}\}. \quad (66)$$

### C. Proof of Theorem 3.2

In this section we prove our main theorem. The proof requires a growth estimate of the subgradients of  $f^{**}$  which we give in the following lemmas.

**Lemma C.1.** *If  $\mathbf{z} \in \partial f^{**}(\mathbf{x})$  and  $\mathbf{z}' \in \partial f^{**}(\mathbf{x}')$  and  $d \leq 1$  then*

$$\langle \mathbf{z}' - \mathbf{z}, \mathbf{x}' - \mathbf{x} \rangle > d \|\mathbf{x}' - \mathbf{x}\|^2, \quad (67)$$

if

$$\langle \pi \tilde{\mathbf{z}}' - \tilde{\mathbf{z}}, \pi \tilde{\mathbf{x}}' - \tilde{\mathbf{x}} \rangle > d \|\pi \tilde{\mathbf{x}}' - \tilde{\mathbf{x}}\|^2, \quad (68)$$

for all permutation matrices  $\pi$ .

*Proof.* We have that (67) can be written

$$C - \langle \mathbf{z}' - d\mathbf{x}', \mathbf{x} \rangle - \langle \mathbf{z} - d\mathbf{x}, \mathbf{x}' \rangle > 0, \quad (69)$$

where

$$C = \langle \mathbf{z}' - d\mathbf{x}', \mathbf{x}' \rangle + \langle \mathbf{z} - d\mathbf{x}, \mathbf{x} \rangle. \quad (70)$$

Note that since the elements of  $\mathbf{z}'$  and  $\mathbf{x}'$  have the same signs and  $\tilde{z}'_i \geq \tilde{x}'_i$  for all  $i$  the term  $C$  is independent of signs. For fixed magnitudes and permutations the term  $\langle \mathbf{z}' - d\mathbf{x}', \mathbf{x} \rangle + \langle \mathbf{z} - d\mathbf{x}, \mathbf{x}' \rangle$  is clearly maximized when  $\mathbf{z}'$  and  $\mathbf{z}$  have the same signs. In which case we have

$$\langle \mathbf{z}' - \mathbf{z}, \mathbf{x}' - \mathbf{x} \rangle = \langle \pi \tilde{\mathbf{z}}' - \tilde{\mathbf{z}}, \pi \tilde{\mathbf{x}}' - \tilde{\mathbf{x}} \rangle \quad (71)$$

and  $\|\mathbf{x}' - \mathbf{x}\|^2 = \|\pi \tilde{\mathbf{x}}' - \tilde{\mathbf{x}}\|^2$  for some permutation  $\pi$ .  $\square$

In the matrix case we have  $Z \in \partial f^{**}(X)$  and  $Z' \in \partial f^{**}(X')$ . Recall that here  $\tilde{\mathbf{x}}, \tilde{\mathbf{z}}, \tilde{\mathbf{x}}', \tilde{\mathbf{z}}'$ , are the singular values of the matrices  $X, Z, X', Z'$  respectively. The corresponding statement is then that

$$\langle Z' - Z, X' - X \rangle > d \|X' - X\|^2 \quad (72)$$

holds whenever (69) holds. The proof is however more complicated than the vector case. We therefore refer the reader to Proposition 4.5 [11] from which it is clear that the above statement holds.

We are now ready to establish the growth estimates on the directional derivatives needed to prove Theorem 3.2. We will first consider directional derivatives between points where the relaxation is tight, that is  $f_g(\mathbf{x}) = f_g^{**}(\mathbf{x})$ . In the subsequent result we then relax this assumption to only be valid for one of the points (namely the stationary point we want to prove is unique).

**Lemma C.2.** *Suppose that  $2\mathbf{z} \in \partial f^{**}(\mathbf{x})$  and  $2\mathbf{z}' \in \partial f^{**}(\mathbf{x}')$ , and that neither  $\tilde{\mathbf{x}}$  nor  $\tilde{\mathbf{x}}'$  have values in  $(0, \sqrt{g_i})$ . If the elements of  $\tilde{\mathbf{z}}$  fulfill*

$$\tilde{z}_i \notin \left[ (1-d)\sqrt{g_k}, \frac{\sqrt{g_k}}{(1-d)} \right] \text{ and } \tilde{z}_{k+1} < (1-2d)\tilde{z}_k, \quad (73)$$

where  $k$  is defined so that  $\tilde{x}_i \geq \sqrt{g_i}$  for  $i \leq k$  and  $\tilde{x}_i = 0$  if  $i > k$ , then

$$\langle \mathbf{z}' - \mathbf{z}, \mathbf{x}' - \mathbf{x} \rangle > d \|\mathbf{x}' - \mathbf{x}\|^2. \quad (74)$$

*Proof.* We need show that

$$\langle \pi \tilde{\mathbf{z}}' - \tilde{\mathbf{z}}, \pi \tilde{\mathbf{x}}' - \tilde{\mathbf{x}} \rangle > d \|\pi \tilde{\mathbf{x}}' - \tilde{\mathbf{x}}\|^2, \quad (75)$$

where  $\pi$  is a permutation. For ease of notation let  $\mathbf{z}' = \pi \tilde{\mathbf{z}}'$  and  $\mathbf{x}' = \pi \tilde{\mathbf{x}}'$ . We let the  $\tilde{I} = \{i; \tilde{x}_i \neq 0\} = \{i; i \leq k\}$  and  $I' = \{i; x'_i \neq 0\}$ . Then

$$\begin{aligned} \langle \mathbf{z}' - \tilde{\mathbf{z}}, \mathbf{x}' - \tilde{\mathbf{x}} \rangle &= \sum_{i \in \tilde{I}, i \in I'} (x'_i - \tilde{x}_i)^2 + \sum_{i \in \tilde{I}, i \notin I'} \tilde{x}_i (\tilde{x}_i - z'_i) \\ &\quad + \sum_{i \notin \tilde{I}, i \in I'} x'_i (x'_i - \tilde{z}_i). \end{aligned} \quad (76)$$

Note that

$$d \|\mathbf{x}' - \tilde{\mathbf{x}}\|^2 = \sum_{i \in \tilde{I}, i \in I'} d(x'_i - \tilde{x}_i)^2 + \sum_{i \in \tilde{I}, i \notin I'} d\tilde{x}_i^2 + \sum_{i \notin \tilde{I}, i \in I'} dx_i'^2. \quad (77)$$

We first consider pairs of terms from the second and third sums of (76). If  $i \in \tilde{I}, i \notin I'$  and  $j \notin \tilde{I}, j \in I'$  we have

$$\tilde{x}_i (\tilde{x}_i - z'_i) + x'_j (x'_j - \tilde{z}_j) = \tilde{x}_i^2 + x_j'^2 - \tilde{x}_i z'_i - x_j' \tilde{z}_j. \quad (78)$$

Since  $j \notin \tilde{I}$  and  $i \in \tilde{I}$  we have  $\tilde{z}_j < (1-2d)\tilde{z}_k \leq (1-2d)\tilde{z}_i = (1-2d)\tilde{x}_i$ . Similarly, since  $j \in I'$  and  $i \notin \tilde{I}$  we have  $z'_i \leq z'_j = x'_j$ . Therefore

$$\tilde{x}_i z'_i \leq \tilde{x}_i x'_j \leq \frac{\tilde{x}_i^2 + x_j'^2}{2} \quad (79)$$

$$\tilde{x}_j z_j < (1-2d)\tilde{x}_j x_i \leq (1-2d) \frac{x_i^2 + \tilde{x}_i^2}{2} \quad (80)$$

which gives

$$\tilde{x}_i(\tilde{x}_i - z'_i) + x'_j(x'_j - \tilde{z}_j) < d(\tilde{x}_i^2 + x'_j{}^2). \quad (81)$$

If the number of elements in  $\tilde{I}$  and  $I'$  are the same the two last sums of (76) have the same number of terms. Then (81) shows that (76) larger is than (77) since clearly  $(x'_i - \tilde{x}_i)^2 > d(x'_i - \tilde{x}_i)^2$ . It therefore remains to consider the two cases when  $I$  has more elements than  $I'$  and vice versa.

Let  $k'$  be the number of elements in  $I'$ . Suppose first that  $\tilde{I}$  has more elements than  $I'$ , that is,  $k > k'$ . Then the middle sums of (76) and (77) have more terms than the third ones. Therefore we need to show that

$$\tilde{x}_i(\tilde{x}_i - z'_i) > d\tilde{x}_i^2, \quad (82)$$

for  $i \in \tilde{I}$  and  $i \notin I'$ . Suppose that  $\pi_{ij} = 1$ , that is element  $j$  of  $\tilde{\mathbf{z}}'$  is moved to element  $i$  of  $\mathbf{z}'$  by the permutation  $\pi$ . By Corollary A.4 we have that  $z'_i = \tilde{z}'_j \leq \tilde{x}'_{k'} \leq \sqrt{g_{k'}} \leq \sqrt{g_k}$ . Since  $i < k$  we have by assumption (73) that  $\tilde{x}_i = \tilde{z}_i > \frac{\sqrt{g_k}}{(1-d)}$ . Therefore

$$\tilde{x}_i - z'_i = (1-d)\tilde{x}_i + d\tilde{x}_i - z'_i > \sqrt{g_k} + d\tilde{x}_i - z'_i \geq d\tilde{x}_i, \quad (83)$$

which gives (82).

Now suppose instead that  $\tilde{I}$  has fewer elements than  $I'$ , that is,  $k' > k$ . Then we need to show that

$$x'_i(x'_i - \tilde{z}_i) \geq dx_i^2, \quad (84)$$

for  $i \notin \tilde{I}$  and  $i \in I'$ . Suppose again that  $\pi_{ij} = 1$ , that is element  $j$  of  $\tilde{\mathbf{z}}'$  is moved to element  $i$  of  $\mathbf{z}'$  by the permutation  $\pi$ . By Corollary A.4 we have  $x'_i = \tilde{z}'_j > \tilde{z}'_{k'} = \tilde{x}'_{k'} \geq \sqrt{g_{k'}}$ . Furthermore by assumption (73) we have  $\tilde{z}_i < (1-d)\sqrt{g_k} \leq (1-d)\sqrt{g_{k'}}$  since  $k < k'$ . Therefore

$$x'_i - \tilde{z}_i = (1-d)x'_i + dx'_i - \tilde{z}_i \geq (1-d)\sqrt{g_{k'}} + dx'_i - \tilde{z}_i > dx'_i. \quad (85)$$

which gives (84).  $\square$

**Lemma C.3.** *Suppose that  $\mathbf{x}$  fulfills the assumptions of Lemma C.2. If  $2\mathbf{z}' \in \partial f^{**}(\mathbf{x}')$  (without any additional assumptions on the values of  $\mathbf{x}'$  or  $\mathbf{z}'$ ) then (74) holds.*

*Proof.* By Lemma A.1 we have  $f^{**}(\mathbf{x}') = \sum_j \lambda_j f^{**}(\mathbf{x}^j)$ , where  $\mathbf{x}^j$  are points where  $f^{**}(\mathbf{x}^j) = f(\mathbf{x}^j)$ , that is  $\mathbf{x}^j$  has no elements in  $(0, \sqrt{g_i})$ . Then by Lemma C.2, for any  $\mathbf{z}' \in \partial f^{**}(\mathbf{x}') \subset \bigcap_j \partial f^{**}(\mathbf{x}^j)$  we have

$$\langle \mathbf{z}' - \mathbf{z}, \mathbf{x}^j - \mathbf{x} \rangle > d\|\mathbf{x}^j - \mathbf{x}\|^2. \quad (86)$$

By convexity of  $\|(\cdot) - X\|^2$  we now get

$$\begin{aligned} \langle \mathbf{z}' - \mathbf{z}, \mathbf{x}' - \mathbf{x} \rangle &> d \sum_j \lambda_j \|\mathbf{x}^j - \mathbf{x}\|^2 \\ &\geq d \|\sum_j \lambda_j \mathbf{x}^j - \mathbf{x}\|^2 = d\|\mathbf{x}' - \mathbf{x}\|^2. \end{aligned} \quad (87)$$

$\square$

*Proof of Theorem 3.2.* We will show that  $\nabla h(\mathbf{x}') = 2(I - A^T A)\mathbf{x}' + 2A^T b \notin \partial f^{**}(\mathbf{x}')$ . Suppose that  $2\mathbf{z}' \in \partial f^{**}(\mathbf{x}')$ . Since  $\mathbf{x}$  is stationary we have  $2\mathbf{z} + \nabla h(\mathbf{x}) = 0$

$$\langle \mathbf{z}' + \nabla h(\mathbf{x}'), \mathbf{x}' - \mathbf{x} \rangle = \langle \mathbf{z}' - \mathbf{z}, \mathbf{x}' - \mathbf{x} \rangle + \langle \nabla h(\mathbf{x}') - \nabla h(\mathbf{x}), \mathbf{x}' - \mathbf{x} \rangle. \quad (88)$$

For the second term we have

$$\begin{aligned} \langle \nabla h(\mathbf{x}') - \nabla h(\mathbf{x}), \mathbf{x}' - \mathbf{x} \rangle &= \|\mathbf{x}' - \mathbf{x}\|_F^2 - \|A(\mathbf{x}' - \mathbf{x})\|^2 \\ &\leq \delta_r \|\mathbf{x} - \mathbf{x}'\|^2, \end{aligned} \quad (89)$$

if  $\text{rank}(\mathbf{x} - \mathbf{x}') \leq r$ , which clearly holds if  $\text{card}(\mathbf{x}') \leq r - k$ . On the other hand we also have by Lemma C.2 that

$$\langle \mathbf{z}' - \mathbf{z}, \mathbf{x}' - \mathbf{x} \rangle > \delta_r \|\mathbf{x}' - \mathbf{x}\|^2, \quad (90)$$

and therefore (88) is positive and  $\mathbf{x}'$  cannot be a stationary point.

Suppose now that  $\mathbf{x}$  is a point that has  $\text{card}(\mathbf{x}) < \frac{r}{2}$ . We will consider the directional derivatives along the line  $\mathbf{x} + t\mathbf{v}$ , where  $\mathbf{v} = \frac{\mathbf{x}' - \mathbf{x}}{\|\mathbf{x}' - \mathbf{x}\|_F}$ . Since  $f^{**}$  is convex (and finite) the directional derivative of the objective function exists and is given by

$$\sup_{\mathbf{z}' \in \partial f^{**}(\mathbf{x} + t\mathbf{v})} \langle \mathbf{z}' + \nabla h(\mathbf{x} + t\mathbf{v}), \mathbf{v} \rangle. \quad (91)$$

Since the  $\text{card}(\mathbf{v}) \leq r$  it is clear by the arguments above that this is positive.  $\square$

## D. Proof of Theorem 3.3

*Proof.* We will let  $\mathbf{x}$  be a global solution to  $\min_{\text{card}(\mathbf{x}) \leq k} \|A\mathbf{x} - \mathbf{b}\|$  and show that this point will be stationary under the conditions above. To do this we need to show that  $2\mathbf{z} \in \partial f(\mathbf{x})$  for  $\mathbf{z} = (I - A^T A)\mathbf{x} + A^T \mathbf{b}$ . We first note that since  $\|A\| < 1$  the vector  $\mathbf{x}$  will be the global minimizer of (7) for the fixed-cardinality relaxation, that is, the special case  $g_i = 0$  if  $i \leq k$  and  $g_i = \infty$  if  $i > k$ . This shows that  $\mathbf{x}$  is stationary in (7) for this particular choice of  $g$ . In particular  $\mathbf{x} = D_s \pi \tilde{\mathbf{x}}$  and  $\mathbf{z} = D_s \pi \tilde{\mathbf{z}}$  with the same  $s$  and  $\pi$ . (In the matrix case the corresponding statement is that the SVD's of  $X$  and  $Z$  have the same  $U$  and  $V$  matrices.) Furthermore, since  $\text{card}(\mathbf{x}) \leq k$  and  $g_i = 0$  when  $i \leq k$  it is clear from Lemma A.2 that the  $\tilde{z}_i = \tilde{x}_i$  for  $i \leq k$ .



To show that  $\mathbf{x}$  is stationary for a general choice of  $g$  fulfilling (18) it is enough to show that  $\sqrt{g_i} \leq \tilde{x}_i$  for  $i \leq k$  and  $\sqrt{g_i} \geq \tilde{x}_k$  for  $i > k$  by Lemma A.2. This is however implied by the stricter constraints (14) and we therefore proceed by proving these directly.

First we show that  $\tilde{x}_i$  is close to  $\tilde{y}_i$ . Since  $\|A\mathbf{x} - \mathbf{b}\| \leq \|A\mathbf{y} - \mathbf{b}\| = \|\boldsymbol{\epsilon}\|$  we have

$$\begin{aligned} \sqrt{1 - \delta_{2k}} \|\mathbf{x} - \mathbf{y}\| &\leq \|A(\mathbf{x} - \mathbf{y})\| \\ &\leq \|A\mathbf{x} - \mathbf{b}\| + \|A\mathbf{y} - \mathbf{b}\| \quad (92) \\ &\leq 2\|\boldsymbol{\epsilon}\|. \end{aligned}$$

Therefore

$$|x_i - y_i| \leq \frac{2}{\sqrt{1 - \delta_{2k}}} \|\boldsymbol{\epsilon}\|. \quad (93)$$

Furthermore

$$\begin{aligned} \|\mathbf{z} - \mathbf{x}\| &= \|A^T A(\mathbf{x} - \mathbf{y}) - A^T \boldsymbol{\epsilon}\| \\ &\leq \|A^T\| \|A\| \|\mathbf{x} - \mathbf{y}\| + \|A^T\| \|\boldsymbol{\epsilon}\| \\ &\leq \|\mathbf{x} - \mathbf{y}\| + \|\boldsymbol{\epsilon}\| \quad (94) \\ &\leq \frac{3}{\sqrt{1 - \delta_{2k}}} \|\boldsymbol{\epsilon}\|. \end{aligned}$$

And since  $\tilde{x}_{k+1} = 0$  this means that

$$\tilde{z}_{k+1} \leq \frac{3}{\sqrt{1 - \delta_{2k}}} \|\boldsymbol{\epsilon}\|. \quad (95)$$

Now inserting the above estimates in  $\mathbf{z}_{k+1} < (1 - 2\delta_k)\mathbf{z}_k$  shows (after some simplification) that this constraint holds if

$$\tilde{y}_k > \frac{5 - 4\delta_k}{\sqrt{1 - \delta_{2k}}(1 - 2\delta_{2k})} \|\boldsymbol{\epsilon}\|, \quad (96)$$

which is implied by (17) since  $\delta_k \geq \delta_{2k} > 0$ . Furthermore since  $\sqrt{g_i}$  is non-decreasing (18) and (93) implies that  $\tilde{z}_i = \tilde{x}_i > \frac{\sqrt{g_k}}{1 - \delta_k}$  for  $i \leq k$ , while (18) and (95) implies that  $\mathbf{z}_i < (1 - \delta_k)\sqrt{g_i}$  for  $i > k$ .  $\square$