

Probabilistic Iterative Hard Thresholding for Sparse Learning *

Matteo Bergamaschi[†], Andrea Cristofari[‡], Vyacheslav Kungurtsev[§], and Francesco Rinaldi[†]

Abstract. For statistical modeling wherein the data regime is unfavorable in terms of dimensionality relative to the sample size, finding hidden sparsity in the ground truth can be critical in formulating an accurate statistical model. The so-called “ ℓ_0 norm”, which counts the number of non-zero components in a vector, is a strong reliable mechanism of enforcing sparsity when incorporated into an optimization problem. However, in big data settings wherein noisy estimates of the gradient must be evaluated out of computational necessity, the literature is scant on methods that reliably converge. In this paper we present an approach towards solving expectation objective optimization problems with cardinality constraints. We prove convergence of the underlying stochastic process, and demonstrate the performance on two Machine Learning problems.

Key words. cardinality constraint, stochastic optimization

MSC codes. 68Q25, 68R10, 68U05

1 Introduction In this paper we consider the cardinality constrained expectation objective problem,

$$(1.1) \quad \begin{aligned} \min_{x \in \mathbb{R}^n} \quad & f(x) := \mathbb{E}[F(x, \xi)] \\ \text{s.t.} \quad & \|x\|_0 \leq K, \end{aligned}$$

where $f(\cdot)$ is $L(f)$ continuously differentiable. We say that $x \in C_K$ if $\|x\|_0 \leq K$ and thus a feasible x corresponds to $x \in C_K$.

This optimization problem is particularly important in applications of data science. In particular, the expectation objective serves to quantify the minimization of some empirical loss function that enforces the fit of a statistical model fit to empirical data. Cardinality constraints enforce sparsity in the model, enabling the discovery of the most salient features as far as prediction accuracy.

Cardinality constraints present a significant challenge to optimization solvers. The so-called (as it is not, formally) zero norm is a discontinuous function that results in a highly nonconvex and disconnected feasible set, as well as an unusual topology of stationary points and minimizers [21, 22]. Algorithmic development has been, as similar to many such problems, a parallel endeavor from the mathematical optimization and the machine learning communities. When dealing with a deterministic objective function, procedures attuned to the structure of the problem and seeking stationary points of various strength are presented, for instance, in [3]. Methods for deterministic optimization problems with sparse symmetric sets are proposed in, e.g., [4, 18], while methods for deterministic optimization problems with both cardinality and nonlinear constraints are described in, e.g., [8, 9, 10, 14, 24, 23, 25]. Simultaneously works appearing in machine learning conferences, e.g., [31, 30, 27, 19], exhibit weak theoretical convergence guarantees, but appear to scale more adequately as far as numerical experience. Thus, an algorithm that enjoys both reliable performance together with strong

*Submitted to the editors September 2, 2024.

Funding: This work was funded by the European Union’s Horizon Europe research and innovation programme under grant agreement No. 101084642.

[†]Department of Mathematics “Tullio Levi-Civita”, University of Padua (bergamas@math.unipd.it, rinaldi@math.unipd.it),

[‡]Department of Civil Engineering and Computer Science Engineering, University of Rome “Tor Vergata” (andrea.cristofari@uniroma2.it),

[§]Czech Technical University in Prague (kunguvya@fel.cvut.cz),

38 theoretical guarantees, as sought for the high dimensional high data volume model fitting
 39 problems in contemporary data science, is as of yet unavailable.

40 In this paper we attempt to reconcile these two and present an algorithm that is associated
 41 with reasonably strong theoretical convergence guarantees, while at the same time able to
 42 solve large scale problems of interest in statistics and machine learning. To this end we
 43 present a procedure under the framework of *Probabilistic Models*, which can be understood as
 44 a sequential linear Sample Average Approximation (SAA) scheme for solving problems with
 45 statistics in the objective function. First introduced in [2], then rediscovered with extensive
 46 analysis in [1, 15], this approach can exhibit asymptotic (and even worst case complexity)
 47 results to a local minimizer of the original problem, while still allowing the use of Newton-
 48 type second order iterations of subproblem solutions, and thus faster convergence as far as
 49 iteration count. The use of probabilistically accurate estimates within a certain bound in these
 50 methods permit a rather flexible approach to estimating the gradient, including techniques
 51 that introduce bias, while foregoing the necessity of a stepsize asymptotically diminishing to
 52 zero. However, asymptotic accurate convergence still requires increasing the batch size, so the
 53 tradeoffs in precision and certainty relative to computation become apparent, and adaptive
 54 for the user, in deciding at which point to stop the algorithm and return the current iterate
 55 as an estimate of the solution.

56 As contemporary Machine Learning applications, we shall consider Adversarial Attacks
 57 (see, e.g., [11, 16, 26] and references therein for further details), and Probabilistic Graphical
 58 Model training (see, e.g., [5, 28] and references therein for further details). In this paper we
 59 shall see how the use of a stochastic gradient and hard sparsity constraint can improve the
 60 performance and model quality in the considered problems.

61 The paper is organized as follows: In Section 2, we introduce some basic definitions and
 62 preliminary results related to optimality conditions of problem (1.1) that ease the theoretical
 63 analysis. We then describe the details of the proposed algorithmic scheme in Section 3. We
 64 then prove almost sure convergence to suitable stationary points in Section 4. Numerical
 65 results on some relevant Machine Learning applications are reported Section 5. Finally, we
 66 draw some conclusions and discuss some possible extensions in Section 6.

67 **2 Background** Cardinality constrained optimization presents an extensive hierarchy of
 68 stationarity conditions, as due to the geometric complexity of the feasible set. This neces-
 69 sitates specialized notions of projection, and presents complications due to the projection
 70 operation's generic non-uniqueness.

71 *Definitions and Preliminaries* The active and inactive set of a vector $x \in \mathbb{R}^n$ are respectively
 72 denoted by

$$73 \quad I_{\mathcal{A}}(x) := \{i \in \{1, \dots, n\}, x_i = 0\}, \quad I_{\mathcal{I}}(x) := \{i \in \{1, \dots, n\}, x_i \neq 0\}.$$

74 A set T is a *super-support* of $x \in C_K$ if $I_{\mathcal{A}}(x) \subseteq T$ and $|T| = s$. Let the permutation group
 75 of $\{1, \dots, n\}$ be denoted as Σ_n and for a permutation $\sigma \in \Sigma_n$, we write $(x^\sigma)_i = x_{\sigma(i)}$. For a
 76 vector $x \in \mathbb{R}^n$ we denote with $M_i(x)$ the i -th largest absolute-value component of x , thus we
 77 have $M_1(x) \leq M_2(x) \leq \dots \leq M_n(x)$.

78 We finally define the orthogonal projection as

$$79 \quad P_{C_K}(x) = \arg \min\{\|z - x\|^2, z \in C_K\},$$

80 that is, an n -length vector consisting of the s components of x with the largest absolute value.
 81 Such operator, as already highlighted in the previous section, is not single-valued due to the
 82 inherent non-convexity of the set C_K and plays a critical role in the development of algorithms
 83 for sparsity constrained optimization (see, e.g., [3, Section 2] for a discussion on this matter).

84 *Optimality Conditions* Now we define several optimality conditions for (1.1), borrowing
 85 heavily from [3]. Observe that a notable characteristic of cardinality constrained optimization
 86 is the presence of a hierarchy of optimality conditions, that is, a number of conditions that
 87 hold at optimal points that range across levels of restriction.

88 When restricted to a specific support, the "no descent directions" rule still provides a
 89 necessary optimality condition, which is referred to as basic feasibility. For a full support,
 90 this condition aligns with the standard stationarity condition, but only applies within the
 91 support set. If the support is not full, the stationarity condition must hold for any potential
 92 full support set that includes the given support, that is the gradient needs to be zero.

93 **Definition 2.1.** $x^* \in C_K$ is Basic Feasible (BF) for problem (1.1) when

94. $\nabla f(x^*) = 0$, if $\|x^*\|_0 < K$,
 92. $\nabla f_i(x^*) = 0$ for all $i \in I_{\mathcal{I}}(x^*)$, if $\|x^*\|_0 = K$.

96 We thus have that when a point $x^* \in C_K$ is optimal for problem (1.1), then x^* is a BF point
 97 (see Theorem 1 in [3]). The BF property is however a relatively weak necessary condition
 98 for optimality. Consequently, stronger necessary conditions are required to achieve higher
 99 quality solutions. This is why we use L-stationarity, an extension of the stationarity concept
 100 for convex constrained problems.

101 **Definition 2.2.** $x^* \in C_K$ is L - stationary for problem (1.1) when

$$102 \quad (2.1) \quad x^* \in P_{C_K} \left(x^* - \frac{1}{L} \nabla f(x^*) \right).$$

103 An equivalent analytic property of L-stationarity is given by the following lemma.

104 **Lemma 2.3.** [3, Lemma 2.2] L-stationarity at x^* is equivalent to $\|x^*\|_0 \leq K$ and

$$105 \quad |\nabla_i f(x^*)| \begin{cases} \leq LM_K(x^*) & i \in I_{\mathcal{A}}(x^*) \\ = 0 & i \in I_{\mathcal{I}}(x^*) \end{cases} .$$

106 The next result relates L-stationarity and Basic Feasibility:

107 **Corollary 2.4.** [3, Corollary 2.1] Suppose that $x^* \in C_k$ is an L-stationary of problem (1.1)
 108 for some L. Then x^* is BF for problem (1.1).

109 In addition, the likely intuition that the L-stationarity is related to the gradient Lipschitz
 110 constant is correct:

111 **Theorem 2.5.** [3, Theorem 2.2] If x^* is an optimal solution for problem (1.1) then it is
 112 L-stationary for all $L > L(f)$.

113 To see the distinction between BF and L-stationary, we can consider that if the Lipschitz
 114 constant of f is 1, then $x^* = (1, 0)$ with $\nabla f(x^*) = (-10, 1)$ satisfies BF but not L-stationarity.
 115 In particular it is clear from a linearization that $f((0, y)) < f(x^*)$ for y small.

116 In this sense L-stationarity is stronger than a linearized feasible direction stationarity
 117 measure, as constructed in [20]. This is because any feasible path for an active component,
 118 that is a direction from which a zero component becomes non-zero, would require a discrete
 119 jump from another component, that is the assignment of zero to a different component, in
 120 order to maintain the constraint. Thus there is no feasible linearized direction in which a zero
 121 component becomes non-negative on which to consider possible descent when the cardinality
 122 constraint is active. L-stationarity enables a relaxation of this by considering Lipschitz bounds
 123 on how much the function value can change along various directions depending on the gradient
 124 vector components.

125 *Iterative Hard Thresholding.* An important component of particularly machine learning
 126 literature procedures to solve (1.1) is the *Hard Thresholding Operator* (see, e.g., [3, 7] for
 127 further details). Consider the operator $\mathbf{HT}^x(v)$ applied to a vector v as one that projects v
 128 onto the sparsity constraint, i.e.,

$$129 \quad (2.2) \quad \mathbf{HT}^x(v) \in \arg \min_w \{ \|v - w\|, \|w\|_0 \leq K \} := P_{C_K}(v).$$

130 **3 Algorithm**

131 *Rolling Projection Estimator* Recall that, the sparse projection operation $P_{C_K}(v)$ for a
 132 vector v amounts to performing a sorting operation $\sigma \in \tilde{\Sigma}(v)$ on $\sigma(v)$, and then keeping the
 133 K largest magnitude components of v while setting the rest to zero.

134 Observe that an algorithmic iterative descent procedure would involve the negative of the
 135 gradient of f or an estimate thereof. Indeed, as the objective function is an expectation, we
 136 do not have access to the exact value of the $\nabla f(x)$ and hence the magnitude ranking of the
 137 its components. Thus, we must by necessity use noisy gradient estimates $\nabla F(x, \hat{\xi})$ to attempt
 138 to estimate the actual ranking of component magnitudes.

139 Asymptotically, we want to ensure that this sparse projector estimates the true ranking
 140 at any limit point. Given the natural source of asymptotically increasing sample sizes, this
 141 present a natural opportunity to use the Algorithm iterate sequence itself to perform this
 142 estimate, ultimately relying on consistency for statistical guarantees on accurate identification.

143 Let x_k correspond to the current iterate. Now we define our particular sequential estimate
 144 of the ranking of the magnitude of the vector components of the gradient of $f(x_k)$. Specifically,
 145 we are given a noisy evaluation $g_k \approx \nabla f(x_k)$, and an application

$$146 \quad (3.1) \quad \sigma_k(g_k) \in \tilde{\Sigma} \left(x_k - \alpha \min \left\{ 1, \frac{\delta_k}{\alpha \|g_k\|} \right\} g_k \right).$$

147 At the same time, there exists a set of permutations $S_k = \{ \sigma^{(j)} \}_{j \in [J]}$, $\sigma^{(j)} \in \Sigma_n$ with coefficient
 148 weights $\{ \omega^{(j)} \}_{j \in [J]}$, $\omega \in \Delta_J$.

149 We now perform exponential smoothing (exponential moving average) on the estimate,
 150 with smoothing parameter α_s :

$$151 \quad (3.2) \quad \begin{aligned} \sigma_k = \sigma^{(j)} \in S_k, &\implies \begin{cases} \omega \leftarrow (1 - \alpha_s)\omega, \\ \omega^{(j)} \leftarrow \omega^{(j)} + \alpha_s, \end{cases} \\ \sigma_k \notin S_k &\implies \begin{cases} S_k \leftarrow S_k \cup \{ \sigma_k \}, \\ \omega \leftarrow (1 - \alpha_s)\omega, \\ \omega^{(|S_k|)} \leftarrow \alpha_s. \end{cases} \end{aligned}$$

152 This accomplishes the following: We maintain a set of possible permutations with associ-
 153 ated mixture weights. With each new iteration, we sort the components of the noisy gradient
 154 estimate. If this sorting permutation has been found before, then we add to a weight corre-
 155 sponding to that permutation and lower the weights of others. Otherwise, i.e. this is a new
 156 permutation, we add it to the list of options.

157 Now, let $\hat{\sigma}^k$ be such that

$$158 \quad (3.3) \quad \hat{\sigma}^k = \sigma^{(j)} \in S_k \text{ with } \omega^{(j)} = \arg \max_{l \in [S_k]} \omega^{(l)}.$$

159 Thus, rather than taking the maximal components based on the current sorting, we use the
 160 moving average historical estimate. Then, taking

$$161 \quad (3.4) \quad I_k = \text{supp} \left(\max_K \hat{\sigma}^k \right),$$

162 that is, the set of indices whose components are largest, we present the **Pseudo Hard**
 163 **Thresholding** operator corresponding to iteration k , defined as follows:

$$164 \quad (3.5) \quad \mathbf{HT}^{x,\delta,k}(v) \in \arg \min_w \{ \|v - w\|, w_{[n] \setminus I_k} = 0, \|w - x\| \leq \delta \}.$$

165 We take a *clipped* step, wherein we step in the negative direction of the scaled negative gra-
 166 dient $-\alpha \min \left\{ 1, \frac{\delta_k}{\alpha \|g_k\|} \right\} g_k$, with α being a positive constant. The Pseudo-Hard Thresholding
 167 algorithm can be computed in a straightforward closed form expression:

$$168 \quad (3.6) \quad [\hat{x}_k]_i = \begin{cases} 0 & i \notin I_k \\ \left[x_k - \alpha \min \left\{ 1, \frac{\delta_k}{\alpha \|g_k\|} \right\} g_k \right]_i & i \in I_k. \end{cases}$$

169 From (3.6), observe that

$$170 \quad (3.7) \quad \hat{x}_k = P_{I_k} \left(x_k - \alpha \min \left\{ 1, \frac{\delta_k}{\alpha \|g_k\|} \right\} g_k \right).$$

171 This presents an opportunity to get a sort of descent lemma in the context of cardinality
 172 constrained optimization problems. To this end, define

$$173 \quad (3.8) \quad h_k(y) = f(x_k) + g_k^T(y - x_k),$$

174 so that

$$175 \quad (3.9) \quad h_k(\hat{x}_k) - h_k(x_k) = g_k^T(\hat{x}_k - x_k) \leq -\frac{1}{\alpha} \max \left\{ 1, \frac{\alpha \|g_k\|}{\delta_k} \right\} \|\hat{x}_k - x_k\|^2 \leq -\frac{1}{\alpha} \|\hat{x}_k - x_k\|^2,$$

176 where the first inequality follows from known results on the projection operator [6].

177 Accuracy Estimates

178 **Definition 3.1.** Define $s_k = \hat{x}_k - x_k$. The function estimates f_k^0 and f_k^s are ε_f -accurate
 179 estimates of $f(x_k)$ and $f(x_k + s_k)$, respectively, for a given δ_k if

$$180 \quad (3.10) \quad |f_k^0 - f(x_k)| \leq \varepsilon_f \delta_k^2 \quad \text{and} \quad |f_k^s - f(x_k + s_k)| \leq \varepsilon_f \delta_k^2.$$

181 **Definition 3.2.** The model for generating the iterate is κ - δ_k , or (κ_f, κ_g) - δ_k accurate, when

$$182 \quad (3.11) \quad \|\nabla F(y) - g_k\| \leq \kappa_g \delta_k \quad \text{and} \quad |f(y) - f(x_k) - g_k^T(y - x_k)| \leq \kappa_f \|y - x_k\| \delta_k^2$$

183 for all $y \in B(x_k, \delta_k)$.

184 Note that this implies:

$$185 \quad (3.12) \quad \|[\nabla F(y) - g_k]_{I_k}\| \leq \kappa_g \delta_k \quad \text{and} \quad |f(y) - f(x_k) - [g_k]_{I_k}^T [y - x_k]_{I_k}| \leq \kappa_f \|y - x_k\| \delta_k^2$$

186 for all $y \in B(x_k, \delta_k)$.

187 **4 Convergence Theory** Now we develop our argument for justifying the long term con-
 188 vergence of the Algorithm based on classic arguments on probabilistic models given in [15](see
 189 also [1, 13]). To this end, we remark that the iterates, being dependent on random function
 190 and gradient estimates, define a stochastic process X_k . The Algorithm itself is a realization,
 191 thus denoting $x_k = X_k(\omega)$, $\delta_k = \Delta_k(\omega)$, etc. for ω the random element defining the realiza-
 192 tion. Similar as to the original, we can consider a filtration with the sigma algebra \mathcal{F}_k defining
 193 the start of the iteration, and $\mathcal{F}_{k+\frac{1}{2}}$ defining the algebra after the minibatch has been sampled
 194 and g_k computed. This filtration will be implicit in the statements of the convergence results.

195 We begin with a standard assumption on a probability bound on the accuracy of the
 196 conditions given by Definition 3.1 and 3.2. To this end define θ, β to be the probability that
 197 a given sample of g_k .

Algorithm 3.1 Probabilistic Iterative Hard Thresholding

1: **Initialization:** $x_0 \in C_K$, $\delta_0 \in (0, \delta_{max}]$, Parameters $\delta_{max} > 0$, $\gamma \in (0, 1)$.
2: **for** $k = 0, 1, 2, \dots$ **do**
3: Sample a minibatch $\xi_k \sim \Xi$ and compute $g_k = \nabla F(x_k, \xi_k)$
4: Compute σ^k by (3.1) and update ω by (3.2)
5: Compute $\hat{\sigma}^k$ from (3.3) and use it to define I_k by (3.4).
6: Compute \hat{x}_k from the Pseudo-Hard-Thresholding (3.6)
7: Compute stochastic estimates $f_k^s \approx f(\hat{x}_k)$, $f_k^0 \approx f(x_k)$
8: **if** $\frac{f_k^0 - f_k^s}{\|[g_k]_{I_k}\| \delta_k} \geq \eta_1$ and $\|[g_k]_{I_k}\| \geq \eta_2 \delta_k$ **then**
9: Set $\delta_{k+1} = \min\{\gamma \delta_k, \delta_{max}\}$, let $x_{k+1} = \hat{x}_k$
10: **else**
11: Set $\delta_{k+1} = \gamma^{-1} \delta_k$, let $x_{k+1} = x_k$
12: **end if**
13: **end for**

198 **Assumption 4.1.** Given $\theta, \beta \in (0, 1)$ and ε_f , there exist κ_g, κ_f such that the sequence of $\{g_k\}$
199 is such that with probability θ , $\kappa\text{-}\delta_k$ -accuracy holds as per Definition 3.2, and with probability
200 β , ε_f accuracy holds as by Definition 3.1.

201 We can consider that [3, Lemma 3.1] provides for the enforcement of function decrease
202 in the favorable probabilistic cases in the convergence theory. Indeed, one can derive the
203 following lemma which also functionally corresponds to [15, Lemma 4.5].

204 **Lemma 4.2.** If the model for generating the iterate k is $\kappa\text{-}\delta_k$ accurate according to Defini-
205 tion 3.2, with \hat{x}_k and δ_k being such that

$$206 \quad (4.1) \quad \delta_k \leq \frac{1}{2\alpha\kappa_g\delta_{max}} \|x_k - \hat{x}_k\|,$$

207 then

$$208 \quad (4.2) \quad f(x_k) - f(\hat{x}_k) \geq \frac{1}{2\alpha} \|\hat{x}_k - x_k\|^2.$$

209 *Proof.* Using the definition of h_k given in (3.8), we can write

$$\begin{aligned} 210 \quad f(\hat{x}_k) - f(x_k) &= f(\hat{x}_k) - h_k(\hat{x}_k) + h_k(\hat{x}_k) - h_k(x_k) + h_k(x_k) - f(x_k) \\ &= f(\hat{x}_k) - h_k(\hat{x}_k) + h_k(\hat{x}_k) - h_k(x_k) \\ &= f(\hat{x}_k) - f(x_k) - g_k^T(\hat{x}_k - x_k) + g_k^T(\hat{x}_k - x_k) \\ &\leq \kappa_g \|x_k - \hat{x}_k\| \delta_k^2 + g_k^T(\hat{x}_k - x_k), \end{aligned}$$

where the inequality follows from the second condition in (3.11). Using (3.9), we also have that

$$g_k^T(\hat{x}_k - x_k) \leq -\frac{1}{\alpha} \|\hat{x}_k - x_k\|^2.$$

211 Then, we obtain

$$212 \quad (4.3) \quad f(\hat{x}_k) - f(x_k) \leq \kappa_g \|x_k - \hat{x}_k\| \delta_k^2 - \frac{1}{\alpha} \|\hat{x}_k - x_k\|^2 \leq \kappa_g \delta_{max} \|x_k - \hat{x}_k\| \delta_k - \frac{1}{\alpha} \|\hat{x}_k - x_k\|^2,$$

213 where the last inequality follows from the fact that $\delta_k \leq \delta_{max}$. Moreover, (4.1) implies that

$$214 \quad \kappa_g \delta_{max} \|x_k - \hat{x}_k\| \delta_k \leq \frac{1}{2\alpha} \|\hat{x}_k - x_k\|^2.$$

215 Using this inequality in (4.3), the desired result follows. ■

216 Now, taking inspiration from [15, Lemma 4.6], we can bound the decrease with respect to
217 the projected real gradient.

218 **Lemma 4.3.** *If the model for generating the iterate k is $\kappa\text{-}\delta_k$ accurate according to Defini-
219 tion 3.2 and*

$$220 \quad (4.4) \quad \delta_k \leq a \left\| x_k - P_{I_k} \left(x_k - \alpha \min \left\{ 1, \frac{\delta_k}{\alpha \|\nabla f(x_k)\|} \right\} \nabla f(x_k) \right) \right\|,$$

221 where

$$222 \quad (4.5) \quad a = \frac{1}{2\alpha\kappa_g\delta_{max} + 2\sqrt{K}}$$

223 and

$$224 \quad (4.6) \quad \alpha > \frac{\sqrt{K}}{\kappa_g\delta_{max}},$$

225 then

$$226 \quad (4.7) \quad f(x_k) - f(\hat{x}_k) \geq c \left\| x_k - P_{I_k} \left(x_k - \alpha \min \left\{ 1, \frac{\delta_k}{\alpha \|\nabla f(x_k)\|} \right\} \nabla f(x_k) \right) \right\|^2,$$

227 with

$$228 \quad c = \frac{1 - 4a\sqrt{K}}{2\alpha} > 0.$$

229 *Proof.* We can write

$$230 \quad \left\| x_k - P_{I_k} \left(x_k - \alpha \min \left\{ 1, \frac{\delta_k}{\alpha \|\nabla f(x_k)\|} \right\} \nabla f(x_k) \right) \right\| \leq \\ \left\| x_k - \hat{x}_k \right\| + \left\| \hat{x}_k - P_{I_k} \left(x_k - \alpha \min \left\{ 1, \frac{\delta_k}{\alpha \|\nabla f(x_k)\|} \right\} \nabla f(x_k) \right) \right\|.$$

231 Using (3.7), we get

$$232 \quad (4.8) \quad \left\| x_k - P_{I_k} \left(x_k - \alpha \min \left\{ 1, \frac{\delta_k}{\alpha \|\nabla f(x_k)\|} \right\} \nabla f(x_k) \right) \right\| = \\ \left\| x_k - \hat{x}_k \right\| + \left\| \alpha \min \left\{ 1, \frac{\delta_k}{\alpha \|g_k\|} \right\} [g_k]_{I_k} - \alpha \min \left\{ 1, \frac{\delta_k}{\alpha \|\nabla f(x_k)\|} \right\} [\nabla f(x_k)]_{I_k} \right\| = \\ \left\| x_k - \hat{x}_k \right\| + \delta_k \left\| \min \left\{ \frac{\alpha \|g_k\|}{\delta_k}, 1 \right\} \frac{[g_k]_{I_k}}{\|g_k\|} - \min \left\{ \frac{\alpha \|\nabla f(x_k)\|}{\delta_k}, 1 \right\} \frac{[\nabla f(x_k)]_{I_k}}{\|\nabla f(x_k)\|} \right\| \leq \\ \left\| x_k - \hat{x}_k \right\| + 2\sqrt{K}\delta_k,$$

233 where the last inequality follows from the fact that $\|u - v\| \leq \sqrt{K}\|u - v\|_\infty \leq 2\sqrt{K}$ for all
234 $u, v \in \mathbb{R}^K$ such that $\|u\| = \|v\| = 1$. From (4.4), the first term in (4.8) is greater of equal to
235 δ_k/a , leading to

$$236 \quad \frac{\delta_k}{a} \leq \|x_k - \hat{x}_k\| + 2\sqrt{K}\delta_k.$$

237 Using the definition of a given in (4.5), it follows that (4.1) is satisfied and we can apply
 238 Lemma 4.2, obtaining

$$239 \quad (4.9) \quad f(x_k) - f(\hat{x}_k) \geq \frac{1}{2\alpha} \|\hat{x}_k - x_k\|^2.$$

240 Finally, in order to lower bound the right-hand side term in the above inequality, using (4.8)
 241 we can write

$$\begin{aligned} \|x_k - \hat{x}_k\|^2 &\geq \left(\left\| x_k - P_{I_k} \left(x_k - \alpha \min \left\{ 1, \frac{\delta_k}{\alpha \|\nabla f(x_k)\|} \right\} \nabla f(x_k) \right) \right\| - 2\sqrt{K}\delta_k \right)^2 \\ &\geq \left\| x_k - P_{I_k} \left(x_k - \alpha \min \left\{ 1, \frac{\delta_k}{\alpha \|\nabla f(x_k)\|} \right\} \nabla f(x_k) \right) \right\|^2 + \\ 242 \quad &\quad - 4\sqrt{K}\delta_k \left\| x_k - P_{I_k} \left(x_k - \alpha \min \left\{ 1, \frac{\delta_k}{\alpha \|\nabla f(x_k)\|} \right\} \nabla f(x_k) \right) \right\| \\ &\geq (1 - 4a\sqrt{K}) \left\| x_k - P_{I_k} \left(x_k - \alpha \min \left\{ 1, \frac{\delta_k}{\alpha \|\nabla f(x_k)\|} \right\} \nabla f(x_k) \right) \right\|^2, \end{aligned}$$

243 where the last inequality follows from (4.4). From (4.6), it also follows that $c > 0$, thus leading
 244 to the desired result. \blacksquare

245 The next lemma states conditions on δ_k to guarantee that an iteration is successful,
 246 similarly as in [15, Lemma 4.7].

247 **Lemma 4.4.** *If, at iteration k , the estimates f_k^0, f_k^s are ε_f -accurate according to Defini-*
 248 *tion 3.1 and the model is κ - δ_k accurate according to Definition 3.2, with*

$$249 \quad \delta_k \leq \min \left\{ \frac{1}{\eta_2}, \frac{1 - \eta_1}{2\varepsilon_f + \kappa\delta_{max}} \right\} \|[g_k]_{I_k}\|,$$

250 *then the step is accepted.*

251 *Proof.* Define

$$252 \quad \rho_k = \frac{f_k^0 - f_k^s}{\|[g_k]_{I_k}\| \delta_k}.$$

253 Using (3.10) and (3.11), we can write

$$\begin{aligned} \rho_k &= \frac{f_k^0 - f(x_k)}{\|[g_k]_{I_k}\| \delta_k} + \frac{f(x_k) - f(\hat{x}_k)}{\|[g_k]_{I_k}\| \delta_k} + \frac{f(\hat{x}_k) - f_k^s}{\|[g_k]_{I_k}\| \delta_k} \\ 254 \quad &\leq \frac{2\varepsilon_f \delta_k}{\|[g_k]_{I_k}\|} + \frac{[g_k]_{I_k}^T [\hat{x}_k - x_k]_{I_k} + \kappa_g \|\hat{x}_k - x_k\| \delta_k^2}{\|[g_k]_{I_k}\| \delta_k} \\ &\leq \frac{2\varepsilon_f \delta_k}{\|[g_k]_{I_k}\|} + 1 + \frac{\kappa_g \delta_{max} \delta_k}{\|[g_k]_{I_k}\|}, \end{aligned}$$

255 where the last inequality follows from the fact that $\|\hat{x}_k - x_k\| \leq \delta_k$ and $\delta_k \leq \delta_{max}$. Then

$$256 \quad |\rho_k - 1| \leq \frac{(2\varepsilon_f + \kappa_g \delta_{max}) \delta_k}{\|[g_k]_{I_k}\|} \leq 1 - \eta_1,$$

257 where we have used the assumption on δ_k in the last inequality. Hence, $\rho_k \geq \eta_1$. Since we
 258 have also assumed that $\|[g_k]_{I_k}\| \geq \eta_2 \delta_k$, from the instructions of the algorithm (see line 8 of
 259 Algorithm 3.1) it follows that the step is accepted. \blacksquare

260 **Lemma 4.5.** *If the estimates f_k^0, f_k^s at iteration k are ε_f -accurate according to Defini-*
 261 *tion 3.1 with $\varepsilon_f < (\eta_1\eta_2)/2$ and the step is accepted, then*

$$262 \quad f(x_{k+1}) - f(x_k) \leq -C\|\delta_k\|^2,$$

263 *with $C = \eta_1\eta_2 - 2\varepsilon_f > 0$.*

264 *Proof.* Since the step is accepted, from the instructions of the algorithm (see line 8 of
 265 Algorithm 3.1) we can write

$$266 \quad (4.10) \quad f_k^0 - f_k^s \geq \eta_1 \| [g_k]_{I_k} \| \delta_k \geq \eta_1 \eta_2 \delta_k^2.$$

267 Moreover,

$$268 \quad f(x_k + s_k) - f(x_k) = f(x_k + s_k) - f_k^s + f_k^s - f_k^0 + f_k^0 - f(x_k) \leq 2\varepsilon_f \delta_k^2 - \eta_1 \eta_2 \delta_k^2,$$

269 where the inequality follows from (3.10) and (4.10). Then, using the definition of C given in
 270 the assertion, the desired result follows. ■

271 Now we define the stochastic process

$$272 \quad (4.11) \quad \Phi_k := \nu f(x_k) + (1 - \nu) \delta_k^2.$$

273 The next Theorem is along the lines of Theorem 4.11 in [15]. The result requires a
 274 compactness assumption, which we present first.

275 **Assumption 4.6.** *Let \mathcal{L} be the level set of the iterates generated by the algorithm, that is,*

$$276 \quad \mathcal{L} = \{x : f(x) \leq f(x_k)\}, \forall x_k$$

277 *noting that this depends on the stochastic realization of the iterates and gradient estimates.*
 278 *Assume that \mathcal{L} is bounded below and that f is L -Lipschitz and its gradient is L -Lipschitz*
 279 *continuous on \mathcal{L} .*

280 **Theorem 4.7.** *Let $\{x_k\}$ be the sequence of iterates generated by the Probabilistic Iterative*
 281 *Hard Thresholding Algorithm (Algorithm 3.1) under Assumption 4.1, and moreover assume*
 282 *that the function and iterates are such that Assumption 4.6 holds. Also assume that the step*
 283 *acceptance parameter η_2 satisfies*

$$284 \quad (4.12) \quad \eta_2 \geq 3\kappa_f \alpha$$

285 *and the function accuracy parameter ε_f satisfying,*

$$286 \quad (4.13) \quad \varepsilon_f \leq \min \{ \kappa_f, \eta_1 \eta_2 \}.$$

287 *Then it holds that the sequence of trust region radii $\{\delta_k\}$ satisfy the summability condition*
 288

$$289 \quad (4.14) \quad \sum_{k=0}^{\infty} \delta_k^2 < \infty$$

290 *almost surely.*

291 *Proof.* We define the constants ζ together with ν appearing in (4.11) as satisfying,

$$292 \quad (4.15) \quad \zeta \geq \max \left\{ a^{-1}, \kappa_g + \max \left\{ \eta_2, \frac{2\varepsilon_f + \kappa_g \delta_{max}}{1 - \eta_1} \right\} \right\},$$

293 where we recall that

$$294 \quad a = \frac{1}{2\alpha\kappa\delta_{max} + 2\sqrt{K}}$$

295 and

$$296 \quad (4.16) \quad \frac{\nu}{1-\nu} > \max \left\{ \frac{4\gamma^2}{\zeta c}, \frac{4\gamma^2}{\eta_1\eta_2}, \frac{\gamma^2}{\kappa_f} \right\},$$

297 with c defined by Lemma 4.3.

298 We observe that on successful, or accepted, iterations,

$$299 \quad (4.17) \quad \Phi_{k+1} - \Phi_k \leq \nu(f(x_{k+1}) - f(x_k)) + (1-\nu)(\gamma^2 - 1)\delta_k^2$$

300 and on unsuccessful iterations,

$$301 \quad (4.18) \quad \Phi_{k+1} - \Phi_k \leq (1-\nu) \left(\frac{1}{\gamma^2} - 1 \right) \delta_k^2 < 0.$$

302 Let us define the event sequence I_k as the satisfaction of model accuracy according to
303 Definition 3.2:

$$304 \quad \|\nabla F(y) - g_k\| \leq \kappa\delta_k, \quad \text{and} \quad |f(y) - f(x_k) - g_k^T(y - x_k)| \leq \kappa\|y - x_k\|\delta_k^2 \quad \forall y \in B(x_k, \delta_k).$$

305 And J_k is defined as the satisfaction of function evaluation accuracy according to Defini-
306 tion 3.1:

$$307 \quad |f_k^0 - f(x_k)| \leq \varepsilon_f\delta_k^2, \quad \text{and} \quad |f_k^s - f(x_k + s_k)| \leq \varepsilon_f\delta_k^2.$$

308 Now we break down the different cases of an approximate stationarity condition denoted
309 as:

$$310 \quad \|(\nabla f(x_k))_{I_k}\| \leq \varepsilon,$$

311 **Case 1** $\|(\nabla f(x_k))_{I_k}\| \geq \zeta\delta_k$

312 We examine the following subcases based on different events:

313 (a) $I_k \cap J_k$: The model g_k satisfies the κ - δ_k accuracy condition as well as having ε_f accurate
314 function evaluations. Applying (4.15),

$$315 \quad \|(\nabla f(x_k))_{I_k}\| \geq \delta_k/a.$$

316 Rearranging, we obtain

$$317 \quad \delta_k \leq a\|(\nabla f(x_k))_{I_k}\| \leq \frac{a \max \{ \delta_k, \alpha\|(\nabla f(x_k))_{I_k}\| \}}{\alpha}$$

318 Notice that this implies (4.4), that is,

$$319 \quad \delta_k \leq a \left\| x_k - P_{I_k} \left(x_k - \alpha \min \left\{ 1, \frac{\delta_k}{\alpha\|\nabla f(x_k)\|} \right\} \nabla f(x_k) \right) \right\|,$$

320 and so we can apply Lemma 4.3 to conclude that

$$321 \quad f(x_k) - f(\hat{x}_k) \geq \frac{1}{2\alpha} \|\hat{x}_k - x_k\|^2.$$

322 Moreover, due to model accuracy it holds that

$$323 \quad \|g_k\| \geq \|\nabla f(x_k)\| - \kappa_g\delta_k \geq (\zeta - \kappa_g)\delta_k \geq \min \left\{ \frac{1}{\eta_2}, \frac{1 - \eta_1}{2\varepsilon_f + \kappa\delta_{max}} \right\} \delta_k.$$

324 As such, we can apply Lemma 4.5 to conclude that the step is accepted and Lemma 4.3 to
 325 conclude that the stochastic process proceeds as

(4.19)

$$326 \quad \begin{aligned} \Phi_{k+1} - \Phi_k &\leq -\nu c \delta_k \left\| x_k - P_{I_k} \left(x_k - \alpha \min \left\{ 1, \frac{\delta_k}{\alpha \|\nabla f(x_k)\|} \right\} \nabla f(x_k) \right) \right\| + (1-\nu)(\gamma^2 - 1) \delta_k^2 \\ &\leq [-\nu c \zeta + (1-\nu)(\gamma^2 - 1)] \delta_k^2 < 0 \end{aligned}$$

327 where the second inequality uses the case assumption.

328 $I_k \cap J_k^c$: The function values f_k^0, f_k^s do not satisfy the ε_f -accuracy condition, while model
 329 accuracy still holds. In this case the same argument as part a holds, with the caveat that
 330 erroneous function estimates could lead to a step rejection. In that case, the change in the
 331 stochastic process is bounded by (4.18), that is,

$$332 \quad \Phi_{k+1} - \Phi_k = (1-\nu) \left(\frac{1}{\gamma^2} - 1 \right) \delta_k^2 < 0.$$

333 $I_k^c \cap J_k$: If the step is unsuccessful then again we can apply (4.18). Otherwise, with accurate
 334 function estimates, we know from Lemma 4.5 together with (4.13) that in this case

$$335 \quad \Phi_{k+1} - \Phi_k \leq [-\nu \eta_1 \eta_2 + (1-\nu)(\gamma^2 - 1)] \delta_k^2,$$

336 which is still bounded by (4.18) on account of (4.16).

337 $I_k^c \cap J_k^c$: In this case, standard Lipschitz arguments give the following bound on the increase
 338 in the value of Φ :

$$339 \quad \Phi_{k+1} - \Phi_k \leq \nu C_L \|\nabla f(x_k)\|_{I_k} \delta_k + (1-\nu)(\gamma^2 - 1) \delta_k^2, \quad C_L := \left(1 + \frac{3L}{2\zeta} \right).$$

340 We can finally combine these results to obtain, using the definitions of the probabilities θ
 341 and β ,

$$342 \quad \begin{aligned} \mathbb{E} [\Phi_{k+1} - \Phi_k | \mathcal{F}_k] &\leq \theta \beta [-\nu c \|\nabla f(x_k)\|_{I_k}] \Delta_k + (1-\nu)(\gamma^2 - 1) \Delta_k \\ &\quad + [\theta(1-\beta) + (1-\theta)\beta] (1-\nu) \left(\frac{1}{\gamma^2} - 1 \right) \Delta_k^2 \\ 344 \quad &\quad + (1-\theta)(1-\beta) \left[C_L \|\nabla f(x_k)\|_{I_k} \delta_k + (1-\nu)(\gamma^2 - 1) \Delta_k^2 \right]. \end{aligned}$$

345 We can observe that we can proceed along the same lines as the proof of Case 1 in [15,
 346 Theorem 4.11] to conclude that with θ, β chosen to satisfy

$$347 \quad (4.20) \quad \frac{(\theta\beta - 1/2)}{(1-\theta)(1-\beta)} \geq \frac{C_L}{c},$$

348 we can apply (4.16) to obtain that both

$$349 \quad (4.21) \quad \mathbb{E} \left[\Phi_{k+1} - \Phi_k | \mathcal{F}_k, \{ \|\nabla f(x_k)\|_{I_k} \geq \zeta \Delta_k \} \right] \leq -\frac{1}{4} c \nu \|\nabla f(x_k)\| \Delta_k$$

350 and

$$351 \quad (4.22) \quad \mathbb{E} \left[\Phi_{k+1} - \Phi_k | \mathcal{F}_k, \{ \|\nabla f(x_k)\|_{I_k} \geq \zeta \Delta_k \} \right] \leq -\frac{1}{2} (1-\nu)(\gamma^2 - 1) \Delta_k^2.$$

352 **Case 2:** $\|(\nabla f(x_k))_{I_k}\| < \zeta \delta_k$

353 If $\|g_k\| < \eta \delta_k$ then (4.18) holds. Now assume that $\|g_k\| \geq \eta_2 \delta_k$. We again examine the
354 following subcases based on different events:

355 (a) $I_k \cap J_k$: The model g_k satisfies the κ - δ_k accuracy condition as well as having ε_f accurate
356 function evaluations. In this case, since it cannot be ensured that the step is accepted, we
357 can apply the argument of Case 1c to conclude that again (4.18) holds.

358 (b) $I_k \cap J_k^c$: The function values f_k^0, f_k^s do not satisfy the ε_f -accuracy condition, while model
359 accuracy still holds. An unsuccessful iteration yields (4.18) a successful iteration satisfies

$$360 \quad f(x_k) - f(x_{k+1}) = f(x_k) - h_k(x_k) + h_k(x_k) - h_k(\hat{x}_k) + h_k(\hat{x}_k) - f(\hat{x}_k) \leq (\eta_2/\alpha - 2\kappa_f)\delta_k^2 \geq \kappa_f \delta_k^2$$

361 with (4.12) responsible for the last inequality. Finally (4.16) implies (4.18) holds again.

362 (c) $I_k^c \cap J_k$: It is the same as Case 1c .

363 (d) $I_k^c \cap J_k^c$: It is the same as Case 1d.

364 Now, with θ, β chosen such that

$$365 \quad (4.23) \quad (1 - \theta)(1 - \beta) \leq \frac{\gamma^2 - 1}{\gamma^4 - 1 + 2\gamma^2 C_L \zeta \frac{\nu}{1-\nu}},$$

366 we follow similar arguments to obtain

$$367 \quad (4.24) \quad \mathbb{E} \left[\Phi_{k+1} - \Phi_k | \mathcal{F}_k, \{ \|(\nabla f(x_k))_{I_k}\| < \zeta \Delta_k \} \right] \leq -\frac{1}{2}(1 - \nu) \left(1 - \frac{1}{\gamma^2} \right) \Delta_k^2.$$

368 Finally, combining the two cases yields that

$$369 \quad \mathbb{E} [\Phi_{k+1} - \Phi_k | \mathcal{F}_k] \leq -\sigma \Delta_k^2$$

370 with $\sigma > 0$, and the theorem has been proven. ■

371 We may proceed now to the main and final result. The rest of the original convergence
372 argument can be applied directly to $\|(\nabla f(x_k))_{I_k}\|$. However, recall that this is not the object
373 that is of primary interest. We are indeed interested in proving that the proposed algorithm
374 gives us a point satisfying some suitable optimality condition with high probability.

375 **Theorem 4.8.** *Almost surely,*

$$376 \quad (4.25) \quad \lim_{k \rightarrow \infty} \|(\nabla f(x_k))_{I_k}\| = 0.$$

377 Moreover, for θ sufficiently large, if it holds that, almost surely, for any limit point x^* of a
378 realization of iterates $\{x_k\}$ satisfying

$$379 \quad (4.26) \quad |\nabla f(x^*)|_{\sigma(K)} \geq |\nabla f(x^*)|_{\sigma(K+1)} + \chi, \text{ with } \chi > 0,$$

380 it holds that, for some S , for all $k \geq S$,

$$381 \quad (4.27) \quad I_k = I_{\mathcal{I}}(x^*) = I_{\mathcal{I}} \left(x^* - \frac{1}{L} \nabla f(x^*) \right)$$

382 and x^* satisfies L -stationarity. Moreover at least one such limit point exists.

383 *Proof.* The first part of the statement follows directly from the identical arguments in [15,
384 Theorem 16, Lemma 17, Theorem 18].

385 For the second statement: first observe that $\Delta_k \rightarrow 0$ almost surely and thus $\|X_{k+1} - X_k\|$
386 almost surely, and so on a set of dense probability, $\{X_k\}$ is a Cauchy sequence. As such, for

any realization there exists a limit point x^* satisfying $x_k \rightarrow x^*$. Now fix the realization for the remainder of the proof.

We compare the ranking of the gradient components, that is $\sigma(\{|g_i|\})$, $\sigma \in \bar{\Sigma}(\{|[g_k]_i|\})$ to $\sigma(\{|(\nabla f(x^*))_i|\})$. To begin with we see that for the subsequence \mathcal{S}_g wherein the model is $\kappa - \delta$ accurate we have that $k \in \mathcal{S}_g$ iterations satisfy

$$([g_k]_i - [\nabla f(x_k)]_i) + ([\nabla f(x_k)]_i - [\nabla f(x^*)]_i) \rightarrow 0$$

where the first summand goes to zero from $\delta_k \rightarrow 0$ and the second from the continuity of ∇f and the convergence of $x_k \rightarrow x^*$. Thus for sufficiently large \bar{S} , for $k \geq \bar{S}$ and $k \in \mathcal{S}_g$, it holds that

$$|[g_k]_i| > |\nabla f(x^*)|_{\sigma(K+1)} + \chi/2$$

for $i \in I_{\mathcal{I}}(x^*)$, and

$$|[g_k]_i| < |\nabla f(x^*)|_{\sigma(K)} + \chi/2$$

for $i \in I_{\mathcal{A}}(x^*)$. Thus, with probability θ , $\sigma_k \in \bar{\Sigma}_k$ satisfies that $\sigma_k[1 : K] = I_{\mathcal{I}}(x^*)$.

When θ is sufficiently large, it holds that for $k \geq \bar{S}$ sufficiently large, by smoothing properties [17], $\hat{\sigma}^k$ satisfies $\{\hat{\sigma}_{(1)}^k, \dots, \hat{\sigma}_{(K)}^k\} = I_{\mathcal{I}}(x^*)$.

This together with Lemma 2.3 proves the statement (4.27). ■

The restriction on θ is just that $\theta > \frac{1}{2}$ if all the components are separated, i.e.,

$$[|\nabla f(x^*)|]_{\sigma(1)} > [|\nabla f(x^*)|]_{\sigma(2)} > [|\nabla f(x^*)|]_{\sigma(3)} > \dots > [|\nabla f(x^*)|]_{\sigma(n)}$$

A larger θ would be necessary otherwise, in case ties prevent a unique $\hat{\sigma}^k$.

5 Numerical Results In this section, we present two machine learning applications of the algorithm 3.1: adversarial attacks on neural networks and the reconstruction of sparse Gaussian graphical models. The implementation was carried out using the Python programming language, using the NumPy, Keras, Tensorflow, scikit-learn, and Pandas libraries. The hyperparameters were selected as follows: $\eta_1 = 10^{-4}$, $\eta_2 = 10^{-4}$, $\delta_0 = 1$, $\delta_{\max} = 10$, and $\gamma = 2$. All the experiments were conducted on a machine equipped with an 11th Gen Intel(R) Core(TM) i7-1165G7 CPU @ 2.80GHz (1.69 GHz). The code is available at https://github.com/Berga53/Probabilistic_iterative_hard_thresholding.

Both applications involve high-dimensional data, making the use of the Pseudo Hard Thresholding operator, as defined in 3, computationally expensive. For practical implementation, we instead utilize the classic Hard Thresholding operator [3]. However, tests on smaller instances have shown that the two operators perform similarly when a suitable value of α_s is chosen.

5.1 Adversarial Attacks on Neural Networks Adversarial attacks are techniques used to craft imperceptible perturbations that, when added to regular data inputs, induce misclassifications in neural network models. These perturbations are typically designed to evade human detection while successfully fooling the model's classification process. One of the most powerful type of adversarial attack is the Carlini and Wagner [12], characterized by the following formulation:

$$(5.1) \quad \min_{\delta} D(x, x + \delta) + c \cdot f(x + \delta) \\ \text{such that } x + \delta \in [0, 1]^n$$

with δ being the perturbation, D being usually the ℓ_2 or ℓ_0 distance, and

$$f(x) = \left(\max_{i \neq t} (F(x)_i) - F(x)_t \right)^+.$$

428 Using our algorithm, we can incorporate the ℓ_0 penalty directly in the constraint, so our
 429 final formulation of the problem is

$$\begin{aligned}
 & \min_{\|\delta\|_0 \leq K} \|\delta\|_2 + c \cdot f(x + \delta) \\
 & \text{such that } x + \delta \in [0, 1]^n
 \end{aligned}
 \tag{5.2}$$

431 In practice, this allows us to decide how many pixels to perturb during the attack. While
 432 usual attacks are trained against selected samples of the dataset, in this paper, we will demon-
 433 strate a universal adversarial attack: the attack is performed against the entirety of the
 434 dataset, producing only one global perturbation. We will show that, in both targeted and
 435 untargeted attacks, we can significantly lower a model’s accuracy using very few pixels. We
 436 tested the attack on the MNIST dataset, which consists of 60,000 images of handwritten digits
 437 (0-9) that are 28×28 pixels in size. We performed both targeted and untargeted attacks.
 438 In the targeted attack, we aimed to misclassify the images into a specific class, while in the
 439 untargeted attack, we simply aimed to cause any misclassification. However, the untargeted
 440 attack is generally a bit weaker in the context of the Carlini and Wagner Attack. We will show
 441 that, in both targeted and untargeted attacks, we can significantly lower a model’s accuracy
 442 using very few pixels. We gradually increase the sparsity constraint and observe that this
 443 gradually increases the errors made by the model. In particular, in Figure 1, we can see both
 444 the accuracy decreasing and the number of samples predicted as the attack target increasing,
 445 indicating that the attack is performed as desired.

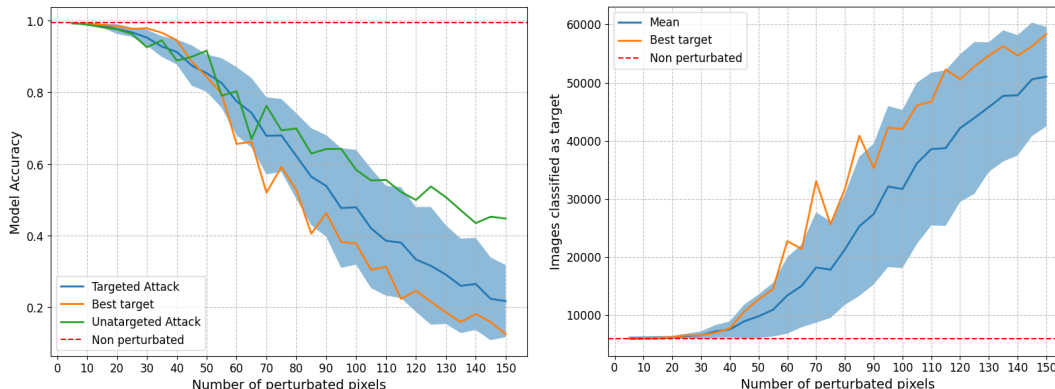


Figure 1. Effect of increasing the sparsity constraint on accuracy and targeted attack predictions.

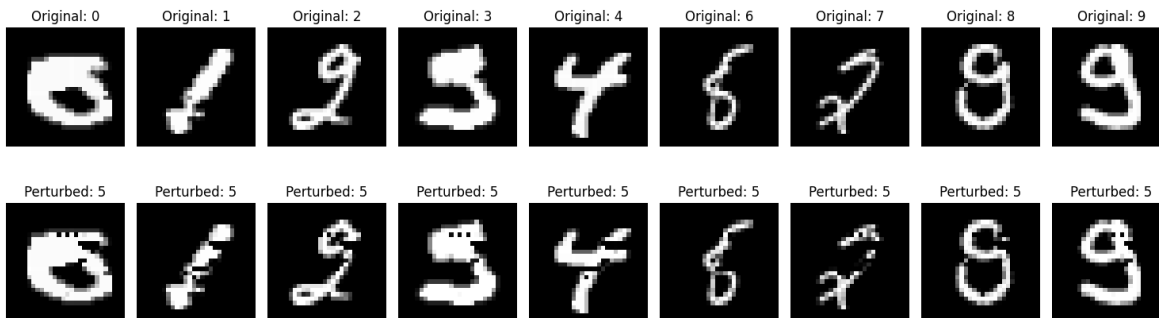


Figure 2. Example of perturbed images with $\|\delta\|_0 = 25$ and target 5

446 **5.2 Sparse Gaussian Graphical Models** Probabilistic Graphical Models are a popular
 447 tool in machine learning to model the relationships between random variables. The Gaussian
 448 Graphical Model is an undirected graph with each edge corresponding to a Gaussian condi-
 449 tional probability of one variable at the end of the edge to another. By learning the adjacency
 450 matrix together with the model weights, we can infer the proximal physical, and possibly
 451 causal, relationships between quantities.

452 This is of special importance in high dimensional settings (see, e.g., [29]). Whereas in many
 453 contemporary “big data” approaches the sample size is many orders of magnitudes larger than
 454 the dimensionality of feature space, there are a number of settings wherein obtaining data
 455 samples is costly, and such a regime cannot be expected to hold. Indeed this is often the case
 456 in medical applications, wherein recruiting volunteers for a clinical trial, or even obtaining
 457 health records, presents formidable costs to significant scaling in sample size. On the other
 458 hand, the precision of instrumentation has led to detailed “omics” data, yielding a very high
 459 dimensional feature space. One associated observation is that in the underdetermined case,
 460 when the dimensionality of the features exceeds the number of samples, some of the guarantees
 461 associated with the ℓ_1 proxy for sparsity are no longer applicable, bringing greater practical
 462 salience to having a reliable algorithm enforcing sparsity explicitly.

463 The recent work [5] presented an integer programming formulation for training sparse
 464 Gaussian graphical models. Prior to redefining the sparsity regularization using binary vari-
 465 ables, their ℓ_0 optimization problem is given as

$$466 \quad (5.3) \quad \min_{\Theta \in \mathbb{S}^p} F_0(\Theta) := \sum_{i=1}^p \left(-\log(\theta_{ii}) + \frac{1}{\theta_{ii}} \|\tilde{X}\theta_i\|^2 \right) + \lambda_0 \|\Theta\|_0 + \lambda_2 \|\Theta\|_2^2$$

467 with $\Theta \in \mathbb{S}^p$ being the weights associated with the graph and $\tilde{X} = \frac{1}{\sqrt{n}}X$ the scaled feature
 468 matrix, with $X \in \mathbb{R}^{p \times n}$ consisting of p measures and n samples. Functionally, Θ_{ij} defines
 469 an edge between node i and j in the graph, with a nonzero indicating the presence of an
 470 active edge, which corresponds to a direct link in the perspective of DAG structure of the
 471 group. The value associated with the edge corresponds to the weight defining the strength of
 472 the interaction between the features i and j . We seek to regularize cardinality for the sake
 473 of encouraging parsimonious models, as well as minimizing the total norm of the weights for
 474 general regularization.

475 Due to the structure of our algorithm, we can modify the formulation of the problem by
 476 incorporating the ℓ_0 constraint. The final formulation of the problem is then expressed as
 477 follows:

$$478 \quad (5.4) \quad \min_{\Theta \in \mathbb{S}^p, \|\Theta\|_0 \leq K} F_0(\Theta) := \sum_{i=1}^p \left(-\log(\theta_{ii}) + \frac{1}{\theta_{ii}} \|\tilde{X}\theta_i\|^2 \right) + \lambda_2 \|\Theta\|_2^2$$

479 We also observed that the ℓ_0 constraint in our formulation is very strong. In practical ap-
 480 plications, we eliminate λ_2 penalty term, as the ℓ_0 constraint was the dominant factor in the
 481 model.

482 We applied the model to the GDS2910 dataset from the Gene Expression Omnibus (GEO).
 483 This dataset consists of gene expression profiles, which naturally yield a high-dimensional
 484 feature space, with 1900 features and 191 samples. Given this feature-to-sample ratio, we can
 485 assume some level of sparsity in the final adjacency matrix. Since there is no ground truth
 486 for the underlying structure, our goal is to investigate how changing the ℓ_0 constraint affects
 487 the results of our method, while also gathering information on the true sparsity nature of
 488 the data. We performed the test by gradually increasing K , the ℓ_0 constraint, from 5000 to

489 15000. This range was previously determined to be optimal based on preliminary tests. Note
 490 that the adjacency matrix we are searching for is of size 1900×1900 , resulting in a total of
 491 3,610,000 entries. To ensure the robustness of the results, for each value of K , we performed
 492 ten runs starting from different randomly chosen feasible points, and the algorithm was given
 493 a total of 1000 iteration for every run. We also decided to set the λ_2 parameter to zero, as
 494 we observed that the strong ℓ_0 constraint was dominant over the ℓ_2 penalty.

495 We also divided the dataset into training and validation sets to determine whether the
 496 reconstructed matrix is a result of overfitting. In Figure 3, we show the effect of varying K ,
 497 which represents the number of nonzero entries that the matrix is allowed to have. The figure
 498 on the left, which shows the average objective value found over the ten runs, demonstrates
 499 that increasing K eventually stops being beneficial to the model's performance. Additionally,
 500 we observe that the number of mean accepted iterations also stops increasing, indicating that
 501 the model cannot extract more information from the data. This suggests that the true sparsity
 502 of the data can be estimated by identifying the point at which further increasing K no longer
 503 improves the model's results. In Figure 4, we present an example from our tests where the
 504 objective function decreases over the successful iterations.

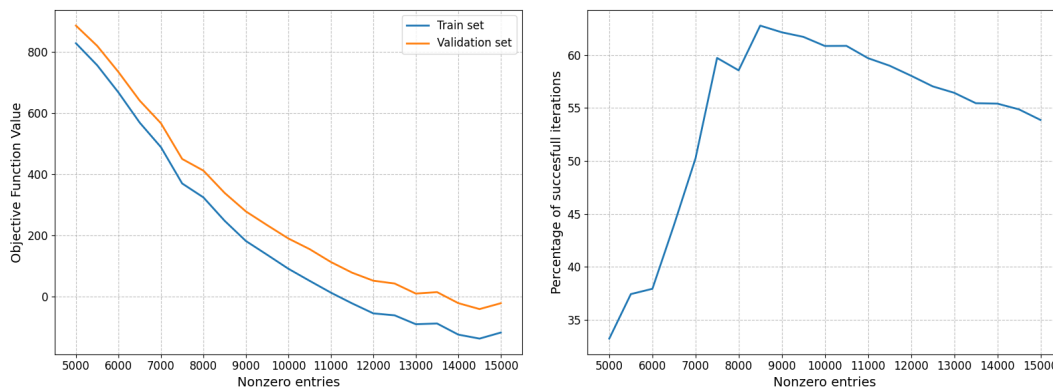


Figure 3. Effect of increasing the sparsity constraint K .

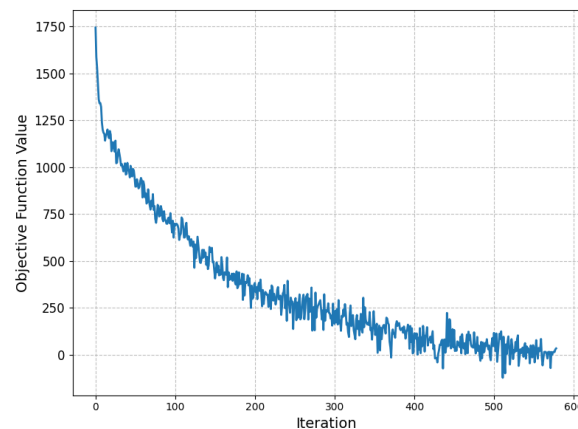


Figure 4. Objective function over the iterations.

505 **6 Conclusions** In this paper, we addressed the stochastic cardinality-constrained opti-
 506 mization problem, providing a well defined algorithm, convergence theory and illustrative
 507 experiments. Many contemporary machine learning applications involve scenarios where spar-

508 sity is crucial for high-dimensional model fitting. We proposed an iterative hard-thresholding
 509 like algorithm based on probabilistic models that nicely balances computational efficiency
 510 and solution precision by allowing flexible gradient estimates while incorporating hard spar-
 511 sity constraints.

512 We analyzed the theoretical properties of the method and proved almost sure convergence
 513 to L-stationary points under mild assumptions. This extends previous work in the optimiza-
 514 tion literature on finding solutions with strong stationarity guarantees together with machine
 515 learning articles that perform iterative hard thresholding with stochastic gradients to achieve
 516 a novel balance between ease of a fast implementation and formal guarantees of performance.
 517 The numerical experiments confirmed the practical effectiveness of our method, showcasing
 518 its potential in machine learning tasks such as adversarial attacks and probabilistic graphi-
 519 cal model training. By enforcing explicit cardinality constraints, our approach was able to
 520 produce models with enhanced sparsity and interpretability in the end.

521 Future work may involve extending the algorithm to accommodate additional nonlinear
 522 constraints, exploring techniques to further improve scalability and performance, as well as
 523 testing the algorithm on some other relevant Machine Learning applications, like, e.g., sparse
 524 Dynamic Bayesian Network training.

525

REFERENCES

- 526 [1] A. S. BANDEIRA, K. SCHEINBERG, AND L. N. VICENTE, *Convergence of trust-region methods based on*
 527 *probabilistic models*, SIAM Journal on Optimization, 24 (2014), pp. 1238–1264.
- 528 [2] F. BASTIN, C. CIRILLO, AND P. L. TOINT, *An adaptive monte carlo algorithm for computing mixed logit*
 529 *estimators*, Computational Management Science, 3 (2006), pp. 55–79.
- 530 [3] A. BECK AND Y. C. ELДАР, *Sparsity constrained nonlinear optimization: Optimality conditions and*
 531 *algorithms*, SIAM Journal on Optimization, 23 (2013), pp. 1480–1509.
- 532 [4] A. BECK AND N. HALLAK, *On the minimization over sparse symmetric sets: projections, optimality*
 533 *conditions, and algorithms*, Mathematics of Operations Research, 41 (2016), pp. 196–223.
- 534 [5] K. BEHDIN, W. CHEN, AND R. MAZUMDER, *Sparse gaussian graphical models with discrete optimization:*
 535 *Computational and statistical perspectives*, arXiv preprint arXiv:2307.09366, (2023).
- 536 [6] D. P. BERTSEKAS, *Nonlinear programming*, Athena Scientific, Belmont, MA, 1999.
- 537 [7] T. BLUMENSATH AND M. E. DAVIES, *Iterative thresholding for sparse approximations*, Journal of Fourier
 538 analysis and Applications, 14 (2008), pp. 629–654.
- 539 [8] M. BRANDA, M. BUCHER, M. ČERVINKA, AND A. SCHWARTZ, *Convergence of a scholtes-type regular-*
 540 *ization method for cardinality-constrained optimization problems with an application in sparse robust*
 541 *portfolio optimization*, Computational Optimization and Applications, 70 (2018), pp. 503–530.
- 542 [9] M. BUCHER AND A. SCHWARTZ, *Second-order optimality conditions and improved convergence results*
 543 *for regularization methods for cardinality-constrained optimization problems*, Journal of Optimization
 544 Theory and Applications, 178 (2018), pp. 383–410.
- 545 [10] O. BURDAKOV, C. KANZOW, AND A. SCHWARTZ, *Mathematical programs with cardinality constraints:*
 546 *Reformulation by complementarity-type conditions and a regularization method*, SIAM Journal on
 547 Optimization, 26 (2016), pp. 397–425.
- 548 [11] N. CARLINI AND D. WAGNER, *Towards evaluating the robustness of neural networks*, in 2017 IEEE sym-
 549 posium on security and privacy (sp), Ieee, 2017, pp. 39–57.
- 550 [12] N. CARLINI AND D. A. WAGNER, *Towards evaluating the robustness of neural networks*, CoRR,
 551 abs/1608.04644 (2016), <http://arxiv.org/abs/1608.04644>, <https://arxiv.org/abs/1608.04644>.
- 552 [13] C. CARTIS AND K. SCHEINBERG, *Global convergence rate analysis of unconstrained optimization methods*
 553 *based on probabilistic models*, Mathematical Programming, 169 (2018), pp. 337–375.
- 554 [14] M. ČERVINKA, C. KANZOW, AND A. SCHWARTZ, *Constraint qualifications and optimality conditions for*
 555 *optimization problems with cardinality constraints*, Mathematical Programming, 160 (2016), pp. 353–
 556 377.
- 557 [15] R. CHEN, M. MENICKELLY, AND K. SCHEINBERG, *Stochastic optimization using a trust-region method*
 558 *and random models*, Mathematical Programming, 169 (2018), pp. 447–487.
- 559 [16] F. CROCE AND M. HEIN, *Sparse and imperceivable adversarial attacks*, in Proceedings of the IEEE/CVF
 560 international conference on computer vision, 2019, pp. 4724–4732.
- 561 [17] E. S. GARDNER JR, *Exponential smoothing: The state of the art*, Journal of forecasting, 4 (1985), pp. 1–
 562 28.

- 563 [18] N. HALLAK, *A path-based approach to constrained sparse optimization*, SIAM Journal on Optimization,
564 34 (2024), pp. 790–816.
- 565 [19] P. JAIN, A. TEWARI, AND P. KAR, *On iterative hard thresholding methods for high-dimensional m-*
566 *estimation*, Advances in neural information processing systems, 27 (2014).
- 567 [20] C. KANZOW, A. B. RAHARJA, AND A. SCHWARTZ, *Sequential optimality conditions for cardinality-*
568 *constrained optimization problems with applications*, Computational Optimization and Applications,
569 80 (2021), pp. 185–211.
- 570 [21] S. LÄMMEL AND V. SHIKHMAN, *On nondegenerate m-stationary points for sparsity constrained nonlinear*
571 *optimization*, Journal of Global Optimization, 82 (2022), pp. 219–242.
- 572 [22] S. LÄMMEL AND V. SHIKHMAN, *Critical point theory for sparse recovery*, Optimization, 72 (2023), pp. 521–
573 549.
- 574 [23] M. LAPUCCI, T. LEVATO, F. RINALDI, AND M. SCIANDRONE, *A unifying framework for sparsity-*
575 *constrained optimization*, Journal of Optimization Theory and Applications, 199 (2023), pp. 663–692.
- 576 [24] M. LAPUCCI, T. LEVATO, AND M. SCIANDRONE, *Convergent inexact penalty decomposition methods*
577 *for cardinality-constrained problems*, Journal of Optimization Theory and Applications, 188 (2021),
578 pp. 473–496.
- 579 [25] Z. LU AND Y. ZHANG, *Sparse approximation via penalty decomposition methods*, SIAM Journal on Op-
580 timization, 23 (2013), pp. 2448–2478.
- 581 [26] A. MODAS, S.-M. MOOSAVI-DEZFOOLI, AND P. FROSSARD, *Sparsefool: a few pixels make a big differ-*
582 *ence*, in Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, 2019,
583 pp. 9087–9096.
- 584 [27] T. MURATA AND T. SUZUKI, *Sample efficient stochastic gradient iterative hard thresholding method for*
585 *stochastic sparse linear regression with limited attribute observation*, Advances in Neural Information
586 Processing Systems, 31 (2018).
- 587 [28] M. M. NEGRI, F. AREND TORRES, AND V. ROTH, *Conditional matrix flows for gaussian graphical models*,
588 Advances in Neural Information Processing Systems, 36 (2023), pp. 25095–25111.
- 589 [29] M. J. WAINWRIGHT, *High-dimensional statistics: A non-asymptotic viewpoint*, vol. 48, Cambridge uni-
590 versity press, 2019.
- 591 [30] B. ZHOU, F. CHEN, AND Y. YING, *Stochastic iterative hard thresholding for graph-structured sparsity*
592 *optimization*, in International Conference on Machine Learning, PMLR, 2019, pp. 7563–7573.
- 593 [31] P. ZHOU, X. YUAN, AND J. FENG, *Efficient stochastic gradient hard thresholding*, Advances in Neural
594 Information Processing Systems, 31 (2018).