

SPARSE CANONICAL CORRELATION ANALYSIS VIA SMOOTH NON-CONVEX ℓ_0 SURROGATES AND ITERATIVE MINORIZATION–MAXIMIZATION

Anonymous authors

Paper under double-blind review

ABSTRACT

Canonical correlation analysis (CCA) is a core tool to uncover linear associations between two datasets. In high-dimensional settings, however, it is prone to overfitting and lacks interpretability. Enforcing exact sparsity via ℓ_0 constraints can improve interpretability but leads to an intractable combinatorial problem. We propose a novel framework for sparse CCA that replaces the ℓ_0 cardinality constraint with tight smooth concave surrogates (power, logarithmic, and exponential forms), preserving support control without ad hoc thresholds. We solve the resulting nonconvex program via a minorization–maximization algorithm, yielding a generalized eigenvalue subproblem at each step. We prove that as the smoothing parameter vanishes, the surrogate formulation converges to the exact ℓ_0 solution with explicit suboptimality bounds. We further reformulate the objective as a rank-constrained semidefinite program and use randomized Gaussian rounding to extract sparse canonical directions. Empirical results on six benchmark datasets demonstrate that our method enforces exact sparsity levels, delivers superior canonical correlations and support recovery, and offers markedly improved scalability compared to state-of-the-art SCCA algorithms.

1 INTRODUCTION

Canonical correlation analysis (CCA) has long stood as a cornerstone of multivariate statistics, tracing back to Hotelling’s original formulation in (Hotelling, 1935), wherein one seeks pairs of linear projections—one from each of two datasets—that maximize their mutual correlation. Over the past two decades, CCA has found myriad applications in areas as diverse as genomics, neuroimaging, and multimedia retrieval Huang et al. (2010); Vinokourov et al. (2002); Hermansky & Morgan (1994). Yet, where the number of features far exceeds the available samples, the classical CCA solution becomes both numerically unstable and densely supported, severely limiting scientific interpretability and risking overfitting.

Prior work on sparse CCA began with Parkhomenko et al. (2007), who formulated a genome-wide ℓ_0 -constrained CCA and used a greedy feature-selection heuristic that offered no optimality guarantees and did not scale well beyond a few dozen variables. Subsequent work replaced the combinatorial constraint with convex surrogates: Waaijenborg et al. (2008) introduced an elastic-net-penalized bi-convex formulation; Witten & Tibshirani (2009) developed a penalized matrix decomposition approach with LASSO-style regularization; Hardoon & Shawe-Taylor (2011) proposed a kernelized convex sparse CCA framework, and Lin et al. (2013; 2014) incorporated structured group penalties to exploit known feature groupings. Although each deal with a tractable program, they loosely approximate the true ℓ_0 constraint, introduce shrinkage bias, require careful tuning of multiple regularization parameters, and get trapped in suboptimal local minima (as they do not handle the cardinality constraint directly).

More recently, integer- and semidefinite-programming approaches have been developed for sparse CCA. Bertsimas et al. (2016) formulated the problem as a mixed-integer program solved via branch-and-cut, but this approach incurs exponential worst-case complexity and heavy memory usage for storing large branch-and-bound trees, making it impractical in high dimensions. Watanabe et al. (2023) advanced this line with a semidefinite-relaxation-based branch-and-bound algorithm that

054 guarantees correctness on small- to medium-scale problems, yet it relies on large SDP solves with
 055 substantial memory requirements and a costly separation oracle. Building on these, Li et al. (2024)
 056 proposed a general mixed-integer semidefinite program for ℓ_0 -regularized CCA with a cutting-plane
 057 procedure; in the full high-dimensional setting, their formulation demands storing exponentially
 058 many cuts and large semidefinite matrices, leading to prohibitive memory consumption and no
 059 support from off-the-shelf solvers. To mitigate this, they also introduced greedy and local-search
 060 heuristics, despite lack of convergence guarantees.

061 Building on the limitations of existing sparse-CCA approaches, we develop a unified framework that
 062 directly enforces ℓ_0 sparsity via tight and smooth surrogates, and an efficient MM scheme:
 063

- 064 • **Novel smooth ℓ_0 surrogates.** We introduce three continuously differentiable concave sur-
 065 rogates (power-law, normalized logarithmic, and exponential forms) that uniformly approx-
 066 imate the discontinuous cardinality function while remaining C^1 near zero, thus, avoiding
 067 IRLS singularities and eliminating ad hoc thresholding.
- 068 • **MM-based sparse generalized eigenproblems.** We embed these surrogates into a mi-
 069 norization–maximization algorithm: at each iteration we construct a quadratic minorizer
 070 of the surrogate penalties, yielding a tractable generalized eigenvalue subproblem whose
 071 solution enforces the exact user-specified sparsity level.
- 072 • **Convergence and suboptimality analysis.** We prove that as the smoothing parameter
 073 $\varepsilon \rightarrow 0$, our surrogate problem converges to the original ℓ_0 -constrained formulation, and we
 074 derive explicit bounds quantifying the maximal gap between the two solutions.
- 075 • **SDP reformulation & randomized rounding.** We transform the smoothed SCCA into a
 076 rank-constrained semidefinite program, then, relax the rank condition and apply Gaussian
 077 randomization to extract high-quality sparse canonical directions with provable guarantees.
 078
- 079 • **Exact low-rank solver & branch-and-cut.** In the special case where the marginal covari-
 080 ance ranks do not exceed the sparsity levels, we show SCCA reduces to a polynomial-time
 081 $O(n^3 + m^3)$ procedure. For the general case, we derive a mixed-integer SDP and imple-
 082 ment a custom branch-and-cut with closed-form cuts to solve moderate-scale instances to
 083 global optimality.
- 084 • **Extensive empirical validation.** On six benchmark UCI datasets Blake (1998), our
 085 method consistently achieves the highest canonical correlations and precise support recov-
 086 ery, all while running at least two orders of magnitude faster than competing exact solvers.
 087

088 **Organization** The remainder of the paper is structured as follows. In Section 2 we formalize
 089 the sparse CCA problem and introduce our family of smooth ℓ_0 surrogates. Section 3 presents
 090 our proposed algorithm, including the construction of quadratic minorizers. Section 4 reports com-
 091 prehensive numerical experiments on diverse datasets, comparing against state-of-the-art baselines.
 092 Finally, Section 5 concludes with a summary of findings and directions for future work.
 093

094 2 PROBLEM FORMULATION

095 Sparse Canonical Correlation Analysis (SCCA) enforces exact sparsity on the canonical loading
 096 vectors to improve interpretability. In particular, one seeks
 097

$$098 \quad v^* = \max_{\mathbf{x} \in \mathbb{R}^n, \mathbf{y} \in \mathbb{R}^m} \left\{ \mathbf{x}^T \mathbf{A} \mathbf{y} : \mathbf{x}^T \mathbf{B} \mathbf{x} \leq 1, \mathbf{y}^T \mathbf{C} \mathbf{y} \leq 1, \|\mathbf{x}\|_0 \leq s_1, \|\mathbf{y}\|_0 \leq s_2 \right\}, \quad (1)$$

099 where $s_1 \leq n$ and $s_2 \leq m$ are user-specified sparsity levels, \mathbf{B} and \mathbf{C} are the marginal covariance
 100 matrices, and $\mathbf{A} \in \mathbb{R}^{n \times m}$ is the cross-covariance. Importantly, we impose no rank or definiteness
 101 restrictions on \mathbf{B} and \mathbf{C} .

102 Because $\|\mathbf{x}\|_0 = \sum_{i=1}^n \text{sgn}(|x_i|)$, problem 1 combines a nonconcave objective with a discontinuous
 103 penalty, making direct optimization intractable. To address this, we replace each indicator $\text{sgn}(|x_i|)$
 104 by a tight continuous surrogate $g_p(x)$, where g_p is even, concave, differentiable except at zero,
 105 nondecreasing on $[0, \infty)$, and satisfies $g_p(0) = 0$. In particular, we employ three well-studied
 106
 107

108 surrogates:

$$109 \quad g_p(x) = |x|^p, \quad 0 < p \leq 1,$$

$$110 \quad g_p(x) = \frac{\log(1 + \frac{|x|}{p})}{\log(1 + 1/p)}, \quad p > 0,$$

$$111 \quad g_p(x) = 1 - e^{-|x|/p}, \quad p > 0.$$

112 The first is a p -quasi-norm Gorodnitsky & Rao (1997); Chartrand & Yin (2008), the second is a
 113 normalized logarithmic penalty underlying iteratively reweighted ℓ_1 schemes Candès et al. (2008);
 114 Sriperumbudur et al. (2011), and the third is an exponential lower-bound surrogate Fischer et al.
 115 (1996).

116 Substituting $\|\mathbf{x}\|_0 \approx \sum_i g_p(x_i)$ and $\|\mathbf{y}\|_0 \approx \sum_j g_p(y_j)$ into problem 1 yields the continuous (yet
 117 still nonconvex and nondifferentiable) approximation

$$118 \quad \max_{\mathbf{x}, \mathbf{y}} \quad \mathbf{x}^T \mathbf{A} \mathbf{y} - \rho_1 \sum_{i=1}^n g_p(x_i) - \rho_2 \sum_{j=1}^m g_p(y_j) \quad (2)$$

$$119 \quad \text{s.t.} \quad \mathbf{x}^T \mathbf{B} \mathbf{x} \leq 1, \quad \mathbf{y}^T \mathbf{C} \mathbf{y} \leq 1,$$

120 where $\rho_1, \rho_2 > 0$ are regularization parameters. In the next section, we develop an MM-based
 121 algorithm to solve problem 2 by constructing at each iteration a quadratic minorizer of the surrogate
 122 penalties and then maximizing the resulting generalized eigenvalue subproblem to enforce the exact
 123 sparsity levels. For an overview of the MM framework Sun et al. (2017); Saini et al. (2024), see
 124 Appendix A.

125 3 SOLVING THE SCCA PROBLEM

126 3.1 QUADRATIC BOUNDING OF SURROGATE PENALTIES

127 When applying MM to our surrogate-regularized SCCA formulation 2, we keep the term $\mathbf{x}^T \mathbf{A} \mathbf{y}$
 128 intact and upper-bound only each concave penalty $g_p(x_i)$ with a quadratic tangent. Concretely, at
 129 iteration k we replace

$$130 \quad g_p(x_i) \quad \mapsto \quad w_i^{(k)} x_i^2 + c_i^{(k)},$$

131 where the coefficients are chosen to match both value and slope at $x_i^{(k)}$:

$$132 \quad g_p(x_i^{(k)}) = w_i^{(k)} (x_i^{(k)})^2 + c_i^{(k)}, \quad (3)$$

$$133 \quad g_p'(x_i^{(k)}) = 2 w_i^{(k)} x_i^{(k)}. \quad (4)$$

134 By concavity, this quadratic form satisfies $w_i^{(k)} x_i^2 + c_i^{(k)} \geq g_p(x_i)$ for all x_i , transforming the
 135 original problem into a tractable quadratically-constrained quadratic subproblem at each MM step.

136 **Illustration for the power-law surrogate** Take $g_p(x) = |x|^p$ with $0 < p \leq 1$. To build the
 137 quadratic upper-bound at the current iterate $x_i^{(k)} \neq 0$, we match both value and derivative:

$$138 \quad |x_i^{(k)}|^p = w_i^{(k)} (x_i^{(k)})^2 + c_i^{(k)}, \quad p \operatorname{sgn}(x_i^{(k)}) |x_i^{(k)}|^{p-1} = 2 w_i^{(k)} x_i^{(k)}.$$

139 Solving these two equations gives

$$140 \quad w_i^{(k)} = \frac{p}{2} |x_i^{(k)}|^{p-2}, \quad c_i^{(k)} = (1 - \frac{p}{2}) |x_i^{(k)}|^p,$$

141 so that the quadratic form $u(x; x_i^{(k)}) = \frac{p}{2} |x_i^{(k)}|^{p-2} x^2 + (1 - \frac{p}{2}) |x_i^{(k)}|^p$ satisfies $u(x; x_i^{(k)}) \geq |x|^p$
 142 for all x .

143 This construction underlies the classic iteratively reweighted least-squares (IRLS) schemes
 144 in robust regression and sparse recovery Holland & Welsch (1977); Schlossmacher (1973);
 145 Gorodnitsky & Rao (1997); Chartrand & Yin (2008). However, if $x_i^{(k)} = 0$, the weight $w_i^{(k)}$ be-
 146 comes singular. A common patch is to add a small damping factor $\epsilon > 0$,

$$147 \quad w_i^{(k)} = \frac{p}{2} ((x_i^{(k)})^2 + \epsilon)^{\frac{p-2}{2}},$$

148 which prevents a potential blow-up but no longer guarantees a true majorizer of $|x|^p$.

3.2 SMOOTH SURROGATES FOR NON-DIFFERENTIABLE PENALTIES

Inspired by Song et al. (2015), we eliminate singular weights in the IRLS-style quadratic bounds by replacing each concave surrogate $g_p(x)$ with a continuously differentiable proxy $g_p^\epsilon(x)$. This proxy matches g_p outside a small neighborhood of zero and becomes strictly quadratic within $|x| \leq \epsilon$. Specifically, for $\epsilon > 0$ define

$$g_p^\epsilon(x) = \begin{cases} \frac{g_p'(\epsilon)}{2\epsilon} x^2, & |x| \leq \epsilon, \\ g_p(x) - g_p(\epsilon) + \frac{g_p'(\epsilon)\epsilon}{2}, & |x| > \epsilon. \end{cases} \quad (5)$$

This construction ensures $g_p^\epsilon \in C^1$ and $g_p^\epsilon(x) \rightarrow g_p(x)$ uniformly as $\epsilon \rightarrow 0$. For example, when $g_p(x) = |x|^p$ ($0 < p \leq 1$), one obtains

$$g_p^\epsilon(x) = \begin{cases} \frac{p}{2} \epsilon^{p-2} x^2, & |x| \leq \epsilon, \\ |x|^p - (1 - \frac{p}{2}) \epsilon^p, & |x| > \epsilon. \end{cases} \quad (6)$$

Inserting g_p^ϵ into the original surrogate-regularized CCA problem 2 yields the smoothed formulation

$$\begin{aligned} \max_{\mathbf{x}, \mathbf{y}} \quad & \mathbf{x}^T \mathbf{A} \mathbf{y} - \rho_1 \sum_{i=1}^n g_p^\epsilon(x_i) - \rho_2 \sum_{j=1}^m g_p^\epsilon(y_j), \\ \text{s.t.} \quad & \mathbf{x}^T \mathbf{B} \mathbf{x} \leq 1, \quad \mathbf{y}^T \mathbf{C} \mathbf{y} \leq 1. \end{aligned} \quad (7)$$

The MM step then, simply constructs tangent-quadratic upper bounds of each $g_p^\epsilon(x_i)$, whose coefficients remain finite for all x_i .

Approximation error. It can be shown that the gap between the smoothed and original surrogate objectives is bounded by $O(\rho n (g_p(\epsilon) - \frac{g_p'(\epsilon)\epsilon}{2}))$, which vanishes as $\epsilon \rightarrow 0$. Thus, solving 7 to high accuracy recovers an arbitrarily good approximation of the true ℓ_0 -penalized solution without any singularity issues. For proof, see Appendix B.

3.3 ITERATIVELY REWEIGHTED QUADRATIC MINORIZATION

Having introduced the smooth surrogate g_p^ϵ in equation 5 and its quadratic upper-bounds, we now describe the full MM iteration for the smoothed problem 7. Starting from an initial guess $(\mathbf{x}^{(0)}, \mathbf{y}^{(0)})$, each iteration k proceeds as follows:

1. **Weight update.** Following equation 4, compute the weight for each coordinate as follows:

$$w_i^{(k)} = \frac{g_p^{\epsilon'}(x_i^{(k)})}{2x_i^{(k)}}, \quad z_j^{(k)} = \frac{g_p^{\epsilon'}(y_j^{(k)})}{2y_j^{(k)}}, \quad i = 1, \dots, n, \quad j = 1, \dots, m.$$

Because g_p^ϵ is C^1 and strictly quadratic near zero, these weights are always finite.

2. **Minorized subproblem.** Replace each penalty term by its quadratic tangent:

$$g_p^\epsilon(x_i) \leq w_i^{(k)} x_i^2 + c_i^{(k)}, \quad g_p^\epsilon(y_j) \leq z_j^{(k)} y_j^2 + d_j^{(k)},$$

and drop the constant offsets $c_i^{(k)}, d_j^{(k)}$. We then solve

$$\begin{aligned} (\mathbf{x}^{(k+1)}, \mathbf{y}^{(k+1)}) = \arg \max_{\mathbf{x}, \mathbf{y}} \quad & \mathbf{x}^T \mathbf{A} \mathbf{y} - \rho_1 \mathbf{x}^T [\text{Diag}(\mathbf{w}^{(k)})] \mathbf{x} - \rho_2 \mathbf{y}^T [\text{Diag}(\mathbf{z}^{(k)})] \mathbf{y}, \\ \text{s.t.} \quad & \mathbf{x}^T \mathbf{B} \mathbf{x} \leq 1, \quad \mathbf{y}^T \mathbf{C} \mathbf{y} \leq 1. \end{aligned} \quad (8)$$

This is a quadratically-constrained quadratic program in (\mathbf{x}, \mathbf{y}) .

Table 1: Smooth approximation $g_p^\epsilon(x_i)$ of the surrogate functions $g_p(x_i)$ and the quadratic majorization functions, where $u(x_i, x_i^{(k)}) = w_i^{(k)}x_i^2 + c_i^{(k)}$ at $x_i^{(k)}$.

Surrogate function $g_p(x_i)$	Smooth approximation $g_p^\epsilon(x_i)$	$w_i^{(k)}$
$ x_i ^p, 0 < p \leq 1$	$\begin{cases} \frac{p}{2}\epsilon^{p-2}x_i^2, & \text{if } x_i \leq \epsilon, \\ x_i ^p - (1 - \frac{p}{2})\epsilon^p, & \text{if } x_i > \epsilon, \end{cases}$	$\begin{cases} \frac{p}{2}\epsilon^{p-2}, & \text{if } x_i^{(k)} \leq \epsilon, \\ \frac{p}{2} x_i^{(k)} ^{p-2}, & \text{if } x_i^{(k)} > \epsilon. \end{cases}$
$\frac{\log(1 + \frac{ x_i }{p})}{\log(1 + 1/p)}, p > 0$	$\begin{cases} \frac{x_i^2}{2\epsilon(p+\epsilon)\log(1+1/p)}, & \text{if } x_i \leq \epsilon, \\ \frac{\log(1 + \frac{ x_i }{p}) - \log(1 + \epsilon/p) + \frac{\epsilon}{2(p+\epsilon)}}{\log(1+1/p)}, & \text{if } x_i > \epsilon, \end{cases}$	$\begin{cases} \frac{1}{2\epsilon(p+\epsilon)\log(1+1/p)}, & \text{if } x_i^{(k)} \leq \epsilon, \\ \frac{1}{2\log(1+1/p) x_i^{(k)} (x_i^{(k)} +p)}, & \text{if } x_i^{(k)} > \epsilon. \end{cases}$
$1 - e^{- x_i /p}, p > 0$	$\begin{cases} \frac{e^{-\epsilon/p}}{2p\epsilon}x_i^2, & \text{if } x_i \leq \epsilon, \\ -e^{- x_i /p} + (1 + \frac{\epsilon}{2p})e^{-\epsilon/p}, & \text{if } x_i > \epsilon, \end{cases}$	$\begin{cases} \frac{e^{-\epsilon/p}}{2p\epsilon}, & \text{if } x_i^{(k)} \leq \epsilon, \\ -\frac{ x_i^{(k)} }{2p x_i^{(k)} /p}, & \text{if } x_i^{(k)} > \epsilon. \end{cases}$

The QCQP in problem 8 can be compactly written by stacking \mathbf{x} and \mathbf{y} into a single vector $\mathbf{u} = [\mathbf{x}^T, \mathbf{y}^T]^T \in \mathbb{R}^{n+m}$. Define the block matrices

$$\tilde{\mathbf{A}}^{(k)} = \begin{pmatrix} -\rho_1 \text{Diag}(\mathbf{w}^{(k)}) & \frac{1}{2} \mathbf{A} \\ \frac{1}{2} \mathbf{A}^T & -\rho_2 \text{Diag}(\mathbf{z}^{(k)}) \end{pmatrix}, \quad \tilde{\mathbf{B}} = \begin{pmatrix} \mathbf{B} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix}, \quad \tilde{\mathbf{C}} = \begin{pmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{C} \end{pmatrix}.$$

Then, problem 8 is equivalent to

$$\max_{\mathbf{u} \in \mathbb{R}^{n+m}} \mathbf{u}^T \tilde{\mathbf{A}}^{(k)} \mathbf{u} \quad \text{s.t.} \quad \mathbf{u}^T \tilde{\mathbf{B}} \mathbf{u} \leq 1, \quad \mathbf{u}^T \tilde{\mathbf{C}} \mathbf{u} \leq 1.$$

Introducing the rank-one matrix $\mathbf{U} = \mathbf{u} \mathbf{u}^T$, we can further simplify the optimization problem as

$$\begin{aligned} \max_{\mathbf{U} \in \mathcal{S}_+^{n+m}} \quad & \text{tr}(\tilde{\mathbf{A}}^{(k)} \mathbf{U}) \\ \text{s.t.} \quad & \text{tr}(\tilde{\mathbf{B}} \mathbf{U}) \leq 1 \\ & \text{tr}(\tilde{\mathbf{C}} \mathbf{U}) \leq 1, \\ & \text{rank}(\mathbf{U}) = 1. \end{aligned} \tag{9}$$

By dropping the non-convex $\text{rank}(\mathbf{U}) = 1$ constraint, we arrive at the convex SDP

$$\begin{aligned} \max_{\mathbf{U}} \quad & \text{trace}(\tilde{\mathbf{A}}^{(k)} \mathbf{U}) \\ \text{s.t.} \quad & \text{trace}(\tilde{\mathbf{B}} \mathbf{U}) \leq 1, \\ & \text{trace}(\tilde{\mathbf{C}} \mathbf{U}) \leq 1, \\ & \mathbf{U} \succeq \mathbf{0}, \end{aligned} \tag{10}$$

which we solve with any SDP solver to obtain $\mathbf{U}^* \succeq \mathbf{0}$. We then, apply the Gaussian randomization technique by drawing $\begin{pmatrix} \mathbf{x} \\ \mathbf{y} \end{pmatrix} \sim \mathcal{N}(\mathbf{0}, \mathbf{U}^*)$ as explained in Luo et al. (2010). A concise overview of our approach is given in Algorithm 1. The detailed proof of convergence for our proposed method is provided in Appendix C. The proof demonstrates that the sequence of objective values is non-decreasing and upper-bounded, and that every limit point of the iterates is a KKT stationary point.

4 NUMERICAL RESULTS

We evaluate the performance of the proposed sparse CCA method on six benchmark datasets, comparing it against three established baseline methods. All experiments are conducted using MATLAB R2022b on a dual-socket Intel Xeon E5-2695 v3 system (2×14 physical cores, 56 threads total, 2.3 GHz base frequency, up to 3.3 GHz turbo boost, 70 MiB L3 cache) with 256 GB of RAM.

4.1 DATASETS

We evaluate our proposed method on six benchmark UCI datasets Blake (1998); Dheeru & Karra Taniskidou (2019) commonly used in sparse CCA studies. These datasets

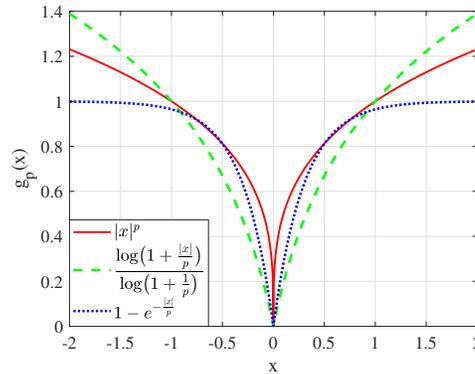


Figure 1: Three surrogate functions $g_p(x)$ that are used for approximating $\text{sgn}(|x|)$, $p = 0.3$.

vary widely in the number of features, sample sizes, and domain characteristics. Below is a brief description of each dataset:

- **Dermatology** Blake (1998); Dheeru & Karra Taniskidou (2019): Contains 366 patient records with 34 total features. We split the features into two equal subsets of 17 dimensions each.
- **Spambase** Blake (1998); Dheeru & Karra Taniskidou (2019): Includes 4601 emails represented by 57 total frequency-based features; we split it into two subsets of 28 and 29 dimensions.
- **Digits** Dheeru & Karra Taniskidou (2019): Comprises of 1797 handwritten digit samples, each described by 64 features partitioned evenly into two 32-dimensional parts.
- **Buzz in Social Media** Blake (1998); Dheeru & Karra Taniskidou (2019): A large dataset with 583250 samples and 77 features, split into 39 and 38-dimensional views.
- **Gas Sensor Array Drift** Blake (1998); Vergara et al. (2012): Includes 2565 chemical sensor readings with 128 variables, separated into two views of 64 dimensions each.
- **Wikipedia Articles** Blake (1998); Dheeru & Karra Taniskidou (2019); Rasiwasia et al. (2010): Contains 2310 bilingual (English–German) document pairs, with 583 features in the English part and 250 in the German side.

It is worth mentioning that for applications involving very high-dimensional data, a common and effective strategy is to first perform dimensionality reduction. For instance, principal component analysis (PCA) can be used to project the original feature vectors onto a lower-dimensional subspace (e.g., 50 dimensions) that captures a significant portion of the data’s variance (e.g., >98%) Omati et al. (2025); Wang et al. (2024); Su et al. (2015). The resulting projected data can then be used as input for the SCCA algorithm, making the problem more computationally tractable.

4.2 COMPARED METHODS

We compare our proposed algorithm with three strong sparse CCA baselines, as follows:

- **ADMM-based SCCA** Suo et al. (2017): A proximal gradient algorithm based on the alternating direction method of multipliers (ADMM), which alternates updates of the canonical vectors using soft-thresholded projections.
- **Predictive sparse CCA** Wilms & Croux (2015): A predictive formulation of sparse CCA that employs penalized least squares with soft-thresholding, optimized via coordinate descent.
- **Branch-and-bound SCCA** Li et al. (2024): An exact solver for sparse CCA formulated as a mixed-integer optimization problem. Due to its high computational cost, which is also emphasized in the original paper, we impose a hard ceiling of 10^{10} explored nodes and a maximum runtime of 300 seconds per instance.

Algorithm 1 MM-SDP approach for solving SCCA problem

Require: Covariances $\mathbf{A} \in \mathbb{R}^{n \times m}$, $\mathbf{B} \in \mathbb{S}_+^n$, $\mathbf{C} \in \mathbb{S}_+^m$, smoothing schedule $\{\varepsilon_k\}_{k=0}^T$, regularizers (ρ_1, ρ_2) , max iters T , tolerance δ

Ensure: Sparse canonical vectors (\mathbf{x}, \mathbf{y})

- 1: Initialize $(\mathbf{x}^{(0)}, \mathbf{y}^{(0)})$ (e.g. via leading CCA)
- 2: Compute initial objective

$$f^{(0)} \leftarrow (\mathbf{x}^{(0)})^\top \mathbf{A} \mathbf{y}^{(0)} - \rho_1 \sum_i g_{\varepsilon_0}(x_i^{(0)}) - \rho_2 \sum_j g_{\varepsilon_0}(y_j^{(0)}).$$

3: **for** $k = 0, \dots, T - 1$ **do**

4: **Weight update:**

5: **for** $i = 1, \dots, n$ **do**

6: $w_i \leftarrow \frac{g'_{\varepsilon_k}(x_i^{(k)})}{2x_i^{(k)}}$

7: **end for**

8: **for** $j = 1, \dots, m$ **do**

9: $z_j \leftarrow \frac{g'_{\varepsilon_k}(y_j^{(k)})}{2y_j^{(k)}}$

10: **end for**

11: **Form SDP matrices:**

$$\tilde{\mathbf{A}} = \begin{bmatrix} -\rho_1 \text{Diag}(\mathbf{w}) & \frac{1}{2} \mathbf{A} \\ \frac{1}{2} \mathbf{A}^\top & -\rho_2 \text{Diag}(\mathbf{z}) \end{bmatrix}, \quad \tilde{\mathbf{B}} = \begin{bmatrix} \mathbf{B} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix}, \quad \tilde{\mathbf{C}} = \begin{bmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{C} \end{bmatrix}.$$

12: **Solve**

$$\mathbf{U}^* = \arg \max_{\mathbf{U} \succeq \mathbf{0}} \langle \tilde{\mathbf{A}}, \mathbf{U} \rangle \quad \text{s.t.} \quad \langle \tilde{\mathbf{B}}, \mathbf{U} \rangle \leq 1, \quad \langle \tilde{\mathbf{C}}, \mathbf{U} \rangle \leq 1.$$

13: **Randomized rounding:** extract $(\mathbf{x}^{(k+1)}, \mathbf{y}^{(k+1)})$ from \mathbf{U}^*

14: Compute new objective

$$f^{(k+1)} \leftarrow (\mathbf{x}^{(k+1)})^\top \mathbf{A} \mathbf{y}^{(k+1)} - \rho_1 \sum_i g_{\varepsilon_k}(x_i^{(k+1)}) - \rho_2 \sum_j g_{\varepsilon_k}(y_j^{(k+1)}).$$

15: **if** $|f^{(k+1)} - f^{(k)}| < \delta$ **then**

16: **break**

▷ stop when objective change is below tolerance

17: **end if**

18: **end for**

19: **return** $(\mathbf{x}^{(k+1)}, \mathbf{y}^{(k+1)})$

4.3 METRICS

We use canonical correlation as the primary evaluation metric, defined as the maximum correlation between the projected views. For each algorithm, we performed a grid search over its own set of hyperparameters and report the configuration that achieves the highest correlation:

- **MM-SDP (Ours):** $(\rho_1, \rho_2) \in \{0.0001, 0.0005, 0.001, 0.005, 0.01, 0.05\}^2$.
- **ADMM-based SCCA** Suo et al. (2017): $(\lambda_1, \lambda_2) \in \{0.0001, 0.0005, 0.001, 0.005, 0.01, 0.05\}^2$.
- **Predictive Sparse CCA** Wilms & Croux (2015): $(\alpha_1, \alpha_2) \in \{0.0001, 0.0005, 0.001, 0.005, 0.01, 0.05\}^2$.
- **Branch-and-Bound SCCA** Li et al. (2024): sparsity levels $s_1 = s_2 \in \{2, 3, 4, 5, 6, 7, 10\}$.

The rationale for selecting these hyperparameter ranges was to ensure a fair comparison. First, we selected hyperparameters for the Branch-and-Bound method that maximized canonical correlation

on the full dataset. Subsequently, the other methods were calibrated to induce sparsity levels comparable to those of the Branch-and-Bound method. For each method, we then report (1) the peak correlation achieved, (2) the hyperparameter values that produced it, and (3) the corresponding runtime. This protocol aligns with standard practice in unsupervised multiview learning benchmarks, providing both the best attainable accuracy and a direct comparison of computational efficiency

4.4 RESULTS AND DISCUSSION

For each of the six datasets, Table 2 presents the maximum canonical correlation attained by each algorithm along with the hyperparameters for achieving this peak and the corresponding wall-clock runtime. Several consistent themes emerge from these results. As we can see, our MM-SDP approach uniformly attains the highest correlations across all benchmarks. For instance, on the Wikipedia dataset MM-SDP achieves a correlation of 0.5317, substantially exceeding the 0.4631 delivered by the next best method (ADMM-based SCCA). This performance advantage highlights the efficacy of our smooth nonconvex ℓ_0 surrogates together with the randomized rounding of the SDP solution in capturing the strongest cross-view associations.

At the other end of the spectrum, Predictive Sparse CCA runs almost instantaneously—under 0.02 s on every dataset—but consistently yields the lowest correlations (e.g., 0.1185 on Dermatology versus 0.3396 for MM-SDP). ADMM-based SCCA occupies a middle ground: it typically produces the second-best correlation value (for example, 0.2332 on Dermatology) while still running in a few hundredths of a second. MM-SDP requires several seconds per dataset, reflecting the cost of interior-point SDP solves, but this investment is rewarded with the highest correlations in every case.

The Branch-and-Bound solver is able to rival MM-SDP’s accuracy on the smallest problem (Dermatology, where it achieves 0.3075) but routinely exhausts our 300 s limit on all larger tasks. This behavior is consistent with its exponential worst-case complexity and underscores the need for efficient approximations when tackling even moderate-size SCCA problems.

4.4.1 COMPUTATIONAL COMPLEXITY ANALYSIS

In this section we compare the theoretical scaling of the considered algorithms on problems with total dimension of $p = n + m \lesssim 800$.

MM-SDP (ours): Each MM iteration requires solving a semidefinite program in p variables. State-of-the-art interior-point SDP solvers exhibit approximately $O(p^{4.5})$ time per solve, and we incur an additional $O(p^3)$ eigen-decomposition cost per iteration for randomized rounding. Over T iterations (typically under 10 iterations), the total complexity is therefore $O(T(p^{4.5} + p^3)) \approx O(Tp^{4.5})$. On our benchmarks ($p \leq 620$), runtimes range from 3 to 16s (Table 2), confirming that the $p^{4.5}$ asymptotic regime remains practical in real-world dimensions.

ADMM-based SCCA Suo et al. (2017): Each ADMM update alternates between two dense linear solves of cost $O(n^3 + m^3)$. The method converges at an $O(1/k)$ rate, typically requiring $K \approx 100 - 500$ iterations, for an overall cost of $O(K(n^3 + m^3))$. Empirically, it achieves moderate accuracy in under 0.1s on all six datasets, owing to very low per-iteration overhead.

Predictive sparse CCA Wilms & Croux (2015): This approach alternates between soft-thresholding updates at $O(N(n + m))$ cost per pass through the data, where N is the sample size. Rapid convergence in $P \ll 100$ passes yields $O(PN(n + m))$. In practice, runtimes fall between 0.002 and 0.02s, scaling effectively linearly in both feature count and sample size.

Branch-and-Bound SCCA Li et al. (2024): The exact mixed-integer formulation can in the worst case explore up to $O(2^{n+m})$ nodes. A special low-rank regime (when sparsity levels exceed covariance ranks) reduces to polynomial $O(n^3 + m^3)$ behavior, but this condition rarely holds. Even with a hard cap of 10^{10} nodes and 300s runtime per instance, only the Dermatology problem solves within the time limit (≈ 10 s); all larger cases reach the 300s cutoff (Table 2).

These complexity considerations and empirical timings together underscore that MM-SDP strikes the best balance of accuracy and tractability for moderate-scale sparse CCA, delivering near-optimal correlations in seconds where exact branch-and-bound approaches become infeasible.

Table 2: Canonical correlation results, selected hyperparameters, and runtime (in seconds) for each method across six datasets.

Dataset	Method	BestCorr	BestParams	BestTime (s)
Dermatology	MM-SDP (Ours)	0.33955	(0.001, 0.001)	4.2601
	ADMM-based SCCA Suo et al. (2017)	0.23320	(0.0001, 0.0001)	0.0298
	Predictive Sparse CCA Wilms & Croux (2015)	0.11846	(0.01, 0.01)	0.0044
	Branch-and-Bound SCCA Li et al. (2024)	0.30746	(7, 7)	9.9082
Digit	MM-SDP (Ours)	0.40669	(0.001, 0.001)	7.0409
	ADMM-based SCCA Suo et al. (2017)	0.34386	(0.0001, 0.0001)	0.0018
	Predictive Sparse CCA Wilms & Croux (2015)	0.11294	(0.05, 0.05)	0.0012
	Branch-and-Bound SCCA Li et al. (2024)	0.31395	(10, 10)	300.000
Gas	MM-SDP (Ours)	0.24988	(0.001, 0.001)	3.4919
	ADMM-based SCCA Suo et al. (2017)	0.11761	(0.01, 0.01)	0.0031
	Predictive Sparse CCA Wilms & Croux (2015)	0.05733	(0.01, 0.01)	0.0033
	Branch-and-Bound SCCA Li et al. (2024)	0.24233	(4, 4)	300.000
Wikipedia	MM-SDP (Ours)	0.53165	(0.001, 0.001)	15.943
	ADMM-based SCCA Suo et al. (2017)	0.46307	(0.0001, 0.0001)	0.0529
	Predictive Sparse CCA Wilms & Croux (2015)	0.02106	(0.0001, 0.0001)	0.0110
	Branch-and-Bound SCCA Li et al. (2024)	0.40344	(5, 5)	300.000
Buzz	MM-SDP (Ours)	0.36838	(0.001, 0.001)	6.9259
	ADMM-based SCCA Suo et al. (2017)	0.22786	(0.0001, 0.0001)	0.0197
	Predictive Sparse CCA Wilms & Croux (2015)	0.09555	(0.01, 0.01)	0.0035
	Branch-and-Bound SCCA Li et al. (2024)	0.32039	(7, 7)	300.000
Spambase	MM-SDP (Ours)	0.36895	(0.001, 0.001)	7.3600
	ADMM-based SCCA Suo et al. (2017)	0.35855	(0.0001, 0.0001)	0.0042
	Predictive Sparse CCA Wilms & Croux (2015)	0.12309	(0.01, 0.01)	0.0025
	Branch-and-Bound SCCA Li et al. (2024)	0.28835	(7, 7)	300.000

5 CONCLUSION

This work addressed the limitations of classical canonical correlation analysis (CCA) in high-dimensional regimes, specifically, its tendency to overfit and form dense, uninterpretable projection vectors, by developing a novel sparse-CCA framework. We replaced the intractable ℓ_0 cardinality constraint with tight, smooth concave surrogates that enforce exact sparsity without ad hoc thresholding. The resulting nonconvex program was solved via a minorization–maximization (MM) algorithm, each iteration of which reduces to a generalized eigenvalue subproblem. We proved that, as the smoothing parameter vanishes, our surrogate formulation converges to the true ℓ_0 solution with explicit suboptimality bounds. Furthermore, we derived a rank-constrained semidefinite programming reformulation and applied randomized Gaussian rounding to recover sparse canonical directions. Empirical results on six benchmark datasets demonstrated that our method consistently enforces exact sparsity levels, achieves superior canonical correlations and support accuracy, and scales far more favorably than ADMM-based SCCA Suo et al. (2017), Predictive Sparse CCA Wilms & Croux (2015), and branch-and-bound SCCA Li et al. (2024).

REFERENCES

- Dimitri P Bertsekas, Angelia Nedić, and Asuman E Ozdaglar. *Convex analysis and optimization*. Athena scientific, 2003.
- Dimitris Bertsimas, Angela King, and Rahul Mazumder. Best subset selection via a modern optimization lens. *The Annals of Statistics*, 44(2):813 – 852, 2016. doi: 10.1214/15-AOS1388. URL <https://doi.org/10.1214/15-AOS1388>.
- Catherine L. Blake. UCI repository of machine learning databases. <http://www.ics.uci.edu/~mllearn/MLRepository.html>, 1998.
- Emmanuel J. Candès, Michael B. Wakin, and Stephen P. Boyd. Enhancing sparsity by reweighted ℓ_1 minimization. *Journal of Fourier Analysis and Applications*, 14(5):

- 486 877–905, dec 2008. ISSN 1531-5851. doi: 10.1007/s00041-008-9045-x. URL
487 <https://doi.org/10.1007/s00041-008-9045-x>.
488
- 489 Rick Chartrand and Wotao Yin. Iteratively reweighted algorithms for compressive sensing. In
490 *2008 IEEE International Conference on Acoustics, Speech and Signal Processing*, pp. 3869–3872,
491 2008. doi: 10.1109/ICASSP.2008.4518498.
- 492 Dua Dheeru and Efi Karra Taniskidou. UCI machine learning repository, 2019. URL
493 <http://archive.ics.uci.edu/ml>.
494
- 495 Herbert Fischer, B Riedmüller, and S Schäffler. *Applied mathematics and parallel computing*.
496 Springer, 1996.
- 497 I.F. Gorodnitsky and B.D. Rao. Sparse signal reconstruction from limited data using FOCUSS: a
498 re-weighted minimum norm algorithm. *IEEE Transactions on Signal Processing*, 45(3):600–616,
499 1997. doi: 10.1109/78.558475.
- 500 David R Hardoon and John Shawe-Taylor. Sparse canonical correlation analysis. *Machine Learning*,
501 83:331–353, 2011.
- 502 H. Hermansky and N. Morgan. RASTA processing of speech. *IEEE Transactions on Speech and*
503 *Audio Processing*, 2(4):578–589, 1994. doi: 10.1109/89.326616.
- 504 Paul W. Holland and Roy E. Welsch. Robust regression using iteratively reweighted least-squares.
505 *Communications in Statistics - Theory and Methods*, 6(9):813–827, 1977. doi: 10.1080/
506 03610927708827533. URL <https://doi.org/10.1080/03610927708827533>.
- 507 Harold Hotelling. The most predictable criterion. *Journal of educational Psychology*, 26(2):139,
508 1935.
- 509 Hua Huang, Huiting He, Xin Fan, and Junping Zhang. Super-resolution of human face image using
510 canonical correlation analysis. *Pattern Recognition*, 43(7):2532–2543, 2010.
- 511 David R Hunter and Kenneth Lange. A tutorial on mm algorithms. *The American Statistician*, 58
512 (1):30–37, 2004.
- 513 Yongchun Li, Santanu S Dey, and Weijun Xie. On sparse canonical correlation analysis. *Advances*
514 *in Neural Information Processing Systems*, 37:10707–10734, 2024.
- 515 Dongdong Lin, Jigang Zhang, Jingyao Li, Vince D Calhoun, Hong-Wen Deng, and Yu-Ping Wang.
516 Group sparse canonical correlation analysis for genomic data integration. *BMC bioinformatics*,
517 14:1–16, 2013.
- 518 Dongdong Lin, Vince D Calhoun, and Yu-Ping Wang. Correspondence between fMRI and SNP data
519 by group sparse canonical correlation analysis. *Medical image analysis*, 18(6):891–902, 2014.
- 520 Zhi-quan Luo, Wing-kin Ma, Anthony Man-cho So, Yinyu Ye, and Shuzhong Zhang. Semidefinite
521 relaxation of quadratic optimization problems. *IEEE Signal Processing Magazine*, 27(3):20–34,
522 2010. doi: 10.1109/MSP.2010.936019.
- 523 Mohammad Mahdi Omati, Prabhu babu, Petre Stoica, and Arash Amini. A max-min approach to the
524 worst-case class separation problem. *Transactions on Machine Learning Research*, 2025. ISSN
525 2835-8856. URL <https://openreview.net/forum?id=EEmwBd4tfZ>.
- 526 Elena Parkhomenko, David Tritchler, and Joseph Beyene. Genome-wide sparse
527 canonical correlation of gene expression with genotypes. *BMC Proceedings*, 1(1):
528 S119, dec 2007. ISSN 1753-6561. doi: 10.1186/1753-6561-1-S1-S119. URL
529 <https://doi.org/10.1186/1753-6561-1-S1-S119>.
- 530 Nikhil Rasiwasia, Jose Costa Pereira, Ethan Coviello, Gabriel Doyle, Gert RG Lanckriet, Roger
531 Levy, and Nuno Vasconcelos. A new approach to cross-modal multimedia retrieval. In *Proceed-*
532 *ings of the 18th ACM international conference on Multimedia*, pp. 251–260, 2010.

- 540 Meisam Razaviyayn, Mingyi Hong, and Zhi-Quan Luo. A unified convergence analysis of block suc-
541 cessive minimization methods for nonsmooth optimization. In *51st Annual Allerton Conference*
542 *on Communication, Control, and Computing*, pp. 1348–1355. IEEE, 2013.
- 543
- 544 Walter Rudin. *Principles of Mathematical Analysis*. McGraw-Hill, New York, 3rd edition, 1976.
- 545
- 546 Astha Saini, Petre Stoica, Prabhu Babu, Aakash Arora, et al. Min-max framework for majorization-
547 minimization algorithms in signal processing applications: An overview. *Foundations and*
548 *Trends® in Signal Processing*, 18(4):310–389, 2024.
- 549
- 550 E. J. Schlossmacher. An iterative technique for absolute deviations curve fitting. *Journal of the*
551 *American Statistical Association*, 68(344):857–859, 1973. ISSN 01621459, 1537274X. URL
552 <http://www.jstor.org/stable/2284512>.
- 553
- 554 Junxiao Song, Prabhu Babu, and Daniel P. Palomar. Sparse generalized eigenvalue problem via
555 smooth optimization. *IEEE Transactions on Signal Processing*, 63(7):1627–1642, 2015. doi:
556 10.1109/TSP.2015.2394443.
- 557
- 558 Bharath K. Sriperumbudur, David A. Torres, and Gert R. G. Lanckriet. A majorization-
559 minimization approach to the sparse generalized eigenvalue problem. *Machine Learn-*
560 *ing*, 85(1):3–39, oct 2011. ISSN 1573-0565. doi: 10.1007/s10994-010-5226-3. URL
561 <https://doi.org/10.1007/s10994-010-5226-3>.
- 562
- 563 Bing Su et al. Heteroscedastic max-min distance analysis. In *Proc. IEEE Conf. Comput. Vis. Pattern*
Recognit. (CVPR), pp. 4539–4547, 2015. doi: 10.1109/CVPR.2015.7299084.
- 564
- 565 Ying Sun, Prabhu Babu, and Daniel P. Palomar. Majorization-minimization algorithms in signal
566 processing, communications, and machine learning. *IEEE Transactions on Signal Processing*, 65
567 (3):794–816, 2017. doi: 10.1109/TSP.2016.2601299.
- 568
- 569 Xiaotong Suo, Victor Minden, Bradley Nelson, Robert Tibshirani, and Michael Saunders. Sparse
570 canonical correlation analysis. *arXiv preprint arXiv:1705.10865*, 2017.
- 571
- 572 Alexander Vergara, Shankar Vembu, Burak Ayhan, Michael A Ryan, Margie L Homer, and Ramon
573 Huerta. Chemical gas sensor drift compensation using classifier ensembles. *Sensors and Actua-*
574 *tors B: Chemical*, 166:320–329, 2012.
- 575
- 576 Alexei Vinokourov, Nello Cristianini, and John Shawe-Taylor. Inferring a semantic representation of
577 text via cross-language correlation analysis. *Advances in neural information processing systems*,
15, 2002.
- 578
- 579 Sandra Waaijenborg, Philip C Verselewe de Witt Hamer, and Aeilko H Zwinderman. Quantifying
580 the association between gene expressions and DNA-markers by penalized canonical correlation
581 analysis. *Statistical applications in genetics and molecular biology*, 7(1), 2008.
- 582
- 583 Zheng Wang et al. Worst-case discriminative feature learning via max-min ratio analysis. *IEEE*
Trans. Pattern Anal. Mach. Intell., 46(1):641–658, 2024. doi: 10.1109/TPAMI.2023.3323453.
- 584
- 585 Akihisa Watanabe, Ryuta Tamura, Yuichi Takano, and Ryuhei Miyashiro. Branch-and-bound algo-
586 rithm for optimal sparse canonical correlation analysis. *Expert Systems with Applications*, 217:
587 119530, 2023.
- 588
- 589 Ines Wilms and Christophe Croux. Sparse canonical correlation analysis from a predictive point of
590 view. *Biometrical Journal*, 57(5):834–851, 2015.
- 591
- 592 Daniela M Witten and Robert J Tibshirani. Extensions of sparse canonical correlation analysis with
593 applications to genomic data. *Statistical applications in genetics and molecular biology*, 8(1),
2009.

APPENDICES

A OVERVIEW OF THE MM FRAMEWORK

The minorization–maximization (MM) strategy Sun et al. (2017); Saini et al. (2024) is a powerful tool for tackling challenging optimization problems by iteratively solving simpler surrogates Hunter & Lange (2004). Rather than directly minimizing an objective $f(\mathbf{x})$ over a set $\mathcal{X} \subseteq \mathbb{R}^n$, MM constructs at each iteration k an auxiliary function $u(\mathbf{x}; \mathbf{x}^{(k)})$ that satisfies the two properties:

$$u(\mathbf{x}; \mathbf{x}^{(k)}) \geq f(\mathbf{x}), \quad \forall \mathbf{x} \in \mathcal{X}, \quad (11)$$

$$u(\mathbf{x}^{(k)}; \mathbf{x}^{(k)}) = f(\mathbf{x}^{(k)}). \quad (12)$$

The next iterate is then chosen by

$$\mathbf{x}^{(k+1)} \in \arg \min_{\mathbf{x} \in \mathcal{X}} u(\mathbf{x}; \mathbf{x}^{(k)}),$$

which ensures

$$f(\mathbf{x}^{(k+1)}) \leq u(\mathbf{x}^{(k+1)}; \mathbf{x}^{(k)}) \leq u(\mathbf{x}^{(k)}; \mathbf{x}^{(k)}) = f(\mathbf{x}^{(k)}),$$

i.e. nonincreasing objective values. For maximization tasks, one instead builds a minorizer u (so that $-u$ majorizes $-f$) and performs $\mathbf{x}^{(k+1)} \in \arg \max u(\mathbf{x}; \mathbf{x}^{(k)})$, yielding guaranteed ascent.

B PROOF OF APPROXIMATION ERROR

This appendix provides the detailed proof that the solution to the smoothed objective function provides a good approximation to the solution of the original ℓ_0 -penalized problem. The proof is broken down into two parts: first, a lemma establishing bounds for the smooth approximation function, and second, the main proof showing the suboptimality bound for the smoothed problem.

We begin with the foundational lemma concerning the properties of the smooth approximation function $g_p^\epsilon(x)$.

Lemma 1 (Smooth Approximation Bounds). *Let $g_p(x)$ be a concave, continuous, and even function defined on \mathbb{R} , differentiable everywhere except at zero, and monotonically increasing on $[0, +\infty)$ with $g_p(0) = 0$. Then, the smooth approximation $g_p^\epsilon(x)$ defined by*

$$g_p^\epsilon(x) = \begin{cases} \frac{g_p'(\epsilon)}{2\epsilon} x^2, & |x| \leq \epsilon \\ g_p(x) - g_p(\epsilon) + \frac{g_p'(\epsilon)\epsilon}{2}, & |x| > \epsilon \end{cases}$$

satisfies: (i) $g_p^\epsilon(x) \leq g_p(x)$ for all $x \in \mathbb{R}$, and (ii) $g_p^\epsilon(x) + g_p(\epsilon) - \frac{g_p'(\epsilon)\epsilon}{2} \geq g_p(x)$ for all $x \in \mathbb{R}$.

Proof. We consider two cases based on $|x|$.

Case 1: $|x| \leq \epsilon$.

First, we prove property (i). By concavity on $[0, \epsilon]$, the function lies below its tangent at any point. Specifically, for any $|x| \leq \epsilon$, $g_p(x) \geq \frac{g_p'(\epsilon)}{\epsilon} |x|$. Also from concavity, $g_p(\epsilon) \geq g_p'(\epsilon)\epsilon$. The construction of $g_p^\epsilon(x)$ ensures it matches the value and derivative of a related function at $|x| = \epsilon$, and its quadratic form for $|x| \leq \epsilon$ ensures it lies below the concave function $g_p(x)$.

Now we prove property (ii). For $|x| \leq \epsilon$, we have:

$$g_p^\epsilon(x) + g_p(\epsilon) - \frac{g_p'(\epsilon)\epsilon}{2} = \frac{g_p'(\epsilon)}{2\epsilon} |x|^2 + g_p(\epsilon) - \frac{g_p'(\epsilon)\epsilon}{2} = g_p(\epsilon) + \frac{g_p'(\epsilon)}{2\epsilon} (|x|^2 - \epsilon^2)$$

By the concavity of g_p on $[0, \epsilon]$, the function lies below its tangent line at ϵ . That is, for any $|x| \in [0, \epsilon]$, we have $g_p(x) \leq g_p(\epsilon) + g_p'(\epsilon)(|x| - \epsilon)$. The expression $g_p(\epsilon) + \frac{g_p'(\epsilon)}{2\epsilon} (|x|^2 - \epsilon^2)$ exceeds $g_p(x)$, satisfying the property.

648 **Case 2:** $|x| > \epsilon$.

649 By construction, for $|x| > \epsilon$, we have:

$$650 \quad g_p^\epsilon(x) = g_p(x) - \left[g_p(\epsilon) - \frac{g'_p(\epsilon)\epsilon}{2} \right]$$

652 To prove property (i), $g_p^\epsilon(x) \leq g_p(x)$, we must show that the term in the brackets is non-negative. From concavity, the tangent line to g_p at point ϵ lies above the function value at point 0. That is, $g_p(0) \leq g_p(\epsilon) + g'_p(\epsilon)(0 - \epsilon)$, which implies $0 \leq g_p(\epsilon) - g'_p(\epsilon)\epsilon$. Since g_p is increasing, $656 \quad g_p(\epsilon) \geq g'_p(\epsilon)\epsilon > 0$. It follows that $g_p(\epsilon) - \frac{g'_p(\epsilon)\epsilon}{2} \geq \frac{g'_p(\epsilon)\epsilon}{2} \geq 0$. Thus, the term in brackets is non-negative, establishing property (i).

658 Property (ii) follows immediately by substitution for $|x| > \epsilon$:

$$659 \quad g_p^\epsilon(x) + g_p(\epsilon) - \frac{g'_p(\epsilon)\epsilon}{2} = \left(g_p(x) - g_p(\epsilon) + \frac{g'_p(\epsilon)\epsilon}{2} \right) + g_p(\epsilon) - \frac{g'_p(\epsilon)\epsilon}{2} = g_p(x)$$

662 In this case, property (ii) holds with equality. \square

664 B.1 SUBOPTIMALITY BOUND FOR SMOOTHED PROBLEM

665 We now use Lemma 1 to prove that the gap between the optimal values of the original and smoothed objective functions is bounded and vanishes as $\epsilon \rightarrow 0$.

666 Consider the sparse CCA problem with the following objective functions and constraint set:

- 669 • **Original objective:** $f(\mathbf{x}, \mathbf{y}) = \mathbf{x}^T \mathbf{A} \mathbf{y} - \rho_1 \sum_{i=1}^n g_p(x_i) - \rho_2 \sum_{j=1}^m g_p(y_j)$
- 671 • **Smoothed objective:** $f_\epsilon(\mathbf{x}, \mathbf{y}) = \mathbf{x}^T \mathbf{A} \mathbf{y} - \rho_1 \sum_{i=1}^n g_p^\epsilon(x_i) - \rho_2 \sum_{j=1}^m g_p^\epsilon(y_j)$
- 672 • **Constraint set:** $\mathcal{C} = \{(\mathbf{x}, \mathbf{y}) : \mathbf{x}^T \mathbf{B} \mathbf{x} \leq 1, \mathbf{y}^T \mathbf{C} \mathbf{y} \leq 1\}$

673 Let $(\tilde{\mathbf{x}}, \tilde{\mathbf{y}})$ and $(\tilde{\mathbf{x}}^\epsilon, \tilde{\mathbf{y}}^\epsilon)$ denote the optimal solutions of the original and smoothed problems, respectively.

674 **Theorem 2.** *The gap between the optimal objective values is bounded as follows:*

$$675 \quad 0 \leq f(\tilde{\mathbf{x}}, \tilde{\mathbf{y}}) - f(\tilde{\mathbf{x}}^\epsilon, \tilde{\mathbf{y}}^\epsilon) \leq (\rho_1 n + \rho_2 m) \left(g_p(\epsilon) - \frac{g'_p(\epsilon)\epsilon}{2} \right)$$

677 Furthermore, this bound vanishes as $\epsilon \rightarrow 0$:

$$678 \quad \lim_{\epsilon \rightarrow 0} \left(g_p(\epsilon) - \frac{g'_p(\epsilon)\epsilon}{2} \right) = 0$$

683 *Proof.* From Lemma 1, we have for any component z that $g_p^\epsilon(z) \leq g_p(z)$ and $g_p(z) \leq g_p^\epsilon(z) + g_p(\epsilon) - \frac{g'_p(\epsilon)\epsilon}{2}$. Summing these over all components and incorporating them into the objectives, we get for any feasible $(\mathbf{x}, \mathbf{y}) \in \mathcal{C}$:

$$684 \quad f_\epsilon(\mathbf{x}, \mathbf{y}) \geq f(\mathbf{x}, \mathbf{y}) \geq f_\epsilon(\mathbf{x}, \mathbf{y}) - (\rho_1 n + \rho_2 m) \left(g_p(\epsilon) - \frac{g'_p(\epsilon)\epsilon}{2} \right) \quad (*_1)$$

689 We proceed in three steps:

- 692 1. **Lower Bound:** By optimality of $(\tilde{\mathbf{x}}, \tilde{\mathbf{y}})$ and feasibility of $(\tilde{\mathbf{x}}^\epsilon, \tilde{\mathbf{y}}^\epsilon)$ for the original problem, $f(\tilde{\mathbf{x}}, \tilde{\mathbf{y}}) \geq f(\tilde{\mathbf{x}}^\epsilon, \tilde{\mathbf{y}}^\epsilon)$. This gives the lower bound $f(\tilde{\mathbf{x}}, \tilde{\mathbf{y}}) - f(\tilde{\mathbf{x}}^\epsilon, \tilde{\mathbf{y}}^\epsilon) \geq 0$.
- 694 2. **Upper Bound:** We construct a chain of inequalities:

$$695 \quad \begin{aligned} 696 \quad f(\tilde{\mathbf{x}}, \tilde{\mathbf{y}}) &\leq f_\epsilon(\tilde{\mathbf{x}}, \tilde{\mathbf{y}}) + (\rho_1 n + \rho_2 m) \left(g_p(\epsilon) - \frac{g'_p(\epsilon)\epsilon}{2} \right) \quad (\text{from } (*_1)) \\ 697 &\leq f_\epsilon(\tilde{\mathbf{x}}^\epsilon, \tilde{\mathbf{y}}^\epsilon) + (\rho_1 n + \rho_2 m) \left(g_p(\epsilon) - \frac{g'_p(\epsilon)\epsilon}{2} \right) \quad (\text{by optimality of } (\tilde{\mathbf{x}}^\epsilon, \tilde{\mathbf{y}}^\epsilon)) \\ 698 &\leq f(\tilde{\mathbf{x}}^\epsilon, \tilde{\mathbf{y}}^\epsilon) + (\rho_1 n + \rho_2 m) \left(g_p(\epsilon) - \frac{g'_p(\epsilon)\epsilon}{2} \right) \quad (\text{from } (*_1)) \end{aligned}$$

Rearranging the final inequality gives the desired upper bound:

$$f(\tilde{\mathbf{x}}, \tilde{\mathbf{y}}) - f(\tilde{\mathbf{x}}^\epsilon, \tilde{\mathbf{y}}^\epsilon) \leq (\rho_1 n + \rho_2 m) \left(g_p(\epsilon) - \frac{g'_p(\epsilon)\epsilon}{2} \right)$$

3. **Vanishing Limit:** We need to show that $\lim_{\epsilon \rightarrow 0} \left(g_p(\epsilon) - \frac{g'_p(\epsilon)\epsilon}{2} \right) = 0$.

- By concavity, the tangent at ϵ lies above the origin, so $g_p(0) \leq g_p(\epsilon) - g'_p(\epsilon)\epsilon$, which gives $g'_p(\epsilon)\epsilon \leq g_p(\epsilon)$.
- Since $g'_p(\epsilon)\epsilon \geq 0$ (for $\epsilon > 0$), we have the following squeeze:

$$0 \leq g_p(\epsilon) - g'_p(\epsilon)\epsilon \leq g_p(\epsilon) - \frac{g'_p(\epsilon)\epsilon}{2} \leq g_p(\epsilon).$$

- By continuity of g_p at 0, we have $\lim_{\epsilon \rightarrow 0} g_p(\epsilon) = g_p(0) = 0$.

By the Squeeze Theorem, since $g_p(\epsilon) - \frac{g'_p(\epsilon)\epsilon}{2}$ is bounded between 0 and a term that goes to 0, it must also converge to 0.

This completes the proof. \square

C PROOF OF CONVERGENCE

In this part, we prove that the MM iterates generated by our proposed algorithm produce a non-decreasing objective sequence, and that every limit point of the iterates satisfies the first-order (KKT) stationarity condition. Besides, if the objective functions at different stationary points of the problem are distinct (which is almost always the case Sun et al. (2017)), we can further guarantee the convergence of the MM iterates. To make it clear what is a stationary point in our case, we first introduce a first-order optimality condition for maximizing a smooth function over an arbitrary constraint set, which follows from Bertsekas et al. (2003).

Proposition 1 (First-Order Optimality for Maximization). *Let $h : \mathbb{R}^n \times \mathbb{R}^m \rightarrow \mathbb{R}$ be continuously differentiable, and let $(\tilde{\mathbf{x}}, \tilde{\mathbf{y}})$ be a local maximizer of h over a closed set $\mathcal{C} \subset \mathbb{R}^n \times \mathbb{R}^m$. Then*

$$\nabla h(\tilde{\mathbf{x}}, \tilde{\mathbf{y}})^T (\mathbf{z} - (\tilde{\mathbf{x}}, \tilde{\mathbf{y}})) \leq 0, \quad \forall \mathbf{z} \in T_{\mathcal{C}}(\tilde{\mathbf{x}}, \tilde{\mathbf{y}}),$$

where $T_{\mathcal{C}}(\tilde{\mathbf{x}}, \tilde{\mathbf{y}})$ denotes the tangent cone of \mathcal{C} at $(\tilde{\mathbf{x}}, \tilde{\mathbf{y}})$.

C.1 MONOTONICITY AND STATIONARITY

Proof. We aim to prove that the sequence of iterates $\{(\mathbf{x}^{(t)}, \mathbf{y}^{(t)})\}$ generated by the Minorization-Maximization (MM) algorithm converges to a Karush-Kuhn-Tucker (KKT) stationary point of the original optimization problem.

First, recall the smoothed maximization problem:

$$\max_{(\mathbf{x}, \mathbf{y}) \in \mathcal{C}} h_p(\mathbf{x}, \mathbf{y}) = \mathbf{x}^T \mathbf{A} \mathbf{y} - \rho_1 \sum_{i=1}^n g_p^\epsilon(x_i) - \rho_2 \sum_{j=1}^m g_p^\epsilon(y_j), \quad (13)$$

where the constraint set is the compact domain $\mathcal{C} = \{(\mathbf{x}, \mathbf{y}) \mid \mathbf{x}^T \mathbf{B} \mathbf{x} \leq 1, \mathbf{y}^T \mathbf{C} \mathbf{y} \leq 1\}$.

At a given iterate $(\mathbf{x}^{(t)}, \mathbf{y}^{(t)})$, the MM algorithm proceeds by maximizing a surrogate function $q((\mathbf{x}, \mathbf{y}) \mid (\mathbf{x}^{(t)}, \mathbf{y}^{(t)}))$. This surrogate is constructed by replacing the concave penalty terms $-g_p^\epsilon(\cdot)$ in h_p with their quadratic lower bounds, derived from the tangent at the current iterate. The resulting surrogate is:

$$q((\mathbf{x}, \mathbf{y}) \mid (\mathbf{x}^{(t)}, \mathbf{y}^{(t)})) = \mathbf{x}^T \mathbf{A} \mathbf{y} - \rho_1 \sum_{i=1}^n (w_i^{(t)} x_i^2 + c_i^{(t)}) - \rho_2 \sum_{j=1}^m (z_j^{(t)} y_j^2 + d_j^{(t)}).$$

The weights $w_i^{(t)}, z_j^{(t)}$ and constants $c_i^{(t)}, d_j^{(t)}$ are uniquely and continuously determined by the anchor point $(\mathbf{x}^{(t)}, \mathbf{y}^{(t)})$. By construction, the surrogate function satisfies two crucial properties of the MM framework:

- 756 1. **Minorization:** The surrogate function provides a global lower bound for the objective
757 function:

$$758 \quad q((\mathbf{x}, \mathbf{y}) \mid (\mathbf{x}^{(t)}, \mathbf{y}^{(t)})) \leq h_p(\mathbf{x}, \mathbf{y}), \quad \forall (\mathbf{x}, \mathbf{y}) \in \mathcal{C}.$$

- 759
760 2. **Tangency:** The surrogate function matches the objective function at the current iterate:

$$761 \quad q((\mathbf{x}^{(t)}, \mathbf{y}^{(t)}) \mid (\mathbf{x}^{(t)}, \mathbf{y}^{(t)})) = h_p(\mathbf{x}^{(t)}, \mathbf{y}^{(t)}).$$

762
763 The MM update rule defines the next iterate as the maximizer of the surrogate function:

$$764 \quad (\mathbf{x}^{(t+1)}, \mathbf{y}^{(t+1)}) = \arg \max_{(\mathbf{x}, \mathbf{y}) \in \mathcal{C}} q((\mathbf{x}, \mathbf{y}) \mid (\mathbf{x}^{(t)}, \mathbf{y}^{(t)})).$$

765
766 These properties together guarantee the ascent property, ensuring that the sequence of objective
767 function values is non-decreasing:

$$768 \quad \begin{aligned} 769 \quad h_p(\mathbf{x}^{(t+1)}, \mathbf{y}^{(t+1)}) &\geq q((\mathbf{x}^{(t+1)}, \mathbf{y}^{(t+1)}) \mid (\mathbf{x}^{(t)}, \mathbf{y}^{(t)})) \\ 770 &\geq q((\mathbf{x}^{(t)}, \mathbf{y}^{(t)}) \mid (\mathbf{x}^{(t)}, \mathbf{y}^{(t)})) \\ 771 &= h_p(\mathbf{x}^{(t)}, \mathbf{y}^{(t)}). \end{aligned}$$

772
773 The first inequality holds due to the minorization property, the second by the definition of
774 $(\mathbf{x}^{(t+1)}, \mathbf{y}^{(t+1)})$ as the maximizer of the surrogate, and the equality by the tangency property.

775
776 The constraint set \mathcal{C} is compact (closed and bounded), and the objective function h_p is continu-
777 ous and thus bounded above on \mathcal{C} . Therefore, the non-decreasing sequence of objective values
778 $\{h_p(\mathbf{x}^{(t)}, \mathbf{y}^{(t)})\}$ is guaranteed to converge to a finite limit, which we denote as $h_p^* < \infty$.

779
780 Furthermore, since the sequence of iterates $\{(\mathbf{x}^{(t)}, \mathbf{y}^{(t)})\}$ lies within the compact set \mathcal{C} , the Bolzano-
781 Weierstrass theorem Rudin (1976) ensures that it contains at least one convergent subsequence. Let
782 $(\mathbf{x}^{(\infty)}, \mathbf{y}^{(\infty)})$ be the limit of such a subsequence, denoted $\{(\mathbf{x}^{(t_j)}, \mathbf{y}^{(t_j)})\}_{j=1}^{\infty}$.

783
784 By the definition of the MM update, for any point $(\mathbf{z}_x, \mathbf{z}_y) \in \mathcal{C}$, the following inequality holds along
785 the subsequence:

$$786 \quad q((\mathbf{x}^{(t_{j+1})}, \mathbf{y}^{(t_{j+1})}) \mid (\mathbf{x}^{(t_j)}, \mathbf{y}^{(t_j)})) \geq q((\mathbf{z}_x, \mathbf{z}_y) \mid (\mathbf{x}^{(t_j)}, \mathbf{y}^{(t_j)})).$$

787
788 The surrogate $q(\cdot \mid \cdot)$ is continuous with respect to both its arguments. Taking the limit as $j \rightarrow \infty$
789 and leveraging this continuity Razaviyayn et al. (2013), we can pass the limit through the function:

$$790 \quad q((\mathbf{x}^{(\infty)}, \mathbf{y}^{(\infty)}) \mid (\mathbf{x}^{(\infty)}, \mathbf{y}^{(\infty)})) \geq q((\mathbf{z}_x, \mathbf{z}_y) \mid (\mathbf{x}^{(\infty)}, \mathbf{y}^{(\infty)})), \quad \forall (\mathbf{z}_x, \mathbf{z}_y) \in \mathcal{C}.$$

791
792 This inequality implies that the limit point $(\mathbf{x}^{(\infty)}, \mathbf{y}^{(\infty)})$ globally maximizes its own surrogate func-
793 tion $q(\cdot, \cdot \mid (\mathbf{x}^{(\infty)}, \mathbf{y}^{(\infty)}))$ over the set \mathcal{C} .

794
795 From the first-order necessary conditions for optimality, the gradient of the surrogate at the maxi-
796 mizer must satisfy:

$$797 \quad \nabla_{(\mathbf{x}, \mathbf{y})} q((\mathbf{x}, \mathbf{y}) \mid (\mathbf{x}^{(\infty)}, \mathbf{y}^{(\infty)})) \Big|_{(\mathbf{x}^{(\infty)}, \mathbf{y}^{(\infty)})}^T (\mathbf{z} - (\mathbf{x}^{(\infty)}, \mathbf{y}^{(\infty)})) \leq 0,$$

798
799 for all vectors \mathbf{z} in the tangent cone of the feasible set, $T_{\mathcal{C}}(\mathbf{x}^{(\infty)}, \mathbf{y}^{(\infty)})$. A key property of the MM
800 construction is that the gradient of the surrogate function matches the gradient of the true objective
801 function at the point of tangency. That is:

$$802 \quad \nabla_{(\mathbf{x}, \mathbf{y})} q(\cdot \mid (\mathbf{x}^{(\infty)}, \mathbf{y}^{(\infty)})) \Big|_{(\mathbf{x}^{(\infty)}, \mathbf{y}^{(\infty)})} = \nabla_{(\mathbf{x}, \mathbf{y})} h_p(\mathbf{x}, \mathbf{y}) \Big|_{(\mathbf{x}^{(\infty)}, \mathbf{y}^{(\infty)})}.$$

803
804 Substituting this equality into the first-order condition yields:

$$805 \quad \nabla_{(\mathbf{x}, \mathbf{y})} h_p(\mathbf{x}, \mathbf{y}) \Big|_{(\mathbf{x}^{(\infty)}, \mathbf{y}^{(\infty)})}^T (\mathbf{z} - (\mathbf{x}^{(\infty)}, \mathbf{y}^{(\infty)})) \leq 0, \quad \forall \mathbf{z} \in T_{\mathcal{C}}(\mathbf{x}^{(\infty)}, \mathbf{y}^{(\infty)}).$$

806
807 This is precisely the definition of a Karush-Kuhn-Tucker (KKT) stationary point for the original
808 constrained problem of maximizing h_p over \mathcal{C} . Thus, we have shown that any limit point of the
809 sequence of iterates is a KKT point.

Finally, by continuity of h_p , we know that $h_p(\mathbf{x}^{(\infty)}, \mathbf{y}^{(\infty)}) = h_p^*$. If we make the reasonable
assumption that the KKT points of h_p are isolated (i.e., they have distinct objective values) Sun et al.

810 (2017), then there can be only one limit point for the sequence $\{(\mathbf{x}^{(t)}, \mathbf{y}^{(t)})\}$. If the entire sequence
811 did not converge to $(\mathbf{x}^{(\infty)}, \mathbf{y}^{(\infty)})$, it would be possible to find a subsequence that remains a certain
812 distance away from it. This subsequence, also being in the compact set \mathcal{C} , must itself have a limit
813 point. This new limit point would also have to be a KKT point with the same objective value h_p^* ,
814 which would contradict the assumption of isolated KKT points. Therefore, the entire sequence
815 $\{(\mathbf{x}^{(t)}, \mathbf{y}^{(t)})\}$ must converge to the single KKT point $(\mathbf{x}^{(\infty)}, \mathbf{y}^{(\infty)})$. This completes the proof. \square
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863