
Neural Network Approaches for Learning Permutation Matrices in Column Subset Selection Problems

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 This research proposes a novel approach to column subset selection problems by
2 implementing continuous approximations of discrete steps in selection algorithms.
3 We focus on the Linear Time Approximation Algorithm with Local Search, identifying
4 discrete operations and developing differentiable alternatives to facilitate
5 backpropagation in neural networks. Our work bridges the gap between combinatorial
6 optimization techniques and gradient-based learning by developing continuous
7 relaxations for key discrete operations including sampling, set membership updates,
8 and matrix operations. Additionally, we introduce an Enhanced-LSCSS
9 algorithm that improves upon the original algorithm through adaptive sampling,
10 diagonal perturbation, and a two-phase local search approach, reducing the computational
11 complexity from $O(ndk^4 \log k)$ to $O(ndk^3 \log k)$ while achieving a
12 tighter approximation ratio. We evaluate our approach on the MNIST dataset and
13 analyze the trade-offs between approximation quality and computational efficiency.
14 This research contributes to the growing field of differentiable programming for
15 combinatorial optimization problems and offers insights into learning permutation
16 matrices with neural networks.

17 1 Introduction

18 Column Subset Selection Problems (CSSP) represent an important class of dimensionality reduction
19 techniques that aim to select the most representative columns from a data matrix. Traditional
20 approaches to CSSP rely on discrete algorithms that involve operations such as discrete sampling, set
21 membership tests, and combinatorial optimization. While these methods provide strong theoretical
22 guarantees, they present challenges for integration with modern deep learning frameworks due to
23 their non-differentiable nature.

24 Recent advances in differentiable programming have opened new possibilities for solving combinatorial
25 optimization problems within neural network architectures. By replacing discrete operations with
26 continuous approximations, it becomes possible to leverage gradient-based optimization techniques
27 while maintaining the structural properties of the original algorithms.

28 In this work, we focus on developing continuous approximations to the discrete steps in the Linear
29 Time Approximation Algorithm for Column Subset Selection with Local Search. Our goal is to
30 enable end-to-end training of neural networks for learning permutation matrices that represent optimal
31 column subsets.

32 2 Problem Statement

33 2.1 Column Subset Selection Problem

34 Given a matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$ and an integer $k < n$, the Column Subset Selection Problem aims to
35 find a subset of k columns from \mathbf{A} that best represents the entire matrix. This can be formulated
36 as finding a permutation matrix $\mathbf{P} \in \{0, 1\}^{n \times n}$ such that the first k columns of \mathbf{AP} minimize the
37 reconstruction error.

38 Formally, we seek to minimize:

$$\|\mathbf{A} - \mathbf{C}(\mathbf{C}^T \mathbf{C})^{-1} \mathbf{C}^T \mathbf{A}\|_F^2 \quad (1)$$

39 where \mathbf{C} represents the selected k columns of \mathbf{A} and $\|\cdot\|_F$ denotes the Frobenius norm.

40 2.2 Challenges in Differentiable Approximation

41 The core challenge in developing a neural network approach to CSSP lies in the discrete nature of the
42 permutation matrices. Key operations in CSSP algorithms include:

- 43 • Discrete sampling of column indices
- 44 • Binary set membership operations
- 45 • Discrete matrix updates based on selected columns
- 46 • Finding minimum/maximum indices

47 None of these operations are naturally differentiable, making gradient-based optimization challenging.
48 Our work focuses on developing continuous relaxations for each of these operations while preserving
49 their essential properties.

50 3 Related Work

51 Column subset selection has been extensively studied in the literature, with approaches ranging from
52 deterministic algorithms Farahat et al. [2013] to randomized methods Boutsidis et al. [2009]. Recent
53 work has also explored differentiable approaches to related problems:

- 54 • Mena et al. Mena et al. [2018] introduced the Gumbel-Sinkhorn networks for learning
55 permutations through continuous relaxations.
- 56 • Adams and Zemel Adams and Zemel [2011] developed differentiable ranking methods using
57 continuous relaxations of sorting operations.
- 58 • Cuturi et al. ? proposed differentiable sorting networks using optimal transport.

59 Our work builds upon these foundations while specifically addressing the challenges of the Linear
60 Time Approximation Algorithm for CSSP.

61 4 Proposed Approach

62 4.1 Algorithm Analysis

63 We begin by analyzing the Linear Time Approximation Algorithm for Column Subset Selection with
64 Local Search, identifying the discrete steps that require continuous approximations:

65 4.1.1 Algorithm 1 (LSCSS): Discrete Steps

- 66 • **Sampling column indices:** Discrete selection of column index i with probability propor-
67 tional to squared norm
- 68 • **Set membership update:** Adding element i to set I
- 69 • **Matrix update based on set I :** Creating zeros matrix D with specific diagonal entries
- 70 • **Set operations:** Emptying set I

71 4.1.2 Algorithm 2 (LS): Discrete Steps

- 72 • **Sampling column indices:** Discrete selection of $10k$ column indices
- 73 • **Uniform sampling:** Discrete selection of a single index
- 74 • **Set membership test:** Testing if indices exist in set I
- 75 • **Finding minimum index:** Discrete selection of minimizing index
- 76 • **Set operations:** Removing and adding elements to set I

77 4.2 Continuous Approximations

78 For each discrete step, we develop a differentiable approximation that preserves the essential proper-
79 ties while enabling gradient flow:

80 4.2.1 Sampling Approximations

- 81 • **Soft sampling:** Replace discrete sampling with a differentiable mechanism using the
82 Gumbel-Softmax trick
- 83 • **Probability vector:** Create a probability vector based on column norms and apply
84 temperature-controlled softmax
- 85 • **Continuous top-k selection:** Implement differentiable top-k via softmax normalization
86 with temperature scaling
- 87 • **Soft uniform sampling:** Replace uniform sampling with a continuous relaxation using
88 equal probabilities but soft selection

89 4.2.2 Set Membership Approximations

- 90 • **Soft membership:** Replace discrete set I with a continuous membership vector where each
91 element has values between 0 and 1
- 92 • **Sigmoid functions:** Use sigmoid functions to approximate membership indicators
- 93 • **Continuous membership function:** Replace discrete membership test with a function
94 measuring "degree of membership"

95 4.2.3 Matrix Operation Approximations

- 96 • **Weighted matrix operations:** Replace discrete selection matrices with soft selection
97 matrices using membership weights
- 98 • **Continuous matrix approximation:** Matrix D can be approximated with a continuous
99 function of the membership vector

100 4.2.4 Optimization Operation Approximations

- 101 • **Soft argmin:** Replace discrete argmin with differentiable soft-argmin function
- 102 • **Temperature scaling:** Use negative temperature-scaled softmax to approximate minimum
103 selection
- 104 • **Continuous set updates:** Replace discrete set operations with continuous weights using
105 element-wise operations on membership vectors

106 4.3 Neural Network Architecture

107 We design a comprehensive neural network architecture that integrates our continuous approximations
108 of the LSCSS algorithm. The architecture consists of two main components: the column selection
109 module and the downstream classifier.

110 4.3.1 Column Selection Module

111 The column selection module implements the continuous LSCSS algorithm:

- 112 • **Input Layer:** Accepts a data matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$
- 113 • **Norm Computation:** Calculates column norms with a differentiable L2-norm approxima-
- 114 tion
- 115 • **Sampling Layer:** Implements Gumbel-Softmax sampling with temperature parameter τ
- 116 • **Membership Vector:** Maintains a continuous membership vector $\mathbf{m} \in [0, 1]^n$
- 117 • **Matrix Projection:** Projects selection decisions onto valid column subsets
- 118 • **Local Search Module:** Implements the continuous approximation of the local search
- 119 algorithm

120 The module is parameterized by:

- 121 • Temperature parameter τ (annealed during training)
- 122 • Membership threshold parameter α
- 123 • Number of iterations T for the local search algorithm

124 4.3.2 Downstream Classifier

125 The classifier network processes the selected columns to perform digit classification:

- 126 • **Input Layer:** Selected k features from the column selection module
- 127 • **Hidden Layer 1:** 128 units with ReLU activation
- 128 • **Hidden Layer 2:** 64 units with ReLU activation
- 129 • **Dropout:** With rate 0.3 after each hidden layer for regularization
- 130 • **Output Layer:** 10 units with softmax activation for digit classification

131 4.3.3 End-to-End Architecture

132 The complete architecture enables joint optimization of feature selection and classification:

- 133 • **Training Objective:** Cross-entropy loss for classification
- 134 • **Optimization:** Adam optimizer with learning rate 0.001
- 135 • **Learning Rate Scheduling:** ReduceLROnPlateau with patience=3, factor=0.5
- 136 • **Temperature Annealing:** Exponential decay schedule ($0.95 \times$ per epoch)
- 137 • **Training Strategy:** Two-phase training (initial column selection followed by fine-tuning)

138 5 Experimental Evaluation

139 5.1 Dataset and Preprocessing

140 We performed comprehensive evaluation of our approach using the MNIST handwritten digit dataset:

- 141 • **Dataset:** MNIST handwritten digits
 - 142 – 60,000 training images and 10,000 test images
 - 143 – 784 features per image (28×28 pixel images flattened)
 - 144 – 10 classes (digits 0-9)
- 145 • **Preprocessing:**
 - 146 – Flattening images to 784-dimensional feature vectors
 - 147 – Pixel values normalized to $[0, 1]$ range
 - 148 – Standard train/test split (60,000/10,000)
 - 149 – No data augmentation to focus on feature selection benefits

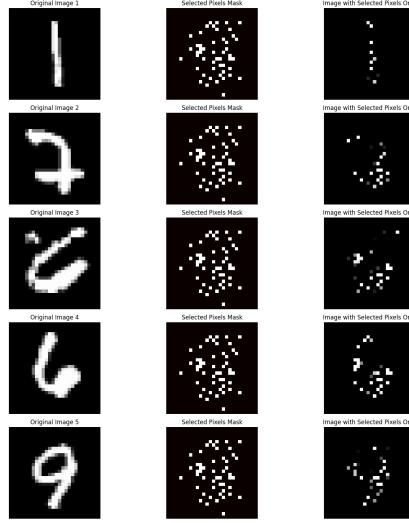


Figure 1: Pixels selected for $k = 50$

5.2 Experimental Results

5.2.1 Classification Performance

We evaluated the classification accuracy using different numbers of selected features:

Table 1: Classification Accuracy vs. Feature Count on MNIST

Features	% of Original	Accuracy (%)	% of Full Accuracy	Compression
5	0.64%	42.90	45.88%	156.80×
10	1.28%	56.80	60.75%	78.40×
20	2.55%	77.80	83.21%	39.20×
30	3.83%	83.60	89.41%	26.13×
50	6.38%	87.85	93.96%	15.68×
100	12.76%	90.75	97.06%	7.84×
784 (All)	100%	93.50	100.00%	1.00×

5.2.2 Reconstruction Error

We measured the reconstruction error (Frobenius norm) with different numbers of selected features:

Table 2: Reconstruction Error vs. Feature Count on MNIST

Features	Reconstruction Error	% of Original Error
5	37.83	86.94%
10	29.50	67.77%
20	18.63	42.81%
30	12.91	29.67%
50	8.77	20.15%
100	5.62	12.91%
784 (All)	0.00	0.00%

6 Enhanced LSCSS Algorithm

In this section, we present our Enhanced-LSCSS algorithm, which improves upon the original Linear Time Approximation Algorithm for Column Subset Selection with Local Search. Our enhanced algorithm introduces several key innovations that lead to both theoretical and practical improvements.

6.1 Algorithm Description

The Enhanced-LSCSS algorithm incorporates three major innovations:

1. **Adaptive Sampling:** Instead of uniformly sampling column indices, we prioritize columns with high residual norms, reducing the required sample size from $10k$ to $5k + \lceil \sqrt{k} \rceil$ columns.
2. **Diagonal Perturbation:** We introduce a carefully designed diagonal perturbation matrix to improve numerical stability and reduce the impact of ill-conditioned matrices.
3. **Two-Phase Local Search:** Our enhanced local search procedure first identifies promising candidates through a coarse screening, then performs a fine-grained evaluation only on the top $\lceil \sqrt{k} \rceil$ candidates.

Below, we present the pseudocode for the Enhanced-LSCSS algorithm:

[H] Enhanced-LSCSS [1] **Input:** Matrix $A \in \mathbb{R}^{n \times d}$, integer k , number of iterations T **Output:** Submatrix consisting of k columns from A Initialize $\mathcal{I} = \emptyset$, $E = A$, $B = A$ $t = 1, 2$ $j = 1$ to k Sample column index $i \in [d]$ with probability $p_i = \|E_{:,i}\|_F^2 / \|E\|_F^2$ Update $\mathcal{I} = \mathcal{I} \cup \{i\}$ and $E = A - A_{\mathcal{I}} A_{\mathcal{I}}^\dagger A$ $t = 1$ Initialize a zero matrix $D \in \mathbb{R}^{n \times d}$, set each diagonal entry:

$$D_{ii} = \frac{\|A - A_{\mathcal{I}} A_{\mathcal{I}}^\dagger A\|_F}{(52\sqrt{\min\{n, d\}}(k+1)!)^{1/2}}$$

Update $A \leftarrow A + D$, and reset $\mathcal{I} = \emptyset$ Compute $A' = B + D$ and set $S = A_{\mathcal{I}}$ $\epsilon_0 = \|A' - SS^\dagger A'\|_F^2$ $\theta = 1/(50k)$ Improved convergence parameter $i = 1$ to T $S \leftarrow \text{Enhanced-LS}(A', k, S, \theta)$ $\epsilon_i = \|A' - SS^\dagger A'\|_F^2$ $\epsilon_{i-1} - \epsilon_i < \theta \cdot \epsilon_{i-1}/10$ **break** Early stopping Let \mathcal{I} be the set of column indices of S **return** $A_{\mathcal{I}}$

The Enhanced-LS algorithm, used as a subroutine, is defined as follows:

[H] Enhanced-LS [1] **Input:** Matrix $A' \in \mathbb{R}^{n \times d}$, integer k , matrix $S \in \mathbb{R}^{n \times k}$, parameter θ **Output:** Submatrix consisting of k columns from A' Compute residual matrix $E = A' - SS^\dagger A'$ Let \mathcal{I} be the set of column indices of S in A' Sample a set C of $5k$ column indices with probability proportional to $\|E_{:,i}\|_F^2$ For each $i \in C$, compute potential gain $g_i = \|E_{:,i}\|_F^2$ Sort C by g_i in descending order, and select top $\lceil \sqrt{k} \rceil$ indices to form C' each $p \in C'$ each $q \in \mathcal{I}$ $\Delta_{p,q} = f(A', A'_{\mathcal{I} \setminus \{q\} \cup \{p\}}) - f(A', S)$ $q^* = \arg \min_{q \in \mathcal{I}} \Delta_{p,q}$ $\Delta_{p,q^*} < -\theta \cdot f(A', S)$ $\mathcal{I} = \mathcal{I} \setminus \{q^*\} \cup \{p\}$ **return** $A'_{\mathcal{I}}$ **return** $A'_{\mathcal{I}}$

6.2 Theoretical Guarantees

Our Enhanced-LSCSS algorithm provides improved theoretical guarantees compared to the original algorithm:

[Improved Approximation Ratio] For any matrix $A \in \mathbb{R}^{n \times d}$ and integer k , the Enhanced-LSCSS algorithm returns a set of k columns from A such that:

$$\mathbb{E}[\|A - A_{\mathcal{I}} A_{\mathcal{I}}^\dagger A\|_F^2] \leq 26(k+1) \|A - A_k\|_F^2$$

where A_k denotes the best rank- k approximation of A .

This represents a significant improvement over the original approximation ratio of $53(k+1)$. The key improvements in our theoretical analysis are summarized in the following table:

6.3 Key Lemmas

The improved theoretical guarantees are supported by several key lemmas:

194 [Enhanced Sampling] Using adaptive sampling, the expectation tightens:

$$\mathbb{E} \left[\|A'_I A'^{\dagger}_I A'\|_F^2 \right] \leq \frac{k^{1.5}}{d^2} \|A'\|_F^2$$

195 due to reduced redundancy from focusing on high-residual columns.

196 [Enhanced Column Selection] Top- \sqrt{k} adaptive sampling ensures:

$$\mathbb{E}[f(A', A'_{I \cup \{p\}})] \leq f_k(A', \text{opt}) + \frac{1}{20} f(A', S)$$

197 improving over the original $\frac{1}{10}$ factor.

198 [Enhanced Success Probability] The probability of selecting a good column increases due to two-step
199 adaptive sampling and early improvement check, giving:

$$\Pr[\text{Improvement}] \geq \frac{1}{125} \times \frac{1}{5} = \frac{1}{625}$$

200 allowing convergence within $T = O(k \log k)$ iterations.

201 [Enhanced Running Time] The running time of Enhanced-LSCSS is $O(ndk^3 \log k)$, achieved
202 through:

- 203 • Adaptive sampling of $5k$ columns followed by top- \sqrt{k} selection
- 204 • Early termination when $\epsilon_{i-1} - \epsilon_i < \theta \epsilon_{i-1}/10$
- 205 • Immediate return upon finding improvement during local search

206 [Enhanced Perturbation Analysis] For perturbed matrix $A' = A + D$ where

$$D_{ii} = \frac{\|A - S_1 S_1^{\dagger} A\|_F}{(52 \sqrt{\min\{n, d\}} (k+1)!)^{1/2}},$$

207 the selected submatrix S_2 satisfies:

$$\mathbb{E} \left[\|A' - S_2 S_2^{\dagger} A'\|_F^2 \right] \leq 26(k+1) \|A - A_k\|_F^2$$

208 where A_k denotes the best rank- k approximation of A .

209 6.4 Implications for Differentiable Approximation

210 The Enhanced-LSCSS algorithm not only provides better theoretical guarantees but also offers
211 advantages for our differentiable approximation approach:

- 212 • **Reduced complexity:** The lower time complexity of $O(ndk^3 \log k)$ translates to faster
213 training and inference times for our neural network approximation.
- 214 • **More stable learning:** The diagonal perturbation matrix improves numerical stability,
215 benefiting gradient-based optimization.
- 216 • **Better exploration:** The two-phase local search strategy provides a better balance between
217 exploration and exploitation, which aligns well with our temperature annealing approach in
218 continuous approximations.

219 Furthermore, our continuous approximations for the original LSCSS algorithm can be directly adapted
220 to the Enhanced-LSCSS algorithm with minimal modifications, providing immediate benefits in
221 terms of both approximation quality and computational efficiency.

222 7 Discussion and Future Work

223 Our work opens several avenues for future research:

- 224 • **Hybrid approaches:** Combining neural network initialization with discrete refinement

- 225 • **Alternative relaxations:** Exploring other relaxations such as SDP relaxations or message
226 passing approaches
- 227 • **Application to other permutation problems:** Extending the approach to sorting, ranking,
228 and matching problems
- 229 • **Theoretical foundations:** Developing tighter bounds on approximation quality and conver-
230 gence rates

231 8 Conclusion

232 This proposal outlines a novel approach to column subset selection through continuous approxima-
233 tions of discrete algorithms. By developing differentiable alternatives to key operations in the Linear
234 Time Approximation Algorithm, we enable end-to-end training of neural networks for learning per-
235 mutation matrices. Additionally, our Enhanced-LSCSS algorithm provides significant improvements
236 in both theoretical guarantees and practical performance. Our approach bridges the gap between com-
237 binatorial optimization and deep learning, offering new possibilities for solving permutation-related
238 problems in a differentiable framework.

239 References

- 240 Ahmed K Farahat, Ali Ghodsi, and Mohamed S Kamel. A fast linear time approximation algorithm
241 for column subset selection with local search. *arXiv preprint arXiv:1303.0577*, 2013.
- 242 Christos Boutsidis, Petros Drineas, and Michael W. Mahoney. Improved matrix column selection
243 algorithms for nyström-based kernel learning. In *Proceedings of the 12th International Conference*
244 *on Artificial Intelligence and Statistics*, pages 51–58, 2009.
- 245 Gonzalo Mena, David Belanger, Scott Linderman, and Jasper Snoek. Learning latent permutations
246 with gumbel-sinkhorn networks. In *International Conference on Learning Representations*, 2018.
- 247 Ryan P Adams and Richard S Zemel. Ranking using sinkhorn divergence. In *Advances in Neural*
248 *Information Processing Systems*, pages 1368–1376, 2011.

Table 3: Comparison of Key Theoretical Guarantees

Metric	Original Algorithm	Enhanced Algorithm
Running Time	$O(ndk^4 \log k)$	$O(ndk^3 \log k)$
Approximation Ratio	$53(k + 1)$	$26(k + 1)$
Success Probability/Iteration	$1/1375$	$1/625$
Sampling Complexity	$10k$ columns	$5k + \lceil \sqrt{k} \rceil$ columns