043

045

046

047

049

050

051

052

053

054

EKM: An Exact, Polynomial-Time Divide-and-Conquer Algorithm for the *K*-Medoids Problem

Anonymous Authors¹

Abstract

The K-medoids problem is a challenging combinatorial clustering task, widely used in data analysis applications. While numerous algorithms have been proposed to solve this problem, none of these are able to obtain an exact (globally optimal) solution for the problem in polynomial time. In this paper, we present EKM: a novel algorithm for solving this problem exactly with worst-case $O(N^{K+1})$ time complexity. The algorithm is provably correct by construction, obtained using formal program derivation steps. We demonstrate the effectiveness of our algorithm by comparing it against various approximate methods and a stateof-the-art, exact branch-and-bound (BnB) algorithm on numerous real-world datasets. Our algorithm can not only provide provably exact solutions, but also consume much less time over all datasets compared with the BnB algorithm. We also show that the wall-clock time of our algorithm aligns with its worst-case time complexity analysis on synthetic datasets. In contrast, a stateof-the-art BnB algorithm not only exhibits exponential time complexity even for fixed K but also frequently produces erroneous solutions. This highlights the importance of employing formal, correct derivation steps when constructing exact algorithms.

1. Introduction

In machine learning (ML), *K*-medoids is the problem of partitioning a given dataset into *K* clusters where each cluster is represented by one of its data points, known as a *medoid*. A plethora of approximate/heuristic methodologies have been proposed to solve this problem, such as PAM (Partitioning Around Medoids) Kaufman (1990), CLARANS (Clustering Large Applications based on randomized Search) (Ng & Han, 2002), and variants of PAM (Schubert & Rousseeuw, 2021; Van der Laan et al., 2003).

While approximate algorithms are computationally efficient and have been scaled up to large dataset sizes, because they are approximate they provide no guarantee that the *exact* (*globally optimal*) solution can be obtained. In *high-stakes* or *safety-critical* applications, where errors are unacceptable or carry significant costs, we want the best possible partition given the specification of the clustering problem. Only an exact algorithm can provide this guarantee.

However, there are relatively few studies on exact algorithms for the K-medoids problem. The use of branchand-bound (BnB) methods predominates research on this problem (Ren et al., 2022; Elloumi, 2010; Christofides & Beasley, 1982; Ceselli & Righini, 2005). An alternative approach is to use off-the-shelf mixed-integer programming solvers (MIP) such as Gurobi (Gurobi Optimization, 2021) or GLPK (GNU Linear Programming Kit) (Makhorin, 2008). These solvers have made significant achievements, for instance, Elloumi (2010); Ceselli & Righini (2005)'s BnB algorithm is capable of processing medium-scale datasets with a very large number of medoids. More recently, Ren et al. (2022) designed another BnB algorithm capable of delivering tight approximate solutions—with an optimal gap of less than 0.1%—on very large-scale datasets, comprising over one million data points with three medoids, although this required a massively parallel computation over 6,000 CPU cores.

These existing studies on the exact *K*-medoids problem have three defects. Firstly, most previous studies on exact algorithms impose a computation time limit, thus exact solutions are rarely actually calculated; the only rigorously exact cases are those computed by (Christofides & Beasley, 1982; Ceselli & Righini, 2005) for small datasets with a maximum size of 150 data items. Secondly, worst-case time and space complexity analyses are not reported in studies of BnB-based algorithms. Such BnB methods have exponential worst-case time/space complexity, yet this critical aspect is rarely discussed or analyzed rigorously. As a result, existing studies on the exact algorithms for the *K*-medoids problem either omit time complexity analysis (Elloumi, 2010;

¹Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

Preliminary work. Under review by the International Conference on Machine Learning (ICML). Do not distribute.

Christofides & Beasley, 1982; Ceselli & Righini, 2005) or overlook essential details (Ren et al., 2022) such as algorithm complexity with respect to cluster size and the impact 058 on performance of upper bound tightness. This omission 059 undermines the reproducibility of findings and hinders ad-060 vancement of knowledge in this domain. Thirdly, proofs of 061 correctness of these algorithms is often omitted. Exact solu-062 tions require such mathematical proof of global optimality, 063 yet many BnB algorithm studies rely on weak assertions 064 or informal explanations that do not hold up under close 065 scrutiny (Fokkinga, 1991).

066 In this report, we take a fundamentally different approach, 067 by combining several modern, broadly applicable algorithm 068 design principles from the theories of constructive algo-069 rithmics (Meertens, 1986; Jeuring & Pekela, 1993; Bird 070 & de Moor, 1996). The derivation of our algorithm is obtained through a rigorous, structured approach known as the 072 *Bird-Meertens formalism* (or the *algebra of programming*) 073 (Meertens, 1986; Jeuring & Pekela, 1993; Bird & de Moor, 074 1996). This formalism enables the development of an effi-075 cient and correct algorithm by starting with an obviously 076 correct but possibly inefficient algorithm and deriving an 077 equivalent, efficient implementation by way of equational 078 reasoning steps. Thus, the correctness of the algorithm is 079 assured, skipping the need for error-prone, post-hoc induction proofs. Our proposed algorithm for the K-medoids 081 problem has worst-case $O(N^{K+1})$ time complexity and 082 furthermore is amenable to optimally efficient vectorization 083 and parallelization.

085 The paper is organized as follows. In Section 2, we ex-086 plain in detail how our efficient EKM algorithm is derived 087 through provably correct equational reasoning steps. Sec-088 tion 3 presents the results of empirical computational com-089 parisons between our proposed algorithms, approximate 090 algorithms, and the state-of-the-art exact algorithm devel-091 oped by Ren et al. (2022). These comparisons were con-092 ducted on datasets from the UCI machine learning repos-093 itory, and a few datasets from Ren et al. (2022). Further-094 more, we provide a detailed time-complexity analysis com-095 paring Ren et al. (2022)'s algorithm with our algorithm. 096 The results demonstrate that, even for fix K, Ren et al. 097 (2022)'s algorithm exhibits exponential time complexity. 098 In contrast, our proposed algorithms consistently achieve 099 predictable polynomial time complexity. In addition, the 100 naive brute-force algorithm for solving this problem also has $O(N^{K+1})$ complexity. To highlight the difference between our algorithm and a brute-force algorithm, we also compared the performance of our algorithm in exhaustive 104 generation and solving the K-medoids problem by compar-105 ing against the classical one-by-one enumeration strategy by 106 using itertools library. Section 4 summarizes the contributions of this study, reviews related work. Finally, Section 5 suggests future research directions. 109

2. Theory

Our novel EKM algorithm for K-medoids is derived using rational algorithm design steps, described in the following subsections. There are two main advantages to this approach of deriving correct (globally optimal) algorithms from specifications. First, since the exhaustive search algorithm is correct for the MIP problem, any algorithm derived through correct equational reasoning steps from this specification is also (provably) correct. Second, the proof is short and elegant: a few generic theorems for shortcut fusion already exist and when the conditions of these shortcut fusion theorems hold, it is simple to apply them to a specific problem such as the K-medoids problem here. This provably correct design step is vitally important in the study of the exact algorithm, as we will see shortly in the experiments, where the state-of-art BnB algorithm produces solutions lower than exact solutions frequently. This highlights the importance of using a correct-by-construction algorithm derivation process rather than relying on ad-hoc BnB algorithms.

2.1. Specifying a mixed-integer program (MIP) for the *K*-medoids problem

We denote a data set \mathcal{D} consists of N data points $x_n =$ $(x_{n1}, x_{n2}, \dots, x_{nD}) \in \mathbb{R}^D, \forall n \in \{1, \dots, N\} = \mathcal{N}, \text{ where }$ D is the dimension of the *feature space*. We assume the data items $\pmb{x}_n \in \mathcal{D}$ are all stored in a ordered sequence (list) $\mathcal{D} = [x_1, x_2, \dots, x_N]$. In clustering problem, we need to find a set of centroids $\mathcal{U} = \{ \boldsymbol{\mu}_k : k \in \mathcal{K} \}$, where $\mathcal{K} = \{1, \ldots, K\}$ in \mathbb{R}^D , and each centroid μ_k associated with a unique *cluster* C_k , which subsumes all data points x are the *closest* to this centroid than other centroids, this closeness is defined by a distance function $d(x, \mu)$, there is no constrains on the distance function in the Kmedoids problem, a common choice is the squared Euclidean distance function $d_2(\boldsymbol{x}, \boldsymbol{\mu}) = \|\boldsymbol{x} - \boldsymbol{y}\|_2^2$. For each set of centroids \mathcal{U} we have a set of *disjoint clusters* $\mathcal{C} = \{C_1, C_2, \dots, C_K\},$ and set \mathcal{K} is then called the *cluster* labels.

The K-medoids problem is usually specified as the following MIP

$$\mathcal{U}^* = \underset{\mathcal{U}}{\operatorname{argmin}} E\left(\mathcal{U}\right)$$

subject to $\mathcal{U} \subseteq \mathcal{D}, |\mathcal{U}| = K,$ (1)

where $E(\mathcal{U}) = \sum_{k \in \mathcal{K}} \sum_{x_n \in C_k} d(x_n, \mu_k)$ is the *objective* function for the K-medoids problem, and \mathcal{U}^* is a set of centroids that optimize the objective function $E(\mathcal{U})$. In our study, we make no assumptions about the objective function, allowing any arbitrary distance function to be applied without affecting the worst-case time complexity while still obtaining the exact solution.

2.2. Representing the solution as an (inefficient) exhaustive search algorithm

110

111

112

113

114

115

116 117

118 119

120 121 122 In the theory of transformational programming (Bird & de Moor, 1996; Jeuring & Pekela, 1993), combinatorial optimization problems such as (1) are solved using the following, generic, *generate-evaluate-select* algorithm,

$$s^* = sel_E\left(eval_E\left(gen\left(\mathcal{D}\right)\right)\right). \tag{2}$$

123 Here, the generator function $gen : \mathcal{D} \to [[\mathbb{R}^D]]$, enumer-124 ates all possible combinatorial configurations $s : [\mathbb{R}^D]$ 125 (here, configurations consist of a list of data items repre-126 senting the K medoids) within the solution (search) space, 127 S (here, this is the set of all possible size K combina-128 tions). For most problems, *gen* is a recursive function so 129 that the input of the generator can be replaced with the 130 index set \mathcal{N} and rewritten $gen(n), \forall n \in \mathcal{N}$. The *evalu*-131 ator $eval_E : [[\mathbb{R}^D]] \to [([\mathbb{R}^D], \mathbb{R})]$ computes the ob-132 jective values r = E(s) for all configurations s generated 133 by qen(n) and returns a list of *tupled configurations* (s, r). 134 Lastly, the selector $sel_E : [([\mathbb{R}^D], \mathbb{R})] \to ([\mathbb{R}^D], \mathbb{R})$ 135 select the best configuration s^* with respect to objective E. 136

Algorithm (2) is an example of exhaustive or brute-force 137 search: by generating all possible configurations in the 138 139 search space S for (1), evaluating the corresponding objective E for each, and selecting an optimal configuration, it 140 is clear that it must solve the problem (1) exactly. Taking 141 a different perspective, (2) can be considered as a generic 142 program for solving the MIP (1). However, program (2) 143 is generally inefficient due to combinatorial explosion; the 144 145 size of $gen(\mathcal{D})$ is often exponential (or worse) in the size of \mathcal{D} . 146

147 To make this exhaustive solution practical, two aspects must 148 be considered to improve the efficiency of the program. The 149 first aspect is the design of an efficient generator. By recog-150 nizing that in many combinatorial problems, generators can 151 be expressed as efficient recursions that take advantage of 152 solving difficult problems by first addressing easier subprob-153 lems, Another principle is known as shortcut fusion. For 154 many problems, it is often possible to fuse sel_E , $eval_E$, and 155 gen into a single, fused, efficient program. This fusion can 156 save substantial amounts of computation because it elim-157 inates the need to generate and store every configuration. 158 However, the fusion is only possible when the generator 159 is defined as a recursive program, this immediately elim-160 inates any one-by-one enumeration approach, as none of 161 them are defined recursively. The remaining paper will be 162 dedicated to explaining an efficient generator for solving the 163 K-medoids problem and the fusion in this problem. 164



Figure 1. A generic recursive generator of all possible sublists for list $[x_1, x_2, x_3, x_4]$, in each recursive stage *n*, the sublists of size *k* (combinations) are automatically grouped into list S_k^n , for all $k \in \{0, 1, ..., n\}$.

2.3. Constructing an efficient, recursive combinatorial configuration generator

In *K*-medoids there are only $N \times (N-1) \times \cdots \times (N-K)$ ways of selecting centroids whose corresponding assignments are potentially distinct. Thus, the solution space *S* for the *K*-medoids problem consists of all size *K* combinations (*K*-sublists) of the data input *D*, denoted as $s : [\mathbb{R}^D]$. If we can design an efficient, recursive combination generator $gen_{combs} : \mathcal{N} \times \mathcal{K} \to [\mathbb{R}^D]$ to enumerate all possible size *K* combinations, at least one set of such centroids corresponds to an optimal value of the clustering objective *E*. An efficient and structured combination generator which achieves this, will be described next.

Denoting S_k^n : $[\mathbb{R}^D]$ as the list that stores all possible size k combinations from list $[x_1, x_2, \dots x_n]$ of length n, and $S^n = [S_k^n | k \in \{0, \dots n\}]$ as the list of all *sublists* of that list. The key to constructing an efficient combination generator can be reduced to solving the following problem: given all sublists S^n for list $\mathcal{D}_1 = [x_1, x_2, \dots x_n]$, and S^m for list $\mathcal{D}_2 = [y_1, y_2, \dots y_m]$, construct all sublists S^{n+m} for list $\mathcal{D}_1 \cup \mathcal{D}_2$. This problem captures the essence of Bellman's *principle of optimality* (Bellman, 1954) i.e. we can solve a problem by decomposing it into smaller *subproblems*, and the solutions for these subproblems are combined to solve the original, larger problem. In this context, constructing all combinations S^{n+m} for list $\mathcal{D}_1 \cup$ \mathcal{D}_2 is the problem; all possible combinations S^n , S^m for lists $\mathcal{D}_1, \mathcal{D}_2$ respectively, are the smaller subproblems.

In previous work of Little et al. (2024), it was demonstrated that the *constraint algebra* (typically a *group* or *monoid*) can be implicitly embedded within a generator semiring through the use of a *convolution algebra*. Subsequently, the filter that incorporates these constraints can be integrated into the generator via the *semiring fusion theorem* (for details of semiring lifting and semiring fusion, see Little et al. (2024)). In our case, we can substitute the convolution algebra used in Little et al. (2024) with the *generator semir*ing ([[T]], \cup , \circ , \emptyset , [[]]) (T stands for type variable), then we have following equality

175 where $l_1 \circ l_2 = [s \cup s' \mid s \in l_1, s' \in l_2]$ is the cross-join 176 operator on lists l_1 and l_2 obtained by concatenating each 177 element s in configuration l_1 with each element s' in l_2 . For 178 instance, $[[1], [2]] \circ [[3], [4]] = [[1, 3], [1, 4], [2, 3], [2, 4]].$ 179 Informally, (3) is true because the size k combinations for 180 the list $\mathcal{D}_1 \cup \mathcal{D}_2$ should be constructed from all possible 181 combinations \mathcal{S}_i^n in \mathcal{D}_1 and \mathcal{S}_j^m in \mathcal{D}_2 such that i + j = k182 for all $0 \le i, j \le k$. In other words, all possible size k 183 combinations should be constructed by joining all possible 184 combinations with a size smaller than k. 185

186 Definition (3) is a special kind of *convolution prod*-187 *uct* for two lists $l_a = [a_0, a_2, \dots, a_{n-1}]$ and $l_b = [b_0, b_2, \dots, b_{m-1}]$, 189

$$conv(f, l_1, l_2, k) = [c_0, c_1, \dots, c_k],$$
 (4)

where c_k is defined as

170

171

190

191

193

195

196

$$c_k = \bigcup_{\substack{i+j=k\\0\le i,j\le k}} f\left(a_i, b_j\right), \quad 0\le k\le n+m-1.$$
(5)

197 Given S^n and S^m , it is not difficult to verify that 198 $\mathcal{S}^{n+m} = conv(\circ, \mathcal{S}^n, \mathcal{S}^m, n+m-1)$, representing all 199 sublists (combinations) for list $\mathcal{D}_1 \cup \mathcal{D}_2$. Furthermore, all 200 combinations with size smaller than k can be obtained by 201 calculating $conv(\circ, S^n, S^m, k)$, denoted as $S^{n+1}_{\leq k}$. Using 202 equation $\mathcal{S}_{\leq k}^{n+1} = conv\left(\circ, \mathcal{S}_{i}^{n}, \mathcal{S}_{j}^{\prime m}, k\right)$, a combination gen-203 erator recursion $gen_{combs} : \mathcal{N} \times \mathcal{K} \to [\mathbb{R}^D]$ can be defined 204 by following *pattern matching*

$$\begin{array}{ll} 206\\ 207\\ gen_{combs}\left(\left[\;\right],k\right) = \left[\left[\left[\;\right]\right]\right]\\ 208\\ gen_{combs}\left(\left[\boldsymbol{x}_{n}\right],k\right) = \left[\left[\left[\;\right]\right],\left[\left[\boldsymbol{x}_{n}\right]\right]\right]\\ 209\\ gen_{combs}\left(xs \cup ys,k\right) =\\ 210\\ conv\left(\circ,gen_{combs}\left(xs,k\right),gen_{combs}\left(ys,k\right),k\right). \end{array}$$

where $\mathcal{D} = xs \cup ys$. The function $gen_{combs}(\mathcal{D}, k)$ generate all combinations with a size k or less. For instance, $gen_{combs}([x_1, x_2, x_3], 2) = S_{\leq 2}^3([x_1, x_2, x_3]) =$ [[[]], $[[x_1], [x_2], [x_3]], [[x_1, x_2], [x_1, x_3], [x_2, x_3]]]$. See Figure 1 for an illustration of the process by which recursive combination generator (6) operates. The generator (6) is an instance of the *divide-and-conquer* (D&C) algorithm, for which the original problem $(xs \cup ys)$ can be divided into arbitrary two disjoint subproblems xs and ys and the subproblems are solved independently.

One of the benefits of being a D&C algorithm is that all D&C algorithms are easy to parallelize. Moreover, our generator is specifically designed for ease of vectorization and parallelization, making it well-suited for modern hardware optimized for array-based operations. By deliberately organizing combinations of the same size into the same list, we enable storage in a pre-allocated array. Operations on these contiguous memories will decrease the possibility of cache misses, which makes our program more efficient. Alternative methods are typically list-based, such as the generator proposed by (He & Little, 2023) and the Gray code generator (Kreher & Stinson, 1999). These methods rely on dynamic memory allocation, making them significantly less efficient than our proposed generator. Alternatively, one-by-one generation approaches, such as lexicographical generation, are primarily designed for random generation and often involve substantial constant factors that are obscured within the big-O notation.

Moreover, the one-by-one enumeration strategy does not permit any form of fusion, rendering future speed-ups for this method unattainable. Consequently, both the best-case and worst-case complexities of the one-by-one enumeration strategy remain $O(N^{K+1})$. In contrast, the recursive structure of our generator enables fusion. Previous research by He & Little (2023) demonstrated an order-of-magnitude speed-up using the simplest bounding techniques. While our focus is on illustrating the principles of program derivation, and we thus omit such optimizations, we show that even the plain version of our algorithm is significantly faster than the state-of-the-art BnB algorithm.

2.4. Applying shortcut fusion equational transformations to derive an efficient, correct implementation

The combination generator gen_{combs} constructed above, gives us an efficient recursive basis to use 2 to derive an efficient algorithm for solving the *K*-medoids problem. Thus, a provably correct algorithm for solving the *K*-medoids problem can be rendered as

$$s^* = sel_E\left(eval_E\left(gen_{combs}\left(xs,k\right)\right)\right). \tag{7}$$

As we mentioned above, if both selector sel_E and evaluator $eval_E$ can be fused into the generator by identifying specific shortcut fusion theorems, then a significant amount of computational effort and memory can be saved.

For any configurations of size k, we can evaluate the objectives of the k-combinations immediately once its was generated. Since adding more medoids provably decreases

the objective value, we do not evaluate combinations of sizes smaller than k. However, this approach can be generalized to solve the k-medoids problem for all $k \leq K$. In other words, the following fusion condition holds:

220

221

222

223

224

225

226

227

228

229

230

231 232 233

239 240

263

264

265

273

274

$$sel_{E} (eval_{E} (gen_{combs} (xs, k))) = \\ sel_{E,k} (conv (\circ_{E,k}, gen_{combs} (xs), gen_{combs} (ys), k)),$$
(8)

where $sel_{E,k}$ select the configuration with lowest objective among all size k combinations, and $\circ_{E,k}$ is the cross-join operator described above augmented with evaluation updates,

$$l_1 \circ_{E,k} l_2 = \begin{cases} [(s_1 \cup s_2, E(s_1 \cup s_2)) \mid \\ (s_1, r) \in l_1, (s_2, r) \in l_2] & \text{if } |s_1 \cup s_2| = k \\ \\ [(s_1 \cup s_2, r) \mid \\ (s_1, r) \in l_1, (s_2, r) \in l_2] & \text{otherwise,} \end{cases}$$
(9)

Due of the universal property (Bird & de Moor, 1996), the solution obtained by recursion

is also a solution to the specification (7), without the proof
 of induction.

250 We now analyze the time and space complexity of the algorithm $EKM(\mathcal{D}, K)$ for a dataset of size N. Since the 251 pairwise distances between data points can be stored in 252 a pre-allocated matrix, accessing any element in this ma-253 trix requires only O(1) time. Evaluating the objective of a 254 255 K-combination requires indexing $K \times N$ entries. Thus evaluating each K-combination takes O(N) time, and there are 256 N257 $= O(N^{K})$ combinations in total. Therefore, the 258 K

overall running time of $EKM(\mathcal{D}, K)$ is $O(N^{K+1})$. +Additionally, $EKM(\mathcal{D}, K)$ only stores partial configurations of sizes smaller than K - 1 resulting in a space complexity of $O(N^{K-1})$.

3. Experiments

In this section, we analyze the computational performance
of our algorithm EKM on both synthetic and real-world data
sets. Our evaluation aims to test the following predictions:
(a) EKM always obtains the best objective value¹; (b) wallclock run-time matches the worst-case polynomial time
complexity analysis; (c) the state-of-art BnB algorithm (Ren

Table 1. Empirical comparison of the EKM algorithm, against widely-used approximate algorithms (PAM, Fast-PAM, and CLARANS) and the state-of-art exact BnB algorithm developed by Ren et al. (2022), for K = 3, in terms of sum-of-squared errors (*E*), smaller is better. For Ren et al. (2022)'s aglorithm, we include both the *upper bound* and the *optimal gap*. The best-performing algorithm is highlighted in **bold**, while incorrect solutions are marked in red. Although some of the upper bounds returned by BnB algorithm are exact, these values are not marked in bold, as the optimal gap has not yet converged to zero. Wall clock time in brackets (seconds), we set a time limit 1.08×10^4 for BnB algorithm.

UCI DATASET	Ν	D	EKM (OURS)	REN'S BNB	PAM	FASTER-PAM	CLARANS
							l
LM	338	3	3.96 × 10-	(2.21 × 10 ¹)	3.99×10^{-3}	4.07×10^{-3}	5.33 × 10 ⁻
			(8.20 × 10 ⁻²)	(2.31 × 10)	(4.02 × 10 -)	(3.01 × 10 · ·)	(0.14)
UKM	403	5	8.36 × 10 ¹	5.51 × 10 ¹	8.44×10^{1}	8.40×10^{1}	1.16×10^{2}
o nan	40.5		(1.37)	(1.62×10^3)	(8.57×10^{-3})	(3.21×10^{-3})	(4.98×10^{1})
				< 0.1%			
LD	345	5	3.31×10^{5}	1.21×10^{1}	3.56×10^{5}	3.31×10^{5}	4.68×10^{5}
			(8.3×10^{-1})	(2.47×10^{1})	(4.11×10^{-3})	(3.87×10^{-3})	(3.40)
			(0.0 // 20)	$\leq 0.1\%$		(01017/100)	
ENERGY	768	8	2.20×10^{6}	2.20×10^{6}	2.28×10^{6}	2.28×10^{6}	2.97×10^{6}
			(1.37×10^{1})	(1.68×10^{1})	(6.95×10^{-3})	(3.94×10^{-3})	(2.71)
				$\leq 0.1\%$			
VC	310	6	3.13×10^{5}	1.50×10^{5}	3.13×10^{5}	3.58×10^{5}	5.27×10^{5}
			(6.82×10^{-1})	(3.83×10^2)	(3.15×10^{-3})	(5.36×10^{-3})	(2.58)
				$\leq 0.1\%$			
WINE	178	13	$2.39 imes 10^6$	1.16×10^{4}	2.39×10^{6}	2.63×10^{6}	6.86×10^{6}
			(2.22×10^{-1})	(5.17×10^{1})	(1.06×10^{-3})	(2.34×10^{-3})	(5.56×10^{-1})
				$\leq 0.1\%$			
YEAST	1484	18	8.37×10^{1}	6.57×10^{1}	8.42×10^{1}	8.42×10^{1}	1.05×10^{2}
			(1.74×10^2)	(1.08×10^4)	(9.54×10^{-2})	(6.08×10^{-2})	(1.73×10^2)
				$\leq 39.19\%$			
IC	3150	0.13	6.9063×10^{9}	6.18×10^{9}	6.9105×10^{9}	6.9063×10^{9}	1.44×10^{10}
			(4.53×10^{9})	(6.37×10^{-9})	(8.68×10^{-1})	(1.91×10^{-1})	$(2.70 \times 10^{*})$
				0			
WDG	5000	21	1.67×10^{-5}	1.60×10^{5}	1.67×10^{5}	1.67×10^{5}	2.77×10^{-5}
			(5.23×10^{-1})	(1.08 × 10 ⁻)	(1.34)	(1.97×10^{-1})	(5.32×10^{-6})
				\$ 785.05%			
IRIS	150	4	8.40 × 10-	8.46 × 10 ⁻	(2.51×10^{-3})	(1.02×10^{-3})	(2.22×10^{-1})
			(1.57×10^{-5})	(2.31 × 10)	(2.51 × 10 -)	(1.03 × 10 ··)	(2.32 × 10)
SEEDS	210	7	5 08 × 10 ²	5 08 × 10 ²	E 08 V 10 ²	5 08 × 10 ²	1.12×10^{3}
SEEDS	210	/	5.98 × 10	(2.42×10^{1})	5.98 X 10	5.98 X 10	(7.82×10^{-1})
			(2.85 × 10)	< 0.1%	(1.14 × 10 ⁻¹)	(3.59 X 10 ·)	(1.02 \ 10)
GLASS	214	9	6.29×10^2	6.29×10^2	6.29×10^{2}	6.29×10^{2}	1.04×10^{3}
012100			(2.90×10^{-1})	(3.13×10^{1})	(1.01×10^{-3})	(1.62×10^{-3})	(2.27)
			(2000)(200))	< 0.1%	(2002.00.20)	(
BM	249	6	8.63×10^{5}	8.63×10^{5}	8.76×10^{5}	8.63×10^{5}	1.33×10^{6}
			(3.96×10^{-1})	(1.19×10^2)	(4.12×10^{-3})	(1.61×10^{-3})	(1.02×10^{1})
				$\leq 0.1\%$			
HF	299	12	7.83×10^{11}	7.83×10^{11}	7.83×10^{11}	7.83×10^{11}	1.88×10^{12}
			(6.26×10^{-1})	(5.20×10^{1})	(1.00×10^{-3})	(4.87×10^{-3})	(6.55×10^{-1})
				0%			
WHO	440	7	8.33×10^{10}	8.33×10^{10}	8.33×10^{10}	8.33×10^{10}	1.21×10^{11}
			(1.98)	(3.71×10^{-1})	(5.62×10^{-3})	(2.81×10^{-3})	(8.12)
				S 0.1%			
UK	258	5	5.08×10^{-1}	5.08×10^{4}	5.08×10^{-2}	5.08×10^{-2}	6.89×10^{4}
			(3.78×10^{-1})	$(1.43 \times 10^{\circ})$	(2.47×10^{-5})	(2.47×10^{-5})	(2.67×10^{-1})
				5 0.1%			
HCV	572	12	2.75 × 10°	$2.75 \times 10^{\circ}$	2.75×10^{-3}	2.75×10^{-3}	4.75 × 10"
			(0.48)	(0.35 × 10)	(5.79 × 10 °)	(2.21×10^{-5})	(0.85 × 10)
ADC	740	10	2 22 V 10 ⁶	2 22 4 100	2.22×10^{6}	2.28 V 10 ⁶	2.06 × 106
ABS	740	19	(1.17×10^{1})	(6.23×10^2)	(2.11×10^{-2})	(5.00×10^{-3})	(7.80×10^{1})
			(1.1) × 10)	< 0.1%	(2.11 × 10)	(0.00 × 10)	(1.00 × 10)
TP	080	10	1.13×10^{3}	1.14×10^{3}	1.13×10^{3}	1.13×10^{3}	1.38×10^{3}
	200		(2.53×10^{1})	(1.08×10^4)	(5.14×10^{-2})	(1.14×10^{-2})	(2.59×10^2)
				< 89%	(0.14 × 10)	(1.14 × 10)	
SGC	1000	21	1.28×10^9	1.28×10^{9}	1.28×10^{9}	1.28×10^{9}	2.52×10^9
			(3.87×10^{1})	(1.75×10^2)	(1.71×10^{-1})	(4.22×10^{-2})	(2.24)
				$\leq 0.1\%$	((
HEMI	1995	7	9.91×10^{6}	9.91×10^{6}	9.91×10^{6}	9.91×10^{6}	1.66×10^{7}
			(9.00×10^2)	(3.92×10^2)	(3.64×10^{-1})	(6.99×10^{-2})	(9.53)
				$\leq 0.1\%$		(
PR2392	2392	2	2.13×10^{10}	2.13×10^{10}	2.13×10^{10}	2.13×10^{10}	3.47×10^{10}
			(1.29×10^3)	$(1.54 \times 10^3) \le$	(3.66×10^{-1})	(8.38×10^{-2})	(1.15×10^2)
				0.1%			

et al., 2022) for solving the K-medoids problem will have exponential time complexity even for fixed K in the worstcase. In our implementation, the matrix operations required at every recursive step are batch processed on a single GPU. We executed all the experiments on an Intel Core i9 CPU, with 24 cores, 2.4-6 GHz, 32 GB RAM and GeForce RTX 4060 Ti GPU.

Remarks, all comparison about the time complexity is executed in the sequential version of our algorithm over CPU only.

3.1. Performance on real-world datasets

We test the performance of our EKM algorithm against the approximate algorithms partition around medoids (PAM),

¹The squares Euclidean distance function was chosen for the experiments, any other proper metrics could also be used.

275 Faster-PAM and Clustering Large Applications based on 276 RANdomized Search (CLARANS)² on 18 datasets from 277 the UCI Machine Learning Repository, two datasets from 278 Ren et al. (2022) (UK, HCV), and two open-source datasets 279 (PR2392, HEMI) from (Wang et al., 2022; Padberg & Ri-280 naldi, 1991; Ren et al., 2022). We show that, as expected, no 281 other algorithms can achieve better objective function values 282 (see Table 1), except in cases where Ren et al. (2022)'s BnB 283 algorithm returned incorrect solutions, which were clearly 284 invalid as they were several orders of magnitude lower than 285 our exact solutions.

286 Our experiments included real-world datasets with a max-287 imum size of N = 5,000. To the best of our knowledge, 288 the largest dataset for which an exact solution has been pre-289 viously obtained is N = 150, as documented by Ceselli 290 & Righini (2005) with K = 3. Existing literature on the 291 K-medoids problem has only reported exact solutions on very small datasets, primarily due to the use of BnB algo-293 rithms. Given their unpredictable run-time and worst-case exponential time complexity, most reported usage of BnB 295 algorithms impose a hard computational time limit to avoid 296 memory overflow or intractable run times. 297

In summary, Ren et al. (2022)'s algorithm returned only 299 approximate solutions (with an optimality gap greater than 300 zero), whereas our algorithm consistently produced prov-301 ably exact solutions in significantly less time. Furthermore, 302 for challenging datasets, such as WDG, Ren et al. (2022)'s 303 algorithm produced a solution with an optimality gap of 304 800% even after running for three hours! Moreover, in 305 nearly all datasets tested in our experiments but not in-306 cluded in Ren et al. (2022) (e.g. IC, Yearst, WDG, wine, 307 LD, VC, UKM and LM), their algorithm produces obvious 308 errorness solutions where upper bounds that were lower 309 than our exact solutions, which is fundamentally incorrect 310 as an upper bound cannot be lower than the exact solution.

311 The only reason that Ren et al. (2022) claim their algorithm 312 can handle datasets with over a million instances is that 313 they test on datasets that are inherently easy to classify-so 314 that even approximate algorithms can obtain exact solutions. 315 As we have demonstrated, almost all the datasets they use 316 can be solved exactly using PAM or Faster-PAM, which 317 achieve exact solutions with significantly fewer resources. 318 In contrast, Ren et al. (2022)'s algorithm requires an ex-319 cessive amount of computational power (6,000 CPU cores) 320 compared to approximate algorithms.

Moreover, we observed that Ren et al. (2022)'s algorithm exhibits **exponential** complexity even when *K* is fixed. This is evident from the **non-polynomial growth** in running time for experiments on datasets such as UK, BM, and Seeds.

328 329 Although these datasets have nearly identical sizes, their running time differs significantly in Ren et al. (2022)'s algorithm. In the following section, we provide an empirical analysis of their algorithm to further investigate this behavior.

3.2. Time complexity analysis without parallelization



Figure 2. Log-log wall-clock run time (seconds) for our algorithm (EKM) tested on synthetic datasets (left panel). The run-time curves from left to right (corresponding to K = 2, 3, 4, 5 respectively), have slopes 3.005, 4.006, 5.018, and 5.995, an excellent match to the predicted worst-case run-time complexity of $O(N^3)$, $O(N^4)$, $O(N^5)$, and $O(N^6)$ respectively. Log-linear wall-clock run-time (seconds) comparing EKM algorithm against Ren et al. (2022)'s algorithm by sampling UK dataset with K = 3 (right panel). On this log-linear scale, exponential run-time appears as a linear function of problem size N, whereas polynomial run-time is a logarithmic function of N.

We test the wall-clock time of our novel EKM algorithm on a synthetic dataset with cluster sizes ranging from K = 2to 5. When K = 2, the data size N ranges from 150 to 2,500, K = 3 ranges from 50 to 530, K = 4 ranges from 25 to 160, and K = 5 ranges from 30 to 200. The worstcase predictions are well-matched empirically (Figure 2, left panel).

As predicted, Ren et al. (2022)'s algorithm exhibits worstcase exponential time complexity even for a fixed K =3 (Figure 2, right panel), whereas our algorithm runs in polynomial time in the worst case.

3.3. Performance compared with one-by-one enumeration

It is easy to confuse our algorithm with the one-by-one enumeration strategy, as both share a worst-case time complexity $O(N^{K+1})$ time complexity in the worst-case. However, this classical approach of modeling algorithm performance using big-O notation is inaccurate in the context of combinatorial generation, as it ignores constant factors, which can significantly impact the actual performance of an algorithm in practice.

To provide a more detailed analysis, it is crucial to examine the wall-clock runtime of different generators in actual implementation. In this section, we compare our approach with

²We set the maximum number of neighbors examined as 4, and the number of iteration as 5.

- the classical combination generator often used in the machine learning community, the itertools package, which
- is based on the lexicographical generation algorithm (Kre-
- her & Stinson, 1999). We demonstrate that although both
- 334 generators have the same asymptotic complexity, the con-
- stant factors hidden in the big-O notation differ significantly,
- leading to a substantial disparity in performance.

Both experiments are based on CPU, but our algorithm is designed for easy implementation on a GPU. The performance difference would likely be even more pronounced with the additional speedup provided by GPU acceleration.

342 In Figure 3, the Python implementation of our algorithm 343 achieves performance improvements of up to 10-fold for the task of exhaustive generation and up to 1000-fold for 345 solving the K-medoids problem, compared to the C-based 346 itertools library on CPU. Furthermore, the performance 347 gap continues to expand as the combinatorial complexity of the problem increases. For exhaustive generation, the 349 performance improvement grows from 3-fold to 10-fold, 350 while for solving the K-medoids problem, it increases from 351 300-fold to 1000-fold. 352

4. Discussion

353

354

384

355 In this paper, we derived EKM, a novel exact algorithm 356 for the K-medoids clustering problem with worst-case 357 $O(N^{K+1})$ time complexity. One immediate objection to 358 our algorithm might be that it shares the same worst-case 359 complexity as the naive one-by-one enumeration strategy. 360 However, our algorithm demonstrates significantly better 361 performance in optimization tasks and offers the potential 362 for further speed-ups through fusion-an advantage that 363 is unattainable with the one-by-one enumeration approach. Since the K-medoids problem is NP-hard, it is unlikely that 365 the K in the exponent can be eliminated in the worst case. Therefore, our focus shifts to designing an efficient recur-367 sive algorithm suitable for modern hardware, such as GPUs. This approach paves the way for future acceleration through 369 bounding techniques based on efficient recursion, rather 370 than attempting to improve existing BnB algorithms, which 371 exhibit exponential worst-case complexity. The bounding 372 techniques used in BnB algorithms cannot improve the 373 worst-case complexity, as no acceleration can be employed 374 in the worst cases.

375 Besides presenting our novel EKM algorithm, we aim to 376 prompt researchers to reconsider the metric for assessing 377 the efficiency of exact algorithms to account for subtleties 378 beyond simple problem scale. In discussing the "goodness" 379 of exact algorithms for ML, it is critical to recognize that 380 focusing solely on the scalability of these algorithms-for 381 instance, their capacity to handle large datasets-does not 382 provide a comprehensive assessment of their utility. This 383

inclination to prioritize scalability when assessing exact algorithms arises from the perceived intractable combinatorics of many ML problems, and most of these problems are classified as NP-hard, so that no known algorithm can solve all instances of the problem in polynomial time.

Although many ML problems may be NP-hard in their most general form, they are often not specified in full generality. In such cases (which are quite common in practice), polynomial-time solutions may be available. Therefore, comparing algorithms solely on the basis of problem scale is neither fair nor accurate, as it ignores variations in actual run-times, such as the tightness of upper bounds, or parameters independent of problem scale.

However, for many ML problems, the problems specified for proving NP-hardness are not the same as their original definitions used in practical ML applications. Apart from the K-medoids problem, polynomial-time algorithms for solving the 0-1 loss classification problem (He & Little, 2023) and other K-clustering problems (Inaba et al., 1994; Tîrnăucă et al., 2018) have also been developed in the literature. If a polynomial-time algorithm does exist for these seemingly intractable problems, overemphasizing scalability can mislead scientific development, diverting attention from important measures such as memory usage, worst-case time complexity, and the practical applicability of the algorithm in real-world scenarios. For example, by setting the cluster size to one, the K-medoids problem can be solved exactly by choosing the closest data item to the mean of the data set, a strategy with O(N) time complexity in the worst-case. While this O(N) time algorithm can efficiently handle very large scale datasets, it does little to advance our understanding of the fundamental principles involved.

Therefore, judging an algorithm implementation solely by the scale of the dataset it can process is not an adequate measure of its effectiveness. Indeed, for large datasets, the use of exact algorithms may often be unnecessary as many highquality approximate algorithms provide very good results, supported by solid theoretical assurances. If the clustering model closely aligns with the ground truth, the discrepancy between approximate and exact solutions should not be significant, provided the dataset is sufficiently large. Thus, it is not surprising that Ren et al. (2022)'s algorithm can achieve excellent approximate solutions with more than one million data points, a typical occurrence in studies involving exact algorithms. Past research has shown that while exact algorithms can quickly find solutions, most of the effort is expended on verifying their optimality (Dunn, 2018; Ustun, 2017). It is possible that the first configuration generated by the algorithm is optimal, but proving its optimality without exploring the entire solution space S is impossible (unless additional information is provided).

For the study of the K-medoids problem, while algorithms

Submission and Formatting Instructions for ICML 2025



Figure 3. The running time comparison between the *itertools* library—a widely used Python library for generating combinations—and our combination generator for the task of exhaustive generation (The panels are arranged from left to right and top to bottom from K = 5 to K=8) and the task of solving the K-medoids problem (K = 3 to K = 6) demonstrates that our Python implementation achieves performance improvements of up to 10-fold and 1000-fold, respectively, over the C-based *itertools* library on the CPU.

presented in Ren et al. (2022); Ceselli & Righini (2005); Elloumi (2010); Christofides & Beasley (1982) are exact in principle, experiments reported by the authors do not demonstrate the actual computation of exact solutions, nor do they provide any theoretical guarantee on the computational time required to achieve satisfactory approximate solutions. If the application of the problem is concerned with only the approximate solution, it may be more beneficial to concentrate on developing more efficient or more robust heuristic algorithms. This could potentially offer more practical value in scenarios where approximate solutions are adequate.

5. Conclusions and future work

As our predictions and experiments show, our novel EKM algorithm clearly outperforms all other algorithms which can be guaranteed to obtain the exact globally optimal solution to the *K*-medoids problem. Whereas the approximate algorithms and the BnB algorithm designed by Ren et al. (2022) can only obtain approximate solutions (or wrong solutions) and use significantly more time.

427 Moreover, precise time and space complexity anlysis of the 428 EKM algorithm precise prediction of time and space re-429 quirements before attempting to solve a clustering problem, 430 in contrast to existing BnB algorithms which in practice 431 require a hard computation time limit to prevent memory 432 overflow or bypass the exponential worst-case run time. In 433 our experiments we were able to process datasets of up to 434 N = 5,000 data items, a considerable increase from the pre-435 vious maximum of around N = 150. We have also shown 436 that the state-of-the-art BnB algorithm, while claimed to 437 be optimal, produces erroneous solutions on many datasets. 438 Additionally, it only returns approximate solutions, requiring significantly more time across all datasets we tested. For difficult datasets, the algorithm by Ren et al. (2022)'s algorithm exhibits exponential complexity, even for fixed K.

The main disadvantage of our algorithm is that its space and time complexity is exponential in K. Thus for problems that involve a large number of medoids, our algorithm quickly becomes intractable. Currently, we only employs the plain form of our algorithm, in other words, we have not yet consider the speed-up by using bounding techniques. In the future, a more sophisticated bounding techniques could be developed; the inherently recursive structure of our combination generator makes this a relatively simple prospect.

With exact solutions for combinatorial ML problems, the memory-computation trade-off is always present. However, with BnB algorithms and MIP solvers, space complexity analysis is often omitted making it difficult to ascertain the actual memory requirements. Specifically, when using off-the-shelf MIP solvers, the memory required just to specify the problem can be substantial. For instance, to describe the *K*-medoids problem in MIP form (1), a constraint matrix³ of size $(N^2 + N + 1) \times N$ is required (Ren et al., 2022; Vinod, 1969). Indeed, memory overflow issues have been reported in almost all practical usage of BnB algorithms (Ceselli & Righini, 2005; Elloumi, 2010; Christofides & Beasley, 1982). Therefore, setting a computation time limit is a necessary restriction for BnB algorithms, a restriction which only applies to EKM at large values of *K* or *N*.

439

³There are N^2 constraints to ensure each point is assigned to exactly one cluster, N constraints to identify which data items are medoids, and one additional constraint for the number of medoids

440 Impact Statement

This paper presents work whose goal is to advance the field of Machine Learning. There are many potential societal consequences of our work, none which we feel must be specifically highlighted here.

References

441

442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485 486

487

488

489

490

491 492

493

494

- Bellman, R. The theory of dynamic programming. *Bulletin* of the American Mathematical Society, 60(6):503–515, 1954.
- Bird, R. and de Moor, O. The algebra of programming. 1996. URL http://www.cs.ox.ac.uk/ publications/Books/algebra/.
- Ceselli, A. and Righini, G. A branch-and-price algorithm for the capacitated p-median problem. *Networks: An International Journal*, 45(3):125–142, 2005.
- Christofides, N. and Beasley, J. E. A tree search algorithm for the p-median problem. *European Journal of Operational Research*, 10(2):196–204, 1982.
- Dunn, J. W. *Optimal trees for prediction and prescription*. PhD thesis, Massachusetts Institute of Technology, 2018.
- Elloumi, S. A tighter formulation of the p-median problem. *Journal of Combinatorial Optimization*, 19(1):69– 83, 2010.
- Fokkinga, M. M. An exercise in transformational programming: Backtracking and branch-and-bound. *Science of Computer Programming*, 16(1):19–48, 1991.
- Gurobi Optimization, L. Gurobi optimizer reference manual. 2021.
- He, X. and Little, M. A. An efficient, provably exact algorithm for the 0-1 loss linear classification problem. *ArXiv.* /*abs*/2306.12344, 2023.
- Inaba, M., Katoh, N., and Imai, H. Applications of weighted voronoi diagrams and randomization to variance-based k-clustering. In *Proceedings of the Tenth Annual Sympo*sium on Computational Geometry, pp. 332–339, 1994.
- Jeuring, J. T. and Pekela, T. *Theories for algorithm calculation*. Universiteit Utrecht, 1993.
- Kaufman, L. Partitioning around medoids (program pam). *Finding groups in data*, 344:68–125, 1990.
- Kreher, D. L. and Stinson, D. R. Combinatorial algorithms: generation, enumeration, and search. 30(1):33–35, 1999.

- Little, M. A., He, X., and Kayas, U. Polymorphic dynamic programming by algebraic shortcut fusion. *Formal Aspects of Computing*, May 2024. ISSN 0934-5043. doi: 10.1145/3664828. URL https://doi.org/10. 1145/3664828. (in press).
- Makhorin, A. GLPK (gnu linear programming kit). http://www.gnu.org/s/glpk/glpk.html, 2008.
- Meertens, L. Algorithmics: Towards programming as a mathematical activity. In *Proceedings of the CWI Symposium on Mathematics and Computer Science*, volume 1, pp. 289–334, 1986.
- Ng, R. T. and Han, J. CLARANS: A method for clustering objects for spatial data mining. *IEEE Transactions on Knowledge and Data Engineering*, 14(5):1003–1016, 2002.
- Padberg, M. and Rinaldi, G. A branch-and-cut algorithm for the resolution of large-scale symmetric traveling salesman problems. *Society for Industrial and Applied Mathematics Review*, 33(1):60–100, 1991.
- Ren, J., Hua, K., and Cao, Y. Global optimal k-medoids clustering of one million samples. *Advances in Neural Information Processing Systems*, 35:982–994, 2022.
- Schubert, E. and Rousseeuw, P. J. Fast and eager k-medoids clustering: O (k) runtime improvement of the pam, clara, and clarans algorithms. *Information Systems*, 101:101804, 2021.
- Tîrnăucă, C., Gómez-Pérez, D., Balcázar, J. L., and Montaña, J. L. Global optimality in k-means clustering. *Information Sciences*, 439:79–94, 2018.
- Ustun, B. T. B. Simple linear classifiers via discrete optimization: learning certifiably optimal scoring systems for decision-making and risk assessment. PhD thesis, Massachusetts Institute of Technology, 2017.
- Van der Laan, M., Pollard, K., and Bryan, J. A new partitioning around medoids algorithm. *Journal of Statistical Computation and Simulation*, 73(8):575–584, 2003.
- Vinod, H. D. Integer programming and the theory of grouping. *Journal of the American Statistical Association*, 64 (326):506–519, 1969.
- Wang, E., Ballachay, R., Cai, G., Cao, Y., and Trajano, H. L. Predicting xylose yield from prehydrolysis of hardwoods: A machine learning approach. *Frontiers in Chemical Engineering*, 4:994428, 2022.