

# RK-MEANS: FAST CLUSTERING FOR RELATIONAL DATA

RYAN CURTIN, BEN MOSELEY, HUNG Q. NGO, XUANLONG NGUYEN, DAN OLTEANU,  
AND MAXIMILIAN SCHLEICH

**ABSTRACT.** Conventional machine learning algorithms cannot be applied until a data matrix is available to process. When the data matrix needs to be obtained from a relational database via a feature extraction query, the computation cost can be prohibitive, as the data matrix may be (much) larger than the total input relation size. This paper introduces Rk-means, or relational  $k$ -means algorithm, for clustering relational data tuples without having to access the full data matrix. As such, we avoid having to run the expensive feature extraction query and storing its output. Our algorithm leverages the underlying structures in relational data. It involves construction of a small *grid coreset* of the data matrix for subsequent cluster construction. This gives a constant approximation for the  $k$ -means objective, while having asymptotic runtime improvements over standard approaches of first running the database query and then clustering. Empirical results show orders-of-magnitude speedup, and Rk-means can run faster on the database than even just computing the data matrix.

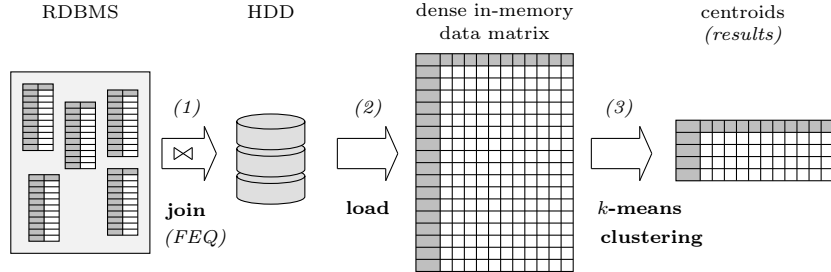
## 1. INTRODUCTION

Clustering is an ubiquitous technique for exploratory data analysis, whether applied to small samples or industrial-scale data. In the latter setting, two steps are typically performed: (1) *data preparation*, or extract-transform-load (ETL) operations, and (2) *clustering the extracted data*—often with a technique such as the popular  $k$ -means algorithm [12, 44]. In this setting, data typically reside in a relational database, requiring a *feature extraction query* (FEQ) to be performed, *joining* involved relations together to form the data matrix: each row corresponds to a data tuple and each column a feature. Then, the data matrix is used as input to a clustering algorithm. Such data matrices can be expensive to compute, and may take up space asymptotically larger than the database itself, which is made of relational tables. Moreover, the join computation time may exceed the time it takes to obtain clusters. It is not uncommon that the exploratory trip into the dataset may be stopped right at the gate.

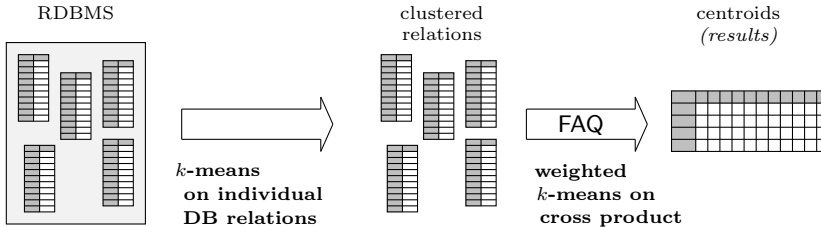
As an example, consider a retailer database consisting of three tables: `product`, which contains data about  $p$  products, `store`, which contains data about  $s$  stores, and `transaction`, which contains the number of transactions for each (product, store) combination on a given day. The table `product` contains information about each of the  $p$  products, `stores` contains information about each of the  $s$  stores, and `transactions` contains the (nonzero) number of transactions for each (product, store) combination on a given day. A practitioner may want to cluster each (product, store) combination as part of an analysis to determine items with related sales patterns across different stores for a given week. To do this, she constructs a data matrix containing all (product, store) combinations (including those with zero sales) for a given week, and additional attributes for each product and store. This can be achieved, for instance, by the following feature extraction query, given in SQLite syntax:

```
SELECT P.id AS i, S.id AS s, P.type AS t, P.price AS p,
       S.yelp_rating AS y, sum(ifnull(T.count, 0)) AS c
FROM product P, store S LEFT JOIN transactions T
ON T.product_id == P.id AND T.store_id == S.id
   AND T.date BETWEEN '2019-05-13' AND '2019-05-20'
GROUP BY P.id, S.id;
```

The result of this query is of size  $\Theta(ps)$ . But the `transaction` table can be significantly smaller than this, since many stores may have zero sales of a particular product in a given week. Thus, the size of the data matrix can be asymptotically greater than the total input relations' sizes. Real-world FEQs possess a similar explosion in both space and time complexity, only at a much larger scale, since they generally involve many more aggregations and tables. In Section 5, we present a real dataset from a large US retailer. The database



(a) Typical  $k$ -means data science workflow. Alternate representations can be used for the data in step (2) for greater computational efficiency (e.g., streaming); and, approximation strategies are known for step (3). However, the dataset often comes from an underlying database system, and in this case the expensive FEQ join (1) is unavoidable.



(b) The Rk-means data science workflow. We avoid ever computing the expensive FEQ by instead clustering each underlying relation (steps 1 and 2, Section 4); we then use FAQs for efficient weighted  $k$ -means of the cross-product of those relations (steps 3 and 4, Section 4). This gives significant empirical and theoretical accelerations, and bounded approximation—*without ever computing the full data matrix*.

FIGURE 1. Conventional  $k$ -means and Rk-means.

has 6 tables of total size 1.5GB. The FEQ result, however, takes up 18GB, and constructing it takes longer than running a learning algorithm on it.

Stripping away the language of databases, a fundamental challenge is how to learn about the joint distribution of a data population given only marginal samples revealed by relational tables. This is possible when the objective function of an underlying model admits some factorization structure similar to conditional independence in graphical models [28]. This insight was exploited recently by database theorists to devise algorithms evaluating a generic class of relational queries called *functional aggregate queries*, or FAQs [4]. The ability to answer FAQs quickly is a building block for a new class of efficient algorithms for training supervised learning models over relational data, *without* having to materialize (i.e. compute) the entire data matrix [3, 2, 37].

The goal of this paper is to devise a method for fast clustering of relational data, without having to materialize the full data matrix. The challenge of unsupervised learning tasks in general and the  $k$ -means algorithm in particular is that the learning objective is not decomposable across marginal samples in relational tables. To enable fast relational computation, we utilize the idea of constructing a *grid coreset*—a small set of points that provide a good summarization of the original (and unmaterialized) data tuples, based on which a provably good clustering can be obtained.

The resulting algorithm, which we call Rk-means, has several remarkable properties. First, Rk-means has a provable constant approximation guarantee relative to the  $k$ -means objective, despite the fact that the algorithm does not require access to the full data matrix. Our approximation analysis is established via a connection of Rk-means to the theory of optimal transport [41]. Second, Rk-means is enhanced by leveraging structures prevalent in relational data: categorical variables, functional dependencies, and the topology of feature extraction queries. These structures lead to exponential reduction in coreset size without incurring loss in the coreset’s approximation error. We show that Rk-means is provably more efficient both in time and space complexity when comparing against the application of the vanilla  $k$ -means to the full data matrix.

Finally, experimental results show significant speedups with little loss in the  $k$ -means objective. We observe orders-of-magnitude improvement in the running time compared to traditional methods. Rk-means is able to operate when other approaches would run out of memory, enabling clustering on truly massive datasets.

## 2. BACKGROUND AND RELATED WORK

**2.1. Background on Database Queries and FAQs.** Recent advancements in the database community have produced new classes of query plans and join algorithms [5, 4, 32, 33] for the efficient evaluation of general database queries. These general algorithms hinge on the expression of a database query as a *functional aggregate query*, or FAQ [4].

Loosely speaking, an FAQ is a collection of *aggregations* (be they sum, max, min, etc.) over a number of functions known as *factors*<sup>1</sup>, in the same sense as that used in graphical models. In particular, if there was only one aggregation (such as sum), then an FAQ is just a sum-product form typically used to compute the partition function. An FAQ is more general as it can involve many marginalization operators, one for each variable, and they can interleave in arbitrary way. Every relational database query can be expressed in this way. Consider the example query of Section 1: for this, the task of the database query evaluator is to compute  $\max(\text{transactions.count})$  for every tuple  $(i, s, t, p, y)$  that exists in the output. We can express this as a function:

$$(1) \quad \phi(i, s, t, p, y) = \max_c \max_i \max_s \psi_P(i, t, p) \psi_T(i, s, c) \psi_S(s, y).$$

In this we have three *factors*  $\psi_P(\cdot)$ ,  $\psi_T(\cdot)$ , and  $\psi_S(\cdot)$ , which correspond to the `product`, `transactions`, and `store` tables, respectively. We define  $\psi_P(i, t, p) = 1$  if the tuple  $(i, t, p)$  exists in the `product` table and 0 otherwise; we define  $\psi_S(s, y)$  similarly. We define  $\psi_T(i, s, c) = c$  if the tuple  $(i, s, c)$  exists in the `transactions` table and 0 otherwise. Thus, given any tuple  $(i, s, t, p, y)$ , we can compute  $\max(\text{transactions.count}) = \phi(i, s, t, p, y)$ .

In order to efficiently solve an FAQ (of which Equation (1) is but one example), the `InsideOut` algorithm of [4] may be used; `InsideOut` is a variable elimination algorithm, inspired from variable elimination in graphical model, with several new twists. One twist is to adapt worst-case optimal join algorithms [ ] to speed up computations by exploiting sparsity in the data. Another twist is that the algorithm has to carefully pick a variable order to minimize the runtime, while at the same time respect the correctness and semantic of the query. Unlike in the case of computing a sum-product where the summation operators are commutative, in a FAQ the operators may not be commutative.

To characterize the runtime of this algorithm, we must first observe that each database query and thus FAQ corresponds to a hypergraph  $\mathcal{H} = \{\mathcal{V}, \mathcal{E}\}$ . The vertices  $\mathcal{V}$  of this hypergraph correspond to the variables of the FAQ expression; in our example, we have  $\mathcal{V} = \{i, s, t, p, y, c\}$ . The hyperedges  $\mathcal{E}$ , then, correspond to each factor  $\psi_P(\cdot)$ ,  $\psi_T(\cdot)$ , and  $\psi_S(\cdot)$ —which in turn correspond to the tables in the database. This hypergraph  $\mathcal{H}$  is shown in Figure 2.1.

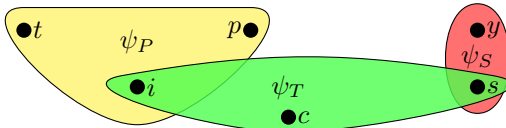


FIGURE 2. Example hypergraph  $\mathcal{H}$  for the example query and FAQ in Equation 1.

Roughly, `InsideOut` proceeds by first selecting a variable ordering  $\sigma$ , reordering the FAQ accordingly, and then solving the inner subproblems repeatedly, in much the same way variable elimination works for inference in graphical models [28]. The runtime of `InsideOut` is dependent on a notion of width of  $\mathcal{H}$  called FAQ-width, or  $\text{faqw}(\cdot)$ . Fully describing this width is beyond the scope of this paper and we encourage readers to refer to [4] for full details. The FAQ-width is a generalized version of *fractional hypertree width* of [19] (denoted by  $\text{fhtw}$ ). When the FAQ query does not have free variables,  $\text{faqw} = \text{fhtw}$ . Given some FAQ with hypergraph  $\mathcal{H}$ , via Section 4.3.4 of [5], `InsideOut` runs in time  $\tilde{O}(N^{\text{faqw}_{\mathcal{H}}(\sigma)} + Z)$ , where we assume that the support of each

<sup>1</sup>A full formal definition of FAQs can be found in [4], but is not required for our work here so we omit it.

factor<sup>2</sup> is no more than  $O(N)$ , and  $Z$  is the number of tuples in the output. As an example, the hypergraph of Figure 2.1 has  $\text{faqw}_{\mathcal{H}}(\sigma) = 1$ . Overall, `InsideOut` gives us the most efficient known way to evaluate problems that can be formulated as FAQs.

**2.2. Coresets for clustering.** From early work on  $k$ -means algorithm [29], ideas emerged for acceleration via coresets [21, 8]. Coresets have become the cornerstone of modern streaming algorithms [20, 11], massively parallel (MPC) algorithms [15, 9], and are used to speed up sequential methods [31, 38].

Unfortunately, existing algorithms for coreset construction do not readily lend themselves to the relational setting; there are several hurdles. First of all, coresets are formed by constructing a set  $S$  of data points (tuples) that represent the entire data set  $X$  well. Typically,  $S$  is a *weighted* representation of the data, where each point in the universe contributes one unit of weight to its closest point in  $X$  [22, 15, 10]. In our relational setting,  $X$  can only be formed by computing the FEQ, but our goal is to avoid materializing  $X$ .

A common challenge for adapting existing coreset constructions given our goal is that most methods construct  $S$  in phases by determining the farthest points from  $S$  [40, 7, 22]. This is difficult without  $X$  fully materialized. Another difficulty is that, even if the points in  $S$  are given, weighting the points in  $S$  is an open problem for relational algorithms [27]. Without  $X$  materialized, again, the points and their attributes are stored across several tables. Fixing a point  $x \in S$  and finding the number of closest points in (unmaterialized)  $X$  is non-trivial. No method, either deterministic or stochastic (e.g., sampling), is known that runs in time asymptotically faster than computing/materializing  $X$ . Our method avoids this by constructing a *grid coreset*  $S$  which can be decomposed over the tables in such a way that computing the weights of the points is a straightforward task.

**2.3. Other Related Work.** Our work draws inspiration from three lines existing work and ideas: coresets for clustering (discussed above), relational algorithms, and optimal transport. Some previous work has focused on the connection to databases; database and disk hardware optimizations have been considered to improve clustering relational data [34, 35]. Recent advances include the work of [3, 37].  $k$ -means has also been connected to optimal transport, which goes back to at least [36] (see also [18]). Recently this connection has received increased interest in the statistics and machine learning communities, resulting in fresh new clustering techniques [14, 24, 45]. To our knowledge, these related lines of work have not been explored together. Motivated by clustering relational data, our attempt at solving a clustering problem formulated as optimal transport in the marginal (projected) spaces to scalably perform  $k$ -means clustering appears to be the first in the literature.

Finally, it is worth noting that despite its popularity, the basic  $k$ -means technique is not always a preferred choice in clustering categorical or high-dimensional data. One may either adopt other clustering techniques [23, 26, 16], or modify the basic  $k$ -means method, e.g., by suitably placing weights on different features of mixed data types and replacing metric  $\ell_2$  by  $\ell_p$  [25], or incorporating a regularizer to combat high dimensionality [43, 39]. As we shall see, the relational techniques and associated theory that we introduce for the basic  $k$ -means extend easily to such improvements.

### 3. RK-MEANS, CORESETS AND OPTIMAL TRANSPORT

Although the Rk-means algorithm is motivated by application to relational databases, its basic idea is also of independent interest and can be easily described without the database language.

First we define the *weighted  $k$ -means* problem, which Rk-means solves (weights are also handy in combining mixed data types [25]). Let  $\mathbf{X}$  be a set of points in  $\mathbb{R}^d$ , and  $\mathbf{Y}$  be a non-empty set of points in the same space. Let  $d(\mathbf{x}, \mathbf{Y}) := \min_{\mathbf{y} \in \mathbf{Y}} \|\mathbf{x} - \mathbf{y}\|$  denote the minimum distance from  $\mathbf{x}$  to an element in  $\mathbf{Y}$ . In some cases, the  $\ell_2$  norm  $\|\cdot\|$  may be replaced by the  $\ell_p$  norm  $\|\cdot\|_p$  for some  $p \geq 1$ . A *weighted  $k$ -means instance* is a pair  $(\mathbf{X}, w)$ , where  $\mathbf{X}$  is a set of points in  $\mathbb{R}^d$  and  $w : \mathbf{X} \rightarrow \mathbb{R}^+$  is a weight function. Without loss of generality, assume  $\sum_{\mathbf{x} \in \mathbf{X}} w(\mathbf{x}) = 1$ . The task is to find a set  $\mathbf{C} = \{\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_k\}$  of  $k$  centroids to minimize the objective  $L(\mathbf{X}, \mathbf{C}, w) = \sum_{\mathbf{x} \in \mathbf{X}} w(\mathbf{x})d(\mathbf{x}, \mathbf{C})^2$ .

That is, we want to solve the problem  $\text{OPT}(\mathbf{X}, w) := \min_{\mathbf{C}} L(\mathbf{X}, \mathbf{C}, w)$ . With Rk-means, we will do this by projecting  $\mathbf{X}$  onto different sets of coordinates, and clustering each projection individually. To this end, let  $[d] = S_1 \cup \dots \cup S_m$  denote an arbitrary *partition* of the dimensions  $[d]$  into non-empty subsets. For every

<sup>2</sup>Or in our case, the number of tuples in the table corresponding to that factor.

---

**Algorithm 1 Rk-means:**  $k$ -means via grid-coreset
 

---

- 1: **Input:** query  $Q$ , number of clusters  $k$
  - 2: **Input:**  $[d] = S_1 \cup \dots \cup S_m$ ,  $\kappa \geq 2$
  - 3: **Output:** centroids  $C \in \mathcal{R}^{k \times d}$
  - 4: **for**  $j = 1$  to  $m$  **do**
  - 5:    $\mathbf{X}_j \leftarrow \{\mathbf{x}_{S_j} \mid \mathbf{x} \in \mathbf{X}\}$
  - 6:    $w_j \leftarrow$  weight function defined in (3)
  - 7:    $C_j \leftarrow \text{wkmeans}_1(\mathbf{X}_j, w_j, \kappa)$  {approx. ratio  $\alpha$ }
  - 8: **end for**
  - 9:  $\mathbf{G} \leftarrow C_1 \times \dots \times C_m$  {the grid coreset}
  - 10:  $w_{\text{grid}} \leftarrow$  weight function defined in (4)
  - 11:  $C \leftarrow \text{wkmeans}_2(\mathbf{G}, w_{\text{grid}}, k)$  {approx. ratio  $\gamma$ }
- 

$\mathbf{x} \in \mathbb{R}^d$ , and  $j \in [m]$ , let  $\mathbf{x}_{S_j}$  denote the projection of  $\mathbf{x}$  onto the coordinates in  $S_j$ . Define the projection set  $\mathbf{X}_j$  and corresponding weight function  $w_j : \mathbb{R}^{S_j} \rightarrow \mathbb{R}$  by

$$(2) \quad \mathbf{X}_j := \{\mathbf{x}_{S_j} \mid \mathbf{x} \in \mathbf{X}\},$$

$$(3) \quad w_j(\mathbf{z}) := \sum_{\mathbf{x} \in \mathbf{X} : \mathbf{x}_{S_j} = \mathbf{z}} w(\mathbf{x}).$$

In words, the  $w_j$  are the *marginal measures* of  $w$  on the subspace of coordinates  $S_j$ . With these notations established, Algorithm 1 presents the high-level description of our algorithm, Rk-means.

For each  $j \in [m]$ , in line 7 we perform  $k$ -means to obtain  $\kappa$  individual clusters on each subspace  $S_j$  for some  $\kappa \geq 2$ . These are solved using some weighted  $k$ -means algorithm denoted by  $\text{wkmeans}_1$  with approximation ratio  $\alpha$ . Then, using the results of these clusterings, we assemble a cross-product weighted grid  $\mathbf{G}$  of centroids, and then perform  $k$ -means clustering on these using the algorithm denoted  $\text{wkmeans}_2$  to reduce down to the desired result of  $k$  centroids. Typically, take  $\kappa = O(k)$ .

Let  $\mathbf{X} := \bigsqcup_{\mathbf{g} \in \mathbf{G}} \mathbf{X}_{\mathbf{g}}$  denote a partition of  $\mathbf{X}$  into  $|\mathbf{G}|$  parts, where  $\mathbf{X}_{\mathbf{g}}$  denote the set of points in  $\mathbf{X}$  closer to  $\mathbf{g}$  than other grid points in  $\mathbf{G}$  (breaking ties arbitrarily). Then, the weight function for line 11 is  $w_{\text{grid}} : \mathbf{G} \rightarrow \mathbb{R}^+$  is defined as

$$(4) \quad w_{\text{grid}}(\mathbf{g}) := \sum_{\mathbf{x} \in \mathbf{X}_{\mathbf{g}}} w(\mathbf{x}).$$

Weighted  $k$ -means and optimal transport. We will analyze the Rk-means algorithm in the language of optimal transport. The connection of  $k$ -means in general, and of our algorithm to optimal transport in particular, provides another interesting insight into our algorithm.

The *optimal transport distance* characterizes the distance between two probability measures, by measuring the optimal cost of transporting mass from one to another [41]. Although this is defined more generally for any two probability measures in abstract spaces, for our purpose it is convenient to consider two discrete probability measures  $P$  and  $P'$  on  $\mathbb{R}^d$ .

Let  $\mathbf{Z}$  and  $\mathbf{Z}'$  be two finite point sets in  $\mathbb{R}^d$ . Let  $\delta$  denote the Dirac measure. Let  $P := \sum_{\mathbf{z} \in \mathbf{Z}} p(\mathbf{z})\delta_{\mathbf{z}}$  and  $P' := \sum_{\mathbf{z}' \in \mathbf{Z}'} p'(\mathbf{z}')\delta_{\mathbf{z}'}$  be two measures with supports  $\mathbf{Z}$  and  $\mathbf{Z}'$ , respectively. The mass transportation plan can be formalized by a *coupling*: a joint distribution  $\mathbf{Q} = (q(\mathbf{z}, \mathbf{z}'))_{(\mathbf{z}, \mathbf{z}') \in \mathbf{Z} \times \mathbf{Z}'}$ , where the marginal constraints  $\sum_{\mathbf{z} \in \mathbf{Z}} q_{\mathbf{z}, \mathbf{z}'} = p'(\mathbf{z}')$  and  $\sum_{\mathbf{z}' \in \mathbf{Z}'} q_{\mathbf{z}, \mathbf{z}'} = p(\mathbf{z})$  hold.

**Definition 3.1.** For any  $p \geq 1$ , the *Wasserstein distance* of order  $p$  is defined by the minimization of  $\mathbf{Q}$  over all possible couplings:  $W_p(P, P') = \min_{\mathbf{Q}} \{\sum_{\mathbf{z}, \mathbf{z}'} q(\mathbf{z}, \mathbf{z}') \|\mathbf{z} - \mathbf{z}'\|_p^p\}^{1/p}$ .

Let  $P^{\text{in}} = \sum_{\mathbf{x} \in \mathbf{X}} w(\mathbf{x})\delta_{\mathbf{x}}$  be the discrete measure associated with the input instance of our weighted  $k$ -means problem; then, this can be expressed precisely as an optimal transport problem:  $M^* = \arg \min W_2^2(M, P^{\text{in}})$ ,

where the optimization is over the space of discrete measures  $M$  that have  $k$  support points (the set  $\mathbf{C}$  of  $k$  centroids). Note that  $\text{OPT}(\mathbf{X}, w) = W_2^2(M^*, P^{\text{in}})$ . Replacing  $\ell_2$  by say  $\ell_1$ , we obtain the  $k$ -median problem, for which the objective becomes  $W_1(M^*, P^{\text{in}})$ .

**Approximation Analysis.** We next analyze the approximation ratio of our Rk-means algorithm working with the  $W_2^2$  objective, provided that  $\text{wkmeans}_1$  has approximation ratio  $\alpha$  and  $\text{wkmeans}_2$  has approximation ratio  $\gamma$ .<sup>3</sup> The reason we might want to invoke different algorithms to solve these sub-problems is because, as we shall show in the next section, we may want to exploit the (relational) structures of the FEQ to construct a “nice” partition  $S_1 \cup \dots \cup S_m$ . We show that the overall approximation ratio of Rk-means is  $(\sqrt{\alpha} + \sqrt{\gamma} + \sqrt{\alpha\gamma})^2$ . In many common cases, the relational database has structure that allows  $\alpha = 1$ , yielding an overall approximation ratio of  $(1 + 2\sqrt{\gamma})^2$ .

For our analysis it is useful to understand Algorithm 1 in the language of optimal transport. For any finite point set  $\mathbf{Y} \subset \mathbb{R}^d$  and a measure  $M = \sum_{\mathbf{y} \in \mathbf{Y}} p(\mathbf{y})\delta_{\mathbf{y}}$  with support  $\mathbf{Y}$ , define the marginal measures  $M_j$  on coordinates  $S_j$  induced by  $M$  in the natural way, i.e.  $M_j := \sum_{\mathbf{z} \in \mathbf{Z}} p_j(\mathbf{z})\delta_{\mathbf{z}}$  where  $p_j$  is defined analogous to  $w_j$  in (3). Under this notation,  $P^{\text{in}}$  induces the marginal measures  $P_j^{\text{in}} := \sum_{\mathbf{z} \in \mathbf{X}_j} w_j(\mathbf{z})\delta_{\mathbf{z}}$ . Then, Algorithm 1 can be described by the following steps:

- (1) For each  $j \in [m]$ , pick  $M_j$  to be the ( $\alpha$ -approximate) minimizer of  $W_2^2(M_j, P_j^{\text{in}})$ , where  $\mathbf{C}_j = \text{supp}(M_j)$  is the support of  $M_j$  and  $|\mathbf{C}_j| = \kappa$  (line 7).
- (2) Collect the  $\kappa^d$  grid points  $\mathbf{G}$  and let probability measure  $Q$  be the one with support in  $\mathbf{G}$  such that  $Q$  minimizes  $W_2^2(Q, P^{\text{in}})$ . (We solve this problem exactly!)
- (3) Finally, return  $P$  which is the measure with exactly  $k$  support points in  $\mathbb{R}^d$  that ( $\gamma$ -approximately) minimizes  $W_2^2(P, Q)$  (line 11).

This is precisely the solution obtained by Algorithm 1. We present next some useful facts.

**Lemma 3.2.** *For any discrete measure  $M$  on  $\mathbb{R}^d$ ,  $W_2^2(M, P^{\text{in}}) \geq \sum_{j=1}^m W_2^2(M_j, P_j^{\text{in}})$ .*

*Proof.* A valid coupling of two measures induces valid marginal couplings of marginal measures. □

**Proposition 3.3.** *The following hold:*

- (a) *If  $\kappa \geq |\text{supp}(M_j^*)| \forall j \in [m]$ , then  $W_2(P^{\text{in}}, P) \leq (\sqrt{\gamma} + \sqrt{\alpha} + \sqrt{\alpha\gamma})W_2(P^{\text{in}}, M^*)$ .*
- (b) *For any  $\kappa \geq 1$ , there exists a distribution  $P^{\text{in}}$  such that*

$$(5) \quad \frac{W_2(P^{\text{in}}, P)}{W_2(P^{\text{in}}, M^*)} \geq \sqrt{1 - e^{-m/(2\kappa)}} \frac{\sqrt{3}k^{3/(2m)}}{2\kappa m^{1/2}}.$$

*Proof.* (a) By the definition of  $Q$ , the optimal transport plan from  $P^{\text{in}}$  to  $Q$  is such that each support point  $s \in S$  is received by all  $x \in \mathbf{X}$  nearest to  $s$  compared to other points in  $S$ . So,

$$(6) \quad W_2^2(P^{\text{in}}, Q) = \sum_{\mathbf{x} \in \mathbf{X}} w(\mathbf{x})d(\mathbf{x}, \mathbf{G})^2$$

$$(7) \quad = \sum_{\mathbf{x} \in \mathbf{X}} w(\mathbf{x}) \sum_{j=1}^m d(\pi_{S_j} \mathbf{x}, \mathbf{C}_j)^2$$

$$(8) \quad = \sum_{j=1}^m \sum_{\mathbf{z} \in \mathbf{X}_j} w_j(\mathbf{z})d(\mathbf{z}, \mathbf{C}_j)^2$$

$$(9) \quad = \sum_{j=1}^m W_2^2(M_j, P_j^{\text{in}})$$

$$(10) \quad \leq \alpha \sum_{j=1}^m W_2^2(M_j^*, P_j^{\text{in}})$$

$$(11) \quad \leq \alpha \cdot W_2^2(M^*, P^{\text{in}}).$$

---

<sup>3</sup>The best known approximation ratio is 6.357 for data in Euclidean space [6].

The second to last inequality is due to the  $\alpha$ -approximation of  $\text{wkmeans}_1$ , and condition that  $|\text{supp}(M_j)| \geq |\text{supp}(M_j^*)|$ . The last inequality follows from Proposition 3.2. By the triangle inequality of  $W_2$ ,

$$\begin{aligned}
(12) \quad & W_2(P^{\text{in}}, P) \leq W_2(P^{\text{in}}, Q) + W_2(Q, P) \\
(13) \quad & \leq W_2(P^{\text{in}}, Q) + \sqrt{\gamma} \cdot W_2(Q, M^*) \\
(14) \quad & \leq W_2(P^{\text{in}}, Q) + \sqrt{\gamma}(W_2(Q, P^{\text{in}}) + W_2(P^{\text{in}}, M^*)) \\
(15) \quad & = (1 + \sqrt{\gamma})W_2(P^{\text{in}}, Q) + \sqrt{\gamma} \cdot W_2(P^{\text{in}}, M^*) \\
(16) \quad & \leq (1 + \sqrt{\gamma})\sqrt{\alpha}W_2(P^{\text{in}}, M^*) + \sqrt{\gamma} \cdot W_2(P^{\text{in}}, M^*) \\
(17) \quad & = (\sqrt{\alpha} + \sqrt{\gamma} + \sqrt{\alpha\gamma}) \cdot W_2(M^*, P^{\text{in}}).
\end{aligned}$$

The second inequality is due to the fact that  $\text{wkmeans}_2$  has approximation ratio  $\gamma$ ; the first and third are the triangle inequality. We conclude the proof.

(b) We only need to construct an example of  $P^{\text{in}}$  for the case  $d = m$ . Although  $P^{\text{in}}$  as an input to the algorithm is a discrete measure, for the purpose of this proof it suffices to take  $P^{\text{in}}$  to be the uniform distribution on  $[0, 1]^m$  (which can be approximated arbitrarily well by a discrete measure). It is simple to verify that if  $k_0 = k^{1/m}$  is a natural number, then  $M^*$  is a uniform distribution on the regular grid of size  $k_0$  in each dimension. It follows that  $W_2^2(P^{\text{in}}, M^*) \leq \frac{m}{12k_0^3} = \frac{m}{12k^{3/m}}$ . The grid points  $\mathbf{G}$  range over the set  $S := [1/(2\kappa), 1 - 1/(2\kappa)]^m$ . Moreover,  $Q$  is a uniform distribution on  $\mathbf{G}$ . Now  $P$  is the outcome of line (11) so the support of  $P$  must lie in the convex hull  $S$  of  $\mathbf{G}$ . The cost of each unit mass transfer from an atom in the complement of set  $[1/(4\kappa), 1 - 1/(4\kappa)]^m$  to one in  $S$  is at least  $(1/4\kappa)^2$ , so  $W_2^2(P^{\text{in}}, P) \geq (1/4\kappa)^2 \cdot [1 - (1 - 1/(2\kappa))^m]$ . We note  $(1 - 1/(2\kappa))^m < e^{-m/2\kappa}$  to conclude the proof.  $\square$

The condition of part (a) is satisfied, for instance, by setting  $\kappa = k$ . In practice,  $\kappa < k$  may suffice. Moreover, part (b) dictates that  $\kappa$  must grow with  $k$  appropriately for our algorithm to maintain a constant approximation guarantee. Since solution  $\mathbf{C}$  has cost  $L(\mathbf{X}, \mathbf{C}, w) = W_2^2(P, P^{\text{in}})$ , and  $\text{OPT}(\mathbf{X}, w) = W_2^2(M^*, P^{\text{in}})$ , the following theorem is immediate from Prop. 3.3(a).

**Theorem 3.4.** *Suppose  $\text{wkmeans}_1$  and  $\text{wkmeans}_2$  have approximation ratios  $\alpha$  and  $\gamma$ . Then by choosing  $\kappa = k$ , the solution  $\mathbf{C}$  given by  $Rk$ -means has the following guarantee:  $L(\mathbf{X}, \mathbf{C}, w) \leq (\sqrt{\gamma} + \sqrt{\alpha} + \sqrt{\gamma\alpha})^2 \text{OPT}(\mathbf{X}, w)$ .*

*Specifically, if both sub-problems are solved optimally ( $\alpha = \gamma = 1$ ),  $Rk$ -means is a 9-approximation.*

Regularized  $Rk$ -means. It is possible to extend our approach to accommodate regularization techniques. This can be useful when the data are very high dimensional [39, 43]. Thus, the clustering formulation can be expressed as a regularized optimal transport problem:  $M^* = \arg \min W_2^2(M, P^{\text{in}}) + \Omega(M)$  where the optimization is over the space of discrete measures  $M$  that have  $k$  support points (the set  $\mathbf{C}$  of  $k$  centroids), and the regularizer  $\Omega(M) \geq 0$  typically decomposes over the  $m$ -partition of variables:  $\Omega(M) = \sum_{j=1}^m \Omega_j(M_j)$ . For instance,  $\Omega_j(M_j)$  may be taken to be a multiple of the  $\ell_1$  norm of  $M_j$ 's supporting atoms (e.g., group lasso penalty). The algorithm has the same three steps as before, with some modification in (1') and (3'):

- (1') For each  $j \in [m]$ , pick  $M_j$  to be the ( $\alpha$ -approximate) minimizer of  $W_2^2(M_j, P_j^{\text{in}}) + \Omega_j(M_j)$ , where  $\mathbf{C}_j = \text{supp}(M_j)$  is the support of  $M_j$  and  $|\mathbf{C}_j| = \kappa$  (line 7).
- (3') Finally, return  $P$  which is the measure with exactly  $k$  support points in  $\mathbb{R}^d$  that ( $\gamma$ -approximately) minimizes  $W_2^2(P, Q) + \Omega(P)$  (line 11).

**Proposition 3.5.** *If  $\kappa \geq |\text{supp}(M_j^*)|$  for all  $j \in [m]$ , then*

$$(18) \quad \frac{W_2^2(P^{\text{in}}, P) + \Omega(P)}{W_2^2(P^{\text{in}}, M^*) + \Omega(M^*)} \leq 2\alpha + 4\gamma + 4\alpha\gamma.$$

*Proof.* As before the optimal transport plan from  $P^{\text{in}}$  to  $Q$  is such that each support point  $s \in S$  is received by all  $x \in \mathbf{X}$  nearest to  $s$  compared to other points in  $S$ . So,

$$(19) \quad W_2^2(P^{\text{in}}, Q) + \Omega(M) = \sum_{j=1}^m W_2^2(M_j, P_j^{\text{in}}) + \Omega(M)$$

$$(20) \quad \leq \alpha \sum_{j=1}^m (\mathbb{W}_2^2(M_j^*, P_j^{\text{in}}) + \Omega_j(M_j^*))$$

$$(21) \quad \leq \alpha (\mathbb{W}_2^2(M^*, P^{\text{in}}) + \Omega(M^*)).$$

The second to last inequality is due to the  $\alpha$ -approximation of (regularized)  $\text{wkmeans}_1$ , and condition that  $|\text{supp}(M_j)| \geq |\text{supp}(M_j^*)|$ . The last inequality follows from Proposition 3.2 and the definition of  $\Omega$ . By the triangle inequality of  $\mathbb{W}_2$ , as before

$$(22) \quad \mathbb{W}_2(P^{\text{in}}, P) \leq \mathbb{W}_2(P^{\text{in}}, Q) + \mathbb{W}_2(Q, P)$$

$$(23) \quad \leq \mathbb{W}_2(P^{\text{in}}, Q) + \sqrt{\gamma \mathbb{W}_2^2(Q, M^*) + \gamma \Omega(M^*)} - \Omega(P)$$

$$(24) \quad \leq \mathbb{W}_2(P^{\text{in}}, Q)$$

$$(25) \quad + \sqrt{2\gamma \mathbb{W}_2^2(P^{\text{in}}, Q) + 2\gamma \mathbb{W}_2^2(P^{\text{in}}, M^*) + \gamma \Omega(M^*)} - \Omega(P).$$

Hence, by Cauchy-Schwarz and combining with (21) we obtain

$$(26) \quad \mathbb{W}_2^2(P^{\text{in}}, P) \leq 2 \left\{ (1 + 2\gamma) \mathbb{W}_2^2(P^{\text{in}}, Q) + 2\gamma \mathbb{W}_2^2(P^{\text{in}}, M^*) + \gamma \Omega(M^*) - \Omega(P) \right\}$$

$$(27) \quad \leq (2\alpha + 4\gamma + 4\alpha\gamma) \mathbb{W}_2^2(P^{\text{in}}, M^*) + (2\alpha + 2\gamma + 4\alpha\gamma) \Omega(M^*) - (2 + 4\gamma) \Omega(M) - 2\Omega(P).$$

The conclusion is immediate by noting that  $\Omega$  is a non-negative function.  $\square$

If both subproblems for regularized  $k$ -means can be solved optimally, our method yields a 10-approximation on the penalized  $\mathbb{W}_2^2$  objective. We conclude by noting that our technique extends easily to the  $\mathbb{W}_p^p$  objective for any  $p \geq 1$ , but the approximation ratio will be changed according to  $p$ .

#### 4. LEVERAGING STRUCTURES IN RELATIONAL DATA

We now explain the “relational” part of the Rk-means algorithm, where we exploit relational structures in the data and the FEQ to achieve significant computational savings. Three classes of relational structures prevalent in RDBMSs are (a) *categorical variables*, (b) *functional dependencies* (FDs), and (c) the topology of the FEQ. We exploit these structures to carefully select the partition  $S_1 \cup \dots \cup S_m$  to use for Rk-means, to compute the marginal sub-problems  $(\mathbf{X}_j, w_j)$ , the components  $\mathbf{C}_j$  of the coreset  $\mathbf{G}$ , and the grid weight  $w_{\text{grid}}$  *without* materializing the entire coreset  $\mathbf{G}$ . When selecting partitions, there are two competing criteria: first, we need a partition so that the approximation ratio  $\alpha$  for  $\text{wkmeans}_1$  is as small as possible. For example, if  $|S_j| = 1$  for all  $j$ , so  $m = d$ , then we can apply the well-known optimal solution for  $k$ -means in 1 dimension using dynamic programming in  $O(n^2k)$  time [42]; this then provides  $\alpha = 1$ . On the other hand, we want the remainder of algorithm to be fast by keeping the size of the grid  $\mathbf{G}$ , namely  $|\mathbf{G}| \leq \kappa^m$ , small.

**4.1. Categorical variables.** Real-world relational database queries typically involve many categorical variables (e.g., color, month, or city). In practice, practitioners may endow non-uniform weights for different categorical variables, or categories [25]. In terms of representation, the most common way to deal with categorical variables is to one-hot encode them, whereby a categorical feature such as city is represented by an indicator vector

$$(28) \quad \mathbf{x}_{\text{city}} = [\mathbf{1}_{\text{city}=c_1} \quad \mathbf{1}_{\text{city}=c_2} \quad \dots \quad \mathbf{1}_{\text{city}=c_L}]$$

where  $\{c_1, \dots, c_L\}$  is the set of cities occurring in the data. The subspace associated with these indicator vectors is known as the *categorical subspace* of a categorical variable. This one-hot representation substantially increases the data matrix size via an increase in the *dimensionality* of the data. For example, a dataset of about 30 mostly categorical features with hundreds or thousands of categories for each feature will have its dimensionality exploded to the order of thousands with one-hot encoding. Fortunately, this is not a problem — by treating each categorical variable as a subset of the partition, the *weighted  $k$ -means* subproblem within a categorical subspace is solvable efficiently and optimally.



This optimal solution can be computed in the same time it takes to find the number of points in each category, which is a vast improvement on either an optimal dynamic program or Lloyd's algorithm. Furthermore, it helps keep  $m$  as low as the number of database attributes in the query.

Consider a weighted  $k$ -means subproblem solved by  $wkmeans_1$  defined on a categorical subspace induced by a categorical feature  $K$  that has  $L$  categories. Then, the instance is of the form  $(\mathbf{I}, v)$ , where  $\mathbf{I}$  is the collection of  $L$  indicator vectors  $\mathbf{1}_e$ , one for each element  $e \in \text{Dom}(K)$ . (One can think of  $\mathbf{I}$  as the identity matrix of order  $L$ .) Define the weight function  $v$  as

$$(29) \quad v(\mathbf{1}_e) = \sum_{\mathbf{x} \in \mathbf{X}, \mathbf{x}_K = e} w(\mathbf{x}).$$

For any set  $F \subseteq \text{Dom}(K)$ , let  $\mathbf{v}_F$  denote the vector  $(v(\mathbf{1}_e))_{e \in F}$ . Also,  $\|\mathbf{v}_F\|_1$  and  $\|\mathbf{v}_F\|_2$  denote the  $\ell_1$  and  $\ell_2$  norm, respectively. It is useful to rewrite the categorical weighted  $k$ -means problem:

**Proposition 4.1.** *The categorical weighted  $k$ -means instance  $(\mathbf{I}, v)$  admits the following optimization objective:*

$$(30) \quad \text{OPT}(\mathbf{I}, v) = \|\mathbf{v}\|_1 - \max_{\mathcal{F}} \sum_{F \in \mathcal{F}} \frac{\|\mathbf{v}_F\|_2^2}{\|\mathbf{v}_F\|_1},$$

where  $\mathcal{F}$  ranges over all partitions of  $\text{Dom}(K)$  into  $k$  parts.

*Proof.* First, consider a subset  $F \subseteq \text{Dom}(K)$  of the categories; the centroid  $\boldsymbol{\mu}$  of (weighted) indicator vectors  $\mathbf{1}_e$ ,  $e \in F$ , can be written down explicitly:

$$(31) \quad \boldsymbol{\mu}_e = \begin{cases} 0 & e \notin F \\ \frac{v_e}{\|\mathbf{v}_F\|_1} & e \in F, \end{cases}$$

The weighted sum of squared distances between  $\mathbf{1}_e$  for all  $e \in F$  to  $\boldsymbol{\mu}$  is

$$\begin{aligned} \sum_{e \in F} (\|\boldsymbol{\mu}\|_2^2 - \mu_e^2 + (\mu_e - 1)^2)v_e &= \frac{\|\mathbf{v}_F\|_2^2}{\|\mathbf{v}_F\|_1} + \sum_{e \in F} ((\mu_e - 1)^2 - \mu_e^2)v_e \\ &= \frac{\|\mathbf{v}_F\|_2^2}{\|\mathbf{v}_F\|_1} + \sum_{e \in F} (-2\mu_e + 1)v_e \\ &= \|\mathbf{v}_F\|_1 - \|\mathbf{v}_F\|_2^2 / \|\mathbf{v}_F\|_1. \end{aligned}$$

Thus, the weighted  $k$ -means objective takes the form

$$(32) \quad \min_{\mathcal{F}} \sum_{F \in \mathcal{F}} \left( \|\mathbf{v}_F\|_1 - \|\mathbf{v}_F\|_2^2 / \|\mathbf{v}_F\|_1 \right) = \|\mathbf{v}\|_1 - \max_{\mathcal{F}} \sum_{F \in \mathcal{F}} \|\mathbf{v}_F\|_2^2 / \|\mathbf{v}_F\|_1,$$

which concludes the proof.  $\square$

In (30), note that  $\|\mathbf{v}\|_1$  is the total weight of input points; hence, we can equivalently solve the inner maximization problem. With the categorical weighted  $k$ -means objective in place, we can derive the optimal clustering. To do so, We next need the following elementary lemma.

**Lemma 4.2.** *Suppose that  $x, a_1, a_2, b_1, b_2 > 0$ ,  $b_1^2 \geq a_1$ ,  $b_2^2 \geq a_2$  and  $x \geq \max\{a_1/b_1, a_2/b_2\}$ . Then  $x + \frac{a_1+a_2}{b_1+b_2} \geq \max\left\{\frac{x^2+a_1}{x+b_1} + \frac{a_2}{b_2}, \frac{x^2+a_2}{x+b_2} + \frac{a_1}{b_1}\right\}$ .*

*Proof.* It suffices to establish  $x + \frac{a_1+a_2}{b_1+b_2} \geq \frac{x^2+a_1}{x+b_1} + \frac{a_2}{b_2}$ , or equivalently

$$x - \frac{x^2 + a_1}{x + b_1} \geq \frac{a_2}{b_2} - \frac{a_1 + a_2}{b_1 + b_2},$$

which can be simplified as

$$(33) \quad x(b_1 + b_2 + a_1/b_1 - a_2/b_2) \geq a_1b_2/b_1 + a_2b_1/b_2.$$

To verify this inequality, consider two cases. If  $a_1/b_1 \geq a_2/b_2$ , then  $LHS \geq x(b_1 + b_2) \geq (a_2/b_2)b_1 + (a_1/b_1)b_2$ . On the other hand, if  $a_2/b_2 > a_1/b_1$ . Since  $b_2 - a_2/b_2 \geq 0$ ,

$$\begin{aligned} LHS &\geq (a_2/b_2)(b_1 + b_2 + a_1/b_1 - a_2/b_2) \\ &= a_2b_1/b_2 + a_2 + a_1a_2/(b_1b_2) - a_2^2/b_2^2 \\ &= a_2b_1/b_2 + a_1b_2/b_1 + (b_2 - a_2/b_2)(a_2/b_2 - a_1/b_1) \\ &\geq a_2b_1/b_2 + a_1b_2/b_1. \end{aligned}$$

Thus the proof is complete.  $\square$

Then, the optimal solution to the categorical  $k$ -means instance is an immediate consequence:

**Corollary 4.3.** *Let  $(e_1, \dots, e_L)$  be a permutation of  $\text{Dom}(K)$  such that  $v_{e_1} \geq v_{e_2} \geq \dots \geq v_{e_L}$ . Then for any  $k \geq 2$  and any  $k$ -partition  $\mathcal{F}$  of  $\text{Dom}(K)$ , there holds*

$$v_{e_1} + \dots + v_{e_{k-1}} + \frac{\sum_{i=k}^L v_i^2}{\sum_{i=k}^L v_i} \geq \sum_{F \in \mathcal{F}} \frac{\|\mathbf{v}_F\|_2^2}{\|\mathbf{v}_F\|_1}.$$

*Proof.* We prove the claim by induction on  $k$ . Let  $F \in \mathcal{F}$  be the set containing the element  $\{e_1\}$ . If there is only one element in  $F$  then we apply the induction hypothesis on the remaining terms. Otherwise,  $F$  contains at least two elements. Let  $G \in \mathcal{F}$  be an arbitrary element of  $\mathcal{F}$  where  $G \neq F$ . Define  $\mathcal{F}'$  to be the partition obtained from  $\mathcal{F}$  by replacing  $(F, G)$  with  $(\{e_1\}, F \cup G - \{e_1\})$ . Then, Lemma 4.2 can be applied to get

$$\sum_{F \in \mathcal{F}} \frac{\|\mathbf{v}_F\|_2^2}{\|\mathbf{v}_F\|_1} \leq \sum_{F \in \mathcal{F}'} \frac{\|\mathbf{v}_F\|_2^2}{\|\mathbf{v}_F\|_1}.$$

Induction on the tail  $k - 1$  terms completes the proof.  $\square$

Theorem 4.4 below follows trivially from the above corollary. Corollary 4.3 and the objective for  $k$ -means on a single attribute in the equation of Proposition 4.1 establishes precisely the structure of the optimal solution for data consisting of a single categorical variable.

**Theorem 4.4.** *Given a categorical weighted  $k$ -means instance, an optimal solution can be obtained by putting each of the first  $k - 1$  highest weight indicator vectors in its own cluster, and the remaining vectors in the same cluster.*

This means that for a categorical variable with  $L$  categories, we can compute the optimal clustering for the sub-problem in only  $O(nL \log L)$  time. The variable gives rise to a categorical subspace of size  $|S_j| = L$ .

**4.2. Functional dependencies.** Next, we address the second call to `wkmeans2`: its runtime is dependent on the size of the grid  $\mathbf{G}$ , which can be up to  $O(k^m)$ , where  $m$  is the number of features from the input. Databases often contain *functional dependencies* (FDs), which we can exploit to reduce the size of  $\mathbf{G}$ . An FD is a dimension whose value depends entirely on the value of another dimension. For example, for a retailer dataset that includes geographic information, one might encounter features such as `storeID`, `zip`, `city`, `state`, and `country`. Here, `storeID` functionally determines `zip`, which determines `city`, which in turns determines `state`, leading to `country`. This common structure is known as an *FD-chain*, and appears often in real-world FEQs.

If we were to apply `Rk-means` without exploiting the FDs, the features `storeID`, `zip`, `city`, `state`, and `country` would contribute a factor of  $k^5$  to the grid size. However, by using the FD structure of the database, we show that only a factor of  $5k$  is contributed to the grid size, because most of the  $k^5$  grid points  $\mathbf{g}$  have  $w_{\text{grid}}(\mathbf{g}) = 0$  (see (4)). More generally, whenever there is an *FD chain* including  $p$  features, their overall contribution to the grid size is a factor of  $O(kp)$  instead of  $O(k^p)$ , and the grid points with non-zero weights can be computed efficiently in time  $O(kp)$ .

**Lemma 4.5.** *Suppose all  $d$  input features are categorical and form an FD-chain. Then, the total number of grid points  $\mathbf{g} \in \mathbf{G}$  with non-zero  $w_{\text{grid}}$  weight is at most  $d(k - 1) + 1$ .*

*Proof.* Suppose the features are  $K_1, \dots, K_d$ , where  $K_i$  functionally determine  $K_{i+1}$ , and  $\text{Dom}(K_i) = \{e_1^i, e_2^i, \dots, e_{n_i}^i\}$ . Without loss of generality, we also assume that the elements in  $\text{Dom}(K_i)$  are sorted in descending order of weights:

$$(34) \quad w(\mathbf{1}_{e_1^i}) \geq w(\mathbf{1}_{e_2^i}) \geq \dots \geq w(\mathbf{1}_{e_{n_i}^i}).$$

From Corollary 4.3, we know the set  $\mathbf{C}_i$  of  $k$  centroids of each of the categorical subspace for  $K_i$ : there is a centroid  $\mu_j^i = \mathbf{1}_{e_j^i}$  for each  $j \in [k-1]$ , and then a centroid  $\mu_k^i$  of the rest of the indicator vectors. The elements  $e_j^i$  for  $j \in [k-1]$  shall be called “heavy” elements, and the rest are “light” elements.

Now, consider an input vector  $\mathbf{x} = (x_1, \dots, x_d)$  where  $x_i \in \text{Dom}(K_i)$ . Under one-hot encoding, this vector is mapped to a vector of indicator vectors  $\mathbf{1}_{\mathbf{x}} := (\mathbf{1}_{x_1}, \dots, \mathbf{1}_{x_d})$ . We need to answer the question: which grid point in  $\mathbf{G} = \mathbf{C}_1 \times \dots \times \mathbf{C}_d$  is  $\mathbf{1}_{\mathbf{x}}$  closest to? Since the  $\ell_2^2$ -distance is decomposable into component sum, we can determine the closest grid point by looking at the closest centroid in  $\mathbf{C}_i$  for  $\mathbf{1}_{x_i}$ , for each  $i \in [d]$ .

If  $x_i \in \{e_1^i, \dots, e_{k-1}^i\}$ , then the corresponding one-hot-encoded version  $\mathbf{1}_{x_i}$  is itself one of the centroids in  $\mathbf{C}_i$ , and thus it is its own closest centroid. Otherwise, the closest centroid to  $\mathbf{1}_{x_i}$  is  $\mu_k^i$ , because  $\|\mathbf{1}_{x_i} - \mu_k^i\|^2 < 2$ , and  $\|\mathbf{1}_{x_i} - \mu_j^i\|^2 = 2$  for every  $j \in [k-1]$ .

Let  $\mu^i(x_i) \in \mathbf{C}_i$  denote the closest centroid in  $\mathbf{C}_i$  to  $\mathbf{1}_{x_i}$ . The closest grid point to  $\mathbf{1}_{\mathbf{x}}$  is completely determined:  $\mathbf{g} = (\mu^1(x_1), \dots, \mu^d(x_d))$ . Furthermore, let  $i \in [d]$  denote the smallest index such that  $x_i$  is heavy. Then, we can write  $\mathbf{g}$  as

$$(35) \quad \mathbf{g} = (\mu_k^1, \dots, \mu_k^{i-1}, \mathbf{1}_{x_i}, \mu^{i+1}(x_{i+1}), \dots, \mu^d(x_d))$$

Note that once  $x_i$  is fixed, due to the FD-chain the *entire* suffix  $(\mathbf{1}_{x_i}, \mu^{i+1}(x_{i+1}), \dots, \mu^d(x_d))$  of  $\mathbf{g}$  is determined. Hence, the number of different  $\mathbf{g}$ s can only be at most  $d(k-1) + 1$ : there are  $d+1$  choices for  $i$  (from 0 to  $d$ ), and  $k-1$  choices for  $x_i$  if  $i > 0$ .  $\square$

Theorem 4.6 below follows trivially from the above lemma, because the  $\ell_2^2$ -distance is the sum over the  $\ell_2^2$ -distances of the subspaces.

**Theorem 4.6.** *Suppose all  $d$  input features can be partitioned into  $m$  FD-chains of size  $d_1, \dots, d_m$ , respectively. Then, the number of grid points  $\mathbf{g} \in \mathbf{G}$  with non-zero  $w_{\text{grid}}$  weight is bounded by  $\prod_{i=1}^m (1 + d_i(k-1))$ . Furthermore, the set of non-zero weight grid points can be computed in time  $\tilde{O}(\prod_{i=1}^m (1 + d_i(k-1)))$ .*

Note that in the above theorem, if there was *no* FD, then  $d$  features each form their own chain of size 1, in which case  $\prod_{i=1}^m (1 + d_i(k-1)) = k^m$ ; thus, the theorem strictly generalizes the no-FD case.

**4.3. Query structure.** Finally, we explain how the FEQ’s structure can be exploited to speed up the computation of subproblems, the grid, and grid weights. In particular, we make use of recent advances in relational query evaluation algorithms [5, 4, 32, 33]. The InsideOut algorithm from the FAQs framework in particular [5] allows us to compute the grid weights without explicitly the grid points.

For concreteness, we describe the steps of Rk-means as implemented in the database, noting the additional speedups we can get over the description in Algorithm 1.

**Step 1** (lines 5 and 6). *Project  $\mathbf{X}$  into each subspace  $S_j$  and compute the weight  $w$  of each point.*

In a relational database, the projected sets  $\mathbf{X}_j$  already exist in normalized form [1], and thus they and their marginal weights can be computed highly efficiently. This step perfectly aligns with our strategy of picking the partition  $S_1 \cup \dots \cup S_m$  to match the database schema!

**Step 2** (line 7). *Find  $\kappa$  centroids in each subspace  $S_j$ .*

If the subspace  $S_j$  corresponds to a single continuous variable, we can solve the one-dimensional  $k$ -means problem quickly and optimally [42], and if the subspace corresponds to a categorical feature, then it is solved trivially (and optimally) using Theorem 4.4.

**Step 3** (lines 9 and 10). *Construct the coreset  $\mathbf{G}$  and the associated weights  $w_{\text{grid}}$ .*

When constructing  $\mathbf{G}$ , it is unnecessary to represent any points in  $\mathbf{G}$  that have zero weight. We use InsideOut [4] to efficiently compute nonzero weights, and then extract only those grid points in  $\mathbf{G}$  with nonzero weight from the database.

**Step 4** (line 11). *Cluster the weighted coreset  $\mathbf{G}$ .*

We use a modified version of Lloyd’s weighted  $k$ -means that exploits the structure of  $\mathbf{G}$  and sparse representation of categorical values to speed up computation

We discuss the optimization and acceleration of Step 4 of the Rk-means implementation in more details here. Recall that the categorical subspace  $k$ -means problem is solved trivially using Theorem 4.4, where we sort all the weights, and the heaviest  $k - 1$  elements form their own centroid, while the remaining vectors are clustered together (the “light cluster”).

If  $S_j$  is a categorical subspace corresponding to a categorical variable  $K$  where  $\text{Dom}(K) = \{e_1, \dots, e_L\}$ . Without loss of generality, assume  $w(\mathbf{1}_{e_1}) \geq \dots \geq w(\mathbf{1}_{e_L})$ , then the centroid of the light cluster is an  $L$ -dimensional vector  $\mathbf{c} = (s_e)_{e \in \text{Dom}(K)}$

$$(36) \quad s_{e_i} := \begin{cases} 0 & i \in [k - 1] \\ \frac{w(\mathbf{1}_{e_i})}{\sum_{j=k}^L w(\mathbf{1}_{e_j})} & i \geq k \end{cases}$$

This encoding is sound and space-inefficient.

Remember also that Step 4 clusters the coreset  $\mathbf{G}$  using a modified version of Lloyd’s weighted  $k$ -means that exploits the structure of  $\mathbf{G}$  and sparse representation of categorical values. We show how to improve the distance computation  $\|\mathbf{c}_j - \boldsymbol{\mu}_j\|^2$  for sub-space  $S_j$ , where  $\mathbf{c}_j$  and  $\boldsymbol{\mu}_j$  are the  $j$ -th components of a grid point and respectively of a centroid for  $\mathbf{G}$ . Since this subspace corresponds to a categorical variable  $K$  with, say,  $L_j$  categories, it is mapped into  $L_j$  sub-dimensions. Let  $\mathbf{c}_j = [s_1, \dots, s_{L_j}]$  and  $\boldsymbol{\mu}_j = (t_1, \dots, t_{L_j})$ . Using the explicit one-hot encoding of its categories, we would need  $O(L_j)$  time to compute  $\|\mathbf{c}_j - \boldsymbol{\mu}_j\|^2 = \sum_{\ell \in [L_j]} (s_\ell - t_\ell)^2$ . We can instead achieve  $O(1)$  time as shown next. There are  $k$  distinct values for  $\mathbf{c}_j$  by our coreset construction, each represented by a vector of size  $L_j$  with one non-zero entry for  $k - 1$  of them and  $L_j - k + 1$  non-zero entries for one of them.

If  $\mathbf{c}_j = \mathbf{1}_e$  is an indicator vector for some element  $e \in K$  ( $e$  is one of the  $k - 1$  heavy categories), then

$$(37) \quad \|\mathbf{c}_j - \boldsymbol{\mu}_j\|^2 = \|\mathbf{1}_e - \boldsymbol{\mu}_j\|^2 = 1 - 2t_e + \|\boldsymbol{\mu}_j\|^2.$$

If  $\mathbf{c}_j$  is a light cluster centroid,

$$(38) \quad \|\mathbf{c}_j - \boldsymbol{\mu}_j\|^2 = \|\mathbf{c}_j\|^2 + \|\boldsymbol{\mu}_j\|^2 - 2\langle \mathbf{c}_j, \boldsymbol{\mu}_j \rangle.$$

In (37), by pre-computing  $\|\boldsymbol{\mu}_j\|^2$  we only spend  $O(1)$ -time per heavy element  $e$ . In (38), by also pre-computing  $\|\mathbf{c}_j\|^2$  and  $\langle \mathbf{c}_j, \boldsymbol{\mu}_j \rangle$ , and by noticing that  $\mathbf{c}_j$  is  $(L_j - k + 1)$ -sparse, we spend  $O(L_j - k)$ -time here. Overall, we spend time  $O(L_j)$  for computing  $\|\mathbf{c}_j - \boldsymbol{\mu}_j\|^2$  per categorical dimension, modulo the precomputation time.

Step 4 thus requires  $O(|\mathbf{G}|mk + \sum_{j \in [m]} L_j k) = O(|\mathbf{G}|mk + Dkm)$  per iteration, whereas a generic approach would take time  $O(\sum_{j \in [m]} |\mathbf{G}|kL_j) = O(|\mathbf{G}|Dkm)$ . Our modified weighted  $k$ -means algorithm thus saves a factor proportional to the total domain sizes of the categorical variables, which may be as large as  $D$ .

**4.4. Runtime analysis.** We compare Rk-means to the standard setting of first extracting the matrix  $\mathbf{X}$  from the database and then perform clustering on  $\mathbf{X}$  directly. The precise runtime statement requires defining a few parameters such as “fractional hypertree width” and “fractional edge cover number” of the FEQ, which we briefly covered in Section 2.1. Hence, we state the main thrust of our runtime result:

**Theorem 4.7.** *There are classes of feature extraction queries (FEQs) for which the runtime of Rk-means is asymptotically less than  $|\mathbf{X}|$ , and the ratio between  $|\mathbf{X}|$  and the runtime of Rk-means can be a polynomial in  $N$ , the size of the largest input relation.*

*Proof of Theorem 4.7.* Let  $N$  denote the maximum number of tuples in any input relation of the FEQ,  $|\mathbf{X}|$  the number of tuples in the data matrix,  $\text{ftw}$  the *fractional hypertree width* of the FEQ  $t$  the number of iterations of Lloyd’s algorithm,  $d$  denote the number of features pre-one-hot encoding,  $r$  number of input relations to the FEQ,  $D$  the real dimensionality of the problem after one-hot-encoding.

We analyze the time complexity for each of the four steps of the Rk-means algorithm.

Step 1 projects  $\mathbf{X}$  into each subspace  $S_j$  and compute the total weight of each projected point:

$$(39) \quad \forall j \in [d] : w_j(\mathbf{x}_{S_j}) := \sum_{\mathbf{x}_{[d] \setminus \{S_j\}}} \prod_{F \in \mathcal{E}} R_F(\mathbf{x}_F)$$

	Retailer	Favorita	Yelp
Relations	5	6	6
Attributes	39	15	25
One-hot Enc.	95	1470	1617
# Rows in $\mathbf{D}$	84M	125M	8.7M
Size of $\mathbf{D}$	1.5GB	2.5GB	0.2GB
# Rows in $\mathbf{X}$	84M	127M	22M
Size of $\mathbf{X}$	18GB	7GB	2.4GB
# Rows in Coreset $\mathbf{G}$			
$\kappa = 5$	1.43M	14.94K	2.69M
$\kappa = 10$	9.58M	85.88K	11.71M
$\kappa = 20$	38.16M	632.5K	11.89M
$\kappa = 50$	73.75M	7.87M	12.46M

TABLE 1. Statistics for the input database  $\mathbf{D}$ , data matrix  $\mathbf{X}$ , and coresets  $\mathbf{G}$  for the three dataset.

Each of the  $d$  FAQs (39) in Step 1 can be computed in time  $\tilde{O}(rd^2N^{\text{ftw}})$  using InsideOut, as we have reviewed in Section 2.1.

In Step 2, the optimal clustering in each dimension takes time  $\tilde{O}(L_j)$  for each categorical variable  $j$  (whose domain size is  $L_j$ , and  $O(kN^2)$  for each continuous variable, with an overall runtime of  $O(kdN^2)$ .

Step 3 constructs  $\mathbf{G}$ , whose size is bounded by  $|\mathbf{X}|$  and by the FD result of Theorem 4.6. In practice, this number can be much smaller since we skip the data points whose weights are zero. To perform this step we construct a tree decomposition of FEQ with with equal  $\text{ftw}$  (this step is data-independent, only dependent on the size of FEQ). Then, from each value  $x_j$  of an input variable  $X_j$ , we determine its centroid  $c(x_j)$  which was computed in step 2. By conditioning on combinations of  $(c_1, \dots, c_j)$ , we can compute  $w_{\text{grid}}$  one for each combination in  $\tilde{O}(dN^{\text{ftw}})$ -time, for a total run time of  $\tilde{O}(rd|\mathbf{G}|N^{\text{ftw}})$ .

Step 4 – as analyzed in Section 4.4 – clusters  $\mathbf{G}$  in time  $O((|\mathbf{G}|+D)kmt)$ , where  $t$  is the number of iterations of  $k$ -means used in Step 4. The most expensive computation is due to the one-dimensional clustering for the continuous variables and the computation of the coreset.

To compare the total runtime with  $|\mathbf{X}|$ , we only need to note that  $|\mathbf{X}|$  can be as large as  $N^{\rho^*}$ , where  $\rho^*$  is the *fractional edge covering number* of the FEQ’s hypergraph [32]. Depending on the query,  $\rho^*$  is always at least 1, and can be as large as the number of features  $d$ . Furthermore, there are classes of queries where  $\text{ftw}$  is bounded by a constant, yet  $\rho^*$  is unbounded [30]. This means, for classes of FEQs where  $\rho^* > \max\{\text{ftw}, 2\}$  the ratio between  $|\mathbf{X}|$  and Rk-means’s runtime will be  $\tilde{O}(\text{omega}(N^{\rho^* - \max\{\text{ftw}, 2\}}/t))$ , which is unbounded.  $\square$

The key insight to read from this theorem is that Rk-means can, in principle, run faster than simply exporting the data matrix, without even running *any* clustering algorithm (be it sampling-based, streaming, etc.). Of course, the result only concerns a class of FEQs “on paper”. Section 5 examines real FEQs, which also demonstrate Rk-means’s runtime superiority.

For reference, we compare the asymptotic runtime of Rk-means to the standard implementation of Lloyd’s algorithm. The standard implementation contains two steps: (1) compute the one-hot-encoded data matrix  $\mathbf{X}$ , and (2) run Lloyd’s algorithm on  $\mathbf{X}$ . The first step, materializing  $\mathbf{X}$ , takes time  $\tilde{O}(rd^2N^{\text{ftw}} + D|\mathbf{X}|)$ . The second step, running Lloyd algorithm, takes time  $\tilde{O}(tkD|\mathbf{X}|)$ , as is well known. Thus, the standard approach takes time  $\tilde{O}(rd^2N^{\text{ftw}} + tkD|\mathbf{X}|)$ .

## 5. EXPERIMENTAL RESULTS

We empirically evaluate the performance of Rk-means on three real datasets for three sets of experiments: (1) we break down and analyze the performance of each step in Rk-means; (2) we benchmark the performance and approximation of Rk-means against mlpack [13] (v. 3.1.0), a fast C++ machine learning library; and (3) we evaluate the performance and approximation of Rk-means for setting  $\kappa < k$ ; i.e., different number of clusters for Steps 2 and 4.

The experiments show that the coresets of Rk-means are often significantly smaller than the data matrix. As a result, Rk-means can scale easily to large datasets, and can compute the clusters with a much lower memory footprint than mpack. When  $\kappa = k$ , Rk-means is orders-of-magnitude faster than the end-to-end computation for mpack—up to  $115\times$ . Typically, the approximation level is very minor. In addition, setting  $\kappa < k$  can lead to further performance speedups with only a moderate increase in approximation, giving over  $200\times$  speedup in some cases.

**Experimental Setup.** We prototyped Rk-means as part of an engine designed to compute multiple FAQ expressions efficiently. Rk-means is implemented in multithreaded C++11 and compiled with -O3 optimizations; this makes mpack a comparable implementation. All experiments were performed on an AWS `x1e.8xlarge` system, which has 1 TiB of RAM and 32 vCPUs. All relations given were sorted by their join attributes.

To construct the data matrix that forms the input to mpack, we use PostgreSQL (`psql`) v. 10.6 to evaluate the FEQ. The seminal  $k$ -means++ algorithm [7] is used for initializing the  $k$ -means cluster. We run Rk-means and mpack + psql five times and report the average approximation and runtime. The timeout for all experiments was set to six hours (21,600 seconds) per trial. Our runtime results omit data loading/saving times. Note that for mpack + psql, psql must export  $\mathbf{X}$  to disk, and then mpack must then read it from disk. Rk-means has no need to do this, and thus the runtime numbers are skewed in mpack’s favor. This skew may be significant: loading and saving a large CSV file may take hours in some cases.

**Datasets.** We use three real datasets: (1) *Retailer* is used by a large US retailer for sales forecasting; (2) *Favorita* [17] is a public dataset for retail forecasting; and (3) *Yelp* is from the public Yelp Dataset Challenge [46] and used to predict users’ ratings of businesses. Table 1 presents key statistics for the three datasets, including the size of data matrix  $\mathbf{X}$  and the coreset  $\mathbf{G}$  for each dataset and different  $\kappa$ -values.  $|\mathbf{G}|$  is highly data dependent. For *Favorita*,  $\mathbf{G}$  is orders-of-magnitude smaller than the data matrix. For *Retailer*, when  $k = 20$  and  $k = 50$ ,  $|\mathbf{G}|$  approaches  $|\mathbf{X}|$ , but Rk-means is still able to provide a speedup. Additional dataset details are given below.

*Retailer* has five relations: *Inventory* stores the number of inventory units for each date, location, and stock keeping unit (sku); *Location* keeps for each store: its zipcode, the distance to the closest competitors, and the type of the store; *Census* provides 14 attributes that describe the demographics of a given zipcode, including population size or average household income; *Weather* stores statistics about the weather condition for each date and store, including the temperature and whether it rained; *Items* keeps track of the price, category, subcategory, and category cluster of each sku.

*Favorita* has six relations: *Sales* stores the number of units sold for items for a given date and store, and an indicator whether or not the unit was on promotion at this time; *Items* provides additional information about the skus, such as the item class and price; *Stores* keeps additional information on stores, like the city they are located it; *Transactions* stores the number of transaction for each date and store; *Oil* provides the oil price for each date; and *Holiday* indicates whether a given day is a holiday. The original dataset gave the `units_sold` attribute with a precision of three decimal places. This resulted in a very many distinct values for this attribute, which has a significant impact on the Step 2 of the Rk-means algorithm. We decreased the precision for this attribute to two decimal places, which decreases the number of distinct values by a factor of four. This modification has no effect on the final clusters or their accuracy.

*Yelp* has five relations: *Review* gives the review rating that a user gave to a business and the date of the review; *User* provides information about the users, including how many reviews they made, when they join, and how many fans they have; *Business* provides information about the businesses that are reviewed, such as their location and average rating; *Category* provide information about the categories, i.e. Restaurant, and respectively attributes of the business, *Attributes* is an aggregated relation, which stores the number of attributes (i.e., open late) that have been assigned to a business. A business can be categorized in many ways, which is the main reason why the size of the join is significantly larger than the underlying relations.

**Breakdown of Rk-means.** Figure 3 shows the time it takes Rk-means to cluster the three datasets for different values of  $k$  with  $\kappa = k$ . The total time is broken down into the four steps of the algorithm from Section 4. We provide the time it takes psql to compute  $\mathbf{X}$  as reference (gray bar). In many cases, Rk-means can cluster *Retailer* and *Favorita* faster than it takes psql to even compute the data matrix. The relative performance of the four steps is data dependent. For *Retailer*, most of the time is spent on constructing  $\mathbf{G}$  in Step 3, which is relatively large. For *Favorita*, however, Step 2 takes the longest, since it contains one

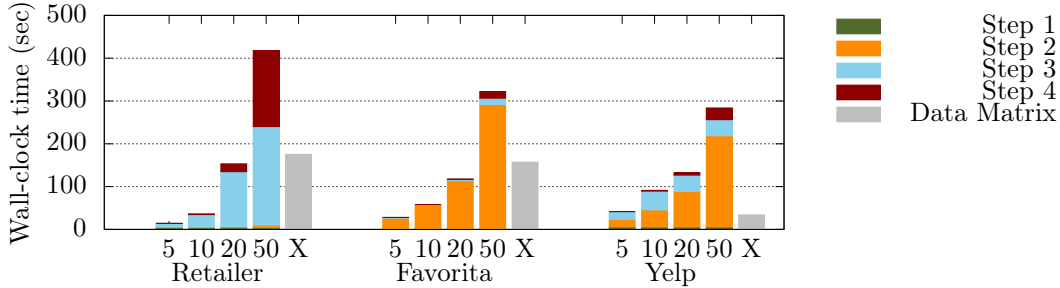


FIGURE 3. Breakdown of the compute time of Rk-means for each step of the algorithm with  $\kappa = k$ . The time to compute  $\mathbf{X}$  is provided as reference.

<b>Retailer</b>	k = 5	k = 10	k = 20	k = 50	k=20, $\kappa = 10$	k = 50, $\kappa = 20$
Compute $\mathbf{X}$ (psql)	175.47	175.47	175.47	175.47	175.47	175.47
Clustering (mlpack)	65.41	158.81	385.67	1,453.88	385.67	1,453.88
Rk-means	15.66	54.59	230.17	650.20	63.51	344.31
Relative Speedup	15.38×	6.12×	2.44×	2.51×	8.84×	4.73×
Relative Approx.	0.20	0.08	0.03	0.00	0.03	0.02
<b>Favorita</b>	k = 5	k = 10	k = 20	k = 50	k=20, $\kappa = 10$	k = 50, $\kappa = 20$
Compute $\mathbf{X}$ (psql)	156.86	156.86	156.86	156.86	156.86	156.86
Clustering (mlpack)	1,002.54	6,449.32	11,794.49	>21,600.00	11,794.49	>21,600
Rk-means	27.95	57.72	118.36	334.65	57.65	120.77
Relative Speedup	41.49×	114.59×	100.98×	>64.55×	207.30×	>178.86×
Relative Approx.	2.99	0.35	0.12	–	1.93	–
<b>Yelp</b>	k = 5	k = 10	k = 20	k = 50	k=20, $\kappa = 10$	k = 50, $\kappa = 20$
Compute $\mathbf{X}$ (psql)	33.83	33.83	33.83	33.83	33.83	33.83
Clustering (mlpack)	210.59	640.43	2,107.83	11,474.24	2,107.83	11,474.24
Rk-means	43.37	107.71	195.22	405.11	114.34	241.34
Relative Speedup	5.64×	6.26×	10.97×	28.41×	18.73×	47.68×
Relative Approx.	0.37	0.26	0.13	0.05	0.27	0.20

TABLE 2. End-to-end runtime and approximation comparison of Rk-means and mlpack on each dataset. The first four columns use different  $\kappa = k$  values; the last two show results for setting  $\kappa < k$ .

continuous variable with many distinct values, and the DP algorithm for clustering runs in time quadratic in the number of distinct values. The runtime for *Favorita* could be improved by clustering this dimension with a different  $k$ -means algorithm instead, but this may increase the approximation.

**Comparison with mlpack.** The left columns of Table 2 compares the runtime and approximation of Rk-means against mlpack on the three datasets for different  $k$  values with  $\kappa = k$ . The approximation is given relative to the objective value obtained by mlpack. Speedup is given by comparing the end-to-end performance of Rk-means and mlpack (ignoring disk I/O time), which for mlpack includes the time needed by psql to materialize  $\mathbf{X}$ . Overall, Rk-means often outperforms even just the clustering step from mlpack, and when end-to-end computation is considered, Rk-means gives up to 115× speedup. mlpack timed out after six hours for *Favorita* with  $k = 50$ . In addition, Rk-means has a much smaller memory footprint than mlpack: for instance, on the *Favorita* dataset with  $k = 20$ , mlpack uses over 900GiB of RAM to cluster the dataset, whereas Rk-means only requires 18GiB. In our simulations, the approximation level is moderate, and consistently well below the 9-approximation bound from Theorem 3.4.

**Setting  $\kappa < k$  for Step 2.** We next evaluate the effect of setting  $\kappa$  to a smaller value than the number of clusters  $k$ . This exploits the speed/approximation tradeoff: smaller  $\kappa$  helps reduce the size of  $\mathbf{G}$ , at the cost of more approximation. Table 2 presents for each dataset the results for setting  $k = 20, \kappa = 10$  and  $k = 50, \kappa = 20$ , and compares them to the relative performance and approximation over computing  $k$  clusters in mlpack.

By setting  $\kappa < k$ , Rk-means can compute the  $k$  clusters up to  $208\times$  faster than mlpack and  $3.6\times$  faster than Rk-means with  $\kappa = k$ , while the relative approximation remains moderate. Our results are data dependent—but as queries and databases scale, our speedups will be even more significant.

## 6. CONCLUSION

We introduce Rk-means, a method to construct  $k$ -means clustering coresets on relational data directly from the database. Rk-means gives a provably good clustering of the entire dataset, without ever materializing the data set; this also yields asymptotic improvements in running time. Experimentally, we observe that the coreset has size up to  $180\times$  smaller than the size of the data matrix and this compression results in orders-of-magnitude improvements in the running time, while still providing empirically good clusterings. Although our work here primarily focuses on  $k$ -means clustering, we believe our construction of grid coresets and the accompanied theory may be useful for other unsupervised learning tasks and plan to explore such possibilities in future work.

## REFERENCES

- [1] S. ABITEBOUL, R. HULL, AND V. VIANU, *Foundations of Databases*, Addison-Wesley, 1995.
- [2] M. ABO KHAMIS, H. Q. NGO, X. NGUYEN, D. OLTEANU, AND M. SCHLEICH, *Ac/dc: In-database learning thunderstruck*, in 2nd Workshop on Data Mgt for End-To-End ML, DEEM'18, 2018, pp. 8:1–8:10.
- [3] M. ABO KHAMIS, H. Q. NGO, X. NGUYEN, D. OLTEANU, AND M. SCHLEICH, *In-database learning with sparse tensors*, in PODS, 2018, pp. 325–340.
- [4] M. ABO KHAMIS, H. Q. NGO, AND A. RUDRA, *FAQ: questions asked frequently*, in PODS, 2016, pp. 13–28.
- [5] M. ABO KHAMIS, H. Q. NGO, AND A. RUDRA, *Juggling functions inside a database*, SIGMOD Rec., 46 (2017), pp. 6–13.
- [6] S. AHMADIAN, A. NOROUZI-FARD, O. SVENSSON, AND J. WARD, *Better guarantees for  $k$ -means and euclidean  $k$ -median by primal-dual algorithms*, in FOCS, 2017, pp. 61–72.
- [7] D. ARTHUR AND S. VASSILVITSKII,  *$k$ -means++: The advantages of careful seeding*, in SODA, 2007, p. 1027–1035.
- [8] O. BACHEM, M. LUCIC, AND A. KRAUSE, *Scalable  $k$ -means clustering via lightweight coresets*, in SIGKDD, 2018, pp. 1119–1127.
- [9] B. BAHMANI, B. MOSELEY, A. VATTANI, R. KUMAR, AND S. VASSILVITSKII, *Scalable  $k$ -means++*, PVLDB, 5 (2012), pp. 622–633.
- [10] M. BALCAN, S. EHRLICH, AND Y. LIANG, *Distributed  $k$ -means and  $k$ -median clustering on general communication topologies*, in Advances in Neural Information Processing Systems 26: 27th Annual Conference on Neural Information Processing Systems 2013. Proceedings of a meeting held December 5-8, 2013, Lake Tahoe, Nevada, United States., 2013, pp. 1995–2003.
- [11] V. BRAVERMAN, G. FRAHLING, H. LANG, C. SOHLER, AND L. F. YANG, *Clustering high dimensional dynamic data streams*, in Proceedings of the 34th International Conference on Machine Learning, ICML 2017, Sydney, NSW, Australia, 6-11 August 2017, 2017, pp. 576–585.
- [12] F. CADY, *The Data Science Handbook*, John Wiley & Sons, 2017.
- [13] R. R. CURTIN, M. EDEL, M. LOZHNIKOV, Y. MENTEKIDIS, S. GHAISAS, AND S. ZHANG, *mlpack 3: a fast, flexible machine learning library*, J. Open Source Software, 3 (2018), p. 726.
- [14] E. DEL BARRIO, J. A. CUESTA-ALBERTOS, C. MATRAN, AND A. MAYO-ISCAR, *Robust clustering tools based on optimal transportation*, Statistics and Computing, (2017), pp. 1–22.
- [15] A. ENE, S. IM, AND B. MOSELEY, *Fast clustering using mapreduce*, in Proceedings of the 17th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, San Diego, CA, USA, August 21-24, 2011, 2011, pp. 681–689.
- [16] B. EVERITT, S. LANDAU, M. LEESE, AND D. STAHL, *Cluster Analysis*, Wiley & Sons, 2011.
- [17] FAVORITA CORP., *Corporacion Favorita Grocery Sales Forecasting: Can you accurately predict sales for a large grocery chain?* <https://www.kaggle.com/c/favorita-grocery-sales-forecasting/>, 2017.
- [18] S. GRAF AND H. LUSCHGY, *Foundations of quantization for probability distributions*, Springer-Verlag, New York, 2000.
- [19] M. GROHE AND D. MARX, *Constraint solving via fractional edge covers*, ACM Trans. Algorithms, 11 (2014), pp. 4:1–4:20.
- [20] S. GUHA, A. MEYERSON, N. MISHRA, R. MOTWANI, AND L. O'CALLAGHAN, *Clustering data streams: Theory and practice*, IEEE Trans. Knowl. Data Eng., 15 (2003), pp. 515–528.
- [21] S. HAR-PELED AND S. MAZUMDAR, *On coresets for  $k$ -means and  $k$ -median clustering*, in SODA, 2004.
- [22] S. HAR-PELED AND S. MAZUMDAR, *Coresets for  $k$ -means and  $k$ -median clustering and their applications*, CoRR, abs/1810.12826 (2018).
- [23] J. A. HARTIGAN, *Clustering algorithms*, Wiley, New York, 1975.
- [24] N. HO, X. NGUYEN, M. YUROCHKIN, H. H. BUI, V. HUYNH, AND D. PHUNG, *Multilevel clustering via Wasserstein means*, in ICML, 2017, pp. 1501–1509.
- [25] Z. HUANG, *Extensions to the  $k$ -means algorithm for clustering large data sets of categorical values*, Data mining and Knowledge discovery, 2 (1998), pp. 283–304.
- [26] L. KAUFMAN AND P. J. ROUSSEW, *Finding Groups in Data - An Introduction to Cluster Analysis*, Wiley & Sons, 1990.
- [27] M. A. KHAMIS, R. CURTIN, B. MOSELEY, H. NGO, X. NGUYEN, D. OLTEANU, AND M. SCHLEICH, *On functional aggregate queries with additive inequalities*, in PODS, 2019.



- [28] D. KOLLER AND N. FRIEDMAN, *Probabilistic graphical models*, Adaptive Computation and Machine Learning, MIT Press, Cambridge, MA, 2009. Principles and techniques.
- [29] S. P. LLOYD, *Least squares quantization in PCM*, IEEE Trans. Inf. Theory, 28 (1982), pp. 129–136.
- [30] D. MARX, *Tractable hypergraph properties for constraint satisfaction and conjunctive queries*, J. ACM, 60 (2013), pp. 42:1–42:51.
- [31] A. MEYERSON, L. O’CALLAGHAN, AND S. A. PLOTKIN, *A k-median algorithm with running time independent of data size*, Machine Learning, 56 (2004), pp. 61–87.
- [32] H. Q. NGO, *Worst-case optimal join algorithms: Techniques, results, and open problems*, in PODS, 2018, pp. 111–124.
- [33] D. OLTEANU AND M. SCHLEICH, *Factorized databases*, SIGMOD Rec., 45 (2016), pp. 5–16.
- [34] C. ORDONEZ, *Integrating k-means clustering with a relational DBMS using SQL*, IEEE Trans. Knowl. Data Eng., 18 (2006), pp. 188–201.
- [35] C. ORDONEZ AND E. OMIECINSKI, *Efficient disk-based k-means clustering for relational databases*, IEEE Trans. Knowl. Data Eng., 16 (2004), pp. 909–921.
- [36] D. POLLARD, *Quantization and the method of k-means*, IEEE Trans. Inf. Theory, 28 (1982), pp. 199–204.
- [37] M. SCHLEICH, D. OLTEANU, AND R. CIUCANU, *Learning linear regression models over factorized joins*, in SIGMOD, 2016, pp. 3–18.
- [38] C. SOHLER AND D. P. WOODRUFF, *Strong coresets for k-median and subspace approximation: Goodbye dimension*, in 59th IEEE Annual Symposium on Foundations of Computer Science, FOCS 2018, Paris, France, October 7-9, 2018, 2018, pp. 802–813.
- [39] W. SUN, J. WANG, AND Y. FANG, *Regularized k-means clustering of high-dimensional data and its asymptotic consistency*, Electronic Journal of Statistics, 9 (2012), pp. 148–167.
- [40] M. THORUP, *Quick k-median, k-center, and facility location for sparse graphs*, in Automata, Languages and Programming, 28th International Colloquium, ICALP 2001, Crete, Greece, July 8-12, 2001, Proceedings, 2001, pp. 249–260.
- [41] C. VILLANI, *Optimal transport*, vol. 338 of Grundlehren der Mathematischen Wissenschaften [Fundamental Principles of Mathematical Sciences], Springer-Verlag, Berlin, 2009.
- [42] H. WANG AND M. SONG, *Ckmeans.1d.dp: Optimal k-means clustering in one dimension by dynamic programming*, The R Journal, 3 (2011), p. 29.
- [43] D. WITTEN AND R. TIBSHIRANI, *A framework for feature selection in clustering*, Journal of the American Statistical Association, 105 (2010), pp. 713–726.
- [44] X. WU, V. KUMAR, J. R. QUINLAN, J. GHOSH, Q. YANG, H. MOTODA, G. J. MCLACHLAN, A. F. M. NG, B. LIU, P. S. YU, Z. ZHOU, M. STEINBACH, D. J. HAND, AND D. STEINBERG, *Top 10 algorithms in data mining*, Knowl. Inf. Syst., 14 (2008), pp. 1–37.
- [45] J. YE, P. WU, J. WANG, AND J. LI, *Fast discrete distribution clustering using barycenter with sparse support*, IEEE Trans. Signal Proc., 65 (2017), pp. 2317–2332.
- [46] YELP, *Yelp dataset challenge*, <https://www.yelp.com/dataset/challenge/>, 2017.

RELATIONALAI

CARNEGIE MELLON UNIVERSITY

RELATIONALAI

UNIVERSITY OF MICHIGAN

OXFORD UNIVERSITY

OXFORD UNIVERSITY