

MINI-BATCH KERNEL k -MEANS

Anonymous authors

Paper under double-blind review

ABSTRACT

We present the first mini-batch kernel k -means algorithm, offering an order of magnitude improvement in running time compared to the full batch algorithm. A single iteration of our algorithm takes $\tilde{O}(kb^2)$ time, significantly faster than the $O(n^2)$ time required by the full batch kernel k -means, where n is the dataset size and b is the batch size. Extensive experiments demonstrate that our algorithm consistently achieves a 10-100x speedup with minimal loss in quality, addressing the slow runtime that has limited kernel k -means adoption in practice. We further complement these results with a theoretical analysis under an early stopping condition, proving that with a batch size of $\tilde{\Omega}(\max\{\gamma^4, \gamma^2\} \cdot k\epsilon^{-2})$, the algorithm terminates in $O(\gamma^2/\epsilon)$ iterations with high probability, where γ bounds the norm of points in feature space and ϵ is a termination threshold. Our analysis holds for any reasonable center initialization, and when using k -means++ initialization, the algorithm achieves an approximation ratio of $O(\log k)$ in expectation. For normalized kernels, such as Gaussian or Laplacian it holds that $\gamma = 1$. Taking $\epsilon = O(1)$ and $b = \Theta(k \log n)$, the algorithm terminates in $O(1)$ iterations, with each iteration running in $\tilde{O}(k^3)$ time.

1 INTRODUCTION

Mini-batch methods are among the most successful tools for handling huge datasets for machine learning. Notable examples include Stochastic Gradient Descent (SGD) and mini-batch k -means (Sculley, 2010). Mini-batch k -means is one of the most popular clustering algorithms used in practice (Pedregosa et al., 2011).

While k -means is widely used due to its simplicity and fast running time, it requires the data to be *linearly separable* to achieve meaningful clustering. Unfortunately, many real-world datasets do not have this property. One way to overcome this problem is to project the data into a high, even *infinite*, dimensional space (where it is hopefully linearly separable) and run k -means on the projected data using the “kernel-trick”. A toy example is given in Figure 1 and a more realistic example is given in Figure 2.

Kernel k -means achieves significantly better clustering compared to k -means in practice. However, its running time is considerably slower. Surprisingly, prior to our work there was no attempt to speed up kernel k -means using a mini-batch approach.

Problem statement We are given an input (dataset), $X = \{x_i\}_{i=1}^n$, of size n and a parameter k representing the number of clusters. A kernel for X is a function $K : X \times X \rightarrow \mathbb{R}$ that can be realized by inner products. That is, there exists a Hilbert space \mathcal{H} and a map $\phi : X \rightarrow \mathcal{H}$ such that $\forall x, y \in X, \langle \phi(x), \phi(y) \rangle = K(x, y)$. We call \mathcal{H} the *feature space* and ϕ the *feature map*.

In kernel k -means the input is a dataset X and a kernel function K as above. Our goal is to find a set \mathcal{C} of k centers (elements in \mathcal{H}) such that the following goal function is minimized: $\frac{1}{n} \sum_{x \in X} \min_{c \in \mathcal{C}} \|c - \phi(x)\|^2$. Equivalently we may ask for a partition of X into k parts, keeping \mathcal{C} implicit.¹

¹A common variant of the above is when every $x \in X$ is assigned a weight $w_x \in \mathbb{R}^+$ and we aim to minimize $\sum_{x \in X} w_x \cdot \min_{c \in \mathcal{C}} \|c - \phi(x)\|^2$. Everything that follows, including our results, can be easily generalized to the weighted case. We present the unweighted case to improve readability.

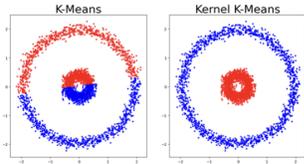
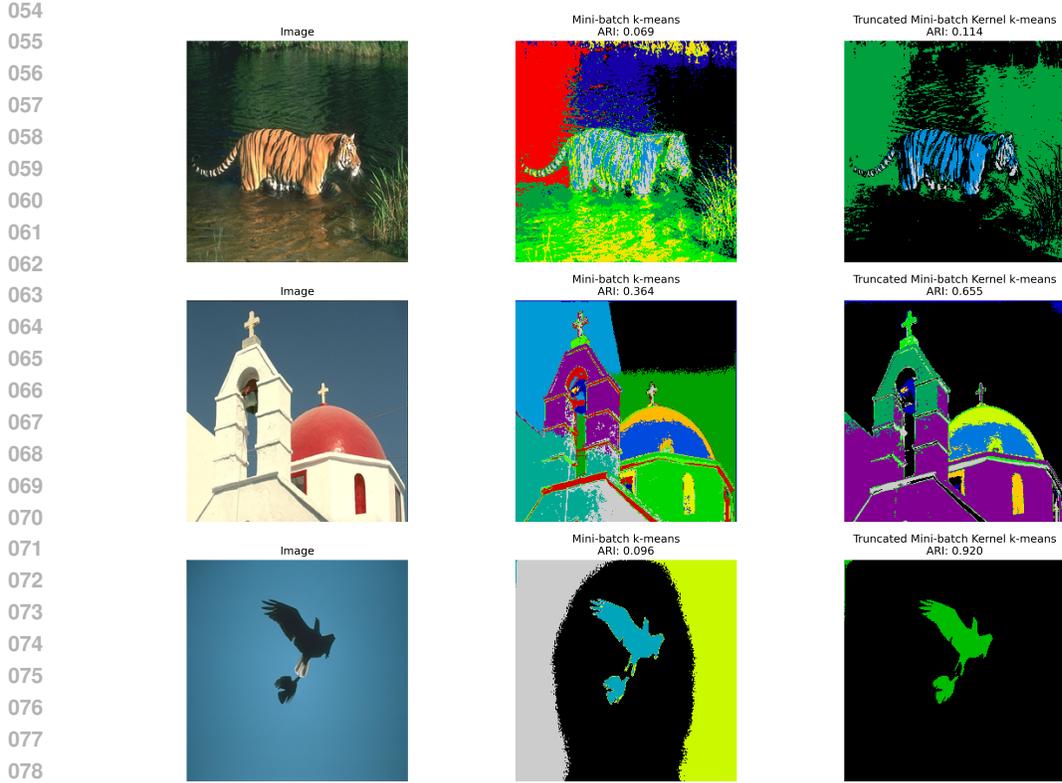


Figure 1: Kernel k -means perfectly clusters the dataset, while k -means cannot.



080 Figure 2: Qualitative comparison of mini-batch k -means and our algorithm (truncated mini-batch
081 kernel k -means) on selected images from the Berkeley Segmentation Data Set (BSDS)(Arbelaez
082 et al., 2010) using the Gaussian kernel. ARI is the Adjusted Rand Index (Rand, 1971).
083
084

085 **Lloyd’s algorithm** The most popular algorithm for (non kernel) k -means is Lloyd’s algorithm,
086 often referred to as the k -means algorithm (Lloyd, 1982). It works by randomly initializing a set of k
087 centers and performing the following two steps: (1) Assign every point in X to the center closest to it.
088 (2) Update every center to be the mean of the points assigned to it. The algorithm terminates when no
089 point is reassigned to a new center. This algorithm is extremely fast in practice but has a worst-case
090 exponential running time (Arthur & Vassilvitskii, 2006; Vattani, 2011).

091 **Mini-batch k -means** To update the centers, Lloyd’s algorithm must go over the entire input at
092 every iteration. This can be computationally expensive when the input data is extremely large. To
093 tackle this, the mini-batch k -means method was introduced by Sculley (2010). It is similar to Lloyd’s
094 algorithm except that steps (1) and (2) are performed on a batch of b elements sampled uniformly
095 at random with repetitions, and in step (2) the centers are updated slightly differently. Specifically,
096 every center is updated to be the weighted average of its current value and the mean of the points (in
097 the batch) assigned to it. The parameter by which we weigh these values is called the *learning rate*,
098 and its value differs between centers and iterations. The larger the learning rate, the more a center
099 will drift towards the new batch cluster mean.

100 **Lloyd’s algorithm in feature space** Implementing Lloyd’s algorithm in feature space is challeng-
101 ing as we cannot explicitly keep the set of centers \mathcal{C} . Luckily, we can use the kernel function together
102 with the fact that centers are always set to be the mean of cluster points to compute the distance from
103 any point $x \in X$ in feature space to any center $c = \frac{1}{|A|} \sum_{y \in A} \phi(y)$ as follows:

104
105
106
107

$$\begin{aligned} \|\phi(x) - c\|^2 &= \langle \phi(x) - c, \phi(x) - c \rangle = \langle \phi(x), \phi(x) \rangle - 2\langle \phi(x), c \rangle + \langle c, c \rangle \\ &= \langle \phi(x), \phi(x) \rangle - 2\langle \phi(x), \frac{1}{|A|} \sum_{y \in A} \phi(y) \rangle + \langle \frac{1}{|A|} \sum_{y \in A} \phi(y), \frac{1}{|A|} \sum_{y \in A} \phi(y) \rangle, \end{aligned}$$

where A can be any subset of the input X . While the above can be computed using only kernel evaluations, it makes the update step significantly more costly than standard k -means. Specifically, the complexity of the above may be quadratic in n (Dhillon et al., 2004).

Mini-batch kernel k -means Applying the mini-batch approach for kernel k -means is even more difficult because the assumption that cluster centers are always the mean of some subset of X in feature space no longer holds.

In Section 4 we first derive a recursive expression that allows us to compute the distances of all points to current cluster centers (in feature space). Using a simple dynamic programming approach that maintains the inner products between the data and centers in feature space, we achieve a running time of $O(n(b+k))$ per iteration compared to $O(n^2)$ for the full-batch algorithm. However, a true mini-batch algorithm should have a running time sublinear in n , preferably only polylogarithmic. We show that the recursive expression can be *truncated*, achieving a fast update time of $\tilde{O}(kb^2)$ while only incurring a small additive error compared to the untruncated version².

While our main contribution is practical — achieving an order-of-magnitude speedup for kernel k -means — we also provide theoretical guarantees for our algorithm (deferred to Appendix B). This is somewhat tricky for mini-batch algorithms due to their stochastic nature, as they may not even converge to a local-minima. To overcome this hurdle, we take the approach of Schwartzman (2023) and answer the question: how long does it take truncated mini-batch kernel k -means to terminate with an *early stopping condition*. Specifically, we terminate the algorithm when the improvement on the batch drops below some user provided parameter, ϵ . Early stopping conditions are very common in practice (e.g., sklearn (Pedregosa et al., 2011)). We show that applying the k -means++ initialization scheme (Arthur & Vassilvitskii, 2007) for our initial centers implies we achieve the same approximation ratio, $O(\log k)$ in expectation, as the full-batch algorithm.

While our general approach is similar to (Schwartzman, 2023), we must deal with the fact that \mathcal{H} may have an *infinite* dimension. The guarantees of (Schwartzman, 2023) depend on the dimension of the space in which k -means is executed, which is unacceptable in our case. We overcome this by parameterizing our results by a new parameter $\gamma = \max_{x \in X} \|\phi(x)\|$. We note that for normalized kernels, such as the popular Gaussian and Laplacian kernels, it holds that $\gamma = 1$. We also observe that it is often the case that $\gamma \ll 1$ for various other kernels used in practice (see Appendix C). We show that if the batch size is $\Omega(\max\{\gamma^4, \gamma^2\} k \epsilon^{-2} \log^2(\gamma n / \epsilon))$ then w.h.p. our algorithm terminates in $O(\gamma^2 / \epsilon)$ iterations. Our theoretical results are summarised in Theorem 1.1 (where Algorithm 2 is presented in Section 4).

Theorem 1.1. *The following holds for Algorithm 2: (1) Each iteration takes $O(kb^2 \log^2(\gamma / \epsilon))$ time, (2) If $b = \Omega(\max\{\gamma^4, \gamma^2\} k \epsilon^{-2} \log^2(\gamma n / \epsilon))$ then it terminates in $O(\gamma^2 / \epsilon)$ iterations w.h.p, (3) When initialized with k -means++ it achieve a $O(\log k)$ approximation ratio in expectation.*

Our result improves upon (Schwartzman, 2023) significantly when a normalized kernel is used since Theorem 1.1 doesn't depend on the input dimension. Our algorithm copes better with non linearly separable data and requires a smaller batch size ($\tilde{\Omega}(1/\epsilon^2)$ vs $\tilde{\Omega}((d/\epsilon)^2)$)³ for normalized kernels. This is particularly apparent with high dimensional datasets such as MNIST (LeCun, 1998) where the dimension squared is already nearly ten times the number of datapoints.

The learning rate we use, suggested in (Schwartzman, 2023), differs from the standard learning rate of sklearn in that it does not go to 0 over time. Unfortunately, this new learning rate is non-standard and (Schwartzman, 2023) did not present experiments comparing their learning rate to that of sklearn. We fill the experimental gap left in (Schwartzman, 2023) by evaluating (non-kernel) mini-batch k -means with their new learning rate compared to that of sklearn. Following our experimental evaluation, the sklearn team accepted a pull request implementing this learning rate in future versions.

In Section 5 we extensively evaluate our results experimentally both with the learning rate of (Schwartzman, 2023) and that of sklearn. To allow a fair empirical comparison, we run each algorithm for a fixed number of iterations without stopping conditions. Our results are as follows: 1) Truncated mini-batch kernel k -means is significantly faster than full-batch kernel k -means, while achieving solutions of similar quality, which are superior to the non-kernel version, 2) The learning

²Where \tilde{O} hides factors that are polylogarithmic in $n, 1/\epsilon, \gamma$.

³In (Schwartzman, 2023) the tilde notation hides factors logarithmic in d instead of γ .

rate of (Schwartzman, 2023) results in solutions with better quality both for truncated mini-batch kernel k -means and (non-kernel) mini-batch k -means.

2 RELATED WORK

Until recently, mini-batch k -means was only considered with a learning rate going to 0 over time. This was true both in theory (Tang & Monteleoni, 2017; Sculley, 2010) and practice (Pedregosa et al., 2011). Recently, (Schwartzman, 2023) proposed a new learning which does not go to 0 over time, and showed that if the batch is of size $\tilde{\Omega}(k(d/\epsilon)^2)^4$, mini-batch k -means must terminate within $O(d/\epsilon)$ iterations with high probability, where d is the dimension of the input, and ϵ is a threshold parameter for termination.

A popular approach to deal with the slow running time of kernel k -means is constructing a *coreset* of the data. A coreset for kernel k -means is a weighted subset of X with the guarantee that the solution quality on the coreset is close to that on the entire dataset up to a $(1 + \epsilon)$ multiplicative factor. There has been a long line of work on coresets for k -means and kernel k -means (Schmidt, 2014; Feldman et al., 2020; Barger & Feldman, 2020), and the current state-of-the-art for kernel k -means is due to (Jiang et al., 2021). They present a coreset algorithm with a nearly linear (in n and k) construction time which outputs a coreset of size $\text{poly}(k\epsilon^{-1})$.

In (Chitta et al., 2011) the authors only compute the kernel matrix for uniformly sampled set of m points from X . Then they optimize a variant of kernel k -means where the centers are constrained to be linear combinations of the sampled points. The authors do not provide worst case guarantees for the running time or approximation of their algorithm.

Another approach to speed up kernel k -means is by computing an approximation for the kernel matrix. This can be done by computing a low dimensional approximation for ϕ (without computing ϕ explicitly) (Rahimi & Recht, 2007; Chitta et al., 2012; Chen & Phillips, 2017), or by computing a low rank approximation for the kernel matrix (Musco & Musco, 2017; Wang et al., 2019).

Kernel sparsification techniques construct sparse approximations of the full kernel matrix in sub-quadratic time. For smooth kernel functions such as the polynomial kernel, (Quanrud) presents an algorithm for constructing a $(1 + \epsilon)$ -spectral sparsifier for the full kernel matrix with a nearly linear number of non-zero entries in nearly linear time. For the gaussian kernel, (Macgregor & Sun, 2024) show how to construct a weaker, cluster preserving sparsifier using a nearly linear number of kernel density estimation queries.

We note that our results are *complementary* to coresets, dimensionality reduction, and kernel sparsification, in the sense that we can compose our method with these techniques.

To the best of our knowledge, the only approach which cannot be directly composed with our work is *kernel sketching* (Liu, 2021; Yin et al., 2022). Here the kernel matrix is used to compute an embedding of the points into a low dimensional Euclidean space, followed by running the standard (non-kernel) k -means algorithm. We compare our algorithm with the state of the art results (Yin et al., 2022) and observe that our algorithm achieves solutions of superior quality for most datasets.

3 PRELIMINARIES

Throughout this paper we work with ordered tuples rather than sets, denoted as $Y = (y_i)_{i \in [\ell]}$, where $[\ell] = \{1, \dots, \ell\}$. To reference the i -th element we either write y_i or $Y[i]$. It will be useful to use set notations for tuples such as $x \in Y \iff \exists i \in [\ell], x = y_i$ and $Y \subseteq Z \iff \forall i \in [\ell], y_i \in Z$. When summing we often write $\sum_{x \in Y} g(x)$ which is equivalent to $\sum_{i=1}^{\ell} g(Y[i])$.

We borrow the following notation from (Kanungo et al., 2004) and generalize it to Hilbert spaces. For every $x, y \in \mathcal{H}$ let $\Delta(x, y) = \|x - y\|^2$. We slightly abuse notation and also write $\Delta(x, y) = \|\phi(x) - \phi(y)\|^2$ when $x, y \in X$ and $\Delta(x, y) = \|\phi(x) - y\|^2$ when $x \in X, y \in \mathcal{H}$ (similarly when $x \in \mathcal{H}, y \in X$). For every finite tuple $S \subseteq X$ and a vector $x \in \mathcal{H}$ let $\Delta(S, x) = \sum_{y \in S} \Delta(y, x)$. Let

⁴The original paper of (Schwartzman, 2023) states the batch size as $\tilde{\Omega}((d/\epsilon)^2)$, however there is a mistake in the calculations which requires an additional k factor. We explain the issue in the proof of Lemma B.12.

us denote $\gamma = \max_{x \in X} \|\phi(x)\|$. Let us define for any finite tuple $S \subseteq X$ the center of mass of the tuple as $cm(S) = \frac{1}{|S|} \sum_{x \in S} \phi(x)$.

Kernel k -means We are given an input $X = (x_i)_{i=1}^n$ and a parameter k . Our goal is to (implicitly) find a tuple $\mathcal{C} \subseteq \mathcal{H}$ of k centers such that the following goal function is minimized: $\frac{1}{n} \sum_{x \in X} \min_{C \in \mathcal{C}} \Delta(x, C)$.

Let us define for every $x \in X$ the function $f_x : \mathcal{H}^k \rightarrow \mathbb{R}$ where $f_x(\mathcal{C}) = \min_{C \in \mathcal{C}} \Delta(x, C)$. We can treat \mathcal{H}^k as the set of k -tuples of vectors in \mathcal{H} . We also define the following function for every tuple $A = (a_i)_{i=1}^\ell \subseteq X$: $f_A(\mathcal{C}) = \frac{1}{\ell} \sum_{i=1}^\ell f_{a_i}(\mathcal{C})$. Note that f_X is our original goal function.

We make extensive use of the notion of *convex combination*:

Definition 3.1. We say that $y \in \mathcal{H}$ is a *convex combination* of X if $y = \sum_{x \in X} p_x \phi(x)$, such that $\forall x \in X, p_x \geq 0$ and $\sum_{x \in X} p_x = 1$.

4 OUR ALGORITHM

We start by presenting a slower algorithm that will set the stage for our truncated mini-batch algorithm and will be useful during the analysis. We present our pseudo-code in Algorithm 1. It requires an initial set of cluster centers such that every center is a convex combination of X . This guarantees that all subsequent centers are also a convex combination of X . Note that if we initialize the centers using the kernel version of k -means++, this is indeed the case.

Algorithm 1 proceeds by repeatedly sampling a batch of size b (the batch size is a parameter). For the i -th batch the algorithm (implicitly) updates the centers using the learning rate α_j^i for center j . Note that the learning rate may take on different values for different centers, and may change between iterations. Finally, the algorithm terminates when the progress on the batch is below ϵ , a user provided parameter. While our termination guarantees (Appendix B) require a specific learning rate, it does not affect the running time of a single iteration, and we leave it as a parameter for now.

Input: Dataset $X = (x_i)_{i=1}^n$, batch size b , early stopping parameter ϵ . Initial centers $(\mathcal{C}_1^j)_{j=1}^k$ where \mathcal{C}_1^j is a convex combination of X for all $j \in [k]$.

for $i = 1$ **to** ∞ **do**

Sample b elements, $B_i = (y_1, \dots, y_b)$, uniformly at random from X (with repetitions)

for $j = 1$ **to** k **do**

$B_i^j = \{x \in B_i \mid \arg \min_{\ell \in [k]} \Delta(x, \mathcal{C}_i^\ell) = j\}$

$\alpha_i^j = \sqrt{|B_i^j|} / b$ is the learning rate for the j -th cluster in iteration i

$\mathcal{C}_{i+1}^j = (1 - \alpha_i^j) \mathcal{C}_i^j + \alpha_i^j \cdot cm(B_i^j)$

end for

if $f_{B_i}(\mathcal{C}_{i+1}) - f_{B_i}(\mathcal{C}_i) < \epsilon$ **return** \mathcal{C}_{i+1}

end for

Algorithm 1: Mini-batch kernel k -means with early stopping

Recursive distance update rule Unlike k -means, the center updates and assignment of points to clusters is tricky for kernel k -means and even harder for mini-batch kernel k -means. Specifically, how do we overcome the challenge that we do not maintain the centers explicitly?

To assign points to centers in the $(i + 1)$ -th iteration, it is sufficient to know $\|\phi(x) - \mathcal{C}_{i+1}^j\|^2$ for every j . This is because we are interested in the closest center to x in kernel space. If we can keep track of this quantity through the execution of the algorithm, we are done. Let us derive a *recursive* expression for the distances: $\|\phi(x) - \mathcal{C}_{i+1}^j\|^2 = \langle \phi(x), \phi(x) \rangle - 2\langle \phi(x), \mathcal{C}_{i+1}^j \rangle + \langle \mathcal{C}_{i+1}^j, \mathcal{C}_{i+1}^j \rangle$.

Let us expand $\langle \phi(x), \mathcal{C}_{i+1}^j \rangle$ and $\langle \mathcal{C}_{i+1}^j, \mathcal{C}_{i+1}^j \rangle$:

$$\begin{aligned} \langle \phi(x), \mathcal{C}_{i+1}^j \rangle &= \langle \phi(x), (1 - \alpha_i^j) \mathcal{C}_i^j + \alpha_i^j \text{cm}(B_i^j) \rangle = (1 - \alpha_i^j) \langle \phi(x), \mathcal{C}_i^j \rangle + \alpha_i^j \langle \phi(x), \text{cm}(B_i^j) \rangle. \\ \langle \mathcal{C}_{i+1}^j, \mathcal{C}_{i+1}^j \rangle &= \langle (1 - \alpha_i^j) \mathcal{C}_i^j + \alpha_i^j \text{cm}(B_i^j), (1 - \alpha_i^j) \mathcal{C}_i^j + \alpha_i^j \text{cm}(B_i^j) \rangle \\ &= (1 - \alpha_i^j)^2 \langle \mathcal{C}_i^j, \mathcal{C}_i^j \rangle + 2\alpha_i^j (1 - \alpha_i^j) \langle \mathcal{C}_i^j, \text{cm}(B_i^j) \rangle + (\alpha_i^j)^2 \langle \text{cm}(B_i^j), \text{cm}(B_i^j) \rangle. \end{aligned}$$

The above is all we need to compute the distances. Furthermore, it is possible to use dynamic programming to update the center for every iteration in $O(n(b+k))$ time and $O(nk)$ space (proof deferred to Appendix A). This is a considerable speedup compared to the best known quadratic update time. Next, we go a step further and show that it is possible to get an update time with only polylogarithmic dependence on n .

4.1 TRUNCATING THE CENTERS

The issue with the above approach is that each center is written as a linear combination of potentially all points in X . We now present a simple way to overcome this issue. We maintain \mathcal{C}_{i+1}^j as an explicit sparse linear combination of X . Let us expand the recursive expression of \mathcal{C}_{i+1}^j for t terms, assuming $t < i$:

$$\mathcal{C}_{i+1}^j = (1 - \alpha_i^j) \mathcal{C}_i^j + \alpha_i^j \text{cm}(B_i^j) = \mathcal{C}_{i-t}^j \prod_{\ell=0}^t (1 - \alpha_{i-\ell}^j) + \sum_{\ell=0}^t \alpha_{i-\ell}^j \text{cm}(B_{i-\ell}^j) \prod_{z=i-\ell+1}^i (1 - \alpha_z^j).$$

The idea behind our truncation technique is that when t is sufficiently large, the term $\mathcal{C}_{i-t}^j \prod_{\ell=0}^t (1 - \alpha_{i-\ell}^j)$ becomes very small and can be discarded. The rate by which this term decays depends on the learning rates, which in turn depend on the number of elements assigned to the cluster in each of the previous iterations.

Let us start with some definitions. Let us denote $b_i^j = |B_i^j|$. We would like to trim the recursive expression such that every cluster center is represented using about τ points, where τ is a parameter to be set later. We define Q_i^j to be the set of indices from i to $i-t$, where t is the smallest integer such that $\sum_{\ell \in Q_i^j} b_\ell^j \geq \tau$ holds. If no such integer exists then $Q_i^j = \{i, i-1, \dots, 1\}$. It is the case that $\sum_{\ell \in Q_i^j} b_\ell^j \leq \tau + b$. Intuitively, Q_i^j is the most recent window of updates to cluster j that contains enough points (at least τ) to serve as a sufficient approximation of the current cluster center.

Next we define the *truncated centers*, for which the contributions of older points to the centers are forgotten after about τ points have been assigned to the center:

$$\widehat{\mathcal{C}}_{i+1}^j = \begin{cases} \sum_{\ell \in Q_i^j} \alpha_\ell^j \text{cm}(B_\ell^j) \prod_{\ell \in Q_i^j \setminus \{i\}} (1 - \alpha_\ell^j), & \min Q_i^j > 1 \\ \mathcal{C}_{i+1}^j & \text{otherwise.} \end{cases} \quad (1)$$

From the above definition it is always the case that either $\widehat{\mathcal{C}}_{i+1}^j = \mathcal{C}_{i+1}^j$ or $\sum_{\ell \in Q_i^j} b_\ell^j \geq \tau$. The following lemma shows that when τ is sufficiently large $\|\widehat{\mathcal{C}}_{i+1}^j - \mathcal{C}_{i+1}^j\|$ is small. Intuitively, this implies that the truncated algorithm should achieve results similar to the untruncated version (we formalize this intuition in Appendix B).

Lemma 4.1. *Setting $\tau = \lceil b \ln^2(28\gamma/\epsilon) \rceil$ it holds that $\forall i \in \mathbb{N}, j \in [k], \|\widehat{\mathcal{C}}_{i+1}^j - \mathcal{C}_{i+1}^j\| \leq \epsilon/28$.*

Proof. We assume that $\sum_{\ell \in Q_i^j} b_\ell^j \geq \tau$, as otherwise the claim trivially holds.

$$\|\widehat{\mathcal{C}}_{i+1}^j - \mathcal{C}_{i+1}^j\| = \|\mathcal{C}_{\min\{Q_i^j\}}^j \prod_{\ell \in Q_i^j} (1 - \alpha_\ell^j)\| \leq \|\mathcal{C}_{\min\{Q_i^j\}}^j\| e^{-\sum_{\ell \in Q_i^j} \alpha_\ell^j}$$

We have $\sum_{\ell \in Q_i^j} \alpha_\ell^j = \sum_{\ell \in Q_i^j} \sqrt{b_\ell^j/b} \geq \sqrt{\sum_{\ell \in Q_i^j} b_\ell^j/b} \geq \sqrt{\tau/b} \geq \ln(28\gamma/\epsilon)$. Plugging this back into the exponent, we get that: $\|\mathcal{C}_{\min\{Q_i^j\}}^j\| e^{-\sum_{\ell \in Q_i^j} \alpha_\ell^j} \leq \gamma e^{\ln(\epsilon/28\gamma)} \leq \epsilon/28$. \square

Algorithm implementation and runtime To implement this, we simply need to swap \mathcal{C}_i^j in Algorithm 1 with $\widehat{\mathcal{C}}_i^j$ (Lines 7 and 8). As before, the main bottleneck of each iteration is assigning points in the batch to their closest center. Once this is done, updating the truncated centers is straightforward by simply adjusting the coefficients in equation 1, removing the last element from the sum and adding a new element to the sum⁵. If $\min \{Q_i^j\}$ is 1, then we also need to add $\mathcal{C}_1^j \prod_{\ell \in Q_i^j} (1 - \alpha_\ell^j)$ which guarantees that $\widehat{\mathcal{C}}_i^j = \mathcal{C}_i^j$. The pseudo code is provided in Algorithm 2.

Input: Dataset $X = (x_i)_{i=1}^n$, batch size b , early stopping parameter ϵ . Initial centers $(\mathcal{C}_1^j)_{j=1}^k$ where \mathcal{C}_1^j is a convex combination of X and $\widehat{\mathcal{C}}_1^j = \mathcal{C}_1^j$ for all $j \in [k]$.

for $i = 1$ **to** ∞ **do**
 Sample b elements, $B_i = (y_1, \dots, y_b)$, uniformly at random from X (with repetitions)
 for $j = 1$ **to** k **do**
 $B_i^j = \{x \in B_i \mid \arg \min_{\ell \in [k]} \Delta(x, \widehat{\mathcal{C}}_i^\ell) = j\}$
 α_i^j is the learning rate for the j -th cluster in iteration i
 $\widehat{\mathcal{C}}_{i+1}^j = \sum_{\ell \in Q_i^j} \alpha_\ell^j \cdot cm(B_\ell^j) \prod_{\ell \in Q_i^j} (1 - \alpha_\ell^j)$
 if $\min \{Q_i^j\} = 1$ **then**
 $\widehat{\mathcal{C}}_{i+1}^j = \widehat{\mathcal{C}}_{i+1}^j + \mathcal{C}_1^j \prod_{\ell \in Q_i^j \setminus \{j\}} (1 - \alpha_\ell^j)$
 end if
 end for
 if $f_{B_i}(\widehat{\mathcal{C}}_{i+1}) - f_{B_i}(\widehat{\mathcal{C}}_i) < \epsilon$ **then**
 Return: $\widehat{\mathcal{C}}_{i+1}$
 end if
end for

Algorithm 2: Truncated Mini-batch kernel k -means with early stopping

As before, let us consider assigning all points in the $(i + 1)$ iteration to their closest centers. Unlike the previous approach, when computing distances between points in B_{i+1} and $\widehat{\mathcal{C}}_{i+1}$ we can do this directly (without recursion) and it is now sufficient to consider a much smaller set of inner products.

As before, the terms we are interested in computing are: $\langle \phi(x), \widehat{\mathcal{C}}_{i+1}^j \rangle$ and $\langle \widehat{\mathcal{C}}_{i+1}^j, \widehat{\mathcal{C}}_{i+1}^j \rangle$. However, there are several differences to the previous approach. We no longer need $\langle \phi(x), \widehat{\mathcal{C}}_{i+1}^j \rangle$ for all $x \in X$, but only for $x \in B_{i+1}$. Furthermore, $\widehat{\mathcal{C}}_{i+1}^j$ can be simply written as a weighted sum of at most $\sum_{\ell \in Q_i^j} b_\ell^j \leq \tau + b$ terms. Summing over all element in B_{i+1} and k centers we get $O(kb(b + \tau))$ time to compute $\langle \phi(x), \widehat{\mathcal{C}}_{i+1}^j \rangle$. For $\langle \widehat{\mathcal{C}}_{i+1}^j, \widehat{\mathcal{C}}_{i+1}^j \rangle$ using the bound on the number of terms we directly get $O(k(\tau + b)^2)$ time. We conclude that every iteration of Algorithm 2 requires $O(k(\tau + b)^2) = \widetilde{O}(kb^2)$ time. The additional space required is $O(k\tau) = \widetilde{O}(kb)$.

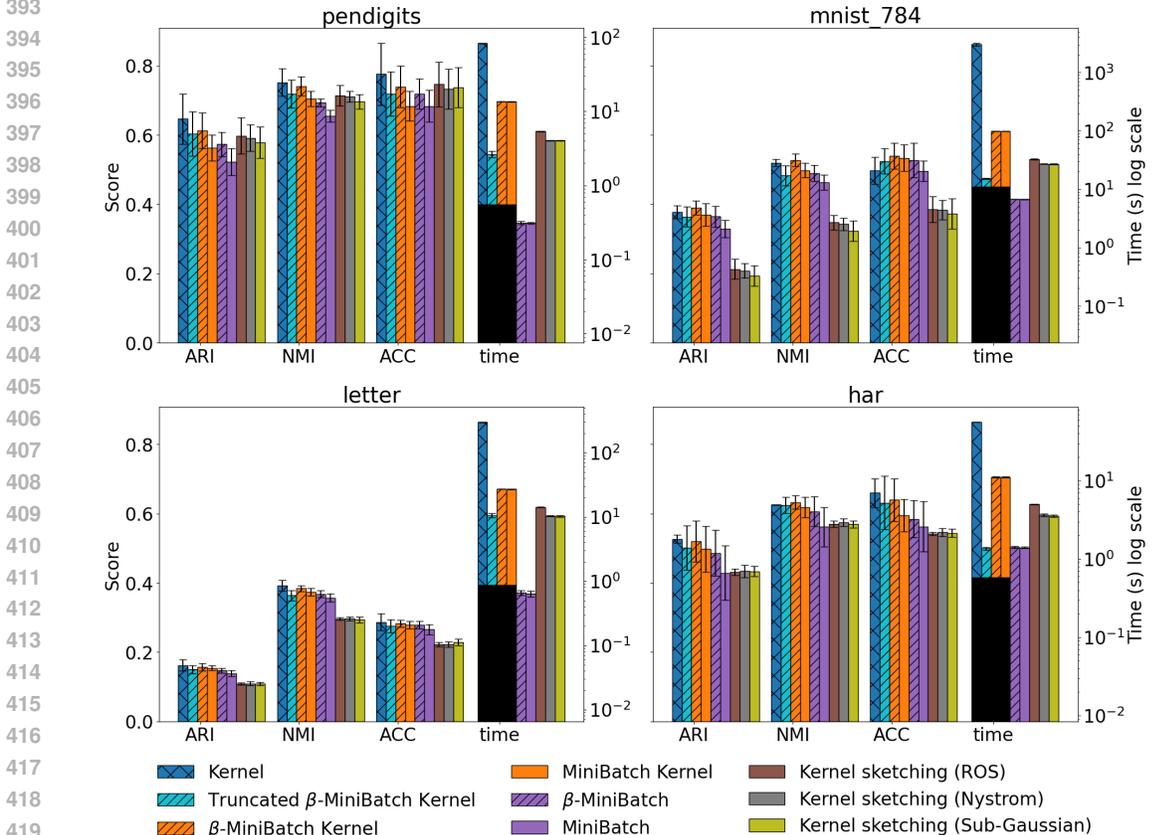
5 EXPERIMENTS

We evaluate our algorithms on the following datasets:

MNIST: The MNIST dataset (LeCun, 1998) has 70,000 grayscale images of handwritten digits (0 to 9), each image being 28x28 pixels. When flattened, this gives 784 features. **PenDigits:** The PenDigits dataset (Alpaydin & Alimoglu, 1998) has 10992 instances, each represented by an 16-dimensional vector derived from 2D pen movements. The dataset has 10 labelled clusters, one for each digit. **Letters:** The Letters dataset (Slate, 1991) has 20,000 instances of letters from ‘A’ to ‘Z’, each represented by 16 features. The dataset has 26 labelled clusters, one for each letter. **HAR:** The HAR dataset (Anguita et al., 2013) has 10,299 instances collected from smartphone sensors, capturing human activities like walking, sitting, and standing. Each instance is described by 561 features. It has 6 labeled clusters, corresponding to different types of physical activities.

⁵In our code we use an efficient sliding window implementation to store and update the coefficients representing each cluster center.

378 We compare the following algorithms: full-batch kernel k -means, truncated mini-batch kernel k -
 379 means, and mini-batch k -means (both kernel and non-kernel) with learning rates from (Schwartzman,
 380 2023) and sklearn. We also implement the three kernel sketching algorithms of Yin et al (Yin et al.,
 381 2022) that use either sub-Gaussian, randomized orthogonal system (ROS), or Nyström sketches.
 382 After sketching, we run k -means. We set the dimension of the sketch to 150, the same as in the
 383 experiments of (Yin et al., 2022). We evaluate our results with batch sizes: 2048, 1024, 512, 256
 384 and $\tau : 50, 100, 200, 300$. We execute every algorithm for 200 iterations. For the results below, we
 385 apply the Gaussian kernel: $K(x, y) = e^{-\|x-y\|^2/\kappa}$, where the κ parameter is set using the heuristic
 386 of (Wang et al., 2019) followed by some manual tuning (exact values appear in the supplementary
 387 materials). We also run experiments with heat and knn kernels in Appendix C. We repeat every
 388 experiment 10 times and present the average Adjusted Rand Index (ARI) (Gates & Ahn, 2017; Rand,
 389 1971), Normalized Mutual Information (NMI) (Lancichinetti et al., 2009) and Accuracy (ACC)⁶
 390 scores for every dataset. All experiments were conducted using an AMD Ryzen 9 7950X CPU with
 391 128GB of RAM and a Nvidia GeForce RTX 4090 GPU. We present partial results in Figure 3 and the
 392 full results in Appendix C. Error bars in the plot measure the standard deviation.



421 Figure 3: Our results for a batch size of size 1024 and $\tau = 200$ using the Gaussian kernel. We use the
 422 β prefix to denote that the algorithm uses the learning rate of (Schwartzman, 2023). Black denotes
 423 the time required to compute the kernel.

425 **Discussion** Throughout our results, we consistently observe that the truncated algorithm achieves
 426 performance on par with the non-truncated version with a running time which is often an order of
 427 magnitude faster. Surprisingly, this often holds for tiny values of τ (e.g., 50) far below the theoretical
 428 threshold (i.e., $\tau \ll b$). We also achieve considerably better quality solutions on most datasets
 429 compared to kernel sketching. We believe that our approach achieves a good balance between speed
 430 and performance, and is a valuable addition to the tool-box of clustering algorithms.

431 ⁶We use the Hungarian algorithm to match labels to clusters such that the accuracy is maximized.

REFERENCES

- 432
433
434 E. Alpaydin and Fevzi. Alimoglu. Pen-Based Recognition of Handwritten Digits. UCI Machine
435 Learning Repository, 1998. DOI: <https://doi.org/10.24432/C5MG6K>. License: CC BY 4.0 DEED,
436 available at <https://creativecommons.org/licenses/by/4.0/>.
- 437 Davide Anguita, Alessandro Ghio, Luca Oneto, Xavier Parra, Jorge Luis Reyes-Ortiz, et al. A public
438 domain dataset for human activity recognition using smartphones. In *Esann*, volume 3, pp. 3,
439 2013. License: CC BY-NC-SA 4.0 DEED, available at [https://creativecommons.org/
440 licenses/by-nc-sa/4.0/](https://creativecommons.org/licenses/by-nc-sa/4.0/).
- 441 Pablo Arbelaez, Michael Maire, Charless Fowlkes, and Jitendra Malik. Contour detection and
442 hierarchical image segmentation. *IEEE transactions on pattern analysis and machine intelligence*,
443 33(5):898–916, 2010.
- 444 David Arthur and Sergei Vassilvitskii. How slow is the k -means method? In *SCG*, pp. 144–153.
445 ACM, 2006.
- 446 David Arthur and Sergei Vassilvitskii. k -means++: the advantages of careful seeding. In *SODA*, pp.
447 1027–1035. SIAM, 2007.
- 448 Artem Barger and Dan Feldman. Deterministic coresets for k -means of big sparse data. *Algorithms*,
449 13(4):92, 2020.
- 450 Di Chen and Jeff M Phillips. Relative error embeddings of the gaussian kernel distance. In
451 *International Conference on Algorithmic Learning Theory*, pp. 560–576. PMLR, 2017.
- 452 Radha Chitta, Rong Jin, Timothy C. Havens, and Anil K. Jain. Approximate kernel k -means: solution
453 to large scale kernel clustering. In *KDD*, pp. 895–903. ACM, 2011.
- 454 Radha Chitta, Rong Jin, and Anil K Jain. Efficient kernel clustering using random fourier features. In
455 *2012 IEEE 12th International Conference on Data Mining*, pp. 161–170. IEEE, 2012.
- 456 Fan RK Chung. *Spectral graph theory*, volume 92. American Mathematical Soc., 1997.
- 457 Inderjit S. Dhillon, Yuqiang Guan, and Brian Kulis. Kernel k -means: spectral clustering and normal-
458 ized cuts. In *Proceedings of the Tenth ACM SIGKDD International Conference on Knowledge
459 Discovery and Data Mining, KDD '04*, pp. 551–556, New York, NY, USA, 2004. Association
460 for Computing Machinery. ISBN 1581138881. doi: 10.1145/1014052.1014118. URL
461 <https://doi.org/10.1145/1014052.1014118>.
- 462 Dan Feldman, Melanie Schmidt, and Christian Sohler. Turning big data into tiny data: Constant-size
463 coresets for k -means, pca, and projective clustering. *SIAM Journal on Computing*, 49(3):601–657,
464 2020.
- 465 Alexander J Gates and Yong-Yeol Ahn. The impact of random models on clustering similarity.
466 *Journal of Machine Learning Research*, 18(87):1–28, 2017.
- 467 Wassily Hoeffding. Probability inequalities for sums of bounded random variables. *Journal of the
468 American Statistical Association*, 58(301):13–30, 1963.
- 469 Shaofeng H.-C. Jiang, Robert Krauthgamer, Jianing Lou, and Yubo Zhang. Coresets for kernel
470 clustering. *CoRR*, abs/2110.02898, 2021.
- 471 Tapas Kanungo, David M. Mount, Nathan S. Netanyahu, Christine D. Piatko, Ruth Silverman, and
472 Angela Y. Wu. A local search approximation algorithm for k -means clustering. *Comput. Geom.*,
473 28(2-3):89–112, 2004.
- 474 Andrea Lancichinetti, Santo Fortunato, and János Kertész. Detecting the overlapping and hierarchical
475 community structure in complex networks. *New journal of physics*, 11(3):033015, 2009.
- 476 Yann LeCun. The mnist database of handwritten digits. <http://yann.lecun.com/exdb/mnist/>, 1998.
477 License: CC0 1.0 DEED CC0 1.0 Universal, available at [https://creativecommons.org/
478 publicdomain/zero/1.0/](https://creativecommons.org/publicdomain/zero/1.0/).

- 486 Yong Liu. Refined learning bounds for kernel and approximate k -means. *Advances in neural*
487 *information processing systems*, 34:6142–6154, 2021.
- 488
- 489 Stuart P. Lloyd. Least squares quantization in PCM. *IEEE Trans. Inf. Theory*, 28(2):129–136, 1982.
- 490
- 491 Peter Macgregor and He Sun. Fast approximation of similarity graphs with kernel density estimation.
492 *Advances in Neural Information Processing Systems*, 36, 2024.
- 493
- 494 Cameron Musco and Christopher Musco. Recursive sampling for the nystrom method. *Advances in*
495 *neural information processing systems*, 30, 2017.
- 496
- 497 Assaf Naor. On the banach-space-valued azuma inequality and small-set isoperimetry of alon-
498 roichman graphs. *Combinatorics, Probability and Computing*, 21(4):623–634, 2012.
- 499
- 500 F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Pretten-
501 hofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and
502 E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*,
503 12:2825–2830, 2011.
- 504
- 505 Kent Quanrud. *Spectral Sparsification of Metrics and Kernels*, pp. 1445–1464. doi: 10.
506 1137/1.9781611976465.87. URL [https://epubs.siam.org/doi/abs/10.1137/1.](https://epubs.siam.org/doi/abs/10.1137/1.9781611976465.87)
507 [9781611976465.87](https://epubs.siam.org/doi/abs/10.1137/1.9781611976465.87).
- 508
- 509 Ali Rahimi and Benjamin Recht. Random features for large-scale kernel machines. *Advances in*
510 *neural information processing systems*, 20, 2007.
- 511
- 512 William M Rand. Objective criteria for the evaluation of clustering methods. *Journal of the American*
513 *Statistical association*, 66(336):846–850, 1971.
- 514
- 515 Melanie Schmidt. Coresets and streaming algorithms for the k -means problem and related clustering
516 objectives. 2014.
- 517
- 518 Gregory Schwartzman. Mini-batch k -means terminates within $O(d/\epsilon)$ iterations. In *ICLR*, 2023.
- 519
- 520 D. Sculley. Web-scale k -means clustering. In *WWW*, pp. 1177–1178. ACM, 2010.
- 521
- 522 David Slate. Letter Recognition. UCI Machine Learning Repository, 1991. DOI:
523 <https://doi.org/10.24432/C5ZP40>. License: CC BY 4.0 DEED, available at [https://](https://creativecommons.org/licenses/by/4.0/)
524 creativecommons.org/licenses/by/4.0/.
- 525
- 526 Cheng Tang and Claire Monteleoni. Convergence rate of stochastic k -means. In *AISTATS*, volume 54
527 of *Proceedings of Machine Learning Research*, pp. 1495–1503. PMLR, 2017.
- 528
- 529 Andrea Vattani. k -means requires exponentially many iterations even in the plane. *Discret. Comput.*
530 *Geom.*, 45(4):596–616, 2011.
- 531
- 532 Shusen Wang, Alex Gittens, and Michael W Mahoney. Scalable kernel k -means clustering with
533 nystrom approximation: Relative-error bounds. *Journal of Machine Learning Research*, 20(12):
534 1–49, 2019.
- 535
- 536 Rong Yin, Yong Liu, Weiping Wang, and Dan Meng. Randomized sketches for clustering: Fast and
537 optimal kernel k -means. *Advances in Neural Information Processing Systems*, 35:6424–6436,
538 2022.
- 539

A OMITTED PROOFS AND ALGORITHMS FOR SECTION 4

Runtime analysis of Algorithm 1 Assuming that $\langle \mathcal{C}_i^j, \mathcal{C}_i^j \rangle$ and $\langle \phi(x), \mathcal{C}_i^j \rangle$ are known for all $j \in [k]$ and for all $x \in X$, we can compute $\langle \mathcal{C}_{i+1}^j, \mathcal{C}_{i+1}^j \rangle$ and $\langle \phi(x), \mathcal{C}_{i+1}^j \rangle$ for all $j \in [k]$ and $x \in X$, which implies we can compute the distances from any point in the batch to all centers.

We now bound the running time of a single iteration of the outer loop in Algorithm 1. Let us denote $b_i^j = |B_i^j|$ and recall that $cm(B_i^j) = \frac{1}{b_i^j} \sum_{y \in B_i^j} \phi(y)$. Therefore, computing $\langle \phi(x), cm(B_i^j) \rangle = \frac{1}{b_i^j} \sum_{y \in B_i^j} \langle \phi(x), \phi(y) \rangle$ requires $O(b_i^j)$ time. Similarly, computing $\langle cm(B_i^j), cm(B_i^j) \rangle$ requires $O((b_i^j)^2)$ time. Let us now bound the time it requires to compute $\langle \phi(x), \mathcal{C}_{i+1}^j \rangle$ and $\langle \mathcal{C}_{i+1}^j, \mathcal{C}_{i+1}^j \rangle$.

Assuming we know $\langle \phi(x), \mathcal{C}_i^j \rangle$ and $\langle \mathcal{C}_i^j, \mathcal{C}_i^j \rangle$, updating $\langle \phi(x), \mathcal{C}_{i+1}^j \rangle$ for all $x \in X, j \in [k]$ requires $O(n(b+k))$ time. Specifically, the $\langle \phi(x), \mathcal{C}_i^j \rangle$ term is already known from the previous iteration and we need to compute $\alpha_i^j \langle \phi(x), cm(B_i^j) \rangle$ for every $x \in X, j \in [k]$ which requires $n \sum_{j \in [k]} b_i^j = nb$ time. Finally, updating $\langle \phi(x), \mathcal{C}_{i+1}^j \rangle$ for all $x \in X, j \in [k]$ requires $O(nk)$ time.

Updating $\langle \mathcal{C}_{i+1}^j, \mathcal{C}_{i+1}^j \rangle$ requires $O(b^2 + kb)$ time. Specifically, $\langle \mathcal{C}_i^j, \mathcal{C}_i^j \rangle$ is known from the previous iteration and computing $\langle cm(B_i^j), cm(B_i^j) \rangle$ for all $j \in [k]$ requires $O(\sum_{j \in [k]} (b_i^j)^2) = O(b^2)$ time. Computing $\langle \mathcal{C}_i^j, cm(B_i^j) \rangle$ for all $j \in [k]$ requires time $O(b)$ using $\langle \phi(x), \mathcal{C}_i^j \rangle$ from the previous iteration. Therefore, the total running time of the update step (assigning points to new centers) is $O(n(b+k))$. To perform the update at the $(i+1)$ -th step we only need $\langle \phi(x), \mathcal{C}_i^j \rangle, \langle \mathcal{C}_i^j, \mathcal{C}_i^j \rangle$, which results in a space complexity of $O(nk)$. This completes the first claim of Theorem 1.1.

B TERMINATION GUARANTEE

In this section we prove the second claim of Theorem 1.1. For most of the section we analyze Algorithm 1, and towards the end we use the fact that the centers of the two algorithms are close throughout the execution to conclude our proof.

Section preliminaries We introduce the following definitions and lemmas to aid our proof of the second claim of Theorem 1.1.

Lemma B.1. *For every y which is a convex combination of X it holds that $\|y\| \leq \gamma$.*

Proof. The proof follows by the triangle inequality: $\|y\| = \|\sum_{x \in X} p_x \phi(x)\| \leq \sum_{x \in X} \|p_x \phi(x)\| = \sum_{x \in X} p_x \|\phi(x)\| \leq \sum_{x \in X} p_x \gamma = \gamma$. \square

Lemma B.2. *For any tuple of k centers $\mathcal{C} \subset \mathcal{H}^d$ which are a convex combination of points in X , it holds that $\forall A \subseteq X, f_A(\mathcal{C}) \leq 4\gamma^2$.*

Proof. It is sufficient to upper bound f_x . Combining that fact that every $C \in \mathcal{C}$ is a convex combination of X with the triangle inequality, we have that

$$\begin{aligned} \forall x \in X, f_x(\mathcal{C}) &\leq \max_{C \in \mathcal{C}} \Delta(x, C) = \Delta(x, \sum_{y \in X} p_y \phi(y)) \\ &= \|\phi(x) - \sum_{y \in X} p_y \phi(y)\|^2 \leq (\|\phi(x)\| + \|\sum_{y \in X} p_y \phi(y)\|)^2 \leq 4\gamma^2. \quad \square \end{aligned}$$

We state the following simplified version of an Azuma bound for Hilbert space valued martingales from (Naor, 2012), followed by a standard Hoeffding bound.

Theorem B.3 ((Naor, 2012)). *Let \mathcal{H} be a Hilbert space and let Y_0, \dots, Y_m be a \mathcal{H} -valued martingale, such that $\forall 1 \leq i \leq m, \|Y_i - Y_{i-1}\| \leq a_i$. It holds that $Pr[\|Y_m - Y_0\| \geq \delta] \leq e^{-\frac{\delta^2}{\sum_{i=1}^m a_i^2}}$.*

Theorem B.4 ((Hoeffding, 1963)). *Let Y_1, \dots, Y_m be independent random variables such that $\forall 1 \leq i \leq m, E[Y_i] = \mu$ and $Y_i \in [a_{min}, a_{max}]$. Then $Pr\left(\left|\frac{1}{m} \sum_{i=1}^m Y_k - \mu\right| \geq \delta\right) \leq 2e^{-2m\delta^2 / (a_{max} - a_{min})^2}$.*

The following lemma provides concentration guarantees when sampling a *batch*.

Lemma B.5. *Let B be a tuple of b elements chosen uniformly at random from X with repetitions. For any fixed tuple of k centers, $\mathcal{C} \subseteq \mathcal{H}$ which are a convex combination of X , it holds that:*
 $Pr[|f_B(\mathcal{C}) - f_X(\mathcal{C})| \geq \delta] \leq 2e^{-b\delta^2/8\gamma^4}$.

Proof. Let us write $B = (y_1, \dots, y_b)$, where y_i is a random element selected uniformly at random from X with repetitions. For every such y_i define the random variable $Z_i = f_{y_i}(\mathcal{C})$. These new random variables are IID for any fixed \mathcal{C} . It also holds that $\forall i \in [b], E[Z_i] = \frac{1}{n} \sum_{x \in X} f_x(\mathcal{C}) = f_X(\mathcal{C})$ and that $f_B(\mathcal{C}) = \frac{1}{b} \sum_{x \in B} f_x(\mathcal{C}) = \frac{1}{b} \sum_{i=1}^b Z_i$.

Applying the Hoeffding bound (Theorem B.4) with parameters $m = b, \mu = f_X(\mathcal{C}), a_{max} - a_{min} \leq 4\gamma^2$ (due to Lemma B.2) we get that: $Pr[|f_B(\mathcal{C}) - f_X(\mathcal{C})| \geq \delta] \leq 2e^{-b\delta^2/8\gamma^4}$. \square

For any tuple $S \subseteq X$ and some tuple of cluster centers $\mathcal{C} = (\mathcal{C}^\ell)_{\ell \in [k]} \subset \mathcal{H}$, \mathcal{C} implies a *partition* $(S^\ell)_{\ell \in [k]}$ of the points in S . Specifically, every S^ℓ contains the points in S closest to \mathcal{C}^ℓ (in \mathcal{H}) and every point in S belongs to a single \mathcal{C}^ℓ (ties are broken arbitrarily). We state the following useful observation:

Observation B.6. Fix some $A \subseteq X$. Let \mathcal{C} be a tuple of k centers, $S = (S^\ell)_{\ell \in [k]}$ be the partition of A induced by \mathcal{C} and $\bar{S} = (\bar{S}^\ell)_{\ell \in [k]}$ be any other partition of A . It holds that $\sum_{j=1}^k \Delta(S^j, \mathcal{C}^j) \leq \sum_{j=1}^k \Delta(\bar{S}^j, \mathcal{C}^j)$.

Recall that \mathcal{C}_i^j is the j -th center in the beginning of the i -th iteration of Algorithm 1 and $(B_i^\ell)_{\ell \in [k]}$ is the partition of B_i induced by \mathcal{C}_i . Let $(X_i^\ell)_{\ell \in [k]}$ be the partition of X induced by \mathcal{C}_i .

We now have the tools to analyze Algorithm 1 with the learning rate of (Schwartzman, 2023). Specifically, we assume that the algorithm executes for at least t iterations, the learning rate is $\alpha_i^j = \sqrt{b_i^j/b}$, where $b_i^j = |B_i^j|$, and the batch size is $b = \Omega(\max\{\gamma^4, \gamma^2\} k \epsilon^{-2} \log(nt))$. We show that the algorithm must terminate within $t = O(\gamma^2/\epsilon)$ steps w.h.p. Plugging t back into b , we get that a batch size of $b = \Omega(\max\{\gamma^4, \gamma^2\} k \epsilon^{-2} \log^2(\gamma n/\epsilon))$ is sufficient. We assume that ϵ is chosen such that $\gamma^2/\epsilon > 1/4$. Otherwise, the stopping condition immediately holds due to Lemma B.2.

Proof outline We note that when sampling a batch it holds w.h.p that $f_{B_i}(\mathcal{C}_i)$ is close to $f_{X_i}(\mathcal{C}_i)$ (Lemma B.5). This is due to the fact that B_i is sampled after \mathcal{C}_i is fixed. If we could show that $f_{B_i}(\mathcal{C}_{i+1})$ is close $f_{X_i}(\mathcal{C}_{i+1})$ then combined with the fact that we make progress of at least ϵ on the batch we can conclude that we make progress of at least some constant fraction of ϵ on the entire dataset.

Unfortunately, as \mathcal{C}_{i+1} depends on B_i , getting the above guarantee is tricky. To overcome this issue we define the auxiliary value $\bar{\mathcal{C}}_{i+1}^j = (1 - \alpha_i^j)\mathcal{C}_i^j + \alpha_i^j cm(X_i^j)$. This is the j -th center at step $i + 1$ if we were to use the entire dataset for the update, rather than just a batch. Note that this is only used in the analysis and not in the algorithm. Note that $\bar{\mathcal{C}}_{i+1}$ is *almost* independent of B_i . Every $\bar{\mathcal{C}}_{i+1}^j$ depends only on \mathcal{C}_i^j, X_i^j and α_i^j . While \mathcal{C}_i^j, X_i^j are independent of B_i , the learning α_i^j is *not*. Nevertheless, the number of possible values of $\{\alpha_i^j\}_{j \in [k]}$ is sufficiently small, and we can overcome this issue by showing concentration for every possible learning rate configuration followed by a union bound. This allows us to use $\bar{\mathcal{C}}_{i+1}$ instead of \mathcal{C}_{i+1} in the above analysis outline. We show that for our choice of learning rate it holds that $\bar{\mathcal{C}}_{i+1}, \mathcal{C}_{i+1}$ are sufficiently close, which implies that $f_X(\mathcal{C}_{i+1}), f_X(\bar{\mathcal{C}}_{i+1})$ and $f_{B_i}(\mathcal{C}_{i+1}), f_{B_i}(\bar{\mathcal{C}}_{i+1})$ are also sufficiently close. That is, $\bar{\mathcal{C}}_{i+1}$ acts as a proxy for \mathcal{C}_{i+1} . Combining everything together we get our desired result for Algorithm 1.

We start with the following useful observation, which will allow us to use Lemma B.1 to bound the norm of the centers by γ throughout the execution of the algorithm.

Observation B.7. If $\forall j \in [k], \mathcal{C}_1^j$ is a convex combination of X then $\forall i > 1, j \in [k], \mathcal{C}_i^j, \bar{\mathcal{C}}_i^j$ are also a convex combinations of X .

We state the following useful lemma. Although the original proof is for Euclidean spaces, it goes through for Hilbert spaces.

Lemma B.8 ((Kanungo et al., 2004)). *For any set $S \subseteq X$ and any $C \in \mathcal{H}$ it holds that $\Delta(S, C) = \Delta(S, cm(S)) + |S| \Delta(C, cm(S))$.*

Proof.

$$\begin{aligned}
\Delta(S, C) &= \sum_{x \in S} \Delta(x, C) = \sum_{x \in S} \langle x - C, x - C \rangle \\
&= \sum_{x \in S} \langle (x - cm(S)) + (cm(S) - C), (x - cm(S)) + (cm(S) - C) \rangle \\
&= \sum_{x \in S} \Delta(x, cm(S)) + \Delta(C, cm(S)) + 2 \langle x - cm(S), cm(S) - C \rangle \\
&= \Delta(S, cm(S)) + |S| \Delta(C, cm(S)) + \sum_{x \in S} 2 \langle x - cm(S), cm(S) - C \rangle \\
&= \Delta(S, cm(S)) + |S| \Delta(C, cm(S)),
\end{aligned}$$

where the last step is due to the fact that

$$\begin{aligned}
\sum_{x \in S} \langle x - cm(S), cm(S) - C \rangle &= \langle \sum_{x \in S} x - |S| cm(S), cm(S) - C \rangle \\
&= \langle \sum_{x \in S} x - \frac{|S|}{|S|} \sum_{x \in S} x, cm(S) - C \rangle = 0.
\end{aligned}$$

□

We use the above to state the following useful lemma.

Lemma B.9. *For any $S \subseteq X$ and $C, C' \in \mathcal{H}$ which are convex combinations of X , it holds that: $|\Delta(S, C') - \Delta(S, C)| \leq 4\gamma |S| \|C - C'\|$.*

Proof. Using Lemma B.8 we get that $\Delta(S, C) = \Delta(S, cm(S)) + |S| \Delta(cm(S), C)$ and that $\Delta(S, C') = \Delta(S, cm(S)) + |S| \Delta(cm(S), C')$. Thus, it holds that $|\Delta(S, C') - \Delta(S, C)| = |S| \cdot |\Delta(cm(S), C') - \Delta(cm(S), C)|$. Let us write

$$\begin{aligned}
&|\Delta(cm(S), C') - \Delta(cm(S), C)| \\
&= |\langle cm(S) - C', cm(S) - C' \rangle - \langle cm(S) - C, cm(S) - C \rangle| \\
&= |-2 \langle cm(S), C' \rangle + \langle C', C' \rangle + 2 \langle cm(S), C \rangle - \langle C, C \rangle| \\
&= |2 \langle cm(S), C - C' \rangle + \langle C' - C, C' + C \rangle| \\
&= |\langle C - C', 2cm(S) - (C' + C) \rangle| \\
&\leq \|C - C'\| \|2cm(S) - (C' + C)\| \leq 4\gamma \|C - C'\|.
\end{aligned}$$

Where in the last transition we used the Cauchy-Schwartz inequality, the triangle inequality, and the fact that $C, C', cm(S)$ are convex combinations of X and therefore their norm is bounded by γ . □

When centers are sufficiently close, these lemmas imply their values are close for any f_A .

Lemma B.10. *Fix some $A \subseteq X$ and let $(C^j)_{j \in [k]}, (\bar{C}^j)_{j \in [k]} \subset \mathcal{H}$ be arbitrary centers such that $\forall j \in [k], \|C^j - \bar{C}^j\| \leq \epsilon/28\gamma$. It holds that $\forall i \in [t], |f_A(\bar{C}_{i+1}) - f_A(C_{i+1})| \leq \epsilon/7$.*

Proof. Let $S = (S^\ell)_{\ell \in [k]}$, $\bar{S} = (\bar{S}^\ell)_{\ell \in [k]}$ be the partitions induced by $\mathcal{C}, \bar{\mathcal{C}}$ on A . Let us expand the expression

$$\begin{aligned} f_A(\bar{\mathcal{C}}) - f_A(\mathcal{C}) &= \frac{1}{|A|} \sum_{j=1}^k \Delta(\bar{S}^j, \bar{\mathcal{C}}^j) - \Delta(S^j, \mathcal{C}^j) \\ &\leq \frac{1}{|A|} \sum_{j=1}^k \Delta(S^j, \bar{\mathcal{C}}^j) - \Delta(S^j, \mathcal{C}^j) \leq \frac{1}{|A|} \sum_{j=1}^k 4\gamma |S^j| \|\bar{\mathcal{C}}^j - \mathcal{C}^j\| \leq \frac{1}{|A|} \sum_{j=1}^k |S^j| \epsilon/7 = \epsilon/7. \end{aligned}$$

The first inequality is due to Observation B.6, the second is due Lemma B.9 and finally we use the assumption about the distances between centers together with the fact that $\sum_{j=1}^k |S^j| = |A|$. Using the same argument we also get that $f_A(\mathcal{C}) - f_A(\bar{\mathcal{C}}) \leq \epsilon/7$, which completes the proof. \square

Now we show that due to our choice of learning rate, \mathcal{C}_{i+1}^j and $\bar{\mathcal{C}}_{i+1}^j$ are sufficiently close.

Lemma B.11. *It holds w.h.p that $\forall i \in [t], j \in [k], \|\mathcal{C}_{i+1}^j - \bar{\mathcal{C}}_{i+1}^j\| \leq \frac{\epsilon}{28\gamma}$.*

Proof. Note that $\mathcal{C}_{i+1}^j - \bar{\mathcal{C}}_{i+1}^j = \alpha_i^j (cm(B_i^j) - cm(X_i^j))$. Let us fix some iteration i and center j . To simplify notation, let us denote: $X' = X_i^j, B' = B_i^j, b' = b_i^j, \alpha' = \alpha_i^j$. Although b' is a random variable, in what follows we treat it as a fixed value (essentially conditioning on its value). As what follows holds for *all* values of b' it also holds without conditioning due to the law of total probabilities.

For the rest of the proof, we assume $b' > 0$ (if $b' = 0$ the claim holds trivially). Let us denote by $\{Y_\ell\}_{\ell=1}^{b'}$ the sampled points in B' . Note that a randomly sampled element from X is in B' if and only if it is in X' . As batch elements are sampled uniformly at random with repetitions from X , conditioning on the fact that an element is in B' means that it is distributed uniformly over X' . Note that $\forall \ell, E[\phi(Y_\ell)] = \frac{1}{|X'|} \sum_{x \in X'} \phi(x) = cm(X')$ and $E[cm(B')] = \frac{1}{b'} \sum_{\ell=1}^{b'} E[\phi(Y_\ell)] = cm(X')$.

Let us define the following martingale: $Z_r = \sum_{\ell=1}^r (\phi(Y_\ell) - E[\phi(Y_\ell)])$. Note that $Z_0 = 0$, and when $r > 0$, $Z_r = \sum_{\ell=1}^r \phi(Y_\ell) - r \cdot cm(X')$. It is easy to see that this is a martingale:

$$E[Z_r | Z_{r-1}] = E\left[\sum_{\ell=1}^r \phi(Y_\ell) - r \cdot cm(X') \mid Z_{r-1}\right] = Z_{r-1} + E[\phi(Y_r) - cm(X') \mid Z_{r-1}] = Z_{r-1}.$$

We bound the differences: $\|Z_r - Z_{r-1}\| = \|\phi(Y_r) - cm(X')\| \leq \|\phi(Y_r)\| + \|cm(X')\| \leq 2\gamma$.

Now we may use Azuma's inequality: $Pr[\|Z_{b'} - Z_0\| \geq \delta] \leq e^{-\Theta(\frac{\delta^2}{\gamma^2 b'})}$. Let us now divide both sides of the inequality by b' and set $\delta = \frac{b' \epsilon}{28\gamma \alpha'}$. We get $Pr[\|cm(B') - cm(X')\| \geq \frac{\epsilon}{28\gamma \alpha'}] = Pr[\|\frac{1}{b'} \sum_{\ell=1}^{b'} \phi(Y_\ell) - cm(X')\| \geq \frac{\epsilon}{28\gamma \alpha'}] \leq e^{-\Theta(\frac{b' \epsilon^2}{\gamma \alpha'^2})}$. Using the fact that $\alpha' = \sqrt{b'/b}$ together with the fact that $b = \Omega(\max\{\gamma^4, \gamma^2\} k \epsilon^{-2} \log(nt))$ (for an appropriate constant) we get that the above is $O(1/ntk)$. Finally, taking a union bound over all t iterations and all k centers per iteration completes the proof. \square

Let us state the following useful lemma.

Lemma B.12. *It holds w.h.p that for every $i \in [t]$,*

$$f_X(\bar{\mathcal{C}}_{i+1}) - f_X(\mathcal{C}_{i+1}) \geq -\epsilon/7 \quad (2)$$

$$f_{B_i}(\mathcal{C}_{i+1}) - f_{B_i}(\bar{\mathcal{C}}_{i+1}) \geq -\epsilon/7 \quad (3)$$

$$f_X(\mathcal{C}_i) - f_{B_i}(\mathcal{C}_i) \geq -\epsilon/7 \quad (4)$$

$$f_{B_i}(\bar{\mathcal{C}}_{i+1}) - f_X(\bar{\mathcal{C}}_{i+1}) \geq -\epsilon/7 \quad (5)$$

Proof. The first two inequalities follow from Lemma B.10. The third is due to Lemma B.5 by setting $\delta = \epsilon/7, B = B_i$:

$$\begin{aligned} \Pr[|f_{B_i}(\mathcal{C}_i) - f_X(\mathcal{C}_i)| \geq \delta] &\leq 2e^{-b\delta^2/8\gamma^4} = e^{-\Theta(b\epsilon^2/\gamma^4)} \\ &= e^{-\Omega(\log(nt))} = O(1/nt). \end{aligned}$$

The last inequality is a bit more involved⁷. Let $\vec{\ell} \in \mathbb{N}^k$ be a vector whose entries sum to b . For every $\vec{\ell}$ we can define $\bar{\mathcal{C}}_{i+1}(\vec{\ell})$ such that $\bar{\mathcal{C}}_{i+1}^j(\vec{\ell}) = \mathcal{C}_i^j(1 - \sqrt{\ell_j/b}) + \sqrt{\ell_j/b} \cdot \text{cm}(X_i^j)$. For every choice of $\vec{\ell}$ it holds that $\bar{\mathcal{C}}_{i+1}(\vec{\ell})$ is independent of B_i and we can apply Lemma B.5 for every possible $\bar{\mathcal{C}}_{i+1}(\vec{\ell})$ by setting $\delta = \epsilon/7, B = B_i$

$$\begin{aligned} \Pr[|f_{B_i}(\bar{\mathcal{C}}_{i+1}(\vec{\ell})) - f_X(\bar{\mathcal{C}}_{i+1}(\vec{\ell}))| \geq \delta] &\leq 2e^{-b\delta^2/8\gamma^4} \\ &= e^{-\Theta(b\epsilon^2/\gamma^4)} = e^{-\Omega(k \log(nt))} = O(1/(nt)^k), \end{aligned}$$

where last inequality is due to the fact that $b = \Omega(\max\{\gamma^4, \gamma^2\} k\epsilon^{-2} \log(nt))$ (for an appropriate constant). Finally, we take a union bound over all possible vectors $\vec{\ell}$, a total of $\binom{b+k-1}{k-1} \leq \left(\frac{(b+k-1)\cdot e}{k-1}\right)^{k-1} = O(n^{k-1})$. As $\bar{\mathcal{C}}_{i+1}$ corresponds to at least one $\bar{\mathcal{C}}_{i+1}(\vec{\ell})$ we are done.

Taking a union bound over t iterations, we obtain the result. \square

Putting everything together We wish to lower bound $f_X(\mathcal{C}_i) - f_X(\mathcal{C}_{i+1})$. We write the following, where the \pm notation means we add and subtract a term:

$$\begin{aligned} f_X(\mathcal{C}_i) - f_X(\mathcal{C}_{i+1}) &= f_X(\mathcal{C}_i) \pm f_{B_i}(\mathcal{C}_i) - f_X(\mathcal{C}_{i+1}) \\ &\geq f_{B_i}(\mathcal{C}_i) - f_X(\mathcal{C}_{i+1}) - \epsilon/7 \\ &= f_{B_i}(\mathcal{C}_i) \pm f_{B_i}(\mathcal{C}_{i+1}) - f_X(\mathcal{C}_{i+1}) - \epsilon/7 \\ &\geq f_{B_i}(\mathcal{C}_{i+1}) - f_X(\mathcal{C}_{i+1}) + 6\epsilon/7 \\ &= f_{B_i}(\mathcal{C}_{i+1}) \pm f_{B_i}(\bar{\mathcal{C}}_{i+1}) \\ &\quad \pm f_X(\bar{\mathcal{C}}_{i+1}) - f_X(\mathcal{C}_{i+1}) + 6\epsilon/7 \geq 3\epsilon/7. \end{aligned}$$

Where the first inequality is due to inequality 4 in Lemma B.12 ($f_X(\mathcal{C}_i) - f_{B_i}(\mathcal{C}_i) \geq -\epsilon/7$), the second is due to the stopping condition of the algorithm ($f_{B_i}(\mathcal{C}_i) - f_{B_i}(\mathcal{C}_{i+1}) > \epsilon$), and the last is due to the remaining inequalities in Lemma B.12. The above holds w.h.p over all of the iterations of the algorithms. Using these guarantees for Algorithm 1 we can easily derive our main result for the truncated version.

Truncated termination Using Lemma 4.1 together with Lemma B.10 and the fact that $f_X(\mathcal{C}_i) - f_X(\mathcal{C}_{i+1}) \geq 3\epsilon/7$ we get that: $f_X(\hat{\mathcal{C}}_i) - f_X(\hat{\mathcal{C}}_{i+1}) \geq f_X(\mathcal{C}_i) - f_X(\mathcal{C}_{i+1}) - 2\epsilon/7 \geq \epsilon/7$. We conclude that when $b = \Omega(\max\{\gamma^4, \gamma^2\} k\epsilon^{-2} \log^2(\gamma n/\epsilon))$, w.h.p. Algorithm 2 terminates within $t = O(\gamma^2/\epsilon)$ iterations. This completes the second claim of Theorem 1.1. The final claim of Theorem 1.1 is due to the following lemma.

Lemma B.13. *The expected approximation ratio of the solution returned by Algorithm 2 is at least the approximation guarantee of the initial centers provided to the algorithm.*

Proof. Let $p = 1 - O(\epsilon/n\gamma^2) = 1 - O(1/n)$ be the success probability of a single iteration. By "success" we mean that all inequalities in Lemma B.12 hold. The value of p is due to the fact that we take $t = O(\gamma^2/\epsilon)$ and that $\gamma^2/\epsilon \geq 1/4$.

With probability at least p , it holds that $f_X(\mathcal{C}_{i+1}) \leq f_X(\mathcal{C}_i) - 2\epsilon/7$. On the other hand, f_X is upper bounded by $4\gamma^2$. Let us denote $Z = f_X(\mathcal{C}_i) - f_X(\mathcal{C}_{i+1})$ the change in the goal function after the i -th iteration. Consider the following:

$$E[Z] = E[Z \mid Z \geq \epsilon/7] \Pr[Z \geq \epsilon/7] + E[Z \mid Z < \epsilon/7] \Pr[Z < \epsilon/7]$$

⁷In (Schwartzman, 2023) this case is treated the same as the third inequality, which is incorrect. Using our approach the analysis can be fixed, with an additional multiplicative k factor in the batch size.

We show that $E[Z] = E[f_X(C_i) - f_X(C_{i+1})] \geq 0$ which implies that $E[f_X(C_{i+1})] \leq E[f_X(C_i)]$ and completes the proof. Note that if $E[Z | Z < \epsilon/7] > 0$ then we are done as we simply have a linear combination of two positive terms which is greater than 0. Let us focus on the case where $E[Z | Z < \epsilon/7] < 0$.

$$\begin{aligned} E[Z] &= E[Z | Z \geq \epsilon/7]Pr[Z \geq \epsilon/7] + E[Z | Z < \epsilon/7]Pr[Z < \epsilon/7] \\ &\geq p\epsilon/7 + E[Z | Z < \epsilon/7](1-p) \geq p\epsilon/7 - 4\gamma^2(1-p) \\ &= (1 - O(1/n))\epsilon/7 - 4\gamma^2O(\epsilon/\gamma^2n) = (1 - O(1/n))\epsilon/7 - O(\epsilon/n) > 0 \end{aligned}$$

Where the first inequality is due to the definition of p and the fact that $E[Z | Z < \epsilon/7] < 0$, the second is due to the upper bound on f_X , and the last inequality is by assuming n is sufficiently large.

□

C FULL EXPERIMENTAL RESULTS

We list our full experimental results in this section. We use the β prefix to denote that the algorithm uses the learning rate of (Schwartzman, 2023). τ denotes the maximum number of data points used to represent each truncated cluster center. We investigate 3 kernel functions: 1) The Gaussian kernel, as presented in Section 5, 2) The k-nearest-neighbor (k-nn) kernel, where the kernel matrix is $D^{-1}AD^{-1}$, A is a k-nn adjacency matrix of the data and D is the corresponding degree matrix, and 3) the heat kernel (Chung, 1997) where the kernel matrix is $\exp(-tD^{-1/2}AD^{-1/2})$ for some $0 < t < \infty$, A is a k-nn adjacency matrix and D is the corresponding degree matrix. All parameter settings can be found in the supplementary material.

Unlike for the Gaussian kernel where $\gamma = 1$; We observe empirically that for both the k-nn and heat kernels, $\gamma \ll 1$. In this case, the dependence on $\max\{\gamma^4, \gamma^2\}$ in the batch size required for Theorem 1.1 actually helps us. We found the parameters for these kernels to be easier to tune in practise than the Gaussian kernel parameter σ . For each kernel, we recorded the empirical value of gamma as follows:

Dataset	Kernel Type	γ
pendigits	knn	0.00100
pendigits	heat	0.0477
pendigits	gaussian	1
har	knn	0.000500
har	heat	0.0468
har	gaussian	1
mnist_784	knn	0.00220
mnist_784	heat	0.0612
mnist_784	gaussian	1
letter	knn	0.00100
letter	heat	0.0399
letter	gaussian	1

Table 1: γ values for various datasets and kernel types, rounded to 3 significant figures.

864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917

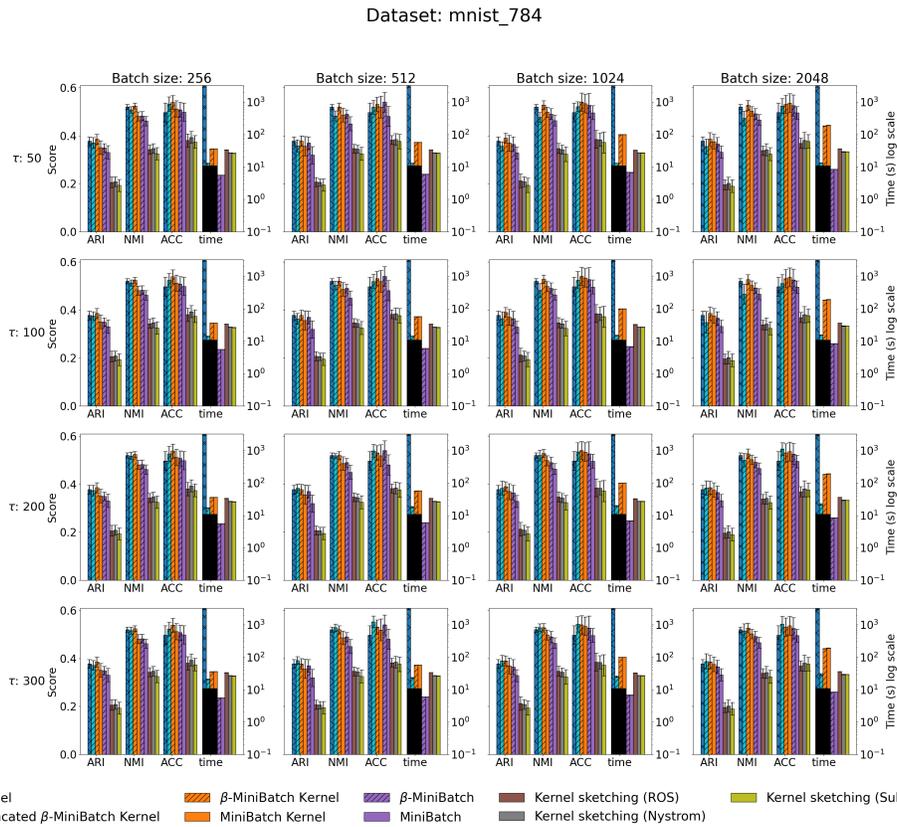


Figure 4: Experimental results on the MNIST dataset where the kernel algorithms use the Gaussian kernel.

918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971

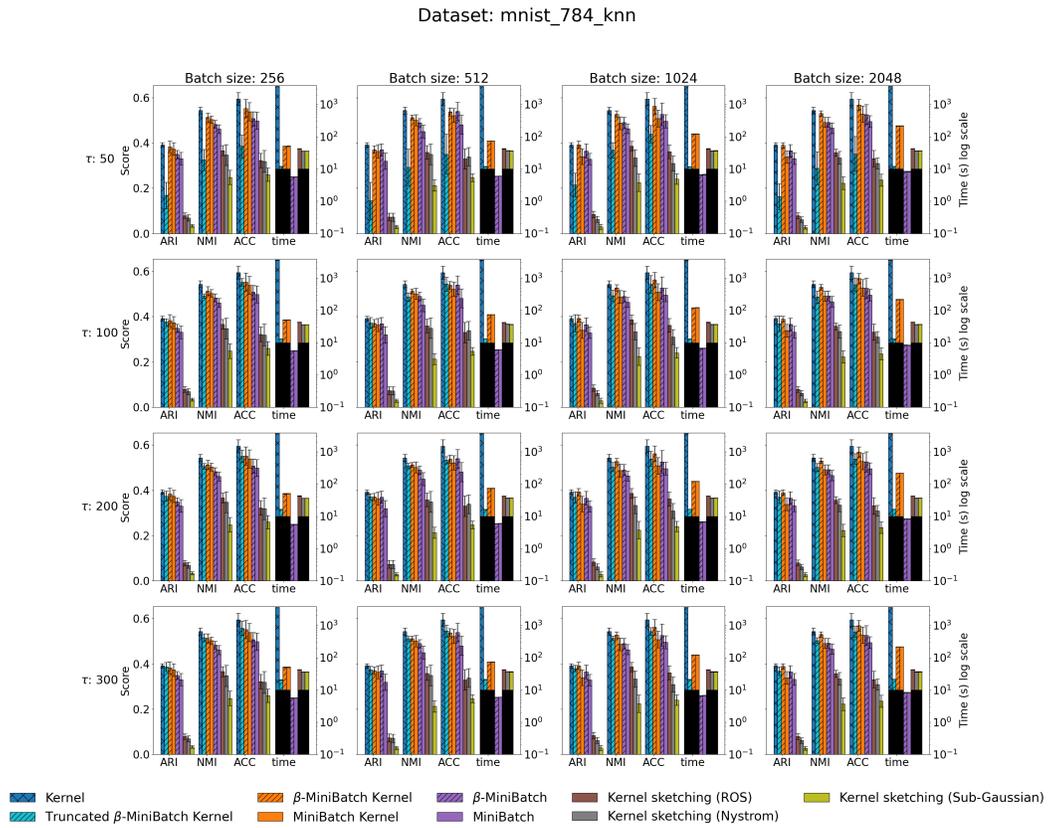


Figure 5: Experimental results on the MNIST dataset where the kernel algorithms use the k-nn kernel.

972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1020
1021
1022
1023
1024
1025

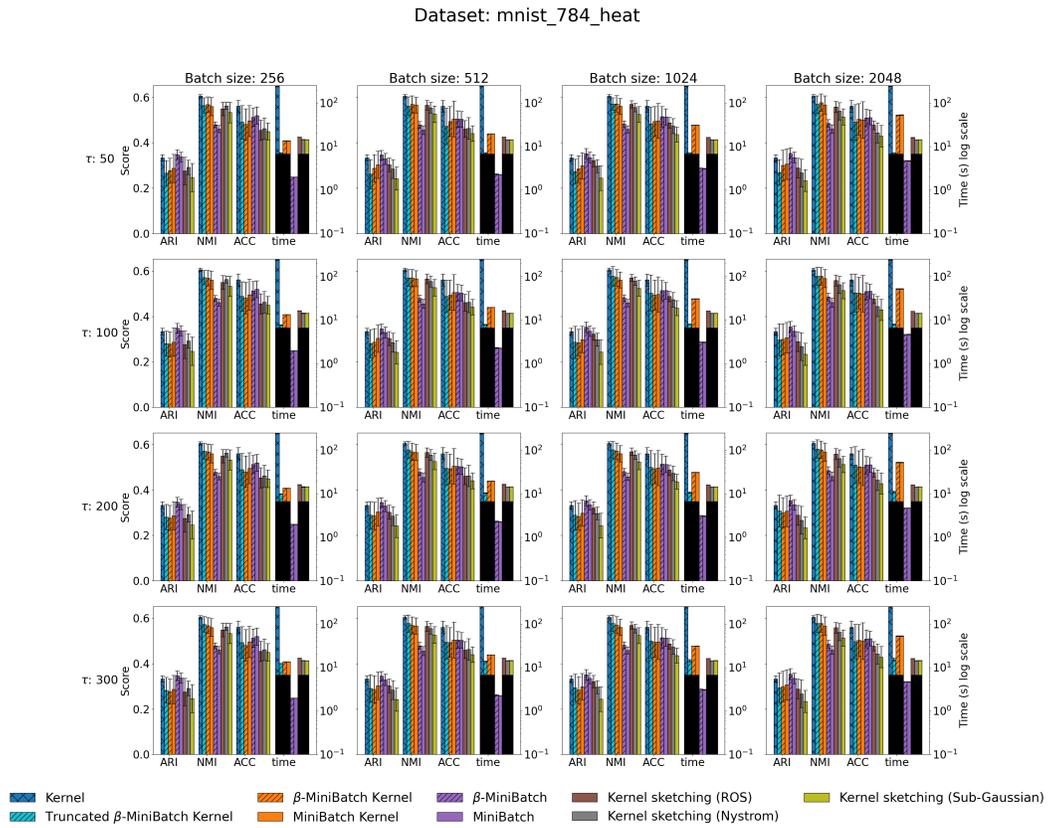


Figure 6: Experimental results on the MNIST dataset where the kernel algorithms use the Heat kernel.

1026
 1027
 1028
 1029
 1030
 1031
 1032
 1033
 1034
 1035
 1036
 1037
 1038
 1039
 1040
 1041
 1042
 1043
 1044
 1045
 1046
 1047
 1048
 1049
 1050
 1051
 1052
 1053
 1054
 1055
 1056
 1057
 1058
 1059
 1060
 1061
 1062
 1063
 1064
 1065
 1066
 1067
 1068
 1069
 1070
 1071
 1072
 1073
 1074
 1075
 1076
 1077
 1078
 1079

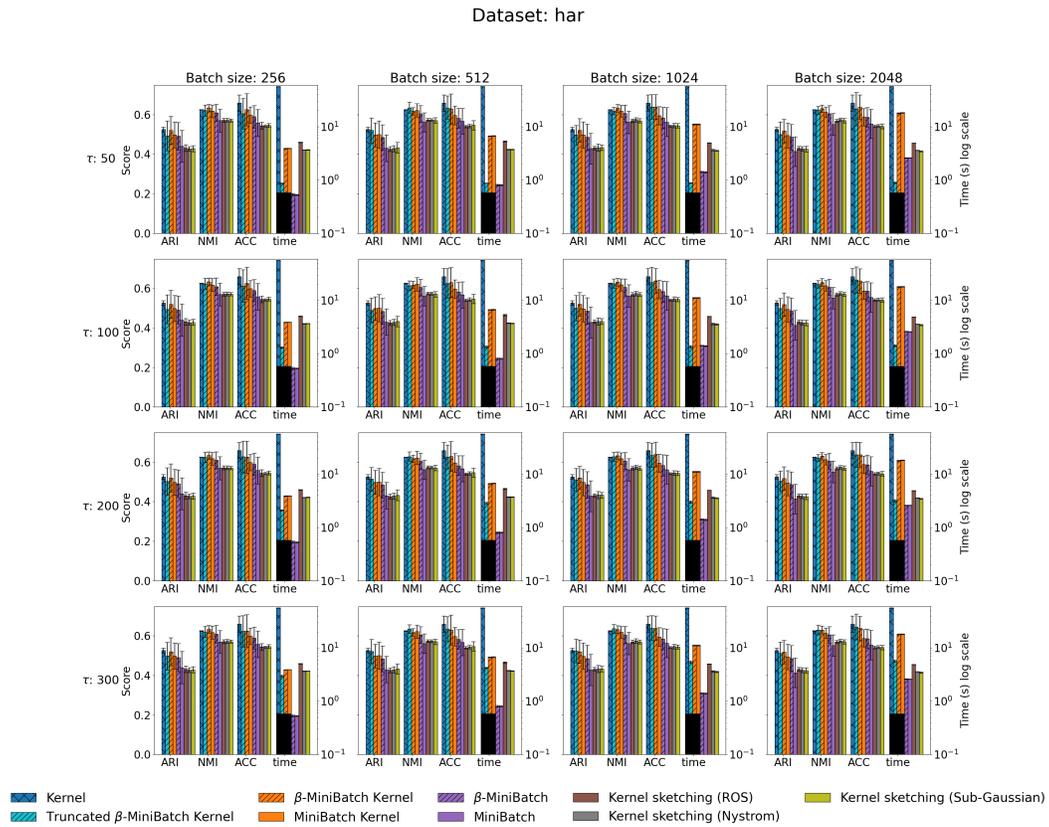


Figure 7: Experimental results on the Har dataset where the kernel algorithms use the Gaussian kernel.

1080
 1081
 1082
 1083
 1084
 1085
 1086
 1087
 1088
 1089
 1090
 1091
 1092
 1093
 1094
 1095
 1096
 1097
 1098
 1099
 1100
 1101
 1102
 1103
 1104
 1105
 1106
 1107
 1108
 1109
 1110
 1111
 1112
 1113
 1114
 1115
 1116
 1117
 1118
 1119
 1120
 1121
 1122
 1123
 1124
 1125
 1126
 1127
 1128
 1129
 1130
 1131
 1132
 1133

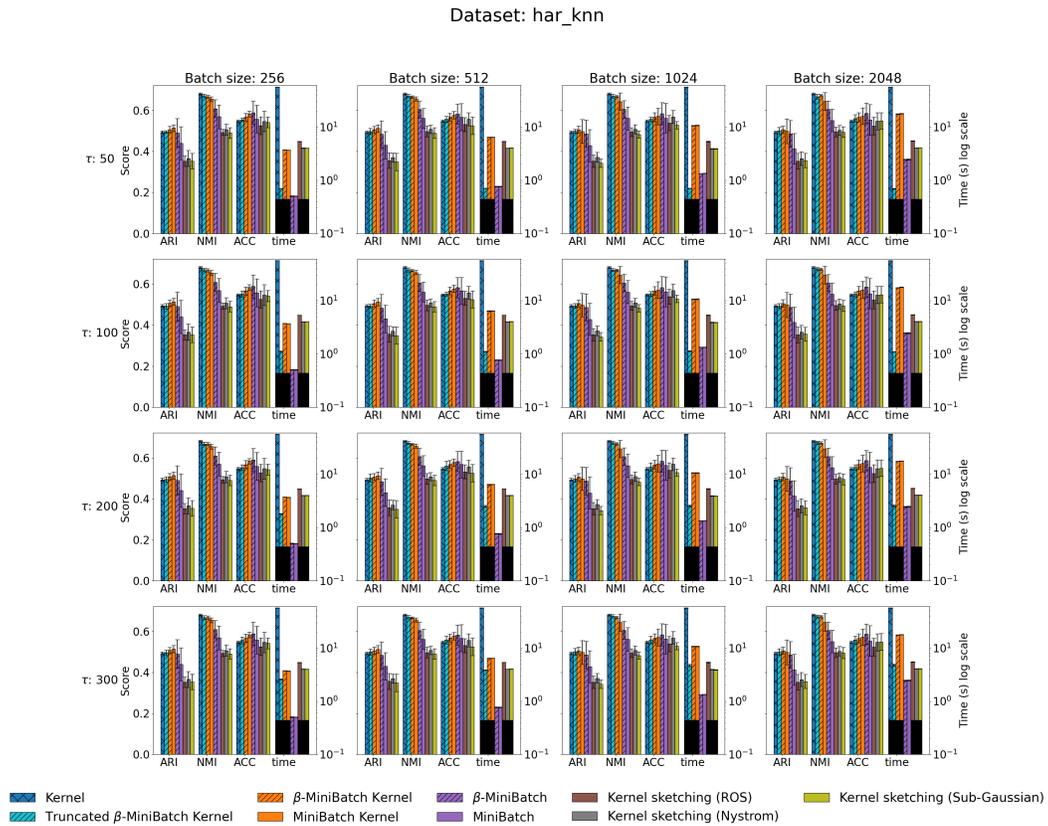


Figure 8: Experimental results on the Har dataset where the kernel algorithms use the k-nn kernel.

1134
1135
1136
1137
1138
1139
1140
1141
1142
1143
1144
1145
1146
1147
1148
1149
1150
1151
1152
1153
1154
1155
1156
1157
1158
1159
1160
1161
1162
1163
1164
1165
1166
1167
1168
1169
1170
1171
1172
1173
1174
1175
1176
1177
1178
1179
1180
1181
1182
1183
1184
1185
1186
1187

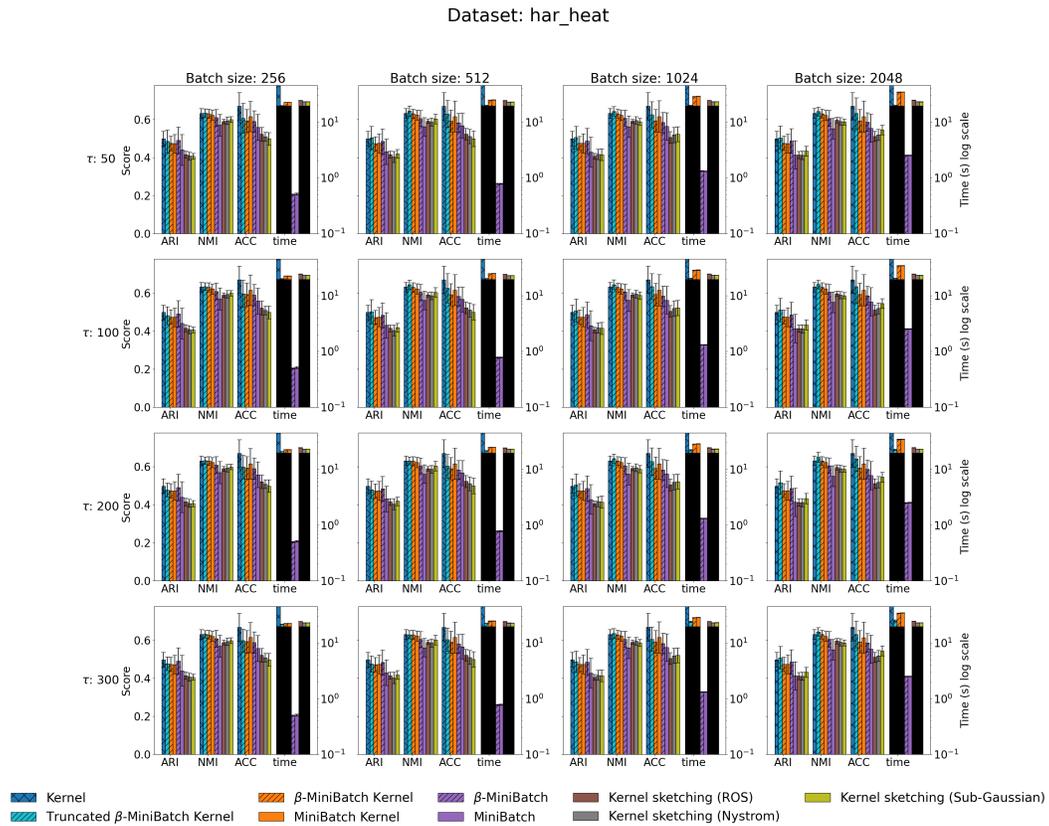


Figure 9: Experimental results on the Har dataset where the kernel algorithms use the Heat kernel.

1188
 1189
 1190
 1191
 1192
 1193
 1194
 1195
 1196
 1197
 1198
 1199
 1200
 1201
 1202
 1203
 1204
 1205
 1206
 1207
 1208
 1209
 1210
 1211
 1212
 1213
 1214
 1215
 1216
 1217
 1218
 1219
 1220
 1221
 1222
 1223
 1224
 1225
 1226
 1227
 1228
 1229
 1230
 1231
 1232
 1233
 1234
 1235
 1236
 1237
 1238
 1239
 1240
 1241

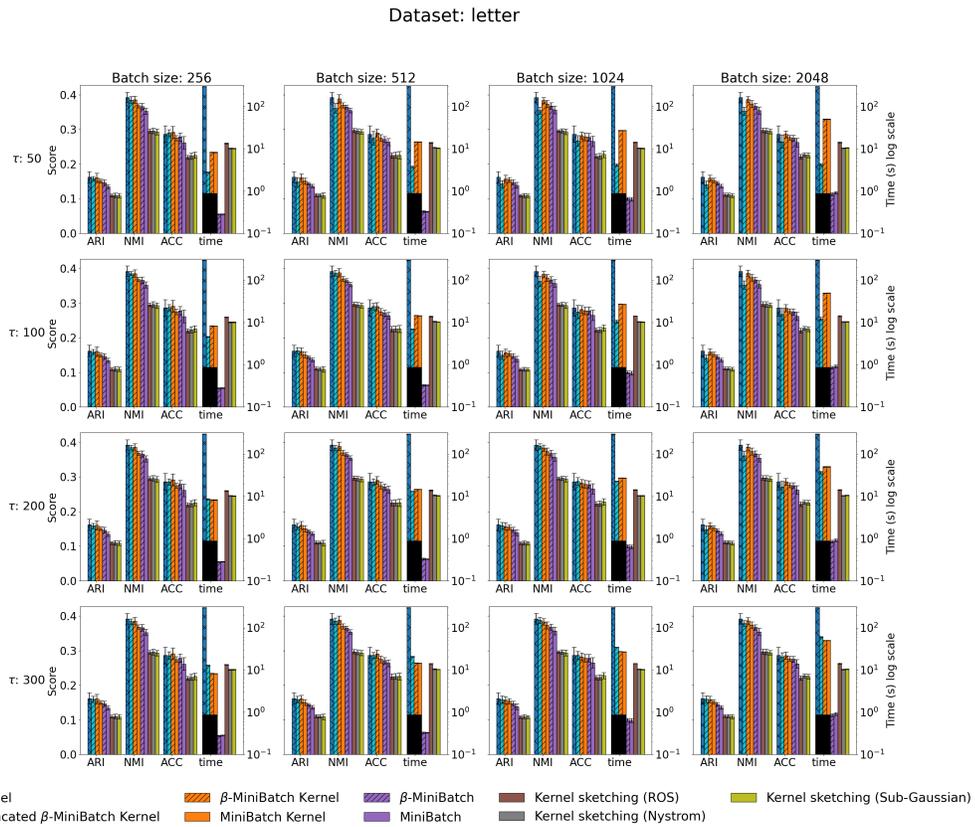


Figure 10: Experimental results on the Letter dataset where the kernel algorithms use the Gaussian kernel.

1242
 1243
 1244
 1245
 1246
 1247
 1248
 1249
 1250
 1251
 1252
 1253
 1254
 1255
 1256
 1257
 1258
 1259
 1260
 1261
 1262
 1263
 1264
 1265
 1266
 1267
 1268
 1269
 1270
 1271
 1272
 1273
 1274
 1275
 1276
 1277
 1278
 1279
 1280
 1281
 1282
 1283
 1284
 1285
 1286
 1287
 1288
 1289
 1290
 1291
 1292
 1293
 1294
 1295

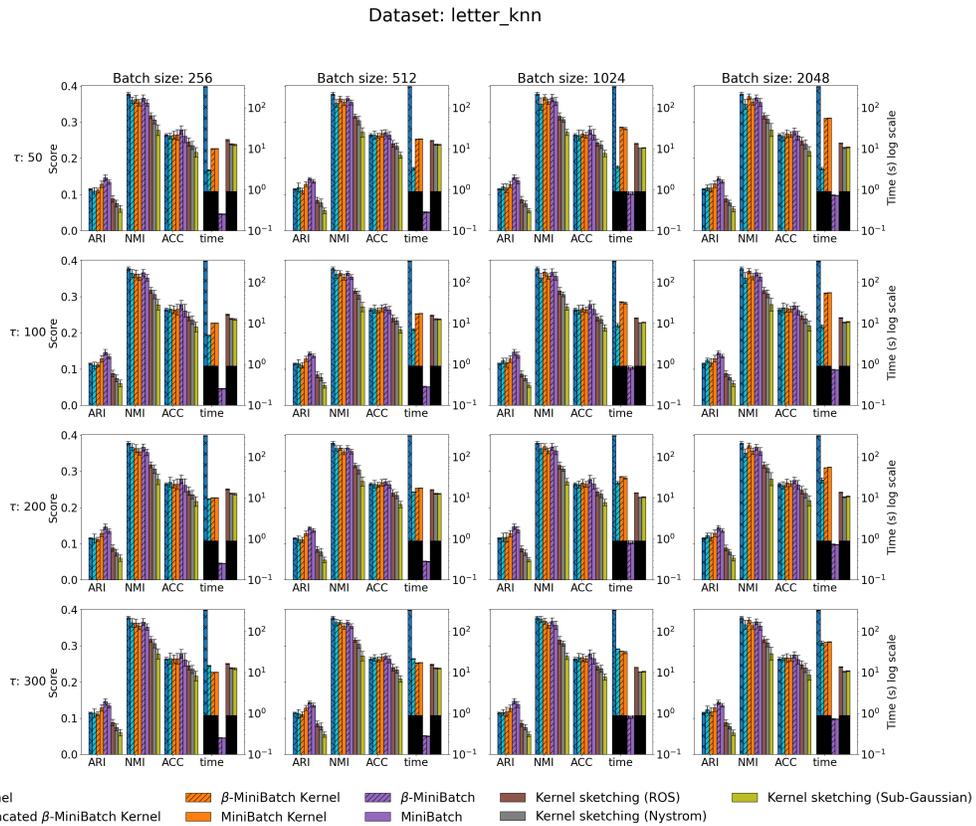


Figure 11: Experimental results on the Letter dataset where the kernel algorithms use the k-nn kernel.

1296
 1297
 1298
 1299
 1300
 1301
 1302
 1303
 1304
 1305
 1306
 1307
 1308
 1309
 1310
 1311
 1312
 1313
 1314
 1315
 1316
 1317
 1318
 1319
 1320
 1321
 1322
 1323
 1324
 1325
 1326
 1327
 1328
 1329
 1330
 1331
 1332
 1333
 1334
 1335
 1336
 1337
 1338
 1339
 1340
 1341
 1342
 1343
 1344
 1345
 1346
 1347
 1348
 1349

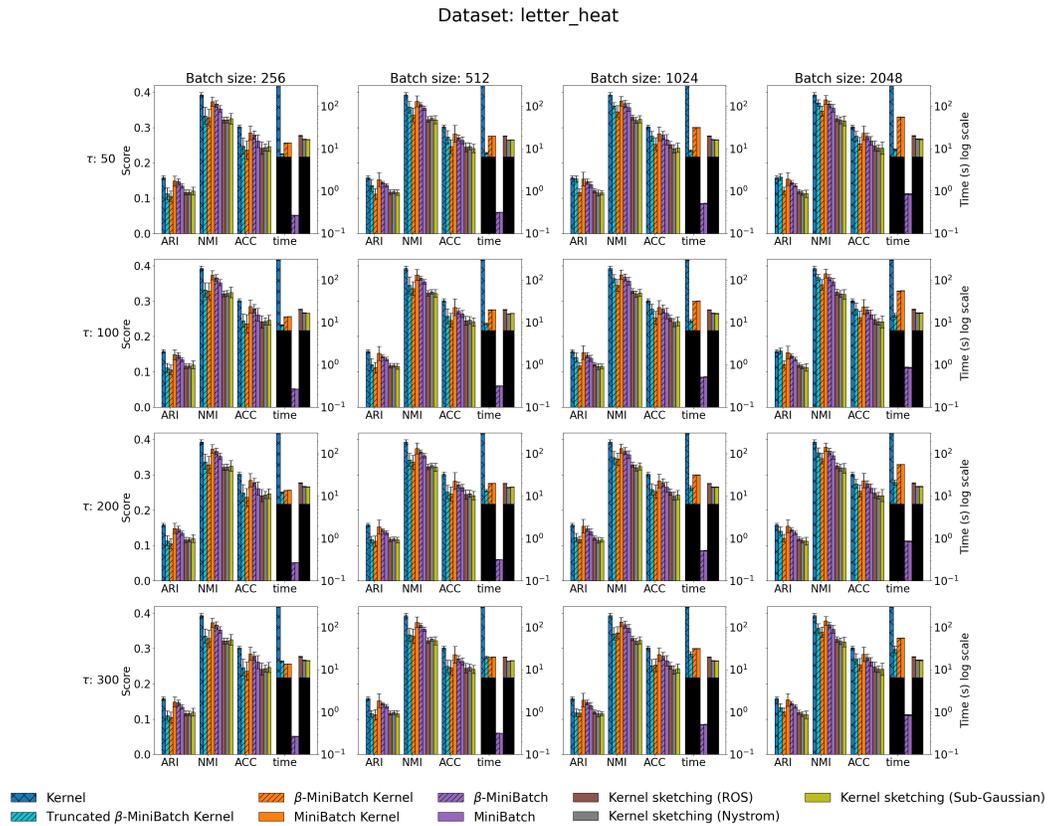


Figure 12: Experimental results on the Letter dataset where the kernel algorithms use the Heat kernel.

1350
 1351
 1352
 1353
 1354
 1355
 1356
 1357
 1358
 1359
 1360
 1361
 1362
 1363
 1364
 1365
 1366
 1367
 1368
 1369
 1370
 1371
 1372
 1373
 1374
 1375
 1376
 1377
 1378
 1379
 1380
 1381
 1382
 1383
 1384
 1385
 1386
 1387
 1388
 1389
 1390
 1391
 1392
 1393
 1394
 1395
 1396
 1397
 1398
 1399
 1400
 1401
 1402
 1403

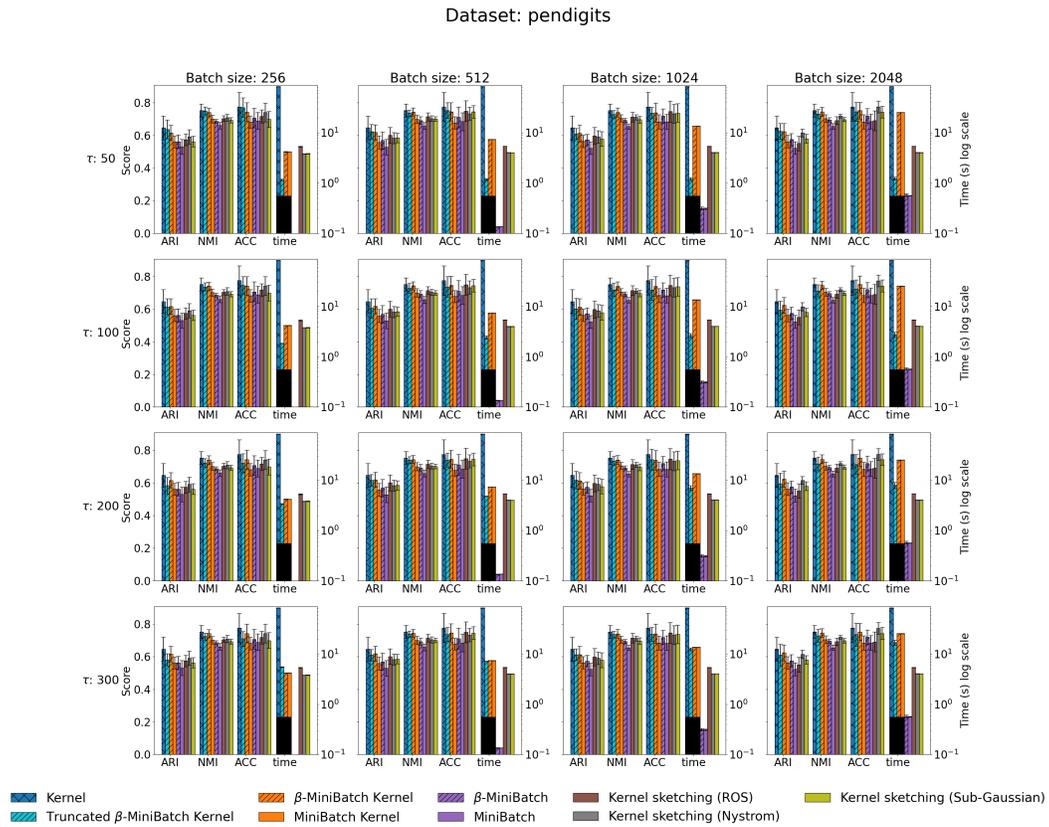


Figure 13: Experimental results on the Pendigits dataset where the kernel algorithms use the Gaussian kernel.

1404
 1405
 1406
 1407
 1408
 1409
 1410
 1411
 1412
 1413
 1414
 1415
 1416
 1417
 1418
 1419
 1420
 1421
 1422
 1423
 1424
 1425
 1426
 1427
 1428
 1429
 1430
 1431
 1432
 1433
 1434
 1435
 1436
 1437
 1438
 1439
 1440
 1441
 1442
 1443
 1444
 1445
 1446
 1447
 1448
 1449
 1450
 1451
 1452
 1453
 1454
 1455
 1456
 1457

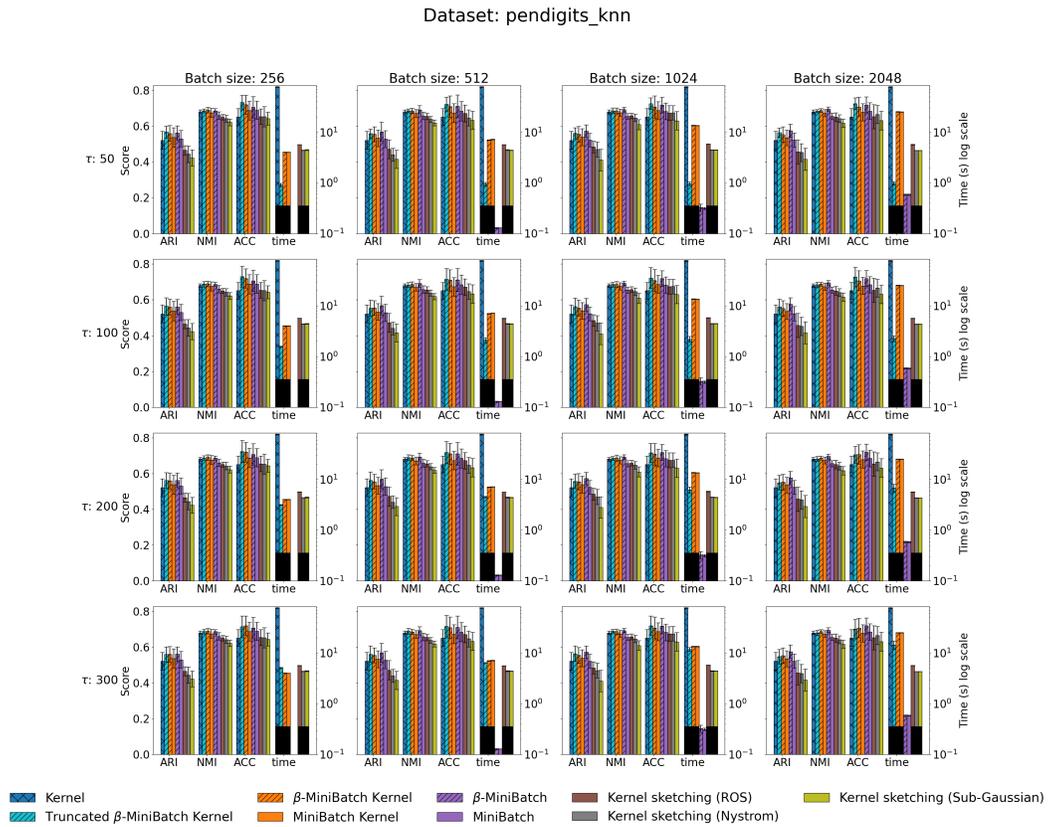


Figure 14: Experimental results on the Pendigits dataset where the kernel algorithms use the k-knn kernel.

1458
 1459
 1460
 1461
 1462
 1463
 1464
 1465
 1466
 1467
 1468
 1469
 1470
 1471
 1472
 1473
 1474
 1475
 1476
 1477
 1478
 1479
 1480
 1481
 1482
 1483
 1484
 1485
 1486
 1487
 1488
 1489
 1490
 1491
 1492
 1493
 1494
 1495
 1496
 1497
 1498
 1499
 1500
 1501
 1502
 1503
 1504
 1505
 1506
 1507
 1508
 1509
 1510
 1511

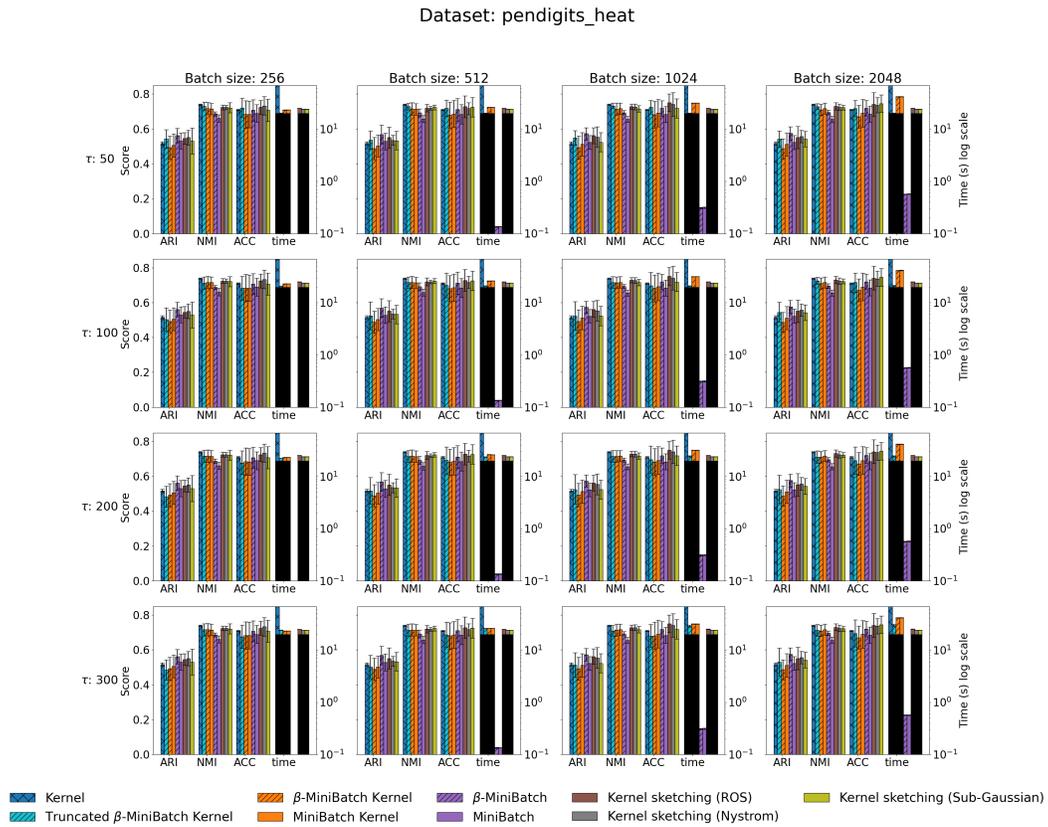


Figure 15: Experimental results on the Pendigits dataset where the kernel algorithms use the Heat kernel.

1512 D LLM STATEMENT

1513
1514
1515
1516
1517
1518
1519
1520
1521
1522
1523
1524
1525
1526
1527
1528
1529
1530
1531
1532
1533
1534
1535
1536
1537
1538
1539
1540
1541
1542
1543
1544
1545
1546
1547
1548
1549
1550
1551
1552
1553
1554
1555
1556
1557
1558
1559
1560
1561
1562
1563
1564
1565

LLMs were used to assist with typesetting, code autocomplete and to polish writing.