

k -variates++: more pluses in the k -means++

Richard Nock

Nicta & The Australian National University

`richard.nock@nicta.com.au`

Raphaël Canyasse

Ecole Polytechnique & The Technion

`raphael.canyasse@polytechnique.edu`

Roksana Boreli

Nicta & The University of New South Wales

`roksana.boreli@nicta.com.au`

Frank Nielsen

Ecole Polytechnique & Sony Computer Science Laboratories, Inc.

`Frank.Nielsen@acm.org`

January 22, 2018

Abstract

k -means++ seeding has become a de facto standard for hard clustering algorithms. In this paper, our first contribution is a two-way generalisation of this seeding, k -variates++, that includes the sampling of general densities rather than just a discrete set of Dirac densities anchored at the point locations, *and* a generalisation of the well known Arthur-Vassilvitskii (AV) approximation guarantee, in the form of a *bias+variance* approximation bound of the *global* optimum. This approximation exhibits a reduced dependency on the “noise” component with respect to the optimal potential — actually approaching the statistical lower bound. We show that k -variates++ *reduces* to efficient (biased seeding) clustering algorithms tailored to specific frameworks; these include distributed, streaming and on-line clustering, with *direct* approximation results for these algorithms. Finally, we present a novel application of k -variates++ to differential privacy. For either the specific frameworks considered here, or for the differential privacy setting, there is little to no prior results on the direct application of k -means++ and its approximation bounds — state of the art contenders appear to be significantly more complex and / or display less favorable (approximation) properties. We stress that our algorithms can still be run in cases where there is *no* closed form solution for the population minimizer. We demonstrate the applicability of our analysis via experimental evaluation on several domains and settings, displaying competitive performances vs state of the art.

1 Introduction

Arthur-Vassilvitskii’s (AV) k -means++ algorithm has been extensively used to address the hard membership clustering problem, due to its simplicity, experimental performance and guaranteed approximation of the *global* optimum; the goal being the k -partitioning of a dataset so as to minimize the sum of within-cluster squared distances to the cluster center (Arthur & Vassilvitskii, 2007), *i.e.*, a centroid or a *population minimizer* (Nock et al., 2016).

The k -means++ non-uniform seeding approach has also been utilized in more complex settings, including tensor clustering, distributed, data stream, on-line and parallel clustering, clustering with non-metric distortions and even clustering with distortions not allowing population minimizers in closed form (Ailon et al., 2009; Balcan et al., 2013; Jegelka et al., 2009; Liberty et al., 2014; Nock et al., 2008; Nielsen & Nock, 2015). However, apart from the non-uniform seeding, all these algorithms are distinct and (seemingly) do not share many common properties.

Finally, the application of k -means++ in some scenarios is still an open research topic, due to the related constraints – e.g., there is limited prior work in a differentially private setting (Nissim et al., 2007; Wang et al., 2015).

Our contribution — In a nutshell, we describe a generalisation of the k -means++ seeding process, k -variates++, which still delivers an efficient approximation of the global optimum, and can be used to obtain *and* analyze efficient algorithms for a wide range of settings, including: distributed, streamed, on-line clustering, (differentially) private clustering, etc. . We proceed in two steps.

First, we describe k -variates++ and analyze its approximation properties. We leverage two major components of k -means++: (i) data-dependent *probes* (specialized to observed data in the k -means++) are used to compute the weights for selecting centers, and (ii) selection of centers is based on an *arbitrary* family of densities (specialized to Diracs in the k -means++). Informally, the approximation properties (when only (ii) is considered), can be shown as:

$$\text{expected_cost}(k\text{-variates++}) \leq (2 + \log k) \cdot \Phi, \text{ with}$$

$\Phi \doteq 6 \cdot \text{optimal_noise-free_cost} + 2 \cdot \text{noise}(\text{bias} + \text{variance})$, where “noise” refers to the family of densities (note that constants are explicit in the bound). The dependence on these densities is arguably smaller than expectable (factor 2 for noise vs 6 for global optimum). There is also not much room for improvement: we show that the guarantee approaches the Fréchet-Cramér-Rao-Darmois lowerbound.

Second, we use this general algorithm in two ways. We use it directly in a differential privacy setting, addressing a conjecture of (Nissim et al., 2007) with weaker assumptions. We also demonstrate the use of this algorithm for a *reduction* to other biased seeding algorithms for distributed, streamed or on-line clustering, and obtain the approximation bounds for these algorithms. This simple reduction technique allows us to analyze lightweight algorithms that compare favorably to the state of the art in the related domains (Ailon et al., 2009; Balcan et al., 2013; Liberty et al., 2014), from the approximation, assumptions and / or complexity aspects. Experiments against state of the art for the distributed and differentially private settings display that solid performance improvement can be obtained.

The rest of this paper is organised as follows: Section 2 presents k -variates++. Section 3 presents approximation properties for distributed, streamed and on-line clustering that use a reduction from k -variates++. Section 4 presents direct applications of k -variates++ to differential

Algorithm 0 *k*-variates++

Input: data $\mathcal{A} \subset \mathbb{R}^d$ with $|\mathcal{A}| = m$, $k \in \mathbb{N}_*$, densities $\{p(\mu_a, \theta_a), a \in \mathcal{A}\}$, probe functions $\wp_t : \mathcal{A} \rightarrow \mathbb{R}^d$ ($t \geq 1$);

Step 1: Initialise centers $\mathcal{C} \leftarrow \emptyset$;

Step 2: **for** $t = 1, 2, \dots, k$

2.1: randomly sample $a \sim_{q_t} \mathcal{A}$, with $q_1 \doteq u_m$ and, for $t > 1$,

$$q_t(a) \doteq D_t(a) \left(\sum_{a' \in \mathcal{A}} D_t(a') \right)^{-1}, \text{ where } D_t(a) \doteq \min_{x \in \mathcal{C}} \|\wp_t(a) - x\|_2^2; \quad (1)$$

2.2: randomly sample $x \sim p(\mu_a, \theta_a)$;

2.3: $\mathcal{C} \leftarrow \mathcal{C} \cup \{x\}$;

Output: \mathcal{C} ;

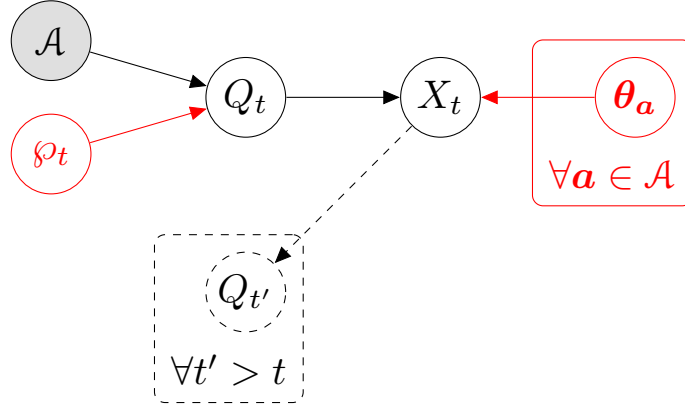


Figure 1: Graphical model for the *k*-means++ seeding process (black) and our generalisation (black + red, best viewed in color).

privacy. Section 5 presents experimental results. Last Section discusses extensions (to more distortion measures) and conclude. In order not to laden the paper’s body, an Appendix, starting page 18, provides all proofs, extensive experiments and additional remarks on the paper’s content.

2 *k*-variates++

We consider the hard clustering problem (Banerjee et al., 2005; Nock et al., 2016): given set $\mathcal{A} \subset \mathbb{R}^d$ and integer $k > 0$, find centers $\mathcal{C} \subset \mathbb{R}^d$ which minimizes the L_2^2 potential to the centers (here, $c(a) \doteq \arg \min_{c \in \mathcal{C}} \|a - c\|_2^2$):

$$\phi(\mathcal{A}; \mathcal{C}) \doteq \sum_{a \in \mathcal{A}} \|a - c(a)\|_2^2, \quad (2)$$

Algorithm 0 describes k -variates++. u_m denotes the uniform distribution over \mathcal{A} ($|\mathcal{A}| = m$). The parenthood with k -means++ seeding, which we name “ k -means++” for short¹ (Arthur & Vassilvitskii, 2007) can be best understood using Figure 1 (the red parts in Figure 1 are pinpointed in Algorithm 0). k -means++ is a random process that generates cluster centers from observed data \mathcal{A} . It can be modelled using a two-stage generative process for a mixture of Dirac distributions: the first stage involves random variable $Q_t \sim \text{Mult}(m, \boldsymbol{\pi}_t)$ whose parameters $\boldsymbol{\pi}_t \in \Delta_m$ (the m -dim probability simplex) are computed from the data and previous centers; sampling Q_t chooses the Dirac distribution, which is then “sampled” for one center (and the process iterates). All the crux of the technique is the design of $\boldsymbol{\pi}_t$, which, under *no* assumption of the data, yield in expectation a k -means potential for the centers chosen that is within $8(2 + \log k)$ of the global optimum (Arthur & Vassilvitskii, 2007).

k -variates++ generalize the process in two ways: first, the update of $\boldsymbol{\pi}_t$ depends on data and previous *probes*, using a sequence of *probe functions* $\wp_t : \mathcal{A} \rightarrow \mathbb{R}^d$ ($\wp = \text{Id}, \forall t$ in k -means++). Second, Diracs are replaced by arbitrary but fixed *local* (sometimes also called *noisy*) distributions with parameters² $(\boldsymbol{\mu}_a, \boldsymbol{\theta}_a)$ that depend on \mathcal{A} .

Let $\mathcal{C}_{\text{opt}} \subset \mathbb{R}^d$ denote the set of k centers minimizing (2) on \mathcal{A} . Let $\mathbf{c}_{\text{opt}}(\mathbf{a}) \doteq \arg \min_{\mathbf{c} \in \mathcal{C}_{\text{opt}}} \|\mathbf{a} - \mathbf{c}\|_2^2$ ($\mathbf{a} \in \mathcal{A}$), and

$$\phi_{\text{opt}} \doteq \sum_{\mathbf{a} \in \mathcal{A}} \|\mathbf{a} - \mathbf{c}_{\text{opt}}(\mathbf{a})\|_2^2, \quad (3)$$

$$\phi_{\text{bias}} \doteq \sum_{\mathbf{a} \in \mathcal{A}} \|\boldsymbol{\mu}_a - \mathbf{c}_{\text{opt}}(\mathbf{a})\|_2^2, \quad (4)$$

$$\phi_{\text{var}} \doteq \sum_{\mathbf{a} \in \mathcal{A}} \text{tr}(\Sigma_a). \quad (5)$$

ϕ_{opt} is the optimal **noise-free** potential, ϕ_{bias} is the bias of the noise³, and ϕ_{var} its variance, with $\Sigma_a \doteq \mathbb{E}_{\mathbf{x} \sim p_a}[(\mathbf{x} - \boldsymbol{\mu}_a)(\mathbf{x} - \boldsymbol{\mu}_a)^\top]$ the covariance matrix of p_a . Notice that when $\boldsymbol{\mu}_a = \mathbf{a}$, $\phi_{\text{bias}} = \phi_{\text{opt}}$. Otherwise, it *may* hold that $\phi_{\text{bias}} < \phi_{\text{opt}}$, and even $\phi_{\text{bias}} = 0$ if expectations coincide with \mathcal{C}_{opt} . Let \mathcal{C}_{opt} denote the partition of \mathcal{A} according to the centers in \mathcal{C}_{opt} . We say that probe function \wp_t is η -stretching if, informally, replacing points by their probes does not distort significantly the observed potential of an optimal cluster, with respect to its actual optimal potential. The formal definition follows.

Definition 1 *Probe functions \wp_t are said η -stretching on \mathcal{A} , for some $\eta \geq 0$, iff the following holds: for any cluster $A \in \mathcal{C}_{\text{opt}}$ and any $\mathbf{a}_0 \in A$ such that $\phi(\wp_t(A); \{\wp_t(\mathbf{a}_0)\}) \neq 0$, for any set of at most k centers $\mathcal{C} \subset \mathbb{R}^d$,*

$$\frac{\phi(A; \mathcal{C})}{\phi(A; \{\mathbf{a}_0\})} \leq (1 + \eta) \cdot \frac{\phi(\wp_t(A); \mathcal{C})}{\phi(\wp_t(A); \{\wp_t(\mathbf{a}_0)\})}, \forall t. \quad (6)$$

¹Both approaches can be completed with the same further local monotonous optimization steps like Lloyd or Hartigan iterations; furthermore, it is the biased seeding which holds the approximation properties of k -means++.

²Because expectations are the major parameter for clustering, we split the parameters in the form of $\boldsymbol{\mu}_a$ (expectation) and $\boldsymbol{\theta}_a$ (other parameters, e.g. covariance matrix).

³We term it *bias* by analogy with supervised classification, considering that the expectations of the densities could be used as models for the cluster centers (Kohavi & Wolpert, 1996).

Since $\phi(A; \mathcal{C}_{\text{opt}}) = \sum_{\mathbf{a}_0 \in A} \phi(A; \{\mathbf{a}_0\})$ (Arthur & Vassilvitskii, 2007) (Lemma 3.2), Definition 1 roughly states that the potential of an optimal cluster with respect to a set of cluster centers, relatively to its potential with respect to the optimal set of centers, does not blow up through probe function \wp_t . The identity function is trivially 0-stretching, for any \mathcal{A} . Many local transformations would be eligible for η -stretching probe functions with η small, including local translations, mappings to core-sets (Har-Peled & Mazumdar, 2004), mappings to Voronoi diagram cell centers (Boissonnat et al., 2010), etc. Notice that ineq. (6) has to hold only for optimal clusters and *not* any clustering of \mathcal{A} . Let $\mathbb{E}[\phi(\mathcal{A}; \mathcal{C})] \doteq \int \phi(\mathcal{A}|\mathcal{C}) dp(\mathcal{C})$ denote the expected potential over the random sampling of \mathcal{C} in k -variates++.

Theorem 2 *For any dataset \mathcal{A} , any sequence of η -stretching probe functions \wp_t and any density $\{p_{\mathbf{a}}, \mathbf{a} \in \mathcal{A}\}$, the expected potential of k -variates++ satisfies:*

$$\mathbb{E}[\phi(\mathcal{A}; \mathcal{C})] \leq (2 + \log k) \cdot \Phi, \quad (7)$$

with $\Phi \doteq (6 + 4\eta)\phi_{\text{opt}} + 2\phi_{\text{bias}} + 2\phi_{\text{var}}$.

(Proof in page 19) Five remarks are in order. First, we retrieve the result of (Arthur & Vassilvitskii, 2007) in their setting ($\eta = \phi_{\text{var}} = 0$, $\phi_{\text{bias}} = \phi_{\text{opt}}$). Second, in the case where $\phi_{\text{bias}} < \phi_{\text{opt}}$, we *may* beat AV’s bound. This is not due to an improvement of the algorithm, but to a finer analysis which shows that special settings may “naturally” favor the improvement. We shall see one example in the distributed clustering case. Third, apart from being η -stretching, there is no constraint on the choice of probe functions \wp_t : it can be randomized, iteration dependent, etc. Fourth, the algorithm can easily be generalized to the case where points are weighted. Last, as we show in the following Lemma, the dependence in noise in ineq. (7) can hardly be improved in our framework.

Lemma 3 *Suppose each point in \mathcal{A} is replaced (i.i.d.) by a point sampled in $p_{\mathbf{a}}$ with $\Sigma_{\mathbf{a}} = \Sigma$. Then any clustering algorithm suffers: $\mathbb{E}[\phi(\mathcal{A}; \mathcal{C})] = \Omega(|\mathcal{A}| \text{tr}(\Sigma))$.*

(Proof in page 22) We make use of k -variates++ in two different ways. First, we show that it can be used to prove approximation properties for algorithms operating in different clustering settings: distributed clustering, streamed clustering and on-line clustering. The proof involves a *reduction* (see page 23) from k -variates++ to each of these algorithms. By reduction, we mean there exists distributions and probe functions (even non poly-time computable) for which k -variates++ yields the same result in expectation as the other algorithm, thus directly yielding an approximability ratio of the global optimum for this latter algorithm via Theorem 2. Second, we show how k -variates++ can *directly* be specialized to address settings for which no efficient application of k -means++ was known.

3 Reductions from k -variates++

Despite tremendous advantages, k -means++ has a serious downside: it is difficult to parallelize, distribute or stream it under relevant communication, space, privacy and/or time resource constraints (Bahmani et al., 2012). Although extending k -means clustering to these settings has been a major research area in recent years, there has been no obvious solution to tailoring k -means++ (Ackermann et al., 2010; Ailon et al., 2009; Bahmani et al., 2012; Balcan et al., 2013; Liberty et al., 2014; Shindler et al., 2011) (and others).

	Ref.	Property	Them	Us
(1)	(Bahmani et al., 2012)	Communication complexity	$O(n^2 \ell \cdot \log \phi_1)$ (expected)	$O(n^2 k)$
(2)	(Bahmani et al., 2012)	# data to compute one center	m	$\leq \max_{i \in [n]} (m/m_i)$
(3)	(Bahmani et al., 2012)	Data points shared	$O(\ell \cdot \log \phi_1)$ (expected)	k
(4)	(Bahmani et al., 2012)	Approximation bound	$O((\log k) \cdot \phi_{\text{opt}})$	$(2 + \log k) \cdot (10\phi_{\text{opt}} + 6\phi_s^F)$
(I)	(Balcan et al., 2013)	Communication complexity	$\Omega((nkd/\varepsilon^4) + n^2 k \ln(nk))$	$O(n^2 k)$
(II)	(Balcan et al., 2013)	Data points shared	$\Omega((kd/\varepsilon^4) + nk \ln(nk))$	k
(III)	(Balcan et al., 2013)	Approximation bound	$(2 + \log k)(1 + \varepsilon) \cdot 8\phi_{\text{opt}}$	$(2 + \log k) \cdot (10\phi_{\text{opt}} + 6\phi_s^F)$
(i)	(Ailon et al., 2009)	Time complexity (outer loop)	— identical —	— identical —
(ii)	(Ailon et al., 2009)	Approximation bound	$(2 + \log k)(1 + \eta) \cdot 32\phi_{\text{opt}}$	$(2 + \log k) \cdot ((8 + 4\eta)\phi_{\text{opt}} + 2\phi_s^F)$
(a)	(Liberty et al., 2014)	Knowledge required	Lowerbound $\phi^* \leq \phi_{\text{opt}}$	None
(b)	(Liberty et al., 2014)	Approximation bound	$O(\log m \cdot \phi_{\text{opt}})$	$(2 + \log k) \cdot (4 + (32/\zeta^2)) \phi_{\text{opt}}$
(A)	(Nissim et al., 2007)	Knowledge required	$\lambda(\phi_{\text{opt}})$	None
(B)	(Nissim et al., 2007)	Noise variance (σ)	$O(\lambda k R/\epsilon)$	$O(R/(\epsilon + \log m))$
(C)	(Nissim et al., 2007)	Approximation bound	$O^*(\phi_{\text{opt}} + m\lambda^2 k R^2/\epsilon^2)$	$O(\log k(\phi_{\text{opt}} + mR^2/(\epsilon + \log m)^2))$
(α)	(Wang et al., 2015)	Assumptions on ϕ_{opt}	Several (separability, size of clusters, etc.)	None
(β)	(Wang et al., 2015)	Approximation bound	$O^*(\phi_{\text{opt}} + km \log(m) R^2/\epsilon^2)$	$O(\log k(\phi_{\text{opt}} + mR^2/(\epsilon + \log m)^2))$

Table 1: Comparison with state of the art approaches for distributed clustering (1-4, I-III), streamed clustering (i, ii), on-line clustering (a, b) and differential privacy (A-C, α , β). Notations used for the "Them" column are as follows. ϕ_1 is the expected potential of a clustering with a single cluster over the *whole* data and ℓ is in general $\Omega(k)$ (Bahmani et al., 2012). ε is the coreset approximation factor in (Balcan et al., 2013). η is the approximation factor of the optimum in (Ailon et al., 2009). λ is the separability factor in Definition 5.1 in (Nissim et al., 2007).

Algorithm 1 Dk -means++ (// PDk -means++)

Input: Forgy nodes $(F_i, \mathcal{A}_i), i \in [n]$,
for $t = 1, 2, \dots, k$
 Round 1 : N^* picks $i^* \sim_{q_t^D} [n]$ and asks F_{i^*} for a center;
 Round 2 : F_{i^*} picks $\mathbf{a} \sim_{u_{i^*}} \mathcal{A}_{i^*}$ and sends \mathbf{a} to $F_i, \forall i$;
 // PDk -means++: F_{i^*} sends $\mathbf{x} \sim p_{(\mu_{\mathbf{a}}, \theta_{\mathbf{a}})}$ to $F_i, \forall i$;
 Round 3 : $\forall i, F_i$ updates $D_t(\mathcal{A}_i)$ and sends it to N^* ;
Output: \mathcal{C} = set of broadcasted \mathbf{a} s (or \mathbf{x} s);

Distributed clustering We consider horizontally partitioned data among *peers*, in line with (Bahmani et al., 2012), and a setting significantly more restrictive than theirs: each peer can only locally run the standard operations of Forgy initialisation (that is, uniform random seeding) on its own data, unlike for example the biased distributions of (Bahmani et al., 2012). This is consistent with the notion that data handling peers are not necessarily computationally intensive resources. Additionally, due to privacy constraints, we limit the data sharing between nodes. We denote the nodes handling the data *Forgy nodes*. We have n such nodes, $(F_i, \mathcal{A}_i), i \in [n]$, where \mathcal{A}_i is the dataset held by F_i . To enable more complex operations necessary to implement k -variates++, we introduce a special node, N^* , that has high computation power, *but* is not allowed to handle *any* data (points) from the Forgy nodes. We therefore split the location of the computational power from the location of the data. We also prevent the Forgy nodes from exchanging *any* data between themselves, with the sole exception of cluster centers. We note that none of the algorithms of (Ailon et al., 2009; Balcan et al., 2013; Bahmani et al., 2012) would be applicable to this setting without non-trivial modifications affecting their properties.

Algorithm 1 defines the mechanism that is consistent with our setting. It includes two variants: a protected version Dk -means++ where Forgy nodes directly share local centers and a private

version PD k -means++ where the nodes share noisy centers, such as to ensure a differentially private release of centers (with relevant noise calibration). Notations used in Algorithm 1 are as follows. Let $D_t(\mathcal{A}_i) \doteq \sum_{\mathbf{a} \in \mathcal{A}_i} D_t(\mathbf{a})$ and $q_{ti}^D \doteq D_t(\mathcal{A}_i) \cdot (\sum_j D_t(\mathcal{A}_j))^{-1}$ if $t > 1$ and $q_{ti}^D \doteq 1/n$ otherwise. Also, u_i is uniform distribution on $[m_i]$, with $m_i \doteq |\mathcal{A}_i|$.

Theorem 4 *Let $\phi_s^F \doteq \sum_{i \in [n]} \sum_{\mathbf{a} \in \mathcal{A}_i} \|c(\mathcal{A}_i) - \mathbf{a}\|_2^2$ be the total spread of the Forgy nodes ($c(\mathcal{A}_i) \doteq (1/m_i) \cdot \sum_{\mathbf{a} \in \mathcal{A}_i} \mathbf{a}$). At iteration k , the expected potential on the **total** data $\mathcal{A} \doteq \cup_i \mathcal{A}_i$ satisfies ineq. (7) with*

$$\Phi \doteq \begin{cases} 10\phi_{\text{opt}} + 6\phi_s^F & (\text{Dk-means++}) \\ 10\phi_{\text{opt}} + 4\phi_s^F + 2\phi_{\text{var}} & (\text{PDk-means++}) \end{cases} \quad (8)$$

Here, ϕ_{opt} is the optimal potential on **total** data \mathcal{A} .

(Proof in page 23) We note that the optimal potential is defined on the total data. The dependence on ϕ_s^F , which is just the peer-wise variance of data, is thus rather intuitive. A positive point is that ϕ_s^F is weighted by a factor smaller than the factor that weights the optimal potential. Another positive point is that this parameter *can* be computed from data, and among peers, without disclosing more data. Hence, it may be possible to estimate the loss against the centralized, k -means++ setting, taking as reference eq. (8). To gain insight in the leverage that Theorem 4 provides, Table 1 compares D k -means++ to (Balcan et al., 2013)’s (ε is the coreset approximation parameter), even though the latter approach would not be applicable to our restricted framework. To be fair, we assume that the algorithm used to cluster the coreset in (Balcan et al., 2013) is k -means++. We note that, considering the communication complexity and the number of data points shared, Algorithm 1 is a clear winner. In fact, Algorithm 1 can also win from the approximability standpoint. The dependence in ε prevents to fix it too small in (Balcan et al., 2013). Comparing the bounds in row (III) shows that if $\varepsilon > 1/4$, then we can also be better from the approximability standpoint if the spread satisfies $\phi_s^F = O(\phi_{\text{opt}})$. While this may not be feasible over arbitrary data, it becomes more realistic on several real-world scenarii, when Forgy nodes aggregate “local” data with respect to features, *e.g.*, state-wise insurance data, city-wise financial data, etc. When n increases, this also becomes more realistic.

Streaming clustering We have access to a stream S , with an assumed finite size: S is a sequence of points $\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_m$. We authorise the computation / output of the clustering at the end of the stream, *but* the memory n allowed for all operations satisfies $n < m$, such as $n = m^\alpha$ with $\alpha < 1$ in (Ailon et al., 2009). We assume for simplicity that each point can be stored in one storage memory unit. Algorithm 2 (Sk -means++) presents our approach. It relies on the standard “trick” of summarizing massive datasets via compact representations (synopses) before processing them (Indyk et al., 2014). The approximation properties of Sk -means++, proven using a reduction from k -variates++, hold regardless of the way synopses are built. They show that two key parameters may guide its choice: the spread of the synopses, analogous to the spread of Forgy nodes for distributed clustering, and the stretching properties of the synopses used as centers.

Theorem 5 *Let $\wp(\mathbf{a}) \doteq \arg \min_{\mathbf{s}' \in S} \|\mathbf{a} - \mathbf{s}'\|_2^2, \forall \mathbf{a} \in S$. Let $\phi_s^\wp \doteq \sum_{\mathbf{a} \in S} \|\wp(\mathbf{a}) - \mathbf{a}\|_2^2$ be the spread of \wp on synopsis set S . Let $\eta > 0$ such that \wp is η -stretching on S . Then the expected*

Algorithm 2 Sk -means++

Input: Stream S Step 1: $\mathcal{S} \doteq \{(s_j, m_j), i \in [n]\} \leftarrow \text{SYNOPSIS}(S, n)$;Step 2: **for** $t = 1, 2, \dots, k$ 2.1: **if** $t = 1$ then let $s_j \sim_{u_n} S$ **else** $s_j \sim_{q_t^S} S$ s.t.

$$q_t^S(s_j) \doteq m_j D_t(s_j) \left(\sum_{j' \in [n]} m_{j'} D_t(s_{j'}) \right)^{-1}; \quad (9)$$

 // $D_t(s_j) \doteq \min_{c \in \mathcal{C}} \|s_j - c\|_2^2$; 2.2: $\mathcal{C} \leftarrow \mathcal{C} \cup \{s_j\}$;**Output:** Cluster centers \mathcal{C} ;

Algorithm 3 OLk -means++

Input: Minibatch S_j , current weighted centers \mathcal{C} ;Step 1: **if** $j = 1$ then let $s \sim_{u_1} S_1$ **else** $s \sim_{q_j^O} S_j$ s.t.

$$q_j^O(s) \doteq D_t(s) \left(\sum_{s' \in S_j} D_t(s') \right)^{-1}; \quad (10)$$

 // $D_t(s) \doteq \min_{c \in \mathcal{C}} \|s - c\|_2^2$;Step 2: $\mathcal{C} \leftarrow \mathcal{C} \cup \{s\}$;

potential of Sk -means++ on stream S satisfies ineq. (7) with

$$\Phi \doteq (8 + 4\eta)\phi_{\text{opt}} + 2\phi_s^{\mathcal{P}},$$

*Here, ϕ_{opt} is the optimal potential on **stream** S .*

(Proof in page 25) It is not surprising to see that Sk -means++ looks like a generalization of (Ailon et al., 2009) and almost matches it (up to the number of centers delivered) when $k' \gg k$ synopses are learned from k' -means#. Yet, we rely on a different — and more general — analysis of its approximation properties. Table 1 compares properties of Sk -means++ to (Ailon et al., 2009) (η relates to approximation of the k -means objective in inner loop).

On-line clustering This setting is probably the farthest from the original setting of the k -means++ algorithm. Here, points arrive in a sequence, finite, but of unknown size and too large to fit in memory (Liberty et al., 2014). We make no other assumptions — the sequence can be random, or chosen by an adversary. Therefore, the expected analysis we make is only with respect to the internal randomisation of the algorithm, *i.e.*, for the fixed stream sequence as it is observed. We do not assume a feedback for learning (common for supervised learning); so, we do not assume that the algorithm has to predict a cluster for each point that arrives, yet it has to be easily modifiable to do so.

Our approach is summarized in Algorithm 3 (OL k -means++), a variation of k -means++ which consists of splitting the stream S into minibatches S_j for $j = 1, 2, \dots$, each of which is used to sample one center. u_1 denotes the uniform distribution with support S_1 . Let $R \doteq \max_{\mathbf{a}, \mathbf{a}' \in S} \|\mathbf{a} - \mathbf{a}'\|_2 (\ll \infty)$ be the diameter of S .

Theorem 6 *Let $\varsigma > 0$ be the largest real such that the following conditions are met (for any $A \in C_{\text{opt}}, j \geq 1$): for any set of at most k centers \mathcal{C} , $\sum_{\mathbf{a}, \mathbf{a}' \in A} \|\mathbf{a} - \mathbf{a}'\|_2^2 \geq \varsigma \cdot \binom{|A|}{2} R^2$ and $\sum_{\mathbf{a} \in A \cap S_j} \|\mathbf{a} - \mathbf{c}(\mathbf{a})\|_2^2 \geq \varsigma \cdot \sum_{\mathbf{a} \in A} \|\mathbf{a} - \mathbf{c}(\mathbf{a})\|_2^2$ (with $\mathbf{c}(\mathbf{a})$ defined in eq. (2)). Then the expected potential of OL k -means++ on stream S satisfies ineq. (7) with*

$$\Phi \doteq \left(4 + \frac{32}{\varsigma^2}\right) \cdot \phi_{\text{opt}} ,$$

where ϕ_{opt} is the optimal potential on **stream** S .

(Proof in page 26) Notice that loss function $\phi(S, \mathcal{C})$ in eq. (2) implies the finiteness of S , and the existence of $\varsigma > 0$; also, the second condition implies $\varsigma \leq 1$. In (Liberty et al., 2014), the clustering algorithm is required to have space and time at most polylog in the length of the stream. Hence, each minibatch can be reasonably large with respect to the stream — the larger they are, the larger ς . The knowledge of ς is not necessary to run OL k -means++; it is just a part of the approximation bound which quantifies the loss in approximation due to the fact that centers are computed from the *partial* knowledge of the stream. Table 1 compares properties of OL k -means++ to (Liberty et al., 2014) (we picked the fully on-line, non-heuristic algorithm). To compare the bounds, suppose that batches have the same size, b , so that $\log k = \log(m/b)$. If batches are at least polylog size, up to what is hidden in the big-Oh notation, our approximation can be quite competitive when ς is large, *e.g.*, if d is large and optimal clusters are not too small.

4 Direct use of k -variates++

The most direct application domain of k -variates++ is differential privacy. Several algorithms have independently emphasised the idea that powerful mechanisms may be amended via a carefully designed noise mechanism to broaden their scope with new capabilities, without overly challenging their original properties. Examples abound (Hardt & Price, 2014; Kalai & Vempala, 2005; Chaudhuri et al., 2011; Chichignoud & Lousteau, 2014), etc. Few approaches are related to clustering, yet noise injected is big — the existence of a smaller, sufficient noise, was conjectured in (Nissim et al., 2007) — and approaches rely on a variety of assumptions or knowledge about the optimum (See Table 1) (Nissim et al., 2007; Wang et al., 2015). To apply k -variates++, we consider that $\wp_t = \text{Id}, \forall t$, and assume $0 < R \ll \infty$ s.t. $\max_{\mathbf{a}, \mathbf{a}' \in A} \|\mathbf{a} - \mathbf{a}'\|_2 \leq R$ (a current assumption in the field (Dwork & Roth, 2014)).

A general likelihood ratio bound for k -variates++ We show that the likelihood ratio of the same clustering for two “close” instances is governed by two quantities that rely on the neighborhood function. Most importantly for differential privacy, when densities $p(\mu_a, \theta_a)$ are carefully chosen, this ratio *always* $\rightarrow 1$ as a function of m , which is highly desirable for differential privacy. We let $\text{NN}_{\mathcal{N}}(\mathbf{a}) \doteq \arg \min_{\mathbf{a}' \in \mathcal{N}} \|\mathbf{a} - \mathbf{a}'\|_2$ denote the nearest neighbour of \mathbf{a} in \mathcal{N} , and let $\mathbf{c}(A) \doteq (1/|A|) \cdot \sum_{\mathbf{a} \in A} \mathbf{a}$.

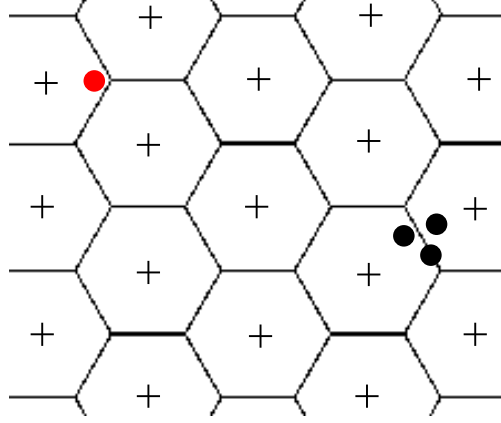


Figure 2: Checking that δ_s is small, for \mathcal{N} the set of crosses (+). Any set A of points close to each other, such as the black dots (•), would be \mathcal{N} -packed (pick $\mathbf{x} = \mathbf{c}(A)$ in this case), but would fail to be \mathcal{N} -packed if too spread (e.g., red dot (•) plus black dots). Segments depict the Voronoi diagram of \mathcal{N} . Best viewed in color.

Definition 7 We say that neighborhood in \mathcal{A} is δ_w -spread for some $\delta_w > 0$ iff for any $\mathcal{N} \subseteq \mathcal{A}$ with $|\mathcal{N}| = k - 1$, and any $\mathcal{B} \subseteq \mathcal{A}$ with $|\mathcal{B}| = |\mathcal{A}| - 1$,

$$\sum_{\mathbf{a} \in \mathcal{B}} \|\mathbf{a} - \text{NN}_{\mathcal{N}}(\mathbf{a})\|_2^2 \geq \frac{R^2}{\delta_w}. \quad (11)$$

Definition 8 We say that neighborhood in \mathcal{A} is δ_s -monotonic for some $\delta_s > 0$ iff the following holds. $\forall \mathcal{N} \subseteq \mathcal{A}$ with $|\mathcal{N}| \in \{1, 2, \dots, k - 1\}$, for any $A \subseteq \mathcal{A} \setminus \mathcal{N}$ which is \mathcal{N} -packed, we have:

$$\begin{aligned} \sum_{\mathbf{a} \in A} \|\mathbf{a} - \text{NN}_{\mathcal{N}}(\mathbf{a})\|_2^2 \\ \leq (1 + \delta_s) \cdot \sum_{\mathbf{a} \in A} \|\mathbf{a} - \text{NN}_{\mathcal{N} \cup \{\mathbf{c}(A)\}}(\mathbf{a})\|_2^2. \end{aligned} \quad (12)$$

Set A is said \mathcal{N} -packed iff there exists $\mathbf{x} \in \mathbb{R}^d$ satisfying $\mathbf{x} = \arg \min_{\mathbf{c} \in \mathcal{N} \cup \{\mathbf{x}\}} \|\mathbf{a} - \mathbf{c}\|_2^2, \forall \mathbf{a} \in A$.

It is worthwhile remarking that as long as $k < |\mathcal{A}| \ll \infty$, both $0 < \delta_w \ll \infty$ and $0 < \delta_s \ll \infty$ always exist. Informally, δ_w brings that the sum of squared distances to any subset of $k - 1$ centers in \mathcal{A} must not be negligible against the diameter R . δ_s yields a statement a bit more technical, but it roughly reduces to stating that adding one center to any set of at most $k - 1$ points that are already close to each other should not decrease significantly the overall potential to the set of centers. Figure 2 provides a schematic view of the property, showing that the modifications of the potential can be very local, thus yielding small δ_s in ineq. (12). The following Theorem uses the definition of neighbouring samples: samples \mathcal{A} and \mathcal{A}' are neighbours, written $\mathcal{A} \approx \mathcal{A}'$, iff they differ by one point. We also define $\mathbb{P}[\mathcal{C}|\mathcal{A}]$ to be the density of output \mathcal{C} given input data \mathcal{A} .

Theorem 9 Fix $\wp_t = \text{Id} (\forall t)$ and densities $p_{(\mu, \theta)}$ having the same support Ω in k -variates++. Suppose there exists $\varrho(R) > 0$ such that densities $p_{(\mu, \theta)}$ satisfy the following pointwise likelihood

ratio constraint:

$$\frac{p(\mu_{\mathbf{a}'}, \theta_{\mathbf{a}'}) (\mathbf{x})}{p(\mu_{\mathbf{a}}, \theta_{\mathbf{a}}) (\mathbf{x})} \leq \varrho(R) , \forall \mathbf{a}, \mathbf{a}' \in \mathcal{A}, \forall \mathbf{x} \in \Omega . \quad (13)$$

Then, there exists a function $f(\cdot)$ such that, for any $\delta_w, \delta_s > 0$ such that \mathcal{A} is δ_w -spread and δ_s -monotonic, for any $\mathcal{A}' \approx \mathcal{A}$, for any $k > 0$ and any $\mathcal{C} \subset \Omega$ of size k output by Algorithm `k-variates++` on whichever of \mathcal{A} or \mathcal{A}' , the likelihood ratio of \mathcal{C} given \mathcal{A} and \mathcal{A}' is upperbounded as:

$$\frac{\mathbb{P}[\mathcal{C}|\mathcal{A}']}{\mathbb{P}[\mathcal{C}|\mathcal{A}]} \leq (1 + \delta_w)^{k-1} + f(k) \cdot \delta_w \cdot (1 + \delta_s)^{k-1} \cdot \varrho(R) . \quad (14)$$

(Proof in page 28) Notice that Theorem 9 makes just one assumption (13) about the densities, so it can be applied in fairly general settings, such as for regular exponential families (Banerjee et al., 2005). These are a key choice because they extensively cover the domain of distortions for which the average is the population minimiser.

An (almost) distribution-free $1 + o(1)$ likelihood ratio We now show that if \mathcal{A} is sampled i.i.d. from any distribution \mathcal{D} which satisfies the mild assumption that it is locally bounded everywhere (or almost surely) in a ball, then with high probability the right-hand side of ineq. (14) is $1 + o(1)$ where the little-oh vanishes with m . The proof, of independent interest, involves an explicit bound on δ_w and δ_s .

Theorem 10 Suppose \mathcal{A} with $|\mathcal{A}| = m > 1$ sampled i.i.d. from distribution \mathcal{D} whose support contains a L_2 ball $\mathcal{B}_2(\mathbf{0}, R)$ with density inside in between $\epsilon_m > 0$ and $\epsilon_M \geq \epsilon_m$. Let $\rho_{\mathcal{D}} \doteq \epsilon_M / \epsilon_m (\geq 1)$. For any $0 < \delta < 1/2$, if (i) $\mathcal{A} \subset \mathcal{B}(\mathbf{0}, R)$ and (ii) the number of clusters k meets:

$$k \leq \frac{\delta^2}{4\rho_{\mathcal{D}}} \cdot \sqrt{m} , \quad (15)$$

then there is probability $1 - \delta$ over the sampling of \mathcal{A} that `k-variates++`, instantiated as in Theorem 9, satisfies $\mathbb{P}[\mathcal{C}|\mathcal{A}'] / \mathbb{P}[\mathcal{C}|\mathcal{A}] \leq 1 + \rho_{\mathcal{D}}^k \cdot g(m, k, d, R)$, $\forall \mathcal{A}' \approx \mathcal{A}$, with

$$g(m, k, d, R) \doteq \frac{4}{m^{\frac{1}{4} + \frac{1}{d+1}}} + \left(\frac{64}{k^{\frac{2}{d}}} \right)^k \cdot \frac{\varrho(2R)}{m} . \quad (16)$$

(Proof in page 34) The key informal statement of Theorem 10 is that one may obtain with high probability some “good” datasets \mathcal{A} , i.e., for which δ_w, δ_s are small, under very weak assumptions about the domain at hand. The key point is that if one has access to the sampling, then one can resample datasets \mathcal{A} until a good one comes.

Applications to differential privacy Let \mathcal{M} be any algorithm which takes as input \mathcal{A} and k , and returns a set of k centers \mathcal{C} . Let $\mathbb{P}_{\mathcal{M}}[\mathcal{C}|\mathcal{A}]$ denote the probability, over the internal randomisation of \mathcal{M} , that \mathcal{M} returns \mathcal{C} given \mathcal{A} and k (k , fixed, is omitted in notations). Following is the definition of differential privacy (Dwork et al., 2006), tailored for conciseness to our clustering problem.

Definition 11 \mathcal{M} is ϵ -differentially private (DP) for k clusters iff for any neighbors $\mathcal{A} \approx \mathcal{A}'$, set \mathcal{C} of k centers,

$$\mathbb{P}_{\mathcal{M}}[\mathcal{C}|\mathcal{A}']/\mathbb{P}_{\mathcal{M}}[\mathcal{C}|\mathcal{A}] \leq \exp \epsilon . \quad (17)$$

A relaxed version of ϵ -DP is (ϵ, δ) -DP, in which we require $\mathbb{P}_{\mathcal{M}}[\mathcal{C}|\mathcal{A}'] \leq \mathbb{P}_{\mathcal{M}}[\mathcal{C}|\mathcal{A}] \cdot \exp \epsilon + \delta$; thus, ϵ -DP = $(\epsilon, 0)$ -DP (Dwork & Roth, 2014). We show that low noise may be affordable to satisfy ineq. (17) using Laplace distribution, $Lap(\sigma/\sqrt{2})$. We refer to the *Laplace mechanism* as a popular mechanism which adds to the output of an algorithm a sufficiently large amount of Laplace noise to be ϵ -DP. We refer to (Dwork et al., 2006) for details, and assume from now on that data belong to a L_1 ball $\mathcal{B}_1(\mathbf{0}, R)$.

Theorem 12 Using notations and setting of Theorem 9, let

$$\tilde{\epsilon} \doteq \log \left(\frac{\exp(\epsilon) - (1 + \delta_w)^{k-1}}{f(k) \cdot \delta_w \cdot (1 + \delta_s)^{k-1}} \right) . \quad (18)$$

Then, k -variates++ with $p_{(\mu, \theta)}$ a product of $Lap(\sigma_1/\sqrt{2})$, for $\sigma_1 \doteq 2\sqrt{2}R/\tilde{\epsilon}$, both meets ineq. (17) **and** its expected potential satisfies ineq. (7) with

$$\Phi = \Phi_1 \doteq 8 \cdot \left(\phi_{\text{opt}} + \frac{mR^2}{\tilde{\epsilon}^2} \right) . \quad (19)$$

On the other hand, if we opt for $\sigma_2 \doteq 2\sqrt{2}kR/\epsilon$, then k -variates++ is an instance of the Laplace mechanism **and** its expected potential satisfies ineq. (7) with

$$\Phi = \Phi_2 \doteq 8 \cdot \left(\phi_{\text{opt}} + \frac{mk^2R^2}{\epsilon^2} \right) . \quad (20)$$

(Proof in page 41) A question is how do σ_1 (resp. Φ_1) and σ_2 (resp. Φ_2) compare with each other, and how do they compare to the state of the art (Nissim et al., 2007; Wang et al., 2015) (we only consider methods with provable approximation bounds of the global optimum). The key fact is that, if m is sufficiently large, then it happens that we can fix $\delta_w = O(1/m)$ and $\delta_s = O(1)$. The proof of Theorem 10 (page 34) and the experiments (page 44) display that such regimes *are* indeed observed. In this case, it is not hard to show that $\tilde{\epsilon} = \Omega(\epsilon + \log m)$, granting $\sigma_1 = o(\sigma_2)$ since

$$\sigma_1 = O \left(\frac{R}{\epsilon + \log(m)} \right) , \quad (21)$$

i.e. the noise guaranteeing ineq. (17) *vanishes* at $1/\log(m)$ rate. Consequently, in this regime, Φ_1 in eq. (19) becomes:

$$\Phi_1 = \tilde{O} \left(\phi_{\text{opt}} + \frac{mR^2}{(\epsilon + \log m)^2} \right) , \quad (22)$$

ignoring all factors other than those noted. Thus, the noise dependence grows *sublinearly* in m . Since in this setting, unless all datapoints are the same, δ_w and δ_s for \mathcal{A} and *any* possible neighbor

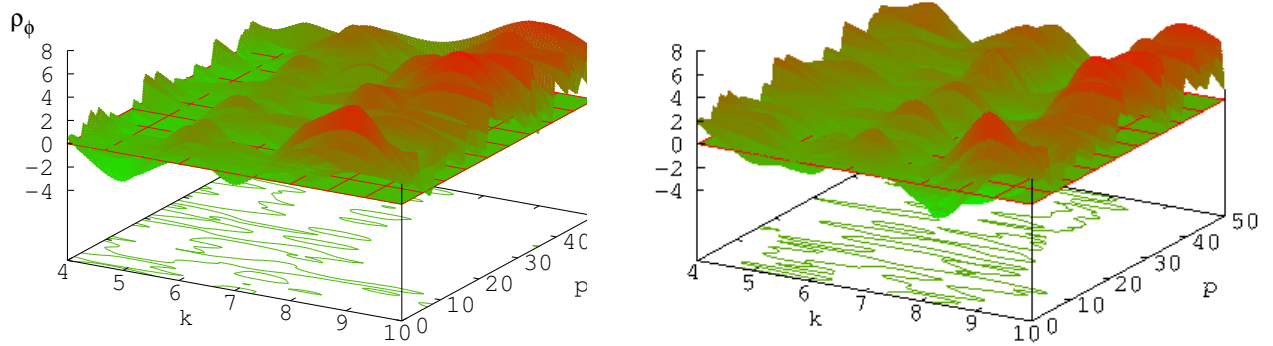


Figure 3: Plot of $\rho_\phi(\mathcal{H}) = f(k, p)$ (points below $z = 0$ — isocontour shown — correspond to superior performances for Dk -means++). Left: $\mathcal{H}=k$ -means++; right: $\mathcal{H}=k$ -means $_{\parallel}$ (best viewed in color).

\mathcal{A}' are within $1 + o(1)$, it is also possible to overestimate δ_w and δ_s to still have $\delta_w = O(1/m)$ and $\delta_s = O(1)$ **and** grant ϵ -DP for k -variates++. Otherwise, the setting of Theorem 10 can be used to grant (ϵ, δ) -DP without any tweak. Table 1 compares k -variates++ to (Nissim et al., 2007; Wang et al., 2015) in this large sample regime, which is actually a prerequisite for (Nissim et al., 2007; Wang et al., 2015). Notation O^* removes all dependencies in their model parameters (assumptions, model parameters, and δ for the (ϵ, δ) -DP in (Wang et al., 2015)), and λ is the separability assumption parameter (Nissim et al., 2007)⁴. The approximation bounds in (Nissim et al., 2007) consider Wasserstein distance between (estimated / optimal) centers, and not the potential involving data points like us. To obtain bounds that can be compared, we have used the simple trick that the observed potential is, up to a constant, no more than the optimal potential plus a function of the distance between (estimated / optimal) centers. This somewhat degrades the bound, but not enough for the observed discrepancies with our bound to reverse or even vanish. It is clear from the bounds that the noise dependence is significantly in our favor, and our bound is also significantly better at least when k is not too large.

5 Experiments

The experiments carried out are provided *in extenso* in the Appendix (from page 44).

Dk -means++ vs k -means++ and k -means $_{\parallel}$ (Bahmani et al., 2012) To address algorithms that can be reduced from k -variates++ (Section 3), we have tested Dk -means++ vs state of the art approach k -means $_{\parallel}$; to be fair with Dk -means++, we use k -means++ seeding as the reclustering algorithm in k -means $_{\parallel}$. Parameters are in line with (Bahmani et al., 2012). To control the spread of Forgy nodes ϕ_s^F (Theorem 4), each peer’s initial data consists of points uniformly sampled in a random hyperrectangle in a space of $d = 50$ (expected number of peers points $m_i = 500, \forall i$). We sample

⁴ λ is named ϕ in (Nissim et al., 2007). We use λ to avoid confusion with clustering potentials.

Dataset	m	d	\bar{k}	$(\bar{\epsilon}/\epsilon)$	$\overline{\rho'_\phi(\text{F-DP})}$	$\overline{\rho'_\phi(\text{GUPT})}$
LifeSci	26 733	10	3	4.5	163.0	0.7
Image	34 112	3	2.5	7.9	188.5	2.9
EuropeDiff	169 308	2	5	13.0	2857.1	40.4

Table 2: k -variates++ vs F-DP and GUPT (see text).

peers until a total of $m \approx 20000$ point is sampled. Then, each point moves with $p\%$ chances to a uniformly sampled peer. We checked that ϕ_s^F blows up with p , *i.e.*, >20 times for $p = 50\%$ with respect to $p = 0$. A remarkable phenomenon was the fact that, even when the number of peers n is quite large (dozens on average), Dk -means++ is able to *beat* both k -means++ and k -means $_{||}$, even for large values of p , as computed by ratio $\rho_\phi(\mathcal{H}) \doteq 100 \cdot (\phi(Dk\text{-means++}) - \phi(\mathcal{H}))/\phi(\mathcal{H})$ for $\mathcal{H} \in \{k\text{-means++}, k\text{-means}_{||}\}$ (Figure 3). Another positive point is that the amount of data to compute a center for Dk -means++ is in average $\approx n$ times smaller than k -means $_{||}$.

The fact that Dk -means++, which locally implements the biased seeding, may be able to beat k -means++, which globally implements this seeding technique, is not surprising, and in fact may come from the leverage brought by the compartmentalization of distributed data: as discussed in deeper details in page 47, this may even improve the approximability ratio of Dk -means++ so that it *beats* the AV bound.

k -variates++ vs Forgy-DP and GUPT To address algorithms that can be obtained via a direct use of k -variates++ (Section 4), we have tested it in a differential privacy framework vs state of the art approach GUPT (Mohan et al., 2012). We let $\epsilon = 1$ in our experiments. We also compare it to Forgy DP (F-DP), which is just Forgy initialisation in the Laplace mechanism, with noise rate (standard dev.) $\propto kR/\epsilon$. In comparison, the noise rate for GUPT is $\propto kR/(\ell\epsilon)$ at the end of its aggregation process, where ℓ is the number of blocks. Table 2 gives results for the average (over the choices of k) parameters used, \bar{k} , $\bar{\epsilon}$, and ratio $\overline{\rho'_\phi}$ where $\rho'_\phi(\mathcal{H}) \doteq \phi(\mathcal{H})/\phi(k\text{-variates++})$ — values above 1 indicate better results for k -variates++. We use $\bar{\epsilon}$ as the equivalent ϵ for k -variates++, *i.e.* the value that guarantees ineq. (17). From Theorem 12, when $\bar{\epsilon} > \epsilon$, this brings a smaller noise magnitude, desirable for clustering. The obtained results show that k -variates++ becomes more of a contender with increasing m , but its relative performance tends to decrease with increasing k . This is in accordance with the “good” regime of Theorem 12. Results on synthetic domains display the same patterns, along with the fact that relative performances of k -variates++ improves with d , making it a relevant choice for “big” domains.

In fact, extensive experiments on synthetic data (page 44) show that intuitions regarding the sublinear noise regime in eq. (22) are experimentally observed, and furthermore they may happen for quite small values of m .

6 Discussion and Conclusion

We first show in this paper that the k -means++ analysis of Arthur and Vassilvitskii can be carried out on a significantly more general scale, aggregating various clustering frameworks of interest and for which no trivial adaptation of k -means++ was previously known. Our contributions stand at two levels: (i) we provide the “meta” algorithm, k -variates++, and two key results, one on its

approximation abilities of the *global* optimum, and one on the *likelihood ratio* of the centers it delivers. We do expect further applications of these results, in particular to address several other key clustering problems: stability, generalisation and smoothed analysis (Arthur et al., 2011; von Luxburg, 2010); (ii) we provide two examples of application. The first is a reduction technique from k -variates++, which shows a way to obtain straight approximability results for other clustering algorithms, some being efficient proxies for the generalisation of existing approaches (Ailon et al., 2009). The second is a direct application of k -variates++ to differential privacy, exhibiting a noise component significantly better than existing approaches (Nissim et al., 2007; Wang et al., 2015).

We have not discussed here the possibility to replace the L_2^2 distortion which computes the potential by elements from large and interesting classes — clustering being a huge practical problem, it is indeed reasonable to tailor the distortion to the application at hand. One example are Bregman divergences, that fail simple metric transforms (Acharyya et al., 2013). Another example are total divergences, that fail the simple computation of the population minimizers (Nock et al., 2016; Liu et al., 2012). Some do not even admit population minimizers in closed form (Nielsen & Nock, 2015). It turns out that k -variates++, and its good approximation properties, *can* be extended to such cases (see page 42) for total Jensen divergence (Nielsen & Nock, 2015).

7 Acknowledgments

Thanks are due to Stephen Hardy, Guillaume Smith, Wilko Henecka and Max Ott for stimulating discussions and feedback on the subject. Nicta is funded by the Australian Government through the Department of Communications and the Australian Research Council through the ICT Center of Excellence Program.

References

- Acharyya, S., Banerjee, A., and Boley, D. Bregman divergences and triangle inequality. In *Proc. of the 13th SIAM International Conference on Data Mining*, pp. 476–484, 2013.
- Ackermann, M.-R., Lammersen, C., Mörtens, M., Raupach, C., Sohler, C., and Swierkot, K. Streamkm++: A clustering algorithms for data streams. In *12th ALENEX*, pp. 173–187, 2010.
- Ailon, N., Jaiswal, R., and Monteleoni, C. Streaming k -means approximation. In *NIPS*22*, pp. 10–18, 2009.
- Arthur, D. and Vassilvitskii, S. k -means++ : the advantages of careful seeding. In *19th SODA*, pp. 1027 – 1035, 2007.
- Arthur, D., Manthey, B., and Röglin, H. Smoothed analysis of the k -means method. *JACM*, 58:19, 2011.
- Bahmani, B., Moseley, B., Vattani, A., Kumar, R., and Vassilvitskii, S. Scalable k -means++. In *38th VLDB*, pp. 622–633, 2012.
- Balcan, M.-F., Ehrlich, S., and Liang, Y. Distributed k -means and k -median clustering on general communication topologies. In *NIPS*26*, pp. 1995–2003, 2013.

- Banerjee, A., Merugu, S., Dhillon, I., and Ghosh, J. Clustering with Bregman divergences. *JMLR*, 6:1705–1749, 2005.
- Boissonnat, J.-D., Nielsen, F., and Nock, R. Bregman voronoi diagrams. *DCG*, 44(2):281–307, 2010.
- Chaudhuri, K., Monteleoni, C., and Sarwate, A.-D. Differentially private empirical risk minimization. *JMLR*, 12:1069–1109, 2011.
- Chichignoud, M. and Lousteau, S. Adaptive noisy clustering. *IEEE Trans. IT*, 60:7279–7292, 2014.
- Dwork, C. and Roth, A. The algorithmic foundations of differential privacy. *Found. & Trends in TCS*, 9:211–407, 2014.
- Dwork, C., McSherry, F., Nissim, K., and Smith, A. Calibrating noise to sensitivity in private data analysis. In *3rd TCC*, pp. 265–284, 2006.
- Har-Peled, S. and Mazumdar, S. On coresets for k -means and k -median clustering. In *37th ACM STOC*, pp. 291–300, 2004.
- Hardt, M. and Price, E. The noisy power method: a meta algorithm with applications. In *NIPS*27*, pp. 2861–2869, 2014.
- Indyk, P., Mahabadi, S., Mahdian, M., and Mirrokni, V.-S. Composable core-sets for diversity and coverage maximization. In *33rd ACM PODS*, pp. 100–108, 2014.
- Jegelka, S., Sra, S., and Banerjee, A. Approximation algorithms for tensor clustering. In *20th ALT*, pp. 368–383, 2009.
- Kalai, A. and Vempala, S. Efficient algorithms for online decision problems. *J. Comp. Syst. Sc.*, pp. 291–307, 2005.
- Kohavi, R. and Wolpert, D. Bias plus variance decomposition for zero-one loss functions. In *13th ICML*, pp. 275–283, 1996.
- Liberty, E., Sriharsha, R., and Sviridenko, M. An algorithm for online k -means clustering. *CoRR*, abs/1412.5721, 2014.
- Liu, M., Vemuri, B.-C., i Amari, S., and Nielsen, F. Shape retrieval using hierarchical total bregman soft clustering. *IEEE Trans. PAMI*, 34(12):2407–2419, 2012.
- McSherry, F. Privacy integrated queries: an extensible platform for privacy-preserving data analysis. *Communications of the ACM*, 53(9):89–97, 2010.
- Mohan, P., Thakurta, A., Shi, E., Song, D., and Culler, D.-E. GUPT: privacy preserving data analysis made easy. In *38th ACM SIGMOD*, pp. 349–360, 2012.
- Nielsen, F. and Nock, R. Total Jensen divergences: definition, properties and clustering. In *40th IEEE ICASSP*, pp. 2016–2020, 2015.

- Nissim, K., Raskhodnikova, S., and Smith, A. Smooth sensitivity and sampling in private data analysis. In *40th ACM STOC*, pp. 75–84, 2007.
- Nock, R., Luosto, P., and Kivinen, J. Mixed Bregman clustering with approximation guarantees. In *19th ECML*, pp. 154–169, 2008.
- Nock, R., Nielsen, F., and Amari, S.-I. On conformal divergences and their population minimizers. *IEEE Trans. IT*, 62:1–12, 2016.
- Shindler, M., Wong, A., and Meyerson, A. Fast and accurate k -means for large datasets. In *NIPS*24*, pp. 2375–2383, 2011.
- von Luxburg, U. Clustering stability: an overview. *Found. & Trends in ML*, 2(3):235–274, 2010.
- Wang, Y., Wang, Y.-X., and Singh, A. Differentially private subspace clustering. In *NIPS*28*, 2015.

Appendix — Table of contents

Appendix on proofs	Pg 19
Proof of Theorem 2	Pg 19
Proof of Lemma 3	Pg 22
Comments on Table 1	Pg 22
Proofs of Theorems 4, 5 and 6	Pg 23
\hookrightarrow Proof of Theorem 4	Pg 23
\hookrightarrow Proof of Theorem 5	Pg 25
\hookrightarrow Proof of Theorem 6	Pg 26
Proof of Theorem 9	Pg 28
Proof of Theorem 10	Pg 34
Proof of Theorem 12	Pg 41
Extension to non-metric spaces	Pg 42
 Appendix on experiments	 Pg 44
Experiments on Theorem 12 and the sublinear noise regime	Pg 44
Experiments with Dk -means++, k -means++ and k -means	Pg 47
Experiments with k -variates++ and GUPT	Pg 51

8 Appendix on Proofs

Several proofs rely on properties of the k -means++ algorithm that are not exploited in the proof of (Arthur & Vassilvitskii, 2007). We assume here the basic knowledge of the proof technique of (Arthur & Vassilvitskii, 2007).

Proof of Theorem 2

Let A denote a subset of \mathcal{A} , and $\mathbf{c}(A) \doteq (1/|A|) \cdot \sum_{\mathbf{a} \in A} \mathbf{a}$ the barycenter of A . It is well known that $\mathbf{c}(A) = \arg \min_{\mathbf{a}' \in \mathbb{R}^d} \sum_{\mathbf{a} \in A} \|\mathbf{a} - \mathbf{a}'\|_2^2$, so the potential of A ,

$$\phi(A) \doteq \sum_{\mathbf{a} \in A} \|\mathbf{a} - \mathbf{c}(A)\|_2^2 \quad (23)$$

is just the optimal potential of A if A defines a cluster in the optimal clustering. We also define the noisy potential of A as:

$$\phi^N(A) \doteq \sum_{\mathbf{a} \in A} \int_{\Omega_{\mathbf{a}}} \|\mathbf{x} - \mathbf{c}(A)\|_2^2 d p_{\mathbf{a}}(\mathbf{x}) . \quad (24)$$

The proof of Theorem 2 follows the same path as the proof of Theorem 3.1 in (Arthur & Vassilvitskii, 2007). Instead of reproducing the proof, we shall assume basic knowledge of the original proof and will just provide the side Lemmata that are sufficient for our more general result. The first Lemma is a generalization of Lemma 3.2 in (Arthur & Vassilvitskii, 2007).

Lemma 13 *Let C_{opt} denotes the optimal partition of \mathcal{A} according to eq. (2). Let A be an arbitrary cluster in C_{opt} . Let C be a single-cluster clustering whose center is chosen at random by one step of Algorithm k -variates++ (i.e. for $t = 1$). Then*

$$\mathbb{E}[\phi(A)] = \phi_{\text{opt}}(A) + \phi_{\text{opt}}^N(A) . \quad (25)$$

Proof The expected potential of cluster A is

$$\begin{aligned}
& \mathbb{E}[\phi(A; \mathcal{C} = \emptyset)] \\
&= \frac{1}{|A|} \cdot \sum_{\mathbf{a}_0 \in A} \int_{\Omega_{\mathbf{a}_0}} \sum_{\mathbf{a} \in A} \|\mathbf{a} - \mathbf{x}\|_2^2 dp_{\mathbf{a}_0}(\mathbf{x}) \\
&= \frac{1}{|A|} \cdot \sum_{\mathbf{a}_0 \in A} \int_{\Omega_{\mathbf{a}_0}} \sum_{\mathbf{a} \in A} \|\mathbf{a} - \mathbf{c}(A) + \mathbf{c}(A) - \mathbf{x}\|_2^2 dp_{\mathbf{a}_0}(\mathbf{x}) \\
&= \frac{1}{|A|} \cdot \sum_{\mathbf{a}_0 \in A} \left(\sum_{\mathbf{a} \in A} \|\mathbf{a} - \mathbf{c}(A)\|_2^2 + |A| \cdot \int_{\Omega_{\mathbf{a}_0}} \|\mathbf{x} - \mathbf{c}(A)\|_2^2 dp_{\mathbf{a}_0}(\mathbf{x}) \right. \\
&\quad \left. + 2 \sum_{\mathbf{a} \in A} \langle \mathbf{a} - \mathbf{c}(A), \mathbf{c}(A) - \int_{\Omega_{\mathbf{a}_0}} \mathbf{x} dp_{\mathbf{a}_0}(\mathbf{x}) \rangle \right) \\
&= \frac{1}{|A|} \cdot \sum_{\mathbf{a}_0 \in A} \left(\sum_{\mathbf{a} \in A} \|\mathbf{a} - \mathbf{c}(A)\|_2^2 + |A| \cdot \int_{\Omega_{\mathbf{a}_0}} \|\mathbf{x} - \mathbf{c}(A)\|_2^2 dp_{\mathbf{a}_0}(\mathbf{x}) \right. \\
&\quad \left. + 2 \underbrace{\left\langle \sum_{\mathbf{a} \in A} \mathbf{a} - |A| \mathbf{c}(A), \mathbf{c}(A) - \mathbf{a}_0 \right\rangle}_{=0} \right) \\
&= \sum_{\mathbf{a} \in A} \|\mathbf{a} - \mathbf{c}(A)\|_2^2 + \sum_{\mathbf{a} \in A} \int_{\Omega_{\mathbf{a}_0}} \|\mathbf{x} - \mathbf{c}(A)\|_2^2 dp_{\mathbf{a}}(\mathbf{x}) \\
&= \phi_{\text{opt}}(A) + \phi_{\text{opt}}^N(A) ,
\end{aligned}$$

as claimed. ■

When $p_{\mathbf{a}}$ is a Dirac anchored at \mathbf{a} , we recover Lemma 3.2 in (Arthur & Vassilvitskii, 2007). The following Lemma generalizes Lemma 3.3 in (Arthur & Vassilvitskii, 2007).

Lemma 14 *Suppose that the optimal clustering C_{opt} is η -probe approximable. Let A be an arbitrary cluster in C_{opt} , and let C be an arbitrary clustering with centers \mathcal{C} . Suppose that the reference point \mathbf{a} chosen according to (1) in Step 2.1 is in A . Then the random point \mathbf{x} picked in Step 2.2 brings an expected potential that satisfies*

$$\mathbb{E}[\phi(A)] \leq (6 + 4\eta) \cdot \phi_{\text{opt}}(A) + 2 \cdot \phi_{\text{opt}}^N(A) . \quad (26)$$

Proof Let us denote $\mathbf{c}^*(\mathbf{u}) \doteq \arg \min_{\mathbf{c} \in \mathcal{C}} \|\mathbf{u} - \mathbf{c}\|_2^2$ (since $C \neq C_{\text{opt}}$ in general, $\mathbf{c}^*(\mathbf{u}) \neq \mathbf{c}_{\text{opt}}(\mathbf{u})$), and $D(\mathbf{a}) \doteq \|\mathbf{a} - \mathbf{c}^*(\mathbf{a})\|_2^2$ the contribution of $\mathbf{a} \in A$ to the k -means potential defined by \mathcal{C} . We have, using Lemma 3.3 in (Arthur & Vassilvitskii, 2007) and Lemma 13,

$$\mathbb{E}_{\mathbf{x}}[\phi(A; \mathcal{C} \cup \{\mathbf{x}\})] = \sum_{\mathbf{a}_0 \in A} \frac{D_t(\mathbf{a}_0)}{\sum_{\mathbf{a} \in A} D_t(\mathbf{a})} \cdot \sum_{\mathbf{a} \in A} \int_{\Omega_{\mathbf{a}_0}} \min\{D(\mathbf{a}), \|\mathbf{a} - \mathbf{x}\|_2^2\} dp_{\mathbf{a}_0}(\mathbf{x}) . \quad (27)$$

The triangle inequality gives, for any $\mathbf{a} \in A$,

$$\begin{aligned}
\sqrt{D_t(\mathbf{a}_0)} &\doteq \|\wp_t(\mathbf{a}_0) - \mathbf{c}^*(\wp_t(\mathbf{a}_0))\|_2 \\
&\leq \|\wp_t(\mathbf{a}_0) - \mathbf{c}^*(\wp_t(\mathbf{a}))\|_2 \\
&\leq \|\wp_t(\mathbf{a}_0) - \wp_t(\mathbf{a})\|_2 + \|\wp_t(\mathbf{a}) - \mathbf{c}^*(\wp_t(\mathbf{a}))\|_2 ;
\end{aligned} \quad (28)$$

since $(a + b)^2 \leq 2a^2 + 2b^2$, then $D_t(\mathbf{a}_0) \leq 2\|\wp_t(\mathbf{a}_0) - \wp_t(\mathbf{a})\|_2^2 + 2D_t(\mathbf{a})$, and so, after averaging over A ,

$$D_t(\mathbf{a}_0) \leq \frac{2}{|A|} \sum_{\mathbf{a} \in A} \|\wp_t(\mathbf{a}_0) - \wp_t(\mathbf{a})\|_2^2 + \frac{2}{|A|} \sum_{\mathbf{a} \in A} D_t(\mathbf{a}) , \quad (29)$$

and eq. (27) can be upperbounded as:

$$\begin{aligned}
\mathbb{E}_{\mathbf{x}}[\phi(A; \mathcal{C} \cup \{\mathbf{x}\})] &\leq \frac{2}{|A|} \sum_{\mathbf{a}_0 \in A} \frac{\sum_{\mathbf{a} \in A} \|\wp_t(\mathbf{a}_0) - \wp_t(\mathbf{a})\|_2^2}{\sum_{\mathbf{a} \in A} D_t(\mathbf{a})} \cdot \sum_{\mathbf{a} \in A} \int_{\Omega_{\mathbf{a}_0}} \min\{D(\mathbf{a}), \|\mathbf{a} - \mathbf{x}\|_2^2\} dp_{\mathbf{a}_0}(\mathbf{x}) \\
&\quad + \frac{2}{|A|} \sum_{\mathbf{a}_0 \in A} \frac{\sum_{\mathbf{a} \in A} D_t(\mathbf{a})}{\sum_{\mathbf{a} \in A} D_t(\mathbf{a})} \cdot \sum_{\mathbf{a} \in A} \int_{\Omega_{\mathbf{a}_0}} \min\{D(\mathbf{a}), \|\mathbf{a} - \mathbf{x}\|_2^2\} dp_{\mathbf{a}_0}(\mathbf{x}) \\
&\leq \underbrace{\frac{2}{|A|} \sum_{\mathbf{a}_0 \in A} \frac{\sum_{\mathbf{a} \in A} D(\mathbf{a})}{\sum_{\mathbf{a} \in A} D_t(\mathbf{a})} \cdot \sum_{\mathbf{a} \in A} \|\wp_t(\mathbf{a}_0) - \wp_t(\mathbf{a})\|_2^2}_{\doteq P_1} \\
&\quad + \underbrace{\frac{2}{|A|} \sum_{\mathbf{a}_0 \in A} \sum_{\mathbf{a} \in A} \int_{\Omega_{\mathbf{a}_0}} \|\mathbf{a} - \mathbf{x}\|_2^2 dp_{\mathbf{a}_0}(\mathbf{x})}_{\doteq P_2} . \tag{30}
\end{aligned}$$

We bound the two potentials P_1 and P_2 separately, starting with P_1 . Fix any $\mathbf{a}_0 \in A$. If $\sum_{\mathbf{a} \in A} \|\wp_t(\mathbf{a}) - \wp_t(\mathbf{a}_0)\|_2^2 = 0$, then trivially

$$\left(\sum_{\mathbf{a} \in A} D(\mathbf{a}) \right) \cdot \left(\sum_{\mathbf{a} \in A} \|\wp_t(\mathbf{a}_0) - \wp_t(\mathbf{a})\|_2^2 \right) \leq (1 + \eta) \cdot \left(\sum_{\mathbf{a} \in A} D_t(\mathbf{a}) \right) \cdot \left(\sum_{\mathbf{a} \in A} \|\mathbf{a}_0 - \mathbf{a}\|_2^2 \right) \tag{31}$$

since the right-hand side cannot be negative. If $\sum_{\mathbf{a} \in A} \|\wp_t(\mathbf{a}) - \wp_t(\mathbf{a}_0)\|_2^2 \neq 0$, then since \wp_t is η -stretching, we have:

$$\frac{\sum_{\mathbf{a} \in A} \|\mathbf{a} - \mathbf{c}^*(\mathbf{a})\|_2^2}{\sum_{\mathbf{a} \in A} \|\mathbf{a} - \mathbf{a}_0\|_2^2} \leq (1 + \eta) \cdot \frac{\sum_{\mathbf{a} \in A} \|\wp_t(\mathbf{a}) - \mathbf{c}^*(\wp_t(\mathbf{a}))\|_2^2}{\sum_{\mathbf{a} \in A} \|\wp_t(\mathbf{a}) - \wp_t(\mathbf{a}_0)\|_2^2} , \tag{32}$$

which is exactly ineq. (31) after rearranging the terms. Ineq (31) implies

$$\begin{aligned}
P_1 &\leq 2(1 + \eta) \cdot \frac{1}{|A|} \sum_{\mathbf{a}_0 \in A} \sum_{\mathbf{a} \in A} \|\mathbf{a}_0 - \mathbf{a}\|_2^2 \\
&= 4(1 + \eta) \cdot \phi_{\text{opt}}(A) , \tag{33}
\end{aligned}$$

where the equality follows from (Arthur & Vassilvitskii, 2007), Lemma 3.2. Also, Lemma 13 brings

$$\begin{aligned}
P_2 &= 2 \cdot \frac{1}{|A|} \sum_{\mathbf{a}_0 \in A} \int_{\Omega_{\mathbf{a}_0}} \sum_{\mathbf{a} \in A} \|\mathbf{a} - \mathbf{x}\|_2^2 dp_{\mathbf{a}_0}(\mathbf{x}) \\
&= 2\phi_{\text{opt}}(A) + 2\phi_{\text{opt}}^N(A) . \tag{34}
\end{aligned}$$

We therefore get

$$\mathbb{E}_{\mathbf{x}}[\phi(A; \mathcal{C} \cup \{\mathbf{x}\})] \leq (6 + 4\eta) \cdot \phi_{\text{opt}}(A) + 2 \cdot \phi_{\text{opt}}^N(A) , \tag{35}$$

as claimed. ■

Again, we recover Lemma 3.3 in (Arthur & Vassilvitskii, 2007) when p_a is a Dirac and the probe function $\wp = \text{Id}$. The rest of the proof of Theorem 2 consists of the same steps as Theorem 3.1 in (Arthur & Vassilvitskii, 2007), after having remarked that $\phi_{\text{opt}}^N(A)$ can be simplified:

$$\begin{aligned}
\phi_{\text{opt}}^N(A) &= \sum_{a \in A} \int_{\Omega_{a_0}} \|\mathbf{x} - \mathbf{c}(A)\|_2^2 dp_a(\mathbf{x}) \\
&= \sum_{a \in A} \int_{\Omega_{a_0}} \|\mathbf{x}\|_2^2 dp_a(\mathbf{x}) - 2\langle \mathbf{c}(A), \boldsymbol{\mu}_a \rangle + \|\mathbf{c}(A)\|_2^2 \\
&= \sum_{a \in A} \int_{\Omega_{a_0}} \|\mathbf{x} - \boldsymbol{\mu}_a\|_2^2 dp_a(\mathbf{x}) + \|\boldsymbol{\mu}_a\|_2^2 - 2\langle \mathbf{c}(A), \mathbf{a} \rangle + \|\mathbf{c}(A)\|_2^2 \\
&= \sum_{a \in A} \{ \text{tr}(\Sigma_a) + \|\boldsymbol{\mu}_a - \mathbf{c}(A)\|_2^2 \} \\
&= \phi_{\text{bias}}(A) + \phi_{\text{var}}(A) .
\end{aligned} \tag{36}$$

Proof of Lemma 3

The proof is a simple application of the Fréchet-Cramér-Rao-Darmois bound. Consider the simple case $k = 1$ and a spherical Gaussian noise for p with a single point in \mathcal{A} . Renormalize both sides of (7) by $m \doteq |\mathcal{A}|$ so that $(1/m) \sum_{a \in \mathcal{A}} \text{tr}(\Sigma_a) = \text{tr}(\Sigma)$. One sees that the left hand side of ineq. (7) is just an estimator of the variance of p_a , which, by Fréchet-Darmois-Cramér-Rao bound, has to be at least the inverse of the Fisher information, that is in this case, the trace of the covariance matrix, *i.e.* $\text{tr}(\Sigma)$.

Comments on Table 1

(Wang et al., 2015) are concerned with approximating subspace clustering, and so they are using a very different potential function, which is, between two subspaces \mathcal{S} and \mathcal{S}' , $d(\mathcal{S}, \mathcal{S}') = \|\mathbf{U}\mathbf{U}^\top - \mathbf{U}'\mathbf{U}'^\top\|_F$, where \mathbf{U} (resp. \mathbf{U}') is an *orthonormal* basis for \mathcal{S} (resp. \mathcal{S}'). To obtain an idea of the approximation on the k -means clustering problem that their technique yields, we compute ϕ in the projected space, using the fact that, because of the triangle inequality and the fact that projections are linear and do not increase norms,

$$\begin{aligned}
\|\text{proj}_{\mathbf{U}}(\mathbf{a}) - \text{proj}_{\mathbf{U}'}(\mathbf{a}')\|_2 &= \|(\text{proj}_{\mathbf{U}}(\mathbf{a}) - \text{proj}_{\mathbf{U}}(\mathbf{a}')) + (\text{proj}_{\mathbf{U}}(\mathbf{a}') - \text{proj}_{\mathbf{U}'}(\mathbf{a}'))\|_2 \tag{37} \\
&\leq \|\text{proj}_{\mathbf{U}}(\mathbf{a}) - \text{proj}_{\mathbf{U}}(\mathbf{a}')\|_2 + \|\text{proj}_{\mathbf{U}}(\mathbf{a}') - \text{proj}_{\mathbf{U}'}(\mathbf{a}')\|_2 \tag{38} \\
&\leq \|\text{proj}_{\mathbf{U}}(\mathbf{a}) - \text{proj}_{\mathbf{U}}(\mathbf{a}')\|_2 + 2\|\mathbf{a}'\|_2 . \tag{39}
\end{aligned}$$

To account for the approximation in the inequalities, we then discard the rightmost term, replacing therefore $\|\text{proj}_{\mathbf{U}}(\mathbf{a}) - \text{proj}_{\mathbf{U}'}(\mathbf{a}')\|_2$ by $\|\text{proj}_{\mathbf{U}}(\mathbf{a}) - \text{proj}_{\mathbf{U}}(\mathbf{a}')\|_2$, which amounts, in the approximation bounds, to remove the dependence in the dimension. At this price, and using the trick to transfer the wasserstein distance between centers to L_2^2 potential between points to cluster centers, we obtain the approximation bound in (β) of Table 1. While it has to be used with care, its main interest is in showing that the price to pay because of the noise component is in fact not decreasing in m .

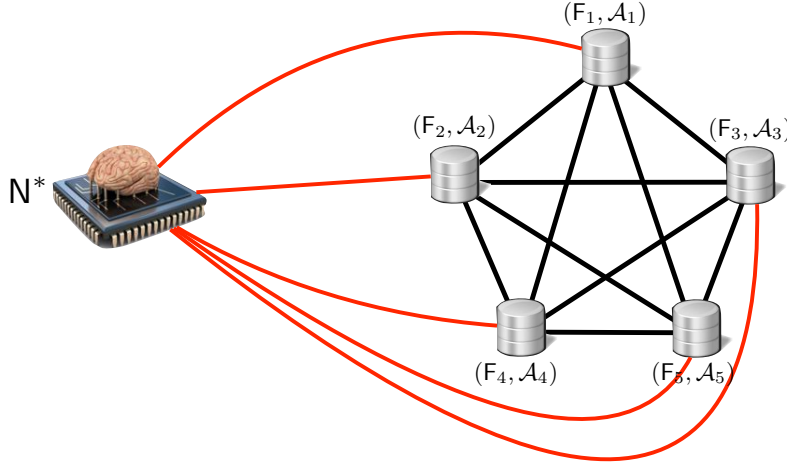


Figure 4: Message passing between peers / nodes in the Dk -means++/PD k -means++ framework. Black edges and red arcs denote message passing between peers / nodes. On each black edge circulates at most k data points; on each red arcs circulates k total potentials.

Proofs of Theorems 4, 5 and 6

The proof of these Theorems uses a *reduction* from k -variates++ to the corresponding algorithms, meaning that there exists particular probe functions and densities for which the set of centers delivered by k -variates++ is the same as the one delivered by the corresponding algorithms.

Definition 15 Let \mathcal{H} (parameters omitted) be any hard membership k -clustering algorithm. We say that k -variates++ **reduces** to \mathcal{H} iff there exists data, densities and probe functions depending on the instance of \mathcal{H} such that, in expectation over the internal randomisation of \mathcal{H} , the set of centers delivered by \mathcal{H} are the same as the ones delivered by k -variates++. We note it

$$k\text{-variates++} \succeq \mathcal{H} . \quad (40)$$

Hence, whenever k -variates++ $\succeq \mathcal{H}$, Theorem 2 immediately gives a guarantee for the approximation of the global optimum in expectation for \mathcal{H} , but this requires the translation of the parameters involved in Φ in ineq. (7) to involve only parameters from \mathcal{H} . In all our examples, this translation poses no problem at all.

Proof of Theorem 4

Figure 4 presents the architecture of message passing in the Dk -means++/PD k -means++ framework. We first focus on the protected scheme, Dk -means++. We reduce k -variates++ to Algorithm

1 using identity probe functions: $\wp_t = \text{Id}, \forall t$. The trick in reduction relies on the densities. We let p_{μ_a, θ_a} be uniform over the subset \mathcal{A}_i to which \mathbf{a} belongs. Thus, the support of densities is discrete, and \mathcal{C} is a subset of \mathcal{A} ; furthermore, the probability $q_t(\mathbf{a})$ that $\mathbf{a} \in \mathcal{A}_i$ is chosen at iteration t in k -variates++ actually simplifies to a convenient expression:

$$q_t(\mathbf{a}) = q_{ti}^D \cdot u_i, \quad (41)$$

where we recall that

$$q_{ti}^D \doteq \begin{cases} D_t(\mathcal{A}_i) \cdot (\sum_j D_t(\mathcal{A}_j))^{-1} & \text{if } t > 1 \\ (1/n) & \text{otherwise} \end{cases}. \quad (42)$$

Hence, picking \mathbf{a} can be equivalently done by first picking \mathcal{A}_i using q_t^D , and then, given the i chosen, sampling uniformly at random \mathbf{a} in \mathcal{A}_i , which is what *Forgy* nodes do. We therefore get the equivalence between Algorithm 1 and k -variates++ as instantiated.

Lemma 16 *With data, densities and probes defined as before, k -variates++ \succeq Dk -means++.*

To get the approximability ratio of Dk -means++, we translate the parameters of Φ in ineq. (7). First, since $(a + b)^2 \leq 2a^2 + 2b^2$,

$$\begin{aligned} \phi_{\text{bias}} &\doteq \sum_{\mathbf{a} \in \mathcal{A}} \|\mu_{\mathbf{a}} - \mathbf{c}_{\text{opt}}(\mathbf{a})\|_2^2 \\ &= \sum_{i \in [n]} \sum_{\mathbf{a} \in \mathcal{A}_i} \|\mathbf{c}(\mathcal{A}_i) - \mathbf{c}_{\text{opt}}(\mathbf{a})\|_2^2 \end{aligned} \quad (43)$$

$$\begin{aligned} &= \sum_{i \in [n]} \sum_{\mathbf{a} \in \mathcal{A}_i} \|\mathbf{c}(\mathcal{A}_i) - \mathbf{a} + \mathbf{a} - \mathbf{c}_{\text{opt}}(\mathbf{a})\|_2^2 \\ &\leq 2 \sum_{i \in [n]} \sum_{\mathbf{a} \in \mathcal{A}_i} \|\mathbf{c}(\mathcal{A}_i) - \mathbf{a}\|_2^2 + 2 \sum_{\mathbf{a} \in \mathcal{A}} \|\mathbf{a} - \mathbf{c}_{\text{opt}}(\mathbf{a})\|_2^2 \\ &= 2\phi_s^F + 2\phi_{\text{opt}}. \end{aligned} \quad (44)$$

Furthermore,

$$\begin{aligned} \phi_{\text{var}} &\doteq \sum_{\mathbf{a} \in \mathcal{A}} \text{tr}(\Sigma_{\mathbf{a}}) \\ &= \sum_{\mathbf{a} \in \mathcal{A}} \int_{\Omega_{\mathbf{a}}} \|\mathbf{x} - \mu_{\mathbf{a}}\|_2^2 d p_{\mathbf{a}}(\mathbf{x}) \\ &= \sum_{i \in [n]} \sum_{\mathbf{a} \in \mathcal{A}_i} \sum_{\mathbf{a}' \in \mathcal{A}_i} \frac{1}{m_i} \cdot \|\mathbf{a}' - \mathbf{c}(\mathcal{A}_i)\|_2^2 \\ &= \sum_{i \in [n]} \sum_{\mathbf{a} \in \mathcal{A}_i} \|\mathbf{a} - \mathbf{c}(\mathcal{A}_i)\|_2^2 = \phi_s^F. \end{aligned} \quad (45)$$

There remains to plug ineq. (44) and eq. (45) in Theorem 2, along with $\eta = 0$ (since $\wp = \text{Id}$), to get $\mathbb{E}[\phi(\mathcal{A}; \mathcal{C})] \leq (2 + \log k) \cdot (10\phi_{\text{opt}} + 6\phi_s)$, as in Theorem 4.

The private version, PDk -means++, follows immediately by leaving ϕ_{var} in Φ instead of carrying eq. (45). This ends the proof of Theorem 4.

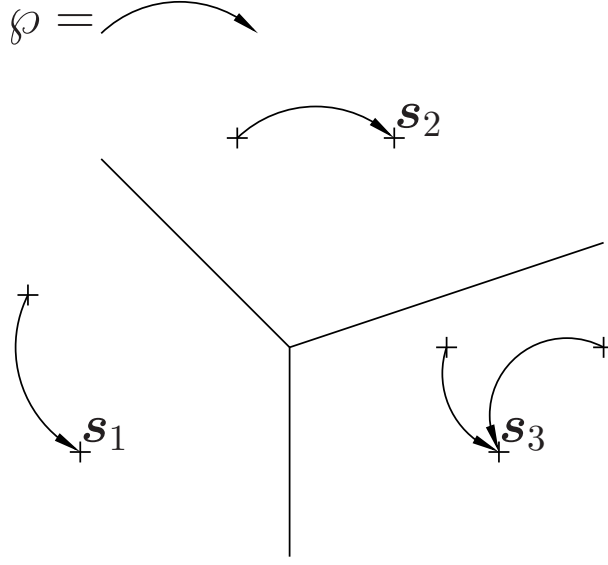


Figure 5: Computation of the probe function \wp for the reduction from k -variates++ to sk -means++. Segments display parts of the Voronoi diagram of \mathcal{S} .

Proof of Theorem 5

The proof proceeds in the same way as for Theorem 4. The probe function (the same for every iteration, $\wp_t = \wp, \forall t$) is already defined in the statement of Theorem 5, from the definition of synopses. The distributions p_{μ_a, θ_a} are Diracs anchored at the *probe* (synopses) locations. The centers chosen in k -variates++ are thus synopses, and it is not hard to check that the probability to pick a synopsis s_j at iteration t factors in the same way as in the definition of q_t^S in eq. (9). We therefore get the equivalence between Algorithm 2 and k -variates++ as instantiated.

Lemma 17 *With data, densities and probes defined as before, k -variates++ $\succeq sk$ -means++.*

The proof of the approximation property of sk -means++ then follows from the fact that $\phi_{\text{var}} = 0$ (Diracs) and

$$\begin{aligned}
 \phi_{\text{bias}} &\doteq \sum_{\mathbf{a} \in \mathcal{A}} \|\mu_{\mathbf{a}} - \mathbf{c}_{\text{opt}}(\mathbf{a})\|_2^2 \\
 &= \sum_{\mathbf{a} \in \mathcal{A}} \|\wp(\mathbf{a}) - \mathbf{c}_{\text{opt}}(\mathbf{a})\|_2^2 \\
 &= \sum_{\mathbf{a} \in \mathcal{A}} \|\wp(\mathbf{a}) - \mathbf{a} + \mathbf{a} - \mathbf{c}_{\text{opt}}(\mathbf{a})\|_2^2 \\
 &\leq 2 \sum_{\mathbf{a} \in \mathcal{A}} \|\wp(\mathbf{a}) - \mathbf{a}\|_2^2 + 2 \sum_{\mathbf{a} \in \mathcal{A}} \|\mathbf{a} - \mathbf{c}_{\text{opt}}(\mathbf{a})\|_2^2 \\
 &= 2 \sum_{\mathbf{a} \in \mathcal{S}} \|\wp(\mathbf{a}) - \mathbf{a}\|_2^2 + 2 \sum_{\mathbf{a} \in \mathcal{S}} \|\mathbf{a} - \mathbf{c}_{\text{opt}}(\mathbf{a})\|_2^2 = 2\phi_s^\wp + 2\phi_{\text{opt}} \tag{46}
 \end{aligned}$$

(using again $(a + b)^2 \leq 2a^2 + 2b^2$). Using Theorem 2, this brings the statement of the Theorem.

Setting	Algorithm	Probe functions \wp_t	Densities $p(\mu, \theta)$
Batch	k -means++ (Arthur & Vassilvitskii, 2007)	Identity	Diracs
Distributed	D k -means++	Identity	Uniform on data subsets
Distributed	PD k -means++	Identity	Non uniform, compact support
Streaming	S k -means++	synopses	Diracs
On-line	OL k -means++	point (batch not hit) / closest center (batch hit)	Diracs

Table 3: Synthesis of the parameters for the reductions from k -variates++. We indicate k -means++ as the batch clustering solution (Arthur & Vassilvitskii, 2007).

Figure 5 shows that the "quality" of the probe function (spread ϕ_s^\wp , stretching factor η) stem from the quality of the Voronoi diagram induced by the synopses in \mathcal{S} .

Proof of Theorem 6

The proof proceeds in the same way as for Theorem 4. The the reduction from k -variates++ to OL k -means++ relies on two things: first, the uniform choice of the first center in k -means++ can be replaced by picking the center uniformly in *any* subset of the data: it does not change the expected approximation properties of the algorithm (this comes from Lemma 3.4 in (Arthur & Vassilvitskii, 2007)); therefore, the choice $q_1 \doteq u_m$ in k -variates++ can be replaced with $q_1 \doteq u_1$ (uniform with support \mathcal{A}_1). Second, a particular probe function needs to be devised, sketched in Figure 6. Basically, all probe functions of a minibatch are the same: each point in the minibatch is probed to itself, while points occurring outside the minibatch are probed to their closest center. The reduction proceeds in the following steps: we first let \mathcal{A} be the complete set of points in the stream \mathcal{S} . Then, we let \mathcal{A}_j denote the set of points of minibatch \mathcal{S}_j . Remark that minibatch \mathcal{A}_j occurs in the stream before $\mathcal{A}_{j'}$ for $j < j'$, and minibatches induce a partition of \mathcal{A} . Let $j(t)$ denote the batch related to iteration t in k -variates++. We define the following probe function $\wp_t(\mathbf{a})$ in k -variates++, letting \mathcal{A}_j the minibatch to which \mathbf{a} belongs (we do not necessarily have $j = j(t)$):

- if $j = j(t)$, then $\wp_t(\mathbf{a}) \doteq \mathbf{a}$;
- else $\wp_t(\mathbf{a}) \doteq \arg \min_{\mathbf{c} \in \mathcal{C}} \|\mathbf{a} - \mathbf{c}\|_2^2$ (remark that $|\mathcal{C}| \geq 1$ in this case).

Finally, densities $p(\mu, \theta)$ are Diracs anchored at selected points, like in k -means++. We get the equivalence between Algorithm 3 and k -variates++ as instantiated.

Lemma 18 *With data, densities and probes defined as before, k -variates++ \succeq OL k -means++.*

The proof is immediate, since each minibatch is hit by a center exactly once in OL k -means++, and when one subset \mathcal{A}_j is hit by a center, then the probe function makes that *no* other center can be sampled again from \mathcal{A}_j (all contributions to the density q_t are then zero in \mathcal{A}_j). We now finish the proof of Theorem 6 by showing the same approximability ratio for k -variates++ as reduced. Because *optimal* clusters are ς -wide with respect to stream \mathcal{S} , we have

$$\frac{1}{|\mathcal{A}|} \cdot \sum_{\mathbf{a}, \mathbf{a}' \in \mathcal{A}} \|\mathbf{a} - \mathbf{a}'\|_2^2 \geq \varsigma \cdot R .$$

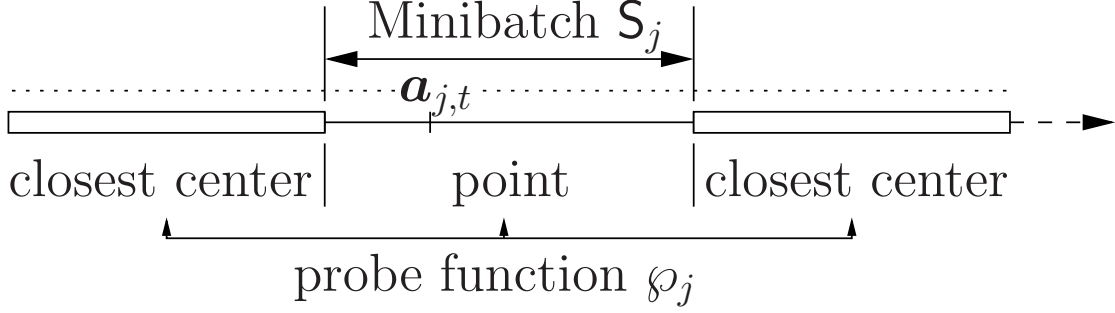


Figure 6: Computation of the probe function \wp_t for the reduction from OL k -means++ to k -variates++, depending on each minibatch stream S_j .

Recall that $\mathbf{c}(A) \doteq (1/|A|) \cdot \sum_{\mathbf{a} \in A} \mathbf{a}$. For any $\mathbf{a}_0 \in A$, it holds that:

$$\frac{1}{|A| - 1} \cdot \sum_{\mathbf{a} \in A} \|\mathbf{a} - \mathbf{a}_0\|_2^2 \geq \frac{1}{|A| - 1} \cdot \sum_{\mathbf{a} \in A} \|\mathbf{a} - \mathbf{c}(A)\|_2^2 \quad (47)$$

$$= \frac{1}{|A| - 1} \cdot \left(\frac{1}{2|A|} \cdot \sum_{\mathbf{a}, \mathbf{a}' \in A} \|\mathbf{a} - \mathbf{a}'\|_2^2 \right) \quad (48)$$

$$= \frac{1}{4} \cdot \frac{2}{|A|(|A| - 1)} \cdot \sum_{\mathbf{a}, \mathbf{a}' \in A} \|\mathbf{a} - \mathbf{a}'\|_2^2 \geq \frac{\varsigma}{4} \cdot R. \quad (49)$$

Ineq. (47) holds because $\mathbf{c}(A)$ is the population minimizer for optimal cluster A (see *e.g.*, (Arthur & Vassilvitskii, 2007), Lemma 2.1). Since probes are points of A ,

$$\begin{aligned} \phi(\wp_j(A); \{\wp_j(\mathbf{a}_0)\}) &\leq |A| \cdot R \\ &\leq \frac{4|A|}{\varsigma(|A| - 1)} \cdot \sum_{\mathbf{a} \in A} \|\mathbf{a} - \mathbf{a}_0\|_2^2. \end{aligned} \quad (50)$$

On the other hand, we have:

$$\phi(\wp_t(A); \mathcal{C}) = \sum_{\mathbf{a} \in A \cap S_j} \|\mathbf{a} - \mathbf{c}(\mathbf{a})\|_2^2, \quad (51)$$

but since minibatches are ς accurate, $\sum_{\mathbf{a} \in A \cap S_j} \|\mathbf{a} - \mathbf{c}(\mathbf{a})\|_2^2 \geq \varsigma \cdot \sum_{\mathbf{a} \in A} \|\mathbf{a} - \mathbf{c}(\mathbf{a})\|_2^2$. Therefore, for any $\mathbf{a}_0 \in A$,

$$\begin{aligned} \frac{\phi(\wp_t(A); \mathcal{C})}{\phi(\wp_t(A); \{\wp_t(\mathbf{a}_0)\})} &\geq \left(\frac{\varsigma^2(|A| - 1)}{4|A|} \right) \cdot \frac{\sum_{\mathbf{a} \in A} \|\mathbf{a} - \mathbf{c}(\mathbf{a})\|_2^2}{\sum_{\mathbf{a} \in A} \|\mathbf{a} - \mathbf{a}_0\|_2^2} \\ &= \left(\frac{\varsigma^2(|A| - 1)}{4|A|} \right) \cdot \frac{\phi(A; \mathcal{C})}{\phi(A; \{\mathbf{a}_0\})}. \end{aligned} \quad (52)$$

In other words, probe functions are η -stretching, for any η satisfying:

$$\eta \geq \frac{4|A|}{\varsigma^2(|A| - 1)} - 1, \quad (53)$$

and they are therefore η -stretching for $\eta = 8/\zeta^2 - 1$. There remains to check that, because of the densities chosen,

$$\phi_{\text{bias}} = \phi_{\text{opt}} , \quad (54)$$

$$\phi_{\text{var}} = 0 . \quad (55)$$

This ends the proof of Theorem 6.

Proof of Theorem 9

To simplify notations in the proof, we let $p_a(\mathbf{x})$ denote the value of density $p(\mu_a, \theta_a)$ on some $\mathbf{x} \in \Omega$. Let us denote $\text{Seq}(n : k)$ the number of sequences of integers in set $\{1, 2, \dots, n\}$ having exactly k elements, whose cardinal is $|\text{Seq}(n : k)| = n!/(n-k)!$. For any sequence $I \in \text{Seq}(n : k)$, we let I_i denote its i^{th} element. For any set $\mathcal{C} \doteq \{\mathbf{c}_1, \mathbf{c}_2, \dots, \mathbf{c}_k\}$ returned by Algorithm *k-variates++* with input instance set $\mathcal{A} \doteq \{\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_n\} \subset \Omega$, the density of \mathcal{C} given \mathcal{A} is:

$$\mathbb{P}[\mathcal{C}|\mathcal{A}] = \sum_{\sigma \in S_k} \sum_{I \in \text{Seq}(n:k)} p(\sigma, I, \mathcal{C}|\mathcal{A}) , \quad (56)$$

where S_k denotes the symmetric group on k elements, and the following shorthand is used:

$$p(\sigma, I, \mathcal{C}|\mathcal{A}) \doteq \prod_{i=1}^k q_i(\mathbf{a}_{I_i}) p_{\mathbf{a}_{I_i}}(\mathbf{c}_{\sigma(i)}) , \quad (57)$$

where q_i is computed using eq. (1) and taking into account the modification due to the choice of each I_j for $j < i$ in the sequence I .

In the following, we let \mathcal{A} and \mathcal{A}' denote two sets of points that differ from one a (they have the same size), say $\mathbf{a}_n \in \mathcal{A}$ and $\mathbf{a}'_n \in \mathcal{A}'$, $\mathbf{a}_n \neq \mathbf{a}'_n$. We analyze:

$$\frac{\mathbb{P}[\mathcal{C}|\mathcal{A}']}{\mathbb{P}[\mathcal{C}|\mathcal{A}]} = \frac{\sum_{\sigma \in S_k} \sum_{I \in \text{Seq}(n:k)} p(\sigma, I, \mathcal{C}|\mathcal{A}')}{\sum_{\sigma \in S_k} \sum_{I \in \text{Seq}(n:k)} p(\sigma, I, \mathcal{C}|\mathcal{A})} . \quad (58)$$

Using the definition of $q(\cdot)$, we refine $p(\sigma, I, \mathcal{C}|\mathcal{A})$ as

$$p(\sigma, I, \mathcal{C}|\mathcal{A}) = \frac{N(I)}{\prod_{i=1}^k M(I^i|\mathcal{A})} \cdot \prod_{i=1}^k p_{\mathbf{a}_{I_i}}(\mathbf{c}_{\sigma(i)}) , \quad (59)$$

where

$$N(I) \doteq \prod_{i=2}^j \|\mathbf{a}_{I_i} - \text{NN}_{I^i}(\mathbf{a}_{I_i})\|_2^2 , \quad (60)$$

$$M(I^i|\mathcal{A}) \doteq \begin{cases} n & \text{if } i = 1 \\ \sum_{j=1}^n \|\mathbf{a}_j - \text{NN}_{I^i}(\mathbf{a}_j)\|_2^2 & \text{otherwise} \end{cases} , \quad (61)$$

and I^i is the prefix sequence I_1, I_2, \dots, I_{i-1} , and $\text{NN}_{I^i}(a) \doteq \arg \min_{j \leq i-1} \|a - \mathbf{a}_{I_j}\|_2$ is the nearest neighbor of a in the prefix sequence. Notice that there is a factor $1/m$ for $q(\cdot)$ at the first iteration that we omit in $N(I)$ since it disappears in the ratio in eq. (58).

We analyze separately each element in (59), starting with $N(I)$. We define the *swapping* operation $s_\ell(I)$ that returns the sequence in which \mathbf{a}_{I_ℓ} and $\mathbf{a}_{I_{\ell+1}}$ are permuted, for $1 \leq \ell \leq k-1$. This incurs non-trivial modifications in $N(s_\ell(I))$ compared to $N(I)$, since the nearest neighbors of \mathbf{a}_{I_ℓ} and $\mathbf{a}_{I_{\ell+1}}$ may change in the permutation:

$$\begin{aligned}
N(s_\ell(I)) &= \prod_{i=2}^{\ell-1} \|\mathbf{a}_{I_i} - \text{NN}_{I^i}(\mathbf{a}_{I_i})\|_2^2 \\
&\quad \cdot \underbrace{\|\mathbf{a}_{I_{\ell+1}} - \text{NN}_{I^\ell}(\mathbf{a}_{I_{\ell+1}})\|_2^2 \cdot \|\mathbf{a}_{I_\ell} - \text{NN}_{I^\ell \cup \{I_{\ell+1}\}}(\mathbf{a}_{I_\ell})\|_2^2}_{\neq \|\mathbf{a}_{I_\ell} - \text{NN}_{I^\ell}(\mathbf{a}_{I_\ell})\|_2^2 \cdot \|\mathbf{a}_{I_{\ell+1}} - \text{NN}_{I^{\ell+1}}(\mathbf{a}_{I_{\ell+1}})\|_2^2} \\
&\quad \cdot \prod_{i=\ell+2}^k \|\mathbf{a}_{I_i} - \text{NN}_{I^i}(\mathbf{a}_{I_i})\|_2^2
\end{aligned} \tag{62}$$

($I \cup \{j\}$ indicates that element j is put at the end of the sequence). We want to quantify the maximal increase in $N(s_\ell(I))$ compared to $N(I)$. The following Lemma shows that the maximal increase ratio is actually a constant, and thus does not depend on the data.

Lemma 19 *The following holds true:*

$$N(s_1(I)) = N(I) , \tag{63}$$

$$N(s_\ell(I)) \leq (1 + \eta)^2 N(I) , \forall 2 \leq \ell \leq k-1 . \tag{64}$$

Here, $0 \leq \eta \leq 3$ is a constant.

The proof stems directly from the following Lemma.

Lemma 20 *For any non-empty $\mathcal{N} \subseteq \mathcal{A}$ and $x \in \Omega$, let $\text{NN}_{\mathcal{N}}(x)$ denote the nearest neighbor of x in \mathcal{N} . There exists a constant $0 \leq \eta \leq 3$ such that for any $\mathbf{a}_i, \mathbf{a}_j \in \mathcal{A}$ and any nonempty subset $\mathcal{N} \subseteq \mathcal{A} \setminus \{\mathbf{a}_i, \mathbf{a}_j\}$,*

$$\frac{\|\mathbf{a}_i - \text{NN}_{\mathcal{N}}(\mathbf{a}_i)\|_2}{\|\mathbf{a}_i - \text{NN}_{\mathcal{N} \cup \{\mathbf{a}_j\}}(\mathbf{a}_i)\|_2} \leq (1 + \eta) \cdot \frac{\|\mathbf{a}_j - \text{NN}_{\mathcal{N}}(\mathbf{a}_j)\|_2}{\|\mathbf{a}_j - \text{NN}_{\mathcal{N} \cup \{\mathbf{a}_i\}}(\mathbf{a}_j)\|_2} . \tag{65}$$

Proof Since $\|\mathbf{a}_j - \text{NN}_{\mathcal{N} \cup \{\mathbf{a}_i\}}(\mathbf{a}_j)\|_2 \leq \|\mathbf{a}_j - \text{NN}_{\mathcal{N}}(\mathbf{a}_j)\|_2$, the proof is true for $\eta = 0$ when $\text{NN}_{\mathcal{N}}(\mathbf{a}_i) = \text{NN}_{\mathcal{N} \cup \{\mathbf{a}_j\}}(\mathbf{a}_i)$. So suppose that $\text{NN}_{\mathcal{N}}(\mathbf{a}_i) \neq \text{NN}_{\mathcal{N} \cup \{\mathbf{a}_j\}}(\mathbf{a}_i)$, implying $\text{NN}_{\mathcal{N} \cup \{\mathbf{a}_j\}}(\mathbf{a}_i) = \mathbf{a}_j$. We distinguish two cases.

Case 1/2, if $\text{NN}_{\mathcal{N} \cup \{\mathbf{a}_i\}}(\mathbf{a}_j) = \mathbf{a}_i$, then we are reduced to showing that $\|\mathbf{a}_i - \text{NN}_{\mathcal{N}}(\mathbf{a}_i)\|_2 \leq (1 + \eta) \|\mathbf{a}_j - \text{NN}_{\mathcal{N}}(\mathbf{a}_j)\|_2$ under the conditions (C) that $\mathcal{N} \cap B(\mathbf{a}_i, \|\mathbf{a}_i - \mathbf{a}_j\|_2) = \emptyset$ and $\mathcal{N} \cap B(\mathbf{a}_j, \|\mathbf{a}_i - \mathbf{a}_j\|_2) = \emptyset$. Here, $B(\mathbf{a}, r)$ denotes the open ball of center \mathbf{a} and radius R . The triangle inequality and conditions (C) bring

$$\begin{aligned}
\|\mathbf{a}_i - \text{NN}_{\mathcal{N}}(\mathbf{a}_i)\|_2 &\leq \|\mathbf{a}_i - \mathbf{a}_j\|_2 + \|\mathbf{a}_j - \text{NN}_{\mathcal{N}}(\mathbf{a}_i)\|_2 \\
&\leq \|\mathbf{a}_j - \text{NN}_{\mathcal{N}}(\mathbf{a}_j)\|_2 + \|\mathbf{a}_j - \text{NN}_{\mathcal{N}}(\mathbf{a}_i)\|_2 .
\end{aligned} \tag{66}$$

If $\text{NN}_{\mathcal{N}}(\mathbf{a}_i) = \text{NN}_{\mathcal{N}}(\mathbf{a}_j)$ then the inequality holds for $\eta = 1$. Otherwise, suppose that $\|\mathbf{a}_j - \text{NN}_{\mathcal{N}}(\mathbf{a}_i)\|_2 > 3 \|\mathbf{a}_j - \text{NN}_{\mathcal{N}}(\mathbf{a}_j)\|_2$. The triangle inequality yields again $\|\mathbf{a}_j - \text{NN}_{\mathcal{N}}(\mathbf{a}_i)\|_2 \leq \|\mathbf{a}_j - \mathbf{a}_i\|_2 + \|\mathbf{a}_i - \text{NN}_{\mathcal{N}}(\mathbf{a}_i)\|_2$, and so we have the inequality:

$$3 \|\mathbf{a}_j - \text{NN}_{\mathcal{N}}(\mathbf{a}_j)\|_2 < \|\mathbf{a}_j - \mathbf{a}_i\|_2 + \|\mathbf{a}_i - \text{NN}_{\mathcal{N}}(\mathbf{a}_i)\|_2 , \tag{67}$$

and since (C) holds, $\|\mathbf{a}_j - \text{NN}_{\mathcal{N}}(\mathbf{a}_j)\|_2 \geq \|\mathbf{a}_j - \mathbf{a}_i\|_2$ which implies

$$\|\mathbf{a}_j - \text{NN}_{\mathcal{N}}(\mathbf{a}_j)\|_2 < \frac{1}{2} \cdot \|\mathbf{a}_i - \text{NN}_{\mathcal{N}}(\mathbf{a}_i)\|_2 . \quad (68)$$

On the other hand, the triangle inequality brings again

$$\begin{aligned} \|\mathbf{a}_i - \text{NN}_{\mathcal{N}}(\mathbf{a}_j)\|_2 &\leq \|\mathbf{a}_i - \mathbf{a}_j\|_2 + \|\mathbf{a}_j - \text{NN}_{\mathcal{N}}(\mathbf{a}_j)\|_2 \\ &\leq 2 \cdot \|\mathbf{a}_j - \text{NN}_{\mathcal{N}}(\mathbf{a}_j)\|_2 \end{aligned} \quad (69)$$

$$< 2 \cdot \frac{1}{2} \cdot \|\mathbf{a}_i - \text{NN}_{\mathcal{N}}(\mathbf{a}_i)\|_2 = \|\mathbf{a}_i - \text{NN}_{\mathcal{N}}(\mathbf{a}_i)\|_2 , \quad (70)$$

a contradiction since $\|\mathbf{a}_i - \text{NN}_{\mathcal{N}}(\mathbf{a}_i)\|_2 \leq \|\mathbf{a}_i - \mathbf{a}_l\|_2, \forall \mathbf{a}_l \in \mathcal{N}$ by definition. Ineq. (69) uses (C) and ineq. (70) uses ineq. (68). Hence, if $\text{NN}_{\mathcal{N}}(\mathbf{a}_i) \neq \text{NN}_{\mathcal{N}}(\mathbf{a}_j)$ then since $\|\mathbf{a}_j - \text{NN}_{\mathcal{N}}(\mathbf{a}_i)\|_2 \leq 3\|\mathbf{a}_j - \text{NN}_{\mathcal{N}}(\mathbf{a}_j)\|_2$, ineq. (66) brings $\|\mathbf{a}_i - \text{NN}_{\mathcal{N}}(\mathbf{a}_i)\|_2 \leq 4 \cdot \|\mathbf{a}_j - \text{NN}_{\mathcal{N}}(\mathbf{a}_j)\|_2$, and the inequality holds for $\eta = 3$.

Case 2/2, if $\text{NN}_{\mathcal{N} \cup \{\mathbf{a}_i\}}(\mathbf{a}_j) \neq \mathbf{a}_i$, then it implies $\text{NN}_{\mathcal{N} \cup \{\mathbf{a}_i\}}(\mathbf{a}_j) = \text{NN}_{\mathcal{N}}(\mathbf{a}_j)$ and so

$$\exists \mathbf{a}_* \in \mathcal{N} : \|\mathbf{a}_j - \mathbf{a}_*\|_2 \leq \|\mathbf{a}_j - \mathbf{a}_i\|_2 . \quad (71)$$

Ineq. (65) reduces to proving

$$\|\mathbf{a}_i - \text{NN}_{\mathcal{N}}(\mathbf{a}_i)\|_2 \leq (1 + \eta) \cdot \|\mathbf{a}_i - \mathbf{a}_j\|_2 , \quad (72)$$

but $\|\mathbf{a}_i - \mathbf{a}_*\|_2 \leq \|\mathbf{a}_i - \mathbf{a}_j\|_2 + \|\mathbf{a}_j - \mathbf{a}_*\|_2 \leq 2\|\mathbf{a}_i - \mathbf{a}_j\|_2$, and since $\mathbf{a}_* \in \mathcal{N}$, $\|\mathbf{a}_i - \text{NN}_{\mathcal{N}}(\mathbf{a}_i)\|_2 \leq \|\mathbf{a}_i - \mathbf{a}_*\|_2 \leq 2\|\mathbf{a}_i - \mathbf{a}_j\|_2$, and (72) is proved for $\eta = 1$. This achieves the proof of Lemma 20. ■

Let I be any sequence not containing the index of \mathbf{a}'_n , and let $I(i)$ denote the sequence in which we replace \mathbf{a}_{I_i} by the index of \mathbf{a}'_n . The sequence of swaps

$$I(k) = (s_{k-1} \circ \dots \circ s_{i+1} \circ s_i)(I(i)) \quad (73)$$

produces a sequence $I(k)$ in which all elements different from \mathbf{a}'_n are in the same relative order as they are in I with respect to each other, and \mathbf{a}'_n is pushed to the end of the sequence in k^{th} rank. We also have

$$N(I(i)) \leq (1 + \eta)^{2(k-i)} N(I(k)) . \quad (74)$$

All the properties we need on $N(\cdot)$ are now established. We turn to the analysis of $M(I^i|\mathcal{A})$.

Lemma 21 *For any $\delta_s > 0$ such that \mathcal{A} is δ_s -monotonic, the following holds. For any $\mathcal{N} \subseteq \mathcal{A}$ with $|\mathcal{N}| \in \{1, 2, \dots, k-1\}$, $\forall \mathbf{x}, \mathbf{x}' \in \Omega$, we have:*

$$\sum_{\mathbf{a} \in \mathcal{A}} \|\mathbf{a} - \text{NN}_{\mathcal{N} \cup \{\mathbf{x}\}}(\mathbf{a})\|_2^2 \leq (1 + \delta_s) \cdot \sum_{\mathbf{a} \in \mathcal{A}} \|\mathbf{a} - \text{NN}_{\mathcal{N} \cup \{\mathbf{x}'\}}(\mathbf{a})\|_2^2 . \quad (75)$$

Proof Since adding a point to \mathcal{N} cannot increase the potential $\sum_{\mathbf{a} \in \mathcal{A}} \|\mathbf{a} - \text{NN}_{\mathcal{N} \cup \{\mathbf{x}\}}(\mathbf{a})\|_2^2$, it comes

$$\sum_{\mathbf{a} \in \mathcal{A}} \|\mathbf{a} - \text{NN}_{\mathcal{N} \cup \{\mathbf{x}\}}(\mathbf{a})\|_2^2 \leq \sum_{\mathbf{a} \in \mathcal{A}} \|\mathbf{a} - \text{NN}_{\mathcal{N}}(\mathbf{a})\|_2^2 , \forall \mathbf{x} \in \Omega . \quad (76)$$

Consider any $\mathbf{x}' \in \Omega$ such that $\sum_{\mathbf{a} \in \mathcal{A}} \|\mathbf{a} - \text{NN}_{\mathcal{N} \cup \{\mathbf{x}'\}}(\mathbf{a})\|_2^2 = \sum_{\mathbf{a} \in \mathcal{A}} \|\mathbf{a} - \text{NN}_{\mathcal{N}}(\mathbf{a})\|_2^2$, i.e., all points of \mathcal{A} are closer to a point in \mathcal{N} than they are from \mathbf{x}' . In this case, we obtain from ineq. (76),

$$\sum_{\mathbf{a} \in \mathcal{A}} \|\mathbf{a} - \text{NN}_{\mathcal{N} \cup \{\mathbf{x}\}}(\mathbf{a})\|_2^2 \leq \sum_{\mathbf{a} \in \mathcal{A}} \|\mathbf{a} - \text{NN}_{\mathcal{N} \cup \{\mathbf{x}'\}}(\mathbf{a})\|_2^2, \quad (77)$$

and since $\delta_s > 0$, the statement of the Lemma holds.

More interesting is the case where $\mathbf{x}' \in \Omega$ is such that $\sum_{\mathbf{a} \in \mathcal{A}} \|\mathbf{a} - \text{NN}_{\mathcal{N} \cup \{\mathbf{x}'\}}(\mathbf{a})\|_2^2 < \sum_{\mathbf{a} \in \mathcal{A}} \|\mathbf{a} - \text{NN}_{\mathcal{N}}(\mathbf{a})\|_2^2$, implying $\mathbf{x}' \notin \mathcal{N}$. In this case, let $A \doteq \{\mathbf{a} \in \mathcal{A} : \text{NN}_{\mathcal{N} \cup \{\mathbf{x}'\}}(\mathbf{a}) = \mathbf{x}'\}$, which is then non-empty. Let us denote for short $\mathbf{c}(A) \doteq (1/|A|) \cdot \sum_{\mathbf{a} \in A} \mathbf{a}$. Since $\mathbf{x}' \notin \mathcal{N}$, $A \cap \mathcal{N} = \emptyset$, and since \mathcal{A} is δ_s -monotonic, then it comes from ineq. (76)

$$\sum_{\mathbf{a} \in \mathcal{A}} \|\mathbf{a} - \text{NN}_{\mathcal{N} \cup \{\mathbf{x}\}}(\mathbf{a})\|_2^2 \leq (1 + \delta_s) \cdot \sum_{\mathbf{a} \in \mathcal{A}} \|\mathbf{a} - \text{NN}_{\mathcal{N} \cup \{\mathbf{c}(A)\}}(\mathbf{a})\|_2^2. \quad (78)$$

We have:

$$\begin{aligned} \sum_{\mathbf{a} \in \mathcal{A}} \|\mathbf{a} - \text{NN}_{\mathcal{N} \cup \{\mathbf{c}(A)\}}(\mathbf{a})\|_2^2 &= \sum_{\mathbf{a} \in \mathcal{A} \setminus A} \|\mathbf{a} - \text{NN}_{\mathcal{N} \cup \{\mathbf{c}(A)\}}(\mathbf{a})\|_2^2 + \sum_{\mathbf{a} \in A} \|\mathbf{a} - \text{NN}_{\mathcal{N} \cup \{\mathbf{c}(A)\}}(\mathbf{a})\|_2^2 \\ &\leq \sum_{\mathbf{a} \in \mathcal{A} \setminus A} \|\mathbf{a} - \text{NN}_{\mathcal{N} \cup \{\mathbf{c}(A)\}}(\mathbf{a})\|_2^2 + \sum_{\mathbf{a} \in A} \|\mathbf{a} - \mathbf{c}(A)\|_2^2 \\ &\leq \sum_{\mathbf{a} \in \mathcal{A} \setminus A} \|\mathbf{a} - \text{NN}_{\mathcal{N} \cup \{\mathbf{c}(A)\}}(\mathbf{a})\|_2^2 + \sum_{\mathbf{a} \in A} \|\mathbf{a} - \mathbf{x}'\|_2^2. \end{aligned} \quad (79)$$

Eq. (79) holds because the arithmetic average is the population minimizer of L_2^2 . Because of the definition of A ,

$$\begin{aligned} \sum_{\mathbf{a} \in \mathcal{A} \setminus A} \|\mathbf{a} - \text{NN}_{\mathcal{N} \cup \{\mathbf{c}(A)\}}(\mathbf{a})\|_2^2 &\leq \sum_{\mathbf{a} \in \mathcal{A} \setminus A} \|\mathbf{a} - \text{NN}_{\mathcal{N}}(\mathbf{a})\|_2^2 \\ &= \sum_{\mathbf{a} \in \mathcal{A} \setminus A} \|\mathbf{a} - \text{NN}_{\mathcal{N} \cup \{\mathbf{x}'\}}(\mathbf{a})\|_2^2, \end{aligned} \quad (80)$$

and, still because of the definition of A ,

$$\sum_{\mathbf{a} \in A} \|\mathbf{a} - \mathbf{x}'\|_2^2 = \sum_{\mathbf{a} \in A} \|\mathbf{a} - \text{NN}_{\mathcal{N} \cup \{\mathbf{x}'\}}(\mathbf{a})\|_2^2, \quad (81)$$

so we get from (80) and (81) $\sum_{\mathbf{a} \in \mathcal{A} \setminus A} \|\mathbf{a} - \text{NN}_{\mathcal{N} \cup \{\mathbf{c}(A)\}}(\mathbf{a})\|_2^2 + \sum_{\mathbf{a} \in A} \|\mathbf{a} - \mathbf{x}'\|_2^2 \leq \sum_{\mathbf{a} \in \mathcal{A}} \|\mathbf{a} - \text{NN}_{\mathcal{N} \cup \{\mathbf{x}'\}}(\mathbf{a})\|_2^2$, and finally from ineq. (79),

$$\sum_{\mathbf{a} \in \mathcal{A}} \|\mathbf{a} - \text{NN}_{\mathcal{N} \cup \{\mathbf{c}(A)\}}(\mathbf{a})\|_2^2 \leq \sum_{\mathbf{a} \in \mathcal{A}} \|\mathbf{a} - \text{NN}_{\mathcal{N} \cup \{\mathbf{x}'\}}(\mathbf{a})\|_2^2, \quad (82)$$

which, using ineq. (78), completes the proof of Lemma 21. ■

Lemma 22 *The following holds true, for any $i > 1$, any $\mathcal{A}' \approx \mathcal{A}$, any $\delta_w, \delta_s > 0$:*

$$\mathcal{A} \text{ is } \delta_w\text{-spread} \Rightarrow (n \notin I^i \Rightarrow M(I^i|\mathcal{A}) \leq (1 + \delta_w) \cdot M(I^i|\mathcal{A}')) , \quad (83)$$

$$\mathcal{A} \text{ is } \delta_s\text{-monotonic} \Rightarrow (n \in I^i \Rightarrow M(I^i|\mathcal{A}) \leq (1 + \delta_s) \cdot M(I^i|\mathcal{A}')) . \quad (84)$$

Proof Suppose first that $n \notin I^i$. In this case, since \mathcal{A} is δ_w -spread,

$$\begin{aligned} M(I^i|\mathcal{A}) &= \sum_{j=1}^n \|\mathbf{a}_j - \text{NN}_{I^i}(\mathbf{a}_j)\|_2^2 \\ &= \sum_{j=1}^{n-1} \|\mathbf{a}_j - \text{NN}_{I^i}(\mathbf{a}_j)\|_2^2 + \|\mathbf{a}_n - \text{NN}_{I^i}(\mathbf{a}_j)\|_2^2 \\ &\leq \sum_{j=1}^{n-1} \|\mathbf{a}_j - \text{NN}_{I^i}(\mathbf{a}_j)\|_2^2 + R^2 \\ &\leq (1 + \delta_w) \cdot \sum_{j=1}^{n-1} \|\mathbf{a}_j - \text{NN}_{I^i}(\mathbf{a}_j)\|_2^2 \end{aligned} \quad (85)$$

$$\begin{aligned} &\leq (1 + \delta_w) \cdot \left(\sum_{j=1}^{n-1} \|\mathbf{a}_j - \text{NN}_{I^i}(\mathbf{a}_j)\|_2^2 + \|\mathbf{a}'_n - \text{NN}_{I^i}(\mathbf{a}'_n)\|_2^2 \right) \\ &= (1 + \delta_w) \cdot M(I^i|\mathcal{A}') , \end{aligned} \quad (86)$$

as indeed computing the nearest neighbors do not involve the n^{th} element of the sets, *i.e.* \mathbf{a}_n or \mathbf{a}'_n . We have used in ineq. (85) the fact that \mathcal{A} is δ_w -spread.

When $n \in I^i$, eq. (84) is an immediate consequence of Lemma 21 in which the distinct elements of \mathcal{A} and \mathcal{A}' play the role of \mathbf{x} and \mathbf{x}' . ■

Lemma 23 *For any $\delta_w > 0$, if \mathcal{A} is δ_w -spread, then for any $\mathcal{N} \subseteq \mathcal{A}$ with $|\mathcal{N}| = k - 1$, $\forall \mathbf{x} \in \Omega$, it holds that $\|\mathbf{x} - \text{NN}_{\mathcal{N}}(\mathbf{x})\|_2^2 \leq \delta_w \sum_{\mathbf{a} \in \mathcal{A}} \|\mathbf{a} - \text{NN}_{\mathcal{N}}(\mathbf{a})\|_2^2$.*

Proof Follows directly from the fact that $\|\mathbf{x} - \text{NN}_{\mathcal{N}}(\mathbf{x})\|_2^2 \leq R^2$ by assumption. ■

Letting $I(k)$ denote a sequence containing element n pushed to the end of the sequence, we get:

$$\begin{aligned} &\sum_{\boldsymbol{\sigma} \in S_k} \sum_{I \in \text{Seq}^+(n:k)} p(\boldsymbol{\sigma}, I, \mathcal{C}|\mathcal{A}') \\ &= \sum_{\boldsymbol{\sigma} \in S_k} \sum_{I \in \text{Seq}^+(n:k)} \frac{N(I)}{\prod_{i=1}^k M(I^i|\mathcal{A}')} \cdot p_{\mathbf{a}'_n}(\mathbf{c}_{\boldsymbol{\sigma}(i)}) \cdot \prod_{i=1: I_i \neq n}^k p_{\mathbf{a}_{I_i}}(\mathbf{c}_{\boldsymbol{\sigma}(i)}) \\ &\leq (1 + \eta)^{2(k-2)} \\ &\quad \cdot \sum_{\boldsymbol{\sigma} \in S_k} \sum_{I \in \text{Seq}^+(n:k)} \frac{N(I(k))}{\prod_{i=1}^k M(I^i|\mathcal{A}')} \cdot p_{\mathbf{a}'_n}(\mathbf{c}_{\boldsymbol{\sigma}(i)}) \cdot \prod_{i=1: I_i \neq n}^k p_{\mathbf{a}_{I_i}}(\mathbf{c}_{\boldsymbol{\sigma}(i)}) . \end{aligned} \quad (87)$$

Now, take any element $I \in Seq_+(n : k)$ with \mathbf{a}'_n in position k , and change \mathbf{a}'_n by some $\mathbf{a} \in \mathcal{A}$. Any of these changes generates a different element $I' \in Seq^-(n : k)$, and so using Lemma 23 and the following two facts:

- the fact that

$$p_{\mathbf{a}'_n}(\mathbf{c}_{\sigma(i)}) \leq \varrho(R) \cdot p_{\mathbf{a}}(\mathbf{c}_{\sigma(i)}) , \quad (88)$$

for any $\mathbf{a} \in \mathcal{A}$,

- the fact that, if \mathcal{A} is δ_s -monotonic,

$$M(I^i_{\mathbf{a}}|\mathcal{A}) \leq (1 + \delta_s) \cdot M(I^i|\mathcal{A}) , \quad (89)$$

for any $\mathbf{a} \in \mathcal{A}$ not already in the sequence, where $I_{\mathbf{a}}$ denotes the sequence I in which \mathbf{a}'_n has been replaced by \mathbf{a} ,

we get from ineq. (87),

$$\begin{aligned} & \sum_{\sigma \in S_k} \sum_{I \in Seq^+(n:k)} p(\sigma, I, \mathcal{C}|\mathcal{A}') \\ & \leq (1 + \eta)^{2(k-2)} \cdot (1 + \delta_s)^{k-1} \cdot \delta_w \\ & \quad \cdot \varrho(R) \cdot \sum_{\sigma \in S_k} \sum_{I \in Seq^-(n:k)} \frac{N(I)}{\prod_{i=1}^k M(I^i|\mathcal{A})} \cdot \prod_{i=1}^k p_{\mathbf{a}_{I_i}}(\mathbf{c}_{\sigma(i)}) . \end{aligned} \quad (90)$$

Lemma 24 *For any $\delta_w, \delta_s > 0$ such that \mathcal{A} is δ_w -spread and δ_s -monotonic, for any $\mathcal{A}' \approx \mathcal{A}$, we have:*

$$\frac{\mathbb{P}[\mathcal{C}|\mathcal{A}']}{\mathbb{P}[\mathcal{C}|\mathcal{A}]} \leq (1 + \delta_w)^{k-1} \cdot \left(1 + \delta_w \cdot \left(\frac{1 + \delta_s}{1 + \delta_w} \right)^{k-1} \cdot (1 + \eta)^{2(k-2)} \cdot \varrho(R) \right) . \quad (91)$$

Proof We get from the fact that \mathcal{A} is δ_w -spread,

$$\sum_{\sigma \in S_k} \sum_{I \in Seq^-(n:k)} p(\sigma, I, \mathcal{C}|\mathcal{A}') \leq (1 + \delta_w)^{k-1} \cdot \sum_{\sigma \in S_k} \sum_{I \in Seq^-(n:k)} p(\sigma, I, \mathcal{C}|\mathcal{A}) , \quad (92)$$

and furthermore ineq. (90) yields:

$$\begin{aligned}
\frac{\mathbb{P}[\mathcal{C}|\mathcal{A}']}{\mathbb{P}[\mathcal{C}|\mathcal{A}]} &= \frac{\sum_{\sigma \in S_k} \sum_{I \in \text{Seq}(n:k)} p(\sigma, I, \mathcal{C}|\mathcal{A}')}{\sum_{\sigma \in S_k} \sum_{I \in \text{Seq}(n:k)} p(\sigma, I, \mathcal{C}|\mathcal{A})} \\
&\leq \frac{\left((1 + \delta_w)^{k-1} \cdot \sum_{\sigma \in S_k} \sum_{I \in \text{Seq}_-(n:k)} p(\sigma, I, \mathcal{C}|\mathcal{A}) \right.}{\sum_{\sigma \in S_k} \sum_{I \in \text{Seq}_+(n:k)} p(\sigma, I, \mathcal{C}|\mathcal{A}') + \left. \sum_{\sigma \in S_k} \sum_{I \in \text{Seq}(n:k)} p(\sigma, I, \mathcal{C}|\mathcal{A}) \right)} \\
&\leq (1 + \delta_w)^{k-1} \cdot \frac{\left(\sum_{\sigma \in S_k} \sum_{I \in \text{Seq}_-(n:k)} p(\sigma, I, \mathcal{C}|\mathcal{A}) \right.}{\delta_w \cdot \left(\frac{1+\delta_s}{1+\delta_w} \right)^{k-1} \cdot (1 + \eta)^{2(k-2)} \cdot \varrho(R) \cdot \sum_{\sigma \in S_k} \sum_{I \in \text{Seq}_-(n:k)} p(\sigma, I, \mathcal{C}|\mathcal{A}') + \left. \sum_{\sigma \in S_k} \sum_{I \in \text{Seq}(n:k)} p(\sigma, I, \mathcal{C}|\mathcal{A}) \right)} \\
&= (1 + \delta_w)^{k-1} \cdot \left(1 + \delta_w \cdot \left(\frac{1 + \delta_s}{1 + \delta_w} \right)^{k-1} \cdot (1 + \eta)^{2(k-2)} \cdot \varrho(R) \right) \\
&\quad \cdot \underbrace{\frac{\sum_{\sigma \in S_k} \sum_{I \in \text{Seq}_-(n:k)} p(\sigma, I, \mathcal{C}|\mathcal{A})}{\sum_{\sigma \in S_k} \sum_{I \in \text{Seq}(n:k)} p(\sigma, I, \mathcal{C}|\mathcal{A})}}_{\leq 1}.
\end{aligned}$$

This ends the proof of Lemma 24. ■

Since

$$\begin{aligned}
&(1 + \delta_w)^{k-1} \cdot \left(1 + \delta_w \cdot \left(\frac{1 + \delta_s}{1 + \delta_w} \right)^{k-1} \cdot (1 + \eta)^{2(k-2)} \cdot \varrho(R) \right) \\
&= (1 + \delta_w)^{k-1} + (1 + \eta)^{2(k-2)} \cdot \delta_w \cdot (1 + \delta_s)^{k-1} \cdot \varrho(R),
\end{aligned}$$

and $\eta \leq 3$ from Lemma 19, we get Theorem 9 with

$$f(k) \doteq 4^{2k-4}. \quad (93)$$

Proof of Theorem 10

Assume that density \mathcal{D} contains a L_2 ball $\mathcal{B}_2(\mathbf{0}, R)$ of radius R , centered without loss of generality in $\mathbf{0}$. Fix $0 < \kappa < m - 1$. For any $\alpha \in (0, 1)$ and $\mathcal{N} \subseteq \mathcal{A}$ with $|\mathcal{N}| \in \{1, 2, \dots, \kappa\} \doteq [\kappa]_*$, let $\mathcal{N} \oplus \alpha \doteq \cup_{\mathbf{x} \in \mathcal{N}} \mathcal{B}_2(\mathbf{x}, \alpha \cdot R)$ be the union of all small balls centered around each element of \mathcal{N} , each of radius $\alpha \cdot R$. An important quantity is

$$q_* \doteq \min_{\mathcal{N} \subseteq \mathcal{A}, |\mathcal{N}| \in [\kappa]_*} \frac{\mu(\mathcal{B}_2(\mathbf{0}, R) \setminus \mathcal{N} \oplus \alpha)}{\mu(\mathcal{B}_2(\mathbf{0}, R))} \quad (94)$$

the minimal mass of $\mathcal{B}_2(\mathbf{0}, R) \setminus \mathcal{N} \oplus \alpha$ relatively to $\mathcal{B}_2(\mathbf{0}, R)$ as measured using \mathcal{D} . As depicted in Figure 7, q_* is a minimal value of the probability to escape the neighborhoods of $\mathcal{N} \oplus \alpha$ when

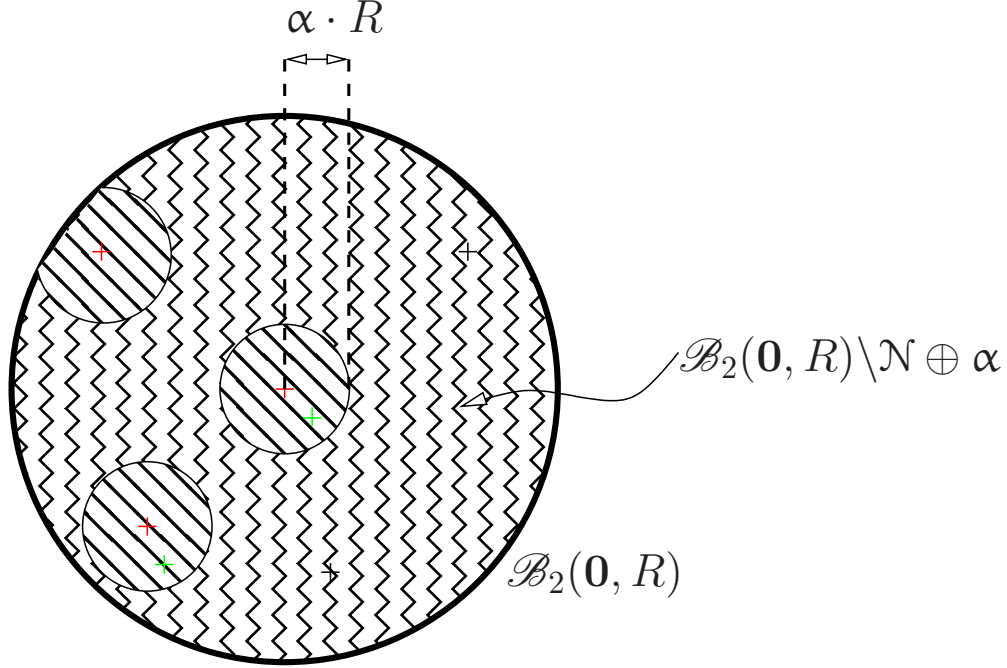


Figure 7: q_* in eq. (94) measures the probability the a point drawn in $\mathcal{B}_2(\mathbf{0}, R)$ escapes the neighborhoods of $\mathcal{N} \oplus \alpha$. In this example, two points in black escape the neighborhoods (defined by three points in red), while two in green do not.

sampling points according to \mathcal{D} in ball $\mathcal{B}_2(\mathbf{0}, R)$. If, for some α that shall depend upon the dimension d and κ , q_* is large enough, then the spread of points drawn shall guarantee "small" values for δ_w and δ_s .

This is formalized in the following Theorem, which assumes $\epsilon_m = \epsilon_M = 1$, *i.e.* the ball has uniform density. Theorem 10 is a direct consequence of this Theorem.

Theorem 25 Suppose $\mathcal{A} \subset \mathcal{B}_2(\mathbf{0}, R)$. For any $\delta \in (0, 1)$, if

$$m \geq 3 \left(\frac{\kappa}{q_* \delta^2} \right)^2, \quad (95)$$

then there is probability $\geq 1 - \delta$ over its sampling that \mathcal{A} is δ_w -spread and δ_s -monotonic for the following values of δ_w, δ_s :

$$\delta_w = \frac{1}{q_*(1 - \delta)(m - \kappa - 1)\alpha^2}, \quad (96)$$

$$\delta_s = \frac{m}{m - \kappa} \cdot \left(\frac{2}{\min \left\{ \frac{1}{4}, q_*(1 - \delta) \right\} \cdot \alpha} \right)^2 - 1. \quad (97)$$

Proof We first prove the following Lemma.

Lemma 26 Suppose $\mathcal{A} \subset \mathcal{B}_2(\mathbf{0}, R)$. Let q_* be defined as in eq. (94). Then for any $\delta \in (0, 1)$, if m meets ineq. (95), then there is probability $\geq 1 - \delta$ that

$$|(\mathcal{B}_2(\mathbf{0}, R) \setminus \mathcal{N} \oplus \alpha) \cap (\mathcal{A} \setminus \mathcal{N})| \geq q_*(1 - \delta)(m - \kappa), \forall \mathcal{N} \subseteq \mathcal{A}, |\mathcal{N}| \in [\kappa]_* . \quad (98)$$

Proof Since we assume $\mathcal{A} \subset \mathcal{B}_2(\mathbf{0}, R)$, Chernoff bounds imply that for any fixed $\mathcal{N} \subseteq \mathcal{A}$ with $|\mathcal{N}| \in [\kappa]_*$,

$$\mathbb{P}_{\mathcal{D}} \left[\frac{|(\mathcal{B}_2(\mathbf{0}, R) \setminus \mathcal{N} \oplus \alpha) \cap (\mathcal{A} \setminus \mathcal{N})|}{|\mathcal{A} \setminus \mathcal{N}|} \leq q_*(1 - \delta) \right] \leq \exp(-\delta^2 q_* |\mathcal{A} \setminus \mathcal{N}| / 2) . \quad (99)$$

Now, remark that

$$\sum_{j=1}^{\kappa} \binom{m}{j} \leq m^{\kappa}, \forall m, \kappa \geq 1 . \quad (100)$$

This can be proven by induction, m being fixed: it trivially holds for $\kappa = 1$ and $\kappa = 2$, and furthermore

$$\begin{aligned} \sum_{j=1}^{\kappa} \binom{m}{j} &= \sum_{j=1}^{\kappa-1} \binom{m}{j} + \binom{m}{\kappa} \\ &\leq m^{\kappa-1} + \frac{m!}{(m - \kappa)! \kappa!} , \end{aligned} \quad (101)$$

by induction at rank $\kappa - 1$. To prove that the right-hand side of (101) is no more than m^{κ} , we just have to remark that

$$\begin{aligned} \frac{m!}{(m - \kappa)! \kappa! m^{\kappa-1}} &< \frac{m}{\kappa!} \\ &\leq m - 1 , \end{aligned} \quad (102)$$

as long as $\kappa > 1$ and $m > 1$. So, the property at rank $\kappa - 1$ for $\kappa > 1$ implies property at rank κ , which concludes the induction.

So, we have at most m^{κ} choices for \mathcal{N} , so relaxing the choice of \mathcal{N} , we get

$$\begin{aligned} \mathbb{P}_{\mathcal{D}} \left[\exists \mathcal{N} \subseteq \mathcal{A}, |\mathcal{N}| = \kappa : \frac{|(\mathcal{B}_2(\mathbf{0}, R) \setminus \mathcal{N} \oplus \alpha) \cap \mathcal{A}_{\mathcal{N}}|}{|\mathcal{A}_{\mathcal{N}}|} \leq q_*(1 - \delta) \right] \\ \leq m^{\kappa} \exp \left(-\frac{\delta^2 q_*(m - \kappa)}{2} \right) . \end{aligned} \quad (103)$$

We want to compute the minimal m such that the right-hand side is no more than δ , this being equivalent to

$$\delta^2 q_* m \geq 2 \log \left(\frac{m^{\kappa}}{\delta} \right) + \kappa \delta^2 q_* ,$$

which, since $\delta \in (0, 1)$, is ensured if

$$\delta^2 q_* m \geq 2\kappa \log \left(\frac{m}{\delta} \right) + \kappa \delta^2 q_* . \quad (104)$$

Suppose

$$m = 3 \left(\frac{\kappa}{q_* \delta^2} \right)^2 .$$

Since we trivially have $\kappa^2 / (q_* \delta^2)^2 \geq \kappa \delta^2 q_*$ ($\kappa \geq 1, q_* \in (0, 1), \delta \in (0, 1)$), it is sufficient to prove:

$$\frac{2\kappa}{q_* \delta^2} \geq 2 \log 3 + 2 \log \left(\frac{\kappa^2}{q_*^2 \delta^5} \right) , \quad (105)$$

which, again observing that $\delta \in (0, 1)$, holds if we can prove

$$\frac{\kappa}{q_* \delta^2} \geq \log 2 + \frac{3}{2} \cdot \log \left(\frac{\kappa}{q_* \delta^2} \right) , \quad (106)$$

which is equivalent to showing $x \geq (3/2) \log x + \log 2$ for $x \geq 1$, which indeed holds (end of the proof of Lemma 26). \blacksquare

The consequence of Lemma 26 is the following: if $\mathcal{A} \subset \mathcal{B}_2(\mathbf{0}, R)$ and m satisfies (95), then for any $\mathcal{N} \subseteq \mathcal{A}$ with $|\mathcal{N}| = k - 1$, and any $\mathcal{B} \subseteq \mathcal{A}$ with $|\mathcal{B}| = |\mathcal{A}| - 1$,

$$\sum_{\mathbf{a} \in \mathcal{B}} \|\mathbf{a} - \text{NN}_{\mathcal{N}}(\mathbf{a})\|_2^2 \geq q_*(1 - \delta)(m - \kappa - 1)\alpha^2 \cdot R^2 , \quad (107)$$

and so from Definition 7 \mathcal{A} is δ_s -spread for:

$$\delta_w = \frac{1}{q_*(1 - \delta)(m - \kappa - 1)\alpha^2} . \quad (108)$$

Now, suppose we add a single point \mathbf{x}_* in \mathcal{N} . If, for some fixed $\alpha_* \in (0, \alpha/2]$,

$$\mathbf{x}_* \notin \mathbf{a} \oplus \alpha_* , \forall \mathbf{a} \in \mathcal{A} , \quad (109)$$

then because of (107),

$$\sum_{\mathbf{a} \in \mathcal{A}} \|\mathbf{a} - \text{NN}_{\mathcal{N} \cup \{\mathbf{x}_*\}}(\mathbf{a})\|_2^2 \geq (m - \kappa) \cdot \min \{ \alpha_*^2, q_*(1 - \delta)\alpha^2 \} \cdot R^2 . \quad (110)$$

Otherwise, consider one \mathbf{a}_* for which $\mathbf{x}_* \in \mathbf{a}_* \oplus \alpha_*$. If we replace \mathbf{a}_* by \mathbf{x}_* in all \mathcal{N} in which \mathbf{a}_* belongs to in Lemma 26, then because $\mathbf{x}_* \oplus \alpha_* \subset \mathbf{a}_* \oplus \alpha$, it comes from Lemma 26:

$$\sum_{\mathbf{a} \in \mathcal{A}} \|\mathbf{a} - \text{NN}_{\mathcal{N} \cup \{\mathbf{x}_*\}}(\mathbf{a})\|_2^2 \geq \frac{1}{4} \cdot (m - \kappa) \cdot q_*(1 - \delta)\alpha^2 \cdot R^2 . \quad (111)$$

We thus get in all cases

$$\sum_{\mathbf{a} \in \mathcal{A}} \|\mathbf{a} - \text{NN}_{\mathcal{N} \cup \{\mathbf{c}(\mathcal{A})\}}(\mathbf{a})\|_2^2 \geq \min \left\{ \frac{\alpha^2}{4}, \alpha_*^2, q_*(1 - \delta)\alpha^2 \right\} (m - \kappa) \cdot q_*(1 - \delta) \cdot R^2 \quad (112)$$

where $\mathbf{c}(A)$ is the arithmetic average computed according to the definition of δ_s -monotonicity, of any $A \subseteq \mathcal{A} \setminus \mathcal{N}$. Since $\mathcal{N} \subseteq \mathcal{A} \subset \mathcal{B}_2(\mathbf{0}, R)$, we have $\sum_{\mathbf{a} \in \mathcal{A}} \|\mathbf{a} - \text{NN}_{\mathcal{N}}(\mathbf{a})\|_2^2 \leq 4mR^2$, and so

$$\sum_{\mathbf{a} \in \mathcal{A}} \|\mathbf{a} - \text{NN}_{\mathcal{N}}(\mathbf{a})\|_2^2 \leq \frac{4m}{\min\left\{\frac{\alpha^2}{4}, \alpha_*^2, q_*(1-\delta)\alpha^2\right\} (m-\kappa) \cdot q_*(1-\delta)} \cdot \sum_{\mathbf{a} \in \mathcal{A}} \|\mathbf{a} - \text{NN}_{\mathcal{N} \cup \{\mathbf{c}(A)\}}(\mathbf{a})\|_2^2 \quad (113)$$

implying from Definition 8 that δ_s -monotonicity holds with:

$$\delta_s = \frac{m}{m-\kappa} \cdot \frac{4}{\min\left\{\frac{\alpha^2}{4}, \alpha_*^2, q_*(1-\delta)\alpha^2\right\} \cdot q_*(1-\delta)} - 1. \quad (114)$$

The statement of the Theorem follows with $\alpha_* = \alpha/2$ (end of the proof of Theorem 25). \blacksquare

We finish the proof of Theorem 10. We have

$$q_* \geq 1 - \kappa\alpha^d, \quad (115)$$

where the lowerbound corresponds to the case where all neighborhoods in $\mathcal{N} \oplus \alpha$ are distinct and included in $\mathcal{B}_2(\mathbf{0}, R)$. So we have, for any fixed choice of $\alpha \in (0, 1)$,

$$\delta_w \leq \frac{1}{\alpha^2 \cdot (1 - \kappa\alpha^d)(1-\delta)(m-\kappa-1)}. \quad (116)$$

To minimize this upperbound, we pick α to maximize $\alpha^2 \cdot (1 - \kappa\alpha^d)$ with $\alpha \in (0, 1)$, which is easily achieved picking

$$\alpha = \left(\frac{1}{\kappa(d+1)} \right)^{\frac{1}{d}}, \quad (117)$$

and yields

$$\begin{aligned} \delta_w &\leq \left(1 + \frac{1}{d}\right) \cdot \frac{1}{(\kappa(d+1))^{\frac{2}{d}}(1-\delta)(m-\kappa-1)} \\ &\leq \left(1 + \frac{1}{d}\right) \cdot \frac{1}{\kappa^{\frac{2}{d}}(1-\delta)(m-\kappa-1)}. \end{aligned} \quad (118)$$

But we have for this choice, $1 - \kappa\alpha^d = d/(d+1) \geq 1/2$, so as long as

$$\delta < 1/2, \quad (119)$$

we shall have $q_*(1-\delta) > 1/4$ and so we shall have

$$\begin{aligned} \delta_s + 1 &= 64 \cdot \frac{m}{m-\kappa} \cdot \frac{1}{\alpha^2} \\ &\leq 64 \cdot \frac{m}{m-\kappa} \cdot \frac{1}{\kappa^{\frac{2}{d}}}. \end{aligned} \quad (120)$$

We now go back to ineq. (14), which reads:

$$\frac{\mathbb{P}[\mathcal{C}|\mathcal{A}']}{\mathbb{P}[\mathcal{C}|\mathcal{A}]} \leq \varrho_1 + \varrho_2, \quad (121)$$

with

$$\varrho_1 \doteq (1 + \delta_w)^{k-1}, \quad (122)$$

$$\varrho_2 \doteq f(k) \cdot \delta_w \cdot (1 + \delta_s)^{k-1} \cdot \varrho(R). \quad (123)$$

We upperbound separately both terms.

Lemma 27 *Suppose ineqs (119) and (15) are met. Then*

$$\varrho_1 \leq 1 + \frac{4}{m^{\frac{1}{4} + \frac{1}{d+1}}}. \quad (124)$$

Proof Since $d \geq 1$ and $\delta < 1/2$, we get from ineq. (118) (using $\kappa = k$)

$$\begin{aligned} (1 + \delta_w)^{k-1} &\leq \left(1 + \left(1 + \frac{1}{d} \right) \cdot \frac{1}{k^{\frac{2}{d}}(1 - \delta)(m - k - 1)} \right)^{k-1} \\ &\leq \left(1 + \frac{2}{k^{\frac{2}{d}}(1 - \delta)(m - k - 1)} \right)^{k-1} \\ &\leq \left(1 + \frac{4}{k^{\frac{2}{d}}(m - k - 1)} \right)^{k-1}. \end{aligned} \quad (125)$$

Let $h(k)$ be the right-hand side of ineq. (125). $h(1)$ trivially meets ineq. (124). When $k \geq 2$, h decreases until $k = 2(m - 1)/(d + 2)$ and then increases. We thus just need to check ineq. (124) for $k = 2$ and $k = \sqrt{m}$ from ineq. (15). We get $h(2) = 1 + 4/(4^{1/d}(m - 3))$. For ineq. (124) to be satisfied, we need to have $4^{1/d}(m - 3) \geq m^{\frac{1}{4} + \frac{1}{d+1}}$, which holds if $m \geq 3 + m^{3/4}$ ($d \geq 1$), that is, $m \geq 8$. But since ineqs (119) and (15) are satisfied, we have $m \geq 16k^2/\delta^2 \geq 64k^2 \geq 64$, and so $h(2)$ satisfies ineq. (124).

There remains to check ineq. (124) for $k = \sqrt{m}$. We have

$$\begin{aligned} h(\sqrt{m}) &= \left(1 + \frac{4}{m^{\frac{1}{d}}(m - \sqrt{m} - 1)} \right)^{\sqrt{m}-1} \\ &\leq \left(1 + \frac{4}{m^{\frac{1}{d}}(m - \sqrt{m})} \right)^{\sqrt{m}} \\ &\leq \left(1 + \frac{2}{\sqrt{m} \cdot m^{\frac{1}{4} + \frac{1}{d}}} \right)^{\sqrt{m}}, \end{aligned} \quad (126)$$

since any $m \geq 64$, we have $m - \sqrt{m} \geq 2m^{3/4}$. To conclude, ineq (126) yields

$$\begin{aligned} h(\sqrt{m}) &\leq \left(1 + \frac{2}{\sqrt{m} \cdot m^{\frac{1}{4} + \frac{1}{d}}} \right)^{\sqrt{m}} \\ &\leq \exp \left(\frac{2}{m^{\frac{1}{4} + \frac{1}{d}}} \right) \\ &\leq 1 + \frac{4}{m^{\frac{1}{4} + \frac{1}{d}}}. \end{aligned} \quad (127)$$

The penultimate ineq. comes from $1 + x \leq \exp x$, and the last one comes from the fact that $\exp(2x) \leq 1 + 4x$ for $x \leq 1$. Since $m^{\frac{1}{4}+\frac{1}{d}} \geq m^{\frac{1}{4}+\frac{1}{d+1}}$, we obtain the statement of the Lemma for $h(\sqrt{m})$. This concludes the proof of Lemma 27. \blacksquare

Lemma 28 *Suppose ineqs (119) and (15) are met. Then*

$$\varrho_2 \leq \left(\frac{64}{k^{\frac{2}{d}}} \right)^k \cdot \frac{\varrho(2R)}{m}. \quad (128)$$

Proof We fix $\kappa = k$, use $f(k) = 4^{2k-4}$ (eq. 93), so we get

$$\begin{aligned} \varrho_2 &= 4^{2k-2} \cdot \left(1 + \frac{1}{d} \right) \cdot \frac{1}{k^{\frac{2}{d}}(1-\delta)(m-k-1)} \cdot \left(64 \cdot \frac{m}{m-k} \cdot \frac{1}{k^{\frac{2}{d}}} \right)^{k-1} \cdot \varrho(2R) \\ &\leq 2 \cdot 64^{k-1} \cdot \left(1 + \frac{1}{d} \right) \cdot \frac{1}{k^{\frac{2k}{d}}(m-k-1)} \cdot \left(1 + \frac{k}{m-k} \right)^{k-1} \cdot \varrho(2R) \end{aligned} \quad (129)$$

$$\leq \underbrace{4 \cdot \frac{1}{(m-k-1)} \cdot \left(1 + \frac{k}{m-k} \right)^{k-1}}_{\doteq \varrho_3} \cdot 64^{k-1} \cdot \frac{1}{k^{\frac{2k}{d}}} \cdot \varrho(2R), \quad (130)$$

using the fact that $\delta < 1/2$ and $d \geq 1$. Now, we also have

$$\left(1 + \frac{k}{m-k} \right)^{k-1} \leq \exp \left(\frac{k^2}{m-k} \right) \quad (131)$$

$$\leq e, \quad (132)$$

as long as $k \leq (1/16) \cdot \sqrt{m}$, and furthermore, since $m \geq 64$ (see the proof of Lemma 27), we also have $1/(m-k-1) \leq 5/m$. We thus obtain

$$\begin{aligned} \varrho_3 &\leq \frac{20e}{m} \\ &\leq \frac{64}{m}, \end{aligned} \quad (133)$$

which yields

$$\varrho_2 \leq \left(\frac{64}{k^{\frac{2}{d}}} \right)^k \cdot \frac{\varrho(2R)}{m}, \quad (134)$$

as claimed. \blacksquare

Putting altogether Lemmata 27 and 28, we get:

$$\frac{\mathbb{P}[\mathcal{C}|\mathcal{A}']}{\mathbb{P}[\mathcal{C}|\mathcal{A}]} \leq 1 + \frac{4}{m^{\frac{1}{4}+\frac{1}{d+1}}} + \left(\frac{64}{k^{\frac{2}{d}}} \right)^k \cdot \frac{\varrho(2R)}{m}, \quad (135)$$

as claimed. There remains to check that, with our choice of α , the constraint on m in (95) is satisfied if

$$m \geq \frac{12k^2}{\delta^4} \quad (136)$$

since $q_* \geq d/(d+1)$. We obtain the sufficient constraint on k :

$$k \leq \frac{\delta^2}{4} \cdot \sqrt{m} \ , \quad (137)$$

which proves Theorem 10 when $\epsilon_m = \epsilon_M = 1$.

When the density do not satisfy $\epsilon_m = \epsilon_M = 1$ we just have to remark that the lowerbound on q_* is now

$$q_* \leq \frac{\epsilon_m}{\epsilon_M} \cdot (1 - \kappa \alpha^d) \ . \quad (138)$$

Ineq. (118) becomes

$$\delta_w \leq \frac{\epsilon_M}{\epsilon_m} \cdot \left(1 + \frac{1}{d}\right) \cdot \frac{1}{\kappa^{\frac{2}{d}}(1-\delta)(m-\kappa-1)} \ , \quad (139)$$

ineq. (120) becomes

$$\delta_s + 1 \leq \frac{\epsilon_M}{\epsilon_m} \cdot 64 \cdot \frac{m}{m-\kappa} \cdot \frac{1}{\kappa^{\frac{2}{d}}} \ . \quad (140)$$

So, the only difference with the $\epsilon_m = \epsilon_M = 1$ is the ratio $\epsilon_M/\epsilon_m (\geq 1)$ which multiplies all quantities of interest, and yields, in lieu of ineq. (135),

$$\frac{\mathbb{P}[\mathcal{C}|\mathcal{A}']}{\mathbb{P}[\mathcal{C}|\mathcal{A}]} \leq 1 + \left(\frac{\epsilon_M}{\epsilon_m}\right)^k \cdot \left(\frac{4}{m^{\frac{1}{4}+\frac{1}{d+1}}} + \left(\frac{64}{k^{\frac{2}{d}}}\right)^k \cdot \frac{\varrho(2R)}{m}\right) \ , \quad (141)$$

which is the statement of Theorem 10.

Proof of Theorem 12

When $p(\mu_a, \theta_a)$ is a product of Laplace distributions $Lap(b)$ (b being the *scale* parameter of the distribution (Dwork & Roth, 2014)), condition in ineq. (13) becomes:

$$\begin{aligned} \frac{p(\mu_{a'}, \theta_{a'}) (\mathbf{x})}{p(\mu_a, \theta_a) (\mathbf{x})} &\leq \exp\left(\frac{\|\mathbf{a} - \mathbf{a}'\|_1}{b}\right) \\ &= \exp\left(\frac{\sqrt{2}\|\mathbf{a} - \mathbf{a}'\|_1}{\sigma_1}\right) \\ &\leq \exp\left(\frac{2\sqrt{2}R}{\sigma_1}\right) \ , \forall \mathbf{a}, \mathbf{a}' \in \mathcal{A}, \forall \mathbf{x} \in \Omega \ , \end{aligned} \quad (142)$$

assuming $\mathcal{A} \subset \mathcal{B}_1(\mathbf{0}, R)$. Let us fix $\varrho(R) \doteq \exp(2\sqrt{2}R/\sigma_1)$. Since $\mathcal{B}_1(\mathbf{0}, R) \subset \mathcal{B}_2(\mathbf{0}, R)$ (the L_2 ball), we now want $(1 + \delta_w)^{k-1} + f(k) \cdot \delta_w \cdot (1 + \delta_s)^{k-1} \cdot \varrho(R) = \exp(\epsilon)$. Solving for σ_1 yields:

$$\sigma_1 = \frac{2\sqrt{2}R}{\log\left(\frac{\exp(\epsilon) - (1 + \delta_w)^{k-1}}{f(k) \cdot \delta_w \cdot (1 + \delta_s)^{k-1}}\right)}, \quad (143)$$

as claimed. The proof that k -variates++ meets ineq. (7) with

$$\Phi = \Phi_1 \doteq 8 \cdot \left(\phi_{\text{opt}} + \frac{mR^2}{\tilde{\epsilon}^2} \right) \quad (144)$$

comes from a direct application of Theorem 2, with

$$\begin{aligned} \eta &= 0, \\ \phi_{\text{bias}} &= \phi_{\text{opt}}, \\ \phi_{\text{var}} &= m \cdot \left(\frac{2\sqrt{2}R}{\tilde{\epsilon}} \right)^2. \end{aligned}$$

The statements for σ_2 and Φ_2 are direct applications of the Laplace mechanism properties (Dwork & Roth, 2014; Dwork et al., 2006).

Extension to non-metric spaces

Since its inception, the k -means++ seeding technique has been successfully adapted to various distortion measures $D(\cdot\|\cdot)$ to handle non-Euclidean features (Jegelka et al., 2009; Nock et al., 2008, 2016). Similarly, our extended seeding technique can be adapted to these scenarios: this boils down to putting the distortion as a free parameter of the algorithm, replacing $D_t(\mathbf{a})$ (eq. (1)) by $D_t(\mathbf{a}) \doteq \min_{\mathbf{a}' \in \mathcal{P}} D(\mathbf{a}\|\mathbf{a}')$. For example, by noticing that the squared Euclidean distance is merely an example of Bregman divergences (the well-known canonical divergences in information geometry of dually flat spaces), k -variates++ can be extended to that family of dissimilarities (Nock et al., 2008). But more interesting examples now appear, that build on constraints that distortions have to satisfy for certain problems, like the invariance to rotations of the coordinate space. This is all the more challenging in practice for clustering since sometimes **no-closed form** solution are available for some of these divergences. Because it bypasses the construction of the population minimisers, k -variates++ offers an elegant solution to the problem. Such hard distortions include the skew Jeffreys α -centroids (Nock et al., 2016). This also include the recent class of total Bregman/Jensen divergences that are examples of conformal divergences (Nielsen & Nock, 2015; Nock et al., 2016). We give an example of the extension of k -variates++ to the total Jensen divergence, to show that k -variates++ can approximate the optimal clustering even without closed form solutions for the population minimisers (Nielsen & Nock, 2015). For any convex function $\varphi: \mathbb{R}^d \rightarrow \mathbb{R}$ and $\alpha \in (0, 1)$, the skew Jensen divergence is

$$J_\alpha(\mathbf{a}, \mathbf{a}') \doteq \alpha\varphi(\mathbf{a}) + (1 - \alpha)\varphi(\mathbf{a}') - \varphi(\alpha\mathbf{a} + (1 - \alpha)\mathbf{a}'), \quad (145)$$

and the total Jensen divergence is

$$tJ_\alpha(\mathbf{a}, \mathbf{a}') \doteq \frac{1}{\sqrt{1 + U^2}} \cdot J_\alpha(\mathbf{a}, \mathbf{a}'), \quad (146)$$

where $U \doteq (\varphi(\mathbf{a}) - \varphi(\mathbf{a}'))/\|\mathbf{a} - \mathbf{a}'\|_2$. There is no closed form solution for the population minimiser of tJ_α , yet we can prove the following Theorem, which builds upon Theorem 3 in (Nielsen & Nock, 2015).

Theorem 29 *In k -variates++, replace $D_t(\mathbf{a})$ (eq. (1)) by $D_t(\mathbf{a}) \doteq \min_{\mathbf{a}' \in \mathcal{P}} tJ_\alpha(\mathbf{a}, \mathbf{a}')$ and suppose for simplicity that probe functions are identity: $\wp_t = \text{Id}, \forall t$. Denote ϕ_{opt} the optimal noise-free potential of the clustering problem using tJ_α as distortion measure. Then there exists a **constant** $\omega > 0$ such that for **any** choice of densities $p_{\mu, \theta}$, the expected tJ_α -potential ϕ of k -variates++ satisfies:*

$$\mathbb{E}[\phi(\mathcal{A}; \mathcal{C})] \leq \omega \cdot \log k \cdot (6\phi_{\text{opt}} + 2\phi_{\text{bias}} + 2\phi_{\text{var}}) \quad , \quad (147)$$

where ϕ_{var} is defined in Theorem 2 and ϕ_{bias} is defined in eq. (4).

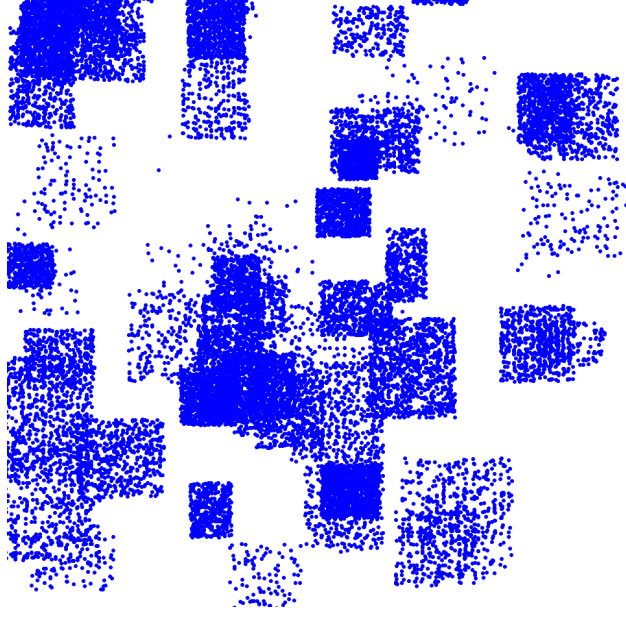


Figure 8: Final dataset for the experiments in Table 4 (plot of the two first coordinates ($d = 10$)).

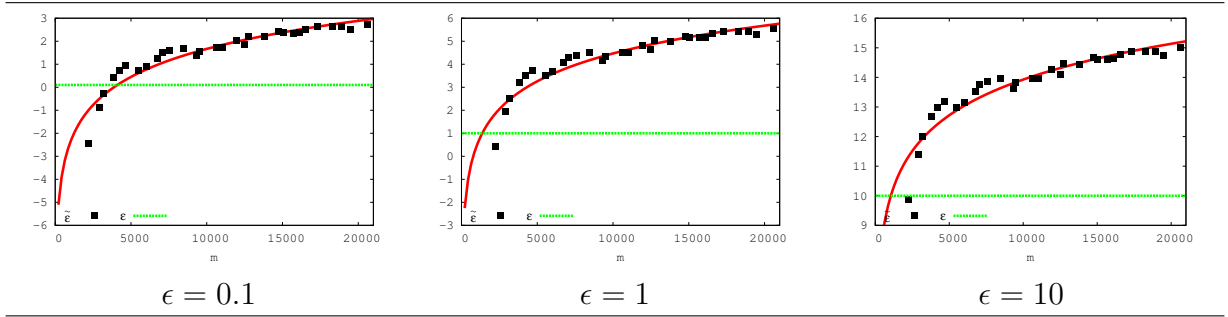


Table 4: Case $d = 10, k = 3$ — Plot of $\tilde{\epsilon}$ as in Theorem 12 (see also eq. (148) below) and best fit for model $\tilde{\epsilon} = a + b \log m$. Figure 8 displays the final dataset obtained (see text).

9 Appendix on Experiments

Experiments on Theorem 12 and the sublinear noise regime

\hookrightarrow **comments on $\tilde{\epsilon}$** An important parameter of Theorem 12 is $\tilde{\epsilon}$, which replaces ϵ in the computation of the noise standard deviation in σ_1 : the larger it is compared to ϵ , the less noise we can put while still ensuring $\mathbb{P}[\mathcal{C}|\mathcal{A}']/\mathbb{P}[\mathcal{C}|\mathcal{A}] \leq \exp \epsilon$ in Definition 11. Recall its formula:

$$\tilde{\epsilon} \doteq \log \left(\frac{\exp(\epsilon) - (1 + \delta_w)^{k-1}}{f(k) \cdot \delta_w \cdot (1 + \delta_s)^{k-1}} \right). \quad (148)$$

The experimental setting is the following one: we repeatedly sample clusters that are uniform in a subset of the domain (with limited, random size), taken to be a d -dimensional hyperrectangle of randomly chosen edge lengths. Each cluster contains a randomly picked number of points between

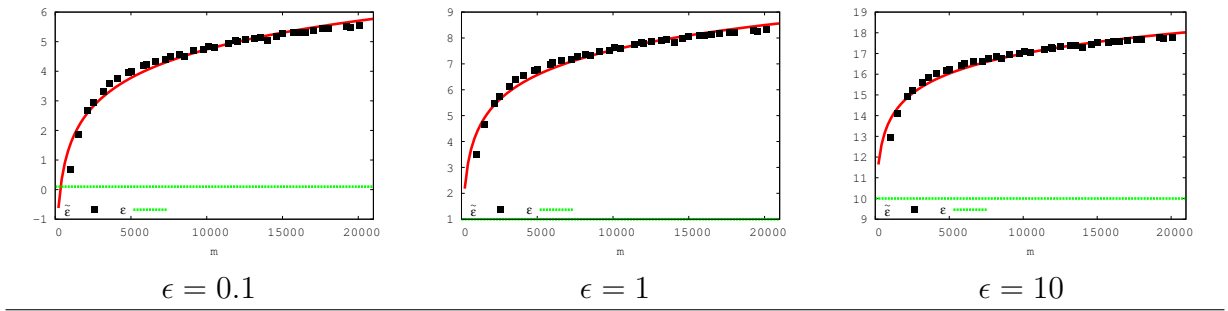


Table 5: Case $d = 50, k = 3$ — Plot of $\tilde{\epsilon}$ as in Theorem 12 (see also eq. (148) below) and best fit for model $\tilde{\epsilon} = a + b \log m$. All other parameters are the same as for Table 4.

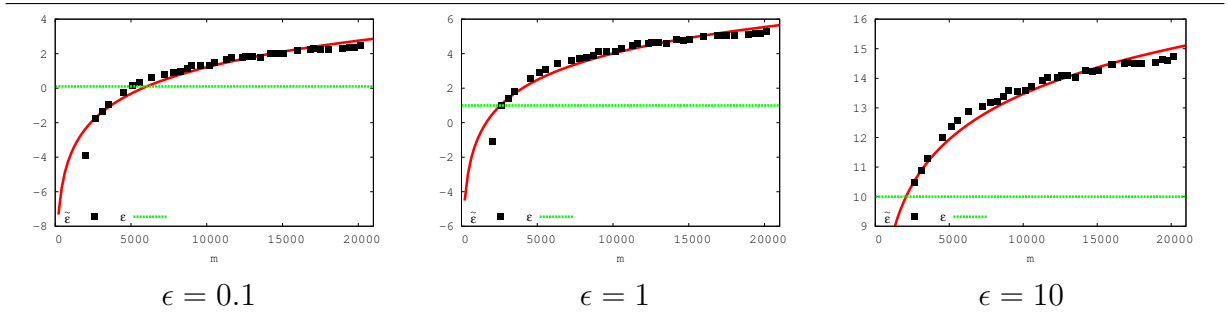


Table 6: Case $d = 50, k = 4$ — Plot of $\tilde{\epsilon}$ as in Theorem 12 (see also eq. (148) below) and best fit for model $\tilde{\epsilon} = a + b \log m$. All other parameters are the same as for Table 4.

1 and 1000. After each cluster is picked, we updated an *estimation* of δ_w and δ_s :

- we compute δ_w by randomly picking \mathcal{B} and \mathcal{N} for a total number of n_{est} iterations, with $n_{\text{est}} = 5000$;
- we compute δ_s by randomly picking \mathcal{N} for a total number of n_{est} iterations. Instead of computing A then \mathbf{x} , we opt for the fast proxy which consists in replacing $c(A)$ by a random data point, thus *without* making the \mathcal{N} -packed test. This should reasonably overestimate δ_s and thus slightly loosen our approximation bounds.

Figure 8 shows the dataset obtained for $d = 10$ at the end of the process. Predictably, the distribution on the whole space looks like a highly non-uniform cover by locally uniform clusters. Tables 4, 5 and 6 display results obtained for three different values of ϵ and three different values for the couple (d, k) . To test the large sample regime intuition and the fact that the noise dependence grows sublinearly in m , we have regressed in each plot $\tilde{\epsilon}$ as a function of m for

$$\tilde{\epsilon}(m) = a + b \log m . \quad (149)$$

The plots obtained confirm a good approximation of this intuition, but they also display some more good news. The smaller ϵ , the larger can be $\tilde{\epsilon}$ relatively to ϵ , by order of magnitudes if ϵ is small. Hence, despite the fact that we eventually overestimate δ_s , we still get large $\tilde{\epsilon}$. Furthermore, if k is small, this "large sample" regime in which $\tilde{\epsilon} > \epsilon$ actually happens for quite small values of m .

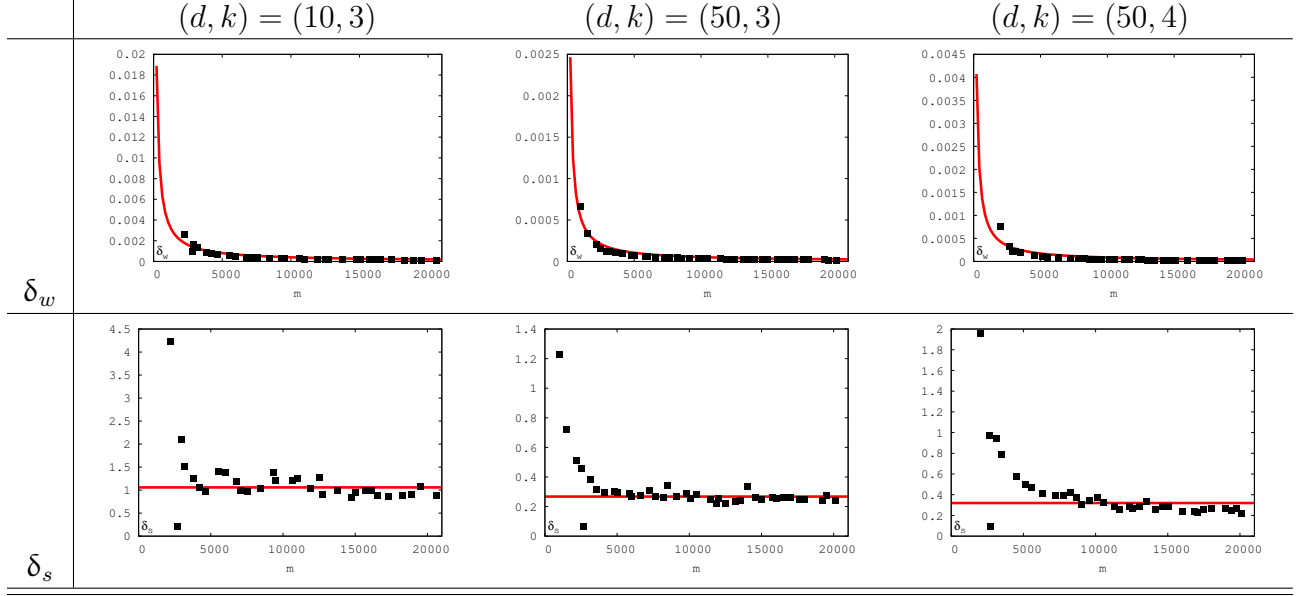


Table 7: Estimations of δ_w (top row) and δ_s (bottom row) as a function of m , for three values of (d, k) . We also indicate the best fit for $\delta_w(m) = a/m$ (top row) and $\delta_s(m) = b$ (for $m \geq 4000$, bottom row).

Also, one may remark that the curves all look like an approximate translation of the same curve. This is not surprising, since we can reformulate

$$\tilde{\epsilon} = \epsilon + \log \left(1 - \frac{U}{\epsilon} \right) + g(m) , \quad (150)$$

whene $U \doteq (1 + \delta_w)^{k-1}$ and g do not depend on ϵ . It happens that δ_w quickly decreases to very small values (bringing also a separate empirical validation of its behavior as computed in ineq. (139) in the proof of Theorem 10). Hence, we rapidly get for small m some $\tilde{\epsilon}$ that looks like

$$\begin{aligned} \tilde{\epsilon} &\approx \epsilon + \log \left(1 - \frac{1 + o(1)}{\epsilon} \right) + g(m) \\ &\approx h(\epsilon) + g(m) , \end{aligned} \quad (151)$$

which may explain what is observed experimentally.

We can summarise the global picture for $\tilde{\epsilon}$ vs ϵ by saying that it becomes more and more in favor of $\tilde{\epsilon}$ as data size (d or m) increase, but become less in favor of $\tilde{\epsilon}$ as the number of clusters k increases (predictably).

\hookrightarrow **comments on δ_w and δ_s** Table 7 presents the estimated values of δ_w and δ_s for the settings of Tables 4, 5 and 6. We wanted to test the intuition as to whether, for m sufficiently large, it would hold that $\delta_w = O(1/m)$ while $\delta_s = O(1)$. The essential part is on δ_w , since such a behaviour would be sufficient for the sublinear growth of the noise dependence. One can check that such behaviours are indeed observed, and more: δ_w converges very rapidly to zero, at least for all settings in which

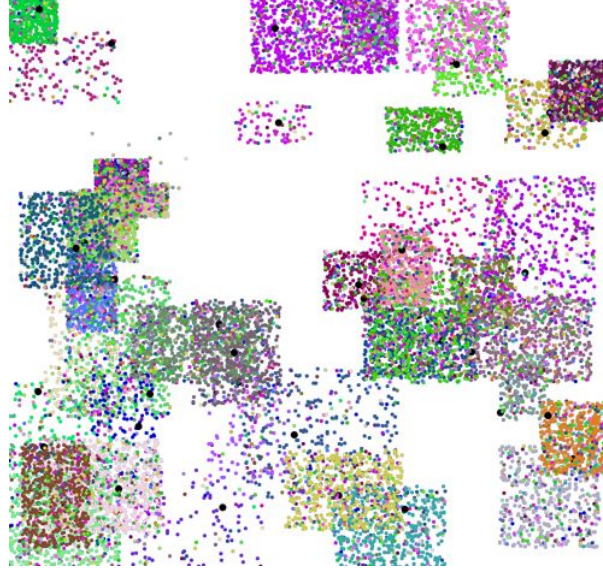


Figure 9: Example dataset obtained for $p = 50\%$ ($d = 50$). Each color represents the points held by a peer (Forgy node) after the process of moving each point from a true cluster to another cluster with probability $p = 0.5$. Big black dots are the datapoints that are the closet to the true cluster centers.

we have tested data generation. Another quite good news, is that δ_s seems indeed to be $\theta(1)$, but for an actual value which is also not large, so the denominator of eq. (148) is actually driven by $f(k)$, even when, as we already said, we may have a tendency to overestimate δ_s with our randomized procedure.

Experiments with Dk -means++, k -means++ and k -means_{||}

★ **Experiments on synthetic data** We have generated a set of $m \approx 20\,000$ points using the same kind of clusters as in the experiments related to Theorem 12: we add "true" clusters until the total number of points exceeds 20 000. To simulate the spread of data among peers (Forgy nodes) and evaluate the influence of the spread of Forgy nodes (ϕ_s^F) for Dk -means++, we have devised the following protocol: let us name "true" clusters the hyperrectangle clusters used to build the dataset. Each true cluster corresponds to the data held by a peer. Then, for some $p \in [0, 100]$ (%), *each* point in *each* true cluster moves into another cluster, with probability p . The choice of the target cluster is made uniformly at random. Thus, as p increases, we get a clustering problem in which the data held by peers is more and more spread, and for which the spread of Forgy nodes

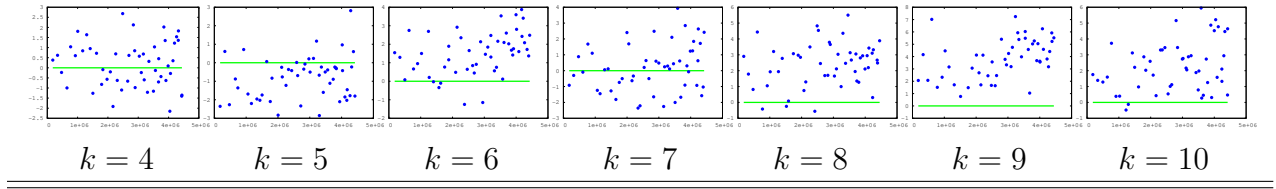


Table 8: Simulated data — Plot of ratio $\rho_\phi(k\text{-means}++)$ in eq. (152) as a function of ϕ_s^F . Points *below* the green line correspond to (average) runs in which $Dk\text{-means}++$ *beats* $k\text{-means}++$.

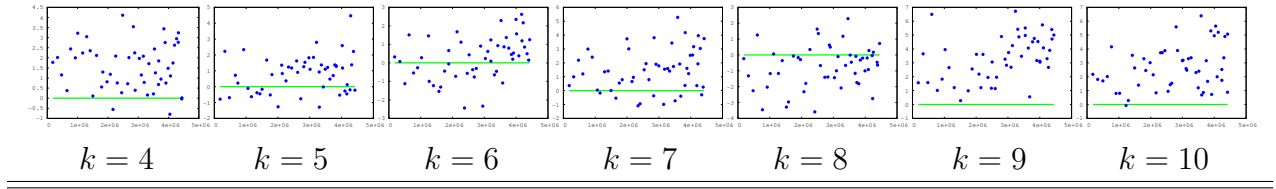


Table 9: Simulated data — Plot of ratio $\rho_\phi(k\text{-means}_\parallel)$ in eq. (152) as a function of ϕ_s^F . Points *below* the green line correspond to (average) runs in which $Dk\text{-means}++$ *beats* $k\text{-means}_\parallel$.

ϕ_s^F increases. Figure 9 presents a typical example of the spread for $p = 50\%$. Notice that in this case many Forgery nodes have data spreading through a much larger domain than the initial, true clusters. Figure 10 displays that this happens indeed, as ϕ_s^F is multiplied by a factor exceeding 20 (compared to ϕ_s^F at $p = 0$) for the largest values of p .

We have compared $Dk\text{-means}++$ to $k\text{-means}++$ and $k\text{-means}_\parallel$ (Bahmani et al., 2012). In the case of that latter algorithm, we follow the paper’s statements and pick the number of outer iterations to be $\lceil \log \phi_1 \rceil$, where ϕ_1 is the potential for one Forgery-chosen center. We also pick $\ell = 2k$, considering that it is a value which gives some of the best experimental results in (Bahmani et al., 2012). Finally, we recluster the points at the end of the algorithm using $k\text{-means}++$. For each algorithm $\mathcal{H} \in \{k\text{-means}++, k\text{-means}_\parallel\}$, we run it on the complete dataset and its results are averaged over 10 runs. We run $Dk\text{-means}++$ for each $p \in \{0\%, 1\%, \dots, 50\%\}$. More precisely, for each p , we average the results of $Dk\text{-means}++$ over 10 runs. We use as metric the relative increase in the potential of $Dk\text{-means}++$ compared to \mathcal{H} :

$$\rho_\phi(\mathcal{H}) \doteq \frac{\phi(Dk\text{-means}++) - \phi(\mathcal{H})}{\phi(\mathcal{H})} \cdot 100 . \quad (152)$$

that we plot as a function of ϕ_s^F , or surface plot as a function of (k, p) . The intuition for the former plot is that the larger ϕ_s^F , the larger should be this ratio, since the data held by peers spreads across the domain and each peer is constrained to pick its centers with uniform seeding.

$\hookrightarrow Dk\text{-means}++$ vs $k\text{-means}++$ Figure 8 presents results for $\rho_\phi(k\text{-means}++) = f(\phi_s^F)$ obtained for various k . First, the intuition is indeed confirmed for $k = 8, 9, 10$, but an interesting phenomenon appears for $k = 5$: $Dk\text{-means}++$ almost consistently beats $k\text{-means}++$. The decrease in the average potential ranges up to 3%. Furthermore, this happens even for large values of ϕ_s^F . Finally, for *all but one* value of k , there exists spread values for which $Dk\text{-means}++$ beats $k\text{-means}++$.

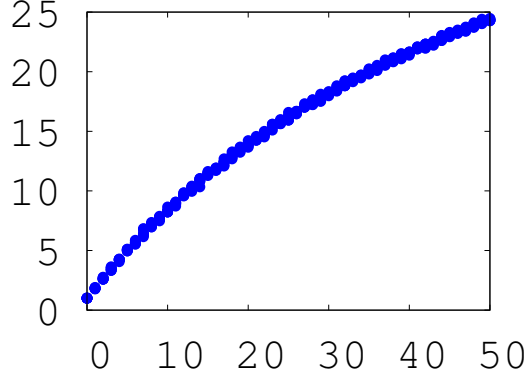


Figure 10: Simulated data — Relative increase of spread, $\phi_s^F(p)/\phi_s^F(0)$, through the runs, as a function of p .

The surface plot in Figure 3 displays that superior performances of Dk -means++ are probably not random. One possible explanation to this phenomenon relies on the expression of ϕ_{bias} given in the proof of Theorem 4 (eq. (43)), recalled here:

$$\begin{aligned}\phi_{\text{bias}} &\doteq \sum_{\mathbf{a} \in \mathcal{A}} \|\boldsymbol{\mu}_{\mathbf{a}} - \mathbf{c}_{\text{opt}}(\mathbf{a})\|_2^2 \\ &= \sum_{i \in [n]} \sum_{\mathbf{a} \in \mathcal{A}_i} \|\mathbf{c}(\mathcal{A}_i) - \mathbf{c}_{\text{opt}}(\mathbf{a})\|_2^2 .\end{aligned}\tag{153}$$

Recall that ϕ_{bias} can be $< \phi_{\text{opt}}$, and it can even be zero, in which case Theorem 2 says that the approximation bound may actually be *better* than that of k -means++ in (Arthur & Vassilvitskii, 2007) (furthermore, $\eta = 0$ for Dk -means++). Hence, what happens is probably that in several cases, there exists a union of peers data (the number of peers is larger than k) that gives a at least reasonably good approximation of the global optimum. In all our experiments indeed, we obtained a number of peers larger than 30.

\hookrightarrow **Dk -means++ vs k -means $_{\parallel}$** Figure 3 appear to display performances for Dk -means++ that are even more in favor of Dk -means++, compared to k -variates++. Figure 9 presents results for $\rho_{\phi}(k\text{-means}_{\parallel}) = f(\phi_s^F)$ obtained for various k . The fact that each of them is a vertical translation of a picture in Figure 8 comes from the fact that the results of k -means $_{\parallel}$ and k -means++ do not depend on the spread of the neighbors ϕ_s^F .

★ **Experiments on real world data** We consider the EuropeDiff dataset⁵ (Dataset characteristics provided in Table 10). Figures 11 and 12 give the results for the equivalent settings of the experimental data. To simulate N peers with real data, reasonably spread geographically, we have

⁵<http://cs.joensuu.fi/sipu/datasets/>

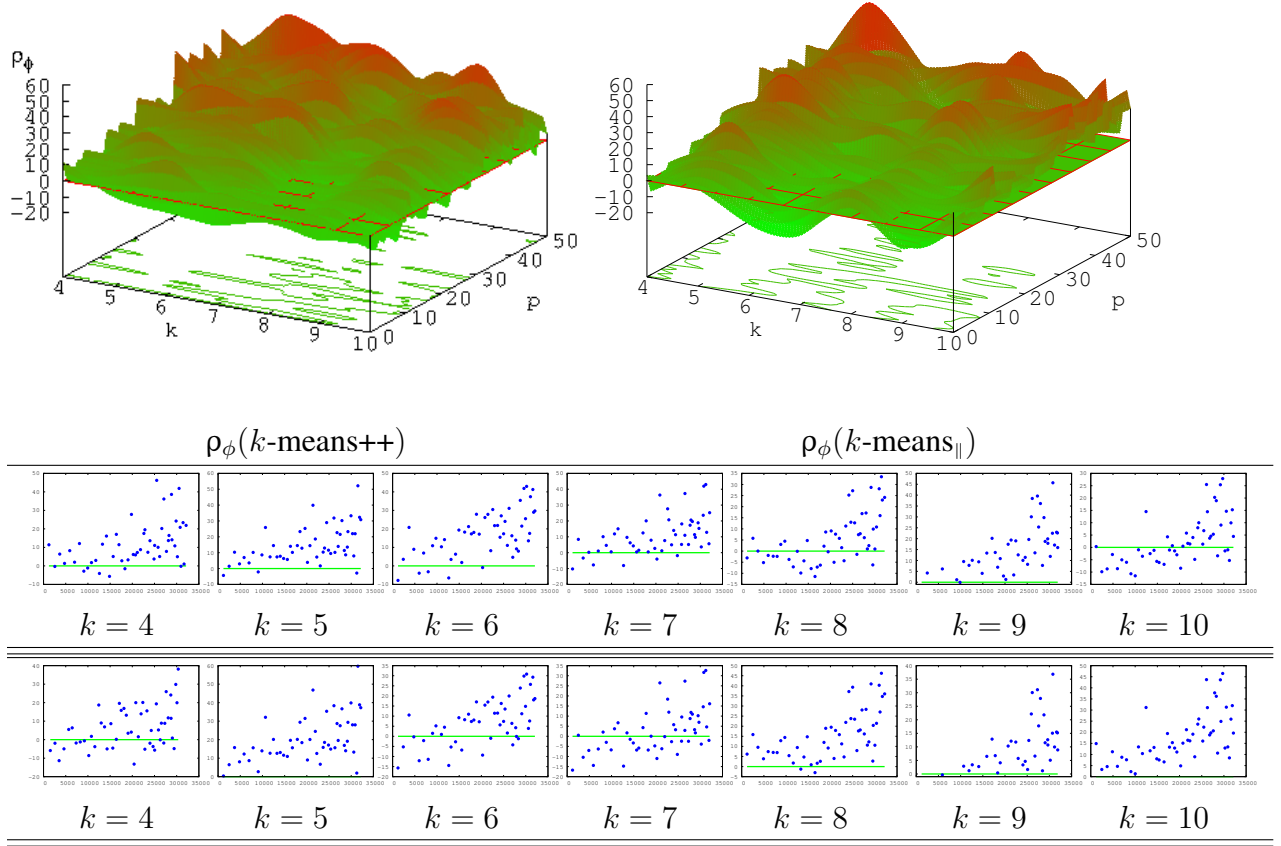


Figure 11: Experiments on real world data "EuropeDiff" with $N = 30$ simulated peers. Top plot: Plots corresponding to Figure 3. Middle and bottom plot ranges: plots corresponding respectively to Figures 8 and 9.

sampled N points ("peer centers") with k -means++ seeding in data and then aggregated for each peer the subset of data in the corresponding Voronoi 1-NN cell. We then simulate the spread for parameter p as in the simulated data. Figures 11 and 12 globally display (and confirm) the same trends as for the simulated data. They, however, clearly emphasize this time that the spread of Forgy nodes ϕ_s^F is one key parameter that drives the performances of Dk -means++. Notice also that Dk -means++ remains on this dataset competitive up to $p \geq 30\%$, which means that it remains competitive when a significant proportion of peers' data is scattered without any constraint.

To further address the way the spread of Forgy nodes affects results, we have used another real world data with highly non-uniform distribution, Mopsi-Finland locations⁵ ($m = 13467, d = 2$). We have sampled peers using two different schemes for the peer centers: k -means++ and Forgy. In this latter initialisation, we just pick peer centers at random. In the former k -means++ initialisation, the initial peer centers are much more evenly geographically spread before we complete the peers data with the closest points. They remain more spread after the $p\%$ uniform displacement of data between peers, as shown on the top plots of Figure 13. What is interesting about this data is that it displays that if peers' data are indeed geographically located, then Dk -means++ is competitive up to quite reasonable values of $p \leq 20\%$ (depending on k). That, is Dk -means++ works well

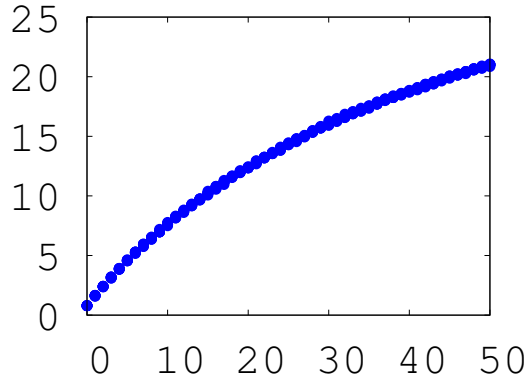


Figure 12: Experiments on real world data "EuropeDiff" with $N = 40$ simulated peers. Plot corresponding to Figure 10.

when each peer aggregates 80 % data which is reasonably "localized in the domain" and 20 % data which can be located everywhere in the domain.

Experiments with k -variates++ and GUPT

Among the state-of-the-art approaches against which we could compare k -variates++, there are two major contenders, PINQ (McSherry, 2010) and GUPT (Mohan et al., 2012). Even when PINQ is a broad system, we switched our preferences to GUPT for the following reasons. The performance of k -means based on PINQ relies on two principal factors: the initialisation (like in the non differentially private version) and the number of iterations. To compete against heavily tuned specific applications, like k -variates++, this scheme requires substantial work for its optimisation. For example, if one allocates part of the privacy budget to release a differential private initialisation, the noise has to be proportional to the domain width, which would release poor centers. Also, generating points uniformly at random from the domain, to obtain data-independent initial centers, yields to a poor initialisation. Finally, the number of iterations has to be tuned very carefully: if too small, the algorithm keeps poor solutions; if too large, the number of iteration increase the added noise for privacy and harms PINQ's final accuracy. We thus chose GUPT. k -means implemented in the GUPT proceeds the following way: the dataset is cut in a certain number of blocks ℓ (following (Mohan et al., 2012), we fix $\ell = m^{0.4}$ in our experiments), the usual k -means algorithm is performed on each block. Before releasing the final centroids, results are aggregated and a noise is applied. Finally, we also compare against the vanilla approach of *Forgy Initialisation* using the Laplace mechanism. The noise rate (*i.e.*, standard deviation) is then proportional to $\propto kR/\epsilon$ (we do not run k -means afterwards, hence the privacy budget remains "small"). In comparison, GUPT adds noise $\propto kR/(\ell\epsilon)$ at the end of this aggregation process. Note that we disregard the fact that our data are multidimensional, which should require a finer-grained tuning of ℓ , and choose to rely on the $\ell = m^{0.4}$ suggestion from (Mohan et al., 2012).

↔ **Comparison on real world domains** Our domains consist of 3 real-world datasets⁵. Lifesci contains the value of the top 10 principal components for a chemistry or biology experiment. Image

Dataset	m	d	k	$\tilde{\epsilon}$	$\rho'_\phi(\text{F-DP})$	$\rho'_\phi(\text{GUPT})$
LifeSci	26733	10	2	8.5	311	1.6
			3	4.4	172	0.4
			4	0.6	6	0.02
Image	34112	3	2	12.6	300	4.8
			3	3.2	77	0.9
EuropeDiff	169308	2	2	19.0	1200	46.1
			3	21.0	3120	66.5
			4	18.0	3750	55.0
			5	14.0	4000	51.0
			6	10.4	5000	36.0
			7	6.6	2600	26.0
			8	1.8	350	2.0

Table 10: Comparison of k -variates++, Forgy Initialisation differentially private (F-DP) and GUPT on the real world domains. On each domain, we compute ratio ρ'_ϕ of the clustering potential of the contender to that of k -variates++, a value > 1 indicating that k -variates++ is better. The potential of each algorithm has been averaged over 30 runs. $\tilde{\epsilon}$ is given in eq. (18).

is a 3D dataset with RGB vectors, and finally EuropeDiff is the differential coordinates of Europe map.

Table 10 presents the extensive results obtained, that are averaged in the paper’s body. We have fixed $\epsilon = 1$ in the differentially privacy parameters. The column $\tilde{\epsilon}$ (eq. (18)) provides the differential privacy parameter which is equivalent from the protection standpoint, but exploits the computation of δ_w, δ_s (which we compute exactly, and not in a randomized way like in the experiments on Theorem 12 above) and ineq. (80). Therefore, each time $\tilde{\epsilon} > \epsilon (=1$ in our applications), it means that our analysis brings a sizeable advantage over “raw protection” by Laplace mechanism (in our application we chose for $p_{\mu_\alpha, \theta_\alpha}$ a Laplace distribution). R is computed from the data by an upperbound of the smallest enclosing ball radius. The results display several interesting patterns. First, the largest the domain, the better we compare with respect to the other algorithms. On EuropeDiff for example, we often have the ratio of the potentials $\phi(\text{GUPT})/\phi(k\text{-variates++})$ of the order of *dozens*. Also, the performances of k -variates++ degrade if k increases, which is again consistent with the “good” regime of Theorem 10.

↔ **Comparison on synthetic domains** The synthetic datasets contain points uniformly sampled on a unit d -ball, in low dimension $d = 2$ and higher dimension $d = 15$, we generated datasets with size in $\{10^5, 10^6\}$.

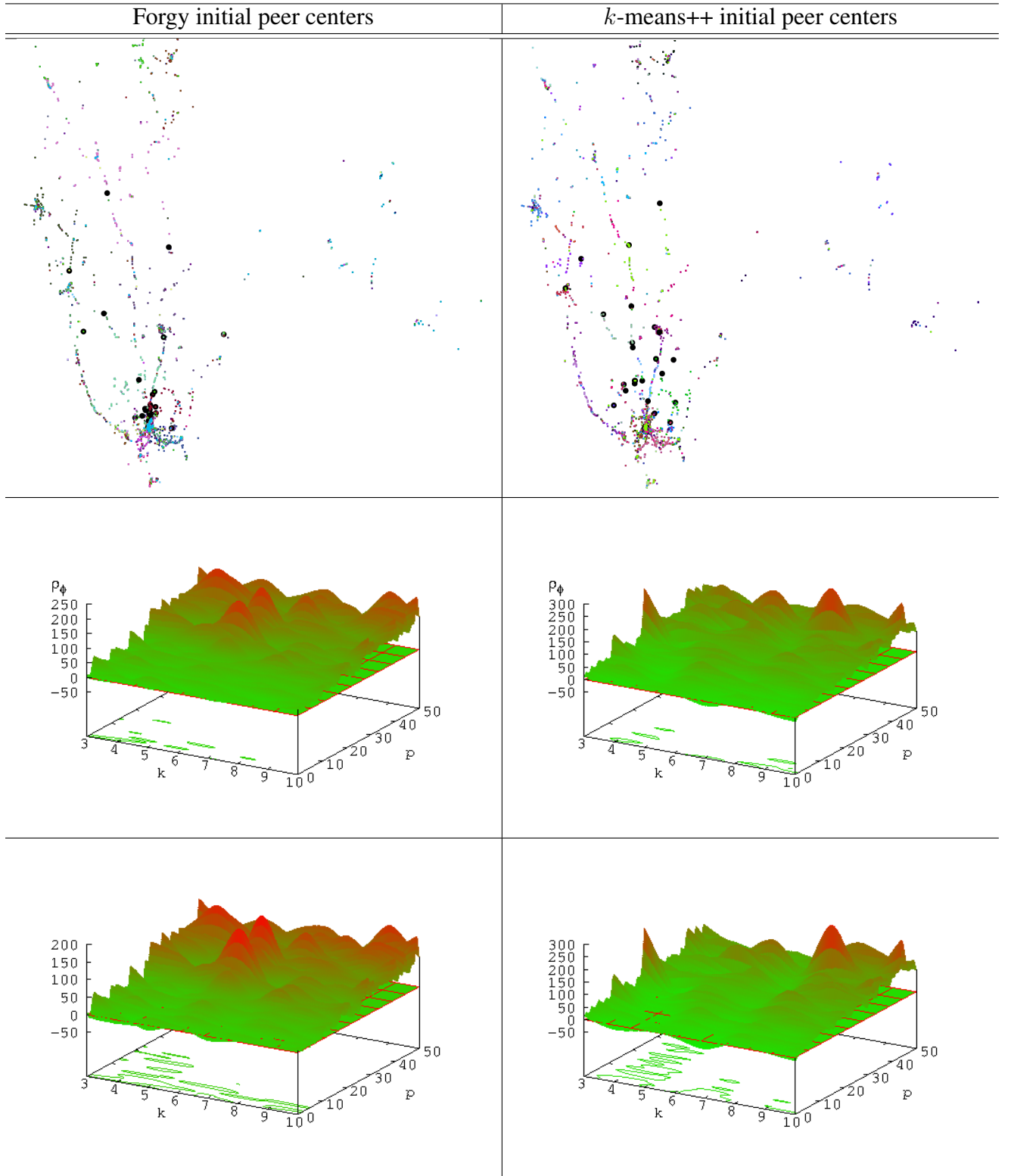


Figure 13: Mopsi-Finland locations data — Top: peer centers (big black dots) after $p = 50\%$ moving probability changes in data. Remark from the right plot (k -means++ initial peer centers) that peer data are less "attracted" towards the highest density regions. Center: plots of $\rho_\phi(k\text{-means}++)$. Bottom: plots of $\rho_\phi(k\text{-means}_\parallel)$.

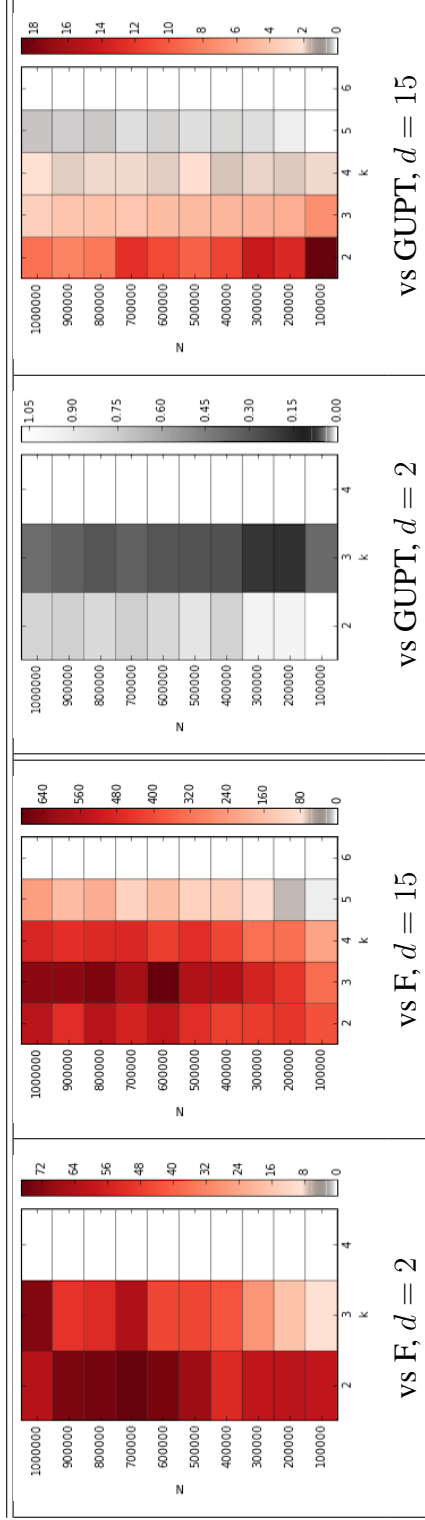


Figure 14: k -variates++ vs Forge initialization differentially private and GUPT. We use ratio ρ'_ϕ between the potential of the contender in (F-DP, GUPT) over the potential of k -variates++ (potentials are averaged 30 times). The more red, the better is k -variates++ with respect to the contender. Grey values indicate less positive outcomes for k -variates++; white values indicate that k -variates++ does not manage to find an ϵ' larger than ϵ , and thus does not manage to put smaller noise rate than in the Laplace mechanism.