
Coresets for Near-Convex Functions

Murad Tukan
muradtuk@gmail.com

Alaa Maalouf
alaamalouf12@gmail.com

Dan Feldman
dannyf.post@gmail.com

The Robotics and Big Data Lab,
Department of Computer Science,
University of Haifa,
Haifa, Israel

Abstract

Coreset is usually a small weighted subset of n input points in \mathbb{R}^d , that provably approximates their loss function for a given set of queries (models, classifiers, etc.). Coresets become increasingly common in machine learning since existing heuristics or inefficient algorithms may be improved by running them possibly many times on the small coreset that can be maintained for streaming distributed data. Coresets can be obtained by sensitivity (importance) sampling, where its size is proportional to the total sum of sensitivities. Unfortunately, computing the sensitivity of each point is problem dependent and may be harder to compute than the original optimization problem at hand. We suggest a generic framework for computing sensitivities (and thus coresets) for wide family of loss functions which we call near-convex functions. This is by suggesting the f -SVD factorization that generalizes the SVD factorization of matrices to functions. Example applications include coresets that are either new or significantly improves previous results, such as SVM, Logistic regression, M-estimators, and ℓ_z -regression. Experimental results and open source are also provided.

1 Introduction

In common machine learning problems, we are given a set of input points $P \subseteq \mathbb{R}^d$ (training data), and a loss function $f : P \times \mathbb{R}^d \rightarrow [0, \infty)$, where the goal is to solve the optimization problem of finding a query (model, classifiers, centers) x^* that minimizes the sum of fitting errors $\sum_{p \in P} f(p, x)$ over every query x in a given (usually infinite) set. For example, in k -median (or k -mean) clustering, each query is a set of k centers and the loss function is the distance (or squared distance) of a point to its nearest center. In linear regression or SVM, every input point includes a label, and the loss function is the fitting error between the classification of p via a given query to the actual label of p . Empirical risk minimization (ERM) may be used to generalize the result from train to test data.

Modern machine learning. In practice, many of these optimization or learning problems are usually hard even to approximate. Instead, practical heuristics with no provable guarantees may be used to solve them. Even for well understood problems, which have close optimal solution, such as linear regression or classes of convex optimization, in the era of big data we may wish to maintain the solution in other computation models such as: streaming input data (“on-the-fly”) that provably uses small memory, parallel computations on distributed data (on the cloud, network or GPUs) as well as deletion of points, constrained optimization (e.g. sparse classifiers). Cross validation [34] or hyper-parameter tuning techniques such as AutoML [30, 32] need to evaluate many queries for different subsets of the data, and different constraints.

Coresets. One approach is to redesign existing machine learning algorithms for faster, approximate solutions and these new computation models. A different approach that is to use data summarization techniques. *Coresets* in particular were first used to solve problems in computational geometry [1] and 34th Conference on Neural Information Processing Systems (NeurIPS 2020), Vancouver, Canada.

got increasing attention [3, 4, 5, 6, 8, 20, 26, 27, 38] over the recent years; see surveys in [23, 47, 50]. Informally, coreset is a small weighted subset of the input points (unlike e.g. sketches, or dimension-reduction techniques) that approximates the loss of the input set P for *every* feasible query x , up to a provable bound of $1 \pm \varepsilon$ for a given error parameter $\varepsilon \in (0, 1)$. The size of the coreset is usually polynomial in $1/\varepsilon$ but independent or near-logarithmic in the size of the input. Since such a coreset approximates every query (and not just the optimal one), it supports constraint optimization, and the above computation models using merge-and-reduce trees; see details in [23]. Moreover, coresets may be computed in time that is near-linear in the input, even for NP-hard optimization problems. Existing heuristic or inefficient algorithms may then be applied many times on the small coreset to obtain improved or faster models in such cases.

Example coresets in machine learning include SVM [33, 57, 58, 59, 60], ℓ_z -regression [18, 21, 54], clustering [2, 16, 24, 31, 37, 42, 53], logistic regression [35, 47], LMS solvers and SVD [28, 44, 45, 52], where all of these works have been dedicated to suggest a coreset for a specific problem.

A generic framework for constructing coresets was suggested in [25, 40]. It states that, with high probability, non-uniform sampling from the input set yields a coreset. Each point should be sampled i.i.d. with a probability that is proportional to its importance or sensitivity, and assigned a multiplicative weight which is inverse proportional to this probability, so that the expected original sum of losses over all the points will be preserved. Here, the sensitivity of an input point $p \in P$ is defined to be the maximum of its relative fitting loss $s(p) = f(p, x) / \sum_{q \in P} f(q, x)$ over every possible query x . The size of the coreset is near-linear in the total (sum) t of these sensitivities; see Theorem 3 for details. It turns out in the recent years that many classical and hard machine learning problems [7, 43, 55] have total sensitivity that is near-logarithmic or independent of the input size $|P|$ which implies small coresets via sensitivity sampling.

Paper per problem. The main disadvantage of this framework is that the sensitivity $s(p)$, as defined above, is problem dependent: namely on the loss function f and the feasible set of queries. Moreover, maximizing $s(p) = f(p, x) / \sum_{q \in P} f(q, x)$ is equivalent to minimizing the inverse $\sum_{q \in P} f(q, x) / f(p, x)$. Unfortunately, minimizing the enumerator is usually the original optimization problem which motivated the coreset in the first place. The denominator may make the problem harder, in addition to the fact that now we need to solve this optimization problem for each and every input point in P . While approximations of the sensitivities usually suffice, sophisticated and different approximation techniques are frequently tailored in papers of recent machine learning conferences for each and every problem.

1.1 Problem Statement

To this end, the goal of this paper is to suggest a framework for sensitivity bounding of a *family* of functions, and not for a specific optimization problem. This approach is inspired by convex optimization: while we do not have a single algorithm to solve any convex optimization, we do have generic solutions for family of convex functions. E.g., linear programming, Semi-definite programming, and so on.

We choose the following family of near-convex loss functions, with example supervised and unsupervised applications that include support vector machines, logistic regression, ℓ_z -regression for any $z \in (0, \infty)$, and functions that are robust to outliers. In the Supplementary Material we suggest a more generalized version that handles a bigger family of functions; see Definition 13, and hope that this paper will inspire the research of more and larger families.

Definition 1 (Near-convex functions). *Let $P \subseteq \mathbb{R}^d$ be a set of n points, and let $f : P \times \mathbb{R}^d \rightarrow [0, \infty)$ be a loss function. We call f a near-convex function if there are a convex loss function $g : P \times \mathbb{R}^d \rightarrow [0, \infty)$ (see Definition 12 at Supplementary Material), a function $h : P \times \mathbb{R}^d \rightarrow [0, \infty)$, and a scalar $z > 0$ satisfying:*

(i) *There exist $c_1, c_2 > 0$ such that for every $p \in P$, and $x \in \mathbb{R}^d$,*

$$c_1 (g(p, x)^z + h(p, x)^z) \leq f(p, x) \leq c_2 (g(p, x)^z + h(p, x)^z).$$

(ii) *For every $p \in P$, $x \in \mathbb{R}^d$ and $b > 0$, we have $g(p, bx) = b \cdot g(p, x)$.*

(iii) *For every $p \in P$ and $x \in \mathbb{R}^d$, we have $\frac{h(p, x)^z}{\sum_{q \in P} h(q, x)^z} \leq \frac{2}{n}$.*

(iv) The set $\mathcal{X}_g = \left\{ x \in \mathbb{R}^d \mid \sum_{p \in P} g(p, x)^{\max\{1, z\}} \leq 1 \right\}$ is centrally symmetric, i.e., for every $x \in \mathcal{X}_g$ we have $-x \in \mathcal{X}_g$, and there exist $R, r \in (0, \infty)$ such that $B(0_d, r) \subset \mathcal{X}_g \subset B(0_d, R)$, where $B(0_d, y)$ denotes a ball of radius $y > 0$, centered at 0_d .

We denote by \mathcal{F} , the union of all functions f with the above properties.

The intuition behind Definition 1. Properties (i)-(iii) are used to reduce the problem to dealing with a “simpler” pair of functions where the first is a convex function “g” that is linear in its argument x and the second function “h” being independent of the input points. Property (iv) ensures that the ellipsoid which encloses the level set of g (the convex function) exists and is centered at the origin to avoid dealing with the center. By combining the properties associated with the level set of g (the convex function) and Properties (i)-(iv), we manage to bound the loss function from above and below by the mahalanobis distance with respect to the enclosing ellipsoid. This is due to the fact that the level set encloses a contracted version of the ellipsoid which encloses the level set of g .

We are interested in a generic algorithm that would get a set of input points, and a loss function as above, and compute a sensitivity for each point, based on the parameters of the given loss function. In addition, we wish to use worst-case analysis and prove that for every input the total sensitivity (and thus size of coreset) would be small, depending on the “hardness” of the loss function that is encapsulated in the above parameters z, R , etc.

2 Related Work

Logistic Regression. A coreset construction algorithms for the problem of logistic regression were suggested by [35], [56], and [47]. All of these works handled variations of the problem, e.g., they all lack the incorporation of the bias term (intercept) in their loss function. Specifically speaking, both [35] and [47] didn’t account for the regularization term and its parameter. Furthermore, the coreset’s size established by [47], was dependant on the structure of the input data. As for [56], the coreset only succeed for a small subset of queries (a ball in \mathbb{R}^d of radius r , where the coreset’s size is near linear in r). Contrary to previous works, our coreset approximates the logistic regression loss function including the bias parameter (intercept) and the regularization term for every possible query. This is the loss function that is usually used in practice, e.g., see Sklearn library in [49]. Finally, our coreset’s size is independent of the structure of the data.

SVM. [11, 57, 58] addressed the problem of coreset construction for SVM, yet they used squared hinge loss to enforce the SVM cost function to be strongly convex. At [60], the coreset is constructed with respect to the hinge loss which most used form of SVM in practice (see Sklearn library at [49]). However for the coreset to be constructed, a (sub-)optimal solution was required for the problem itself. In addition, the coreset size depended heavily on on the ratio between the variance of each class of points. In this paper, we also address a coreset with respect to the hinge loss, yet we don’t require any (sub-)optimal solution to construct the coreset, and our coreset’s size depends on the ratio between the number of points of each class (see Corollary 9).

ℓ_z -Regression. A notable line of work [10, 18, 21, 54, 65] addressed the construction of coresets and sketches in this area, however, all such papers addressed the case of $z \geq 1$. Most of these works used tools similar to the well-conditioned basis which was first suggested at [21] to compute such coresets. Intuitively it can be thought of as a generalization of the *SVD* factorization of an input set with respect to the loss function of ℓ_z -regression for any $z \geq 1$. In our framework we generalize this factorization in order to compute coresets for the near-convex functions. To our knowledge, we suggest the first coreset for the problem of ℓ_z -regression for any $z \in (0, 1)$.

Outlier resistant functions (similar to M -estimators). Among such functions, is the ℓ_z -regression for any $z \in (0, 1]$ that is mentioned above, *Huber* loss function [13], *Tukey* loss functions [12], and many more [14]. However, to our knowledge, we present the first coreset for the problem formulation which is given at Corollary 10.

3 Our contribution

In this paper, we suggest an ε -coreset construction algorithm with respect to any near-convex function. Specifically speaking, we provide:

(i) A generalization of the well conditioned bases of [21] to a broader family of functions, i.e., not just for ℓ_z -Regression problems where $z \geq 1$. This informally describes a factorization of the input data with respect to a given near-convex loss function. We call such factorization the f -SVD of P (see Definition 4).

(ii) A framework for bounding the sensitivity of each point in an input set with respect to any near-convex function. The heart of the framework relies on computing the f -SVD factorization described in (i); see Lemma 5 and Algorithm 1.

(iii) By (ii), we provide the first ε -coreset for the problem of ℓ_z -regression where $z \in (0, 1)$, and the first ε -coreset for certain outlier resistant functions. We also unify existing works of coreset construction for the problems of logistic regression and SVM; see Section 6.

(iv) Experimental results on real-world and synthetic datasets for common machine learning solvers (supported by our framework) of Scikit-learn library [49], assessing the practicability and efficacy of our algorithm.

(v) An open source code implementation of our algorithm, for reproducing our results and future research [61].

3.1 Novelty

f -SVD factorization. In this work, we suggest a novel factorization technique of an input dataset with respect to a specific loss function f , we call it the f -SVD factorization. Roughly speaking, the heart of the f -SVD factorization lies in finding a diagonal matrix $D \in [0, \infty)^{d \times d}$ and an orthogonal matrix $V \in \mathbb{R}^{d \times d}$ such that the total loss $\sum_{p \in P} f(p, x)$ for any query $x \in \mathbb{R}^d$ can be bounded from above by $\sqrt{d} \|DV^T x\|_2$ and from below by $\|DV^T x\|_2$. In some sense, this can be thought of as a $(1 - 1/\sqrt{d})$ -coreset (or a sketch) since it approximates the total loss for any query in \mathbb{R}^d up to a multiplicative factor of $(1 - 1/\sqrt{d})$. In order to obtain such factorization, we forge a link between the Löwner ellipsoid [36] and the properties of near-convex functions; see Fig. 1 for a detailed illustrative explanation, Definition 4 and Lemma 16 for the formal details.

Note that SVD factorization is a special case of f -SVD due to that fact that SVD handles functions of the form $\sqrt{\sum_{p \in P} |p^T x|^2}$ and attempts to achieve the same purpose. The f -SVD factorization is a generalization of the *well-conditioned bases* of [21].

From f -SVD to sensitivity bounds. With the lower bound on the total loss that is guaranteed by the f -SVD, we show how to bound the sensitivity of each point in the dataset. On the other hand, the upper bound on the total loss provided by the f -SVD factorization, helps us in bounding the total sensitivity. Having this being said, we use the f -SVD factorization to suggest a sensitivity bounding framework for a set of points with respect to any near-convex function $f \in \mathcal{F}$; see Lemma 5.

4 Preliminaries

Notations. For integers $n, d \geq 2$, we denote by 0_d the origin of \mathbb{R}^d , and by $[n]$ the set $\{1, \dots, n\}$. The set $\mathbb{R}^{n \times d}$ denotes the union over every $n \times d$ real matrix, and $I_d \in \mathbb{R}^{d \times d}$ denotes the identity matrix. We say that a matrix $A \in \mathbb{R}^{d \times d}$ is *orthogonal* if and only if $A^T A = A A^T = I_d$. Finally, throughout the paper, vectors are addressed as column vectors, and $\mathcal{W} : \mathbb{R}^d \rightarrow 1$ is a weight function.

In what follows, we provide formally the notion of ε -coreset in our context.

Definition 2 (ε -coreset). *Let $P \subseteq \mathbb{R}^d$ be a set of n points, $f : P \times \mathbb{R}^d \rightarrow [0, \infty)$ be a near-convex function, and let $\varepsilon \in (0, 1)$. An ε -coreset for P with respect to f , is a pair (S, v) where $S \subseteq P$, $v : S \rightarrow (0, \infty)$ is a weight function, such that for every $x \in \mathbb{R}^d$, $\left| 1 - \frac{\sum_{q \in S} v(q) f(q, x)}{\sum_{p \in P} f(p, x)} \right| \leq \varepsilon$.*

The following theorem formally describes how to construct an ε -coreset via the sensitivity framework.

Theorem 3 (Restatement of Theorem 5.5 in [7]). *Let $P \subseteq \mathbb{R}^d$ be a set of n points, and let $f : P \times \mathbb{R}^d \rightarrow [0, \infty)$ be a loss function. For every $p \in P$ define the sensitivity of p as $\sup_{x \in \mathbb{R}^d} \frac{f(p, x)}{\sum_{q \in P} f(q, x)}$, where the sup is over every $x \in \mathbb{R}^d$ such that the denominator is non-zero. Let $s : P \rightarrow [0, 1]$ be*

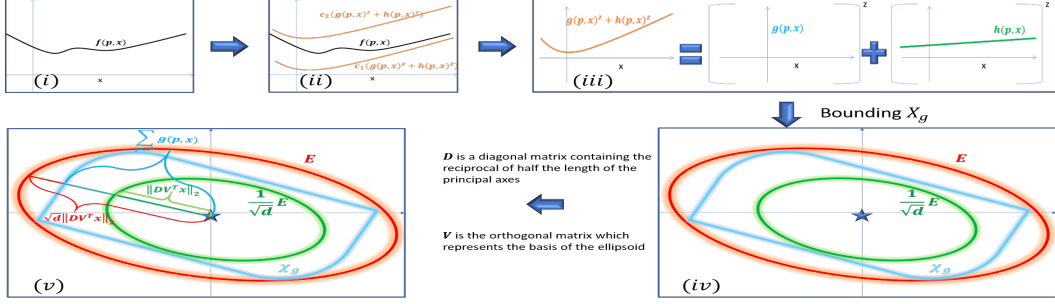


Figure 1: **How to compute f -SVD:** (i) Given a set $P \subseteq \mathbb{R}^2$, and a function $f : P \times \mathbb{R}^2 \rightarrow [0, \infty)$, (ii) find a function which can bound $f(p, \cdot) \times \mathbb{R}^2 \rightarrow [0, \infty)$ from above and below for every $p \in P$, (iii) decompose this function into two functions g and h where for every $p, q \in P$ and $x \in \mathbb{R}^2$, $g(p, \cdot)$ is a convex function (e.g., $g(p, x) = |p^T x|^4$), and $h(p, x) \approx h(q, x)$ (e.g., $h(p, x) = \|x\|_1 + 10$), here $z = 1$. (iv) Since g is convex, we find the Löwner ellipsoid E which contains $\mathcal{X}_g = \{x \in \mathbb{R}^2 \mid \sum_{p \in P} g(p, x) \leq 1\}$, and the contracted ellipsoid $1/\sqrt{d}E$ is inscribed in \mathcal{X}_g . Following this, we compute a diagonal matrix $D \in [0, \infty)^{2 \times 2}$ and an orthogonal matrix $V \in \mathbb{R}^{2 \times 2}$ such that $E = \{x \in \mathbb{R}^2 \mid \|DV^T x\|_2 \leq 1\}$. (v) By properties of the Löwner ellipsoid, we show that the total loss of g (cyan line) for any query $x \in \mathbb{R}^2$ is in the range $[\|DV^T x\|_2, \sqrt{d}\|DV^T x\|_2]$ (green and red lines). When combined with the bounds on f , we obtain an upper bound on the sensitivity of each point in P and on the total sensitivity.

a function such that $s(p)$ is an upper bound on the sensitivity of p . Let $t = \sum_{p \in P} s(p)$ and d' be the VC dimension of the quadruple $(P, \mathbf{1}, \mathbb{R}^d, f)$; see Definition 15. Let $c \geq 1$ be a sufficiently large constant, $\varepsilon, \delta \in (0, 1)$, and let S be a random sample of $|S| \geq \frac{ct}{\varepsilon^2} (d' \log t + \log \frac{1}{\delta})$ i.i.d points from P , such that every $p \in P$ is sampled with probability $s(p)/t$. Let $v(p) = \frac{t}{s(p)|S|}$ for every $p \in S$. Then, with probability at least $1 - \delta$, (S, v) is an ε -coreset for P with respect to f .

5 Coreset for near-convex loss functions

For brevity purposes, proofs of the technical results have been omitted from this manuscript; we refer the reader to the supplementary material for the proofs. In addition, for simplicity of notation, we assume that the weight of each point in the input set is 1, while in the supplementary material, we handle the general case where each point may have any nonnegative weight. We also discuss generalized versions of Definition 1 and Definition 4.

5.1 Bounding the sensitivity

The following provides the generalization of the well-conditioned basis of [21], which will be used to bound the sensitivities.

Definition 4 (f -SVD of P). Let $P \subseteq \mathbb{R}^d$ be a set of n points, $f \in \mathcal{F}$ be a near-convex loss function (see Definition 1), and let g, h, c_1, z be defined as in the context of Definition 1 with respect to f . Let $D \in [0, \infty)^{d \times d}$ be a diagonal matrix, and let $V \in \mathbb{R}^{d \times d}$ be an orthogonal matrix, such that for every $x \in \mathbb{R}^d$, $c_1 \left((\|DV^T x\|_2)^z + \sum_{p \in P} h(p, x)^z \right) \leq \sum_{p \in P} f(p, x)$, and let $\alpha \in \Theta(\sqrt{d})$ such that for every $x \in \mathbb{R}^d$, $\sum_{p \in P} g(p, x)^{\max\{1, z\}} \leq (\alpha \|DV^T x\|_2)^{\max\{1, z\}}$. Define $U : P \rightarrow \mathbb{R}^d$ such that $U(p) = (VD)^{-1} p$ for every $p \in P$. The tuple (U, D, V) is the f -SVD of P .

Note that (i) such factorization exists for any set of points P and any near-convex loss function $f : P \times \mathbb{R}^d \rightarrow [0, \infty)$ satisfying Definition 1, and (ii) the matrix VD is invertible due to the fact that D is of full rank which is a result of Property (iv) of Definition 4. Both (i)-(ii) hold by using Löwner ellipsoid; see Fig. 1 for intuitive explanation, and Lemma 16 at the Supplementary Material for formal proof.

In what follows, we proceed to bound the sensitivity of each point and the total sensitivity, with respect to a loss function $f \in \mathcal{F}$. This is by using the f -SVD of P .

Lemma 5. Let $P \subseteq \mathbb{R}^d$ be a set of n points, and let $f \in \mathcal{F}$ be a near-convex loss function as in Definition 1. Let g, h, c_1, c_2, z be defined as in the context of Definition 1 with respect to f , (U, D, V)

be the f -SVD of P , and let $\alpha \in \Theta(\sqrt{d})$ which satisfies the conditions in Definition 4. Suppose that there exists a set $\{v_j\}_{j=1}^{O(d)} \subseteq \mathbb{R}^d$ of $O(d)$ unit vectors and $c > 0$, such that for every unit vector $y \in \mathbb{R}^d$ and $p \in P$, $g(p, (DV^T)^{-1}y)^z \leq c \sum_{j=1}^{O(d)} g(p, (DV^T)^{-1}v_j)^z$. Then, for every $p \in P$, the sensitivity of p is bounded by $s(p) \leq \frac{2c_2}{c_1 n} + \frac{cc_2}{c_1} \sum_{j=1}^{O(d)} \left(g(p, (DV^T)^{-1}v_j)\right)^z$, and the total sensitivity is bounded by $\sum_{p \in P} s(p) \in \frac{2c_2}{c_1} + \frac{cc_2}{c_1} \max\{n^{1-z}, 1\} \alpha^z O(d)$.

The existence of the set $\{v_j\}_{j=1}^{O(d)}$ is discussed in details at the supplementary material at Section D.

5.2 The coresets construction

Algorithm 1 receives as input, a set P of n points in \mathbb{R}^d , a loss function $f \in \mathcal{F}$ (see Definition 1), and a sample size $m > 0$. As Theorem 6 states, if the sample size m is sufficiently large, then Algorithm 1 outputs a pair (S, v) that is with high probability, an ε -coreset for P with respect to f .

First, we set d' to be VC dimension of the quadruple $(P, \mathbb{1}, \mathbb{R}^d, f)$; See Definition 15. The crux of our algorithm lies in generating the importance sampling distribution via efficiently computing upper bound on the sensitivity of each point (Lines 5–7). To do so, we compute the f -SVD of P at Lines 3–4, and we use it to bound the sensitivity of each $p \in P$ as stated in Lemma 5; see Line 6. Now we have all the needed ingredients to use Theorem 3 in order to obtain an ε -coreset, i.e., we sample i.i.d m points from P based on their sensitivity bounds (see Line 9), and assign a new weight for every sampled point at Line 10.

Algorithm 1: CORESET(P, f, m)

Input: A set $P \subseteq \mathbb{R}^d$ of n points, a near-convex loss function $f : P \times \mathbb{R}^d \rightarrow [0, \infty)$, and a sample size $m \geq 1$.

Output: A pair (S, v) that satisfies Theorem 6.

- 1 Set $d' :=$ the VC dimension of quadruple $(P, \mathbb{1}, \mathbb{R}^d, f)$ // See Definition 15
 - 2 Set g and $\{z, c_1, c_2\}$ to be a function and a set of real positive numbers respectively, satisfying Property (i) and (ii) of Definition 1 with respect to f
 - 3 Set $c > 0$ and $\{v_1, \dots, v_d\}$ to be positive scalar and a set of d unit vectors in \mathbb{R}^d respectively satisfying Lemma 5
 - 4 Set (U, D, V) to be the f -SVD of (P, w) // See Definition 1
 - 5 **for** every $p \in P$ **do**
 - 6 Set $s(p) := \frac{cc_2}{c_1} \sum_{j=1}^d g(p, (DV^T)^{-1}v_j)^z + \frac{2c_2}{c_1 n}$
 // the bound of the sensitivity of p as in Lemma 5
 - 7 Set $t := \sum_{p \in P} s(p)$
 - 8 Set $\tilde{c} \geq 1$ to be a sufficiently large constant // Can be determined from Theorem 6
 - 9 Pick an i.i.d sample S of m points from P , where each $p \in P$ is sampled with probability $\frac{s(p)}{t}$.
 - 10 set $v : \mathbb{R}^d \rightarrow [0, \infty]$ to be a weight function such that for every $q \in S$, $v(q) = \frac{t}{s(q) \cdot m}$.
 - 11 **return** (S, v)
-

Theorem 6. Let $P \subseteq \mathbb{R}^d$ be set of n points, and $f \in \mathcal{F}$ be a near-convex function. Let $R, r > 0$ be a pair of positive scalars as in Definition 1 with respect to f , and let c, c_1, c_2, α be defined as in the context of Lemma 5 with respect to f . Let $\varepsilon, \delta \in (0, 1)$ be an error parameter and a probability of failure respectively, and let d' be the VC dimension of the triplet (P, f, \mathbb{R}^d) . Let $t = \frac{2c_2}{c_1} + \frac{cc_2}{c_1} \max\{n^{1-z}, 1\} \alpha^z d$, $m \in O\left(\frac{t}{\varepsilon^2} (d' \log(t) + \log(\frac{1}{\delta}))\right)$, and let (S, v) be the output of a call to CORESET(P, f, m). Then, (i) with probability at least $1 - \delta$, (S, v) is an ε -coreset of size m for P with respect to f ; see Definition 2. (ii) The overall time for constructing (S, v) is bounded by $O\left(T(n, d) d^4 \log\left(\frac{R}{r}\right)\right)$, where $T(n, d)$ is a bound on the time it takes to compute a gradient of $\sum_{p \in P} f(p, x)$ with respect to any query $x \in \mathbb{R}^d$.

Poly-logarithmic coresets size. We provide an analysis that shows how to obtain a coresets of size poly-logarithmic in the input size n ; see Algorithm 2 and Lemma 17 at the Supplementary Material.

6 Applications

In what follows, we provide various applications for our framework, .e.g, SVM, Logistic Regression, ℓ_z for $z \in (0, 1)$, outlier resistant functions (similar to Tukey in behavior). For additional problems supported by our framework, we refer the reader to Section G at the Supplementary Material.

Table 1: Results: The table below presents the coresets size and the time needed for constructing it with respect to a specific set of problems, where the input is a set of n points in \mathbb{R}^d denoted by P . In the table, $\text{nnz}(P)$ denotes the total number of nonzero entries in the set P , \tilde{C} denotes the ratio between the number of positive and negative labeled points (in practice, it's a constant number), $\lambda = \sqrt{n}$ is the given regularization parameter for the problems, $\gamma \geq 1$ is defined as in Corollary 10, ε is the error parameter, and δ is the probability of failure.

Problem type	Coreset's size	Construction time ¹
Logistic regression	$O\left(\frac{d\sqrt{n}}{\varepsilon^2} \left(d \log(d\sqrt{n}) + \log\left(\frac{1}{\delta}\right)\right)\right)$	$O(nd^2)$
ℓ_z -Regression for $z \in (0, 1)$	$O\left(\frac{n^{1-z}d^{\frac{z}{2}+1}}{\varepsilon^2} \left(d \log(n^{1-z}d^{\frac{z}{2}+1}) + \log\left(\frac{1}{\delta}\right)\right)\right)$	$O(\text{nnz}(P) \log n + d^{O(1)})$
SVM	$O\left(\frac{d\sqrt{n} + \frac{\tilde{C}^2+1}{\tilde{C}}}{\varepsilon^2} \left(d \log\left(d\sqrt{n} + \frac{\tilde{C}^2+1}{\tilde{C}}\right) + \log\left(\frac{1}{\delta}\right)\right)\right)$	$O(nd^2)$
Restricted ℓ_z -regression	$O\left(\frac{\gamma d^{2+\left \frac{1}{2}-\frac{1}{z}\right }}{\varepsilon^2} \left(d \log\left(\gamma d^{2+\left \frac{1}{2}-\frac{1}{z}\right }\right) + \log\left(\frac{1}{\delta}\right)\right)\right)$	$O(\text{nnz}(P) \log n + d^{O(1)})$

Corollary 7 (Logistic Regression). *Let $P \subseteq \mathbb{R}^d$ be a set of n points such that for every $p \in P$, $\|p\|_2 \leq 1$, $y : P \rightarrow \{-1, 1\}$ be a labeling function, $\lambda \geq 1$ be a regularization parameter such that for every $p \in P$, $x \in \mathbb{R}^d$ and $b \in \mathbb{R}$, $f_{\text{LOG}}\left(p, \begin{bmatrix} x \\ b \end{bmatrix}\right) = \frac{1}{\lambda} \ln\left(1 + e^{p^T x + y(p) \cdot b}\right) + \frac{1}{2n} \|x\|_2^2$.*

Let $\varepsilon, \delta \in (0, 1)$ be an error parameter and a probability of failure respectively, $m \in O\left(\frac{dn}{\lambda \varepsilon^2} \left(d \log\left(\frac{dn}{\lambda}\right) + \log\left(\frac{1}{\delta}\right)\right)\right)$, and let (S, v) be the output of a call to CORESET (P, f_{LOG}, m) . Then, with probability at least $1 - \delta$, (S, v) is an ε -coreset (of size m) for P with respect to f_{LOG} .

Corollary 8 (ℓ_z -Regression where $z \in (0, 1)$). *Let $P \subseteq \mathbb{R}^d$ be a set of n points, $z \in (0, 1)$ and let $f_{\text{NCL}_z} : P \times \mathbb{R}^d$ be a loss function such that for every $x \in \mathbb{R}^d$, and $p \in P$, $f_{\text{NCL}_z}(p, x) = |p^T x|^z$.*

Let $\varepsilon, \delta \in (0, 1)$, $m \in O\left(\frac{n^{1-z}d^{\frac{z}{2}+1}}{\varepsilon^2} \left(d \log(n^{1-z}d^{\frac{z}{2}+1}) + \log\left(\frac{1}{\delta}\right)\right)\right)$, and let (S, v) be the output of a call to CORESET (P, f_{NCL_z}, m) . Then, with probability at least $1 - \delta$, (S, v) is an ε -coreset (of size m) for P with respect to f_{NCL_z} .

We now show how our framework can be used to compute an ε -coreset for some query spaces where the involved loss functions are not from the family \mathcal{F} . The coreset construction algorithms are hidden in the constructive proofs of the following corollaries.

Corollary 9 (Support Vector Machines). *Let $P \subseteq \mathbb{R}^d$ be a set of n points such that for every $p \in P$, $\|p\| \leq 1$. Let $y : P \rightarrow \{1, -1\}$ be a labelling function, $\lambda \geq 1$ be a regularization parameter such that for every $p \in P$, $x \in \mathbb{R}^d$, and $b \in \mathbb{R}$, $f_{\text{SVM}}\left(p, \begin{bmatrix} x \\ b \end{bmatrix}\right) = \lambda \max\{0, 1 - (p^T x + y(p) \cdot b)\} + \frac{1}{2n} \|x\|_2^2$. Let $P_+ = \{p | p \in P, y(p) = 1\}$, $P_- = P \setminus P_+$, $\tilde{C} = \frac{|P_+|}{|P_-|}$.*

Then, there exists an algorithm that gets the set P as an input, and returns a pair (S, v) , such that (i) with probability at least $1 - \delta$, (S, v) is an ε -coreset for P with respect to f_{SVM} , and (ii) the size of the coreset is $|S| \in O\left(\frac{1}{\varepsilon^2} \left(\frac{dn}{\lambda} + \frac{\tilde{C}^2+1}{\tilde{C}}\right) \left(d \log\left(\frac{dn}{\lambda} + \frac{\tilde{C}^2+1}{\tilde{C}}\right) + \log\left(\frac{1}{\delta}\right)\right)\right)$.

Corollary 10 (Outlier resistant functions). *Let $P \subseteq \mathbb{R}^d$ be a set of n points, and let $f_{\text{RES}_z} : P \times \mathbb{R}^d \rightarrow [0, \infty)$ be loss function such that for every $x \in \mathbb{R}^d$, and $p \in P$, $f_{\text{RES}_z}(p, x) = \min\{|p^T x|, \|x\|_z\}$.*

Then, there exists an algorithm that gets the set P as an input, and returns a pair (S, v) , such that (i) with probability at least $1 - \delta$, (S, v) is an ε -coreset for P with respect to f_{RES_z} , and (ii) the size of the coreset is $O\left(\frac{\gamma d^{2+\left|\frac{1}{2}-\frac{1}{z}\right|}}{\varepsilon^2} \left(d \log\left(\gamma d^{2+\left|\frac{1}{2}-\frac{1}{z}\right|}\right) + \log\left(\frac{1}{\delta}\right)\right)\right)$, where γ is defined in the proof.

¹Problems which are reduced to ℓ_z -regression problems for any $z \geq 1$, are easier to be dealt with in term of coreset construction time due to the existence of randomized algorithm of computing the Löwner ellipsoid by [15]; see Section H for detailed description.

7 Experimental Results

In what follows we evaluate our coreset against uniform sampling on real-world datasets, with respect to the SVM problem, Logistic regression problem and ℓ_z -regression problem for $z \in (0, 1)$. Additional details of our setup can be found at Section H of the Supplementary Material.

Software/Hardware. Our algorithms were implemented in Python 3.6 [63] using “Numpy” [48], “Scipy” [64] and “Scikit-learn” [49]. Tests were performed on 2.59GHz i7-6500U (2 cores total) machine with 16GB RAM.

Datasets. The following datasets were used for our experiments mostly from UCI machine learning repository [22]: (i) **HTRU** [22] — 17, 898 radio emissions of the Pulsar star each consisting of 9 features. (ii) **Skin** [22] — 245, 057 random samples of R,G,B from face images consisting of 4 dimensions. (iii) **Cod-rna** [62] — consists of 59, 535 samples, 8 features, which has two classes (i.e. labels), describing RNAs. (iv) **Web** dataset [9] – 49, 749 web pages records where each record is consists of 300 features. (v) **3D spatial networks** [22] – 3D road network with highly accurate elevation information (+20cm) from Denmark used in eco-routing and fuel/Co2-estimation routing algorithms consisting of 434, 874 records where each record has 4 features.

Evaluation against uniform sampling. At Fig. 2a–2f, we have chosen 20 sample sizes, starting from 50 till 500, at Figures 2g–2h, we have chosen 20 sample sizes starting from 4000 till 16, 000. At each sample size, we generate two coresets, where the first is using uniform sampling and the latter is using Algorithm 1. For each coreset (S, v) , we find $x^* \in \arg \min_{x \in \mathbb{R}^d} \sum_{p \in S} v(p) f(p, x)$, and the approximation error ε is set to be $(\sum_{p \in P} f(p, x^*)) / (\min_{x \in \mathbb{R}^d} \sum_{p \in P} f(p, x)) - 1$. The results were averaged across 40 trials, while the shaded regions correspond to the standard deviation.

Evaluation against prior work. We can not have a fair comparison between our coreset to prior coresets for Logistic regression[47, 56] due to the fact that our formulation of the problem is different. As for support vector machines, we compared our efficacy against [60], the same way that we have compared against uniform sampling. Although not in all cases our approach outperforms [60] in terms of relative error (i.e., ε), our approach is much faster than that of [60]; see Figure 3.

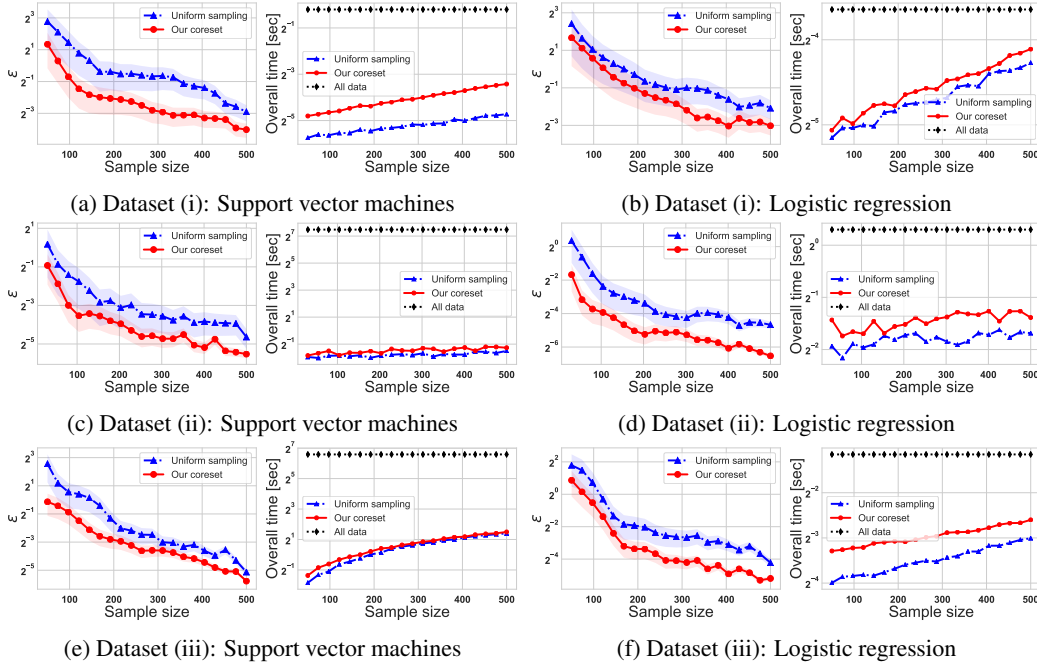


Figure 2: Experimental results

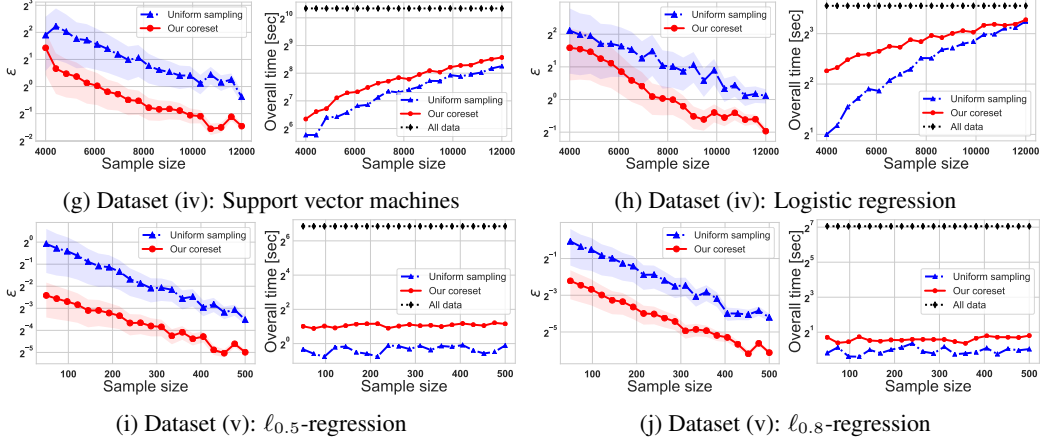


Figure 2: Experimental results

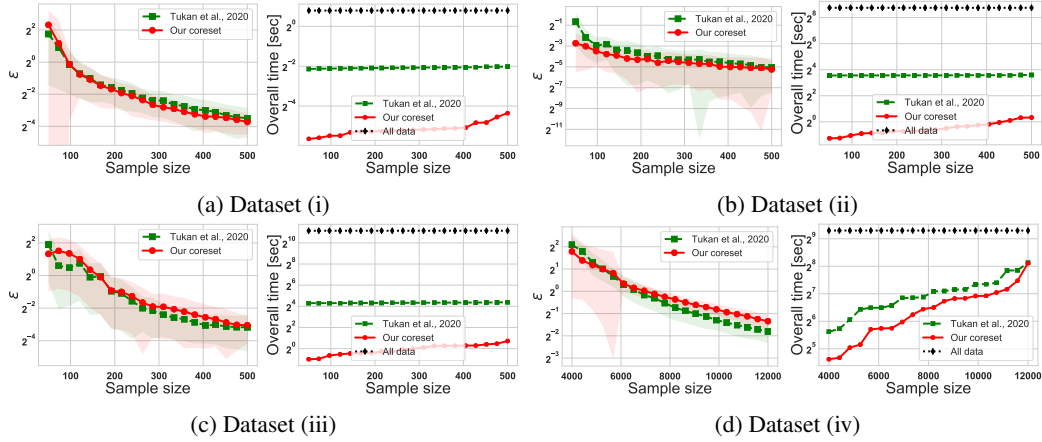


Figure 3: Comparison against prior work in the context of SVMs

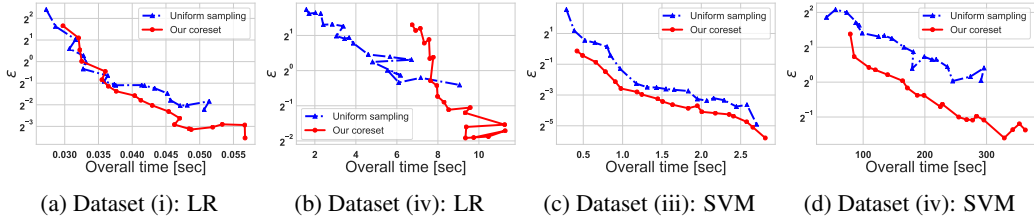


Figure 4: Relative error as a function of the total running time. Here, LR stands for Logistic regression.

8 Conclusions and open problems

In this paper, we have provided what we call the f -SVD of P with respect a given near-convex loss function $f \in \mathcal{F}$, as well as sensitivity bounding framework using the f -SVD. What interests us is to draw back forcing f to have a centrally symmetric level set as well as embedding the center of the Löwner ellipsoid into the sensitivity bound. This is crucial step for generalizing the framework towards a much broader family of functions, e.g., loglog-Lipschitz functions [29]. We are aware that for ℓ_z -regression problems where $z \geq 1$, Lewis weights have been used by [17] and are considered to be the state of the art coreset for these problems. We aim to generalize the applicability of Lewis weights and other sketching techniques towards different functions, and as far as we know, we consider the above issues to be open problems.

Broader Impact

Our work provides a strong theoretical result, where we have suggested a generic framework for bounding the sensitivity with respect to broad family of functions. Practically, this family imposes widely used applications such as *SVM*, *Logistic regression*, ℓ_z -Regression and more.

Although, Broader Impact discussion is not directly applicable, our work can be used to accelerate many known machine learning solvers under various settings such as distributed, streaming, etc.

References

- [1] P. Agarwal, S. Har-Peled, and K. Varadarajan. Approximating extent measures of points. *Journal of the ACM*, 51(4):606–635, 2004.
- [2] O. Bachem, M. Lucic, and A. Krause. Coresets for nonparametric estimation-the case of dp-means. In *ICML*, pages 209–217, 2015.
- [3] O. Bachem, M. Lucic, and A. Krause. Scalable k-means clustering via lightweight coresets. In *KDD’18 Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pages 1119–1127. ACM, 2018.
- [4] O. Bachem, M. Lucic, and S. Lattanzi. One-shot coresets: The case of k-clustering. In A. Storkey and F. Perez-Cruz, editors, *Proceedings of the Twenty-First International Conference on Artificial Intelligence and Statistics*, volume 84 of *Proceedings of Machine Learning Research*, pages 784–792, Playa Blanca, Lanzarote, Canary Islands, 09–11 Apr 2018. PMLR.
- [5] M. Bădoiu and K. L. Clarkson. Optimal core-sets for balls. *Computational Geometry*, 40(1):14–22, 2008.
- [6] M.-F. F. Balcan, S. Ehrlich, and Y. Liang. Distributed k -means and k -median clustering on general topologies. In *Advances in Neural Information Processing Systems*, pages 1995–2003, 2013.
- [7] V. Braverman, D. Feldman, and H. Lang. New frameworks for offline and streaming coreset constructions. *arXiv preprint arXiv:1612.00889*, 2016.
- [8] V. Braverman, S. H.-C. Jiang, R. Krauthgamer, and X. Wu. Coresets for ordered weighted clustering. In K. Chaudhuri and R. Salakhutdinov, editors, *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pages 744–753, Long Beach, California, USA, 09–15 Jun 2019. PMLR.
- [9] C.-C. Chang and C.-J. Lin. LIBSVM: A library for support vector machines. *ACM Transactions on Intelligent Systems and Technology*, 2:27:1–27:27, 2011. Software available at <http://www.csie.ntu.edu.tw/~cjlin/libsvm>.
- [10] K. L. Clarkson. Subgradient and sampling algorithms for l_1 regression. In *Proceedings of the sixteenth annual ACM-SIAM symposium on Discrete algorithms*, pages 257–266. Society for Industrial and Applied Mathematics, 2005.
- [11] K. L. Clarkson. Coresets, sparse greedy approximation, and the frank-wolfe algorithm. *ACM Transactions on Algorithms (TALG)*, 6(4):63, 2010.
- [12] K. L. Clarkson, R. Wang, and D. P. Woodruff. Dimensionality reduction for tukey regression. *arXiv preprint arXiv:1905.05376*, 2019.
- [13] K. L. Clarkson and D. P. Woodruff. Sketching for m -estimators: A unified approach to robust regression. In *Proceedings of the twenty-sixth annual ACM-SIAM symposium on Discrete algorithms*, pages 921–939. SIAM, 2014.
- [14] K. L. Clarkson and D. P. Woodruff. Input sparsity and hardness for robust subspace approximation. In *2015 IEEE 56th Annual Symposium on Foundations of Computer Science*, pages 310–329. IEEE, 2015.

- [15] K. L. Clarkson and D. P. Woodruff. Low-rank approximation and regression in input sparsity time. *Journal of the ACM (JACM)*, 63(6):1–45, 2017.
- [16] M. B. Cohen, S. Elder, C. Musco, C. Musco, and M. Persu. Dimensionality reduction for k-means clustering and low rank approximation. In *Proceedings of the forty-seventh annual ACM symposium on Theory of computing*, pages 163–172, 2015.
- [17] M. B. Cohen and R. Peng. Lp row sampling by lewis weights. In *Proceedings of the forty-seventh annual ACM symposium on Theory of computing*, pages 183–192. ACM, 2015.
- [18] M. B. Cohen and R. Peng. Lp row sampling by lewis weights. In *Proceedings of the forty-seventh annual ACM symposium on Theory of computing*, pages 183–192, 2015.
- [19] R. M. Corless, G. H. Gonnet, D. E. Hare, D. J. Jeffrey, and D. E. Knuth. On the lambertw function. *Advances in Computational mathematics*, 5(1):329–359, 1996.
- [20] R. Curtain, S. Im, B. Moseley, K. Pruhs, and A. Samadian. On coresets for regularized loss minimization. *arXiv preprint arXiv:1905.10845*, 2019.
- [21] A. Dasgupta, P. Drineas, B. Harb, R. Kumar, and M. W. Mahoney. Sampling algorithms and coresets for ℓ_p regression. *SIAM Journal on Computing*, 38(5):2060–2078, 2009.
- [22] D. Dua and C. Graff. UCI machine learning repository, 2017.
- [23] D. Feldman. Core-sets: Updated survey. In *Sampling Techniques for Supervised or Unsupervised Tasks*, pages 23–44. Springer, 2020.
- [24] D. Feldman, M. Faulkner, and A. Krause. Scalable training of mixture models via coresets. In *Advances in neural information processing systems*, pages 2142–2150, 2011.
- [25] D. Feldman and M. Langberg. A unified framework for approximating and clustering data. In *Proceedings of the forty-third annual ACM symposium on Theory of computing*, pages 569–578. ACM, 2011.
- [26] D. Feldman, M. Monemizadeh, C. Sohler, and D. P. Woodruff. Coresets and sketches for high dimensional subspace approximation problems. In *Proceedings of the twenty-first annual ACM-SIAM symposium on Discrete Algorithms*, pages 630–649. SIAM, 2010.
- [27] D. Feldman, G. Rossman, M. Volkov, and D. Rus. Coresets for k-segmentation of streaming data. In *NIPS*, 2014.
- [28] D. Feldman, M. Schmidt, and C. Sohler. Turning big data into tiny data: Constant-size coresets for k-means, pca and projective clustering. In *Proceedings of the twenty-fourth annual ACM-SIAM symposium on Discrete algorithms*, pages 1434–1453. SIAM, 2013.
- [29] D. Feldman and L. J. Schulman. Data reduction for weighted and outlier-resistant clustering. In *Proceedings of the twenty-third annual ACM-SIAM symposium on Discrete Algorithms*, pages 1343–1354. SIAM, 2012.
- [30] M. Feurer, A. Klein, K. Eggensperger, J. Springenberg, M. Blum, and F. Hutter. Efficient and robust automated machine learning. In *Advances in neural information processing systems*, pages 2962–2970, 2015.
- [31] L. Gu. A coreset-based semi-supervised clustering using one-class support vector machines. In *Control Engineering and Communication Technology (ICCECT), 2012 International Conference on*, pages 52–55. IEEE, 2012.
- [32] I. Guyon, K. Bennett, G. Cawley, H. J. Escalante, S. Escalera, T. K. Ho, N. Macià, B. Ray, M. Saeed, A. Statnikov, et al. Design of the 2015 chlearn automl challenge. In *2015 International Joint Conference on Neural Networks (IJCNN)*, pages 1–8. IEEE, 2015.
- [33] S. Har-Peled, D. Roth, and D. Zimak. Maximum margin coresets for active and noise tolerant learning. In *IJCAI*, pages 836–841, 2007.

- [34] T. Hastie, R. Tibshirani, and J. Friedman. *The elements of statistical learning: data mining, inference, and prediction*. Springer Science & Business Media, 2009.
- [35] J. Huggins, T. Campbell, and T. Broderick. Coresets for scalable bayesian logistic regression. In *Advances In Neural Information Processing Systems*, pages 4080–4088, 2016.
- [36] F. John. Extremum problems with inequalities as subsidiary conditions. In *Traces and emergence of nonlinear programming*, pages 197–215. Springer, 2014.
- [37] I. Jubran, M. Tukan, A. Maalouf, and D. Feldman. Sets clustering. *arXiv preprint arXiv:2003.04135*, 2020.
- [38] Z. Karnin and E. Liberty. Discrepancy, coresets, and sketches in machine learning. In *Conference on Learning Theory*, pages 1975–1993, 2019.
- [39] J.-C. Kuang. Applied inequalities. *Shandong Science and Technology Press, Jinan, China*, 3, 2004.
- [40] M. Langberg and L. J. Schulman. Universal ε -approximators for integrals. In *Proceedings of the twenty-first annual ACM-SIAM symposium on Discrete Algorithms*, pages 598–607. SIAM, 2010.
- [41] L. Lovász. *An algorithmic theory of numbers, graphs and convexity*. SIAM, 1986.
- [42] M. Lucic, O. Bachem, and A. Krause. Strong coresets for hard and soft bregman clustering with applications to exponential family mixtures. In A. Gretton and C. C. Robert, editors, *Proceedings of the 19th International Conference on Artificial Intelligence and Statistics*, volume 51 of *Proceedings of Machine Learning Research*, pages 1–9, Cadiz, Spain, 09–11 May 2016. PMLR.
- [43] M. Lucic, M. Faulkner, A. Krause, and D. Feldman. Training gaussian mixture models at scale via coresets. *The Journal of Machine Learning Research*, 18(1):5885–5909, 2017.
- [44] A. Maalouf, I. Jubran, and D. Feldman. Fast and accurate least-mean-squares solvers. In *Advances in Neural Information Processing Systems*, pages 8305–8316, 2019.
- [45] A. Maalouf, A. Statman, and D. Feldman. Tight sensitivity bounds for smaller coresets. *arXiv preprint arXiv:1907.01433*, 2019.
- [46] C. D. Meyer. *Matrix analysis and applied linear algebra*, volume 71. Siam, 2000.
- [47] A. Munteanu, C. Schwiegelshohn, C. Sohler, and D. Woodruff. On coresets for logistic regression. In *Advances in Neural Information Processing Systems*, pages 6561–6570, 2018.
- [48] T. E. Oliphant. *A guide to NumPy*, volume 1. Trelgol Publishing USA, 2006.
- [49] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830, 2011.
- [50] J. M. Phillips. Coresets and sketches. *arXiv preprint arXiv:1601.00617*, 2016.
- [51] W. Rudin. *Functional Analysis*. International series in pure and applied mathematics. McGraw-Hill, 1991.
- [52] T. Sarlos. Improved approximation algorithms for large matrices via random projections. In *2006 47th Annual IEEE Symposium on Foundations of Computer Science (FOCS'06)*, pages 143–152. IEEE, 2006.
- [53] M. Schmidt, C. Schwiegelshohn, and C. Sohler. Fair coresets and streaming algorithms for fair k-means. In *International Workshop on Approximation and Online Algorithms*, pages 232–251. Springer, 2019.
- [54] C. Sohler and D. P. Woodruff. Subspace embeddings for the l_1 -norm with applications. In *Proceedings of the forty-third annual ACM symposium on Theory of computing*, pages 755–764, 2011.

- [55] C. Sohler and D. P. Woodruff. Strong coresets for k-median and subspace approximation: Goodbye dimension. In *2018 IEEE 59th Annual Symposium on Foundations of Computer Science (FOCS)*, pages 802–813. IEEE, 2018.
- [56] E. Tolochinsky and D. Feldman. Generic coreset for scalable learning of monotonic kernels: Logistic regression, sigmoid and more. *arXiv preprint arXiv:1802.07382*, 2018.
- [57] I.-H. Tsang, J.-Y. Kwok, and J. M. Zurada. Generalized core vector machines. *IEEE Transactions on Neural Networks*, 17(5):1126–1140, 2006.
- [58] I. W. Tsang, J. T. Kwok, and P.-M. Cheung. Core vector machines: Fast svm training on very large data sets. *Journal of Machine Learning Research*, 6(Apr):363–392, 2005.
- [59] I. W. Tsang, J. T.-Y. Kwok, and P.-M. Cheung. Very large svm training using core vector machines. In *AISTATS*, 2005.
- [60] M. Tukan, C. Baykal, D. Feldman, and D. Rus. On coresets for support vector machines. *arXiv preprint arXiv:2002.06469*, 2020.
- [61] M. Tukan, A. Maalouf, and D. Feldman. Open source code for all the algorithms presented in this paper, 2019. [Link for open-source code](#).
- [62] A. V. Uzilov, J. M. Keegan, and D. H. Mathews. Detection of non-coding rnas on the basis of predicted secondary structure formation free energy change. *BMC bioinformatics*, 7(1):173, 2006.
- [63] G. Van Rossum and F. L. Drake. *Python 3 Reference Manual*. CreateSpace, Scotts Valley, CA, 2009.
- [64] P. Virtanen, R. Gommers, T. E. Oliphant, M. Haberland, T. Reddy, D. Cournapeau, E. Burovski, P. Peterson, W. Weckesser, J. Bright, S. J. van der Walt, M. Brett, J. Wilson, K. Jarrod Millman, N. Mayorov, A. R. J. Nelson, E. Jones, R. Kern, E. Larson, C. Carey, Í. Polat, Y. Feng, E. W. Moore, J. Vand erPlas, D. Laxalde, J. Perktold, R. Cimrman, I. Henriksen, E. A. Quintero, C. R. Harris, A. M. Archibald, A. H. Ribeiro, F. Pedregosa, P. van Mulbregt, and S. . . Contributors. SciPy 1.0: Fundamental Algorithms for Scientific Computing in Python. *Nature Methods*, 2020.
- [65] D. Woodruff and Q. Zhang. Subspace embeddings and ℓ_p -regression using exponential random variables. In *Conference on Learning Theory*, pages 546–567, 2013.