



**HAL**  
open science

# Large-Scale High-Dimensional Clustering with Fast Sketching

Antoine Chatalic, Rémi Gribonval, Nicolas Keriven

► **To cite this version:**

Antoine Chatalic, Rémi Gribonval, Nicolas Keriven. Large-Scale High-Dimensional Clustering with Fast Sketching. ICASSP 2018 - IEEE International Conference on Acoustics, Speech and Signal Processing, Apr 2018, Calgary, Canada. pp.4714-4718, 10.1109/ICASSP.2018.8461328 . hal-01701121

**HAL Id: hal-01701121**

**<https://inria.hal.science/hal-01701121>**

Submitted on 5 Feb 2018

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# LARGE-SCALE HIGH-DIMENSIONAL CLUSTERING WITH FAST SKETCHING

Antoine Chatalic<sup>\*</sup>, Rémi Gribonval<sup>†</sup> and Nicolas Keriven<sup>\*</sup>

<sup>\*</sup> Université de Rennes 1, France    <sup>†</sup> Inria Rennes, France

## ABSTRACT

In this paper, we address the problem of high-dimensional k-means clustering in a large-scale setting, i.e. for datasets that comprise a large number of items. Sketching techniques have already been used to deal with this “large-scale” issue, by compressing the whole dataset into a single vector of random nonlinear generalized moments from which the  $k$  centroids are then retrieved efficiently. However, this approach usually scales quadratically with the dimension; to cope with high-dimensional datasets, we show how to use fast structured random matrices to compute the sketching operator efficiently. This yields significant speed-ups and memory savings for high-dimensional data, while the clustering results are shown to be much more stable, both on artificial and real datasets.

**Index Terms**— Sketching, Sketched Learning, Fast Transforms, Structured Matrices, k-means, Random Fourier Features.

## 1. INTRODUCTION

Let  $\mathcal{X} = \{x_1, \dots, x_n\} \subset \mathbb{R}^d$  be a set of  $n$   $d$ -dimensional points. We consider the problem of k-means clustering, which consists in finding  $k$  centroids  $\mathcal{C} = \{c_1, \dots, c_k\} \subset \mathbb{R}^d$  minimizing the sum of squared errors (SSE):

$$\text{SSE}(\mathcal{X}, \mathcal{C}) = \sum_{i=1}^n \min_j \|x_i - c_j\|^2. \quad (1)$$

We are interested in the case where the size of the dataset  $n$ , the dimension  $d$ , and possibly the number of clusters  $k$  are large. The standard iterative k-means heuristic of Lloyd [1] is widely used because of its simplicity but scales in  $\Theta(ndk)$  per iteration, cannot be easily distributed, and requires to load all the dataset in memory, which limits its usability in this context.

A framework [2] has been proposed to deal with large collections by compressing the whole dataset into a single  $m$ -dimensional vector  $\hat{z}$  of random generalized moments calculated as follows:

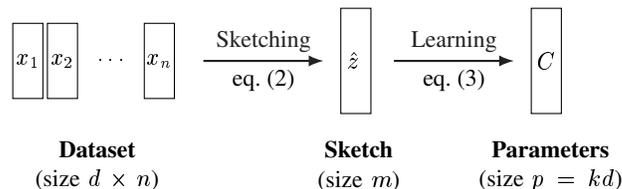
$$\hat{z} = \frac{1}{n} \sum_{i=1}^n \Phi(x_i), \text{ where } \Phi(x) = [e^{-i\omega_1^T x}, \dots, e^{-i\omega_m^T x}]^T. \quad (2)$$

The  $\omega_i$  are here frequency vectors drawn i.i.d. from an isotropic distribution  $\Lambda$  [3]. The sketch is therefore simply made of  $m$  random samples of the empirical characteristic function. Note that the sketching process can be very easily performed in a distributed manner, or even on a data stream. As depicted in Figure 1, the centroids  $\mathcal{C}$  and associated weights  $\alpha$  can then be estimated efficiently from this sketch, and without using the initial data, as one solution of:

$$\mathcal{C}, \alpha \in \arg \min_{c, \alpha} \left\| \hat{z} - \sum_{i=1}^k \alpha_i \Phi(c_i) \right\|_2. \quad (3)$$

This problem has many similarities with compressive sensing, and theoretical guarantees have been obtained in this context [4].

The optimization problem is non-convex, but approximate solutions can be found using greedy heuristics such as CL-OMPR [3]



**Fig. 1:** Overview of the general workflow. In practice, the sketch size  $m$  should be of the order of  $p$  to get good results, and  $p = kd$ .

which is inspired from orthogonal matching pursuit. Some other algorithms, bearing similarities with CL-OMPR, have been proposed in other contexts for solving such sparse inverse problems [5, 6].

If  $X = [x_1, \dots, x_n]$  denotes the dataset in a matrix form, and  $W = [\omega_1, \dots, \omega_m]$  the dense  $d \times m$  matrix of frequency vectors, computing the sketch involves computing the matrix product  $W^T X$ . Previous empirical studies [2] suggested that the size  $m$  of the sketch should be of the order of the number of parameters to learn (i.e.  $m \approx p = kd$ ) to obtain a quality of clustering that is similar to k-means. Under this assumption, the cost of both sketching and learning phases is dominated by such matrix products and scales quadratically with  $d$ . Multiplications by large random matrices appear in various contexts, and multiple works have proposed to replace them by random structured matrices, which behave similarly but have less degrees of freedom, and for which the matrix product can be computed efficiently.

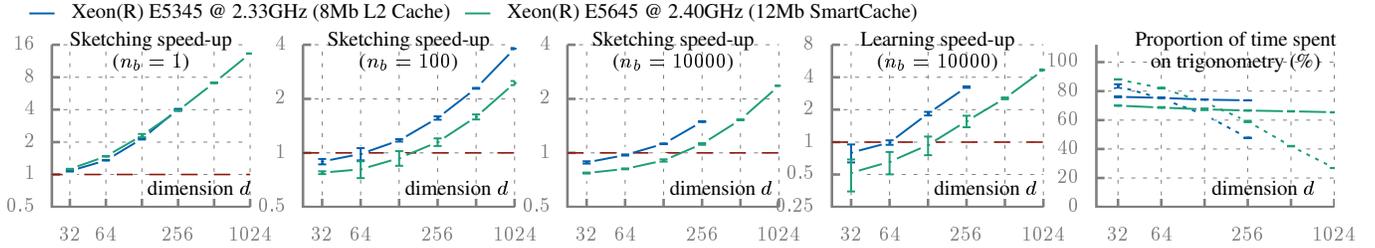
We show how to combine in a single framework the scalability of the sketching approach — thus allowing us to cope with very large and distributed collections (large  $n$ ) — with the computational efficiency of such fast transforms (large  $d$ ). In our approach, both sketching and learning phases now scale with  $d$  in  $\Theta(d \log_2(d))$ , which makes it possible to work with high-dimensional data. When the number of centroids  $k$  is large as well, we propose to leverage a hierarchical learning algorithm previously introduced for Gaussian mixtures [3], thus reducing the complexity of the learning phase with respect to  $k$  from  $\Theta(k^3)$  to  $\Theta(k^2 \log k)$ .

We present some related works in Section 2, detail the proposed method (Section 3) and give experimental results (Section 4); the fast transforms are shown to give significant speed-ups in high dimension, and perhaps surprisingly much more *stable* clustering results both on artificial and real datasets.

## 2. RELATED WORK

Multiple approaches have been used for large-scale high-dimensional k-means. The dimension of the data can be reduced using feature selection [7], sparsification [8], or geometry-preserving dimensionality-reduction techniques [9] according to the Johnson-Lindenstrauss lemma [10]. Coresets methods [11], on the other side, reduce the size of the dataset but not the dimension. Sketching using random Fourier sampling [12] has been used not only for k-means [2], but





**Fig. 2:** Sketching speed-ups (i.e. ratios of running times without/with fast transforms) for three values of the batch size  $n_b$ , speed-up of the learning phase, and proportion of the sketching time spent on computing the nonlinearity (sines and cosines). Results obtained on 30 experiments. In the rightmost figure, plain and dashed lines denote respectively the usage of fast and dense matrices.

### 3.5. Summary

In the following, KM stands for “k-means” [1], CKM for “Compressive k-means” [2], and FCKM for “Fast Compressive k-means”—i.e. using structured matrices. We refer to the learning algorithm using fast hierarchical initialization as “Hierarchical”, and both the CL-OMPR and Hierarchical algorithms can be used with dense or structured matrices, i.e. in the CKM or FCKM frameworks.

We give in Table 1 a summary of space and time complexities of the different methods assuming  $m = \Theta(kd)$ , as it empirically [2] seems to be a necessary condition to get good clustering results. One should keep in mind that the sketching time can be drastically reduced by relying on distributed computing, which is not the case when using Lloyd’s k-means.

## 4. EXPERIMENTAL VALIDATION

We first give some implementation details (Section 4.1), and then present experiments on randomly generated (Section 4.2) and real data (Sections 4.3 and 4.4).

### 4.1. Implementation details

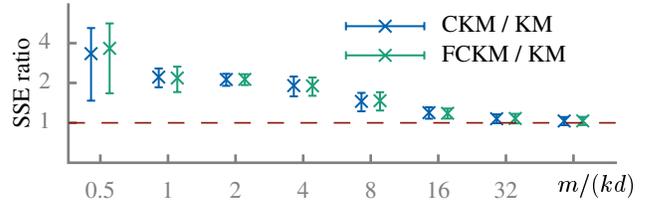
We implemented the fast transform proposed in Section 3 as a contribution to the SketchMLbox Matlab Toolbox [21], which already includes the sketching procedure, the CL-OMPR heuristic and the Hierarchical algorithm for Gaussian mixtures, but always using a dense matrix of frequency vectors.

The critical part of the code is written in C and compiled to binary MEX files. We use the adaptive Fast Walsh-Hadamard Transform of the Spiral project [22, 23], which is designed for an optimized usage of the cache hierarchy [24]. One might get higher speed-ups by relying on carefully designed SIMD implementations [25]. In the experiments, KM always refer to the Matlab implementation of k-means using  $I_{max} = 1000$  as the maximum number of iterations, and with uniform initialization—except when using k-means++ in Section 4.3.

### 4.2. Synthetic Data

We first show on artificial data that one can replace the dense matrix by a structured one without degrading the quality of the results. We perform here k-means clustering on  $n = 10000$  data vectors randomly generated according to a mixture of  $k = 10$  Gaussians with identity covariance matrix. The means are drawn with respect to a centered Gaussian with covariance  $1.5k^{1/d}I_d$  to create clusters that are well separated with high probability [3].

The quality of the clustering is measured using SSE (1), and Figure 3 shows the ratios obtained using CKM or FCKM with respect to Matlab’s k-means (KM). These results confirm that the compressive



**Fig. 3:** Ratios of SSE for CKM and FCKM with respect to one run of Matlab’s k-means (KM) as a function of  $m/(kd)$ . Averaged on 30 repetitions for each  $d \in \{8, 16, 32, 64, 128, 256, 512\}$ , and  $k = 10$ .

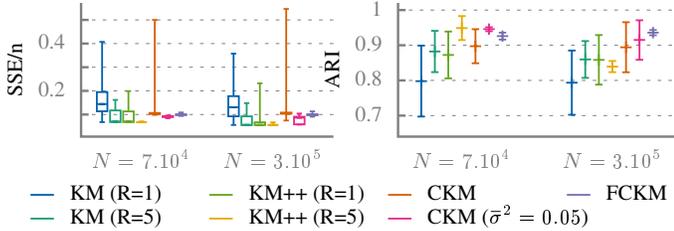
k-means framework gives good clustering results provided that the sketch is large enough, but also that in this case, using fast transforms does not degrade the clustering quality.

Sketching speed-ups obtained using structured rather than dense matrices of frequencies are given in Figure 2. We consider different values of the size  $n_b$  of the batches used for sketching—i.e. we compute products such as  $W^T X_i$ , where  $X_i$  is a batch of size  $d \times n_b$ . The speed-ups are significant for high-dimensional data, and especially in the “streaming” case, i.e. when sketching the vectors one by one; this is very interesting when the amount of available memory is limited. For larger batch sizes, using fast transforms might not be useful for small dimensions (the BLAS matrix-matrix product being too optimized to compete with it), but still allows us to deal with high-dimensional datasets. The learning phase benefits similarly from using fast transforms as depicted on the 4th sub-figure. Note that the computational cost of the complex exponential is very high (5th sub-figure): approximately 70% of the time for all dimensions when using fast transforms. Working with alternative feature maps relying on cheaper non-linear functions could be interesting for future works.

### 4.3. Spectral clustering on MNIST

As long as we work on high-dimensional datasets, one should expect to observe on real datasets the same speed-ups as the ones obtained on random data. In the following, we consider datasets with dimensions for which one should not expect to get significant speed-ups when using a large batch size according to Figure 2; however, it is highly interesting to check whether the fact that structured matrices yield the same clustering quality on random data (see Section 4.2) also holds for real datasets.

We perform here clustering on the MNIST dataset [26] of handwritten digits, which has  $k = 10$  classes. The original dataset contains  $n = 7 \times 10^4$  pictures. Distorted variants of these images have been generated using infMNIST [27], so that one other dataset of size  $n = 3 \times 10^5$  is used for evaluation as well. For every image, we extract dense SIFT descriptors [28], which are concatenated into a single vector. We compute the similarity matrix between these vectors, and the  $k$  first eigenvectors of the associated Laplacian matrix



**Fig. 4:** Min/quartiles/max boxplots and median of SSE (left, the lower the better) and mean & standard deviation of the adjusted Rand index (ARI) (right, the higher the better). Results obtained on 120 experiments,  $m = 1008$ ;  $R$  corresponds to the number of replicates. KM was run with a maximum number of 1000 iterations, using uniform initialization. The legend, when read column by column from KM to FCKM, corresponds to the order of the boxplots from left to right.

in order to get  $n$  spectral features [29] in dimension  $d = k = 10$ .

Figure 4 gives the results in terms of SSE and adjusted Rand index (ARI), for KM, CKM and FCKM. We also consider k-means++ (KM++) [30], and include results with  $R = 5$  replicates, i.e. running the algorithm 5 times and keeping the best results. Note that for CKM and FCKM, the distribution used to draw the frequency vectors involves a parameter  $\bar{\sigma}^2$  that is estimated automatically by first sketching a small subset of the data [3]. For comparison, we also consider using CKM with a fixed value of  $\bar{\sigma}^2$  rather than relying on this estimation. Results obtained with CKM contain roughly 15% of outliers. Using CKM with  $\bar{\sigma}^2 = 0.05$  gives better and highly concentrated results, but one usually does not have this knowledge of the dataset. FCKM turns out to give similar results in terms of SSE and ARI, but without requiring any knowledge on  $\bar{\sigma}^2$ , as it is here again automatically estimated. Results are well concentrated, contrarily to what is obtained with Lloyd’s k-means, even with  $R = 5$  replicates. KM++ with  $R = 5$  replicates gives good results in terms of SSE, but lower ARI for  $n = 3 \times 10^5$ ; one should keep in mind that KM and KM++ require to have the whole dataset in memory.

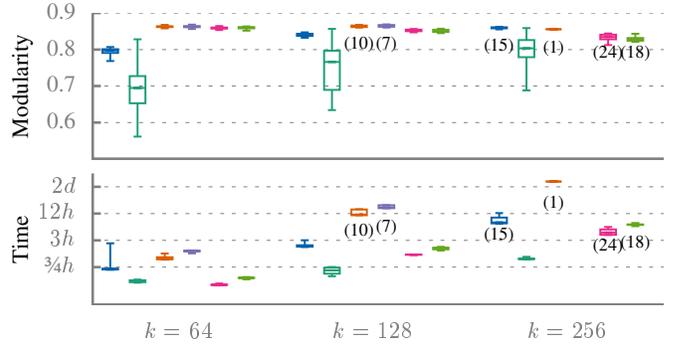
In summary, even in small dimensions where sketching might be faster with dense matrices, the use of structured random matrices is potentially useful in terms of stability. To get the best possible speed, one can rely on the explicit dense representation of such structured matrices if applicable.

#### 4.4. Hierarchical clustering on a co-purchasing graph

Computing spectral features involves computing the eigendecomposition of the Laplacian matrix, which is usually expensive. Tremblay et al. proposed to bypass this step by using as features  $\Theta(\log k)$  filtered random graph signals on the graph associated to the Laplacian matrix [31]. Standard KM is then applied on a random subset of these features, and interpolated to the whole collection. We combine these fast random features with the FCKM framework, thus allowing us to avoid the subsampling step.

We work on the Amazon co-purchasing network [32], which is a graph comprising  $n = 334863$  nodes and  $E = 925872$  edges. As there is no ground truth for this dataset, we used  $k = 64, 128, 256$ . We compare the original spectral clustering (SC), compressive spectral clustering (CSC) [31], and 4 methods using sketching on the same random features: we combine the two types of matrices (dense/structured) with the two learning procedures discussed in Section 3.4 (CL-OMPR/Hierarchical). Please refer to the table of Figure 5 for a summary.

Standard KM is launched with  $R = 2$  replicates. Using compressive spectral clustering, the k-means step is performed only on a



Method	Features	Subs.	Sk. matrix	Clustering
SC	spectral	No	n/a	KM
CSC	random	Yes	n/a	KM
S2C	random	No	Dense	CL-OMPR
FS2C	random	No	Structured	CL-OMPR
HS2C	random	No	Dense	Hierarchical
HFS2C	random	No	Structured	Hierarchical

**Fig. 5:** Boxplots of modularity (the higher the better) and clustering time for  $k = 64, 128, 256$ . Only the learning times are displayed for sketching methods; sketching times are much smaller, even on a single core. All experiments were run on Intel(R) Xeon(R) CPUs E5640 and repeated 30 times (or less for experiments that were too long; the number of iterations is indicated below in these cases, and FS2C does not appear for  $k = 256$ ). The table is a summary of the different methods (in the order of the boxplots, from left to right).

subset of the features and is therefore much faster; we used  $R = 20$  replicates for a fair comparison. We used  $m = 10kd$  for the sketch size when using CKM. All initializations were performed uniformly.

The results are presented in Figure 5 (top), where the elapsed times are given and the clustering quality is measured with the modularity metric [33]. As regards the elapsed times, CL-OMPR is not competitive but satisfying results are obtained with the hierarchical algorithm. Similar modularities are obtained with and without structured matrices; the results are slightly lower when using the hierarchical algorithm, but in both cases they are highly concentrated, whereas CSC yields a high variance.

## 5. PERSPECTIVES

We proposed a way of combining the efficiency of fast structured matrices with the scalability of sketching-based approaches and hierarchical learning methods, yielding a k-means framework which can handle large and high-dimensional collections with a limited memory footprint. Experimental validation confirms that significant speedups are obtained in high dimension, and the clustering results seem to be much more stable.

It would of course be interesting to be able to control the error induced by the use of structured matrices. In the theoretical framework [4] which has been proposed for sketched learning, the matrix used for sketching implicitly defines a kernel function, and studying the kernel associated with the structured matrices used in our approach should help to establish theoretical guarantees.

Although the hierarchical algorithm proposed for the approximate minimization allows to deal with a larger number of centroids compared to CL-OMPR, it still scales quadratically with  $k$  when  $m = \Theta(kd)$ , leaving room for improvement; locality could for instance be leveraged, as it has been done for orthogonal matching pursuit [34]. Designing provably-good procedures for this optimization problem is challenging as well.

## 6. REFERENCES

- [1] Stuart Lloyd, “Least squares quantization in pcm,” *IEEE transactions on information theory*, vol. 28, no. 2, pp. 129–137, 1982.
- [2] Nicolas Keriven, Nicolas Tremblay, Yann Traonmilin, and Rémi Gribonval, “Compressive k-means,” 2017.
- [3] Nicolas Keriven, Anthony Bourrier, Rémi Gribonval, and Patrick Pérez, “Sketching for large-scale learning of mixture models,” *To be published in Information and Inference*, vol. abs/1606.02838, 2016.
- [4] Rémi Gribonval, Gilles Blanchard, Nicolas Keriven, and Yann Traonmilin, “Compressive statistical learning with random feature moments,” .
- [5] Kristian Bredies and Hanna Katriina Pikkarainen, “Inverse problems in spaces of measures,” *ESAIM: Control, Optimization and Calculus of Variations*, vol. 19, no. 1, pp. 190–218, 2013.
- [6] Nicholas Boyd, Geoffrey Schiebinger, and Benjamin Recht, “The alternating descent conditional gradient method for sparse inverse problems,” *SIAM Journal on Optimization*, vol. 27, no. 2, pp. 616–639, 2017.
- [7] Christos Boutsidis, Petros Drineas, and Michael W Mahoney, “Unsupervised feature selection for the  $k$ -means clustering problem,” in *Advances in Neural Information Processing Systems*, 2009, pp. 153–161.
- [8] Piotr Indyk, Jiří Matoušek, and Anastasios Sidiropoulos, “Low-distortion embeddings of finite metric spaces,” in *Handbook of discrete and computational geometry*, Csaba D Toth, Joseph O’Rourke, and Jacob E Goodman, Eds. CRC press.
- [9] Christos Boutsidis, Anastasios Zouzias, and Petros Drineas, “Random projections for  $k$ -means clustering,” in *Advances in Neural Information Processing Systems*, 2010, pp. 298–306.
- [10] William B. Johnson and Joram Lindenstrauss, “Extensions of lipschitz mappings into a hilbert space,” vol. 26, no. 189.
- [11] Jeff M. Phillips, “Coresets and sketches,” in *Handbook of discrete and computational geometry*, Csaba D Toth, Joseph O’Rourke, and Jacob E Goodman, Eds., chapter 48. CRC press, 2017.
- [12] Ali Rahimi and Benjamin Recht, “Random features for large-scale kernel machines,” in *Advances in Neural Information Processing Systems (NIPS)*. 2007, pp. 1177–1184, Curran Associates, Inc.
- [13] Alastair R Hall, *Generalized method of moments*, Oxford University Press, 2005.
- [14] Lars Peter Hansen, “Large sample properties of generalized method of moments estimators,” *Econometrica: Journal of the Econometric Society*, pp. 1029–1054, 1982.
- [15] Quoc Le, Tamás Sarlós, and Alex Smola, “Fastfood—approximating kernel expansions in loglinear time,” in *Proceedings of the international conference on machine learning (ICML)*, 2013.
- [16] Krzysztof Choromanski and Vikas Sindhwani, “Recycling randomness with structure for sublinear time kernel expansions,” in *Proceedings of the international conference on machine learning (ICML)*. 2016, vol. 48 of *JMLR Workshop and Conference Proceedings*, pp. 2502–2510, JMLR.org.
- [17] Felix X. Yu, Ananda Theertha Suresh, Krzysztof Marcin Choromanski, Daniel Holtmann-Rice, and Sanjiv Kumar, “Orthogonal random features,” in *Advances in Neural Information Processing Systems (NIPS)*, 2016, pp. 1975–1983.
- [18] Mariusz Bojarski, Anna Choromanska, Krzysztof Choromanski, Francois Fagan, Cedric Gouy-Pailler, Anne Morvan, Nouri Sakr, Tamas Sarlos, and Jamal Atif, “Structured adaptive and random spinners for fast machine learning computations,” in *The 20th International Conference on Artificial Intelligence and Statistics*.
- [19] Krzysztof Choromanski, Mark Rowland, and Adrian Weller, “Random features for large-scale kernel machines,” in *Advances in Neural Information Processing Systems (NIPS)*.
- [20] Nir Ailon and Bernard Chazelle, “The fast johnson–lindenstrauss transform and approximate nearest neighbors,” vol. 39, no. 1, pp. 302–322.
- [21] Nicolas Keriven, Nicolas Tremblay, and Rémi Gribonval, “Sketchmlbox: A matlab toolbox for large-scale mixture learning.” <http://sketchml.gforge.inria.fr/>,” 2016.
- [22] “Spiral project: Wht package.” <http://www.spiral.net/software/wht.html>,” .
- [23] J. Johnson and M. Puschel, “In search of the optimal walsh-hadamard transform,” in *IEEE International Conference on Acoustics, Speech, and Signal Processing. Proceedings*, vol. 6.
- [24] Neungsoo Park and N. K. Prasanna, “Cache conscious walsh-hadamard transform,” in *Proceedings of the International Conference on Acoustics, Speech and Signal Processing 2001 (ICASSP)*, vol. 2, pp. 1205–1208 vol.2.
- [25] Joachim Curto, Irene Zarza, Feng Yang, Alexander Smola, and Luc Van Gool, “F2f: A library for fast kernel expansions,” .
- [26] Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner, “Gradient-based learning applied to document recognition,” vol. 86, no. 11, pp. 2278–2324.
- [27] Gaëlle Loosli, Stéphane Canu, and Léon Bottou, “Training invariant support vector machines using selective sampling,” *Large scale kernel machines*, pp. 301–320, 2007.
- [28] David G Lowe, “Distinctive image features from scale-invariant keypoints,” *International journal of computer vision*, vol. 60, no. 2, pp. 91–110, 2004.
- [29] Andrew Y. Ng, Michael I. Jordan, and Yair Weiss, “On spectral clustering: Analysis and an algorithm,” in *Advances in neural information processing systems*, pp. 849–856.
- [30] David Arthur and Sergei Vassilvitskii, “k-means++: The advantages of careful seeding,” in *Proceedings of the eighteenth annual ACM-SIAM symposium on Discrete algorithms*, 2007.
- [31] Nicolas Tremblay, Gilles Puy, Rémi Gribonval, and Pierre Vandergheynst, “Compressive spectral clustering,” in *ICML 2016*, June, pp. 20–22.
- [32] Jaewon Yang and Jure Leskovec, “Defining and evaluating network communities based on ground-truth,” *Knowledge and Information Systems*, vol. 42, no. 1, pp. 181–213, 2015.
- [33] Mark EJ Newman and Michelle Girvan, “Finding and evaluating community structure in networks,” *Physical review E*, vol. 69, no. 2, pp. 026113, 2004.
- [34] Boris Mailhé, Rémi Gribonval, Frédéric Bimbot, and Pierre Vandergheynst, “LocOMP: algorithme localement orthogonal pour l’approximation parcimonieuse rapide de signaux longs sur des dictionnaires locaux,” in *GRETSI*.