

MACHINE LEARNING FOR DATA STREAMS

CLUSTERING

1. Efficient and Effective Clustering Methods for Spatial Data Mining (1994):

- A cluster is represented by its **medoid**, or the most, centrally located data point in the cluster.
- The clustering process is formalized as **searching a graph in which each node is a K-partition represented by a set of K medoids**, and two nodes are neighbors if they only differ by one medoid.
- Starts with a randomly selected node. For the current node, it checks at most the *maxneighbor* number of neighbors randomly, and if a better neighbor is found, it moves to the neighbor and continues; otherwise it records the current node as a *local minimum*, and restarts with a new randomly selected node to search for another local minimum. It stops after the *numlocal* number of the so-called *local minima* have been found, and returns the best of these.
- Suffers from the same drawbacks as the above IO method wrt. efficiency.
- It may not find a real local minimum due to the searching trimming controlled by *maxneighbor*.
- All of the k-medoid types of approaches, including PAM, CLARA, and CLARANS, are known not to be scalable and thus are not appropriate in a streaming context.

2. BIRCH: An Efficient Data Clustering Method for Very Large Databases (1996):

- Balanced Iterative Reducing and Clustering using Hierarchies.
- Appropriate for **very large datasets**, by making the time and memory constraints explicit.
- **Incrementally and dynamically clusters incoming multidimensional metric data points** to try to produce the best quality clustering with the available resources (i.e., available memory and time constraints).
- It is **local** (as opposed to global) in that each clustering decision is made without scanning all data points or all currently existing clusters.
- It exploits the observation that the data space is **usually not uniformly occupied**, and hence **not every data point is equally important for clustering purposes**. A dense region of points is treated collectively as a **single cluster**. Points in sparse regions are treated as **outliers** and removed optionally.
- It makes full use of available memory to derive the finest possible subclusters (to ensure accuracy) while minimizing I/O costs (to ensure efficiency).

- Can typically find a good clustering with a **single scan of the data**, and improve the quality further with a few additional scans.
- First clustering algorithm proposed in the database area to **handle “noise”** (data points that are not part of the underlying pattern) effectively.
- Considers **metric** attributes, as in most of the Statistics literature (attributes whose values satisfy the requirements of Euclidean space, i.e., self identity and triangular inequality).
- Offers opportunities for **parallelism**, and for **interactive** or **dynamic performance tuning** based on knowledge about the dataset, gained over the course of the execution.
- **The clustering and reducing process is organized and characterized by the use of an in-memory, height-balanced and highly-occupied tree structure.** Due to these features, its running time is **linearly scalable**.
- **Centroid, radius and diameter** as properties of a single cluster, and **Euclidean distance, Manhattan distance, average inter-cluster distance, average intra-cluster distance** (diameter of the merged cluster) and **variance increase distance** as properties between two clusters and state them separately.
- The concepts of **Clustering Feature** and **CF** tree are at the core of *BIRCH*'s incremental clustering. A **Clustering Feature** is a triple summarizing the information that we maintain about a cluster (summarizes a set of points by the sum of the points, the number of points and the sum of the squared lengths of all points).
- Given the **CF** vectors of clusters, the corresponding distances, as well as the usual quality metrics, can all be calculated easily.
- **CF** summary is not only **efficient** because it stores much less than all the data points in the cluster, but also **accurate** because it is sufficient for calculating all the measurements that we need for making clustering decisions in *BIRCH*.
- A **CF** tree is a height-balanced tree with two parameters: branching factor B and threshold T .
- **BIRCH algorithm:**
 - Phase 1: Scan all the data and build an initial in-memory **CF** tree using the given amount of memory and recycling space on disk (an initial memory summary of the data).
 - Phase 2 (optional): Scan the leaf entries in the initial **CF** tree to rebuild a smaller **CF** tree, while removing more outliers and grouping crowded subclusters into larger ones.
 - Phase 3: A global or semi-global algorithm is used to cluster all leaf entries. They adapt an agglomerative hierarchical clustering algorithm by applying it directly to the subclusters represented by their CF vectors. The undesirable effect of the skewed input order, and splitting triggered by page size causes to be unfaithful to the actual clustering patterns in the data. After this phase, we

obtain a set of clusters that captures the major distribution pattern in the data.

- Phase 4 (optional): Entails the cost of additional passes over the data to correct minor and localized inaccuracies because of the rare misplacement problem and refine the clusters further .
 - When using Clustering Features to store a small summary of points, the quality of the summary decreases when storing points together in one CF that should be assigned to different centers later on. If we summarize points in a CF and later on get centers where all these points are closest to the same center, then their clustering cost can be computed with the CF without any error. Thus, the idea of BIRCH is to heuristically identify points that are likely to be clustered together. For this purpose, they use the insertion process defined in the paper.
 - BIRCH works with increasing thresholds when processing the input data. It starts with threshold $t = 0$ and then increases t whenever the number of CFs exceeds a given space bound, calling a rebuilding algorithm to shrink the tree. This algorithm ensures that the number of CFs is decreased sufficiently. Notice that CFs cannot be split again, so the rebuilding might return a different tree than the one computed directly with the new threshold.
 - **BIRCH has no theoretical quality guarantees and does indeed sometimes perform badly in practice.**
- Comparison with CLARANS (1)
 - CLARANS assumes that the memory is enough for holding the whole dataset, so it needs much more memory than BIRCH does.
 - For all three datasets of the base workload, CLARANS is at least 15 times slower than BIRCH, and is sensitive to the pattern of the dataset.
 - The D (with a hyphen on the top) value for the CLARANS clusters is much larger than that for the BIRCH clusters.
 - In conclusion, for the base workload, BIRCH uses much less memory, but is faster, more accurate, and less order-sensitive compared with CLARANS.

RESUMEN BIRCH (STREAM): Birch compresses a large dataset into a smaller one via a CFtree (clustering feature tree). Each leaf of this tree captures sufficient statistics (namely the first and second moments) of a subset of points. Internal nodes capture sufficient statistics of the leaves below them. The algorithm for computing a CFtree repeatedly inserts points into the leaves provided that the radius of the set of points associated with a leaf does not exceed a certain threshold. If the threshold is exceeded then a new leaf is created and the tree is appropriately balanced. If the tree does not fit in main memory then a new threshold is used to create a smaller tree.

3. Streaming-Data Algorithms For High-Quality Clustering (2002)_STREAM:

- Most heuristics, including k-Means, are also infeasible for data streams because they require random access. As a result, several heuristics have been proposed for scaling clustering algorithms. In the database literature, the BIRCH system is commonly considered to provide a competitive heuristic for this problem. There are no guarantees on their SSQ performance.
- The k-Means algorithm and BIRCH are most relevant to their results. Most other previous work on clustering either does not offer the

scalability required for a fast streaming algorithm or does not directly optimize SSQ.

- It mentions [CLARANS \(1\)](#) because choosing a new medoid among all the remaining points is time-consuming and, to address this problem, CLARANS draws a fresh sample of feasible centers before each calculation of SSQ improvement.
- They assume that our data stream is not sorted in any way.
- Their solution for k-Median is obtained via a variant called **facility location**, which does not specify in advance the number of clusters desired, and instead evaluates an algorithm's performance by a combination of SSQ and the number of centers used.
- The streaming algorithm given in Section 3 is shown to enjoy theoretical quality guarantees.
- The SSQ minimization problem is identical to k-Median except that $\mathbf{M} = \mathbf{R}^b$ for some integer b , d is the Euclidean metric, the medians can be arbitrary points in \mathbf{M} , and $d^2(x; c_i)$ replaces $d(x; c_i)$ in the objective function; that is, **they minimize the “sum of squares” (SSQ) rather than the sum of distances.**
- **Facility location** is the same as k-Median except that instead of restricting the number of medians to be at most k they simply impose a cost for each median, or facility. The additive cost associated with each facility is called the **facility cost**. It allows as input a range of number of centers.
- Previous problems are known to be **NP-hard**, and several theoretical approximation algorithms are known. Their algorithm for clustering streaming data uses a subroutine called **LSEARCH**. The streaming algorithm, **STREAM**, is as follows: They cluster the i th chunk X_i using **LSEARCH**, and assign each resulting median a weight equal to the sum of the weights of its members from X_i . They then purge memory, retaining only the k weighted cluster centers, and apply **LSEARCH** to the weighted centers they have retained from X_1, \dots, X_i , to obtain a set of (weighted) centers for the entire stream $X_1 \cup \dots \cup X_i$.
- The CG algorithm does not directly solve k-Median but could be used as a subroutine to a k-Median algorithm, as follows. They first set an initial range for the facility cost z (between 0 and an easy-to-calculate upper bound); they then perform a binary search within this range to find a value of z that gives them the desired number k of facilities; for each value of z that they try, we call Algorithm CG to get a solution.

- Time calling CG as a subroutine of a binary search is prohibitive for large data streams. Therefore, they describe a new local search algorithm that relies on the correctness of the above algorithm but avoids the super-quadratic running time by taking advantage of the structure of local search in certain ways.
- The algorithm is implemented as a continuous version of k-means algorithm which continues to maintain a number of cluster centers which change or merge as necessary throughout the execution of the algorithm. Such an approach is especially risky when the characteristics of the stream evolve over time. This is because the k-means approach is highly sensitive to the order of arrival of the data points. For example, once two cluster centers are merged, there is no way to informatively split the clusters when required by the evolution of the stream at a later stage. (Mentioned by CluStream).
- Comparison with BIRCH (2) and k-Means:
 - They give empirical evidence that the clustering algorithm outperforms the commonly-used k-Means algorithm.
 - They also experimentally demonstrate our streaming algorithm's superior performance to Birch.
 - We found that SSQ for k-means was worse than that for LSEARCH, and that LSEARCH typically found near-optimum (if not the optimum) solution.
 - Over the course of multiple runs, there was a large variance in the performance of k-Means, whereas LSEARCH was consistently good. LSEARCH took longer to run for each trial but for most datasets found a near-optimal answer before k-Means found an equally good solution. On many datasets k-Means never found a good solution.
 - K-means is infeasible for data streams because they require random Access.
 - STREAM and BIRCH have a common method of attack: repeated preclustering of the data. However the preclustering of STREAM is bottom up, where every substep is a clustering process, whereas the preclustering in BIRCH is top down partitioning.
 - Overall, our results point to a cluster quality vs. running time tradeoff. In applications where speed is of the essence, e.g., clustering web search results, BIRCH appears to do a reasonable quick and dirty job. In applications like intrusion detection or target marketing where mistakes can be costly our STREAM algorithm exhibits superior SSQ performance.
 - StreamLS is slower than BIRCH but provides a clustering with much better quality (with respect to the sum of squared errors).
 - Shows better accuracy than BIRCH.

RESUMEN STREAM (STREAMKM++): Partitions the input stream into chunks and computes for each chunk a clustering using a local search algorithm.

4. A Framework for Clustering Evolving Data Streams (2003)_CluStream:

- The idea is divide the clustering process into an **online component** which periodically stores detailed summary statistics and an **offline component** which uses only this summary statistics. The offline component is utilized by the analyst who can use a wide variety of inputs (such as time horizon or number of clusters) in order to provide a quick understanding of the broad clusters in the data stream. This two-phased approach also provides the user the flexibility to explore the nature of the evolution of the clusters over different time periods.
- Their proposal utilize two concepts which are useful for efficient data collection in a fast stream: **microclusters** (They maintain statistical information about the data locality in terms of micro-clusters, defined as a temporal extension of the cluster feature vector from BIRCH) and **Pyramidal Time Frame** (the micro-clusters are stored at snapshots in time which follow a pyramidal pattern, at different levels of granularity depending upon the recency, that provides an effective trade-off between the storage requirements and the ability to recall summary statistics from different time horizons). This summary information in the micro-clusters is used by an offline component which is dependent upon a wide variety of user inputs such as the time horizon or the granularity of clustering.
- The offline macro-clustering algorithm will use the finer level micro-clusters in order to create higher level clusters over specific time horizons.
- The macro-clustering algorithm will use some of the subtractive properties of the micro-clusters stored at snapshots t_c and (t_c-h) in order to find the higher level clusters in a history or time horizon of length h .
- Since it is not possible to store the snapshots at each and every moment in time, it is important to choose particular instants of time at which the micro-clusters are stored. The aim of choosing these particular instants is to ensure that clusters in any user specified time horizon $(t_c - h; t_c)$ can be approximated. In order to achieve this, they use the *pyramidal time frame*).
- The **micro-clustering** phase is the online statistical data collection portion of the algorithm. This process is not dependent on any user input such as the time horizon or the required granularity of the clustering process.
- The micro-clusters generated by the algorithm serve as an intermediate statistical representation which can be maintained in an efficient way even for a data stream of large volume. On the other hand, the macro-clustering process does not use the (voluminous) data stream, but the compactly stored summary statistics of the micro-clusters. Therefore, **it is not constrained by one-pass requirements**.

- It is assumed, that as input to the algorithm, the user supplies the time-**horizon h** , and the **number of higher level clusters k** which he wishes to determine. The choice of the time horizon h determines the amount of history which is used in order to create higher level clusters. The choice of the number of clusters k determines whether more detailed clusters are found, or whether more rough clusters are mined.
- The subtractive property of the cluster vector feature helps considerably in determination of the micro-clusters over a pre-specified time horizon.
- The clusters are determined by using a **modification of a k-means algorithm**. In this technique, the resulting micro-clusters are treated as pseudo-points which are re-clustered in order to determine higher level clusters.
- It is important to note that a given execution of the macro-clustering process only needs to use two (carefully chosen) snapshots from the pyramidal time window of the micro-clustering process. The compactness of this input thus allows the user considerable flexibilities for querying the stored micro-clusters with different levels of granularity and time horizons.
- One novel feature of CluStream is that it can create a set of macro-clusters for any user-specified horizon at any time upon demand.
- First paragraph of the Discussion and Conclusions section.
- Comparison with STREAM and BIRCH:
 - Previous algorithms on clustering data streams such as those discussed in the paper where STREAM is explained assume that the clusters are to be computed over the entire data stream. Such methods simply view the data stream clustering problem as a variant of one-pass clustering algorithms. While such a task may be useful in many applications, a clustering problem needs to be defined carefully in the context of a data stream. This is because a data stream should be viewed as an infinite process consisting of data which continuously evolves with time. As a result, the underlying clusters may also change considerably with time.
 - Uses cluster feature vector from BIRCH.
 - The nature of the maintenance process ensures that a very large number of micro-clusters can be efficiently maintained as compared to the STREAM method.
 - CluStream has very good scalability in terms of stream size, dimensionality, and the number of clusters.
 - The offline clustering algorithms cannot detect such intrusions in real time (dataset from KDD-CUP'99). Even the recently proposed stream clustering algorithms such as BIRCH and STREAM cannot be very effective because the clusters reported by these algorithms are all generated from the entire history of data stream, whereas the current cases may have evolved significantly.

- The surprisingly high clustering quality of CluStream benefits from its good design. On the one hand, the pyramidal time frame enables CluStream to approximate any time horizon as closely as desired. On the other hand, the STREAM clustering algorithm can only be based on the entire history of the data stream.
- CluStream is more reliable than STREAM algorithm. In most cases, no matter how many times we run CluStream, it always
- returns the same (or very similar) results. More interestingly, the fine granularity of the micro-cluster maintenance algorithm helps CluStream in detecting the real attacks.
- CluStream is effective for both evolving and stable streams.
- CluStream can achieve higher accuracy than STREAM due to its registering of more detailed information than the k points used by the k -means approach.

5. StreamKM++: A Clustering Algorithms for Data Streams (2010):

- K-means clustering algorithm for data streams.
- Their algorithm computes a small weighted sample of the data stream and solves the problem on the sample using the **k-means++** algorithm.
- To compute the small sample, they propose two new techniques:
 - First, they use a **non-uniform sampling approach** similar to the k -means++ seeding procedure to obtain small coresets from the data stream. Coreset construction applied in this paper has low dependency on the dimensionality of the data.
 - Second, they propose a new data structure which they call a **coreset tree** (a coreset tree T for a point set P is a binary tree that is associated with a hierarchical divisive clustering for P) that solves the problem concerning the k -means++ seeding procedure. The use of these coreset trees significantly speeds up the time necessary for the non-uniform sampling during our coreset construction. The advantage of the coreset tree is that it enables to compute subsequent sample points by taking only points from a subset of P into account that is significantly smaller than n . The resulting sample set S has essentially the same properties as the original k -means++ seeding. They use the same distribution as the k -means++ algorithm does with the exception that the partition is determined by the coreset tree rather than by assigning each point to the nearest sample point.
- Computes very good solutions and is reasonably fast for small k .
- Arthur and Vassilvitskii developed the k -means++ algorithm, which is a seeding procedure for Lloyd's k -means algorithm that guarantees a solution with certain quality and gives good practical results. However, the k -means++ algorithm (as well as Lloyd's algorithm) needs random access on the input data and is not suited for data streams. In this paper, they develop a new clustering algorithm for data streams that is based on the idea of the **k-means++ seeding procedure**.

- Their streaming algorithm maintains a small sketch of the input using the **merge-and-reduce** technique, i.e., the data is organized in a small number of samples, each representing 2^m input points (for some integer m and a fixed value n). Everytime when two samples representing the same number of input points exist they take the union (**merge**) and create a new sample (**reduce**).
- For the reduce step, they propose a new **coreset construction**. Here, we focus on giving a construction that is suitable for high-dimensional data. Since the k-means++ seeding works well for high-dimensional data, a coreset construction based on this approach seems to be more promising. In order to implement this approach efficiently, they develop a new data structure, which they call the **coreset tree**.
- They prove that, with high probability, **sampling according to the k-means++ seeding procedure gives small coresets, at least in low dimensional spaces**.
- A **coreset** for a set P is a small (weighted) set, such that for any set of k cluster centers the (weighted) clustering cost of the coreset is an approximation for the clustering cost of the original set P with small relative error. The advantage of such a coreset is that we can apply any fast approximation algorithm (for the weighted problem) on the usually much smaller coreset to compute an approximate solution for the original set P more efficiently.
- Their clustering algorithm maintains a small coreset in the data streaming model. Applies only one sequential scan over the data.
- ALGORITHM: First, they extract a small coreset of size m from the data stream by using the merge-and-reduce technique. For the reduce step, they employ our new coreset construction, using the coreset. After that, a k -clustering can be obtained at any point of time by running any k -means algorithm on the coreset of size m because the size of the coreset is much smaller than the size of the data stream, and they are no longer prohibited from algorithms that require random access on their input data. In their implementation, they run the k -means++ algorithm on their coreset five times independently and choose the best clustering result obtained this way.
- Its running time is still far too large for big datasets, especially for large k .
- Comparison with StreamLS (3), BIRCH (2) and non-streaming version of k -MEANS++:
 - StreamLS is significantly outperformed by Stream-KM++ which computes a coreset and then solves the k -means problem on the coreset by applying k -means++.
 - StreamLS uses a batch approach, and Stream-KM++ uses merge & reduce.
 - In terms of quality (sum of squared errors), their algorithm is comparable with StreamLS and significantly better than BIRCH (up to a factor of 2). In terms of running time, their algorithm is slower than BIRCH. Comparing the running time with StreamLS,

it turns out that their algorithm scales much better with increasing number of centers. We conclude that, if the first priority is the quality of the clustering, then their algorithm provides a good alternative to BIRCH and StreamLS, in particular, if the number of cluster centers is large.

- It turns out that their algorithm is slower than BIRCH, but it computes significantly better solutions (in terms of sum of squared errors). In addition, to obtain the desired number of clusters, our algorithm does not require the trial-and-error adjustment of parameters as BIRCH does.
- The quality of the clustering of algorithm StreamLS is comparable to that of their algorithm, but the running time of StreamKM++ scales much better with the number of cluster centers. While for about $k < 10$ centers StreamLS is sometimes faster than their algorithm, for a larger number of centers their algorithm easily outperforms StreamLS.
- In comparison with the standard implementation of k-means++, our algorithm runs much faster on larger datasets and computes solutions that are on a par with k-means++.
- Comparing k-means++ with our streaming algorithm, we find that on all datasets the quality of the clusterings computed by algorithm StreamKM++ is on a par with or even better than the clusterings obtained by algorithm k-means++. They conjecture that this is due to the fact that in the last step of our algorithm they run the k-means++ algorithm five times on the coreset and choose the best clustering result obtained this way.
- Algorithm k-means++ should only be used if the size of the dataset is not too large. For larger datasets, algorithm StreamKM++ computes comparable clusterings in a significantly improved running time.

6. BICO: BIRCH meets Coresets for k-means clustering* (2013):

- A data stream algorithm for the k-means problem that combines the data structure of the SIGMOD Test of Time award winning algorithm **BIRCH** with the theoretical concept of **coresets** for clustering problems.
- A **coreset** is a small weighted set of points S , that ensures that if we compute the weighted clustering cost of S for any given set of centers C , then the result will be a $(1 + \epsilon)$ -approximation of the cost of the original input.
- It computes high quality solutions in a time short in practice.
- First, BICO computes a summary S of the data with a provable quality guarantee: For every center set C , S has the same cost as P up to a $(1 + \epsilon)$ -factor, i. e., S is a coreset. Then, it runs k-means++ on S .
- This paper contributes to the field of interlacing theoretical and practical work to develop an algorithm good in theory and practice.

- The main problem with the insertion procedure of BIRCH is that the decision whether points are added to a CF or not is based on the increase of the radius of the candidate CF.
- Like BIRCH, BICO uses a tree whose nodes correspond to CFs.
- Algorithm 1.
- The aim of the rebuilding algorithm is to create a tree which is similar to the tree which would have resulted from using the new threshold in the first place.
- Comparison with StreamKM++ (4) and StreamLS (3) (approximation algorithms, designed for high quality solutions) and comparison with BIRCH (2) and MacQueen's proposal (very fast heuristics):
 - They achieve the same quality as the approximation algorithms mentioned with a much shorter running time, and they get much better solutions than the heuristics at the cost of only a moderate increase in running time.
 - Approximation algorithms are rather slow in practice. Algorithms fast in practice are usually heuristics and known to compute bad solutions on occasions. The best known one is BIRCH. It also computes a summary of the data, but without theoretical quality guarantee.
 - Stream-KM++ computes a coreset and then solves the k-means problem on the coreset by applying k-means++. It computes very good solutions and is reasonably fast for small k. However, especially for large k, its running time is still far too large for big data sets. BICO also computes a coreset and achieves very good quality in practice, but is significantly faster than Stream-KM++.
 - BIRCH and BICO algorithms compute a summary of the data, but while the summary computed by BIRCH can be arbitrarily bad, they show that BICO computes a coreset S , so for every set of centers C , the cost of the input point set P can be approximated by computing the cost of S .
 - For small k, BICO's running time is only beaten by MacQueen's and in particular, BICO is 5-10 times faster than Stream-KM++ (and more for StreamLS). For larger k, BICO needs to maintain a larger coreset to keep the quality up. However, BICO can trade quality for speed. They do additional testruns showing that with different parameters, BICO still beats the cost of MacQueen and BIRCH in similar running time. They believe that BICO provides the best quality-speed trade-off in practice.
 - BIRCH decides heuristically how to group the points into subclusters. We aim at refining this process such that the reduced set is not only small but is also guaranteed to approximate the original point set. More precisely, we aim at constructing a coreset.
 - Stream-KM++ also aims at a trade-off between quality and speed which makes it most relevant for our work. BIRCH is the most relevant practical algorithm.

- BIRCH and MacQueen both tend to compute much worse solutions.
- BICO is still practical for large k despite the large running times when adjusted.