# k Nearest Neighbour Efficiency
# IT 494 PROJECT

Asish Joel - 202103015
Aaditya Meher - 202103024
Raj Kariya - 202103048

November 2023

# Contents

# 1    Introduction

In recent years, technological advancements have automated and made information gathering cost-effective, leading to a surge in available data from various fields. This abundance poses a challenge for conventional machine learning methods, necessitating updates to handle the volume, diversity, and complexity of the data. The k-Nearest Neighbor algorithm (kNN), a nonparametric model for classification and regression, faces scalability issues in the big data context. This study focuses on classification tasks and explores partial and approximate variants of kNN to address computational challenges, considering approximation error bounds for distance computations. The need for adapting existing learning techniques to cope with substantial data volumes is emphasized.

In response to the escalating volume of available data, recent cloud-based technologies offer a solution, with the MapReduce framework in Hadoop being an early tool for data-intensive applications. However, limitations in Hadoop MapReduce, especially for iterative algorithms, led to the emergence of platforms like Spark, known for faster distributed computing using in-memory primitives. The k-Nearest Neighbor algorithm (kNN) has been explored in the big data context, with some works incorporating kNN in MapReduce, often for partial kNN applications. This paper proposes an iterative MapReduce-based approach for kNN implemented under Apache Spark, leveraging Spark's flexibility and in-memory operations to alleviate consumption costs. The method iteratively addresses chunks of enormous test sets, performing a kNN MapReduce process in each iteration, emitting class labels and distance values. The approach, referred to as kNN with Spark (kNN-IS), aims to efficiently classify large test samples against extensive training datasets, combining scalability with in-memory solutions to enhance performance.

# 2    Understanding the Algorithm

## 2.1    MapReduce

"Large datasets can be processed in parallel and distributed using the MapReduce programming language and processing framework in a Hadoop cluster.Since its introduction by Google in 2003, MapReduce has been a fundamental component of big data processing.The Map function consists of two fundamental parts: Map and Reduce. It takes raw data in the form of key-value pairs and converts it into intermediate pairs, which may be of different kinds. The description of key and value types is part of the user-defined specifications. After then, during the shuffle stage of MapReduce, values connected to the same intermediate key are combined into a list.It enables users to build programs that, by dividing up the burden among several nodes, can process enormous volumes of data.The entire procedure, as seen in, takes place on a scalable, transparent platform that processes data in distributed clusters, alleviating users from technical intricacies such as data transformation, analysis, and pattern recognition in big data applications, making it a fundamental tool in Hadoop's ecosystem for efficiently handling large-scale data processing tasks."[7]
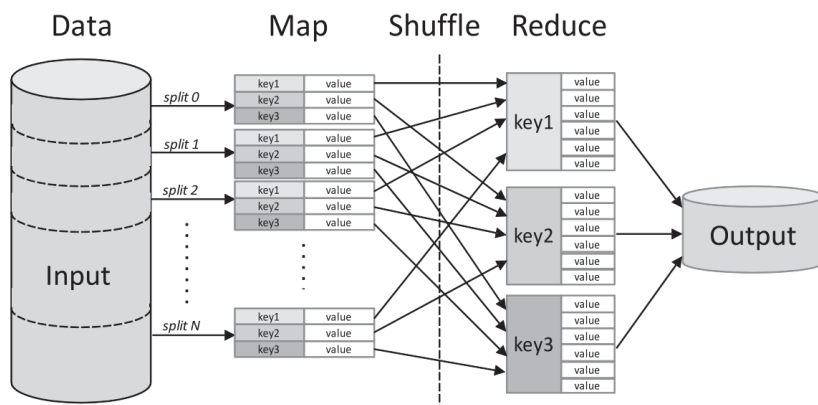


Figure 1: Data Flow Overview of Graph

## 2.2   Hadoop

So, now what exactly is hadoop? Apache Hadoop is an open-source Large datasets can be processed and stored distributed using the Hadoop open-source framework. It is made up of two primary parts: the MapReduce programming model for parallel data processing and the Hadoop Distributed File System (HDFS) for distributed storage. This framework's performance, open source nature, and ease of installation have led to its extensive adoption in numerous fields. Big data can be handled in a scalable and fault-tolerant manner on clusters of commodity technology thanks to Hadoop. It is becoming a vital tool for businesses handling complex data processing and analytics jobs.Despite their widespread use, Hadoop and MapReduce have proven to be ill-suited for many applications, such as online or iterative computation. Its incapacity to recycle data by its inability to reuse data through in-memory primitives makes the application of Hadoop for many machine learning algorithms unfeasible.

## 2.3   Spark

Spark is based on Resilient Distributed Datasets (RDDs), a special type of data structure used to parallelize the computations in a transparent way. Unlike its predecessor, MapReduce (MR), Spark stands out for its in-memory processing capability, which significantly accelerates iterative algorithms such as machine learning and graph processing. Spark's Resilient Distributed Datasets (RDDs) enable fault-tolerant parallel processing, improving overall performance compared to the disk-based approach of MR.Moreover, they also let us manage the partitioning to optimize data placement, and manipulate data using a wide set of transparent primitives. All these features allow users to easily design new data processing pipelines.

## 2.4   HDFS

An essential part of the open-source Hadoop system, the Hadoop Distributed File System (HDFS) is made for distributed storage on reasonably priced hardware. High-throughput data access is made possible by HDFS, an efficient and fault-tolerant system that is especially well-suited for applications utilizing massive datasets in Apache Hadoop. HDFS acts as the file system inside Hadoop, providing access to stored data, in contrast to Hadoop, which is a complete framework for storing, processing, and analyzing data.

HDFS is centered on NameNodes and DataNodes in terms of architecture. As the master server, the NameNode oversees file operations, manages files, and regulates client access. However, DataNodes—which are found on every HDFS cluster node—manage data storage functions like block management, file creation, and replication in response to NameNode commands.

Using hundreds of nodes per cluster, HDFS facilitates the efficient handling of massive datasets, one of its many important purposes. It is particularly good at fault detection, which is important for managing the frequent failures of components in a distributed system. HDFS also improves hardware economy by reducing network traffic and maximizing processing speed—two important factors to take into account when working with large datasets.

## 2.5   KNN - K Nearest Neighbours

"In the k-nearest neighbors (KNN) algorithm, a non-parametric and supervised learning algorithm used for classification and regression tasks.It operates on the principle of proximity, where the classification or prediction of a data point is determined by the majority class or average value of its k-nearest neighbors in the feature space. In classification, the algorithm assigns a data point to the class most common among its k-nearest neighbors, while in regression, it predicts a continuous value based on the average of the target values of those neighbors.

The algorithm works by calculating the distances between the target data point and all other points in the dataset, often using metrics like Euclidean distance. The 'k' in KNN represents the number of neighbors considered. A smaller k value leads to a more sensitive model, while a larger k value provides a smoother decision boundary but might be less sensitive to local variations.

KNN's performance can be influenced by the choice of k and the distance metric, and it may not scale well with large datasets. Regularization techniques and feature scaling can be applied to enhance its effectiveness in various applications. "[8]

### 2.5.1 Working of KNN

"Consider a training dataset (TR) and a test set (TS), both comprising a specified number of samples (n for TR, t for TS). Each sample, denoted as xp, is a tuple containing D features (xp1, xp2, ..., xpD,) where, xp f is the value of the f-th feature of the p-th sample and belongs to a class represented by xp in a D-dimensional space. While the class is known for TR, it remains unknown for TS. In the k-Nearest Neighbors (kNN) algorithm, applied to each xtest in TS, the goal is to find the k closest samples in TR. By calculating Euclidean distances between xtest and all TR samples, the algorithm ranks training samples in ascending order based on these distances. The top k neighbors (neigh1, neigh2, ..., neighk) are then used to determine the most prevalent class label. The choice of k significantly impacts the algorithm's performance and noise tolerance."[3]

### 2.5.2 KNN-advantages in big data

"While the k-Nearest Neighbors (kNN) algorithm has demonstrated exceptional performance across diverse problems, it faces challenges in scalability when dealing with large training (TR) datasets. The primary issues encountered in managing extensive data sets include:

- **Runtime Complexity:** The time complexity to identify the nearest neighbor training example for a single test instance is $O((n \cdot D))$, where 'n' is the number of training instances and 'D' is the number of features. This complexity increases when seeking the k closest neighbors, involving the additional computational expense of sorting computed distances ($O(n \cdot \log(n))$). This process must be reiterated for each test example.

- **Memory Consumption:** To swiftly compute distances, the kNN model necessitates storing the training data in memory. However, when both TR and the TS sets are substantial, they may surpass the available RAM memory, posing challenges in terms of memory consumption.

These limitations drive the exploration of big data technologies to distribute kNN processing across a node cluster. Within existing literature, various approaches exist for executing kNN joins through MapReduce. It's crucial to note that kNN joins differ from kNN classifiers in their anticipated outputs. While a kNN classifier seeks to predict the class, a kNN join outputs the neighbors themselves for a given test instance. Consequently, these methods are not suitable for classification tasks. "[7]

## 2.6 KNN-MR

"For an exact kNN join, two primary alternatives exist. The first, named H-BkNNJ, employs a single round of MapReduce where TR and TS sets are jointly partitioned. Each map task processes a pair TSi and TRi, executing pairwise distance comparisons between training and test splits. With 'm' as the number of used partitions, it generates $m^2$ blocks through a linear scan on both sets. The subsequent reduce task processes computed distances for a given test instance, sorting them in ascending order to output the top k results. The second alternative, H-BNLJ, uses two MapReduce processes to reduce the complexity of the sort phase but still necessitates $m^2$ tasks. Deficiencies include the generation of extra data blocks, exacerbating problem size, quadratic complexity ($m^2$ tasks), and reliance on Hadoop MapReduce, incurring serialization costs. Some models like PGBJ perform preprocessing and distance-based partitioning to reduce tasks to 'm,' albeit introducing additional computational costs. A more recent alternative, Spitfire, employs a distributed procedure distinct from MapReduce. It calculates k nearest neighbors for all set elements by partitioning the search space, replicating the k nearest neighbors in each split, and conducting a final local kNN computation for the ultimate result.

In the context of classification tasks, which also applies to regression, existing methods are more straightforward compared to kNN join approaches, as they don't require providing the neighbors themselves—only their classes. Two main approaches have been introduced, both centered on leveraging the Map phase to partition the training data into 'm' disjoint parts. The first approach, advocates iteratively repeating a MapReduce process (without an explicitly defined reduce function) for each individual test example. However, this proves to be time-consuming in both Hadoop and Spark, as further elaborated in the experiment section. The second, denoted as MR-kNN, employs a single MapReduce process to handle the classification of the sizable test set. Hadoop primitives are utilized to read the test data line by line within the map phase, rendering the model scalable, albeit with potential performance enhancements achievable through in-memory solutions."[7]

## 2.7 KNN- Iterative Spark Design of the KNN classifier

"In this part, we introduce a different way to do big data classification using Spark called KNN-IS.Our main goal in this model is to make the KNN classifier run faster, especially when both training dataset and test datasets are large.

When we are implementing KNN in a parallel system like Spark, there are several factors such as number of map-reduce jobs(denoted by j) or number of Map tasks(denoted by m) and number of Reduce tasks(denoted by r) required, etc affects the Execution time.Therefore to obtain an efficient and accurate modelis a difficult task."[1]

To address the challenges posed by MapReduce Solution for KNN.We will introduce new ideas in this model as given below:

- "A MapReduce process divides the training dataset into m tasks.Unlike KNN-join methods that requires $m^2$ tasks, we have reduced it to m tasks without any preprocessing.

- To handle large test datasets, we use spark to recycle the earlier training dataset with different chunks of the test dataset.Using multiple MR jobs over the same data does not imply significant extra costs in Spark, but we aim to minimize it.Now the MR-KNN approach only needs to do m tasks, regardless of test data size by reading the test dataset line by line within the maps.

- In spark, RDD are the specific type of data structures which are used to parallelize the computation in a transparent way.Every operation will be performed within the RDD objects which makes operations such as normalization,etc scalable to handle large amounts of data."[2]

### 2.7.1 MapReduce for KNN classification within Spark

In this part, we discuss about a MapReduce process which manages the classification of subsets of test data against the whole training dataset.This process is builded upon earlier suggestes MR-KNN method.This process allows use for multiple reducers, monitor the iterations to prevent swaps, and its specifically designed to work with Spark framework.

"This model divides the computation into 2 phases: map phase and reduce phase. In map phase, the training data is divided, and for each chunk the distances and the corresponding classes of the k nearest neighbors are calculated for every test sample. In the reduce stage, the distance of the k nearest neighbors from each map are combined and we create a list of k nearest neighbors.Lastly it follows the typical majority voting procedure of the KNN Algorithm to predict the resulting class."[3]

### 2.7.2 Map Phase

Let's say we have loaded the training set(TR) and a specific subset of test dataset(TS) which are stored as RDD objects after reading it from HDFS.When we read the Training dataset, we broke it into disjoint subsets, each defined by the user, denoted as m.Each part of the Training dataset will be handled by a different task like Map1, Map2, Map3, so on upto Map m.Each tasks will deal against the specific subset of training set, labelled as $TRj$ where $1 \leq j \leq m$. Therefore each map processes a similar number of training instances.

---

**Algorithm 1** Map Phase

**Require:** $TR_j, TS_i$; k

    1. for t=0 to size($TS_i$) do

    2. $CD_{t,j} \leftarrow ComputeKNN(TR_j, TS_i(x), k)$

    3. $result_j \leftarrow (\langle key : t, value : CD_{t,j} \rangle)$

    4. EMIT($result_j$)

---

For each test sample in test dataset (denoted by $TS_i$), every map $j$ will constitute a class-distance vector $CD_{t,j}$ of pairs $\langle class, distance \rangle$ of dimension k.To speed up the later update of the nearest neighbors in the reducers, each class-distance vector(CD) vector have been arranged in the increasing order based on the distance to the sample.

The proposed method will send a key-value pair where key is a test instance identifier (denoted by t) for each output and value is vector $CD_{t,j}$.We enable the method to utilize multiple reducers. Whenever we deal with large

datasets, this method allows the use of multiple reducers which will improve the performance and scalability of the algorithm.

### 2.7.3 Reduce by Key Operation

In reduce phase, we have to collect the closest ones based on the distance from the k nearest neighbors which are provided by map phase. Reduce by key operation will group all the elements which have same key after the completion of map phase. A reducer is run over a list($CD_{t,0}, CD_{t,1}, \ldots, CD_{t,m}$) and determines the k nearest neighbors of this test example t.

---

**Algorithm 2** Reduce by key Operation

---

**Require:** $result_{\text{key}}, k$
1: $cont = 0$
2: **for** $i = 0$ to $k$ **do**
3:     **if** $result_{\text{key}}(cont) < result_{\text{reducer}}(i).Dist$ **then**
4:         $result_{\text{reducer}}(i) = result_{\text{key}}(cont)$
5:         $cont + +$

---

"This process goes through each element of such list one after another.In the Map phase vector are generated and are sorted based distance which allows the update faster during the subsequent phases in the Reduce phase. It involves combining two sorted list to obtain the k closest values, making the worst time complexity $O(k)$.In Spark, this transformation is called Reduce by key operation."[4]
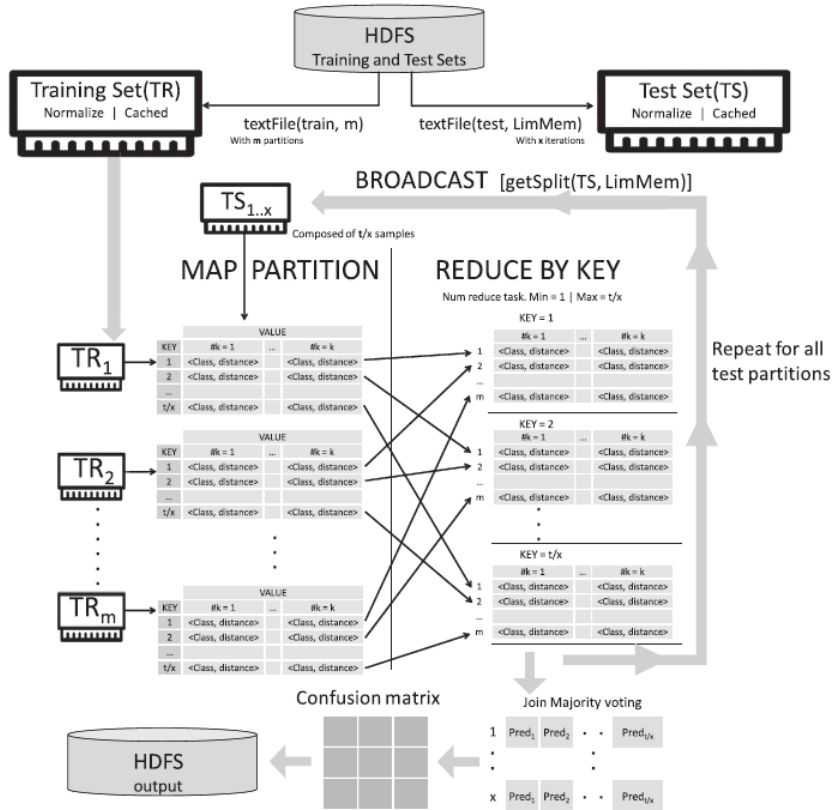
## 2.8 Scheme of KNN-IS



Figure 2: Flowchart of the proposed kNN-IS algorithm.

If the test dataset is too big, then it might not fit into the memory allocated for the map tasks.When this type of situation occurs we have to divide test data into smaller parts and run MapReduce Process multiple times.

---

**Algorithm 3** KNN

---
1: TR - $RDD_{raw}$ ← textfile(TR, #Maps)
2: TS - $RDD_{raw}$ ← textfile(TS).zipWithIndex()
3: TR - $RDD$ ← TR - $RDD_{raw}$.map(normalize).cache
4: TS - $RDD$ ← TS - $RDD_{raw}$.map(normalize).cache
5: #Iter ← $callter(TR - RDD.weight(), TS - RDD.weight, MemAllow)$
6: **for** $i = 0$ to #Iter **do**
7:   $TS_i$ ← broadcast(TS - RDD.getSplit(i))
8:   $resultKNN \leftarrow TR - RDD.mapPartition(TR_j \rightarrow KNN(TR_j, TS_i, k))$
9:   $result \leftarrow resultKNN.reduceByKey(combineResult, \#Reduces).collect$
10:   right-predictedClasses[i] ← calculateRightPredicted(result)
11: **end for**
12: cm ← $calculateConfusionMatrix(right - predictedClasses)$

---

We will start by getting the location of training and test dataset in the HDFS. We also have information about the number of maps(m), reducers(r), the number of neighbors(k) and how much memory each map would use.

We create a Resilent Distributed Dataset for the training set(TR) which is divided into m blocks. The test dataset(TS) is also read as an RDD, but don't specify the number of partitions. The function named zipwithindex() provides a unique key to each test instance based on its position in the dataset.

Normally we use euclidean distances to compute the similarity between instances, we also need to normalize both datasets.This means we are going to adjust the values within specific range, in this case [0, 1]. We do this in parallel for both dataset.

Later on, both datasets are store in cache to quickly access. This is important because it allows us to reuse the data without having to read it from disk again."[5]

In this process, Spark is used to handle data iteratively. However, to improve performance, the number of iterations is minimized. The minimum number of iterations needed is calculated using the size of the training dataset chunks, the size of the test set, and the memory limit for each map. This calculated number of iterations is then used to split the test dataset into subsets with a similar number of samples. The keys previously set in the test set are used for this purpose. Finally, the test dataset is partitioned using the RangePartitioner function.

"The algorithm classifies subsets of the test set in a loop. It first gets the current iteration's split and uses the filterByRange function to efficiently select the corresponding subset. This subset is then broadcasted into the memory of all computing nodes. The main map phase starts next, where the mapPartition function computes the k-nearest neighbors for each partition of the training and test sets and emits a pair RDD. The reduce phase then groups the results by key. As a result, we get the k-nearest neighbors and their classes for each test input."[6]

The last step in the loop collects the right and predicted classes stroing them as an array in every iteration. When loop gets compelted, we compute the confusion matrix and outputs the desired performance measures.

## 2.9 Analysis of Results

### 2.9.1 Comparison with MR-kNN

This section compares MR-kNN, which is thought to be a possibly faster alternative, with kNN-IS. The Susy and PokerHand datasets are used in the evaluation. Owing to constraints, the investigation was unable to achieve findings for the sequential kNN using datasets other than these ones. Since the test datasets fit into the memory of each map task, kNN-IS applied to these datasets only requires one iteration. To create a comparison between closely related MapReduce alternatives, specifically under Hadoop (MR-kNN) and Spark (kNN-IS), the number of reducers in kNN-IS is fixed at 1.

These datasets are first processed using the sequential version of kNN in order to create a baseline. This sequential variant mirrors the MR-kNN strategy by reading the test set line by line, similar to a technique described in [28]. This is a simple fix meant to prevent memory problems. It's crucial to remember that this situation is the worst that can happen with the sequential version. Assuming that big test sets cannot fit into memory together with the training set, the goal here is to compare with the simplest sequential version, even though there may be more sophisticated versions.

Table 5 displays the average accuracy (AccTest) and runtime (in seconds) of the typical kNN algorithm as a

**Table 5**
Sequential kNN performance.

| Dataset | Number of Neighbors | Runtime(s) | AccTest |
|---|---|---|---|
| PokerHand | 1 | 105475.0060 | 0.5019 |
| | 3 | 105507.8470 | 0.4959 |
| | 5 | 105677.1990 | 0.5280 |
| | 7 | 107735.0380 | 0.5386 |
| Susy | 1 | 3258848.8114 | 0.6936 |
| | 3 | 3259619.4959 | 0.7239 |
| | 5 | 3265185.9036 | 0.7338 |
| | 7 | 3325338.1457 | 0.7379 |

Figure 3: Sequential kNN performance.

**Table 6**
Results obtained by MR-kNN and kNN-IS algorithms in PokerHand dataset.

| Dataset | #Map | MR-kNN | | kNN-IS | |
|---|---|---|---|---|---|
| | | AvgRunTime | Speedup | AvgRuntime | Speedup |
| PokerHand | 128 | 804.4560 | 131.1135 | 102.9380 | 1024.6460 |
| | 64 | 1470.9524 | 71.7052 | 179.2381 | 588.4631 |
| | 32 | 3003.3630 | 35.1190 | 327.5347 | 322.0270 |
| | 256 | 12367.9657 | 263.4911 | 1900.0393 | 1715.1481 |
| Susy | 128 | 26438.5201 | 123.2614 | 3163.9710 | 1029.9869 |
| | 64 | 50417.4493 | 64.6373 | 6332.8108 | 514.5975 |

Figure 4: Results obtained by MR-kNN and kNN-IS algorithms in PokerHand dataset.

**Table 7**
Results obtained with Susy dataset.

| k | #Map | AvgMapTime | AvgRedTime | AvgTotalTime |
|---|---|---|---|---|
| 1 | 512 | 730.3893 | 560.4334 | 2042.2533 |
| | 256 | 1531.6145 | 345.5664 | 1900.0393 |
| | 128 | 2975.3013 | 166.3976 | 3163.9710 |
| | 64 | 6210.6177 | 92.8188 | 6332.8108 |
| 3 | 512 | 770.3924 | 736.8489 | 2298.3384 |
| | 256 | 1553.4222 | 410.2235 | 2615.0150 |
| | 128 | 3641.9363 | 253.5656 | 3921.3640 |
| | 64 | 6405.3132 | 152.9890 | 6593.5531 |
| 5 | 512 | 781.3855 | 928.2620 | 2511.8909 |
| | 256 | 1773.3579 | 479.0801 | 2273.6377 |
| | 128 | 3685.3194 | 332.9783 | 4042.1755 |
| | 64 | 6582.0373 | 194.7054 | 6802.8159 |
| 7 | 512 | 782.5756 | 930.5107 | 2516.5011 |
| | 256 | 1827.9189 | 522.6219 | 2372.4100 |
| | 128 | 3401.2547 | 414.2961 | 3838.2360 |
| | 64 | 6637.8837 | 224.7191 | 6890.8242 |

Figure 5: Results obtained with Susy dataset.

function of neighbor count. The outcomes of both approaches are summarized in Table 6 for k=1. The impact of k's value will be discussed in detail in the section that follows. It shows the speed improvement over the sequential version and the average total time (AvgRuntime) for each number of maps (Maps). It is important to note that

8

both approaches guarantee precisely the same average accuracy as shown in Table 5 and correspond to an accurate implementation of the kNN.

As the number of maps increases (k=1), Figure 3 shows speedup comparisons between the two approaches and the sequential version. The following analysis is possible using the data shown in the tables and figures:

The following conclusions can be drawn from the analysis of the tables and figures:

High Runtime in Sequential kNN: Table 5 shows that, for both datasets, the runtime required for the sequential version of the kNN algorithm is significantly high.

Significant Runtime decrease: As the number of maps increases, Table 6 shows a significant runtime decrease for both approaches. This is a significant improvement that demonstrates the effectiveness of parallelization.

Regular Accuracy: As previously indicated, both options regularly offer the same degree of precision as the sequential version. This highlights how the parallelized approaches are successful without sacrificing precision.

As per Figure 3: Hadoop-based kNN with Linear Speedup: Considering that both models read the test dataset line-by-line, the Hadoop-based kNN model exhibits a linear speedup. The reason for the speedup, which is sometimes even superlinear, is that the previous kNN model had problems with memory usage when handling the training set.

Faster Speedup in kNN-IS: With regard to this sequential version, kNN-IS delivers a faster speedup than linear. This is explained by the use of in-memory data structures, which do away with the requirement to read the test data from HDFS line by line.

A comparison of kNN-IS and MR-kNN The findings demonstrate that Spark (kNN-IS) has slashed the necessary runtime by almost ten times when compared to Hadoop (MR-kNN). This highlights the increases in productivity brought about by using Spark.
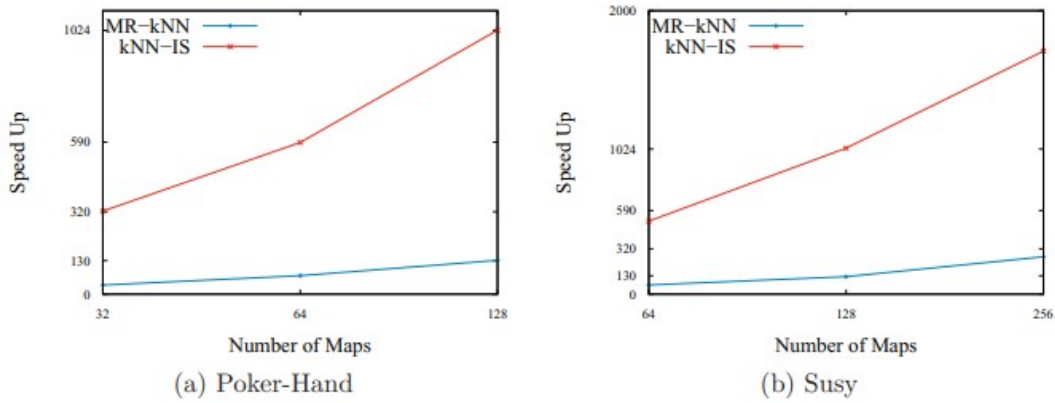


Figure 6: Speedup comparisons between MR-kNN and kNN-IS against the sequential kNN.

### 2.9.2 Influence of neighbors

To thoroughly examine the impact of the number of neighbors, we concentrate on the Susy dataset, maintaining the number of reducer tasks at one. We assess its effects on both the map and reduce phases. Table 7 compiles, for each number of neighbors (#Neigh) and number of maps (#Maps), the average execution time for maps (AvgMap-Time), the average execution time for reduces (AvgRedTime), and the average total runtime (AvgTotalTime). It's important to note that in our cluster, the maximum number of map tasks is capped at 256. Consequently, the total runtime for 512 maps may not exhibit a linear reduction, but the reduction in map runtime is evident. Figure 4 illustrates how the value of k influences map runtimes, displaying the map runtimes concerning the number of maps for k=1, 3, 5, and 7. Additionally, Figure 5 depicts the reduce runtime in relation to the k value and the number of maps.

Following are list of the observations to these tables and plots,

1. Larger values of k do not significantly alter the total runtimes, even if they suggest that more data is transmitted from the maps to the reducers. Table 7 shows that the total runtime generally increases by a small

amount.

2. By contrasting Figs. 4 and 5, we can observe that the reduction runtime appears to be more affected by the number of neighbors than the map phase. This is because the number of neighbors has no bearing on the map phase's primary calculation cost (calculating the distances between test and training instances), but it might have an impact on the reducers' updating procedure because of its $O(k)$ complexity.

In summary, Figure 5 offers a broad understanding that the reduction runtimes significantly increase when more maps are used, which is frequently required to manage huge datasets. The investigation carried out in the next part is motivated by this observation.

# 3   Conclusion

1. **Efficiency and Accuracy:** The developed Iterative MapReduce solution, kNN-IS, achieves the same accuracy as the traditional k-Nearest Neighbors (kNN) algorithm while effectively addressing runtime and memory consumption challenges associated with large-scale datasets.

2. **Spark's Impact:** Leveraging Apache Spark, kNN-IS presents a straightforward and efficient environment for parallelizing the kNN algorithm through an iterative MapReduce process. This choice of technology enables a transparent and simple implementation.

3. **Competitive Runtime Performance:** The experimental study demonstrates that kNN-IS exhibits highly competitive runtime performance across datasets of varying sizes, considering different numbers of features and samples. It outperforms the traditional kNN algorithm in terms of runtime efficiency.

4. **Significant Speedup with Spark:** kNN-IS, implemented in Spark, substantially reduces the runtime by nearly 10 times compared to MR-kNN implemented in Hadoop. This highlights the effectiveness of using Spark for large-scale kNN computations.

5. **Flexibility and Optimization:** The flexibility of adjusting the number of maps and reducers in kNN-IS allows for runtime optimization. Despite potential data transfer challenges between map and reduce phases, the number of neighbors (k) does not drastically impact the overall runtime. This suggests a feasible trade-off between computational efficiency and dataset-specific configurations, enhancing the adaptability of kNN-IS.

# 4   Future Work

In future work, the optimization and scalability of the kNN-IS algorithm can be explored further to enhance its performance on extremely large datasets. Investigating adaptive strategies for dynamically adjusting the number of maps and reducers based on dataset characteristics could provide insights into achieving even more efficient computations. Additionally, the exploration of advanced distance or similarity measures may contribute to refining the accuracy of the kNN-IS algorithm. Integration with emerging technologies, such as distributed computing frameworks beyond Spark, can also be considered to ensure adaptability to evolving computing landscapes. Furthermore, research into extending kNN-IS for applications beyond classification tasks may open avenues for broader and more versatile data analytics solutions.

# 5 Refernces

1. [1] https://www.nature.com/articles/455028a

2. [2] https://www.google.co.in/books/edition/Big_Data_Big_Analytics/HYYaX9dsZsYC?hl=en&gbpv=1&pg=PR13&printsec=frontcover

3. [3] https://scholar.google.com/scholar_lookup?title=Nearest%20neighbor%20pattern%20classification&author=T.M.%20Cover&publication_year=1967&pages=21-27

4. [4] https://scholar.google.com/scholar_lookup?title=Using%20MPI%3A%20portable%20parallel%20programming%20with%20the%20message-passing%20interface&author=W.%20Gropp&publication_year=1999

5. [5] http://refhub.elsevier.com/S0950-7051(16)30175-7/sbref0046

6. [6] https://scholar.google.com/scholar_lookup?title=Graphlab%3A%20a%20new%20parallel%20framework%20for%20machine%20learning&author=Y.%20Low&publication_year=2010

7. [7] https://dl.acm.org/doi/10.14778/1920841.1920881

8. [8] https://www.sciencedirect.com/journal/knowledge-based-systems