

Inclusion Dependencies Identification using Spark

Department of Computer Science and Mathematics
Philipps-Universität Marburg

Course: Big Data Systems
Instructor: Prof. Dr. Thorsten Papenbrock

submitted by Nikita Galaganov (3873245) and Erik Maiterth (3695656)

February 11, 2025

Contents

1	Introduction	2
2	System Design and Architecture	2
3	Sindy Algorithm Overview	2
4	Implementation Details	3
5	Advantages Compared to Alternative Approaches	4

1 Introduction

Ensuring data integrity in large-scale relational databases is crucial. This report presents a system designed to discover Inclusion Dependencies (INDs) using **Apache Spark**, a distributed computing framework. The system processes large CSV datasets to efficiently identify column dependencies.

2 System Design and Architecture

The system follows a structured pipeline for IND discovery:

1. **Data Reading:** CSV files are parsed into Spark DataFrames.
2. **Data Transformation:** Cell values are converted into tuples of the form `(cellValue, Set(columnName))`.
3. **Dependency Generation:** Grouping by cell value merges the column sets, generating candidate dependencies.
4. **Dependency Reduction:** For each column, the intersection of candidate dependency sets is computed.
5. **Result Output:** The final INDs are collected, sorted, and displayed.

3 Sindy Algorithm Overview

The Sindy algorithm is a scalable method for discovering inclusion dependencies in large datasets. It leverages Apache Spark's distributed processing capabilities to identify relationships between columns based on overlapping cell values. In essence, the algorithm operates as follows:

- **Extraction:** Each cell value is paired with its corresponding column name, forming tuples of the type `(cellValue, Set(columnName))`.
- **Candidate Generation:** By grouping these tuples by cell value, the algorithm determines the set of columns in which a particular value appears. For each value, candidate dependency pairs are generated by subtracting the current column from the set.
- **Dependency Reduction:** The algorithm computes the intersection of candidate dependency sets for each column. A non-empty intersection indicates that every value in a column is also present in the intersected columns, establishing an inclusion dependency.

This approach efficiently uncovers dependencies such as `Column_A < Column_B, Column_C`, meaning every value in `Column_A` is also found in `Column_B` and `Column_C`.

4 Implementation Details

The core function `discoverINDs` implements the IND discovery pipeline. The following code listing shows the main steps of the implementation.

```
1 def discoverINDs(inputs: List[String], spark: SparkSession): Unit = {
2   import spark.implicits._
3
4   // Step 1: Read CSV files and convert each row
5   // into (cellValue, Set(columnName)) tuples
6   val extractedData = inputs
7     .map { input =>
8       spark.read
9         .option("inferSchema", "false")
10        .option("header", "true")
11        .option("quote", "\"")
12        .option("delimiter", ";")
13        .csv(input)
14        .flatMap { row =>
15          // For each row and column, produce (cellValue, Set(
16            columnName))
17          row.schema.fieldNames.map { colName =>
18            (row.getAs[Any](colName).toString, Set(colName))
19          }
20        }
21    }.reduce(_ union _)
22
23   // Step 2: Group by the cellValue and combine columns,
24   // yielding (column, otherColumns) dependency pairs
25   val dependencyCandidates = extractedData
26     .groupByKey { case (cellValue, _) => cellValue }
27     .mapGroups { case (_, tuples) =>
28       val mergedCols = tuples.flatMap(_._2).toSet
29       mergedCols.map(col => (col, mergedCols - col))
30     }
31     .flatMap(identity)
32
33   // Step 3: For each column, intersect all candidate sets,
34   // leaving only columns that appear in every set
35   val indCandidates = dependencyCandidates
36     .groupByKey { case (col, _) => col }
37     .mapGroups { case (col, depSets) =>
38       (col, depSets.map(_._2).reduce(_ intersect _))
39     }
40     .filter(_._2.nonEmpty)
41
42   // Step 4: Collect, sort, and print the final INDs
43   val finalINDs = indCandidates
44     .collect()
45     .map { case (col, deps) => (col, deps.toList.sorted) }
46     .sortBy(_._1)
47
48   finalINDs.foreach {
49     case (col, deps) => println(s"$col < ${deps.mkString(", ")")
50   }
51 }
```

5 Advantages Compared to Alternative Approaches

While techniques such as the Longest Common Substring algorithm are useful for string similarity tasks, the Sindy approach offers several advantages in the context of inclusion dependency discovery:

- **Scalability:** Sindy is designed to run on Apache Spark, enabling efficient distributed processing of very large datasets.
- **Applicability:** No like the Longest Common Substring algorithm—which focuses on character-level similarity between two strings—Sindy operates on cell values across multiple columns to identify inclusion relationships.
- **Efficiency:** By leveraging grouping and set intersection operations, Sindy quickly filters out irrelevant candidates and hones in on true inclusion dependencies.
- **Robustness:** The algorithm is tailored for structured data in relational databases, making it more robust to the diverse and noisy data often encountered in real-world datasets.
- **Flexibility:** Sindy’s design allows for easy extension to additional data formats (e.g., JSON) and integration with further optimization strategies such as caching and indexing.