# Distributed Data Management

## spark-homework, Team null

Marcian Seeger & Justus Gillmann

# Code 1

```scala
def discoverINDs(inputs: List[String], spark: SparkSession): Unit = {
  import spark.implicits._

  val valuesWithEmptySet = mutable.HashSet[String]()
  inputs.map(input => readData(input, spark))       // list of files -> list of Spark Datasets
    .map(ds => {                                    // Spark Dataset -> tuples of value and corresponding columnName
      val columns = ds.columns
      ds.flatMap(row => {
        for (i <- columns.indices) yield {
          (row.getString(i), columns(i))
        }
      }) // (value, colName)
    })                                              : List[Dataset[(String, String)]]
    .reduce(_ union _)                              : Dataset[(String, String)]
```

- Filepaths → Spark Datasets

- Dataset of rows → Dataset of (value, columnName)-pairs (for each row and value)

# Code 2

```scala
def discoverINDs(inputs: List[String], spark: SparkSession): Unit = {
  import spark.implicits._

  val valuesWithEmptySet = mutable.HashSet[String]()
  inputs.map(input => readData(input, spark))        // list of files -> list of Spark Datasets
    .map(ds => {                                     // Spark Dataset -> tuples of value and corresponding columnName
      val columns = ds.columns
      ds.flatMap(row => {
        for (i <- columns.indices) yield {
          (row.getString(i), columns(i))
        }
      }) // (value, colName)
    })
    .reduce(_ union _)
    .groupByKey(x => x._1)                           // group by value --> get unique values with corresponding colNames
    .mapGroups((_, it) => it.map(elem => elem._2).toSet)  // throw away value, only keep columnNames as Sets
    .flatMap(set => {                                // get each possible combination of elements in the set (candidates)
      if(set.size == 1)
        valuesWithEmptySet.add(set.head)
      set.map(colName => (colName, set - colName))
    })
```

(value, columnName)
→Dataset[(String, String)]

← save columnNames without a partner for later

: List[Dataset[(String, String)]]
: Dataset[(String, String)]

: Dataset[(String, Set[String])]

- (value, columnName)-Tuples → groups of (value, all existing columnNames for this value)-Tuples
- → Sets of (all existing columnNames for „this" value)
       („this value" is thrown away)
- build all possible combinations of elements in the set (= IND candidates)

# Code 3

(columnName, [more columnNames])
→ Dataset[(String, Set[String])]

```
40      .groupByKey(x => x._1)                          // group by candidate              : KeyValueGroupedDataset[String, (String, ...)]
41      .mapGroups((s, it) =>
42        (s, it
43          .dropWhile(_ => valuesWithEmptySet.contains(s)) // (if not relevant, empty iterator to skip map-reduce)
44          .map(x => x._2)
45          .reduce((a, b) => a.intersect(b)))
46      )                                                                                  : Dataset[(String, Set[String])]
```

• group by key (columnName / IND candidate) → get all Sets for each columnName

• clear iterator if the grouped by key was partnerless before (skip map-reduce)

→ only keep columnNames that are included in every Set for this key (columnName / …)

# Code 4

(columnName, [more columnNames])
→ Dataset[(String, Set[String])]

```scala
47        .filter(elem => elem._2.nonEmpty)          // throw away results with empty sets
48        .collect()                                  // collect to Array
49        .sortBy(x => x._1)                          // sort keys
50        .map(x => (x._1, x._2.toList.sorted))       // sort values
51        .foreach(x => println(x._1 + " < " + x._2.mkString(", ")))              : Unit
52     }
53
```

• filter out results with empty sets

• sort results by key, then sort the sets

• print results

# Code

```scala
19  def discoverINDs(inputs: List[String], spark: SparkSession): Unit = {
20    import spark.implicits._
21
22    val valuesWithEmptySet = mutable.HashSet[String]()
23    inputs.map(input => readData(input, spark))      // list of files -> list of Spark Datasets
24      .map(ds => {                                   // Spark Dataset -> tuples of value and corresponding columnName
25        val columns = ds.columns
26        ds.flatMap(row => {
27          for (i <- columns.indices) yield {
28            (row.getString(i), columns(i))
29          }
30        }) // (value, colName)
31      })                                                                          : List[Dataset[(String, String)]]
32      .reduce(_ union _)                                                          : Dataset[(String, String)]
33      .groupByKey(x => x._1)                         // group by value --> get unique values with corresponding colNames
34      .mapGroups((_, it) => it.map(elem => elem._2).toSet)  // throw away value, only keep columnNames as Sets
35      .flatMap(set => {                              // get each possible combination of elements in the set (candidates)
36        if(set.size == 1)
37          valuesWithEmptySet.add(set.head)
38        set.map(colName => (colName, set - colName))
39      })                                                                          : Dataset[(String, Set[String])]
40      .groupByKey(x => x._1)                         // group by candidate                  : KeyValueGroupedDataset[String, (String, ...)]
41      .mapGroups((s, it) =>
42        (s, it
43          .dropWhile(_ => valuesWithEmptySet.contains(s)) // (if not relevant, empty iterator to skip map-reduce)
44          .map(x => x._2)
45          .reduce((a, b) => a.intersect(b)))
46      )                                                                           : Dataset[(String, Set[String])]
47      .filter(elem => elem._2.nonEmpty)              // throw away results with empty sets
48      .collect()                                     // collect to Array
49      .sortBy(x => x._1)                             // sort keys
50      .map(x => (x._1, x._2.toList.sorted))          // sort values
51      .foreach(x => println(x._1 + " < " + x._2.mkString(", ")))                  : Unit
52  }
53
```