# Data Analysis and Integration

Data Profiling

# Motivation

# Data Profiling

- ## What is data profiling?

  – Analyze the contents of a data source

  – Gather statistics about the data contained therein

    - minimum, maximum, average, range, value distribution, etc.

  – Identify data quality problems

    - missing or incomplete data, errors, inconsistencies, etc.

  – Understand the logic and relationships between data

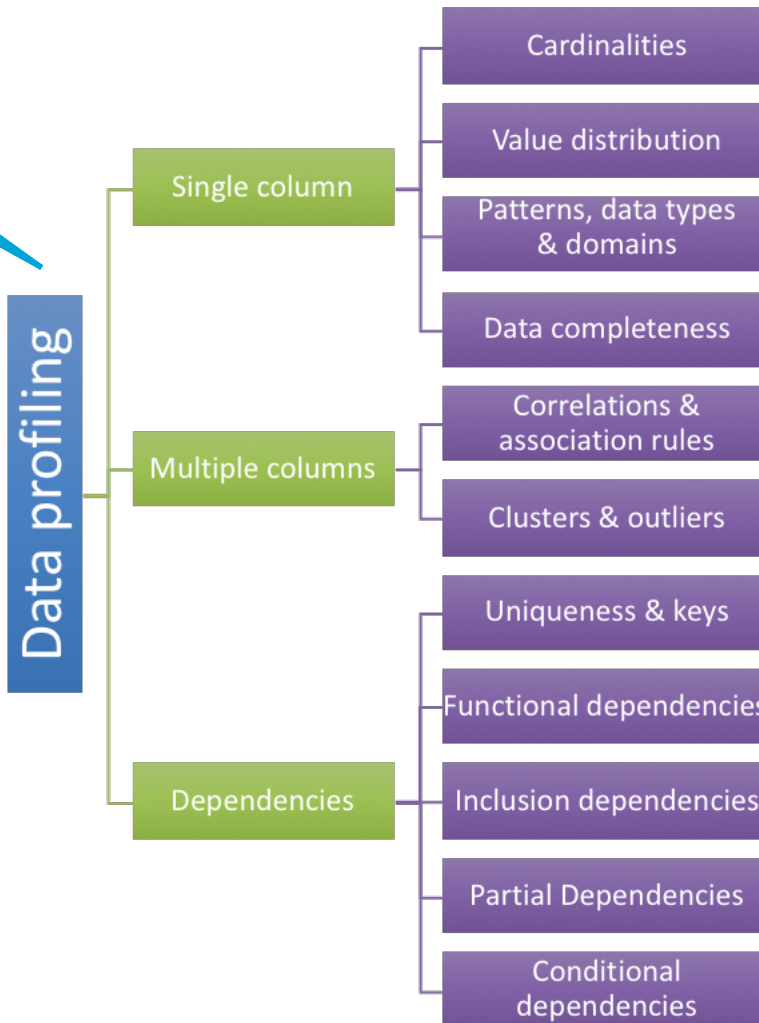    - unique values, keys, foreign keys and other constraints

# Data Profiling

- These tasks (and many more) are instances of data profiling
  - number of rows, number of null values, number of distinct values
  - minimum value, maximum value, minimum length, maximum length
  - single- and multi-column frequency histogram
  - precision of numeric values, length of string values
  - data type discovery
  - uniqueness and constancy
  - single- and multi-column primary key discovery
  - single- and multi-column foreign key discovery
  - value overlap (cross-domain analysis)
  - text field profiling
  - pattern discovery (e.g. phone number patterns)
  - soundex frequency histogram
  - etc.

# Data Profiling

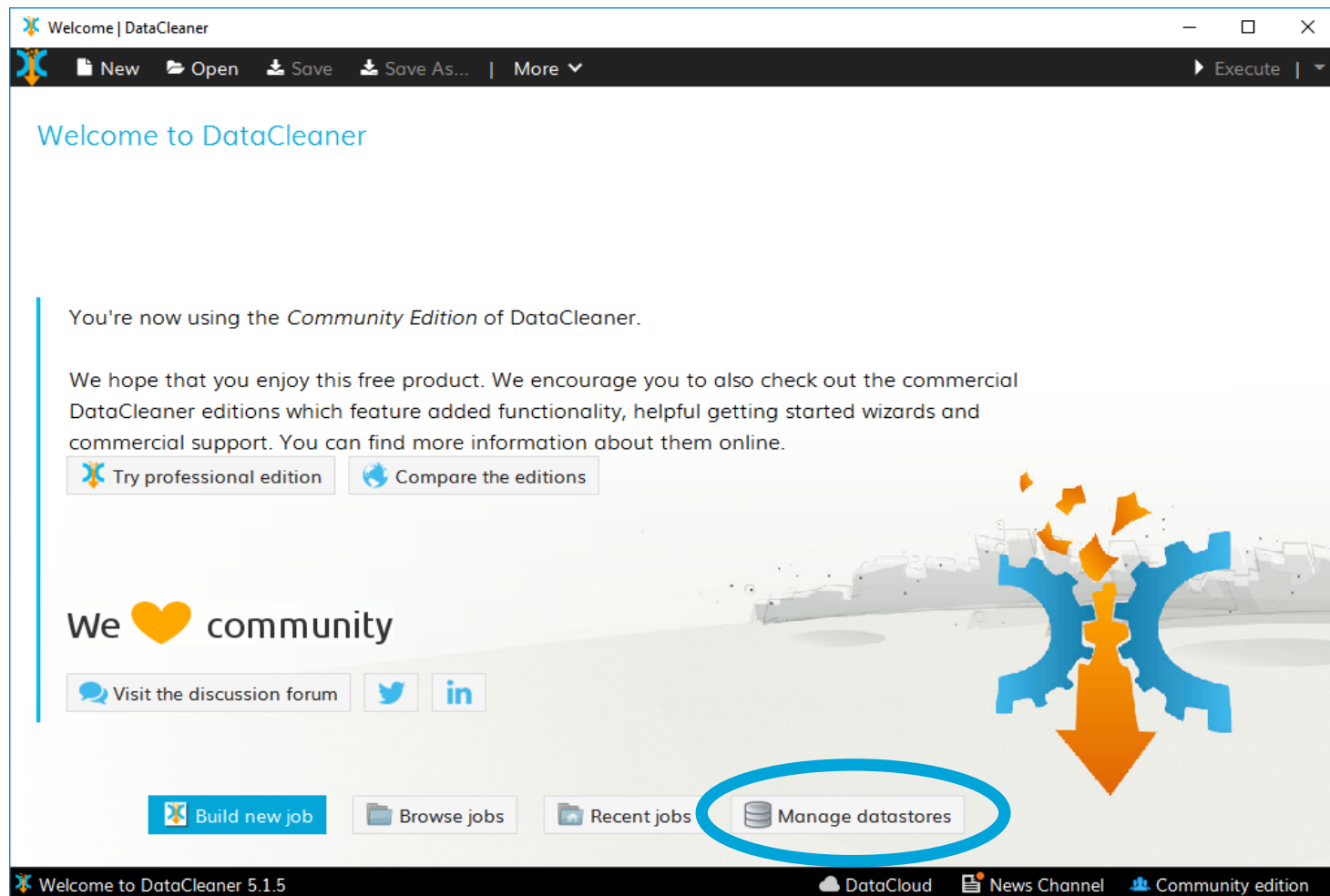Taxonomy of Data Profiling Operations

**Data profiling**

- **Single column**
  - Cardinalities
  - Value distribution
  - Patterns, data types & domains
  - Data completeness
- **Multiple columns**
  - Correlations & association rules
  - Clusters & outliers
- **Dependencies**
  - Uniqueness & keys
  - Functional dependencies
  - Inclusion dependencies
  - Partial Dependencies
  - Conditional dependencies

Z. Abedjan, L. Golab, F. Naumann
*Profiling relational data: a survey*
VLDB Journal, vol. 24, no. 4, 2015

# Data Profiling

- Data Profiling tools



data
quality
problems

**Data
profiling**

**metadata**
- data types
- keys
- dependencies
- …

data
statistics

# Data Profiling Tool Concepts

# Data Profiling

- Data profiling tools

# Data Profiling

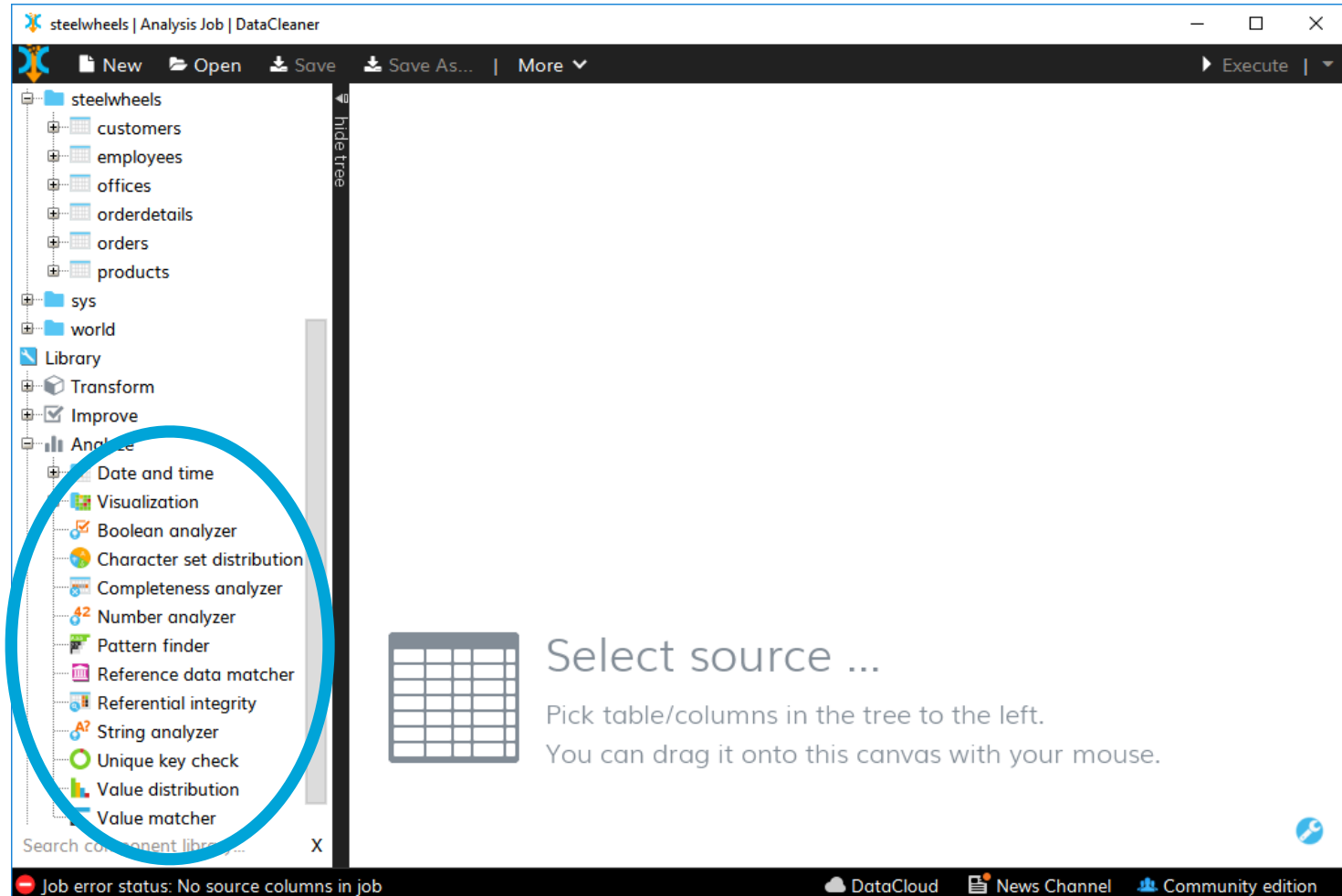- A datastore can be a file, database, etc.

# Data Profiling

- Creating a new datastore

# Data Profiling

- Creating a new datastore

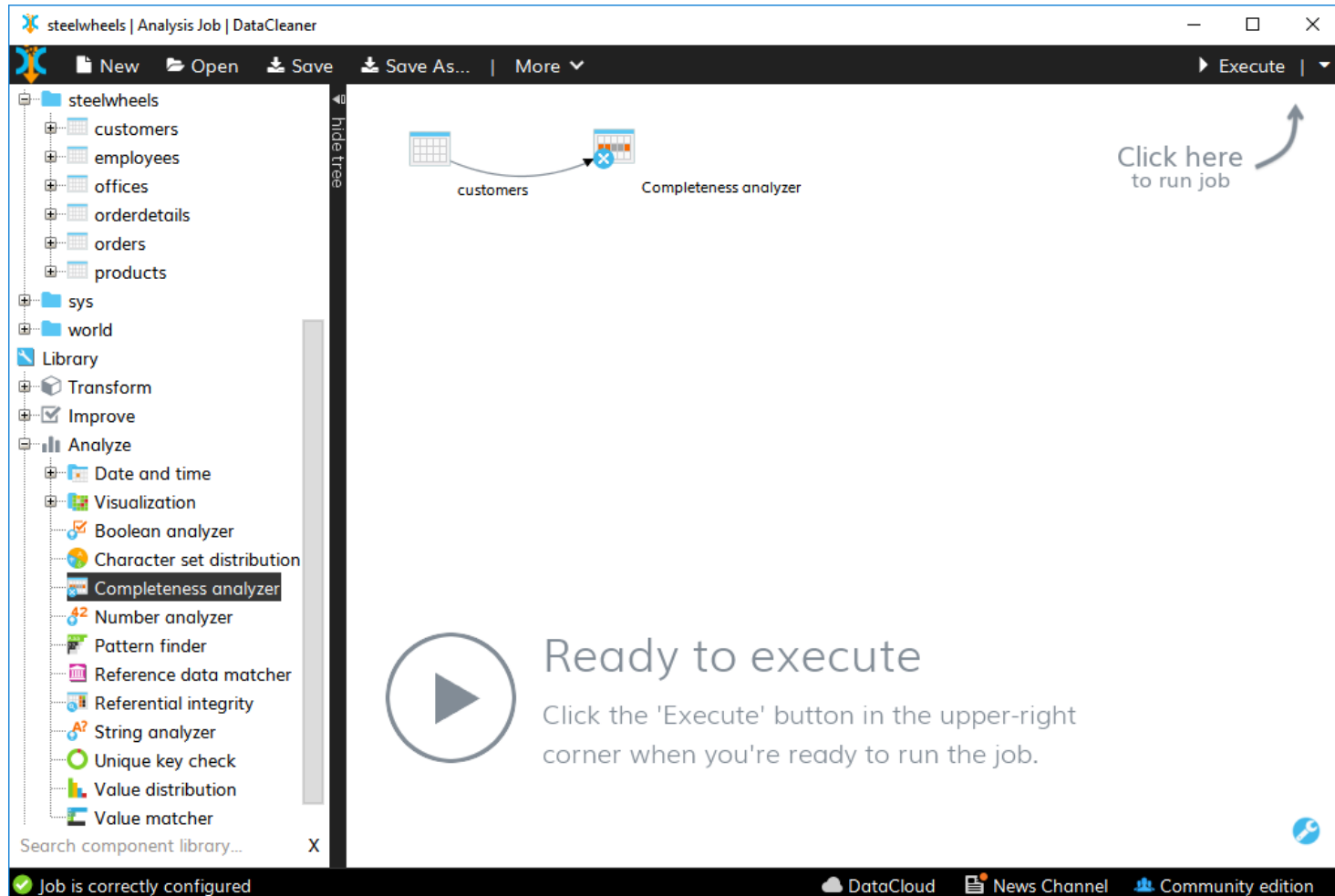# Data Profiling

- Several options for data profiling tasks

# Data Profiling

- Selecting the data source

# Data Profiling

- Completeness analysis

# Data Profiling

- Completeness analysis
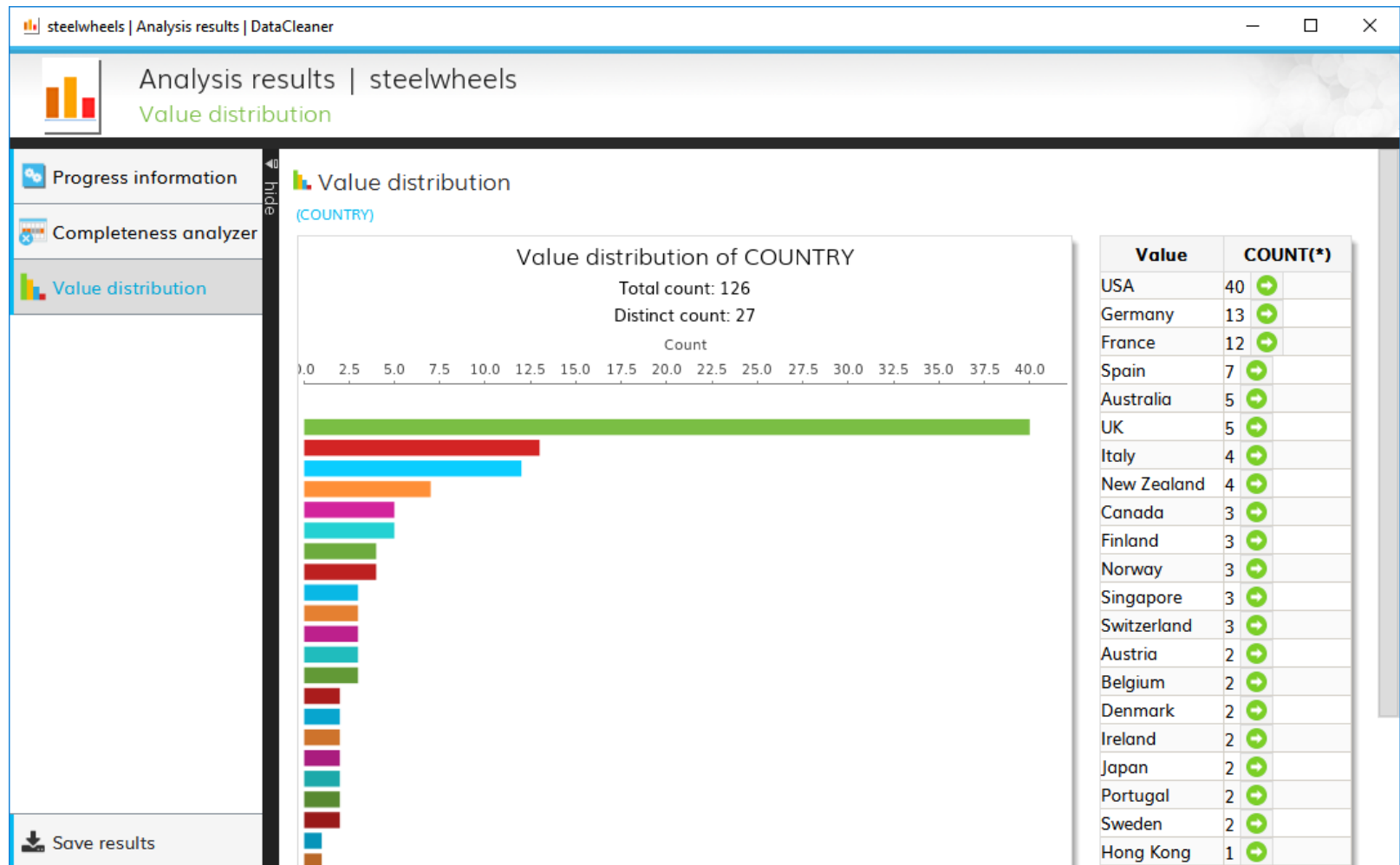
# Data Profiling

- Completeness analysis

# Data Profiling

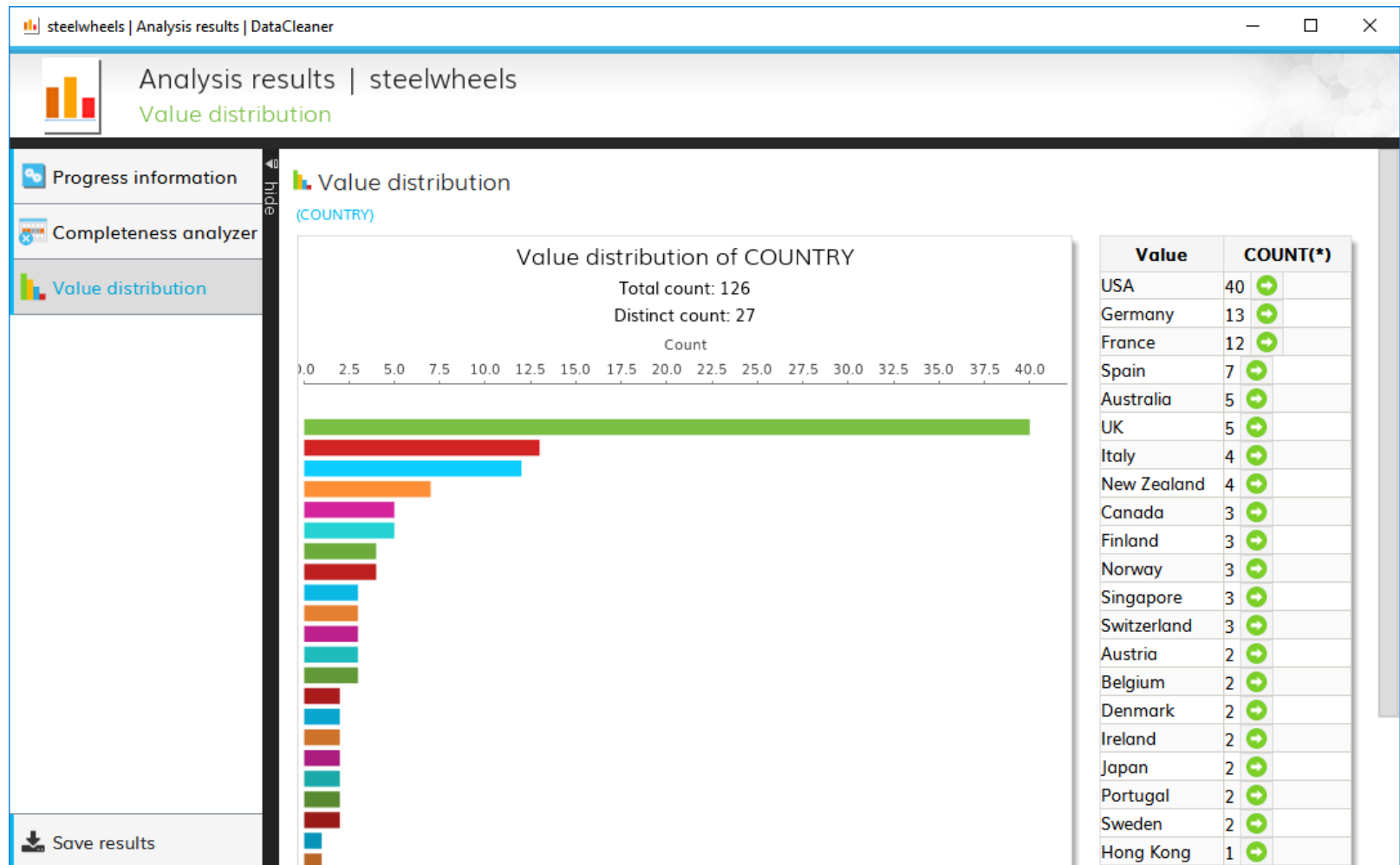- Value distribution

# Data Profiling

- Value distribution

# Data Profiling

- Value distribution

# Data Profiling

- Value distribution

# Data Profiling

- Value distribution – detailed results

# Data Profiling

- String analysis

# Data Profiling

- String analysis

# Data Profiling

- String analysis



Analysis results | steelwheels
String analyzer

| | Progress information |
| | Completeness analyzer |
| | Value distribution |
| | String analyzer |

**String analyzer**
(ADDRESSLINE1,CITY,STATE,COUNTRY)

| | ADDRESSLINE1 | CITY | STATE | COUNTRY |
|---|---|---|---|---|
| Row count | 126 | 126 | 126 | 126 |
| Null count | 0 | 0 | 74 | 0 |
| Blank count | 0 | 0 | 0 | 0 |
| Entirely uppercase count | 1 | 6 | 44 | 45 |
| Entirely lowercase count | 0 | 0 | 0 | 0 |
| Total char count | 2474 | 990 | 153 | 709 |
| Max chars | 46 | 17 | 13 | 12 |
| Min chars | 11 | 3 | 2 | 2 |
| Avg chars | 19.635 | 7.857 | 2.942 | 5.627 |
| Max white spaces | 6 | 2 | 2 | 1 |
| Min white spaces | 1 | 0 | 0 | 0 |
| Avg white spaces | 2.468 | 0.183 | 0.058 | 0.048 |
| Uppercase chars | 293 | 168 | 100 | 217 |
| Uppercase chars (excl. first letters) | 176 | 42 | 47 | 91 |
| Lowercase chars | 1407 | 798 | 49 | 486 |
| Digit chars | 365 | 0 | 0 | 0 |
| Diacritic chars | 9 | 8 | 1 | 0 |
| Non-letter chars | 774 | 24 | 4 | 6 |
| Word count | 435 | 148 | 55 | 132 |
| Max words | 7 | 3 | 3 | 2 |
| Min words | 2 | 1 | 1 | 1 |

| | Save results |

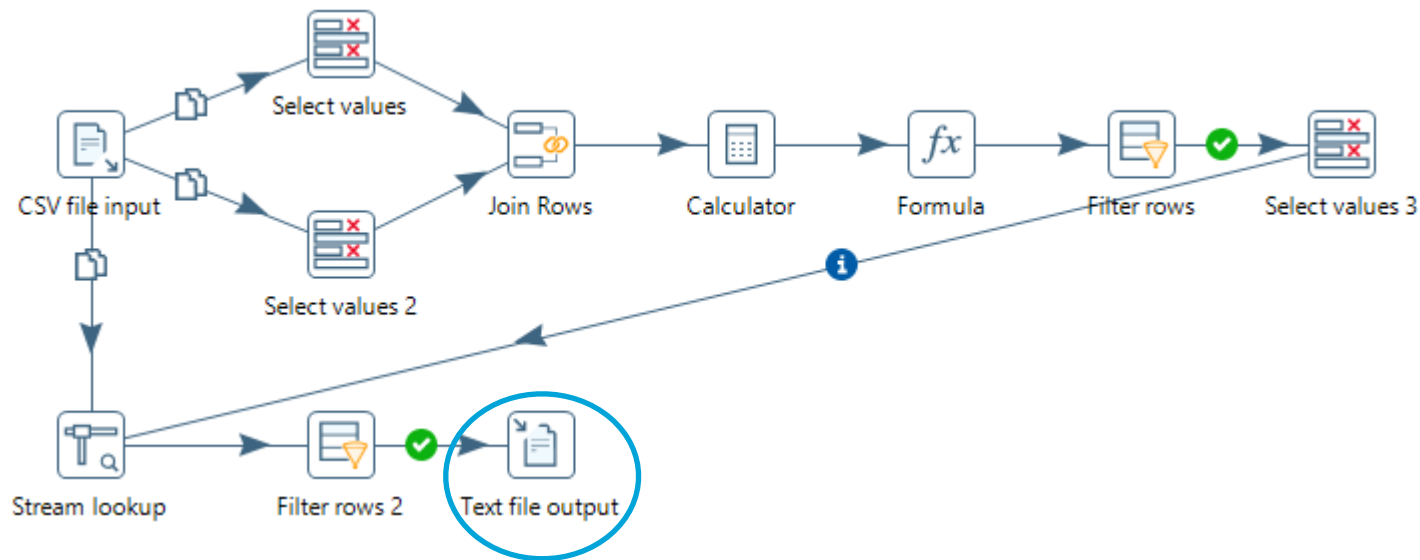# Integration

# Data Profiling Integration

- Integration between tools
  - Pentaho Data Integration (PDI) with DataCleaner plugin
  - The output of any transformation step can be a data source for data profiling
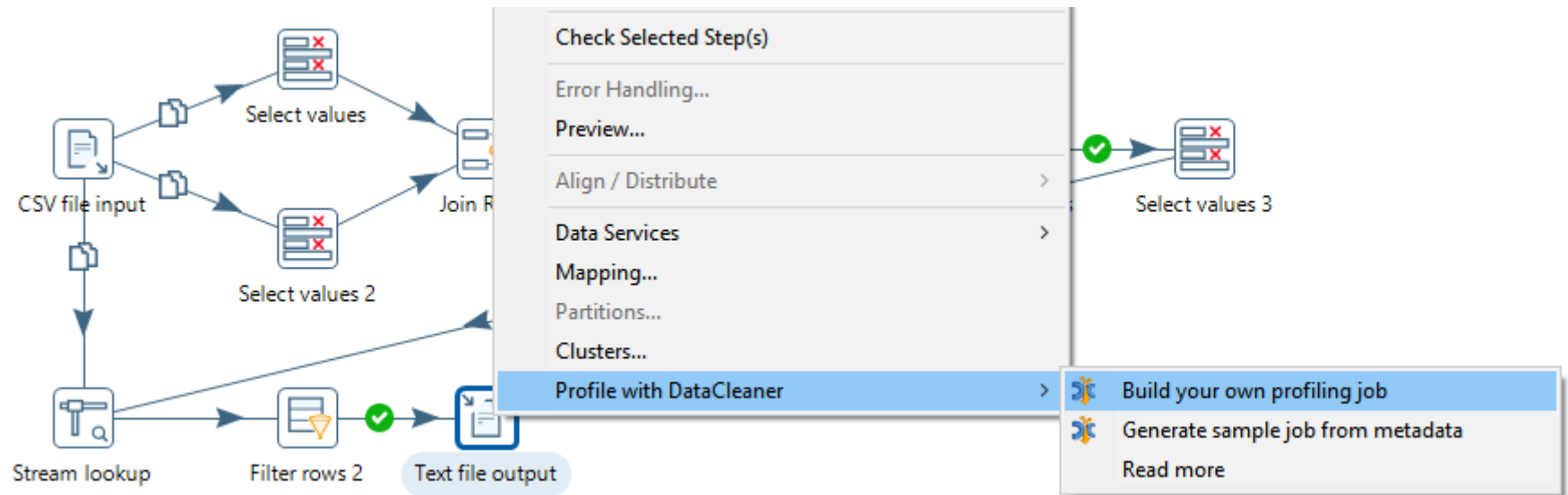
# Data Profiling Integration
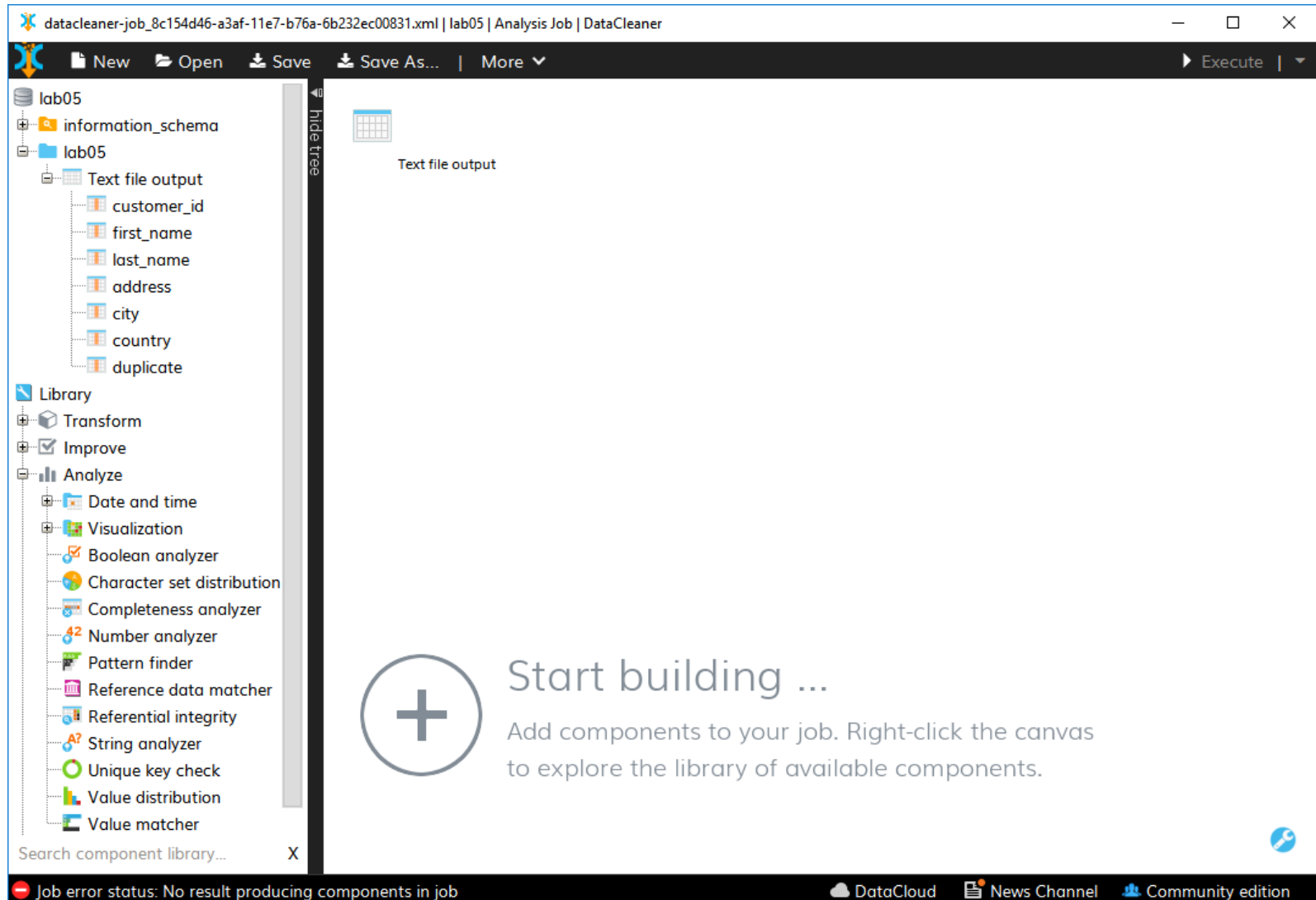


Rows of step: Text file output (64 rows)

| # | customer_id | first_name | last_name | address | city | country | duplicate |
|---|---|---|---|---|---|---|---|
| 1 | 103 | Carine | Schmitt | 54, rue Royale | Nantes | France | <null> |
| 2 | 119 | Janine | Labrune | 67, rue des Cinquante Otages | Nantes | France | <null> |
| 3 | 121 | Jonas | Bergulfsen | Erling Skakkes gate 78 | Stavern | Norway | <null> |
| 4 | 125 | Zbygniew | Piestrzeniewicz | ul. Filtrowa 68 | Warszawa | Poland | <null> |
| 5 | 128 | Roland | Keitel | Lyonerstr. 34 | Frankfurt | Germany | <null> |
| 6 | 141 | Diego | Freyre | c/ Moralzarzal, 86 | Madrid | Spain | <null> |
| 7 | 144 | Christina | Berglund | Berguvsvägen 8 | Luleå | Sweden | <null> |
| 8 | 145 | Jytte | Petersen | Vinbæltet 34 | Kobenhavn | Denmark | <null> |
| 9 | 146 | Mary | Saveley | 2, rue du Commerce | Lyon | France | <null> |

# Data Profiling Integration

- Example
  - Perform data profiling after duplicate elimination

# Data Profiling Integration

# Data Profiling Integration

# Data Profiling Integration

- Value distribution

# Data Profiling Integration

- Value distribution

# Data Profiling Integration

# Data Profiling Integration

- String analysis

# Data Profiling Integration

- String analysis



Analysis results | lab05 | datacleaner-job_67a84087-a42e-11e7-b651-4379f2e87310.xml
String analyzer

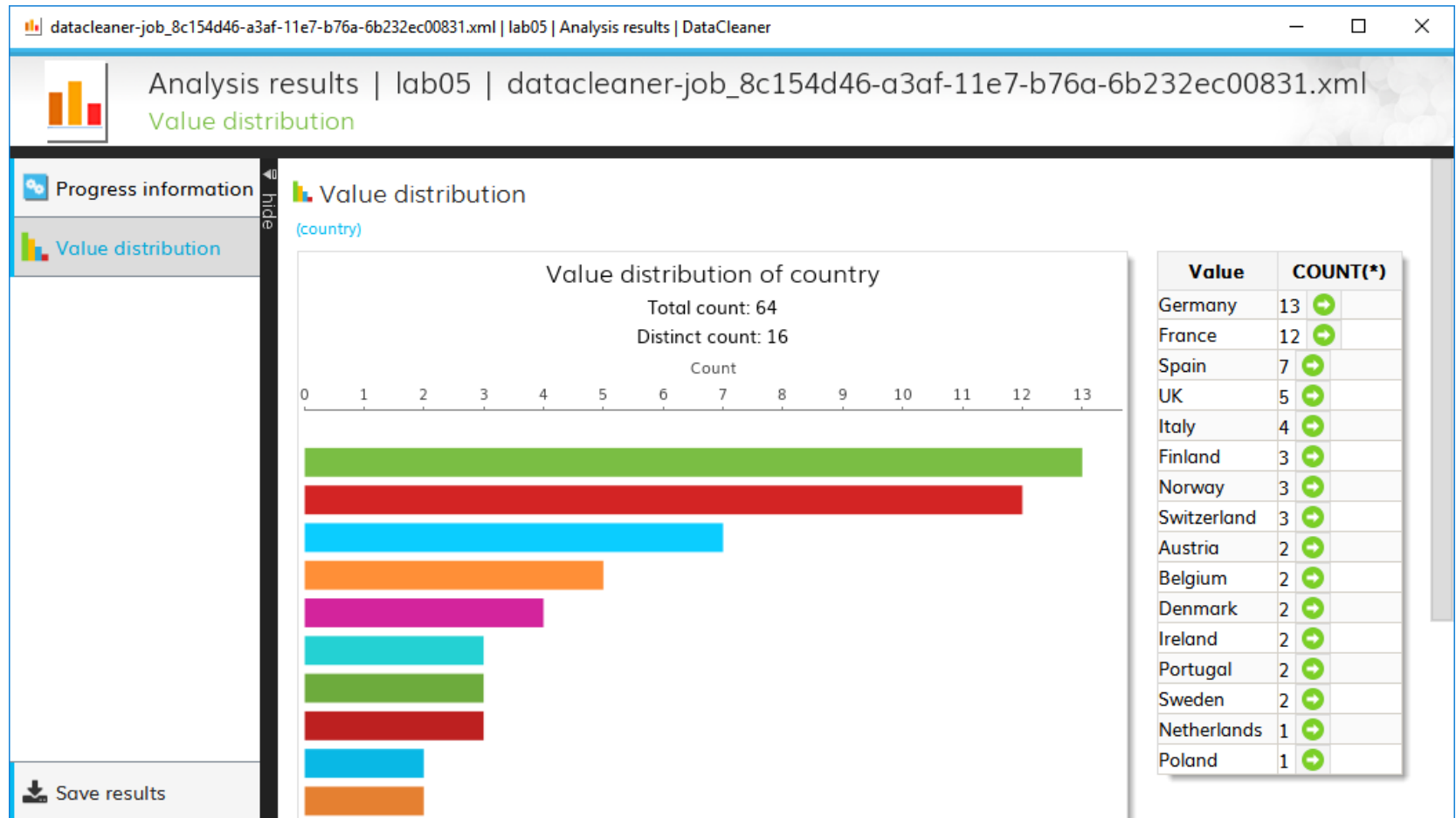| | Progress information | Value distribution | String analyzer | Save results |

### String analyzer
(5 columns)

| | first_name | last_name | address | city | country |
|---|---|---|---|---|---|
| Row count | 64 | 64 | 64 | 64 | 64 |
| Null count | 0 | 0 | 0 | 0 | 0 |
| Blank count | 0 | 0 | 0 | 0 | 0 |
| Entirely uppercase count | 0 | 0 | 0 | 0 | 5 |
| Entirely lowercase count | 0 | 0 | 0 | 0 | 0 |
| Total char count | 397 | 441 | 1171 | 435 | 401 |
| Max chars | 10 | 15 | 38 | 13 | 11 |
| Min chars | 3 | 3 | 11 | 4 | 2 |
| Avg chars | 6.203 | 6.891 | 18.297 | 6.797 | 6.266 |
| Max white spaces | 1 | 1 | 5 | 1 | 0 |
| Min white spaces | 0 | 0 | 1 | 0 | 0 |
| Avg white spaces | 0.016 | 0.031 | 2.188 | 0.016 | 0 |
| Uppercase chars | 65 | 66 | 115 | 65 | 69 |
| Uppercase chars (excl. first letters) | 1 | 3 | 65 | 1 | 5 |
| Lowercase chars | 331 | 373 | 736 | 369 | 332 |
| Digit chars | 0 | 0 | 136 | 0 | 0 |
| Diacritic chars | 4 | 2 | 9 | 7 | 0 |
| Non-letter chars | 1 | 2 | 320 | 1 | 0 |
| Word count | 65 | 66 | 203 | 65 | 64 |
| Max words | 2 | 2 | 6 | 2 | 1 |
| Min words | 1 | 1 | 2 | 1 | 1 |