

Data Analysis and Integration

Data profiling

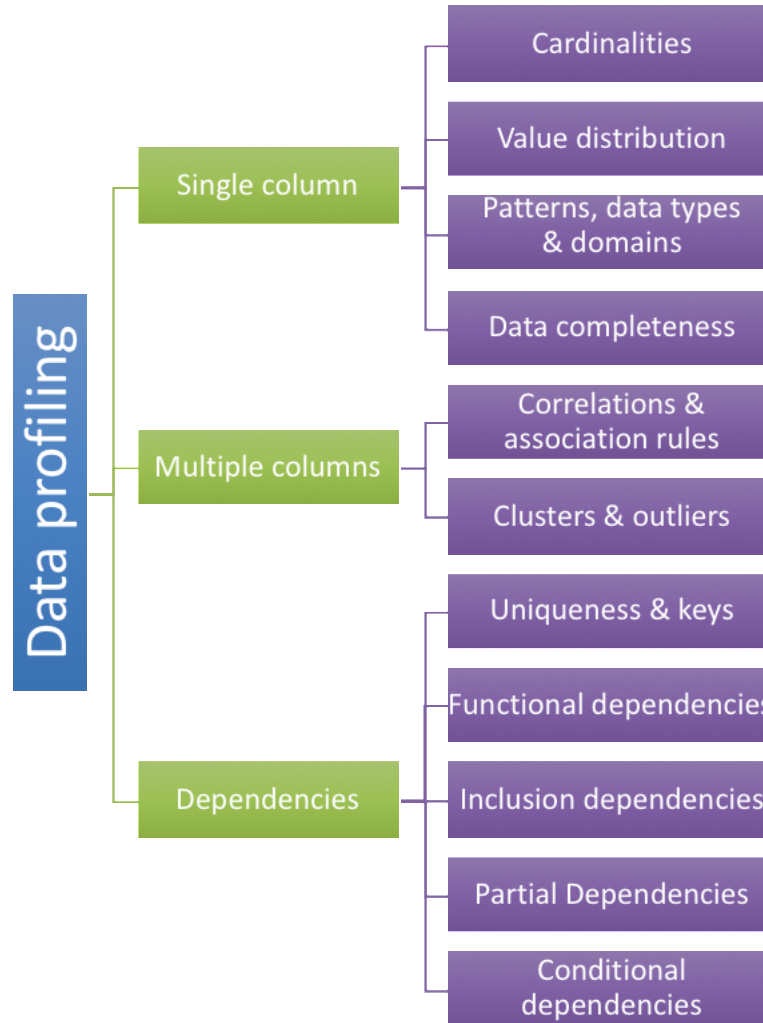
Data profiling

- 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 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.

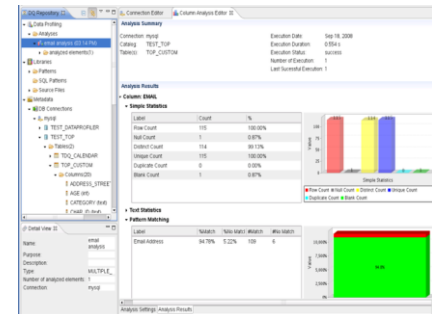
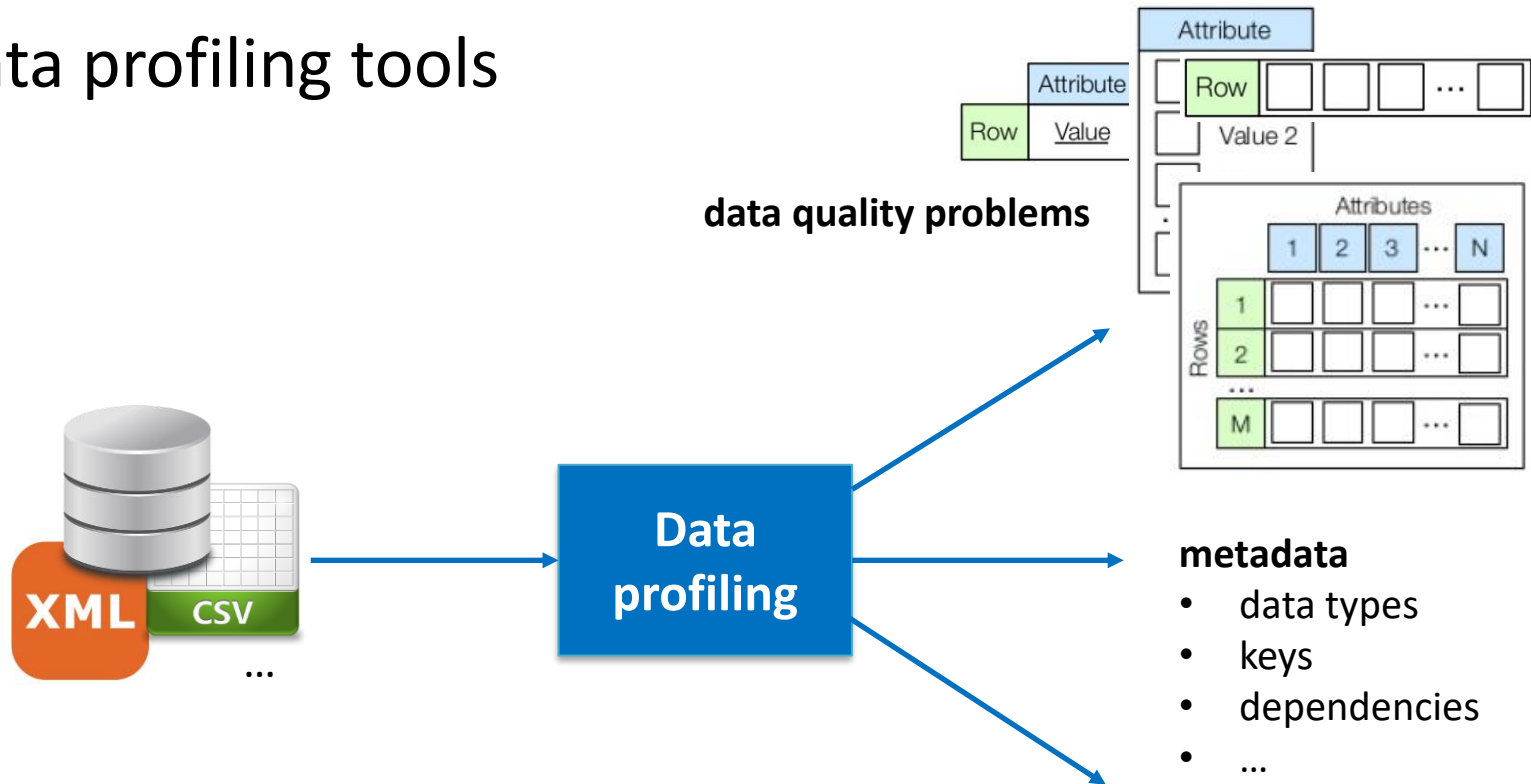
Data profiling



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

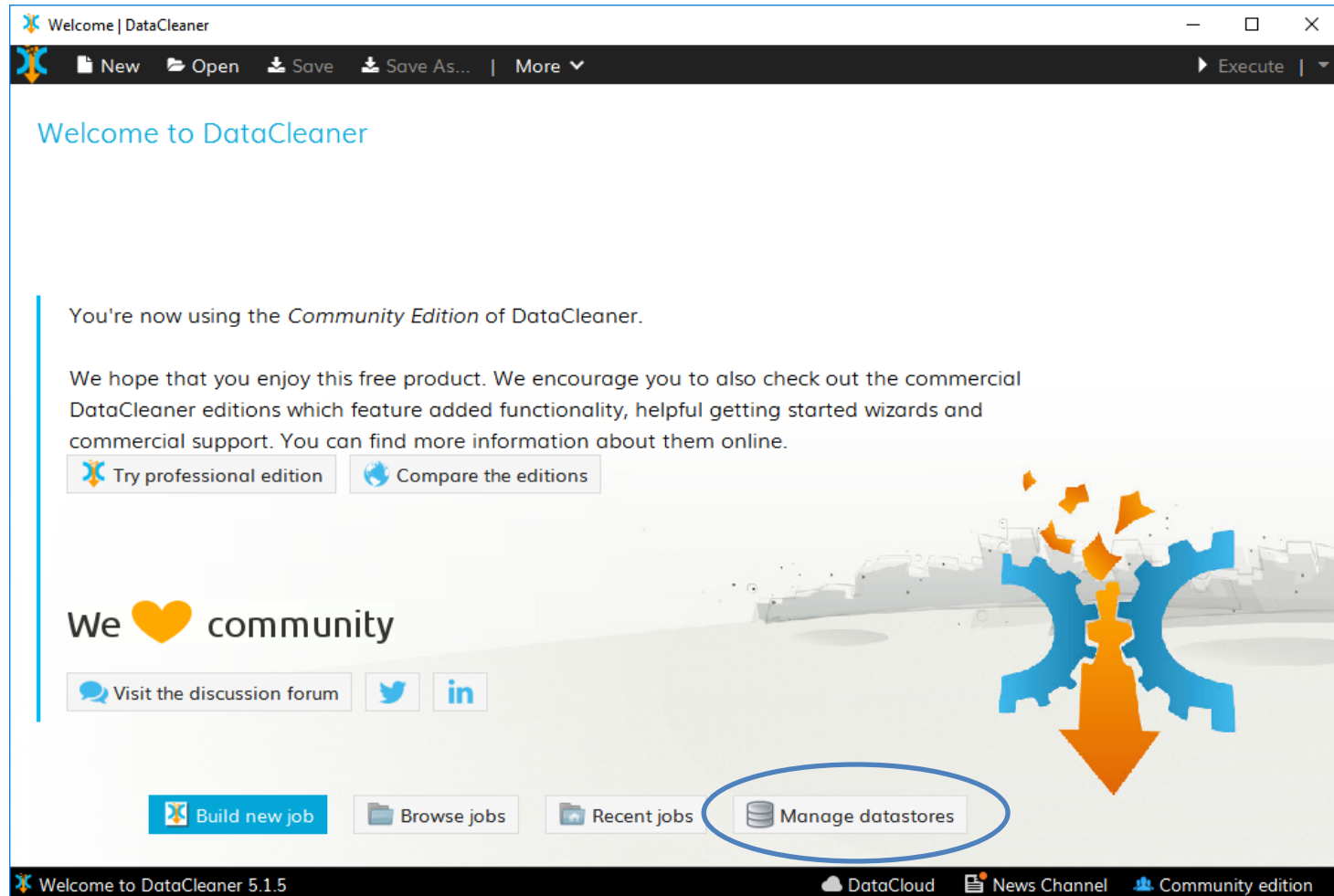
Data profiling

- Data profiling tools



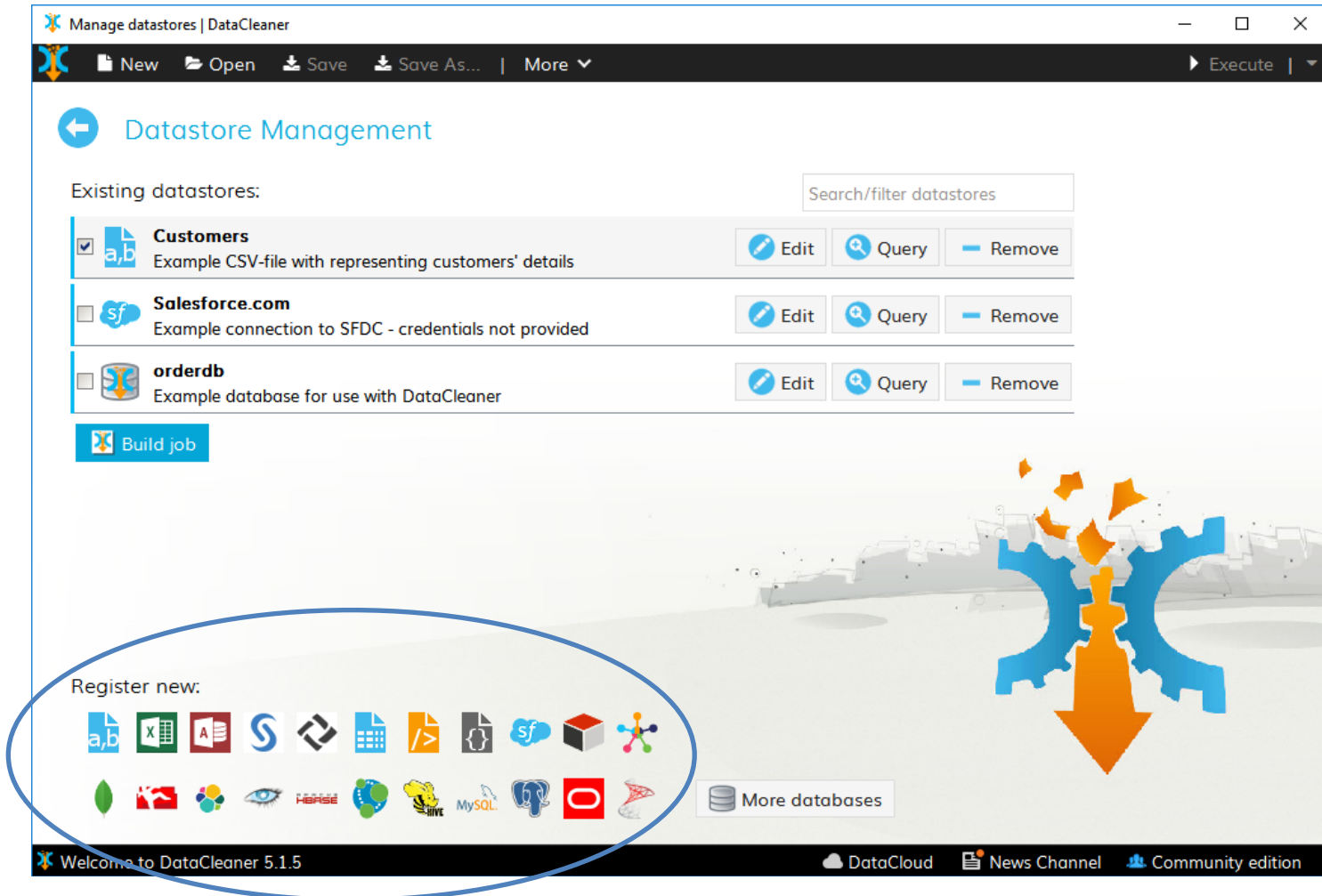
Data profiling

- Data profiling tools



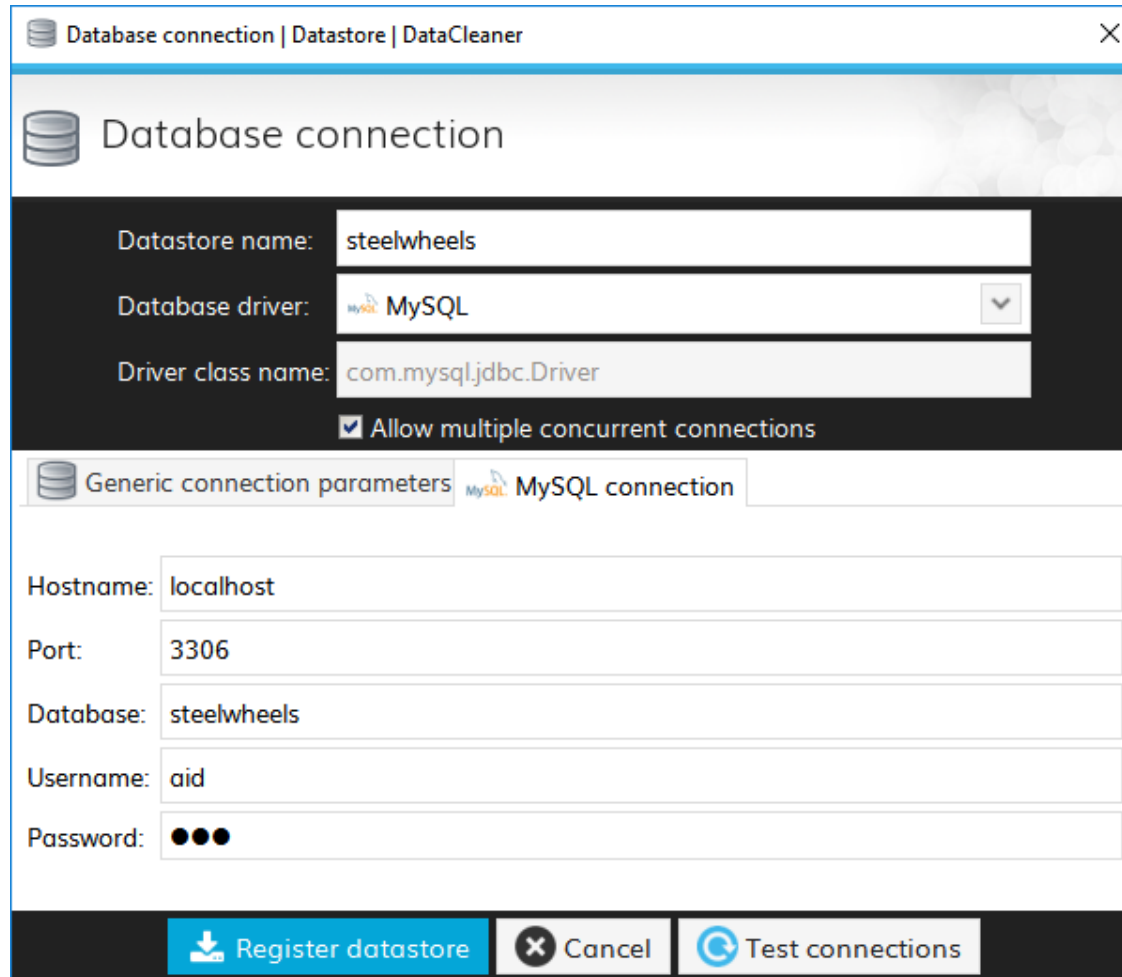
Data profiling

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



Data profiling

- Creating a new datastore



Database connection | Datastore | DataCleaner

Database connection

Datastore name: steelwheels

Database driver: MySQL

Driver class name: com.mysql.jdbc.Driver

☒ Allow multiple concurrent connections

Generic connection parameters MySQL connection

Hostname: localhost

Port: 3306

Database: steelwheels

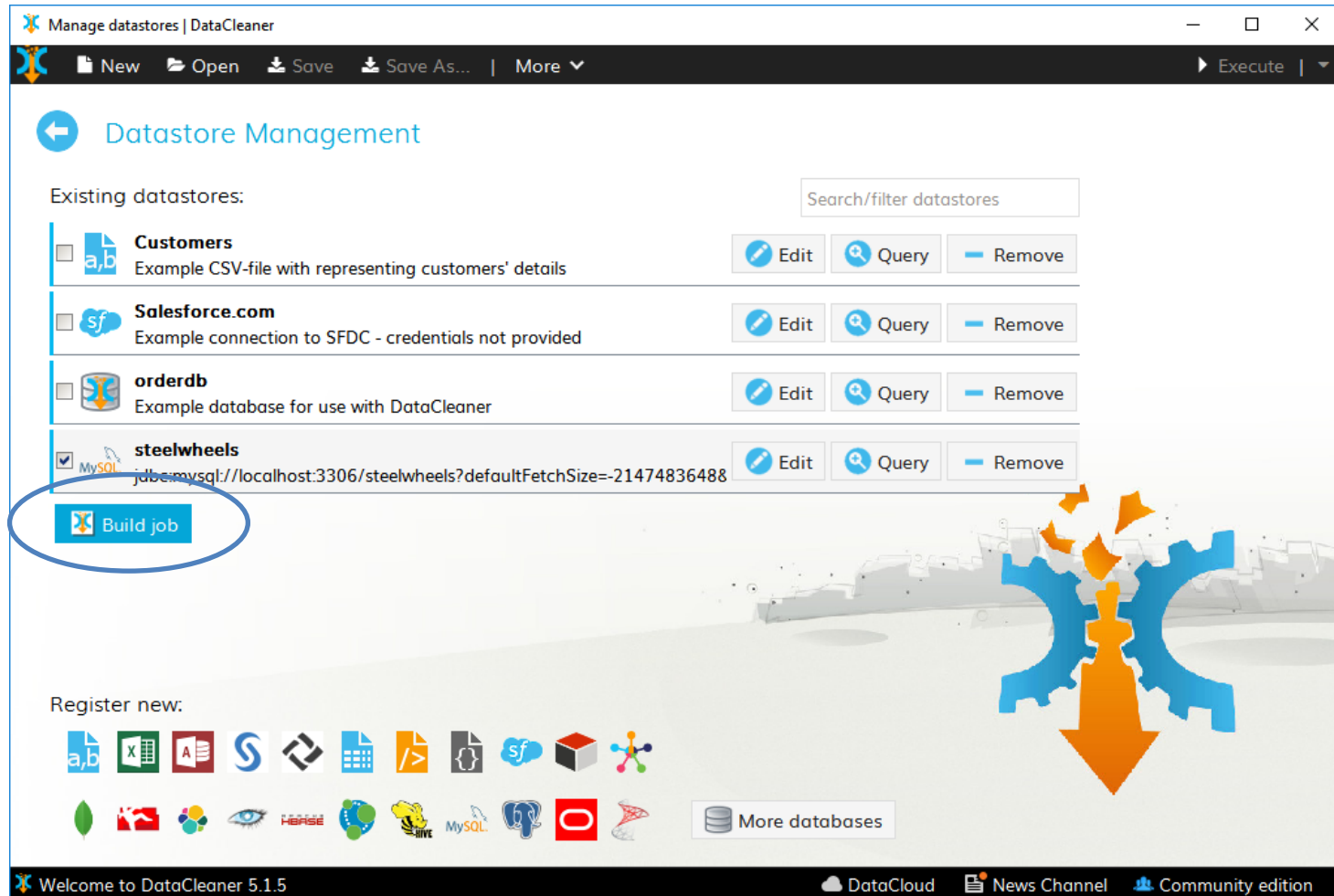
Username: aid

Password: ●●●

Register datastore Cancel Test connections

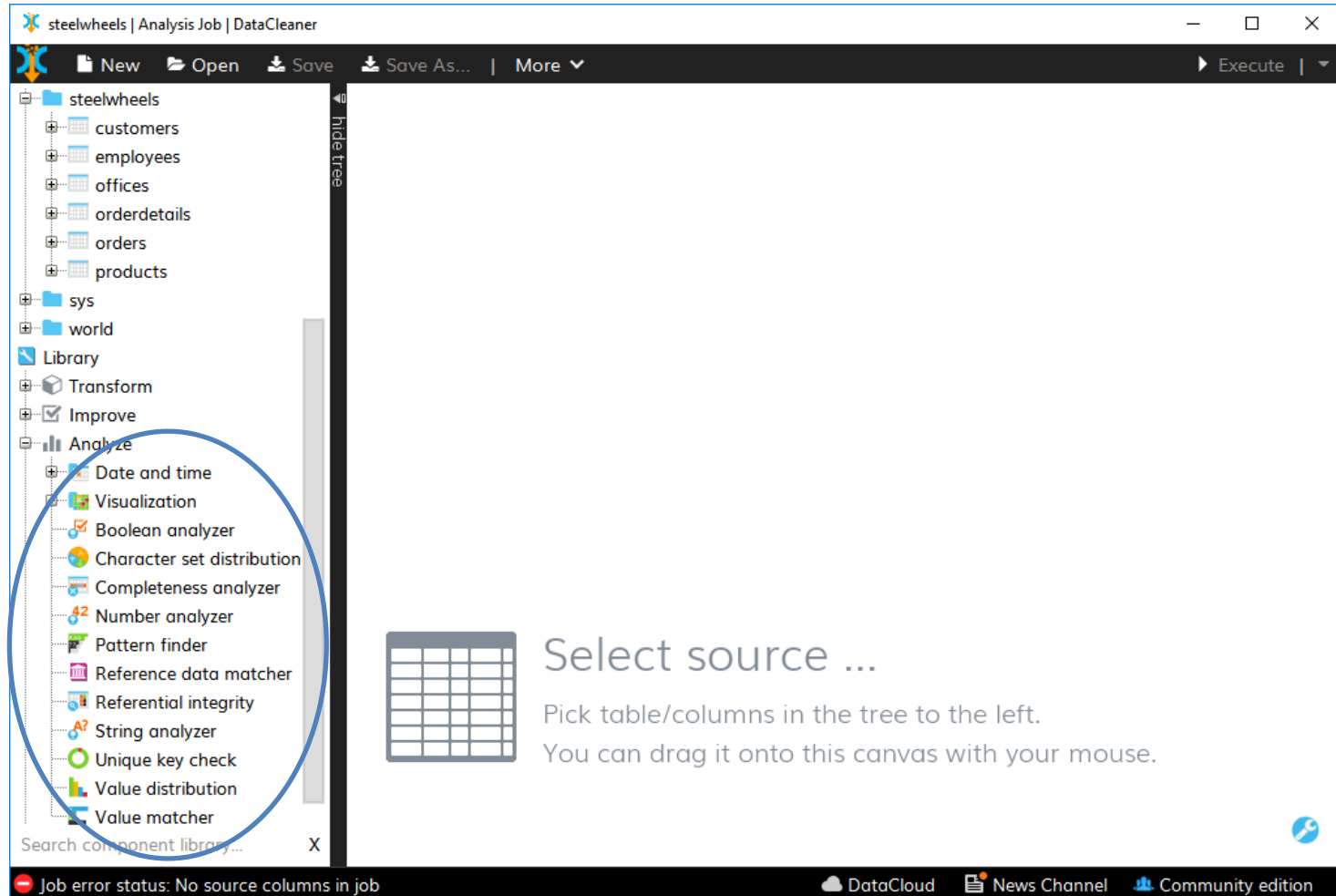
Data profiling

- Creating a new datastore



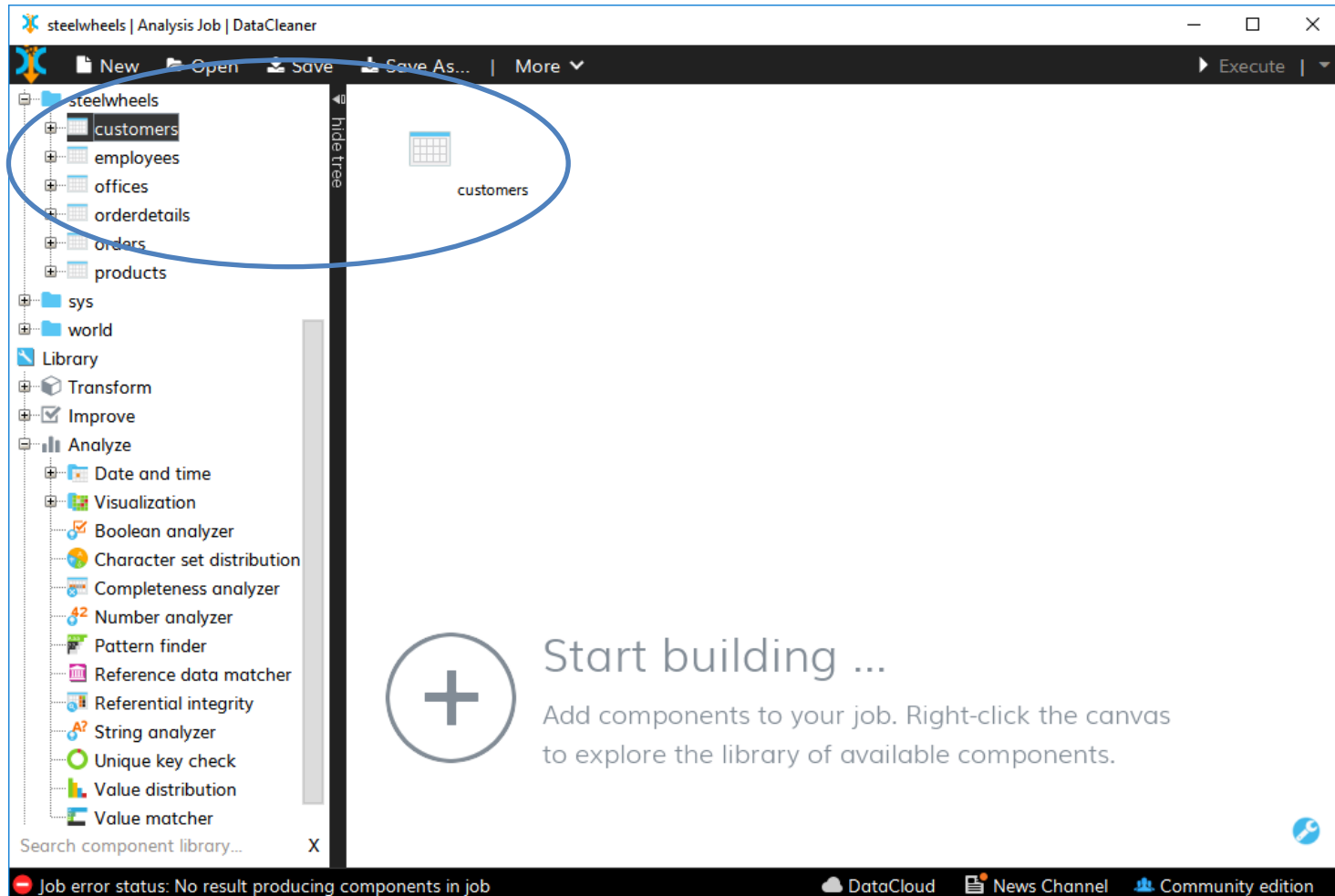
Data profiling

- Several options for data profiling tasks



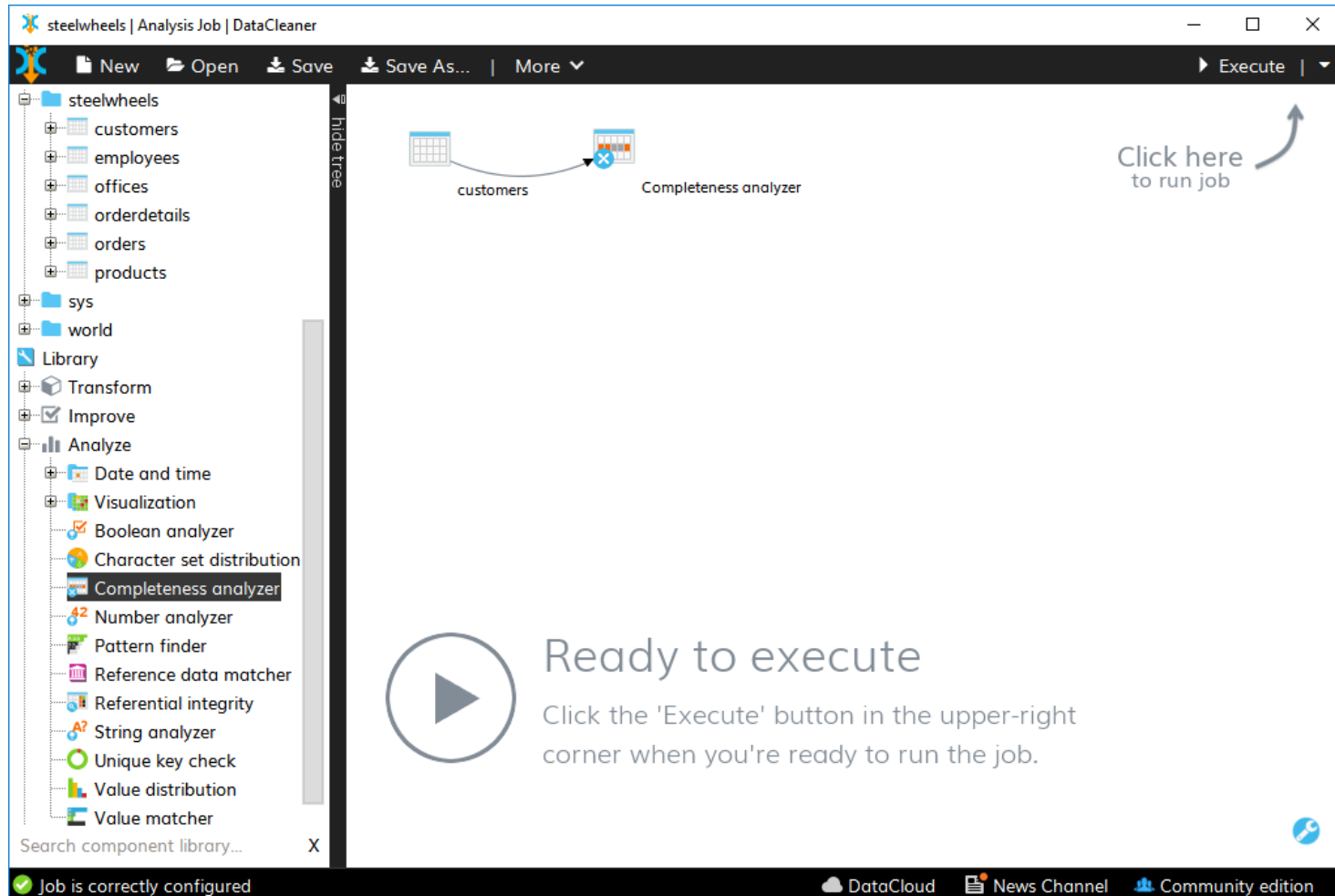
Data profiling

- Selecting the data source



Data profiling

- Completeness analysis



Data profiling

- Completeness analysis

Completeness analyzer | DataCleaner

Completeness analyzer

Documentation Rename

Input columns

Values: Select all Select none

Values to check for completeness

- ☐ CUSTOMERNUMBER
- ☐ CUSTOMERNAME
- ☐ CONTACTLASTNAME
- ☐ CONTACTFIRSTNAME
- ☐ PHONE
- ☐ ADDRESSLINE1
- ☐ ADDRESSLINE2
- ☐ CITY
- ☒ STATE
- ☐ POSTALCODE
- ☐ COUNTRY
- ☐ SALESREPEMPOYEEENUMBER
- ☐ CREDITLIMIT

Not <blank> or <null>

Required properties

Evaluation mode: When any field is incomplete, the record is incomplete

Optional properties (1)

Close

Data profiling

- Completeness analysis

steelwheels | Analysis results | DataCleaner

Analysis results | steelwheels
Completeness analyzer

Progress information
Completeness analyzer

Completeness analyzer
(STATE)

Incomplete records (74)

View detailed rows

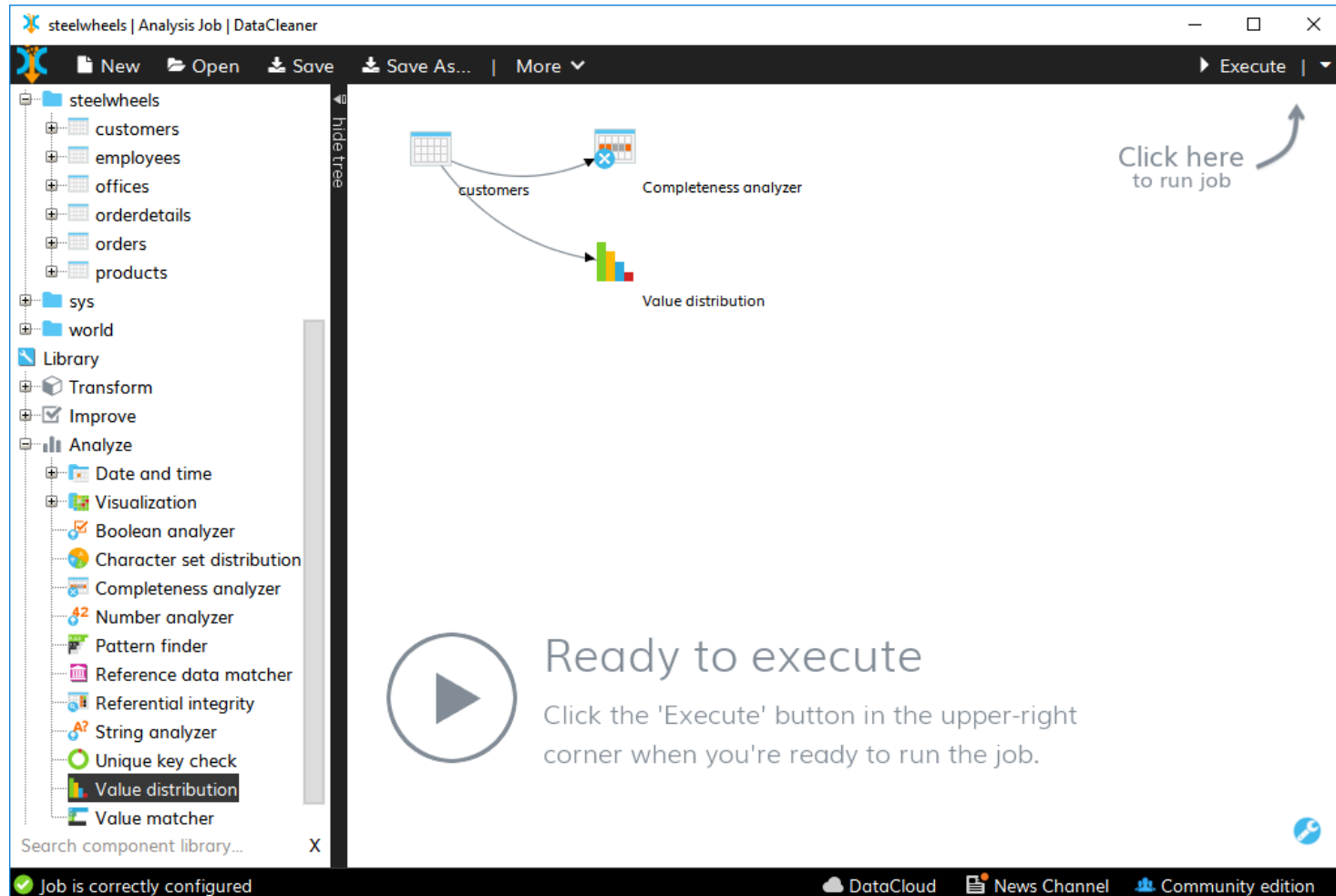
Save dataset

CUSTOMERNUMBER ^	ADDRESSLINE1	CITY	STATE	COUNTRY
103	54, rue Royale	Nantes	<null>	France
119	67, rue des Cinquante Otages	Nantes	<null>	France
121	Erling Skakkes gate 78	Stavern	<null>	Norway
125	ul. Filtrowa 68	Warszawa	<null>	Poland
128	Lyonerstr. 34	Frankfurt	<null>	Germany
141	C/ Moralzarzal, 86	Madrid	<null>	Spain
144	Berguvsvägen 8	Luleå	<null>	Sweden
145	Vinbæltet 34	København	<null>	Denmark
146	2, rue du Commerce	Lyon	<null>	France
148	Bronz Sok., Bronz Apt. 3/6 Tesvikiye	Singapore	<null>	Singapore
166	Village Close - 106 Linden Road Sandown	Singapore	<null>	Singapore
167	Drammen 121, PR 744 Sentrum	Bergen	<null>	Norway
169	Estrada da saúde n. 58	Lisboa	<null>	Portugal
171	184, chaussée de Tournai	Lille	<null>	France
172	265, boulevard Charonne	Paris	<null>	France
186	Keskuskatu 45	Helsinki	<null>	Finland
187	Fauntleroy Circus	Manchester	<null>	UK
189	25 Maiden Lane	Dublin	<null>	Ireland
201	Berkeley Gardens 12 Brewery	Liverpool	<null>	UK
206	Penthouse Level, Suntec Tower Three, 8 Tem...	Singapore	<null>	Singapore
209	24, place Kléber	Strasbourg	<null>	France

Save results

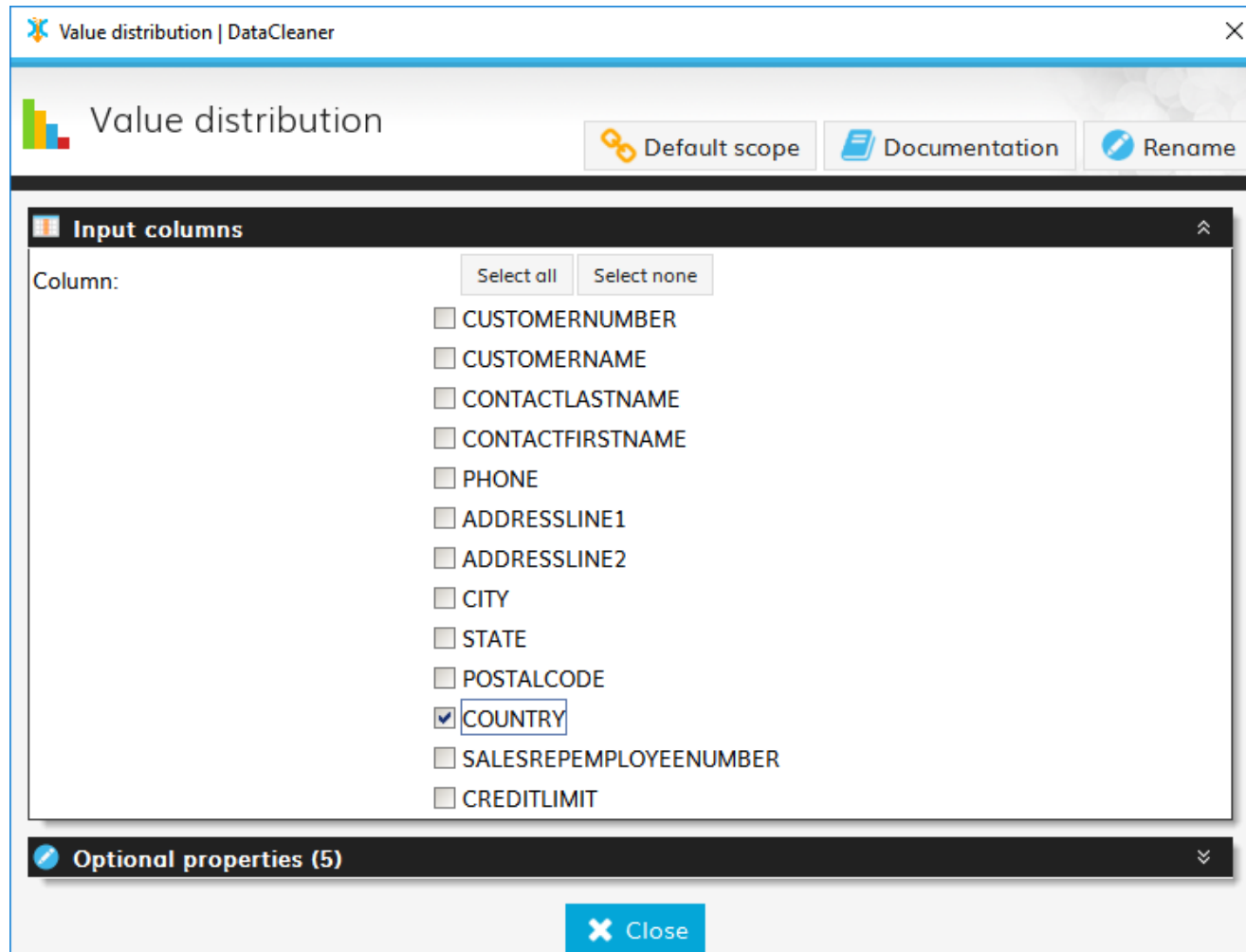
Data profiling

- Value distribution



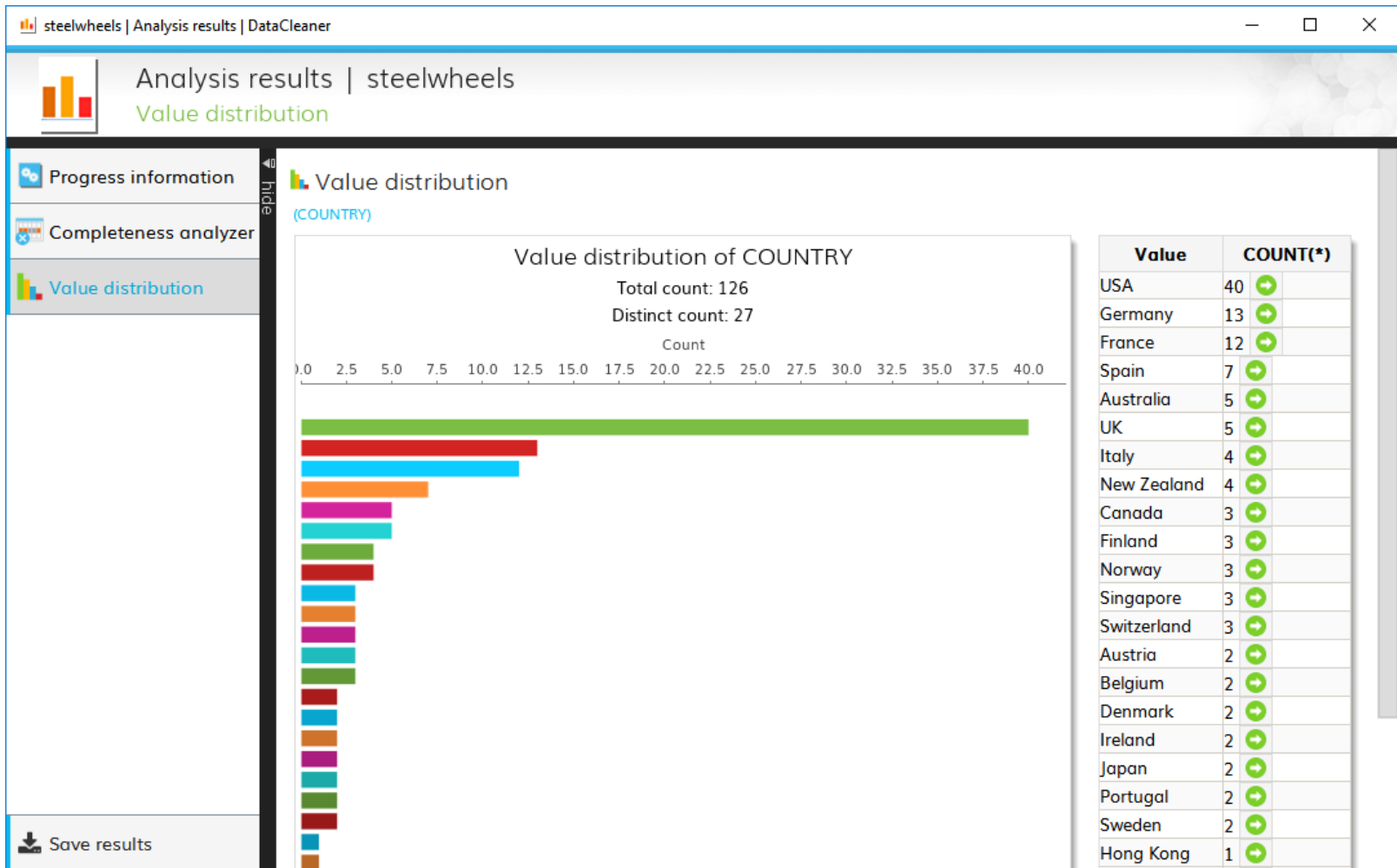
Data profiling

- Value distribution



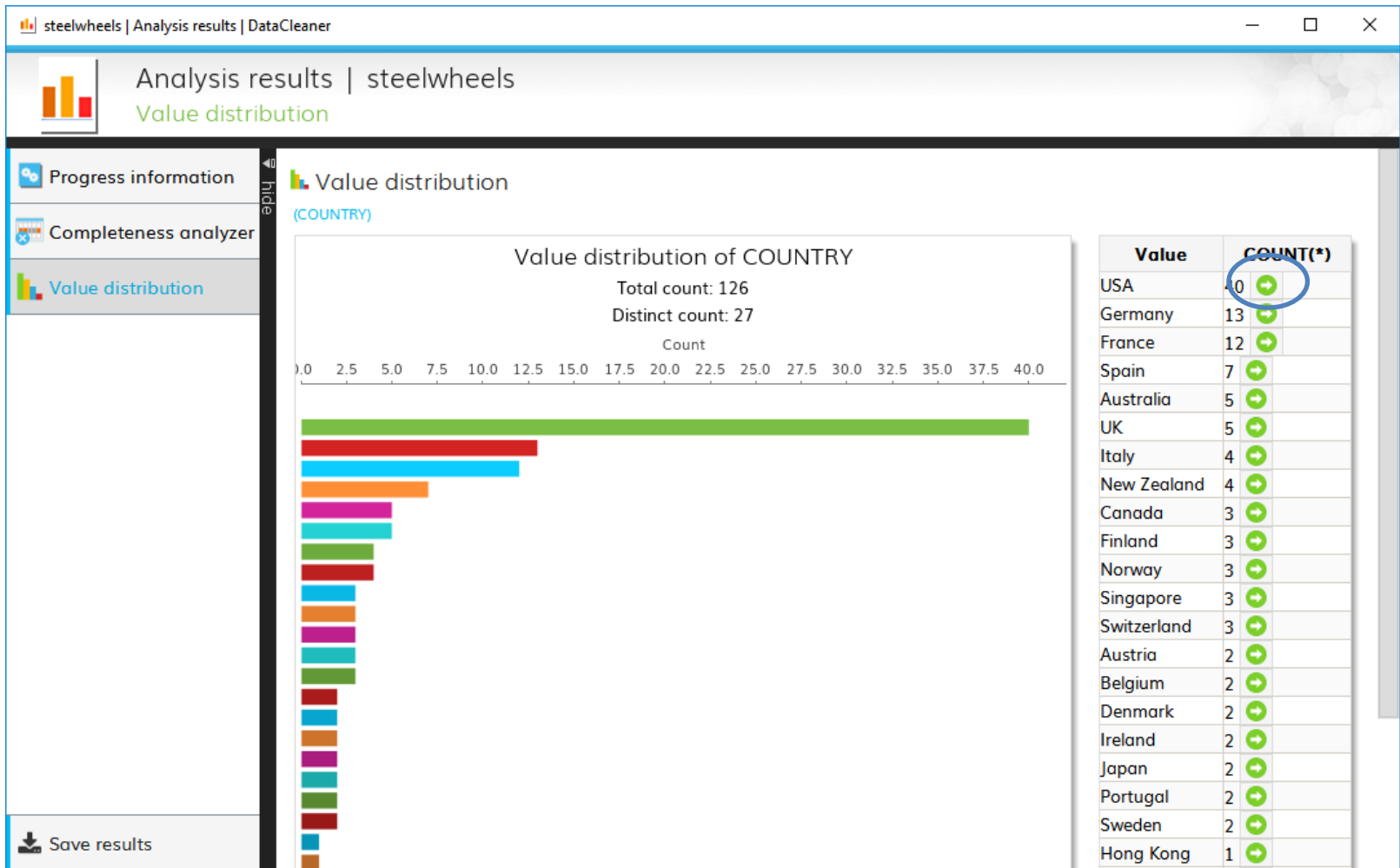
Data profiling

- Value distribution



Data profiling

- Value distribution



Data profiling

- Value distribution – detailed results

Detailed results for [USA] | DataCleaner

Detailed results

Records (40) View detailed rows Save dataset

CUSTOMERNUMBER	ADDRESSLINE1	CITY	STATE	POS...	COUNTRY
97	1302 PONCE DE L...	ST AUGUSTINE	FL	32084	USA
98	123 Sesame Street	ST Cloud	FL	34769	USA
175	25593 South Bay L...	Bridgewater	CT	97562	USA
198	16780 Pompton St.	Brickhaven	MA	58339	USA
204	7635 Spinnaker Dr.	Brickhaven	MA	58339	USA
205	78934 Hillside Dr.	Pasadena	CA	90003	USA
219	4097 Douglas Av.	Glendale	CA	92561	USA
239	361 Furth Circle	San Diego	CA	91217	USA
286	39323 Spinnaker Dr.	Cambridge	MA	51247	USA
319	3758 North Pendl...	White Plains	NY	24067	USA
320	4575 Hillside Dr.	New Bedford	MA	50553	USA
321	7734 Strong St.	San Francisco	CA	94217	USA
328	7476 Moss Rd.	Newark	NJ	94019	USA
339	782 First Street	Philadelphia	PA	71270	USA
347	6047 Douglas Av.	Los Angeles	CA	91003	USA
362	8616 Spinnaker Dr.	Boston	MA	51003	USA
363	2304 Long Airport ...	Nashua	NH	62005	USA
379	7825 Douglas Av.	Brickhaven	MA	58339	USA
424	5905 Pompton St.	NYC	NY	10022	USA
447	2440 Pompton St.	Glendale	CT	97561	USA
450	3086 Ingle Ln.	San Jose	CA	94217	USA

Data profiling

- String analysis

The screenshot displays the DataCleaner software interface. On the left, a 'hide tree' sidebar lists various data sources and analysis components. The 'Analyze' section is expanded, showing a list of analysis tools, with 'String analyzer' highlighted. The main workspace shows a workflow diagram where a 'customers' data source is connected to three analysis components: 'Completeness analyzer', 'Value distribution', and 'String analyzer'. An arrow points from the 'String analyzer' to a large play button icon. In the upper right corner, there is a button labeled 'Execute' with a dropdown arrow. A callout bubble points to this button with the text 'Click here to run job'. At the bottom, a status bar indicates 'Job is correctly configured' and includes links to 'DataCloud', 'News Channel', and 'Community edition'.

steelwheels | Analysis Job | DataCleaner

New Open Save Save As... More

Execute

customers

Completeness analyzer

Value distribution

String analyzer

Click here to run job

Ready to execute

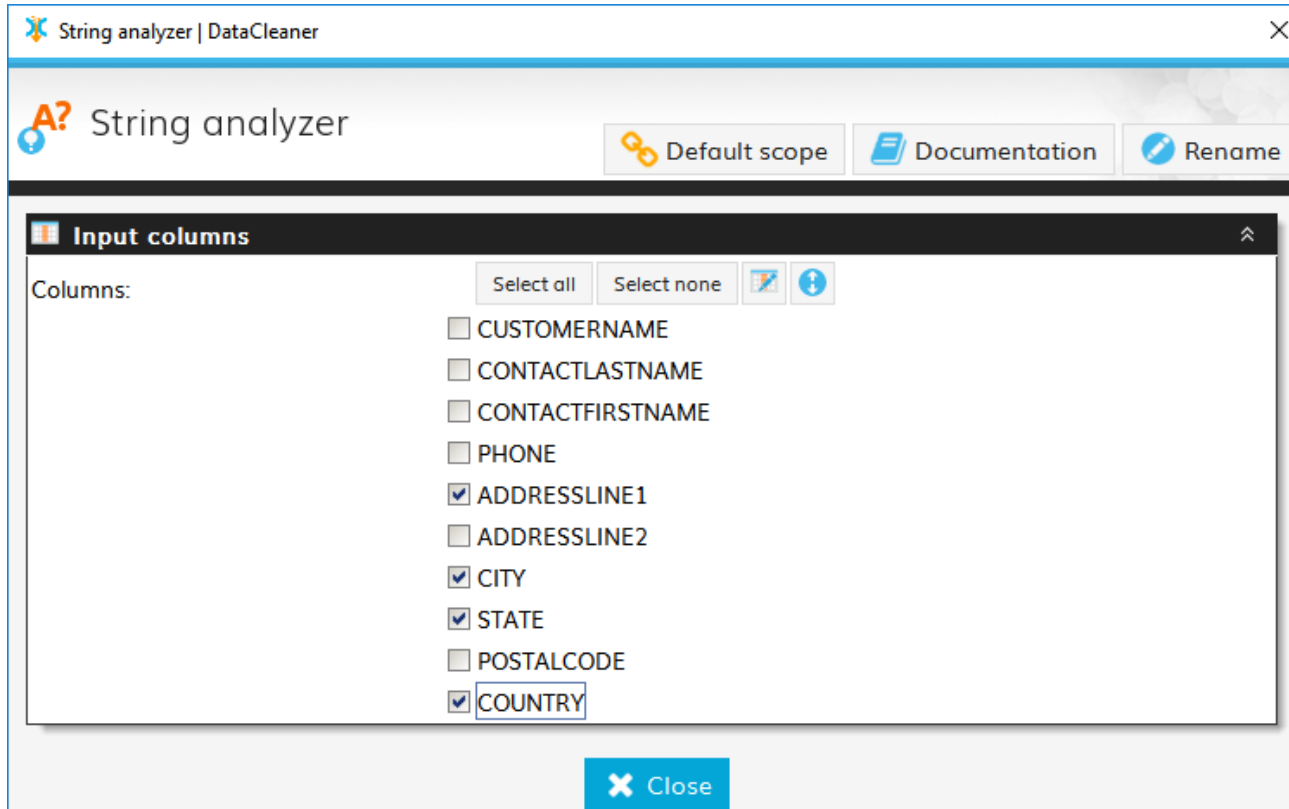
Click the 'Execute' button in the upper-right corner when you're ready to run the job.

Job is correctly configured

DataCloud News Channel Community edition

Data profiling

- String analysis



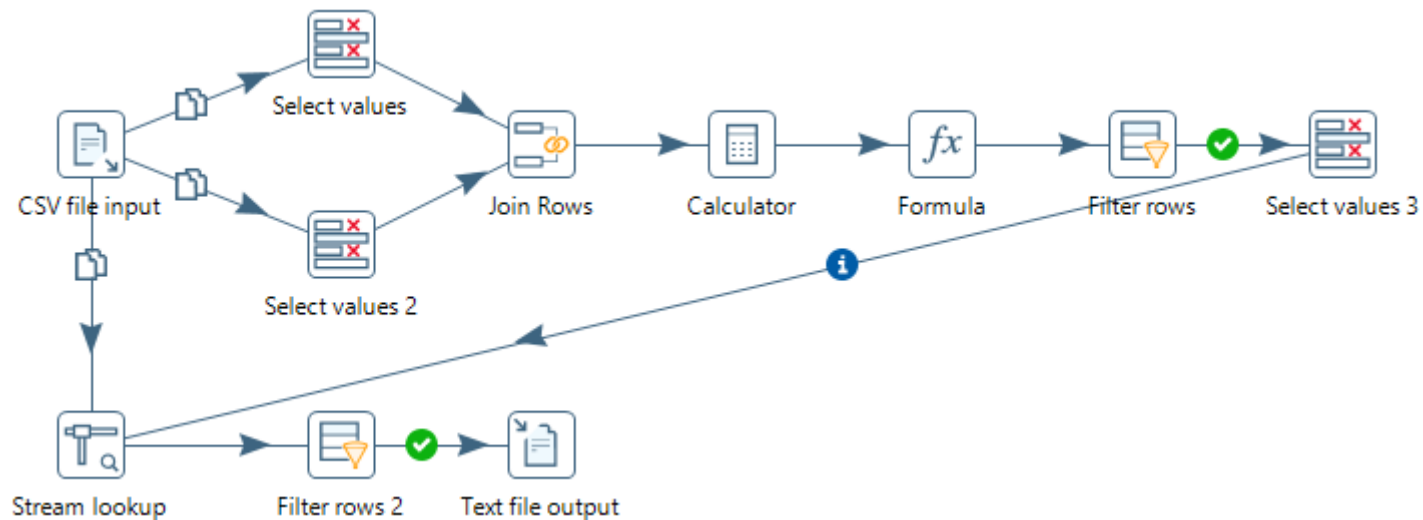
Data profiling

- String analysis

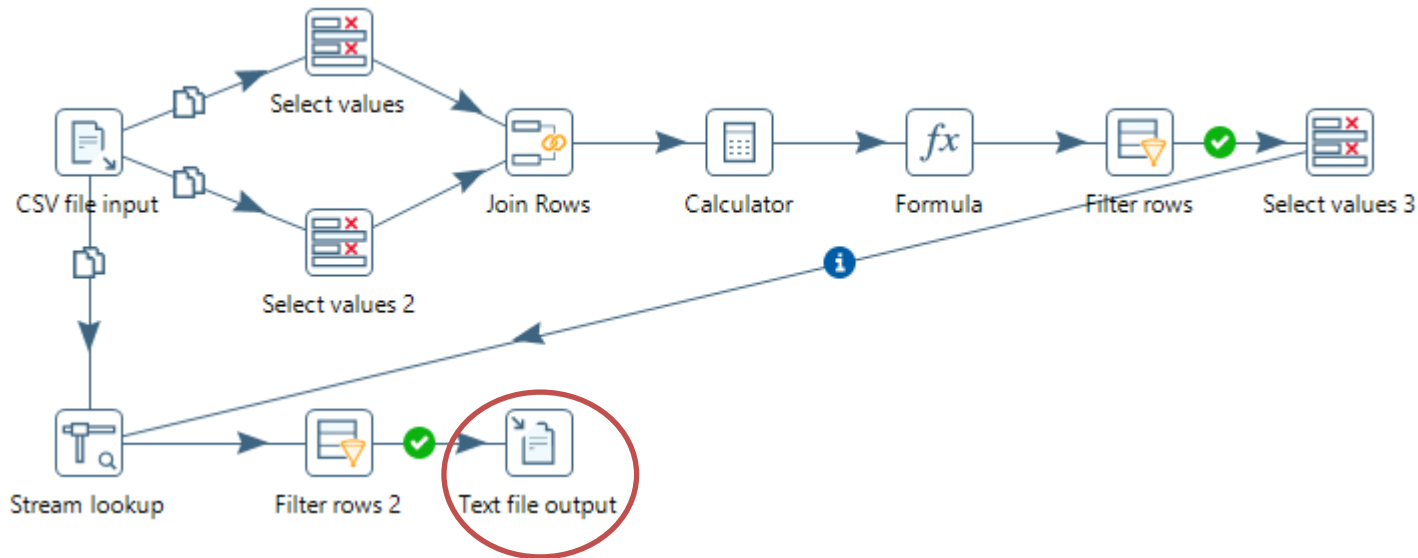
Analysis results steelwheels					
String analyzer					
(ADDRESSLINE1,CITY,STATE,COUNTRY)					
	ADDRESSLINE1	CITY	STATE	COUNTRY	
Row count	126	126	126	126	
Null count	0	0	74	0	
Blank count	0	0	0	0	
Entirely uppercase count	1	6	44	45	
Entirely lowercase count	0	0	0	0	
Total char count	2474	990	153	709	
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	

Data profiling

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



Duplicate detection

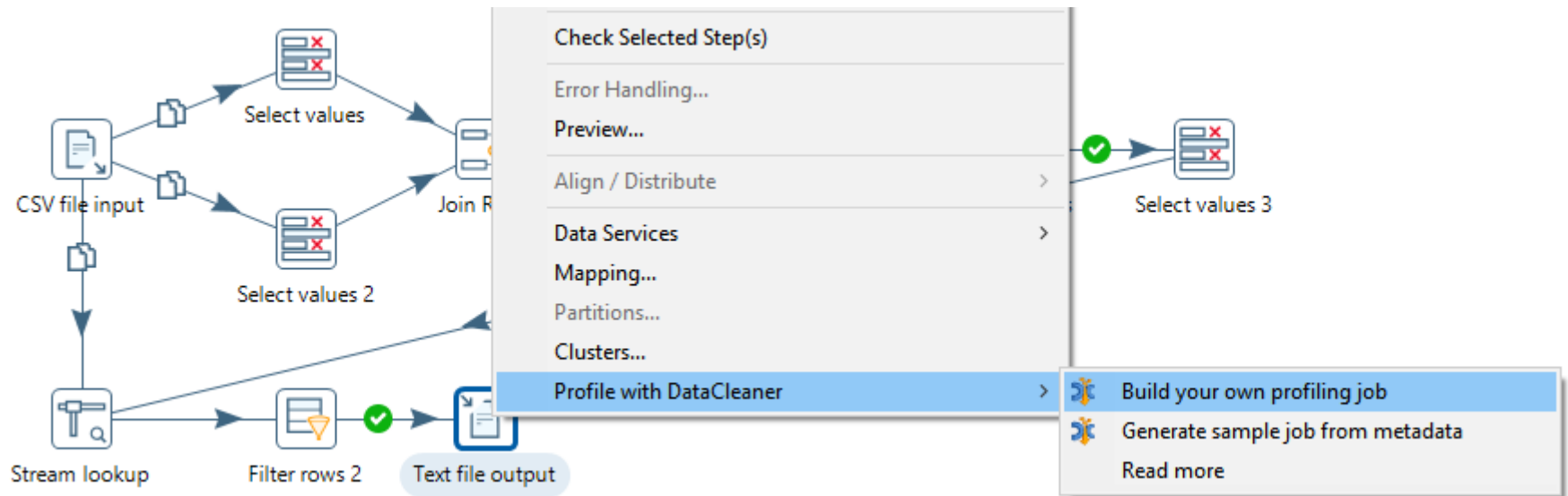


Rows of step: Text file output (64 rows)

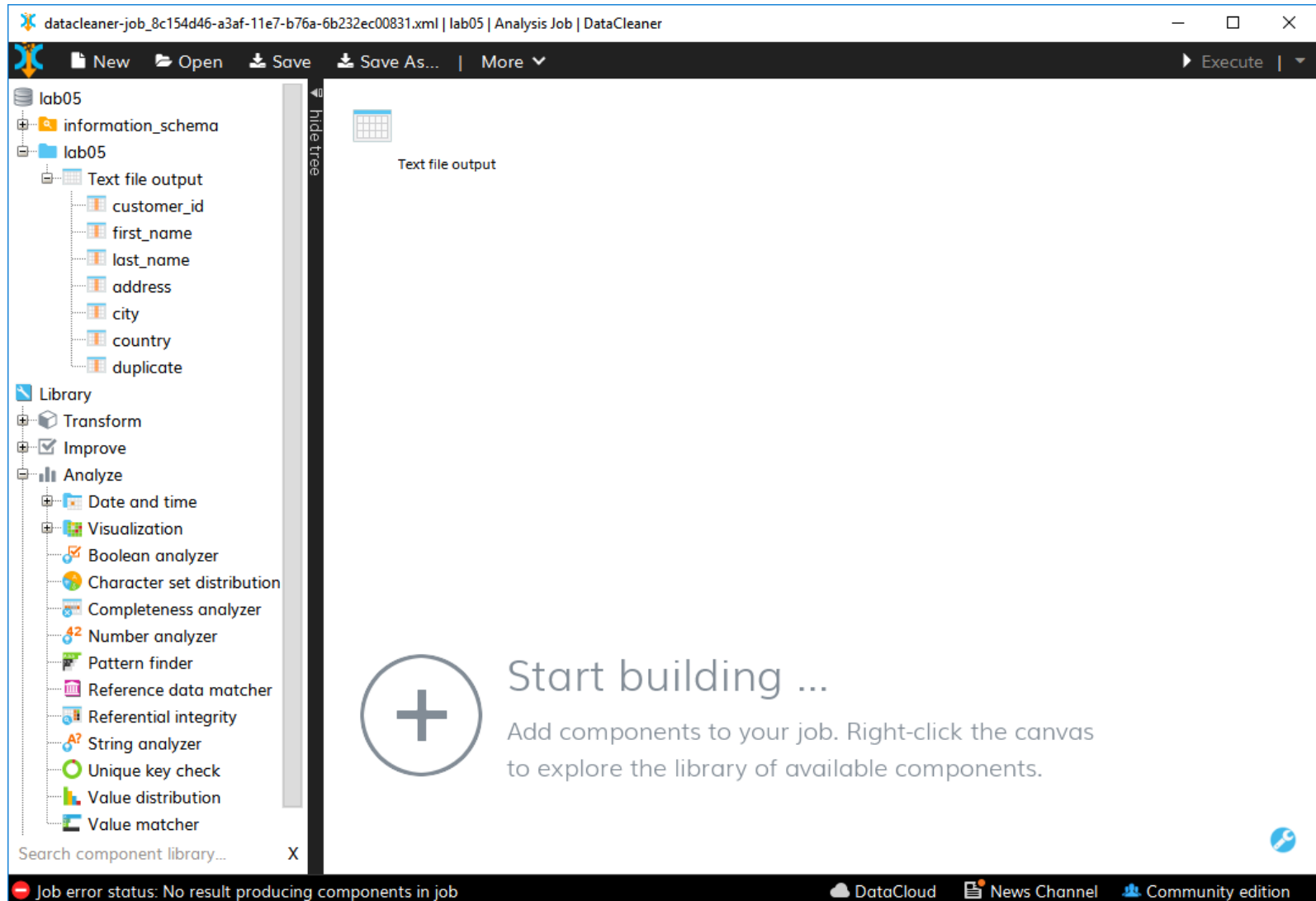
#	customer_id	first_name	last_name	address	city	country	duplicate
1	103	Carine	Schmitt	54, rue Royale	Nantes	France	<null>
2	119	Janine	Labruno	67, rue des Cinquante Otages	Nantes	France	<null>
3	121	Jonas	Bergulfsen	Erling Skakkes gate 78	Stavern	Norway	<null>
4	125	Zbysgniew	Piastzeniewicz	ul. Filtrowa 68	Warszawa	Poland	<null>
5	128	Roland	Keitel	Lyonerstr. 34	Frankfurt	Germany	<null>
6	141	Diego	Freyre	c/ Morazarzal, 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

- Example
 - perform data profiling after duplicate elimination



Data profiling



Data profiling

The screenshot displays the DataCleaner application window. The title bar shows the file path and job name: "datacleaner-job_8c154d46-a3af-11e7-b76a-6b232ec00831.xml | lab05 | Analysis Job | DataCleaner". The interface includes a menu bar with "New", "Open", "Save", "Save As...", and "More". On the left, a "lab05" project tree shows a folder "information_schema" containing a "lab05" folder, which lists fields: "customer_id", "first_name", "last_name", "address", "city", "country", and "duplicate". Below this is a "Library" section with categories: "Transform", "Improve", and "Analyze". The "Analyze" category is expanded, showing various tools like "Date and time", "Visualization", "Boolean analyzer", "Character set distribution", "Completeness analyzer", "Number analyzer", "Pattern finder", "Reference data matcher", "Referential integrity", "String analyzer", "Unique key check", "Value distribution" (highlighted), and "Value matcher". A search bar at the bottom left says "Search component library...". The main workspace shows a workflow diagram with a "Text file output" component connected to a "Value distribution" component. A large play button icon and the text "Ready to execute" are prominently displayed, along with instructions: "Click the 'Execute' button in the upper-right corner when you're ready to run the job." An arrow points to the "Execute" button in the top right corner. The bottom status bar indicates "Job is correctly configured" and includes links for "DataCloud", "News Channel", and "Community edition".

datacleaner-job_8c154d46-a3af-11e7-b76a-6b232ec00831.xml | lab05 | Analysis Job | DataCleaner

New Open Save Save As... More Execute

lab05

- information_schema
 - lab05
 - Text file output
 - customer_id
 - first_name
 - last_name
 - address
 - city
 - country
 - duplicate

Library

- Transform
- Improve
- Analyze
 - Date and time
 - Visualization
 - Boolean analyzer
 - Character set distribution
 - Completeness analyzer
 - Number analyzer
 - Pattern finder
 - Reference data matcher
 - Referential integrity
 - String analyzer
 - Unique key check
 - Value distribution
 - Value matcher

Search component library... X

Text file output Value distribution

Click here to run job

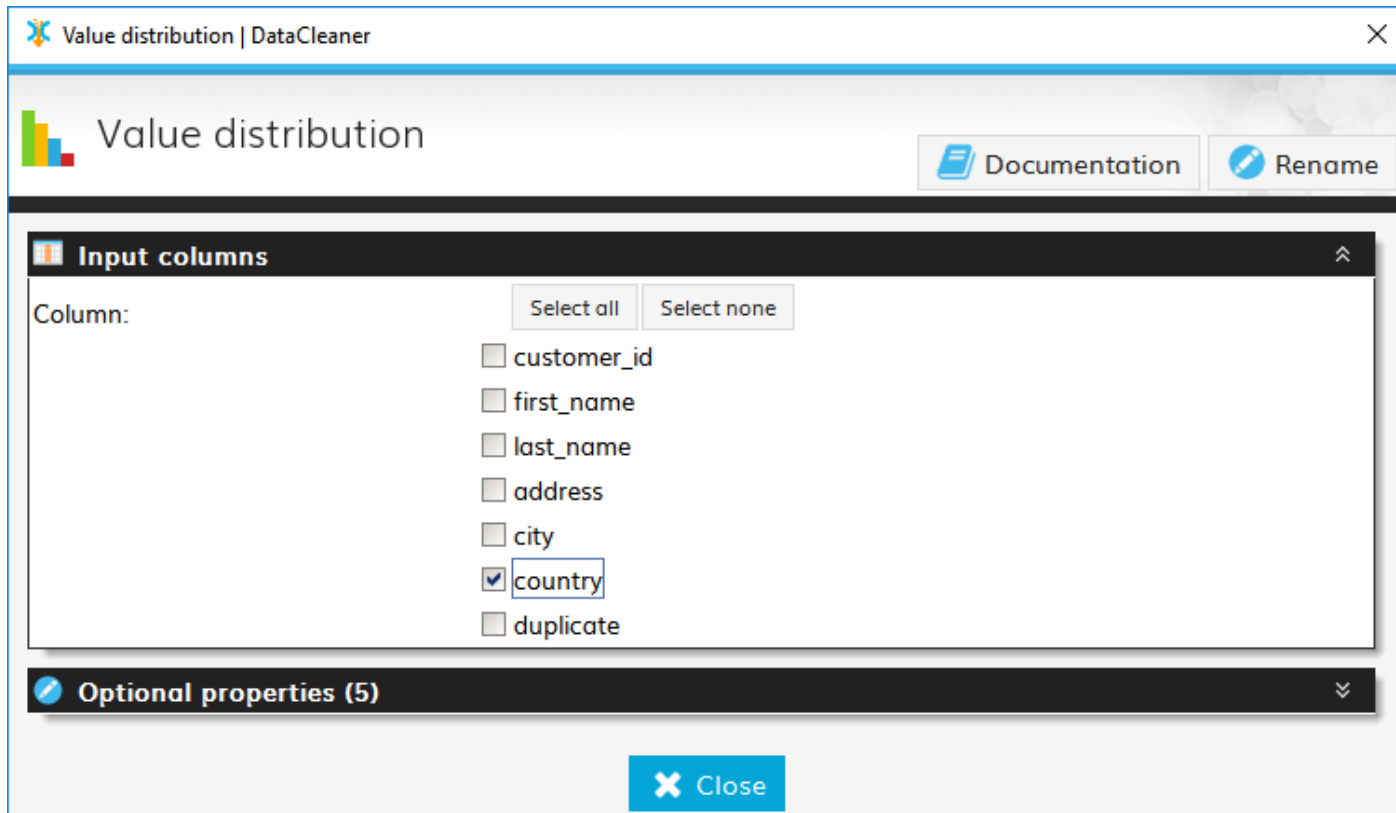
Ready to execute

Click the 'Execute' button in the upper-right corner when you're ready to run the job.

Job is correctly configured DataCloud News Channel Community edition

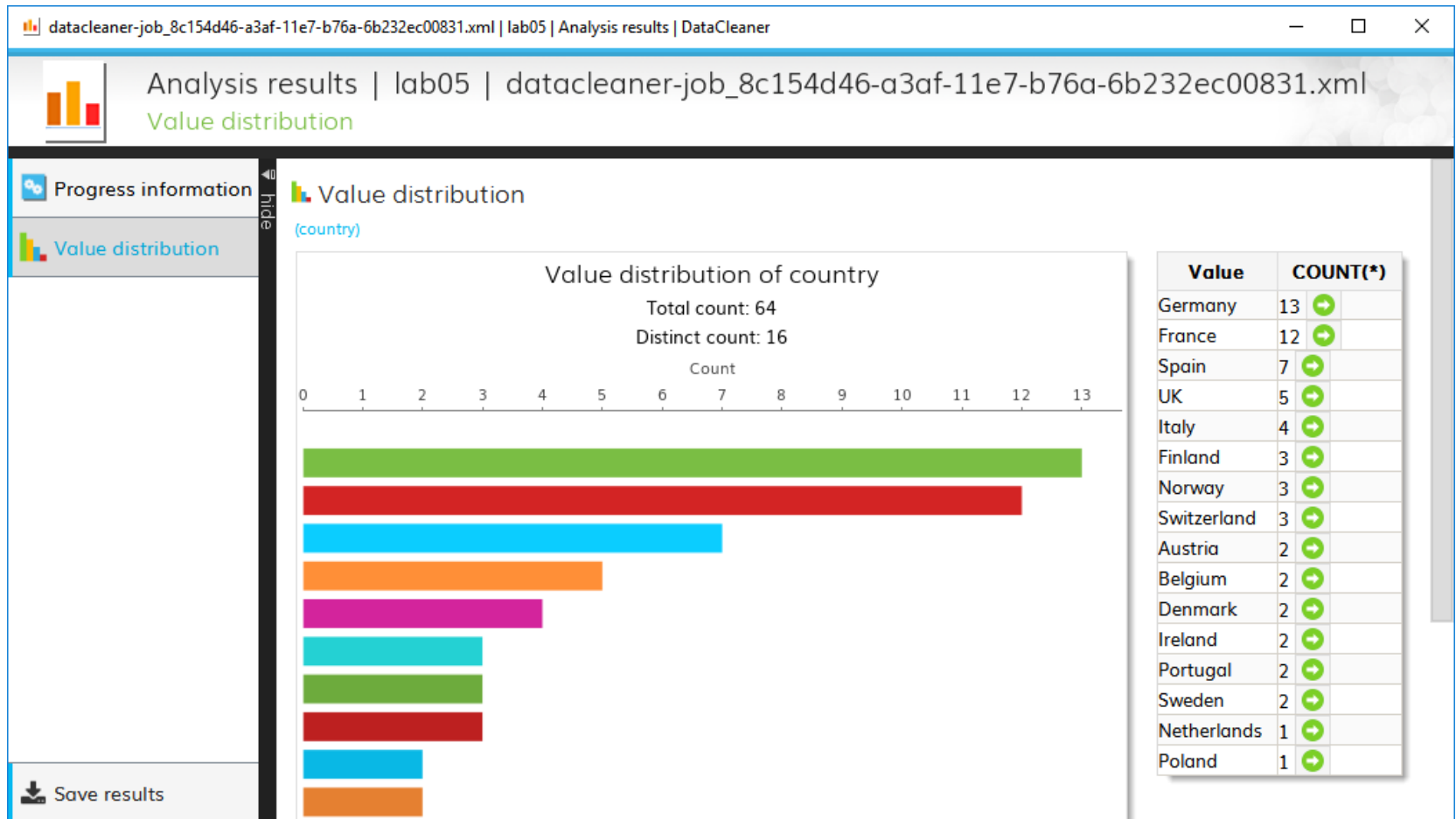
Data profiling

- Value distribution



Data profiling

- Value distribution



Data profiling

The screenshot displays the DataCleaner application window. The title bar reads "datacleaner-job_67a84087-a42e-11e7-b651-4379f2e87310.xml | lab05 | Analysis Job | DataCleaner". The menu bar includes "New", "Open", "Save", "Save As...", and "More". The left sidebar contains a "lab05" tree view with "information_schema" and "lab05" folders. Under "lab05", there is a "Text file output" component and a list of fields: "customer_id", "first_name", "last_name", "address", "city", "country", and "duplicate". Below this is a "Library" section with categories "Transform", "Improve", and "Analyze". The "Analyze" category is expanded, showing various analysis tools: "Date and time", "Visualization", "Boolean analyzer", "Character set distribution", "Completeness analyzer", "Number analyzer", "Pattern finder", "Reference data matcher", "Referential integrity", "String analyzer" (highlighted), "Unique key check", "Value distribution", and "Value matcher". The main workspace shows a workflow diagram with a "Text file output" component connected to "Value distribution" and "String analyzer". A large play button icon and the text "Ready to execute" are prominently displayed. A callout bubble points to the "Execute" button in the top right corner, stating "Click here to run job". At the bottom, a status bar indicates "Job is correctly configured" and provides links to "DataCloud", "News Channel", and "Community edition".

datacleaner-job_67a84087-a42e-11e7-b651-4379f2e87310.xml | lab05 | Analysis Job | DataCleaner

New Open Save Save As... More Execute

lab05

- information_schema
- lab05
 - Text file output
 - customer_id
 - first_name
 - last_name
 - address
 - city
 - country
 - duplicate

Library

- Transform
- Improve
- Analyze
 - Date and time
 - Visualization
 - Boolean analyzer
 - Character set distribution
 - Completeness analyzer
 - Number analyzer
 - Pattern finder
 - Reference data matcher
 - Referential integrity
 - String analyzer
 - Unique key check
 - Value distribution
 - Value matcher

Search component library... X

Text file output

Value distribution

String analyzer

Click here to run job

Ready to execute

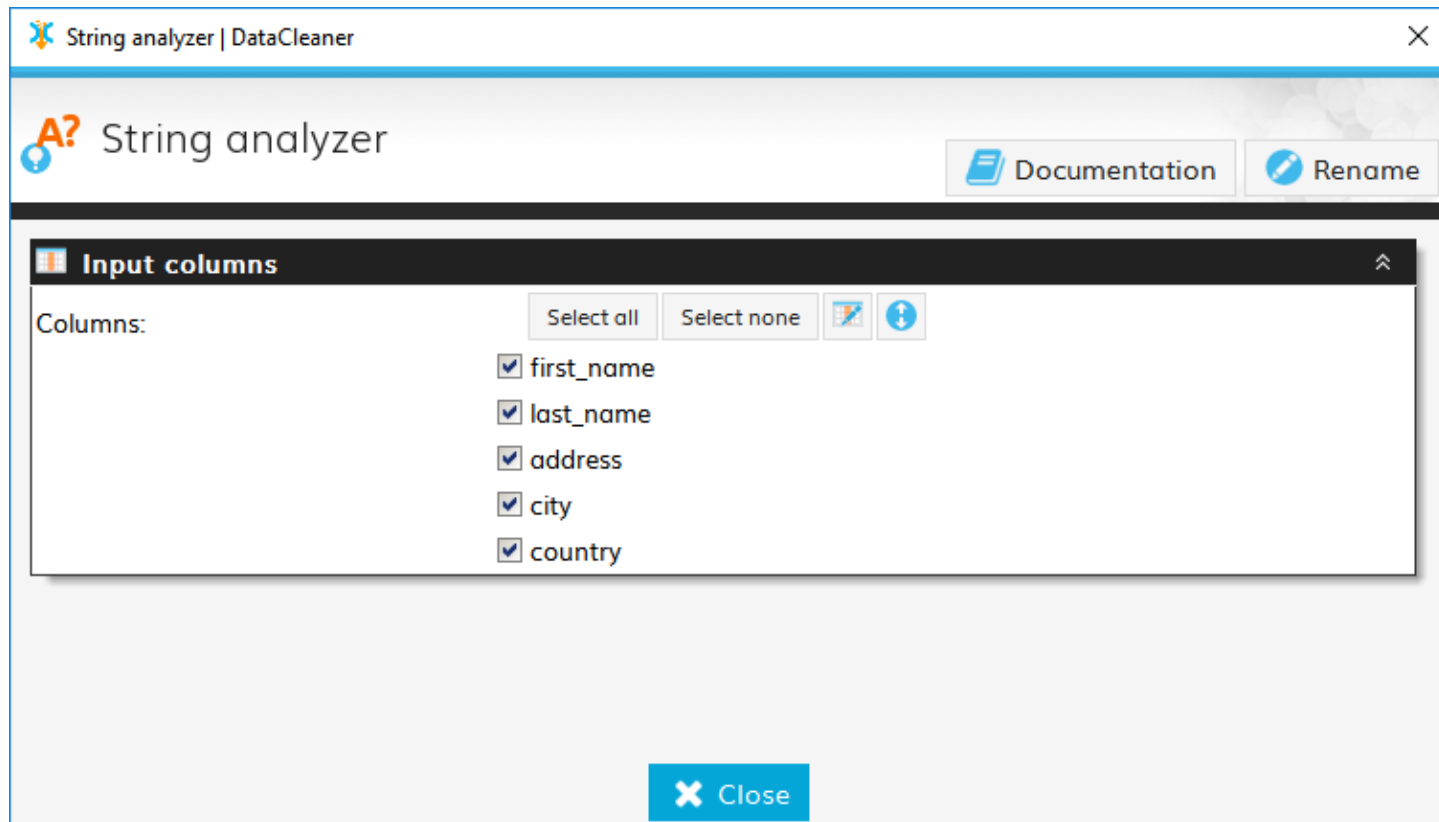
Click the 'Execute' button in the upper-right corner when you're ready to run the job.

Job is correctly configured

DataCloud News Channel Community edition

Data profiling

- String analysis



Data profiling

- String analysis

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