

DataGraft Service for Data Publishing and Transformation: User Guide

This tutorial describes core functionality of the DataGraft platform and gives detailed step-by-step explanation of data publishing and transformation process with help of DataGraft portal.

Platform overview

The most demonstrable way to get an overview of what can be done with help of DataGraft platform is to explore data pages and data transformations that other users of this platform chose to share. The two terms mentioned above, **data pages** and **data transformations**, are two main concepts you work with while using DataGraft platform, therefore it may be useful to understand what each of them means in a context of the service.

Data pages contain cleaned and transformed data you want to publish. Each data page is stored with some metadata, including data page name; short description; keywords, describing a data page; owner; creation date and public/private property (see Figure 1). The latter is defined by data page owner and specifies whether this data page can be explored by other users of a platform or not.

The screenshot shows a web form titled "Data page properties". It contains the following fields and values:

- Name:** My first data page
- Description:** This data page contains information I wish to publish
- Expose as public:** An unchecked checkbox.
- Owner:** newUser
- Creation Date:** 12 Aug 2015
- Keyword:** A row of tags: "keywords", "relevant", "for", and "datapage", each with a close button (x). To the right of these tags is a text input field containing "Keywords...".

Figure 1: Data page properties

Users that have access to the data page (i.e. just owners in case of private pages and everyone else for public pages) can locally download information associated with a data page in a raw tabular format by pressing [EXPORT RAW](#) button; or as RDF by pressing [EXPORT RDF](#). The list of supported RDF formats includes RDF/XML(.rdf), n-triple(.nt), turtle(.ttl), n3(.n3), nquads(.nq), RDF/JSON(.rj). [View portal, SPARQL](#)

Another type of asset that users may create and share in DataGraft is **data transformation**. Before publishing data, in most cases you will need to transform the original dataset – clean messy data, remove unnecessary information, probably add some new data fields and convert tabular data to RDF. This sequence of operations you perform on your data to convert it to desirable form is called data transformation. The greatest thing about data transformations in DataGraft platform is that you may reuse them repeatably on another datasets, share them with other users, modify existing transformations and create new transformations by extending ones that you or other users created and shared. As well as data pages, data transformations are stored with some metadata. The transformation properties include transformation name, description, owner and public/private property. They are defined on a data transformation creation stage(see).

User Registration

In order to create your own data pages and transformations through a platform, you should first sign up for DataGraft account. After registration you are automatically redirected to data page creation service, from where you may start a process of creating your first data page. This process is described in detail in section .

User Dashboard

User dashboard helps to manage data pages and data transformations created by user. The dashboard view gives you an overview of your data pages and transformations with search capabilities.

Publishing data

Publishing data with help of DataGraft platform is rather simple process. You may start by switching to a [Publish](#) tab in a main menu. The first thing you do is uploading data. To do so you just drop your dataset file in a raw CSV or RDF format in a white frame under "Upload your data" label.

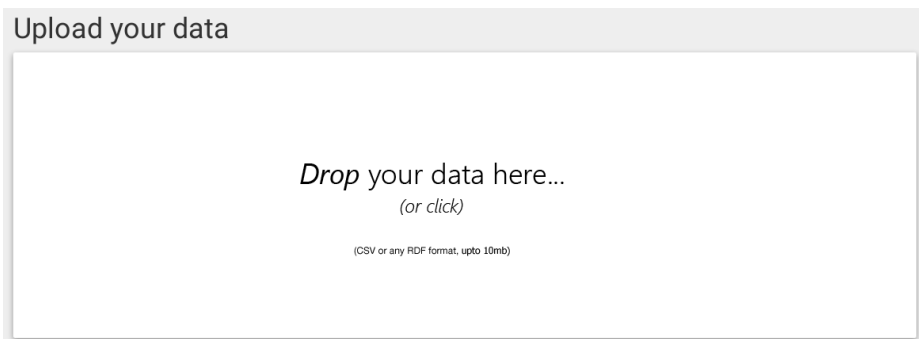


Figure 2: Uploading data

After data is successfully uploaded (this is indicated by green mark in the top right corner of a file icon) you have several options:

1. Create data page from a row data without applying any transformation on it.
2. Transform your data before creating a datapage. This in turn can be done in two ways
 - (a) By creating a new transformation to use on your data
 - (b) By applying an existing transformation to your data

Let's go through the most simple scenario by choosing the first alternative. To do this you just click **CREATE RAW** button. This automatically takes you to the next page where you specify data page properties (see Figure 1). After everything is in order, you simply click **CREATE DATA PAGE**. And that's it, you have just created your very first data page. Now you (and other users in case if you defined this data page as public) have access to the data page, are able to download associated data, add more information and features to the created asset.

However in most cases you still need to process your data before publishing it. In this case you should use the transformation service. By clicking **CREATE USING A NEW TRANSFORMATION** button you may start transforming your data.

Data cleaning and transformation

This section explains how tabular data is transformed in DataGraft platform and gives you the best strategies for data transformation.


In the basis of DataGraft data transformations there lies a Grafter DSL (Domain Specific Language), which in its turn is implemented in Clojure. Therefore, to

take maximum advantage of the service, one should be acquainted with mentioned languages. However number of transformations, depending on their complexity, can be done through intuitive and user-friendly GUI without any coding.

Transformation metadata

The first tab seen in the transformation creation window is "Metadata". Here user defines transformation title and gives a short description of how this transformation processes targeted data. If you wish to share transformation, it is possible to expose it as public. In this case other platform users will be able to explore and use given transformation.

After describing metadata, you may save transformation by clicking "Save"

button  in the top right corner. Transformation may be as well saved later at any moment.

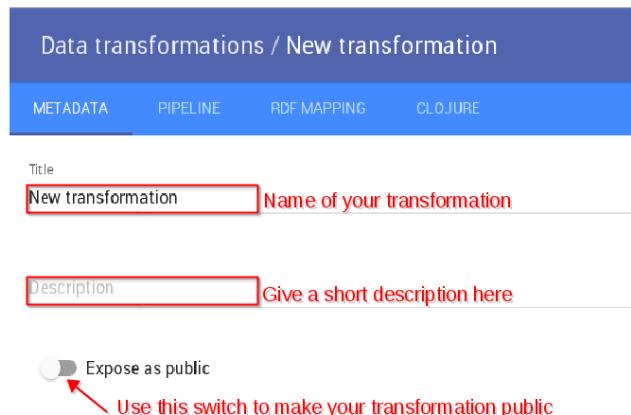



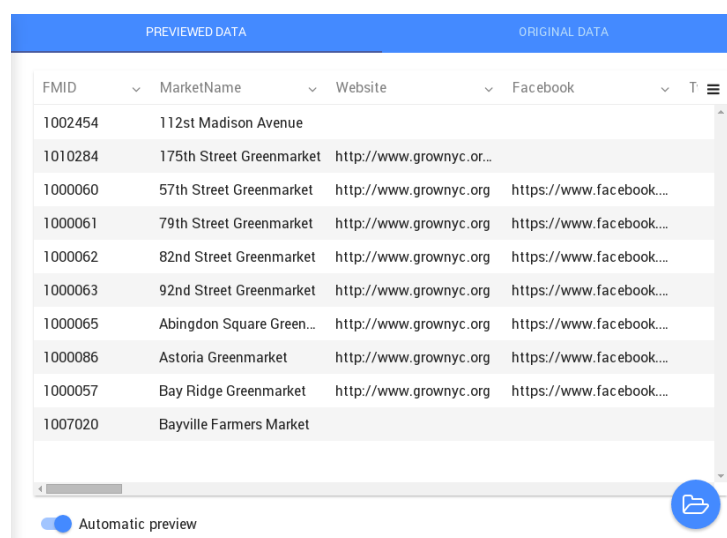


Figure 3: Transformation properties

Transformation preview

After your transformation was saved, in the bottom right corner you may see this icon: . By clicking on it you may immediately apply transformation being created to target data. This may be whether distribution already uploaded to a platform, or alternatively you may upload a new file. For the first option use button depicted as , and for the second – this one . In this way you are able to see the instant effect of created transformation on your data in the preview area.

Preview area is located in the right part of transformation window. You can see two tabs there – one with original data and another with changes made through transformation pipeline. Each time you modify a pipeline, the transformation is applied to the previewed dataset immediately, so you can see the effect of each performed step. You may adjust preview settings to check and evaluate transformation steps you are creating. Thus, it is possible to hide columns and to sort visible data. The changes made through these settings are not part of transformation and affect just previewed data. However at any time you may export tabular data either as it looks in preview or in a format it has at current stage of your transformation.



FMID	MarketName	Website	Facebook	T
1002454	112st Madison Avenue			
1010284	175th Street Greenmarket	http://www.grownyc.or...		
1000060	57th Street Greenmarket	http://www.grownyc.org	https://www.facebook...	
1000061	79th Street Greenmarket	http://www.grownyc.org	https://www.facebook...	
1000062	82nd Street Greenmarket	http://www.grownyc.org	https://www.facebook...	
1000063	92nd Street Greenmarket	http://www.grownyc.org	https://www.facebook...	
1000065	Abingdon Square Green...	http://www.grownyc.org	https://www.facebook...	
1000086	Astoria Greenmarket	http://www.grownyc.org	https://www.facebook...	
1000057	Bay Ridge Greenmarket	http://www.grownyc.org	https://www.facebook...	
1007020	Bayville Farmers Market			

Figure 4: Transformation preview

Transformation pipeline

Data cleaning and transformation in DataGraft platform is performed with help of a “pipeline” concept. To begin with, each single transformation step is defined as a pipe – a function that performs simple data conversion on its input. The greatest fact about these functions is that they may be combined together in a such way, that output of one pipe acts as an input for another. Obviously, the input/output data, that travels through this pipeline is dataset being transformed. This way of composing operations gives a great flexibility and allows to perform rather complex data conversions.

The first implicit step of every transformation pipeline is getting the very first pipeline input. Therefore, each transformation starts from reading a dataset from uploaded file. However you do not need to include this step into your

pipeline manually, since this action is performed automatically for each transformation.

To add a first transformation step click the  button next to the pipeline



Figure 5: Add pipe function

Now you can see the list of functions you may use to modify uploaded dataset. Available functions are logically grouped according to the type of effect they have on data. Consequently, OPERATIONS ON COLUMNS add, remove or modify dataset columns, while OPERATIONS ON ROWS extract certain rows from a dataset based on row numbers or some condition that user defines. Operations "Make dataset" and "Reshape dataset" affect the entire dataset.



Figure 6: List of available functions

For each operation you can see a short documentation with simple illustrated example by clicking "show/hide documentation" button. For every pipeline function you create you may leave a short description note in the "Comment" field. This information helps you and other users of your transformation to understand operations that are performed here. If you ignore this field, the note will be created automatically based on function parameters you have specified.

Once you have added a new function to the pipeline, it will instantly appear in the pipeline view. You are free to change function parameters any time you need it by simply clicking on correspondent function icon. To get a short information about actions performed by this function you may just hover mouse

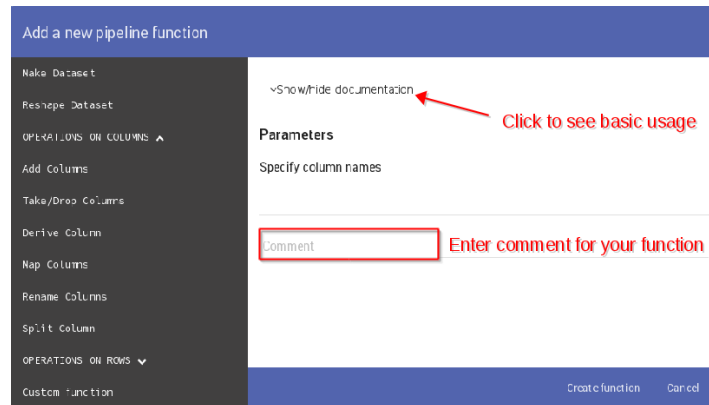



Figure 7: Common operations for all pipe functions

pointer over its name. In many cases function order significantly affects the transformation result. It is very simple to change this order by just dragging function icons along the pipeline. To remove a function click  button next to the function you would like to remove.

The following sections provide you with detailed guidelines for each function usage.

Make Dataset

As its name suggests "Make dataset" operation creates new dataset from its input. If you leave all parameter fields blank new dataset will be created from all the input columns with column names given as simple alphabetic sequence. By checking "move first row to header" option you get all the column names from the first row. First row will be removed from your dataset. You may as well specify column names you wish to see in a new dataset or fetch first n columns.

Reshape Dataset

Reshape dataset "melts" given dataset in a such way, that each row of new dataset represents a unique combination of variables and values for given column array. The best way to explain this function is through an example.

Add Columns

The "Add columns" results in adding a new column(s) to a given dataset. To add a new column you should specify column name and a value for a new column. This value will be copied into every row within dataset. You may as well populate new column with some custom value. This may be current date, dataset filename, row index or custom Clojure code. Note that the "Expression"

field is prioritized, in other words if you define both value and expression, only expression will be used to get a value for a new column. You are able to add as many columns as you need within one operation. To add one more column simply click

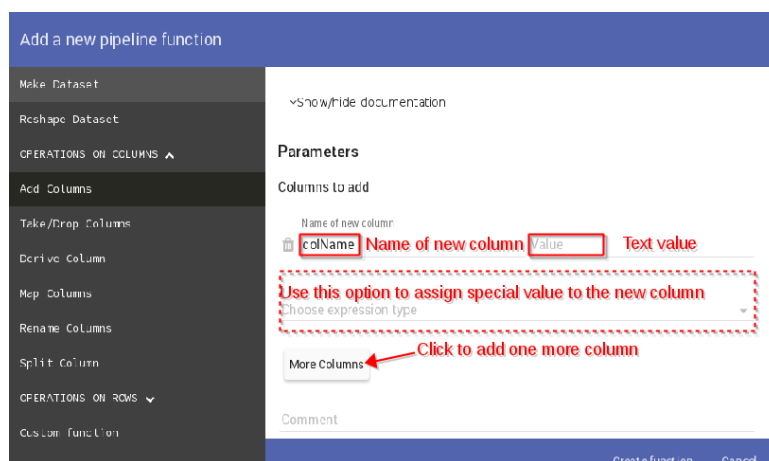


Figure 8: Add column function

Take/Drop Columns

The "Take/Drop Columns" function narrows given dataset. You may explicitly define columns you wish to see in a dataset by listing their names, get first n columns or remove unnecessary columns from a dataset. In second case, dataset will be narrowed to specified number of columns with their names assigned as alphabetical sequence("A", "B", "C" etc.). If more than 26 columns are fetched column names will count "AA", "AB", "AC" ... "BA", "BB", "BC" etc.

Derive Column

This function creates a new column in a dataset by applying some transformation to existing column(or columns). To use this function you should specify name for a new column, define one or several columns you are going to use to obtain new value and specify a function you apply to them. This function can be chosen from a drop-down list, which contains some standard functions and custom utility functions. You may define custom utility functions by yourself using Clojure code(see) or use utility functions written by other users.

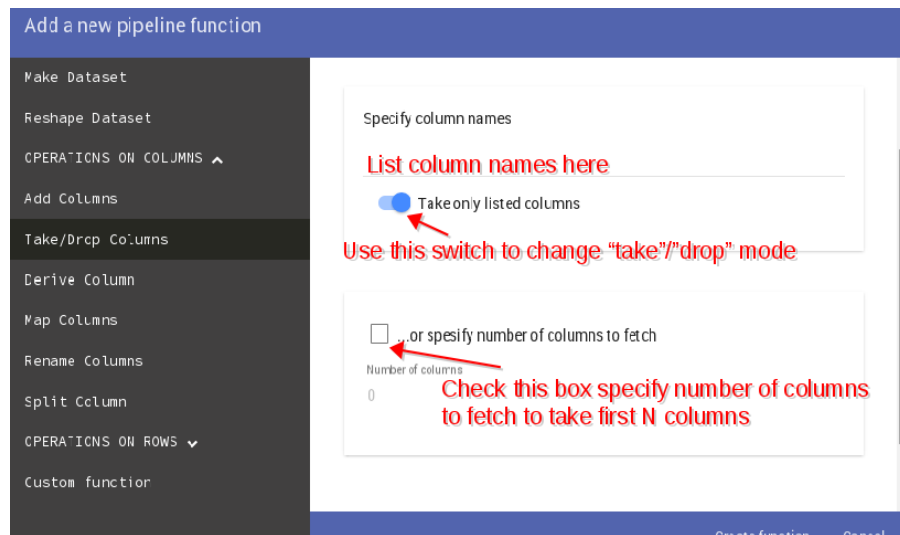



Figure 9: Columns function

One powerful feature of derive-column and similar pipeline functions is possibility to combine functions you apply to columns in the same way you combine functions in pipelines. Essentially, these are pipelines inside your dataset manipulation pipeline. To add one more "internal" function into your pipeline function, click Compose with button. Note, that function order is significant in this case. Functions are composed as they appear from top to down. To remove function from composition pipeline click trash icon  next to the function you wish to delete.

Drop Rows

This function just removes first n rows from a given dataset.


Filter Dataset

This method filters rows in the dataset for matches. This is a rather flexible filtering tool that allows you to filter your data in several different ways.

First field in in "Filter Dataset" parameter list specifies whether filter will be applied to specific column(s) or to all columns in dataset. For latter option just leave this field empty and matching will be performed for each column.

To perform actual data filtering you have several options.

First, it is possible to select only rows containing specified text. For this you enter the text in field marked as "Text to match". By default this text is treated as case-sensitive. To perform a case insensitive filtering check the box "ignore case" next to the text field.

You may as well filter dataset with help of regular expressions. By pressing  button next to the "Regular expression" field you will get short quick start tutorial for pattern usage.

Finally, you may filter dataset by applying utility functions to columns. Note, that the result of function(or combination of functions) will be treated a true/false expression.

The priority of listed option is defined as they appear - from top to down: if "Text to match" field is specified, other fields are ignored, if "Text to match" is ignored, but "Regular expression" is defined – this one will be used to filter your dataset ignoring functions below(if specified any).

Map Columns

This method allows you to apply function transformation to column and put the transformed value back to the same column. Transformation is done cell by cell. In parameter list you must specify a column you wish to change and a function that will be used to perform a transformation. You may add as many column-function pairs as you need.

Rename Columns

This operation allows you to rename columns in dataset. To get new names for columns you may either apply function(or function pipeline) to current column names or assign mapping from old to new column names directly.

Split Column

This operation allows to split column on some separator. Column is splitted once, leaving substring obtained as original text before first instance of separator in the source column and copying text after separator into new column. Separator itself is extracted from text. Name of new column is optional.

Take Rows

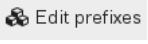
This operation allows you to narrow dataset just to first n rows.

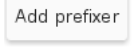
Adding Custom Function to Your Pipeline

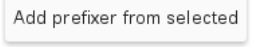
Some complex transformations cannot be done with help of operations described above. In this case you may need to define your own pipeline functions. This can be done with help of "Custom code" option using Clojure language

Creating and editing custom utility functions

Creating and editing prefixes

You may create and edit prefixes in the "Edit prefixes" window. To see this press  button in the pipeline view. Here you can see the list of all prefixes you created for current transformation. You may add a new prefix by

specifying its name and URI and pressing  button. Created prefix will instantly appear in the list of prefixes above. It is possible as well to create prefix by adding some string to the existing one. In this case select a prefix you wish to choose as base one, enter new prefix name and string value and

press . Now, if the base prefix is changed, changes will be as well applied to the child prefixes. To remove prefixes select all prefixes

you wish to remove and press .

Example 1. Here is given a small example, that illustrates how you can transform your data starting from the point where you upload raw csv data and up to the moment where you get the result in desired RDF format.