

# DATUM



**K S Ananth**  
CS22BTECH11029  
IIT Hyderabad

**Nimai Parsa**  
CS22BTECH11044  
IIT Hyderabad

**Kartikeya  
Mandapati**  
CS22BTECH11032  
IIT Hyderabad

**M Bhavesh  
Chowdary**  
CS22BTECH11041  
IIT Hyderabad

**Arugonda  
Srikar**  
CS22BTECH11008  
IIT Hyderabad

**Gunnam Sri  
Satya  
Koushik**  
CS22BTECH11026  
IIT Hyderabad

**Mohammed  
Gufran Ali**  
CS22BTECH11040  
IIT Hyderabad

September 2024

Supervisor **Dr. Ramakrishna Upadrasta**

# Abstract

This whitepaper presents a domain-specific language(DSL) which is designed for tasks related to data manipulation and visualization. With the recent developments in machine learning and data science applications, the importance of data is ever growing and is needed for every model, may it be a large model made by professionals or a small model made by a beginner. Which is why there is a need for a language that is intuitive, brief, and user friendly so that the user can preprocess the data as required before using the same. Therefore our DSL provides these and more by offering features such as intuitive and chainable syntax, strong semantic analysis, and pipeline functions, enabling users to create readable, and less verbose data transformation processes. Possible optimizations of the DSL are lazy evaluation, vectorization, memory-efficient data slicing, and the ability to easily integrate various data transformation functions. The language focuses on making the code intuitive and readable by adding simple English-like syntax. All of the above discussed aspects about the DSL is explained in detail in this whitepaper.

# Acknowledgements

We would like to express our gratitude as a team to our course instructor, **Dr. Ramakrishna Upadrasta**. For the guidance and support throughout this project and course. We also extend our thanks to our Teaching Assistants(TA's) **Mr. Rajapu Jayachandra Naidu** and **Mr. Rajiv Shailesh Chitale** for their insights on the DSL and always being available for any doubts. Finally, we acknowledge the course CS3423 Compilers-II at the Indian Institute of Technology, Hyderabad for providing the opportunity to do this project.

# Contents

<b>Abstract</b>	<b>ii</b>
<b>Acknowledgements</b>	<b>iii</b>
<b>1 Introduction</b>	<b>1</b>
1.1 What is DATUM:	1
1.2 Motivation	1
1.3 Ease of Use:	1
1.4 Completeness:	1
1.5 Intuitive:	2
<b>2 Aim of our DSL</b>	<b>3</b>
2.1 Making Intuitive Syntax	3
2.2 Write Concise Codes with Chainable Syntax	4
2.3 Pipeline Functions	4
2.4 Strong Semantic Analysis	5
<b>3 Program Structure</b>	<b>6</b>
3.1 Function Declarations	7
3.2 Start of Program	7
3.3 Syntax of Functions	7
3.3.1 Pipeline Functions	7
3.3.2 Parameter Functions	7
<b>4 Lexemes</b>	<b>9</b>
4.1 Data types:	9
4.2 Comments:	10

4.2.1	Single Line comments:	10
4.2.2	Multi Line comments:	10
4.3	Operators:	11
4.3.1	Arithmetic:	11
4.3.2	Relational Operators:	12
4.3.3	Logical Operators:	12
4.3.4	Additional Keywords:	12
4.3.5	Built-in Pipeline Functions:	13
4.3.6	Visualization:	15
4.4	Control Flow:	16
4.4.1	Conditional Flow:	16
4.4.2	Iterative flow:	16
5	Possible Optimizations	18
5.1	Lazy Evaluation	18
5.2	Vectorization	18
5.3	Memory-efficient Data Slicing	19
5.4	Pipeline Optimization and Reordering	19
6	Sample Codes	21
6.1	Sample Code 1:	21
6.2	Sample Code 2:	22
6.3	Sample Code 3:	23

# Chapter 1

## Introduction

### 1.1 What is DATUM:

DATUM is a simple and complete domain specific language which circles around data manipulation, visualization, pre-processing and data handling while focusing on ease of use and intuitiveness.

### 1.2 Motivation

As ML and DS continue to evolve, so does the demand for the languages and tools. But many existing languages and tools are often complex, less productive, bulky or lack some processing to do some specific data related tasks. Datum is developed to address some of these issues by providing productive, less complex language with possible optimizations and it also provides data visualization.

### 1.3 Ease of Use:

Programming in Datum is simple. The language provides a special data type ‘dataset’ and features that can be used on them.

### 1.4 Completeness:

The DSL is complete i.e. the user can program their own functions and manipulate the data accordingly.

## 1.5 Intuitive:

The language provides a new operator ‘ $\rightarrow$ ’ which eases the use of language as it suggests the flow of data through features enabling pipelining which helps in understanding the code better and also helps in debugging. It also offers common English keywords more on this is covered later.

## Chapter 2

# Aim of our DSL

There already exist several libraries for data manipulation and visualization. The motivation of creating a DSL when there already exists libraries is covered in the following points:

- Availability of all features at one place. The language offers some features that may not be present in one library, but present in other libraries.
- It is not function oriented rather it is language oriented where the user can come up with their own functions.
- Possible optimization offerings are Lazy Evaluation, Vectorization, Memory-efficient Data Slicing, Pipeline optimization and Reordering, which are explained later.
- The lack of strong semantic analysis is fulfilled by our DSL.
- The syntax is designed in such a way that improves the productivity of a programmer by increasing the readability and making it less verbose.

The way we achieve this is explained in the following sections.

### 2.1 Making Intuitive Syntax

The syntax and keywords in our DSL are designed to be intuitive and self-explanatory, making it easy for users to understand and follow. For instance,

```
1 row(from 1 to 10)
```



clearly expresses the intention to select rows 1 to 10, using plain English that is accessible even to those with limited programming experience. By prioritizing clear, human-readable commands, our DSL allows users to focus on their data rather than decoding complex syntax.

## 2.2 Write Concise Codes with Chainable Syntax

The chainable syntax of our DSL offers a highly intuitive and user-friendly approach to data manipulation, preprocessing, and visualization. Our DSL allows a user to chain multiple pipeline functions together, thus allowing the user to perform data manipulations / transformations in a pipelined fashion.

Example:

```
1 ds->row((from 1 to 5) also (from 10 to 15 step 3))
2   ->col(to "W" step 2)
3   ->filter({elem == 10;})
4   = 100;
```

Our DSL is less verbose and complex than the traditional data manipulation frameworks. This syntax minimizes nested function calls, making data operations feel more natural and readable. Unlike existing languages, which often rely on verbose function nesting, our DSL's clear and concise chaining mechanism makes it easy for users to combine operations without losing track of their data transformations, enhancing productivity and reducing error.

## 2.3 Pipeline Functions

Apart from the basic pipeline functionalities provided by our DSL, we also provide the freedom to the user to create their custom pipeline functions to use them in their custom pipeline to match their specific needs.

Example:

```
1 function (dataset ds) -> fill_null_values(int num) -> dataset {
2   ds->filter(ds == null) = num;
3   return ds;
4 }
5
6 ds->fill_null_values(0)->rows(from 1 to 10)->show();
```

This pipeline function takes a dataset, replaces all the null values in that dataset with the user provided number and returns the dataset.

A pipeline function is declared by using the function keyword. Each pipeline function expects the datatypes of the input of the pipe, datatypes of the output of the pipe, and the list of datatypes of the additional parameters.

## 2.4 Strong Semantic Analysis

Our DSL offers a very good semantic analysis ensuring that the logic of the code and the manipulation of the dataset being performed by the user is correct and also makes sense. We perform semantic checks on various aspects such as data type check while performing operations, pipeline order, behavior of functions, scope and data tracking at each step.

A few other data manipulation libraries already present out there mostly use Python as their primary language which does not offer an extremely impressive semantic analysis as it gives a lot of runtime errors. In our DSL we try to perform the semantic check and reduce the same.

### Some key semantic checks:

1. **Data Type check:** Ensures that all the data types in the DSL are operated only on the data types that they are allowed to operate on. For example, arithmetic operations are only allowed on numerical types such as integer, float and dataset with all data types of type integers and floats, while concatenation is limited to string. By adding strict data types rules we can limit errors such as concatenation of number with a string which happens in languages like python or JavaScript to name a few.
2. **Pipeline Ordering:** The pipeline feature of our DSL allows to chain data transformations in a logical sequence. The semantic analysis checks that each step in the pipeline operates on valid input and produces valid output for the next function in the chain.
3. **Behaviour of Functions:** This semantic check ensures that these functions behaves accordingly with their declared input and output types. This check prevents issues where functions may return incorrect types, leading to pipeline failures.
4. **Data Tracking:** As data moves through different stages of manipulation, the DSL tracks its transformations to ensure consistency. It ensures that the shape and structure remain valid for the upcoming operations.

## Chapter 3

# Program Structure

```
1  $$ Program that drops rows with null values in the first column, and fills all
   the null values in the second column with zeros. Then normalises each column
   of the dataset.
2  function_declarations:
3  function (dataset ds) -> fill_null_values(int num) -> dataset {
4      ds->filter(ds == null) = num;
5      return ds;
6  }
7  function (dataset ds) -> normalize() -> dataset {
8      array(integer) mean_arr = ds->mean(1);
9      loop i from 0 to ds->shape[0] {
10         ds->col(i) -= mean_arr[i];
11     }
12     array(integer) std_arr = ds*ds->sum(1) / ds->shape[0];
13     loop i from 0 to ds->shape[0] {
14         ds->col(i) /= std_arr[i];
15     }
16     return ds;
17 }
18 start:
19 dataset ds = read("./dataset", "csv")
20 ds->col(1)
21     ->fill_null_values(0)
22     ->row({
23         return r[0] == null;
24     })->drop();
25
26 ds->normalize()->write("./dataset.csv", "csv");
```

The above code is a sample code using which we will explain the program structure of the DSL.

## 3.1 Function Declarations

All the function declarations in the code are found under the `function_declarations:` label and above the `start:` . The `function_declarations:` label is optional. We can have programs without any user defined functions.

## 3.2 Start of Program

The program starts executing from the line below `start:` label. The `start:` label is mandatory.

## 3.3 Syntax of Functions

### 3.3.1 Pipeline Functions

A pipeline function is a function that can be part of a pipeline.

A pipeline function is declared by using the `function` keyword. Each pipeline function expects the datatypes of the input of the pipe, datatypes of the output of the pipe, and the list of datatypes of the additional parameters.

```
1 function (datatype1 param1, ...) -> function_name(datatype1 param1, ...) -> (
    datatype1, ...) {
2   $$ Statement 1
3   $$ ...
4   $$ Statement N
5 }
```

### 3.3.2 Parameter Functions

A parameter function is a function that can be passed as a parameter to a pipeline function. It is supported for the built-in pipeline functions like `row()`;

```
1 ... -> function_name({
2     $$ Statement 1
3     $$ ...
4     $$ Statement N
```

```
5  }) -> ...
```

The body of the function is surrounded by '{}'. We don't have to explicitly mention the parameters of the function. The undeclared variables in the body become the function parameters.

Example:

```
1  (ds1, ds2)->join({
2      integer sum1 = r1->sum();
3      integer sum2 = r2->sum();
4      return sum1 != sum2;
5  })
6      ->row(to 10)->show();
7
8  $$ the parameter function here is equivalent to
9
10 (r1, r2) {
11     integer sum1 = r1->sum();
12     integer sum2 = r2->sum();
13     return sum1 != sum2;
14 }
```

The order of the undeclared variables determine the order of the parameters of the function. The datatype of `r1` and `r2` are assumed to be arrays.

## Chapter 4

# Lexemes

DATUM follows some of the standard data types already present in C/C++, But it also offers some new cool data types and keywords which are domain specific and makes it easier for the programmer to use and understand.

### 4.1 Data types:

Datatype	Keyword	Description	Size
<b>Integer</b>	integer	It is the datatype to store signed numeric data values that fit within the 64-bit range.	8 Bytes (64 bits)
<b>Float</b>	float	The float data type is to store floating point values that fit within 64-bit representations.	8 Bytes (64 bits)
<b>Char</b>	char	The char type is used to store a single character using a single byte.	1 Byte (8 bits)
<b>String</b>	string	The string type is used to store a character string or a sequence of characters at contiguous memory locations.	Dynamic

Continued on next page

---

<b>Datatype</b>	<b>Keyword</b>	<b>Description</b>	<b>Size</b>
<b>Boolean</b>	bool	The boolean data type, which represents two possible values: true or false, used to indicate binary states or logical conditions.	1 Byte (8 bits)
<b>Dataset</b>	dataset	The dataset datatype is designed to store data read from CSV or TSV files, for further transformations and manipulations through pipelining and various built-in functions. It is a complex data structure used to store the entire dataset into a variable.	Dynamic
<b>Array</b>	array	The array data type is a dynamic container data structure designed to store a sequence of elements, where each element can be of only one single data type, including integer, float, char, string, boolean, and dataset. These elements are stored in contiguous memory locations, allowing for efficient access and manipulation.	Dynamic

## 4.2 Comments:

### 4.2.1 Single Line comments:

A single line comment starts with '\$\$'.

\$\$ This is a single line comment.

### 4.2.2 Multi Line comments:

Multi line comments start with '\$!' and terminate with '!\$'.

\$! This is a comment.

This comment illustrates multi line comments. !\$

## 4.3 Operators:

All basic numeric, relational and logical operators are supported. The following are the supported operators:

### 4.3.1 Arithmetic:

Operator	Description	Supported Data Types
<b>Increment</b> (++)	Increments the variable by the value 1.	integer, float
<b>Decrement</b> (--)	Decrements the variable by the value 1.	integer, float
<b>Division</b> (/)	Standard division	integer, float
<b>Multiply</b> (*)	Standard multiplication	integer, float
<b>Modulus</b> (%)	Modulus operator	integer
<b>Addition</b> (+)	Standard addition for integers and floats, offers concatenation for strings	integer, float, string (concatenation)
<b>Subtraction</b> (-)	Standard subtraction for integers and floats, also used to get the integer ASCII value difference between two characters or strings.	integer, float (subtraction), char, string (difference in ASCII value)
<b>Assignment</b> (=)	Standard assignment operator	All data types
<b>Multiply and Assign</b> (*=)	Multiply with a new value and assign the result to the same variable	integer, float
<b>Divide and Assign</b> (/=)	Divide with a new value and assign the result to the same variable	integer, float
<b>Modulus and Assign</b> (%=)	Perform modulo operation with a new value and assign the result to the same variable	integer
<b>Add and Assign</b> (+=)	Add with a new value and assign the result to the same variable	integer, float
<b>Subtract and Assign</b> (-=)	Subtract with a new value and assign the result to the same variable	integer, float



#### 4.3.2 Relational Operators:

Operator	Description	Supported Data Types
Greater than ( > )	Returns true when the left value is greater than the right value.	integer, float, char, string
Lesser than ( < )	Returns true when the left value is lesser than the right value.	integer, float, char, string
Greater than or equal to ( >= )	Returns true when the left value is greater than or equal to the right value.	integer, float, char, string
Lesser than or equal to ( <= )	Returns true when the left value is lesser than or equal to the right value.	integer, float, char, string
Equality ( == )	Returns true when the left value is equal to the right value.	integer, float, char, string
Inequality ( != )	Returns true when the left value is not equal to the right value.	integer, float, char, string

#### 4.3.3 Logical Operators:

Operator	Description	Supported Data Types
Logical and ( and )	Takes the logical AND of two statements. Returns true if both statements are true.	integer, float
Logical or ( or )	Takes the logical OR of two statements. Returns true if at least one statement is true.	integer, float

#### 4.3.4 Additional Keywords:

Keyword	Description
step	Used with from to syntax to skip the values in chunks of step value
also	Union of two arrays
from	The initial value

Continued on next page

Keyword	Description
<b>to</b>	The end value, end value is exclusive
<b>function</b>	Used to declare a function

#### 4.3.5 Built-in Pipeline Functions:

input_type → Keyword → output_type	Expected Parameters
<b>ds → row() → ds</b>	<ol style="list-style-type: none"> <li>1. Function returning boolean (parameter - array or array, int)</li> <li>2. From, to, step, also syntax</li> <li>3. Array of int (indices) or boolean (bitmask)</li> <li>4. Single row</li> </ol>
<b>ds → col() → ds</b>	<ol style="list-style-type: none"> <li>1. Function returning boolean (parameter - array or array, int)</li> <li>2. From, to, step, also syntax</li> <li>3. Array of int (indices), string (col names), boolean (bitmask)</li> <li>4. Single column</li> </ol>
<b>ds → filter() → ds</b>	Function returning boolean (parameter - int or int, int)
<b>ds → sum → array</b> <b>ds → sum → int</b> <b>array → sum → int</b>	<ol style="list-style-type: none"> <li>1. Parameter (int - axis when input is dataset)</li> <li>2. No parameter when input is array</li> </ol>
<b>ds → max → array</b> <b>ds → max → int</b> <b>array → max → int</b>	<ol style="list-style-type: none"> <li>1. Parameter (int - axis when input is dataset)</li> <li>2. No parameter when input is array</li> </ol>

Continued on next page

<b>input_type → Keyword → output_type</b>	<b>Expected Parameters</b>
<b>ds → min → array</b> <b>ds → min → int</b> <b>array → min → int</b>	1. Parameter (int - axis when input is dataset) 2. No parameter when input is array
<b>ds → mean → array</b> <b>ds → mean → int</b> <b>array → mean → int</b>	1. Parameter (int - axis when input is ds) 2. No parameter when input is array
<b>(ds, ds) → join → ds</b>	Function returning boolean (parameters are two arrays, rows)
<b>void → read → ds</b>	Read ds from file
<b>ds → write → void</b>	Write ds to file
<b>ds → unique → ds</b>	Unique rows, no parameters
<b>ds → split → array(ds)</b>	Parameter - array of proportions
<b>array → sort → array</b> <b>ds → sort → ds</b>	1. int or string (to sort according to the values of a column) 2. No parameters
<b>array → shuffle → array</b> <b>ds → shuffle → ds</b>	Shuffle the elements. No parameters
<b>ds → add(. . .) → ds</b>	The parameters in the add feature must be provided in a comma separated values pattern with the right data type as that of the of the respective column
<b>ds → shape → array</b>	Returns [rows, cols]
<b>ds → drop()</b>	1. Parameter - 0 or 1, 0 is for row and 1 is for column 2. Returns void

#### 4.3.6 Visualization:

Keyword/Function	Description	Usage
<code>ds → col(i) → show_bar()</code>  <code>ds → col(i) → show_bar(array)</code>	<ol style="list-style-type: none"> <li>1. Displays a barchart of frequencies plotted on the input column.</li> <li>2. The displayed values can be constrained to the ones present in the input array.</li> </ol>	Supported for discrete values with limited cardinality.
<code>(ds → col(i), ds → col(j)) → show_scatter()</code>  <code>ds → col(i) → show_scatter(array)</code>	<ol style="list-style-type: none"> <li>1. Creates a scatter plot to visualize the relationship between two columns.</li> <li>2. The range of the values can be constrained to [l,r].</li> </ol>	Supported for numerical data.
<code>(ds → col(i), ds → col(j)) → show_line()</code>  <code>ds → col(i) → show_line([l1, r1], [l2, r2])</code>	<ol style="list-style-type: none"> <li>1. Displays a line chart between two columns so visualize the relationship.</li> <li>2. The range of the values can be constrained to [l,r].</li> </ol>	Supported for numerical data.
<code>(ds → col(i), ds → col(j)) → show_box()</code>  <code>ds → col(i) → show_box()</code>	<ol style="list-style-type: none"> <li>1. Displays box plots to visualize the distribution of data across different categories.</li> <li>2. The range of the values can be constrained to [l,r].</li> </ol>	Supported between numerical and categorical variable type attributes.

## 4.4 Control Flow:

### 4.4.1 Conditional Flow:

#### **if:**

Executes the code block if the specified condition evaluates to true. Usage:

```
1 if ( condition ){  
2     Statement1;  
3     Statement2;  
4 }
```

#### **else if:**

Used after an if condition to specify a new condition if the previous if or another else if condition present above was false. It proceeds to execute the code block if this new condition is satisfied.

Usage:

```
1 if ( condition1 ){  
2     Statement1;  
3 }  
4 ...(cascaded else if conditions)...  
5 else if ( condition2 ){  
6     Statement2;  
7 }
```

#### **else:**

Provides a default code block that executes if none of the preceding if or else if conditions are true.

```
1 ...( if followed by (else if)* conditions )...  
2 else {  
3     Statement3;  
4 }
```

### 4.4.2 Iterative flow:

#### **loop:**

Iterates over a sequence or range, executing a specified block of code for each element in the sequence or iteration. Usage:

```

1 integer temp = 0;
2 loop (i from 1 to 10 step 3){
3     temp++;
4 }
5 $$ temp value is 3 after the loop

```

### break:

Exits from the nearest enclosing loop, immediately terminating the iteration and continuing execution with the code that follows the loop. Usage:

```

1 integer temp = 0;
2 loop (i from 5 to 10 step 4){
3     temp++;
4     if( condition ) break;
5 }

```

### continue:

Skips the rest of the current iteration of a loop and proceeds with the next iteration. Usage:

```

1 integer temp = 0;
2 loop (i from 1 to 100 step 19){
3     temp++;
4     if( condition ) continue;
5 }

```

### return:

Returns the data type as defined in the function. Usage:

```

1 function (dataset ds) -> fill_null_values(int num) -> dataset {
2     ds->filter(ds == null) = num;
3     return ds;
4 }

```

## Chapter 5

# Possible Optimizations

This chapter summarises the main outcomes and conclusions resulting from this body of work.

### 5.1 Lazy Evaluation

- Instead of executing operations immediately, we can collect all the commands in a pipeline and evaluate them only when the final result is needed. (Example: When we have to write the dataset back to the disk)
- In this way, we can improve the performance by minimizing unnecessary computations, especially when chained operations may override each other.
- In the Example given, before the second unique instruction is executed, the rows of `ds` are already unique. So the performance of this code can be optimized by the skipping redundant computation.

```
1 dataset ds = read("./dataset.csv", "csv");  
2 ds = ds->unique();  
3 ds = ds->row(to 10)->col(to 10)->unique();  
4 ds->write("./dataset.csv", "csv");
```

### 5.2 Vectorization

- When performing operations to a dataset, like adding two columns or adding a scalar to a dataset, we can apply operations to entire columns or rows at once instead of iterating element-by-element.

- In this way we can speed our program by leveraging the power of multi-processing / multi-threading or GPU programming.
- In the Example given, say, the user tries to convert height from meters to centimeters. Instead of iterating over every row, we can utilize the GPU to perform a single operation on multiple data.

```
1 dataset ds = read("./dataset.csv", "csv");
2 ds.col('height') *= 100 $$ can be parallelized
```

### 5.3 Memory-efficient Data Slicing

- When new datasets are declared to store a result of a data transformation, we can use views or references instead of copying data.
- Copying large dataset are both time and memory intensive. Therefore this optimization reduces memory overhead and improves slicing performance.
- In the above example, three other datasets are declared. If copies were created instead, the program would take way more time and space.

```
1 dataset ds = read("./dataset.csv", "csv");
2 dataset ds1 = ds->col(to 10)
3 dataset ds2 = ds->col(from 10 to 20)
4 dataset ds3 = ds->col(from 20 to 30)
```

### 5.4 Pipeline Optimization and Reordering

- We can analyze the pipeline and automatically reorder operations in a pipeline for an efficient performance.
- In this way, we can increase the performance by minimizes the amount of data processed in later stages of the pipeline.
- In the example given above, if filter is applied first to the dataset, then it has to be applied to all rows and columns of the dataset. But on reordering the operations in the pipeline, we can avoid the expensive computation.

```
1 dataset ds = read("./dataset.csv", "csv");
2
3 ds->filter({
```



```
4         return elem = null;
5     })
6     ->row(from 5 to 10) = 0;
7
8 $$ can be reordered to
9
10 ds->row(from 5 to 10)
11     ->filter({
12         return elem = null;
13     }) = 0;
```

## Chapter 6

# Sample Codes

### 6.1 Sample Code 1:

```
1  $$ Program that drops all the rows whose sum is not 1 and replaces the values in
   the quality column with 1 if it is greater than or equal to 7 and 0 if it
   is less than 7
2
3  function_declarations:
4
5  function (dataset ds) -> convert_values_to_binary(int num, string colName) ->
   dataset {
6      ds->col(colName)->filter({e < num}) = 0;
7      ds->col(colName)->filter({e >= num}) = 1;
8  }
9
10 start:
11
12 dataset ds = read("./dataset.csv", "csv")
13
14 ds->convert_values_to_binary(7, "quality")
15   ->row({
16       return r.sum() != 1;
17   })
18   ->drop(0);
19
20 ds->show_bar();
21 ds->write("./dataset-updated.csv", "csv");
```

## 6.2 Sample Code 2:

```
1  $$ Program that drops rows with null values in the first column, and fills all
   the null values in the second column with zeros. Then normalizes each column
   of the dataset.
2
3  function_declarations:
4
5  function (dataset ds) -> fill_null_values(int num) -> dataset {
6    ds->filter(ds == null) = num;
7    return ds;
8  }
9
10 function (dataset ds) -> normalize() -> dataset {
11   array(int) mean_arr = ds->mean(1);
12   loop i from 0 to ds->shape[0] {
13     ds->col(i) -= mean_arr[i];
14   }
15   array(int) std_arr = ds*ds->sum(1) / ds->shape[0];
16   loop i from 0 to ds->shape[0] {
17     ds->col(i) /= std_arr[i];
18   }
19   return ds;
20 }
21
22 start:
23
24 dataset ds = read("./dataset", "csv")
25 ds->col(1)
26   ->fill_null_values(0)
27   ->row({
28     return r[0] == null;
29   })->drop(0);
30
31 ds->normalize()->write("./dataset.csv", "csv");
```

## 6.3 Sample Code 3:

```
1  $$ First convert each row with their respective probabilities row-wise and split
   to train and test sets with ratios 0.8 and 0.2 respectively.
2
3  function_declarations:
4
5  function (dataset ds1) -> train_test_split(float testSize) -> (dataset, dataset)
6  {
7      ds->shuffle();
8      float trainSize = 1 - testSize;
9      array(ds) splitDs = ds->split([trainSize, testSize]);
10     return splitDs[0], splitDs[1];
11 }
12
13 function (dataset ds) -> convert_rows_to_orobabilities() -> (dataset)
14 {
15     ds->row({r = r/r->sum()});
16     return ds;
17 }
18
19 start:
20
21 dataset ds = read("./dataset.csv", "csv")
22
23 train, test = ds->convert_rows_to_orobabilities()->train_test_split(0.2);
24
25 X_train = train->col(to 10)
26 y_train = train->col(10)
27 X_test = test->col(to 10)
28 y_test = test->col(10)
29
30 X_train -> write("./x-train", "csv");
31 y_train -> write("./y-train", "csv");
32 X_test -> write("./x-test", "csv");
33 y_test -> write("./y-test", "csv");
```