

Paper SAS0000-2021

Azure Databricks and the SAS® Viya QKB for Better Data Quality and Entity Resolution

Arnold Toporowski, SAS Institute (Canada) Inc

ABSTRACT

Azure Databricks Python coders can now leverage the power of the SAS® Viya Quality Knowledge Base (QKB) and dramatically improve their data quality and data matching results. This session explores the capabilities now available to Python coders and gives coding examples and demonstrations showing how to leverage SAS Quality Knowledge Base capabilities such as parsing, standardization, and match-coding to better prepare data for analytics. Techniques for entity resolution and duplicate elimination are also explored.

INTRODUCTION

Data quality is a pervasive problem. There are some Python packages for data quality, but they are mostly about detecting or reporting on data quality, not for improving data quality. If you are using Python in an Azure Databricks notebook today to improve data quality, you are probably writing a lot of your own regular expressions and `np.where()` code.

For over twenty years SAS users have been able to use SAS® Data Quality and the SAS Quality Knowledge Base (QKB) to quickly and easily improve data quality, enrich data, and to facilitate duplicate elimination and entity resolution using fuzzy matching techniques.

Now, with SAS® Data Quality on SAS® Viya®, Python coders have access to the rules-based AI contained in the QKB and can leverage that same rich set of data quality capabilities.

Those capabilities include the following:

- Identification Analysis (to highlight data that is in the wrong place)
- Parsing (to take apart data into its constituent parts)
- Standardization (to correct inconsistent formatting)
- Gender Analysis (to get gender values from name information)
- Match Coding (to generate consistent codes for similar data values)
- And more! (Casing, Extraction, Pattern Analysis, Locale Guess, Language Guess)

The SAS QKB currently has support for 43 locales in 29 different languages, including languages using non-latin characters, such as Chinese, Japanese, Arabic, and Russian.

This paper will use a typical data quality problem and show you how the first five capabilities of the SAS QKB mentioned above can be leveraged from an Azure Databricks notebook to transform data exhibiting a poor level of data quality into a higher level of data quality that you would require for downstream analytics.

EXAMPLE PROBLEM

The CSV file shown below contains the data that we will be using for our example.

	A	B	C	D	E	F	G	H	I
1	ID	Name	Address	City	Prov	PostCode	Amount	Phone	Email
2	201	Mr. Jacques Plante	14 Denis Road	Cantley	QC	J8V 3J5	50	(819)-555-2334	
3	202	CGI	140 Grand-Allee Est Rue Unit 200	Quebec	QC	G1R 5P7	500		
4	203	CGI Ltée	140 Grand-Allee Est, Bureau 200	Ville de Québec	QC	G1R 5P7	400	info@cgi.ca	
5	204	Amar Singh	5264 Joel Avenue	Burlington, ON	L7L3Y7		300		
6	205	Jack Plant					50	555-2334	JPlante@gmail.com
7	206	Mr. Arnold Toporowski	38 Metropol, Unit 1605	Ottawa ON K1Z 1E9		613 755-2313	90		ArnoldT@sas.com
8	207	Ms MJ Belanger	4500 Sherbrooke St W	Montreal	QC	H3Z 1E6	950	514 799 9239	
9	208	CGI Canada Inc	200-140 Grand Allee	Quebec	Que	G1R 5P7	500	418.627.2227	
10	209	Plant, Jack	201-14 Denis Rd	Gatineau		J8V3J5	150	555-2334	
11	210	Mme Marie-Josée Bélanger	4500, rue Sherbrook Ouest	Montréal	PQ		900		MJBelanger@bell.ca
12	211	Ms. M.-J. Bélanger	4500 Sherbrooke O	Mon.	QC	H3Z 1E6	100	799-9239	
13	212	Jacques Plante	14, Chemin Denis, app 201	Cantley	Quebec		40		
14	213	Amar Singh	5264 Joel Av	Burlington	ON		100	(905) 637 5119	amar.singh@lost.com
15	214	Arnie Toporowski	38, privé Metropole, app 1605	Ottawa	Ont	K1Z1E9	400		
16	215	Jacques Plante	14 Denise Unit 201	Cantley, QC, J8V 3J5			10		
17	216	JF Tremblay	P.O. Box 123	St-Marc-du-Lac-Long	QC		200	819-555-4545	
18	217	Jean-Francois Tremblay	CP 123	Saint Marc QC	G0L 1T0		90		JFTremblay@bell.ca
19	218	Tremblay, JF	CP 123	St-Marc	Quebec	G0L 1T0	50	555-4545	
20	219	A. Toporowski	38 Metropole Private, Unit 1605	Ottawa	Ontario	K1Z 1E9	100	7552313	ArnoldT@sas.com
21	220	C.G.I.	140 Grand-Allee East, Unit 200	Quebec City	QC	G1R 5P7	600	627-2227	info@cgi.ca

Figure 1. CSV Data That Exhibits Poor Data Quality

This data contains multiple rows per individual/organization. Say for the type of analysis we want to do with this data we would like to have just six rows, one for each entity, with a total Amount per entity, and the other data cleansed and consolidated. Also, we would like to have a Gender field (Male/Female/Unknown) added for individuals.

Here are the data quality challenges that we are facing:

- no unique identifier per entity, and different types of entities (individuals/organizations)
- various formats of names of entities, including nick-names, initials, and legal terms
- various formats of phone and address information.
- missing data and incomplete data (e.g., some phone numbers missing area codes).
- data in the wrong places (e.g., Email in the Phone column).
- concatenated data that needs parsing (e.g., "Cantley QC J8V 3J5" all in the City field).
- typos ("Toporowski" vs "Toperowski", PostCode "G0L1T0" vs "GOL1T0", etc.).

The simplified business rules that we will use for this example are that we would like to get the longest Name, Address, City, Phone, and Email for each individual, along with total Amount, standardized Province, PostCode, Phone, and a generated gender code, like this:

	A	B	C	D	E	F	G	H	I
1	Name	gender	Address	City	Prov	PostCode	Amount	Phone	Email
2	Mr. Jacques Plante	M	14, Chemin Denis, app 201	Gatineau	QC	J8V 3J5	300	(819) 555 2334	JPlante@gmail.com
3	CGI Canada Inc	-	140 Grand-Allee Est Rue Unit 200	Ville de Québec	QC	G1R 5P7	2000	(418) 627 2227	info@cgi.ca
4	Amar Singh	M	5264 Joel Avenue	Burlington	ON	L7L 3Y7	400	(905) 637 5119	amar.singh@lost.com
5	Mr. Arnold Toporowski	M	38 Metropole Private, Unit 1605	Ottawa	ON	K1Z 1E9	590	(613) 755-2313	ArnoldT@sas.com
6	Mme Marie-Josée Bélanger	F	4500, rue Sherbrook Ouest	Montréal	QC	H3Z 1E6	1950	(514) 799 9239	MJBelanger@bell.ca
7	Jean-Francois Tremblay	M	P.O. Box 123	St-Marc-du-Lac-Long	QC	G0L 1T0	340	(819) 555 4545	JFTremblay@bell.ca

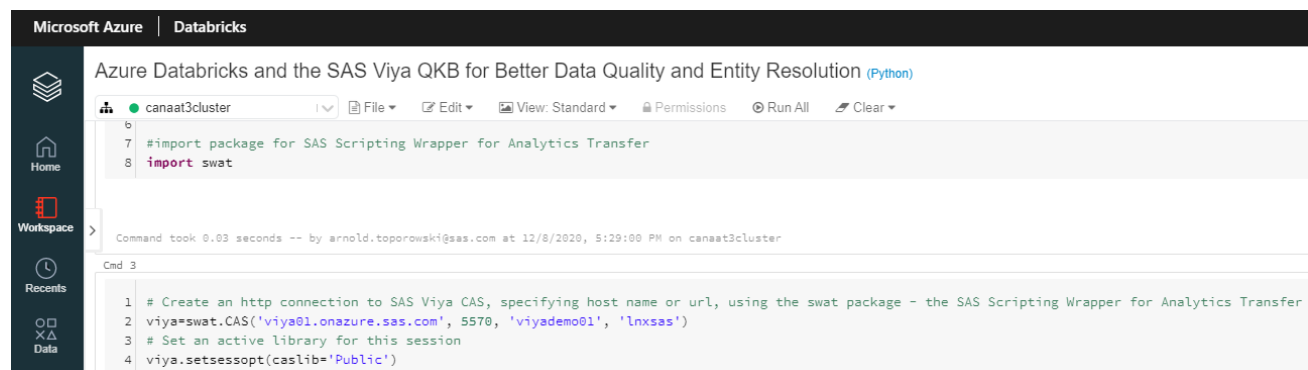
Figure 2. Results of the Data Quality Processing That We Would Like to Achieve

LOADING DATA INTO SAS® CLOUD ANALYTIC SERVICES

SAS Cloud Analytic Services (CAS) is the cloud-based server running on the SAS Viya high-performance, fault-tolerant analytics architecture. The smallest unit of work for the CAS server is a CAS action. CAS actions can load data, transform data, perform analytics, and create output.

To use CAS in a Databricks notebook you need to make sure that you have the SAS Scripting Wrapper for Analytics Transfer (SWAT) installed. See <https://github.com/sassoftware/python-swat> for details.

From a Databricks notebook we can then import the SAS SWAT package (**import swat**) and connect to SAS CAS (**swat.CAS**), and then set the active library with the **setessopt** action.



```
Microsoft Azure | Databricks

Azure Databricks and the SAS Viya QKB for Better Data Quality and Entity Resolution (Python)

6
7 #import package for SAS Scripting Wrapper for Analytics Transfer
8 import swat

Command took 0.03 seconds -- by arnold.toporowski@sas.com at 12/8/2020, 5:29:00 PM on canaat3cluster

Cmd 3

1 # Create an http connection to SAS Viya CAS, specifying host name or url, using the swat package - the SAS Scripting Wrapper for Analytics Transfer
2 viya=swat.CAS('viya01.onazure.sas.com', 5570, 'viyademo01', 'lnxsas')
3 # Set an active library for this session
4 viya.setessopt(caslib='Public')
```

Display 1. Setting Up the Environment and Connecting to CAS

Then we need to upload the data, in this case we used the **read_excel** method. To display the uploaded data, we can use **CASTable**, a swat DataFrame-like object:

```
1 # Load data to Viya either using a databricks client side load
2 testdata = viya.read_csv('/dbfs/FileStore/tables/CDN_ACCOUNTS/CDN_ACCOUNTS2.csv', casout=dict(name='CDN_Accounts2',caslib='Public', promote='true'))
3 # or load using a SAS Viya server side load
4 #
5 # use the swat CASTable object to treat a CAS Table like a pandas DataFrame
6 testdata = viya.CASTable('CDN_Accounts2')
7 testdata.head(20)
```

Display 2. Uploading Data and Then Fetching That Data Back for Display

NOTE: Cloud Analytic Services made the uploaded file available as table CDN_ACCOUNTS2 in caslib Public.
NOTE: The table CDN_ACCOUNTS2 has been created in caslib Public from binary data uploaded to Cloud Analytic Services.
Out[9]:

Selected Rows from Table CDN_ACCOUNTS2									
	ID	Name	Address	City	Prov	PostCode	Amount	Phone	Email
0	201.0	Mr. Jacques Plante	14 Denis Road	Cantley	QC	J8V 3J5	50.0	(819)-555-2334	
1	217.0	Jean-Francois Tremblay	CP 123	Saint Marc QC	G0L 1T0		90.0		JFTremblay@bell.ca
2	202.0	CGI	140 Rue Grand-Allee Est Unit 200	Quebec	QC	G1R 5P7	500.0		
3	218.0	Tremblay, JF	CP 123	St-Marc	Quebec	GOL 1T0	50.0	555-4545	
4	203.0	CGI Ltée	140 Grand-Allee Est, Bureau 200	Ville de Québec	QC	G1R 5P7	400.0	info@cgi.ca	
5	219.0	A. Toporowski	38 Metropole Private, Unit 1605	Ottawa	Ontario	K1Z 1E9	100.0	7552313	ArnoldT@sas.com
6	204.0	Amar Singh	5264 Joel Avenue	Burlington,ON	L7L3Y7		300.0		

Display 3. Output From the Display Data Action

DATA QUALITY AND DATA ENRICHMENT OPERATIONS

IDENTIFICATION ANALYSIS

The first data quality operation we will invoke is the **dqIdentify** function, through the **dataStep.runCode** CAS action. This inspects the City, Prov, PostCode, and Phone using the “Field Content” definition, in the English Canadian locale (ENCAN) of the QKB.

The results that come back show where these fields are empty, or where we are getting data different than we expected, using the newly created `_Ident` fields:

- The Phone_Ident field shows where Phone is empty or looks like an email.
- The Post_Ident field shows where PostCode is empty or looks like a phone number.
- The Prov_Ident field show where the Prov Field is empty or contains Postal Code.
- The City_Ident field shows where the City fields contains more than just City.
- The Name_Ident field shows whether the name is either an Individual or Organization

```
1 viya.dataStep.runCode(  
2     code=''' data public.test_dq_from_Python ;  
3         set public.CDN_Accounts2 ;  
4             City_Ident = dqIdentify(City,'Field Content','ENCAN');  
5             Prov_Ident = dqIdentify(Prov,'Field Content','ENCAN');  
6             Post_Ident = dqIdentify(PostCode,'Field Content','ENCAN');  
7             Phone_Ident = dqIdentify(Phone,'Field Content','ENCAN');  
8             Name_Ident = dqIdentify(Name,'Field Content','ENCAN');  
9         run;''')  
10 # let's look at just the first six rows of our data ...  
11 dq = viya.CASTable('test_dq_from_Python')  
12 dq.head(6)
```

Out[10]:

Selected Rows from Table TEST_DQ_FROM_PYTHON

Address	City	Prov	PostCode	Amount	Phone	Email	City_Ident	Prov_Ident	Post_Ident	Phone_Ident	Name_Ident
14 Denis Road	Cantley	QC	J8V 3J5	50.0	(819)-555-2334		CITY	STATE/PROVINCE	POSTAL CODE	PHONE	INDIVIDUAL
CP 123	Saint Marc QC	G0L 1T0		90.0	JFTremblay@bell.ca		INDIVIDUAL	POSTAL CODE	EMPTY	EMPTY	INDIVIDUAL
140 Rue Grand-Allee Est Unit 200	Quebec	QC	G1R 5P7	500.0			STATE/PROVINCE	STATE/PROVINCE	POSTAL CODE	EMPTY	UNKNOWN
CP 123	St-Marc	Quebec	GOL 1T0	50.0	555-4545		UNKNOWN	STATE/PROVINCE	POSTAL CODE	PHONE	UNKNOWN
140 Grand-Allee Est, Bureau 200	Ville de Québec	QC	G1R 5P7	400.0	info@cgl.ca		CITY	STATE/PROVINCE	POSTAL CODE	E-MAIL	ORGANIZATION
38 Metropole Private, Unit 1605	Ottawa	Ontario	K1Z 1E9	100.0	7552313	ArnoldT@sas.com	CITY	STATE/PROVINCE	POSTAL CODE	PHONE	INDIVIDUAL

Display 4. Commands and Output from the Identification Analysis

For a real-world problem, you might want to invoke the `dqIdentify` function on more fields, or even all fields, depending on the extent of the problems. In this case, we are only working on these five fields since we know from previous data profiling work on the CSV file that these data issues are limited to these five fields.

RIGHT-FIELDING

Next, we want to use the intelligence gained with the use of the **dqIdentify** function within SAS DATA step code that is invoked with the **dataStep.runCode** CAS action.

In the simple example code below, we find those cases where Email, Phone and Postal Code are in the wrong spot, and move them to the correct spot (and update the _Ident indicator fields as well).

Cmd 6

```

1 # using the above results of the SAS QKB Identity Analysis,
2 # run dataStep code to move stray phone, email, Postal code info to their correct places AKA "Right-Fielding"
3 viya.dataStep.runCode(
4     code=''' data public.test_dq_from_Python;
5         set public.test_dq_from_Python;
6         if Phone_Ident = 'E-MAIL' then do;
7             Email = Phone;
8             Phone = '';
9             Phone_Ident = 'EMPTY';
10        end;
11        if Post_Ident = 'PHONE' then do;
12            Phone = PostCode;
13            PostCode = '';
14            Phone_Ident = 'PHONE';
15            Post_Ident = 'EMPTY';
16        end;
17        if Prov_Ident = 'POSTAL CODE' then do;
18            PostCode = Prov;
19            Prov = '';
20            Post_Ident = 'POSTAL CODE';
21            Prov_Ident = 'EMPTY';
22        end;
23        run;'''
24    dq.head(6)

```

Out[11]:

Selected Rows from Table TEST_DQ_FROM_PYTHON

	ID	Name	Address	City	Prov	PostCode	Amount	Phone	Email	City_Ident	Prov_Ident	Post_Ident	Phone_Ident
0	201.0	Mr. Jacques Plante	14 Denis Road	Cantley	QC	J8V 3J5	50.0	(819)-555-2334		CITY	STATE/PROVINCE	POSTAL CODE	PHONE
1	217.0	Jean-Francois Tremblay	CP 123	Saint Marc QC		G0L 1T0	90.0	JFTremblay@bell.ca		INDIVIDUAL	EMPTY	POSTAL CODE	EMPTY
2	202.0	CGI	140 Rue Grand-Allee Est Unit 200	Quebec	QC	G1R 5P7	500.0			STATE/PROVINCE	STATE/PROVINCE	POSTAL CODE	EMPTY
3	218.0	Tremblay, JF	CP 123	St-Marc	Quebec	GOL 1T0	50.0	555-4545		UNKNOWN	STATE/PROVINCE	POSTAL CODE	PHONE
4	203.0	CGI Ltée	140 Grand-Allee Est, Bureau 200	Ville de Québec	QC	G1R 5P7	400.0		info@cgi.ca	CITY	STATE/PROVINCE	POSTAL CODE	EMPTY
5	219.0	A. Toporowski	38 Metropole Private, Unit 1605	Ottawa	Ontario	K1Z 1E9	100.0	7552313	ArnoldT@sas.com	CITY	STATE/PROVINCE	POSTAL CODE	PHONE

Command took 0.11 seconds -- by arnold.toporowski@sas.com at 12/8/2020, 5:32:58 PM on canaat3cluster

Display 5. The Code Used for Right-Fielding and the Results

PARSING

Next, we will use the **dqParse** and **dqParseTokenGet** functions to fix those rows where City, Province, and Postal Code have been concatenated together.

Display 6. The Code Used for Parsing Apart the City, Province, and Postal Code, and the Results

Now that we have all information in the correct place, we can move on to standardizations, corrections, and enrichment.

STANDARDIZATION AND ENRICHMENT

We want to standardize the Prov, PostCode, and Phone fields using the **dqStandardize** function. You want to use the correct standardization definition on each data type. (for example, the "State/Province (Postal Standard)" definition on the Prov field). We also specify the locale as "ENCAN", so that the data gets standardized to "English, Canadian" standards (note that the ENCAN and FRCAN definitions handle data in both French and English). If we had some USA data, we would want to use the "ENUSA" definition on those rows.

Here we also invoke the **dqGender** function to generate the gender field. It will get a value of M, F, or U, depending on what it finds in the Name field (unless Name_Ident indicates that it is an Organization). It considers name prefixes (Mr., Ms, Mme, etc.) and uses a lookup table of given names that skew toward a specific gender.

Cmd 8

```

1 # use the SAS QKB Standardize definitions for Province, PostCode, and Phone to standardize those columns in place
2 # use the SAS QKB Gender Analysis definition to enrich the data with a gender field, based on the Name
3 viya.dataStep.runCode(
4     code=''' data public.test_dq_2(drop=Name_Ident);
5               set public.test_dq_2;
6               Prov   = dqStandardize(Prov,'State/Province (Postal Standard)','ENCAN');
7               Phone  = dqStandardize(Phone,'Phone');
8               PostCode= dqStandardize(PostCode,'Postal Code','ENCAN');
9               gender  = dqGender(Name,'Name','ENCAN') ;
10              if Name_Ident = 'ORGANIZATION' then gender = '-';
11              run;'''
12 dq2.to_frame()

```

Out[13]:

Selected Rows from Table TEST_DQ_2

	ID	Name	Address	City	Prov	PostCode	Amount	Phone	Email	gender
0	201.0	Mr. Jacques Plante	14 Denis Road	Cantley	QC	J8V 3J5	50.0	(819) 555 2334		M
1	217.0	Jean-Francois Tremblay	CP 123	Saint Marc	QC	G0L 1T0	90.0		JFTremblay@bell.ca	M
2	202.0	CGI	140 Rue Grand-Allee Est Unit 200	Quebec	QC	G1R 5P7	500.0			U
3	218.0	Tremblay, JF	CP 123	St-Marc	QC	G0L 1T0	50.0	555 4545		U
4	203.0	CGI Ltée	140 Grand-Allee Est, Bureau 200	Ville de Québec	QC	G1R 5P7	400.0		info@cgi.ca	-
5	219.0	A. Toporowski	38 Metropole Private, Unit 1605	Ottawa	ON	K1Z 1E9	100.0	755 2313	ArnoldT@sas.com	U
6	204.0	Amar Singh	5264 Joel Avenue	Burlington	ON	L7L 3Y7	300.0			M
7	220.0	C.G.I.	140 Grand-Allee East, Unit 200	Quebec City	QC	G1R 5P7	600.0	627 2227	info@cgi.ca	-
8	205.0	Jack Plant					50.0	555 2334	JPlante@gmail.com	M
9	206.0	Mr. Arnold Toporowski	38 Metropol, Unit 1605	Ottawa	ON	K1Z 1E9	90.0	(613) 755 2313	ArnoldT@sas.com	M
10	207.0	Ms MJ Belanger	4500 Sherbrooke St W	Montreal	QC	H3Z 1E6	950.0	(514) 799 9239		F
11	208.0	CGI Canada Inc	200-140 Grand Allee	Quebec	QC	G1R 5P7	500.0	(418) 627 2227		-
12	209.0	Plant, Jack	201-14 Denis Rd	Gatineau		J8V 3J5	150.0	555 2334		M
13	210.0	Mme Marie-Josée Bélanger	4500, rue Sherbrook Ouest	Montréal	QC		900.0		MJBelanger@bell.ca	F
14	211.0	Ms. M.-J. Bélanger	4500 Sherbrooke O	Mon.	QC	H3Z 1E6	100.0	799 9239		F
15	212.0	Jacques Plante	14, Chemin Denis, app 201	Cantley	QC		40.0			M
16	213.0	Amar Singh	5264 Joel Av	Burlington	ON		100.0	(905) 637 5119	amar.singh@lost.com	M
17	214.0	Arnie Toporowski	38, privé Metropole, app 1605	Otawa	ON	K1Z 1E9	400.0			M
18	215.0	Jacques Plante	14 Denise Unit 201	Cantley	QC	J8V 3J5	10.0			M
19	216.0	JF Tremblay	P.O. Box 123	St-Marc-du-Lac-Long	QC		200.0	(819) 555 4545		U

Display 7. Standardization of Prov, Phone, and PostCode, and Generating Gender

ENTITY RESOLUTION OPERATIONS

MATCH-CODING

Next is the **dqMatch** function. Sensitivity is the third parameter, which determines how much “fuzziness” is allowed in the matchcode. The standard setting for sensitivity is 85, but can be set higher to require closer matches, or set lower to permit “fuzzier” matches. Both Canadian locales (ENCAN and FRCAN) will handle bilingual English and French data and generate the same match code for similar data no matter the language. If we had some USA data, we would probably want to use the “ENUSA” definition on those rows.

Cmd 9

```

1 # using the SAS QKB matchcode definitions create matchcodes on Name, Address, City.
2 viya.dataStep.runCode(
3   code='' data public.test_dq_3;
4     set public.test_dq_2;
5     Name_matchcode55 = dqMatch(Name,'Name',55,'ENCAN');
6     Org_matchcode85 = dqMatch(Name,'Organization',85,'ENCAN');
7     Address_matchcode70 = dqMatch(Address,'Address',70,'ENCAN');
8     City_matchcode85 = dqMatch(City,'City',85,'ENCAN');
9     if length(Phone) > 7 then Phone7 = substr(Phone,length(Phone)-7,8);
10    run;''')
11 # let's look at results for rows 3 to 14 ...
12 dq3 = viya.CASTable('test_dq_3')
13 dq3[['ID','Name','Address','City','Name_matchcode55','Org_matchcode85','Address_matchcode70','City_matchcode85','Phone7']].fetch(from_=3,to=14)

```

Out[22]:

§ Fetch

Selected Rows from Table TEST_DQ_3

ID	Name	Address	City	Name_matchcode55	Org_matchcode85	Address_matchcode70
2 202.0	CGI	140 Rue Grand-Allee Est Unit 200	Quebec	3F7\$	JF7\$	\$Z9FYP~\$\$\$\$\$\$\$\$...
3 218.0	Tremblay, JF	CP 123	St-Marc	~YBM\$	~YEBMWERC\$	\$N@MZH\$...
4 203.0	CGI Ltée	140 Grand-Allee Est, Bureau 200	Ville de Québec	W-7\$	JF7\$	\$Z9FYP~\$\$\$\$\$\$\$\$...
5 219.0	A. Toporowski	35 Metropole Private, Unit 1605	Ottawa	~MYL\$	~FNEYEL437\$	\$SKB~YNWNV~\$\$\$\$\$\$\$\$...
6 204.0	Amar Singh	5264 Joel Avenue	Burlington	4BF\$	8B&Y47PF\$	\$5H6CW\$\$\$\$\$\$\$\$\$\$...
7 220.0	C.G.I.	140 Grand-Allee East, Unit 200	Quebec City	7\$	JF7\$	\$Z9FYP~\$\$\$\$\$\$\$\$...
8 205.0	Jack Plant		NWB	~\$	C&J3NWEP~\$	\$Z8P\$\$\$\$\$\$\$\$\$\$...
9 206.0	Mr. Arnold Toporowski	38 Metropol, Unit 1605	Ottawa	~MYL\$	BY&YP7~FNEYEL437\$	\$SKB~YNW\$\$\$\$\$\$\$\$\$\$...
10 207.0	Ms MJ Belanger	4500 Sherbrooke St W	Montreal	MWBF\$	B4BCMEW&PFEY\$	\$S042YMY3\$\$\$\$\$\$\$\$...
11 208.0	CGI Canada Inc	200-140 Grand Allee	Quebec	8B3\$	JF7\$	\$Z9FYP~\$\$\$\$\$\$\$\$...
12 209.0	Plant, Jack	201-14 Denis Rd	Gatineau	NWB~\$	NWEP~C&J3\$	\$Z8P\$\$\$\$\$\$\$\$\$\$...
13 210.0	Mme Marie-Josée Belanger	4500, rue Sherbrook Ouest	Montréal	MWBF\$	BB_B&YRC@4_MEWEPFEY\$	\$S042YMY3\$\$\$\$\$\$\$\$...

Display 8. Similar Data Values Get the Exact Same Matchcode Using the dqMatch Function

CLUSTERING

Next we need to load the Entity Resolution CAS action set, and use the **entityres.match** CAS action to cluster together records based on the matching rules that we specify. In this example, we are using the matchcodes to match on the following rules:

- Name & Street Address & PostalCode - **OR** -
- Name & City & Province - **OR** -
- Name & Phone (last 7 digits) - **OR** -
- Name & Email - **OR** -
- Org Name & PostalCode

```
> 1 #load Entity Resolution CAS action set
2 viya.loadactionset(actionset="entityRes")
```

NOTE: Added action set 'entityRes'.
Out[15]:

```
§ actionset
entityRes
```

Cmd 11

```
1 # use the entityres.match CAS action to match rows on: (Name & Address & Postal Code) |
2 # OR (Name & City & Province) OR (Name & Phone) OR (Name & Email) OR (Org & Postcode)
3 #
4 dq_clustered = viya.CASTable('test_Clustered', groupBy='CLUSTERID', replace=True)
5 viya.entityres.match(clusterid='CLUSTERID',
6     intable=dq3,
7     matchrules=[{'rule':{'columns':['Name_matchcode55','Address_matchcode70','Postcode']},},
8                 {'rule':{'columns':['Name_matchcode55','City_matchcode85','Prov']},},
9                 {'rule':{'columns':['Name_matchcode55','Phone7']},},
10                {'rule':{'columns':['Name_matchcode55','Email']},},
11                {'rule':{'columns':['Org_matchcode85','Postcode']},},
12                ],
13     outtable=dq_clustered)
14
15 # display cluster results. CLUSTERID is a 24-byte character string
16 dq_clustered[['CLUSTERID','ID','Name','Address','City','Phone','Email']].sort_values('CLUSTERID').to_frame()
```

Out[16]:

Selected Rows from Table TEST_CLUSTERED

CLUSTERID	ID	Name	Address	City	Phone	Email
0	0.0 215.0	Jacques Plante	14 Denise Unit 201	Cantley		
1	0.0 212.0	Jacques Plante	14, Chemin Denis, app 201	Cantley		
2	0.0 205.0	Jack Plant			555 2334	JPlante@gmail.com
3	0.0 201.0	Mr. Jacques Plante	14 Denis Road	Cantley	(819) 555 2334	
4	0.0 209.0	Plant, Jack	201-14 Denis Rd	Gatineau	555 2334	
5	1.0 216.0	JF Tremblay	P.O. Box 123	St-Marc-du-Lac-Long	(819) 555 4545	
6	1.0 217.0	Jean-Francois Tremblay	CP 123	Saint Marc		JFTremblay@bell.ca
7	1.0 218.0	Tremblay, JF	CP 123	St-Marc	555 4545	
8	2.0 202.0	CGI	140 Rue Grand-Allee Est Unit 200	Quebec		
9	2.0 208.0	CGI Canada Inc	200-140 Grand Allee	Quebec	(418) 627 2227	
10	2.0 203.0	CGI Ltée	140 Grand-Allee Est, Bureau 200	Ville de Québec		info@cgi.ca
11	2.0 220.0	C.G.I.	140 Grand-Allee East, Unit 200	Quebec City	627 2227	info@cgi.ca

Display 9. Clustering Rules Using Matchcodes, and the Resulting CLUSTERID

SURVIVORSHIP

Now that we have all the records for our six people grouped together as we had hoped we would, all that's left to do is sum the amounts for each person, and choose the best information that we have for each person on each of the other columns. This is called Survivorship.

Below is some example SAS DATA step code to create a surviving record for each person. In this example, we are simply summing up the Amount for each cluster, and selecting the Name, Address, City, Postal, Phone, and Email with the longest string lengths in the cluster. Real world survivorship rules are usually more robust than this, so please don't consider this code to be a "best-practices" survivorship coding example.

After this step is done, you only need to download one record per individual to your Databricks client. Or leave the data on the server and leverage SAS Analytics in CAS from Python as well.

```
Cmd 12
>
1 # run datastep code to create a "golden record" for each cluster with best info, and total amount. AKA "Survivorship"
2 viya.dataStep.runCode(
3     code=''' data public.test_done
4         (drop= max_Name min_gender max_Address max_City max_Postal tot_Amount max_Phone max_Email
5              ID Name_matchcode55 Address_matchcode70 City_matchcode85 Phone7);
6         set public.test_Clustered;
7         by CLUSTERID;
8
9         retain max_Name min_gender max_Address max_City max_Postal tot_Amount max_Phone max_Email;
10        if first.CLUSTERID then do;
11            max_Name = Name;
12            min_gender = gender;
13            max_Address = Address;
14            max_City = City;
15            max_Postal = PostCode;
16            tot_Amount = Amount;
17            max_Phone = Phone;
18            max_Email = Email;
19        end;
20        else do;
21            if length(Name) < length(max_Name) then Name = max_Name; else max_Name = Name;
22            if gender > min_gender then gender = min_gender; else min_gender = gender;
23            if length(Address) < length(max_Address) then Address = max_Address; else max_Address = Address;
24            if length(City) < length(max_City) then City = max_City; else max_City = City;
25            if length(PostCode) < length(max_Postal) then PostCode = max_Postal; else max_Postal = PostCode;
26            tot_Amount = Amount + tot_Amount;
27            Amount = tot_Amount;
28            if length(Phone) < length(max_Phone) then Phone = max_Phone; else max_Phone = Phone;
29            if length(Email) < length(max_Email) then Email = max_Email; else max_Email = Email;
30        end;
31        if last.CLUSTERID then output;
32    run;''')
33 dq_done = viya.CASTable('test_done')
34 dq_done[['CLUSTERID','Name','gender','Address','City','Prov','PostCode','Amount','Phone','Email']].sort_values('CLUSTERID').head(20)
```

Out[17]:

	CLUSTERID	Name	gender	Address	City	Prov	PostCode	Amount	Phone	Email
0	0.0	Mr. Jacques Plante	M	14, Chemin Denis, app 201	Gatineau	QC	J8V 3J5	300.0	(819) 555 2334	JPlante@gmail.com
1	1.0	Jean-Francois Tremblay	M	P.O. Box 123	St-Marc-du-Lac-Long	QC	G0L 1T0	340.0	(819) 555 4545	JFTremblay@bell.ca
2	2.0	CGI Canada Inc	-	140 Rue Grand-Allee Est Unit 200	Ville de Québec	QC	G1R 5P7	2000.0	(418) 627 2227	info@cgi.ca
3	3.0	Mr. Arnold Toporowski	M	38 Metropole Private, Unit 1605	Ottawa	ON	K1Z 1E9	590.0	(613) 755 2313	ArnoldT@sas.com
4	4.0	Amar Singh	M	5264 Joel Avenue	Burlington	ON	L7L 3Y7	400.0	(905) 637 5119	amar.singh@lost.com
5	5.0	Mme Marie-Josée Bélanger	F	4500, rue Sherbrook Ouest	Montréal	QC	H3Z 1E6	1950.0	(514) 799 9239	MJBelanger@bell.ca

Display 10. Surviving Just One Record Per Individual/Organization

CONCLUSION

This paper has shown that the rules-based AI capabilities in the SAS QKB are a powerful way for you, as a Python coder in Databricks, to achieve better data quality more quickly and easily than you could by trying to write your own code to achieve the same results.

What's next? I recommend you try out these capabilities on a data quality problem of your own, or take the example shown here and take it further. What if you wanted to standardize the formatting of the Names and remove name prefixes, or add missing name prefixes? (hint: there are Name Standardization and Name Parsing definitions in the QKB).

Finally, you might be starting to think about how to move beyond just an interactive data science session in a Databricks notebook and thinking about how to operationalize this data quality process. Will you want to schedule this as a regular batch job? Or turn it into a web service, callable in real-time by other applications? Or embed it in a process ingesting streaming data? All these deployment options for data quality are available with SAS.

ACKNOWLEDGMENTS

Thanks go out to Kevin D. Smith for his guidance on coding style, and to colleagues Ron Yee, André Lafreniere, and Guy Bourassa for content suggestions for this paper. Also, thanks to my colleagues Uttam Kumar, Cesar Guerra, Jocelyn Gascon-Giroux, Artem Glazkov and Kevin D. Smith for helping me to get over the technical hurdles around Azure networking, Databricks SSL libraries, and SAS® Viya® SSL certificates.

RECOMMENDED READING

- *SAS® Blogs: "Using Python to work with SAS Viya and CAS", by Chris Hemedinger*
- *SAS® Viya®: The Python Perspective, by Kevin D. Smith and Xiangxiang Meng*
- *SAS® Data Quality: Getting Started*
- *SAS® Quality Knowledge Base for Contact Information: Online Help*
- *SAS® Data Quality: Language Reference* (the "Functions supported in CAS" section)
- *SAS® Viya®: System Programming Guide* (the Python syntax examples)
- *SAS® Data Quality: CAS Action Programming Guide*
- *SAS® Cloud Analytics Services: CASL Programmer's Guide*

CONTACT INFORMATION

Your comments and questions are valued and encouraged. Contact the author at:

Arnold Toporowski
SAS Institute (Canada) Inc
+1 (613) 755-2313
Arnold.Toporowski@sas.com

SAS and all other SAS Institute Inc. product or service names are registered trademarks or trademarks of SAS Institute Inc. in the USA and other countries. ® indicates USA registration.

Other brand and product names are trademarks of their respective companies.