

# Fuzzy Matching with SAS® PROC SQL



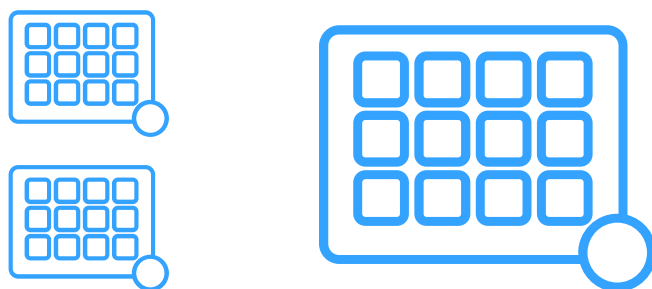
Grace Barnhill  
SAS HQ (Cary, NC)  
Email: [Grace.Barnhill@sas.com](mailto:Grace.Barnhill@sas.com)  
Linkedin: <http://linkedin.com/in/grace-barnhill/>



Copyright © SAS Institute Inc. All rights reserved.

In this demonstration, you'll learn about several techniques you can use for fuzzy matching data values that aren't strictly identical but represent the same entity.

## The Match Game



Copyright © SAS Institute Inc. All rights reserved.



We often need to assemble data from separate tables into a single table for further processing, analysis, and reporting.



UUID	Name
123e1212-e12b-12d3-a123-123123123001	Able
123e1212-e12b-12d3-a123-123123123002	Baker
123e1212-e12b-12d3-a123-123123123003	Charlie



UUID	Age
123e1212-e12b-12d3-a123-123123123001	20
123e1212-e12b-12d3-a123-123123123002	31
123e1212-e12b-12d3-a123-123123123003	42

Copyright © SAS Institute Inc. All rights reserved.



Ideally, the individual tables would contain a common unique identifier field, making for a simple and precise join. Of course, we don't live in a perfect world, so this scenario isn't as common as we'd like...

Name	Birthday	Date	Result
Jacqueline Smith	1994-06-30	2023-02-21	Good
John Smith	1994-06-30	2023-02-22	Bad
Joseph Jay	2000-01-01	2023-02-23	Avg

Name	Birthday	PatientID	Age
Jackie Smith	1994-06-30	0021345	25
John Smythe	1994-06-30	0021354	52
Joseph J.	2000-01-01	0021368	66



Copyright © SAS Institute Inc. All rights reserved.



More commonly, we will need to use the information in multiple columns to match records. If the data all comes from the same system, this often is also a simple join. But often, variations in how the data is entered can make this task trickier than it seems, requiring the use of fuzzy matching techniques to successfully join the records.

CustomerID	FirstName	LastName	StreetAddress	City	State	ZipCode
100000001	Katherin	Thomas	3005 Brookside Drive	Guin	AL	35563
100000002	Anthony	Feldmann	4026 Maple Lane	Huntsville	AL	35802
100000003	Sam	Tennyson	2715 Mulberry Avenue	Donaldson	AR	71941
100000004	Kristiana	Curry	399 Masonic Hill Road	Little Rock	AR	72212
100000016	Margaret	Robinson	622 Chestnut Street	Tampa	FL	33610

SubscriberID	Name	Address1	City	State	Zip
10010060	CHRISTY CURRY	399 MASONIC HILL RD	LITTLE ROCK	AR	72212-2122
10010127	JENNIFER RYAN	4974 WATER ST	FREMONT	CA	94539-9354
10010188	ANNA CHRISTIAN	3323 MORGAN ST	FT WALTON BEACH	FL	32548-8452
10010264	PEGGY ROBINSON	622 CHESTNUT ST	TAMPA	FL	33610-0163
10010305	MICHEAL TILLMAN	4854 POPLAR ST	CALUMET CITY	IL	60409-9040

Copyright © SAS Institute Inc. All rights reserved.



For example, say we are the owners of Cats Adore, a cat food company, and we'd like to identify potential new customers. To help with this endeavor, we purchased the Cat Lovers Magazine mailing list.

Now, while there is no ID column that provides a direct link between the Customers and CatLovers tables, we do have the individual's name and address information in both tables. This could be enough for matching. There are some obvious issues we'll have to deal with. For example:

In the Customer table, data is entered in proper case, while the mailing list data is all upper case. In the customer table, ZipCode is a numeric column containing 5-digit zip codes, but in the CatLovers table, Zip is a character column containing ZIP+4 codes. And finally, the names. Not every instance of a name is spelled the same, and some folks even use nicknames...

```

proc sql;
create table xref as
  select CustomerID, SubscriberID
  from catyums.customers as c
  inner join
    catyums.catlovers as cat
  on scan(cat.Name,-1,' ')=upcase(LastName)
     and scan(cat.Name,1,' ') =upcase(FirstName)
     and cat.State=c.State
     and cat.City =upcase(c.City)
     and put(Zipcode,z5.)=scan(Zip,1,'-')
;
quit;

```

NOTE: Table WORK.XREF created, with 2 rows and 2 columns

Copyright © SAS Institute Inc. All rights reserved.



Let's look at a sample join of this data. This join accounts for differences in text casing and zip code formats. The code runs without errors, but surely there must be more than 2 cat lovers in our customer database! This calls for a little fuzzy matching magic.

CustomerID	FirstName	LastName	StreetAddress	City	State	ZipCode
100000001	Katherin	Thomas	3005 Brookside Drive	Guin	AL	35563
100000002	Anthony	Feldmann	4026 Maple Lane	Huntsville	AL	35802
100000003	Sam	Tennyson	2715 Mulberry Avenue	Donaldson	AR	71941
100000004	Kristiana	Curry	399 Masonic Hill Road	Little Rock	AR	72212
100000016	Margaret	Robinson	622 Chestnut Street	Tampa	FL	33610

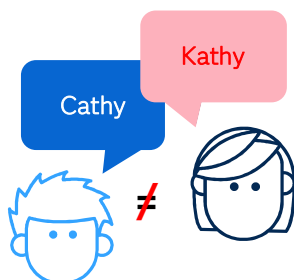
SubscriberID	Name	Address1	City	State	Zip
10010060	CHRISTY CURRY	399 MASONIC HILL RD	LITTLE ROCK	AR	72212-2122
10010127	JENNIFER RYAN	4974 WATER ST	FREMONT	CA	94539-9354
10010188	ANNA CHRISTIAN	3323 MORGAN ST	FT WALTON BEACH	FL	32548-8452
10010264	PEGGY ROBINSON	622 CHESTNUT ST	TAMPA	FL	33610-0163
10010305	MICHEAL TILLMAN	4854 POPLAR ST	CALUMET CITY	IL	60409-9040

Copyright © SAS Institute Inc. All rights reserved.



Fuzzy matching, sometimes referred to as fuzzy name matching or fuzzy string matching, uses techniques to identify text elements in a data set that are similar, but not necessarily identical. While both the SAS 9 Data Quality Server and SAS Viya provide the most powerful and configurable techniques for this type of work, there are several functions available in base SAS that can help tremendously.

## SOUNDEX()



```
proc sql;  
select FirstName, LastName  
       ,soundex(FirstName)  
from catyums.customers  
where CustomerID in  
       (100000004,000100000008);  
select Name  
       ,soundex(scan(Name,1,' '))  
from catyums.catlovers  
where SubscriberID in  
       (10010060,10010127);  
quit;
```

FirstName	LastName	Score
Kristiana	Curry	K6235
Jennifer	Ryan	J516

Name	Score
CHRISTY CURRY	C623
JENNIFER RYAN	J516

Copyright © SAS Institute Inc. All rights reserved.



Soundex is a phonetic algorithm that compares the sound of two text segments as they would be pronounced in English. It is often used for fuzzy matching based on names. Most SQL-based databases include a SOUNDEX function, and so does base SAS. SOUNDEX accepts a single argument consisting of a text string. In addition to the English bias, this algorithm is particularly sensitive to spelling diversions early in the text, so while it is very easy to use, it often misses matches that you or I would find exceptionally easy to detect. There is no way to control the scoring with SOUNDEX – so scores either match or don't.



```
proc sql;
create table xref as
  select CustomerID, SubscriberID
  from catyums.customers as c
  inner join
    catyums.catlovers as cat
  on SOUNDEX(scan(cat.Name,-1,' '))=SOUNDEX(LastName)
    and SOUNDEX(scan(cat.Name,1,' '))=SOUNDEX(FirstName)
    and cat.State=c.State
    and cat.City =upcase(c.City)
    and put(Zipcode,z5.)=scan(Zip,1,'-')
;
quit;
```

NOTE: Table WORK.XREF created, with 3 rows and 2 columns



Copyright © SAS Institute Inc. All rights reserved.

Using the SOUNDEX function in our join criteria identifies an additional cat lover among our customers. But I'm sure that more cat lovers are lurking in there!

# SPEDIS

**SPEDIS(compare word, base word)**

```
proc sql;  
select Name, FirstName, LastName  
      ,sum(SPEDIS(scan(name, 1, ' '),uppercase(FirstName))  
      ,SPEDIS(scan(name,-1, ' '),uppercase(LastName))) as Score  
from catyums.customers  
   ,catyums.catlovers  
where calculated score <75  
;  
quit;
```



Name	FirstName	LastName	Score
JENNIFER RYAN	Jennifer	Ryan	0
ANNA CHRISTIAN	Anna	Christian	0
MICHEAL TILLMAN	Micheal	Tillman	0
JIM CARPENTER	James	Carpenter	66
MICHEAL KAUFMAN	Micheal	Tillman	71
MICHEAL KAUFMAN	Michael	Kaufman	7
MICKEY KAUFMAN	Michael	Kaufman	41

Copyright © SAS Institute Inc. All rights reserved.



The asymmetric spelling distance algorithm is used to calculate the likelihood that two words match. SPEDIS is a base SAS function which finds the "distance" between two words by using a sequence of operations to make the strings match, each with a specific base "cost". In other words, how many operations do you have to do to make the comparison word match the base word, and what is the weight of those operations? Deleting or appending a letter is going to have a smaller cost than inserting or replacing a letter. The total cost is divided by the length of the string to produce the final score. The resulting score falls between 0 and 200, with 0 indicating a perfect match. So, the closer to zero, the better the match. You choose the cutoff point to maximize true matches while minimizing false positives. We'll discuss two more fuzzy matching functions, COMPGED and COMPLEV. Of the three, SPEDIS is generally the slowest processing.

## COMPGED

**COMPGED(string-1, string-2 <,cutoff><,modifiers>)**

```
proc sql;  
select Name, FirstName, LastName  
      ,sum(compged(scan(name, 1, ' '),FirstName,'n'))  
      ,compged(scan(name,-1, ' '),LastName,'n')) as Score  
from catyums.customers  
   ,catyums.catlovers  
where calculated score <250  
;  
quit;
```



Name	FirstName	LastName	Score
JENNIFER RYAN	Jennifer	Ryan	0
ANNA CHRISTIAN	Anna	Christian	0
BOB KELLEY	Robert	Kelley	230
MICHEAL TILLMAN	Micheal	Tillman	0
JIM CARPENTER	James	Carpenter	120

Copyright © SAS Institute Inc. All rights reserved.



The generalized edit distance function, or COMPGED, is a generalization of Levenshtein edit distance. Like with our last function, it's making edits to transform string 1 into string 2 and it's calculating the cost of each edit. The cost of each edit is additive, with a maximum cost of 200 and 0 indicating a perfect match. This time, the cost is not divided by the number of letters in the string. The COMPGED modifiers allow some control over processing and are applied in the order listed. Here the 'n' modifier removes quotation marks and ignores the case for both arguments. When choosing a cutoff point, the values returned by COMPGED are much larger than SPEDIS, making choosing a cutoff point to maximize true matches while minimizing false positives a bit easier. COMPGED processes faster than SPEDIS, but slower than COMPLEV.

## COMPLEV

**COMPLEV(string-1, string-2 <,cutoff><,modifiers>)**

```
proc sql;  
select Name, FirstName, LastName  
      ,sum(COMPLEV(scan(name, 1, ' '),FirstName,'n')  
          ,COMPLEV(scan(name,-1, ' '),LastName,'n')) as Score  
from catyums.customers  
   ,catyums.catlovers  
where calculated score <=3  
;  
quit;
```



Name	FirstName	LastName	Score
JENNIFER RYAN	Jennifer	Ryan	0
ANNA CHRISTIAN	Anna	Christian	0
MICHEAL TILLMAN	Micheal	Tillman	0
JIM CARPENTER	James	Carpenter	3
MICHEAL KAUFMAN	Michael	Kaufman	2
MICKEY KAUFMAN	Michael	Kaufman	3

Copyright © SAS Institute Inc. All rights reserved.



The Levenshtein edit distance function, or COMPLEV, produces a score based on the number and type of edits required to make two strings match. Once again, a score of 0 indicates a perfect match. The COMPLEV modifiers are the same as those for COMPGED, but the values returned are very small integers because it's only counting the number of edits required. This can make it tricky to choose a cutoff point to maximize true matches while minimizing false positives. But on the plus side, COMPLEV is faster than both COMPGED and SPEDIS.



## Fuzzy Matching with PROC SQL in base SAS

This demonstration illustrates using base SAS functionality with PROC SQL for fuzzy matching rows in tables without a common unique identifier.

Fuzzy Matching - 1 - Why I need it.sas  
Fuzzy Matching - 2 - SOUNDEX.sas  
Fuzzy Matching - 3 - SPEDIS.sas  
Fuzzy Matching - 4 - COMPLEV.sas  
Fuzzy Matching - 5 - COMPGED.sas



Copyright © SAS Institute Inc. All rights reserved.

Code and data available in the ZIP file.

#### Customers

Name	Street Address	City	State	Zip
Catherine M. Thomas	3005 Brookside Drive	Guin	AL	35563
János L. von Neumann	4026 Maple Lane	Huntsville	AL	35802
Sam C. Tennyson	2715 Mulberry Avenue	Donaldson	AR	71941
Kristi R. Curry	399 Masonic Hill Road	Little Rock	AR	72212
Wanda R. Smith	2190 Cedar Street	Pine Bluff	AR	71601

#### Cat Lover's Magazine Mailing List

Name	Address	City	State	Zip
BEAUTIFUL KITTY	3005 BROOKSIDE DR	GUIN	AL	35563
JOHN VON NEUMANN	4026 MAPLE LN	HUNTSVILLE	AL	35802
SAM CHARLES	3750 MULBERRY AVE	GROTON	CT	71941
CHRISTY CURRY	399 MASONIC HILL RD	LITTLE ROCK	AR	72212
TIGER MAN	2190 CEDAR ST	PINE BLUFF	AR	71601

Copyright © SAS Institute Inc. All rights reserved.



Sometimes the problems are a lot tougher. This calls for either a ton of coding, or the use of some power tools. Here's an example of some of these problems. Let's say that we're matching our customers with the list of subscribers for Cat Lover's Magazine again, but this time some subscribers have entered their cat's name as the recipient. The computer wouldn't match any of the rows in these tables, yet as a human, I can clearly see several matches:

- It looks like our customer Catherine Thomas has subscribed her cat, Beautiful Kitty, to Cat Lover's magazine.
- Our customer *Yanosh* von Neumann is listed as the cat lover *John* von Neumann.
- Customer Kristi R. Curry and the subscriber Christy Curry are most probably the same person.
- And it looks like Wanda Smith has also subscribed her cat, Tiger Man, to the magazine.

We may have to ignore names altogether when matching the table in this example.

#### Customers

Name	Street Address	City	State	Zip
Catherine M. Thomas	3005 Brookside Drive	Guin	AL	35563
János L. von Neumann	4026 Maple Lane	Huntsville	AL	35802
Sam C. Tennyson	2715 Mulberry Avenue	Donaldson	AR	71941
Kristi R. Curry	399 Masonic Hill Road	Little Rock	AR	72212
Wanda R. Smith	2190 Cedar Street	Pine Bluff	AR	71601

#### Cat Lover's Magazine Mailing List

Name	Address	City	State	Zip
BEAUTIFUL KITTY	3005 BROOKSIDE DR	GUIN	AL	35563
JOHN VON NEUMANN	4026 MAPLE LN	HUNTSVILLE	AL	35802
SAM CHARLES	3750 MULBERRY AVE	GROTON	CT	71941
CHRISTY CURRY	399 MASONIC HILL RD	LITTLE ROCK	AR	72212
TIGER MAN	2190 CEDAR ST	PINE BLUFF	AR	71601

```
Address=tranwrd(Address, 'ST', 'Street');
Address=tranwrd(Address, 'RD', 'Road');
Address=tranwrd(Address, 'AVE', 'Avenue');
...
```

Copyright © SAS Institute Inc. All rights reserved.



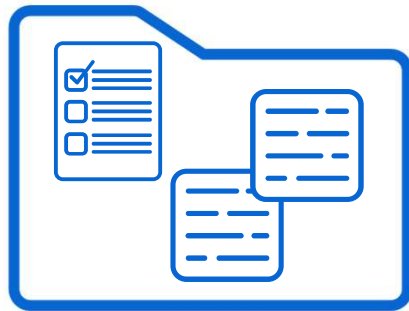
Instead, we could try to match households, using address, city, state, and zip code as the keys. But there are several problems with the text. The mailing list is all upper case and uses standard abbreviations like “ST” for “Street” and “AVE” for “Avenue”. The customer data information is proper cased and spells out all of the words. This could be a coding nightmare...



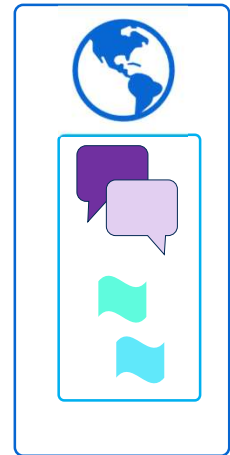
SAS Data Quality provides SAS language elements that perform data quality operations such as matching and standardization to improve data accuracy and consistency. Higher quality data yields higher-value results. SAS Data Quality comes with SAS Viya, but if you are a SAS 9 user, a SAS Data Quality license is required.



## SAS Quality Knowledge Base (QKB) Overview



SAS QKB for Contact  
Information (CI)



Copyright © SAS Institute Inc. All rights reserved.



At the heart of SAS Data Quality is the SAS Quality Knowledge Base or QKB. A QKB is a collection of complex rules and reference data used for data standardization and cleansing. Your SAS Data Quality deployment includes a QKB called the SAS Quality Knowledge Base for Contact Information, supporting an array of SAS processes used to simplify and accelerate profiling, standardizing, modifying, and matching data. The QKB contains

Information in a QKB is hierarchical in nature. The global, language, and locale levels each add more specific data and definitions in addition to those inherited from their parent level.

## QKB Locale Support



Copyright © SAS Institute Inc. All rights reserved.



QKB definitions are applied as specified by a locale. A locale is a unique combination of geography and language. For example, English is spoken in many geographies, but there are considerable differences between them. The same is true of Spanish. It's important to know which locale applies to your data, because the locale affects the results produced by SAS Data Quality functions and procedures.



John Smith

4B&~2\$\$\$\$\$\$\$\$\$C@P\$\$\$\$\$\$\$\$\$

	Name	Match Code
✓	John Smith	4B&~2\$\$\$\$\$\$\$\$\$C@P\$\$\$\$\$\$\$\$\$
✓	John Q. Smith	4B&~2\$\$\$\$\$\$\$\$\$C@P\$\$\$\$\$\$\$\$\$
✓	Jon Q. Smythe	4B&~2\$\$\$\$\$\$\$\$\$C@P\$\$\$\$\$\$\$\$\$
✗	John Q. Public	Q\$\$\$\$\$\$\$\$\$C@P\$\$\$\$\$\$\$\$\$

Copyright © SAS Institute Inc. All rights reserved.



Consider this table. If we try to match John Smith to these rows using traditional criteria to join on **Name**, only the first row of data matches. SAS Data Quality's "secret sauce" for fuzzy matching is the match code. Match codes are an encoded representation of a text string, computed based on a locale-sensitive combination of the data type, number of tokens, and sensitivity level. Using match codes instead of the raw text values, the computer can now see that these additional rows represent the same person, and still reliably decide that "John Q. Public" is not a match.



DQMATCH Function

- Creates match codes for any valid character input value
- Can be used in any SAS programming statement where functions are permitted



PROC DQMATCH

- Creates match codes for character variables from input data set
- Can create cluster IDs for matching records
- Options control output behavior

Copyright © SAS Institute Inc. All rights reserved.



SAS Data Quality provides both a DQMATCH function and a DQMATCH procedure. The DQMATCH function creates match codes for any valid text value. You can create a match code for a constant character value, the contents of a character variable, or the result of any valid character expression. You can use the DQMATCH function anywhere functions are permitted, which makes using match codes very flexible. I've used DQMATCH in the DATA step, PROC SQL, FEDSQL, and even DS2. Alternatively, you can use the DQMATCH procedure to create match codes for character variables in the input data set using multiple criteria, and cluster the data to help identify duplicate records. PROC DQMATCH includes several options for controlling how match codes and clusters are created, and for controlling output behavior.

## DQMATCH Function Syntax and Example

Syntax:

```
DQMATCH(character-value, 'match-definition' <, sensitivity> < , 'locale' >)
```

Example:

```
Matchcode=dqmatch(Contact, 'NAME');
```

Contact	Matchcode
Mark Abbott	&M&~\$BY3\$
Mark Abbott	&M&~\$BY3\$
Mark & Julie Abbott	&M&~\$BY3\$

Copyright © SAS Institute Inc. All rights reserved.



For fuzzy matching in PROC SQL, we'll focus on the DQMATCH function. The two required arguments are a *character-value*, which is the basis for the generated match code, and a *match-definition*, which specifies how the match code is produced. Two optional arguments provide finer control of match code generation. The default value for *Sensitivity* is 85, with valid values ranging from 50 to 95. The lower the sensitivity, the higher risk of false positives. The higher the sensitivity, the higher the risk of false negatives. In my experience, the default 85 works best for most applications. *Locale* specifies the locale to use when generating the match code. The default locale is used if no locale is specified. Both *match-definition* and *locale* values are text strings – hard-coded values must be written in quotes. In this example, the DQMATCH function generates match codes for the Contact variable using the NAME match definition. The output shows the match codes generated for each Contact value. In this example, the match codes indicate that all three rows of output refer to the same person.



## Fuzzy Matching with SAS Data Quality

This demonstration illustrates using SAS Data Quality functions with PROC SQL for fuzzy matching rows in tables without a common unique identifier.

Fuzzy Matching - 6 - Complex Problems DQ Solution.sas

Copyright © SAS Institute Inc. All rights reserved.



Code and data available in the ZIP file.

# Fuzzy Matching with SAS® PROC SQL



Grace Barnhill  
SAS HQ (Cary, NC)  
Email: [Grace.Barnhill@sas.com](mailto:Grace.Barnhill@sas.com)  
Linkedin: <http://linkedin.com/in/grace-barnhill/>



Get the ZIP file with the PDF and code!  
Scan the QR code or go to  
<https://bit.ly/SASFuzzyMatchSQL>



Copyright © SAS Institute Inc. All rights reserved.