$ipprl\_tools$  Documentation: Linkability Measures v1.0

July 2019

# Contents

1	Compatibility			
2	Required Dependencies			
3	Optional Dependencies			
4	4.1 PIP Method (Recommended)	2 2		
5	5.1 Importing the Package	2 3 5 6		
6	6.1			
	$ \begin{array}{cccccccccccccccccccccccccccccccccccc$	0		

# 1 Compatibility

The *ipprl\_tools* package was written using Python 3.6, but should be compatible with any version of Python 3 (Python 3.x).

# 2 Required Dependencies

The following packages are required dependencies for the *ipprl\_tools* package. If you installed *ipprl\_tools* through PIP, these dependencies should be installed automatically.

- Pandas  $\geq$  v0.23
  - https://pandas.pydata.org
- $NumPy \ge v1.16$

- https://www.numpy.org
- $SciPy \ge v1.2$ 
  - https://www.scipy.org

# 3 Optional Dependencies

The following packages are optional dependencies for the *ipprl\_tools* package. These dependencies will not be installed automatically when installing *ipprl\_tools* with PIP, so they must be installed manually if needed.

- Fuzzy  $\geq$  v1.2.2
  - This package is required for the Soundex corruption method. For more information about the package, visit https://pypi.org/project/Fuzzy/.
- Jupyter  $\geq$  v1.0.0
  - This package is required to view and run the tutorial Jupyter notebook. For more information about Jupyter, visit https://jupyter.org/

# 4 Installation

# 4.1 PIP Method (Recommended)

To install the package via PIP run the command:

```
pip install git+git://github.com/cu-recordlinkage/ipprl_tools
```

through a command-line interface.

This command will install the *ipprl\_tools* package into your default Python environment. This command will also install the required dependencies (Pandas, NumPy, SciPy, etc) if they are not already installed.

# 4.2 GitHub Method

: The source code can also be cloned directly from GitHub using the following command from a command-line interface.

```
git clone https://github.com/cu-recordlinkage/ipprl_tools
```

# 5 Usage

#### 5.1 Importing the Package

To use *ipprl\_tools*, first import the metrics submodule.

```
In [ ]: 1 from ipprl_tools import metrics
```

Figure 1: Importing ipprl\_tools

This command will import all of the functions defined in the metrics submodule. If you only need a subset of the functions, you can specify the exact functions to import as well:

```
In [ ]: 1 from ipprl_tools.metrics import convert_data,run_metrics
```

Figure 2: Importing specific functions

# 5.2 Data Prerequisites

The linkability metric functions expect that data will be contained in a Pandas DataFrame. To read in a file using Pandas, first call the appropriate read function. In our case we are reading CSV data, so we use  $pandas.read\_csv()$ , but alternative functions are available for other types of data. For additional ways to import data using Pandas, refer to the Pandas Documentation here:

Pandas Documentation: IO

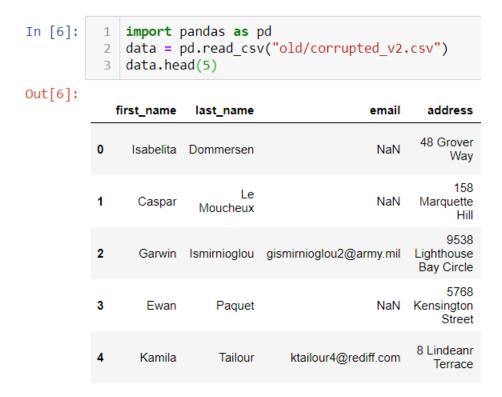


Figure 3: Reading in a CSV file with Pandas

The linkability metric functions provided in *ipprl\_tools* are designed to operate on Pandas DataFrames, but in order for them to work correctly, we must first ensure that the DataFrame

is in the correct format.

By default, Pandas will attempt to parse columns of your input file differently, depending on the type of data in the column.

In [6]:	<pre>data = pd.read_csv("corrupted_3_parts.csv") data.dtypes</pre>		
Out[6]:	id	int64	
	net_id	int64	
	first_name	object	
	last_name	object	
	email	object	
	address1	object	
	ssn	object	
	gender	object	
	city	object	
	zip	float64	
	state	object	
	dob	object	
	phone	object	
	phone2	object	
	phone3	object	
	race .	object	
	pcp_npi	object	
	suffix	object	
	title	object	
	middle_name	float64	
	mothers_maiden	float64	
	address2	float64	
	phone4	float64	
	mrn	float64	
	dtype: object		

Figure 4: Data types of each column, after reading data from a CSV file.

In this case, Pandas has parsed some of the numerical columns (*id*, *zip*, etc.) as different data types. In order for the linkability metrics to work, we first need to convert all columns in the DataFrame to the *string* data type. We also need to ensure that missing data is handled in the correct way. All linkability metrics treat the empty string ("") as a missing value.

In order to make this process easy,  $ipprl\_tools$  contains a function called  $convert\_data()$  (6.2.1), which can be used to automatically convert a DataFrame to the correct format. To use, import  $convert\_data$  from the  $ipprl\_tools.metrics$  submodule, then call it on your DataFrame.

```
In [8]:
          data = convert_data(data)
             data.dtypes
Out[8]: id
                           object
         net id
                           object
                           object
         first_name
         last name
                           object
         email
                           object
         address1
                           object
         ssn
                           object
         gender
                           object
                            object
         city
         zip
                           object
         state
                           object
         dob
                           object
         phone
                           object
                           object
         phone2
                           object
        phone3
         race
                           object
                           object
        pcp npi
         suffix
                           object
         title
                            object
         middle_name
                           object
                           object
         mothers_maiden
         address2
                           object
         phone4
                           object
                            object
         mrn
        dtype: object
```

Figure 5: Example of calling *convert\_data()* and viewing the result.

In the above figure, we call  $metrics.convert\_data()$  on our DataFrame, and view the data types of the output. We can see that the DataFrame was converted correctly, as every column is now of type(object).

# 5.3 Computing Metrics

# 5.3.1 Calling Individual Functions

To call a linkability metric on a DataFrame, first ensure that you've imported the function by following the instructions in Section 5.1.

We can call an individual function on every column in the DataFrame, or specify some subset of the columns to use.

```
In [10]: 1 | metrics.theoretical_maximum_entropy(data,columns=["first_name","last_name"])
Out[10]: {'first_name': 14.306488964728633, 'last_name': 15.712123926676279}
```

Figure 6: Example of calling a linkability metric by manually specifying the columns to operate on.

```
metrics.theoretical maximum entropy(data)
In [11]:
Out[11]: {'id': 17.042599881712917,
           'net_id': 16.609640474436812,
           'first_name': 14.306488964728633,
           'last_name': 15.712123926676279,
           'email': 14.965919530982598,
           'address1': 16.66379434619952,
           'ssn': 16.38407544713785,
            gender': 1.5849625007211563,
           city': 13.272775563286087,
           'zip': 13.45455634396194,
           'state': 10.521600439723727,
           'dob': 12.467350814487581,
           phone': 14.253847484987404,
            phone2': 14.261507309202056,
            phone3': 16.49039651667306,
           'race': 11.326991174900817,
           'pcp_npi': 10.396604781181859,
           'suffix': 2.584962500721156,
           'title': 11.88683970588442,
           'middle_name': -0.0,
           'mothers maiden': -0.0,
           'address2': -0.0,
           'phone4': -0.0,
           'mrn': -0.0}
```

Figure 7: Example of calling a linkability metric by using the default column argument (all columns)

In either case, the output of each linkability metric is a Python *dict* object, which has keys corresponding to the names of the selected columns, and values corresponding to the calculated linkability metric.

# 5.3.2 Using $run\_metrics()$

Although each metric can be run individually, *ipprl\_tools* also provides a helper function to run all of the metrics at once, and provide a formatted output DataFrame.

To use this function, import  $run\_metrics$  from the  $ipprl\_tools.metrics$  submodule, and call it on your DataFrame.

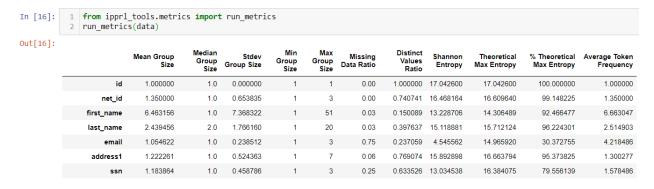


Figure 8: Example of importing and using the run\_metrics function.

The output of the function is a formatted Pandas DataFrame, where each column corresponds to

a linkability metrics, and the rows are the columns from your original DataFrame.

# 6 Function Documentation

# 6.1 Linkability Metrics

## $6.1.1 \quad missing\_data\_ratio(data, \ columns=None)$

This function calculates the Missing Data Ratio, or MDR for each of the specified columns. The Missing Data Ratio is defined as:

$$MDR_{i} = \frac{Number\ of\ records\ with\ missing\ value\ in\ Variable_{i}}{Total\ number\ of\ records\ in\ Variable_{i}} \tag{1}$$

This metric can be used to determine what fraction of the columns values are usable for linkage.

#### **Parameters:**

- data The Pandas DataFrame to be modified. Please ensure that your data is in the correct format (as specified in Section 5.2) before passing as an argument.
- columns An optional list of columns to operate on. These columns should correspond to column names in the DataFrame that you would like to calculate metrics for. The default argument (columns=None) will calculate metrics for each column in data.

#### Return Value:

 values - A dictionary containing a key/value mapping of column\_name → mdr\_value for each column specified in columns.

# 6.1.2 $distinct\_values\_ratio(data, columns=None)$

This function calculates the Distinct Values Ratio (DVR), which is defined as:

$$DVR_i = \frac{Number\ of\ distinct\ values\ for\ Variable_i}{Number\ of\ records\ in\ Variable_i} \tag{2}$$

This metric can be used to determine how many unique values a variable/column has, relative to its size. A high DVR indicates that there are a large number of distinct values in the column, whereas a low DVR indicates that there are fewer distinct values.

#### **Parameters:**

- data The Pandas DataFrame to be modified. Please ensure that your data is in the correct format (as specified in Section 5.2) before passing as an argument.
- columns An optional list of columns to operate on. These columns should correspond to column names in the DataFrame that you would like to calculate metrics for. The default argument (columns=None) will calculate metrics for each column in data.

#### Return Value:

• values - A dictionary containing a key/value mapping of column\_name → dvr\_value for each column specified in columns.

### 6.1.3 group\_size(data, columns=None)

This function calculates and returns the Group Sizes for each distinct value in each *column*. Group Size is defined as:

$$GS_{i,j} = Number\ of\ rows\ for\ Variable_i, Value_j$$
 (3)

This function is used to calculate  $agg\_group\_size$  (6.1.4), but may also be called by the user.

#### **Parameters:**

- data The Pandas DataFrame to be modified. Please ensure that your data is in the correct format (as specified in Section 5.2) before passing as an argument.
- columns An optional list of columns to operate on. These columns should correspond to column names in the DataFrame that you would like to calculate metrics for. The default argument (columns=None) will calculate metrics for each column in data.

#### Return Value:

• values~(dict) - For each column name in columns, returns a Counter object holding key/value mappings of  $distinct\_value \rightarrow num\_of\_occurrences$ 

# $6.1.4 \quad agg\_group\_size(data, agg\_func = np.mean, columns=None)$

This function calculates the aggregate group size for each column specified in *columns*, based on an arbitrary function  $agg\_func$ . The function calls  $group\_size()$  (6.1.3), and uses  $agg\_func$  to reduce the results to a scalar value for each column.

#### Parameters:

- data The Pandas DataFrame to be modified. Please ensure that your data is in the correct format (as specified in Section 5.2) before passing as an argument.
- agg\_func This is an arbitrary reduction function which accepts a list of values, and returns a scalar. Some examples of compatible reduction/aggregation functions are: numpy.amin, numpy.amax, numpy.mean, and numpy.std. The default argument (agg\_func = np.mean) will pass NumPy's mean function to compute the mean group size for each column name specified in columns.
- columns An optional list of columns to operate on. These columns should correspond to column names in the DataFrame that you would like to calculate metrics for. The default argument (columns=None) will calculate metrics for each column in data.

#### Return Value:

• values - A dictionary containing a key/value mapping of  $column\_name \rightarrow agg\_group\_size\_val$  for each column specified in columns.

# $6.1.5 \quad shannon\_entropy(data, columns=None)$

This function calculates the Shannon Entropy for each column name specified in *columns*. Shannon Entropy is defined as:

$$H_x = -\sum_{i=1}^{N} P(i) \log_2 P(i)$$
 (4)

This metric is useful for determining the amount of information containing in a column. Typically, columns with large numbers of evenly distributed unique values will have higher values for Shannon Entropy.

#### **Parameters:**

- data The Pandas DataFrame to be modified. Please ensure that your data is in the correct format (as specified in Section 5.2) before passing as an argument.
- columns An optional list of columns to operate on. These columns should correspond to column names in the DataFrame that you would like to calculate metrics for. The default argument (columns=None) will calculate metrics for each column in data.

#### Return Value:

• values - A dictionary containing a key/value mapping of  $column\_name \rightarrow shannon\_entropy$  for each column specified in columns.

### 6.1.6 theoretical\_maximum\_entropy(data, columns=None)

This function calculates the Theoretical Maximum Entropy (TME) for a column, given its current number of distinct values. The formal definition for Theoretical Maximum Entropy is:

$$TME_x = -log_2(\frac{1}{N}) \tag{5}$$

where N is the number of distinct values in the column.

### Parameters:

- data The Pandas DataFrame to be modified. Please ensure that your data is in the correct format (as specified in Section 5.2) before passing as an argument.
- columns An optional list of columns to operate on. These columns should correspond to column names in the DataFrame that you would like to calculate metrics for. The default argument (columns=None) will calculate metrics for each column in data.

#### Return Value:

• values - A dictionary containing a key/value mapping of  $column\_name \rightarrow theo\_max\_entropy$  for each column specified in columns.

#### 6.1.7 percent\_theoretical\_maximum\_entropy(data, columns=None)

This function will calculate the Percent of Theoretical Maximum Entropy (PTME) for each column name specified in *columns*. The Percent of Theoretical Maximum Entropy is defined as:

$$PTME_x = \frac{H_x}{TME_x} * 100 (6)$$

Where  $H_x$  and  $TME_x$  are the Shannon Entropy and Theoretical Maximum Entropy of the column, respectively. This function will call theoretical\_maximum\_entropy() (6.1.5) and percent\_theoretical\_maximum\_entropy (6.1.6) in order to calculate the result.

#### Parameters:

- data The Pandas DataFrame to be modified. Please ensure that your data is in the correct format (as specified in Section 5.2) before passing as an argument.
- columns An optional list of columns to operate on. These columns should correspond to column names in the DataFrame that you would like to calculate metrics for. The default argument (columns=None) will calculate metrics for each column in data.

### Return Value:

values - A dictionary containing a key/value mapping of column\_name → pct\_theo\_max\_entropy
for each column specified in columns.

# $6.1.8 \quad average\_token\_frequency(data,\ columns=None)$

This function will calculate the Average Token Frequency (ATF) for each column name specified in *columns*. The Average Token Frequency is defined as:

$$ATF = \frac{|V|}{N} \tag{7}$$

Where |V| is the size of the column, and N is the number of unique values in the column.

#### Parameters:

- data The Pandas DataFrame to be modified. Please ensure that your data is in the correct format (as specified in Section 5.2) before passing as an argument.
- columns An optional list of columns to operate on. These columns should correspond to column names in the DataFrame that you would like to calculate metrics for. The default argument (columns=None) will calculate metrics for each column in data.

#### Return Value:

• values - A dictionary containing a key/value mapping of  $column\_name \rightarrow avg\_tok\_freq$  for each column specified in columns.

### 6.2 Utilities

The utilities functions are a set of functions provided alongside the Linkability Metric functions in order to perform common tasks associated with calculating linkability metrics on a dataset.

# $6.2.1 \quad convert\_data(data)$

This function will convert the Pandas DataFrame data into a format suitable for calculating linkability metrics. The function will convert data to a DataFrame of string objects, automatically converting NaN values to the empty string ("").

### Parameters:

• data - The Pandas DataFrame to be converted. If your DataFrame contains columns that are not of type(str), or a value for NaN that is not the empty string (""), you will need to call this function to convert your DataFrame before running any of the linkage metrics.

#### Return Value:

• data - A version of the input DataFrame that has been converted to type(str), and had all NaN values replaced with the empty string.

# $6.2.2 \quad run\_metrics(data)$

This function is a helper function to automatically calculate all linkability metrics on all columns in data. The output will be collected and formatted into a Pandas DataFrame.

#### Parameters:

• data - The Pandas DataFrame to be modified. Please ensure that your data is in the correct format (as specified in Section 5.2) before passing as an argument.

### Return Value:

• metrics\_df - A Pandas DataFrame containing the values of every linkability metric calculated on every column in data. The columns of the DataFrame correspond to the linkability metrics, and the rows correspond to the names of the columns in data.columns.