
Data Refinery Documentation

Release 0.2.9

BBVA-Labs

Feb 08, 2018

Contents

1	Introduction	1
1.1	Motivation	1
1.2	Solutions	1
1.3	Help not substitution	2
1.4	Focus on streaming	2
2	Installation	3
2.1	Requirements	3
2.2	Pip	3
3	The Basics	5
3.1	Field operations	5
3.1.1	How process errors in data	6
3.1.2	List of functions	6
3.1.3	Composing field operations	10
3.2	Event operations	11
3.2.1	List of functions	11
3.2.2	Combine event operations	13
3.3	Pandas use	15
3.4	Review exercises	16
4	Tuple DSL documentation	17
4.1	DSL	17
4.1.1	List of functions	17
5	Python modules	19
5.1	datarefinery	19
5.2	datarefinery.FieldOperations	19
5.3	datarefinery.tuple.Formats	19
5.4	datarefinery.tuple.TupleDSL	19
5.5	datarefinery.TupleOperations	19
5.6	datarefinery.DateFieldOperations	19
5.7	datarefinery.CombineOperations	19
6	Glossary	21
7	Indices and tables	23

1.1 Motivation

Data Analytics is now very popular on modern companies. Nowadays nobody deny the data insights value regarding business development. Also technology allow us to manage huge amount of data, even when look imposible to process. In fact, scientific libraries allow us to perform advaced analytics.

> But i can't stop thinking about move local work to the cluster, and for me looks like a lonely valley.

This work is a must every time that a hypothesis looks promising and must hit the real world. Doing this job it's like walk along a lonely valley plenty of obstacles. The scientific libraries usually are not desingned to scale beyond machine memory (usually this is a very hard task) because modern machines are big enought. But a lot of times I spent time translating from R o Python to Java or Scala; just because clusters software is typically programed on JVM languages.

I think that smaller the transformation code, the better is translation to the cluster software. And also Python nowadays live between the two worlds.

1.2 Solutions

I usually talk about this translation problem with my colleagues on Big Data events. Some of then put cluster tools in hands of their data scientist, like Spark. But this comes with a cost, every test boot up a complete cluster on the local machine, slow down the development process. But not every test must hit the real work, so comes with an unnecessary adaptation to the cluster, instead a quick local test. Also there is an slightly outdated AI algorithms on cluster libraries. Some times people try to do the work directly on real cluster, with real data. Thats awesome but await to all the data be processed, usually slow down development process.

There is no silver bullet for this problem, every solution that you try has its own tradeoffs. Like "i must develop on cluster, but every test run for days" or like "i have my program finished but i must rewrite it in java"

On my opinion the only way to go is more automation and be iterative, so the scientific can use their own tools and the engineer can work less. When we say iterative it's not only about development, it's also about use more data if our hypotesis is still alive. Automation allow us to be more productive and secure.

> Automation reveals the importance of a good data pipeline.

1.3 Help not substitution

If we want to be successful we need steal the better of other tools, not try to reinvent the wheel. Big Data framework move data from one place to another very well, and also can distribute computation easilly.

> But when we want add a field, or split a date field, it's completely on our hands.

Scientific libraries, also, help us loading data and show information, but if we want do a transformation, again it's in our hands.

Transformations are the more tedious work. And are written directly on the language of the cluster or notebook. The data refinery library comes to help here, giving you code that runs on cluster or notebook. This library it's not a complete ETL solution, but proves to be a huge help when things going real.

This can be done thanks it's minimalist interface and contained scope. Library only do one thing very well, transform data. It build a function that you can use anywhere. So if you need run a small script, you can use it, or if you need to scale up the process, you can use Spark.

1.4 Focus on streaming

Real world is not static, never stop, so we focus on streaming solutions. Sometimes people ask me "where i config the data file?" or "how i obtain all the values of the column?". On streaming environment you only have two things, a function and the current row. If you have a file, it's your work to split into rows and pass it to the transform function. Streamings are endless so you can't know all the values of a column.

As we say, this library does not do work that other frameworks do better. Pandas and Spark can load files easily and send results elsewhere. I strongly recommend you use a framework where you find comfortable, and use this library as help.

Data scientist are usually used to work with Pandas. Its interface is based on columns, very efficient on batch, but not on streaming. If you are used to Pandas, maybe this approach will twist your mind.

> Work thinking about row not column.

This approach it's resource friendly, and you can use on any machine, you only need one row in memory.

2.1 Requirements

- Python 3.5 or 3.6

2.2 Pip

To install datarefinery using pip:

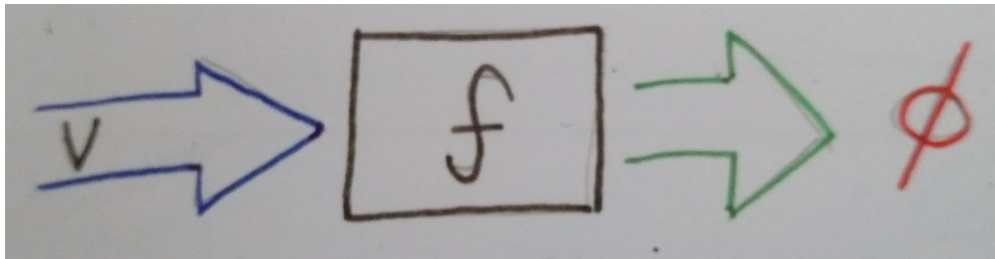
```
pip install datarefinery
```


3.1 Field operations

TODO: explain how we build the operation from field to a complete process

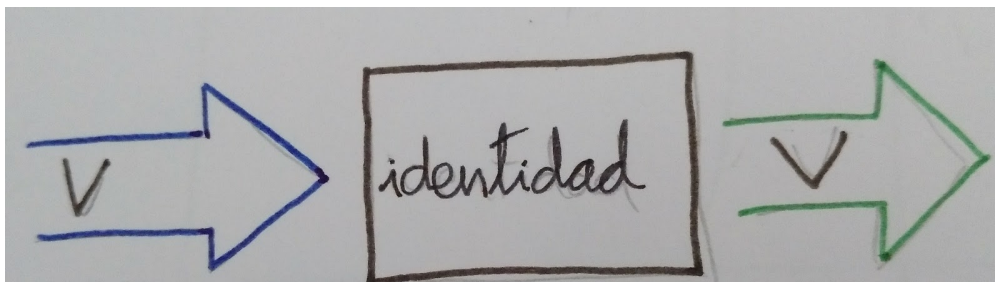
An ETL process is composed by a group of performed operations on a row and rows are composed by groups of columns. So the minimal possible operation is performed on a field.

Suppose a function is useless on a field V :



As the process is composed by operations on fields, if we execute this function to a row with one field V , there won't be output.

Thus we need an usable function, as example a basic function that keeps value of a field. This function returns the input and its name is identity:



3.1.1 How process errors in data

One of the most important situations we have to process is errors on data. A typical issue is when we expect some field but row doesn't have it (unstructured data source).

If any transformation fails we must report precisely without stop the process. All reasons for failure will aggregate by row, so the datascientist will have all the issues regarding a row together, and the process will continue with the next row. Because of this, field error will not stop the row process, and this one will not stop the dataset process. On production environment unexpected data is a common issue, so this data must be logged along their reasons for further research and process improvement.

To fix this problem, Data Refinery transformation functions returns the transformed value and the possible error. This situation takes us to three scenes: - Function returns a value and no error. - Function returns an error and provokes the value doesn't matter. - Function execution is OK, but there is not value.

In the first case we can use the transformed value. In the second one, we don't care about the value because there is an insuperable error. Finally it depends on the current operation, the process rejects the field or fails because there is not value.

3.1.2 List of functions

On data transformation processes, it's common to transform types or values if a condition is triggered, decompose data,.. Data Refinery has defined this kind of functions. Furthermore, if it's necessary we can create our own functions.

More complex situations can be solved chaining together two or more simple functions. This atomic design allows changes on application behaviour with very little user effort.

The majority of ETL's functions require some context parameters in order to produce another, more specialized function. This specialized function can then be used to transform the field. This style of operation improves simple composition given that all functions will have the same input and output instead of a variable number of parameters.

type_enforcer

Expected function in order to cast the value to a type that can rise a cast exception.

Sometimes, data source doesn't contain type or this type is wrong (numbers as string, etc). This function use as context a type changer function. If this function raises an error, type enforcer will return None and the error message raised.

For example, if you want to change the data type to integer:

```
from datarefinery.FieldOperations import type_enforcer

int_enforcer = type_enforcer(lambda x: int(x))
(res, err) = int_enforcer("6")
print(res) # 6
```

min_max_normalization

Typical operation on machine learning. It consists of interpolate a value between 0 (as minimum) and 1 (as maximum). So, this function need a context with the maximum and minimum to run. For example:

```
from datarefinery.FieldOperations import min_max_normalization

normalizator = min_max_normalization(1, 10)
(res, err) = normalizator(5)
print(res) # 0.4444444444444444
```

std_score_normalization

Returns the distance of a particular datum ...

This operation represents how far it's the value from mean, using deviation as scale.

The function requires two input parameters, the column average and the standard deviation. Usage example:

```

from datarefinery.FieldOperations import std_score_normalization

normalizator = std_score_normalization(79, 8)
(res, err) = normalizator(85)
print(res) # 0.75

```

buckets_grouping

Transform a lineal numeric value into a categorical one. For instance it can be used to group users by age.

A minimum of one input value is mandatory. This will produce two groups, the first one from negative infinity to the given value, and the second one from the given value to infinity.

For example, in order to categorize users into three groups (children, adults and elderly) the values 18 and 70 can be passed to the function. This will produce the following groups:

1. From negative infinity to 18.
2. From 18 to 70.
3. From 70 to infinity.

```

from datarefinery.FieldOperations import buckets_grouping

group = buckets_grouping(18, 70)
(res, err) = group(10)
print(res) # 1
(res, err) = group(20)
print(res) # 2
(res, err) = group(73)
print(res) # 3

```

linear_category

Translates the textual value of a field into a numeric value given a list of possible values.

The input value is a list of categories. Keep in mind that this list must always be in the same order to consistently translate the values.

The translated value will be the category index of the list.

As an example, the age can be categorized again, but this time the input a text value instead of a numeric one.

```

from datarefinery.FieldOperations import linear_category

categorizer = linear_category(["niño", "adulto", "jubilado"])
(res, err) = categorizer("adulto")
print(res) # 2

```

column_category

Translates the textual value of a field into a set of columns given a list of possible values. A column will be produced by each one of the members of the input list. This columns will have a value of *0* by default except for the corresponding category that will have a value of *1*. This is known as *one hot vector*.

Example:

```
from datarefinery.FieldOperations import column_category

categorizer = column_category(["niño", "adulto", "jubilado"])
(res, err) = categorizer("niño")
print(res) # OrderedDict([('niño', 1), ('adulto', 0), ('jubilado', 0)])
```

This operation adds new columns, so is usually used along with an event operation of type [append](##Change it).

add_column_prefix

Adds a prefix to the column name. This is useful in a scenario when other function generates a new column with the same name of another already existing.

```
from datarefinery.FieldOperations import add_column_prefix

prefix = add_column_prefix("good")
(res, err) = prefix({"one": "me"})
print(res) # {"good_one": "me"}
```

explode

Flattens a nested data structure even when is made up by a list of objects.

In the case of just one inner object, only the original name prefix will be added.

When multiple objects are present the same prefix will be added and in addition, a numerical suffix (starting on 1) for the second position.

In this example we exploded the field *name*:

```
from datarefinery.FieldOperations import explode
from datarefinery.TupleOperations import append

explode_name = append(["name"], explode("name"))
(res, err) = explode_name({"name": {"first": "Bob", "last": "Dylan"}})
print(res) # {"name_first": "Bob", "name_last": "Dylan"}
```

replace_if

Replaces a value when some condition fulfilled.

Two functions are expected, the former should return a boolean value and the latter should produce a new value in case of the former function returns *True*. Both function will receive the field value.

As an example, if we want to replace by zero all negative values:

```
from datarefinery.FieldOperations import replace_if

change = replace_if(lambda x: x<0, lambda x: 0)
(res, err) = change(-3)
print(res) # 0
```

date_parser - time_parser

Tries to parse a date with the given list of date formats. If none of the formats successfully parses the date then the function returns an error.

The expected formats are Python standard time formats.

```
from datarefinery.DateFieldOperations import date_parser

parser = date_parser(["%Y-%m-%d"])
(res, err) = parser("2017-03-22")
print(res) # <datetime class>
```

There is a similar function to format hours, minutes and seconds.

explode_date - explode_time

Transforms a datetime object to a series of columns with numeric values.

```
from datetime import datetime
from datarefinery.DateFieldOperations import explode_date

(res, err) = explode_date(datetime(2017,3,22))
print(res) # {"year": 2017, "month": 3, "day": 22, "hour":0, "minute": 0, "second": 0}
```

If multiple date exists on the event, please consider using the function [add_prefix](###Prefijo de columna). If no all fields are needed the function [remove column](###Quitando columnas) can be used.

This function is typically used along with *date_parser*.

remove_columns

Removes one or more columns from a set.

This function is usually used along with other functions which generate multiple columns.

In case of not require a column, is better just not to operate it. This non-operated column will be removed automatically.

```
from datetime import datetime
from datarefinery.tuple.TupleDSL import compose
from datarefinery.DateFieldOperations import explode_date
from datarefinery.FieldOperations import remove_columns

only_year_month = compose(explode_date, remove_columns("day", "hour", "minute",
↳ "second"))
(res, err) = only_year_month(datetime(2017,3,22))
print(res) # {"year": 2017, "month": 3}
```

match_dict

Translates values from a table.

```
from datarefinery.FieldOperations import match_dict

d = {"Spain": "ES", "United States of America": "US"}
iso_decoder = match_dict(d)
(res, err) = iso_decoder("Spain")
print(res) # "ES"
```

3.1.3 Composing field operations

This powerful concept from functional programming allow us to build complex applications using simple blocks (functions).

Composition is similar to programming in the sense that a small set of operations can be combined to solve a very large set of problems.

All functions on the library can be combined together using the function *combine*.

This concept is better shown by example.

Normalize Numeric Data

Convert a numeric string to a numeric format and then normalize using min max approximation.

```
from datarefinery.tuple.TupleDSL import compose
from datarefinery.FieldOperations import type_enforcer, min_max_normalization

str_2_min_max = compose(
    type_enforcer(lambda x: int(x)),
    min_max_normalization(0, 100)
)
(res, err) = str_2_min_max("50")
print(res) # 0.5
```

Date data

Explode a date is a typical operation too. It keeps year, month and day as data. Furthermore, it adds a prefix to avoid problems with other fields.

```
from datarefinery.tuple.TupleDSL import compose
from datarefinery.FieldOperations import remove_columns, add_column_prefix
from datarefinery.DateFieldOperations import date_parser, explode_date

complete_date = compose(
    date_parser(["%Y-%m-%d"]),
    explode_date,
    remove_columns("hour", "minute", "second"),
    add_column_prefix("x")
)
(res, err) = complete_date("2017-03-22")
print(res) # {"x_year": 2017, "x_month": 3, "x_day": 22}
```

Day to one hot vector

This example returns one hot vector using a date string and week days.

```
from datarefinery.tuple.TupleDSL import compose
from datarefinery.TupleOperations import wrap
from datarefinery.FieldOperations import match_dict, column_category
from datarefinery.DateFieldOperations import date_parser

week_days={
    0: "Mo", 1: "Tu", 2: "We", 3: "Th", 4: "Fr", 5: "Sa", 6: "Su"
}

def day_of_week(dat):
    return dat.weekday()

day_hot = compose(
    date_parser(["%Y-%m-%d"]),
    wrap(day_of_week),
    match_dict(week_days),
    column_category(week_days.values())
)

(res, err) = day_hot("2017-10-19")
print(res) # OrderedDict([('Mo', 0), ('Tu', 0), ('We', 0), ('Th', 1), ('Fr', 0), ('Sa', 0), ('Su', 0)])
```

3.2 Event operations

Field functions has no affect on the row, so we need Event functions; maybe we need to change the value of a field; or maybe create a new field. Field functions has the same interface. They receives the input, and return result and error. Think in this functions as a tools for dictionaries.

3.2.1 List of functions

keep - Keep files

Keep is the simplest operation, no need of any field function. In essence take the value of a field from the input and put it on the output without change neither the value nor the name of field.

```
from datarefinery.TupleOperations import keep

operation = keep(["greet"])
(res, err) = operation({"greet": "hello", "who": "world"}, {})
print(res) # {"greet": "hello"}
```

If you need keep several similar fields you can use `keep_regexp`.

substitution - Value Substitution

The next operation change the value of a field with the supplied field function. This function will not change the name of the field. By example, given a `to_float` function, you can do this:

```
from datarefinery.TupleOperations import wrap, substitution

operation = substitution(["greet"], wrap(lambda x: len(x)))
(res, err) = operation({"greet": "hello", "who": "world"}, {})
print(res) # {"greet": 5}
```

append - Append new fields

Usually we need add new field or change the name of the field. We can use append to do this, but it expects a field function that return a python dictionary, where every key will be a new field. By example, given a len_cap function that will return the len of a string and the first letter in uppercase:

```
from datarefinery.TupleOperations import wrap, append

operation = append(["greet"], wrap(lambda x: {x: "you", "y": "None"}))
(res, err) = operation({"greet": "hello", "who": "world"}, {})
print(res) # {'hello': 'you', 'y': 'None'}
```

Notice that the field “greet” it’s not in the output. Append only add the result of the function, and the function has no “greet” in their output.

fusion - several fields one output

We already define functions to change the values and add more than one field to the output. But also we can create one field from several inputs.

This function it’s complex, and we have several scenarios. Also need the name of the new field. By example we can sum several numeric fields into total field.

```
from datarefinery.TupleOperations import wrap, fusion

operation = fusion(["a", "b", "c"], "sum_abc", wrap(lambda x: sum(x)))
(res, err) = operation({"a": 1, "b": 2, "c": 3}, {})
print(res) # {'sum_abc': 6}
```

But we can perform some decision with the input fields. Suppose that you need to change the value of money amount field by currency code field. But functions receive only one parameter input, so fusion will put all the values in a list as the input, in the same order that you request in fields field. Knowing that, we can decompose the list into the fields that we need.

By example, given a function to_eur that convert any currency into EUR:

```
def to_eur_wrapped(x):
    [currency, value] = x
    return to_eur(currency, value)
```

We can use this way with fusion:

```
from datarefinery.TupleOperations import wrap, fusion

val_eur_op = fusion(["currency", "value"], "val_eur", wrap(to_eur_wrapped))
(res, err) = val_eur_op({"currency": "USD", "value": 1})
print(res) # {"val_eur": 0.8459}
```


fusion_append - Multiple values in, multiple values out

In fact it's the same that a fusion, but expect that the field function returns a python dict, in the same way than append function.

We can perform the same example that fusion but we can generate new fields with the currency as name of field and the money value as value:

```
from datarefinery.TupleOperations import wrap, fusion_append

def to_eur_cols(x):
    [currency, value] = x
    return {"EUR": to_eur(currency, value), currency: value}

val_eur_op = fusion_append(["currency", "value"], "val_eur", wrap(to_eur_cols))
(res, err) = val_eur_op({"currency": "USD", "value": 1})
print(res) # {"EUR": 0.8459, "USD": 1}
```

filter_tuple - Discard irrelevant

When we don't need all data, the best approach it's discard as soon as possible the rows that we don't need.

With no_none function and filter field operation the generated function will return None as result if the event it's discarded:

```
from datarefinery.TupleOperations import wrap, filter_tuple

no_none = filter_tuple(["value"], wrap(lambda x: x is not None))

(res, err) = no_none({"value": None})
print(res) # None
```

It's up to you not to fail when result it's None, if no error means that event is discarded.

alternative - Plan B

Some times we have several ways to transform the event. If the first approach fail, alternative will try the next, until success or the last fail.

Suppose that you want to multiply by two, but if this operation fail, you want to append value 0.

```
from datarefinery.TupleOperations import wrap, alternative, substitution, append

need_value = alternative(
    substitution(["value"], wrap(lambda x: x*2)),
    append(["name"], wrap(lambda x: {"value": 0}))
)
(res, err) = need_value({"name": "John"})
print(res) # {"value": 0}
```

3.2.2 Combine event operations

Normally, a transformation is a group of different type of functions, not only one type. For example, you want to keep some fields and change the value of a field using a function.

So that's why we need an interface for doing this kind of transformations. Data Refinery has two functions to perform this kind of work: parallel and sequential.

With sequential you can perform transformations that affect the same data several times, like, generate two new fields and then sum together.

With parallel you can change several fields using the same input and with no interactions between them.

For example, if we want to keep a field and replace another field value with a x2 function (multiply by two):

```
from datarefinery.TupleOperations import wrap, keep, substitution
from datarefinery.CombineOperations import parallel

x2 = wrap(lambda x: x*2)

operation = parallel(keep(["name"]), substitution(["value"], x2))
(res, err) = operation({"name": "John", "value": 10})
print(res) # {"name": "John", "value": 20}
```

A world of possibilities

We can save our set of transformations every time we want because it's just a plain function.

Training data is often as production data, but it contains an extra field called "label". This field indicates what machine learning model must learn.

In next example, data transformation is a specified module of your application. You can get it with *etl* function. Then we will add the label logic:

```
from datarefinery.TupleOperations import keep
from datarefinery.CombineOperations import parallel

op = etl()
if training == True:
    operation = parallel(op, keep("label"))
else:
    operation = op
```

By this way, the output will contain the label without knowing other transformations over the data during training phase.

Sequentiality

When you use parallel to create a group of operations, these are executed in one step. In other words, all operations use the same input and write in the same output. So, if we want to modify the value of a field which have been modified yet, even using parallel, all operations take place at the same time and we only see the last transformation.

For example:

```
from datarefinery.TupleOperations import wrap, substitution
from datarefinery.CombineOperations import parallel

x2 = wrap(lambda x: x*2)

operation = parallel(substitution(["value"], x2), substitution(["value"], x2))
(res, err) = operation({"value": 2})
print(res) # {"value": 4}
```

In this example, you expect output value will be 8, but it's not true. When operations run in parallel, the execution is something like this:

input	value (1st time)	value(2nd time)
2	4	4

The input is the same for both operations and the output of the first operation was overwritten by the second one.

If you want to do different operations on the same value field, then you have to use *compose*. This function lets operations run sequentially, as the next example:

```
from datarefinery.TupleOperations import wrap, substitution, compose

x2 = wrap(lambda x: x*2)

operation = substitution(["value"], compose(x2,x2))
(res, err) = operation({"value": 2})
print(res) # {"value": 8}
```

3.3 Pandas use

You can use Pandas Data Frame easily with this function:

```
import pandas as pd

def pandas_dataframe_operator(df, operation):
    if df is None or operation is None:
        return None, None
    raw_df = pd.DataFrame(df.apply(operation, axis=1), columns=["res"])

    raw_df['ok'] = raw_df["res"].apply(lambda x: x[0])
    raw_df['ko'] = raw_df["res"].apply(lambda x: x[1])

    raw_df = pd.concat([df, raw_df[['ok', 'ko']]], axis=1)

    ok_filter = raw_df["ok"].apply(lambda x: x is not None)
    ok_df = raw_df[ok_filter]
    ok_df = pd.DataFrame(ok_df["ok"].apply(pd.Series))
    ko_df = raw_df[~ok_filter]
    ko_df = pd.concat([ko_df, pd.DataFrame(ko_df["ko"].apply(pd.Series))], axis=1)
    del ko_df['ok']
    del ko_df['ko']

    return ok_df, ko_df
```

This function expect a data frame and any tuple operation, and returns two data frames. One it's the correct tuples transformed, the other is the original data frame with error columns added, this way you can easily fix the tuple operation.

You can check the example of use in tests/examples.

3.4 Review exercises

If you want or need to do basic exercises for review all your knowledge, you can execute the next notebook on your Jupyter instance: [thebasics.ipynb](#).

4.1 DSL

DSL let us create any operation over a row, but you know... with great power comes great responsibility.

When you create an operation, you have to spread the input, the modified output and the error if proceeds.

Its usage is very simple. For example, *keep* function is shown below:

```
def keep(fields) -> Callable[[dict, dict, dict], Tuple[dict, dict, dict]]:
    operations = [compose(use_input(), read_field(f), write_field(f)) for f in fields]
    return reduce(compose, map(apply_over_output, operations))
```

In this example you can see a composition that takes different functions for every field. *keep* returns an unique function that contains all of them applied on output.

4.1.1 List of functions

- use_input
- use_output
- read_field
- read_match
- read_fields
- write_field
- write_error_field
- dict_enforcer
- apply_over_output
- compose

You can create your own function too. It has to return a function that takes 3 dictionaries (input, output and error) and returns them modified if it's necessary.

5.1 datarefinery

5.2 datarefinery.FieldOperations

5.3 datarefinery.tuple.Formats

5.4 datarefinery.tuple.TupleDSL

5.5 datarefinery.TupleOperations

5.6 datarefinery.DateFieldOperations

5.7 datarefinery.CombineOperations

ETL

Extract, Transform and Load.

Field

Element of an event which was identified by name. For example, the first column of an excel sheet is identified as A field.

Row

Minimal unit of transformation at ETL.

DSL

Domain Specific Language

Event

Operations over a tuple or row.

CHAPTER 7

Indices and tables

- `genindex`
- `modindex`
- `search`