https://stochasticsolutions.com/pdf/toronto2021-reproducibility-in-safe-haven-slides.pdf

REPRODUCIBILITY IN A SAFE HAVEN



GLOBAL OPEN FINANCE

CENTRE OF EXCELLENCE

Nicholas J. Radcliffe, Ph.D.

The Global Open Finance Centre of Excellence & Department of Mathematics, University of Edinburgh & Stochastic Solutions Limited



Stochastic Solutions

Pei Shan Yu, Ph.D.

The Global Open Finance Centre of Excellence University of Edinburgh













REPRODUCIBLE RESEARCH

(If a tree falls in a forest, and no one is around to hear it, does it make a sound?)

actually REPRODUCED RESULTS



LINUS'S LAW

"Given enough eyeballs, all bugs are shallow"*

Just a claim; not actually true

TEST-DRIVEN DEVELOPMENT (TDD)



- Well-understood goal
- *Many kinds of errors/failures are unmistakable*

TESTS (ADDITION)

2 + 2 = 4

-1 + 1 = 0

 $10^{100} + 10^{100} = 2 \times 10^{100}$

0 2 x 10100

TEST-DRIVEN DEVELOPMENT (TDD)



- Well-understood goal
- *Many kinds of errors/failures are unmistakable*

TEST-DRIVEN DEVELOPMENT (TDD)



- Well-understood goal
- *Many kinds of errors/failures are unmistakable*

Constantly run tests with CI?



CONTINUOUS **INTEGRATION** SYSTEM

Transform the derta Try to understand the data Generate results Formulate an analytical Opproach Drown Sorrows Try to formulate the problem Make sense? Try that approach Eyeball the data -Show to a colleague Segment & profile Discover the approact cloest work Make surce? > Discover you don't understand the data Show to expert Question ann Sounty Malie sense? Discover the dates Deploy 1 Distribute Curse is wrong Question others' sanity REFORMULATE K Make sense? Re-source the data

"Shouldn't we do test-driven data analysis?"

— Patrick Surry, c. 2010

$\mathsf{TDD} \mapsto \mathsf{TDDA}$

- We need to extend TDD's idea of testing for
 - software correctness
 - with the idea of testing for
 - meaningfulness of analysis,
- correctness and validity of input and output data,
 - & correctness of interpretation.

"test-driven data analysis"



- sustainability
- Incubated in the Edinburgh Futures Institute at University of Edinburgh \bullet
- assessment

• Non-profit established as global economic observatory to study the way people and businesses earn, spend and save with a view to improving financial outcomes, equity and

Currently looking at the financial impact of COVID-19 on citizens and businesses in UK

Outputs to national, regional and local governments to guide policy development and





No transfers (file, copy & paste etc.) to/from the user's computer to the Safe Haven Data ingress and egress only though Research Coordinators (information governance)

via secure, encrypted, managed file transfer processes





To date, every output shared with government has been reproduced by two separate software systems / teams in the Safe Haven

CORE ASPECTS

- Data ingestion and verification
 Checks for personal/identifying data and data matches specs
- 2. Replicating outputs and testing analytical code and processes
- 3. Output disclosure control
 Vetting of outputs, audit trails, data provenance & lineage

REPRODUCIBILITY & TDDA

1. Data ingestion and verification

3. Output disclosure control

- Vetting of outputs, audit trails, data provenance & lineage

- Checks for personal data and data matching specs

2. Replicating outputs and testing analytical code and processes

Reproducible Research **TDDA**

KΕ





CONSTRAINT GENERATION, DATA VERIFICATION & ANOMALY DETECTION

Install from PyPI (recommended) pip install tdda or from Github (source) python setup.py install

TDDA LIBRARY

git clone https://github.com/tdda/tdda.git

- Runs on Python 2* & Python 3, Mac, Linux & Windows, under unittest and pytest
- MIT Licensed
- Documentation:
 - Sphinx source in **doc** subdirectory
 - Built copy at http://tdda.readthedocs.io
- Quick reference:

* But no one should still be using Python 2!

TDDA LIBRARY

http://www.tdda.info/pdf/tdda-quickref.pdf

AUTOMATIC CONSTRAINT GENERATION



TRAINING DATA

(believed to *be "good"*)

AUTOMATIC DISCOVERY OF CONSTRAINTS C1: Age ≥ 0 C2: ID is not null C3: CardNumber ~ DDDD DDDD DDDD DDDD

> DISCOVERED **CONSTRAINTS**



```
"creation metadata": {
  "as at": "2021-01-20 13:51:08",
  "local time": "2021-01-20 13:51:08",
  "utc_time": "2021-01-20 13:51:08",
  "creator": "Miro 3.1.18",
  "creation_command": "discover -o /tmp/tbank",
  "creation session log":
      "~/miro/log/2021/01/20/session073.html",
  "source": "banksc-raw",
  "host": "gofcoe01",
  "user": "njr",
  "dataset": banksc-raw",
  "n records": 205723986,
  "n selected": 205723986,
  "tddafile": "~/citizendata/qa/banksc.tdda",
  "last edit": "Nick Radcliffe",
  "version": "1.1",
  "last edit date": "2021-02-25 16:00:00"
},
"fields": {
  "pseudoid": {
    "type": "string",
    "min_length": 40,
    "max_length": 40,
    "max nulls": 0,
    "rex": [
      "^([0-9a-z]{40})$",
  "end_of_previous_period": {
    "type": "date",
    "min": "2018-12-30",
    "max": "2022-01-01",
    "max_nulls": 0
  },
  "end_of_this_period": {
    "type": "date",
    "min": "2019-01-07",
    "max": "2022-01-01",
    "max nulls": 0
  },
```

```
"end_of_this_period:time-before-now": {
  "min": "6 days",
   "transform": "time-before-now"
},
"end_of_this_period:time-before-now": {
  "min": "0 days",
  "transform": "time-before-now"
},
"postal district": {
  "type": "string",
  "min_length": 2,
  "max length": 5,
  "max nulls": 0,
  "rex": [
     "^([A-Z]{1,2})([0-9]{1,2})([A-Z]?)$"
},
"sex": {
  "type": "string",
  "min_length": 1,
  "max_length": 1,
  "max nulls": 0,
  "allowed values": [
    "F",
    " M "
"total credits": {
  "type": "real",
  "min": 0,
  "max": 10000000,
  "sign": "non-negative",
  "max nulls": 0
"total debits": {
  "type": "real",
  "min": -1000000,
  "max": 0,
  "sign": "non-positive",
  "max_nulls": 0
},
```

```
"final balance": {
    "type": "real",
    "min": -1000000,
    "max": 1000000,
    "max nulls": 0
},
"dataset": {
  "required fields": [
    "pseudoid",
    "end_of_previous_period",
    "end_of_this_period",
    "postal district",
    "sex",
    "total_credits",
    "total debits",
    "final_balance"
  ],
  "allowed fields": [
    "pseudoid",
    "end of previous period",
    "end_of_this_period",
    "postal district",
    "sex",
    "total_credits",
    "total debits",
    "final_balance"
  "allowed_all_null_fields": []
},
"field groups": {
  "end_of_previous_period,end_of_this_period": {
    "lt": true
```



OPERATIONAL DATA



C1: ...

C2: ...

C3: ...

CONSTRAINTS









GENERATING CONSTRAINTS & VERIFYING DATA



training data

operational data

REFERENCE TESTS



DATA & PARAMETERS

> Record inputs

Capture as scripted, parameterised executable procedure ("reproducible research")

Develop a verification procedure (diff) *and periodically rerun: do the same inputs (still) produce the same (or equivalent) outputs?*

ANALYTICAL PROCESS

OUTPUTS

DATASETS, NUMBERS, GRAPHS, MODELS, DECISIONS ETC.

Record ("reference") outputs

TDDA LIBRARY SUPPORT FOR REFERENCE TESTS

- hand
- Output artefacts can be large and complex vary even when semantic results don't vary
 - numbers, environment information etc.
- remain consistent does, while allowing variation in the artefacts results change.

• With Data Science, the outputs are often too complex to produce by

- Output artefacts include random IDs, dates, run numbers, version

• TDDA library provides support for testing that everything that should themselves, and regeneration when fixes mean that reference (target)



TDDA LIBRARY FEATURES



GENERATE CONSTRAINTS FROM DATA

VERIFY DATA WITH CONSTRAINTS

TEST SUPPORT FOR CODE WITH COMPLEX, **NOT-ALWAYS-IDENTICAL** OUTPUTS



CLASSES OF ERRORS IDENTIFIED

(more detail available in written transcript)







Inconsistent handling of invalid values



CSV precision (weird behaviour in R library write.csv in utils package in R)

- Rolling values with missing weeks for some customers
- Inconsistent handling of nulls (missing values)





CLASSES OF ERRORS IDENTIFIED

(more detail available in written transcript)



Numerical accuracy

- · Comparisons to nearest half penny
- · Counts above thresholds Expressing in whole pennies and rounding to ints



Bin allocation inconsistent because of numerical (in)accuracies



Incorrect Loop Unrolling (copy/paste error)



Inconsistent sort ordering (and mode tiebreaks)



ORIGINAL ANALYSIS PIPELINE



RAW DATA

ANALYSE

FINAL REPORT

DECOMPOSED ANALYSIS PIPELINE

MEASURE



RAW DATA EXTRACT, TRANSFORM LOAD



ANALYSE

APPLY DISCLOSURE CONTROLS

GRAPHING, MAPPING REPORT CONTRUCTION



OUTPUT VERIFICATION





CORRECTNESS: TESTS REPRODUCED









CHECKS FOR PERSONAL DATA ("PII")





AUDIT TRAIL & DATA LINEAGE

/.../audit/miro/xxxxx-by-pa-enhanced.csv
(hash 2a4cf239fed057281c3a513888492496)
compared to

/.../audit/r/xxxxx-by-pa-enhanced.csv
(hash 59e5540d687daff81f066cc7884efb9c)
validated as equivalent on 2021-01-05T10:34:51Z.
:

All files produce for report validated successfully.

Closing thoughts

The Data Scientist's Version: How is this misleading data misleading me?

Why is this lying bastard lying to me? --- Louis Heren (Journalist, The Times)



*not William Blake!

ABC (Blake's Dictum*)

Always Be Closing

- Blake (Alec Baldwin) Glengarry Glen Ross

The Data Scientist's Version:

Always Be Checking





http://tdda.info



https://github.com/tdda





https://linkedin.com/in/njradcliffe



#tdda*



Correct interpretation: Zero (Error of interpretation: Letter "Oh")

https://stochasticsolutions.com/pdf/toronto2021-reproducibility-in-safe-haven-transcript.pdf https://stochasticsolutions.com/pdf/toronto2021-reproducibility-in-safe-haven-slides.pdf

globalopenfinance.com • stochasticsolutions.com

njr@StochasticSolutions.com Pei-Shan.Yu@ed.ac.uk

* *tweet (DM) us your email address for invitation. Or email me.*