# Preparing Data Files for Preservation with Colectica Datasets

°C
OIL TEMP
ENG

OIL PRESS.
ENG

PSI
FUEL PRESS

GALS
FUEL QTY

°C
OIL TEMP
TRANS

PSI
OIL PRESS
TRANS

LOADMETER

CLIM

UP

DOWN

VOLT
DC

WACLINE
F20301A

MODEL FLD
FS=30 DCV

MONTE CARLO
HEUER

VOR DEM ABSTELLEN
TOT 2 MINUTEN
IN LEERLAUFDREHZAHL
STABILISIEREN

TRANS
CHIP
T/R
CHIP
A/F FUEL
FILTER
ENG
CHIP
DET
ENG
Gen
SPARE
FUEL
FUEL
TRANS
TRANS
GEN
FAIL
PRESS
OUT
THIS HELICOPTER MUST BE
IN COMPLIANCE WITH THE
LIMITATIONS SPECIFIED IN FAR
ROTORCRAFT FLIGHT MANUAL
COCKPIT WEIGHT 200 LBS W
RANGE EXTENDER FUEL C

°C
OIL TEMP
ENG

OIL PRESS.
ENG

PSI
FUEL PRESS

GALS
FUEL QTY

°C
OIL TEMP
TRANS

PSI
OIL PRESS
TRANS

LOADMETER

VOLT
DC

VOR DEM ABSTELLEN
TOT 2 MINUTEN
IN LEERLAUFDREHZAHL
STABILISIEREN

# Context Matters

# Metadata provides context for Data

# Data

# Metadata

| Variable | Name | Type | | Width | Decimals | Label | Value Labels | | Missing Va |
|---|---|---|---|---|---|---|---|---|---|
| 1 | M2ID | Numeric | ... | 5 | 0 | MIDUS 2 ID number | None | ... | None |
| 2 | M2FAMNUM | Numeric | ... | 6 | 0 | MIDUS 2 Family number | None | ... | None |
| 3 | SAMPLMAJ | Numeric | ... | 8 | 0 | Major sample identificatior | {1, MAIN RDD}... | ... | None |
| 4 | C1STATUS | Numeric | ... | 1 | 0 | Completion status of M3 re | {1, COMPLETED M3 CATI ONLY | | None |
| 5 | C1PRAGE | Numeric | ... | 2 | 0 | Respondent's age | None | ... | None |
| 6 | C1PBYEAR | Numeric | ... | 4 | 0 | Respondent's year of birth | None | ... | None |
| 7 | C1PRSEX | Numeric | ... | 1 | 0 | Respondent's sex | {1, MALE}... | ... | None |
| 8 | C1PIDATE_MO | Numeric | ... | 8 | 0 | Interview date - Month | None | ... | None |
| 9 | C1PIDATE_YR | Numeric | ... | 8 | 0 | Interview date - Year | {9997, DON'T KNOW}... | ... | None |
| 10 | C1PAA1 | Numeric | ... | 1 | 0 | Recession began with spec | {1, YES}... | ... | 7, 8 |

# Metadata

| 4 | C1STATUS | Numeric | ... | 1 | 0 | Completion status of M3 re | {1, COMPLET |
|---|----------|---------|-----|---|---|---------------------------|-------------|
| 5 | C1PRAGE | Numeric | ... | 2 | 0 | Respondent's age | None |
| 6 | C1PBYEAR | Numeric | ... | 4 | 0 | Respondent's year of birth | None |

# Statistical Tools have Limited Metadata

- ☐ Data Types
- ☐ Variable Labels
- ☐ Value Labels

# No Metadata ☹

| Variable | Name | Type | | Width | Decimals | Label | Value Labels | | Missing Values | | Co |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | NEWID | String | ... | 8 | | | None | ... | None | ... | 8 |
| 2 | DIRACC | String | ... | 1 | | | None | ... | None | ... | 1 |
| 3 | DIRACC_ | String | ... | 1 | | | None | ... | None | ... | 1 |
| 4 | AGE_REF | Numeric | ... | 8 | 0 | | None | ... | None | ... | 8 |
| 5 | AGE_REF_ | String | ... | 1 | | | None | ... | None | ... | 1 |
| 6 | AGE2 | Numeric | ... | 8 | 0 | | None | ... | None | ... | 8 |
| 7 | AGE2_ | String | ... | 1 | | | None | ... | None | ... | 1 |
| 8 | AS_COMP1 | Numeric | ... | 8 | 0 | | None | ... | None | ... | 8 |
| 9 | AS_C_MP1 | String | ... | 1 | | | None | ... | None | ... | 1 |
| 10 | AS_COMP2 | Numeric | ... | 8 | 0 | | None | ... | None | ... | 8 |
| 11 | AS_C_MP2 | String | ... | 1 | | | None | ... | None | ... | 1 |
| 12 | AS_COMP3 | Numeric | ... | 8 | 0 | | None | ... | None | ... | 8 |
| 13 | AS_C_MP3 | String | ... | 1 | | | None | ... | None | ... | 1 |

Open... | Save | Go To Variable... | Insert Variable | Split File... | Weight Cases... | Value Labels

# The Metadata Problem

- ☐ Metadata is needed to understand data
- ☐ Statistical tools have limited metadata capabilities

# The Goal

- Make a simpler way to add rich metadata to statistical files

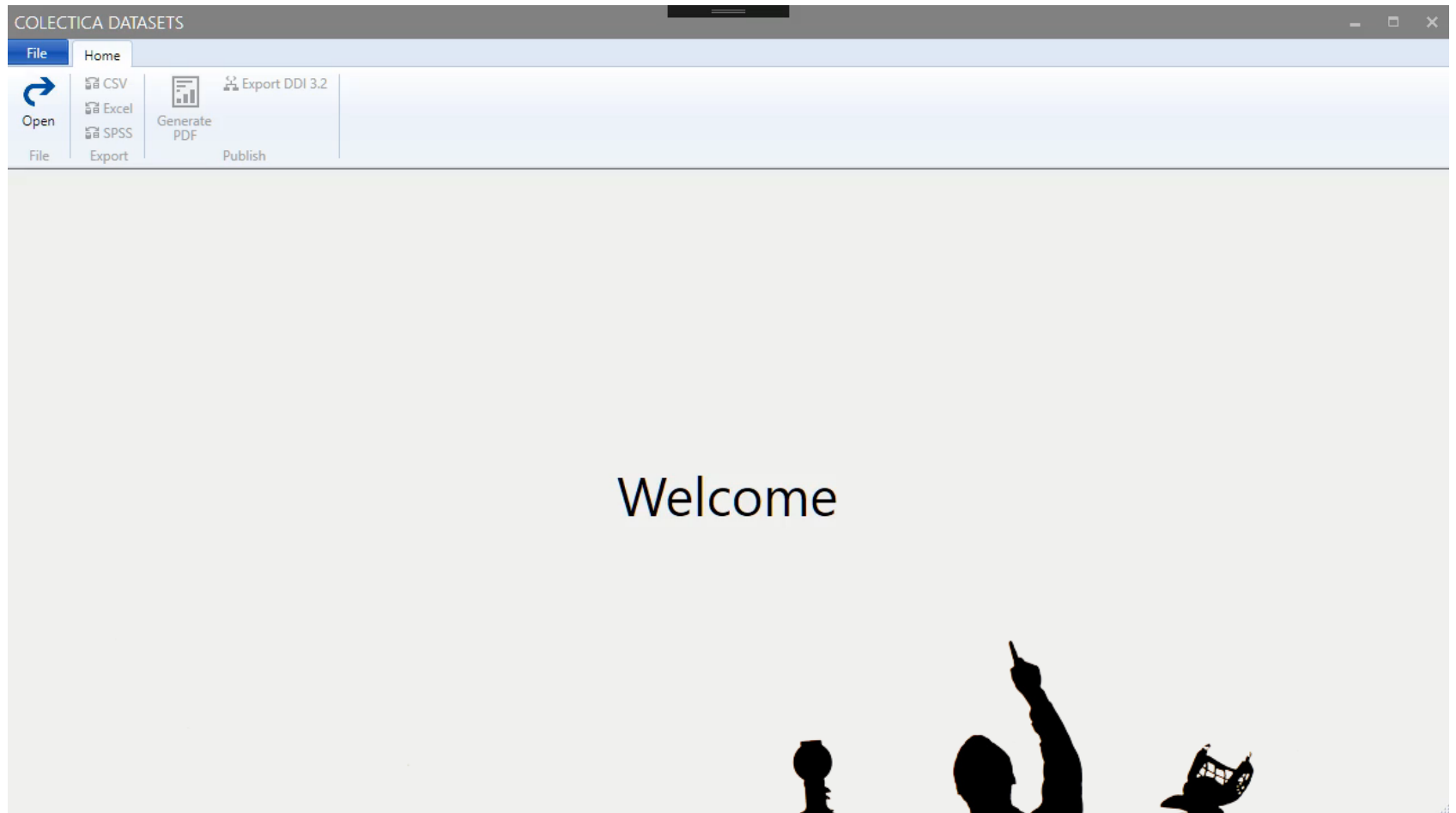# Existing Metadata Tools

- Not quite early enough and easy enough

# Colectica Datasets

- View
- Improve
- Publish

# Explore

# View Data

# Soon: Metadata Quality Checks

□ Data quality report card

- Missing metadata
- Spell checking
- Extensible validations
- Easily fix issues directly in the app

# Publish

# Data Dictionary

M2 and MKE are not directly comparable with M1, MR, MKER: M2 and MKE responses were coded as one single variable, while other waves broke out month and year separately.

| Valid | Invalid | Minimum | Maximum |
|-------|---------|---------|---------|
| 3294 | 0 | 2013 | 2014 |

## C1PAA1

**Label**
Recession began with specific event

**Role**
input

**qstnLit**
Think back to 2008 when the recession first began. Was there a specific event that made you aware that the recession of 2008 had begun?

**Comparability Class**
New Item - Difference exists among waves because some new items were introduced only in particular waves.

**Comparability Notes**
Recession item introduced since MR and MKER.

| Value | Label | Frequency | % |
|-------|-------|-----------|---|
| 1 | YES | 1,970 | 59.8% |
| 2 | NO | 1,308 | 39.7% |
| 7 | DON'T KNOW | 15 | 0.5% |
| 8 | REFUSED | 1 | 0.0% |

# Publish

- Publish to more locations
  - BagIt
  - Figshare
  - Zenodo

# Availability

- 2017
- For Windows and macOS

# Technical Previews

- Get in touch if you'd like to help test and give feedback

colectica

**25** | # Thank you

| | Web | colectica.com |
|---|---|---|
|  | Twitter | @Colectica |
|  | YouTube | youtube.com/colectica |