

Delivering Data Packages for Discovery, Analysis, and Preservation

Research Data Alliance
11th Plenary
Berlin, Germany,
March 21-23, 2018



The United States Department of Transportation (USDOT) Plan to Increase Public Access to the Results of Federally-Funded Scientific Research (PA) requires, in part, "digitally formatted scientific data resulting from unclassified research supported wholly or in part by Federal funding to be stored and publicly accessible for search, retrieval, and analysis." The PA goes on to require that researchers deliver a Data Management Plan (DMP) which identifies practices they will employ to ensure the long-term access and preservation of the project data.

The National Transportation Library (NTL), of USDOT's Bureau of Transportation Statistics (BTS), is applying the

standards and practices required by the PA to new datasets created by BTS. Further, in order to ensure the greatest possible longevity of discovery and preservation, as well as encouraging data interoperability and reuse, NTL is helping BTS staff to create even more robust documentation for datasets. This documentation is collectively known as a "data package." In addition to the final dataset, and the DMP, a data package includes other documentation, which as defined by NTL, is "needed to contextualize the dataset for any and all users."

This poster will explore the various elements of data packages, and look at their initial use within BTS and NTL, using the 2016 National Census of Ferry Operators dataset as an example.

What is a "Data Package"?

A "Data Package" is the dataset, the data management plan (DMP), and all other documentation needed to contextualize the dataset for any and all users and re-users.

At the National Transportation Library, a data package is comprised of 4 required files - or elements - and 2 (or more) optional elements. Each plays a different role in contextualizing a dataset.

Look in the back of truck for the inventory of Data Package Elements.

Element: Data

```
bts_osp_national_census_ferry_operators_2016_Segment_2017_10_26.csv - Notepad
File Edit Format View Help
SEGMENT_ID, SEGMENT_NAME, SEG_TERMINAL1_ID, SEG_TERMINAL2_ID, SEG_TYPE, SERVES_NPS, SURVEYYEAR
222, Ketchikan Terminal - Prince Rupert, 45, 1, 3, 0, 2016
223, Anacortes - Lopez, 613, 534, 1, 0, 2016
224, Friday Harbor - Orcas, 862, 541, 1, 0, 2016
226, Lopez - Shaw, 534, 550, 1, 0, 2016
227, Orcas - Anacortes, 541, 613, 1, 0, 2016
228, Shaw - Orcas, 550, 541, 1, 0, 2016
229, Anacortes - Friday Harbor, 613, 862, 1, 0, 2016
237, Sombra ONT - Marine City, 11, 295, 3, 0, 2016
239, Victoria Inner Harbor - Black Ball Ferry Line, 4, 542, 3, 0, 2016
264, Road Town, Tortola - w. Blyden Marine Terminal, 19, 492, 1, 0, 2016
265, West End, Tortola - w. Blyden Marine Terminal, 20, 492, 1, 0, 2016
276, Enighed Pond - Red Hook, 494, 495, 1, 0, 2016
312, Catano - Old San Juan Pier Two, 447, 452, 1, 0, 2016
```

Best data file format for long-term preservation of tabular data: .CSV

- Open, non-proprietary file format
- Backward and forward compatible
- Long-term preservable

Best Practice Tip:

In order for data to be machine readable and interoperable, it should have only two types of rows:
1) a single header row, and 2) all others rows are data.

Element: Data Management Plan

A data management plan (DMP) is a succinct narrative document which describes the deliberate planning, creation, storage, access, and preservation of data produced from a given investigation.

A robust DMP contains 5 sections:

1. Data Description
2. Standards Employed
3. Access Policies
4. Re-Use, Redistribution, and Derivative Products Policies
5. Archiving and Preservation Plans



Best Practice Tip:

To write a DMP that conforms to the USDOT Public Access Plan, go to the NTL Public Access Guidance website at: <https://ntl.bts.gov/publicaccess/>

Element: Readme.txt

The purpose of the README document is to provide all of that contextual information that could not be included in the data file because it would break the data.

Data Dictionary Example:

AN. Name: SURVEYYEAR
Full name: Survey Year
Variable Definition or Description: Year of survey.
Data format: YYYY
Field length: 4
Representation of Null values: No Null values present; Not Applicable.

NTL ReadMe.txt Template Contents

- A. General Information
- B. Sharing/Access Information
- C. Data and Related File Overview
- D. Methodological Information
- E. Data-Specific Information, Data dictionary
- F. Appendices

Best Practice Tip:

In order for you and future users to know which of the dozens of "README.txt" files on your desktop go with which datasets, use clarifying, human-readable file names, such as:

bts_osp_national_census_ferry_operators_2016_README_2017_10_26.txt

NTL Dataset Data Package Elements

- 1) **Dataset**
 - .csv or other open format
- 2) **Readme.txt**
 - Includes Data Dictionary
 - Notes standards used
 - Defining Zero, Null, and Unknown
 - FAQs and other notes
- 3) **Metadata file**
 - in Project Open Data .json
- 4) **Data Management Plan (DMP)**
- 5) **Code or scripts** used in data analysis
- 6) **Supporting files, tables, etc.**

(Bold = Required; Italics = Optional, or Required if Applicable)

Element: Metadata

Project Open Data Metadata V1.1

Machine-readable; Required by DOT open data policy; Required by data.gov

```
bts_osp_national_census_ferry_operators_2016_Metadata_2017_10_26.json - Notepad
File Edit Format View Help
{
  "conformsTo": "https://project-open-data.cio.gov/v1.1/schema",
  "title": "2016 National Census of Ferry Operators (NCFO)",
  "description": "The 2016 NCFO dataset is comprised of the responses of all operators who completed the 2016 National Census of Ferry Operators (NCFO) survey. The dataset includes information on ferry routes, terminals, vessels, and other information. The data was collected from 2016 to 2016 and is available for download from the National Transportation Library (NTL) website.",
  "keyword": "passenger, freight, transit, travel, maritime, terminals, vessels, ferry",
  "modified": "2017-10-06",
  "publisher": "U.S. Department of Transportation, Bureau of Transportation Statistics, Office of Survey Programs",
  "contactPoint": {
    "type": "vcard:contact",
    "fn": "Janine McFadden",
    "hasEmail": "janine.mcfadden@dot.gov"
  },
  "identifier": null,
  "accessLevel": "public",
  "accessLevelComment": null,
  "bureauCode": "021:04",
  "programCode": "021:053",
  "accessURL": "https://data.transportation.gov/Public-Transit/2016-National-Census-of-Ferry-Operators-NCFO-5K",
  "webService": null,
  "format": "csv",
  "license": "Public Domain",
  "spatial": "United States",
  "temporal": "2015",
  "accrualPeriodicity": "R/P2Y",
  "landingPage": "https://www.bts.gov/product/highlights-2016-national-census-ferry-operators-ncfo"
}
```

Element: Code Book & Supporting Files

Include code books, scripts used during analysis, auxiliary tables, and other supporting files that were created or used to collect, process, clean, or analyze the data.

Best Practice Tip:

Be as complete as possible. The goals of a robust data package include fully documenting your processes to allow a naïve user to replicate your results and understand the full context of the data. This will enable them to decide intelligently whether your data meets their reuse needs.

A Data Package is...

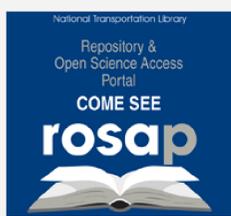
The dataset and all documentation needed to contextualize the dataset for any user.



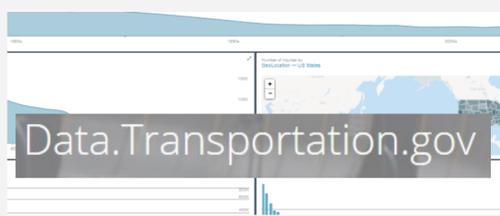
Best Practice Tip:

Use consistent, descriptive, human-readable file names, with date and timestamps to help you keep your data package well organized and up to date.

Where does NTL Deliver Data Package?



<https://rosap.ntl.bts.gov/>



<https://data.transportation.gov/>

Leighton L Christiansen <http://orcid.org/0000-0002-0543-4268>

Data Curator, NTL leighton.christiansen@dot.gov @purpleleighton

Recommended Citation: Christiansen, Leighton L. <http://orcid.org/0000-0002-0543-4268>. 2018. "Delivering Data Packages for Discovery, Analysis, and Preservation." Research Data Alliance 11th Plenary. Berlin, Germany.

Acknowledgements: The author thanks Alpha Wingfield for design assistance, & the National Transportation Library staff for their input, support, and teamwork.

Delivery Truck Image Credits: Font Awesome by Dave Gandy - <http://fortawesome.github.com/Font-Awesome> [CC BY-SA 3.0] (<https://creativecommons.org/licenses/by-sa/3.0/>), via Wikimedia Commons. https://commons.wikimedia.org/wiki/File%3ATruck_font_awesome.svg