# Data Fidelity in Formulation

Trevor Beer
Systems Engineering APX, 312C - Systems Modeling,
Analysis, and Architectures
June 13, 2024

**Jet Propulsion Laboratory**
California Institute of Technology

# About Me

- Completed undergraduate studies at UCSB in Statistics and Data Science
- Graduating with MS in Statistics from UCLA
- Thesis - Generative Data Science: Applications for Early Life Cycle Cost Estimation in the Aerospace Industry
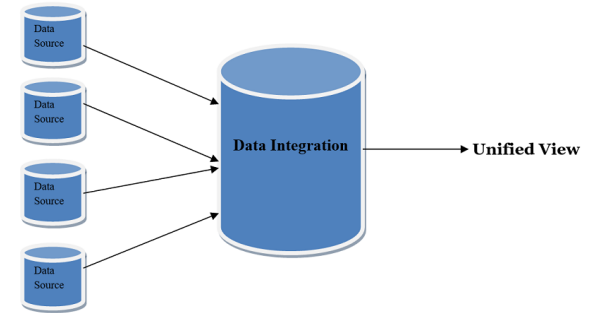- Joined group in April 2023

# The Big Picture

- Group 312C – Systems Modeling, Analysis, and Architectures
- Assist in critical cost estimation tasks during mission formulation activities
- Capabilities range - statistical modeling, data engineering (normalization is a subset of this), web-based tool development
- Dynamic group with members that come from a variety of disciplines – presented many opportunities to broaden my skillset

# My Contributions to the Group

- Supported, modified and built web-apps using Pythonic software to provide user friendly modeling and data visualization tools

- Harvested, merged and processed raw data into easily ingestible, code-friendly datasets

- Normalized incorrect, outdated and/or inconsistent data into clean datasets to ensure high quality data for our modeling efforts

- Documented all of my code, scripts and work into Github or Jupyter notebook to make all work traceable
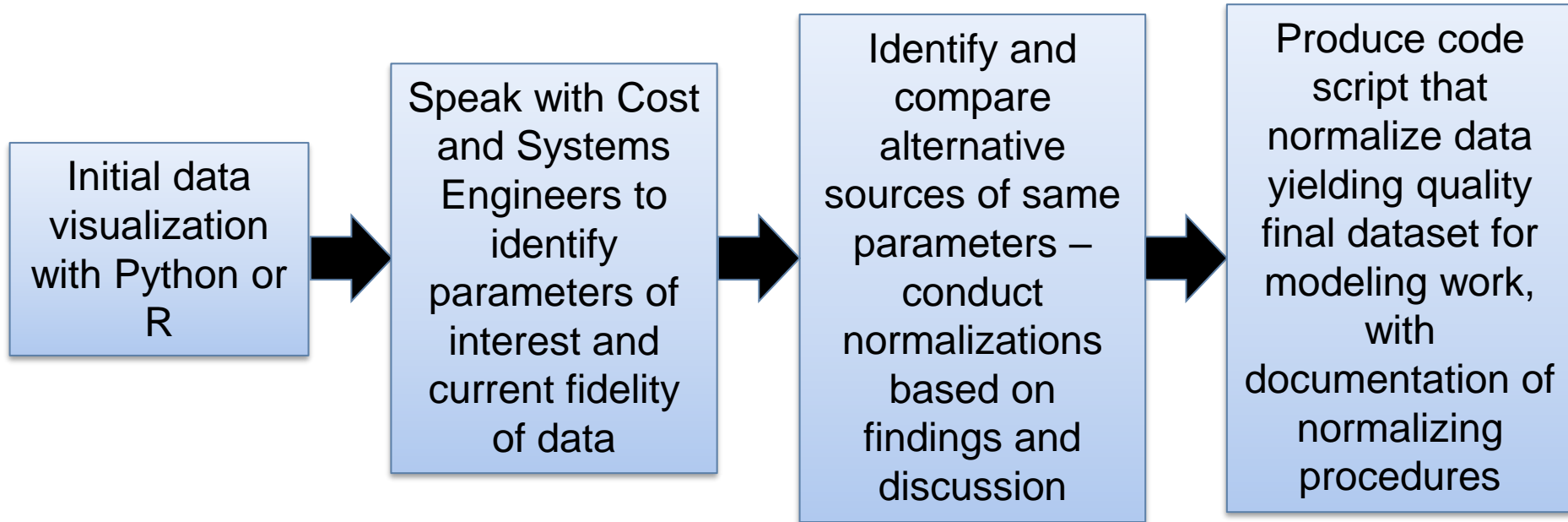
# Relevant Data in Cost Estimation

- Different types of data depending on task
  - FTE (Full-Time Equivalent) – characterizes actual workforce effort
  - Cost – primary variable of interest, includes labor, nonlabor, contributions, and contracts
  - Technical Parameters – by mission, flight element, instrument, include mass, power, data rates
  - Schedule – Stages of lifecycle, viewed in "Phases"
    - Pre-A/A – Study
    - B-D – Development
    - E-F – Operations and Mission Close-Out
- Each comes with unique normalization challenges/constraints
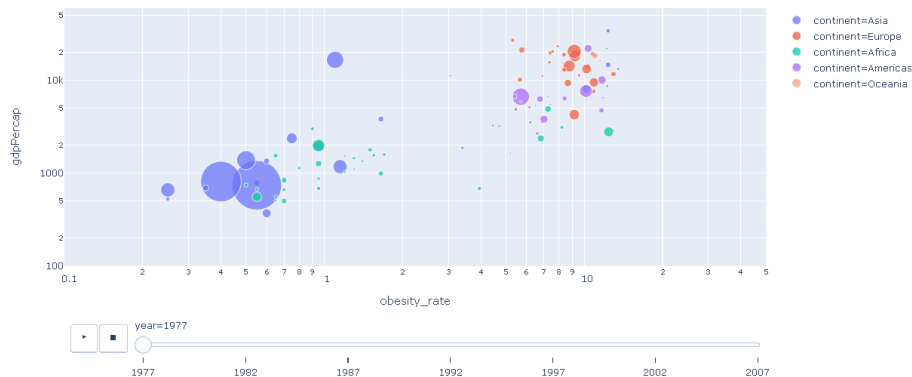
# Data Fidelity

- Pillar upon which all downstream modeling rests upon
- Relevant questions:
  - Where is data coming from?
  - Does the data agree with other sources? If not, why?
  - What (if any) normalizations need to occur to prepare data for analysis?
- Normalization – refers to process of transforming variables
  - Ex. Consistent Fiscal Year for cost, charge accounts mapped to correct WBS
- Iterative procedure to account for as much uncertainty as possible

# General Process

Initial data visualization with Python or R

→

Speak with Cost and Systems Engineers to identify parameters of interest and current fidelity of data

→

Identify and compare alternative sources of same parameters – conduct normalizations based on findings and discussion

→

Produce code script that normalize data yielding quality final dataset for modeling work, with documentation of normalizing procedures

# Processing Tools

- I focused on Python, particularly the Pandas and Plotly packages for data frame manipulation and visualization creation

- Plotly – allows for interactive visuals, elevating users' ability to control their analysis experience

- Both packages scale well with large datasets in terms of visualization and normalization capabilities
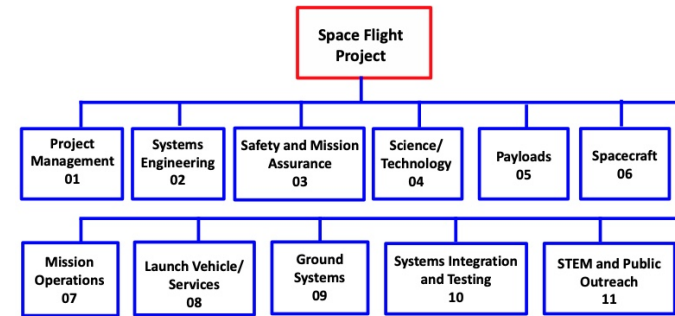
# Examples of Sources

- Oracle Database – Provides granular workforce records, monthly FTE's and cost for various projects

- Work Definition – Project level records and schedule

- CADB – Project level costs by Fiscal year, with information about project phases and WBS mappings

- ZAP – 312C Database intended to serve as warehouse for "gold standard" data based on normalizations conducted

- My work consisted of gathering, merging, and pre-processing data from these sources to generate high-quality datasets for modeling tasks

Reviewed and determined not to contain CUI.

# WBS Normalizations

- Work Breakdown Structure provides common framework for projects to bucket work by major mission components

- JPL uses NASA Standardized WBS, with addition of WBS 12 – Mission Design Navigation (MD Nav)

- Older records may be missing, incomplete, or reflect dated WBS standards which make interpreting through present WBS lens more difficult

- Normalizing Procedures:
  - CADB – Records from 2X that define WBS mappings and account for changes in bookkeeping over time
  - Identify keywords, talk to personnel from projects "grassroots normalization"



Space Flight Project

| Project Management 01 | Systems Engineering 02 | Safety and Mission Assurance 03 | Science/ Technology 04 | Payloads 05 | Spacecraft 06 |

| Mission Operations 07 | Launch Vehicle/ Services 08 | Ground Systems 09 | Systems Integration and Testing 10 | STEM and Public Outreach 11 |

Reviewed and determined not to contain CUI.

# FTE/Cost Normalizations

- FTEs – Monthly or yearly, must determine how to display calculations for interpretable results

- Cost – Standard to inflate to current fiscal year
    - Labor rates do not change 1-1 with inflation, meaning simply inflating may not be sufficient to compare costs temporally

- FTEs are best indicator of raw workforce efforts, but ultimately cost is the variable of interest in most analyses

- Analyzing actuals from formulation lens is difficult:
    - How to back out reserves?
    - Account for slips or other events not typically considered in early life cycle

# ICM Modeling Background

- ICM Model – Institutional Cost Model is a section approved model that serves to quantify cost based on general parameters that guide mission development in formulation
  - Such as planned schedule, anticipated FTE breakouts, and broad technical considerations like the intended target
- Must be signed off by institutional manager – imperative that the calculations and results output represent our best and most accurate efforts to provide high-fidelity cost estimates
- My efforts focused on the model to build WBS 05 and 06 estimates, which correspond to Instrument and Spacecraft management and systems engineering

# ICM Modeling Background

- Speaking to project leads informs us that bookkeeping processes are not perfect during development as indicated in Oracle, so the CADB exists in order to identify and amend bookkeeping errors after they are initially recorded

- This process is critical because the ICM Model must be informed by actual data, and incorrect bookkeeping could lead to over or underestimation of the target WBS

# ICM Modeling Background

- I conducted various normalization procedures on the raw data of over 1.6 million entries, including:
  - Merging CADB mapping of WBS to raw Oracle data to correct known bookkeeping errors
  - Manual WBS normalization based on WBS keywords (such as Payload SE, Flight System SE, etc) to catch unknown mappings
  - Merged mission level characteristics such as Cost Category to allow visualizing groups of missions easily
  - Mapping entries to the proper phase by Project Number and date with CADB and Work Definition information



BIG DATA

# ICM Modeling Example

- The procedures from the previous slide were conducted on both Cost and FTE data to allow further analysis of labor rate calculations and assist in model validation

- I created a Python script to automate as much of this process as possible, merging the unique relational data frames and generating an output CSV that is ready for immediate analysis

- Results stemming from validation with data I normalized serve as key motivator to update model in order to best reflect how JPL conducts business in higher level of detail than previously possible

```python
# ENGINEERING --> ENGINEERING (updated)
raw_data_df.loc[(raw_data_df["Expenditure Type"] == "LJ-ENGINEERING-ASSOC"), "Expenditure Type New"] = "LJ-ENGINEERING-1-2"
raw_data_df.loc[(raw_data_df["Expenditure Type"] == "LJ-ENGINEERING-PRIN"), "Expenditure Type New"] = "LJ-ENGINEERING-5-6"
raw_data_df.loc[(raw_data_df["Expenditure Type"] == "LJ-ENGINEERING-SENIOR"), "Expenditure Type New"] = "LJ-ENGINEERING-4"
raw_data_df.loc[(raw_data_df["Expenditure Type"] == "LJ-ENGINEERING-STAFF"), "Expenditure Type New"] = "LJ-ENGINEERING-3"

raw_data_df.loc[(raw_data_df["OBIEE Wbs Number"] == "99"), "Normalized Wbs Number"] = "01.01"
raw_data_df.loc[(raw_data_df["OBIEE Wbs Number"] == "99"), "Normalized Wbs Name"] = "Proj Mgmt"
```

# Schedule Normalizations

- Schedule and Phases are projected at different checkpoints, such as PDR, SIR, and EoM

- Differences in schedules based on various factors – launch windows, internal/external slips, funding availability

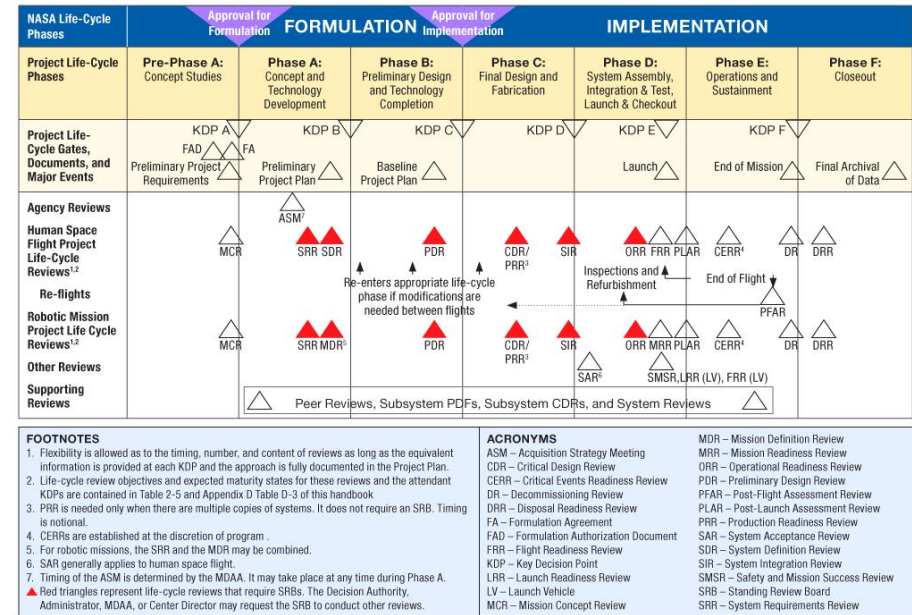- Which form of schedule will best inform models?



FIGURE 3.0-1 NASA Space Flight Project Life Cycle from NPR 7120.5E

# Schedule Normalization Example

- ICM model uses phase lengths to build mission profile to better understand how time drives cost
- Phases B, C, and D (development) correspond to the highest degree of work on Payload and Spacecraft components
- I aided in efforts to document these schedule discrepancies from plans and actuals and identify rules of thumb to breakout sub-schedule durations that are used in mission planning
- Identified source of schedule slips (internal or external) and removed months pertaining the delays caused by non-JPL factors
- Delivered processed and documented schedule data to incorporate in ZAP, the 312C visualization and data warehouse

# Technical Data Normalizations

- May wish to utilize technical parameters from flown missions and studies to deduce cost estimating relationships (CERs)

- Parameters may be given on instrument, flight element, or mission basis

- Studies can include multiple design scenarios with different values – must choose what to accept for analysis
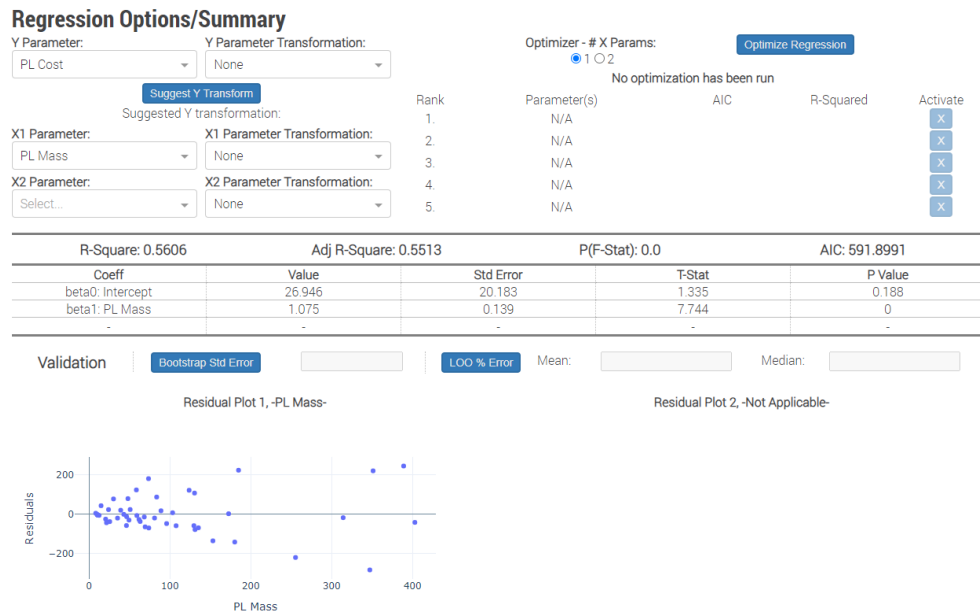
jpl.nasa.gov

# Technical Data Normalizations Example

- MCDB – Mission and Cost Database
- MCDB task aims to unify actual flown data and study to create comprehensive "gold standard" technical data for missions, flight elements, and instruments
- I read and documented study reports and workbooks to intimately understand mission objectives and identify parameters of interest
- Provided data to personnel responsible for updating database



| Mission | Mass (kg) | Power (kW) |
|---------|-----------|------------|
| A | 500 | 26 |
| B | 2000 | 45 |
| C | 1500 | 30 |

# Technical Data Normalizations Example

- Technical data is used to model relationships between physical parameters and cost, such as in the A-Team Regression (ART) tool

- Quantifying these relationships with regression works well within our generally low-data paradigm (more complex models may overfit)

- Allows the user to explore strength of different CERs, such as mass vs cost, data rate vs cost

# Workforce Analysis Visualization and Exploration Tool

- Otherwise known as WAVE, part of 312C Suite of tools
  - https://costing.jpl.nasa.gov/dash/wave/
- I created using Dash – a Pythonic framework for creating powerful data science oriented web apps
- Integrates elements of WBS, FTE, Cost, and Schedule normalizations I conducted in backend data with simple filters and user interface on frontend
  - This data is local right now, with plans to integrate backend directly to ZAP
- Use case: quickly identify most prevalent expenditure categories in WBS 05 and 06 to inform ICM model
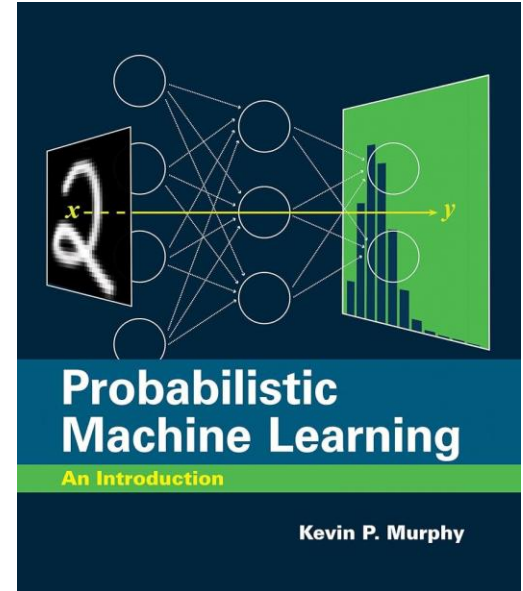
jpl.nasa.gov

# Summary

- Led data normalization procedures for data used in various critical tasks including ICM Model updates and MCDB data collection
- Expanded the capabilities of 312C web presence with the WAVE tool
- Created scripts and templates in Python and R to provide standardized and traceable normalization procedures on various data, namely cost, FTE, schedule, and technical parameters
- Learned about fundamental JPL engineering and business practices and integrated lessons from cost and systems engineers into data normalization process

# Acknowledgements

- Michael Saing and Anto Kolanjian
  - Served as my primary mentors for the duration of my internship
  - Provided invaluable domain knowledge, guidance on tasks, and freedom to explore projects of interest
- Patrick Bjornstad and Patricia Gallagher
  - Collaborated with on a number of projects including web development, ICM, DAG, MCDB

# Additional Acknowledgements

- Weekly Math/Stats reading group – engaging discussions on statistical/ML techniques with eye towards applications in our daily work
- Weekly Tag-Ups – presented opportunity to engage in discussions surrounding other activities performed by group members and understand broader scope of group impact



**Probabilistic Machine Learning**
**An Introduction**

**Kevin P. Murphy**

# Thank You!

Reviewed and determined not to contain CUI.

jpl.nasa.gov

jpl.nasa.gov