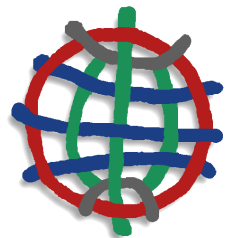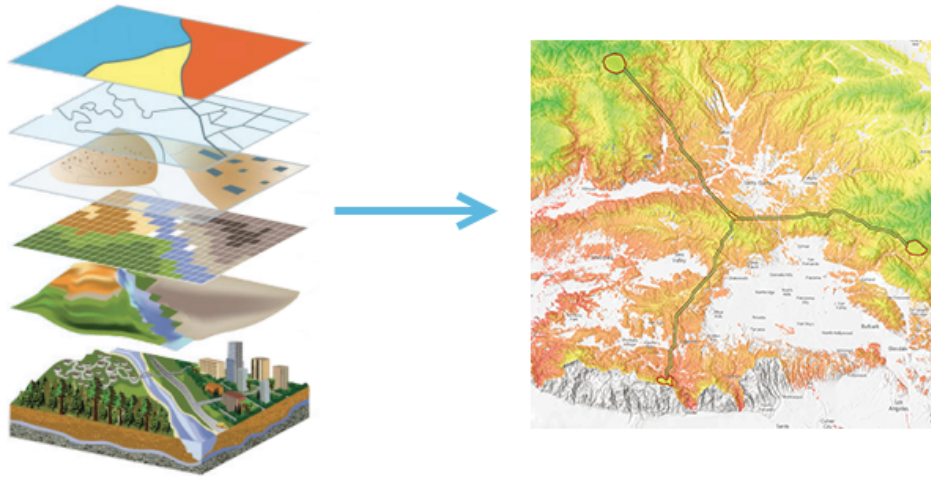# TEST-DRIVEN DATA ANALYSIS

*Do you believe your analytical results?*

Thanks to:
Nicholas Radcliffe
http://tdda.info
njr@StochasticSolutions.com
Dept of Mathematics, University of Edinburgh

OLLIVIER & CO

# GeoPlanner's Suitability Modeler is now part of Web AppBuilder

by Rob Stauder on June 29, 2017
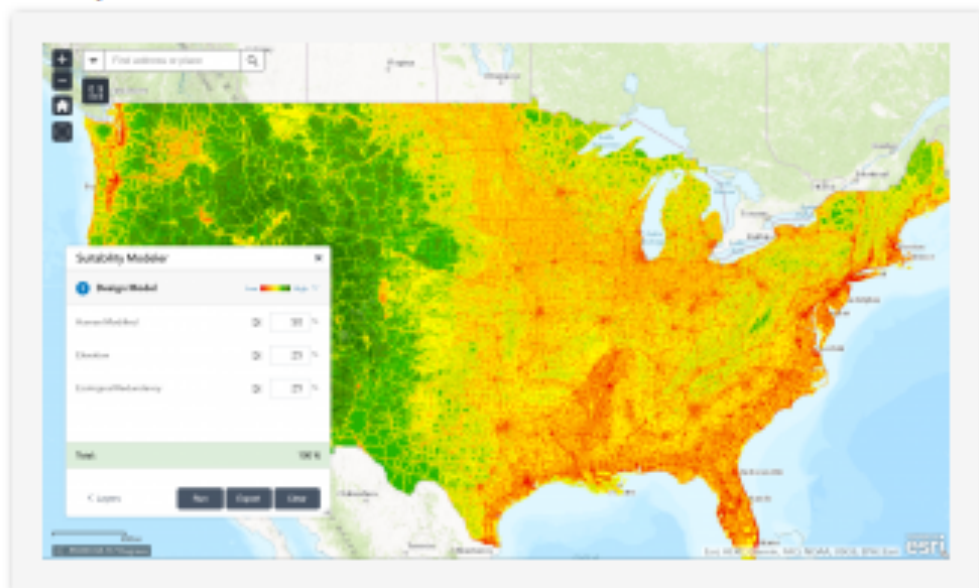
Imagine if, in a few clicks, you could answer multiple-factor spatial questions like *Where are the areas on low angled slopes, in shrubby vegetation and are far from roads*? What if you could do that and emphasize the importance of one of those factors over another?
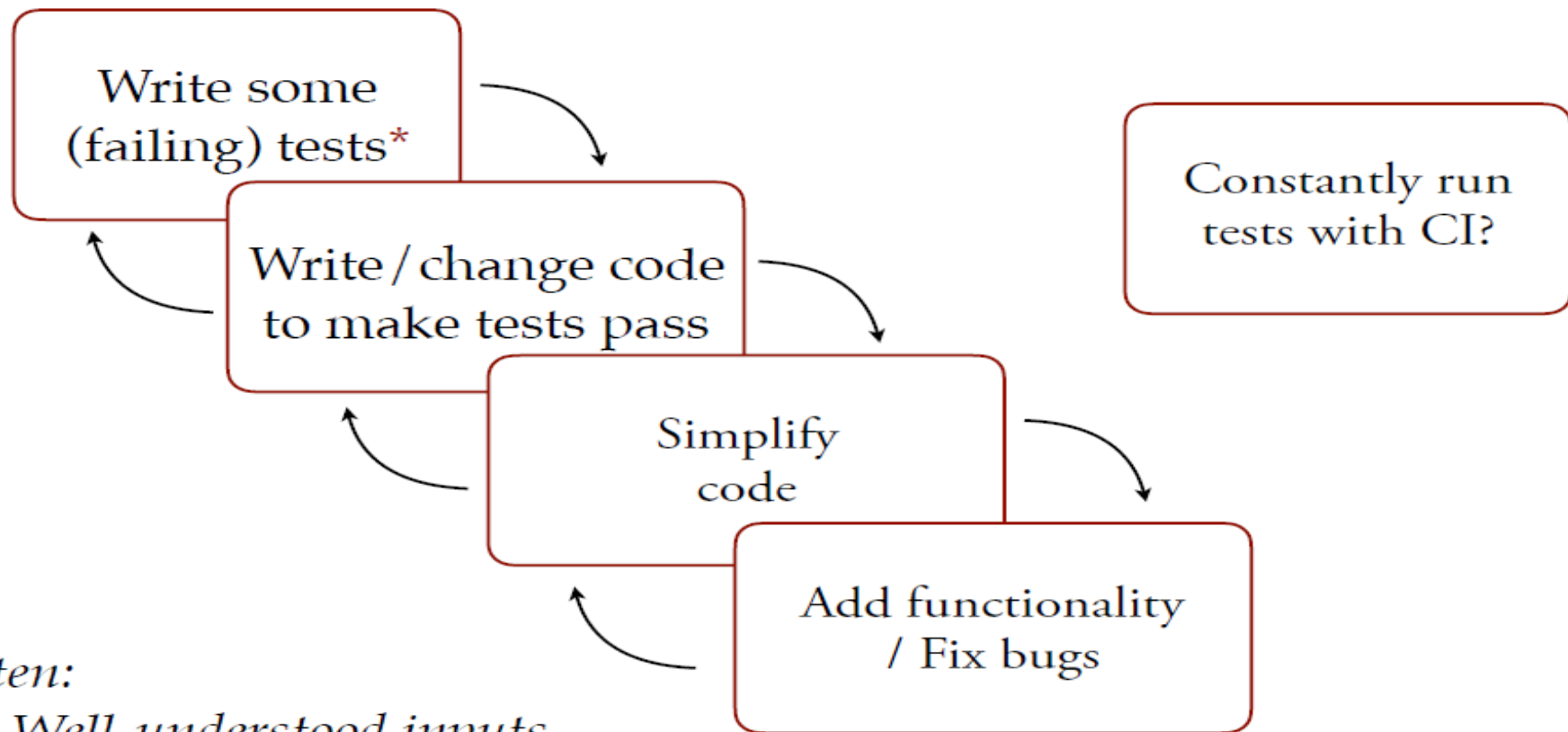
You would be the hero of your workplace!

# The Big Idea

*Transfer the ideas of
test-driven development
from software
development
to data analysis*

# SOFTWARE DEVELOPMENT (WITH TDD)

Write some (failing) tests*

Write/change code to make tests pass

Simplify code

Add functionality / Fix bugs
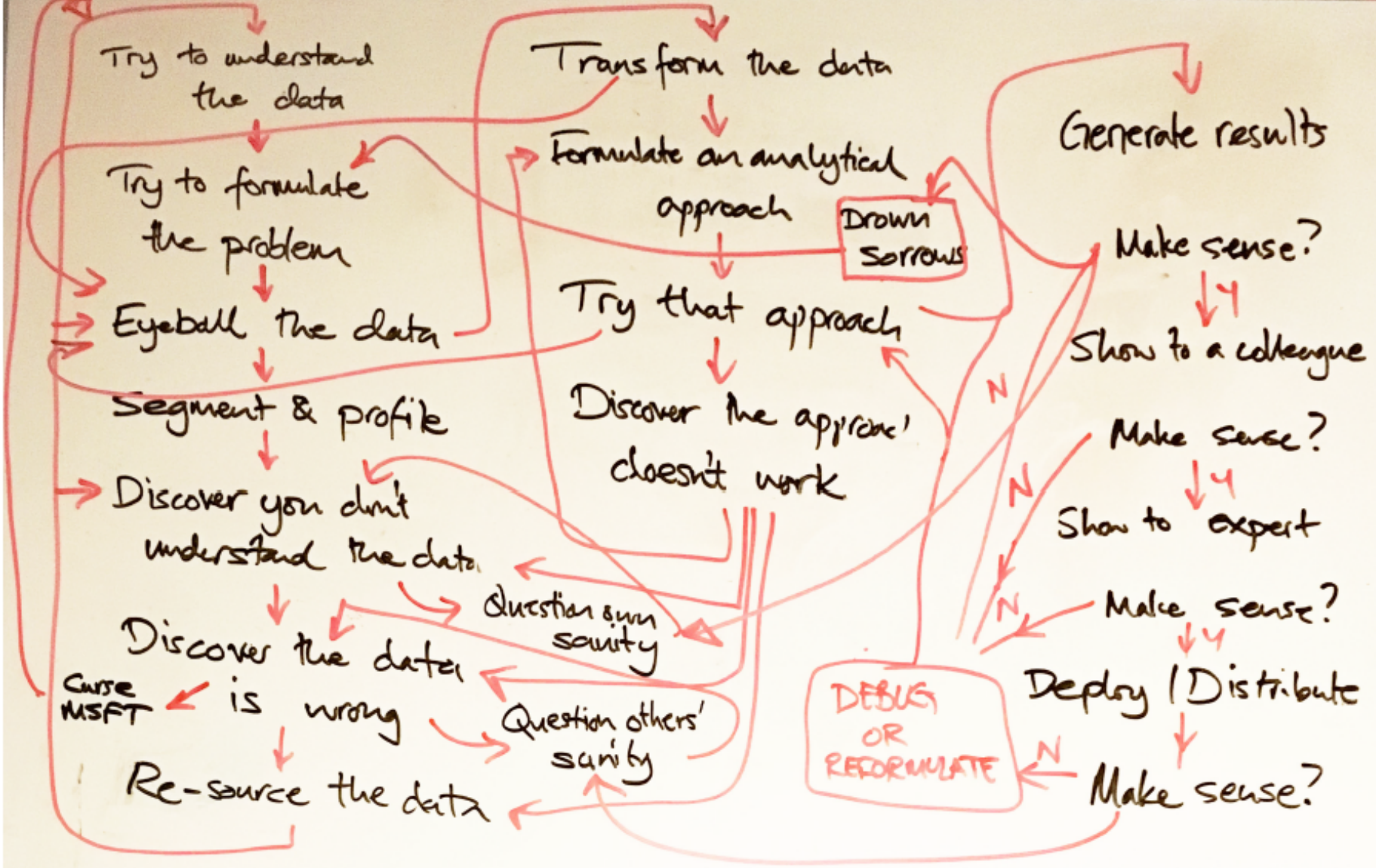
Constantly run tests with CI?

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

* While mocking almost everything
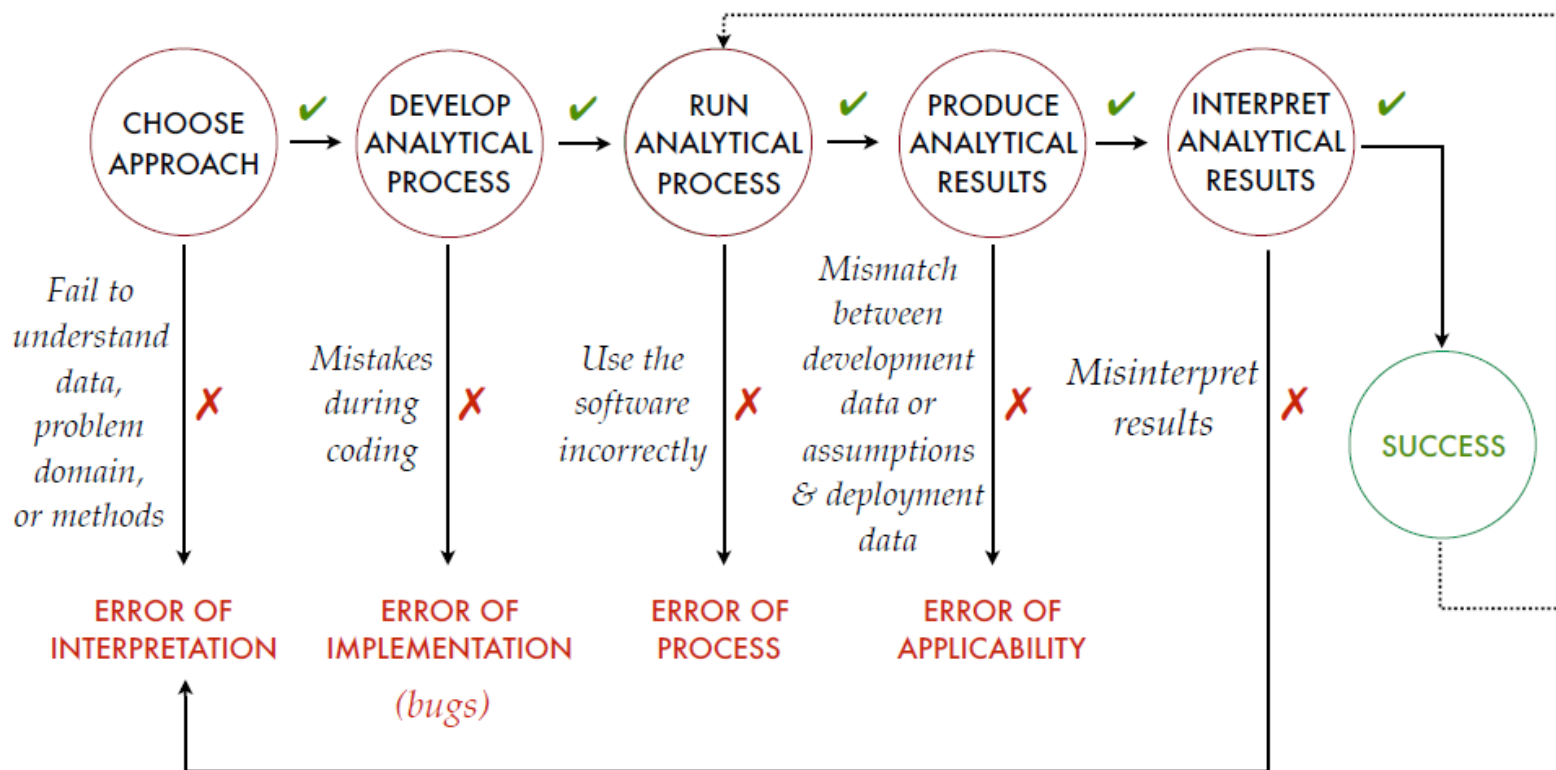
How models are really developed

Try to understand the data

Try to formulate the problem

Eyeball the data

Segment & profile

Discover you don't understand the data

Discover the data is wrong

Curse MSFT

Re-source the data

Transform the data

Formulate an analytical approach

Drown Sorrows

Try that approach

Discover the approach doesn't work

Question own sanity

Question others' sanity

DEBUG OR REFORMULATE

Generate results

Make sense?

Show to a colleague

Make sense?

Show to expert

Make sense?

Deploy / Distribute
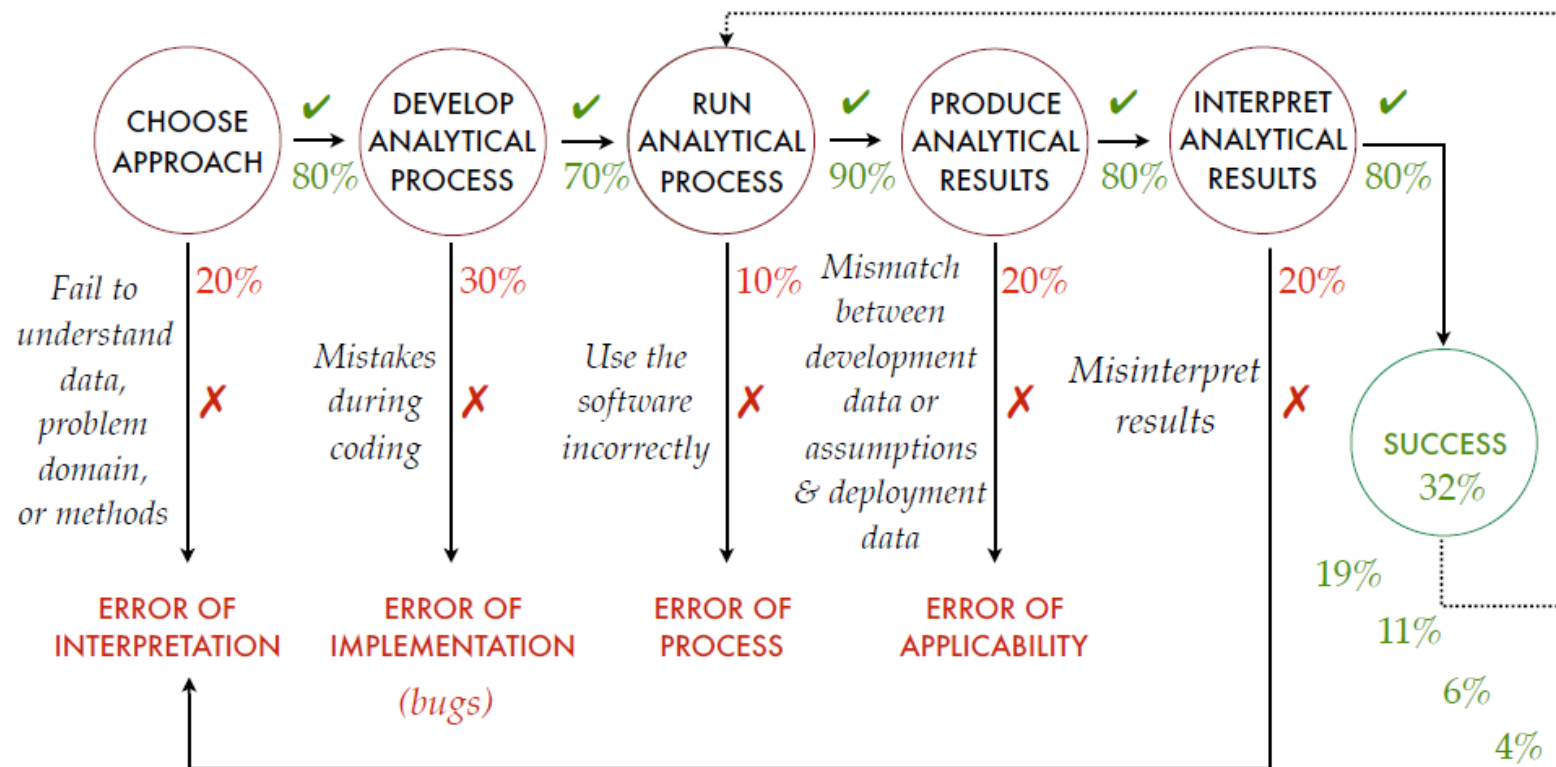
Make sense?

## DEVELOPMENT PHASE

*Using sample/initial datasets & inputs to develop the process*

## OPERATIONAL PHASE

*Using the process with other datasets and inputs, possibly having different characteristics*

**CHOOSE APPROACH** ✔ **DEVELOP ANALYTICAL PROCESS** ✔ **RUN ANALYTICAL PROCESS** ✔ **PRODUCE ANALYTICAL RESULTS** ✔ **INTERPRET ANALYTICAL RESULTS** ✔

*Fail to understand data, problem domain, or methods* ✗

*Mistakes during coding* ✗

*Use the software incorrectly* ✗

*Mismatch between development data or assumptions & deployment data* ✗

*Misinterpret results* ✗

**SUCCESS**

**ERROR OF INTERPRETATION**

**ERROR OF IMPLEMENTATION**

*(bugs)*

**ERROR OF PROCESS**

**ERROR OF APPLICABILITY**

*If you buy into this model, it's sobering to attach probability estimates to each transition and calculate the probability of success after a few runs . . .*

# TDDA: LEVEL ZERO

**INPUTS**

DATA
& PARAMETERS

➡️ **ANALYTICAL PROCESS**

➡️ **OUTPUTS**

DATASETS, NUMBERS,
GRAPHS, MODELS,
DECISIONS ETC.

*Record inputs*

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

*Record ("reference") outputs*

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

# Welcome to Kaggle Competitions

Challenge yourself with real-world machine learning problems

### New to Data Science?

Get started with a tutorial on our most popular competition for beginners, Titanic: Machine Learning from Disaster.

### Build a Model

Get the data & use whatever tools or methods you prefer to make predictions.

### Make a Submission

Upload your prediction file for real-time scoring & a spot on the leaderboard.

Submit »

**Learn more**

# tdda level1:
# CONSTRAINTS

*Look before you leap!*

Checking the data conforms to your assumptions before you start.

- Not just the obvious input, but also intermediate and output sets

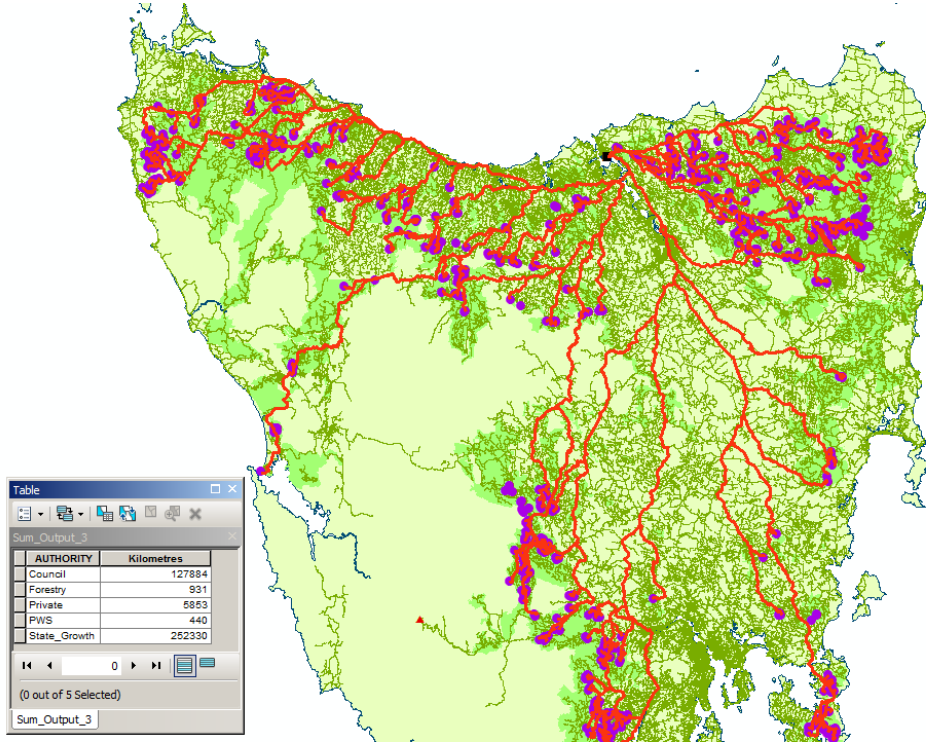This is tedious to generalise so there are tools to help...

# Tools to automate L1 constraint tests

| Tool | Read input | Rules | Action | Flexibility, Complexity |
|---|---|---|---|---|
| ArcGIS | Input Dialog | Existence, schema | Block from starting | Low |
| RDBMS | Table schema | Field constraints, triggers | Rejection | Med |
| TDDA python module | Pandas framework | Regular expression generator | Report | Med |
| FME | Attribute Validator transformer | Choose from built-in rules, custom tests | Report, repair or filter | High |

# EXAMPLE CONSTRAINTS

| SINGLE FIELD CONSTRAINTS | DATASET CONSTRAINTS |
|---|---|
| Age ≤ 150 | The dataset must contain field **CID** |
| type(Age) = int | Number of records must be  118 |
| CID ≠ NULL | One field should be tagged  **O** |
| CID unique | **Date** should be sorted ascending |
| len(CardNumber) = 16 | **MULTI-FIELD CONSTRAINTS** |
| Base in {"C", "G", "A", T"} | StartDate ≤ EndDate |
| Vote ≠ "Trump" | AlmostEqual(F, m * a, 6) |
| StartDate < tomorrow() | sum(Favourite*) = 1 |
| v < 2.97e10 | minVal ≤ medianVal ≤ maxVal |
| Height ~ N(1.8, 0.2) | V ≤ H * w * d |

# TEST-DRIVEN DATA ANALYSIS