# Small Statistics Big Data Curriculum

C.Foster
05 Oct 2018

# Historic

- Database access controlled and complex (Instance / Schema / Table)
  Software Engineer to organize and query data.

- Analysis in proprietary code (FEA/CFD Solvers)
  Analysis Engineer to setup and run solvers

- Statistical Analysis in special programs (JMP/Minitab)
  Statistician run analysis

- Publish results (HTML/Java)
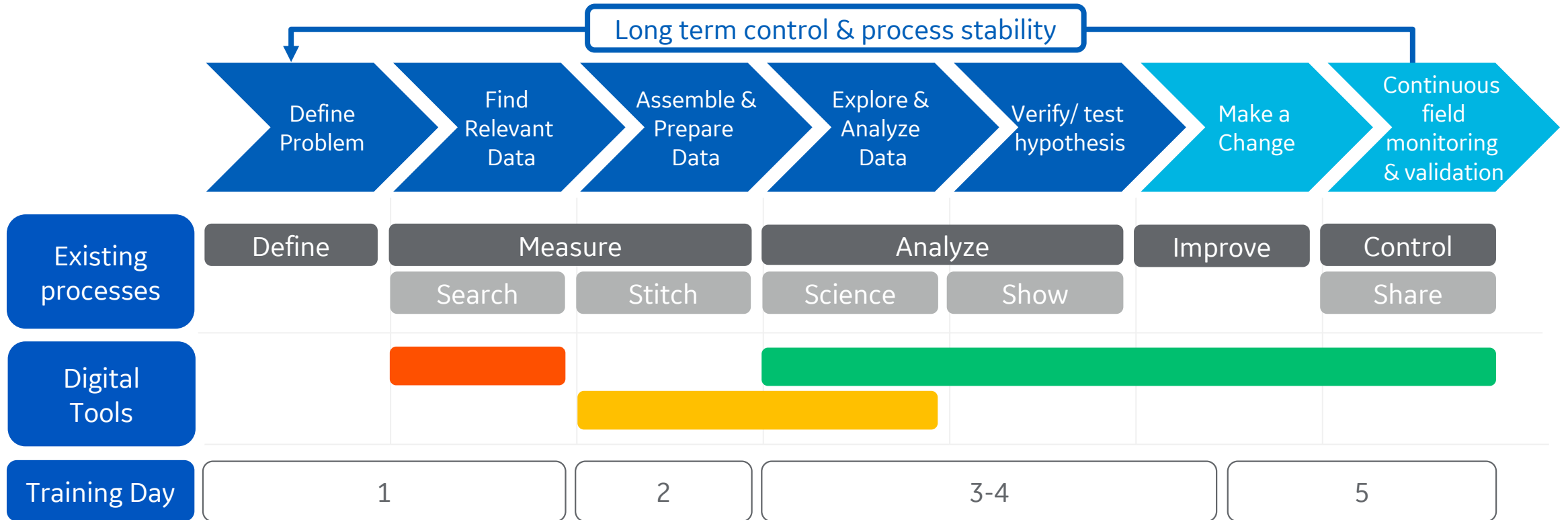  Programmer to create report or dashboard

**No common process or link**

# Process & Tools

- Get Data
- Analyze and model
- Publish results

Target Audience: Engineers / Technicians / Scientists who are not primarily digital.

| | Long term control & process stability | | | | | |
|---|---|---|---|---|---|---|
| Define Problem | Find Relevant Data | Assemble & Prepare Data | Explore & Analyze Data | Verify/ test hypothesis | Make a Change | Continuous field monitoring & validation |

| | Define | Measure | | Analyze | | Improve | Control |
|---|---|---|---|---|---|---|---|
| **Existing processes** | Define | Measure | | Analyze | | Improve | Control |
| | | Search | Stitch | Science | Show | | Share |
| **Digital Tools** | | | | | | | |
| **Training Day** | 1 | | 2 | 3-4 | | | 5 |

Query → Norm / Imputing → Aggregate →

### Data Engineering Process

| Query | Norm / Imputing | Aggregate |
|---|---|---|
| First N - Sampling | Z-Score<br>0 Replacement | Mean |

### Applied Statistics Process

| Stratified Sampling | Min/Max<br>Robust Scaling<br>Imputation | Median<br>IQR |
|---|---|---|

→ Charts → Tests → Models →

### Data Visualization Process

| Charts | Tests | Models |
|---|---|---|
| Scatter / Bar /<br>Line / Combined | Color by… / Line by… /<br>Shape by… | Trendlines |

### Applied Statistics Process

| Quantile / Box Plot | Z-t Test / $\chi2$ / FDR | Linear and Logistic Regression<br>Classification and Regression Trees |
|---|---|---|

# Decision → Implementation → Control / Capability

## Programming Process

| Threshold | Recode and Test initial | Run Reports |

## Applied Statistics Process

| Bootstrapped threshold | DOE Based Validation | Anomaly Detection / Residual Analysis |

**Integrate the Statistical / Data Science methods and the Tools and Process Activities**

**Define the goal in terms of a measurable outcome**

# Four Examples

**Project Charter**

| Project Title: | | |
| --- | --- | --- |
| Problem Statement: | Goal Statement: | VOC: |
| Project Team:<br>Leader:<br>Team member1:<br>Team member2:<br>Team member3: | Project Information:<br>Project start:        Project end:<br>Project approach:<br>Project scope: | Key Metrics:<br>: |
| | | Resources:<br>: |
| Milestones: | | |
| Signatures: _____  _____  _____  _____ | | |

1. N≠All

2. N is big

3. Act

4. Feedback

Models / Ownership

Something is missing

Statistical assumptions breakdown.

What is wrong with the measurements? Are we solving the right problem?

How will the system react and how will it effect the analytic?
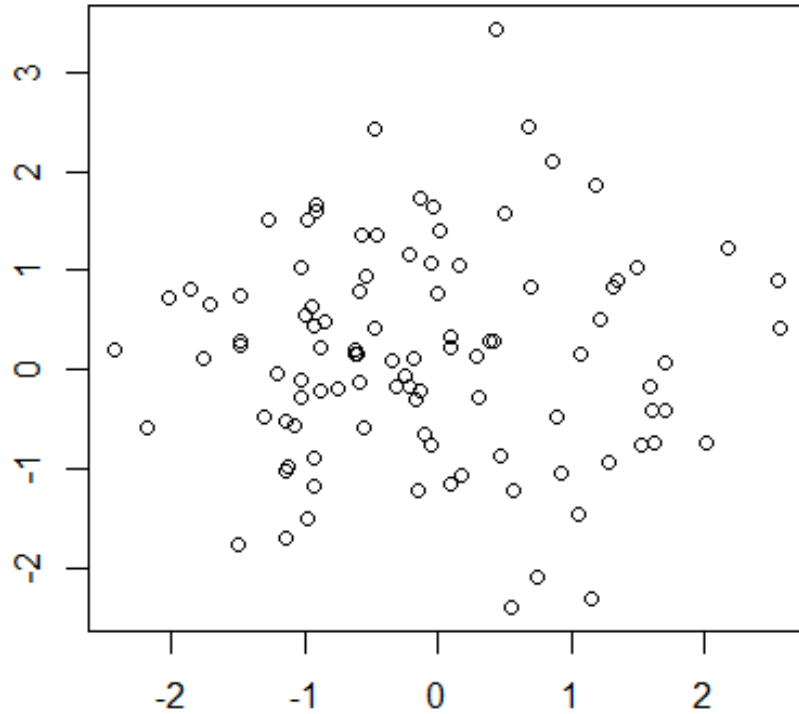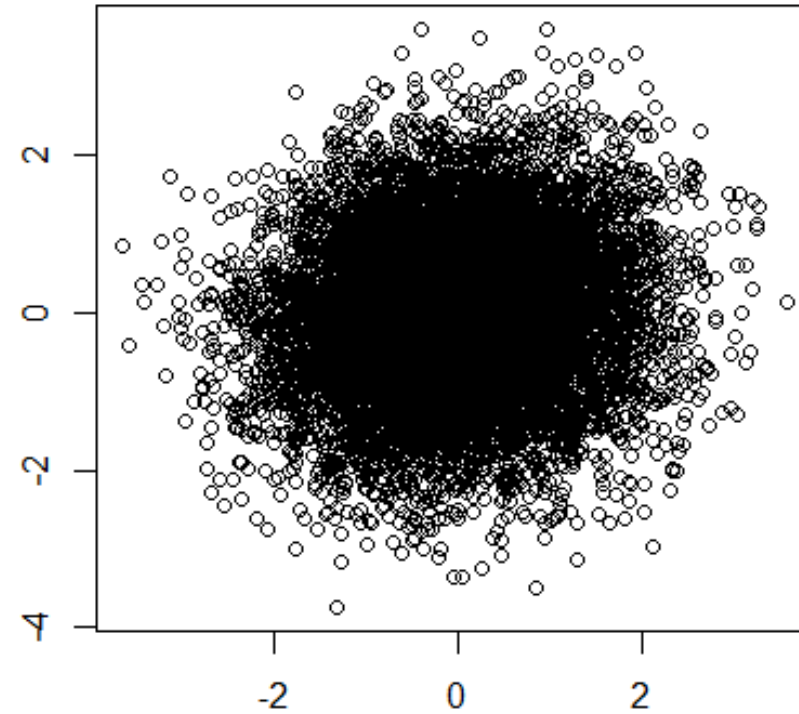
N≠All

N is big

R = .1

100 points p=0.936

Not Significant

10000pts p<2e-16

VERY Significant

Your task is to turn over as few cards as possible to verify whether the following statement is true.

*Every card with a vowel on one side has an even number on the other.*



Which order would you turn over the cards?

1. AB23   4. A2B3
2. AB32   5. A3B2
3. A23B   6. A32B

# Chasing Noise...



## All data are noisy

| Students Have a Poor Landing | Engine has high fuel consumption one month |
|---|---|
| ↓ | |
| Yell at student | ...Tell Operator |
| ↓ | |
| Next Landing Improves | Fuel Economy Improves |

Feedback

Interactions with the System



Calculated Model Output

Time

X is a surrogate for Analytic based on correlations

Customer tunes X and changes correlations to Y

If customer is charged more when they are here

They try to run here
Even if it does not change risk / cost

Control is critical
Not predicting weather

# Methods and Big Data



Medium ANN

Small ANN

Logistic / Random Forest /
Support Vector Machine

Likelihood

Performance

Data

# Gaussian Process
# Random Forests

General use model

Pretty well most of the time

Bayesian Calibration of Computer Models
https://rss.onlinelibrary.wiley.com/doi/10.1111/1467-9868.00294

Fails obviously

Confidence Intervals for Random Forests: The Jackknife and the Infinitesimal Jackknife
http://jmlr.org/papers/volume15/wager14a/wager14a.pdf

Includes Uncertainty

The Analyst:

Owns the proposal

Owns knowing the data - and quality

Gets the data from the lake

Stitches the sources

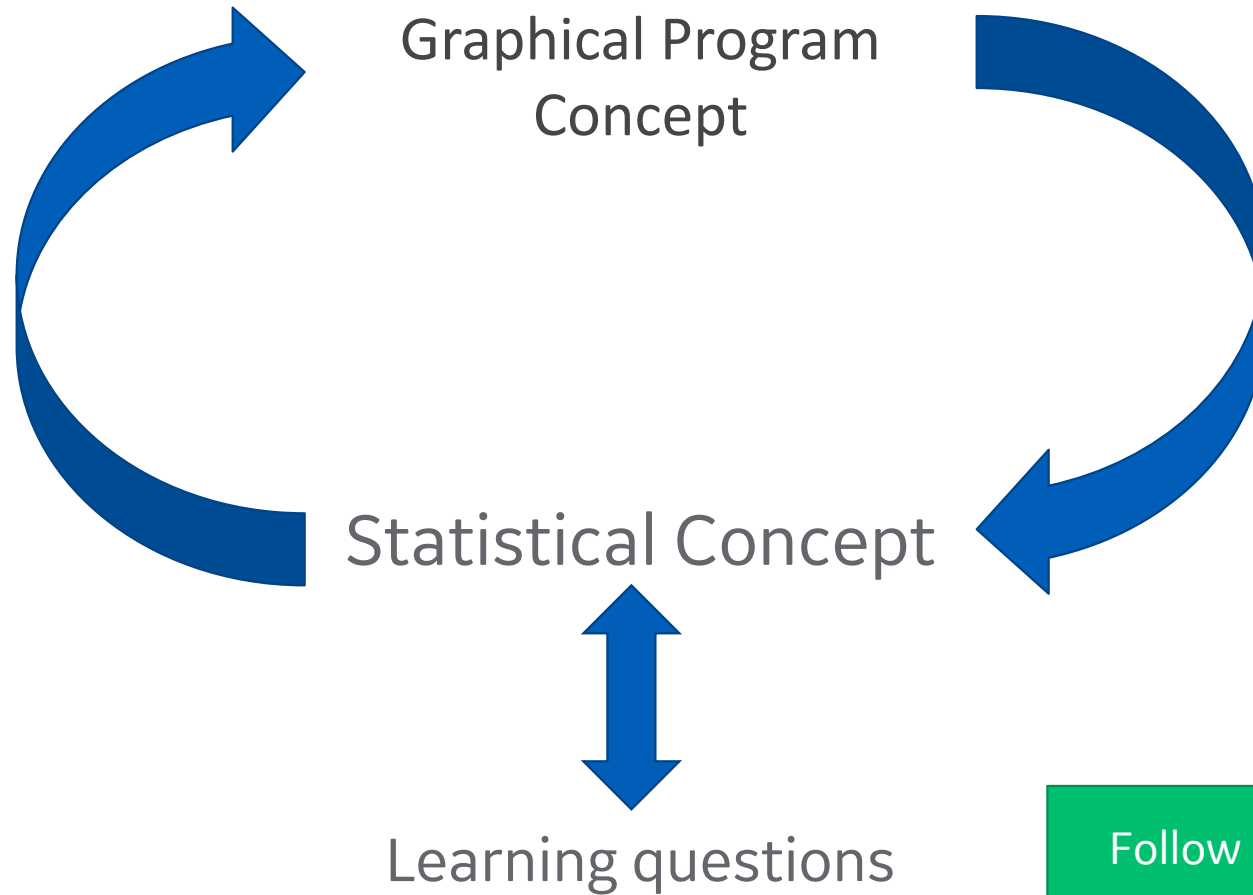Creates the visualization

**Training Process**

Integrate objectives

Graphical Program Concept

Statistical Concept

Learning questions

Follow an available text

https://www.openintro.org/stat/textbook.php

# P-Value Sample

A poll by the National Sleep Foundation found that students on <u>average</u> sleep 7 hours per night.

A sample of 169 students sleep for one night had an <u>average</u> of 6.88 hours and a standard deviation of 0.94 hours.

Assuming that this is a representative random sample, is there sufficient evidence to reject the null hypothesis that students on <u>average</u> sleep 7 hours per night?

Standard Error
= sd/sqrt(samp)

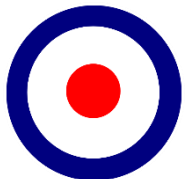What is the p-value for this hypothesis test?

What caveats would you add?

Test
sample mean: 6.88
samp standard dev: 0.94
number samps: 169

test_mean: 7

Z-Score: -1.660

Lower Prob: 4.850005 % (p-value)
Upper Prob: 95.149995 %

https://www.openintro.org/stat/textbook.php

# Student Case Study – West Nile Virus (with parallels to fleet management)
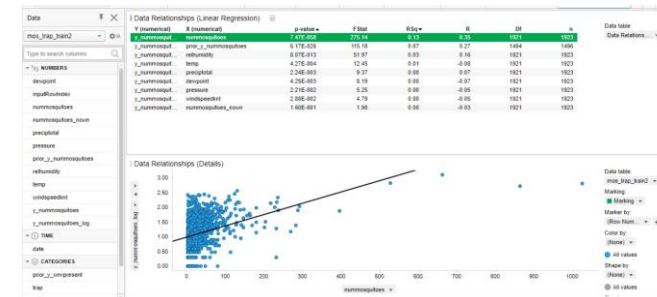
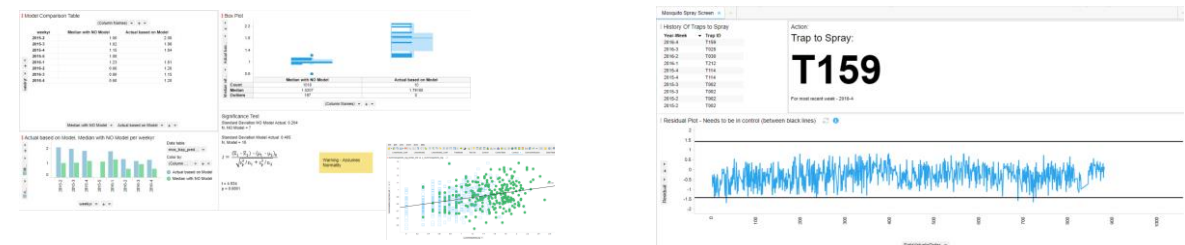Setup project requirements and deliverables

Find and manipulate data

Visualize and build models

Validate and Publish Dashboard

# Complete process with tools

- Engineers and users are capable of entire analytic process
- Basic statistics knowledge and capability
- Integrated with current improvement process

| | Long term control & process stability | | | | | |
|---|---|---|---|---|---|---|
| Define Problem | Find Relevant Data | Assemble & Prepare Data | Explore & Analyze Data | Verify/ test hypothesis | Make a Change | Continuous field monitoring & validation |

| Existing processes | Define | Measure | | Analyze | | Improve | Control |
|---|---|---|---|---|---|---|---|
| | | Search | Stitch | Science | Show | | Share |

| Digital Tools | | | | | | | |
|---|---|---|---|---|---|---|---|

| Training Day | 1 | 2 | 3-4 | 5 |
|---|---|---|---|---|