

Collecting clean and usable data for better reliability modelling

Data Scientist & founder
carles@reliabledynamics.com



carlescg



@carles_



RELIABLE DYNAMICS

- Who we are?
- Modeling 101 - Quick introduction
- So, what can we do?
- Six steps to improve any modeling
- Demo/Case
- Recap

Who we are?

We believe that **reliability engineering** practitioners should focus on delivering value and developing critical thinking in **O&M** environments. **Software should not be a barrier** to apply this techniques in any organization. We take care to eliminate that barrier.

We believe that the competitive advantage is to have **a wide range of apps**, that increase the reliability inside any organization, regardless of **size and budget**.

We develop **basic applications** with **modern technologies**. Trough standard applications or custom solutions. Our products include apps from statistical calculations, reporting, to warranty predictions and predictive maintenance.



RELIABLE DYNAMICS

- Reliability modeling
 - Life data analysis (statistical)
 - Vibration (deterministic)
 - Oil sampling

- Algorithms

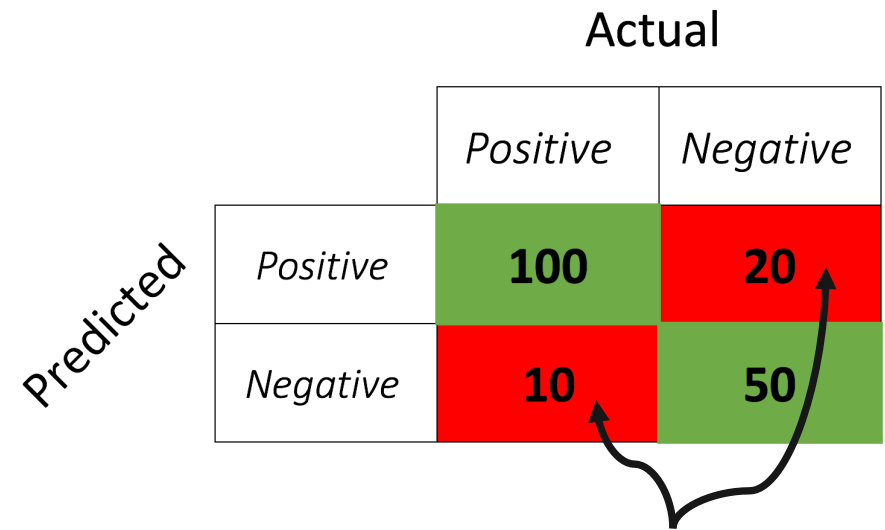
- Machine learning
- Deep learning
- Reinforced learning

- 
- Supervised
 - Unsupervised
 - Semi-supervised

- 
- Classification
 - Regression

Output = function (Input)

- Data driven services require integrated workflows
- Garbage in / garbage out
- No free lunch theorem
- Confusion Matrix

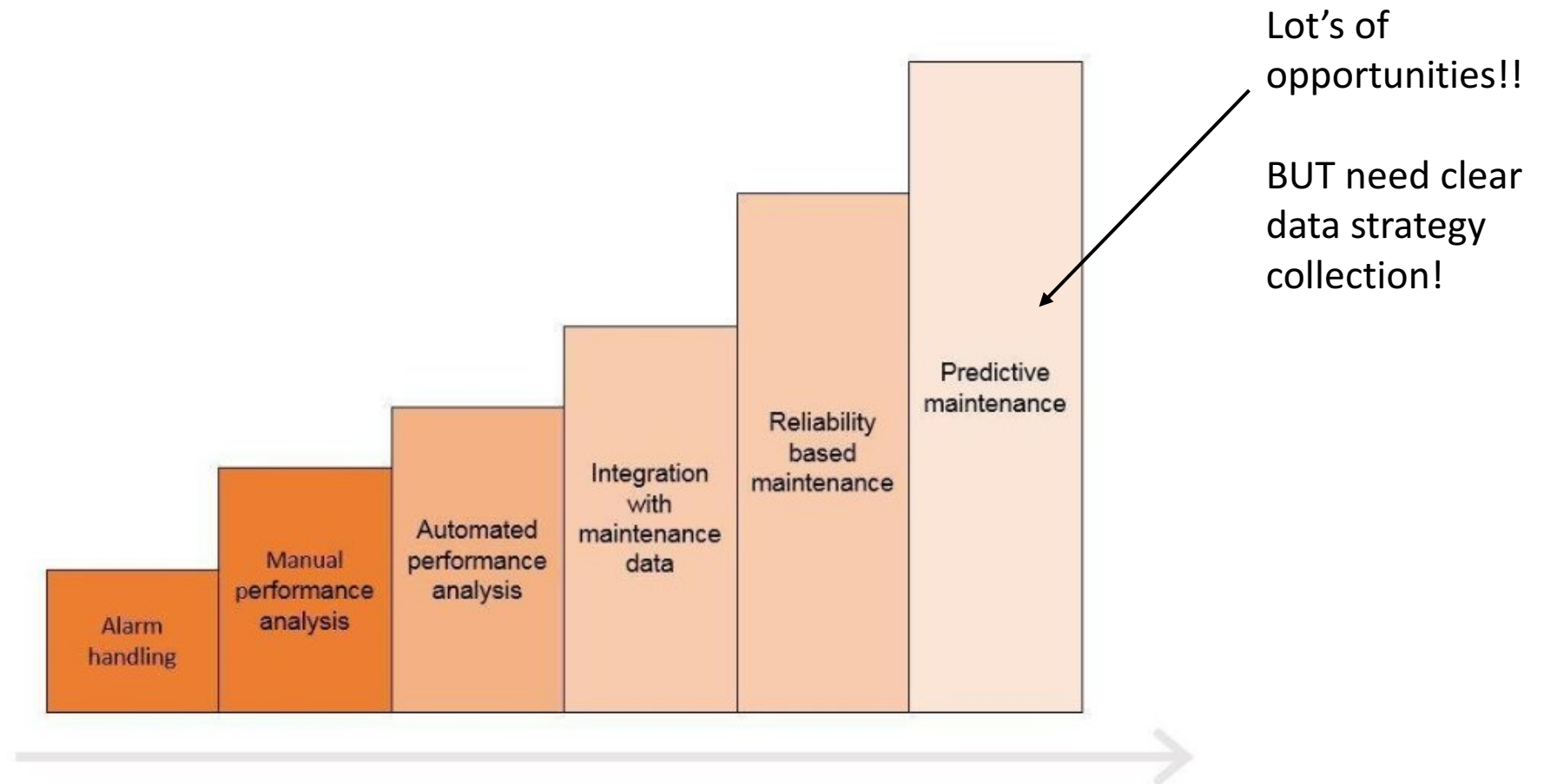


A confusion matrix diagram. The columns are labeled 'Actual' with sub-labels 'Positive' and 'Negative'. The rows are labeled 'Predicted' with sub-labels 'Positive' and 'Negative'. The matrix cells contain counts: 100 (green), 20 (red), 10 (red), and 50 (green). Two curved arrows point from the bottom of the matrix to the red cells (20 and 10), highlighting miss-detections.

		Actual	
		Positive	Negative
Predicted	Positive	100	20
	Negative	10	50

Keep track of miss-detection!

Modeling in Operation & Maintenance



So what can we do??

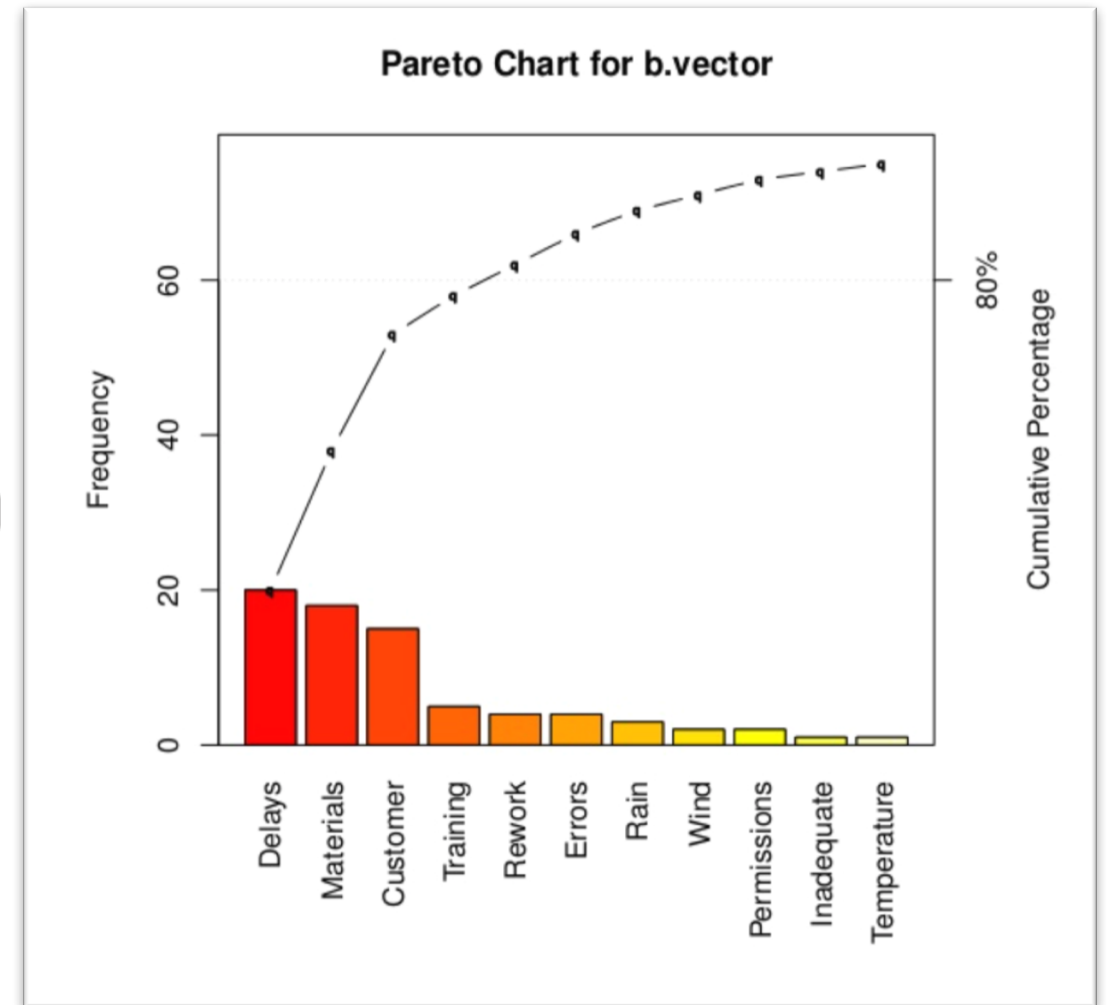
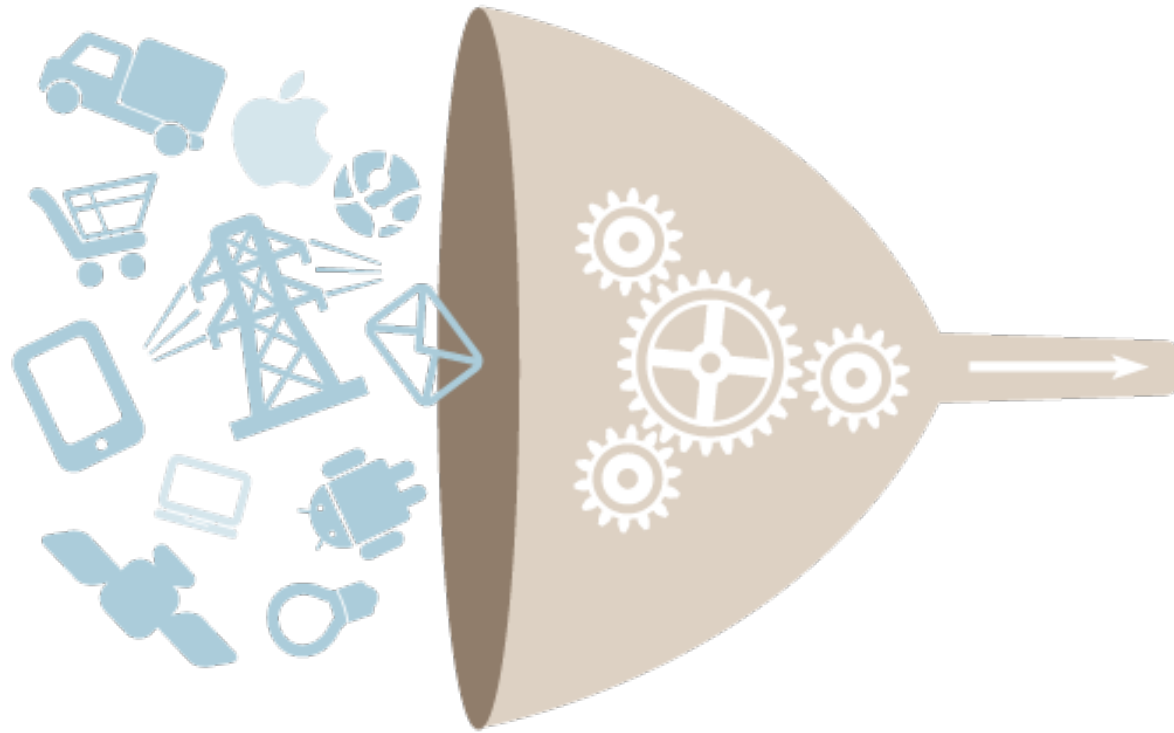


"By sharing their algorithms, Facebook and Google are merely sharing the recipe. Someone has to provide the eggs and flour and provide the baking facilities (which in Google and [Facebook's case](#) are vast data-computation facilities, often located near hydroelectric power stations for cheaper electricity)."

[The Guardian.](#)

"This is probably why Facebook and Google have so freely shared their methodologies: they know that the real value in their companies is the vast quantities of data they retain about each one of us."

[The Guardian.](#)



Six steps to improve your modeling



Photo by Clem Onojeghuo

What is the objective of the modeling?



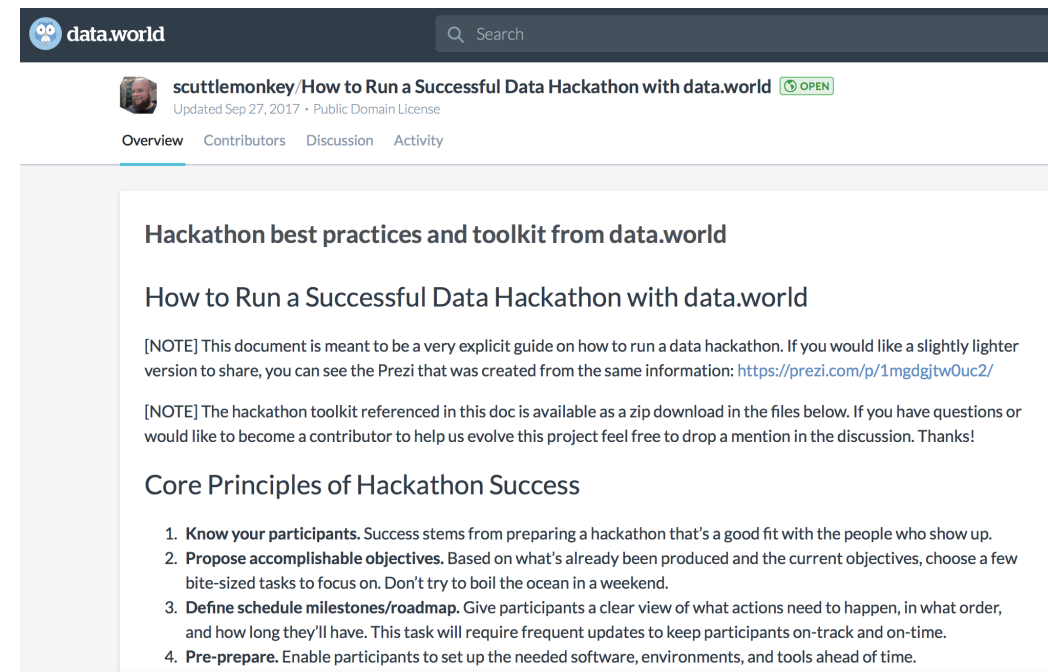
What is the ideal data set?



Step 2 – Ideal data set

Cluster	Parameter
Time	Age of component
	Time
Stress	Full load hours
	Shear modulus
	Deviations
Environment	AMB temperature
	Wind speed
	Wave height
	Wake effect
Maintenance	Crane/non-crane components
	Rate/degree/effort of maintenance
	Human factor

Guide [Data sharing](#) by Jeff Leek



The screenshot shows a data.world project page. At the top, the data.world logo and a search bar are visible. Below the header, the project title "scuttlemonkey/How to Run a Successful Data Hackathon with data.world" is displayed, along with a green "OPEN" button. The page has tabs for "Overview", "Contributors", "Discussion", and "Activity". The main content area features the title "Hackathon best practices and toolkit from data.world" and a subtitle "How to Run a Successful Data Hackathon with data.world". There are two notes: one stating the document is a guide on how to run a data hackathon with a link to a Prezi, and another stating that a toolkit is available as a zip download. Below the notes, the section "Core Principles of Hackathon Success" is listed with five numbered points: 1. Know your participants, 2. Propose accomplishable objectives, 3. Define schedule milestones/roadmap, 4. Pre-prepare, and 5. Build a strong data management system.

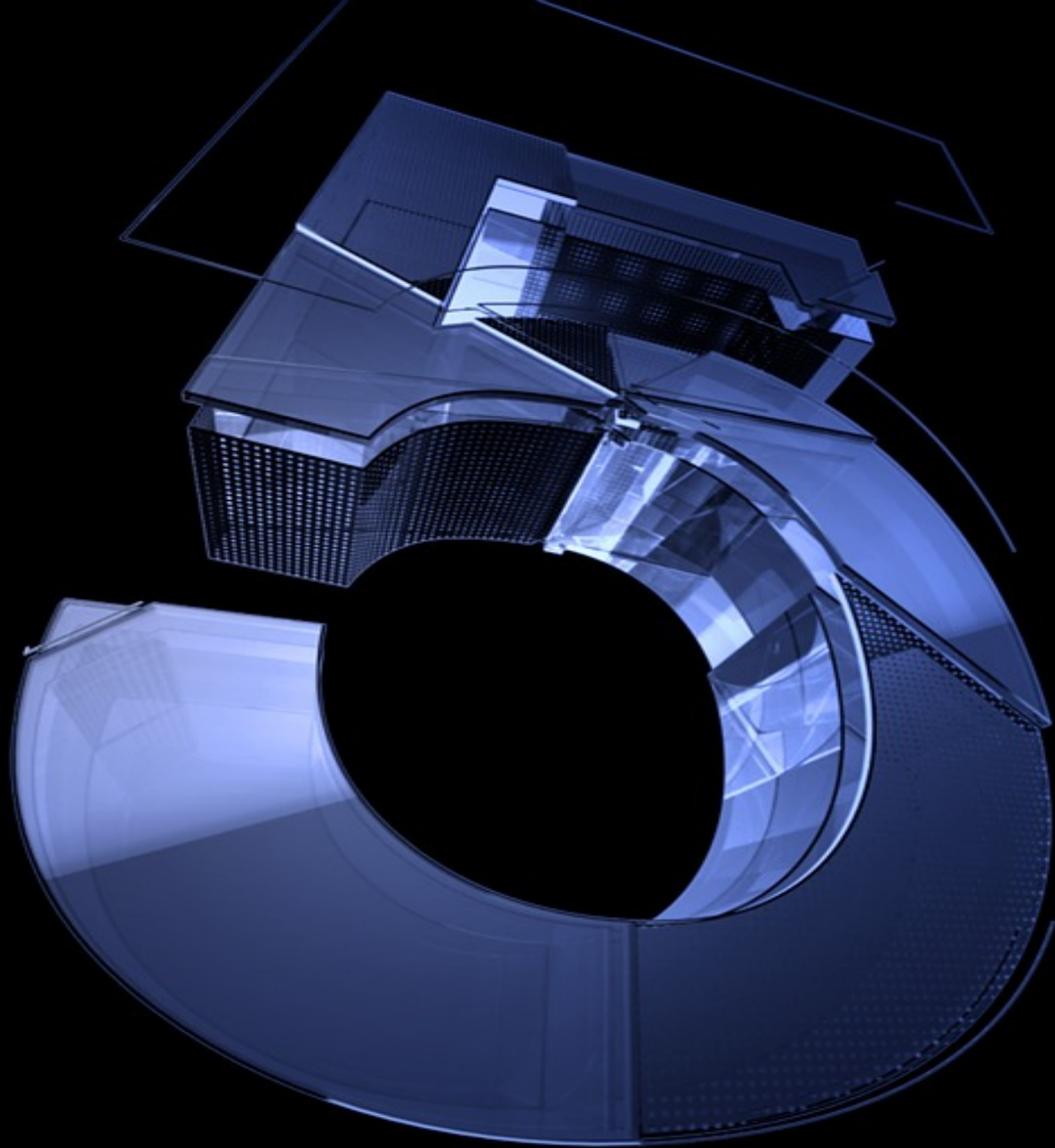
Is there domain
expertise available?

Elevator

3

What data is
available?





Is the data clean
& usable?

Step 5 – Data cleaning

The life of a data scientist

Data scientists, according to interviews and expert estimates, spend from 50 percent to 80 percent of their time mired in this more mundane labor of collecting and preparing unruly digital data, before it can be explored for useful nuggets.

-- "For Big-Data Scientists, 'Janitor Work' Is Key Hurdle to Insight" - The New York Times, 2014

Step 5 – Data cleaning

- Guides

- [Data sharing by Jeff Leek](#)
- [The Quartz guide to bad data](#)
- Reliability Centered Maintenance – Asset data register

country	year	cases	population
Afghanistan	1999	31745	1999071
Afghanistan	2000	2666	2000360
Brazil	1999	31737	17206362
Brazil	2000	80488	17400898
China	1999	212258	127201272
China	2000	210766	128003583

variables

country	year	cases	population
Afghanistan	1999	31745	1999071
Afghanistan	2000	2666	2000360
Brazil	1999	31737	17206362
Brazil	2000	80488	17400898
China	1999	212258	127201272
China	2000	210766	128003583

observations

country	year	cases	population
Afghanistan	1999	31745	1999071
Afghanistan	2000	2666	2000360
Brazil	1999	31737	17206362
Brazil	2000	80488	17400898
China	1999	212258	127201272
China	2000	210766	128003583

values

- Data steward

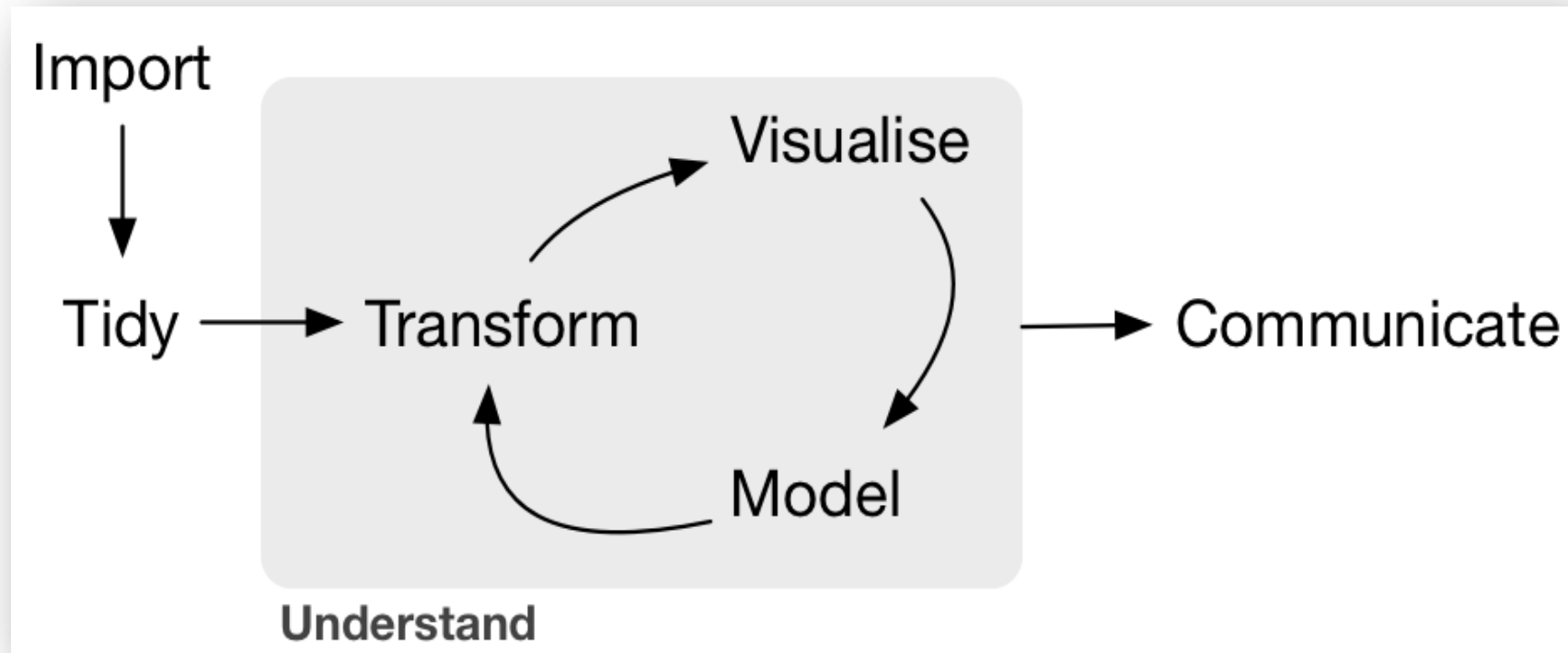
- Is it tidy data?

- Transactional data
- SCADA data
- Failure data (time, energy)

- Each variable in the data set is placed in its own column
- Each observation is placed in its own row
- Each value is placed in its own cell

Discovery Loop

Step 6 – Discovery loop

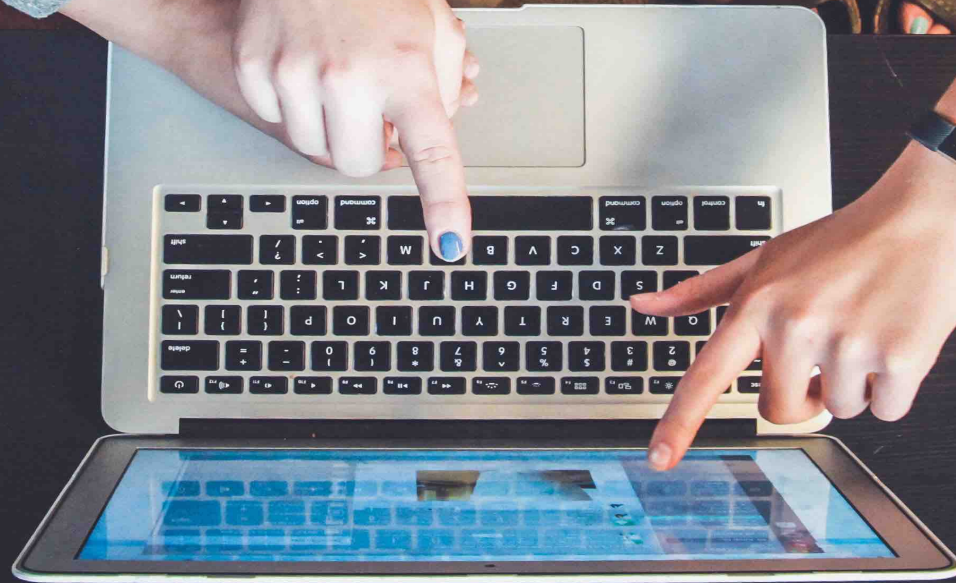


by Hadley Wickham

www.makeaweibull.com

Demo

[#makeaweibull](#)



Main Takeaways

- There is a process!
 - Step 1 – Define objective of the modeling
 - Step 2 – Define ideal data set
 - Step 3 – Add domain expert (iterate step 1 & 2)
 - Step 4 – Document available data
 - Step 5 – Clean data & make usable
 - Step 6 – Discovery loop
- Base models go a long way!

bit.ly/BarcelonaR

Thank you!

Q&A



RELIABLE DYNAMICS

Carles CG

Data Scientist & founder
carles@reliabledynamics.com

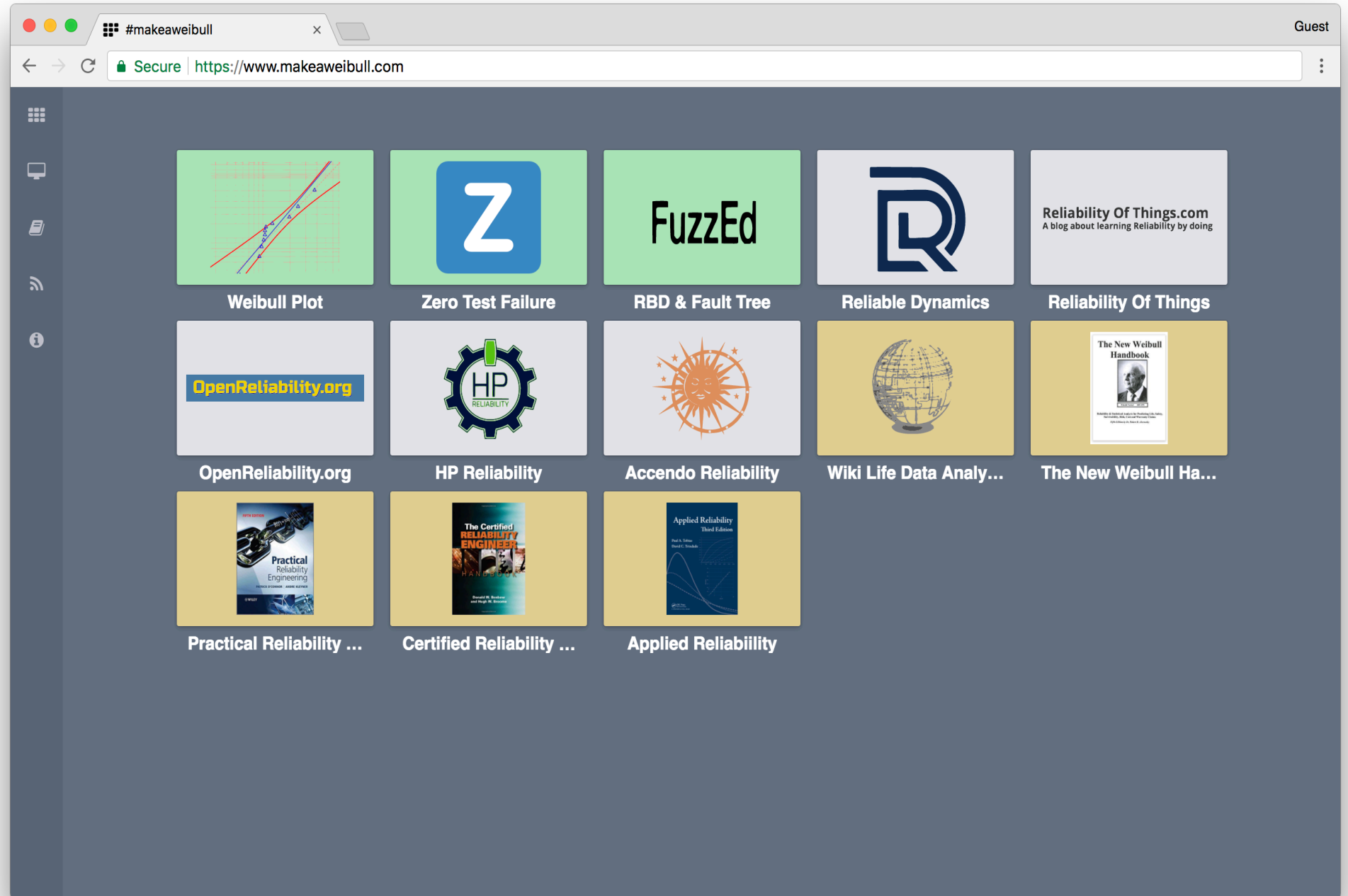


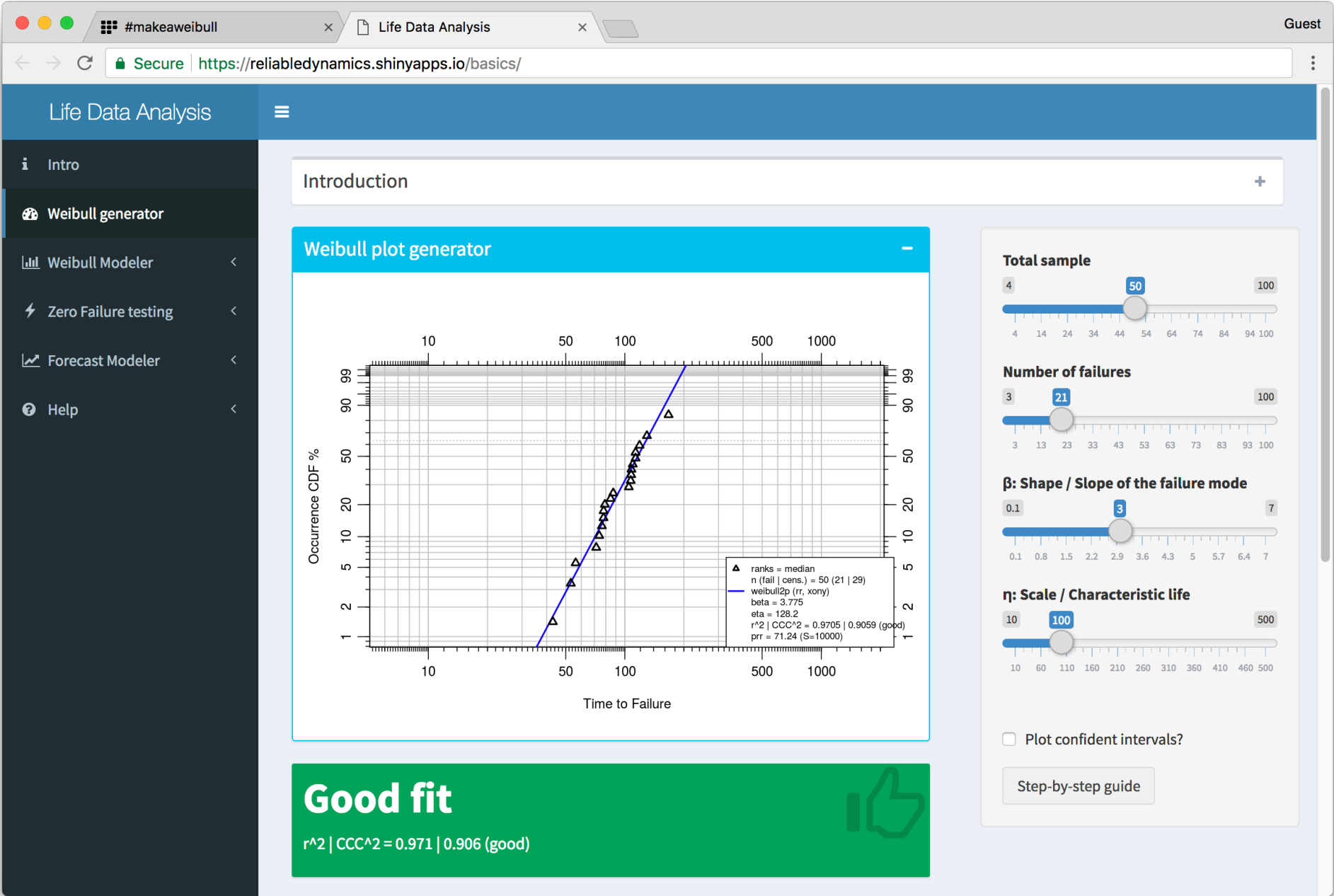
carlescg



@carles_

makeaweibull.com





#makeaweibull

Life Data Analysis

Secure

https://reliabledynamics.shinyapps.io/basics/

Guest

Intro

Weibull generator

Weibull Modeler

Zero Failure testing

Time testing complete

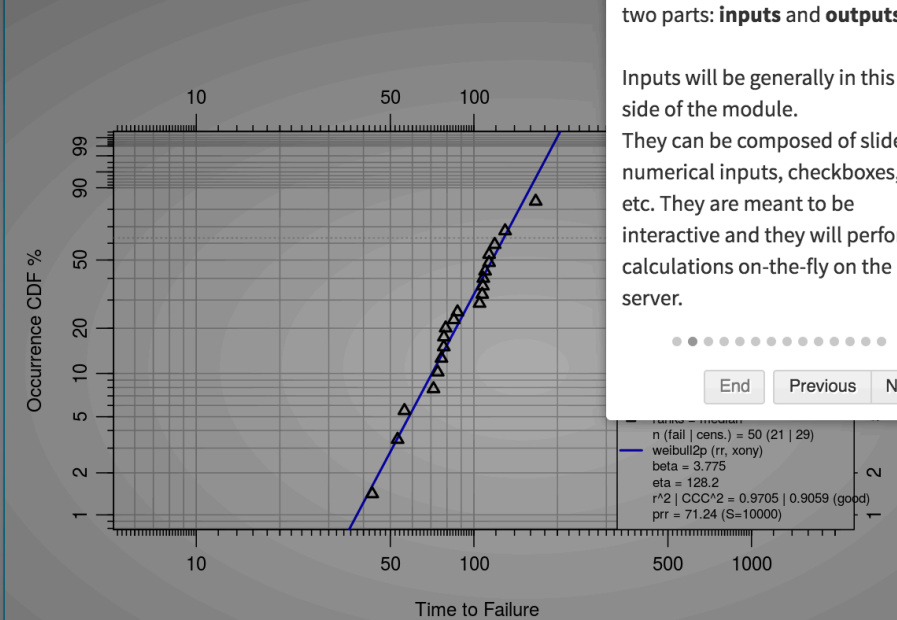
Time testing partial

Forecast Modeler

Help

Introduction

Weibull plot generator



The modules are composed into two parts: **inputs** and **outputs**.

Inputs will be generally in this side of the module.

They can be composed of sliders, numerical inputs, checkboxes, etc. They are meant to be interactive and they will perform calculations on-the-fly on the server.

End

Previous

Next

Good fit

r^2 | CCC^2 = 0.971 | 0.906 (good)

Countour plot

2

Total sample

4

50

100

Number of failures

3

21

100

β : Shape / Slope of the failure mode

0.1

3

7

η : Scale / Characteristic life

10

100

500

☐ Plot confident intervals?

Step-by-step guide

#makeaweibull

BarcelonaR #1

#makeaweibull

Life Data Analysis

Secure

https://reliabledynamics.shinyapps.io/basics/

Guest

Life Data Analysis

Intro

Weibull generator

Weibull Modeler

Upload data

Calculate

Zero Failure testing

Forecast Modeler

Help

Data

Show 25 entries

Search:

row_id	part_id	time	event	failure_mode
1	A10	412	1	Mode A
2	A09	551	1	Mode A
3	A08	858	1	Mode A
4	A08	600	1	Mode A
5	A08	700	1	Mode A
6	A08	100	1	Mode A
7	A01	913	0	None
8	A02	913	0	None
9	A03	913	0	None
10	A04	913	0	None
11	A05	913	0	None
12	A06	913	0	None
13	A07	913	0	None
14	A08	55	0	None

Choose CSV File

Browse...

Weibull_template (2).csv

Upload complete

☒ Header

Separator

☒ Comma

☐ Semicolon

☐ Tab

Quote

☐ None

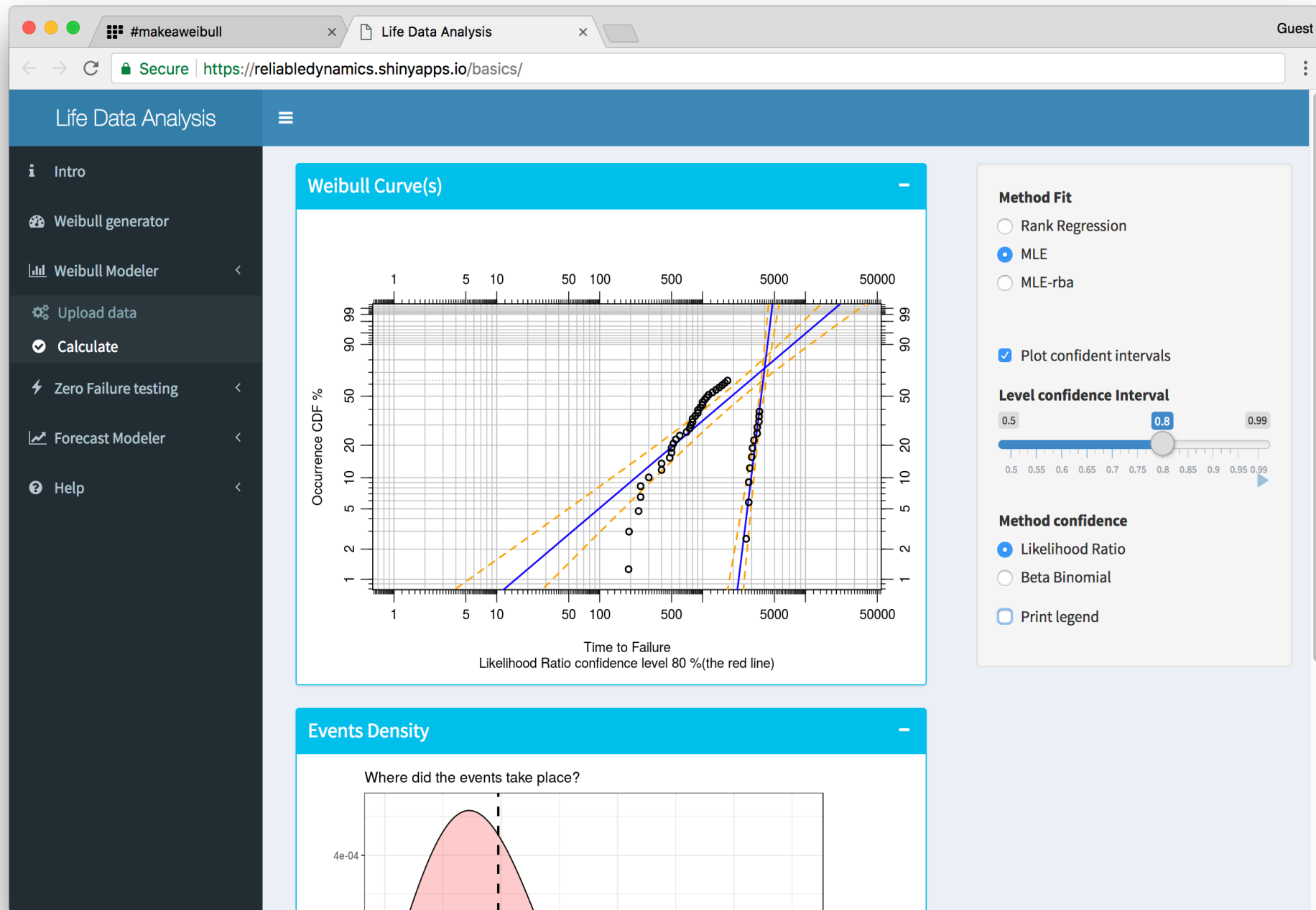
☒ Double Quote

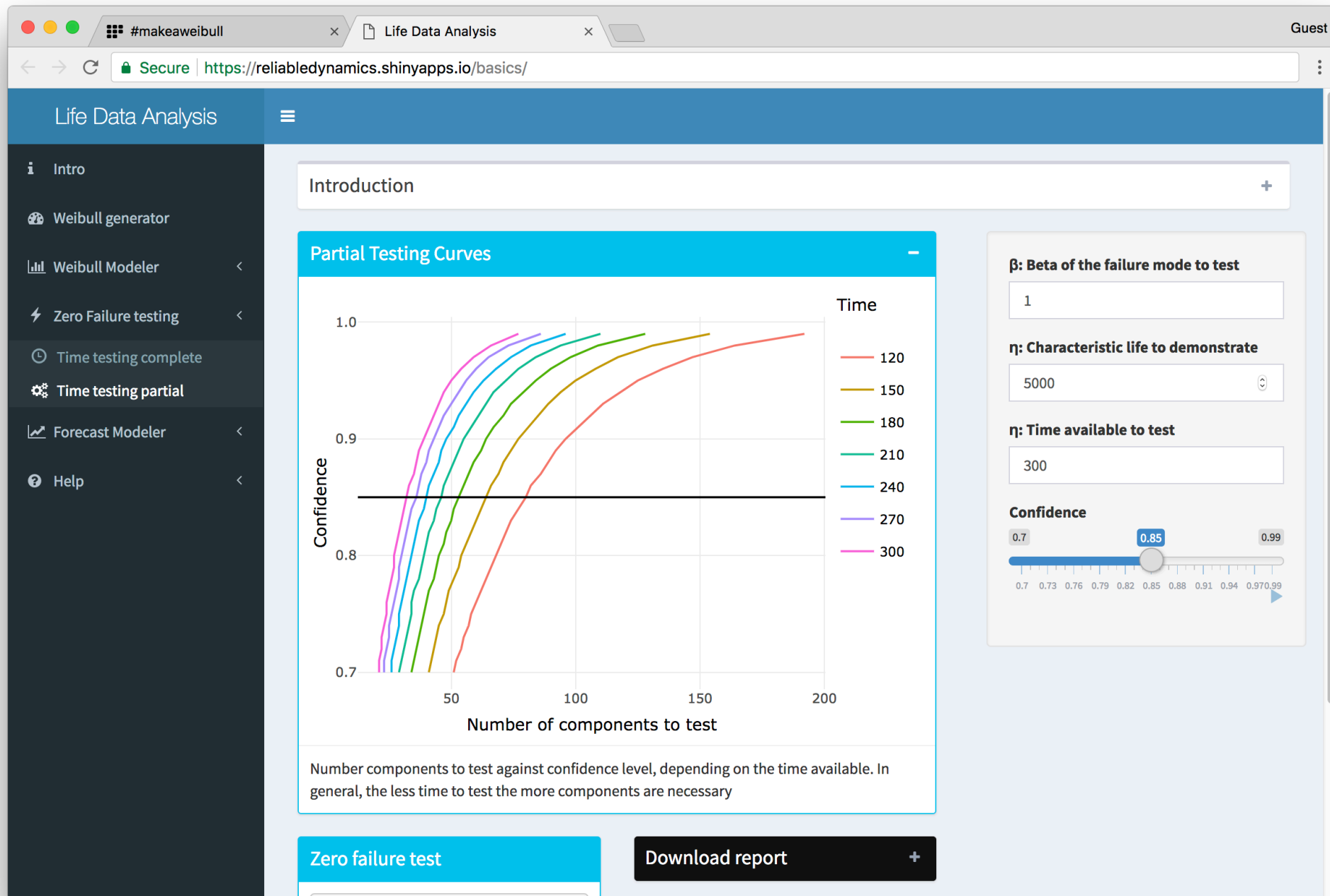
☐ Single Quote

Need a template?

Download

Guide





Demo End