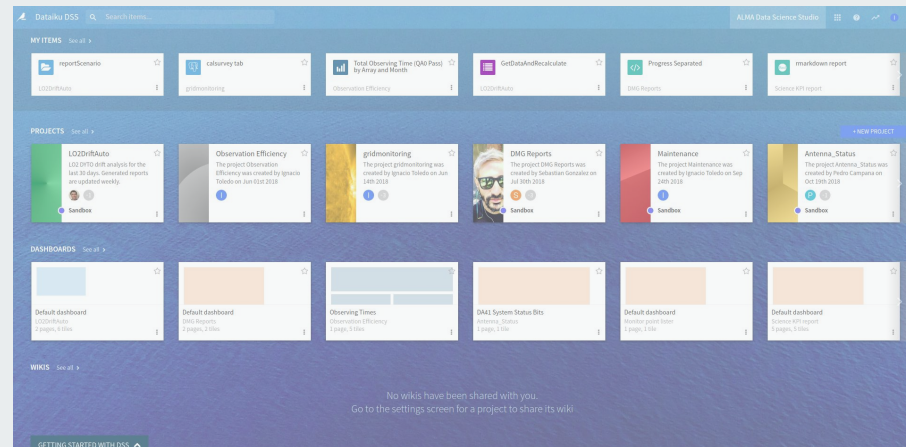# Data Science, not Software Engineering.
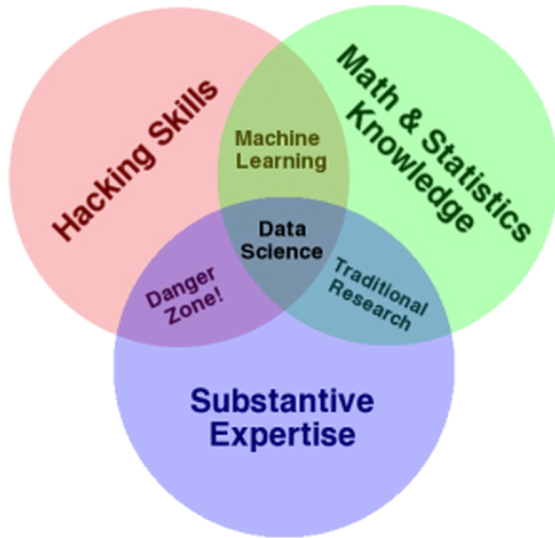
Exploring a workflow for ALMA operations.

# Abstract

In the last few years Data science has emerged as a discipline of its own to address problems where data is usually heterogeneous, complex and abundant. In a nutshell, data science allows to provide answers to situations where a hypothesis can be formulated and later can be either confirmed or rejected following standard scientific methodology using data as raw material.

Data science has been called differently depending of the domain (business intelligence, operational management, astroinformatics) and it has been recently in the center of a hype related to artificial intelligence and machine learning. It has been quickly adopted by the digital industry as the tool to distill information of massive operational data sets. Among the many tools data science requires (mathematics, statistics, domain knowledge of the data sets, …), IT infrastructure and software is by far the most visible and there is at present a whole ecosystem available as open source projects. The downside of this is data science is commonly confused with IT and software development, which creates conflicts between engineering- and scientific- mindsets, and leads to wrongly applying software development methodologies to it neglecting the experimental nature of the problem. In summary, creating the data lab becomes more important than answering questions with it. In the domain of ALMA operations, there are many instances that can be identified and described as data science cases or projects ranging from monitoring array elements to understand performances and predict faults for engineering operations to routine monitoring of calibrators for science operations purposes.
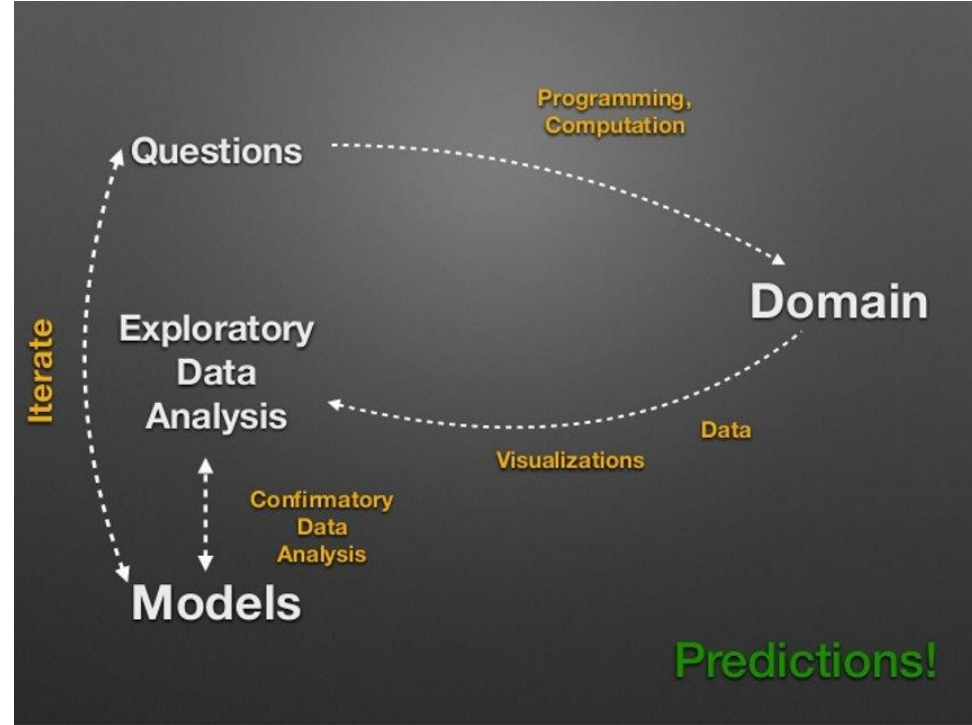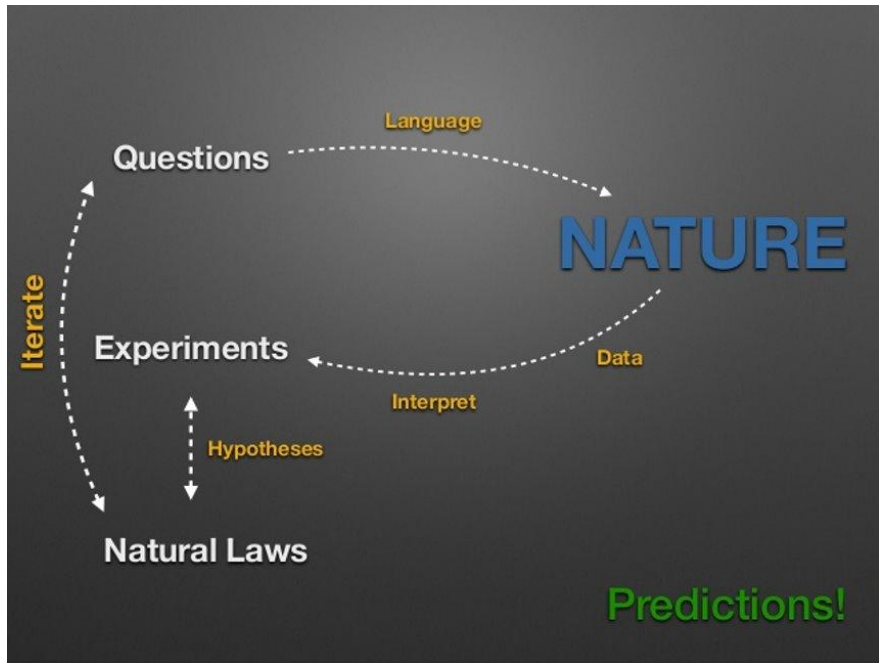
We have identified already around 30 different initial questions (or data science cases) and found that several of them have been addressed through individual efforts. In parallel, several enabling platforms or frameworks have appear in the ecosystem that provides data scientists with both the "laboratory equipment" to conduct their "experiments" as well as enabling tools for collaboration, versioning control, and deploying results in production with a quick turnaround. This talk aims to summarize the results of our exploration to apply data science workflows to resolve ALMA operations issues, identify suitable platforms that are already in use by the industry, share our experience in addressing specific ALMA operations data cases, and discuss the technical and sociological challenges we encountered along the way.

# Overview

- What is Data Science (at least for the next 10 minutes)
- Why focusing in the differences with Software Engineering?
- ALMA - Data iku Experience:
  - Don't reinvent the wheel
  - Take advantage of your strengths
  - Socialize, communicate, collaborate
  - Our challenge: Data Engineering!

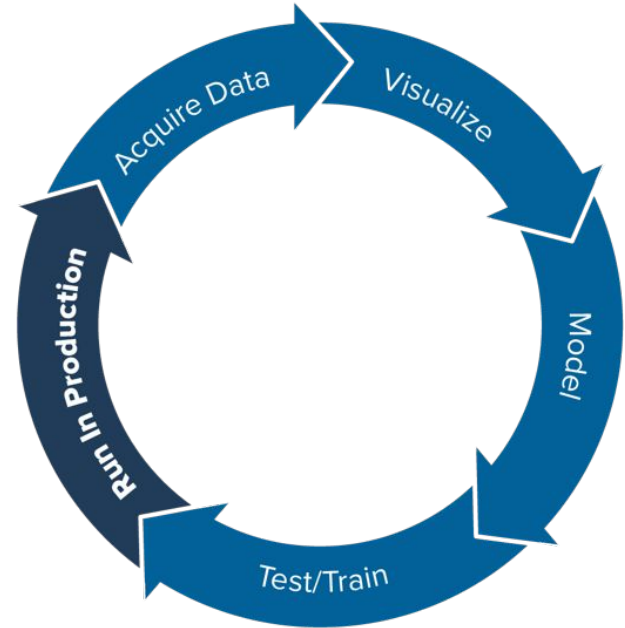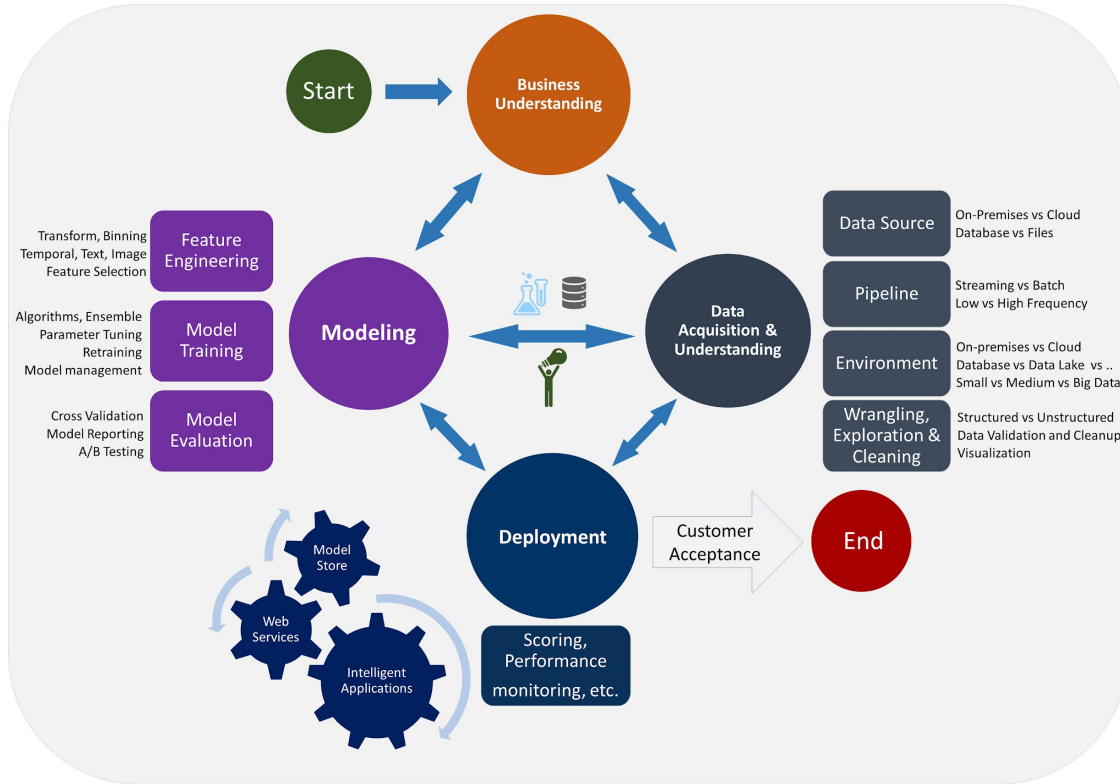*"Data science is the discipline of making data useful."*

*"Data science is science... it just starts with raw data exploration instead of physical observation"*

What is Data Science?

A new methodology for an old process

# Data Science Lifecycle



**Start** → **Business Understanding**

Transform, Binning
Temporal, Text, Image
Feature Selection
**Feature Engineering**

Algorithms, Ensemble
Parameter Tuning
Retraining
Model management
**Model Training**

Cross Validation
Model Reporting
A/B Testing
**Model Evaluation**

**Modeling**

**Data Acquisition & Understanding**

| | |
|---|---|
| Data Source | On-Premises vs Cloud Database vs Files |
| Pipeline | Streaming vs Batch Low vs High Frequency |
| Environment | On-premises vs Cloud Database vs Data Lake vs .. Small vs Medium vs Big Data |
| Wrangling, Exploration & Cleaning | Structured vs Unstructured Data Validation and Cleanup Visualization |

**Deployment**

Model Store
Web Services
Intelligent Applications

Scoring, Performance monitoring, etc.

Customer Acceptance → **End**

A different methodology

Acquire Data — Visualize — Model — Test/Train — Run In Production

How is the observatory time distributed among the different executives?

Can we automate early detection or prediction of hardware failures leading to observation downtime or degradation? Can we automate diagnosis?

There is a need for a general set of tools to explore relationships between disparate arbitrary data streams as a support for problem investigation.

Hardware malfunction fault detection and diagnosis (FDD)

Given the current configuration and the forthcoming, how much pressure do we have of Science Projects, by band and LST?

# Template: experiment name

ALMA Stakeholder: ALMA (department/team/manager). Executed by: Team member(s)

| 1 | What was(is) the question or thesis? | 2 | What data was(is) needed and what science was done? | 3 | What was find out? How was it shared? Any gains? |

Detect and diagnose IFPs amplifiers degradations on any antenna on 2 polarizations each, 4 basebands each polarization. Degradation leads to TP detectors readings calibration drift.

Quantify the dependency of photonic reference failure rate as a function of various factors to identify strategies for improvement

Track Needed Grid Executions, reduction & Ingestions: Observing Log of Grid data. Which sources need to be observed? Which sources/datasets need to be reduced/ingested ? Who is reducing the data ?

Are these software engineering or data science problems? Hint: check the outcome.

# What are we doing at ALMA?

The Data Science initiative

- Not much resources: do not reinvent the wheel!
- Astronomers and engineers with statistical, mathematical, domain knowledge and coding skills. Work with them!
- Collaborate, socialize, communicate!
- Show results.

A data science platform collaboration

ALMA

data iku

# Make data accessible

# Create a laboratory with the tools needed

# Collaborate and socialize



## Summary

Show [ Commits ] for [ Last year ]

|       | Jul 15 | Jul 22 | Jul 29 | Aug 5 | Aug 12 | Aug 19 | Aug 26 | Sep 2 | Sep 9 | Sep 16 | Sep 23 | Sep 30 | Oct 7 | Oct 14 | Oct 21 | Oct 28 | Nov 4 | Nov 11 |
|-------|--------|--------|--------|-------|--------|--------|--------|-------|-------|--------|--------|--------|-------|--------|--------|--------|-------|--------|
| Sun   |        |        |        |       |        |        |        |       |       |        |        |        |       |        |        |        |       |        |
| Mon   |        |        |        |       |        |        |        |       |       |        |        |        |       |        |        |        |       |        |
| Tue   |        |        |        |       |        |        |        |       |       |        |        |        |       |        |        |        |       |        |
| Wed   |        |        |        |       |        |        |        |       |       |        |        |        |       |        |        |        |       |        |
| Thu   |        |        |        |       |        |        |        |       |       |        |        |        |       |        |        |        |       |        |
| Fri   |        |        |        |       |        |        |        |       |       |        |        |        |       |        |        |        |       |        |
| Sat   |        |        |        |       |        |        |        |       |       |        |        |        |       |        |        |        |       |        |

## Contributors activity

During last year, **3 authors** have created **206 commits, 18,605 additions** and **7,275 deletions.**
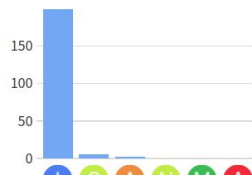
History

- 2 commits on Oct 31, 2018

  C  Updated exposed objects (project:GRIDMONITORING)                    f027515
     cverdugo - 12 days ago

  C  Updated exposed objects (project:GRIDMONITORING)                    9be2d0f
     cverdugo - 12 days ago

- 1 commit on Oct 1, 2018

  I  Updated project permissions (project:GRIDMONITORING)                e28eacd
     itoledo - 1 month ago

- 1 commit on Sep 21, 2018

  I  Saved explore settings of calsurvey_tab                             aebb728
     itoledo - 1 month ago

- 1 commit on Sep 13, 2018

MANAGE

## gridmonitoring

WATCH ▾ 2   ★ STAR  0

The project *gridmonitoring* was created by Ignacio Toledo on Jun 14th 2018

I A C S

imported

No project status yet

| notebook editor for reci... | 3 minutes ago |
| calsurvey_tab | 4 days ago |
| Last Observed | 22 days ago |

### Automation

| 1 | 149 | 14/10 | 21/10 | 28/10 | 4/11 | 11/11 |
| ACTIVE SCENARIOS | RUNS | | | | | |

| Flow | | Lab | | Dashboards | | Wiki | | Tasks | |
|------|--|-----|--|------------|--|------|--|-------|--|
| 🗄 **19** | | </> **16** | | 📊 **0** | | 📖 **0** | | ☑ **0** | |
| DATASETS | | RECIPES | | DASHBOARD | | ARTICLE | | TASK | |

# Show the results!

Days since last observation in Bands 3 and 7 on calsurvey_tab

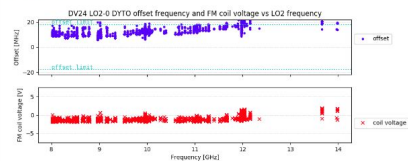| | B3 | B7 | |
|---|---|---|---|
| J0006-0623 | 3 | 4 | 3 |
| J0237+2848 | 2 | 3 | 2 |
| J0238+1636 | 2 | 3 | 2 |
| J0319+4130 | 2 | 3 | 2 |
| J0334-4008 | 0 | 2 | 0 |
| J0423-0120 | 2 | 3 | 2 |
| J0510+1800 | 0 | 2 | 0 |
| J0519-4546 | 0 | 2 | 0 |
| J0522-3627 | 0 | 2 | 0 |
| J0538-4405 | 0 | 2 | 0 |
| J0635-7516 | 0 | 2 | 0 |
| J0725-0054 | 0 | 2 | 0 |
| J0750+1231 | 0 | 2 | 0 |
| J0854+2006 | 21 | 2 | 2 |
| J0904-5735 | 0 | 2 | 0 |
| J1037-2934 | 0 | 2 | 0 |
| J1058+0133 | 0 | 2 | 0 |

Default dashboard

LO2 replacement candidates - Top 5 priorities from frequency offset analysis

| Antenna | BBpr | Frequency offset count | Analysis date |
|---|---|---|---|
| DV24 | 0 | 2339 | 2018-11-12 04:37:25 |
| DA56 | 3 | 1100 | 2018-11-12 04:37:25 |
| CM09 | 1 | 774 | 2018-11-12 04:37:25 |
| DV09 | 3 | 674 | 2018-11-12 04:37:25 |
| DV12 | 0 | 543 | 2018-11-12 04:37:25 |

Default dashboard

priority_0

DV24 LO2-0 DYTO offset frequency and FM coil voltage vs LO2 frequency

priority_1

DA56 LO2-3 DYTO offset frequency and FM coil voltage vs LO2 frequency

priority_2

CM09 LO2-1 DYTO offset frequency and FM coil voltage vs LO2 frequency

priority_3

DV25 LO2-3 DYTO offset frequency and FM coil voltage vs LO2 frequency
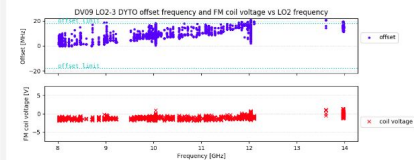
# What is Next? A new challenge.



**Framework to use modern Big Data Software Tools to improve operations at the Paranal Observatory**

Eduardo Pena*, Ricardo Schmutzer, Christian Stephan, Claudio Reinero, Julien Milli, Juan C. Guerra, Juan Osorio European Southern Observatory, Alonso de Córdova 3107, Vitacura, casilla 19001, Santiago, CHILE

# Questions?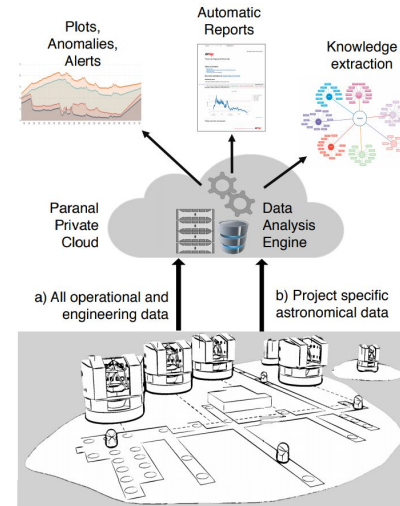