



Human in the Loop: Active Learning for Astronomy

*Michigan Cosmology Summer School - part 2
8 June 2023 - Ann Arbor, USA*

Emille E. O. Ishida

*Laboratoire de Physique de Clermont, CNRS - Université Clermont-Auvergne
Clermont Ferrand, France*



Acknowledgements



Alberto Krone-Martins (UCI)
Rafael S. de Souza (Hertfordshire)
and the entire [Cosmostatistics Initiative \(COIN\)](#)



Konstantin Malanchev (UIUC)
Maria Pruzhinskaya (LPC)
and everyone in the [SNAD collaboration](#)

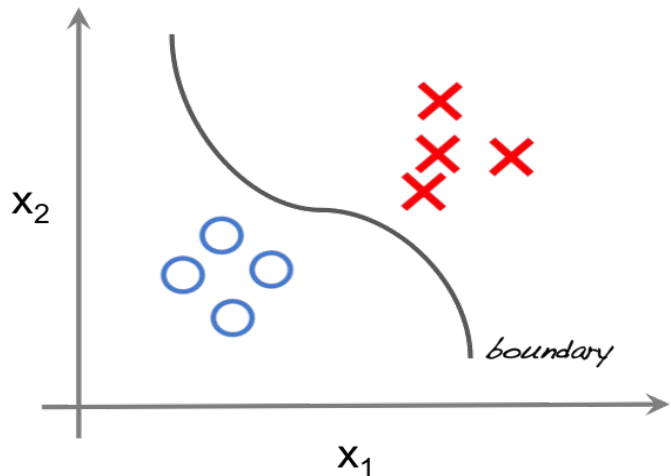


Anais Moller (Swinburne)
Julien Peloton (IJCLab)
and everyone in the [Fink broker](#)



Alex I. Malz (CMU)
Mi Dai (JHU)
Kara Ponder
Amanda Wasserman (UIUC)
and all those working in the [RESSPECT team](#)

Supervised



Training sample:

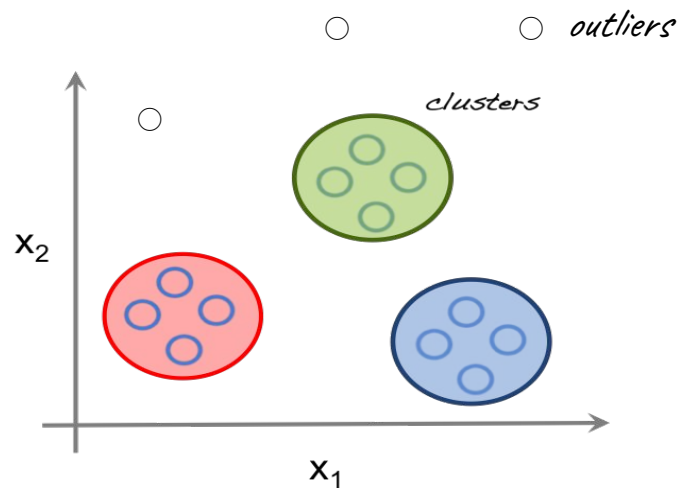
features + labels

Target sample:

features

x

Unsupervised

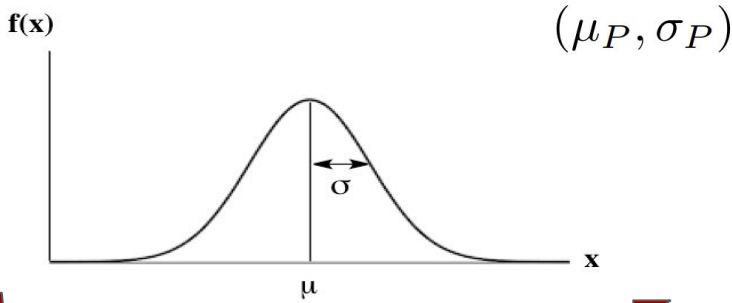


Data sample:

features

Representativeness

Probability distribution, P



This is why it works!

Training

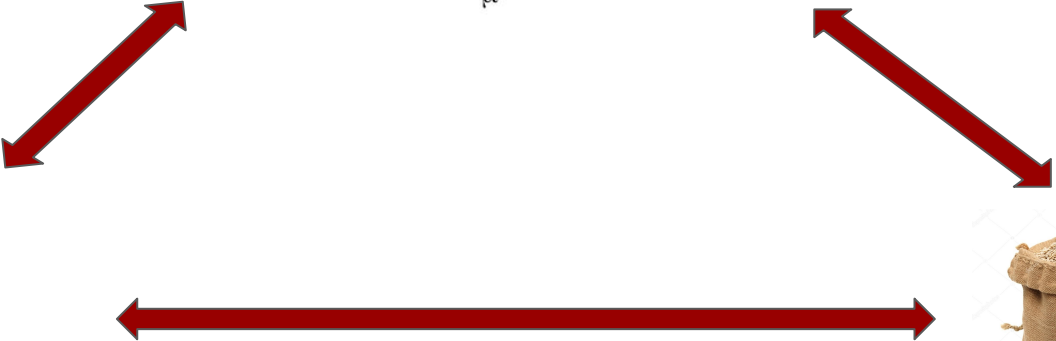


$(\mu_{S_1}, \sigma_{S_1})$

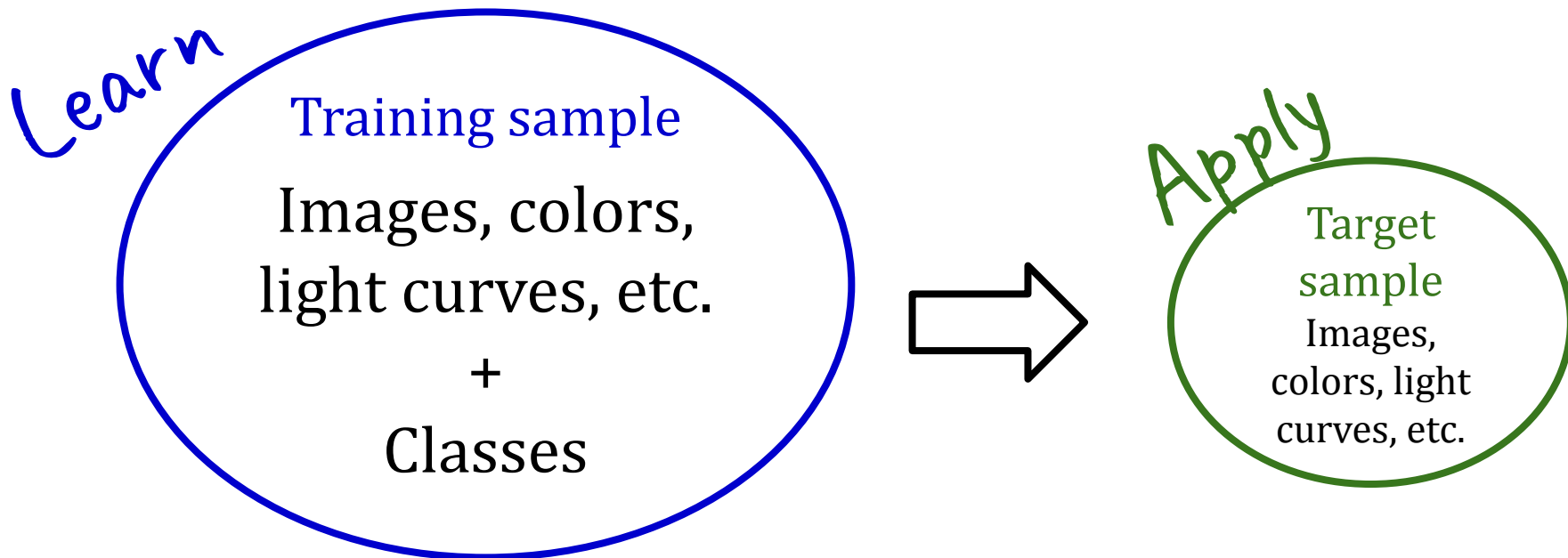
Target



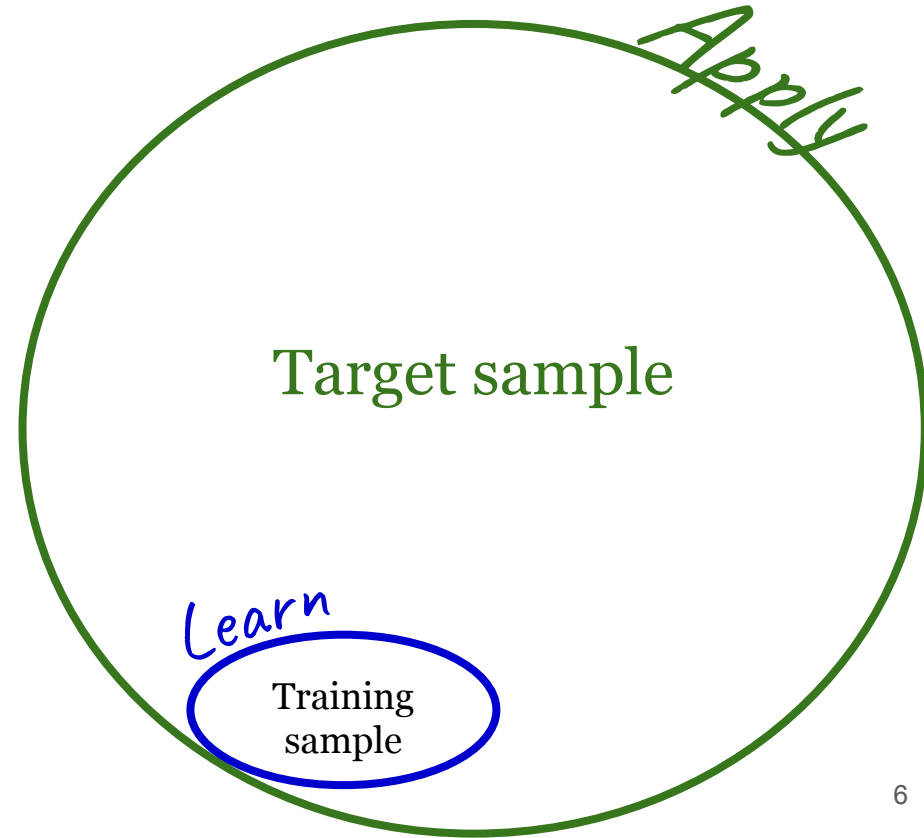
$(\mu_{S_2}, \sigma_{S_2})$



Ideal Supervised learning situation



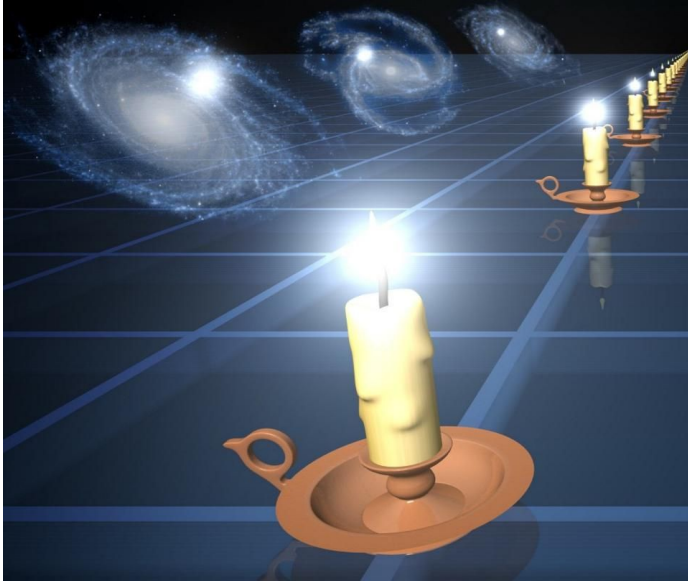
Real astro-learning situation



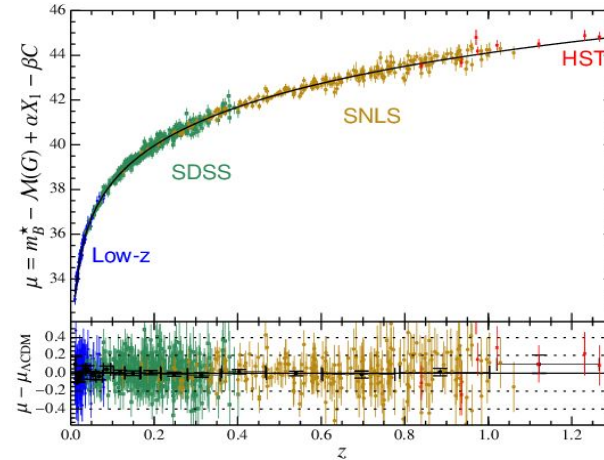
Example science case:

Type Ia supernova cosmology

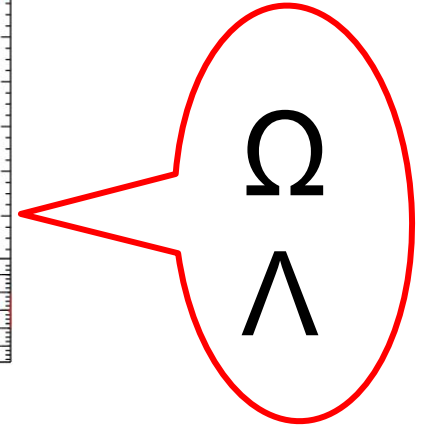
Standard candles used to measure cosmological distances



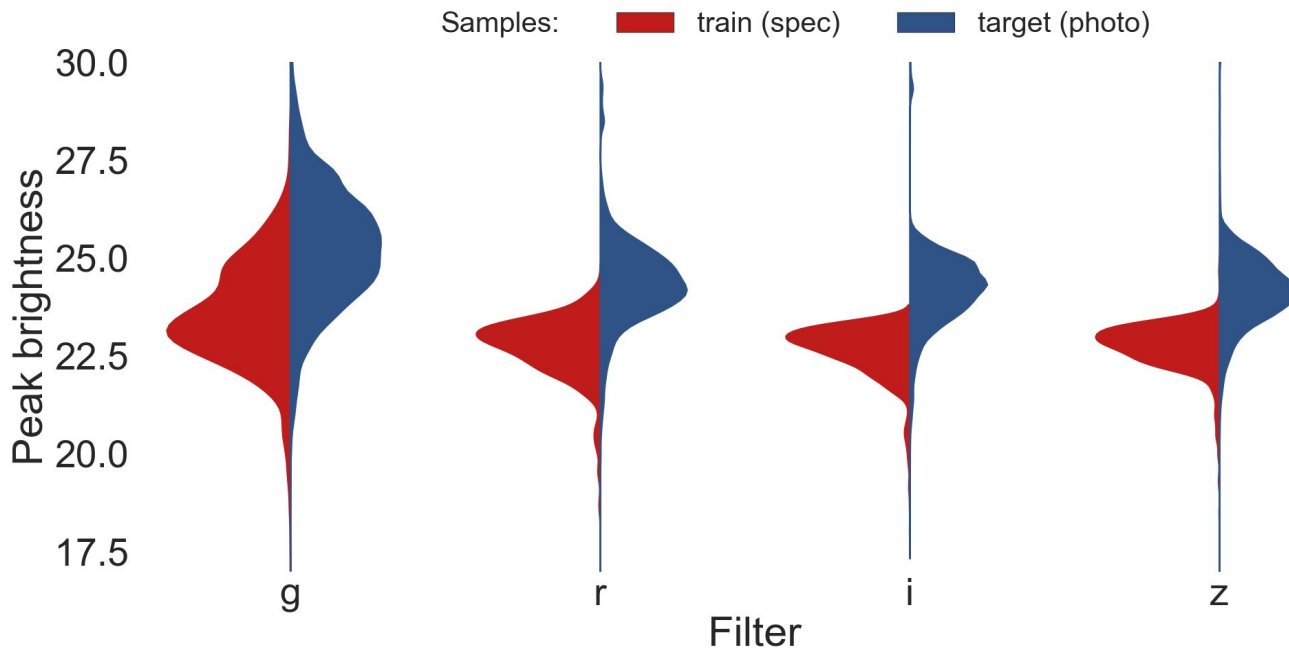
https://supernova.eso.org/exhibition/images/1111_E_549779main_pia14095_full/



http://supernovae.in2p3.fr/sdss_snls_ia/ReadMe.html



Real astro-learning situation



Very common situation

Labels are often far too expensive!



Given limited
resources, we need
recommendation
systems!

amazon

35% OF AMAZON'S REVENUE ARE GENERATED BY IT'S RECOMMENDATION ENGINE.

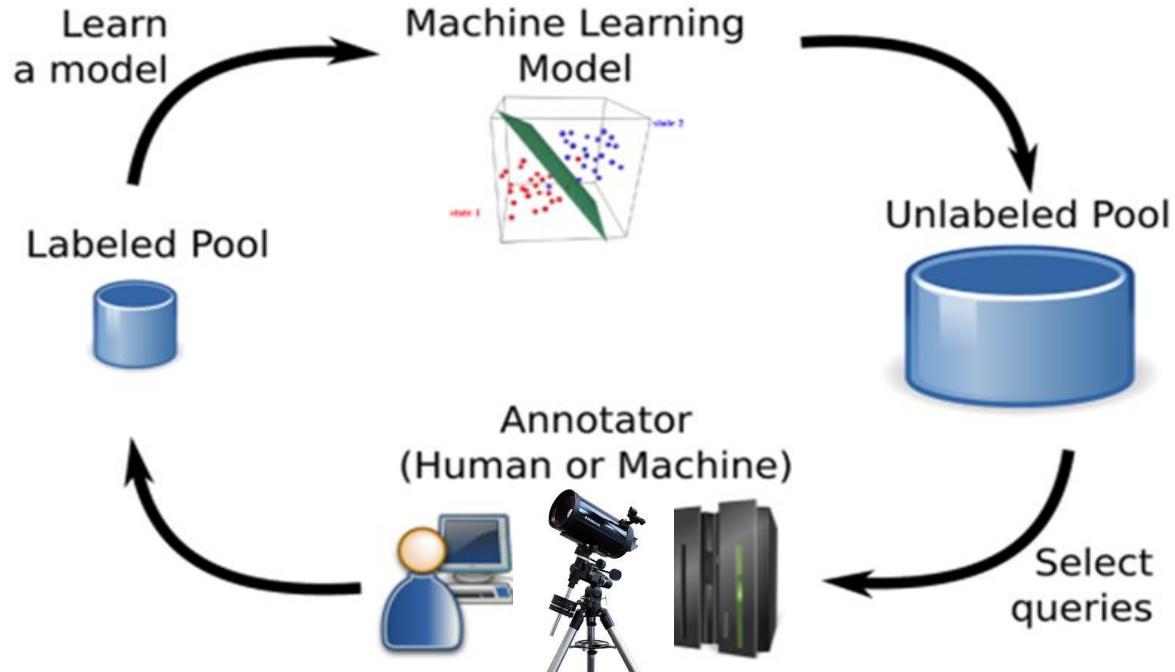
NETFLIX

75% OF USERS SELECT MOVIES BASED ON NETFLIX'S RECOMMENDATIONS.



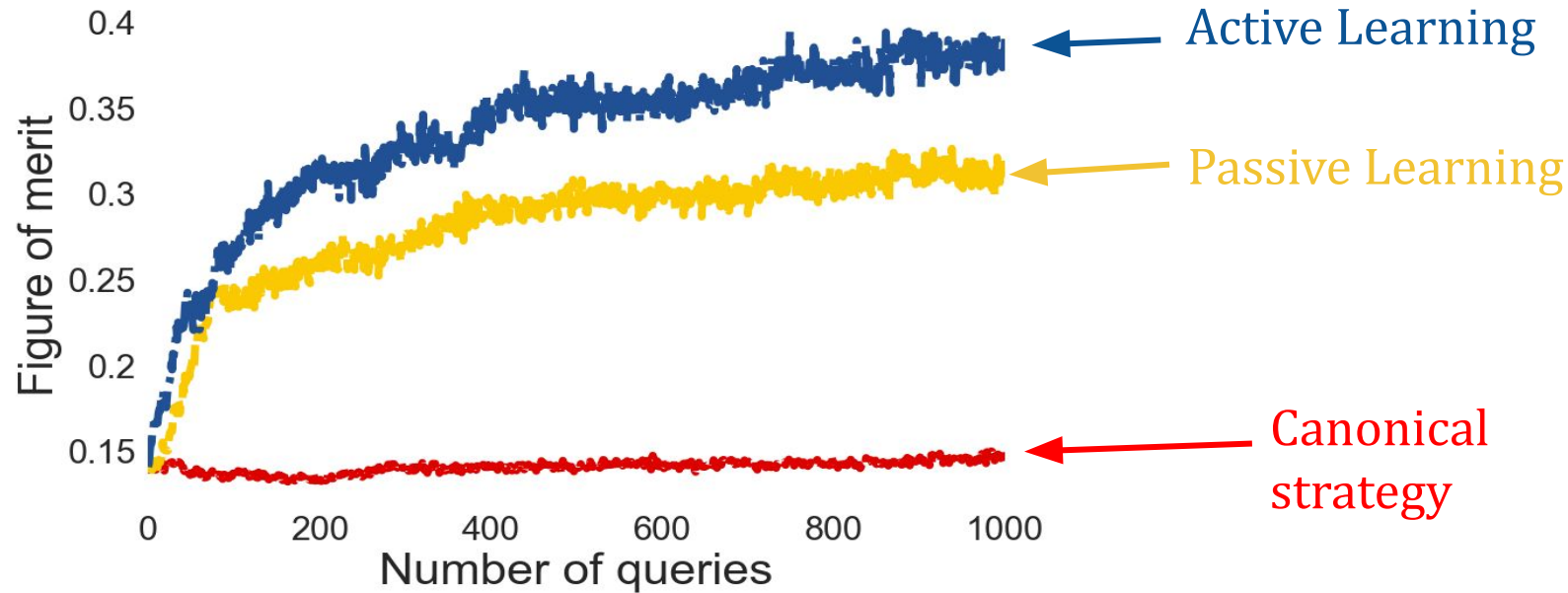
Active Learning

Optimal classification, minimum training



AL for SN classification

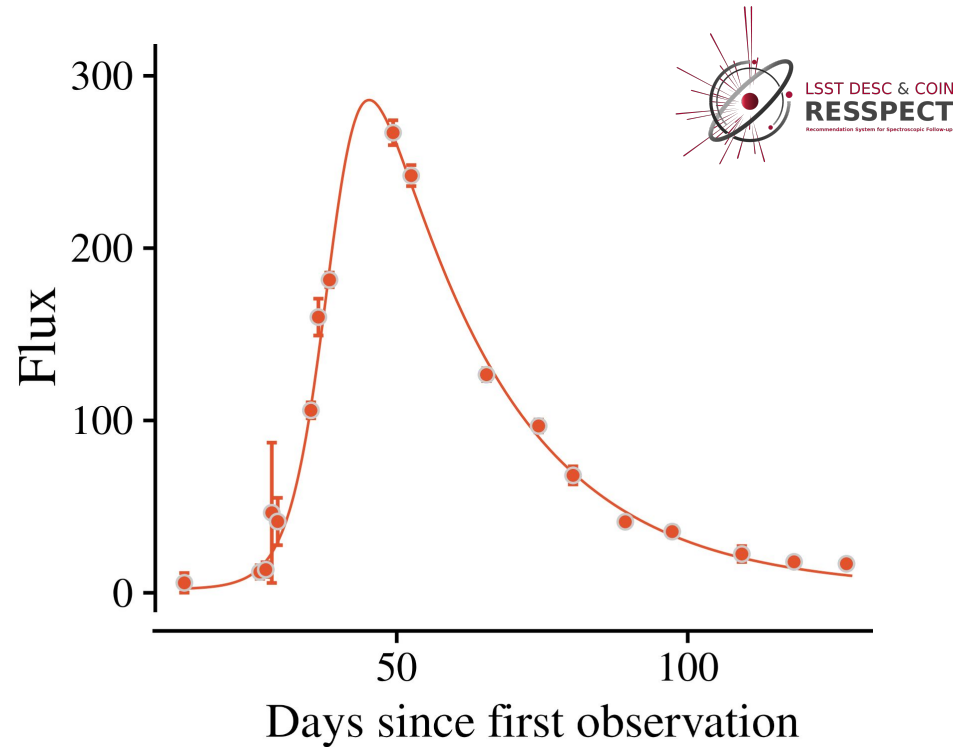
Static results



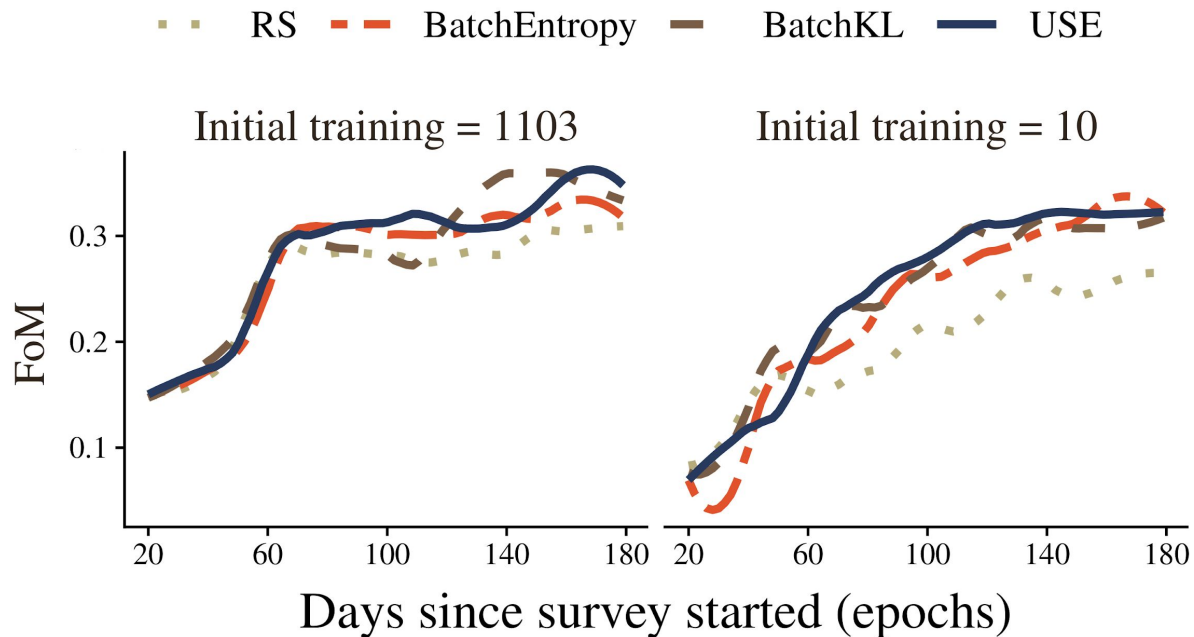
From COIN Residence Program #4, *Ishida et al., 2019, MNRAS, 483 (1), 2–18*

If only it were that simple ...

- Window of Opportunity for Labelling
 - We must make query decisions before we can observe the full LC
- Evolving Samples
 - Other people want to use the telescope
- Multiple Instruments for labelling
- Evolving budget
 - Other people want to use the telescope
- Evolving Costs
 - Observing costs for a given object changes as it evolves.

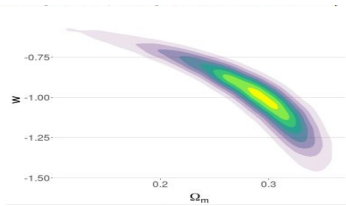


Start from scratch, do not overcomplicate

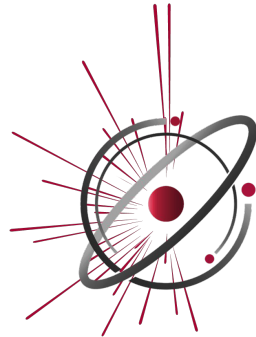


What about science?

Cosmology results from
photometrically classified SN Ia



2. Impacts
this



LSST DESC & COIN
RESSPECT
Recommendation System for Spectroscopic Follow-up

Photo-classified
SN Ia

SN
candidates

Trained
machine
learning
classifier

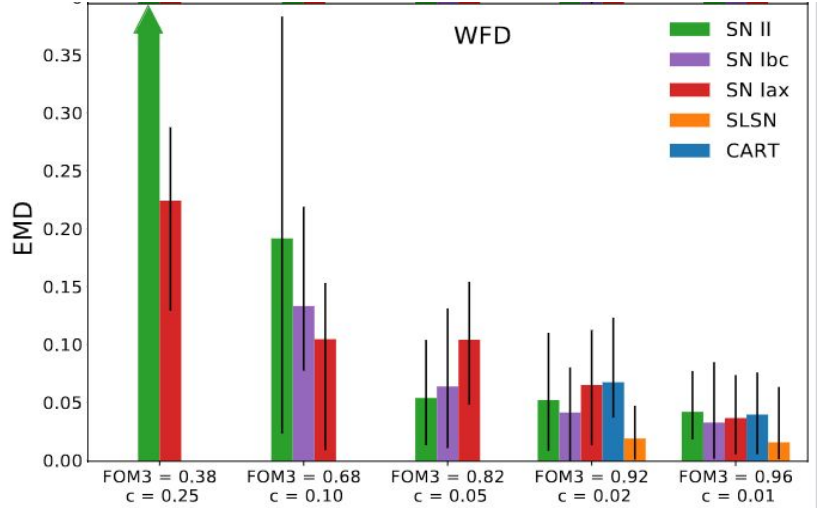
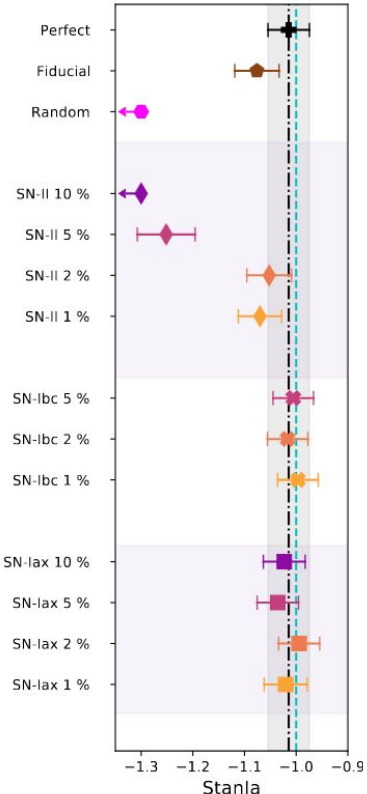
learning
algorithm

+

Training
sample

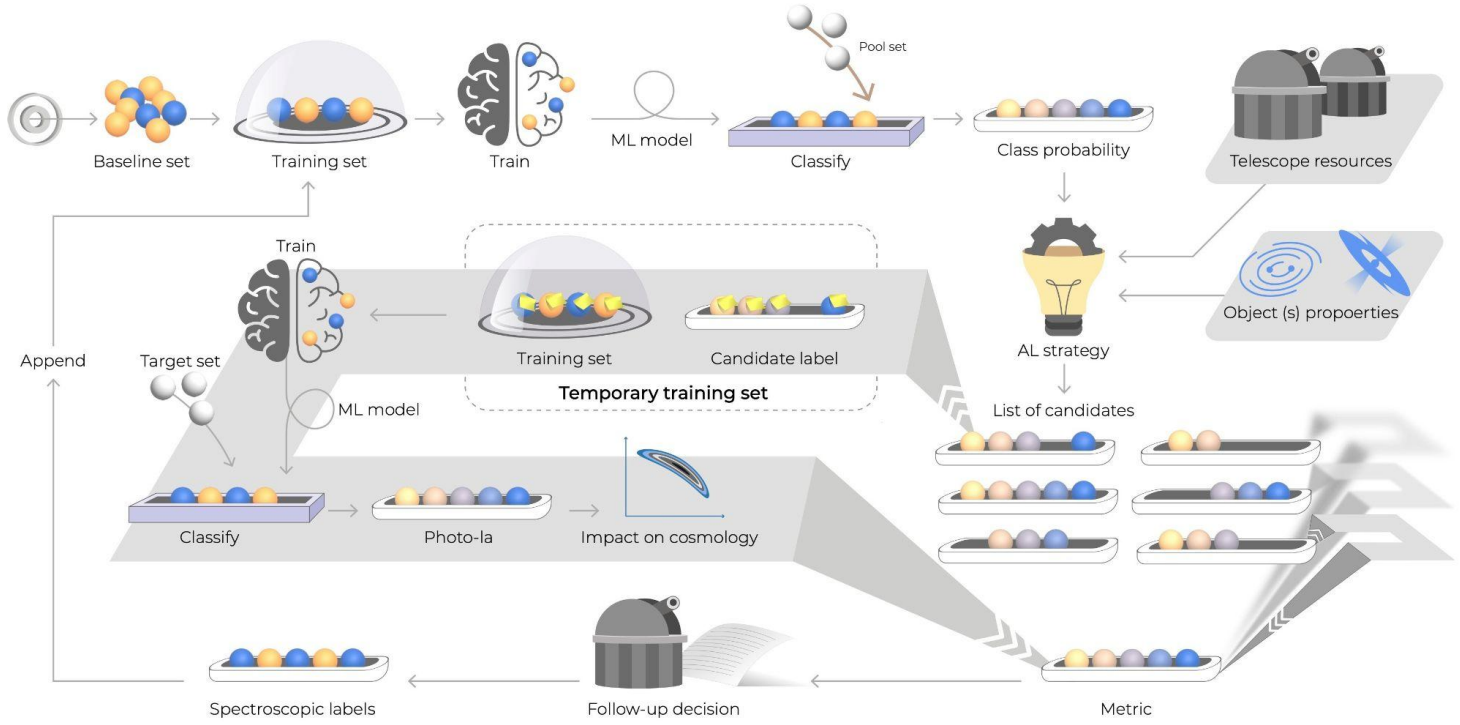
1. Different
choices of
this!

Good classification might not be enough



Malz et al., 2023 - [arXiv:astro-ph/2305.14421](https://arxiv.org/abs/2305.14421) - The RESPECT team: LSST-DESC and COIN, Are classification metrics good proxies for SN Ia cosmological constraining power? -- submitted to A&A

The RESPECT workflow



<https://github.com/COINtoolbox/RESPECT>

The difficult part is data treatment/gathering

- The power of machine learning is in its connection with domain knowledge
- There are caveats in using machine learning and we should avoid off-the-shelf and black boxes applications
- ML for science must be personalized

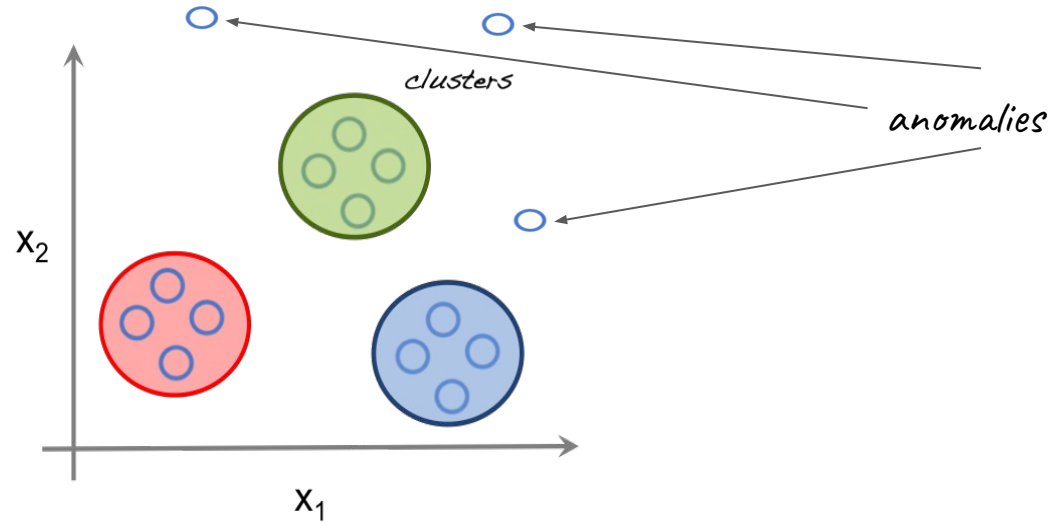
The beauty of an observational science

“... telescopes that merely achieve their stated science goals have probably failed to capture the most important scientific discoveries available to them.”

Norris, R. (2017). Discovering the Unexpected in Astronomical Survey Data. Publications of the Astronomical Society of Australia, 34, E007. doi:10.1017/pasa.2016.63

Statistically,

Anomaly Detection

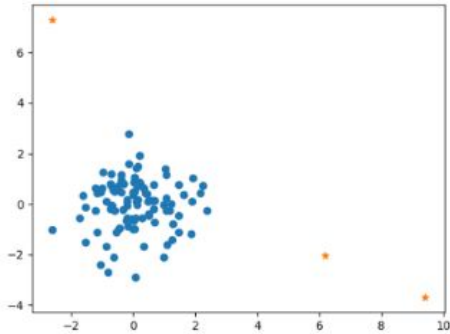


"An anomaly is an observation which deviates so much from the other observations as to arouse suspicions that it was generated by a different mechanism"

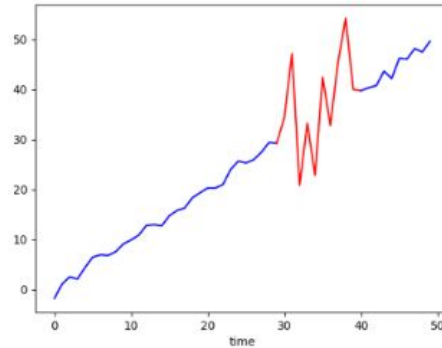
Anomaly Detection

“An anomaly is an observation which deviates so much from the other observations as to arouse suspicions that it was generated by a different mechanism”

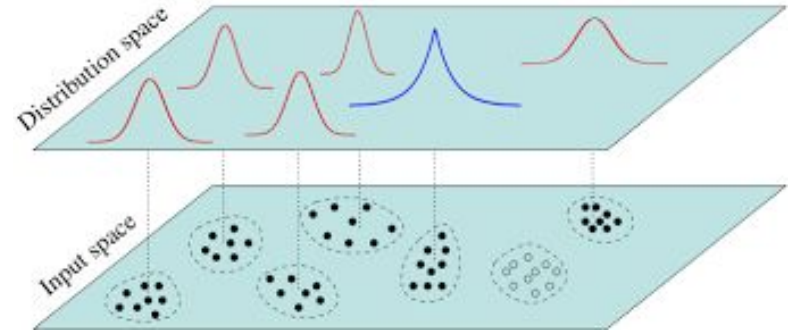
Hawkins, 1980



isolation



behaviour

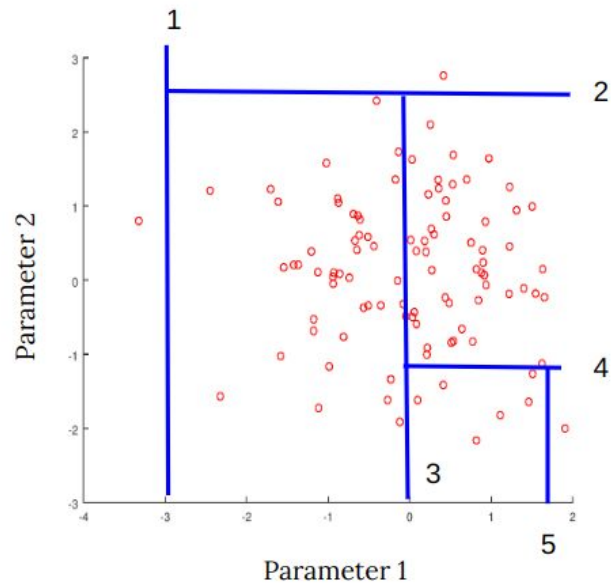
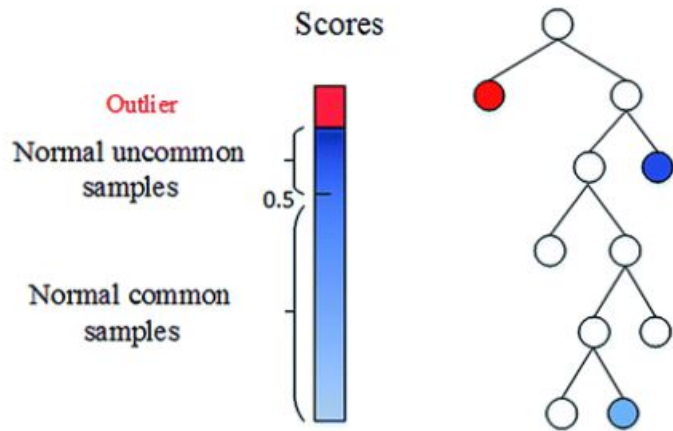


Plot from Muandet and Scholkopf, 2013 - [arXiv:stat.ML/1303.0309](https://arxiv.org/abs/1303.0309)

group

Example of an automatic search for anomalies,

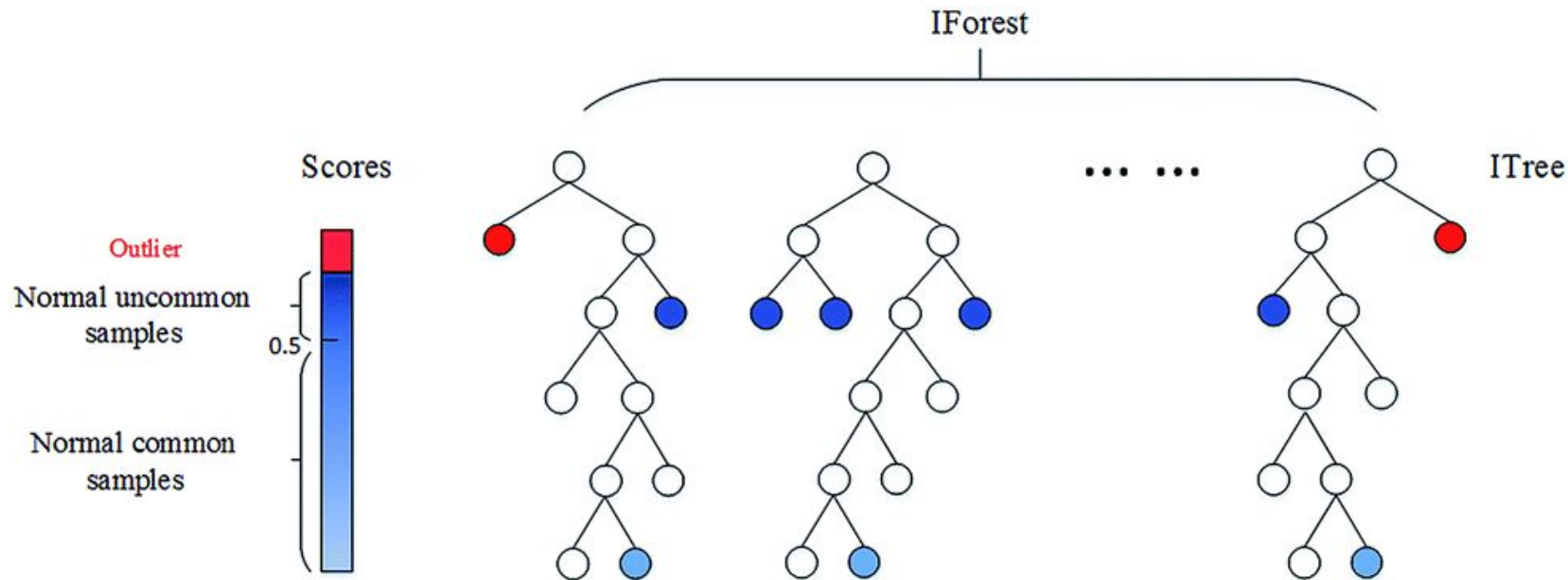
Isolation tree



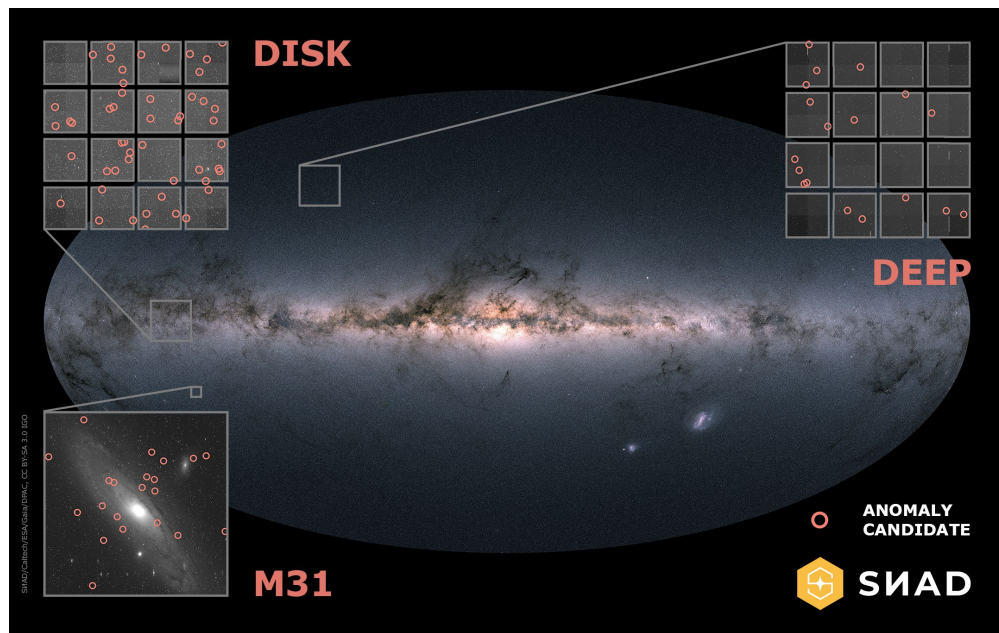
Plot from <https://donghwa-kim.github.io/iforest.html>

Example of an automatic search for anomalies,

Isolation forest



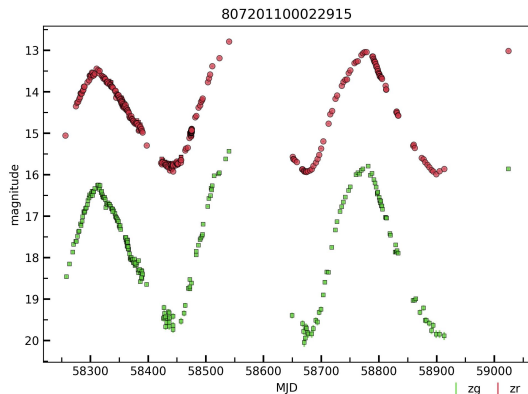
Zwicky Transient Facility DR3



- Survey currently in operation, telescope in California
- 3 fields from Data Release 3 (DR3)

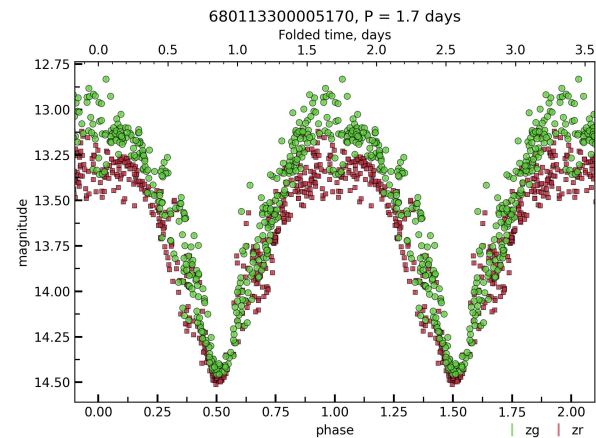
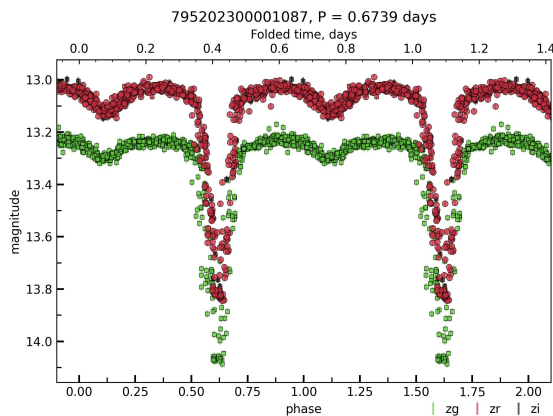
*After selection cuts and feature extraction, **2.25 million objects***

Zwicky Transient Facility DR3



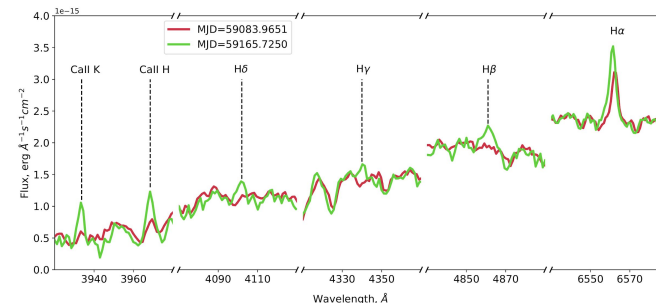
expected to contain stars
and periodic variables
(no transients)

Visualization generated with the SNAD ZTF viewer: <https://ztf.snad.space/>

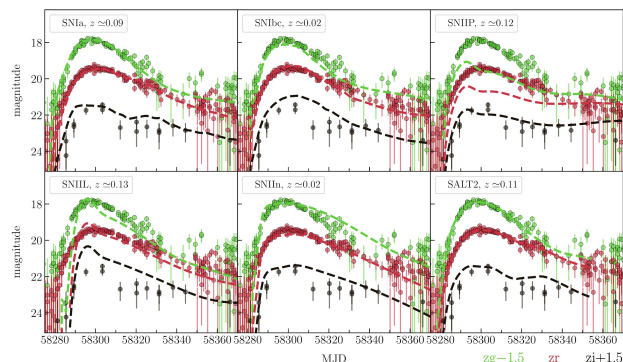


Zwicky Transient Facility DR3

- Feature extraction
- Anomaly detection algorithms:
 - *Isolation Forest*
 - *Local Outlier Factor*
 - *Gaussian Mixture Model*
 - *One-Class Support Vector Machine*
- Initial data: 2.25 million objects
- Expert analysis: 277 objects



- 1 RS Canum Venaticorum star
- 1 red dwarf flare
- 4 Supernova candidates

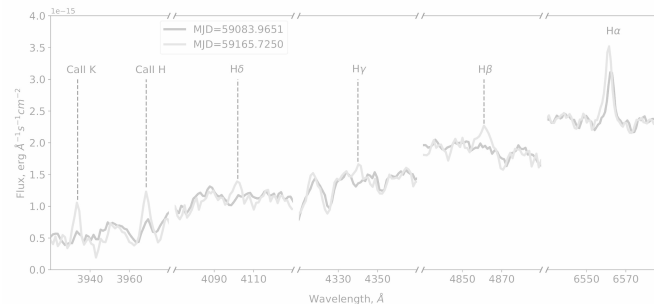


Results:

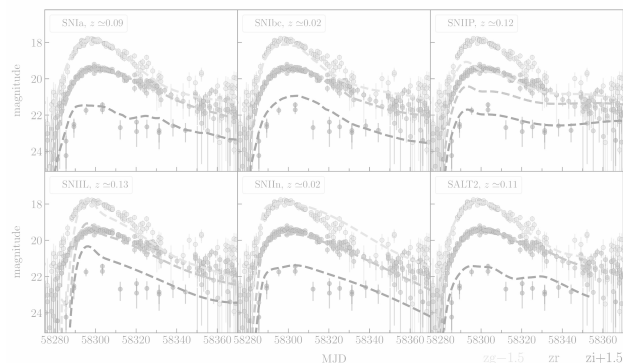
- 68 % (188) - artifacts, bogus
- 24 % (66) - previously cataloged
- 8 % (23) - discoveries

Zwicky Transient Facility DR3

- Feature extraction
- Anomaly detection algorithms:
 - *Isolation Forest*
 - *Local Outlier Factor*
 - *Gaussian Mixture Model*
 - *One-Class Support Vector Machine*
- Initial data: 2.25 million objects
- Expert analysis: 277 objects



- 1 RS Canum Venaticorum star
- 1 red dwarf flare
- 4 Supernova candidates



Results:

- **68 % (188) - artifacts, bogus**
- 24 % (66) - previously cataloged
- 8 % (23) - discoveries

It is about Discovery

“An anomaly is an observation which deviates so much from the other observations as to arouse suspicions that it was generated by a different mechanism”

Stages of discovery in astronomy:

- Detection
- Interpretation
- Understanding
- Acceptance

Which mechanism?
Is it something we are familiar with but fail to proper model or recognise?
Is it something we have never seen before?
Is there something new for us to Learn?

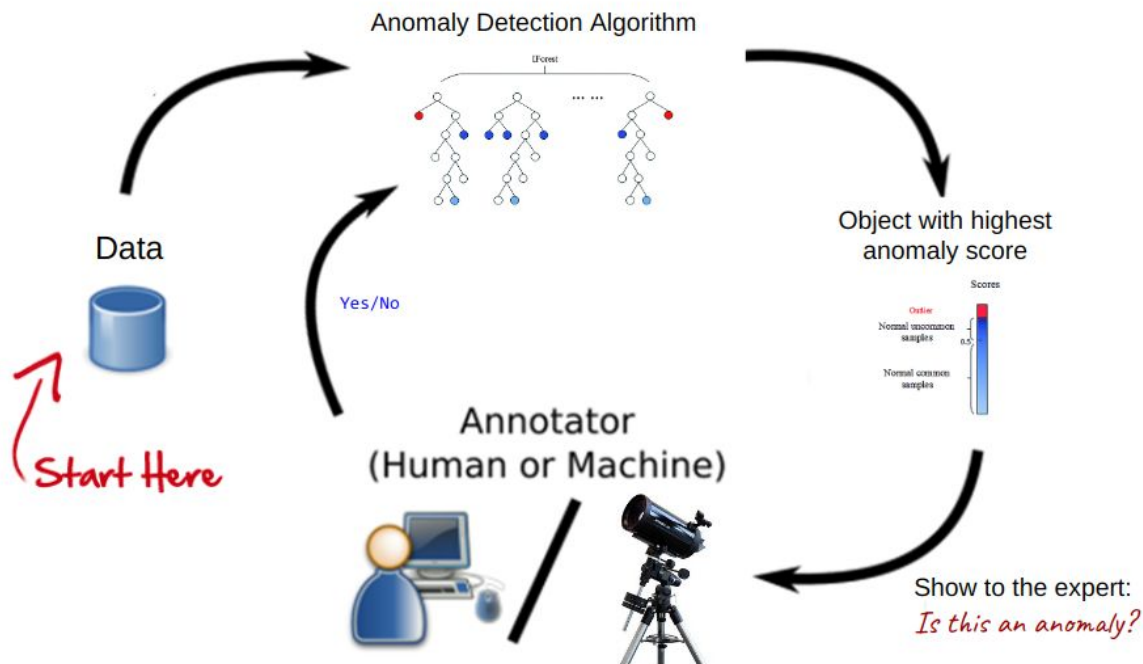


In order to identify the unusual we need to have a clear ideal of what is usual ...

.. and that is a social construct. It changes and adapts with time!



Active Anomaly Detection

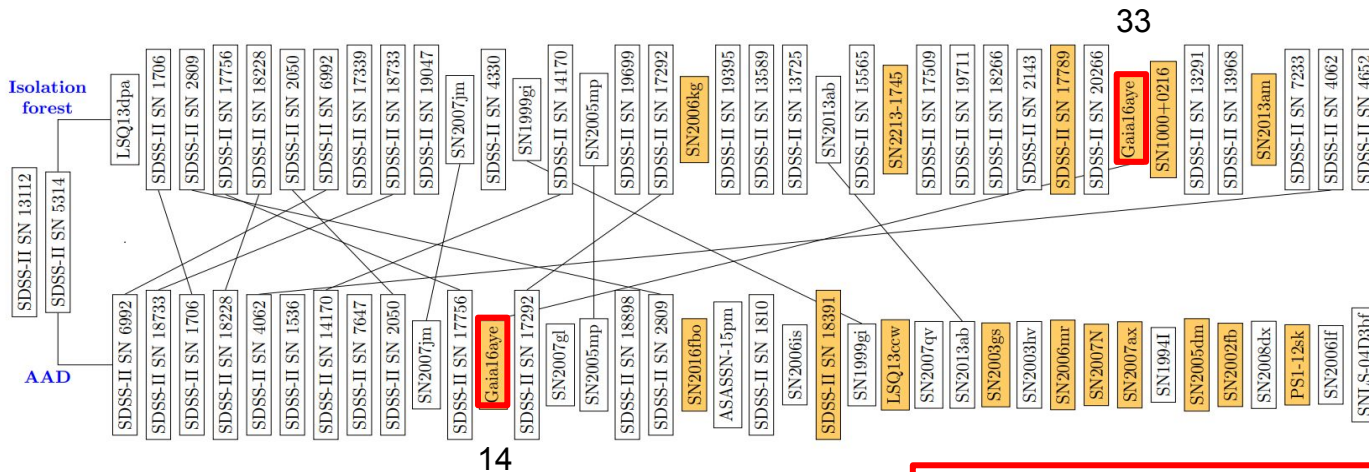


Plot modified from [Chowdhury et al., 2021, SPIE Medical Imaging](#)

Algorithm from Das, S., et al., 2017, in Workshop on Interactive Data Exploration and Analytics (IDEA'17), KDD workshop, [arXiv:cs.LG/1708.09441](#)

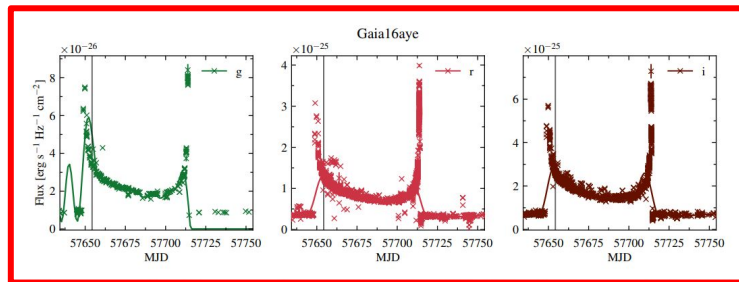
Try the SNAD implementation: <https://coniferest.readthedocs.io/en/latest/tutorial.html>

AAD on real data: The Open Supernova Catalog

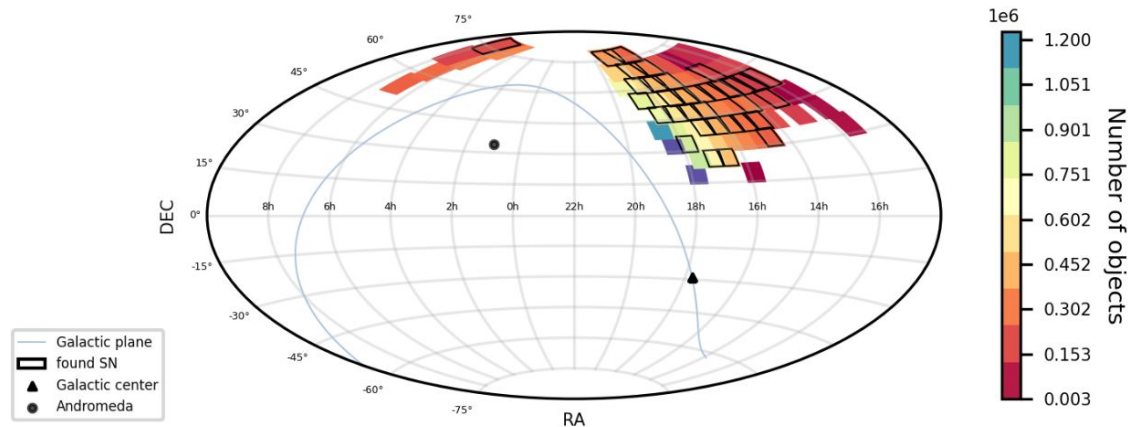


Anomaly

Fast identification of binary microlensing event



AAD on real data: ZTF data releases

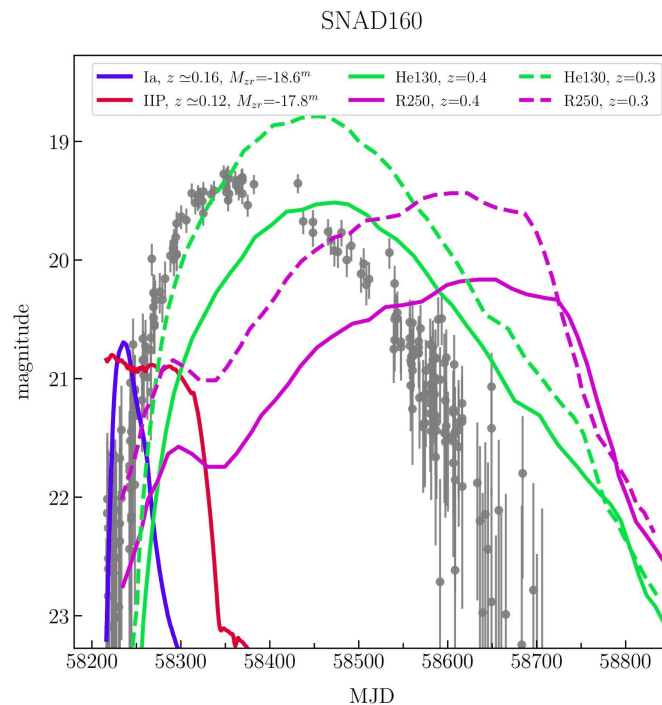
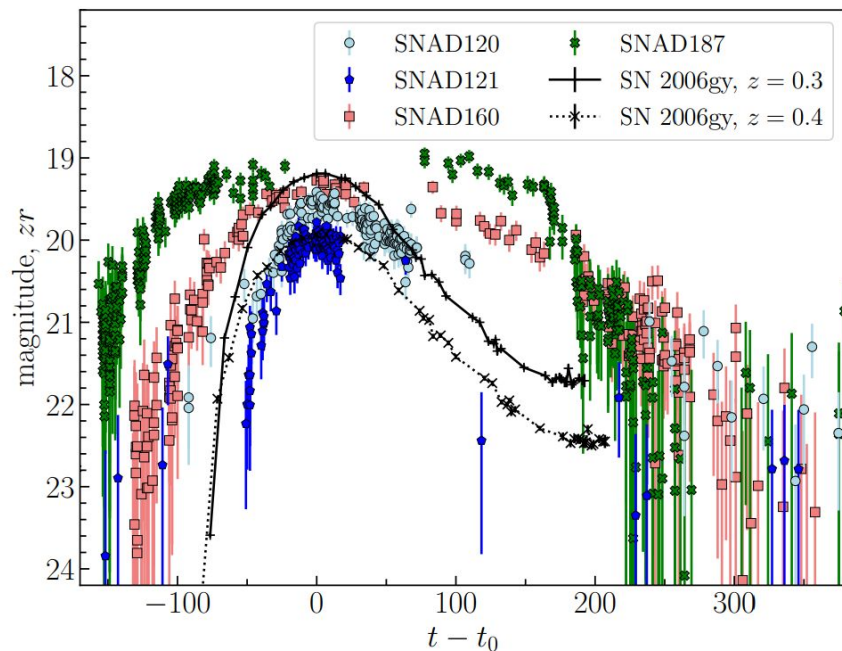


- March - December/2018
- 70 fields
- 30 objects/field
- Total 2100 objects inspected

Found:

- 100 SN-like objects
 - 46 already catalogued
 - 54 newly discovered
- The SNAD catalog: <https://snad.space/catalog/>

Interesting SLSN candidates



AAD is very expensive

Algorithm 2 Active Anomaly Discovery (AAD)

Input: Dataset \mathbf{H} , budget B

Initialize the weights $\mathbf{w}^{(0)} = \{\frac{1}{\sqrt{m}}, \dots, \frac{1}{\sqrt{m}}\}$

Set $t = 0$

Set $\mathbf{H}_A = \mathbf{H}_N = \emptyset$

while $t \leq B$ **do**

$t = t + 1$

 Set $\mathbf{a} = \mathbf{H} \cdot \mathbf{w}$ (i.e., \mathbf{a} is the vector of anomaly scores)

 Let \mathbf{z}_i = instance with highest anomaly score (where $i = \arg \max_i (a_i)$)

 Get feedback {'anomaly'/'nominal'} on \mathbf{z}_i

if \mathbf{z}_i is anomaly **then**

$\mathbf{H}_A = \{\mathbf{z}_i\} \cup \mathbf{H}_A$

else

$\mathbf{H}_N = \{\mathbf{z}_i\} \cup \mathbf{H}_N$

end if

15: $\mathbf{w}^{(t)}$ = compute new weights; normalize $\|\mathbf{w}^{(t)}\| = 1$

end while

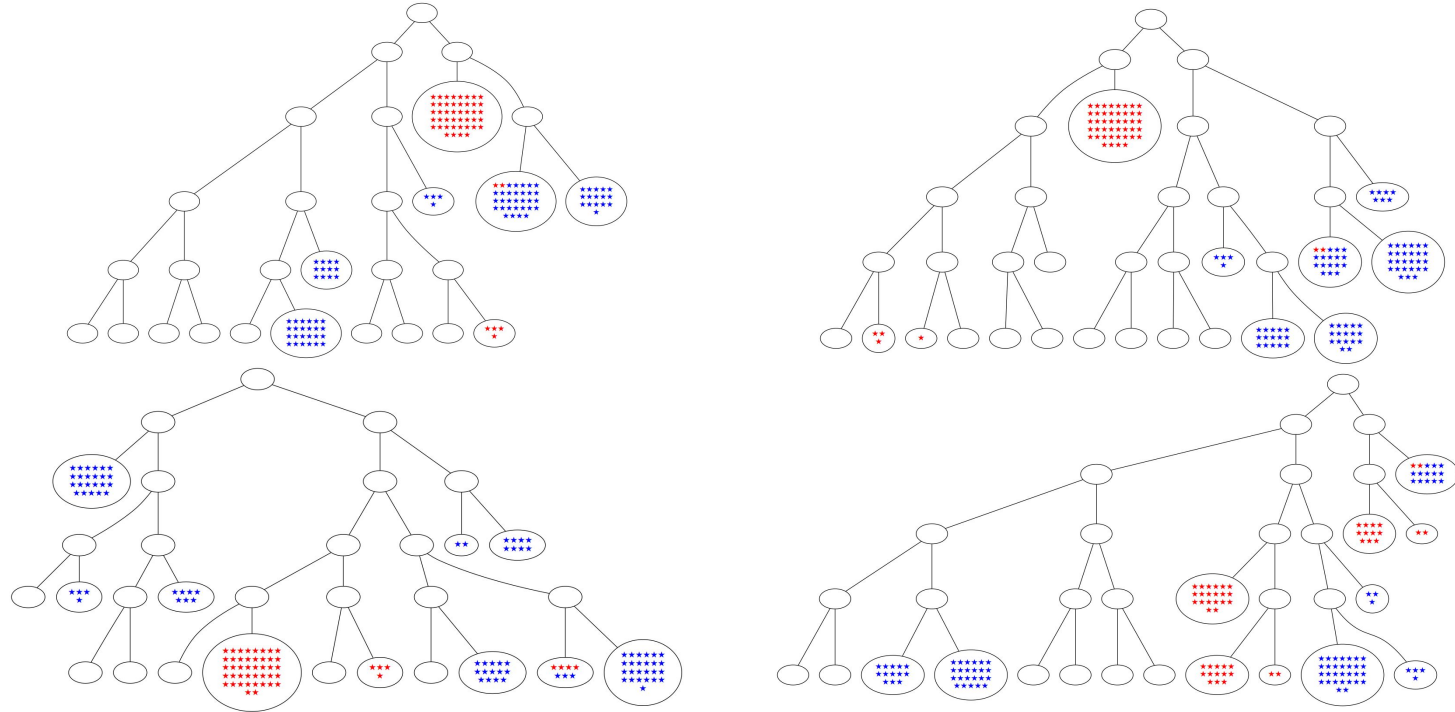
$$l(q, \mathbf{w}; z_i, y_i) = \begin{cases} 0 & \text{if } \mathbf{w} \cdot \mathbf{z}_i \geq q \text{ and } y_i = \text{anomaly} \\ 0 & \text{if } \mathbf{w} \cdot \mathbf{z}_i < q \text{ and } y_i = \text{normal} \\ q - \mathbf{w} \cdot \mathbf{z}_i & \text{if } \mathbf{w} \cdot \mathbf{z}_i < q \text{ and } y_i = \text{anomaly} \\ \mathbf{w} \cdot \mathbf{z}_i - q & \text{if } \mathbf{w} \cdot \mathbf{z}_i \geq q \text{ and } y_i = \text{normal} \end{cases},$$

$$\begin{aligned} \mathbf{w}^{(t)} = \arg \min_{\mathbf{w}, \xi} \frac{C_A}{|\mathbf{H}_A|} & \left(\sum_{z_i \in \mathbf{H}_A} \ell(\hat{q}_\tau(\mathbf{w}^{(t-1)}), \mathbf{w}; (z_i, y_i)) \right) \\ & + \frac{1}{|\mathbf{H}_N|} \left(\sum_{z_i \in \mathbf{H}_N} \ell(\hat{q}_\tau(\mathbf{w}^{(t-1)}), \mathbf{w}; (z_i, y_i)) \right) \\ & + \frac{C_\xi}{|\mathbf{H}_A|} \left(\sum_{z_i \in \mathbf{H}_A} \ell(\mathbf{z}_\tau^{(t-1)} \cdot \mathbf{w}, \mathbf{w}; (z_i, y_i)) \right) \\ & + \frac{C_\xi}{|\mathbf{H}_N|} \left(\sum_{z_i \in \mathbf{H}_N} \ell(\mathbf{z}_\tau^{(t-1)} \cdot \mathbf{w}, \mathbf{w}; (z_i, y_i)) \right) \\ & + \|\mathbf{w} - \mathbf{w}_p\|^2 \end{aligned} \quad (2)$$

Aiming at bigger data

- Requires optimization
- Smooth incorporation of expert knowledge

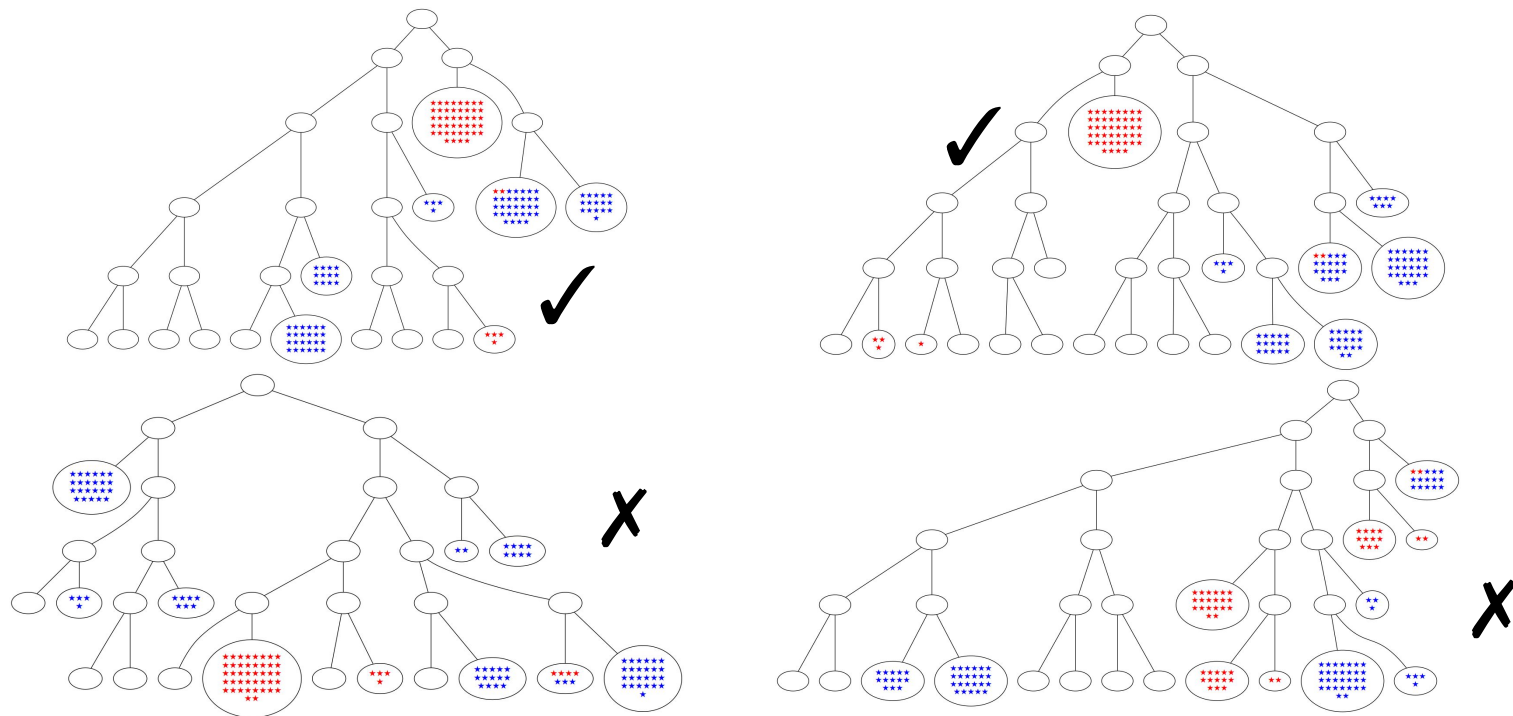
Pine Forest



- Scientifically interesting anomalies
- Statistical outliers

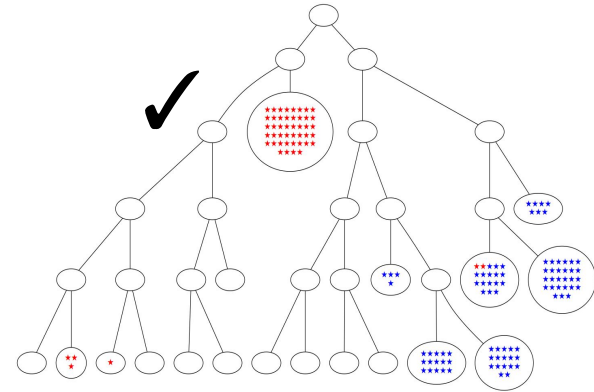
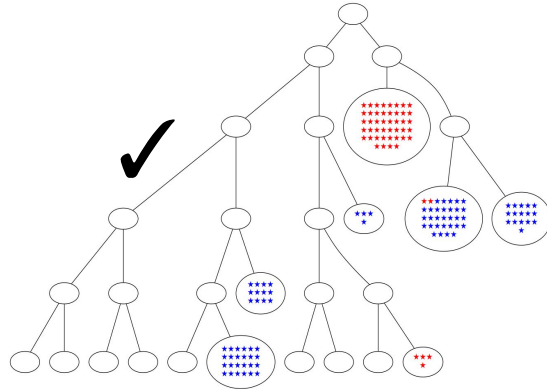
Illustration and algorithm by Vladimir Korolev

Pine Forest



- Scientifically interesting anomalies
- Statistical outliers

Pine Forest



Generate new random trees
and filter again

- Scientifically interesting anomalies
- Statistical outliers

Break

Active Anomaly Detection tutorial

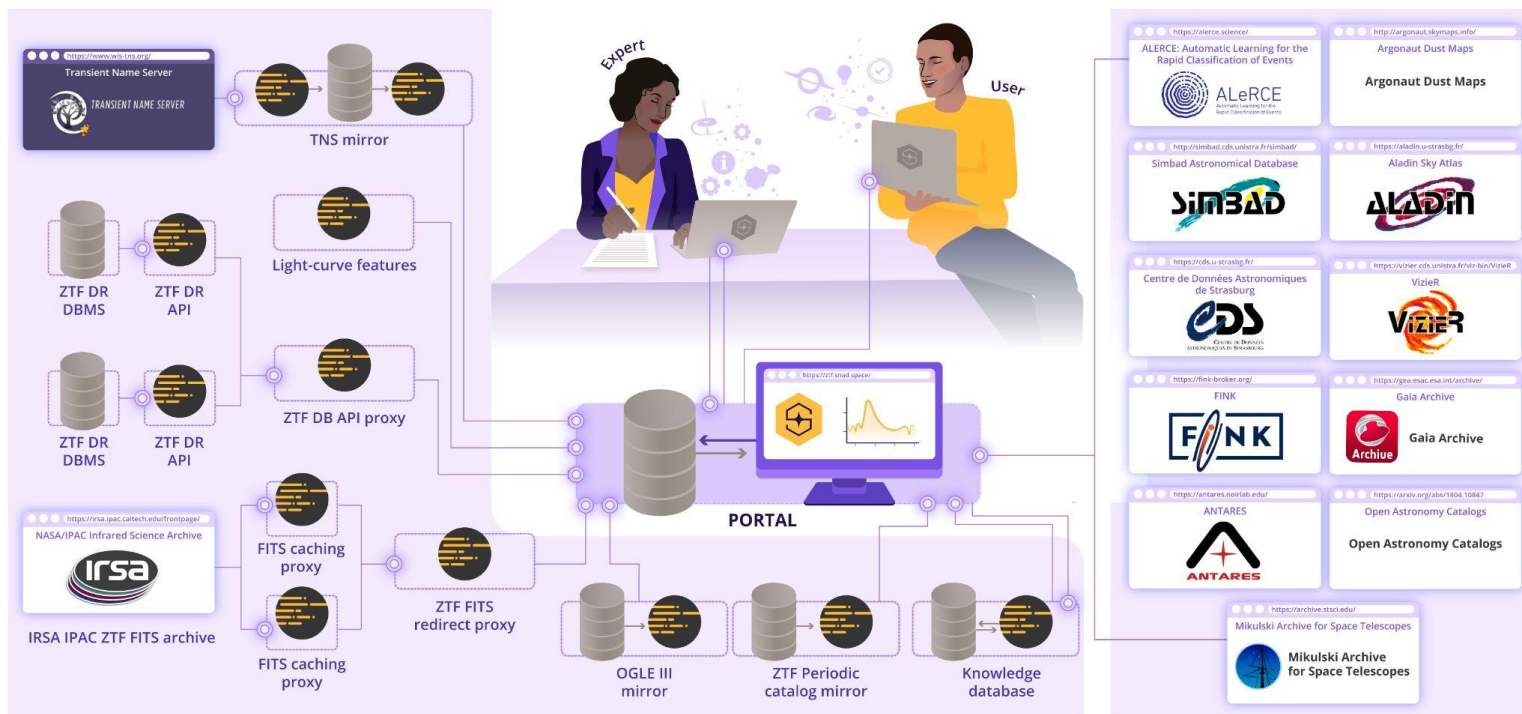
Remember to: File → Save a copy to drive

https://colab.research.google.com/drive/1LvC_a8QE7Q5MECL5u0Adze_xGgRSS4bm?usp=sharing

Active Anomaly Detection tutorial

- The coniferest package: <https://coniferest.readthedocs.io/en/latest/tutorial.html>
- To use with your own data:
 - **Features** (called data in the tutorial) \Rightarrow 2D pandas dataframe
1 line per object, 1 column per feature
 - If you are using light curves, you might want to check the `light_curve` package for optimized feature extraction:
<https://github.com/light-curve/light-curve>
 - **Metadata** \Rightarrow 2D array
1 line per object: it may contain any information that will identify the candidate so you can make a decision.

The SNAD viewer



The knowledge database

SNAD160 — 821207100004043

- artefact column bright_star cosmic defocusing ghost M31 spike track frame_edge
 VAR transient AGN QSO STAR Galaxy
 SN SNIa CCSN SLSN
 Eclipsing EA EB EW
 Pulsating CEP DCEP L LPV M RR RRAB RSG SR
 DSCT
 Cataclysmic AM N UG UGSS UGZ
 Eruptive INS SDOR TTS YSO M_dwarf_flare
 Rotating BY RSCVn
 uncertain non-catalogued 1-point TNS_candidate

Point tag name to see its description. See instructions and tag editor [here](#)

SNAD160

SUBMIT

RESET

Tags	Description	Changed by	Changed at
SN, uncertain	SNAD160	maria	2021-11-11T11:07:27.753Z
AGN, SN, uncertain	SNAD160	maria	2021-11-09T21:42:36.314Z
AGN, SN, uncertain, non-catalogued		maria	2021-11-08T13:33:49.098Z

More than 2000 objects already tagged by experts

Explore the boundaries of your knowledge

- In the era of Rubin, serendipitous discoveries will not happen
- Domain experts **must be included** in the development of new techniques **from the first stages**. They should supervise the first prototypes.

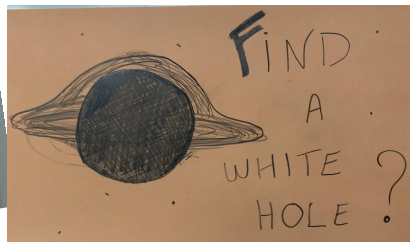
Explore the boundaries of your knowledge

- In the era of Rubin, serendipitous discoveries will not happen
- Domain experts **must be included** in the development of new techniques **from the first stages**. They should supervise the first prototypes.

It is crucial to know what you are looking for

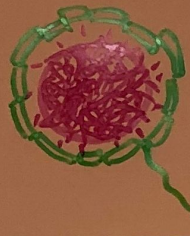
#FinkDreamShots

Build a catalogue
of interstellar asteroids



A classifier
for tidal disruption
events

I WANT TO
FIND A
DYSON
SPHERE



I WANT TO
FIND A LIVE
PISN

~~I wish astronomers~~
~~use REAL units!~~
A switch between
mag. and Lum.

What do you want
to see?

THANK

YOU

