



Machine Learning Tutorial III - Beyond textbook ML

*ADA IX Summer School
20-22 May 2018, Valencia - Spain*

Emille E. O. Ishida

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

In collaboration with
Alexandre Boucaud
Paris-Saclay Center for Data Science



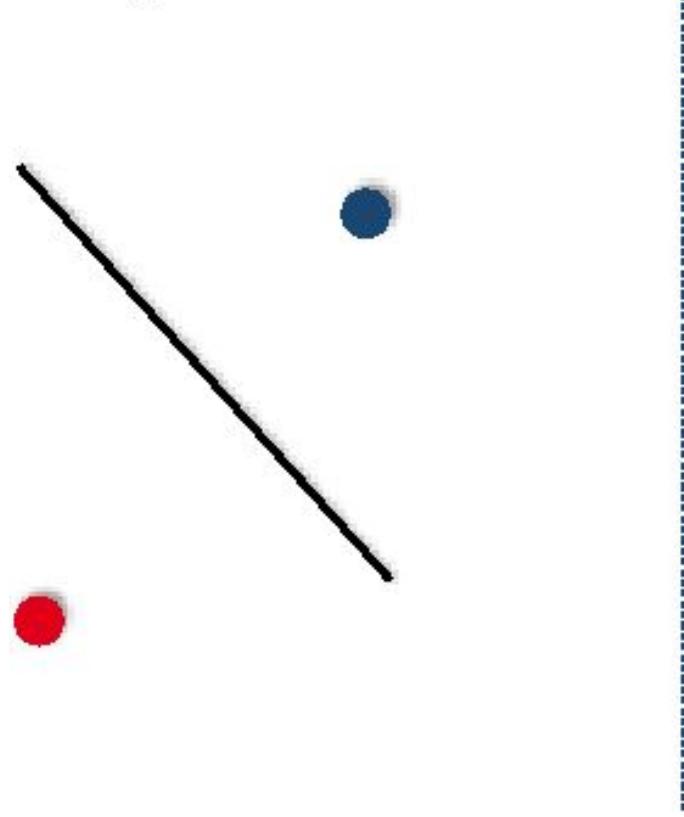
Summary

- I. More than Supervised x Unsupervised
- II. Challenges ahead
- III. Human learning

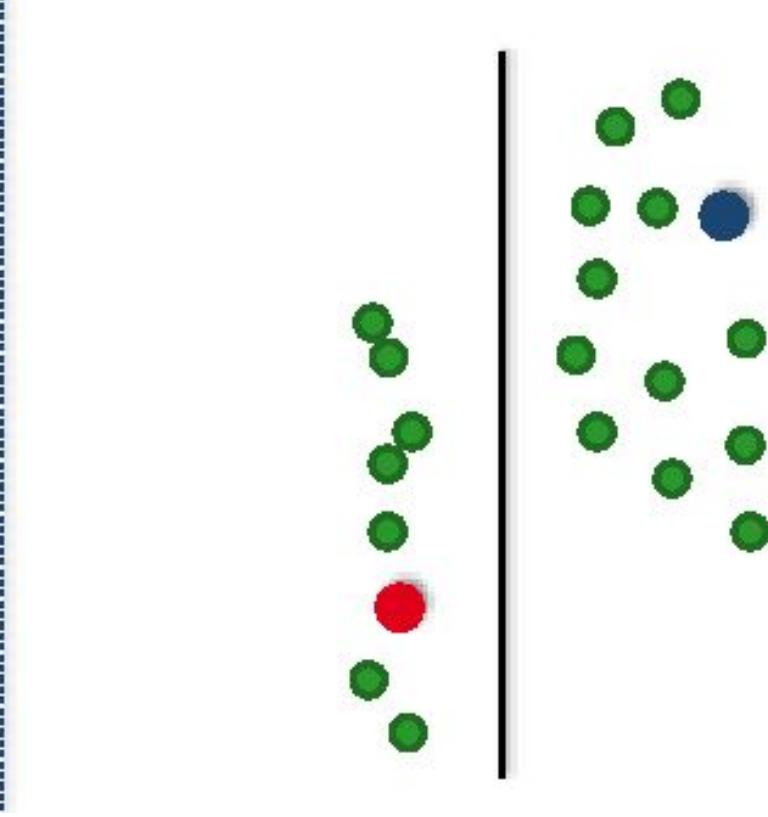
Semi-supervised learning

Getting partial information from the unlabelled sample

only labeled data

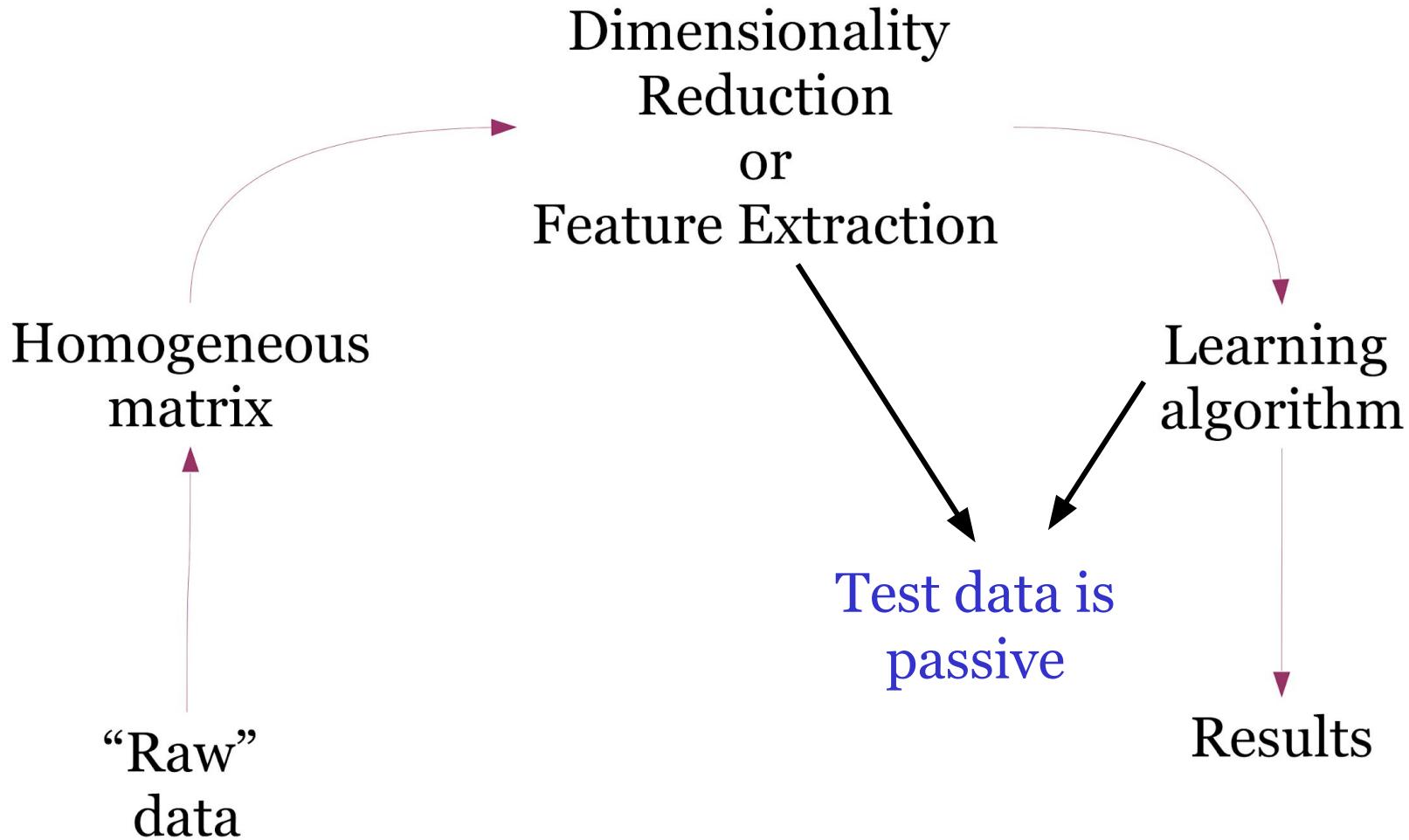


with unlabeled data



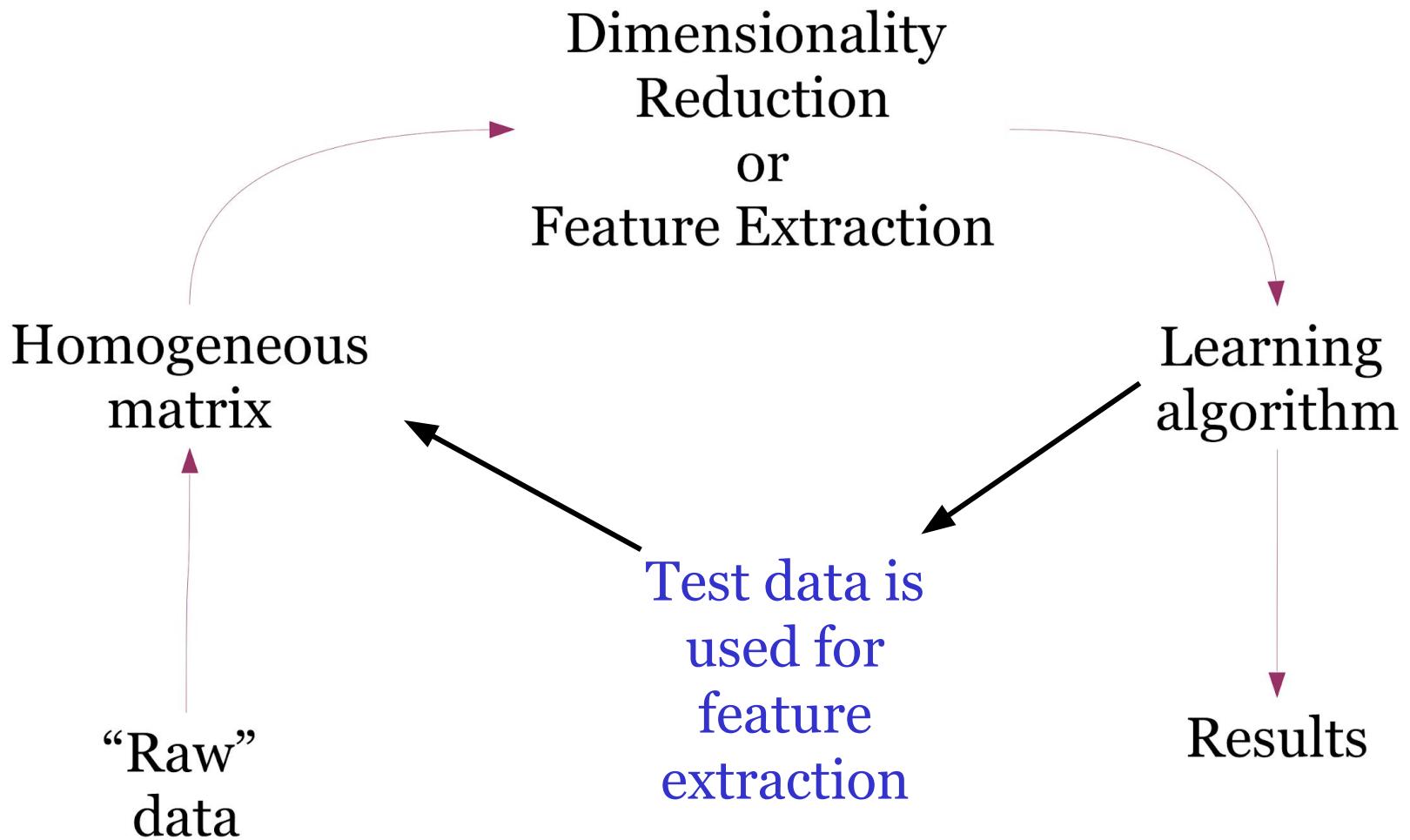
Machine Learning

Work-flow



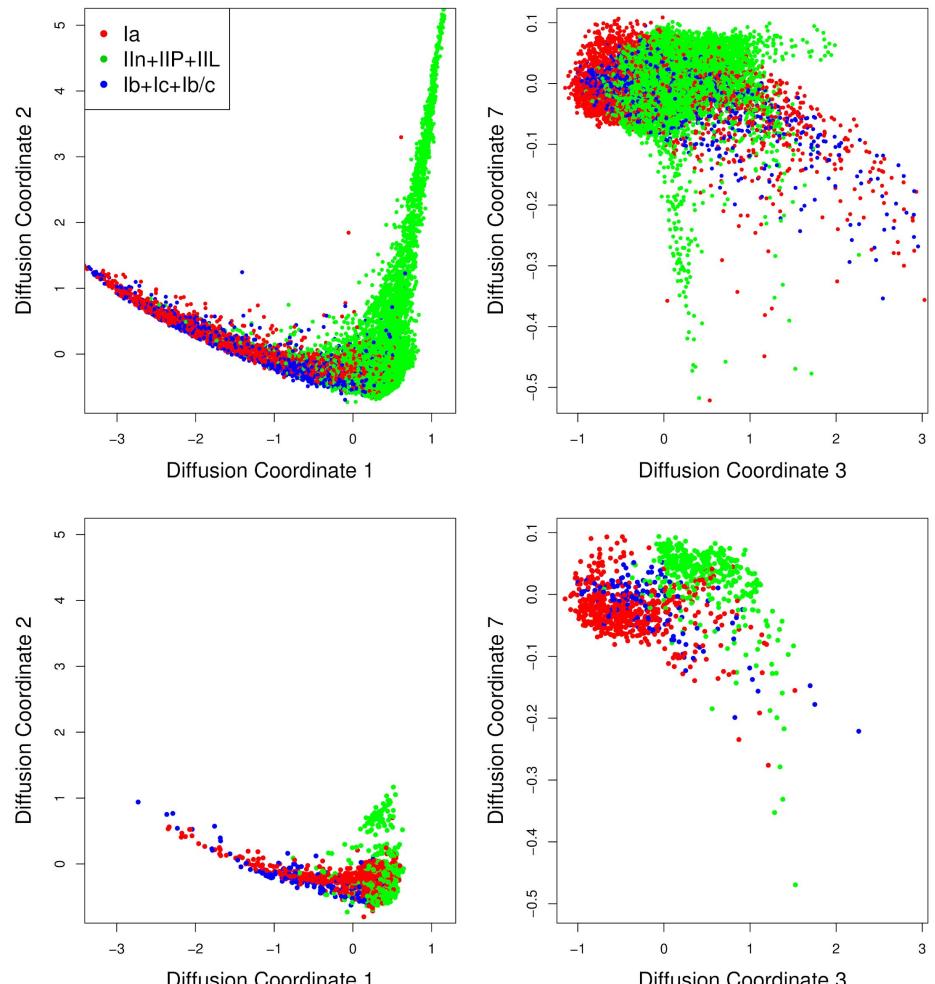
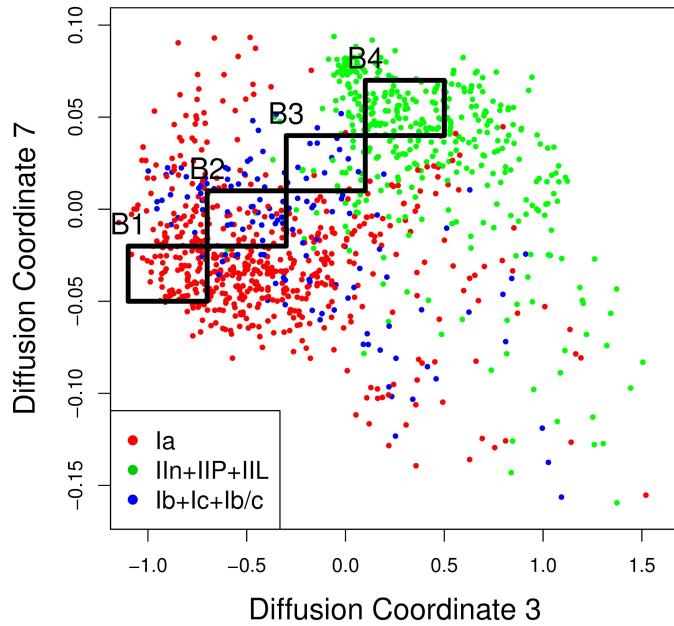
Semi-supervised learning

Work-flow



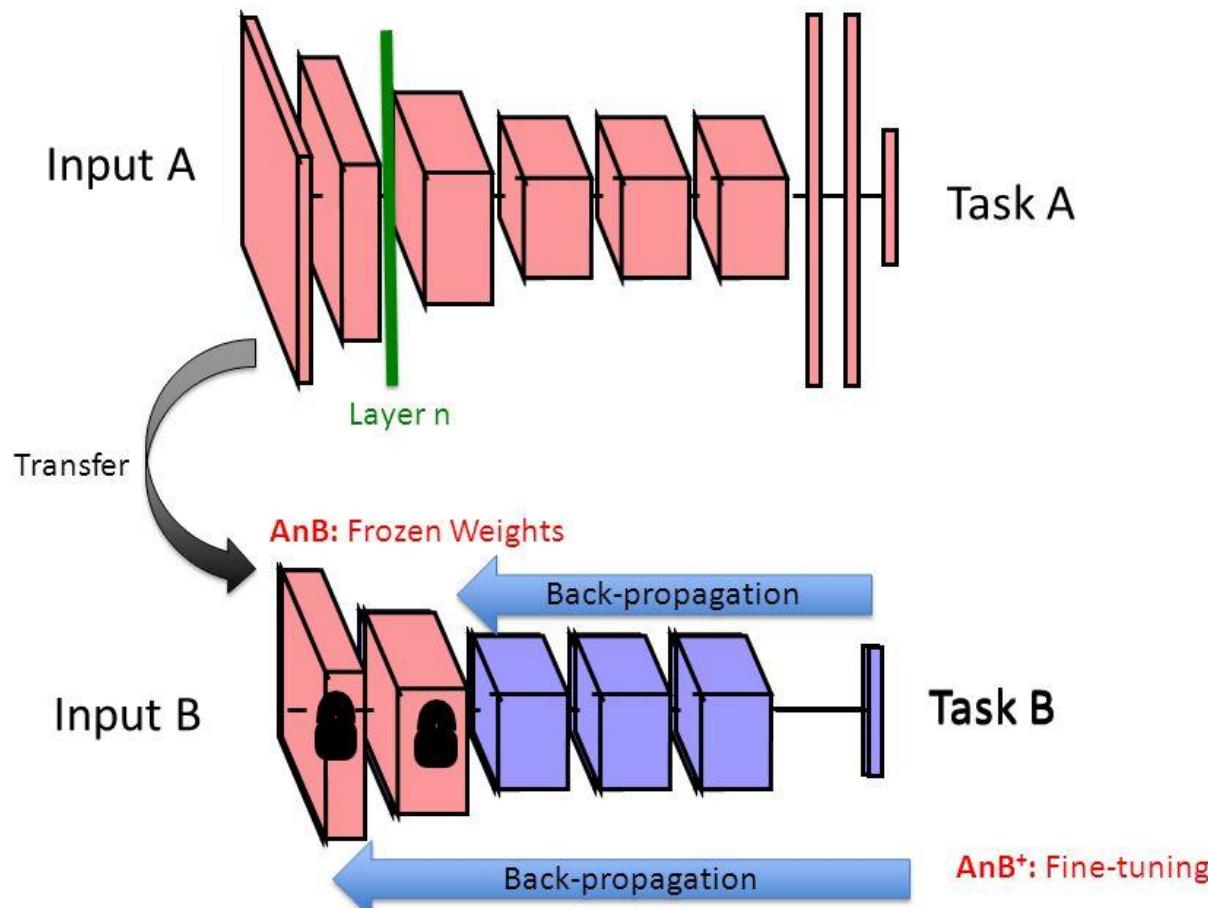
Semi-supervised learning

For Supernova Photometric Classification



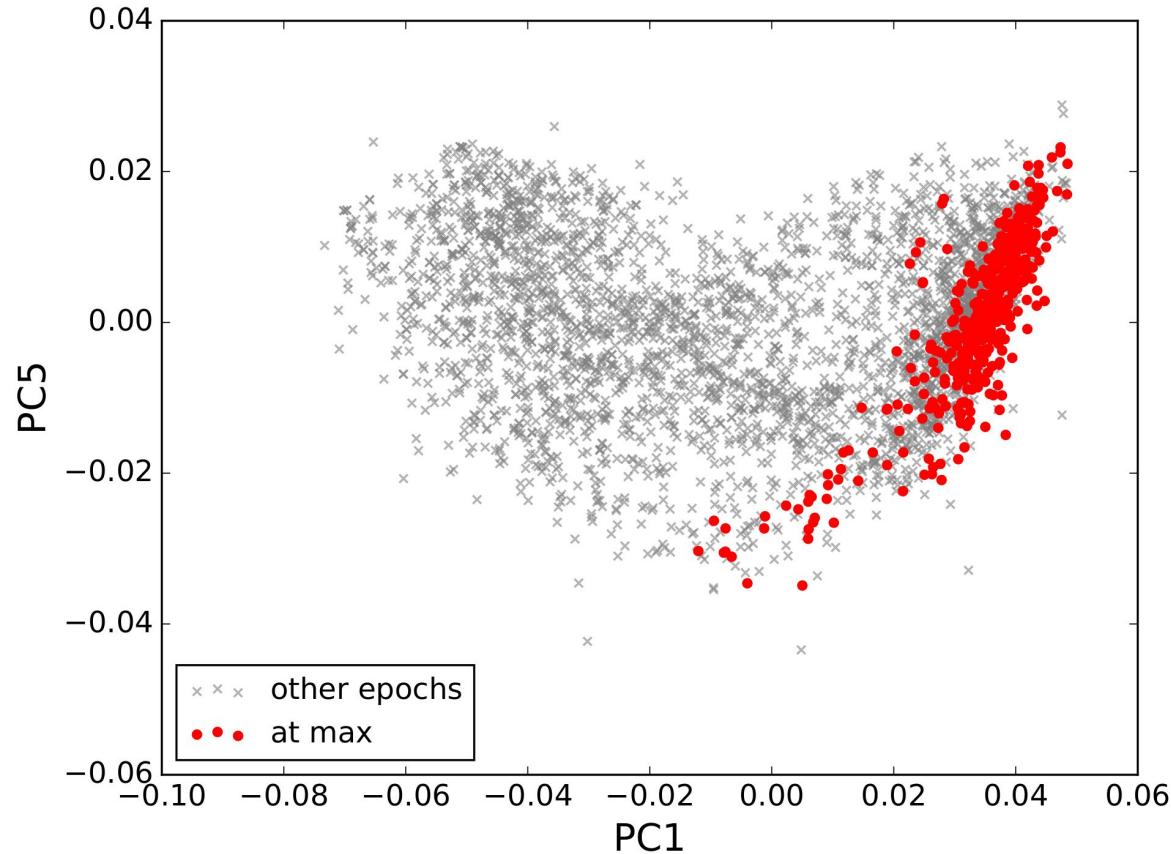
Transfer Learning

Borrowing information from somewhere else



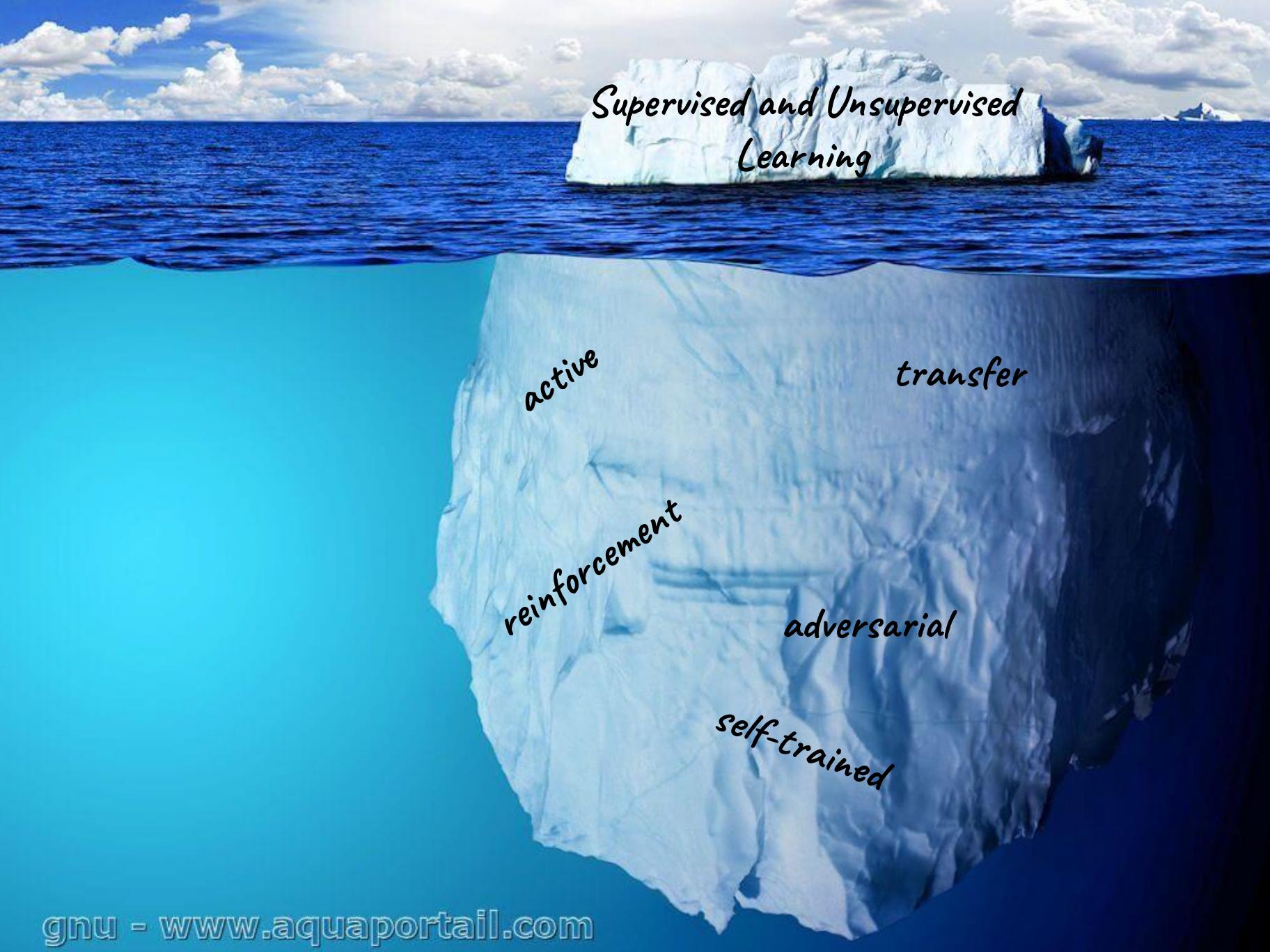
Transfer Learning

In Astronomy





*Supervised and Unsupervised
Learning*



Supervised and Unsupervised Learning

active

transfer

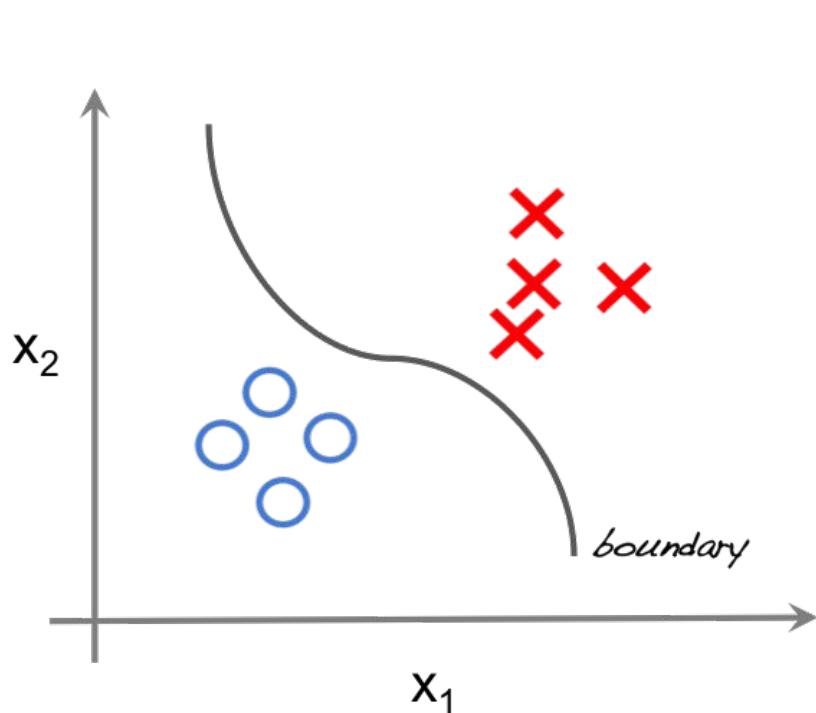
reinforcement

adversarial

self-trained

Categories of Machine Learning:

Supervised x Unsupervised

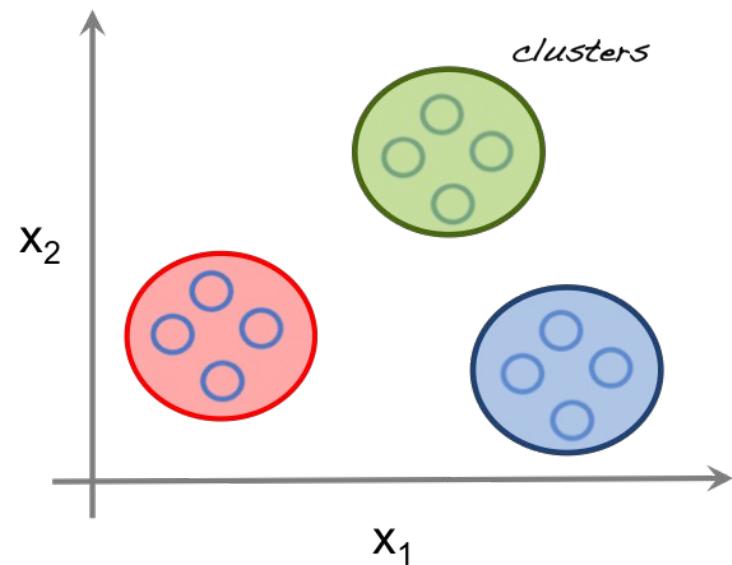


Training sample:

features + labels

Target sample:

features

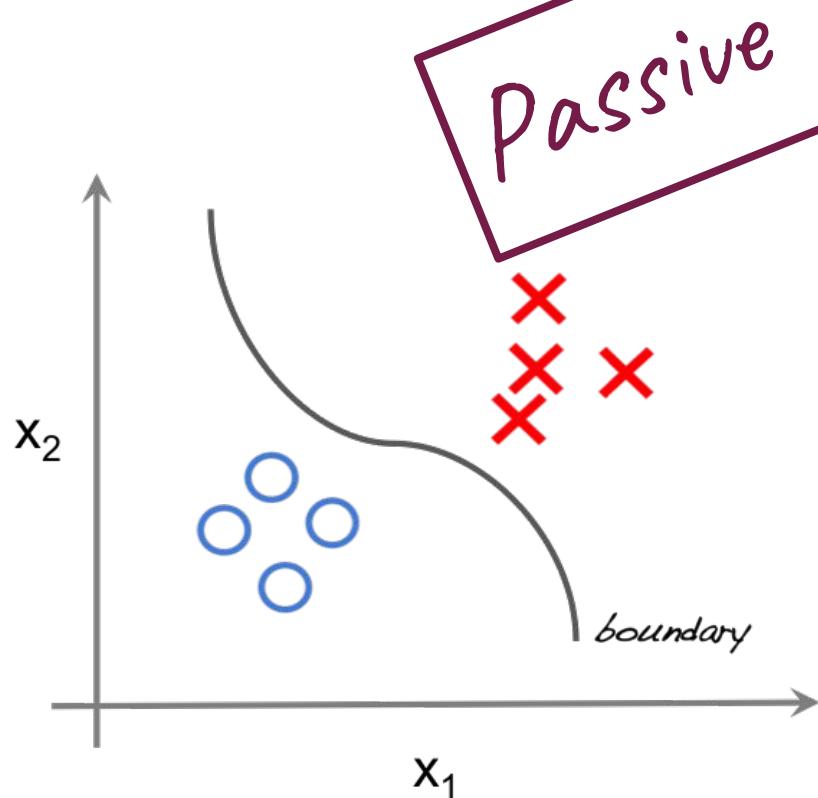


Data sample:

features

Categories of Machine Learning:

Supervised Learning



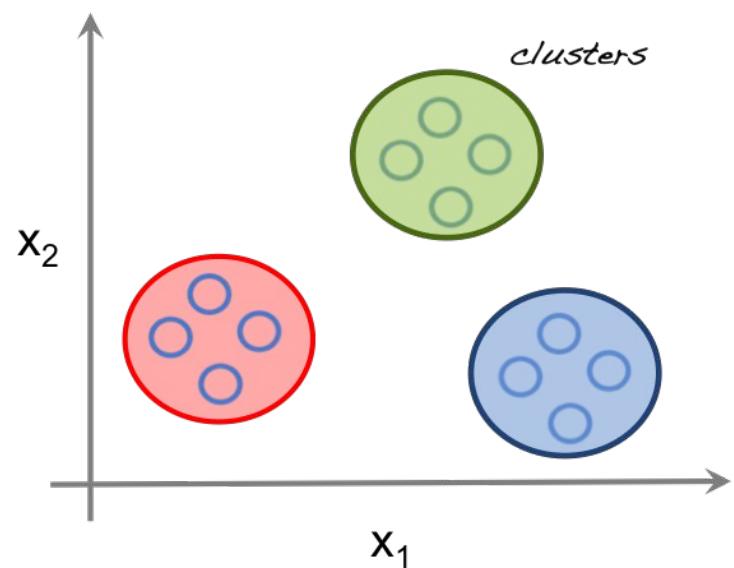
Training sample:

features + labels

Target sample:

features

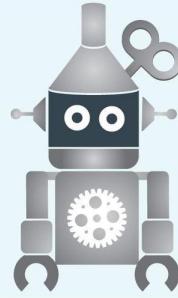
Supervised Learning



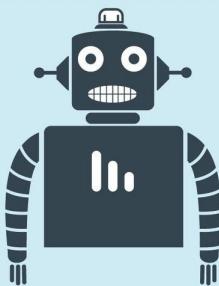
Data sample:

features

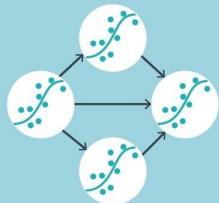
FIRST GENERATION:
Rule-based



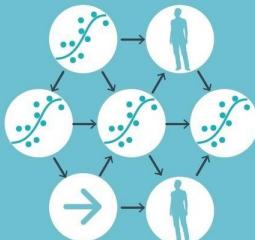
SECOND GENERATION:
Simple machine learning



THIRD GENERATION:
Deep learning

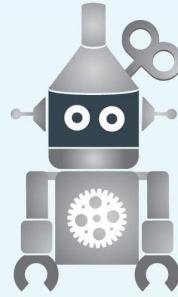


FOURTH GENERATION:
Adaptive learning

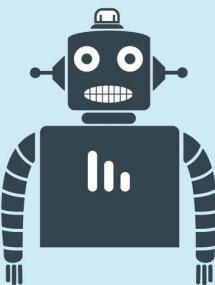


Machines
need to
evolve...

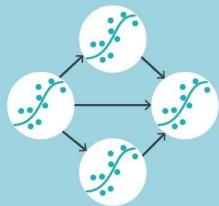
FIRST GENERATION:
Rule-based



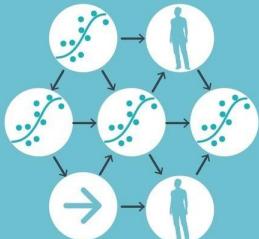
SECOND GENERATION:
Simple machine learning



THIRD GENERATION:
Deep learning



FOURTH GENERATION:
Adaptive learning

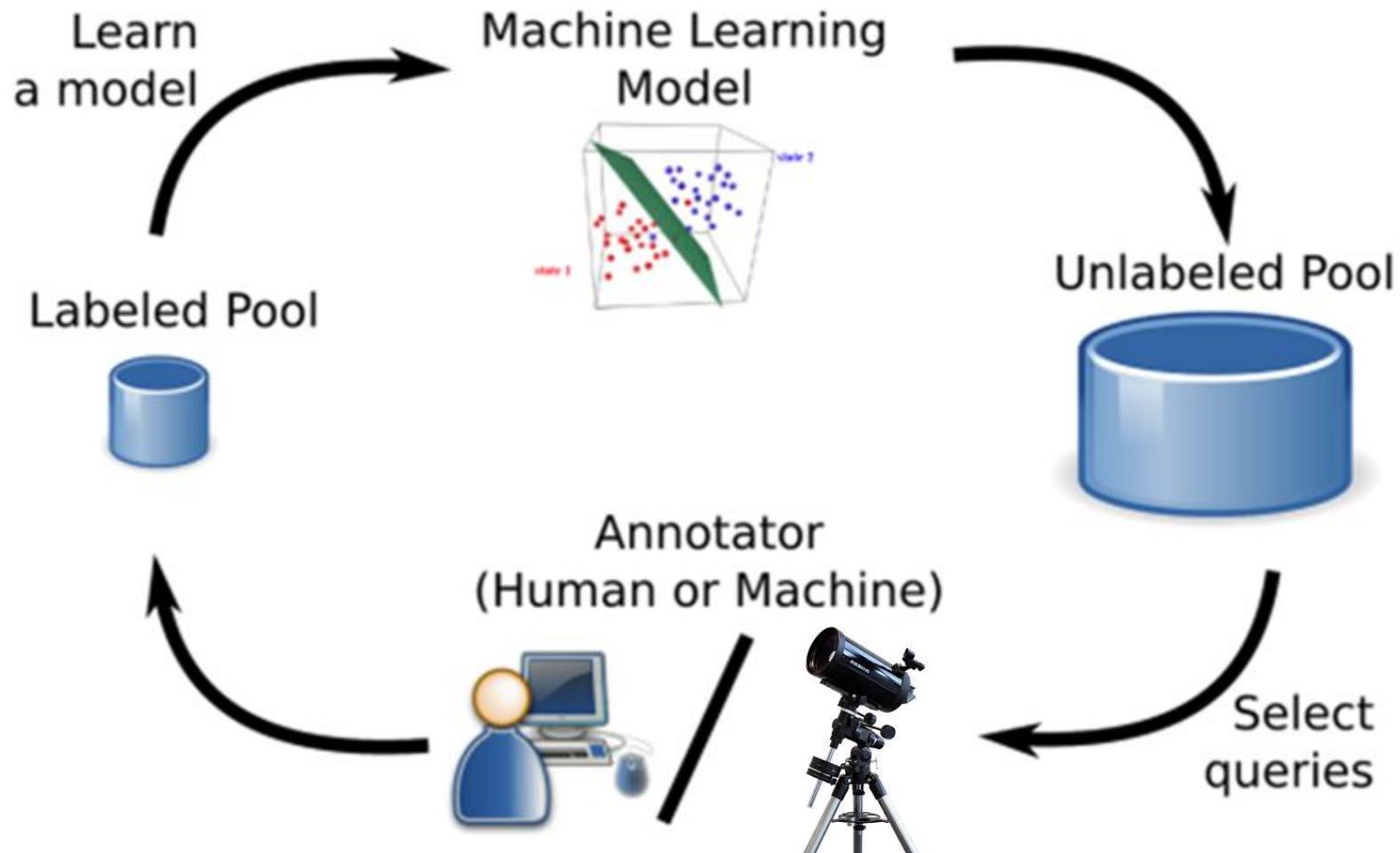


Machines
need to
evolve...

so they
need to
adapt!

Active Learning

Automatizing the construction of the training sample



Active Learning

Automatizing the construction of the training sample

Can machines learn
better, with **fewer**
labelled examples, if
they are carefully
chosen?



Active Learning

Automatizing the construction of the training sample

Can machines learn
better, with **fewer**
labelled examples, if
they are carefully
chosen?

YES!



Active Learning

Automatizing the construction of the training sample



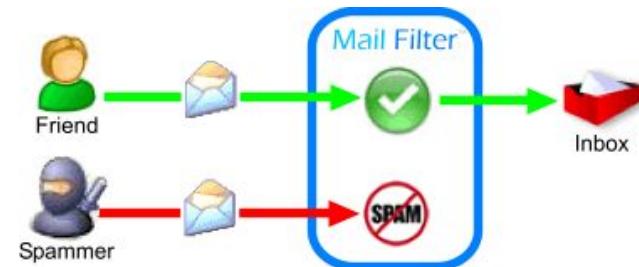
35% OF AMAZON'S REVENUE ARE GENERATED
BY IT'S RECOMMENDATION ENGINE.



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

Can machines learn
better, with **fewer**
labelled examples, if
they are carefully
chosen?

YES!



Active Learning

Automatizing the construction of the training sample



Can machines learn
better, with **fewer**
labelled examples, if
they are carefully
chosen?



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



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

YES!



Active Learning

Automatizing the construction of the training sample



Can machines learn
better, with **fewer**
labelled examples, if
they are carefully
chosen?



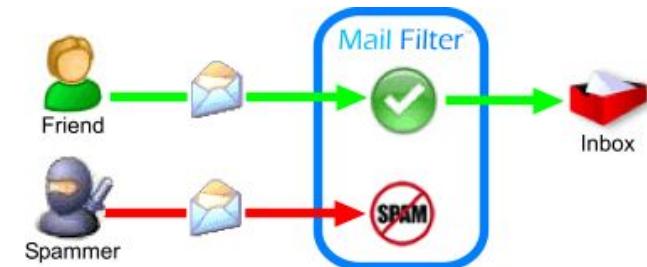
YES!



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

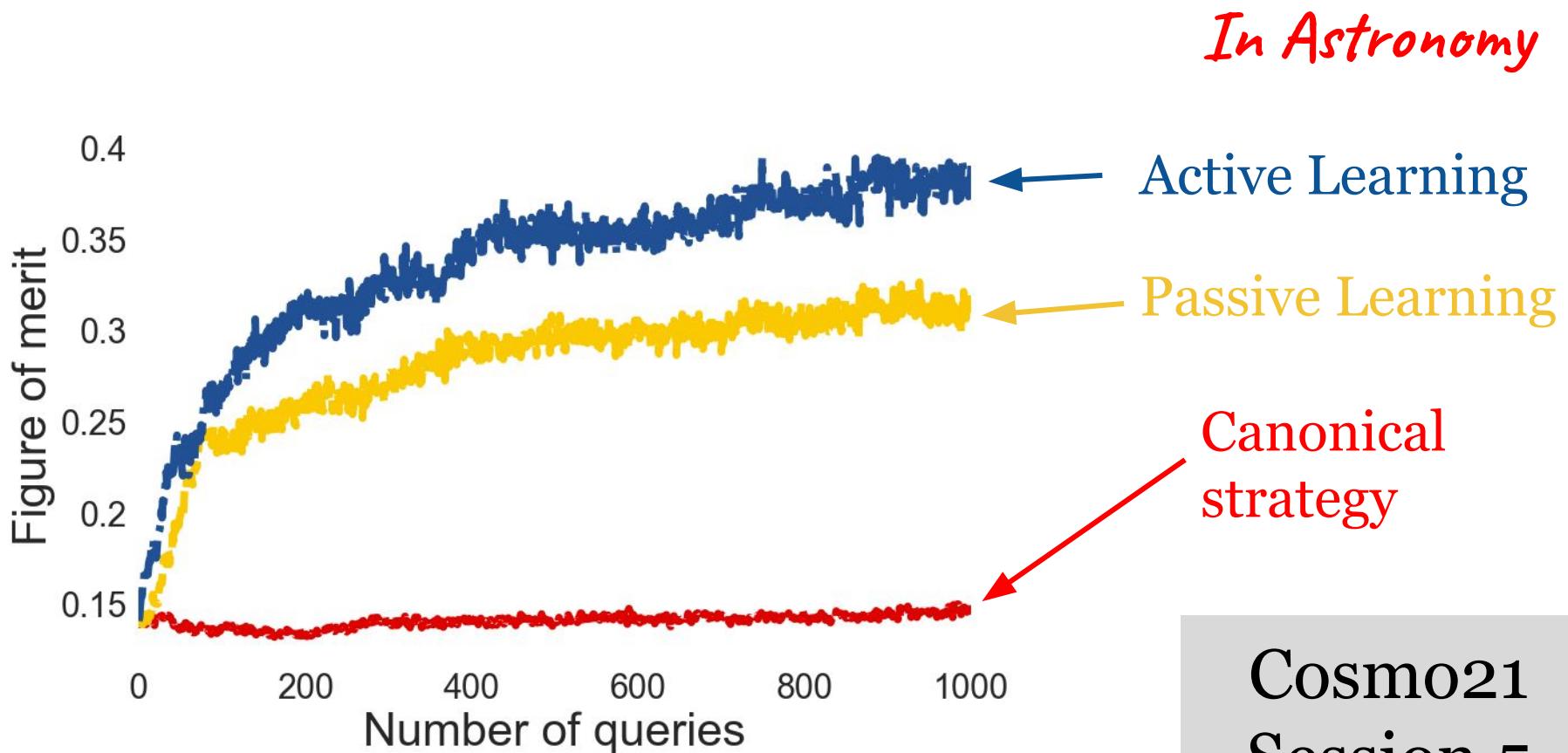


75% OF USERS SELECT MOVIES BASED ON
NETFLIX'S RECOMMENDATIONS.



Active Learning

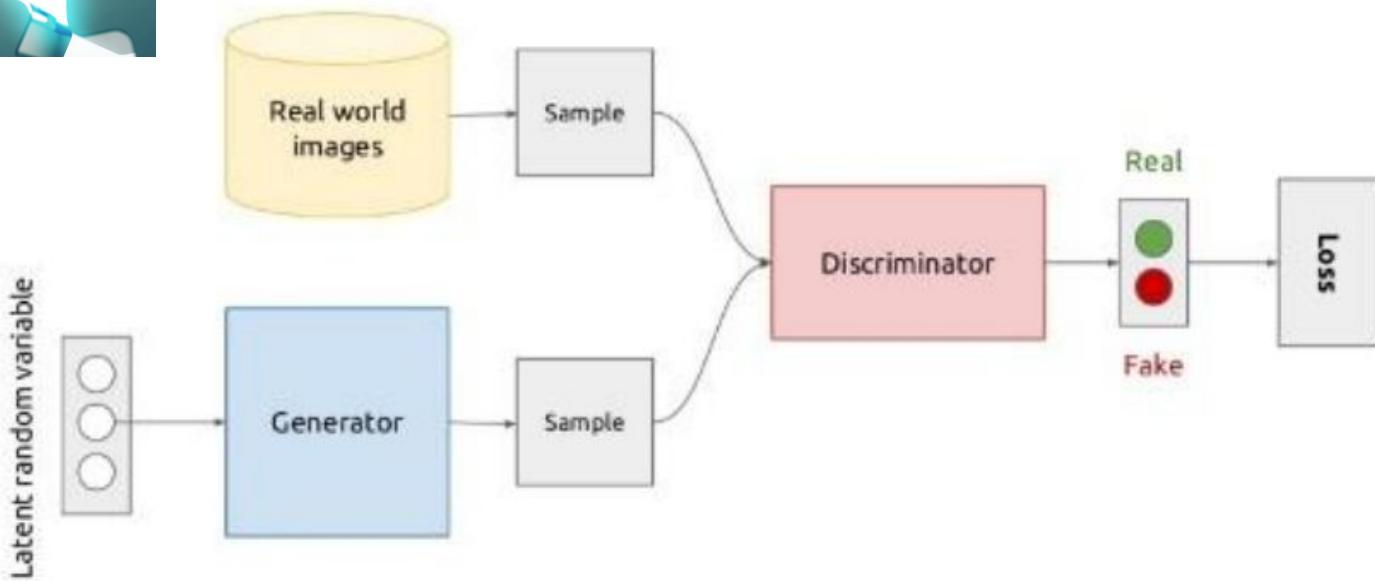
Better performance, fewer training



Cosmo21
Session 5
Wednesday

Adversarial Learning

The benefits of a worthy opponent



<http://www.slideshare.net/xavigiro/deep-learning-for-computer-vision-generative-models-and-adversarial-training-upc-2016>

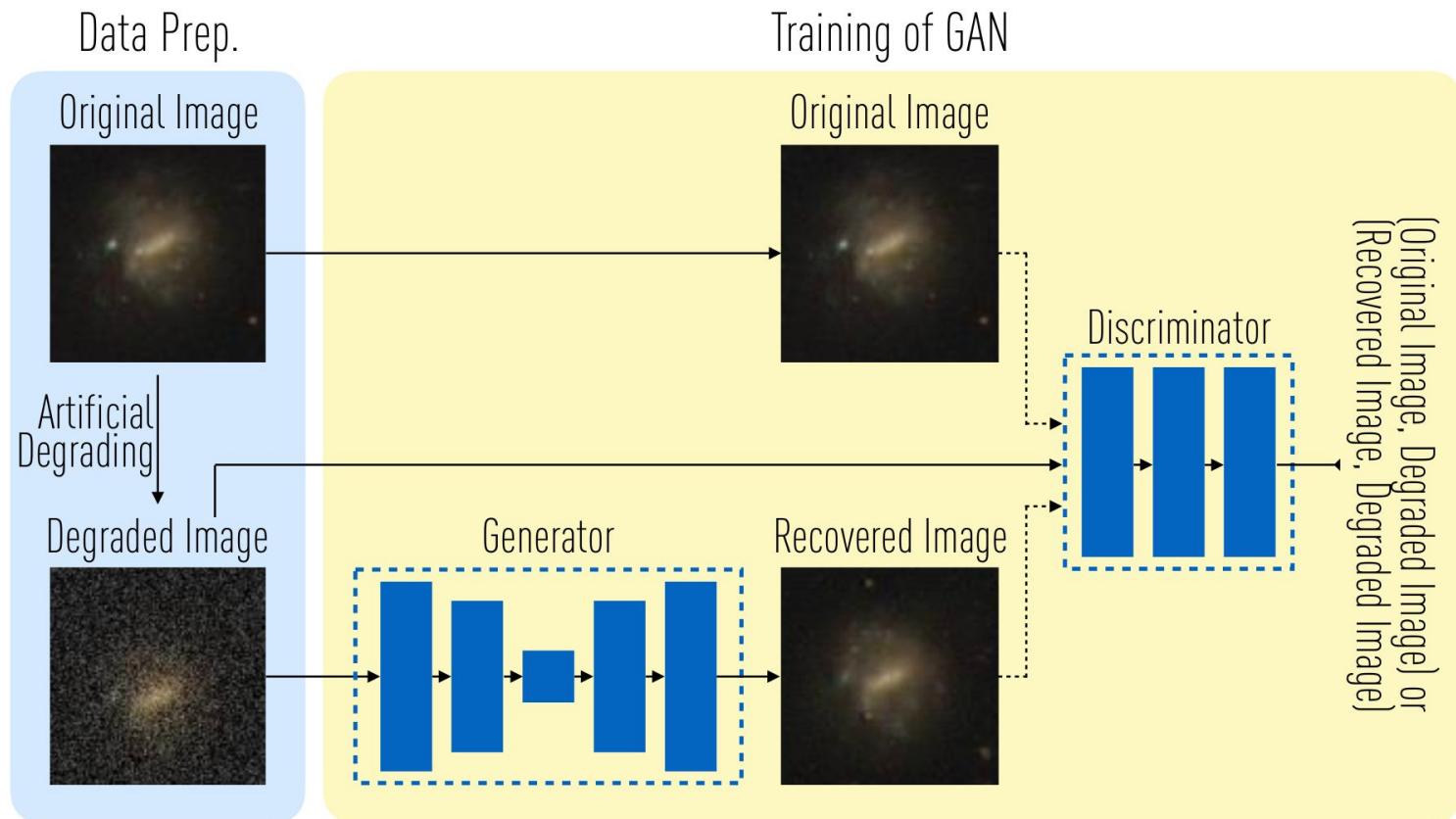
<https://mascherari.press/introduction-to-adversarial-machine-learning/>

Adversarial Learning

The benefits of a worthy opponent

K. Schawinski et al, 2017

In Astronomy

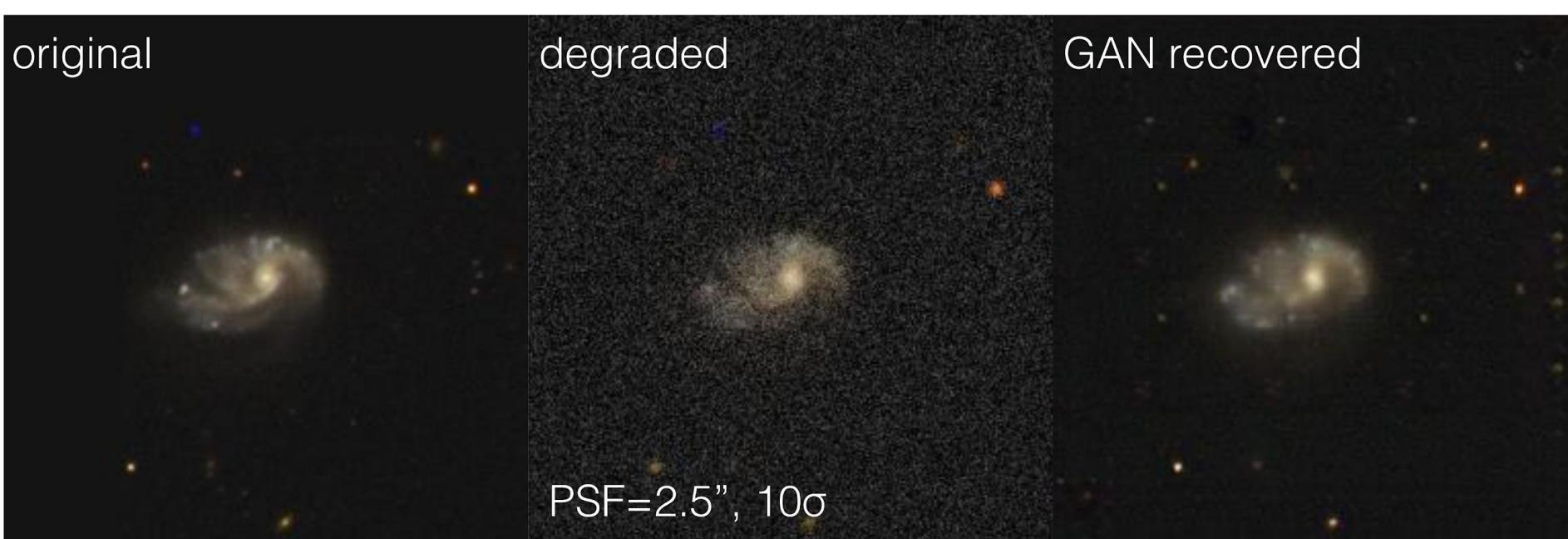


Adversarial Learning

The benefits of a worthy opponent

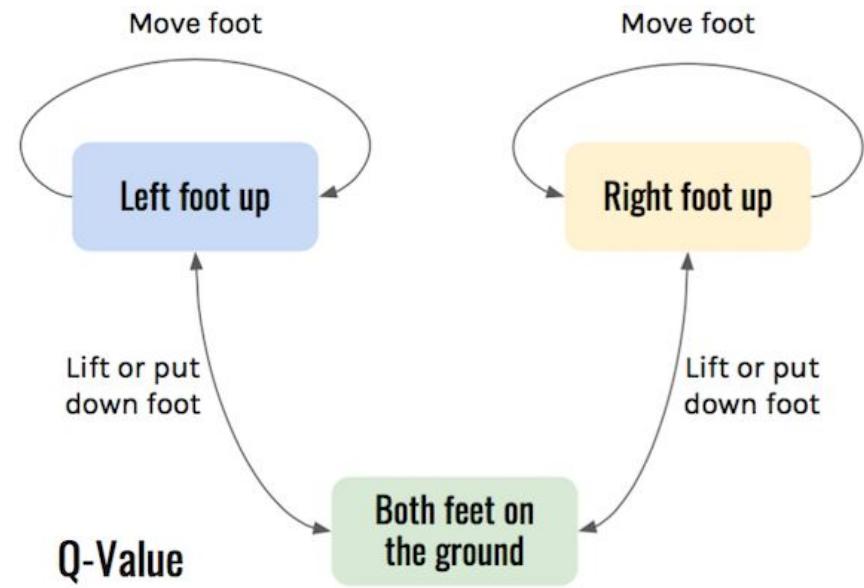
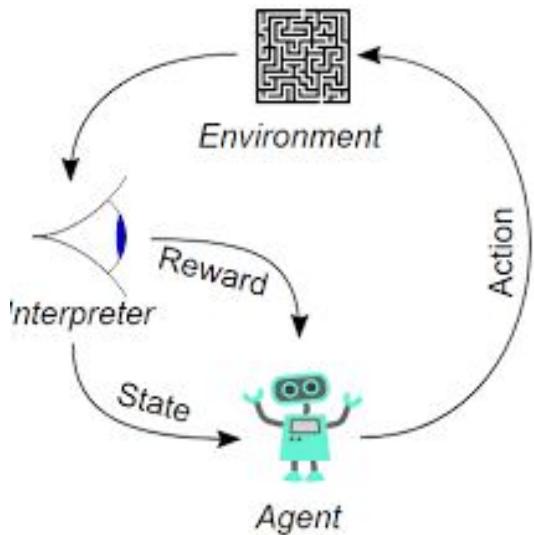
K. Schawinski et al, 2017

In Astronomy



Reinforcement Learning

The importance of feedback



State	Action	Q-Value
Left foot up	Move foot forward	+ 0.5
Left foot up	Put foot down	+ 0.0
Left foot up	Move foot backward	- 0.5
...

This is only the beginning...



... moreover, Astronomical data is tough!

Representativeness

Supervised ML model

data **training**, target

X set of all samples, x

Y set of possible labels, y

h_{train} learner: $y_{est;i} = h_{train}(x_i)$
 L Loss function

Only works if
test is
representative
of training

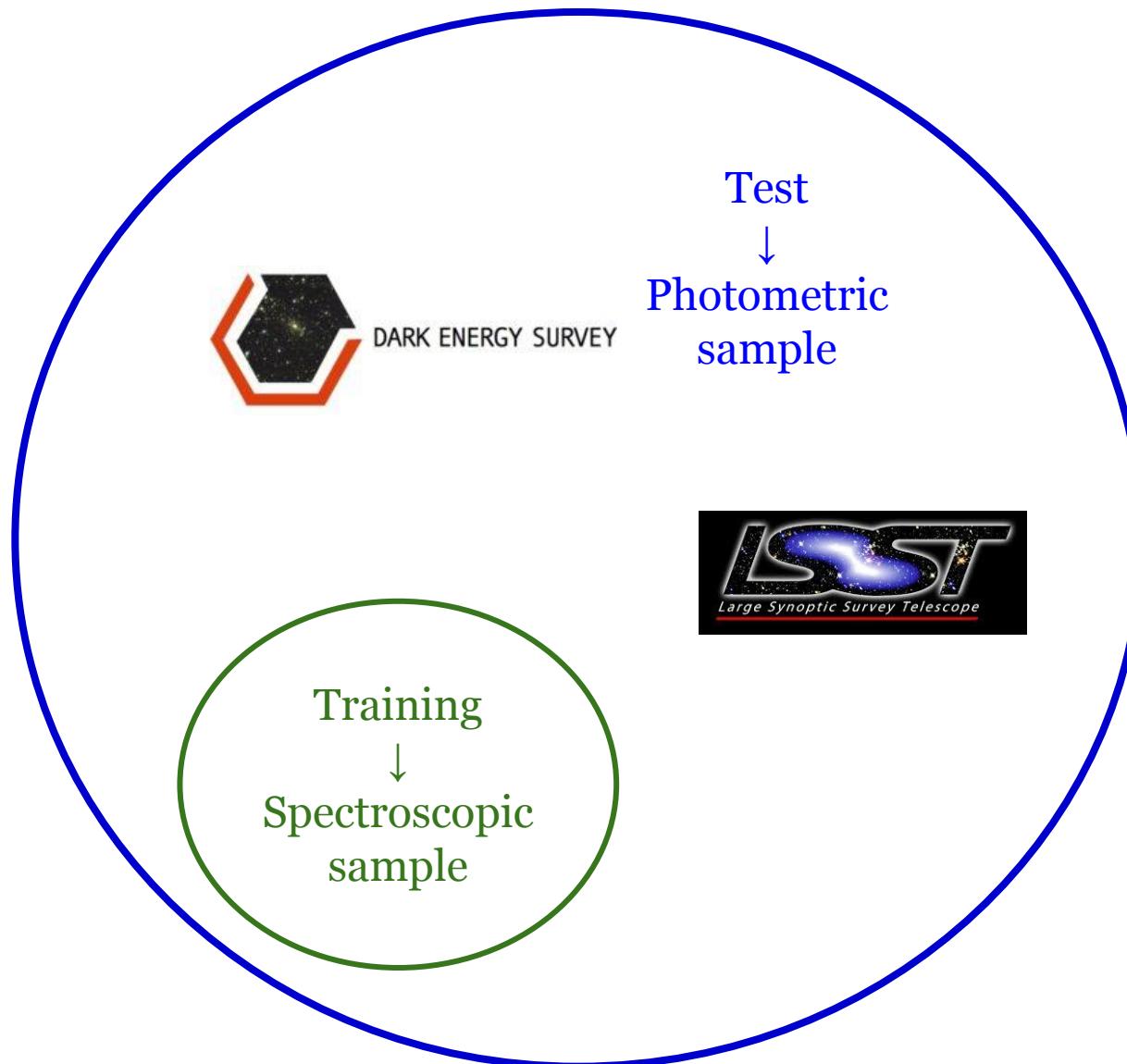
Data generation model:

$$x_i \sim P_X$$

f true labeling function, $y_i = f(x_i)$

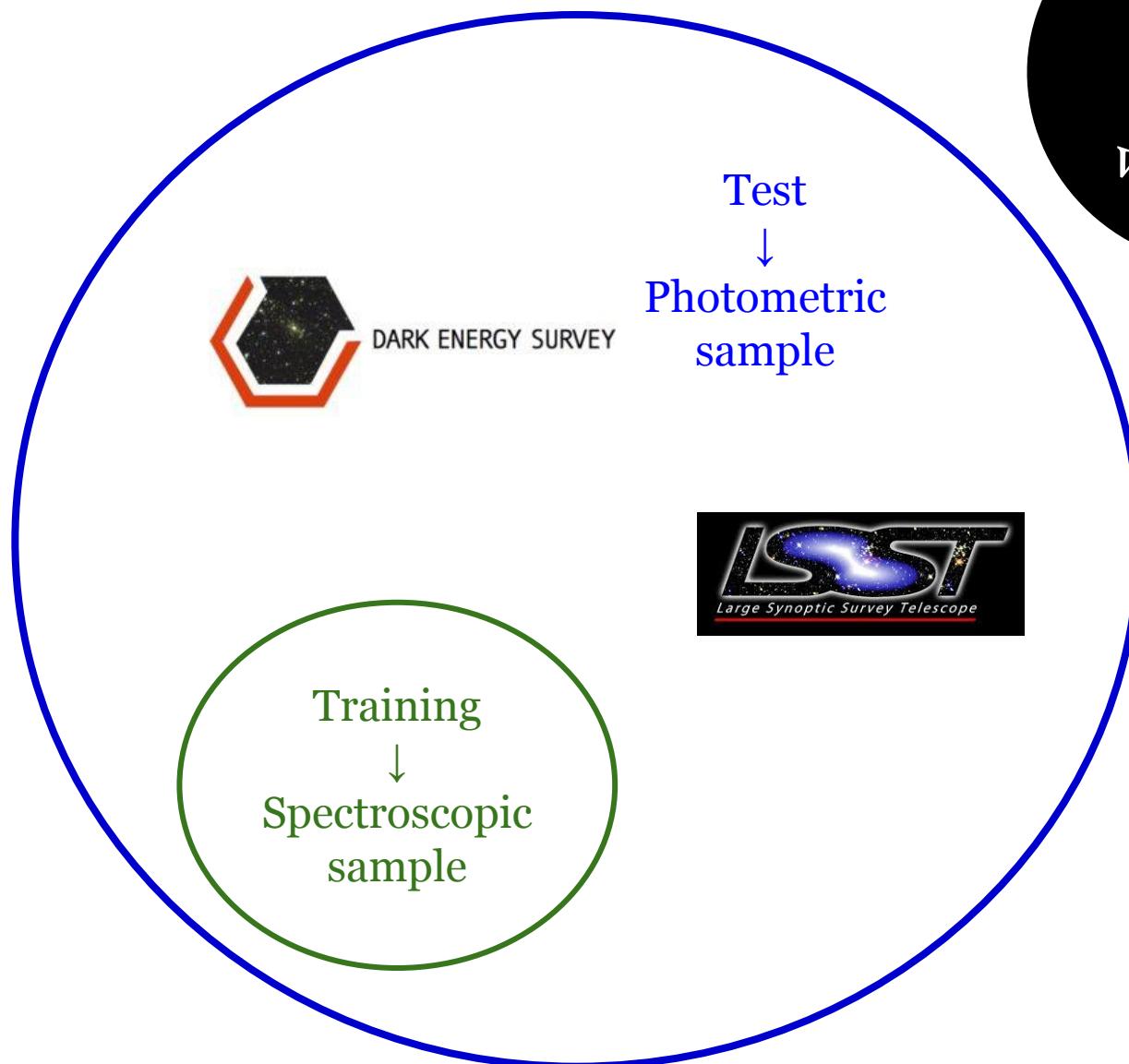
$$L_{data,f}(h) \equiv P_{x \sim data}(h_{train}(x) \neq f(x))$$

In Astronomy ...



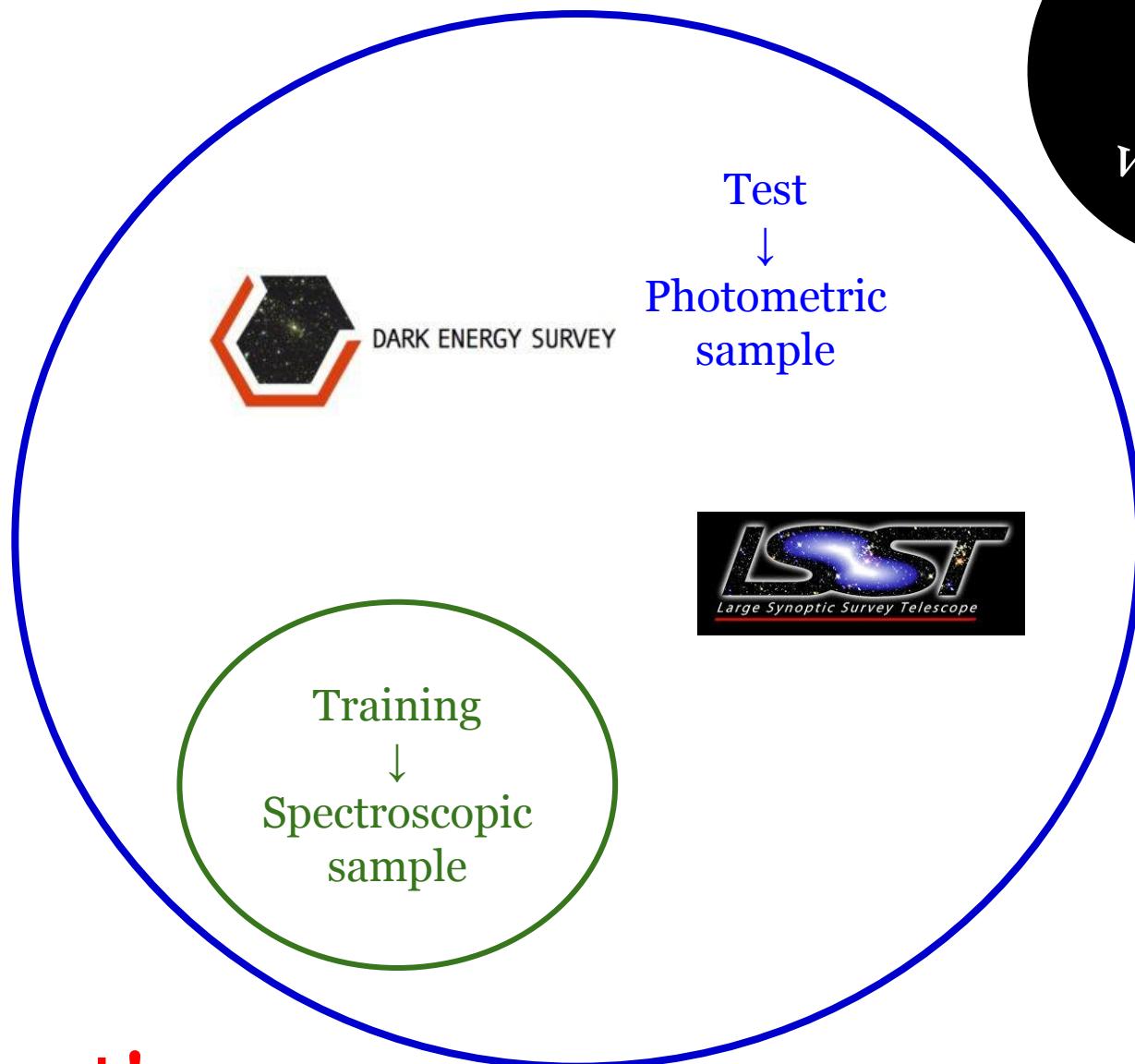
... it is NOT the case!

In Astronomy ...



... it is NOT the case!

In Astronomy ...



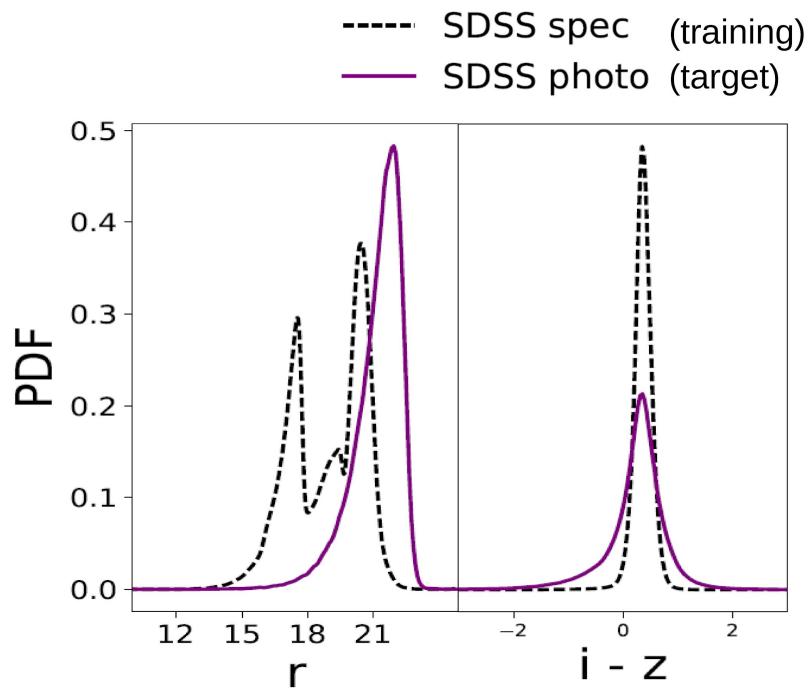
It does not!

... it is NOT the case!

Representativeness

(or the lack of) in Astronomy

Photometric redshift estimation



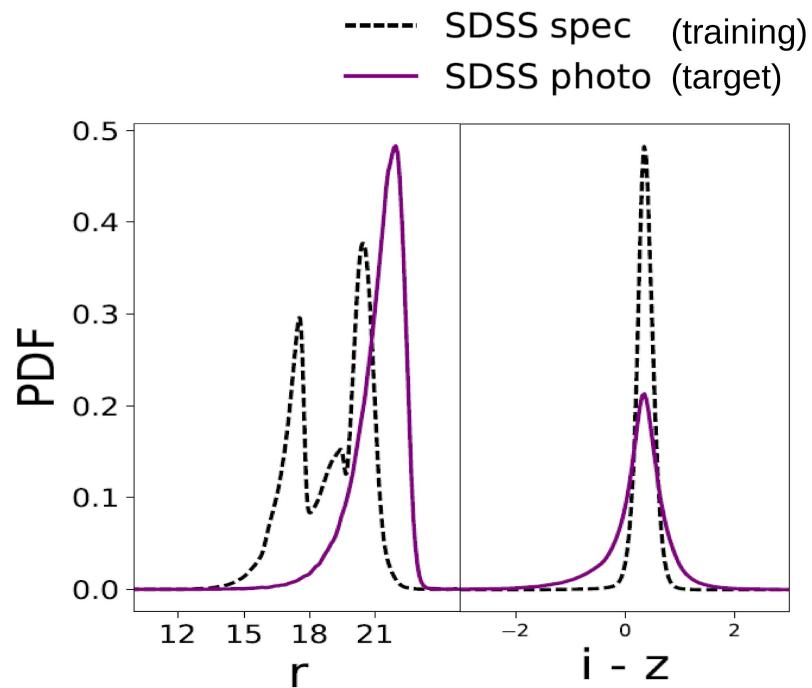
Beck et al., 2017, MNRAS - from CRP #3

Check **Linc's talk**: Cosmo21 - session 10
Thursday afternoon

Representativeness

(or the lack of) in Astronomy

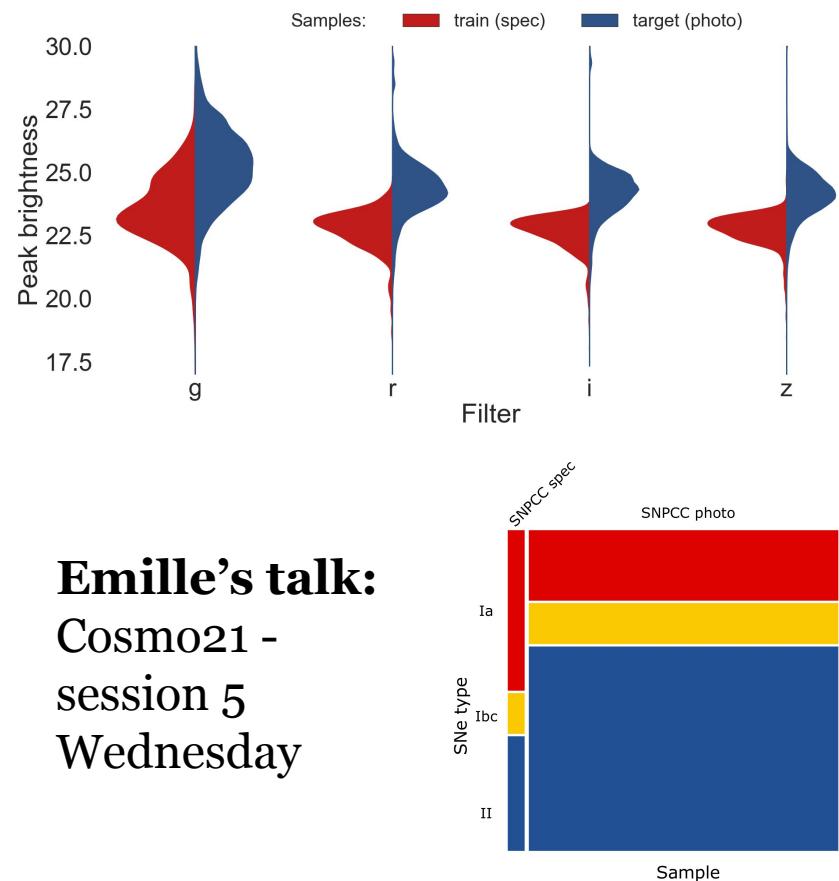
Photometric redshift estimation



Beck et al., 2017, MNRAS - from CRP #3

Check **Linc's talk:** Cosmo21 - session 10
Thursday afternoon

Supernova Classification



Emille's talk:
Cosmo21 -
session 5
Wednesday

Measurement Errors

Bayesian Machine Learning

Time Domain

SNe over LSST sky



year	Number of supernova
1998	42
2014	740
2025	> 10 000

2 million alerts/day
15 TB/day

40 nights of LSST

↓
entire Google database

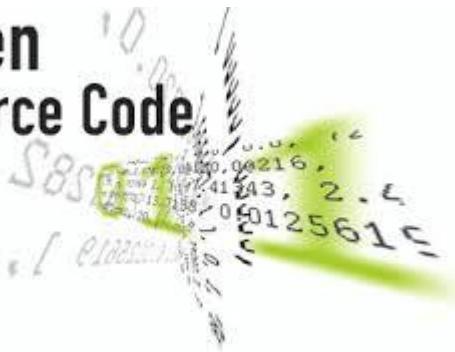
We will have to adapt!



We are getting there...

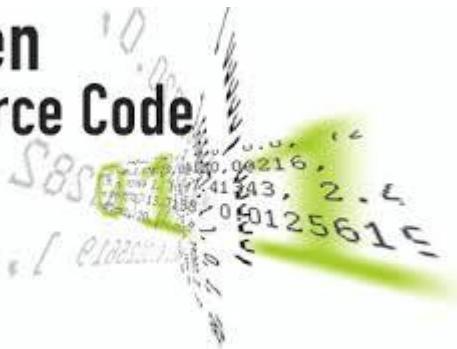
Developments in human learning

Open
Source Code



Developments in human learning

Open
Source Code



Community code
development

Developments in human learning



kaggle

Community code
development

GRAMP

Check **Alex's talk:** Cosmo21 - session 10
Thursday afternoon

PLAsTiCC

Photometric LSST Astronomical Time-series Classification Challenge

A data challenge aimed to prepare
a larger community for the LSST data paradigm

- PI: Renee Hlozek, simulations: Rick Kessler
- SNANA simulations → Light curves in observer-frame (no images!)
- 3 years worth of LSST data, ~ 100 MB
- ~ 10^7 objects
- Around 20 transient models
(galactic and extra-galactic, periodic and non-periodic)
- Please respect model-information policy:
``don't ask, don't tell''
- Not all models will be present in the training sample
- Supervised classification + novelty detection
- Deployment: **kaggle** + **RAMP**

Expected release date:
Summer/Fall 2018



Get ready!

Play with the first SNPCC data!

Here is a peak at the data.

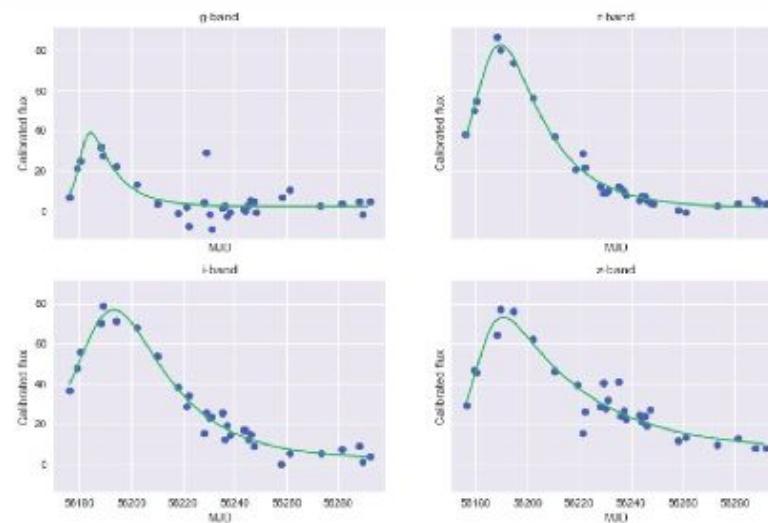
```
In [2]: X, y = read_data('data/des_train.pkl')
```

```
In [3]: # 5 first rows of the dataframe  
X.head()
```

```
Out[3]:
```

	mjd_g	fluxcal_g	fluxcalerr_g	mjd_r	fluxcal_r	fluxcalerr_r	mjd_I	fluxcal_I	fluxerr_I
1186	56283.203,	[3.182,	6.286,	56283.211,	-4.973,	3.326,	56283.219,	-9.615,	-9.284,
	56288.199,	-13.07,	-7.177,	56288.215,	-6.427,	6.181, 4.83,	56288.109,	0.8371,	
	56292.098,	-7.177,	11.79, 11.55,	56292.102,	9.881,	1.657,	56292.191,	0.8885,	
	56304.051,	-0.4645,	2.479,	56304.055,	-1.683,	2.968,	56304.07,	-1.593,	
5...	5...	0.2998,	7.589,	5...	6.724,	2.553,	56...	56...	24.51... 3
	16.63...	6.195, 3.6...	5...	23.82, ...	3.07...	3.07...			
	56177.172,	[27.18,	4.837,	56177.188,	[33.02,	5.375,	56177.203,	49.56,	
	56178.172,	34.03,	2.372,	56179.312,	34.65,	3.109,	56179.329,	42.79,	
	56179.312,	24.15,	1.805,	56179.312,	26.92,	1.365,	56179.329,	30.72,	

```
In [14]: plot lightcurves with fit(2)
```



PLAsTiCC - RAMP: classification of astronomical transients

Alexandre Boucaud (Paris-Saclay Center for Data Science)
Emille E. O. Ishida (Université Clermont Auvergne)

Constructing the data matrix for classification

Now that you can fit each individual light curve, you are ready to build your low dimension representation of the data.

<https://github.com/ramp-kits/supernovae>

Build an interdisciplinary, people-centric community!



The banner features a dark background with a network of white lines connecting black dots, resembling a star map or a complex network. In the center, the COIN logo is displayed with a stylized atom or planet icon integrated into the letter 'O'. Below the logo, the text "Cosmostatistics Initiative" is written in a smaller, light-colored font. At the bottom, the text "COIN Residence Program #5" and "Chania, Greece, 15 - 22 September 2018" is prominently displayed in large, bold, white letters.

COIN Residence Program #5
Chania, Greece, 15 - 22 September 2018

COIN [Home](#) [About](#) [Organizers](#) [Location](#) [Code of Conduct](#) [Apply](#) [Partners](#)

<https://www.cosmostatistics-initiative.org>

Summary

- There is much more than supervised/unsupervised learning
- Not only machines, but humans will be required to adapt
 - Astronomical data can lead to important development in computer science
- We should also think about updating the academic model

Thank you!



Cosmostatistics Initiative

<http://cointoolbox.github.io/>