

3rd ASTERICS-OBELICS International School

8-12 April 2019, Annecy, France.



H2020-Astronomy ESFRI and Research Infrastructure Cluster
(Grant Agreement number: 653477).



Machine Learning Tutorial IV - Beyond textbook ML

*3rd ASTERICS-OBELICS International School
8-12 April 2019, Annecy - France*

Emille E. O. Ishida

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



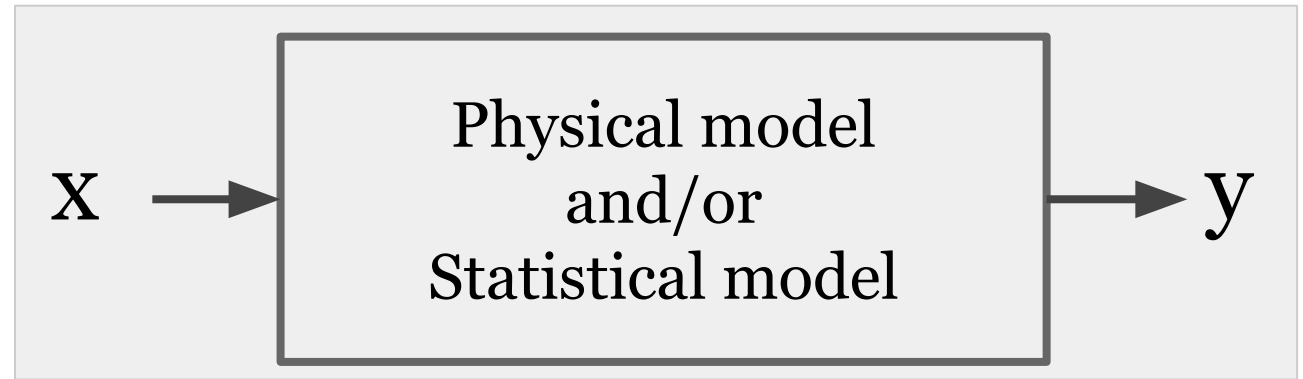
Summary

- I. Quick recap of the week
with a few interesting additions
- II. Representativeness matters
- III. Adaptive Learning Techniques
- IV. The human factor

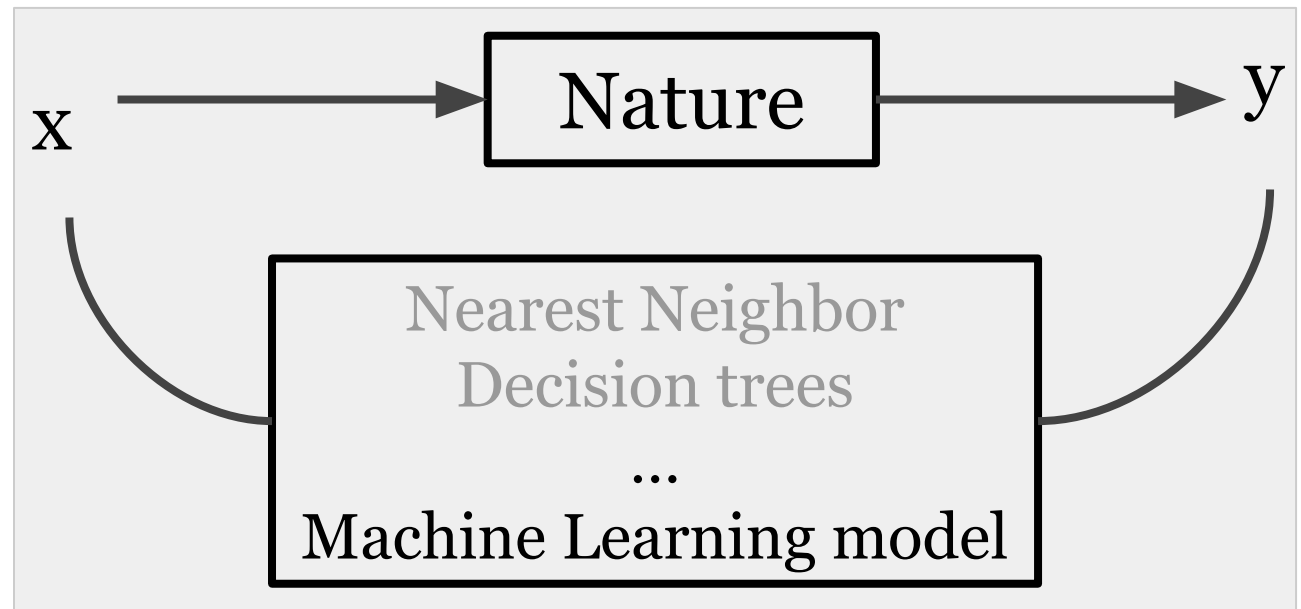
Hypothesis:



Traditional
data
modeling:



Algorithmic
modeling:



Supervised ML model

data **training**, target

\mathcal{X} set of all samples, x

\mathcal{Y} set of possible labels, y

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

L Loss function

Representativeness matters!

Data generation model:

$$x_i \sim P_X$$

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

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

Supervised ML model

data training, target

Machine Learning algorithm

X set of samples, x

Y set of possible labels, y

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

L Loss function

Representativeness matters!

Data generation model:

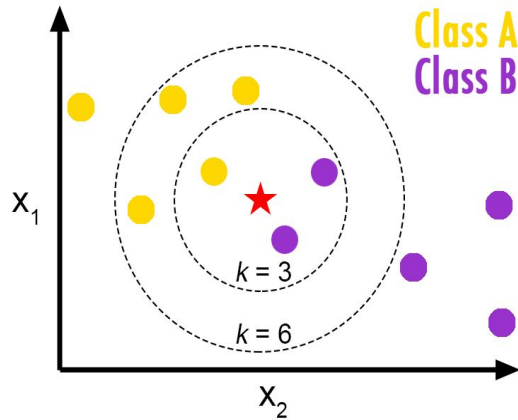
$$x_i \sim P_X$$

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

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

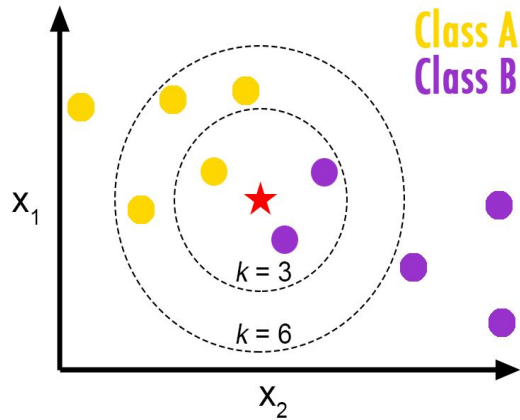
Machine Learning algorithms

K Nearest Neighbor

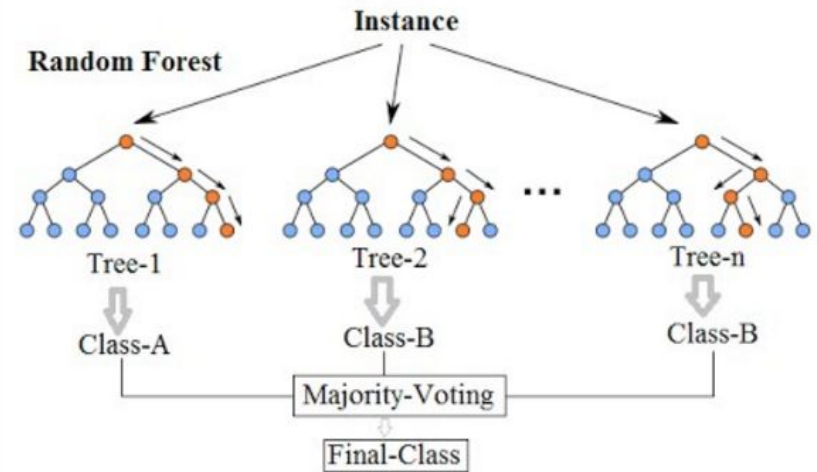


Machine Learning algorithms

K Nearest Neighbor

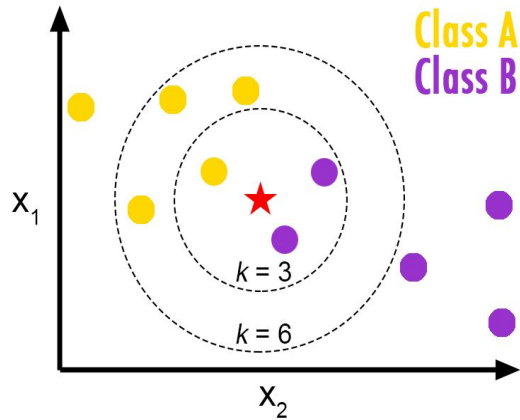


Decision trees and random forests

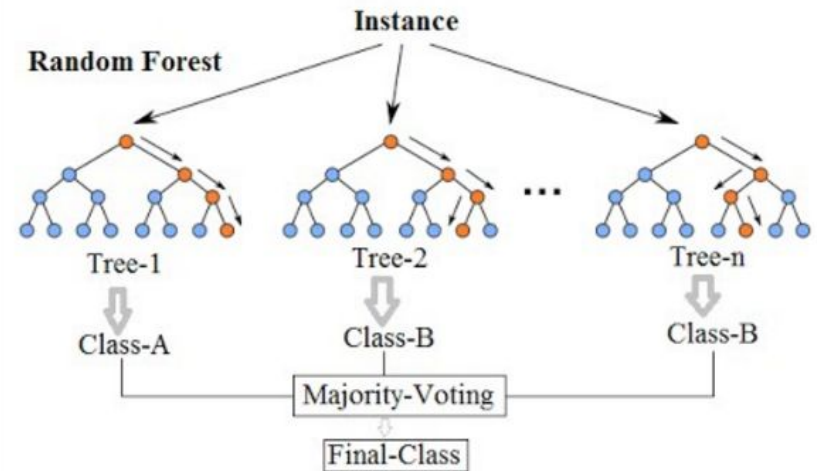


Machine Learning algorithms

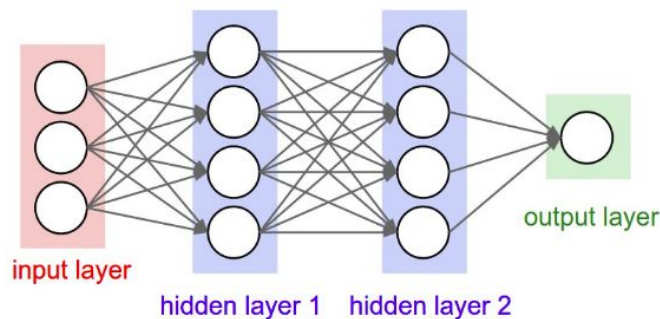
K Nearest Neighbor



Decision trees and random forests

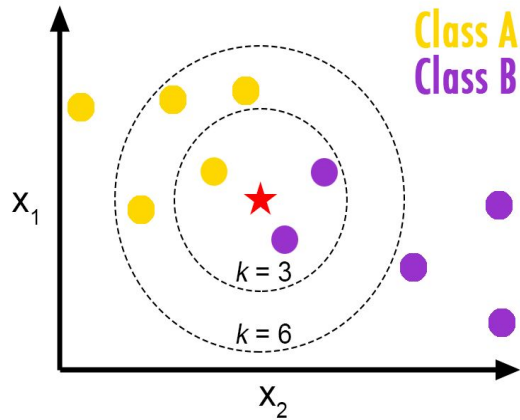


Neural networks and deep learning

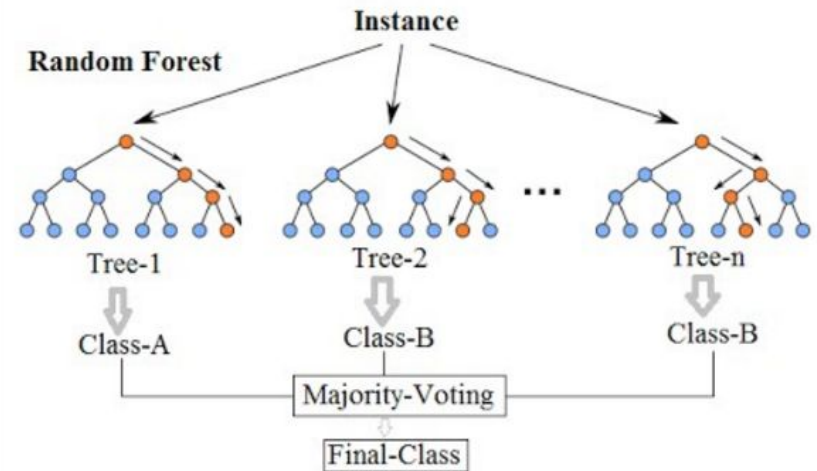


Machine Learning algorithms

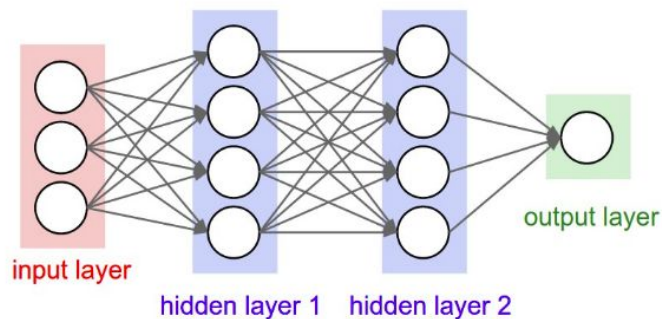
K Nearest Neighbor



Decision trees and random forests

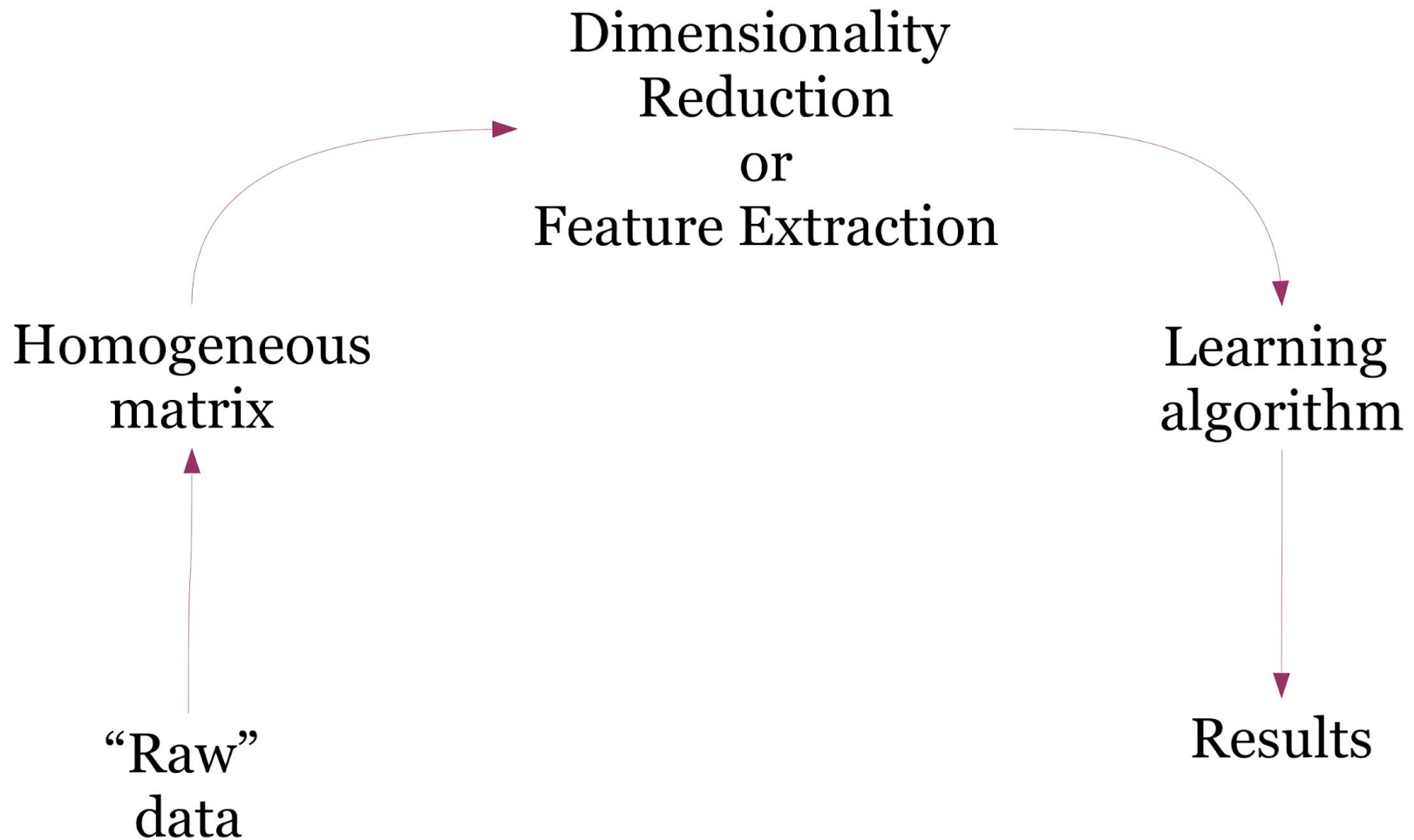


Neural networks and deep learning



A comment on

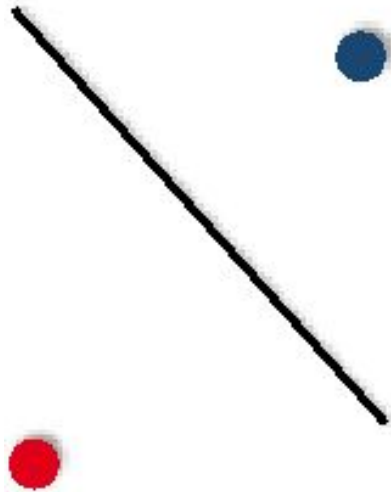
Feature extraction



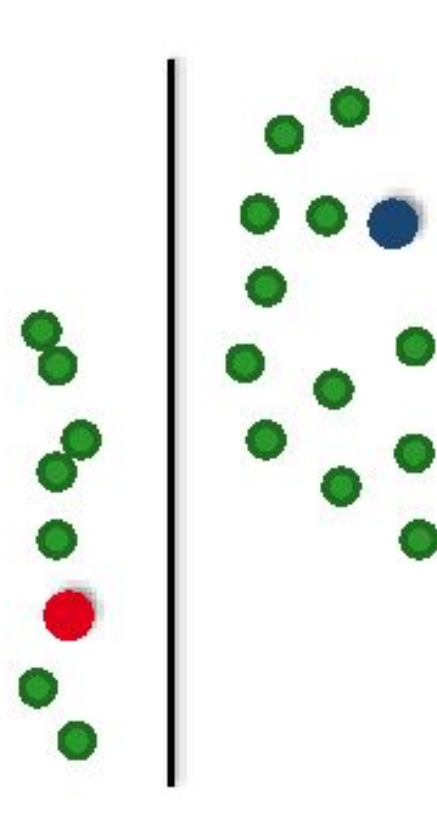
Semi-supervised learning

Getting partial information from the unlabelled sample

only labeled data



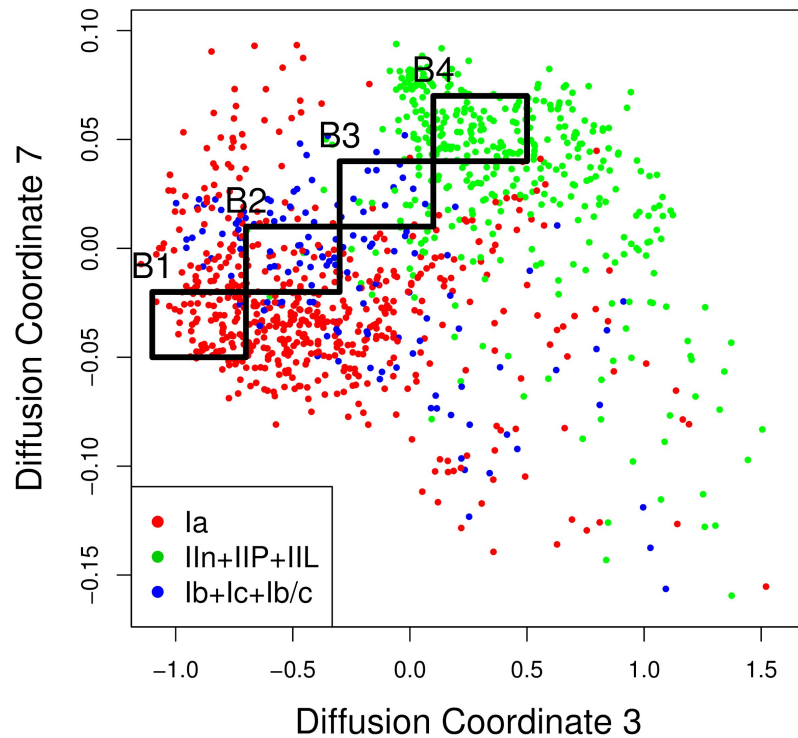
with unlabeled data



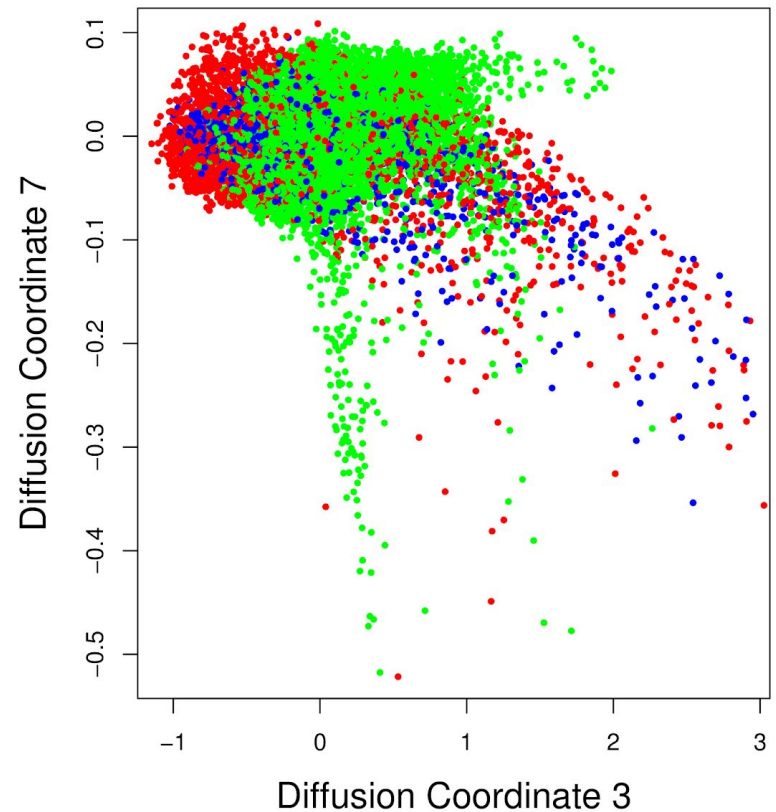
Semi-supervised learning

For Supernova Photometric Classification

Spectroscopic only (training)

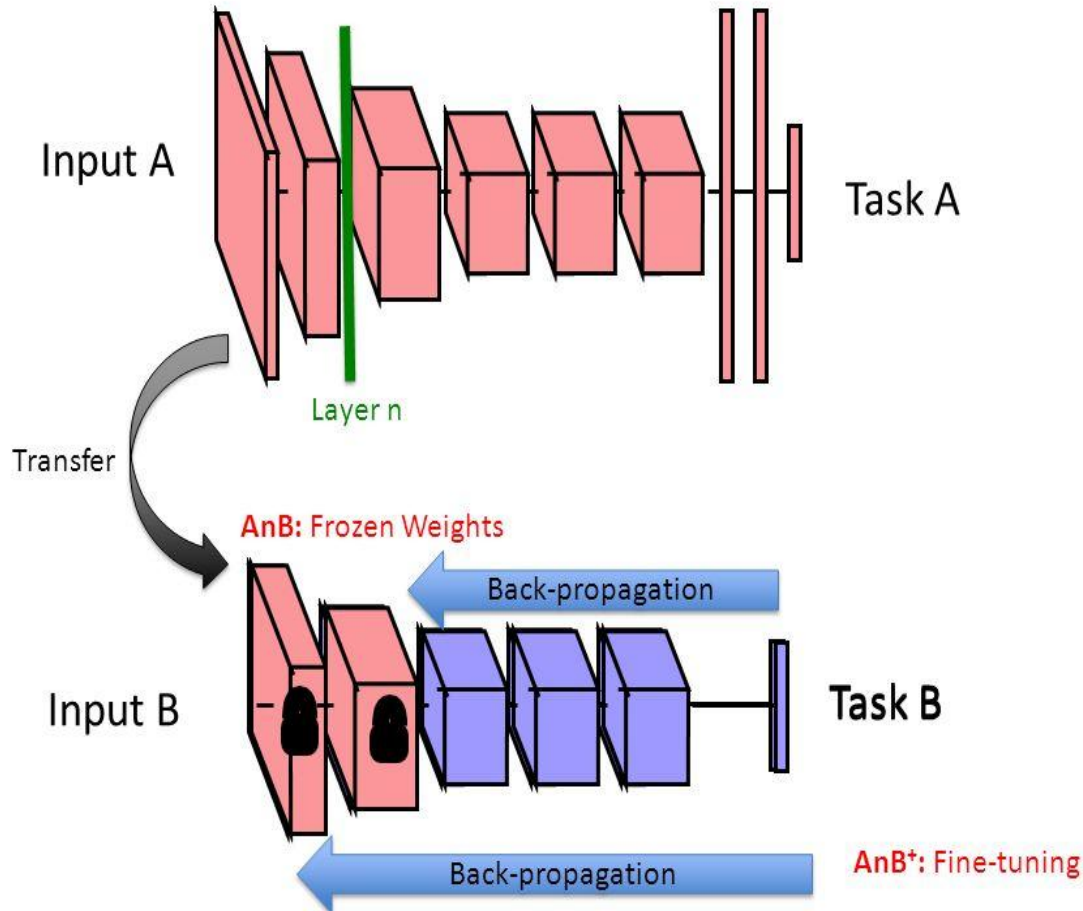


All available data



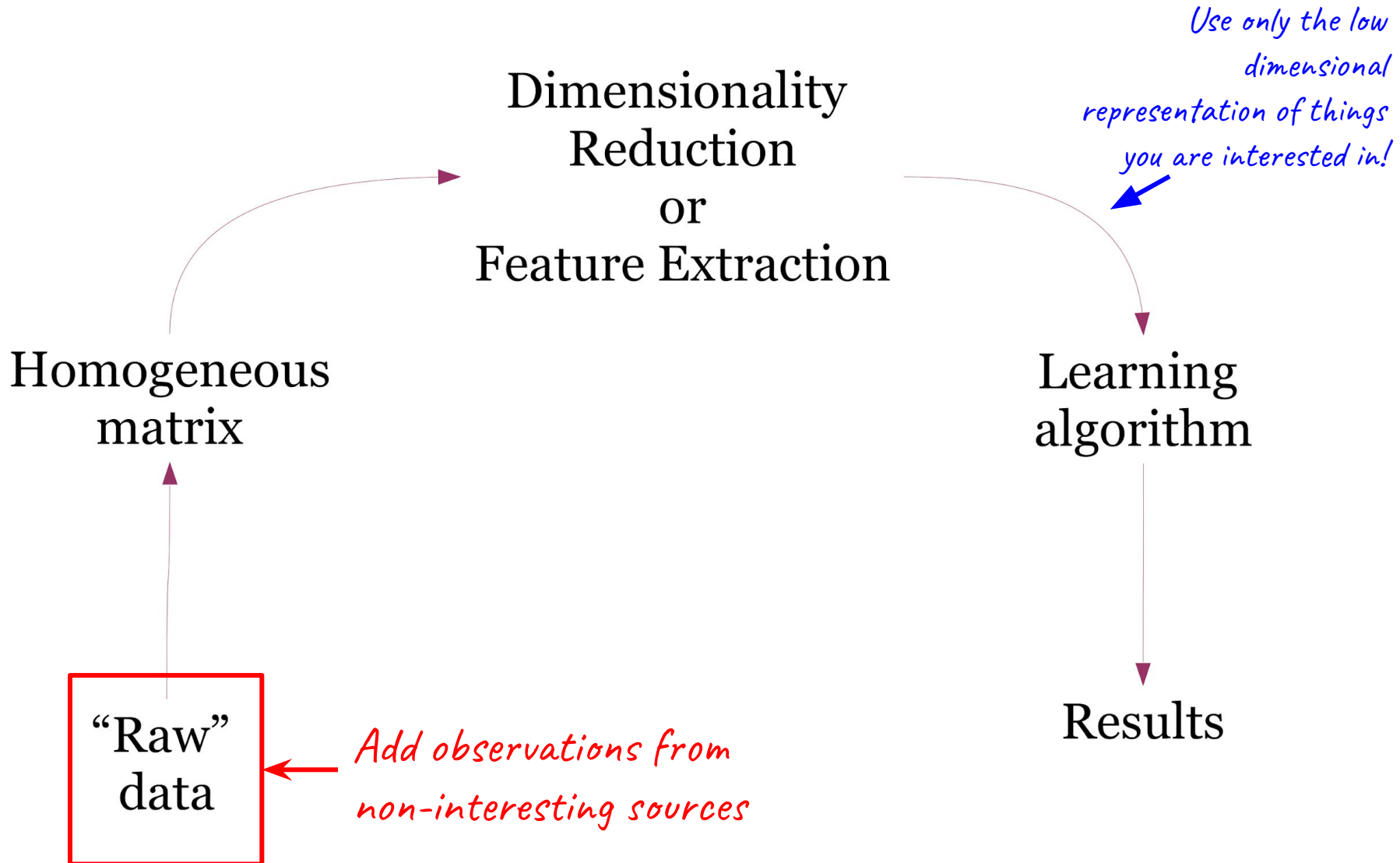
Transfer Learning

Borrowing information from somewhere else



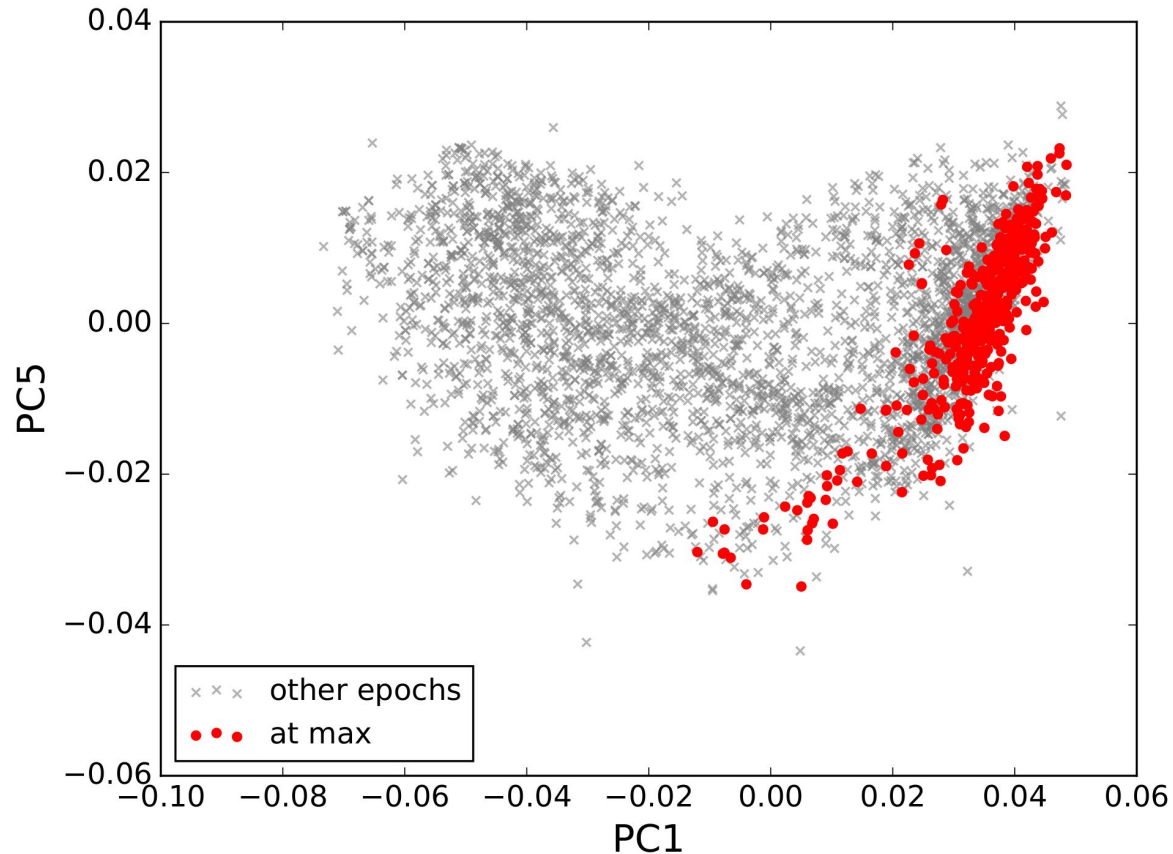
Transfer Learning

Exploiting information from various data sets



Transfer Learning

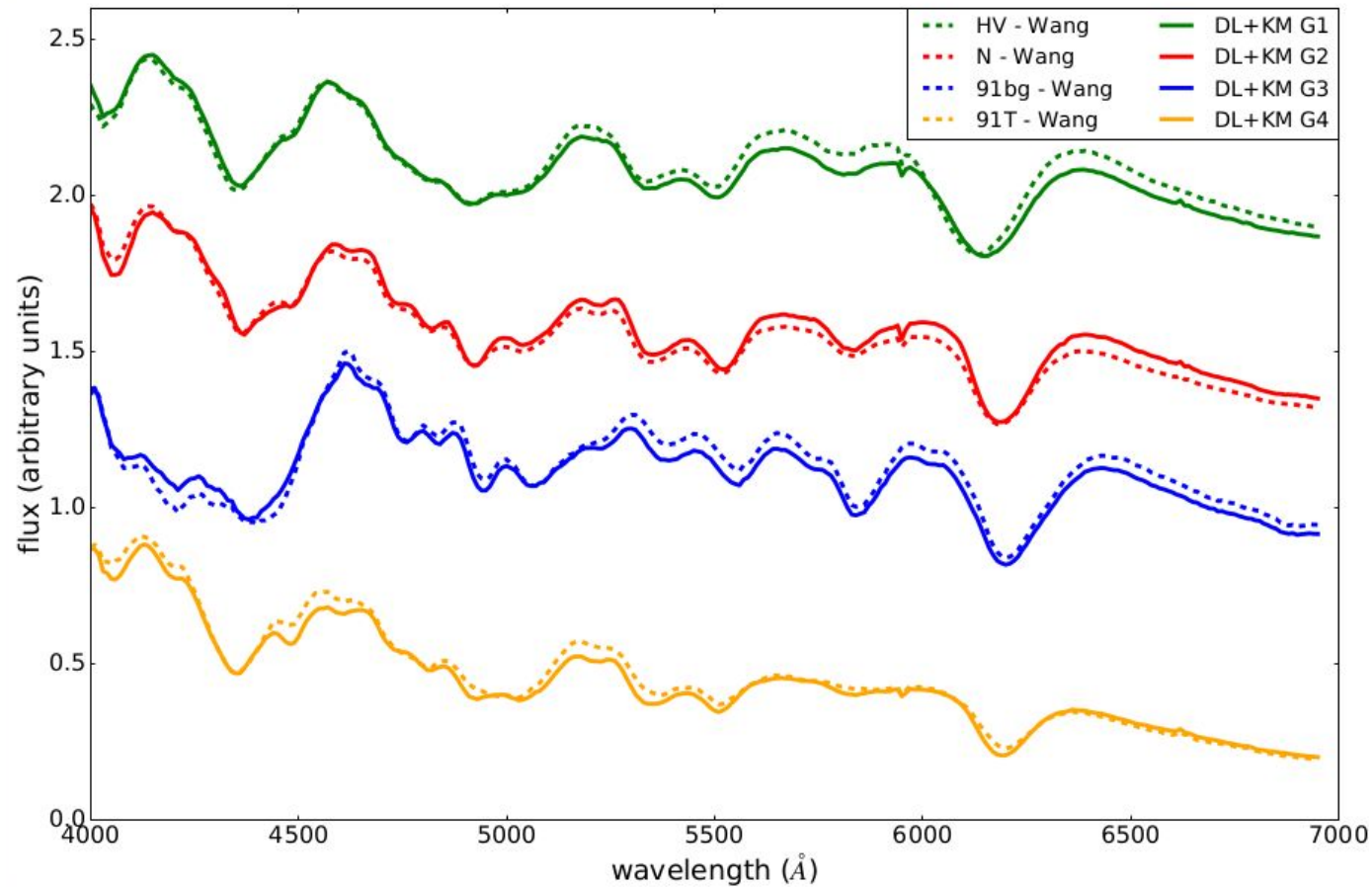
In Astronomy



Sasdelli, Ishida et al., 2016, MNRAS, 461, Issue 2, p.2044, from CRP #2

Unsupervised Clustering

In Astronomy



Sasdelli, Ishida et al., 2016, MNRAS, 461, Issue 2, p.2044, from CRP #2

Neural Network

In astronomy: photometric redshift estimation

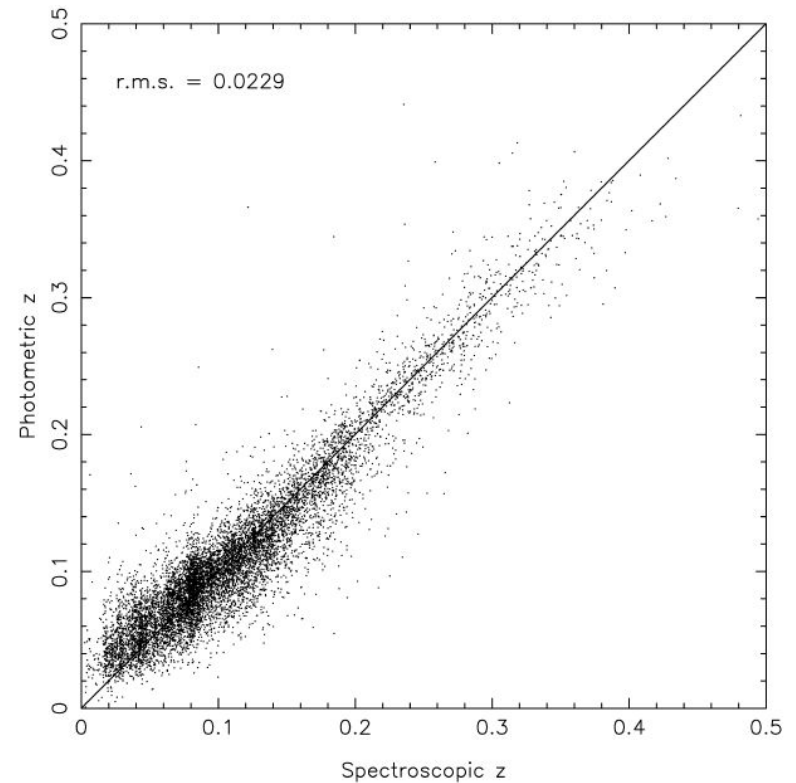
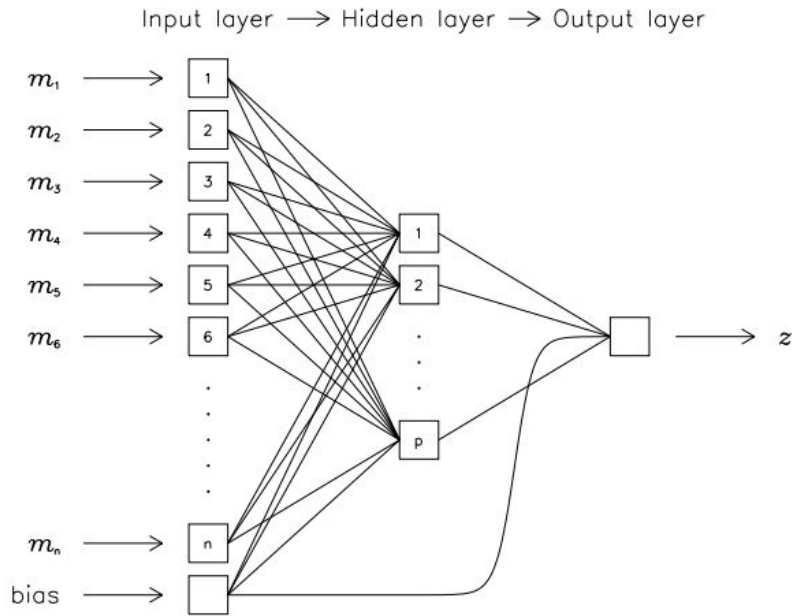
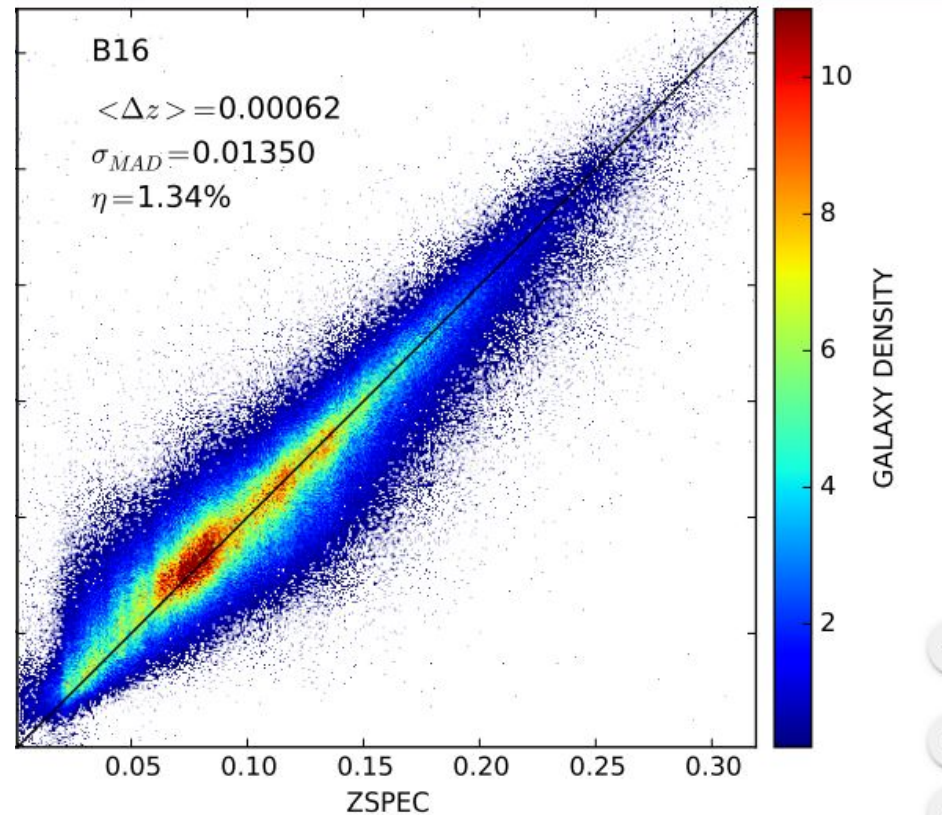
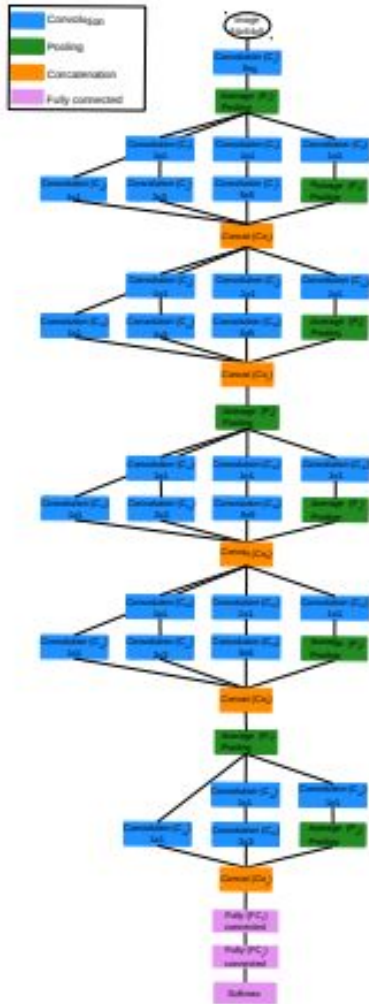


FIG. 2.— Spectroscopic vs. photometric redshifts for ANNz applied to 10,000 galaxies randomly selected from the SDSS EDR.

Neural Network

In astronomy: photometric redshift estimation



Symbolic Regression

Supervised Learning: an extreme regression example

Mathematical atoms:

+ , - , x , / , pow

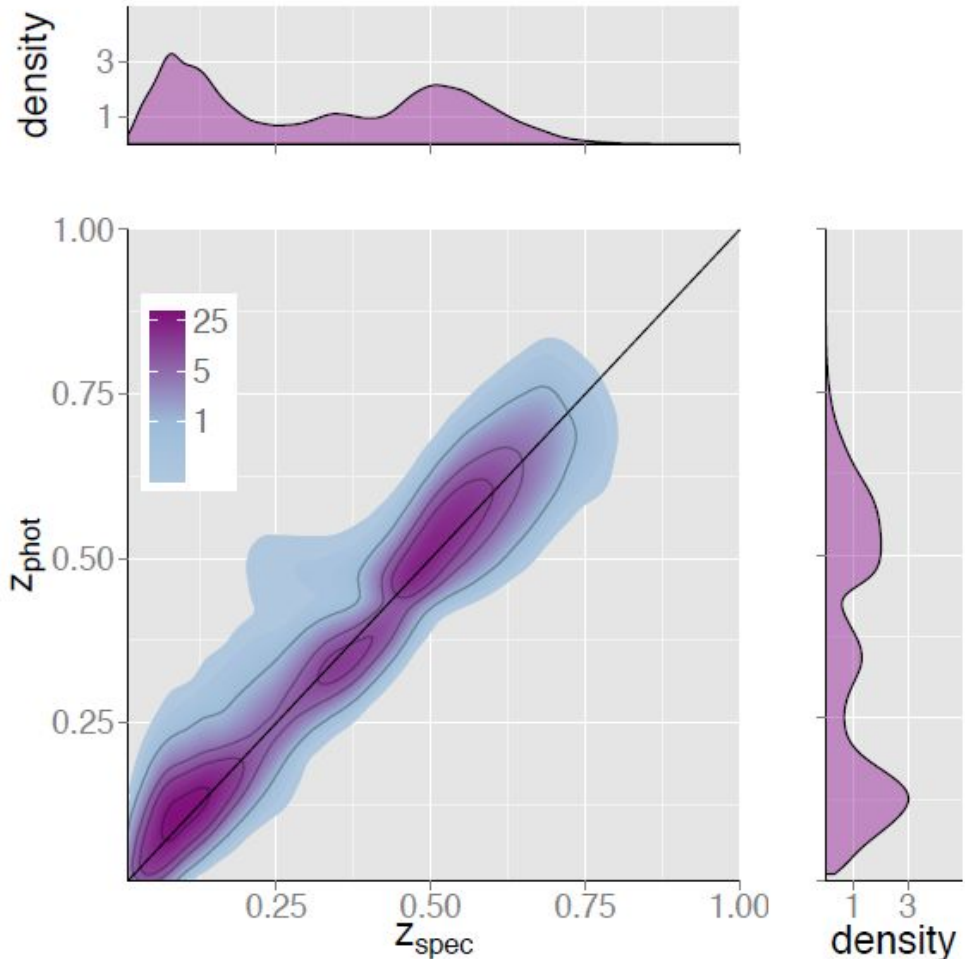
1 - Random construction of an analytical expression

2 - find the best parameters

3 - if result is better than previous keep it, otherwise discard it

Final expression:

$$z_{\text{phot}} = \frac{0.4436r - 8.261}{24.4 + (g - r)^2(g - i)^2(r - i)^2 - g + 0.5152(r - i)}.$$



All of the above relies on
representativeness...

This is a very strong assumption

All of the above relies on
representativeness...

*How often does this hypothesis
hold in astronomy?*

All of the above relies on
representativeness...

*How often does this hypothesis
hold in astronomy?*

What happens when it breaks?

Notebook: Regression2.ipynb

Labels are expensive

We need recommendation systems



Labels are expensive

How to construct optimal training samples?

amazon

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

NETFLIX

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



Labels are expensive

How to construct optimal training samples?



amazon

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

NETFLIX

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



Labels are expensive

How to construct optimal training samples?

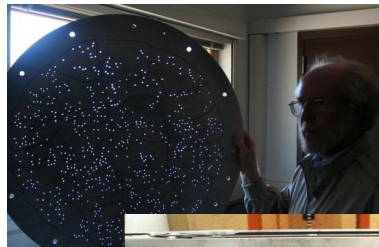


amazon

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

NETFLIX

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



Labels are expensive

How to construct optimal training samples?



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



amazon

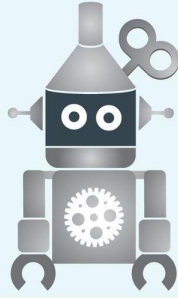
35% OF AMAZON'S REVENUE ARE GENERATED BY IT'S RECOMMENDATION ENGINE.

NETFLIX

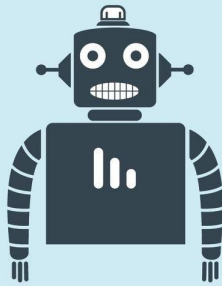
75% OF USERS SELECT MOVIES BASED ON NETFLIX'S RECOMMENDATIONS.



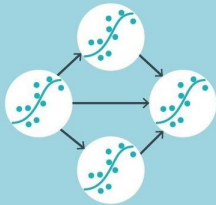
FIRST GENERATION:
Rule-based



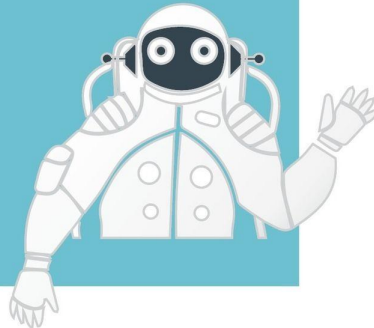
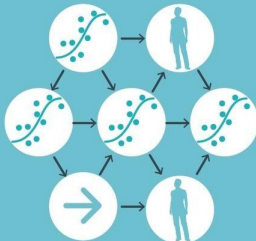
SECOND GENERATION:
Simple machine learning



THIRD GENERATION:
Deep learning



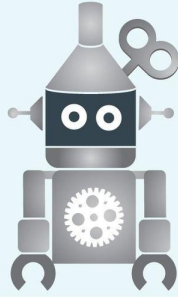
FOURTH GENERATION:
Adaptive learning



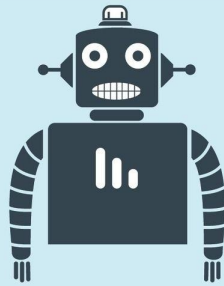
Machines
need to
evolve...

so they
need to
adapt!

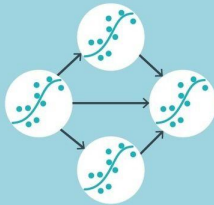
FIRST GENERATION:
Rule-based



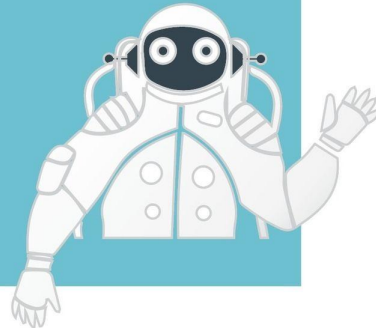
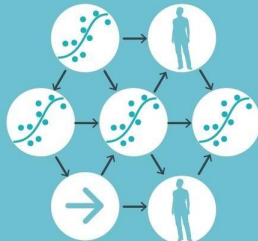
SECOND GENERATION:
Simple machine learning



THIRD GENERATION:
Deep learning



FOURTH GENERATION:
Adaptive learning

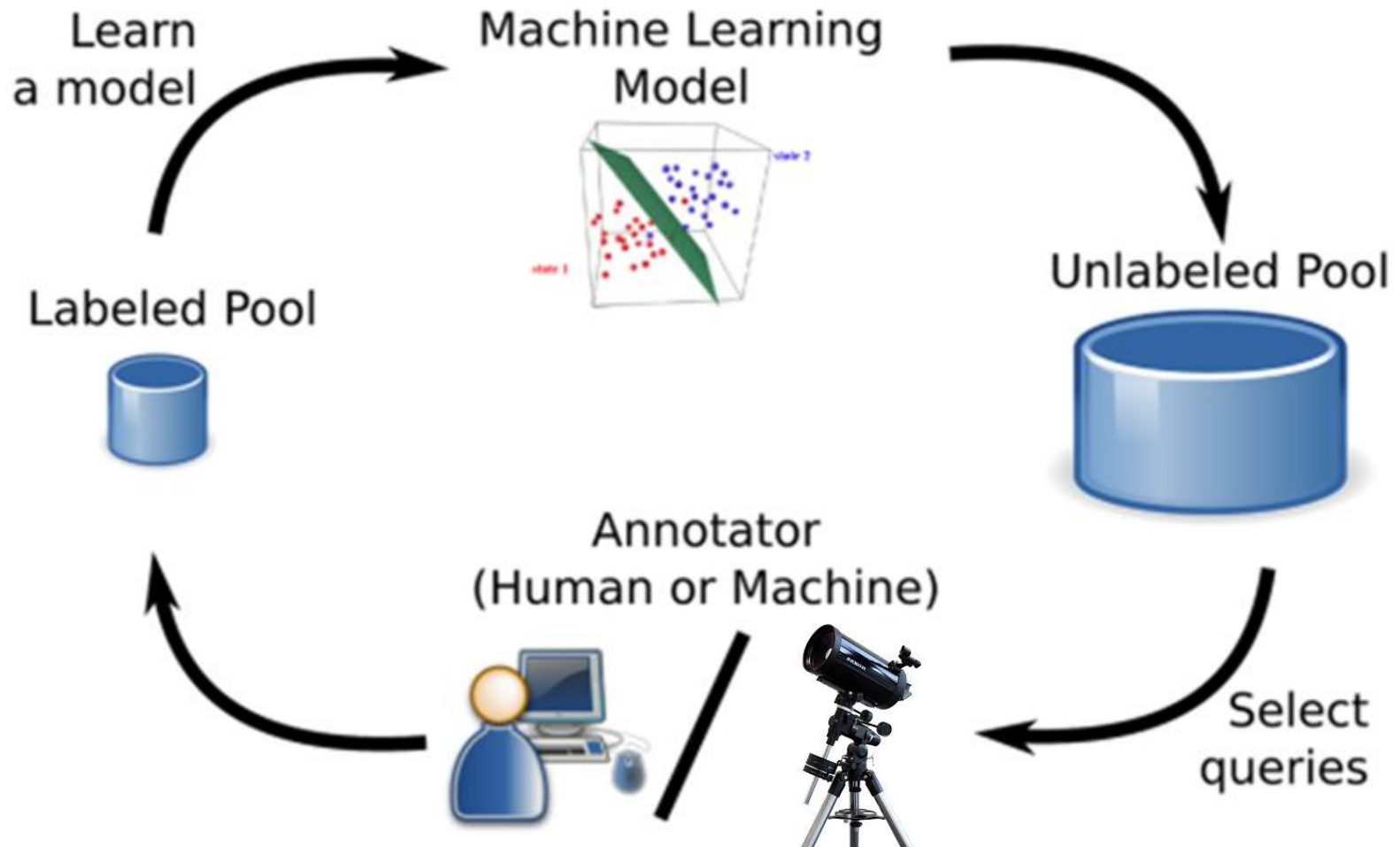


Machines
need to
evolve...

so they
need to
adapt!

Active Learning

Optimal classification, minimum training



Optimal Experiment Design

In Statistics literature

$$PQ_{data,f}(x) \propto P_{x \sim data}(h_{train}(x) \neq f(x) \mid \textit{previous results})$$

- Pool based
- Generative
- Sequential

Active Learning for

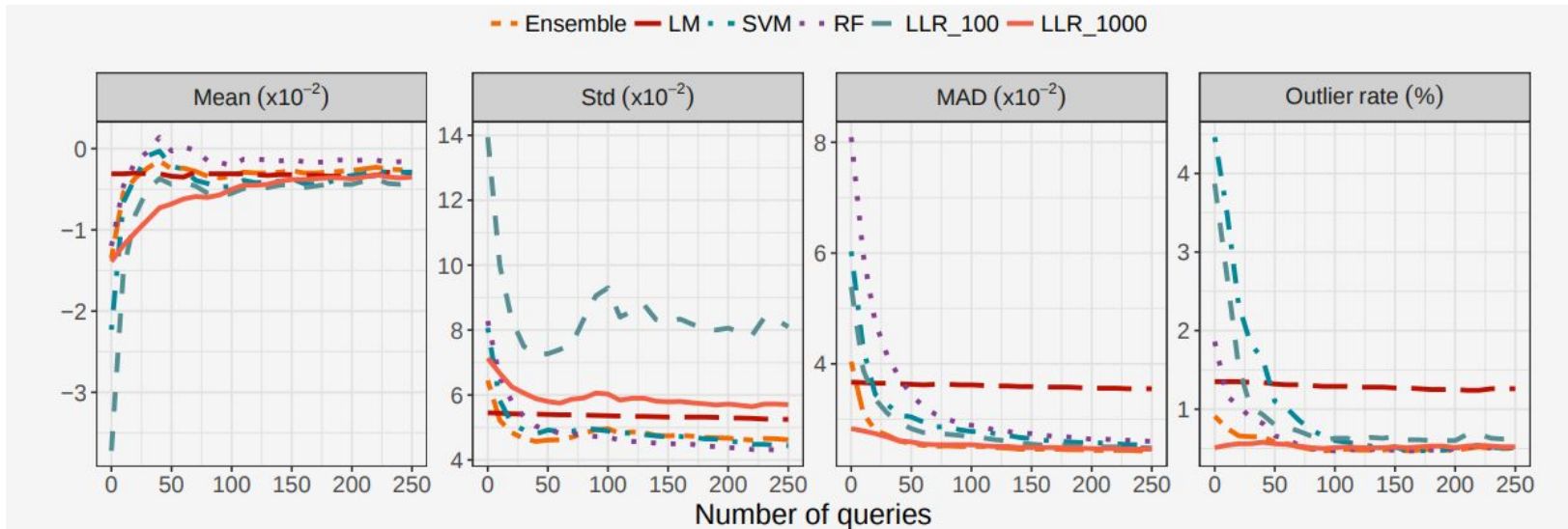
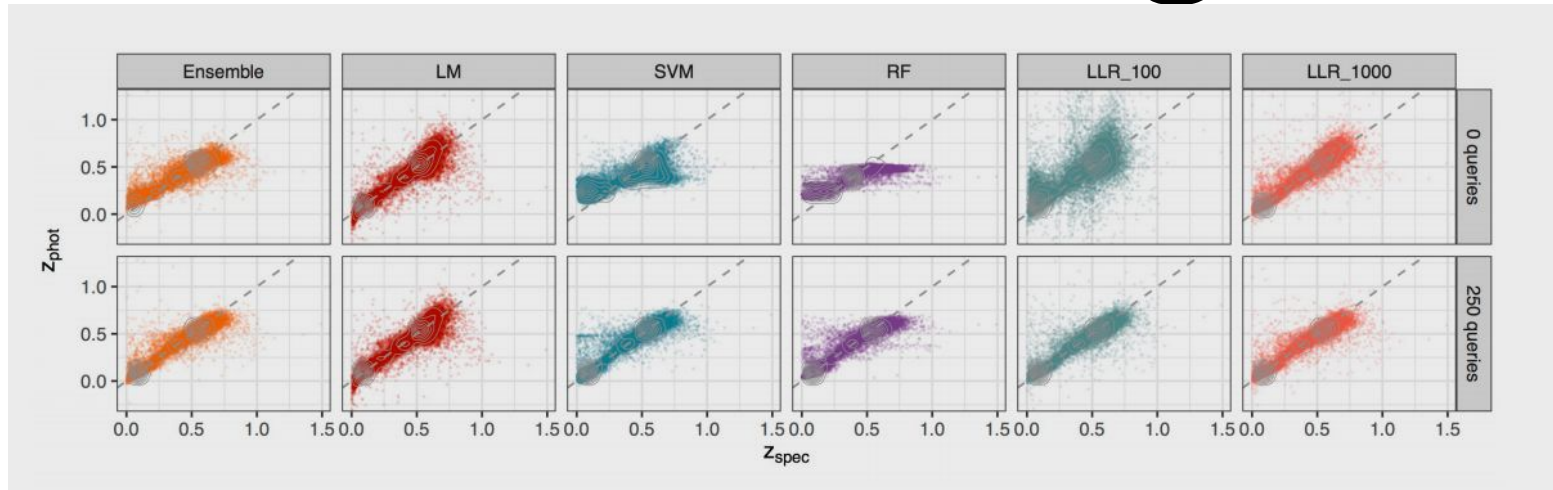
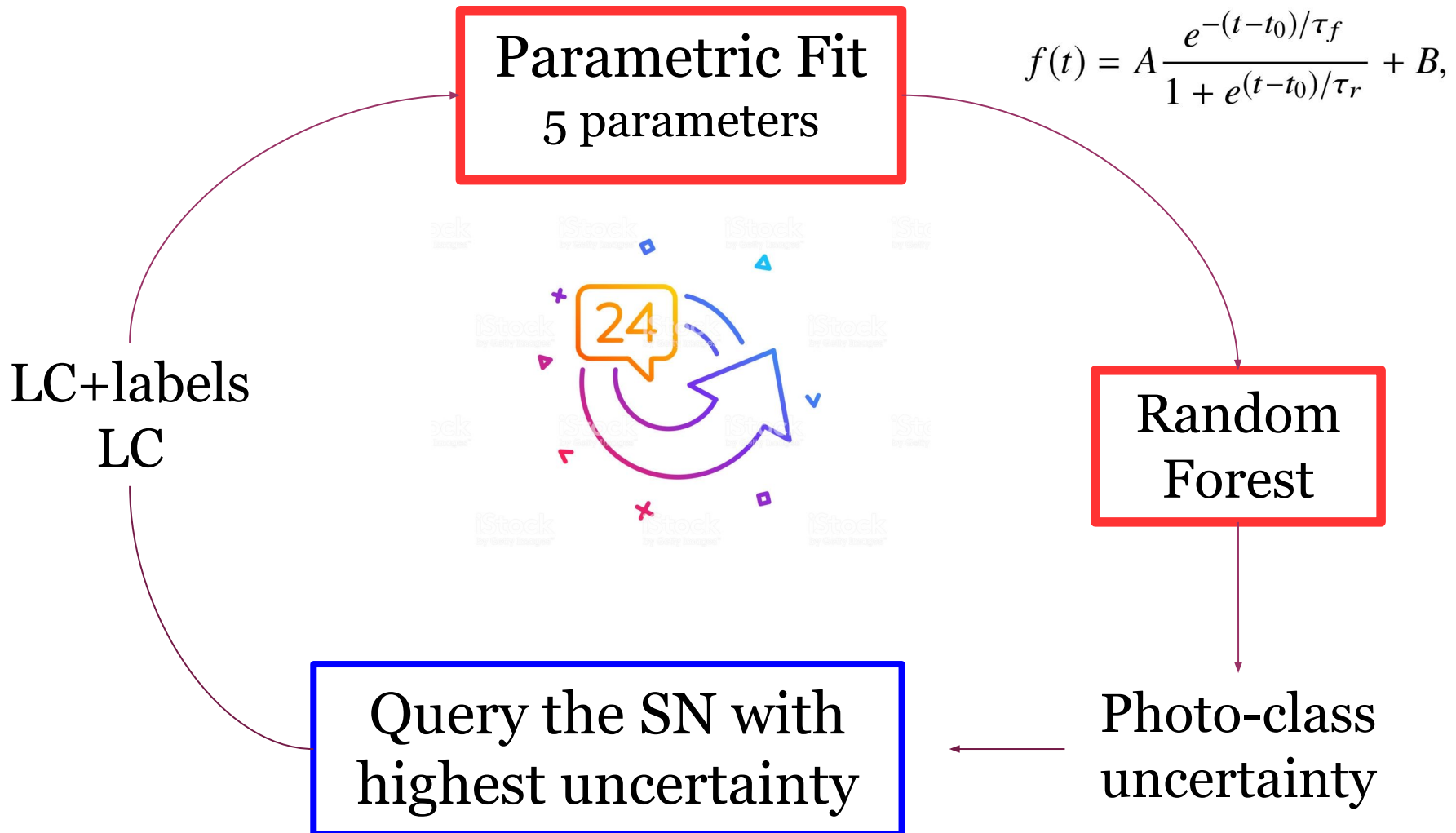


Figure 4. An assessment of the performance of the ensemble model and its constituent models using active learning. Performance diagnostics are shown as a function of the number of queries.

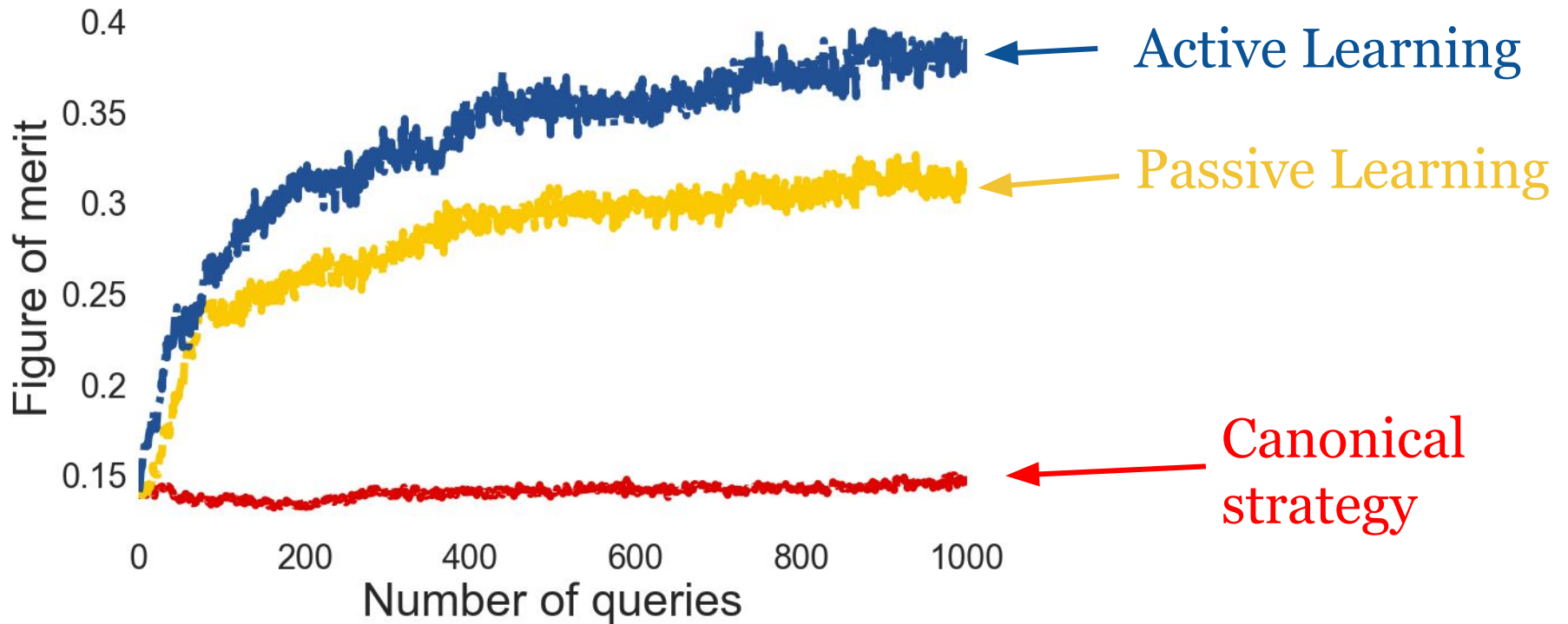
AL for Supernova classification

A strategy



AL for SN classification

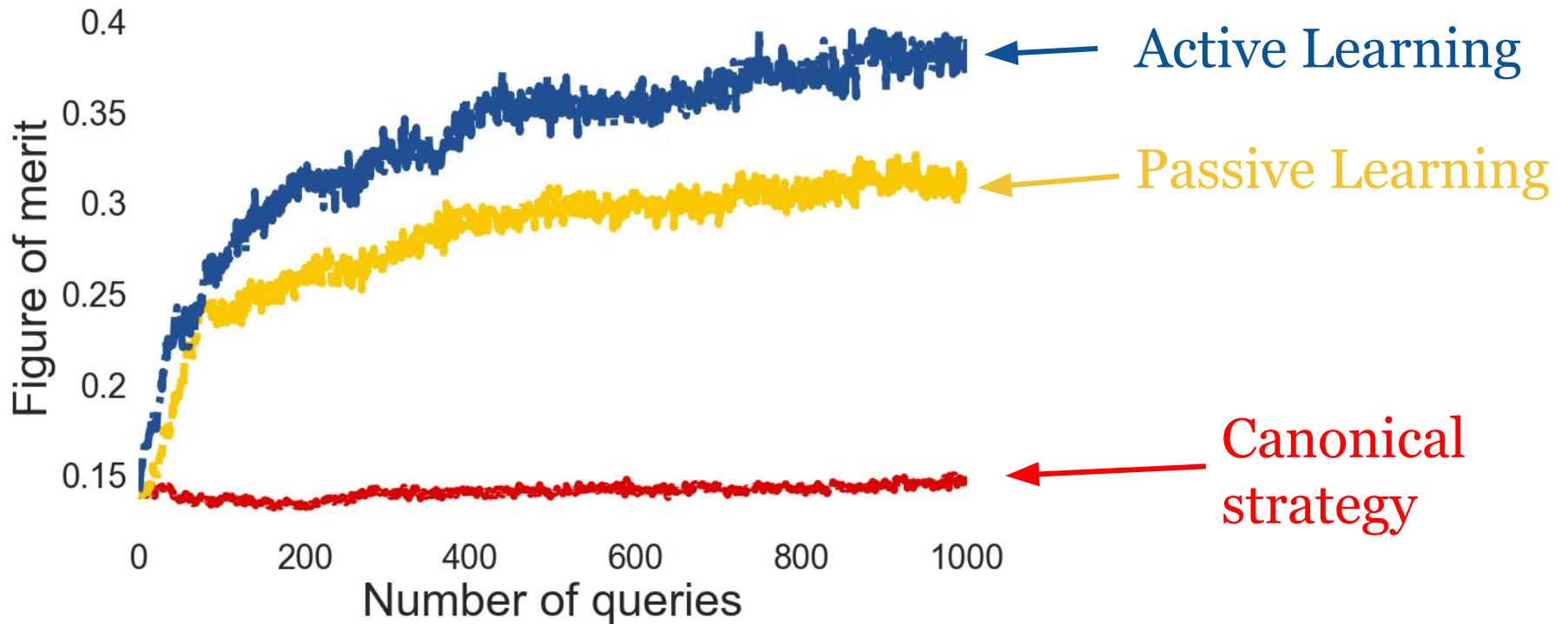
Static results



AL for SN classification

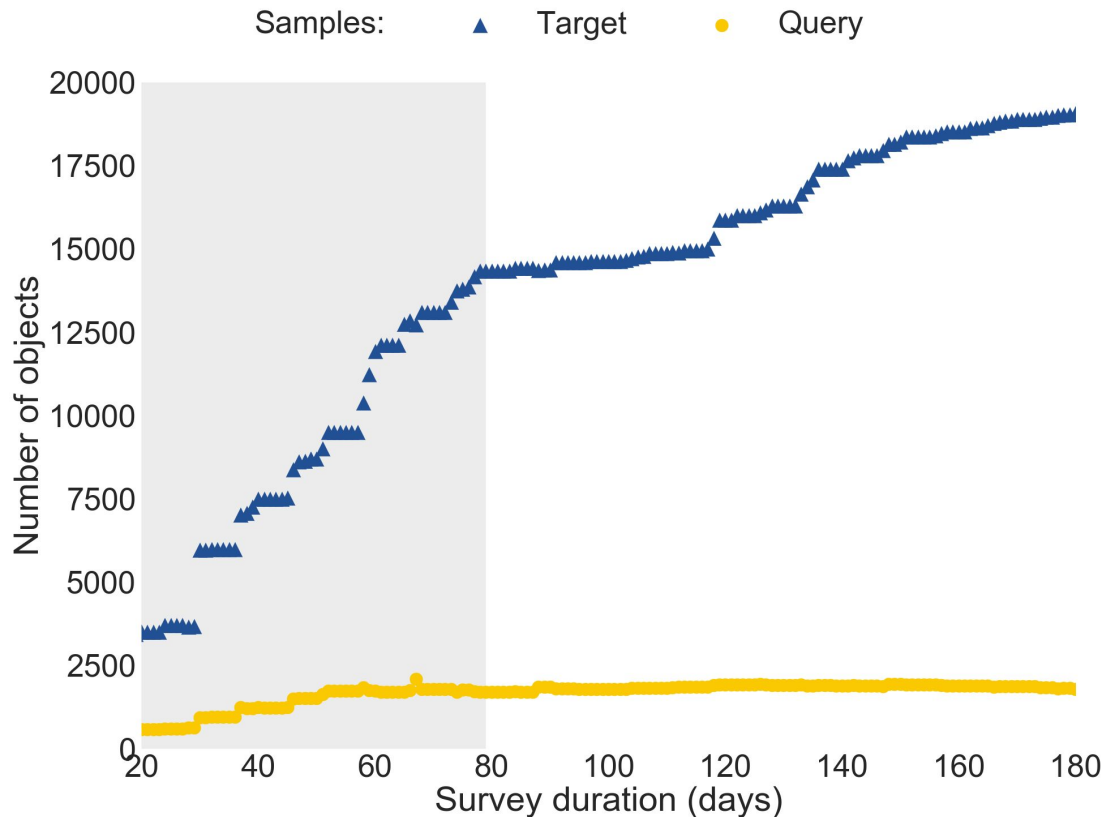
Static results

What astronomical aspect make this setting non-realistic?



SN are transients

Window of opportunity



1. Feature extraction done daily **with available observed epochs until then.**

2. Query sample is also re-defined daily: objects with **r-mag < 24**

Does this solve the
problem completely?

No, it is just the best you can do!

Does this solve the
problem completely?

No, it is just the best you can do!

Is this the only way
of doing it?

No!, it is only one exciting possibility

Summary

“How do we optimize machine learning results for astronomical purposes?”

What we need

What we have

**Adaptive
Learning
designed for
astronomical
data**



Final thoughts ...

We will have to adapt!



We are getting there...

Developments in human learning



Community code
development

Developments in human learning

Open
Source Code



Community code
development



Data challenges

kaggle

SRAMP

Developed by
Paris-Saclay Center for Data
Science

Developments in human learning

Open
Source Code



Community code
development



Data challenges

kaggle

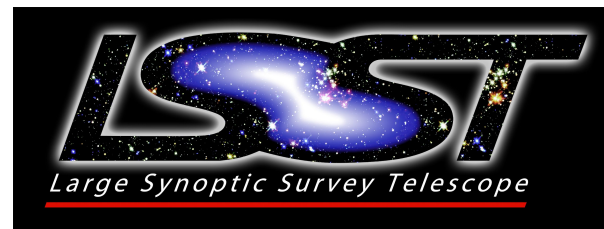
SRAMP

Developed by
Paris-Saclay Center for Data
Science

Large
collaborations



euclid



Good things happen when brains
are connected properly...



C  I N

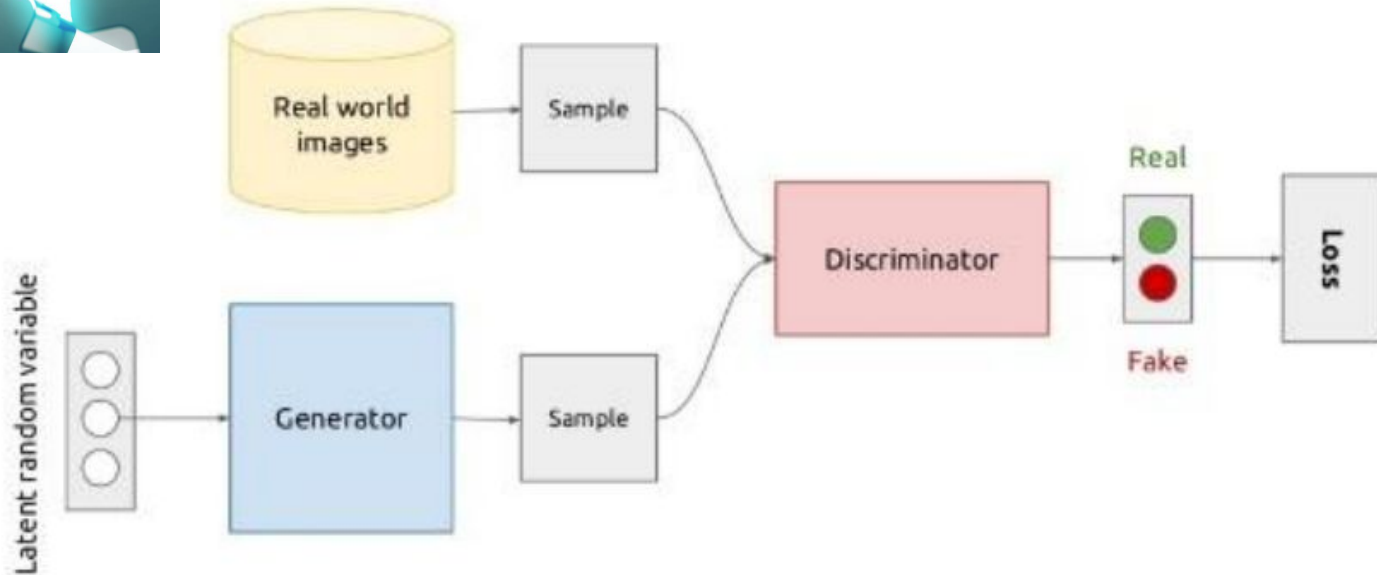
C o s m o s t a t i s t i c s I n i t i a t i v e

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

Extra slides

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>

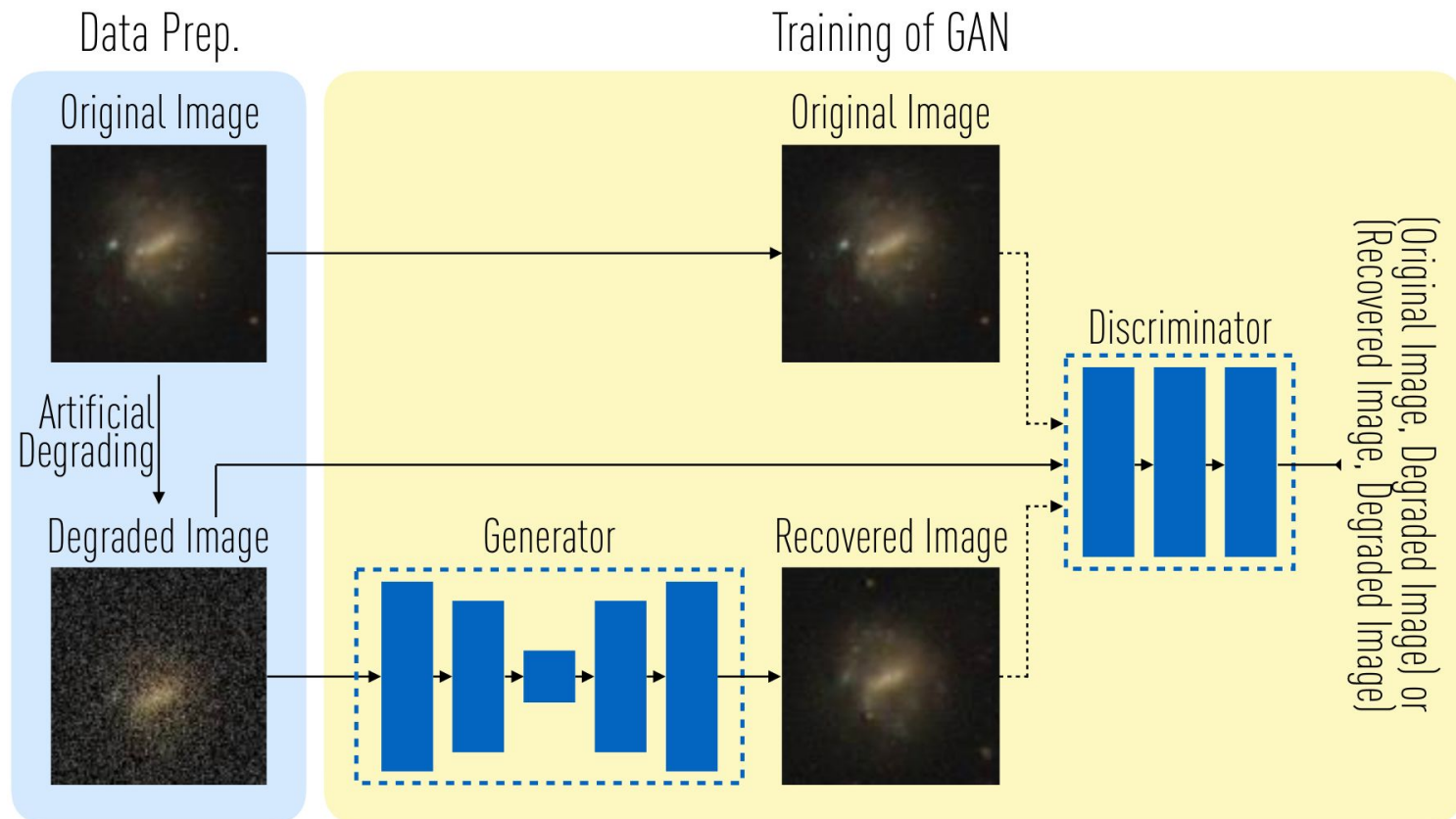
<https://mascherari.press/introduction-to-adversarial-machine-learning/>

Adversarial Learning

The benefits of a worthy opponent

K. Schawinski et al, 2017

In Astronomy

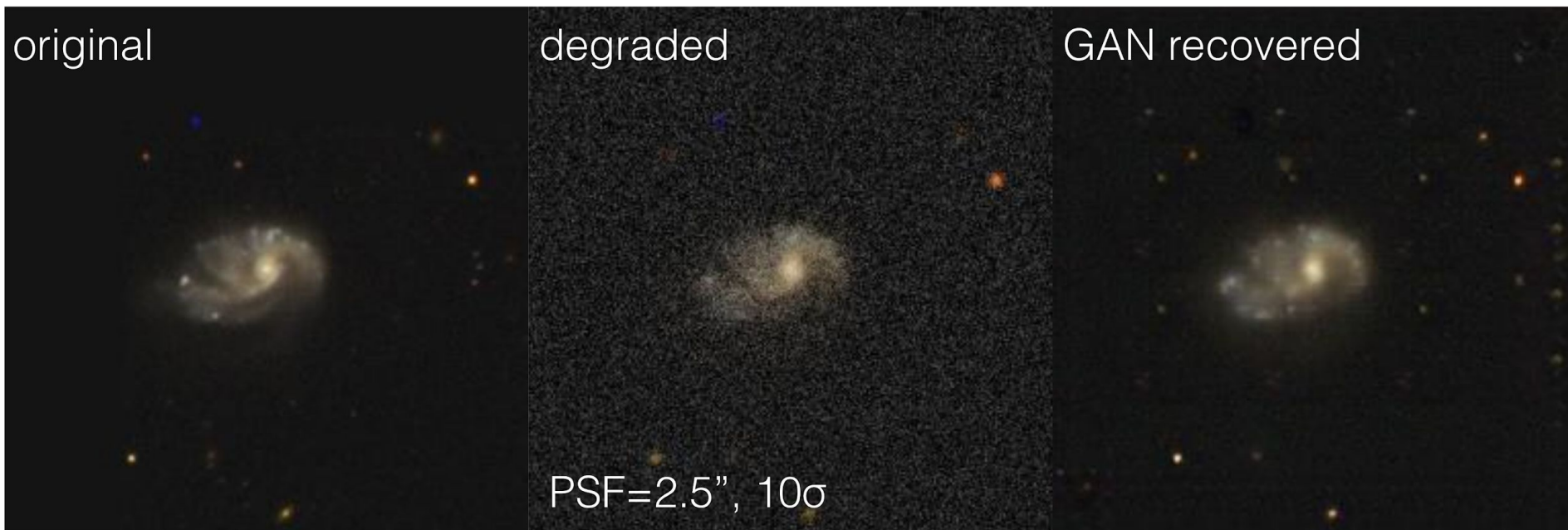


Adversarial Learning

The benefits of a worthy opponent

K. Schawinski et al, 2017

In Astronomy



Adversarial Learning

The benefits of a worthy opponent

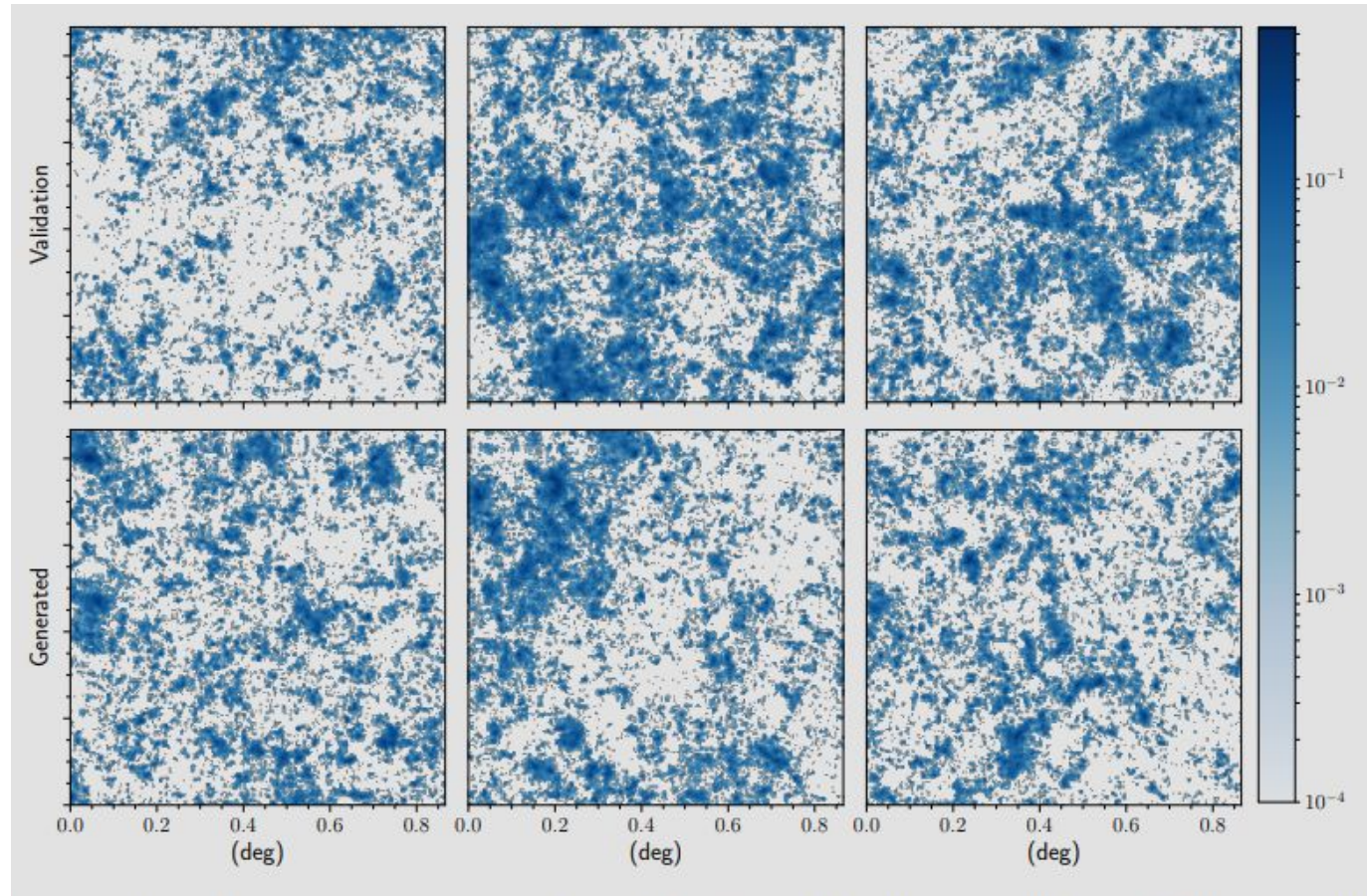
Mustafa et al, 2017 - CosmoGAN

Generating cheap WL maps

In Cosmology

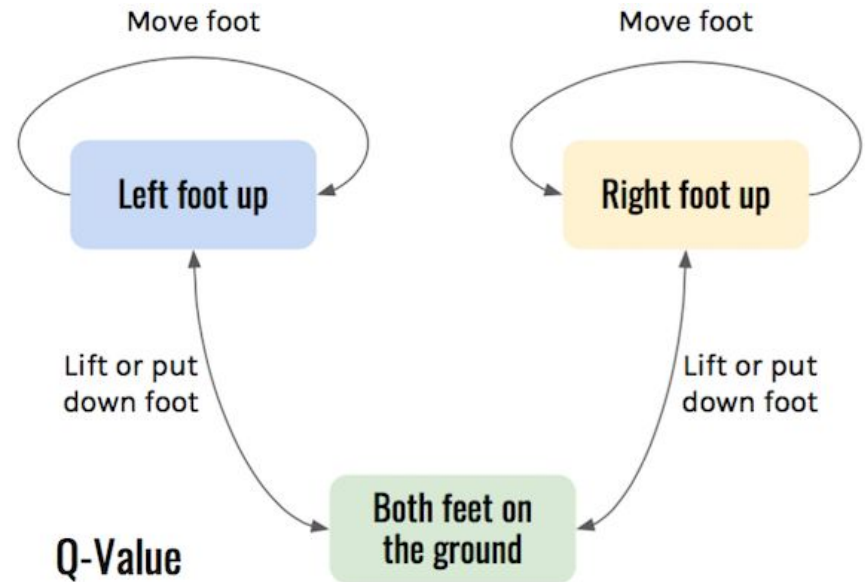
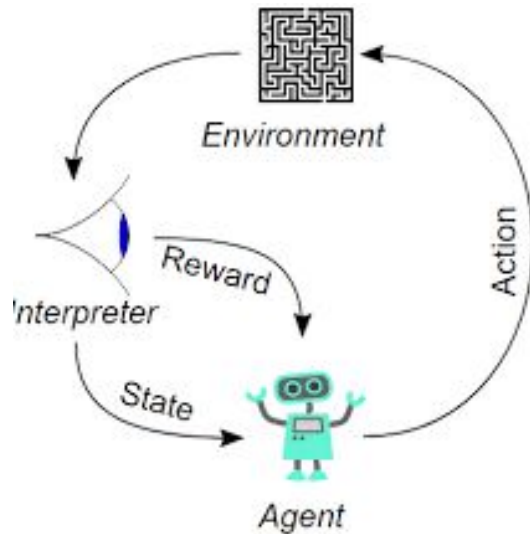
Original

Generated



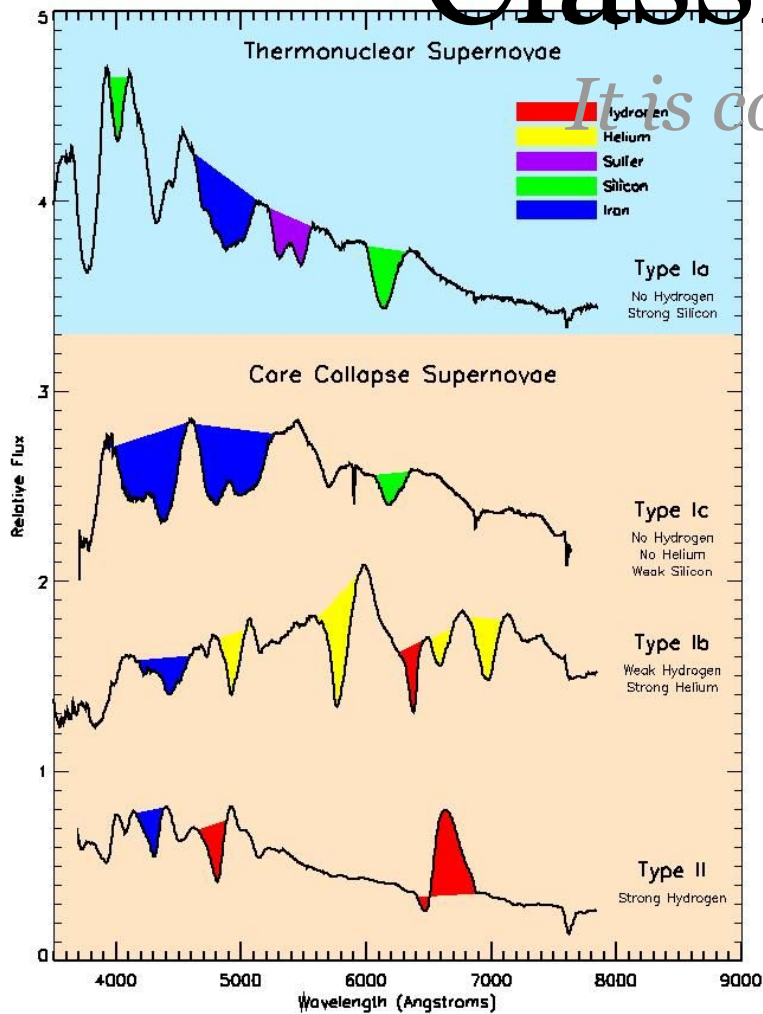
Reinforcement Learning

The importance of feedback

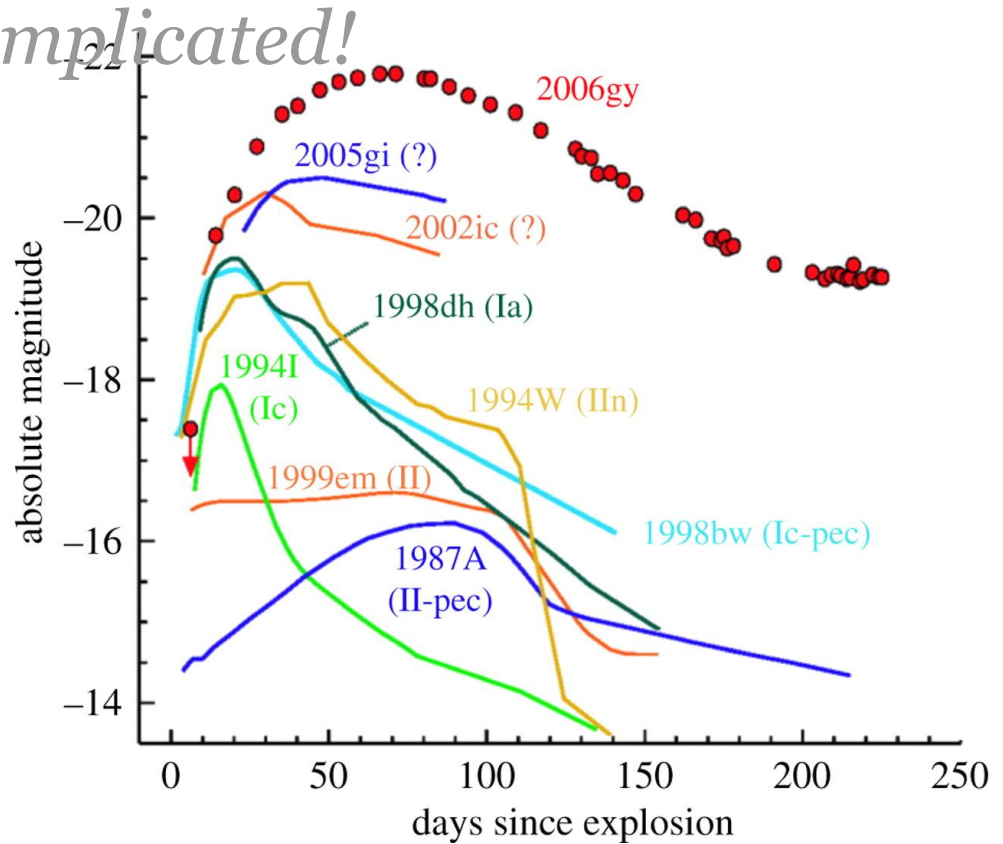


State	Action	Q-Value
Left foot up	Move foot forward	+ 0.5
Left foot up	Put foot down	+ 0.0
Left foot up	Move foot backward	- 0.5
...

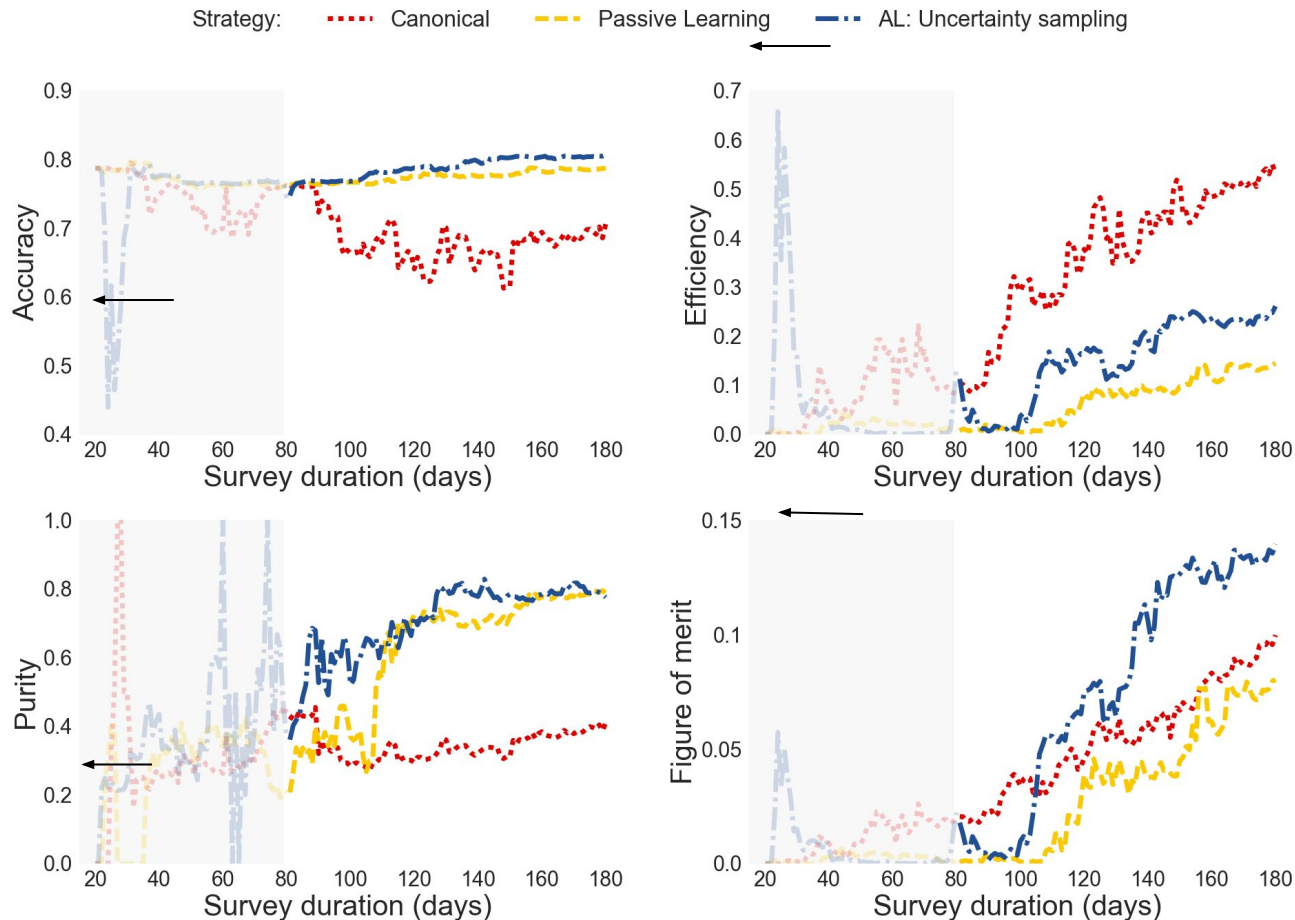
SN Photometric Classification



It is complicated!



No initial training



The arrow shows traditional Full light-curve results with full SNPCC spec

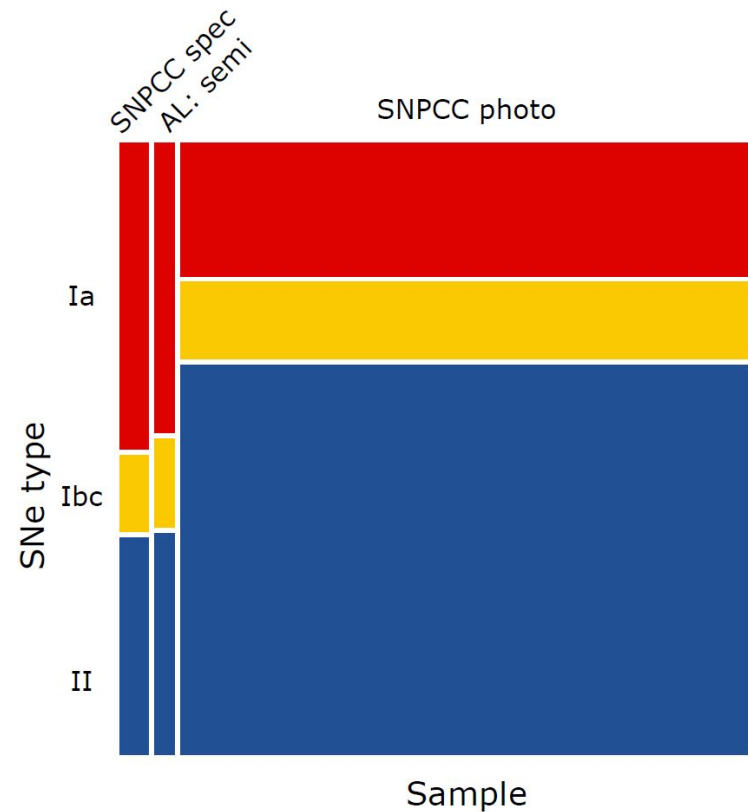
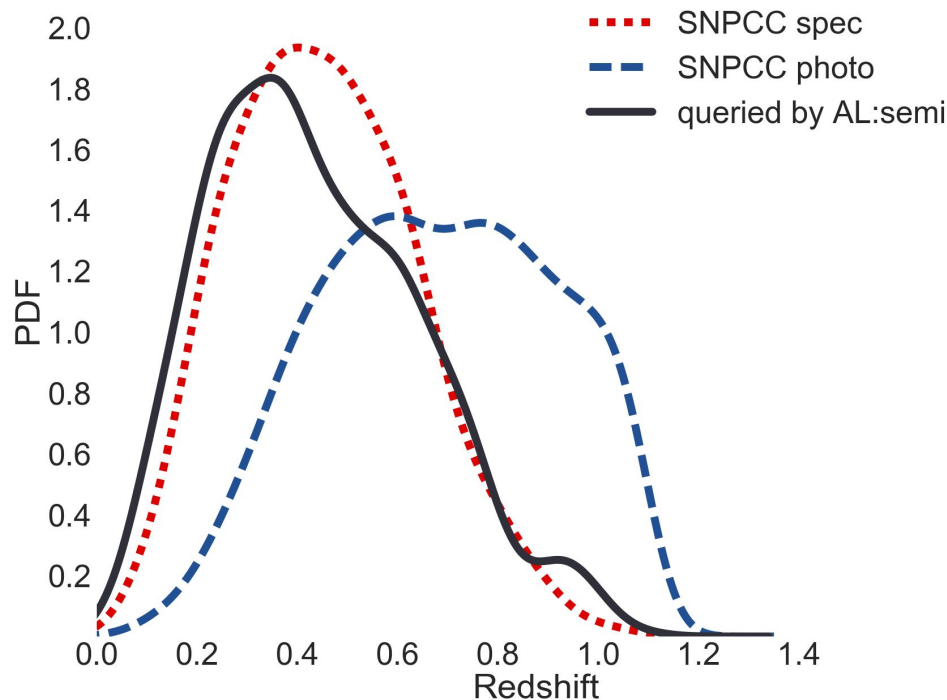
The queried sample

Partial LC, no training, time domain, batch

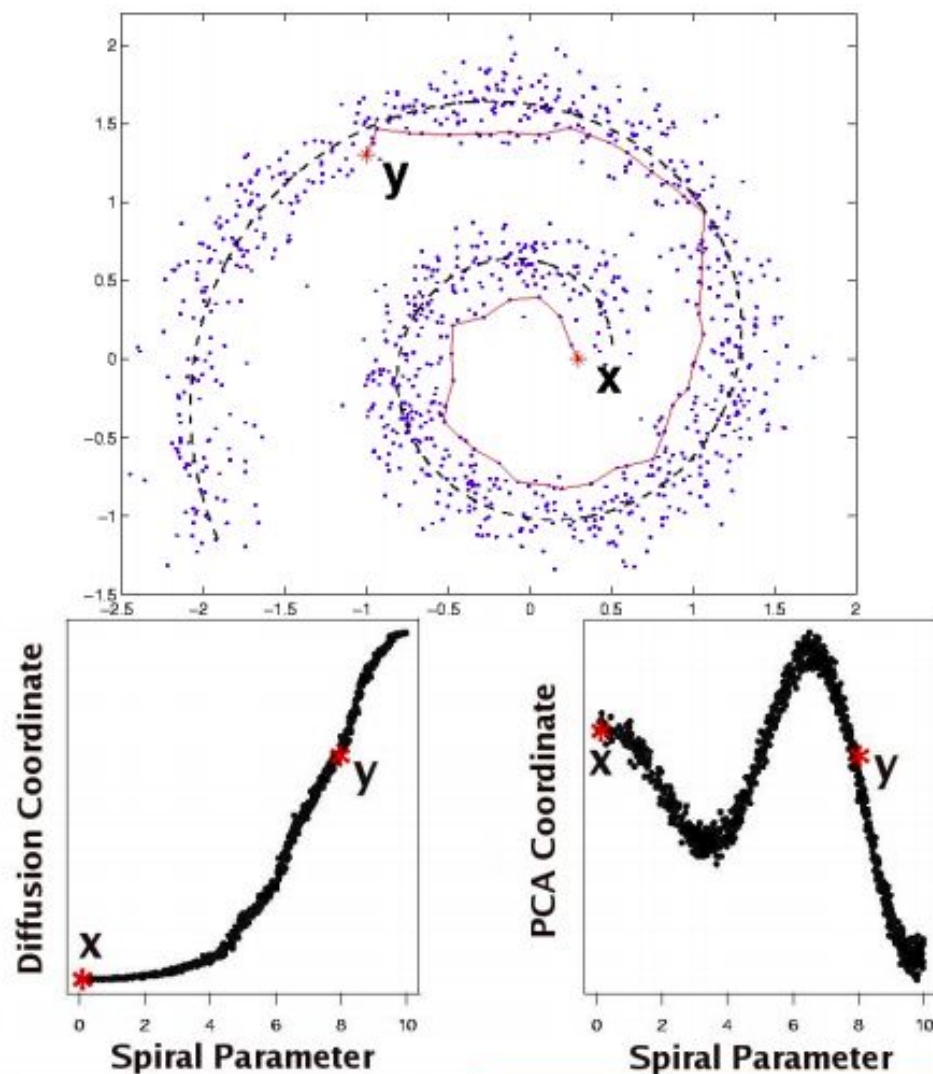
SNPCC spec:
1103 objects

Queried sample:
800 objects

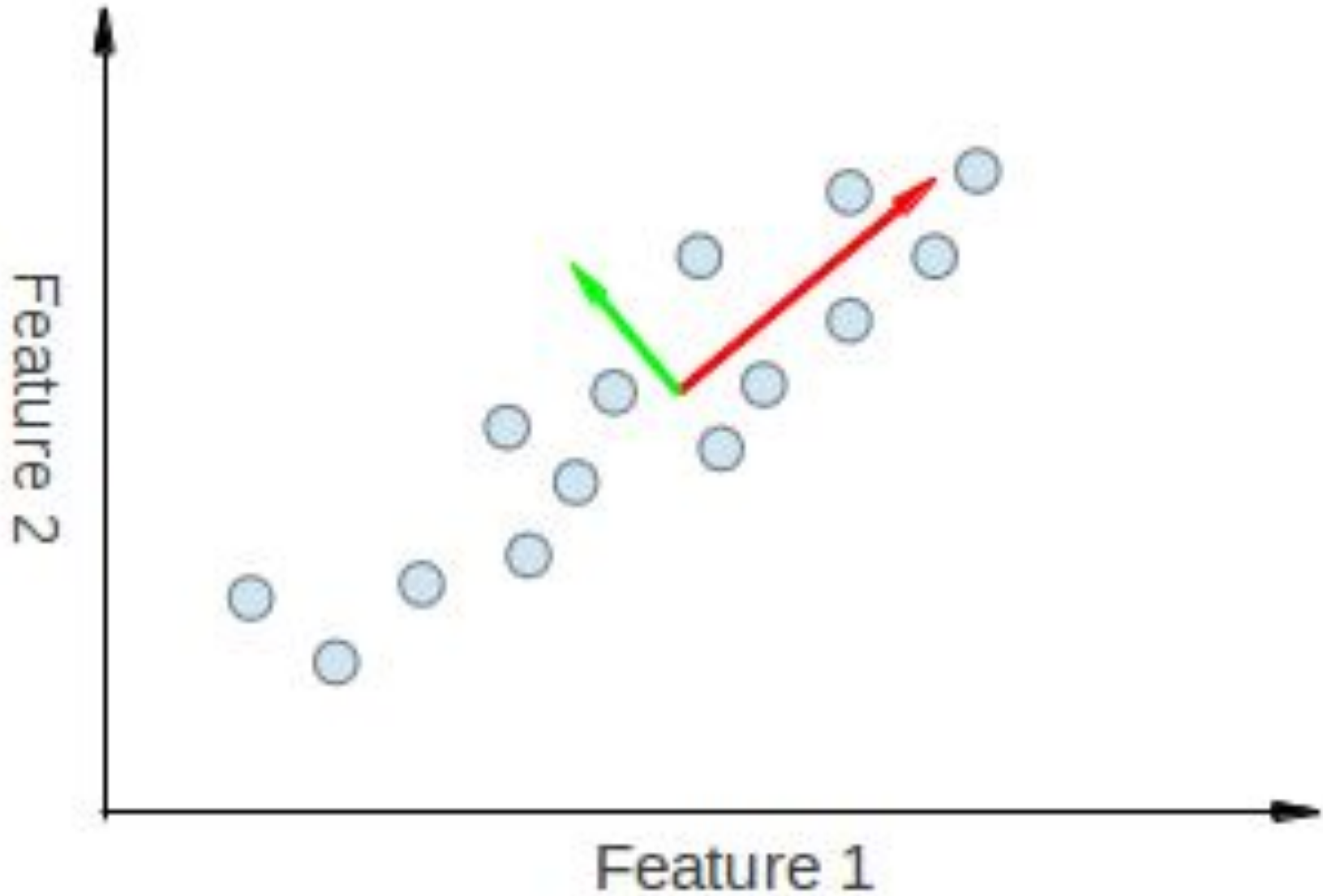
Telescope time:
Queried/spec = 0.999



Diffusion Map: Spiral Example



Principal Component Analysis



Acknowledgement

- H2020-Astronomy ESFRI and Research Infrastructure Cluster (Grant Agreement number: 653477).

