



<https://github.com/COINtoolbox>



The Cosmostatistics Initiative

MAESTRO, June/2017

Emille E. O. Ishida

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

Long term:

Contribute to the establishment of
Astrostatistics and Astroinformatics as full
fledged scientific disciplines

Long term:

Contribute to the establishment of
Astrostatistics and Astroinformatics as full
fledged scientific disciplines **ASAP!**

Long term:

Contribute to the establishment of
Astrostatistics and Astroinformatics as full
fledged scientific disciplines **ASAP!**

Short term:

Make astronomers, statisticians, computer scientists and
data experts understand each other ...
WHILE doing science!



COIN's activities
cannot be merely
pedagogical

Long term:

Contribute to the establishment of
Astrostatistics and Astroinformatics as full
fledged scientific disciplines **ASAP!**

Short term:

Make astronomers, statisticians, computer scientists and
data experts understand each other ...
WHILE doing science!

Directive:

Significantly contribute to the CV of our members

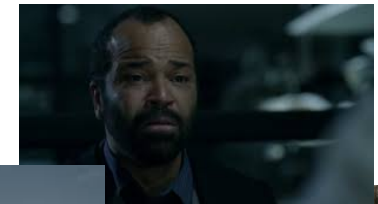
Long term:

Contribute to the establishment of
Astrostatistics and Astroinformatics as full
fledged scientific disciplines **ASAP!**

Short term:

Make astronomers, statisticians, computer scientists and
data experts understand each other ...
WHILE doing science!

Try to remember:
they might work as robots, but they are not!



The COIN Residence Program - CRP

Annual meetings

Conference

Workshop

Hackathon

The COIN Residence Program - CRP

Annual meetings

~~Conference~~

Workshop

Hackathon

The COIN Residence Program - CRP

Annual meetings

~~Conference~~

~~Workshop~~

Hackathon

The COIN Residence Program - CRP

Annual meetings

~~Conference~~

~~Workshop~~

~~Hackathon~~

The COIN Residence Program - CRP

A non-profit start-up?

Annual meetings



John Johnson/HBO



<https://www.theroadtosiliconvalley.com/moving/comparing-sydney-silicon-valley/>

The COIN Residence Program - CRP

A non-profit start-up?

Annual meetings



John Johnson/HBO



CRP #2, UK, 2015



CRP #3, Budapest, 2016



<https://www.theroadtosiliconvalley.com/moving/comparing-sydney-silicon-valley/>





Choosing the participants

The COIN Residence Program



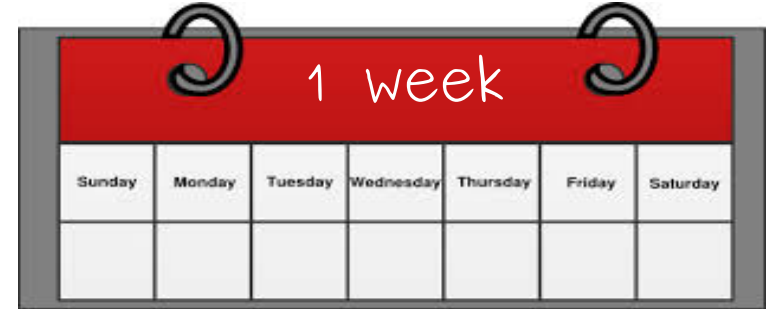
The COIN Residence Program

Once a year



Who wants to collaborate?

What we can guarantee up front



Lots of coffee



paper



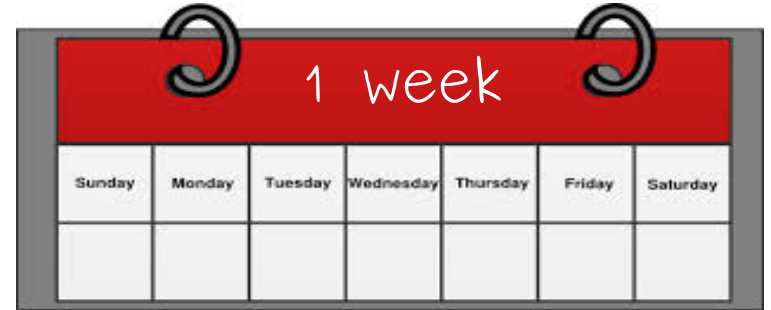
The COIN Residence Program

Once a year



Who wants to collaborate?

What we can guarantee up front



Lots of coffee



paper



What we can NOT guarantee up front



Talks



CATERING



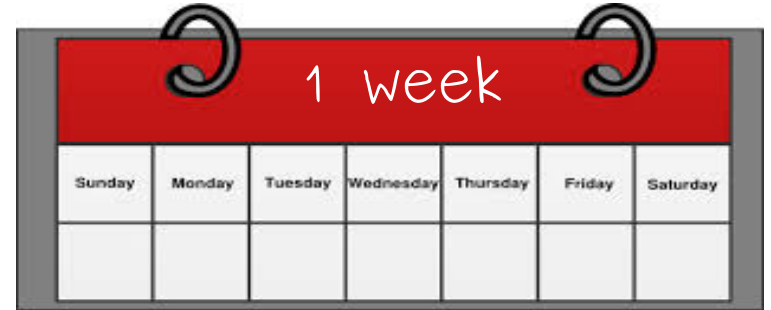
The COIN Residence Program

Once a year



Who wants to collaborate?

What we can guarantee up front



Lots of coffee



paper



What we can NOT guarantee up front



Talks



What we require from participants





Choosing the projects

Questions posed by the organizers

1. What do you know how to do?
2. What do you like to do?
3. What would you like to learn?

Participants can propose projects.
Everyone votes to which project will be selected

Questions posed by the organizers

1. What do you know how to do?
2. What do you like to do?
3. What would you like to learn?

Participants can propose projects.
Everyone votes to which project will be selected

From CRP #2, UK - 2015

<p>1</p> <p>Proj. Name: HBM SNe host galaxies Method: <u>Hierarchical Bayesian Models</u> Object: SNe/GRB Data: TBD Manager : TBD</p> <p>Scientific Question: Can HBM improve our understanding of the bias between type Ia SN hubble residuals and their host galaxy properties?</p> <p>Participants: Heather, Rafael, Emille, Joseph, Mohammad, Paniez, Miguel, Madhura</p>	<p>2</p> <p>Project type: application Method: <u>Spatial Statistics - INLA</u> Object: TBD Data: Spatio-Temporal evolution of chemical abundances in primordial galaxies</p> <p>Scientific Question: Explore Spatial Statistics (spatstat, INLA) capabilities for astronomical research</p> <p>Participants: Marina, Mariana, Sandro, Alberto, Rafael, Ewan, Joseph, Miguel, Mohammad, Eric?</p>	<p>3</p> <p>Proj. Name: Galaxy AGN connection Method: <u>Bayesian Logistic Regression</u> Objects: Galaxy bars, rings, AGNs ... Data: TBD Manager: TBD</p> <p>Scientific Question: How do the properties of the galaxy influence its probability to host an AGN?</p> <p>Participants: Marina, Rafael, Alberto, Joseph, Mohammad, Arlindo, Malu, Paula</p>
<p>4</p> <p>Proj. Name: Review on Variable selection Method: <u>different variable selection algorithms</u> Object: Data: Manager:</p> <p>Scientific Question: What is the state of art methodology to subset the best predictors for multivariate regression in astronomical datasets?</p> <p>Participants: Alberto, Rafael, Zoe, Luke, Paniez, Miguel, Arlindo, Bruce, Alan, Yabebal, Malu</p>	<p>5</p> <p>Proj. Name: MACHine learning in SNe spectra Method: <u>Dim. Reduction + unsupervised learning</u> Object: SNe Data: SN spectral series Manager: Michele</p> <p>Scientific Question: Use unsupervised/ semi-supervised learning to identify subtypes of SNe</p> <p>Participants: Michele, Emille, Rafael, Ricardo, Paolo, Fabian, Arlindo, Paniez</p>	<p>6</p>



Does it work?

COIN products



Rafael S. de Souza
(head) - statistics



Alberto Krone-Martins
astrometry



Emille E. O. Ishida
SN cosmology

60 researchers from 15 countries

Scientific outcomes

In 3 years

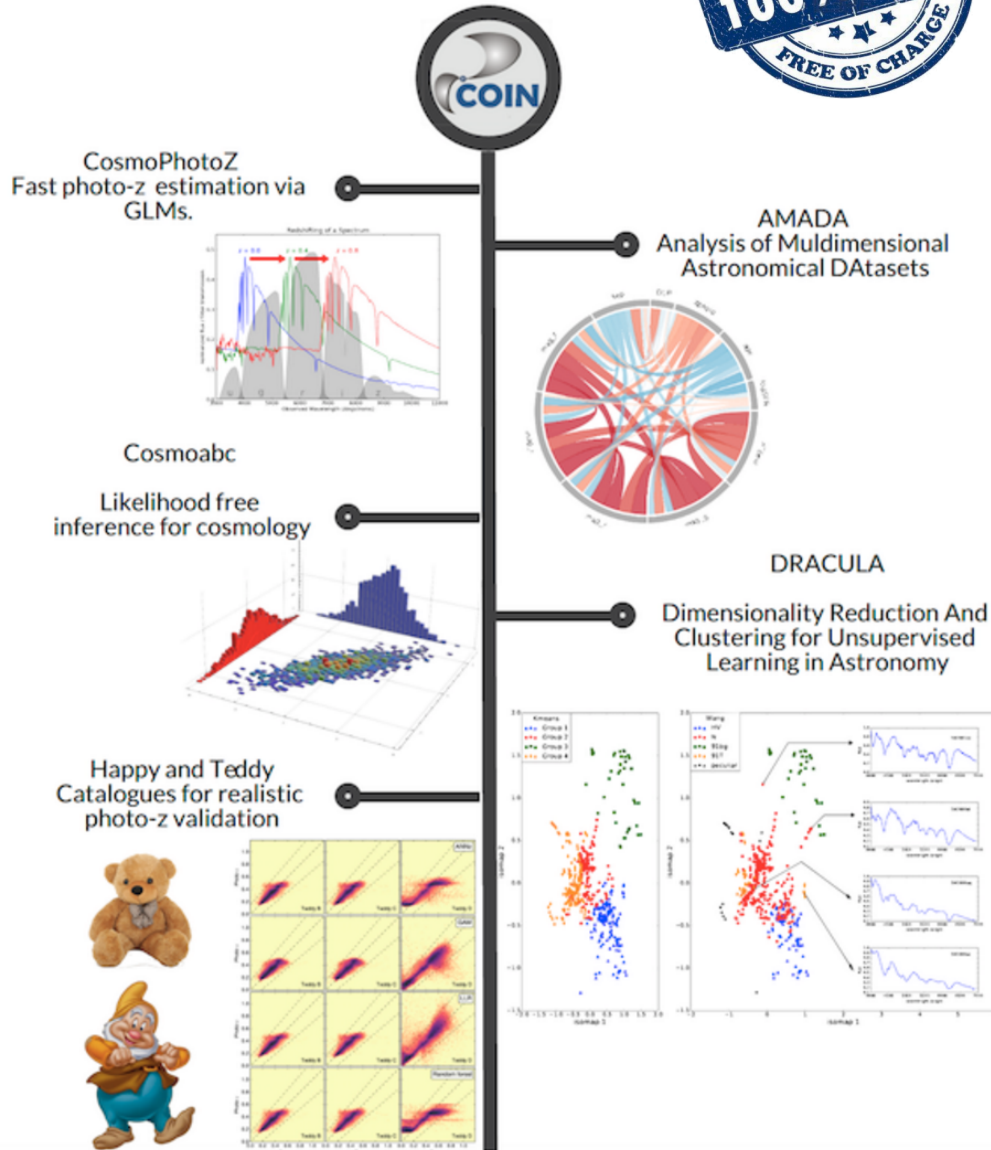


	Paper	Citation
1	GLM I	de Souza <i>et al.</i> , 2015
2	GLM II	Elliott <i>et al.</i> , 2015
3	GLM III	de Souza <i>et al.</i> , 2015
4	AMADA	de Souza & Ciardi, 2015
5	CosmoABC	Ishida <i>et al.</i> , 2015
6	DRACULA	Sasdelli <i>et al.</i> , 2016
7	AGNlogit	de Souza <i>et al.</i> , 2016
8	PhotoZ	Beck <i>et al.</i> , 2017
9	AGNgmm	de Souza <i>et al.</i> , 2017



1	CosmoPhotoZ	de Souza <i>et al.</i> , 2014,
2	AMADA	de Souza & Ciardi, 2015
3	CosmoABC	Ishida <i>et al.</i> , 2015
4	DRACULA	Aguena <i>et al.</i> , 2015

- + 1 galaxy catalog
- + 1 GMM tutorial
- + 2 photoz catalogs



COIN products are open source!



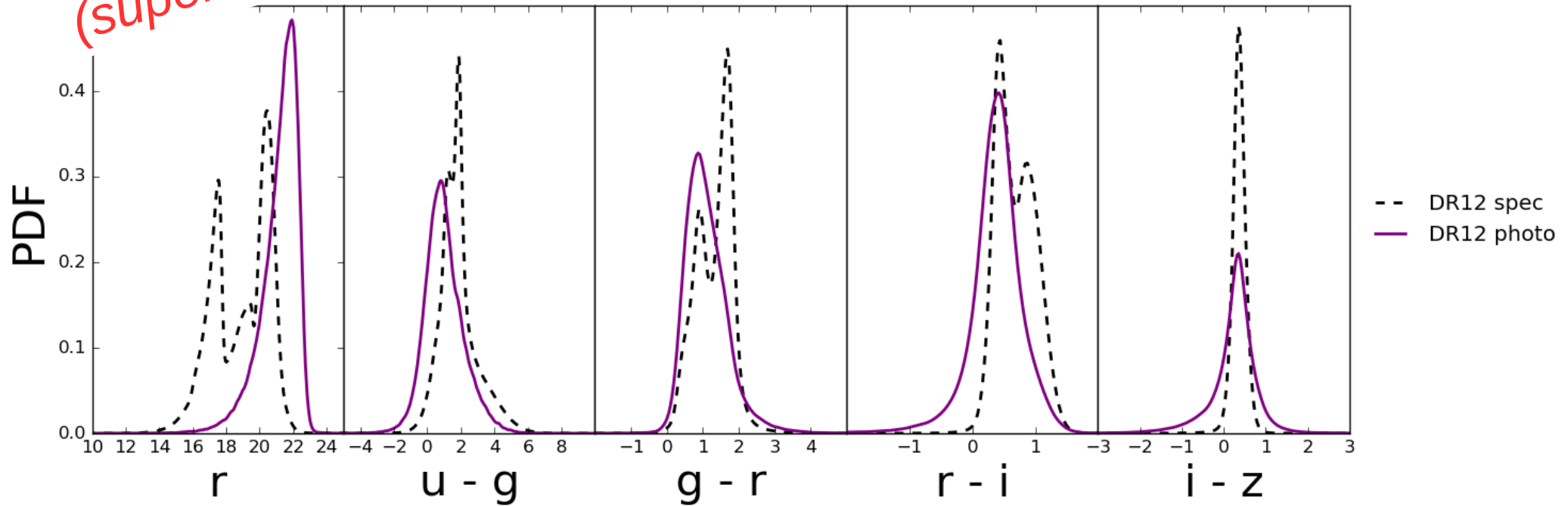
Example from
CRP #4, Budapest
2016

The problem with text-book ML:

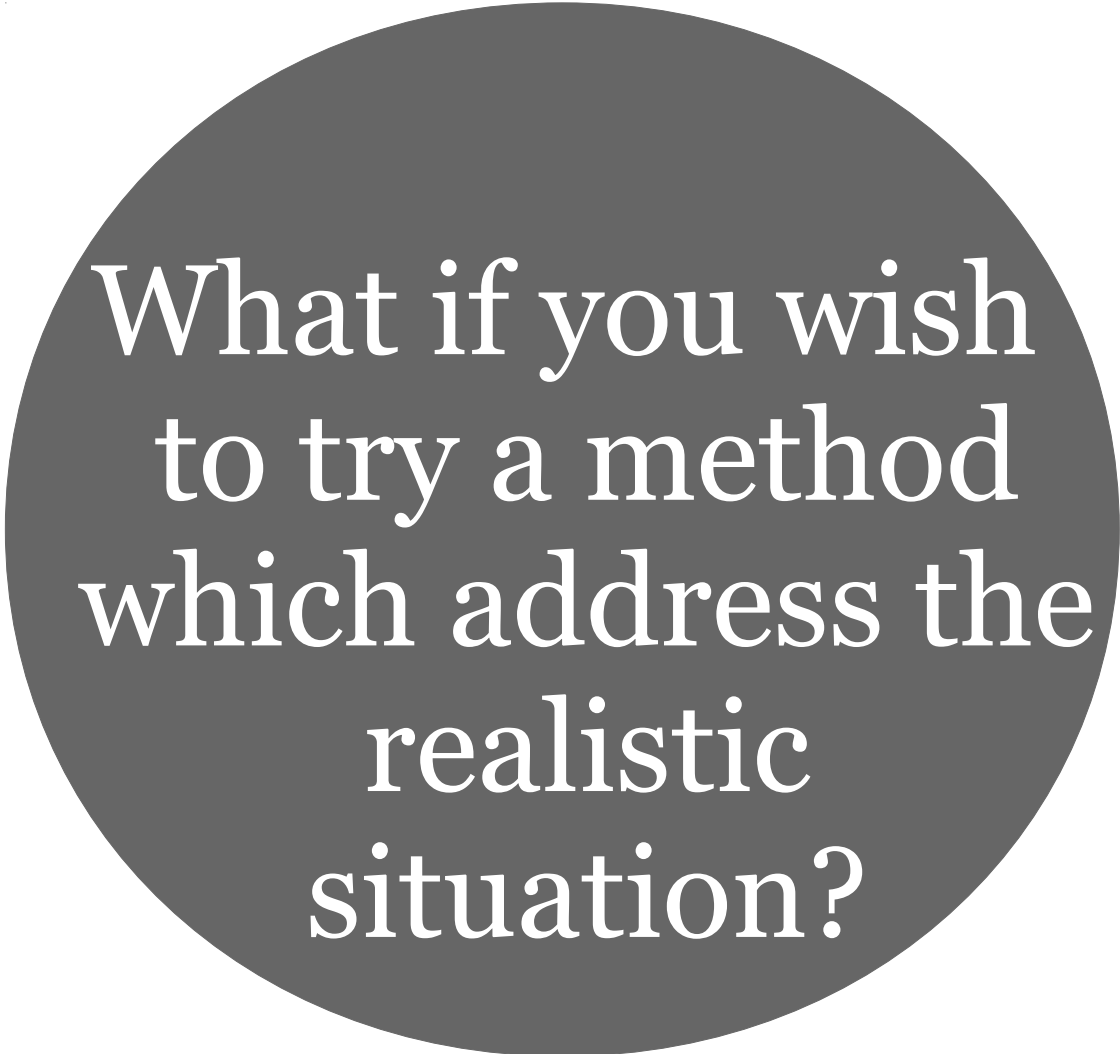
Representativeness

For photo-z
(supervised regression)

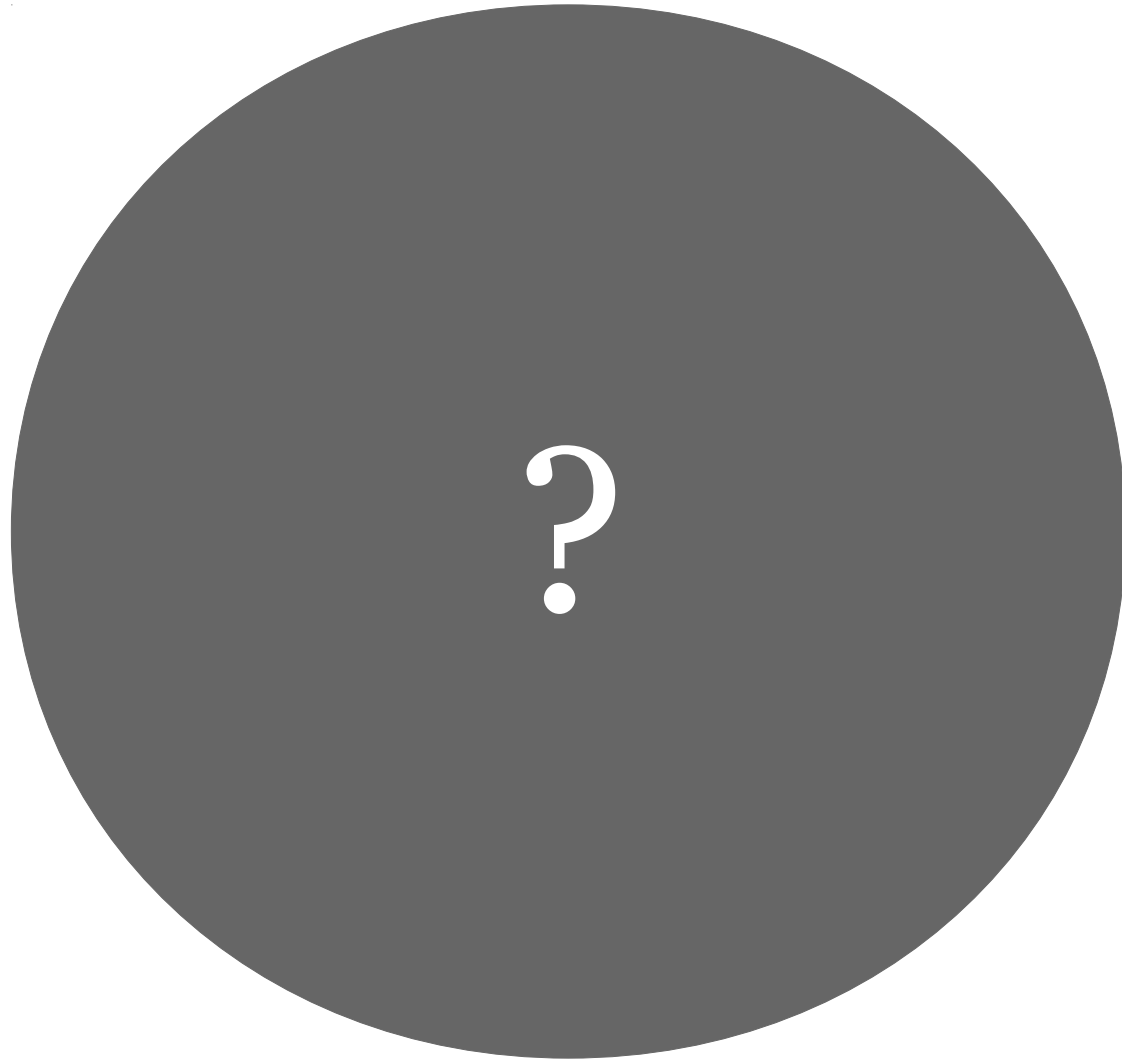
Features (columns)
distributions



Spec → training/validation
Photo → test



What if you wish
to try a method
which address the
realistic
situation?

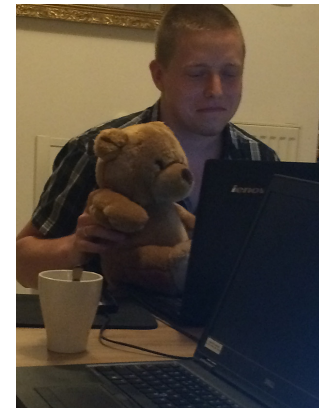


COIN Residence Program #3

<http://iaacoin.wixsite.com/crp2016>

COIN Residence Program 2016

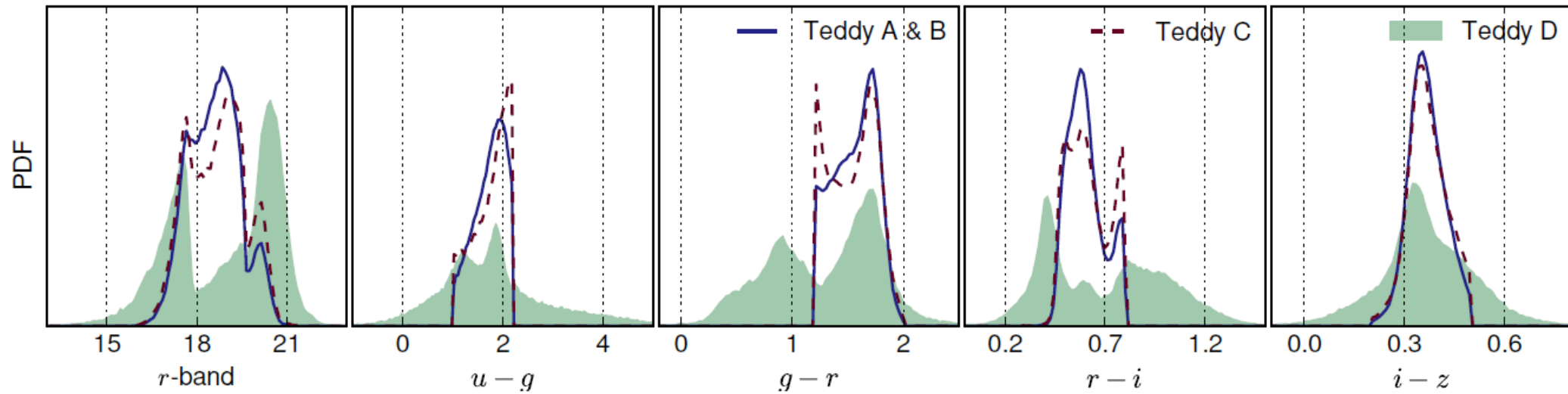
August 21-28 - Budapest, Hungary



Teddy catalogue

The effect of color coverage

A/B follow SDSS spec distribution
B is completely representative of A
C has the same coverage but slightly different shape
D has a wider domain in r-mag and color (no coverage)



Teddy



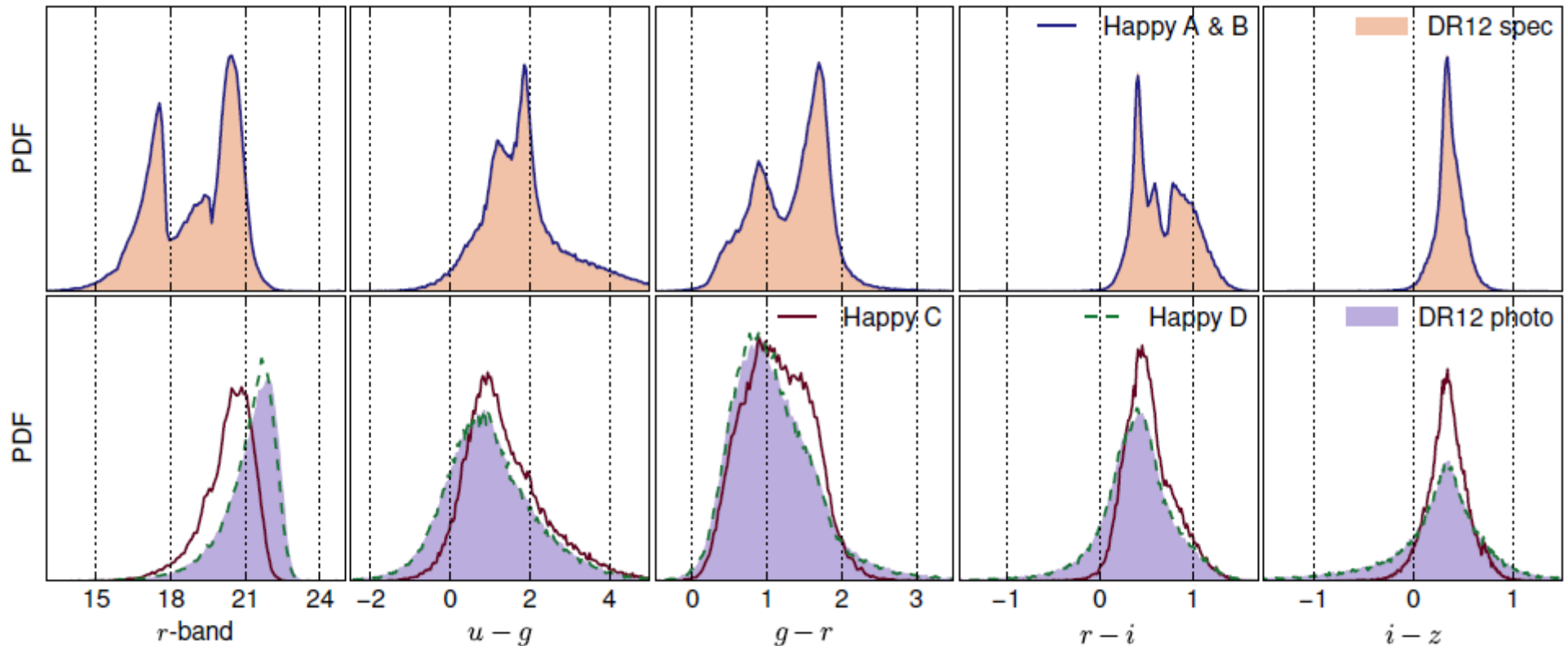
Happy catalogue

The effect of coverage + photometric errors

Photometry from SDSS

Spec-z from many different surveys leads to larger photometric errors and consequently wide domain in r-band and color

- A /B follow SDSS spec distribution
- B is completely representative of A
- C was constructed performing a nearest neighbor between the SDSS-DR12 photo sample and the extended spec sample but with a cut on photometric errors
- D is the same of C but without the photometric error cut.
- Consequently, D follows exactly the SDSS-DR12 photo sample distribution



Happy

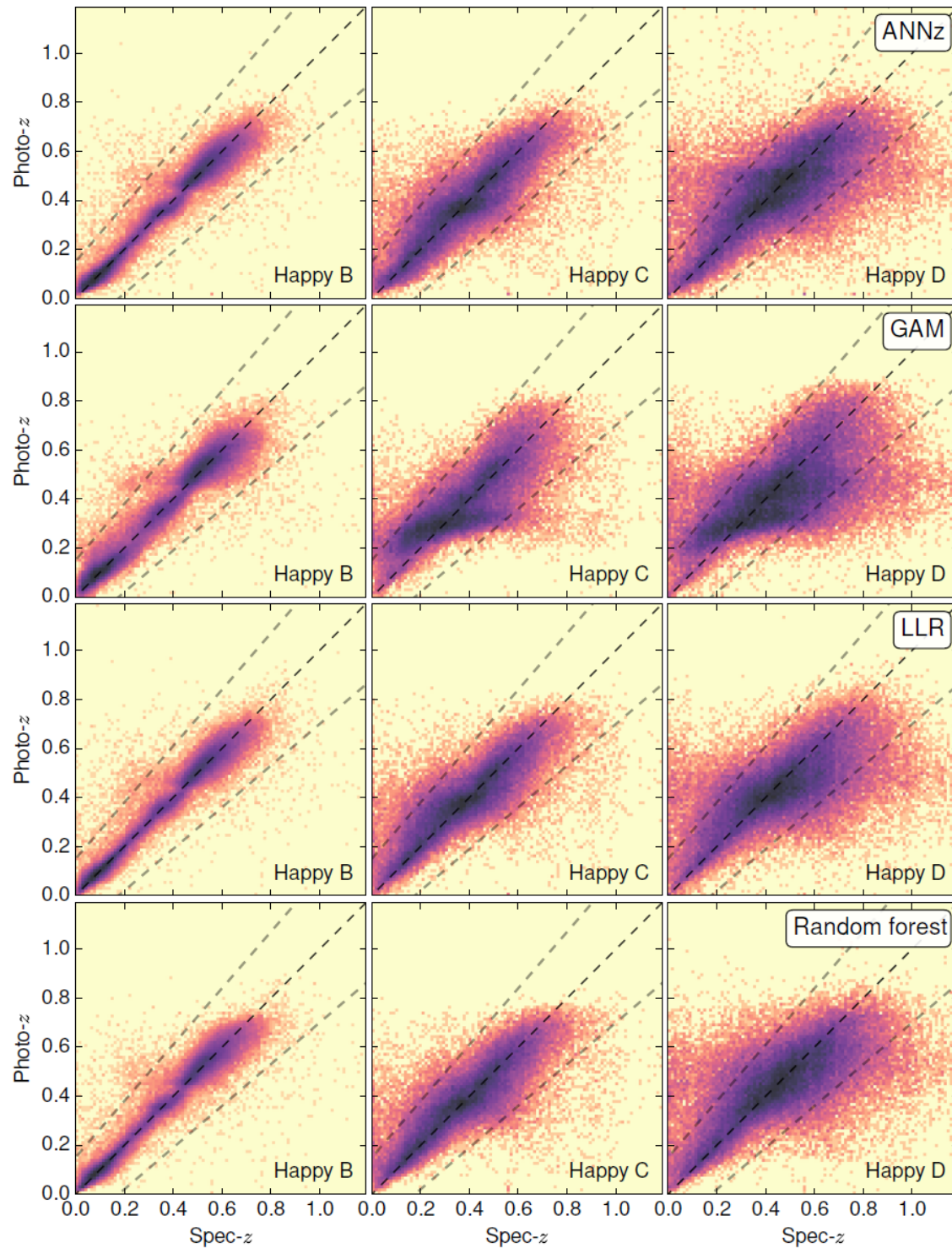


Happy catalogue

*The effect of coverage +
photometric errors*



Empirical methods



Method	Set	Diagnostics			Outlier rate (%)
		Mean ($\times 10^{-2}$)	Std ($\times 10^{-2}$)	MAD ($\times 10^{-2}$)	
ANNz	B	0.04	2.87	1.49	0.99
	C	0.16	5.41	3.60	5.59
	D	-0.52	6.53	5.44	14.01
GAM	B	0.09	3.50	1.95	1.36
	C	0.86	6.34	4.84	7.37
	D	-0.51	7.21	6.70	16.38
LLR	B	0.13	2.81	1.39	1.11
	C	0.52	5.45	3.59	6.07
	D	-0.79	6.62	5.62	14.52
Random Forest	B	0.05	2.82	1.41	1.02
	C	0.34	5.39	3.51	5.58
	D	-0.28	6.51	5.36	14.2

https://github.com/COINtoolbox/photoz_catalogues

README.md

Readme minor

5 months ago

README.md

Teddy and Happy photo-z catalogues

This repository contains the photometric redshift catalogues presented in [Beck et al., 2017](#) - *On the realistic validation of photometric redshifts, or why Teddy will never be Happy.*

This is one of the products of the third edition of the [COIN Residence Program](#), which took place in August/2016 in Budapest (Hungary).

Any questions/suggestions should be sent to iaa.coin@gmail.com.

A general overview of both catalogues is given bellow. Check the individual folders for detailed information on the files presented here.

Teddy

This catalogue was designed to isolate the effect of limited spectroscopic sample coverage in colour/magnitude space.

It is constructed from the [SDSS DR12](#) spectroscopic sample and is maintained by [Chieh-An Lin](#) (CEA, France)

Happy

This catalogue was designed to reproduce the effect of distinct photometric error distributions and their convolution with colour/magnitude space coverage between the spectroscopic and photometric samples.

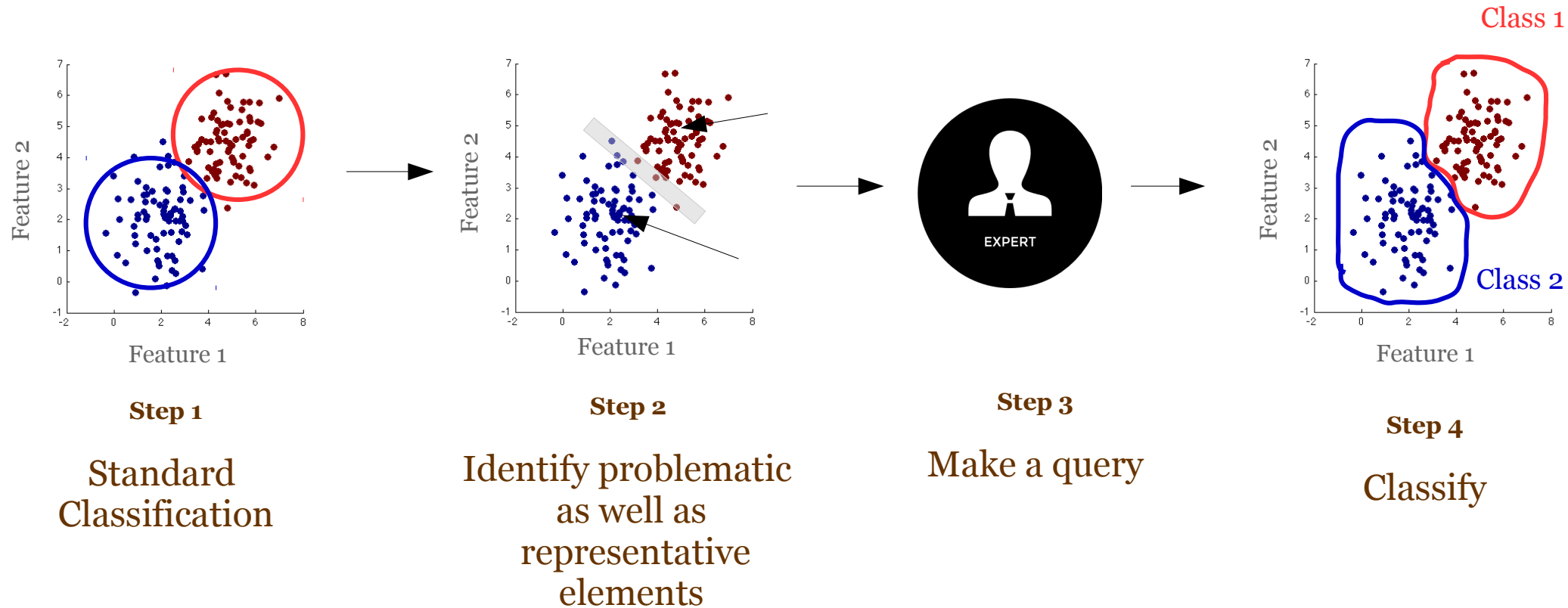
All photometry was taken from [SDSS DR12](#), and spectroscopy was gathered from a set of different sources (see [Beck et al., 2017](#) for further details).

Happy is maintained by [Robert Beck](#) (ELTE, Hungary).



Potential solution: Active Learning

SAMSI & COIN, in prep





The future of
COIN

<https://iaacoin.wixsite.com/crp2017>

#coinCF2017

Home About Organizers Location Conduct Participants Sponsors



COIN Residence Program #4

20-27 August 2017
Clermont Ferrand, France

Find people ...



Find people ...



solve bureaucracy

Find people ...



solve bureaucracy



Find people ...



solve bureaucracy

A screenshot of the website header for "The Cosmostatistics Initiative (COIN) on Overleaf". The header features a dark blue background with a starry night sky and a white curved line. In the center, there is a white square containing the COIN logo, which consists of a stylized white 'C' and the word "COIN" in blue, bold, capital letters. Below the logo, the text "The Cosmostatistics Initiative (COIN) on Overleaf" is displayed in a white, sans-serif font. At the bottom of the header, there are four navigation links: "Overview", "Quick Start", "Templates", and "FAQ & Help", all in a light blue, sans-serif font.

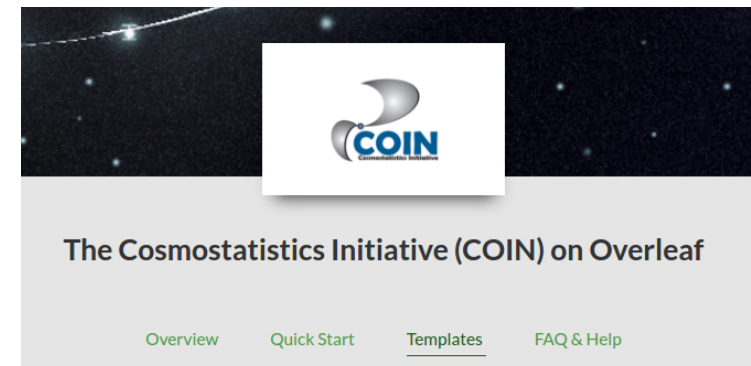
Find people ...



solve bureaucracy



Let it go ...



IAA facebook page

International Astrostatistics Association

Background image by ESO

Organização sem fins lucrativos em Milão
5.0 ★★★★★

Dicas da Página Ver tudo

Como criar publicações eficientes
Publicações pequenas, visuais, criadas para o público certo têm mais êxito.

COIN on twitter

IAA-COIN
@iaa_coin FOLLOWS YOU

COIN promotes the development of novel statistical tools for astronomy. #stats #astrostatistics #cosmology #python #datascience #astronomy #bigdata

Worldwide
goo.gl/alvZYA
Joined June 2015


Tweet to Message

Registrations will open July 7th!

STAT
4
ASTRO

School of Statistics for Astrophysics 2017: Bayesian Methodology
9-13 October 2017, Autrans (France)

Login



Do we need to re-
think the
academic model?

*Thank
you*

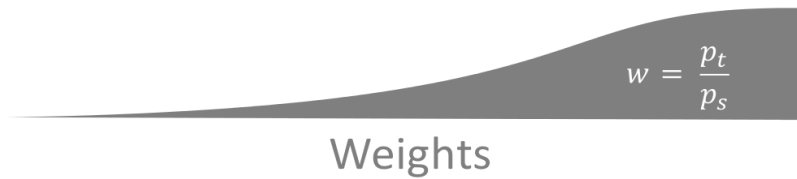
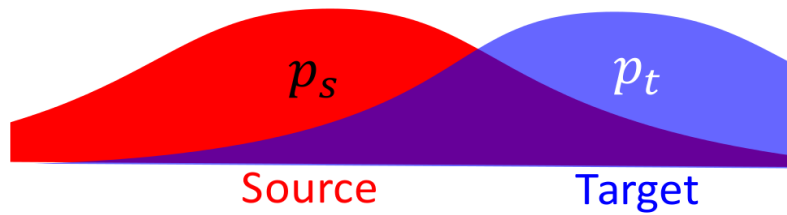
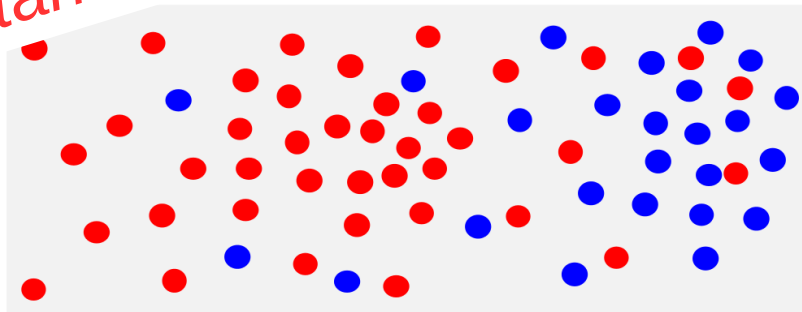


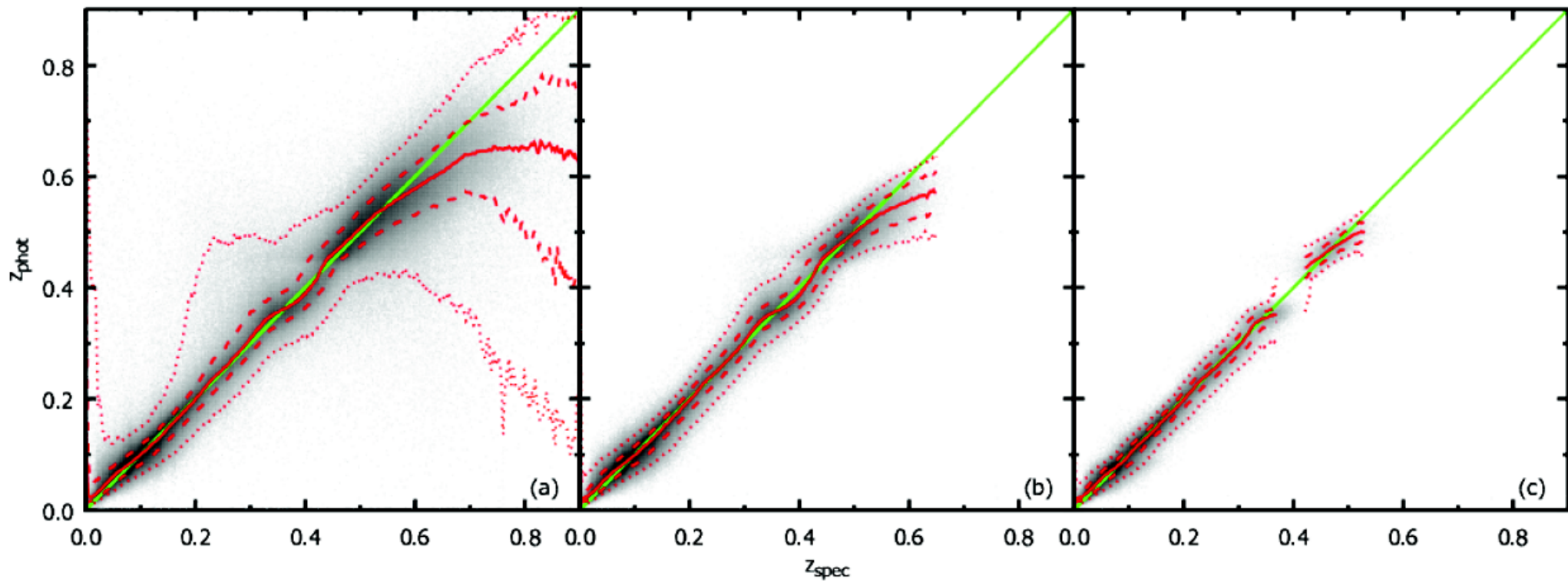


<https://github.com/COINtoolbox>

Attempted solutions: Domain Adaptation

instance weighting

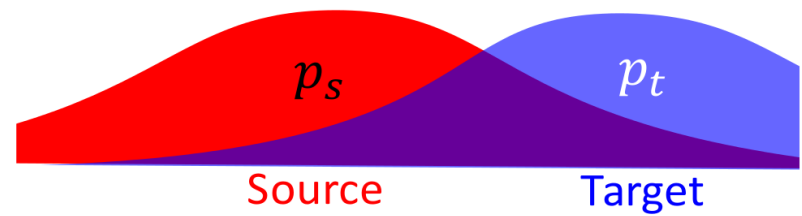
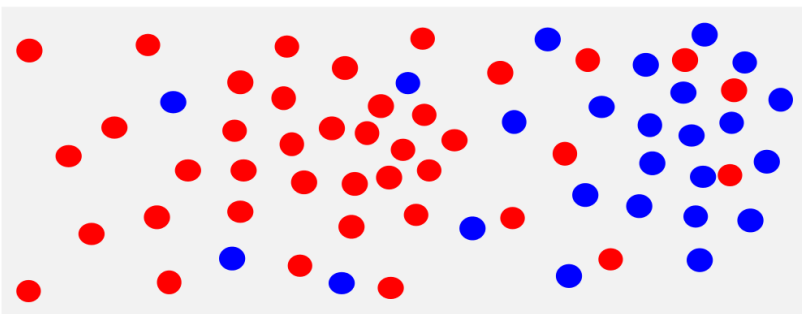




Photometric redshifts for the SDSS Data Release 12

Róbert Beck ✉, László Dobos ✉, Tamás Budavári, Alexander S. Szalay, István Csabai ✉

Mon Not R Astron Soc (2016) 460 (2): 1371-1381.



$$w = \frac{p_t}{p_s}$$

Attempted solution: Domain Adaptation

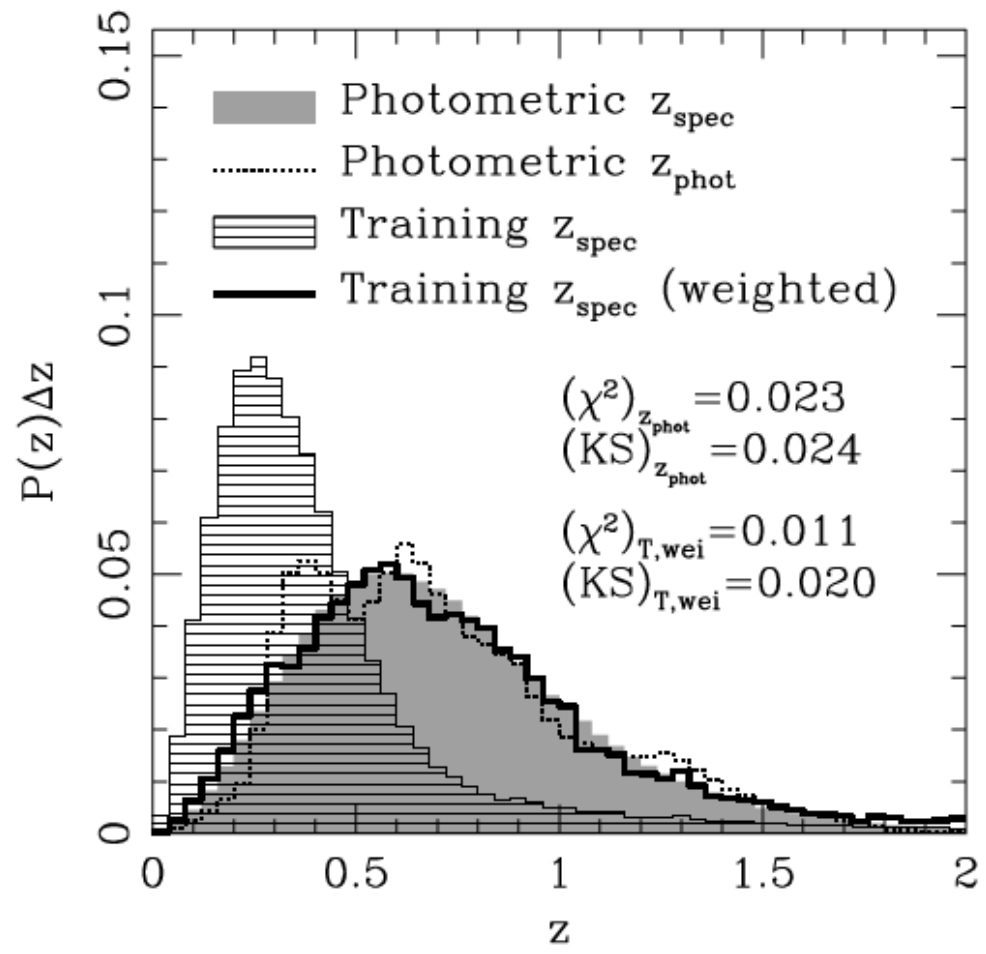
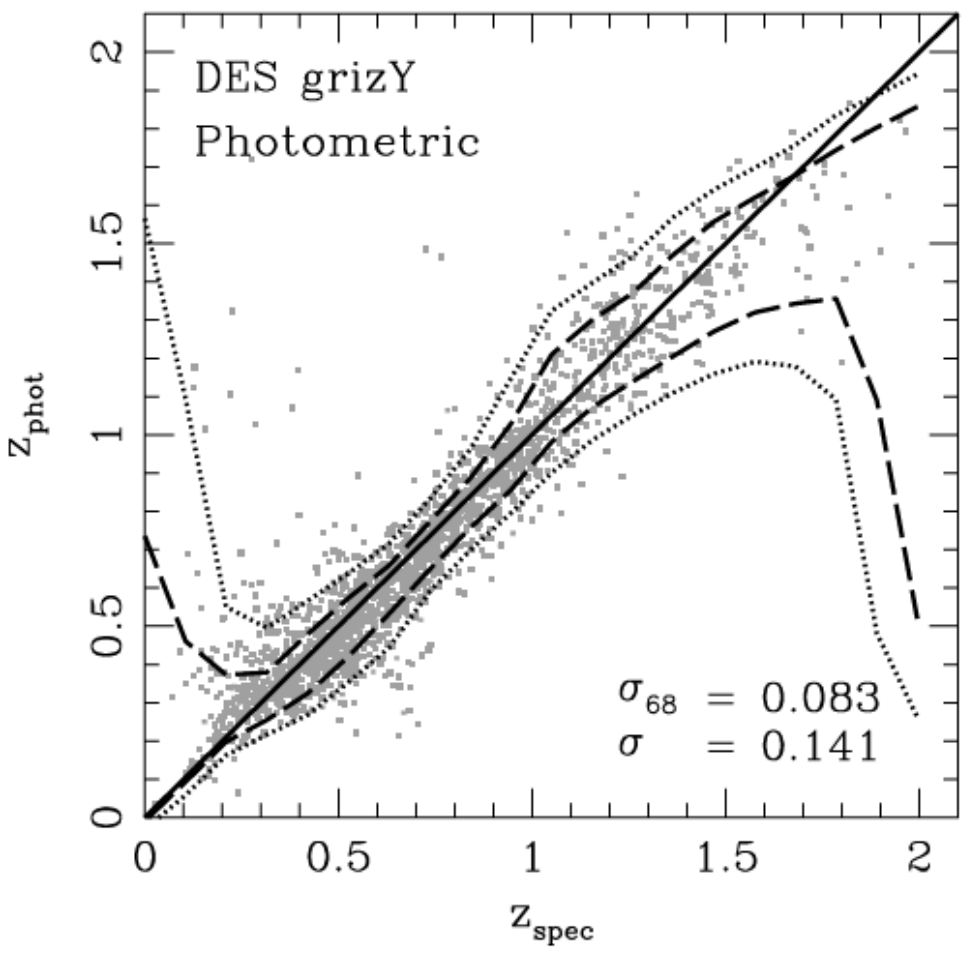
supervised regression

Estimating the redshift distribution of photometric galaxy samples

Marcos Lima,^{1,2*} Carlos E. Cunha,^{2,3} Hiroaki Oyaizu,^{2,3} Joshua Frieman,^{2,3,4}
Huan Lin⁴ and Erin S. Sheldon⁵

Mon. Not. R. Astron. Soc. **390**, 118–130 (2008)

instance weighting



Attempted solutions: Domain Adaptation

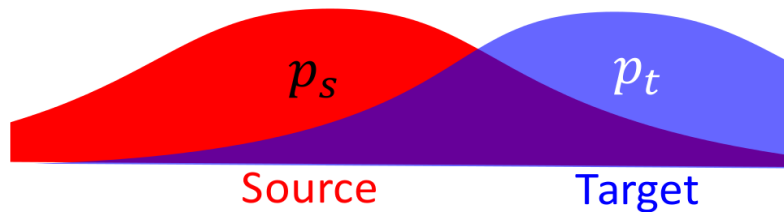
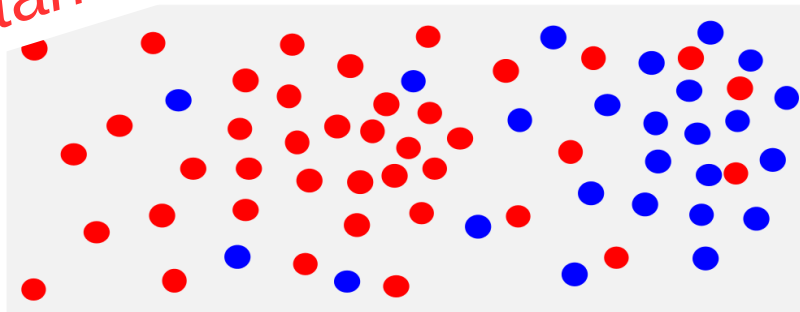
supervised regression

Estimating the redshift distribution of photometric galaxy samples

Marcos Lima,^{1,2*} Carlos E. Cunha,^{2,3} Hiroaki Oyaizu,^{2,3} Joshua Frieman,^{2,3,4}
Huan Lin⁴ and Erin S. Sheldon⁵

Mon. Not. R. Astron. Soc. **390**, 118–130 (2008)

instance weighting



Works well when data is not sparse,
and there is coverage!

Attempted solutions: Domain Adaptation

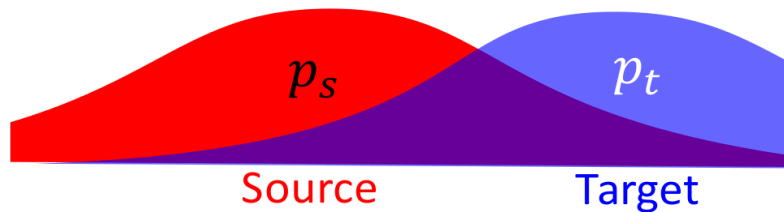
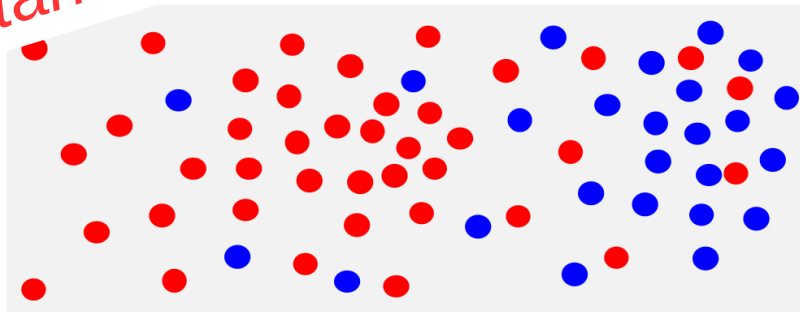
supervised regression

Estimating the redshift distribution of photometric galaxy samples

Marcos Lima,^{1,2*} Carlos E. Cunha,^{2,3} Hiroaki Oyaizu,^{2,3} Joshua Frieman,^{2,3,4}
Huan Lin⁴ and Erin S. Sheldon⁵

Mon. Not. R. Astron. Soc. **390**, 118–130 (2008)

instance weighting



Works well when data is not sparse,
and there is coverage!

Photometric redshifts for the SDSS Data Release 12

Róbert Beck ✉, László Dobos ✉, Tamás Budavári, Alexander S. Szalay, István Csabai ✉

Mon Not R Astron Soc (2016) 460 (2): 1371-1381.

If there is no coverage, identify
problematic areas and discard!

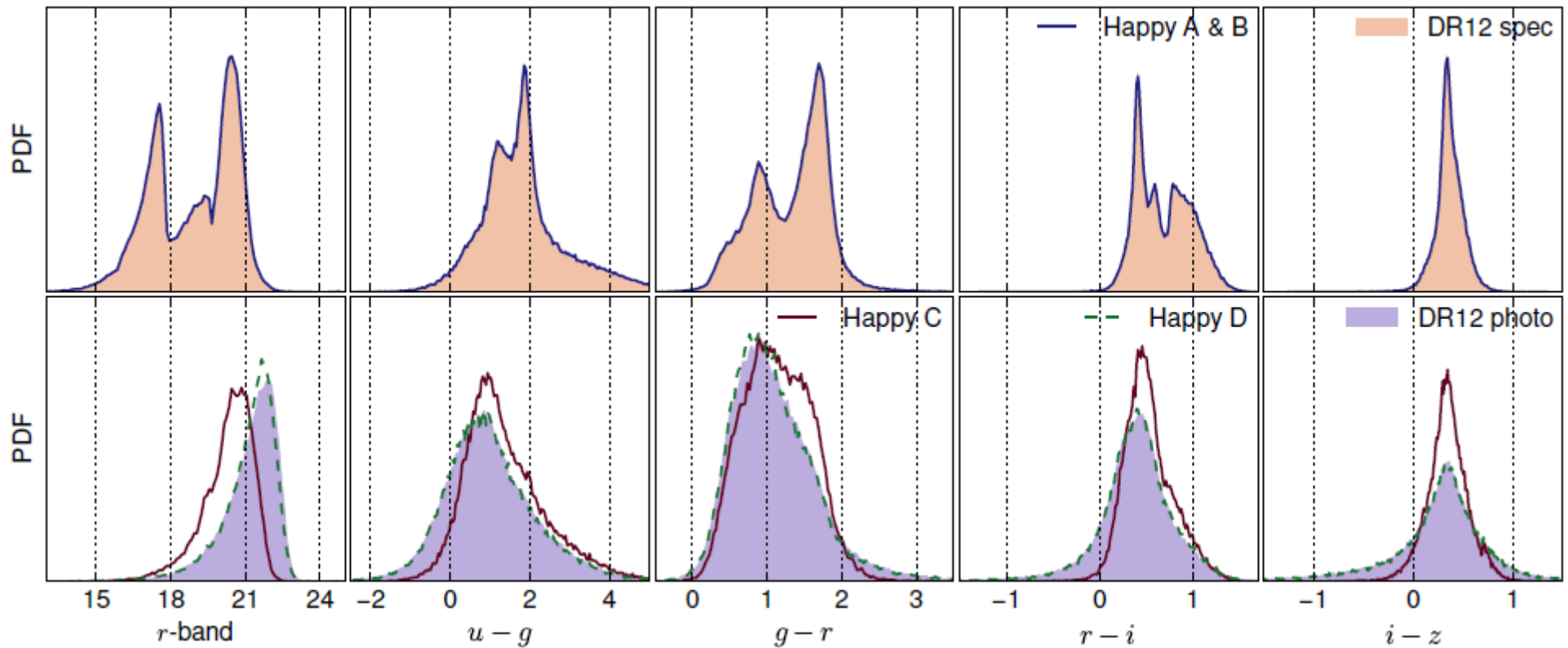
Happy catalogue

The effect of coverage + photometric errors

Photometry from SDSS

Spec-z from many different surveys leads to larger photometric errors and consequently wide domain in r-band and color

- A / B follow SDSS spec distribution
- B is completely representative of A
- C was constructed performing a nearest neighbor between the SDSS-DR12 photo sample and the extended spec sample but with a cut on photometric errors
- D is the same of C but without the photometric error cut.
- Consequently, D follows exactly the SDSS-DR12 photo sample distribution



Happy catalogue

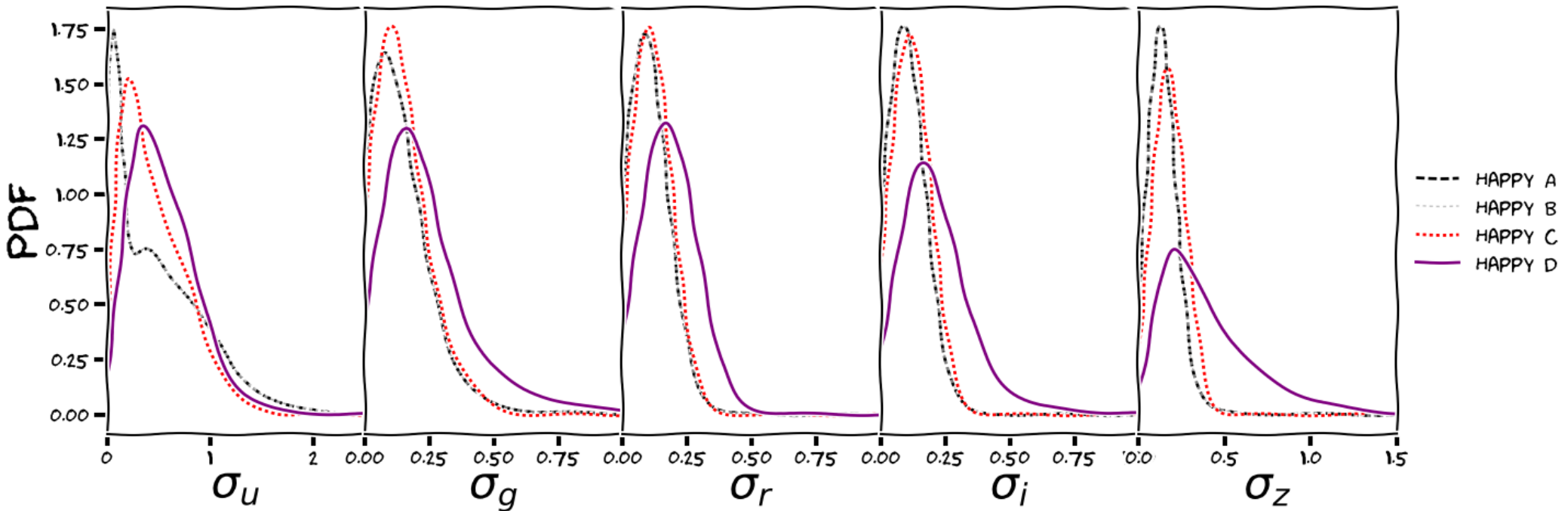
The effect of coverage + photometric errors

Photometry from SDSS

Spec-z from many different surveys leads to larger photometric errors and consequently wide domain in r-band and color

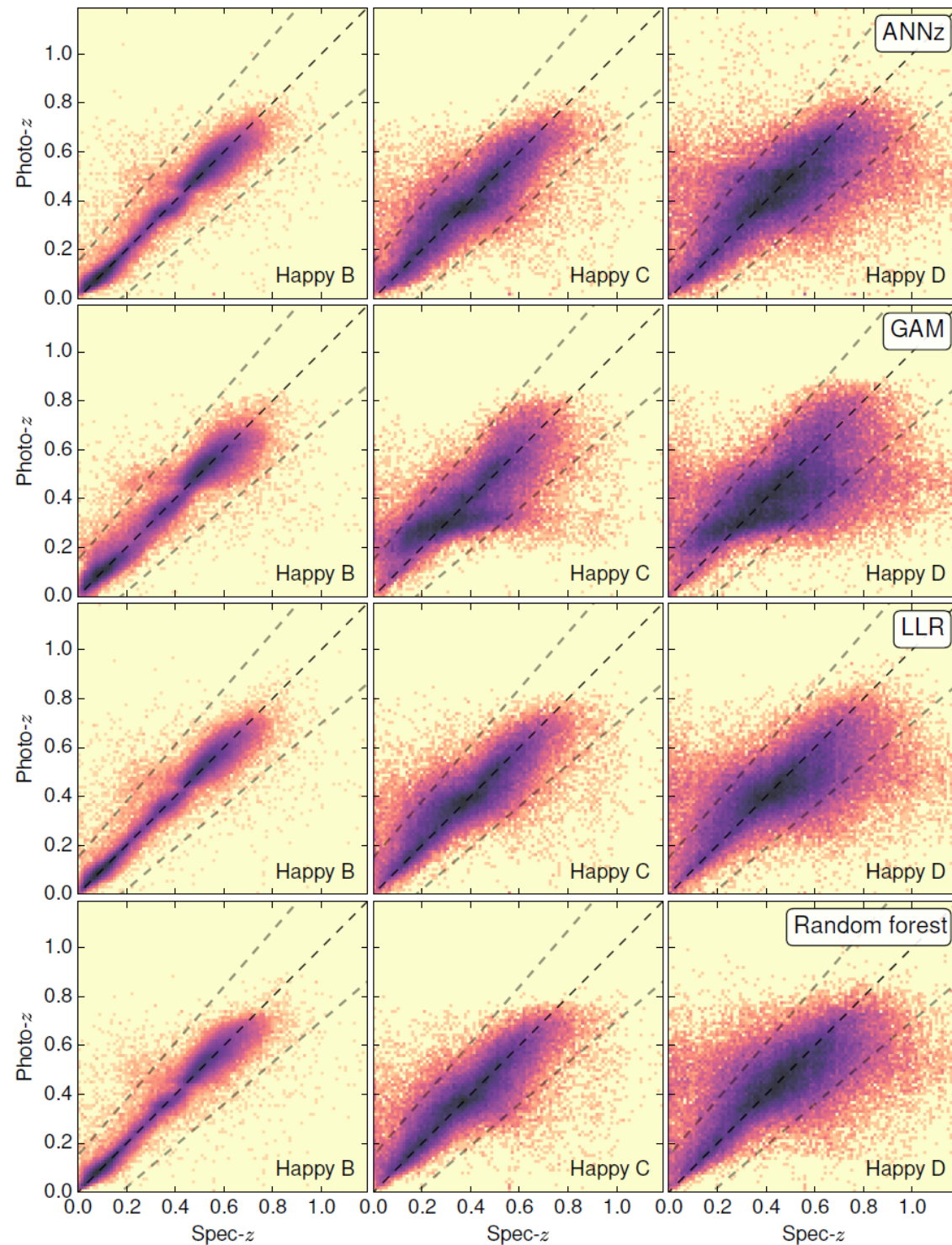
- A /B follow SDSS spec distribution
- B is completely representative of A
- C was constructed performing a nearest neighbor between the SDSS-DR12 photo sample and the extended spec sample but with a cut on photometric errors
- D is the same of C but without the photometric error cut.
- Consequently, D follows exactly the SDSS-DR12 photo sample distribution

Error distributions correlate with features



Happy catalogue

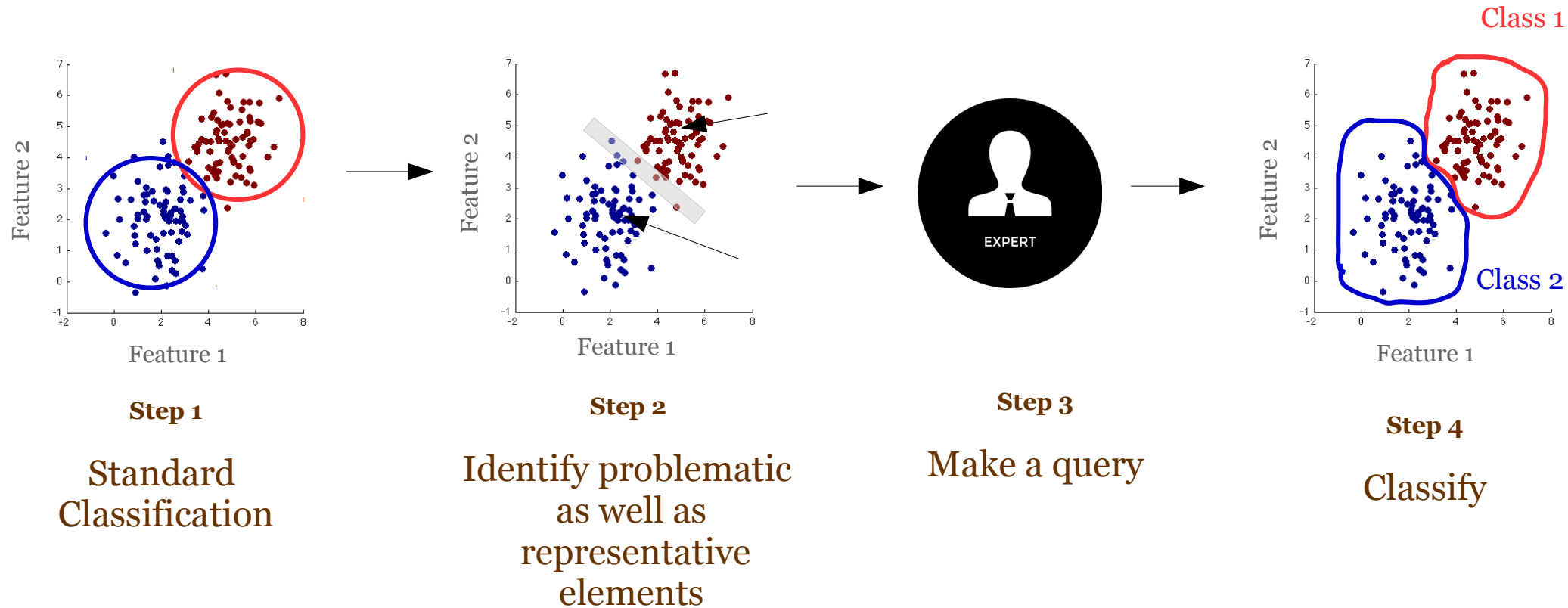
The effect of coverage + photometric errors



Method	Set	Diagnostics			Outlier rate (%)
		Mean ($\times 10^{-2}$)	Std ($\times 10^{-2}$)	MAD ($\times 10^{-2}$)	
ANNz	B	0.04	2.87	1.49	0.99
	C	0.16	5.41	3.60	5.59
	D	-0.52	6.53	5.44	14.01
GAM	B	0.09	3.50	1.95	1.36
	C	0.86	6.34	4.84	7.37
	D	-0.51	7.21	6.70	16.38
LLR	B	0.13	2.81	1.39	1.11
	C	0.52	5.45	3.59	6.07
	D	-0.79	6.62	5.62	14.52
Random Forest	B	0.05	2.82	1.41	1.02
	C	0.34	5.39	3.51	5.58
	D	-0.28	6.51	5.36	14.2

Potential solution: Active Learning

SAMSI & COIN, in prep



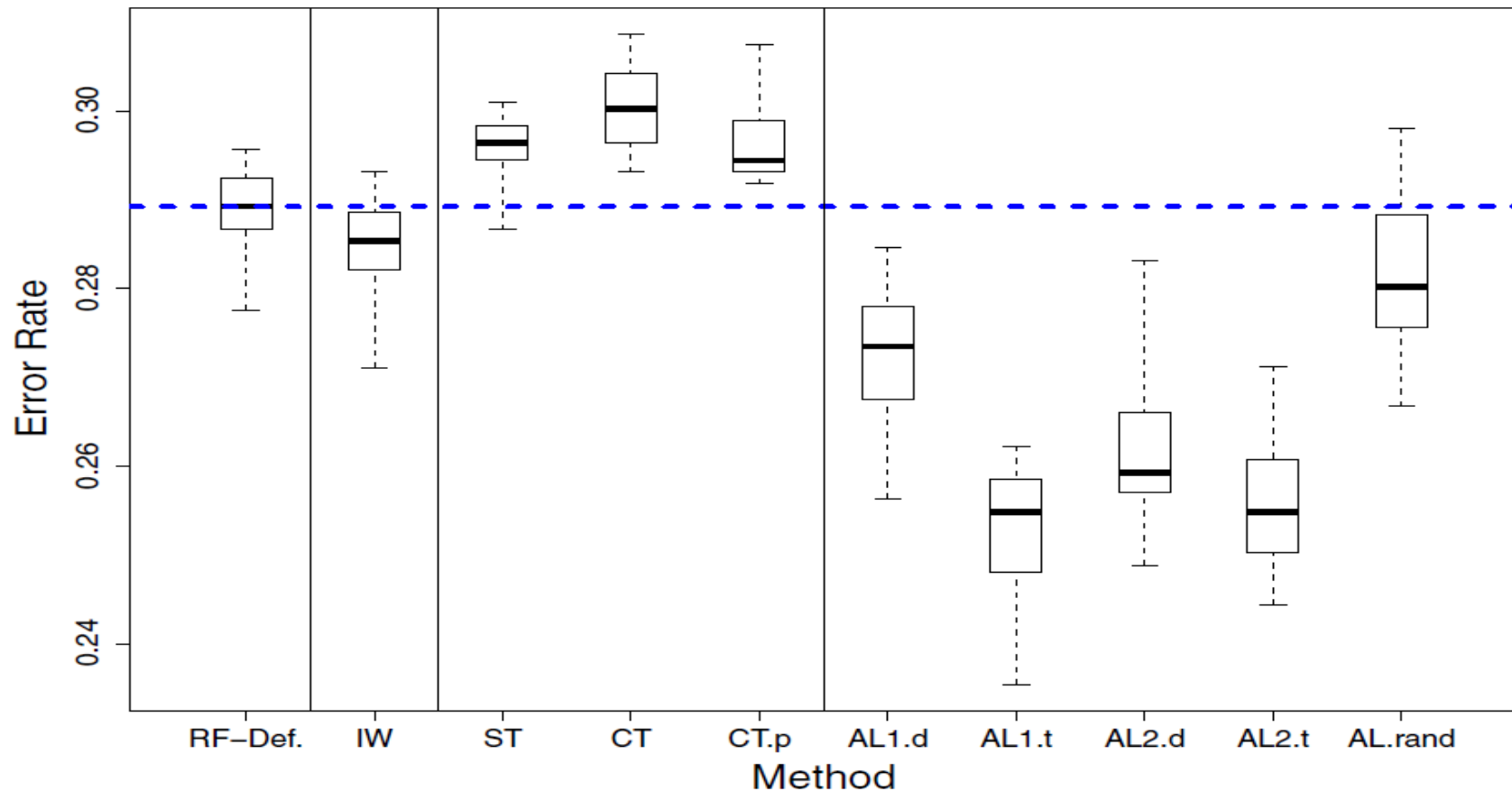
Background: Active Learning in Astronomy

ACTIVE LEARNING TO OVERCOME SAMPLE SELECTION BIAS: APPLICATION TO PHOTOMETRIC VARIABLE STAR CLASSIFICATION

JOSEPH W. RICHARDS^{1,2}, DAN L. STARR¹, HENRIK BRINK³, ADAM A. MILLER¹, JOSHUA S. BLOOM¹,
NATHANIEL R. BUTLER¹, J. BERIAN JAMES^{1,3}, JAMES P. LONG², AND JOHN RICE²

supervised classification

THE ASTROPHYSICAL JOURNAL, 744:192 (19pp), 2012 January 10



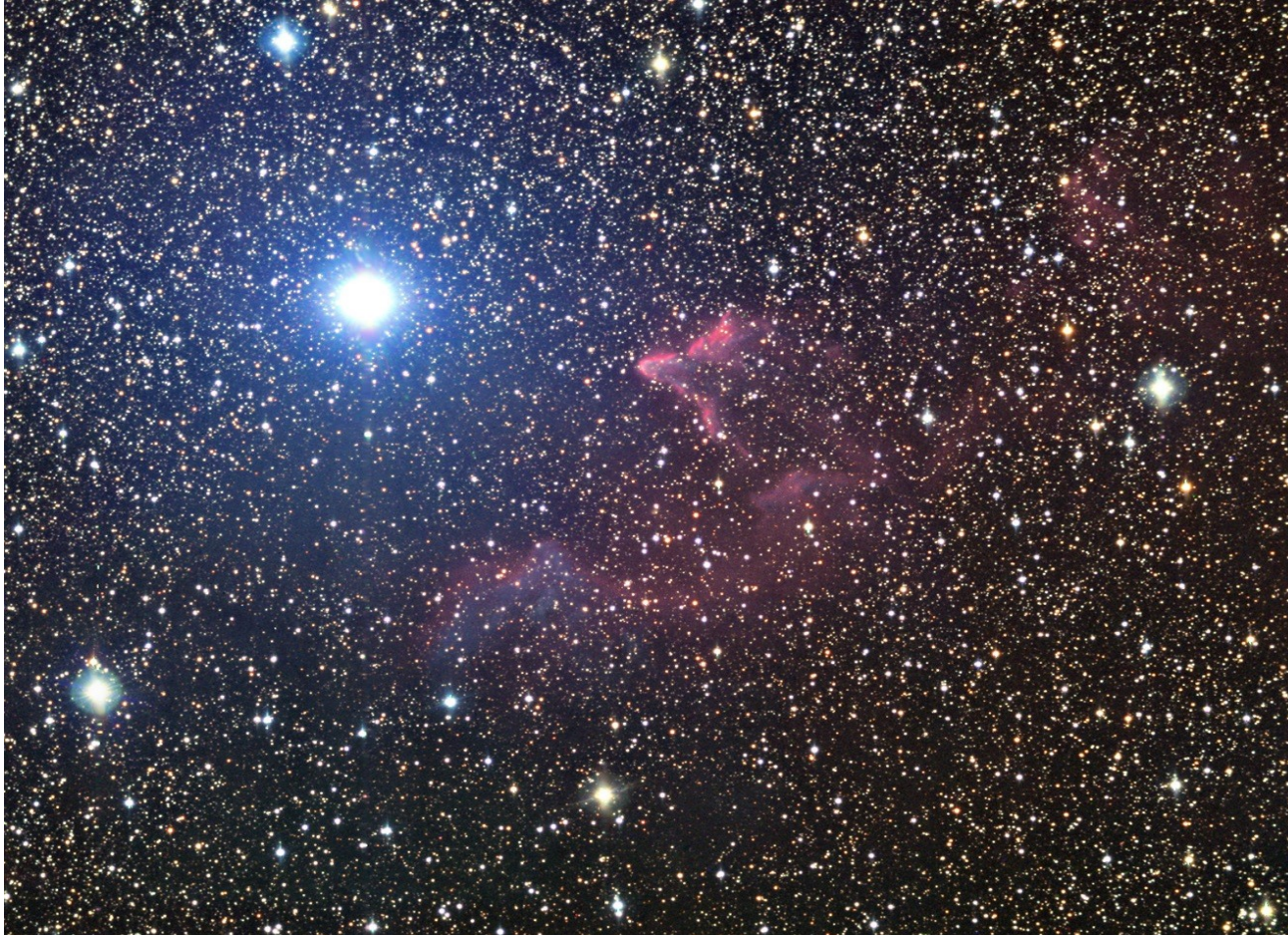
Background: Active Learning in Astronomy

Automated Supernova Ia Classification Using Adaptive Learning Techniques

Kinjal Dhar Gupta*, Renuka Pampana*, Ricardo Vilalta*, Emille E. O. Ishida[†], Rafael S. de Souza[‡]

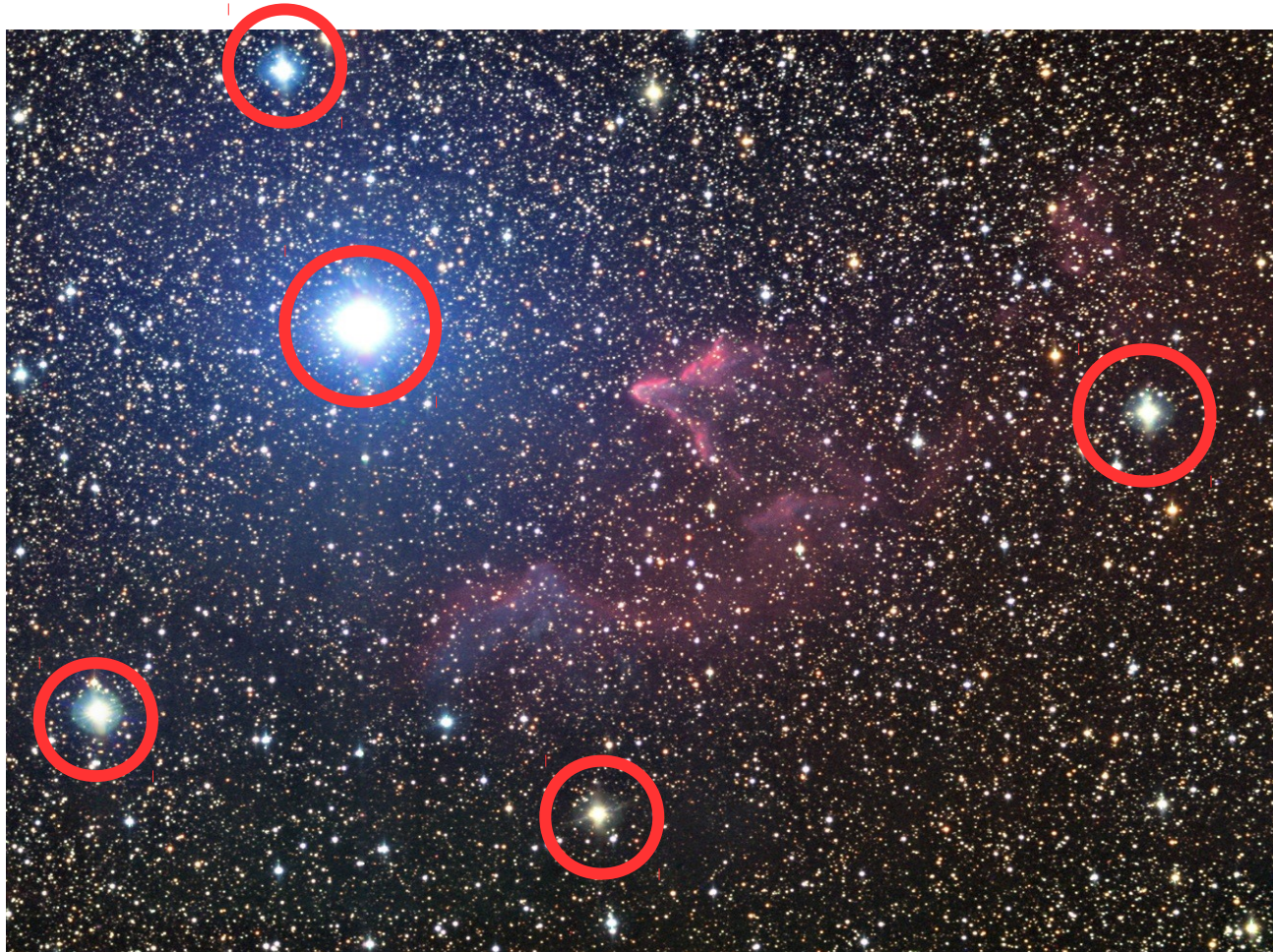
supervised classification

How are spectroscopic sets constructed?



How are spectroscopic sets constructed?

Take spectra for learning and determine everything else



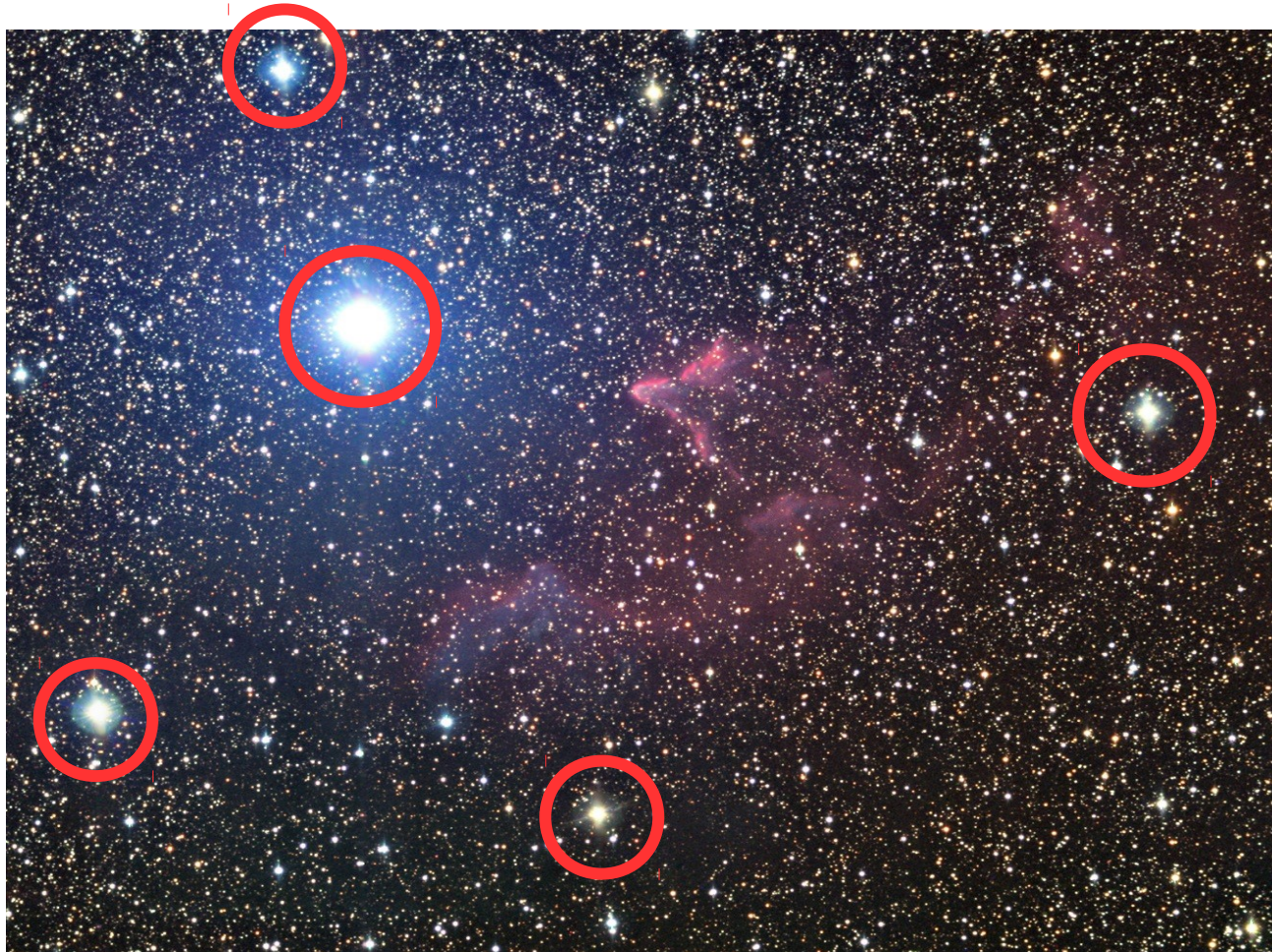
Alternative approach

Landmark selection



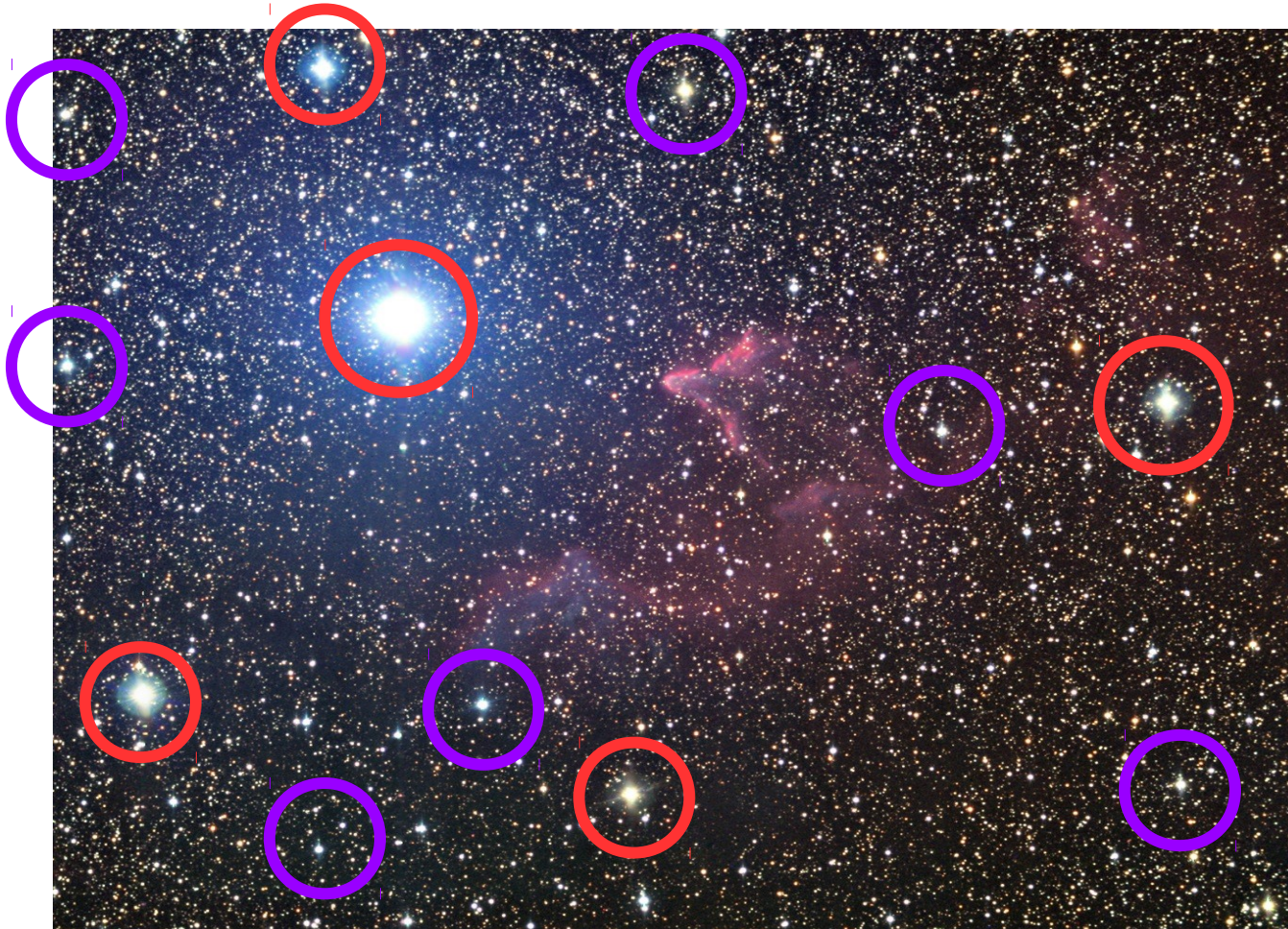
Alternative approach

Landmark selection



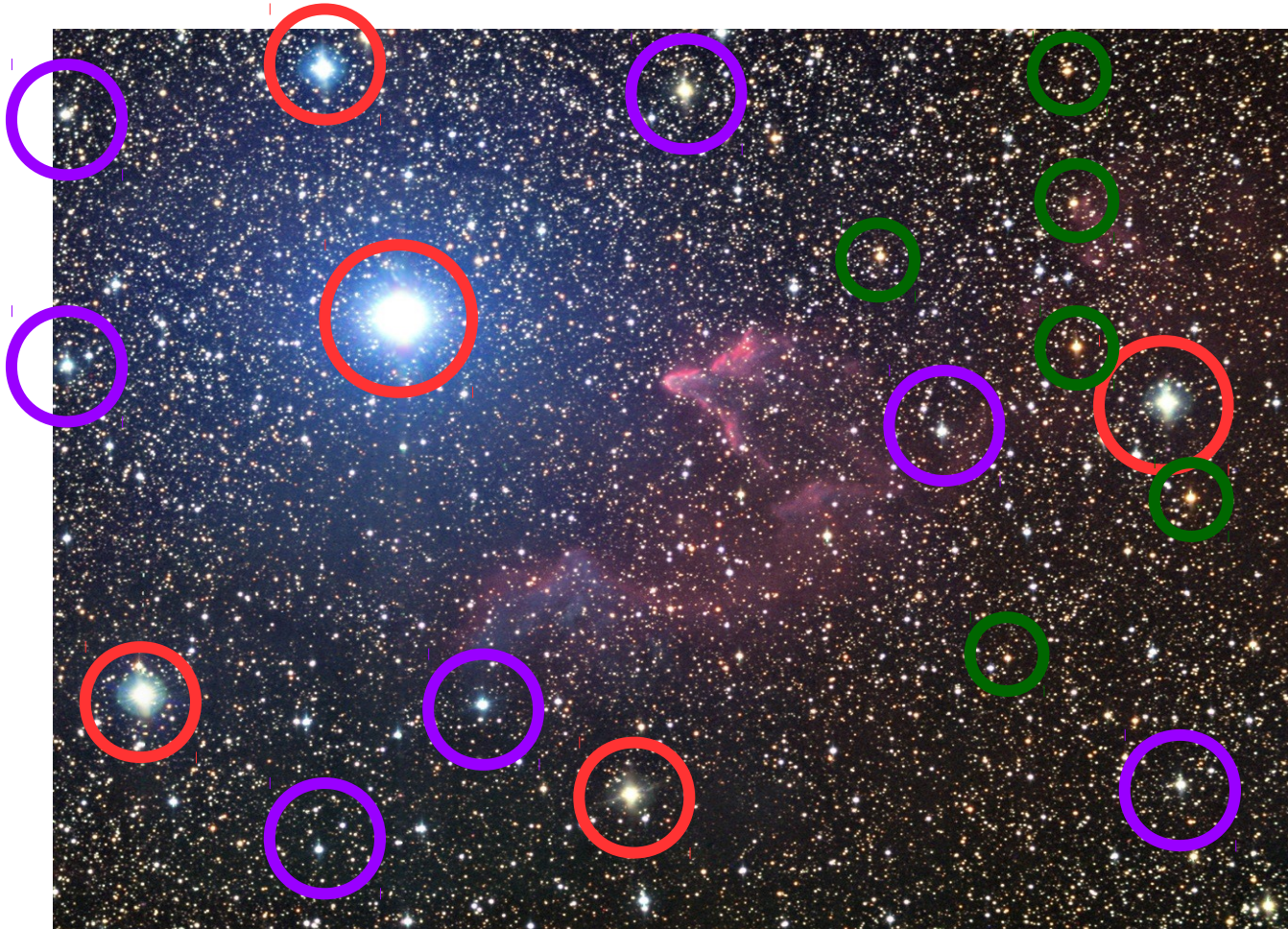
Alternative approach

Landmark selection



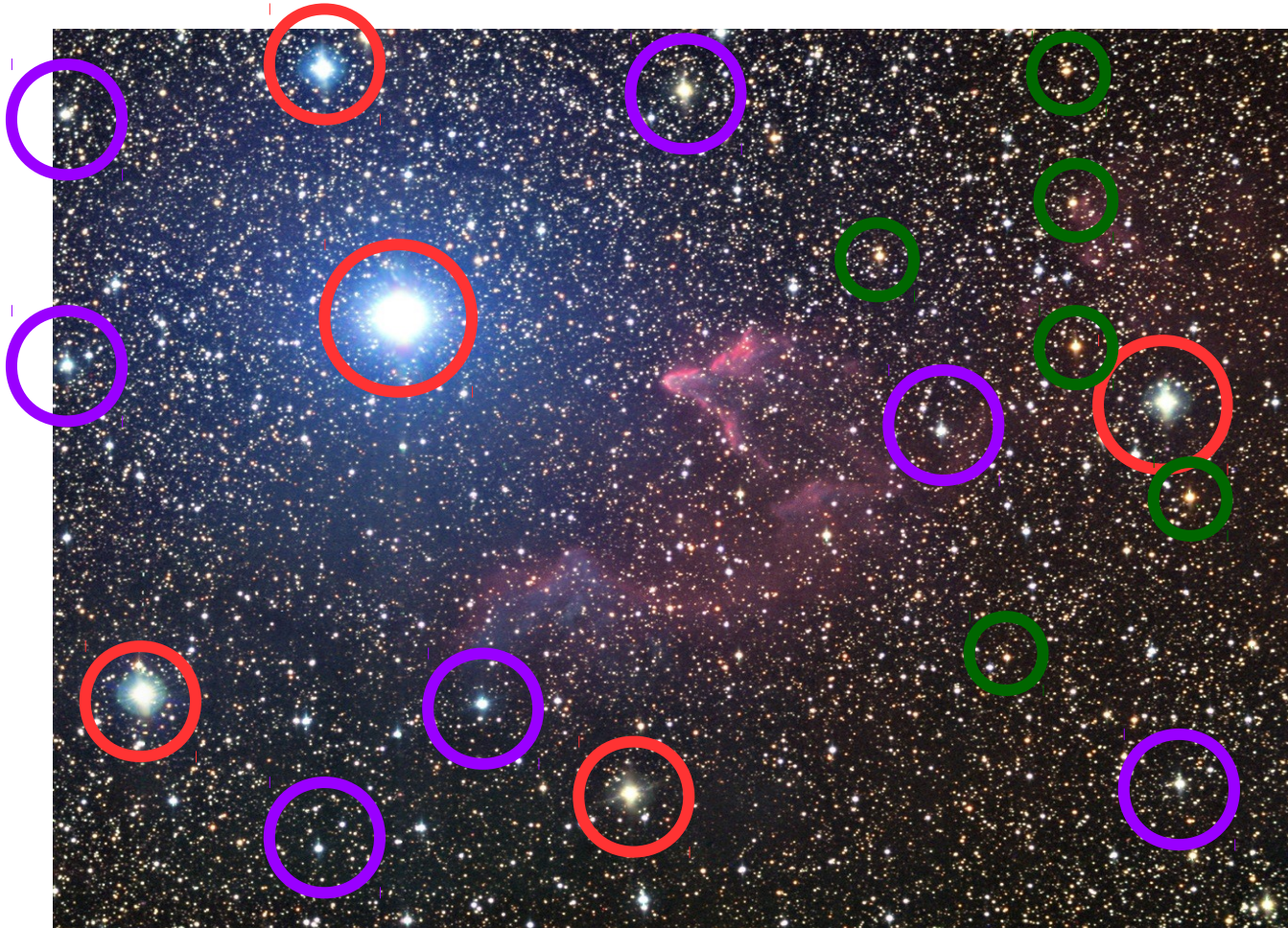
Alternative approach

Landmark selection



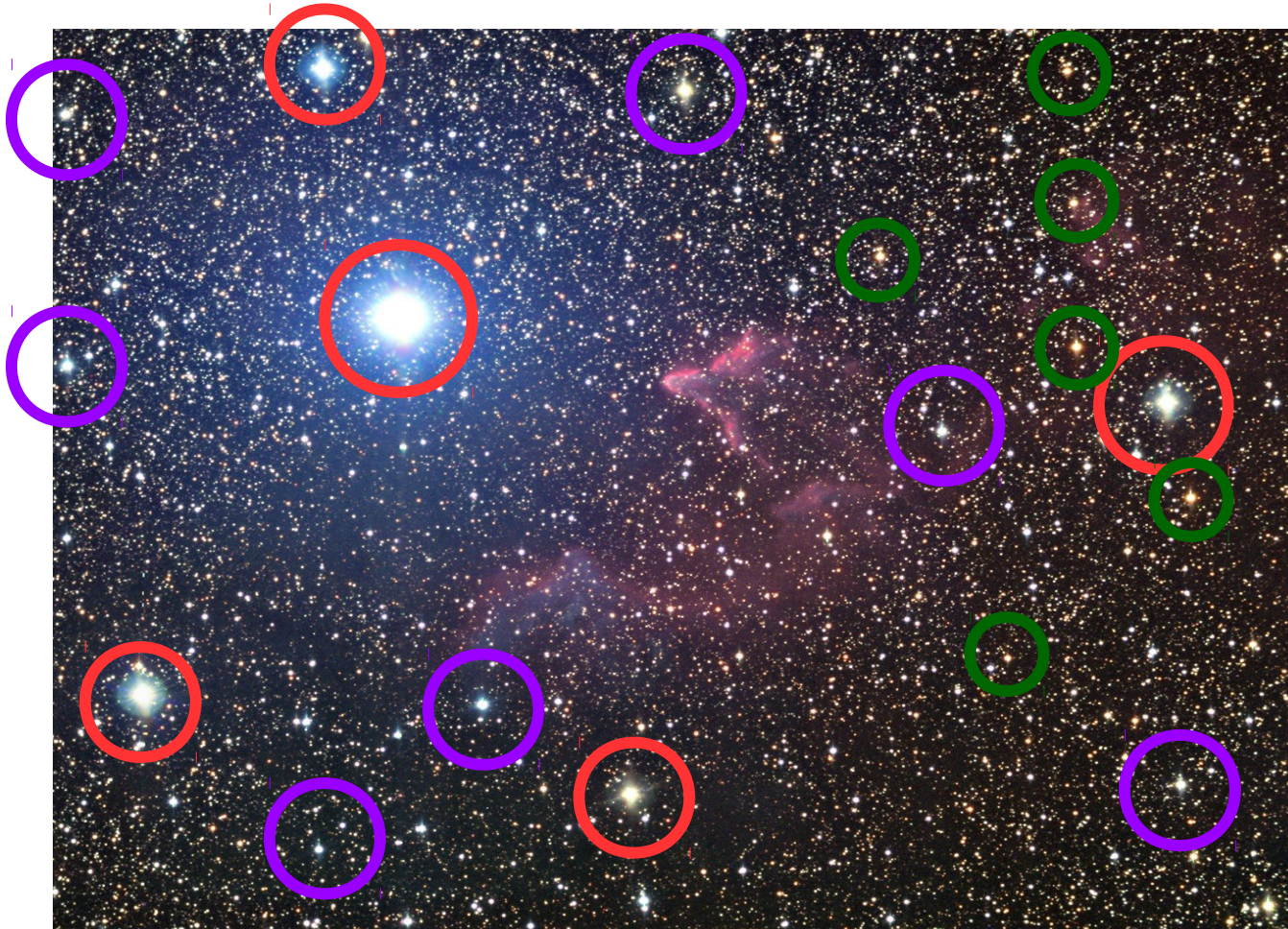
Alternative approach

Landmark selection + Active Learning



Alternative approach

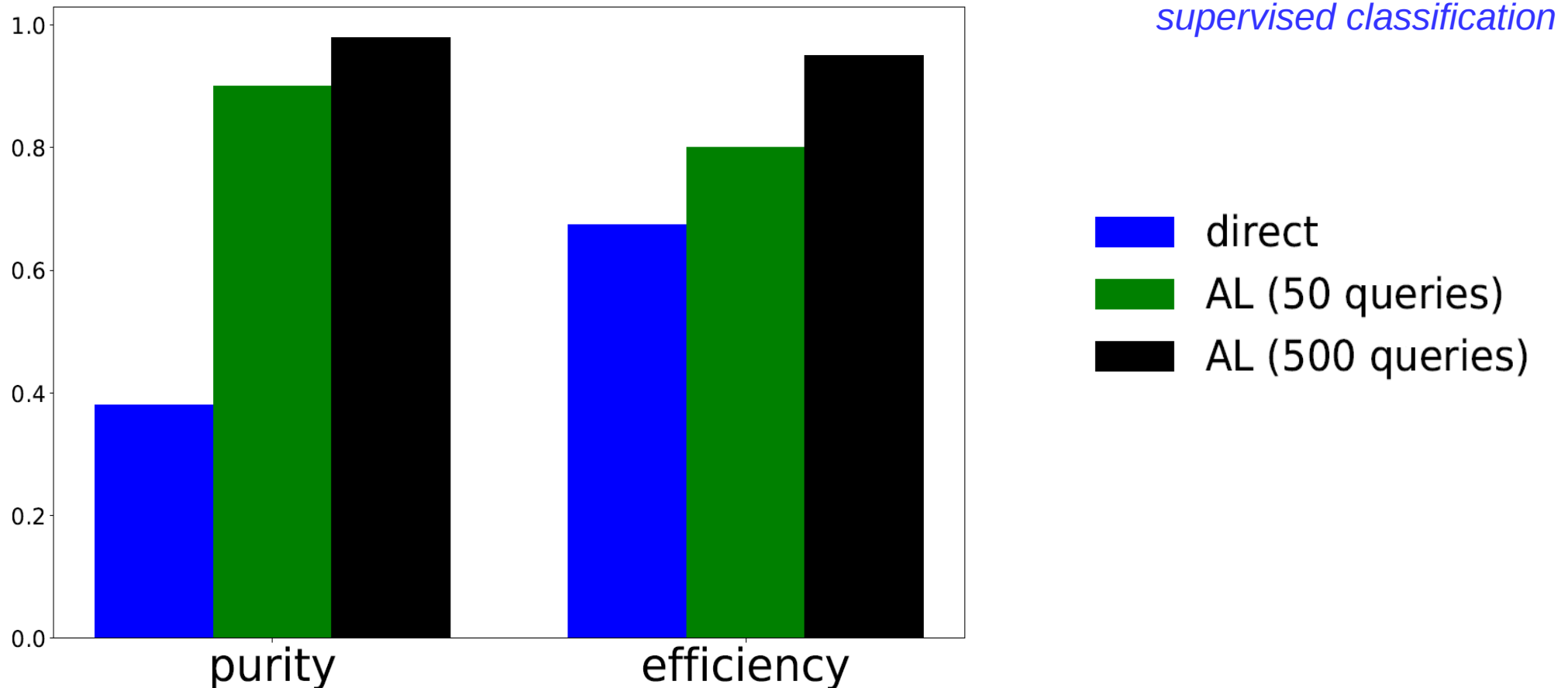
Landmark selection + Active Learning




Background: Active Learning in Astronomy

Automated Supernova Ia Classification Using Adaptive Learning Techniques

Kinjal Dhar Gupta*, Renuka Pampana*, Ricardo Vilalta*, Emille E. O. Ishida[†], Rafael S. de Souza[‡]





Active Learning
for supervised
regression?

Active Learning
for supervised
regression?

**TO BE
CONTINUED...**

Main tasks:

Supervised Learning



Regression:

- photometric redshift
- stellar parameters determination
- ...

Classification:

- detection
- star/galaxy separation
- galaxy morphology
- variable stars
- supernova
- ...

Unsupervised Learning



Clustering:

- SN Ia spectra characterization
- galaxy spectral classification
- ...

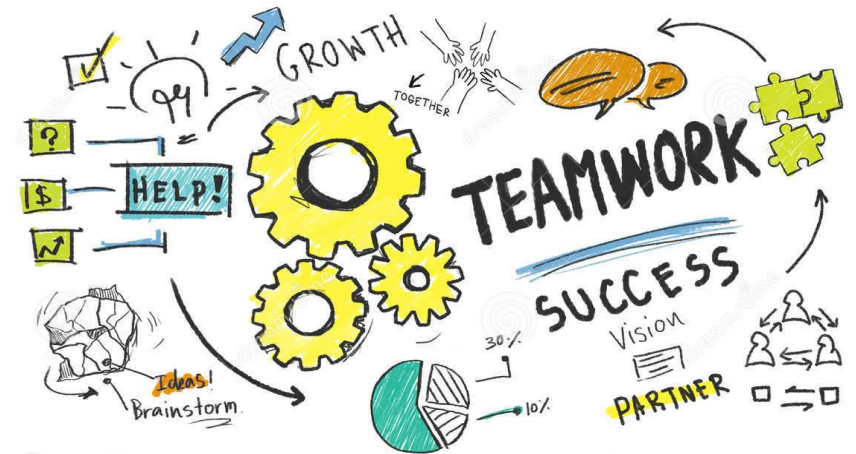
Anomaly/Novelty detection:

- unforeseen new objects
- detection error analysis
- identification of predicted objects
- ...



What about the
future?

urgent: Build a support community



urgent: Build a support community



solve bureaucracy



