



# Astronomy in the 21st century: Drowning in data, Starving for knowledge

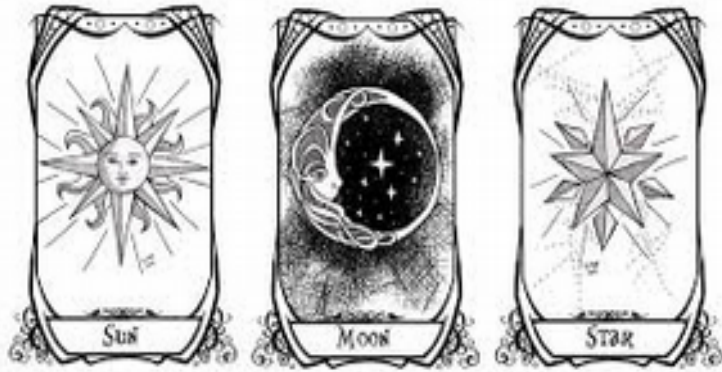
Emille E. O. Ishida

*Laboratoire de Physique Corpusculaire - Université Blaise Pascal  
Clermont Ferrand, France*

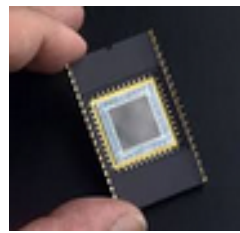


# The Big Picture

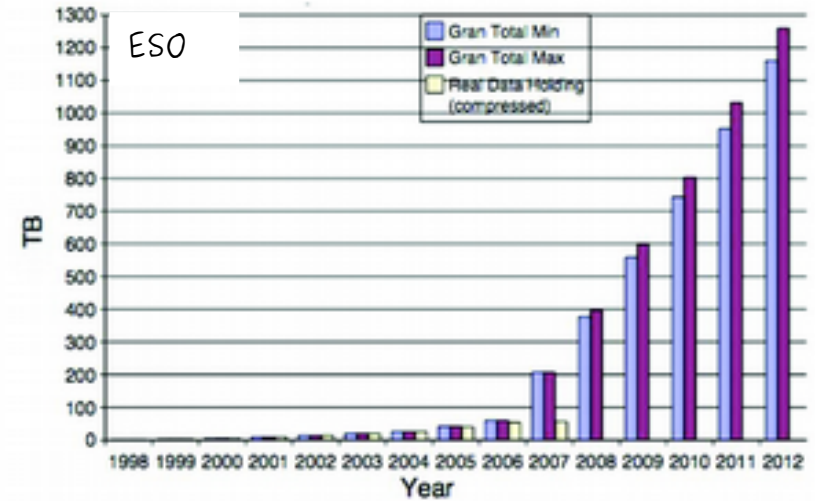
Astronomy began with 3 elements



In 1969, the CCD



Data

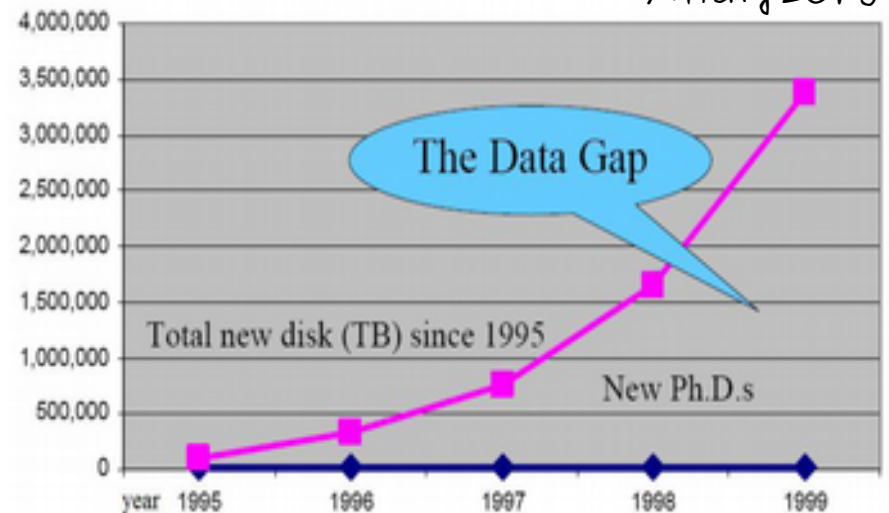


In the 17<sup>th</sup> century...

- 1 telescope
- 1 observer
- 1 object



Analizers



Grossman (2001)

An illustrative example:

# SDSS – Sloan Digital Sky Survey

**1992**

2.5 Terapixels of images

10 TB of raw data

0.5 TB catalogs

1992 – 2000: planning

2001 – 2009: observing

2.5m mirror

New Mexico, USA



An illustrative example:

# SDSS – Sloan Digital Sky Survey

**1992**

2.5 Terapixels of images

10 TB of raw data

0.5 TB catalogs

1992 – 2000: planning

2001 – 2009: observing

2.5m mirror

New Mexico, USA



How to deliver 0.5 TB of useful data to all users?

An illustrative example:

# SDSS – Sloan Digital Sky Survey

## 1992

2.5 Terapixels of images  
10 TB of raw data  
0.5 TB catalogs

1992 – 2000: planning  
2001 – 2009: observing

2.5m mirror  
New Mexico, USA



## 2009

5 Tpx of images  
120TB processed data  
35TB catalogs



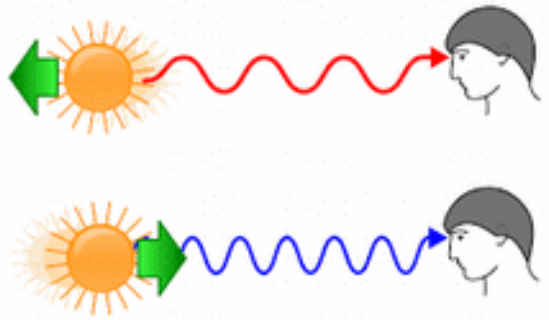
by Ann K. Finkbeiner, 2012

Case study:

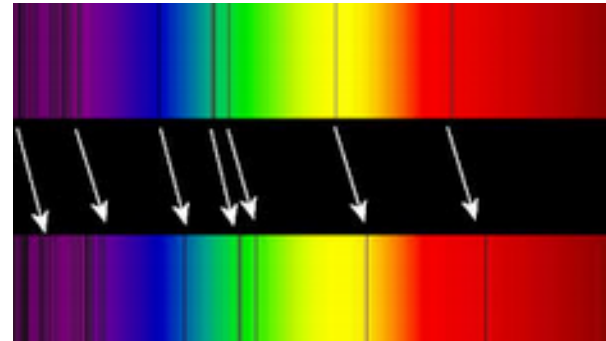
# Photometric Redshifts

# Redshifts (z) ↔ Distances

Idea



Observation



or

Velocities

+

Cosmological  
model

=

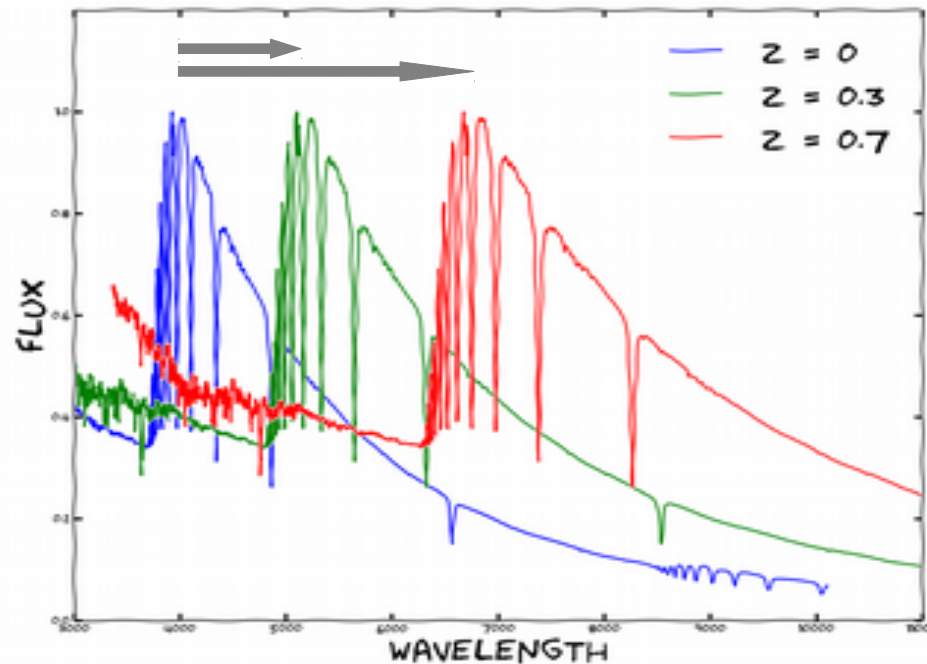
distances

+

finite light velocity



History



Redshifts  
are  
important!

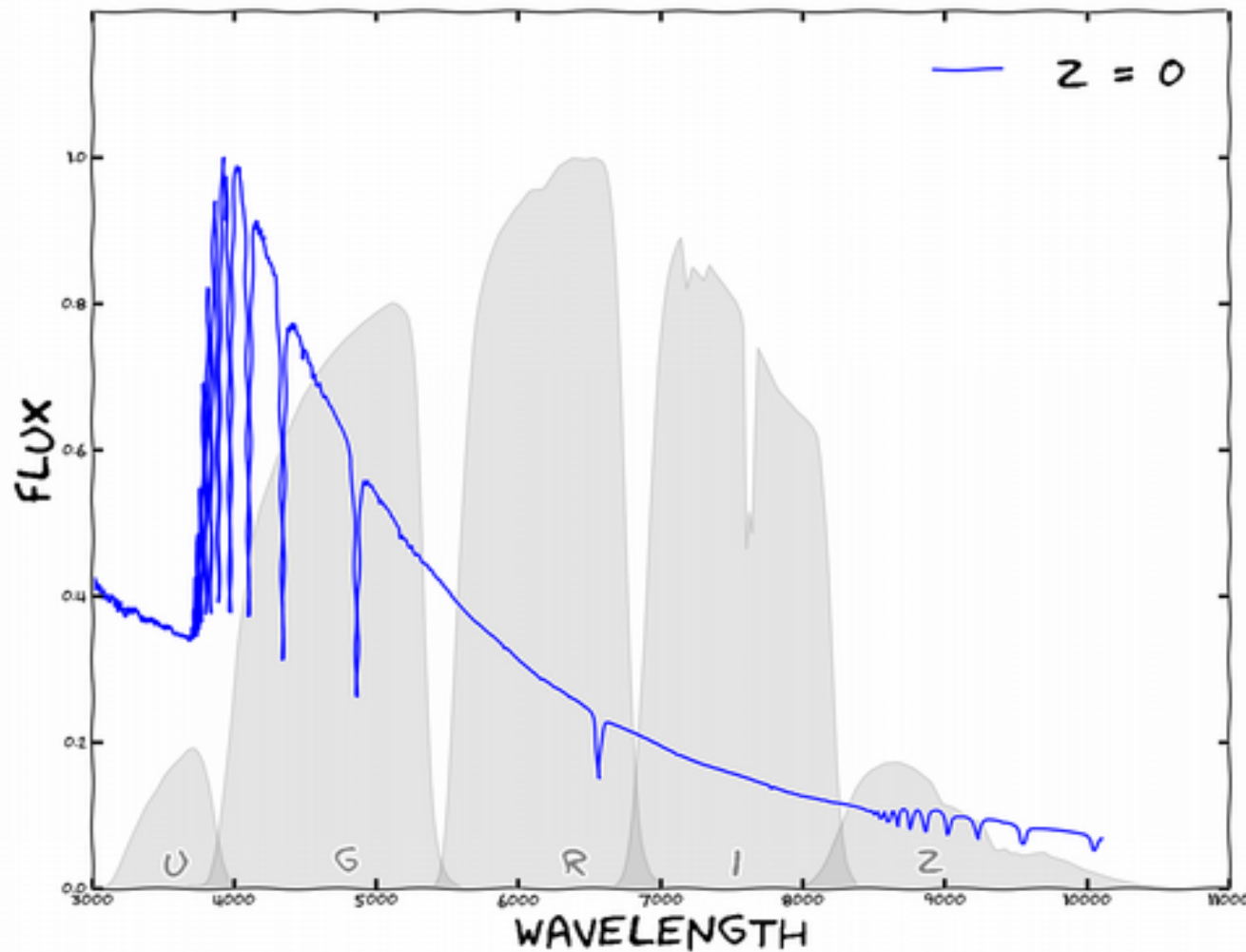


# Photometric Redshifts

Spectra are expensive!

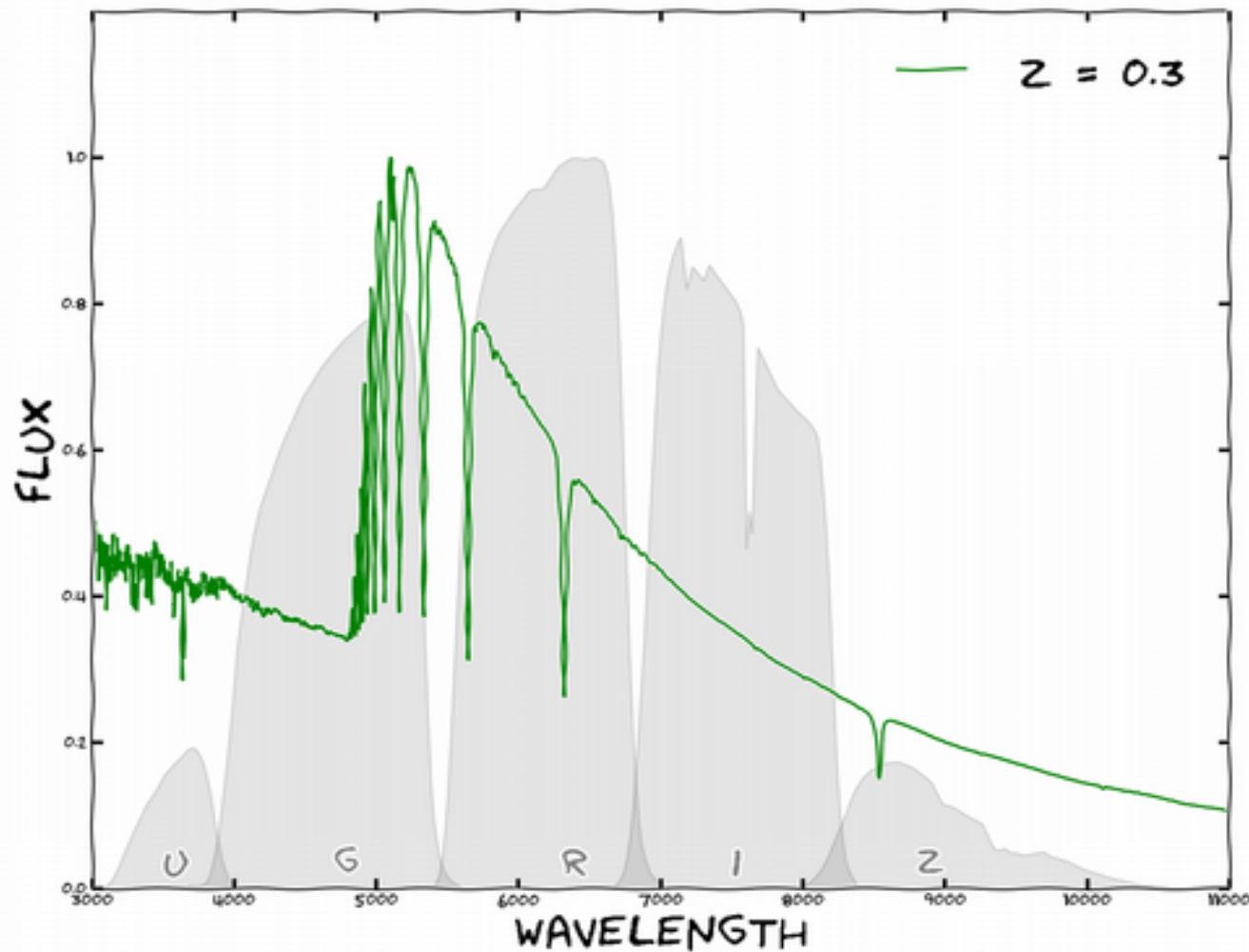
# Photometric Redshifts

Alternative measurement  
300 dim  $\rightarrow$  5 dim



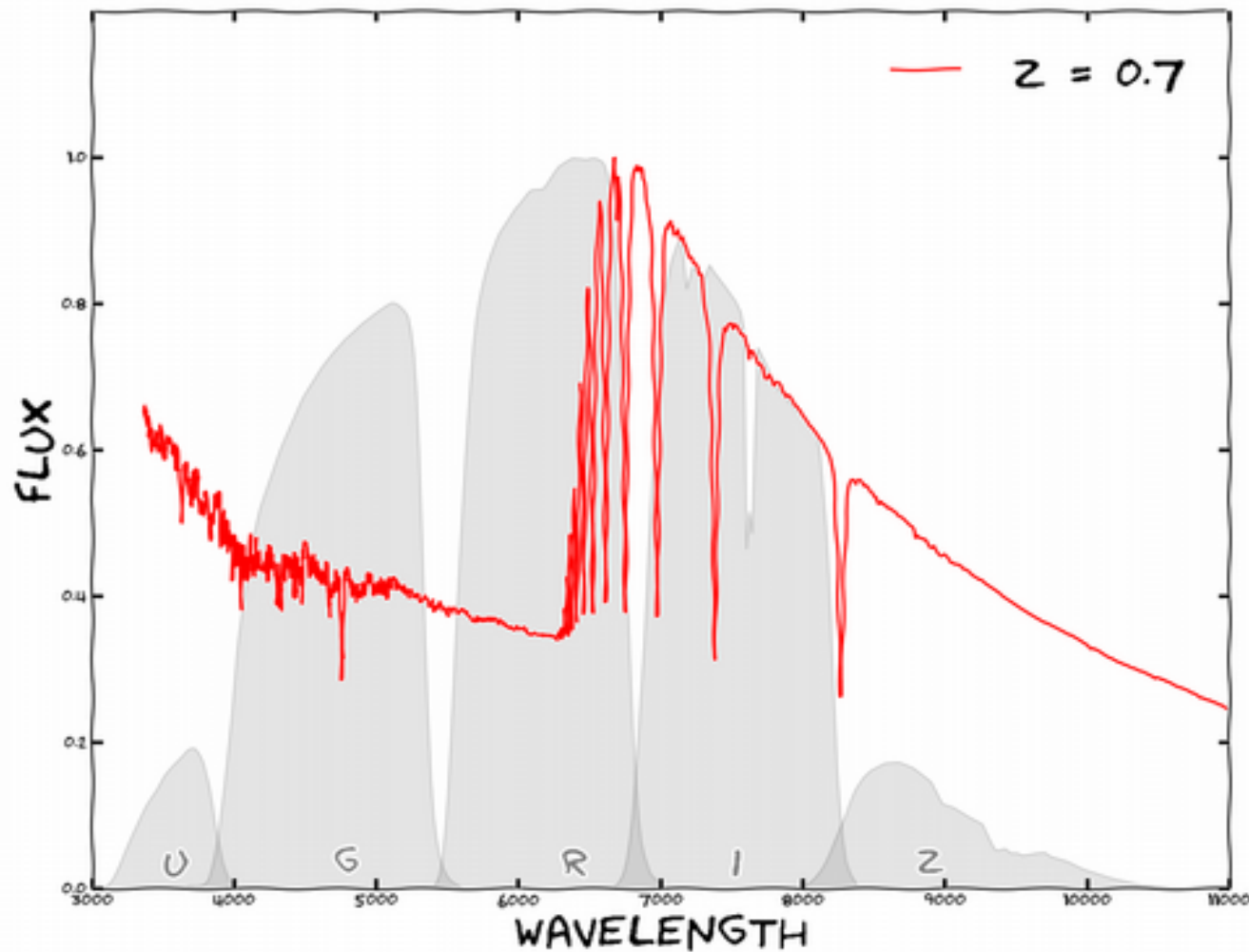
# Photometric Redshifts

Alternative measurement  
300 dim  $\rightarrow$  5 dim



# Photometric Redshifts

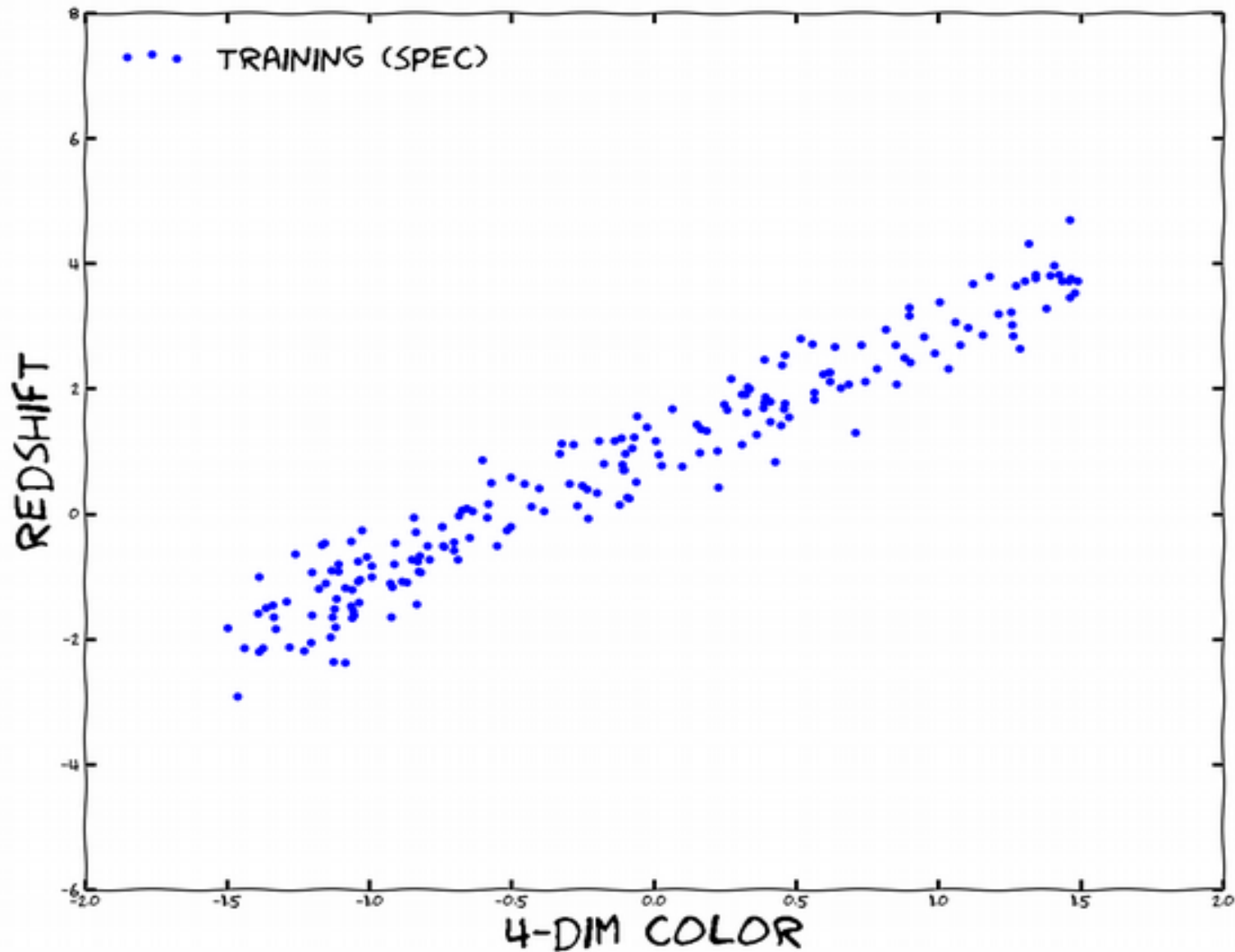
Alternative measurement  
300 dim  $\rightarrow$  5 dim



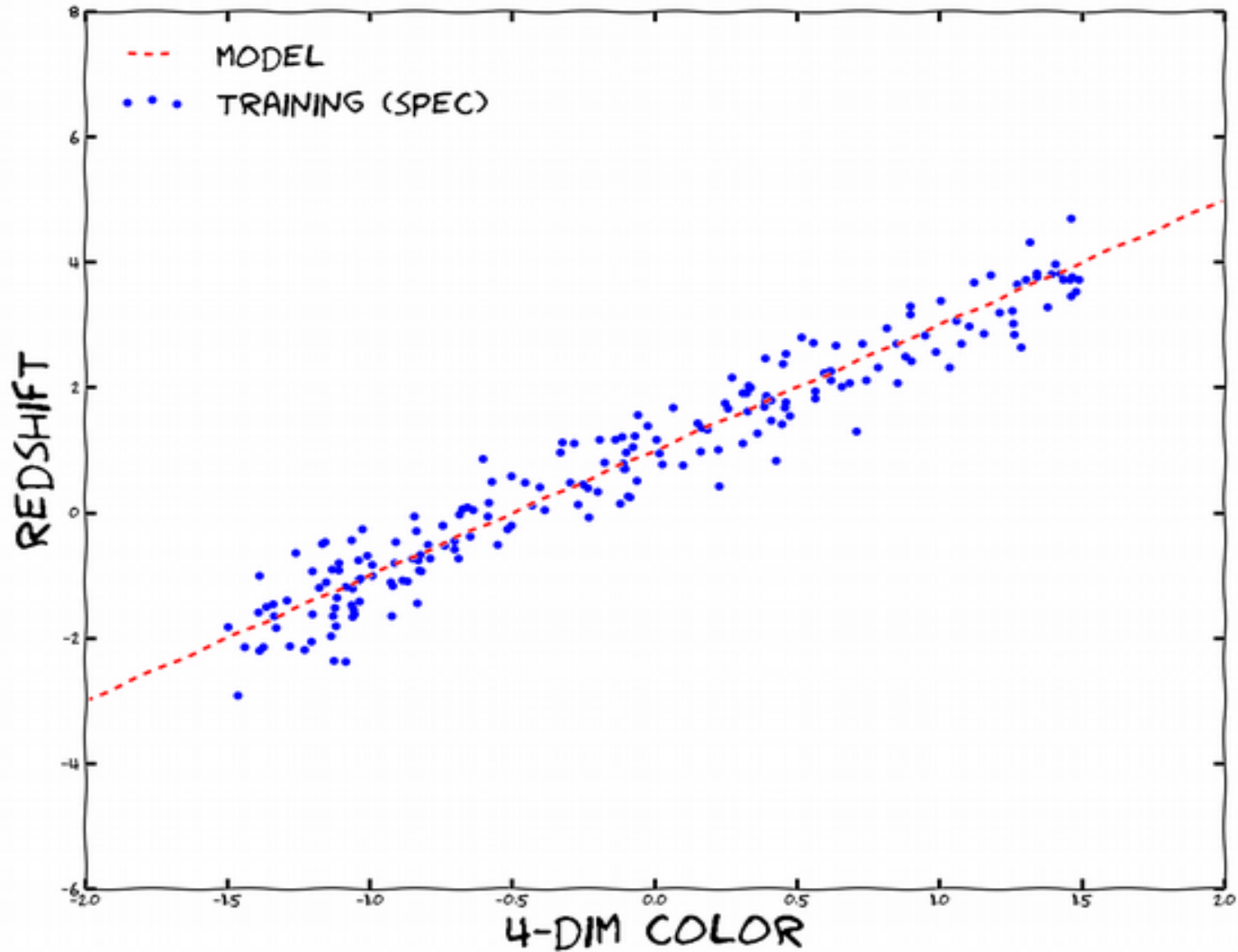
# Photo-z: a regression problem

u-g	g-r	r-i	i-z	redshift
2.07	1.39	0.48	0.27	0.31
1.54	1.58	0.54	0.42	0.34
1.03	1.76	0.67	0.37	0.41
2.17	1.30	0.43	0.30	0.19
...	...	...	...	...
1.36	1.72	0.52	0.36	0.32

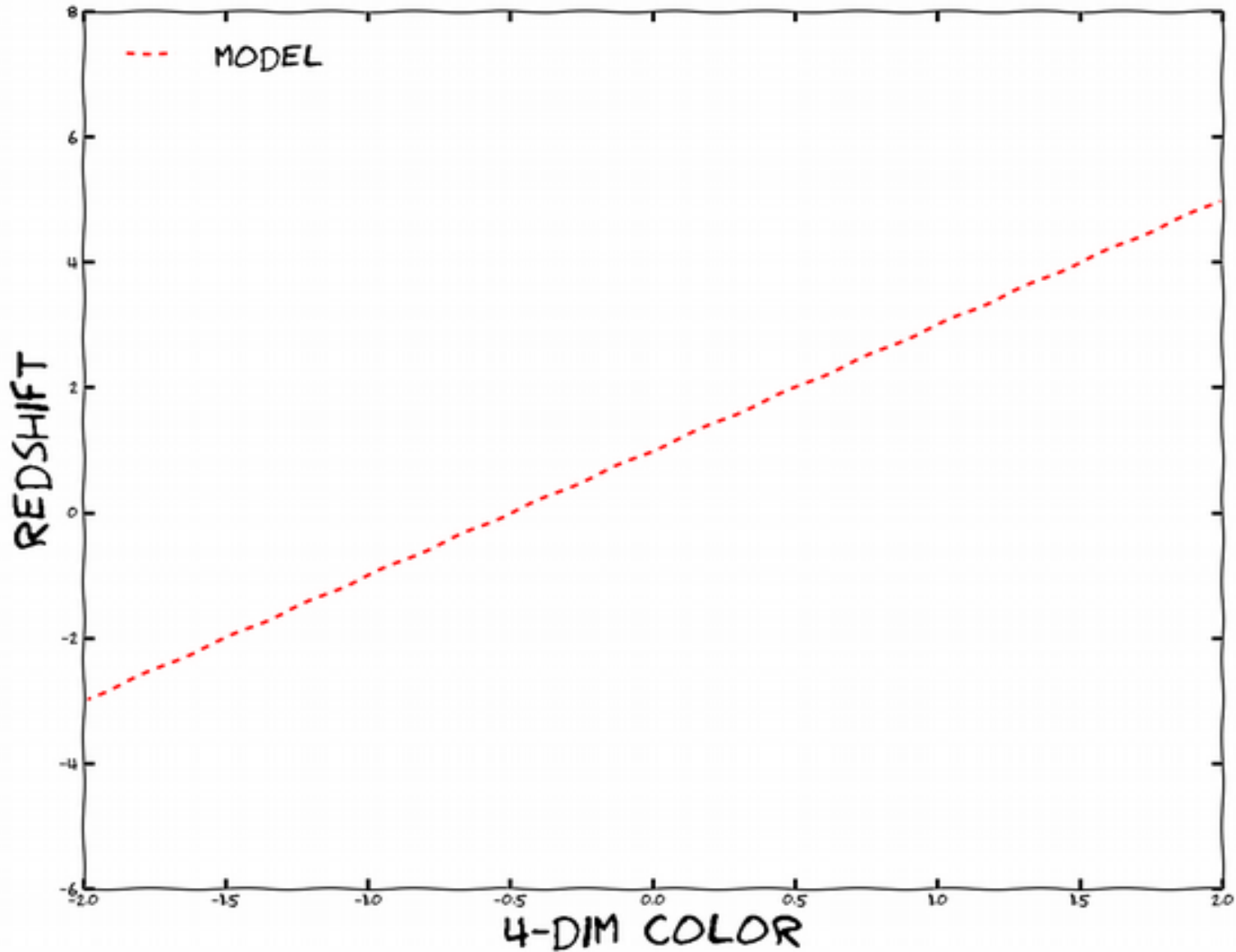
# Photo-z: a regression problem



# Photo-z: a regression problem

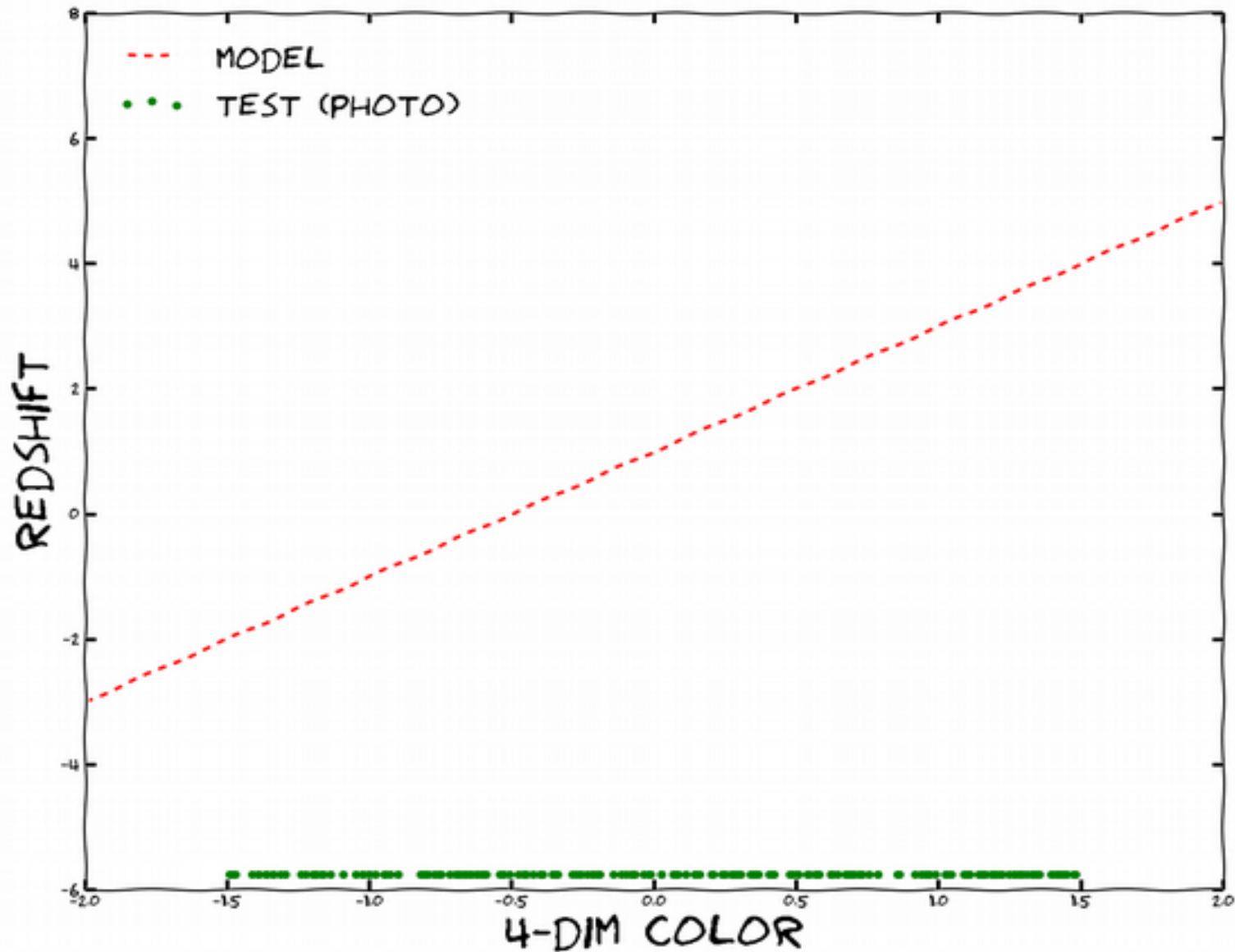


# Photo-z: a regression problem

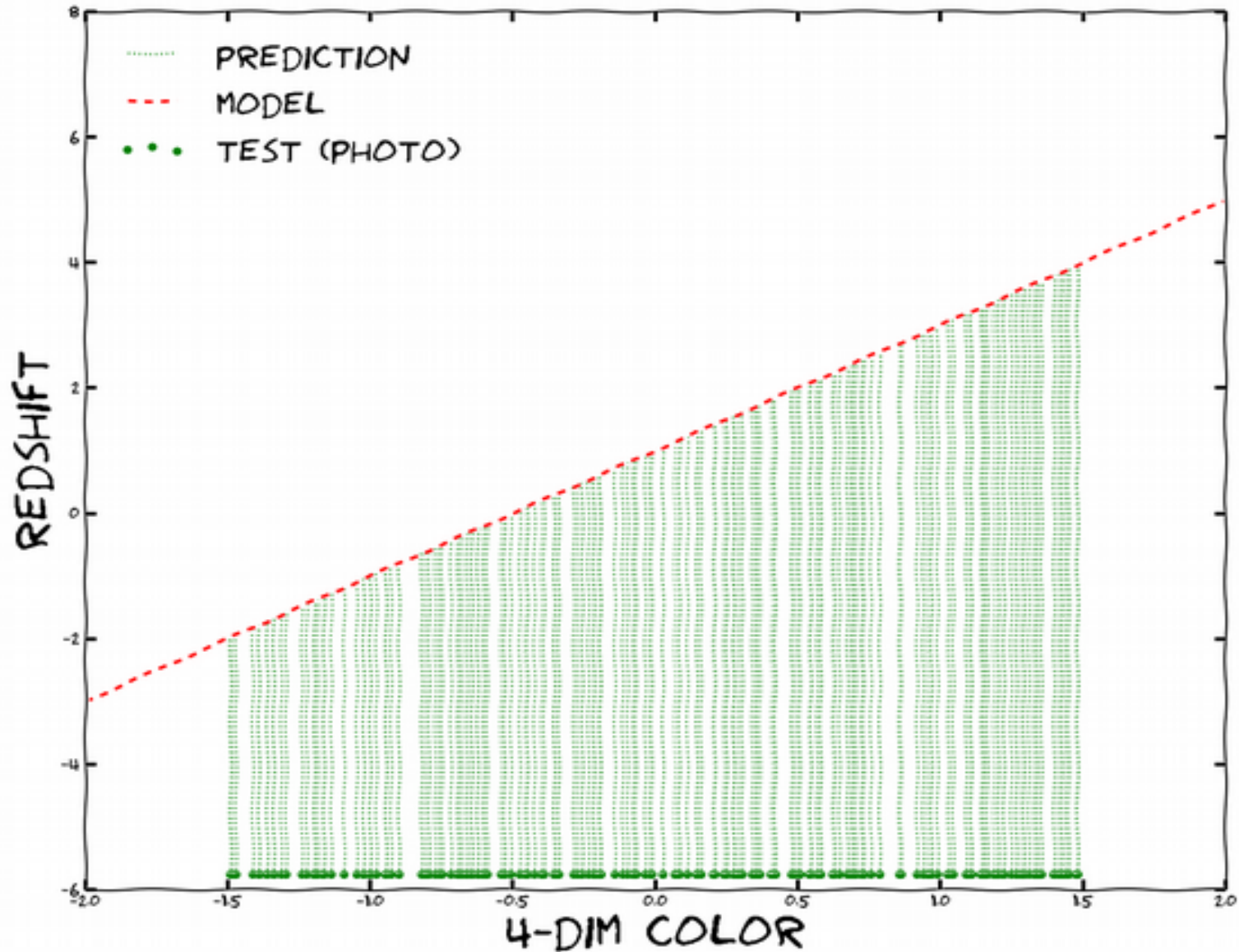




# Photo-z: a regression problem



# Photo-z: a regression problem

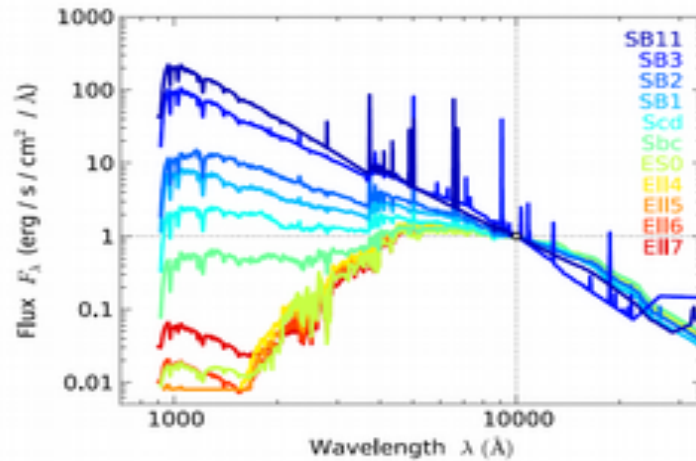
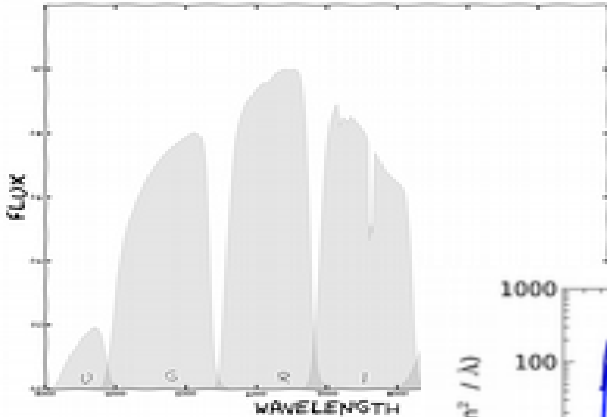


# Photo-z: methods

Hybrid



Template fitting



Machine Learning



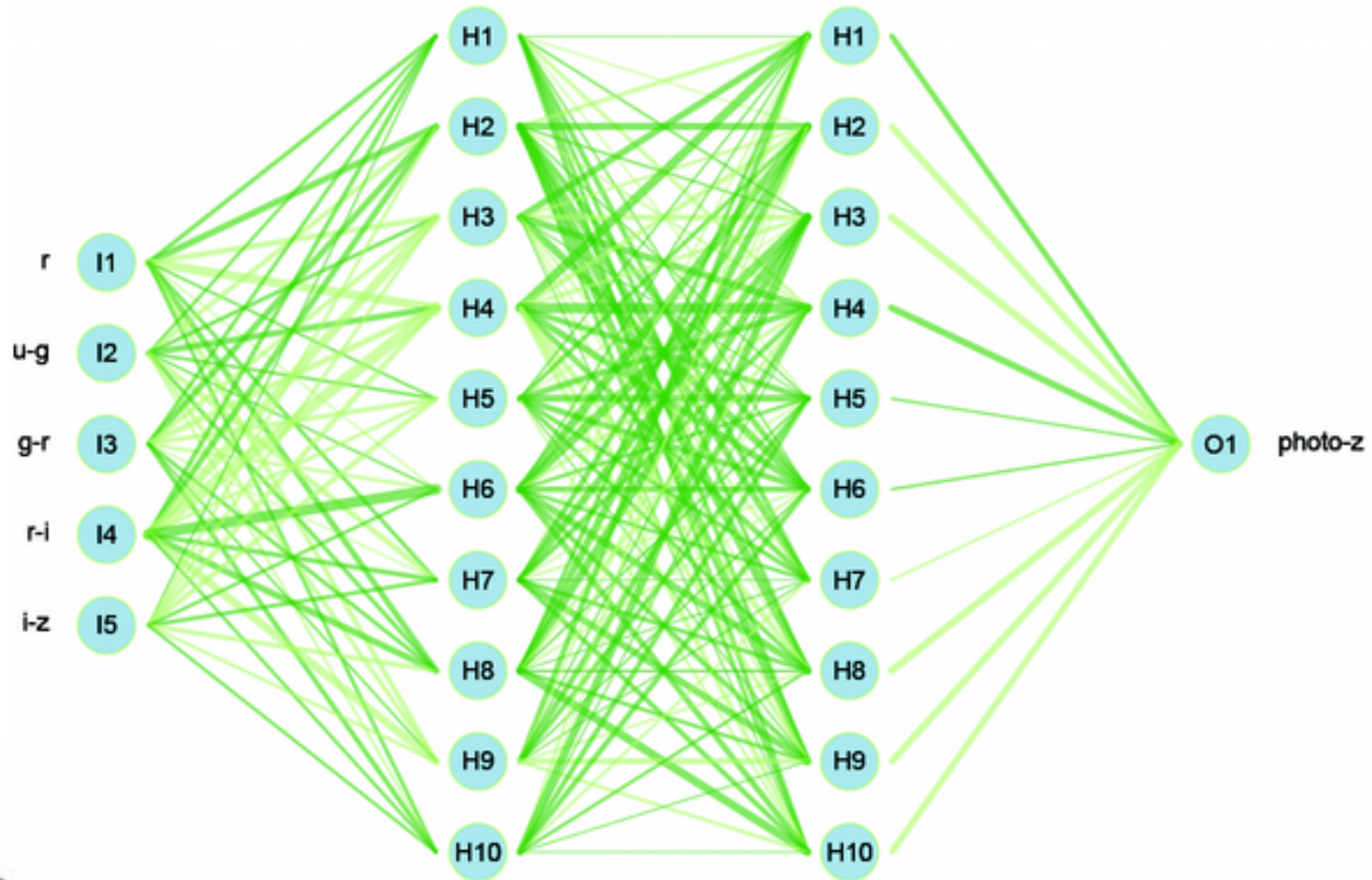
# Photo-z: methods

Machine  
Learning



# Photo-z: Artificial Neural Networks

Used in astronomy since 1990's



Supervised Learning



Plot by Rafael S. de Souza

# Photo-z: Artificial Neural Networks

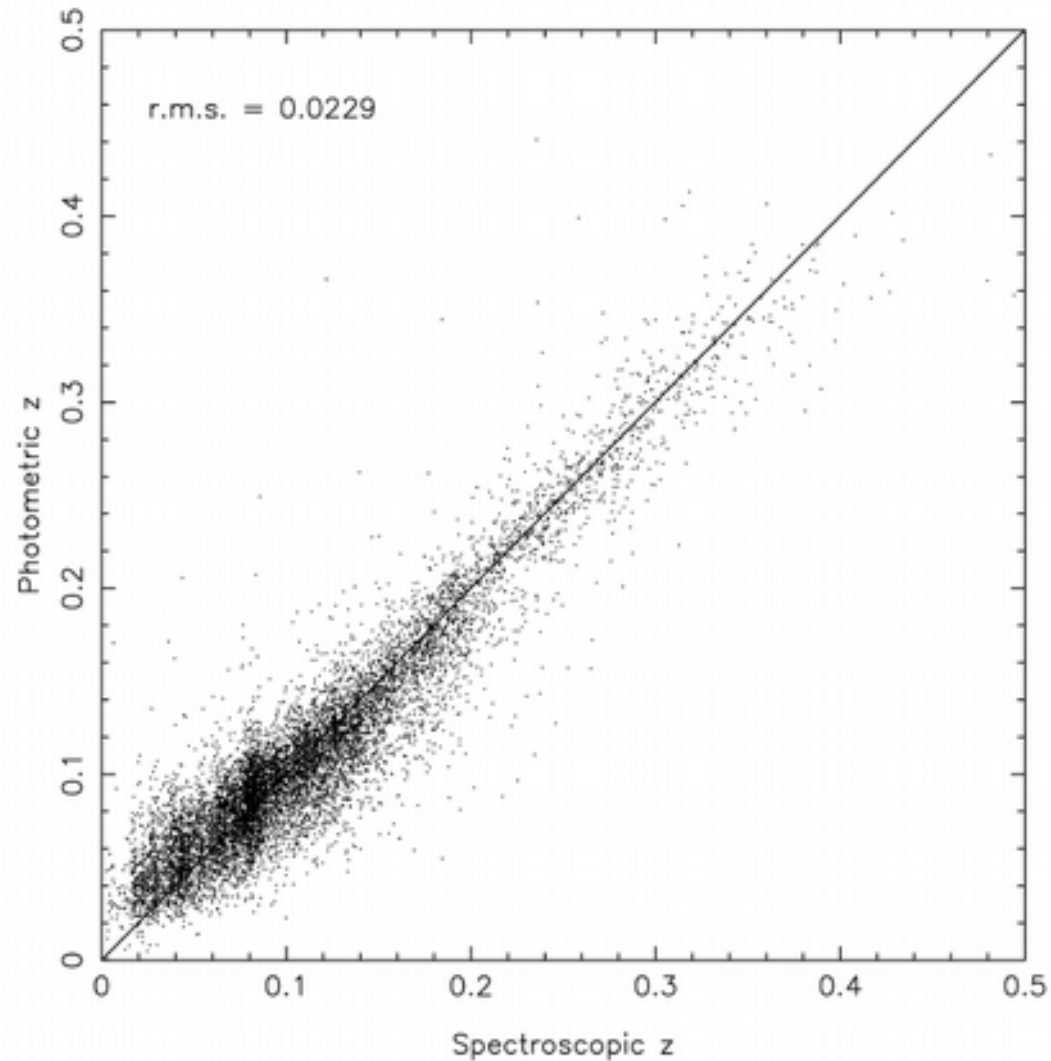
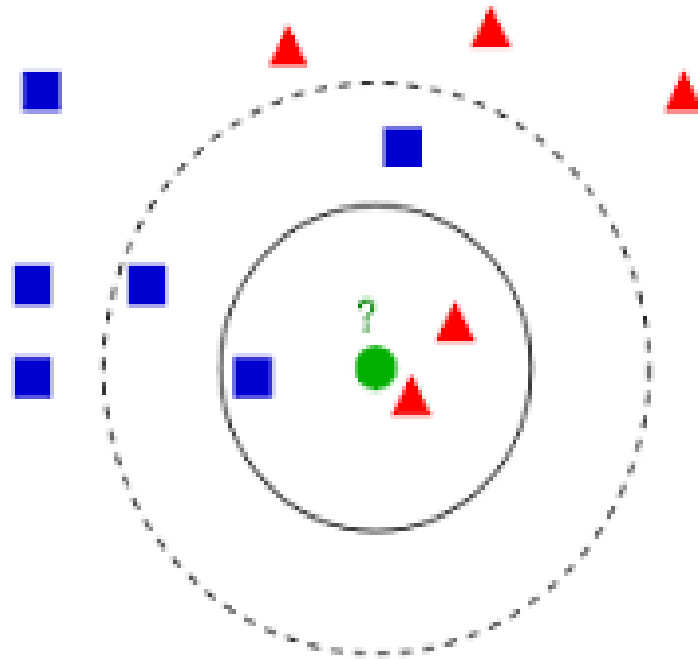


FIG. 2.— Spectroscopic vs. photometric redshifts for ANNz applied to 10,000 galaxies randomly selected from the SDSS EDR.

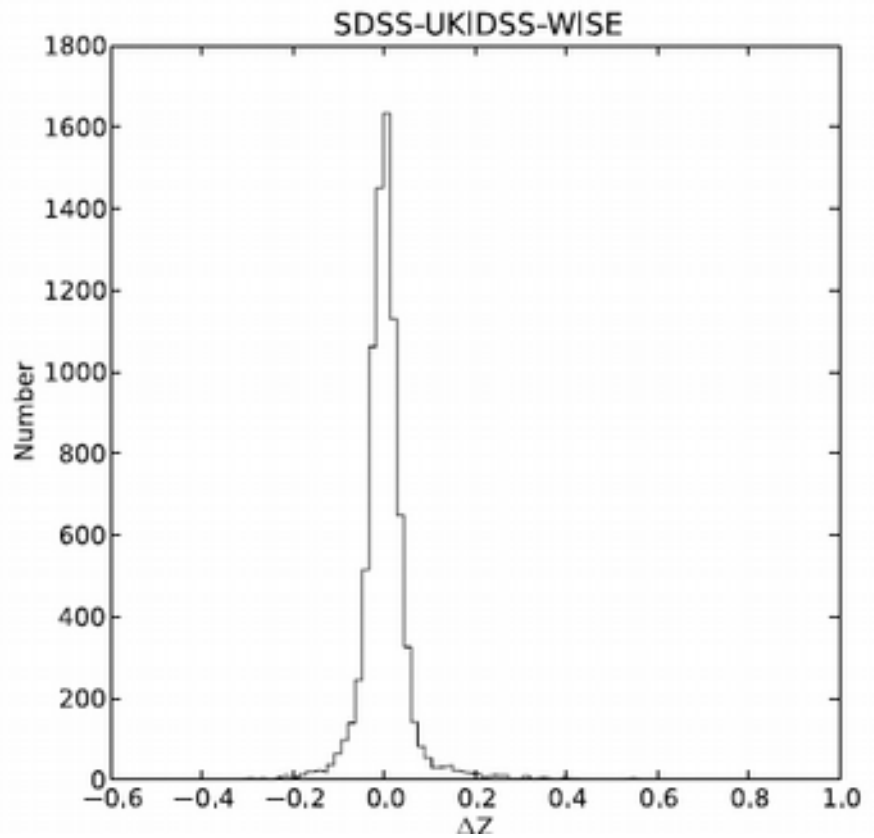
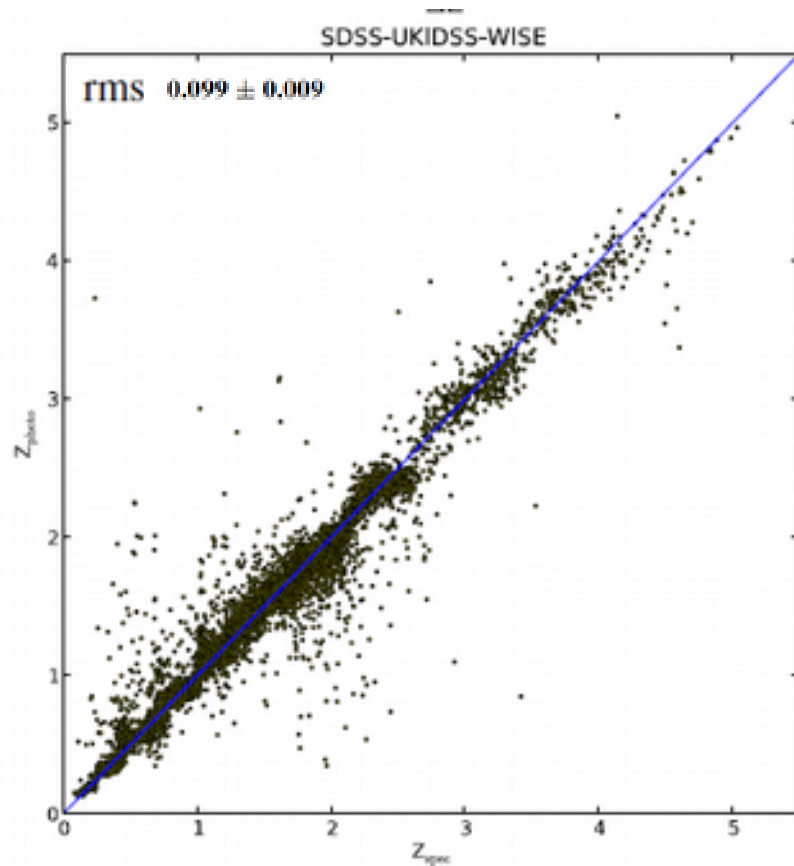
# Photo-z: Nearest Neighbors



Supervised Learning

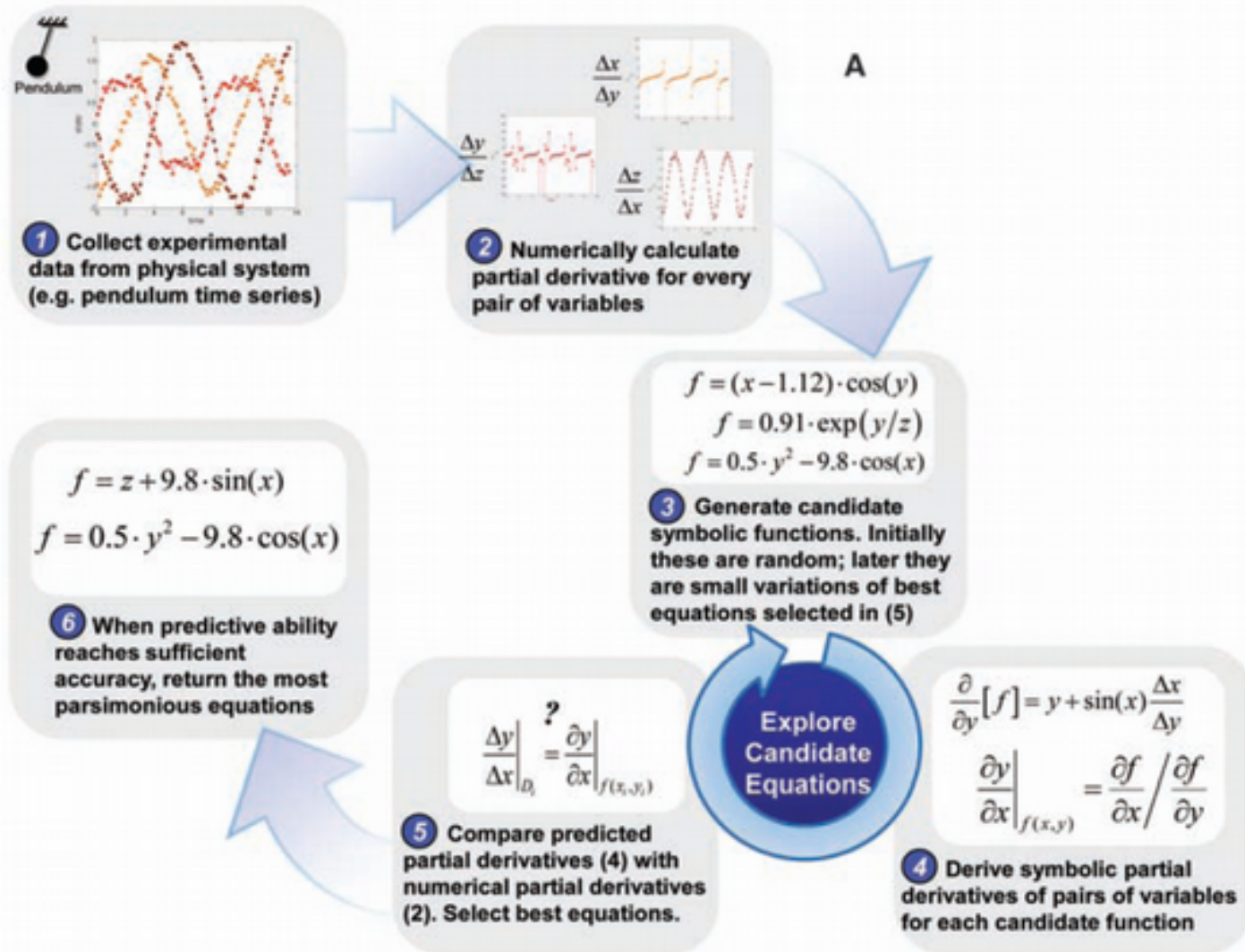


# Photo-z: Nearest Neighbors





# Photo-z: Symbolic Regression



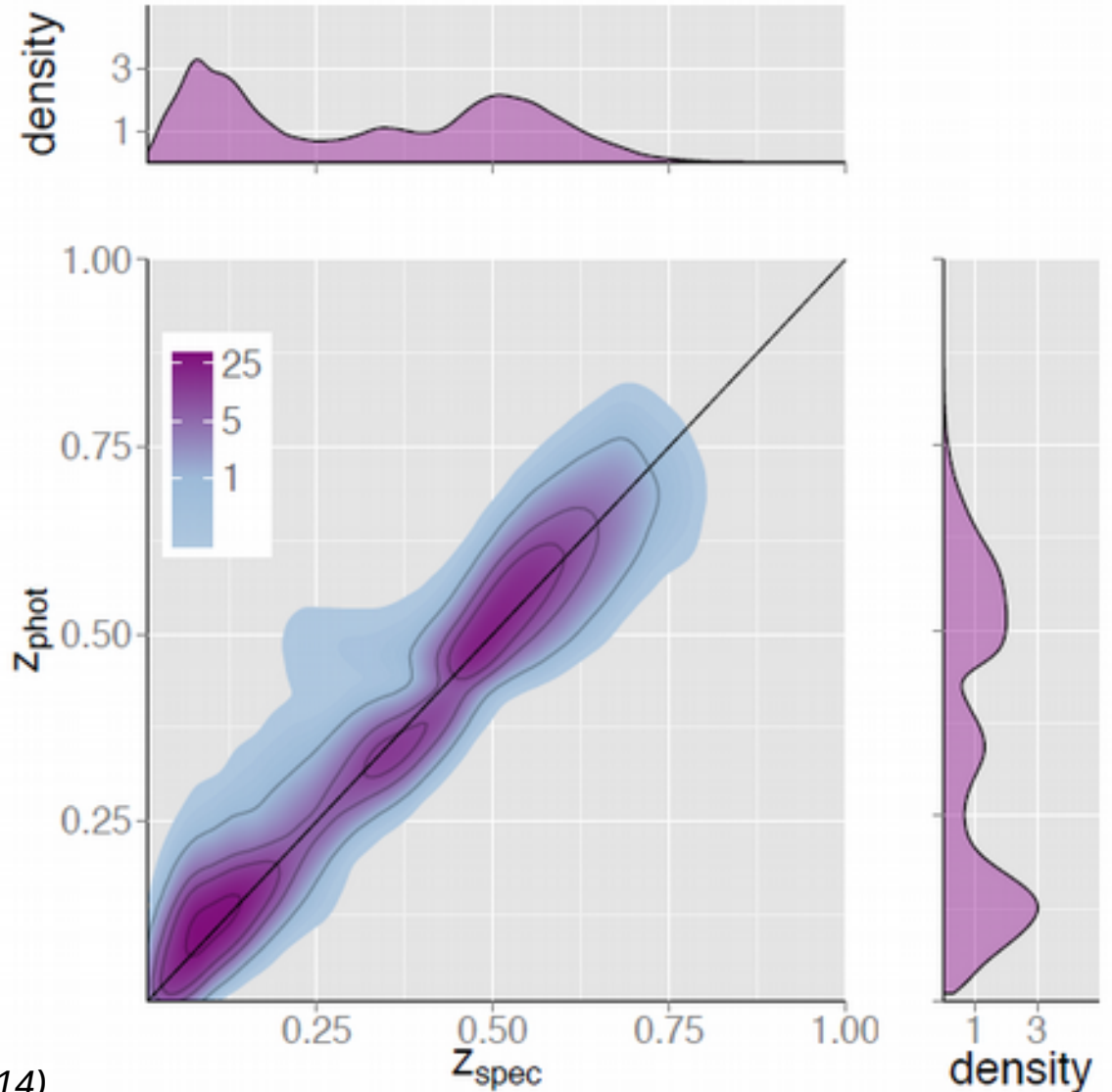
# Photo-z: Symbolic Regression

Final expression:

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

**Parametric model**

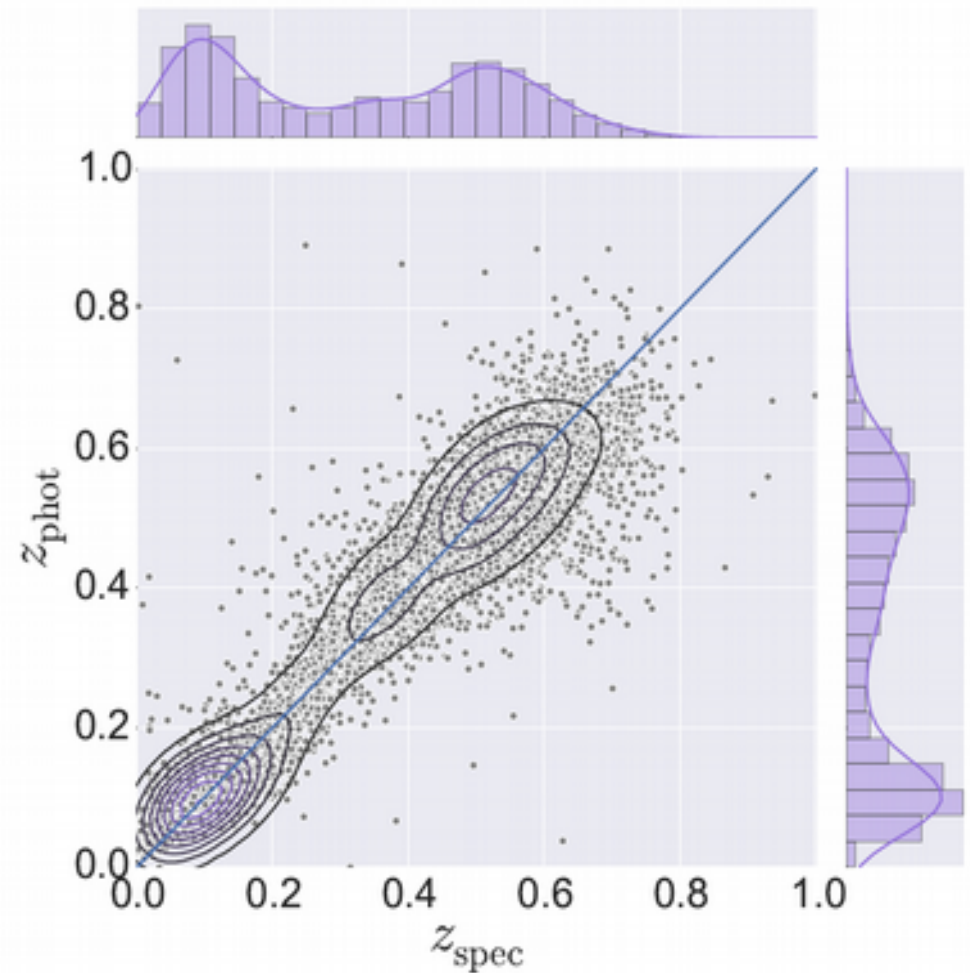
$$\sigma(z_{\text{phot}} - z_{\text{spec}}) / (1 + z_{\text{spec}}) \approx 0.0449$$



# Photo-z: Generalized Linear Models

$$rms(\Delta z) \sim 0.034$$

**Statistical model**

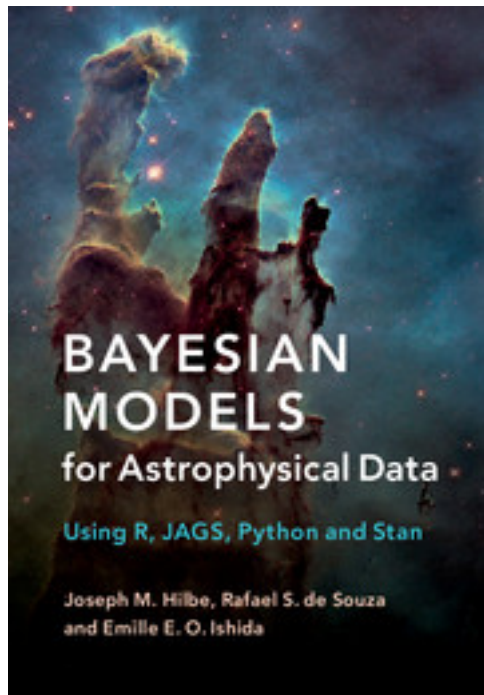


From COIN Residence Program #1:  
*Elliot et al. (incl. Ishida), Astronomy & Computing, 10 (2015)*

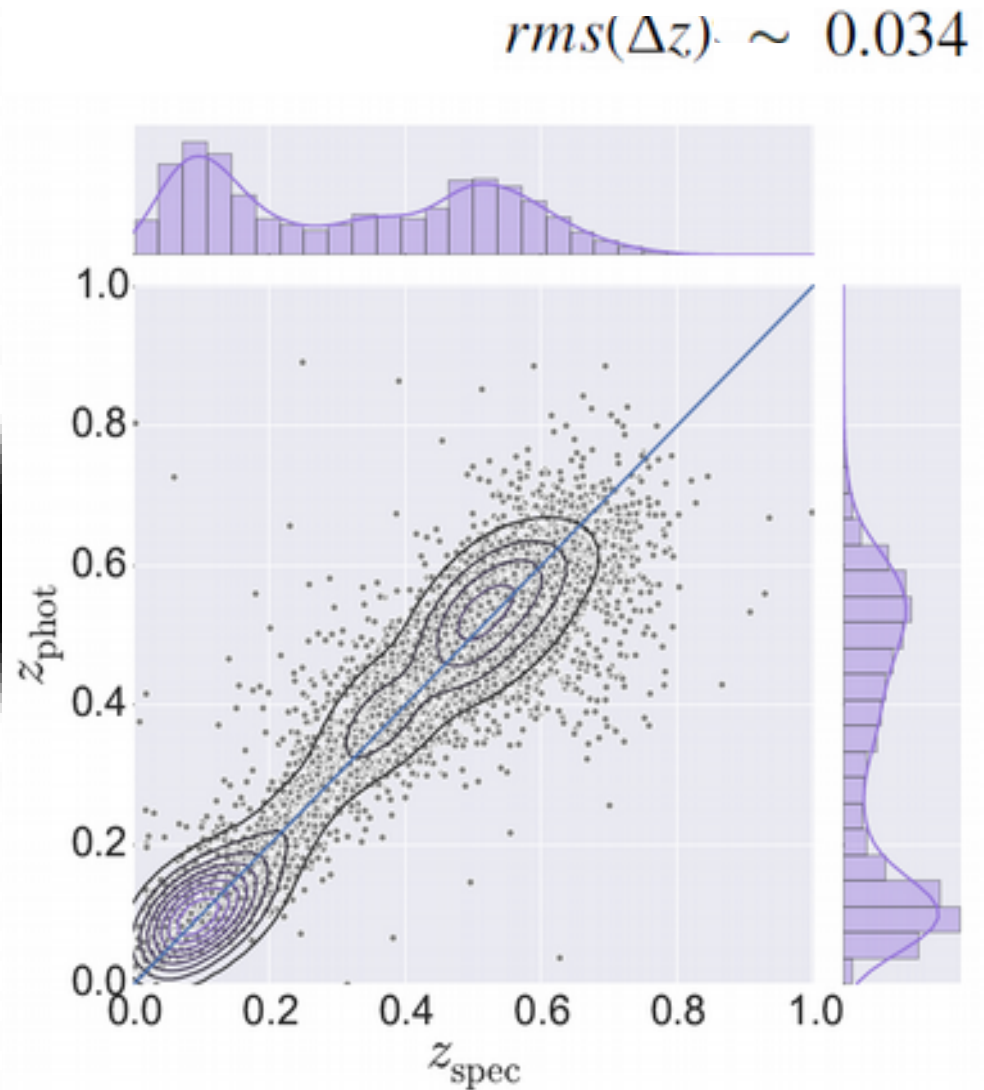
# Photo-z: Generalized Linear Models

More on GLMs  
(Bayesian approach):

<https://github.com/RafaelSdeSouza/ADA8>



Cambridge  
University Press  
May/2017



From COIN Residence Program #1:  
Elliot et al. (incl. **Ishida**), *Astronomy & Computing*, 10 (2015)

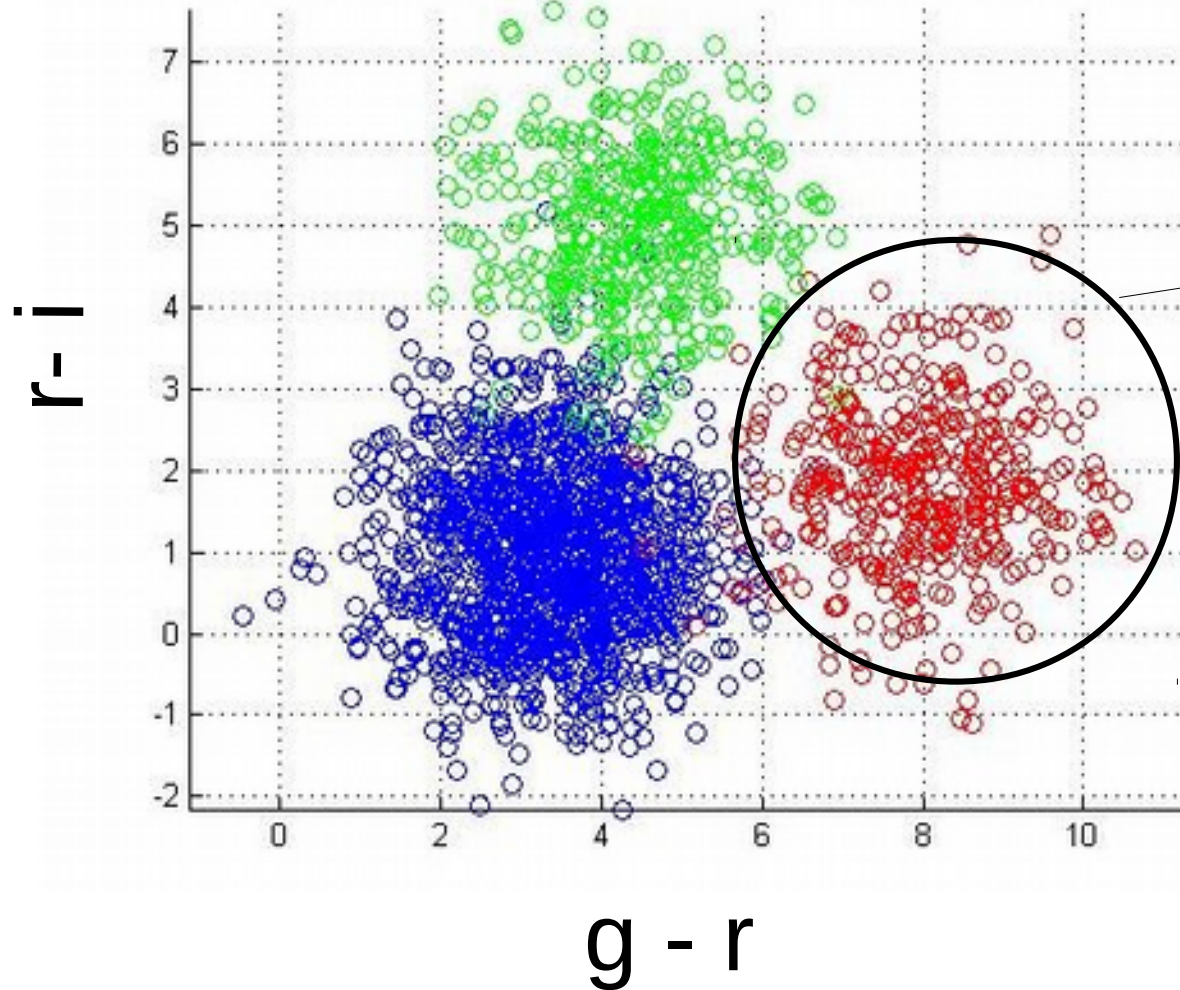
# Photo-z: Local Linear Regression

Unsupervised Learning



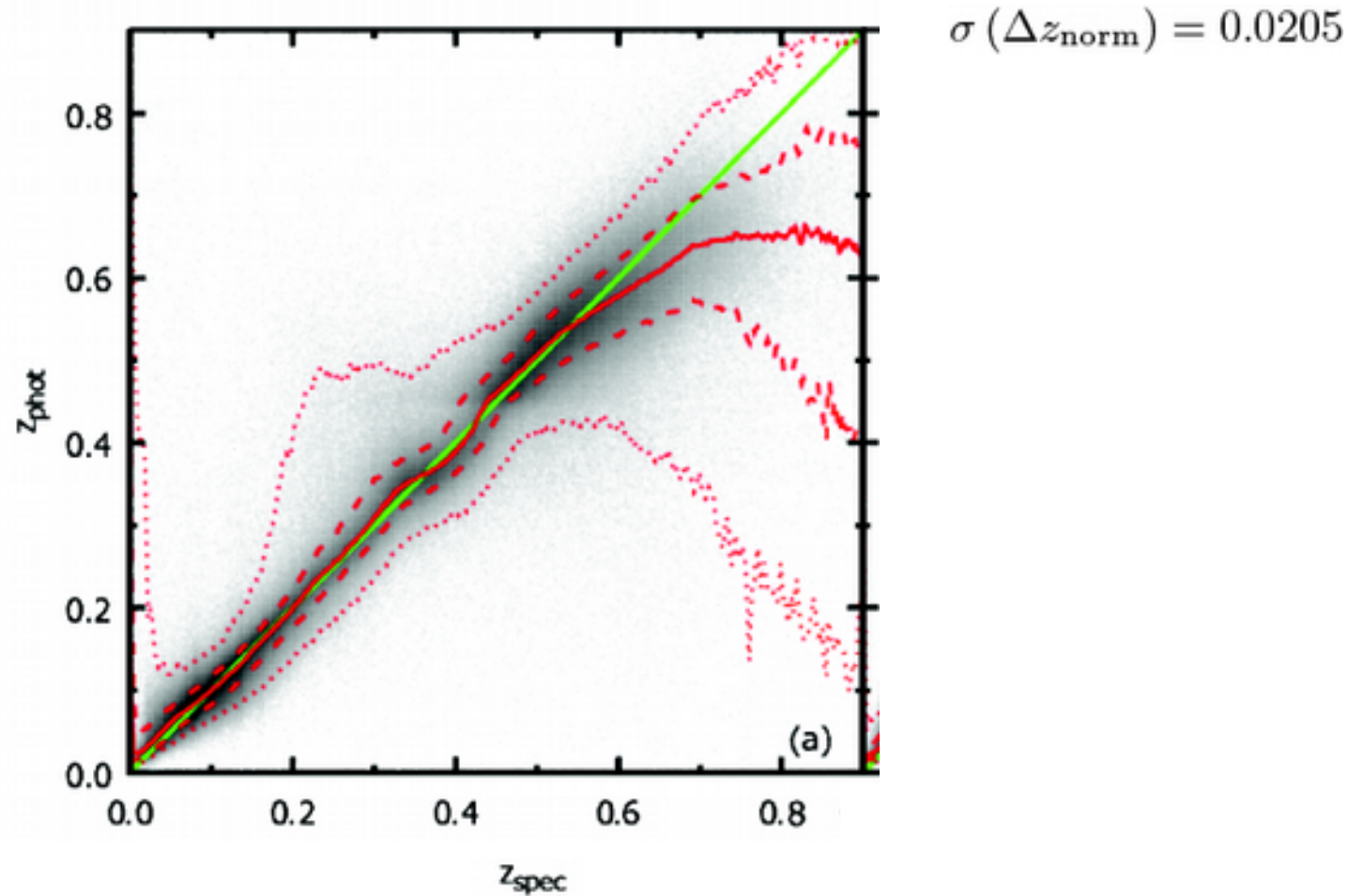
Nearest neighbors

Supervised Learning



# Photo-z: Local Linear Regression

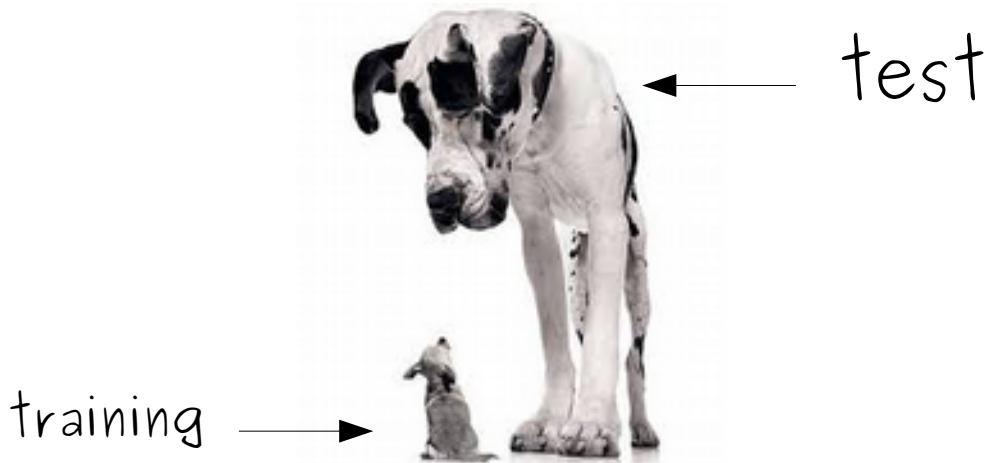
official SDSS DR12 Photoz method



# Summary of results:

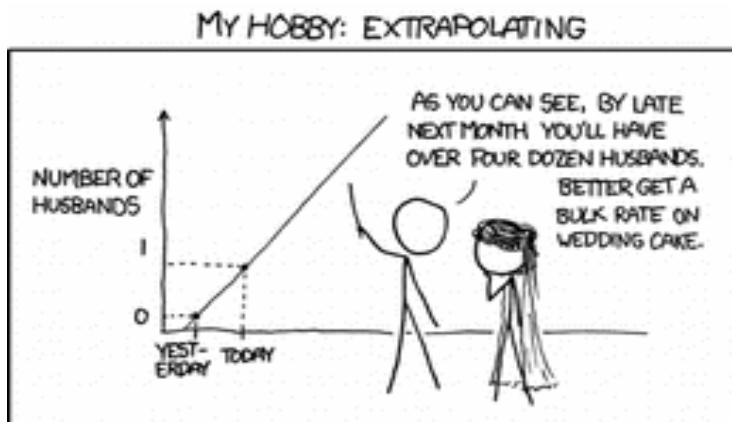


# Challenges



**Measurement errors**

supervised methods cannot extrapolate



Impossible to get a representative training sample

Training will always be:

Brighter

Closer

Higher data quality

Diverse population



# Challenges

Impossible to get a  
representative training  
sample

Training will always be:

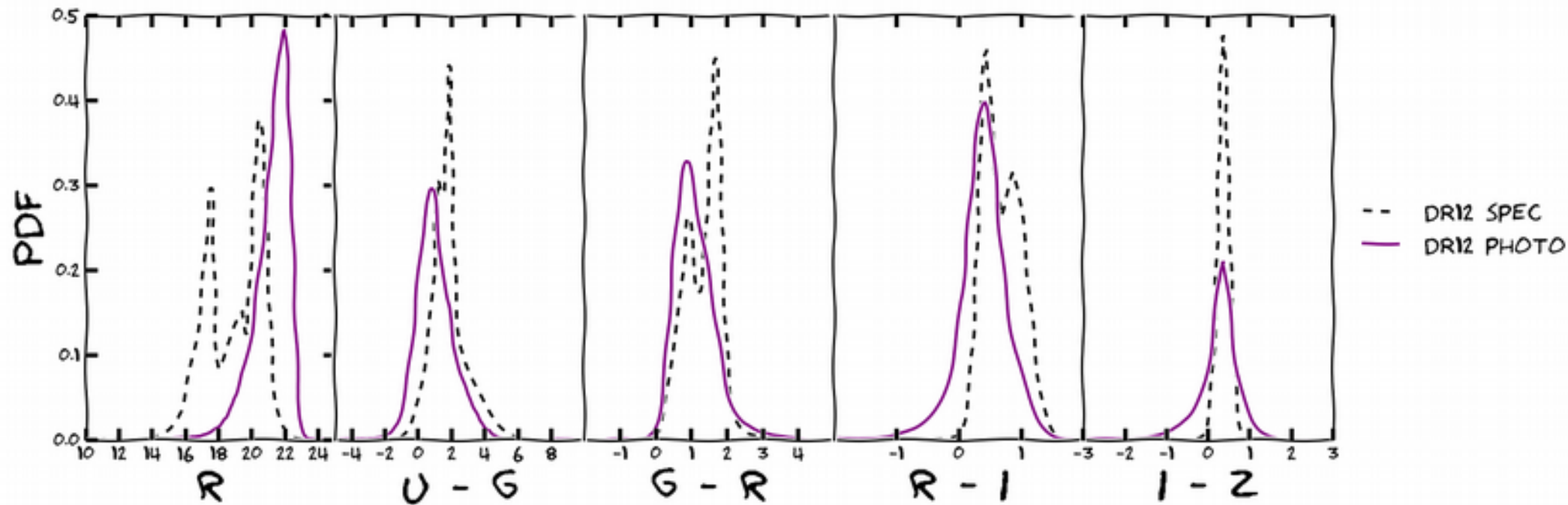
Brighter

Closer

Higher data quality

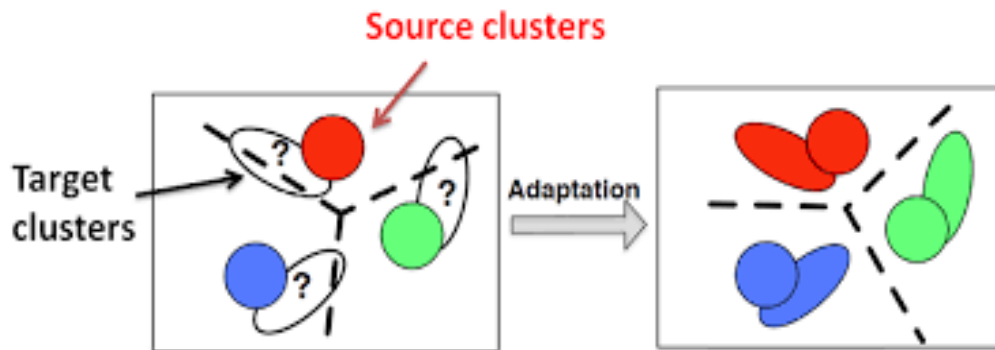
Diverse population

# The quest for representativeness

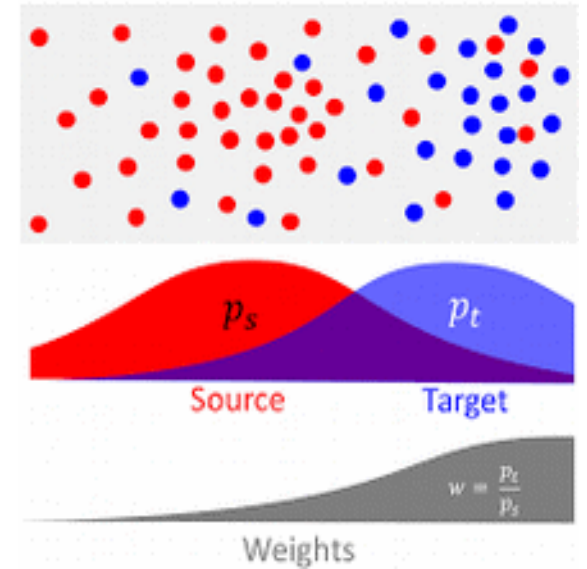


# Domain Adaptation

Define boundaries



Give weights



Feature space transformation



CRP #3 – Budapest, 2016

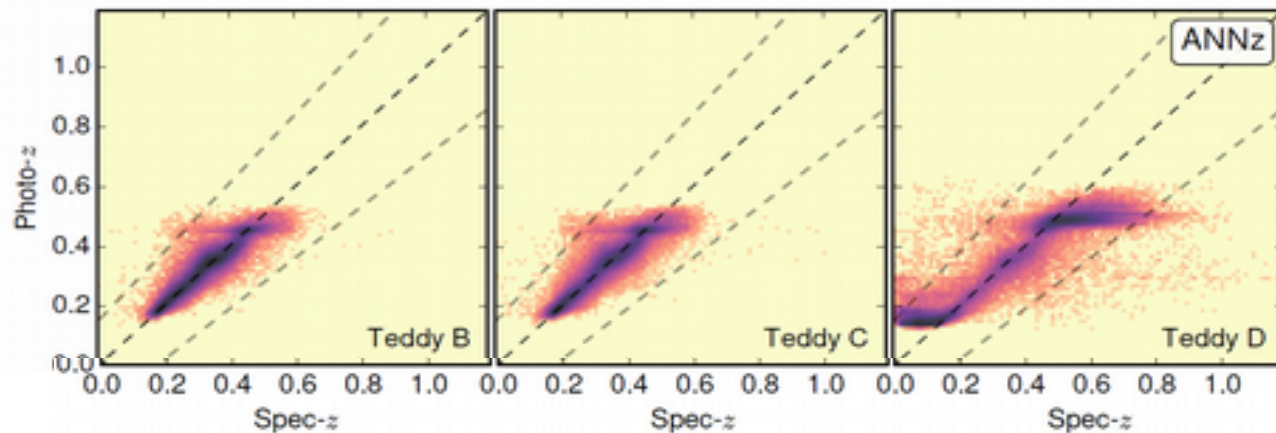
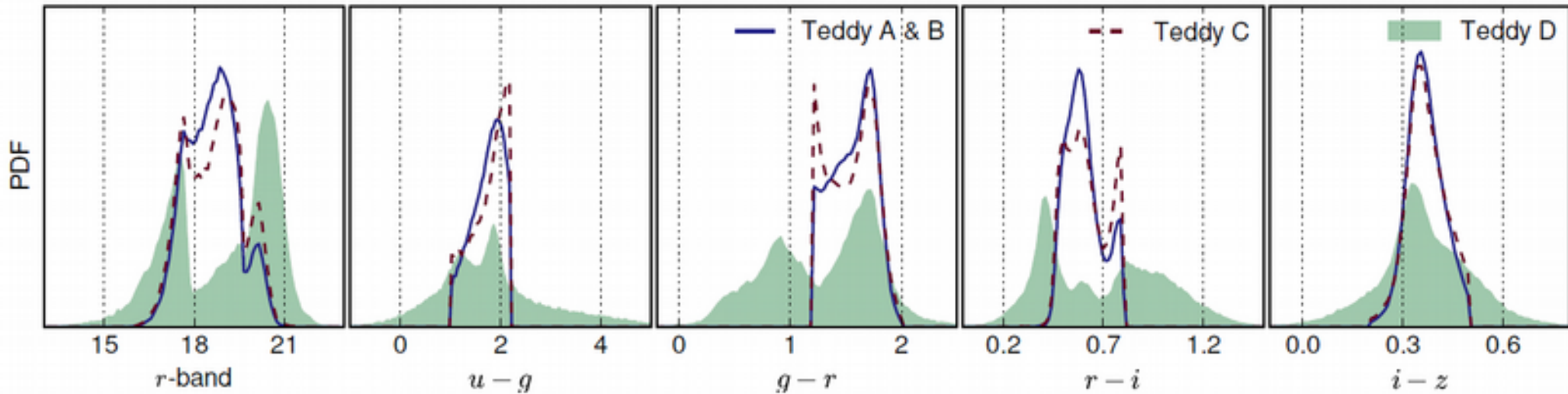


# Teddy catalogue

## Probing the effect of coverage



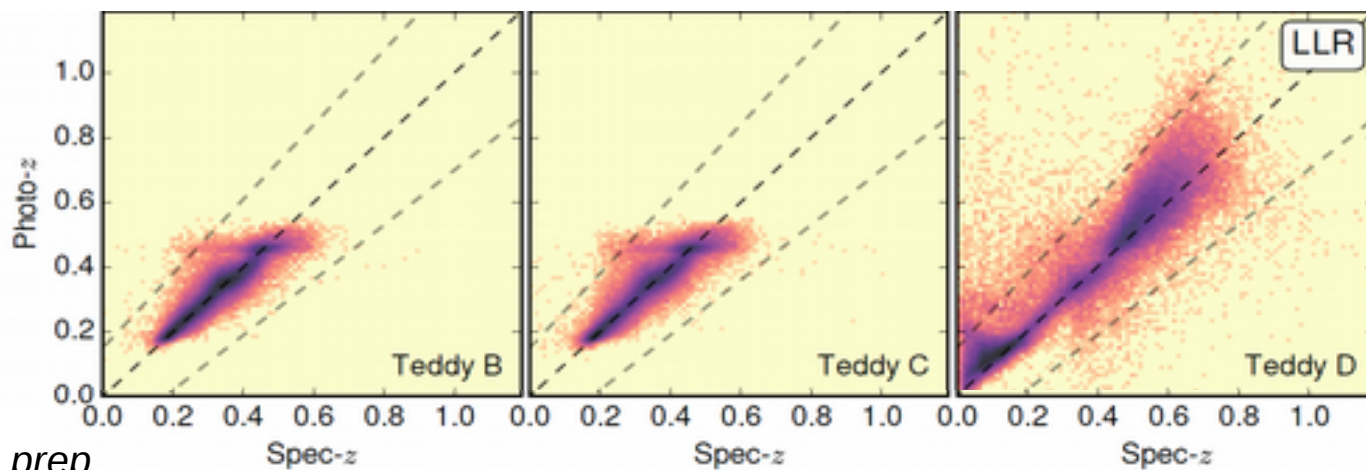
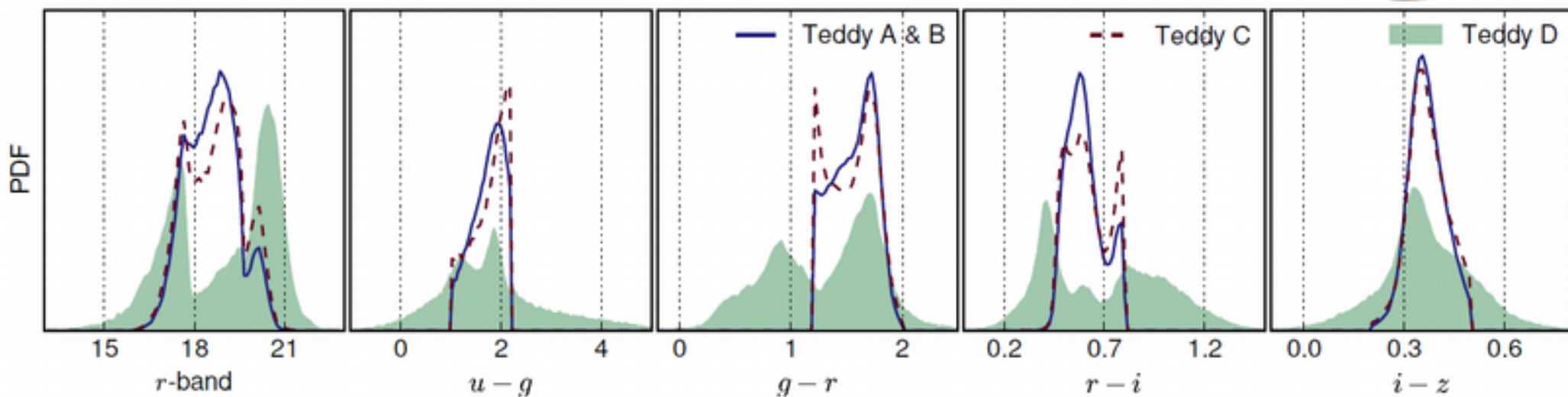
- A/B follow SDSS spec distribution
- B is completely representative of A
- C has the same coverage but slightly different shape
- D has a wider domain in r-mag and color (no coverage)



# Teddy catalogue

## Probing the effect of coverage

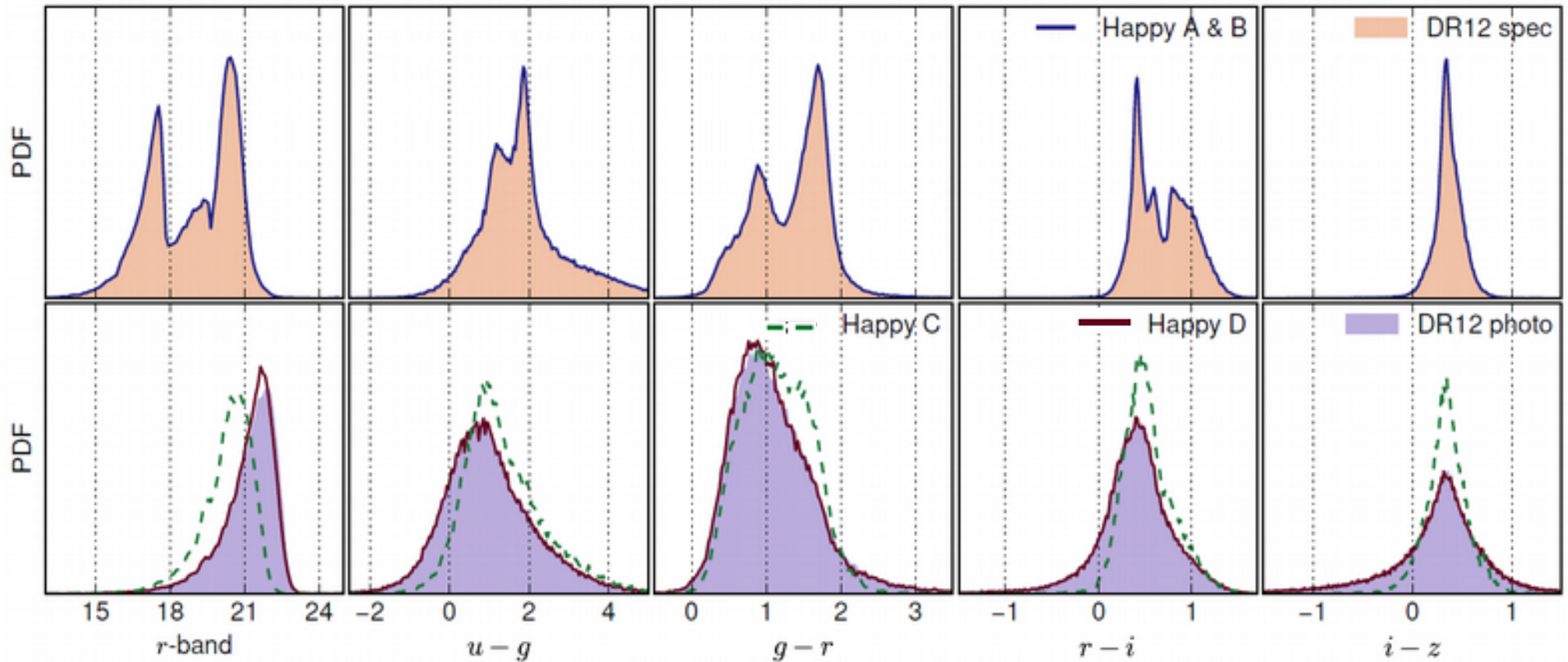
- A/B follow SDSS spec distribution
- B is completely representative of A
- C has the same coverage but slightly different shape
- D has a wider domain in r-mag and color (no coverage)



# Happy catalogue

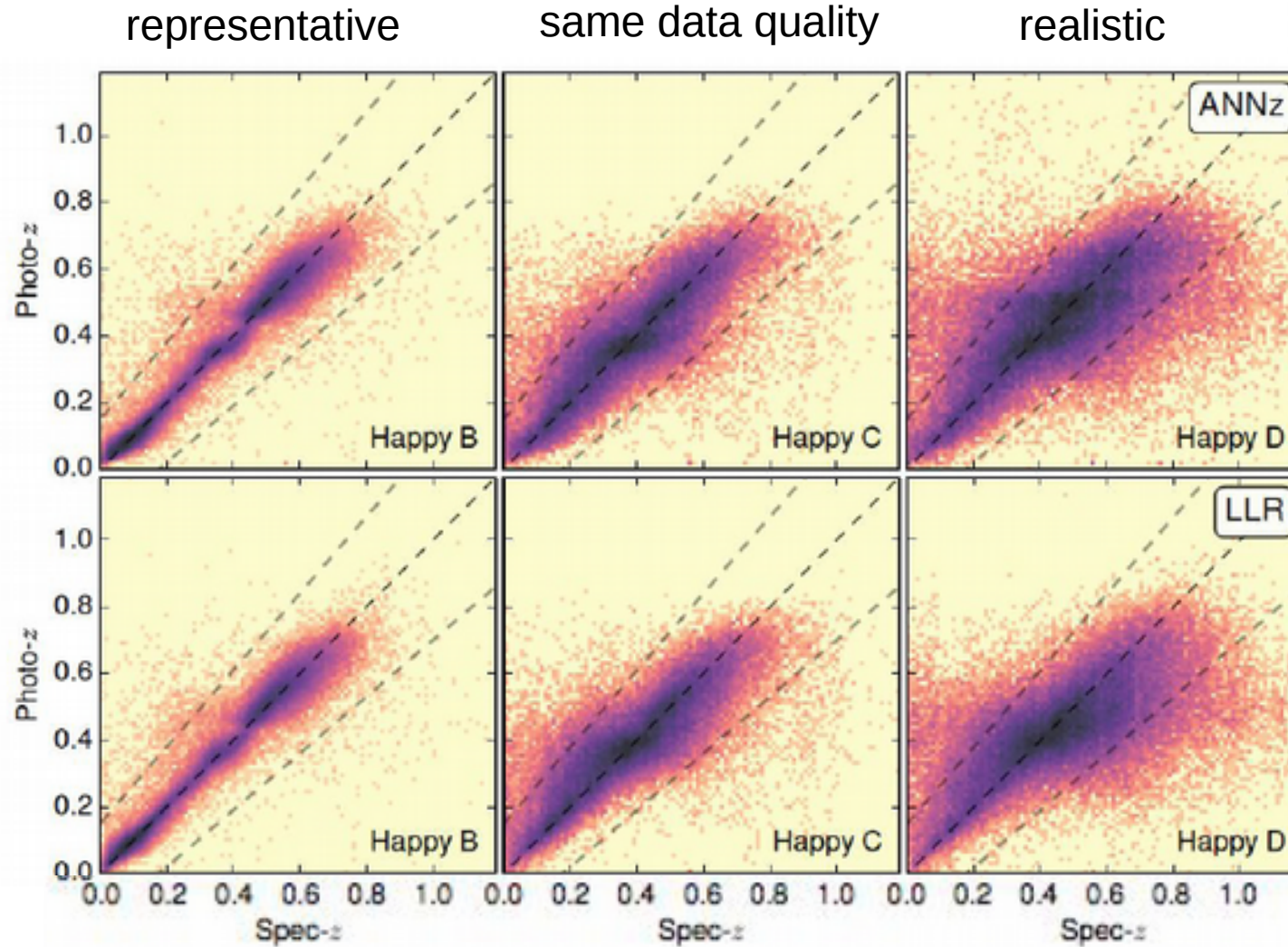
*The effect of coverage + photometric errors*

Happy



# Happy catalogue

*The effect of coverage + photometric errors*



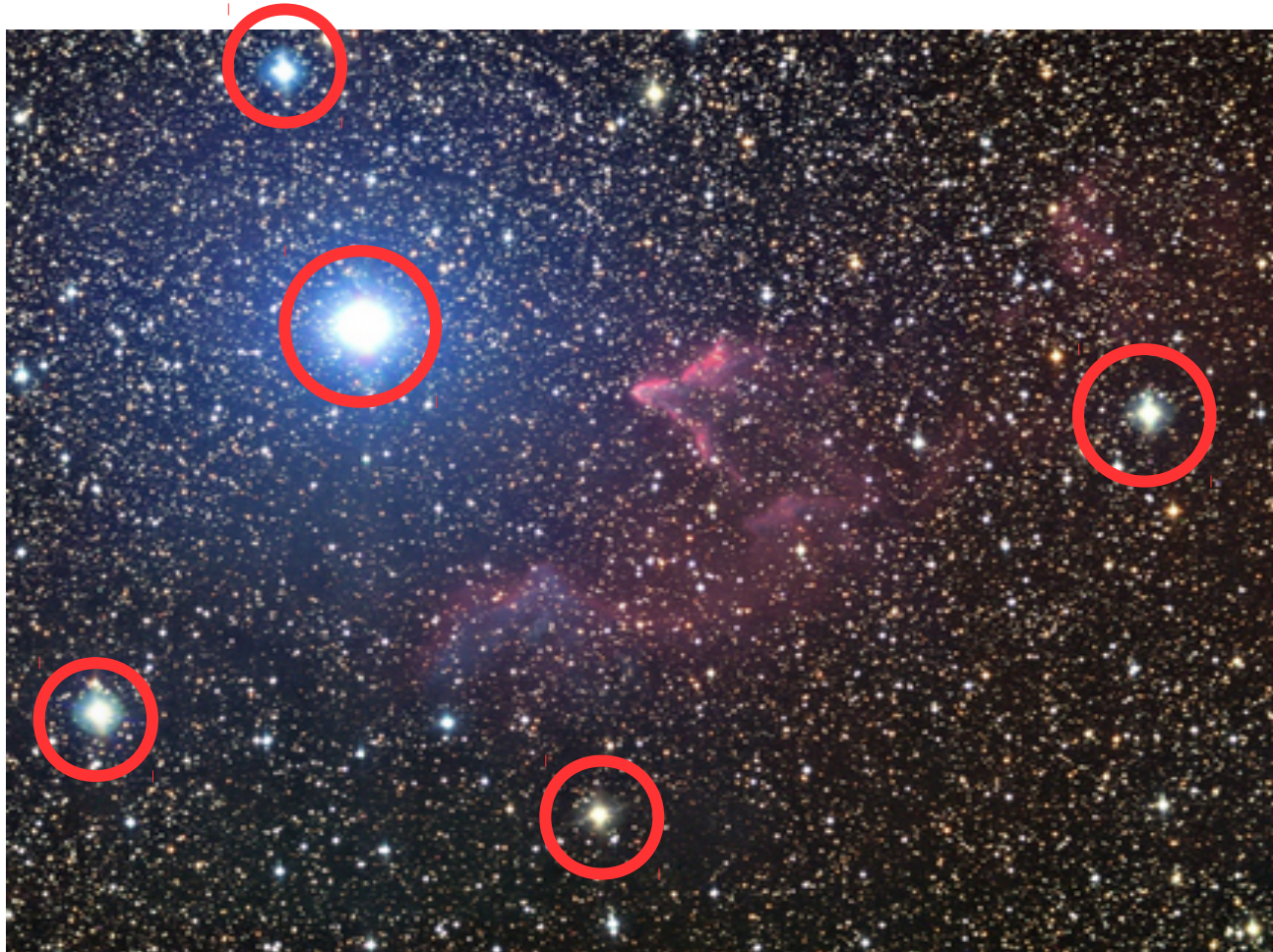
# The Big Picture

(data perspective)



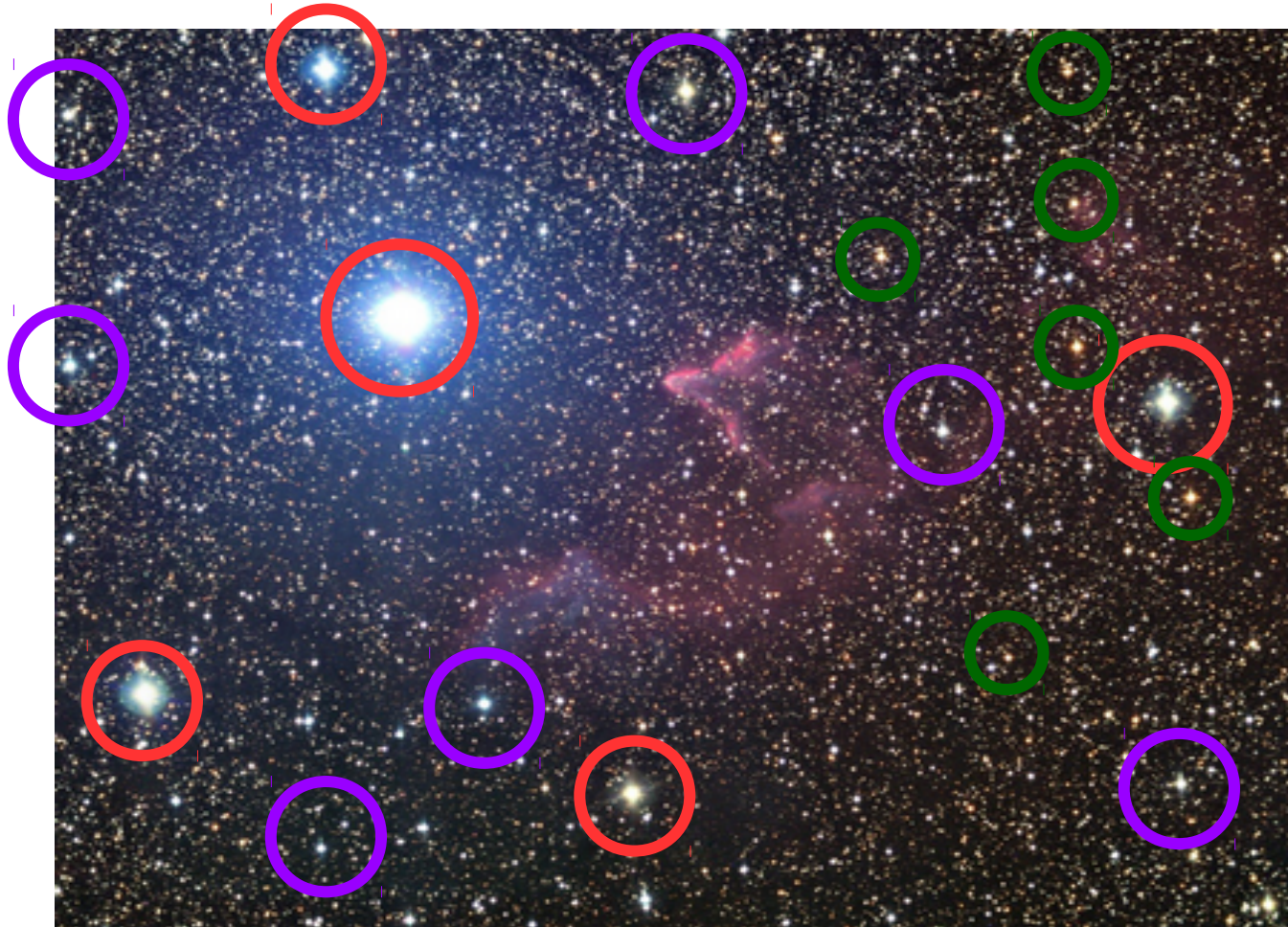
# How are spectroscopic sets constructed?

Take spectra for learning and determine everything else



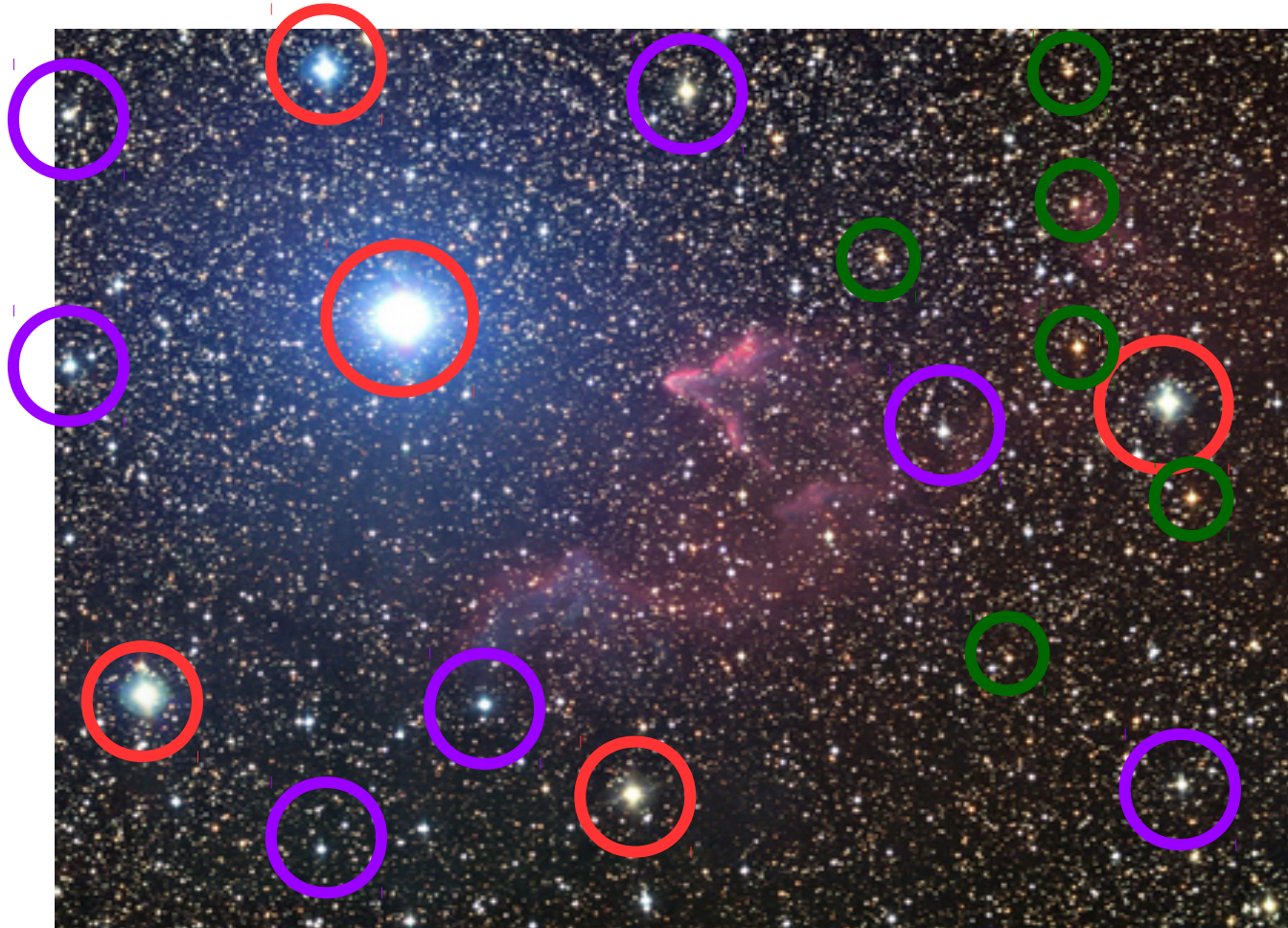
# Alternative approach

Landmark selection + Active Learning



# Alternative approach

Landmark selection + Active Learning



TO BE  
CONTINUED...

Take home message

Astronomy has evolved ...



...there is still a long way to go

Astronomers won't do it alone



The REAL goal is HUMAN learning





THANK YOU!