

# Object classification in SDSS DR12

Farhang Habibi  
LAL - Université Paris SUD-11

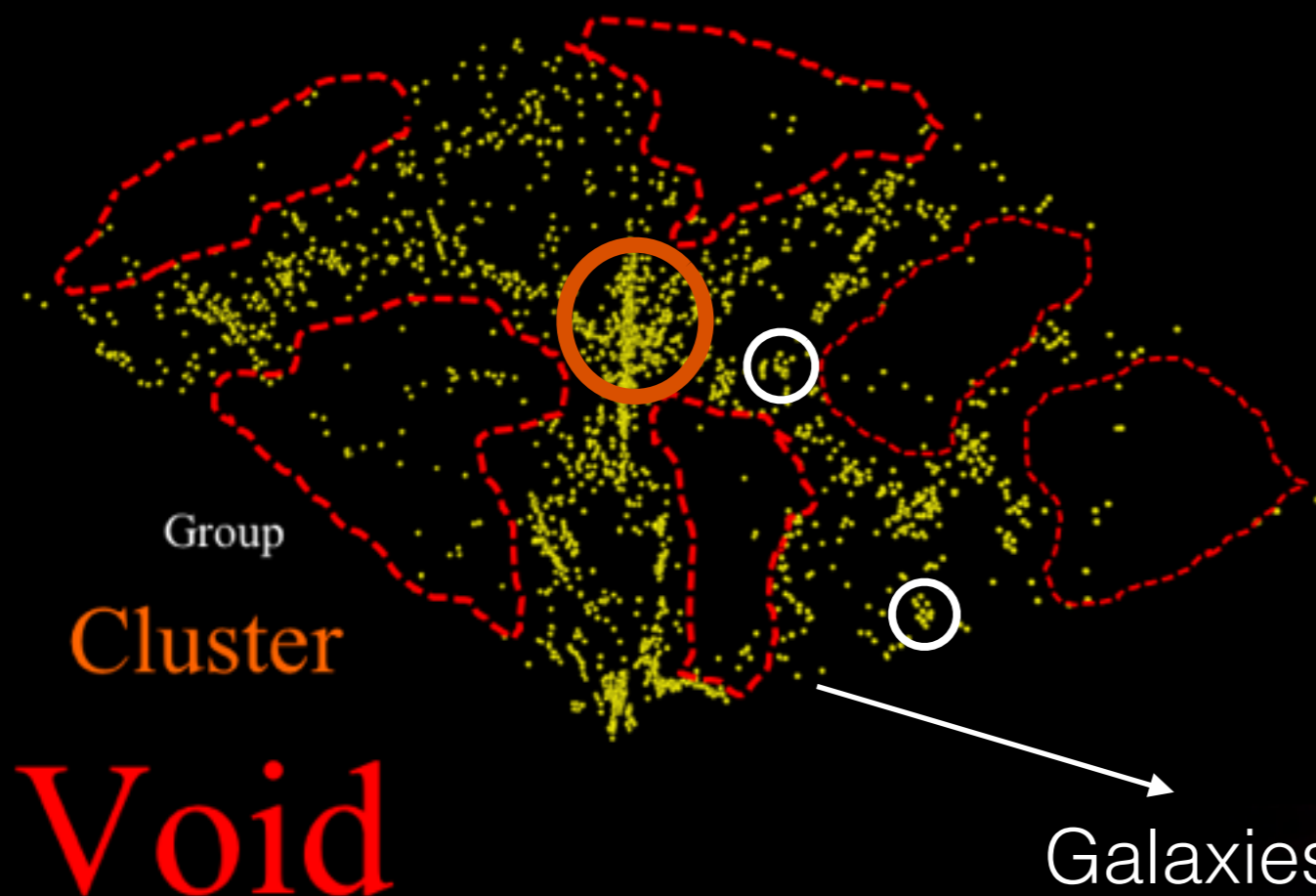
# Aim

**To automatically separate stars, galaxies and Quasars by using the colour indices in the absence of spectroscopic data.**

# Cosmological surveys

**All sky surveys → cosmic structures**

**Deep surveys → structures formation & evolution**



To know about the nature of Dark Matter & Dark Energy



# Object classification

- **Cosmic structures contain galaxies.**
- **Images taken by surveys include galaxies, QSOs and foreground stars.**
- **How to separate these three objects?**

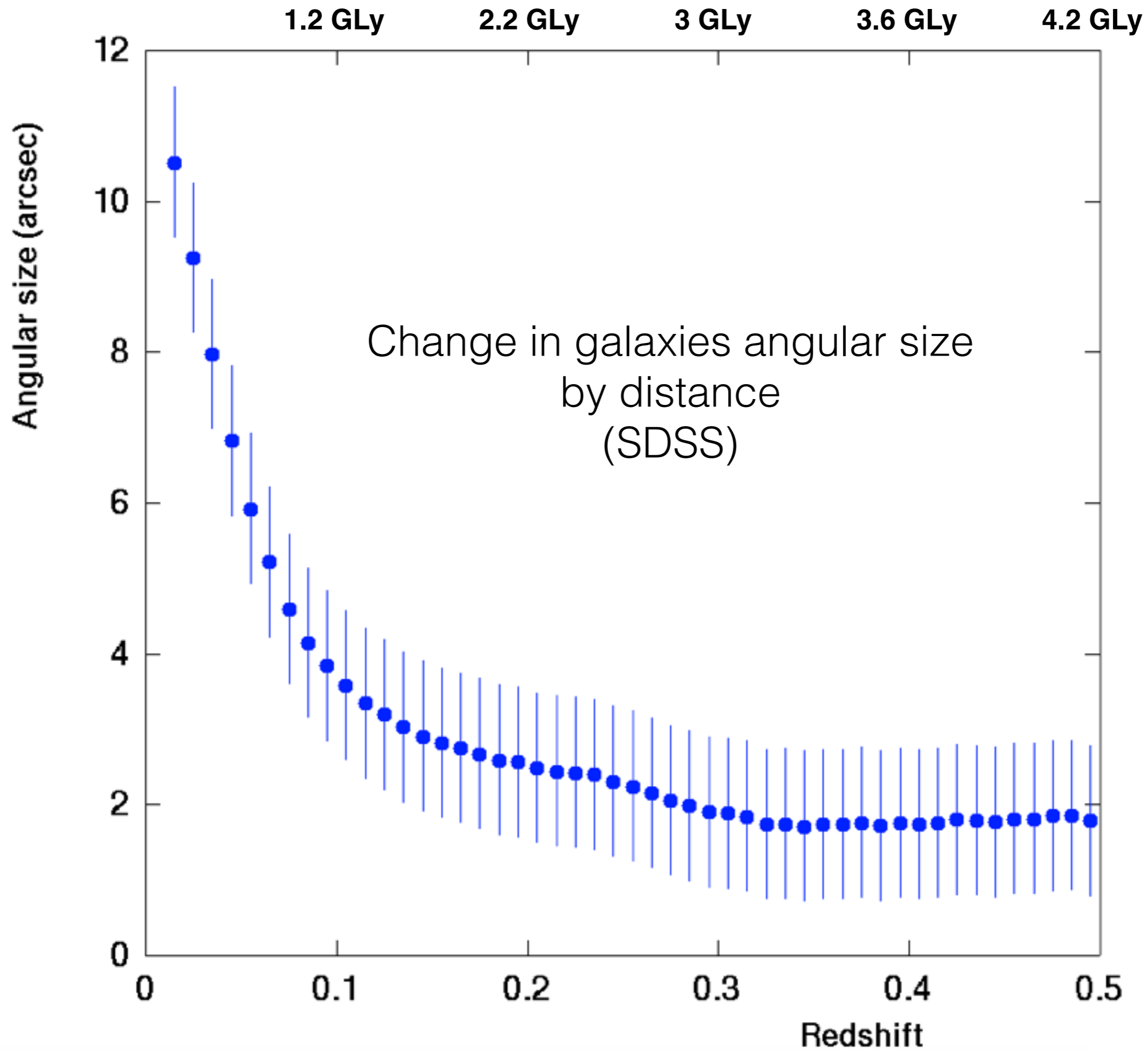


# Nearby galaxies

A deep field image of a starry sky. The background is dark with numerous small, distant stars of various colors (white, yellow, orange, red, blue). In the lower-left quadrant, there is a prominent, elongated, yellowish galaxy with a bright core and a diffuse, glowing envelope. The galaxy is oriented horizontally and appears to be a nearby galaxy.

luminosity spread on CCD:  
stars  $\sim 1$  arcsec  
galaxies  $\sim 10$  arcsec

full moon  $\sim 1800$  arcsec





# Far galaxies

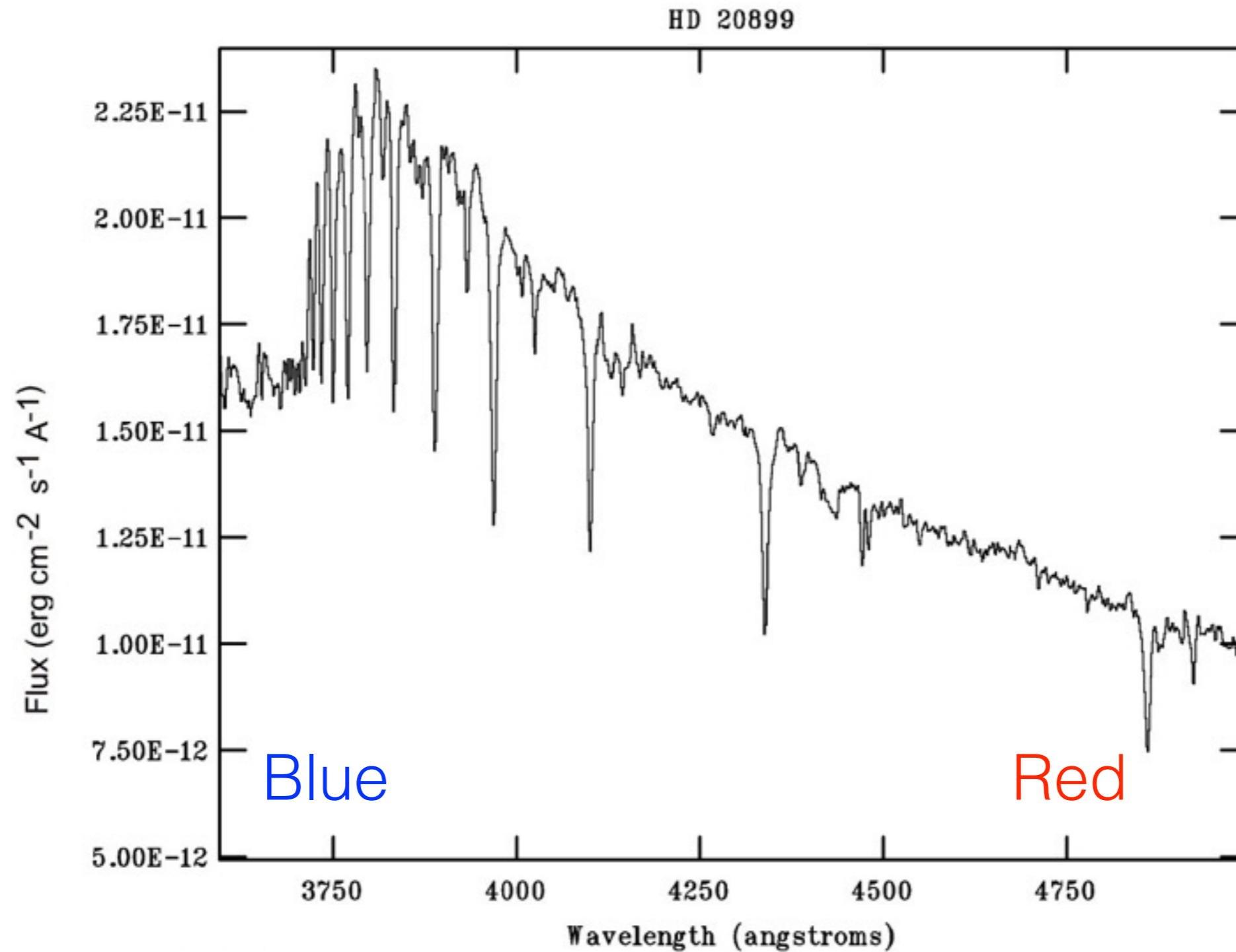
**luminosity spread on CCD:**

**stars, QSOs  $\sim 1$  arcsec**

**galaxies  $\sim 1$  arcsec**

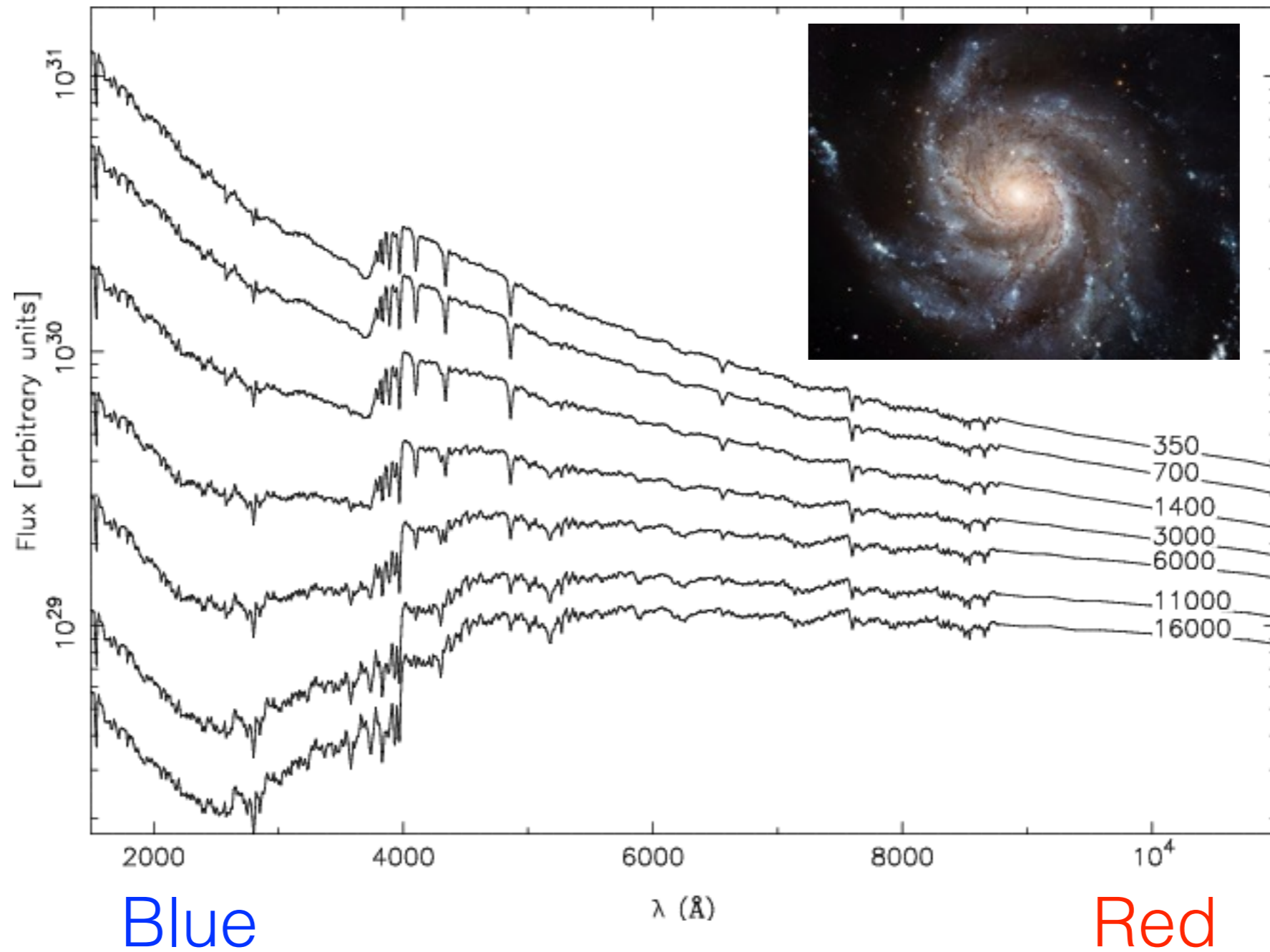


# Stellar energy spectrum

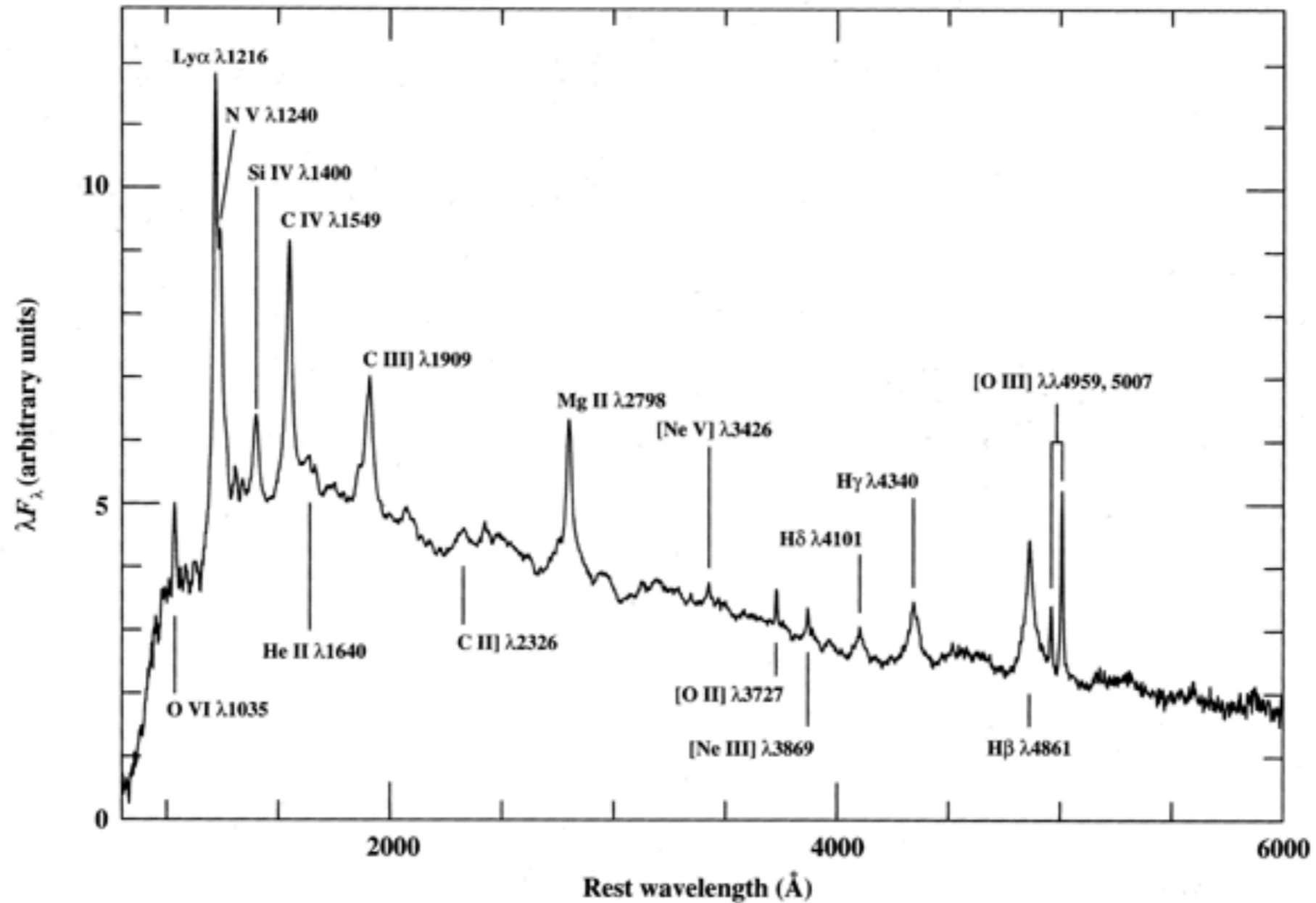




# Galactic energy spectrum



# QSO energy spectrum

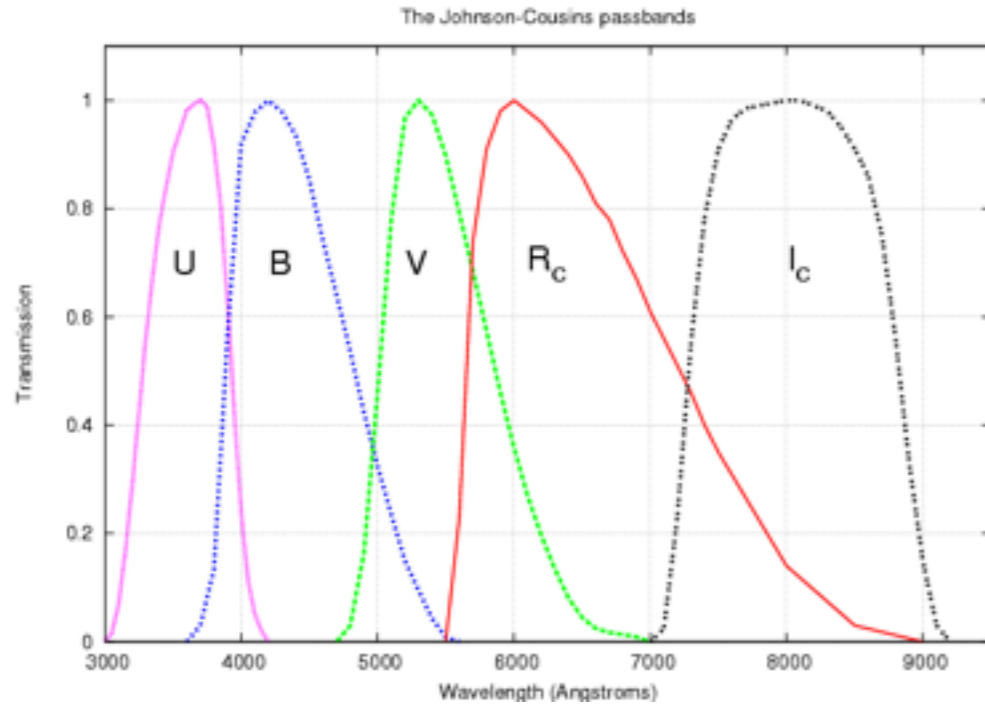


Blue

Red



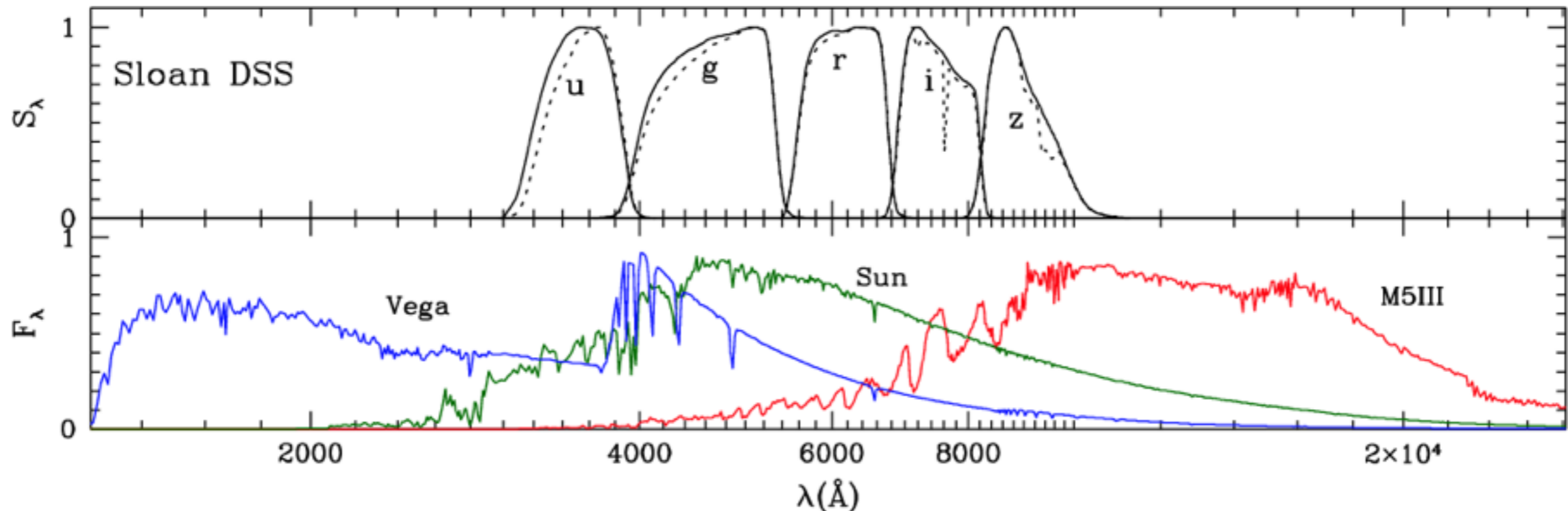
# Photometric systems



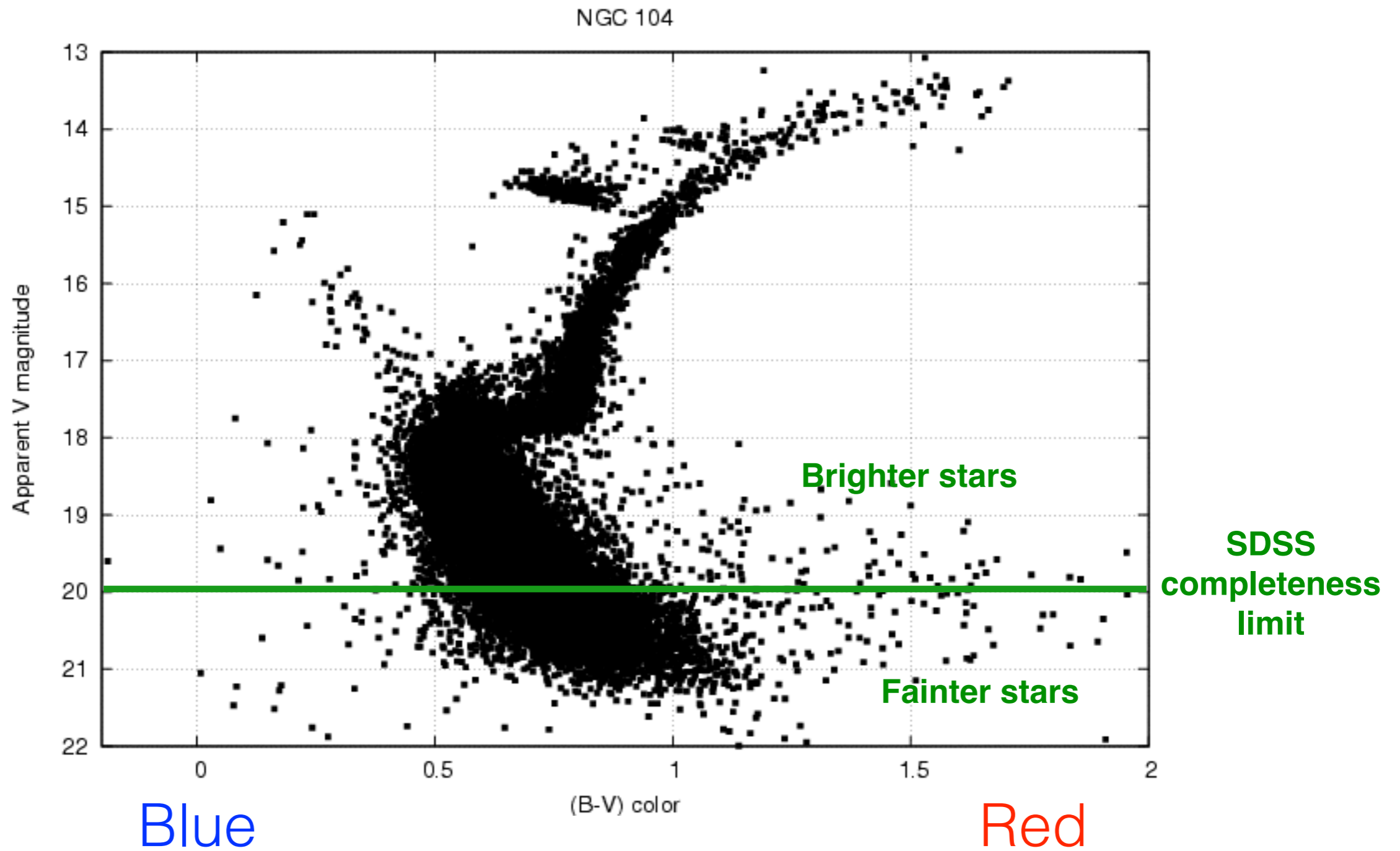
Magnitude in a filter  
 $\sim -\text{Log}(\text{Flux in the same filter})$

Colour = Mag(filter 1) - Mag(filter 2)

**higher colour index: Redder object**

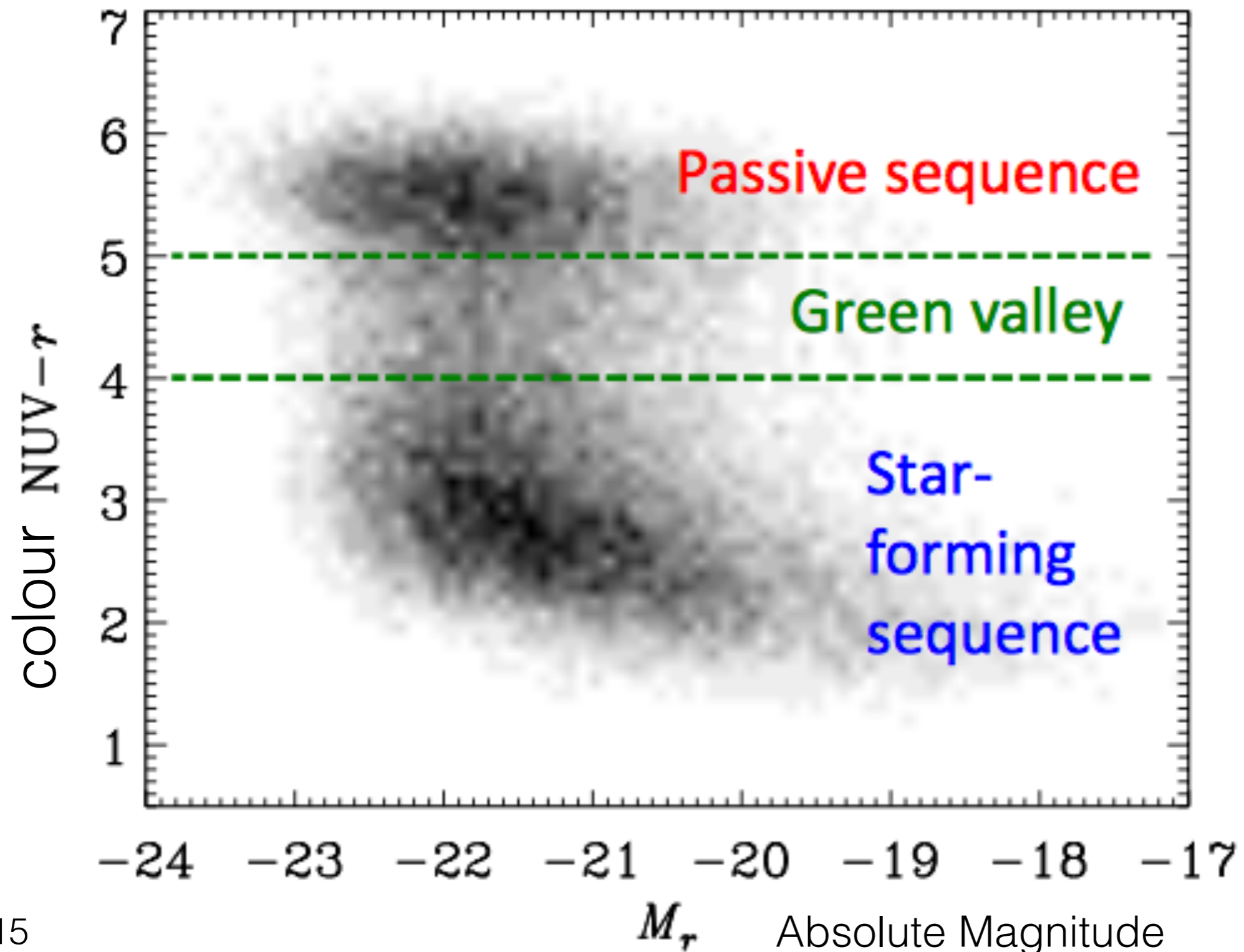


# Stellar colour-magnitude diagram



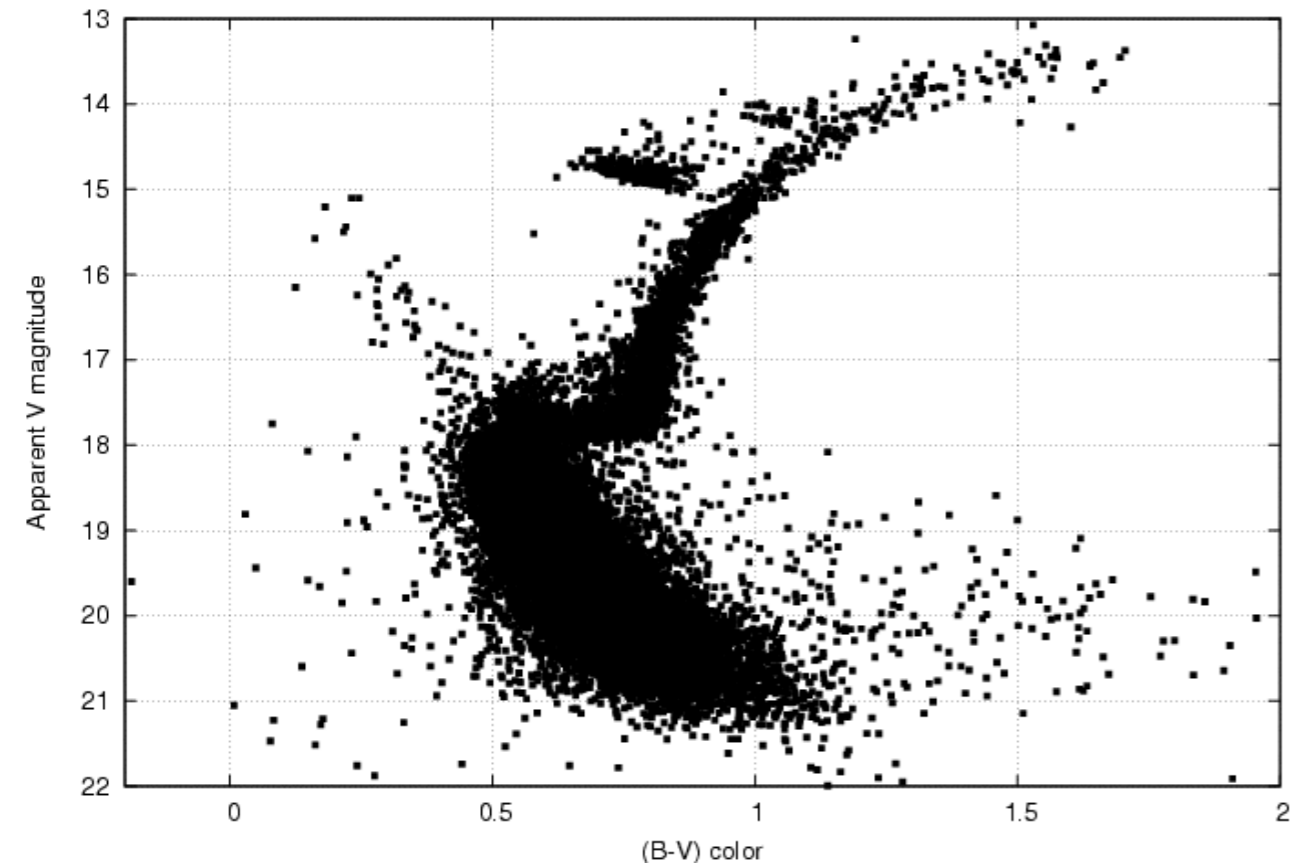


# galactic colour-magnitude diagram

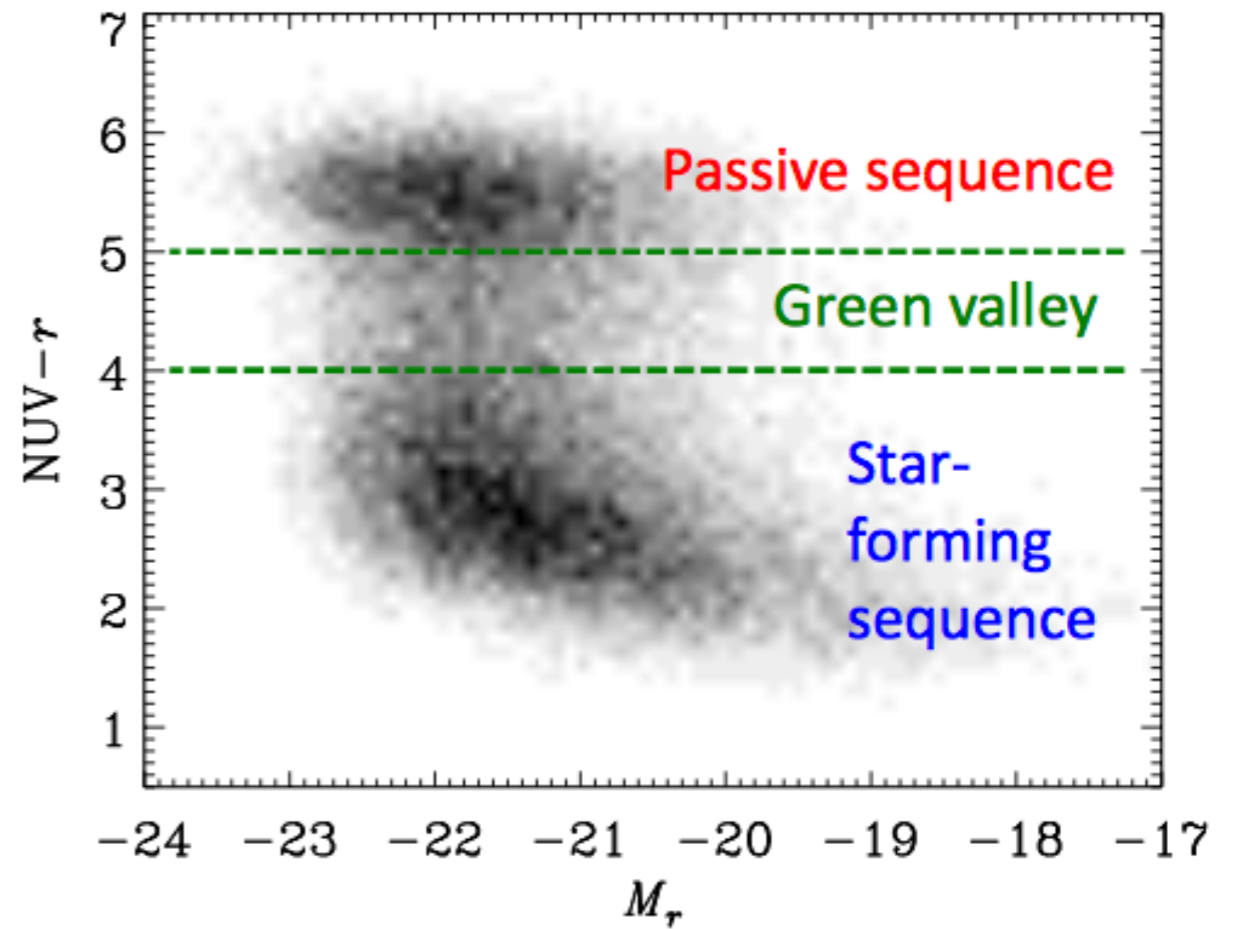


# Colour indices can be used to classify the celestial objects

NGC 104



Stars



Galaxies

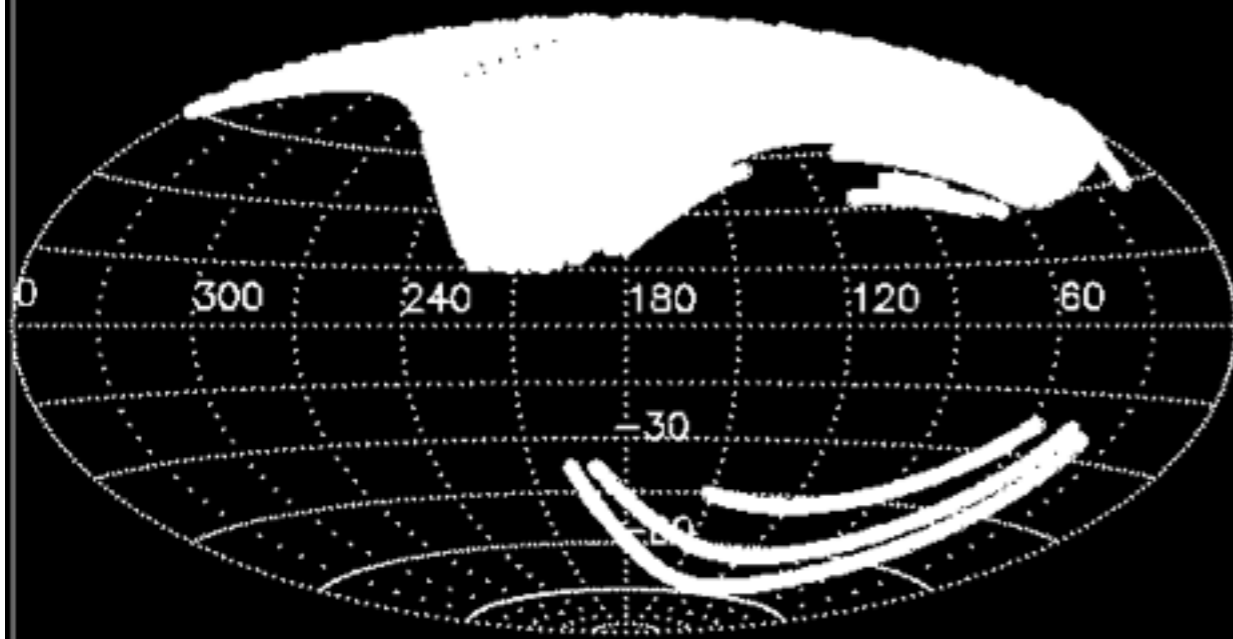


# SDSS survey

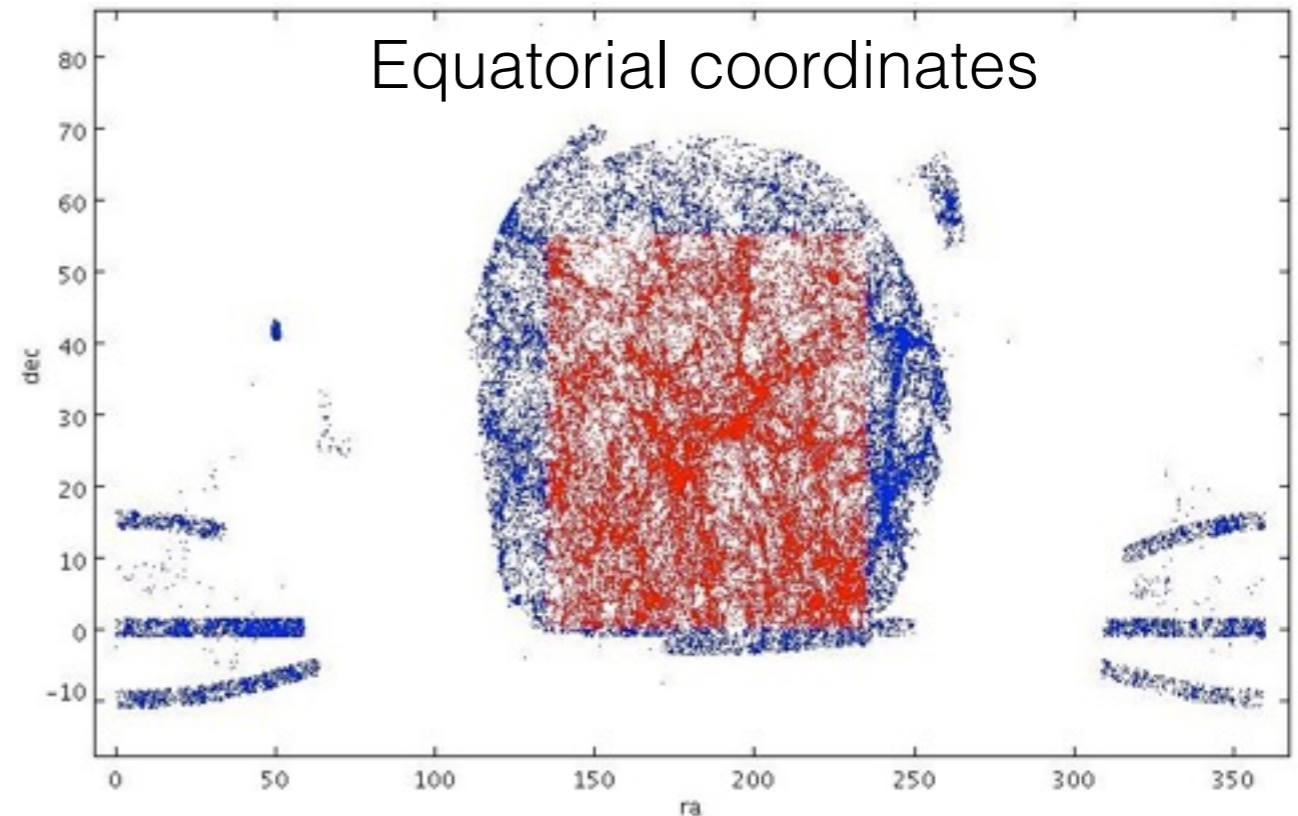
(Sloan digital sky survey)

**.2 m class telescope**  
**.complete up to  $\sim 2.6$  GLy**  
 **$\sim 4$  million spectroscopically**  
**classified objects**

Galactic coordinates

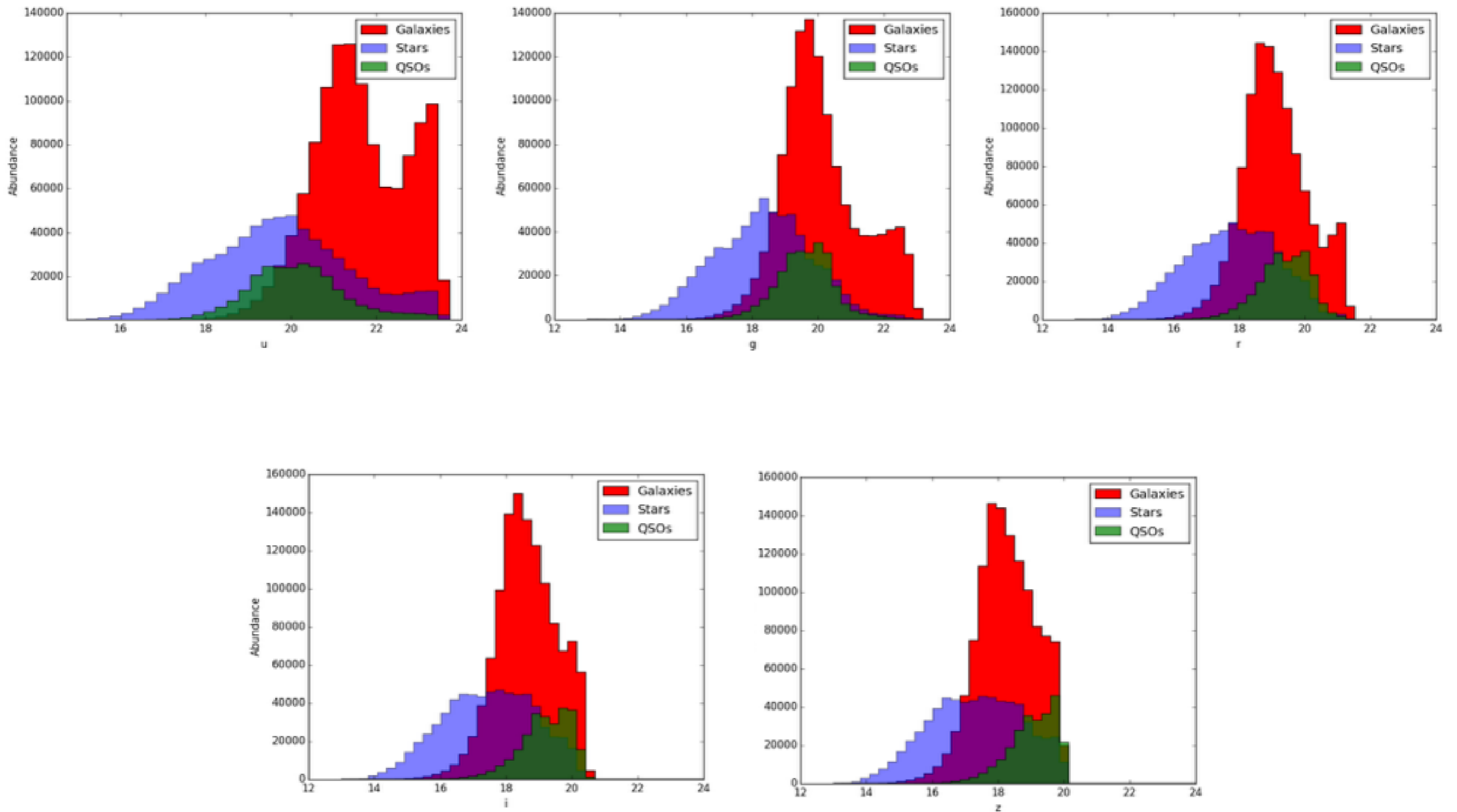


Equatorial coordinates

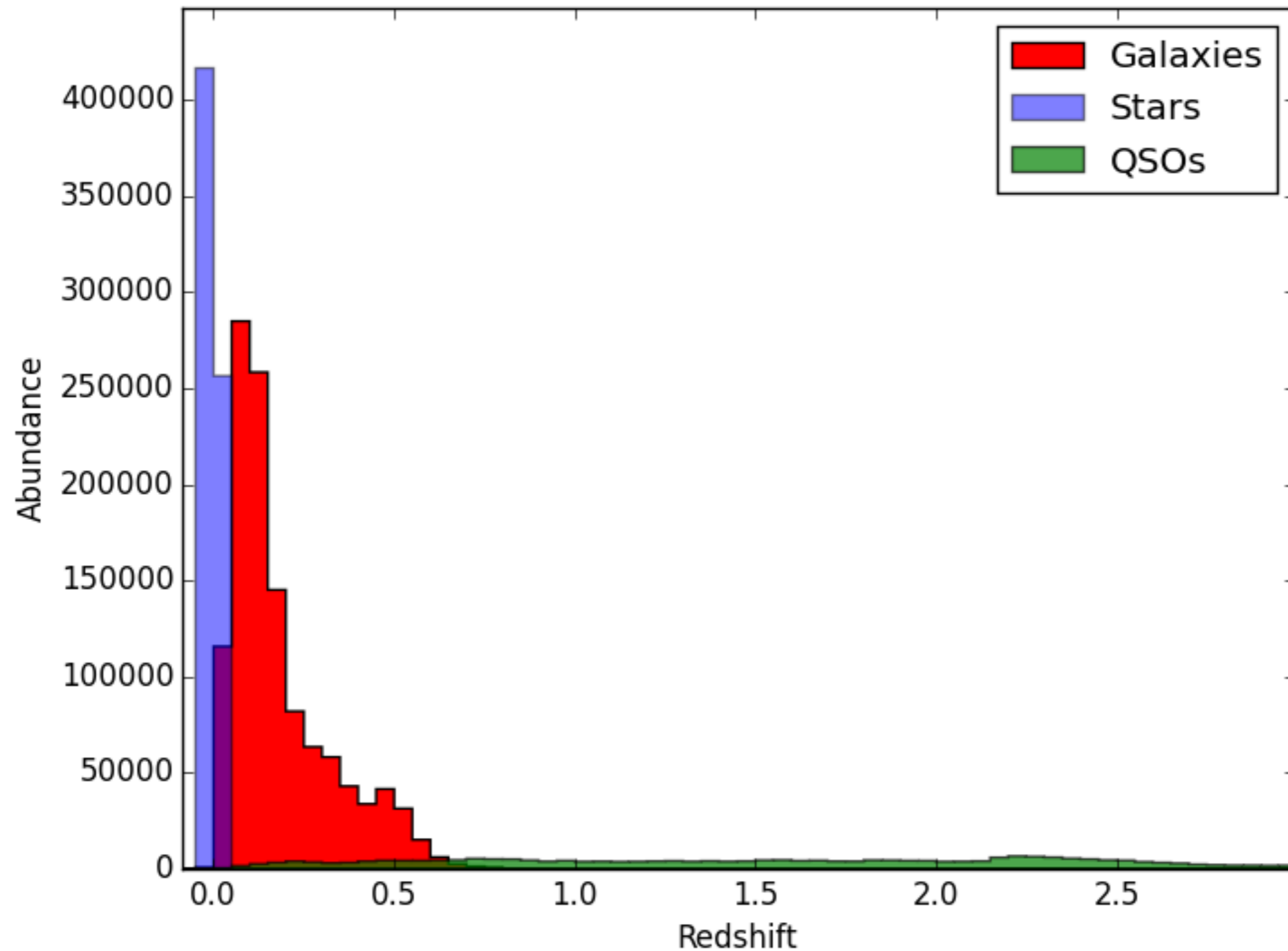


# SDSS DR12 photo-spec sample

~ 2,100,000 objects (after data cleaning)

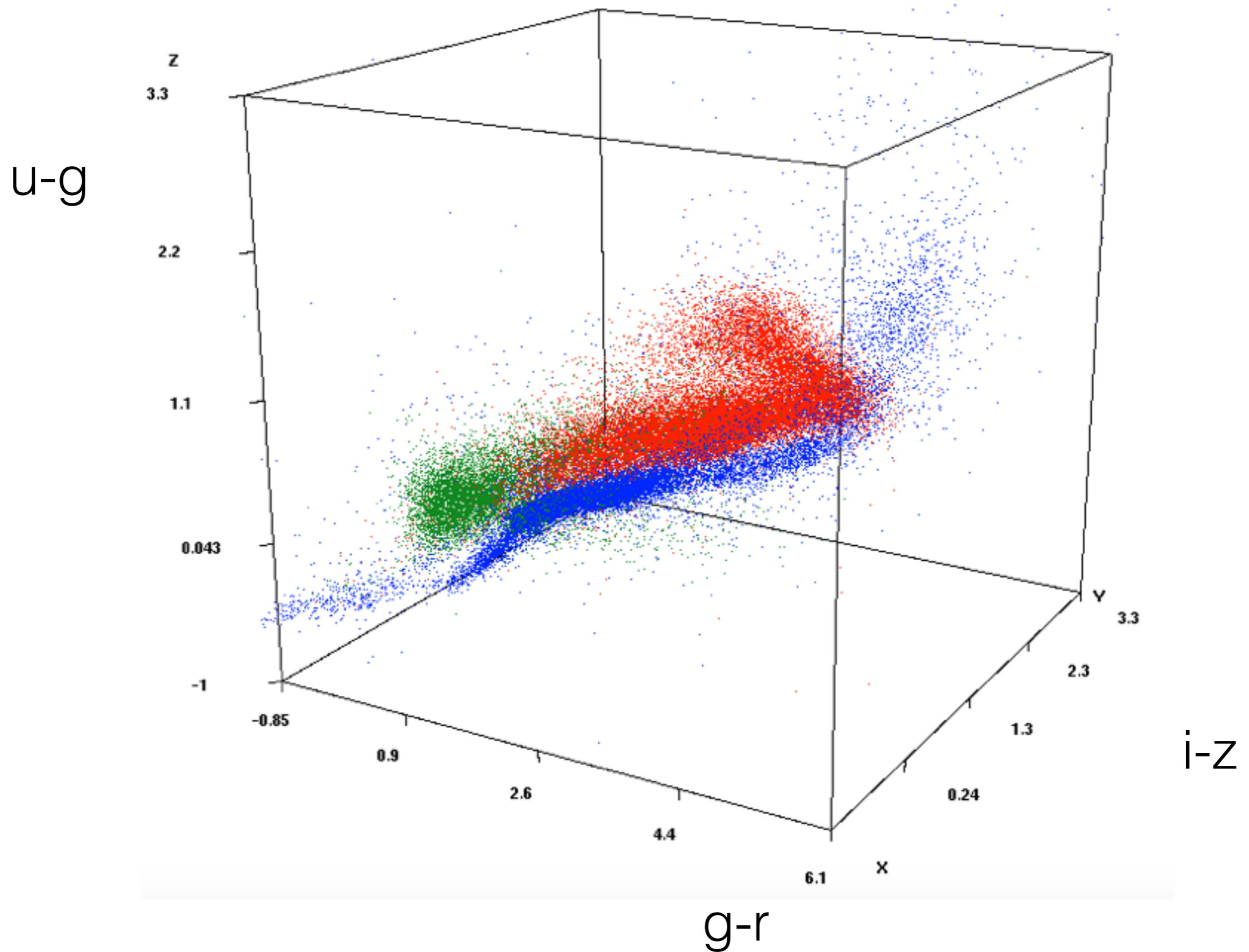


SDSS DR12 photo-spec sample  
~ 2,100,000 objects (after data cleaning)

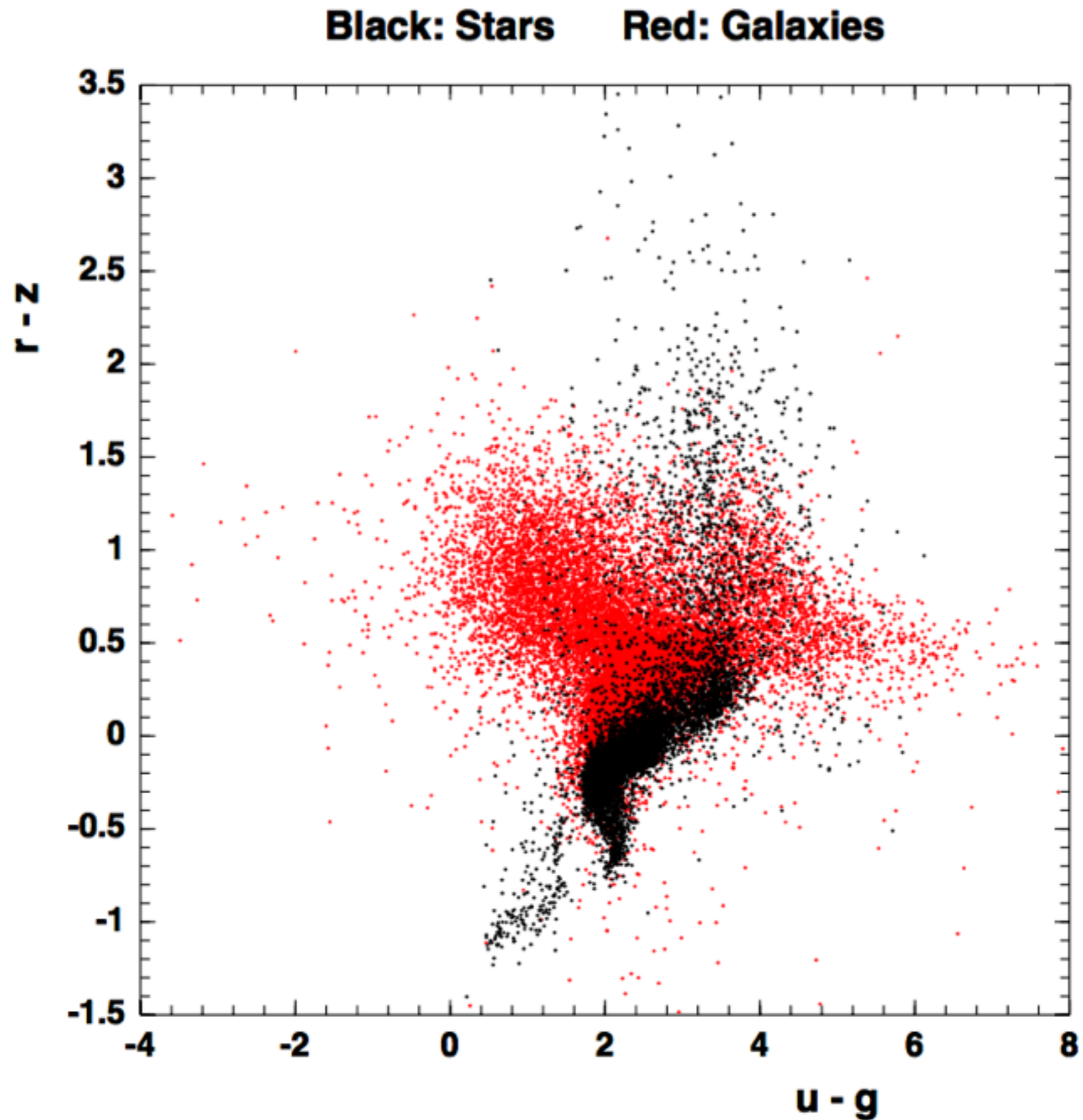




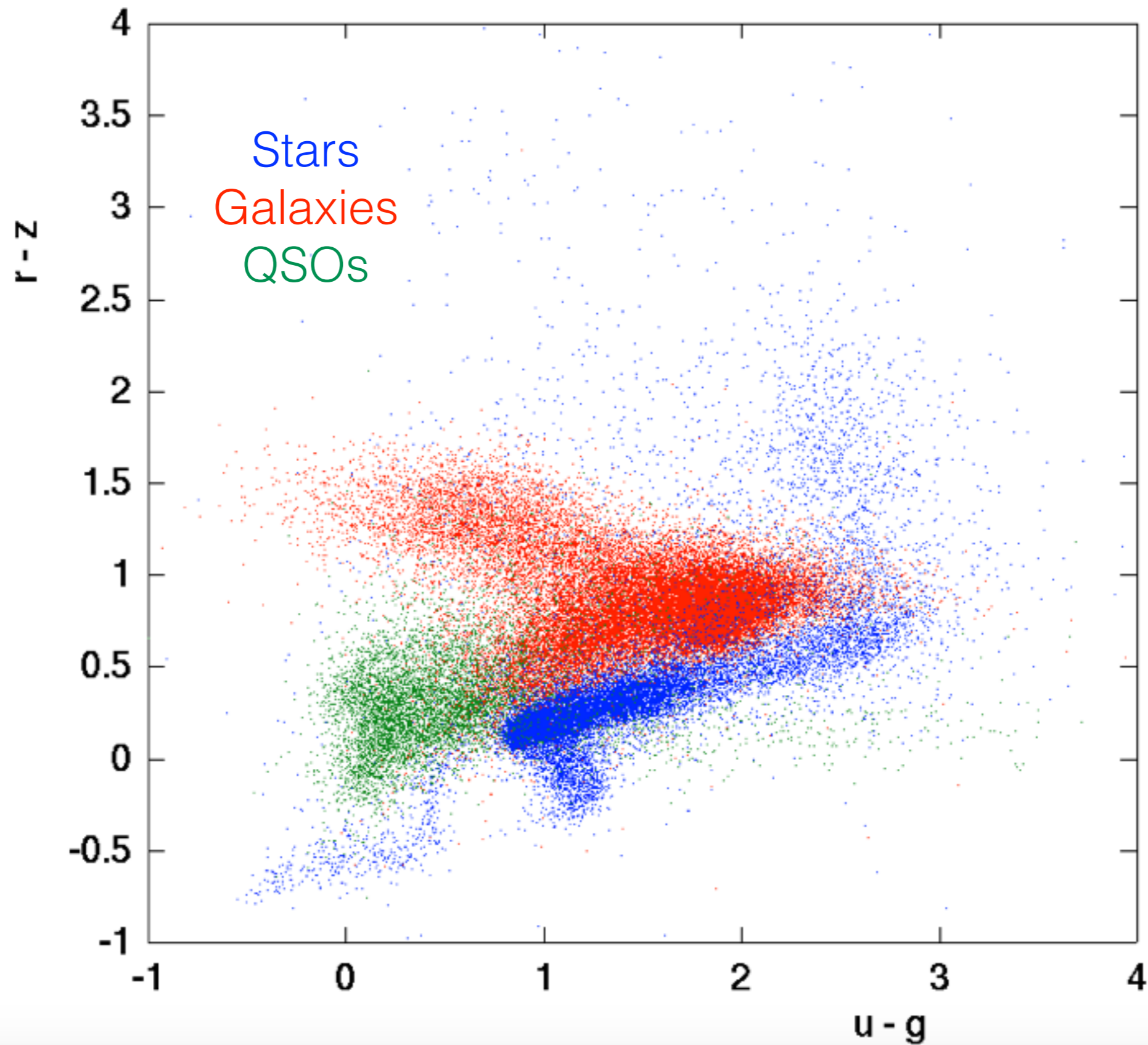
# Colour indices as “features” for classification



# Colour indices as “features” for classification



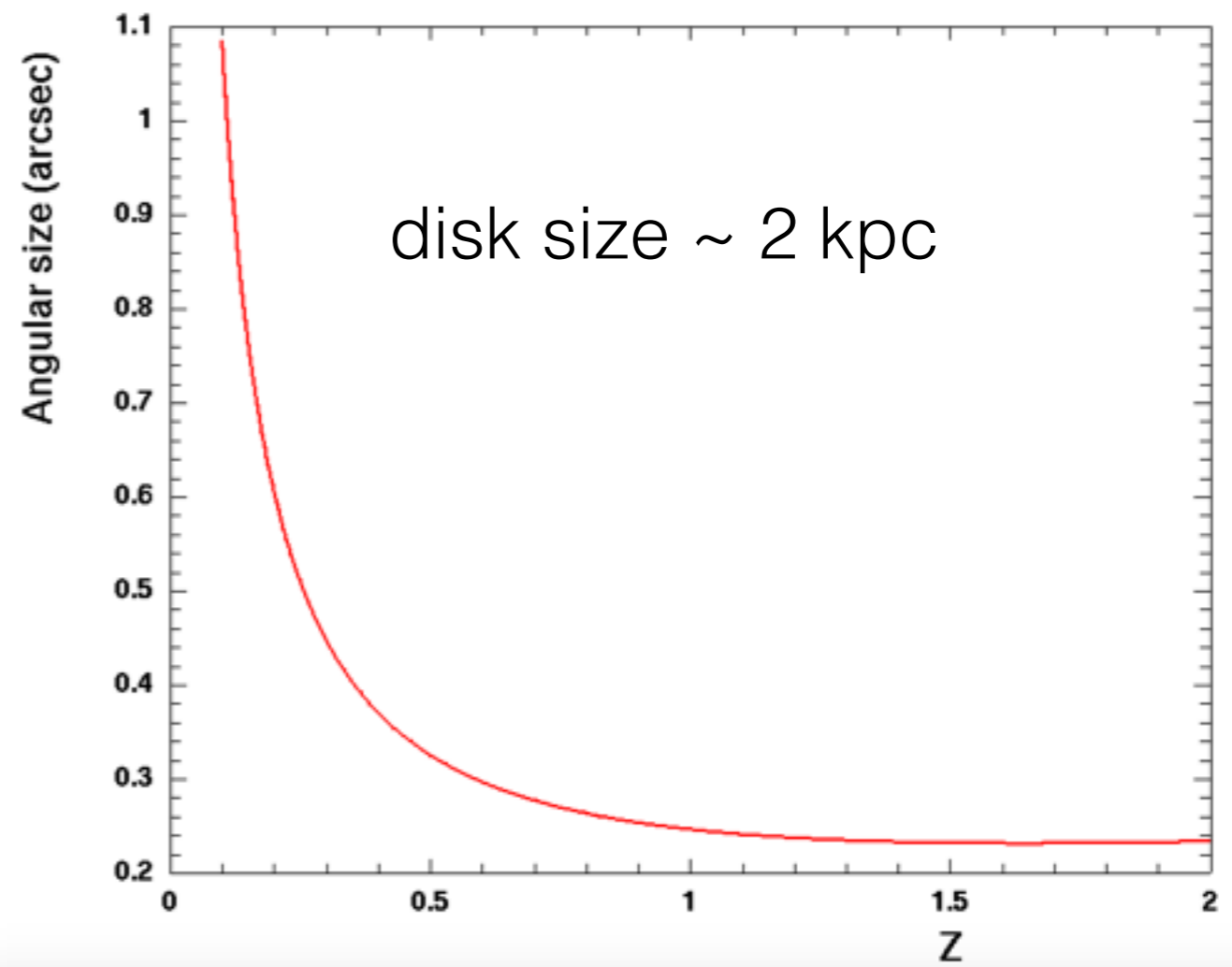
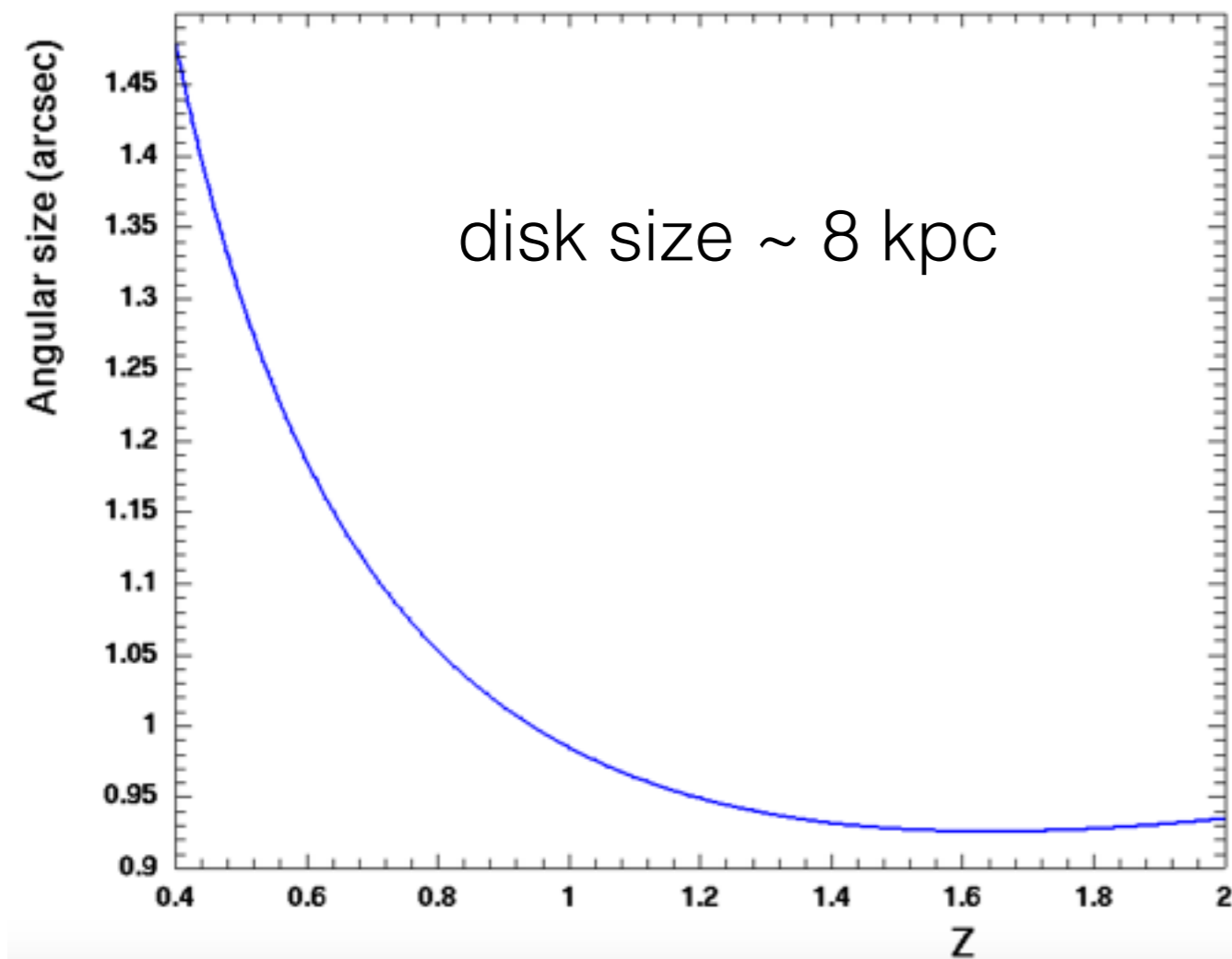
# Colour indices as “features” for classification



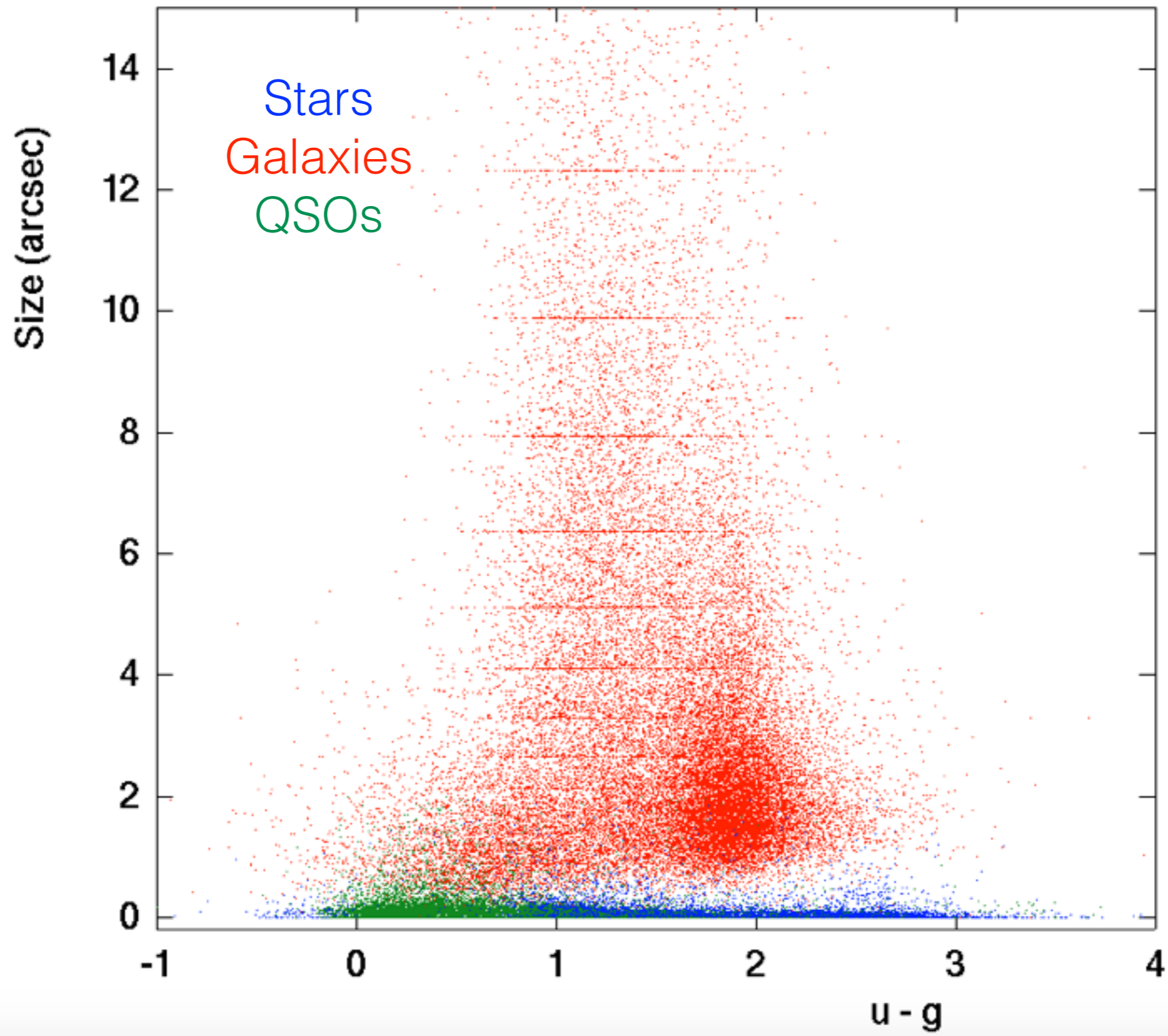


# Object's size as "feature" for classification

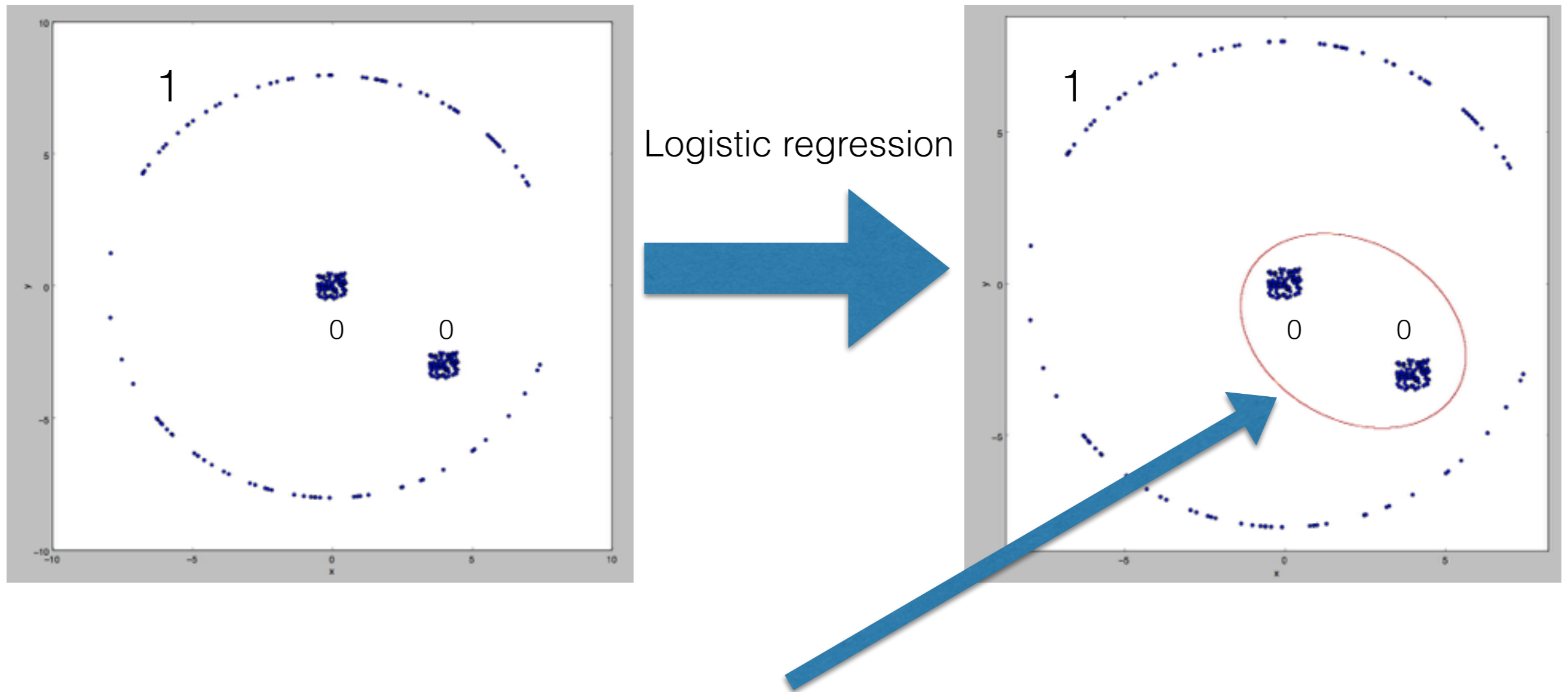
with good seeing  
MW-like galaxies can be resolved by morphology  
but not for faint galaxies (dwarfs)



# Object's size as “feature” for classification



# Supervised Classification



Parameters of the separating curve are derived by the logistic regression method.



# Logistic regression

(thanks to Andrew Ng)

## Cost function to be minimised

$$J(\theta) = -\frac{1}{m} [\sum y^{(i)} \log(h_{\theta}(x^{(i)})) + (1 - y^{(i)}) \log(1 - h_{\theta}(x^{(i)}))] + (1 - y^{(i)}) \log(1 - h_{\theta}(x^{(i)}))]$$

## Sigmoid (logistic) function

$$h_{\theta} = \frac{1}{1 + e^{-\theta^T x}}$$

$m$  : total number of objects in the training set

$i$  : object's index

$x_i$  : vector of features of an object

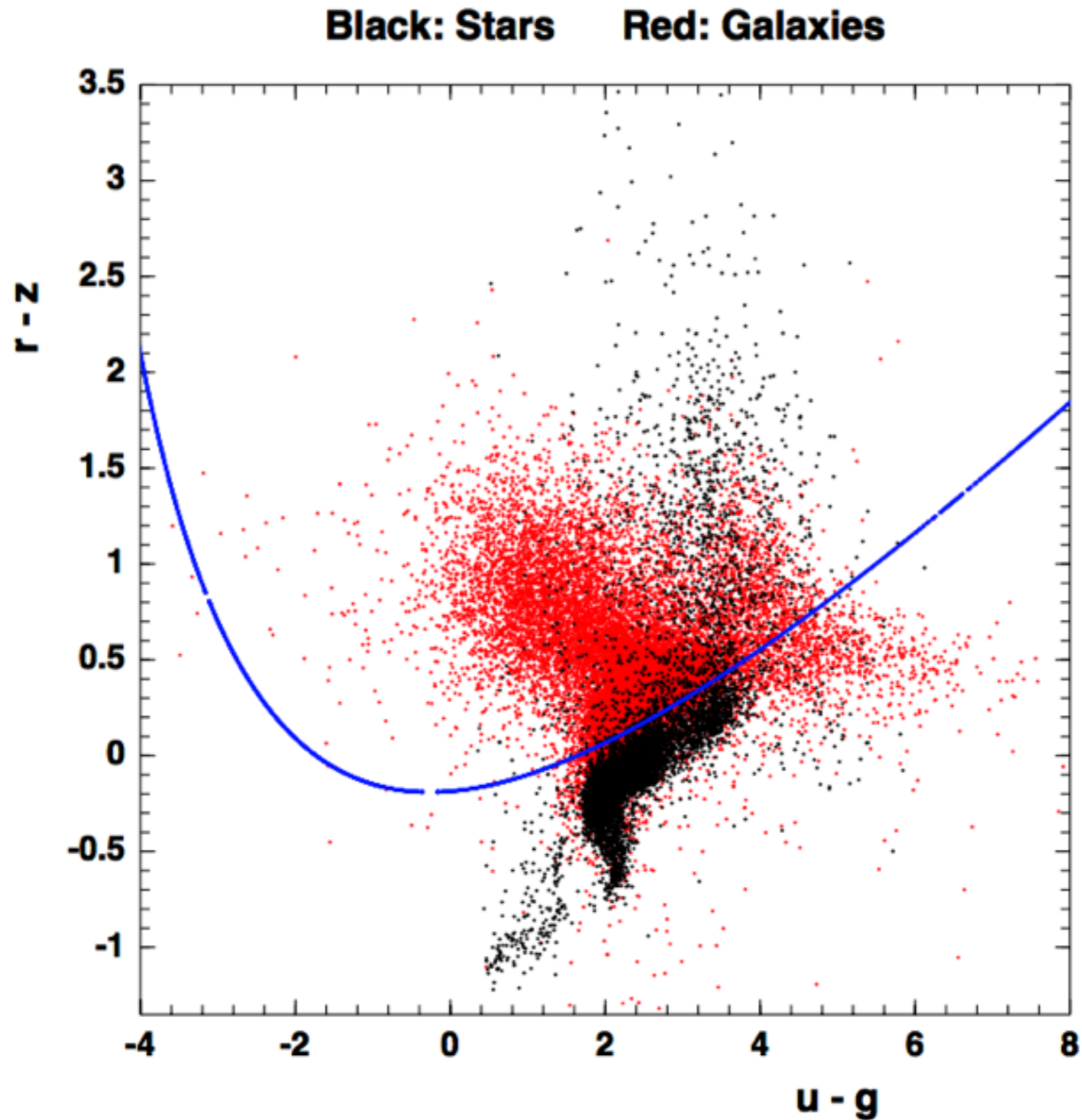
$y_i$  : object's label, 0 for stars, 1 for galaxies

$\theta$  : vector of parameters to be fitted

# Logistic regression

- We take into account size of objects and 10 colours ( $c_1=u-g$ ,  $c_2=u-r$ , ...) plus one magnitude ( $u$ ) and their quadratic function ( $c_i.c_j$ ) to have 77 features.
- The separation region is constrained by a 12 dimension hyper parabola defined by 79 parameters.
- From  $\sim 670,000$  stars,  $\sim 1,100,000$  galaxies and  $250,000$  QSOs we randomly put 20000 from each object into the training sample.

# Results from the L.R. fit

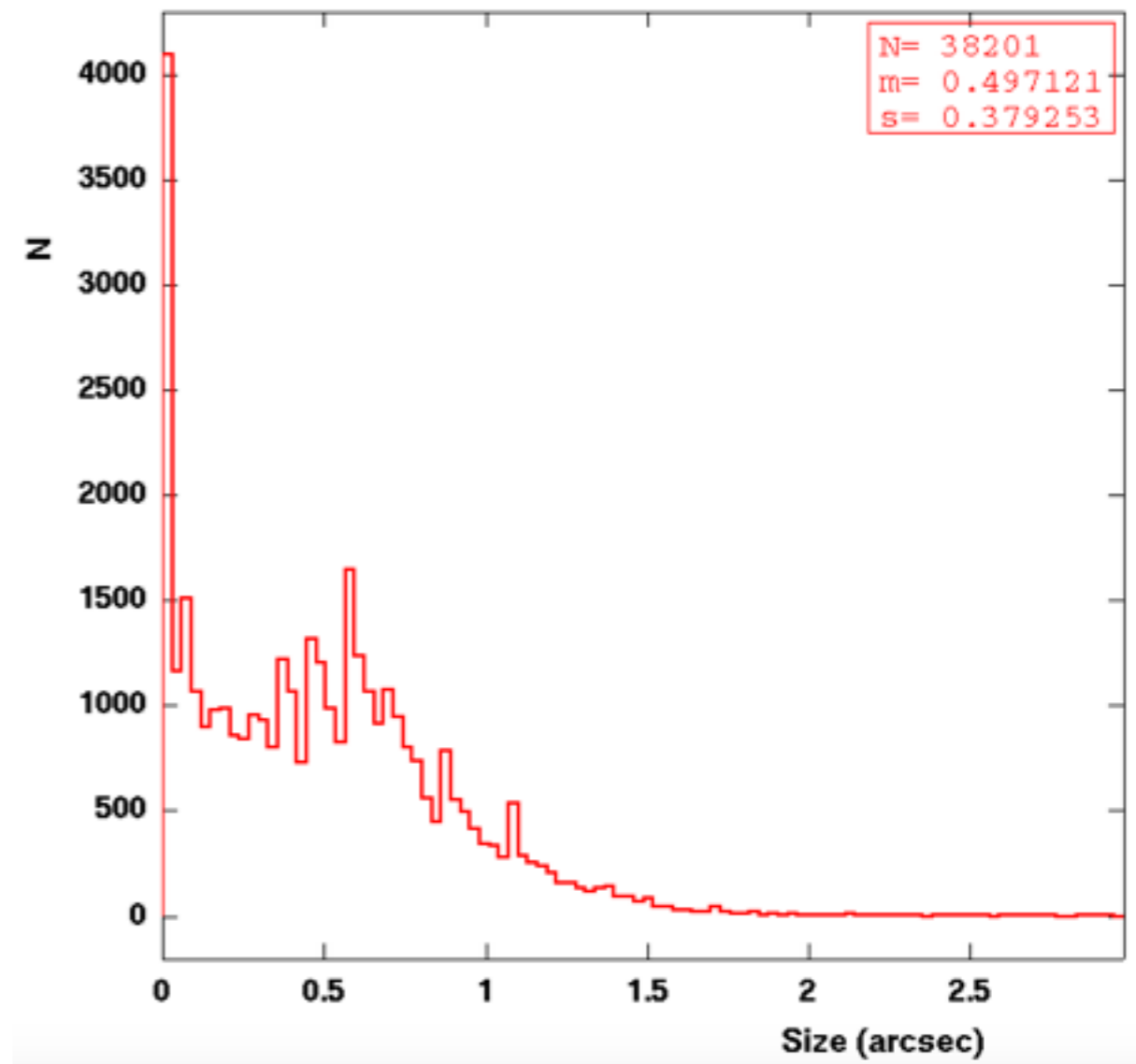
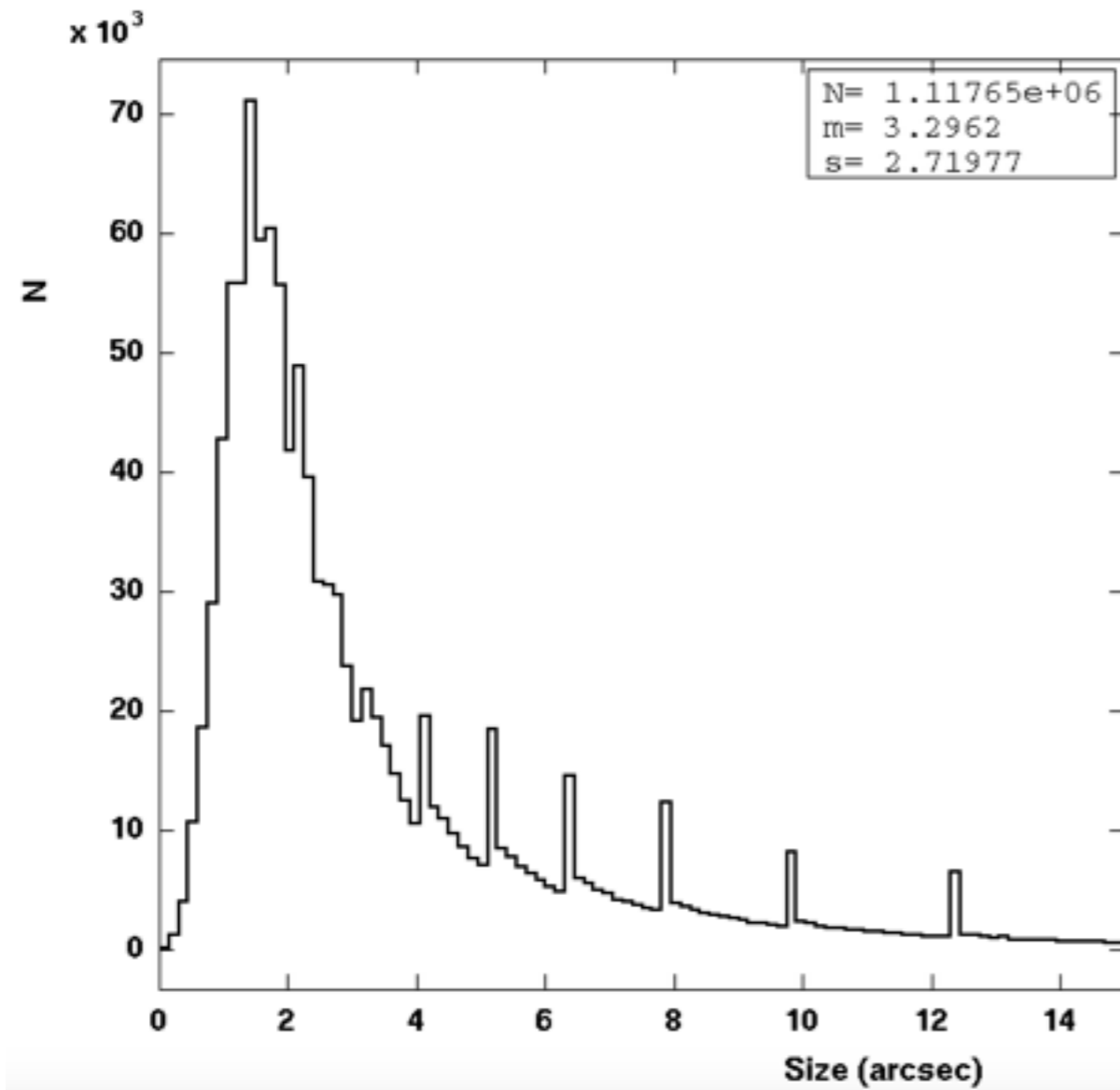




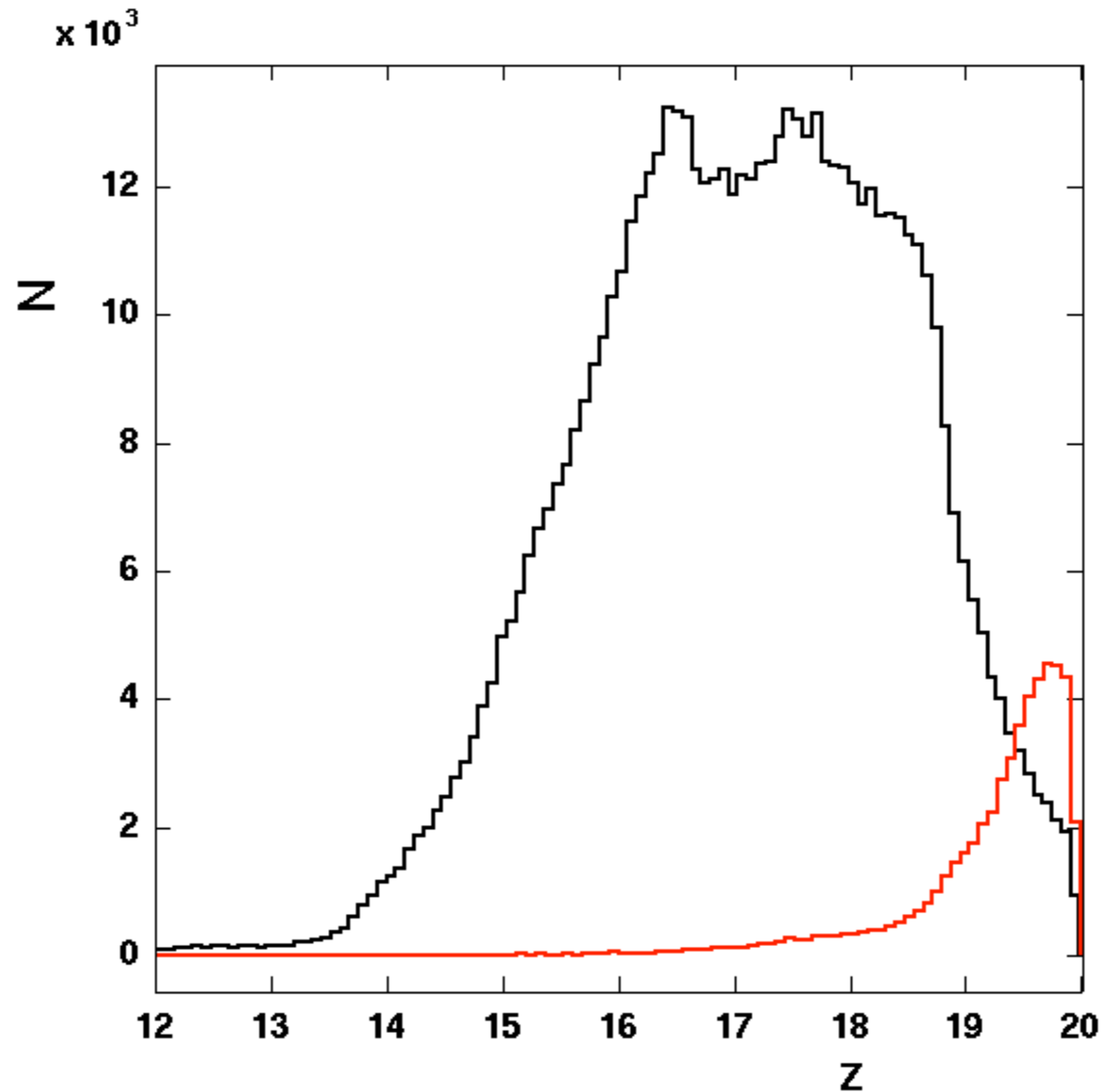
# Results from the classification

- **Classification efficiency** for the whole sample: 94%  
**galaxies**: 96%  
**stars**: 92%  
**QSOs**: 88%
- **Mean size** of the **galaxies** classified wrongly: 0.5 arcsec  
correctly: 3 arcsec
- **Mean magnitude** (extinction corrected) of the **stars** classified wrongly:  $z = 19$  (fainter stars)  
correctly:  $z = 17$
- **Mean redshift** of the **QSOs** classified wrongly: redshift = 2 (further QSOs)  
correctly: redshift = 1.5

# Wrongly and correctly classified galaxies

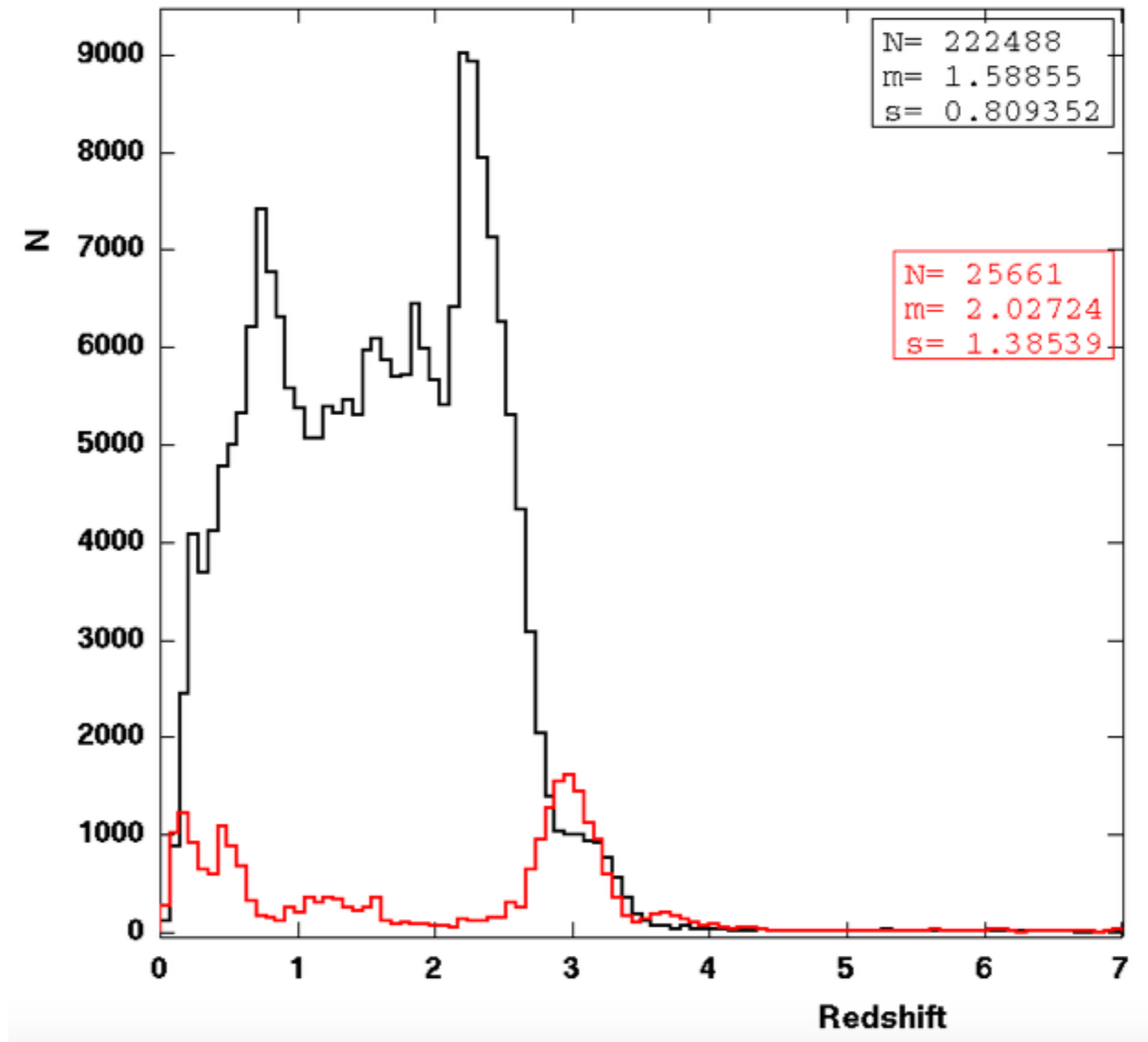


# Wrongly and correctly classified stars





# Wrongly and correctly classified QSOs



# Ongoing work

## Comparison with other classifiers

1. Random forest technique:

Whole sample efficiency ~ 95%

efficiency per object to be investigated  
(special thanks to Mehdi Cherti)

A basic classifier works nicely so far!

2. Go deeper in finding the sources  
of the misclassifications.

# Conclusions & Perspectives

- in SDSS DR12, ~ 94% of galaxies, stars and QSOs can be correctly separated using their colours and size by implementing Logistic Regression.
- 4% of galaxies (small angular size) can be mis-classified as point-like sources.
- 9% of (faint) stars can be mis-classified as galaxy-QSO.
- 12% of (further) QSOs can be mis-classified as galaxy-star.
- Classifying the simulated objects according to the LSST observation ability (higher redshifts and fainter objects).
- What is the effect of misclassified objects on photo-z determination of galaxies and cosmological parameters?