

Astronomical data Mining:

An application to the photometric redshifts of galaxies and QSOs

Giuseppe Longo

University Federico II in Napoli, longo@na.infn.it

In coll. with

M. Brescia, R. D'Abrusco, O. Laurino & the **DAME** team



Ministero degli Affari Esteri



INAF



The company which is making the journey...



University Federico II

- Massimo Brescia (project manager)
- Stefano Cavuoti
- Raffaele D'Abrusco
- Giancarlo D'Angelo (GRID)
- Natalia V. Deniskina
- Michelangelo Fiore (student)
- Mauro Garofalo
- Omar Laurino (project engineer)
- Giuseppe Longo (Principal Investigator)
- Francesco Manna (student)
- Alfonso Nocella
- Civita Vellucci (student)



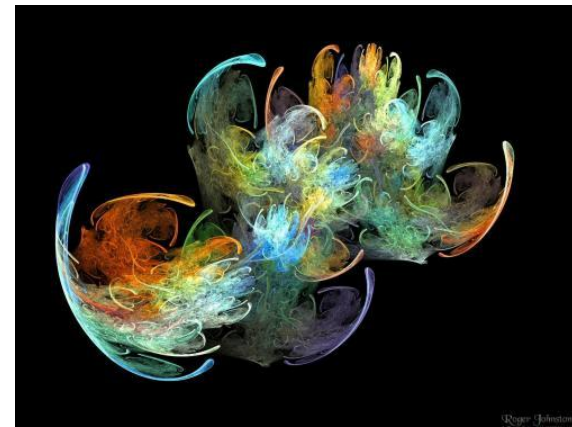
Caltech

- G.S. Djorgovski
- C. Donalek
- A. Mahabal

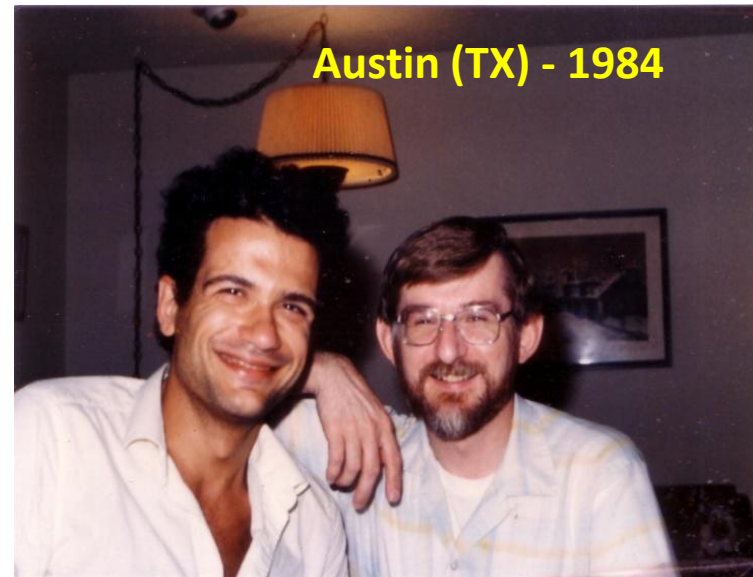
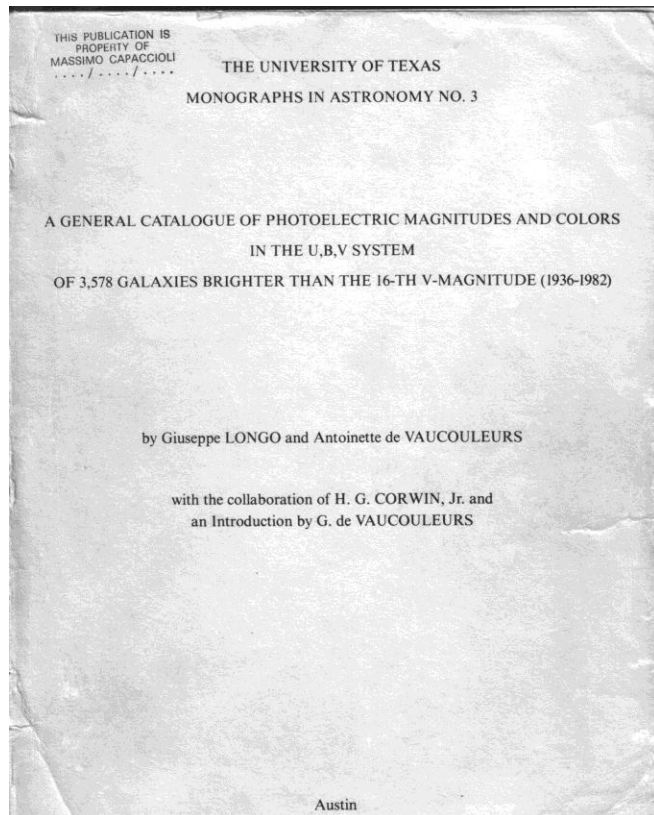


Summary of the talk

- Data Mining and astronomy
 - Why DAME and what is DAME
 - Photometric redshifts and galaxy phot-z's in DAME
 - A DM “pipeline” for QSO's (candidate selection and phot-z's)
 - Some general considerations on the future



Most of us have done it for their whole life



Compilation of photoelectric multiaperture photometry

Through standard luminosity profile curves to derive
“Extrapolation corrections”

.... in order to derive Total Magnitudes of galaxies



TABLE 2. Source contributions*

Source	n(V)	n(B)	n(B-V)	n(U-M)	Source	n(V)	n(B)	n(B-V)	n(U-M)
AAB-72	5	0	0	0	DIB-69	3	3	4	4
AAB-77	8	8	8	0	DIB-70	26	26	26	25
AAB-81	8	8	8	0	DIB-81	32	33	33	32
ABE-71	23	23	23	17	DIB-72	2	2	2	2
ALC-74	37	37	37	37	DIB-77	9	10	9	6
ALC-75	142	141	141	141	DOB-79	47	47	47	47
ANG-74	0	0	0	0	DOI-75	49	49	49	49
ANG-77	22	22	22	21	DOI-81	40	40	40	35
ANG-87	19	19	19	19	DOI-10	39	37	37	37
ARA-77	29	30	30	29	ESD-75	122	122	122	122
ARA-77	13	13	13	13	ESD-81	70	70	70	70
ARS-82	35	35	35	30	FAL-59	1	1	1	1
BAB-68	76	76	76	74	FET-67	6	6	6	3
BAS-84	61	61	61	59	FMA-82	4	4	4	4
BAS-85	290	290	290	290	FOM-76	5	5	5	4
BDE-87	1827	1827	1827	1827	FOT-72	5	5	5	4
BDE-81	7	7	7	7	FRI-75	35	36	36	35
BEG-81	1	1	1	1	FRI-82	3	3	3	3
BFW-78	63	63	64	64	FRI-87	591	591	591	591
BGC-3	0	0	0	0	GAB-83	2	2	2	2
BGC-33	44	44	44	44	GAB-77	0	0	0	0
BGC-54	72	72	72	72	GAI-85	38	37	37	37
BGC-55	16	16	16	0	GAI-87	37	30	30	30
BEG-64	109	107	107	108	GAI-79	2	2	2	2
BHP-80	255	212	212	0	GAV-82	101	99	99	101
BLA-84	14	17	14	0	GAW-89	0	0	0	0
BOR-81	29	29	29	29	GFB-77	115	108	107	0
BOT-82	200	200	200	200	GFS-84	174	174	174	174
BFE-78	89	89	89	78	GFB-80	147	7	7	139
BIM-85	1	1	1	1	GIB-82	141	139	141	139
BRS-85	35	55	55	55	GIB-85	3	3	3	0
BUR-84	101	101	101	101	GIB-86	1	1	1	0
BUT-81	0	0	0	0	GIB-74	1	1	1	0
BUT-82	64	67	67	65	GIB-86	169	169	169	169
BUT-83	18	14	14	0	GIB-78	74	74	74	0
BUT-84	58	58	58	58	GIB-82	718	569	569	569
BUT-84	202	202	202	202	GIB-82	75	41	41	41
CAL-83	31	31	31	31	GIB-85	1	1	1	1
CAL-83	91	91	91	0	GIB-85	4	4	4	6
CAW-97	65	57	57	55	GIB-77	4	4	4	4
CEW-75	13	12	12	12	GIB-87	0	0	0	22
CGR-67	83	83	83	81	GIB-84	81	81	81	81
CGR-80	161	164	164	164	GIB-88	97	96	96	67
CGR-84	171	171	171	171	GIB-87	2	2	2	1
CGR-87	222	222	222	222	GIB-88	1	1	1	2
CFB-81	16	16	16	0	GIB-75	5	2	2	0
CFB-84	396	396	396	0	GIB-82	1	1	1	0
CFA-81	15	15	15	15	GIB-88	287	291	287	0
DAN-78	4	4	4	4	GIB-77	275	275	275	275
DAN-84	827	825	825	816	GIB-79	5	5	5	5
DDP-84	49	46	46	46	GIB-80	4	4	4	4
DDP-78	2	2	2	2	HIF-81	34	34	34	30
DDP-89	2	2	2	2	HIF-82	1	1	1	0
DIB-68	36	36	36	35	JET-82	2	2	2	1

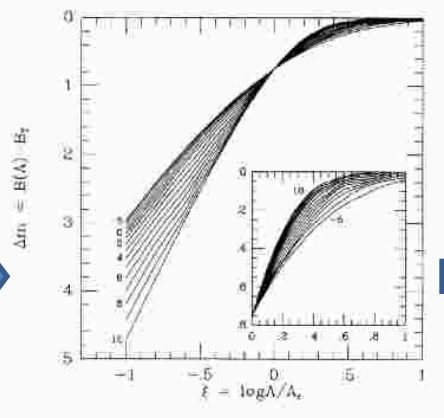
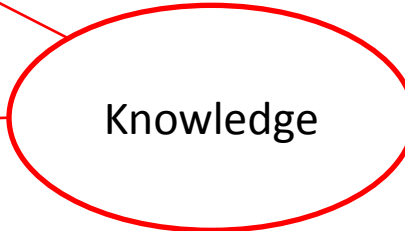
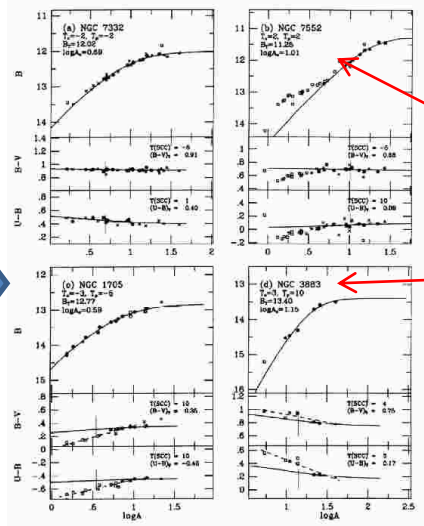


FIG. 1. Revised standard growth curves for B magnitudes. Several curves are labeled with their associated numerical Hubble type index. The upper parts of the curves are shown on an expanded scale in the inset at lower right.



Data

Base of Knowledge (BoK)

Model

Data Mining is not only new astronomy.

In many cases (**but NOT ALL**) it is just the name we give to rather usual stuff when it needs to be performed fast and on billions of records of COMPLEX data



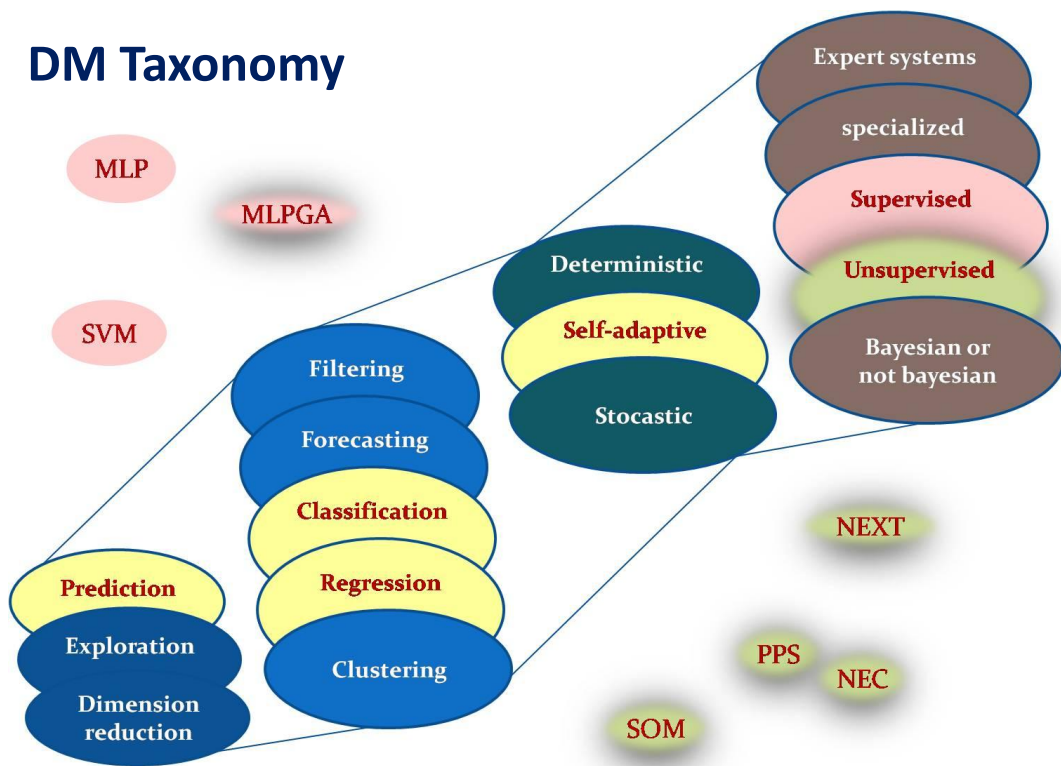
Human brain is not sufficient



Machine learning methods

Data Mining is the activity of extracting **USEFUL** information from **COMPLEX** data using Statistical Pattern Recognition and Machine Learning methods.

DM Taxonomy



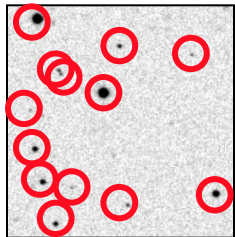
1. To catalogue the known (classification)
2. Characterize the unknown (clustering)
3. Find functional dependencies (regression)
4. Find exceptions (outliers)

Supervised Methods

Patterns are learnt from extensive set of templates (Base of Knowledge = BoK)

Unsupervised Methods

Patterns are discovered using the data themselves



$p = \{\text{isophotal, petrosian, aperture magnitudes, concentration indexes, shape parameters, etc.}\}$

$$p^1 = \{RA^1, \delta^1, t, \{\lambda_1, \Delta\lambda_1, f_1^{1,1}, \Delta f_1^{1,1}, \dots, f_1^{1,m}, \Delta f_1^{1,m}\}, \dots, \{\lambda_n, \Delta\lambda_n, f_n^{1,1}, \Delta f_n^{1,1}, \dots, f_n^{1,m}, \Delta f_n^{1,m}\}\}$$

$$p^2 = \{RA^2, \delta^2, t, \{\lambda_1, \Delta\lambda_1, f_1^{2,1}, \Delta f_1^{2,1}, \dots, f_1^{2,m}, \Delta f_1^{2,m}\}, \dots, \{\lambda_n, \Delta\lambda_n, f_n^{2,1}, \Delta f_n^{2,1}, \dots, f_n^{2,m}, \Delta f_n^{2,m}\}\}$$

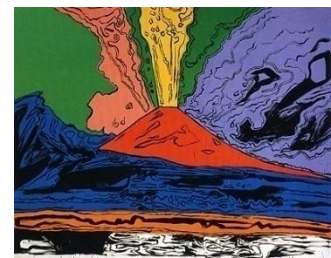
.....

$$p^N = \{RA^N, \delta^N, t, \{\lambda_1, \Delta\lambda_1, f_1^{N,1}, \Delta f_1^{N,1}, \dots, f_1^{N,m}, \Delta f_1^{N,m}\}, \dots\}$$

$$D = 3 + m \times n$$

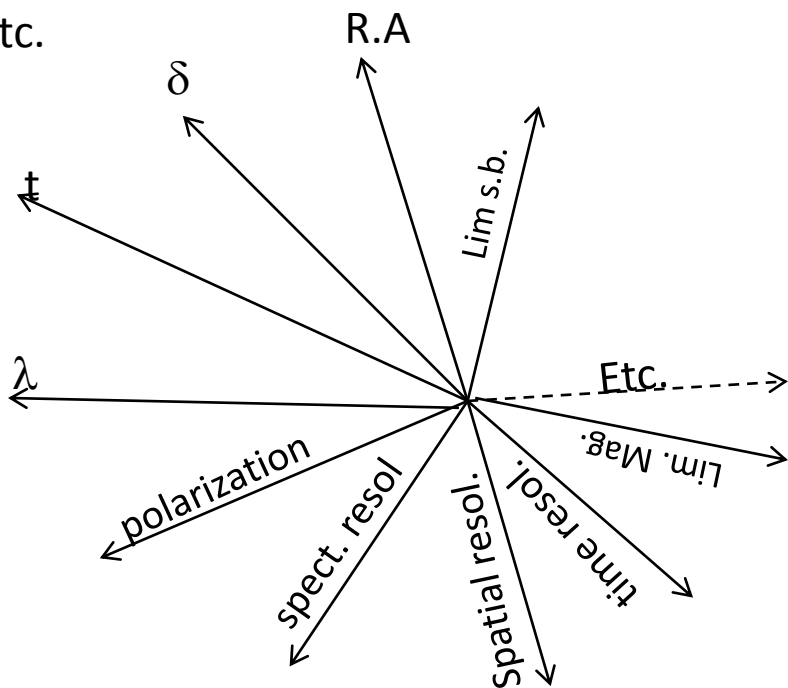
The scientific exploitation of a multi band, multiepoch (K epochs) universe implies to search for patterns, trends, etc. among **N** points in a **DxK** dimensional parameter space:

$N > 10^9, D \gg 100, K > 10$



Any observed (simulated) datum p defines a point (region) in a subset of \mathbb{R}^N . Es:

- RA and dec
- time
- λ
- experimental setup (spatial and spectral resolution, limiting mag, limiting surface brightness, etc.) parameters
- fluxes
- polarization
- Etc.



$$p \in \mathfrak{R}^N \quad N \gg 100$$

The parameter space concept is crucial to:

1. Guide the quest for new discoveries (observations can be guided to explore poorly known regions), ...
2. Find new physical laws (patterns)
3. Etc,



The computational cost of DM:

N = no. of data vectors, D = no. of data dimensions

K = no. of clusters chosen, K_{\max} = max no. of clusters tried

I = no. of iterations, M = no. of Monte Carlo trials/partitions

K-means: $K \times N \times I \times D$

Expectation Maximisation: $K \times N \times I \times D^2$

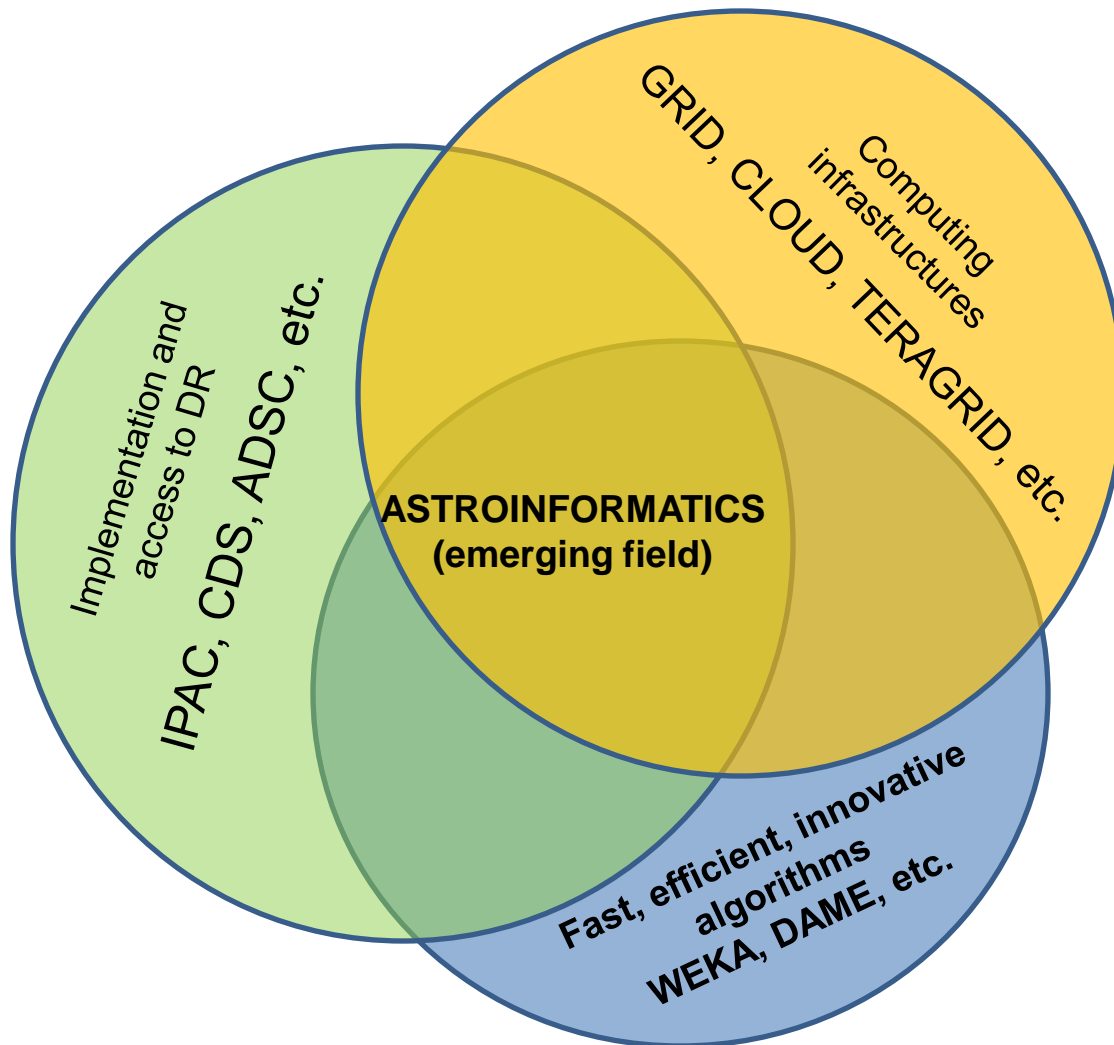
Monte Carlo Cross-Validation: $M \times K_{\max}^2 \times N \times I \times D^2$

Correlations $\sim N \log N$ or N^2 , $\sim D^k$ ($k \geq 1$)

Likelihood, Bayesian $\sim N^m$ ($m \geq 3$), $\sim D^k$ ($k \geq 1$)

SVM $> \sim (N \times D)^3$





Machine Learning problems as “function approximation”

$\mathbf{X} \equiv \{x_1, x_2, x_3, \dots, x_N\}$ input vectors

$\mathbf{Y} \equiv \{x_1, x_2, x_3, \dots, x_M\}$ target vectors $M \ll N$

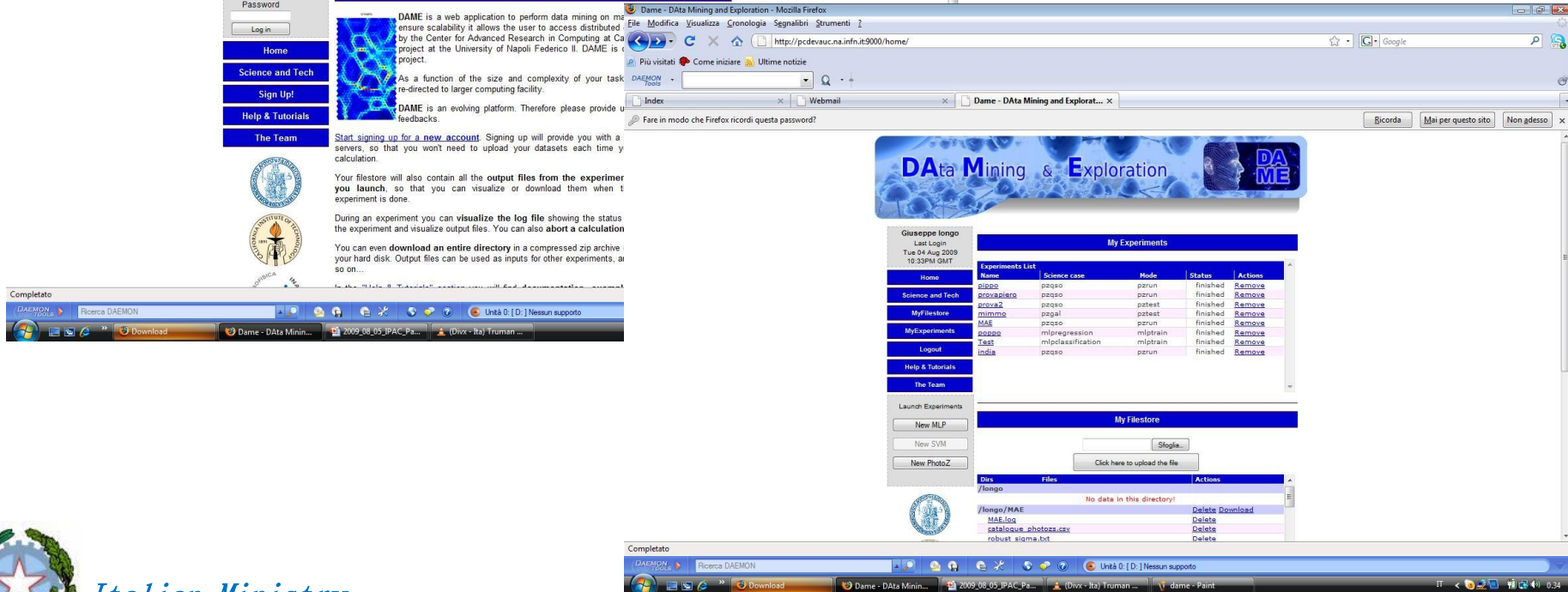
find \hat{f} : $\hat{\mathbf{Y}} = \hat{f}(\mathbf{X})$ is a good approximation of \mathbf{Y}

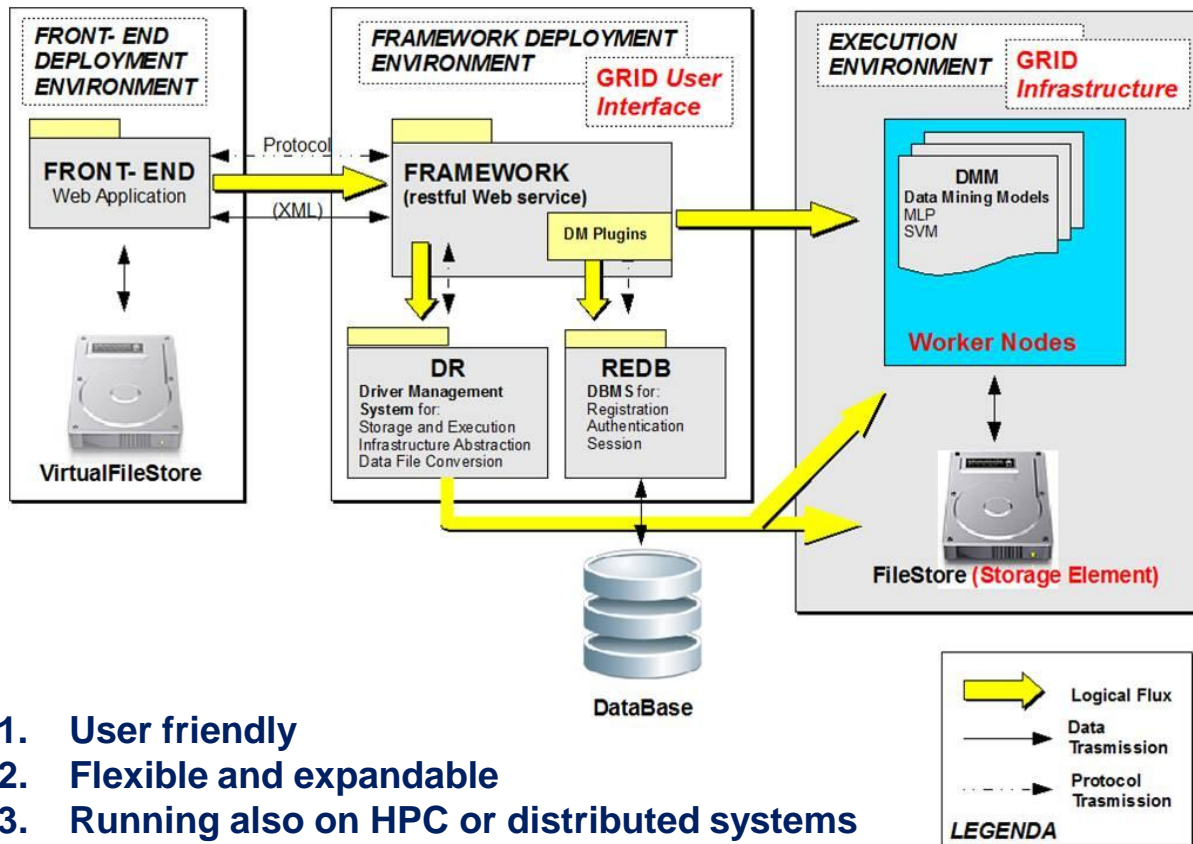
variable	characteristics	Type	Operation/example
Quantitative	Numerical with ordering relationship and possibility to define a metric	Actual measurement	Regression Photometric redshifts
Categorical (non ordered)	Membership into a finite number of classes. No ordering relationship.	Numerical codes (targets) arbitrarily ordered	Classification Search for peculiar objects, QSO's, Star/galaxy, etc.
Ordered categorical	Classes ordered by a relationship but there is no metric	Numerical codes non arbitrarily ordered	Classification Morphological and physical classification of galaxies, etc.



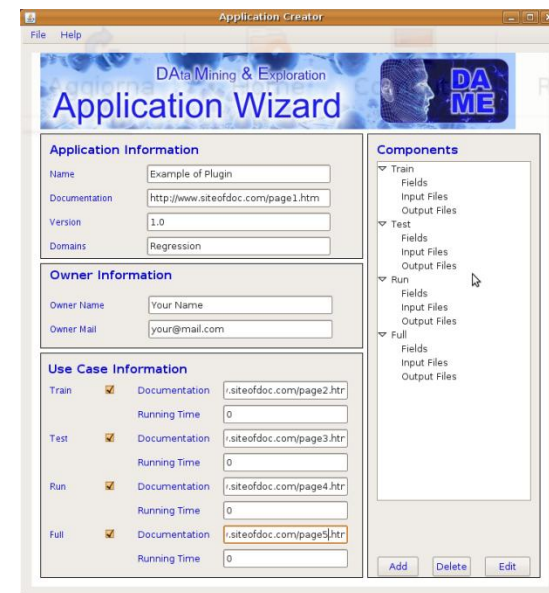
DAME is a joint effort between University Federico II, INAF and Caltech aimed at: implementing (as web application) a suite of data exploration, data mining and data visualization tools.

<http://dame.na.infn.it/>
Web application PROTOTYPE
<http://voneural.na.infn.it/>
Documents





1. User friendly
2. Flexible and expandable
3. Running also on HPC or distributed systems



Will substitute the prototype at the end of October 2009

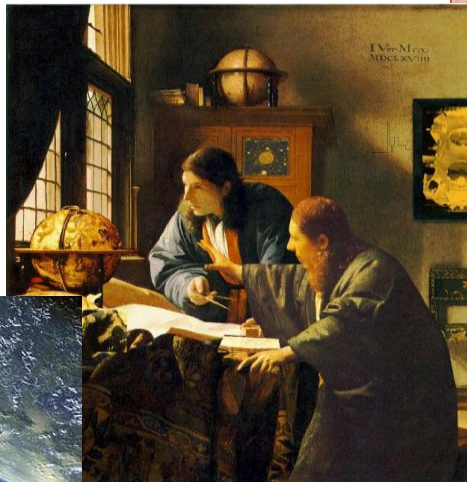
PART II - applications of DAME

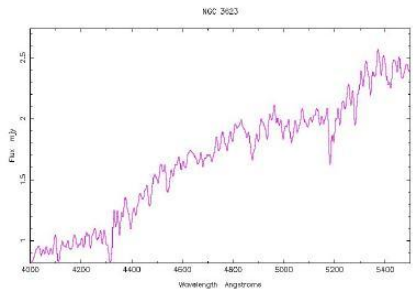
to observational cosmology

Photometric redshifts of galaxies and QSO's

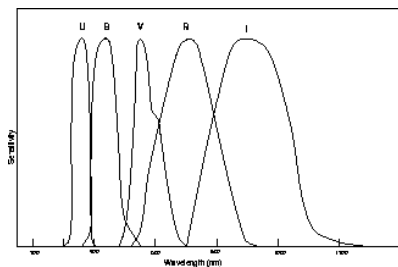
Selection of candidate quasars

- D'Abrusco et al. 2007, ApJ, 663, pp. 752-764
- D'Abrusco et al. 2009, MNRAS, 396, 223-262
- Laurino et al., 2009, Thesis
- Laurino et al., 2009, MNRAS, in preparation





X



=

$$m_U = -2.5 \log_{10} \frac{\int F(\lambda) S_U(\lambda) d\lambda}{\int S_U(\lambda) d\lambda} + c_U$$

$$m_B = -2.5 \log_{10} \frac{\int F(\lambda) S_B(\lambda) d\lambda}{\int S_B(\lambda) d\lambda} + c_B$$

Etc...

Galaxy spectrum - $F(\lambda)$

Photometric system - $S_i(\lambda)$

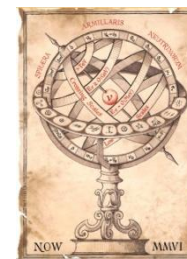
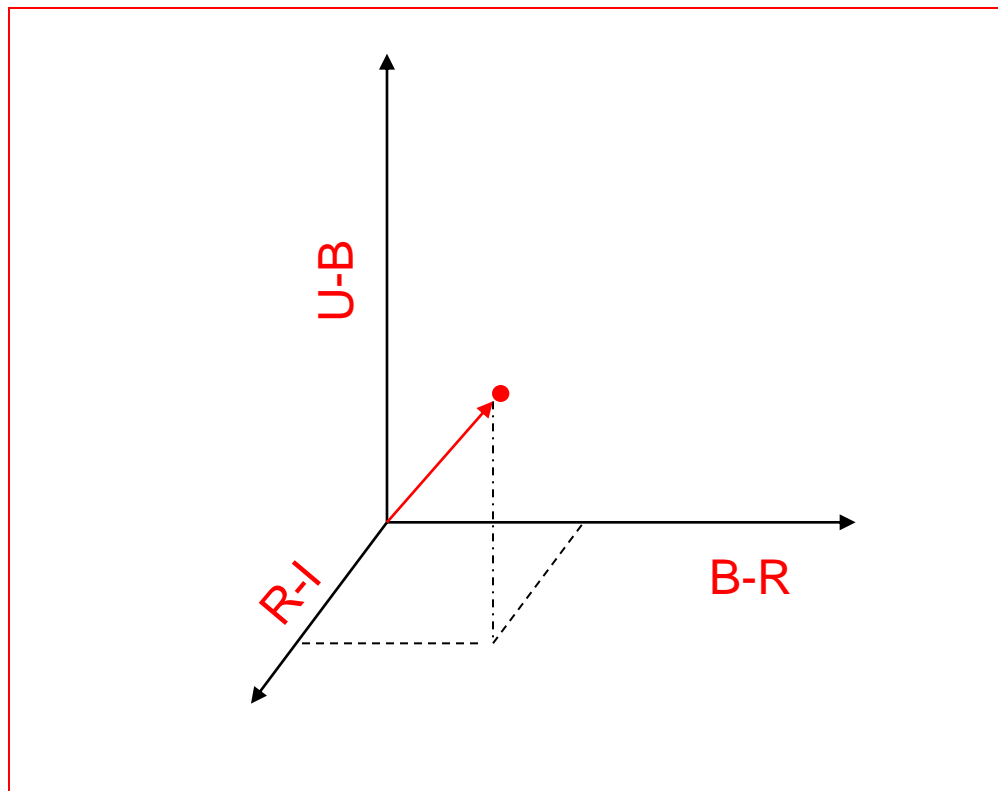


Color indexes

$$U - B \equiv m_U - m_B$$

$$B - R \equiv m_B - m_R$$

etc.



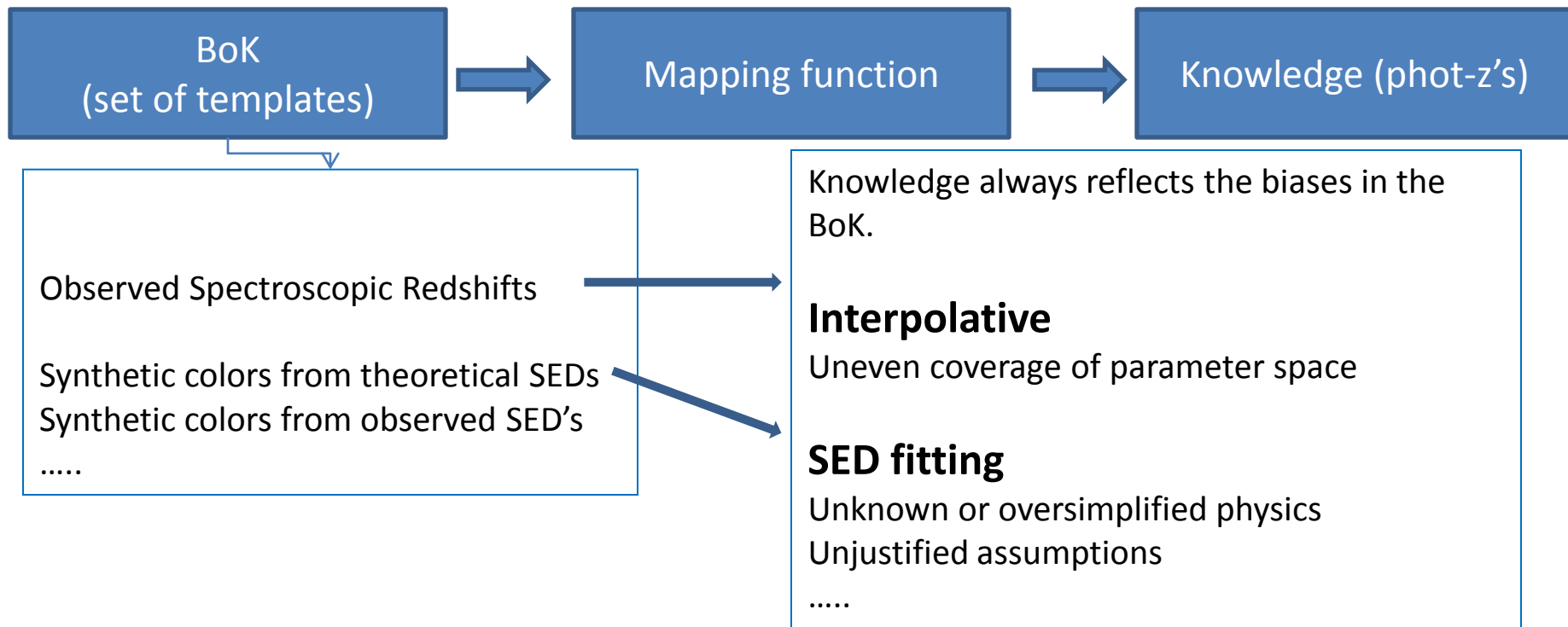


Photometric redshifts are always a function approximation hence a DM problem:

$\mathbf{X} \equiv \{x_1, x_2, x_3, \dots, x_N\}$ **input vectors**

$\mathbf{Y} \equiv \{x_1, x_2, x_3, \dots, x_M\}$ **target vectors** $M \ll N$

find $\hat{f}: \hat{\mathbf{Y}} = \hat{f}(\mathbf{X})$ **is a good approximation of** \mathbf{Y}



Data used in the science cases:

SDSS: 10^8 galaxies in 5 bands;
 BoK: spectroscopic redshifts for 10^6 galaxies
 BoK: incomplete and **biased**.

UKIDDS: overlap with SDSS

GALEX: overlap with SDSS

SDSS

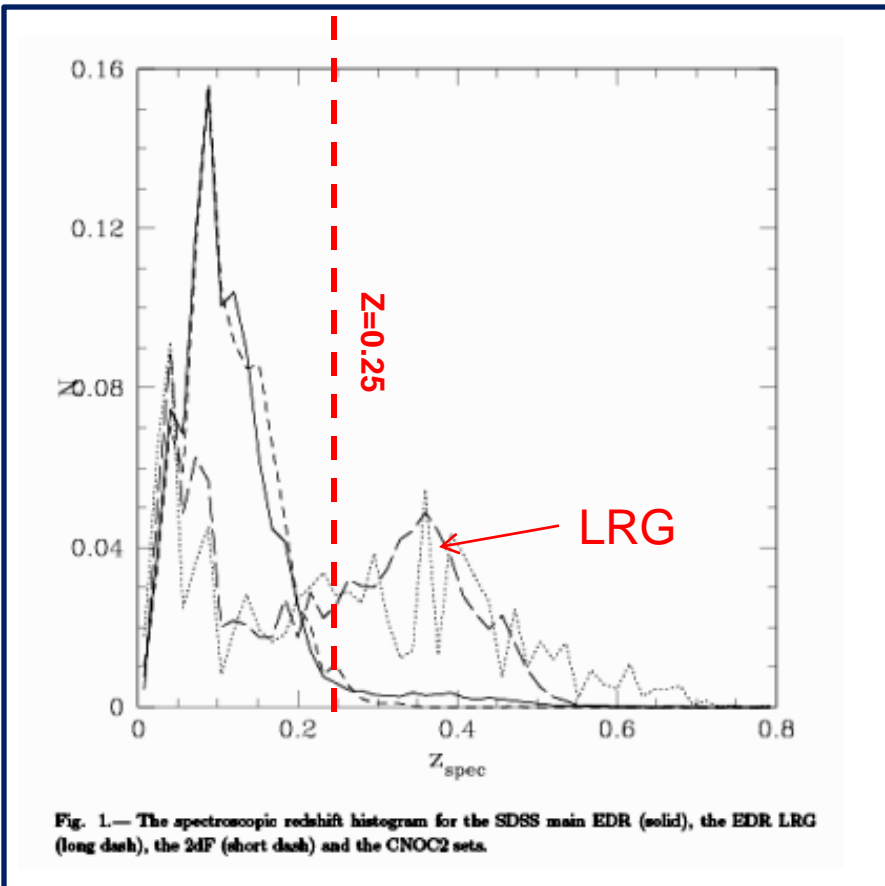
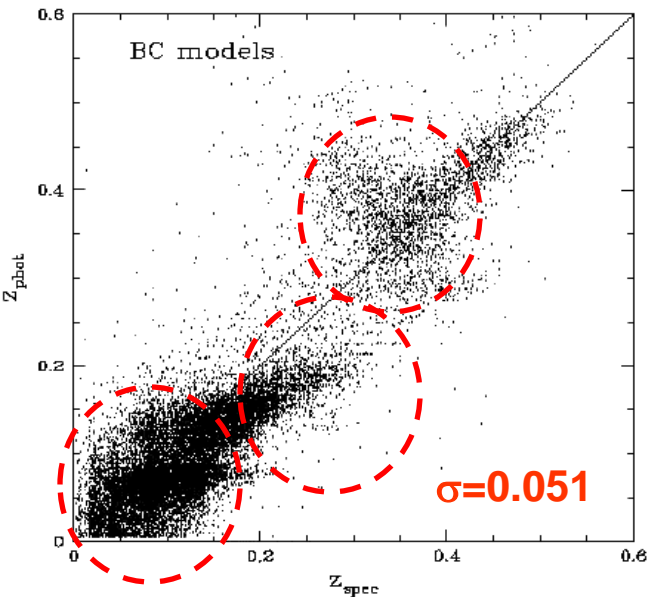
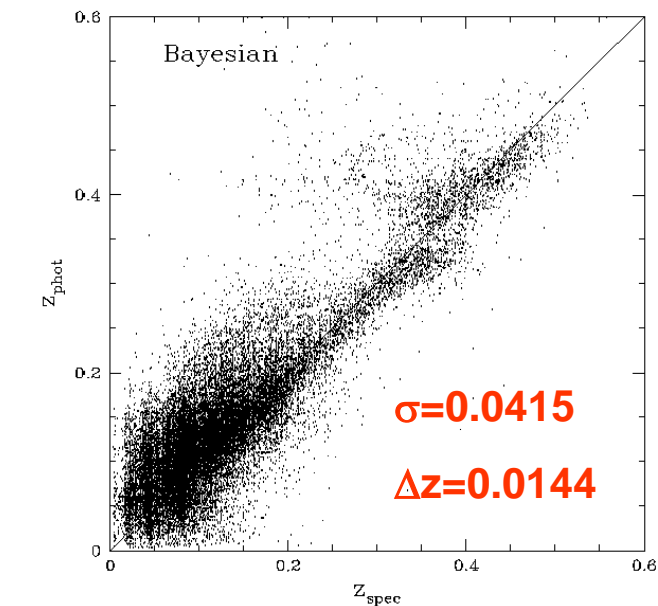


Fig. 1.— The spectroscopic redshift histogram for the SDSS main EDR (solid), the EDR LRG (long dash), the 2dF (short dash) and the CNOC2 sets.



SED fitting

Templates from synthetic colors obtained from theoretical SED's
 Mapping function from simple interpolation



Interpolative

Templates from synthetic colors obtained from theoretical SED's
 Mapping function from Bayesian inference

- the color space is partitioned (KD-tree - a binary search tree) into cells containing the same number of objects from the training set
- In each cell a second order polynomial is fit to BoK.

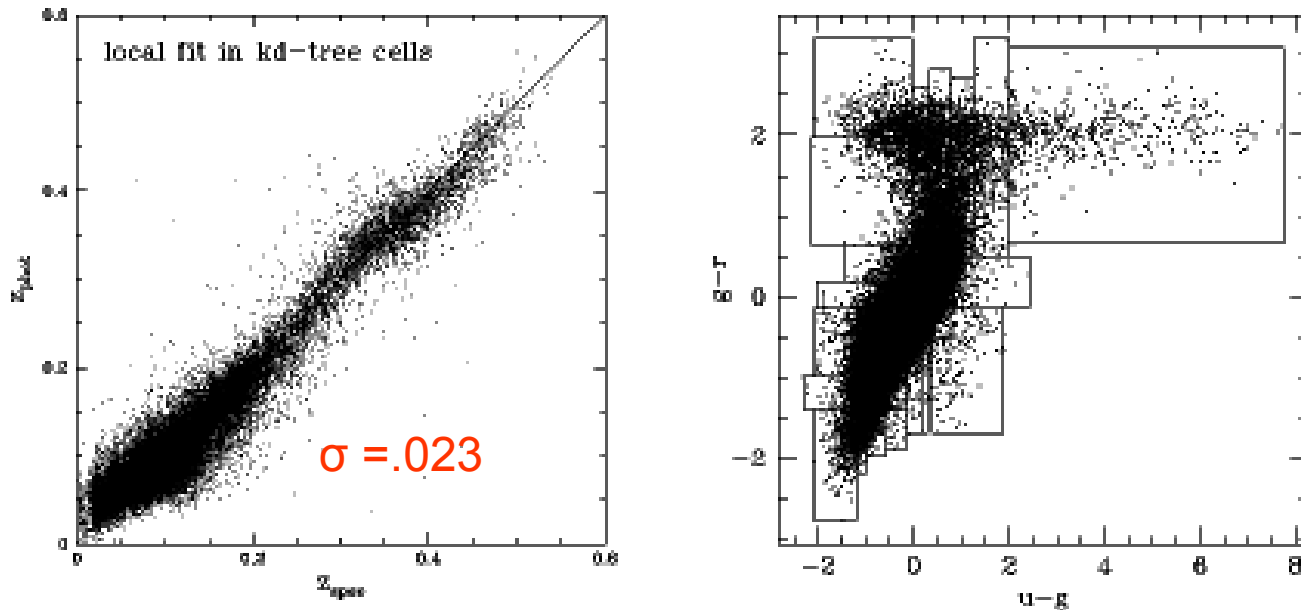
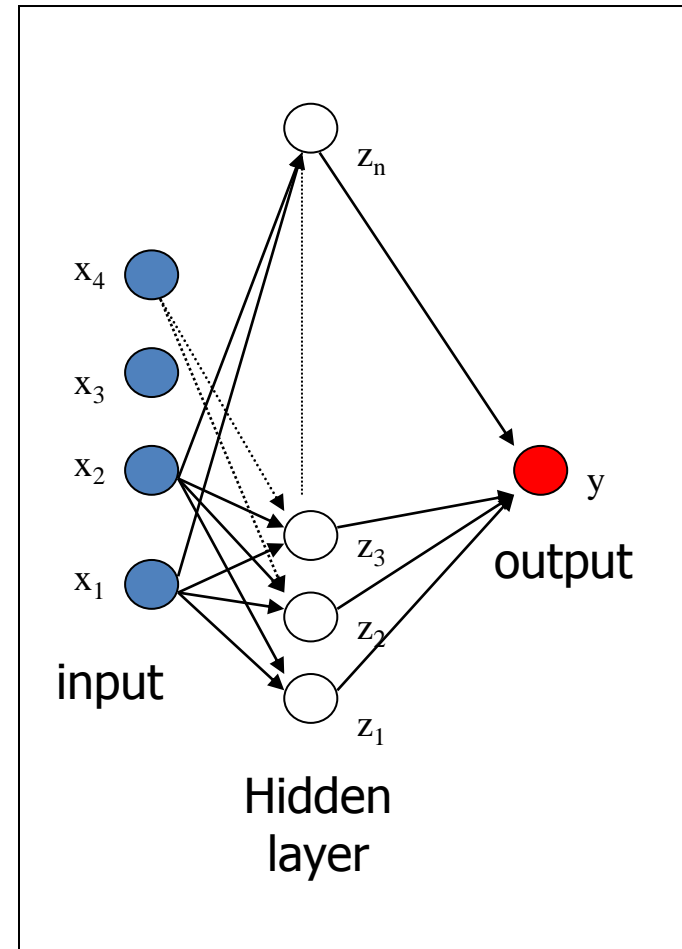
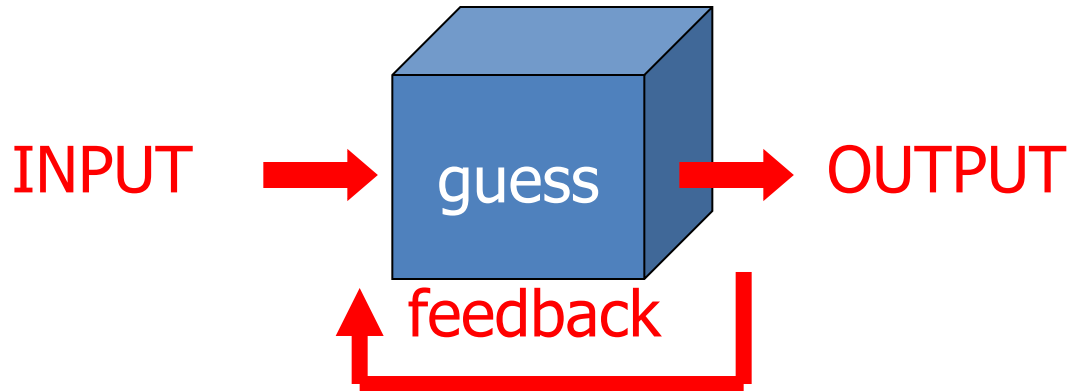


Fig. 4.— On the right we plot a 2 dimensional demonstration of the color space partitioning. In each of these cells we applied the polynomial fitting technique to estimate redshifts. The left figure show the results.

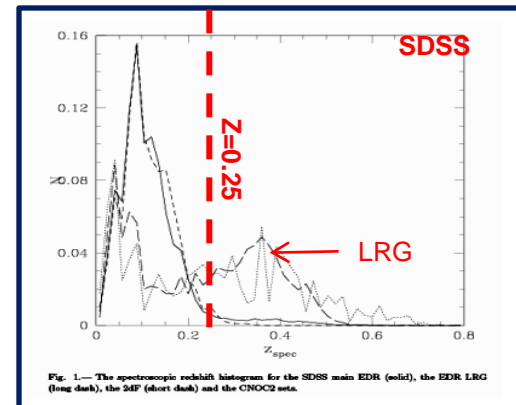
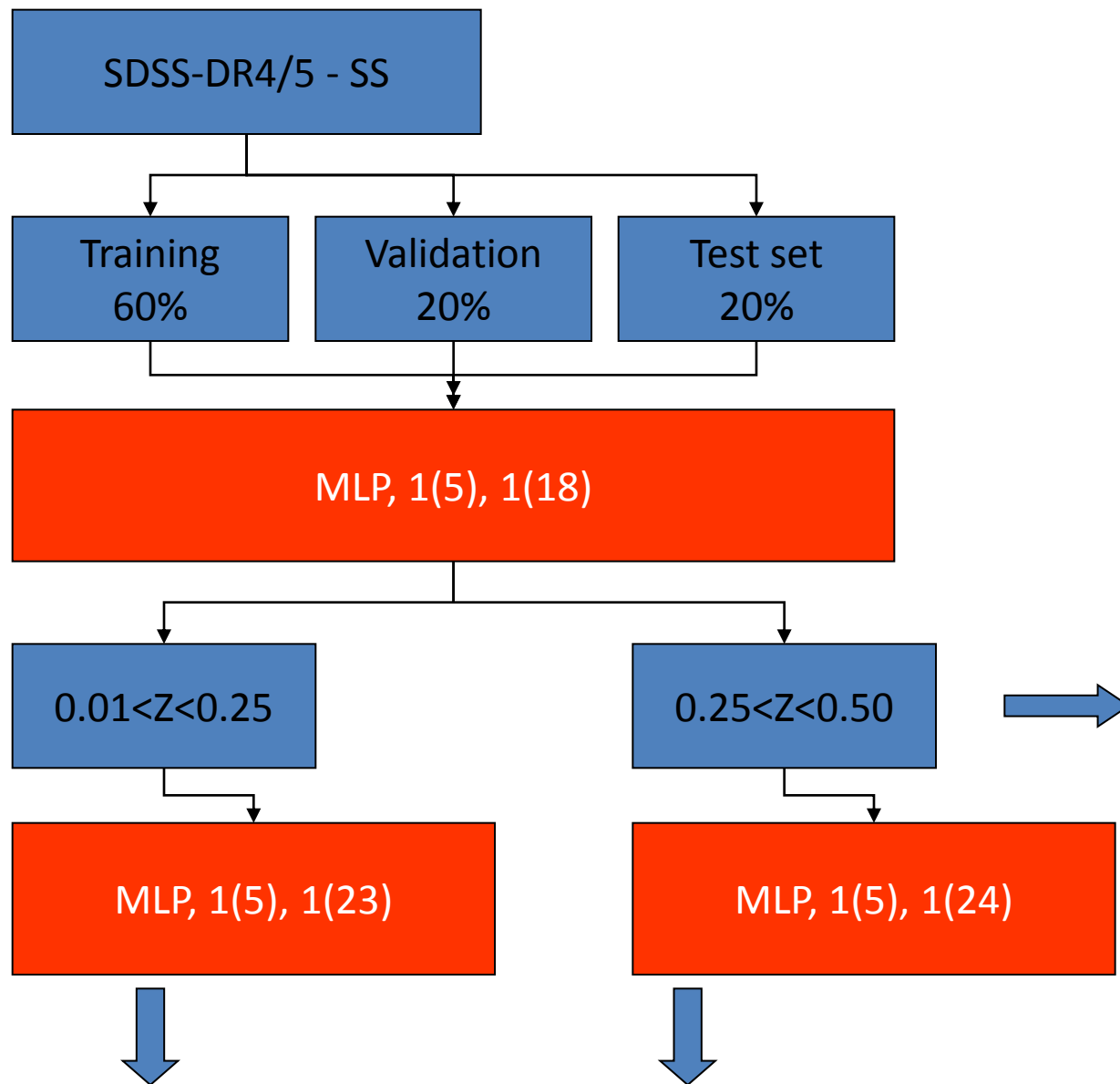
MLP or Multi Layers Perceptron



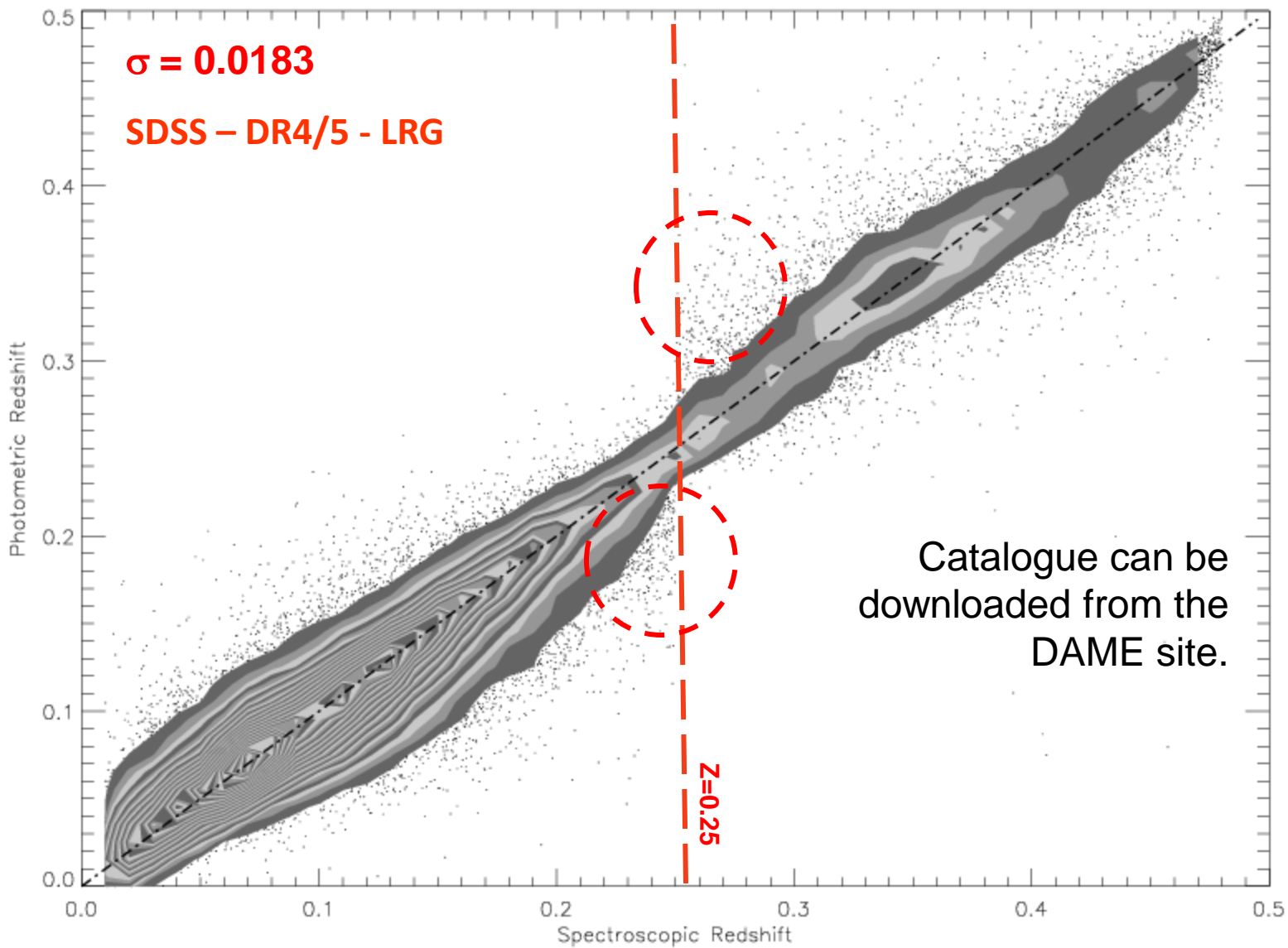
- input layer (n neurons)
- M hidden layer (1 or 2)
- Output layer (n' < n neurons)

Neurons are connected via activation functions

Different NN's given by different topologies, different activation functions, etc.



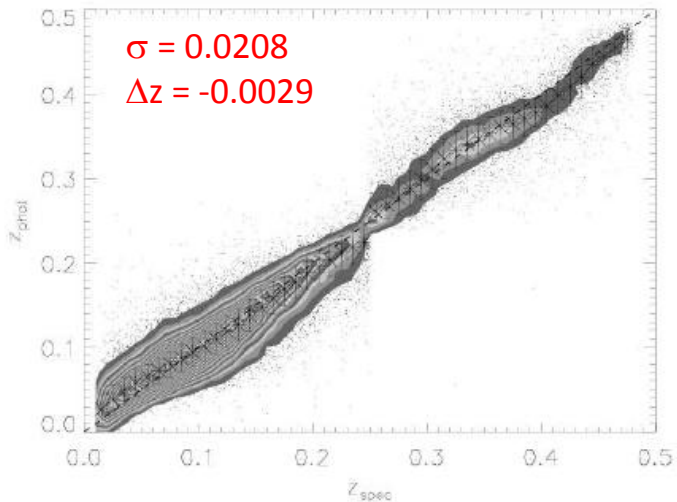
99.6 % accuracy



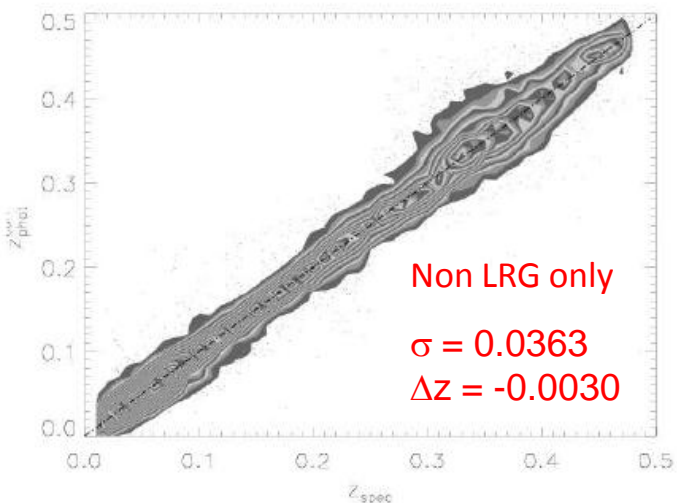
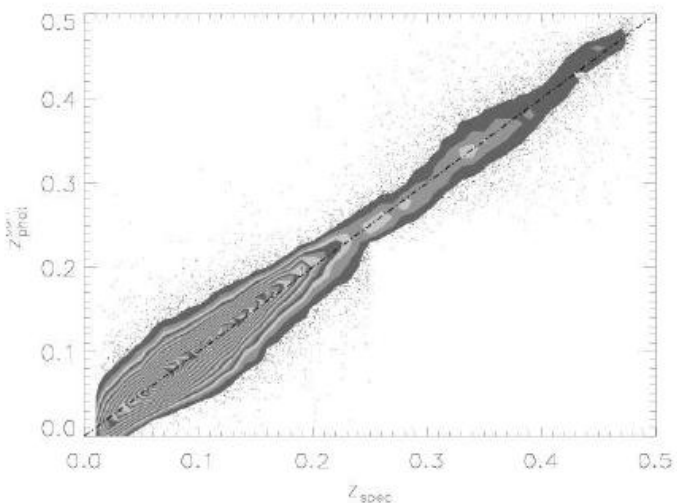
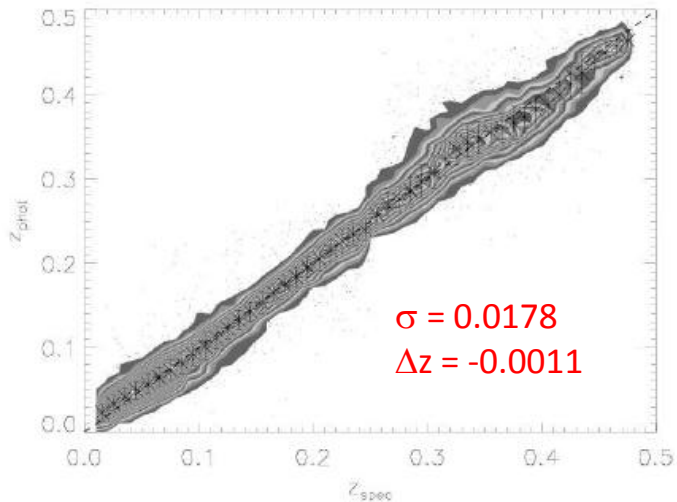


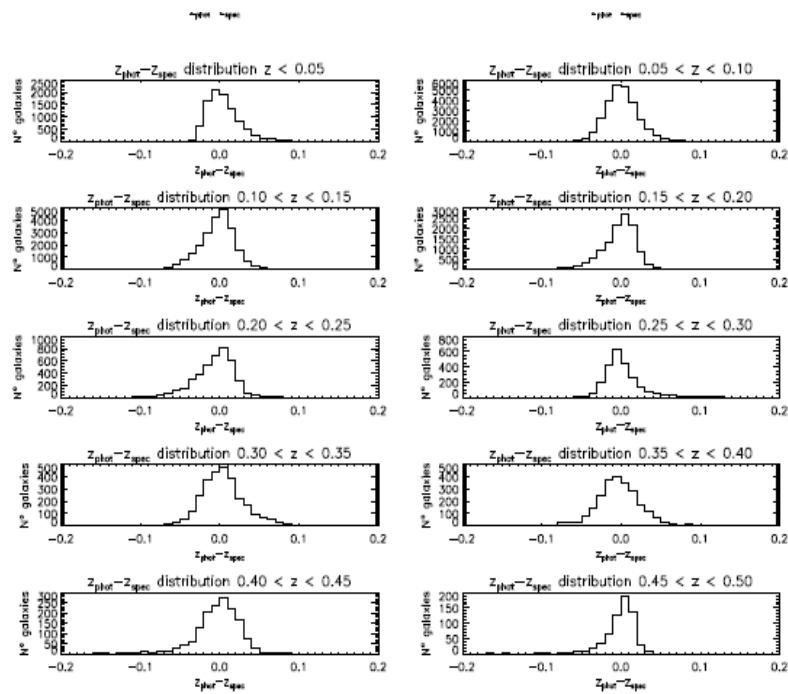
type	method	data	Δz_{rms}	Notes	Reference
SEDF	CWW	EDR	0.0666		(Csabai et al. 2003)
	Bruzual-CHARLOT	EDR	0.0552		(Csabai et al. 2003)
	Interpolated	EDR	0.0451		(Csabai et al. 2003)
	Polyomial	EDR	0.0318		(Csabai et al. 2003)
	KD-tree	EDR	0.0254		(Csabai et al. 2003)
ML	ANNz	EDR	0.0229		(Collister & Lahav 2004)
	SVM	EDR	0.027		(Wadadekar 2004)

General galaxy sample

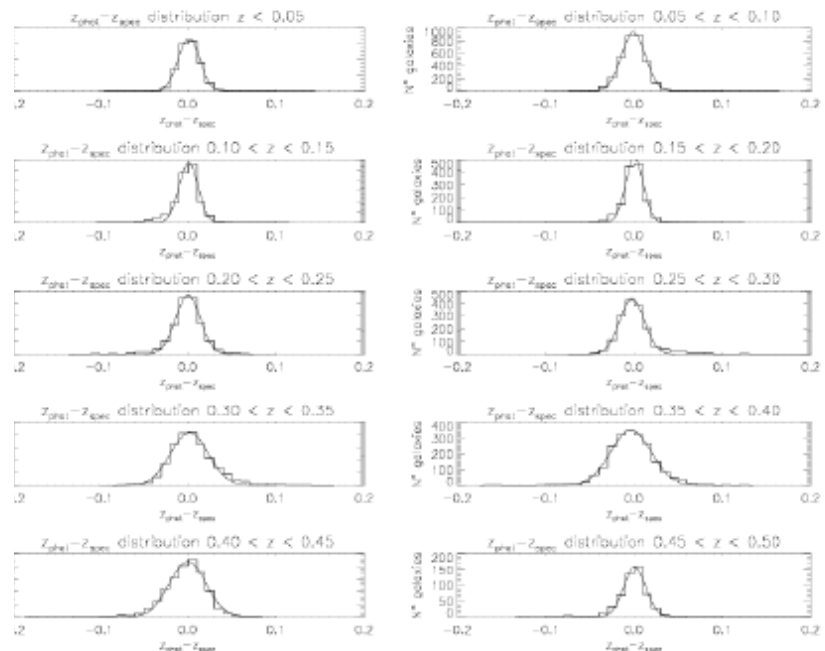


LRG sample





General galaxy sample

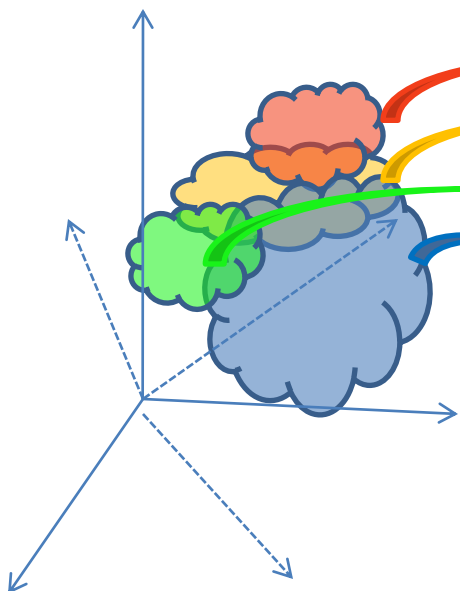


LRG sample

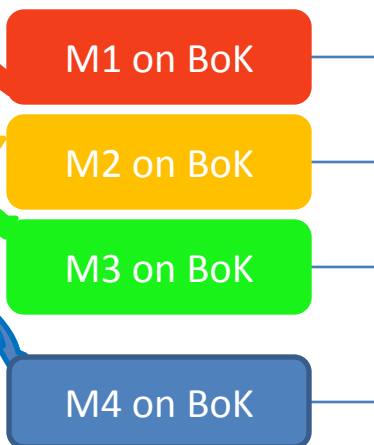
What do we learn if the BoK is biased:

- At high z LRG dominate and interpolative methods are not capable to “generalize” rules
- An unique method optimizes its performances on the parts of the parameter space which are best covered in the BoK

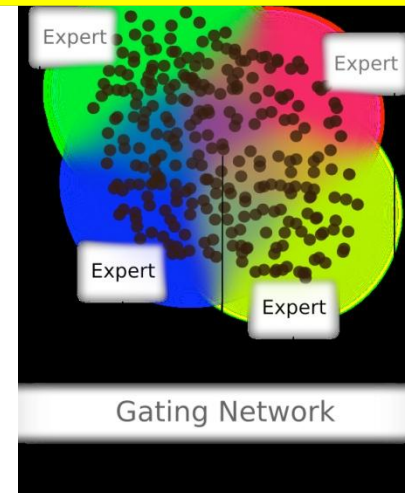
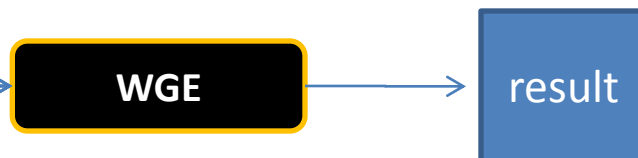
Step 1:
unsupervised clustering in
parameter space



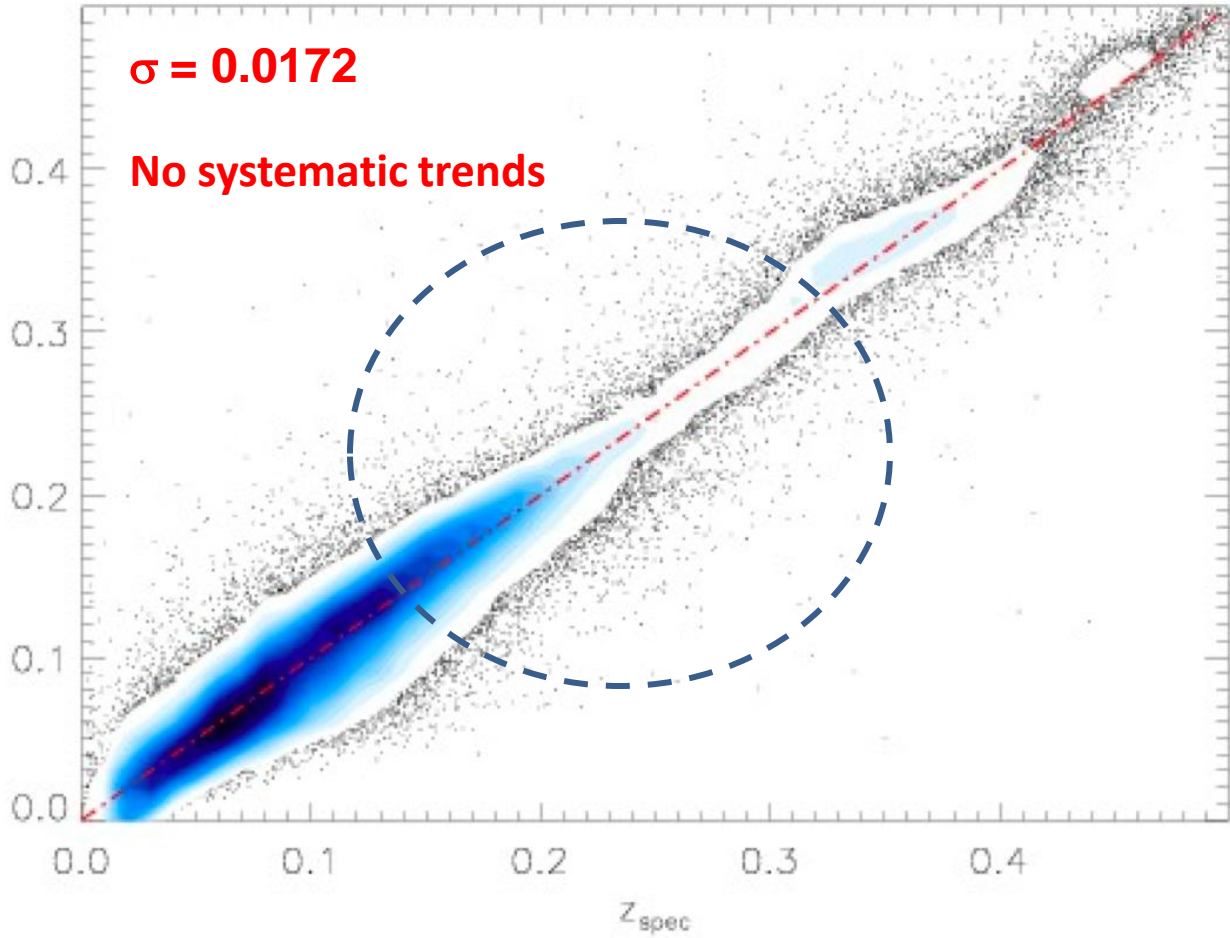
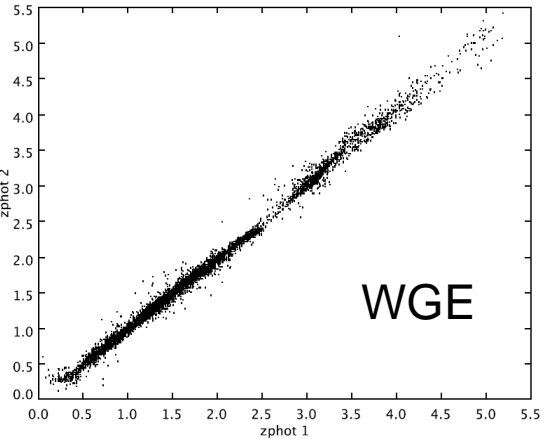
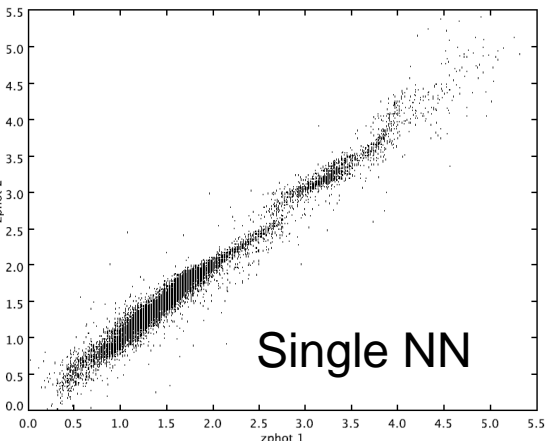
Step 2:
supervised training of
different NN for each cluster



Step 3:
output of all NN go to WGE
which learns the correct
answer



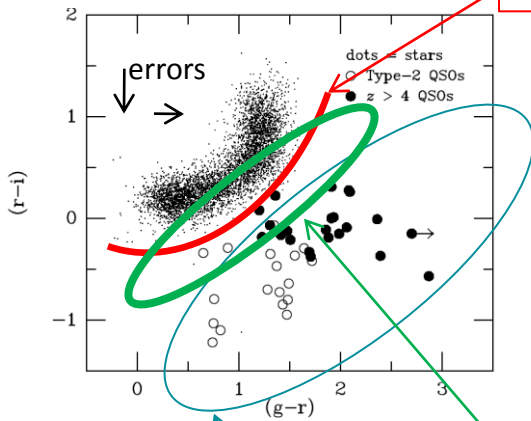
Laurino et al. 2009a,2009b



PART II - applications to observational cosmology

Photometric selection of candidate QSO's (as a clustering problem)

Traditional way to look for candidate QSO in 3 band survey



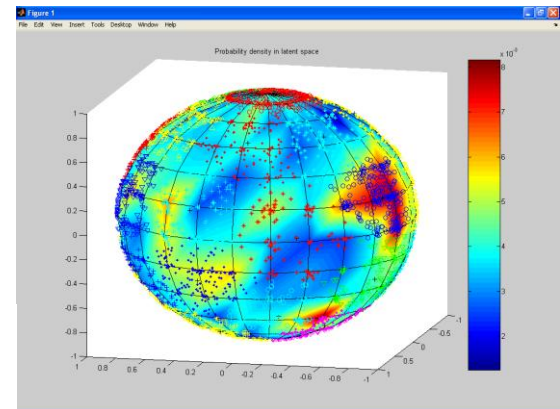
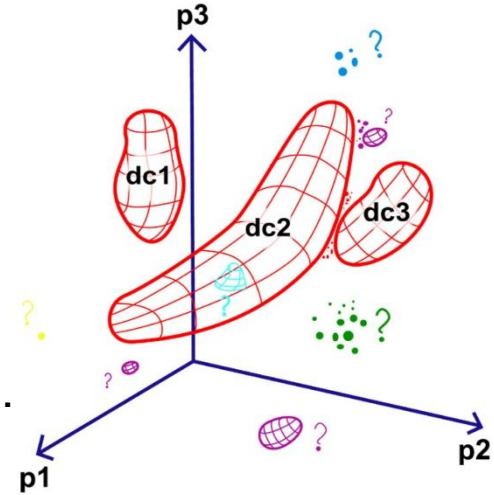
Cutoff line

Candidate QSOs for spectroscopic follow-up's

Ambiguity zone

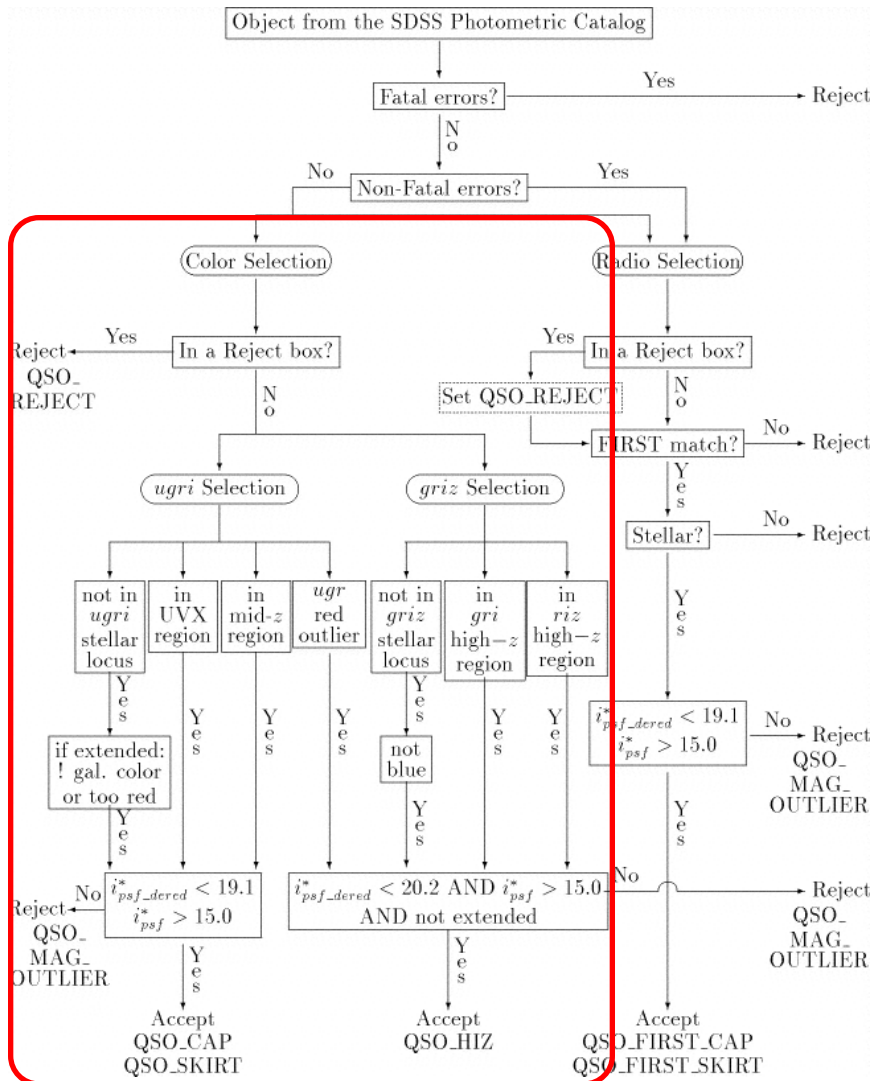
Adding one feature improves separation...

A Generic Machine-Assisted Discovery Problem: Data Mapping and a Search for Outliers



PPS projection of a 21-D parameter space showing as blue dots the candidate quasars. Notice better disentanglement

SDSS QSO candidate selection algorithm (Richards et al, 2002) targets star-like objects as QSO candidate according to their position in the SDSS colours space (u-g,g-r,r-i,i-z), if one of these requirements is satisfied:

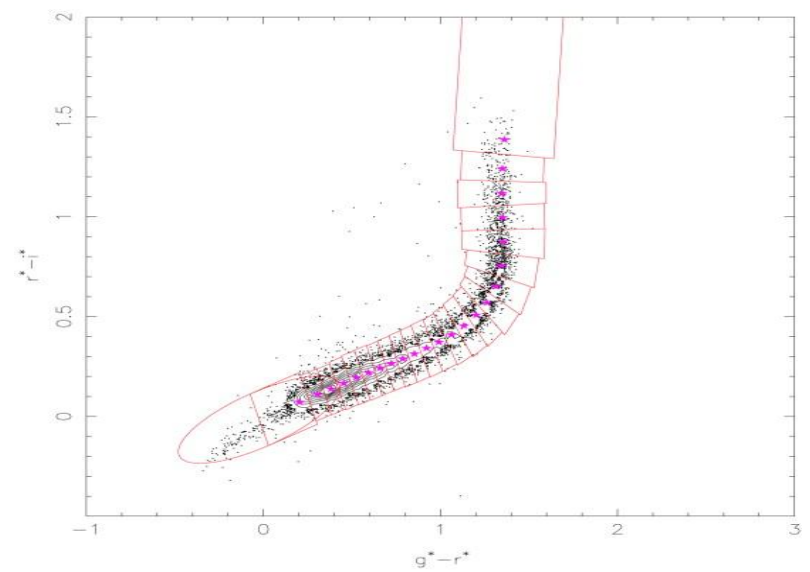
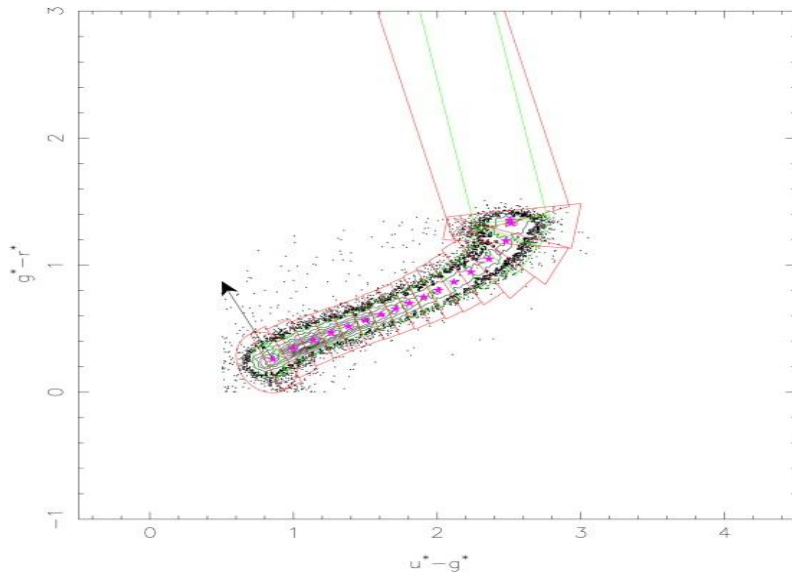


- QSOs are supposed to be placed $>4\sigma$ far from a cylindrical region containing the “stellar locus” (S.L.), where σ depends on photometric errors.

OR

- QSOs are supposed to be placed inside the inclusion regions, even if not meeting the previous requirement.

$c = 95\%$, $e = 65\%$
locally less

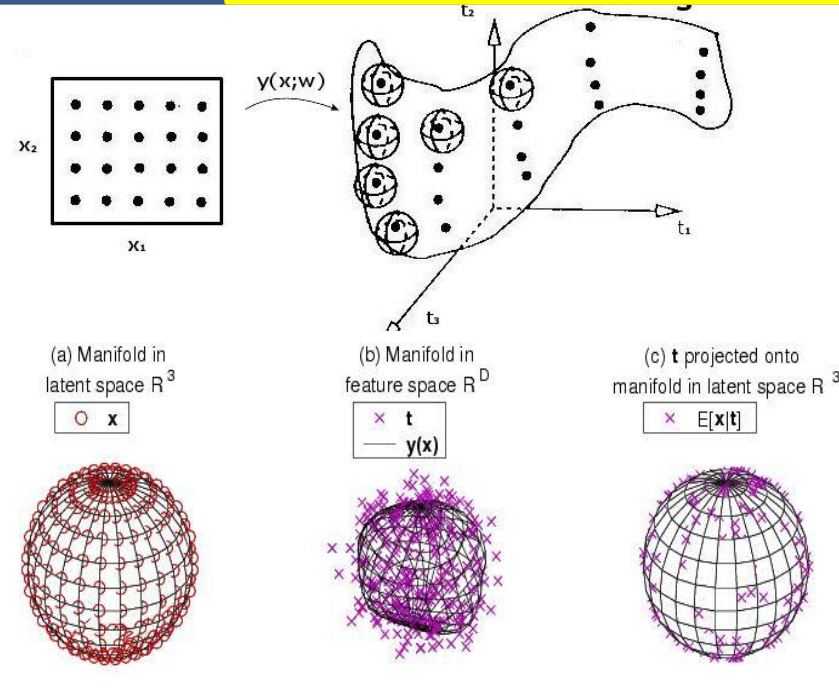


- 1. inclusion regions** are regions where S.L. meets QSO's area (due to absorption from Ly α forest entering the SDSS filters, which changes continuum power spectrum power law spectral index). All objects in these areas are selected so to sample the [2.2, 3.0] redshift range (where QSO density is also declining), but at the cost of a worse efficiency (Richards et al, 2001).
- 2. exclusion regions** are those regions outside the main "stellar locus" clearly populated by stars only (usually WDs). All objects in these regions are discarded.

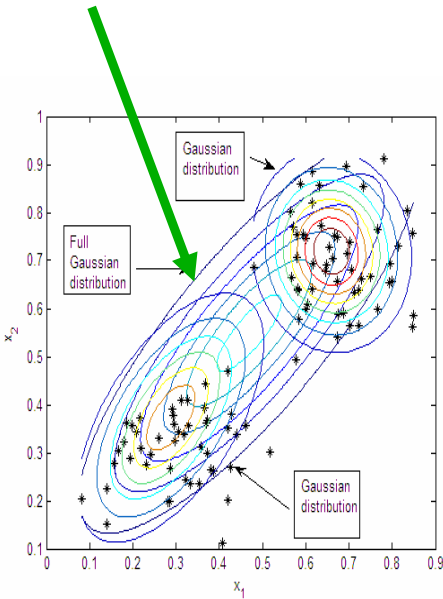
Overall performance of the algorithm: completeness $c = 95\%$, efficiency $e = 65\%$, but locally (in colours and redshift) much less.

Step 1: Unsupervised clustering

PPS determines a large number of distinct groups of objects: nearby clusters in the colours space are mapped onto the surface of a sphere.

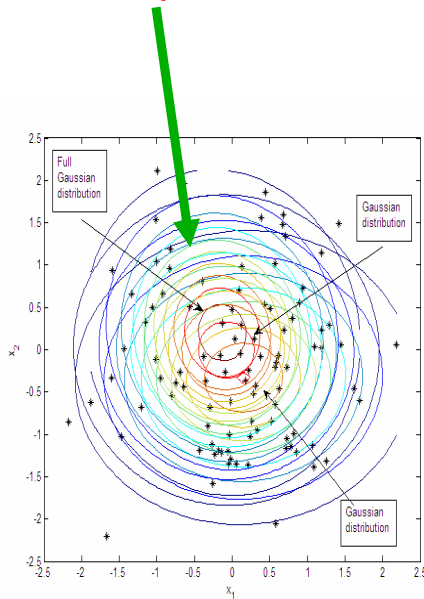


Not replaced!



NegE=750

Replaced!

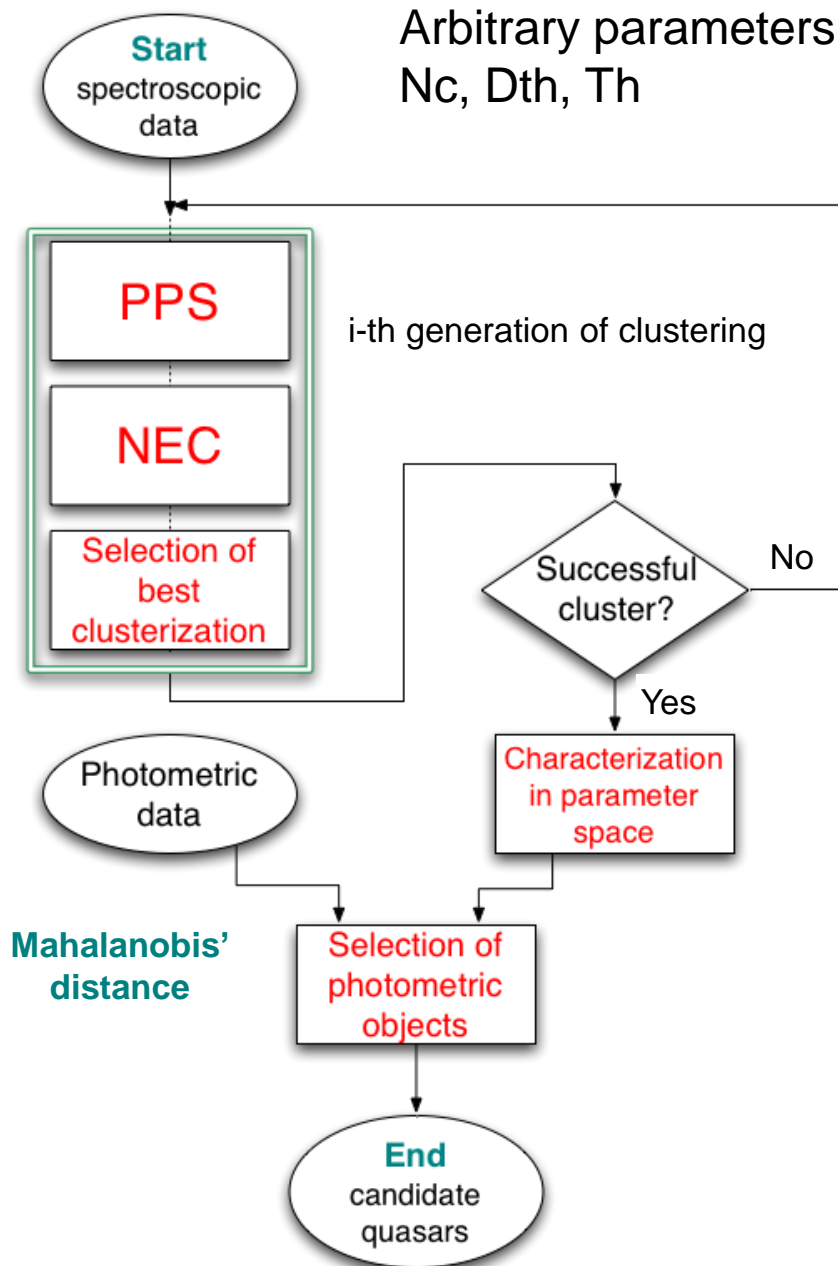


NegE=4

Step 2: Cluster agglomeration

NEC aggregates clusters from PPS to a (a-priori unknown) number of final clusters.

- Plateau analysis:** final number of clusters $N(D)$ is calculated over a large interval of D , and critical value(s) D_{th} are those for which a plateau is visible.
- Dendrogram analysis:** the stability threshold(s) D_{th} can be determined observing the number of branches at different levels of the graph.



To determine the critical dissimilarity D_{th} threshold we rely not only on a stability requirement.

A cluster is successful if fraction of confirmed QSO is higher than assumed fractionary value (Th)

D_{th} is required to maximize **NSR**

$$NSR = \frac{\text{Number of successful clusters}}{\text{Number of total clusters}}$$

The process is recursive: feeding merged unsuccessful clusters in the clustering pipeline until no other successful clusters are found.

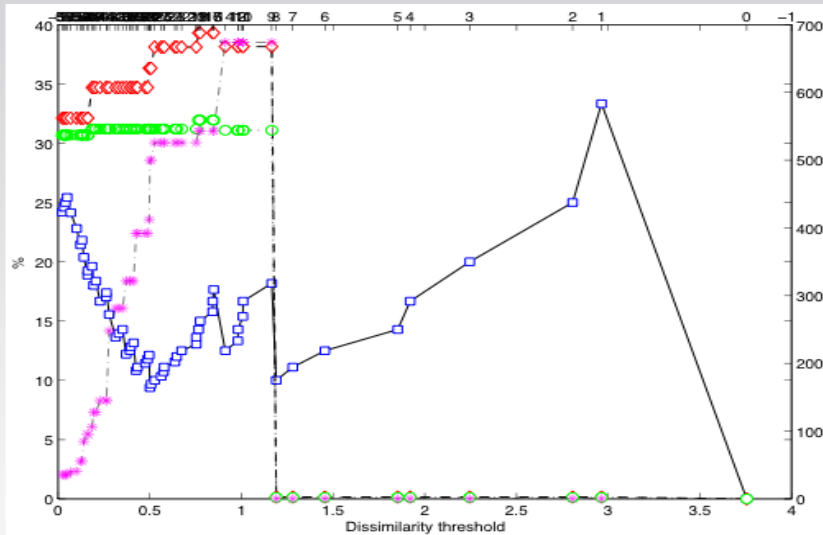
The overall efficiency of the process e_{tot} is the sum of weighed efficiencies e_i for each generation:

$$e_{tot} = \sum_{i=1}^n e_i$$

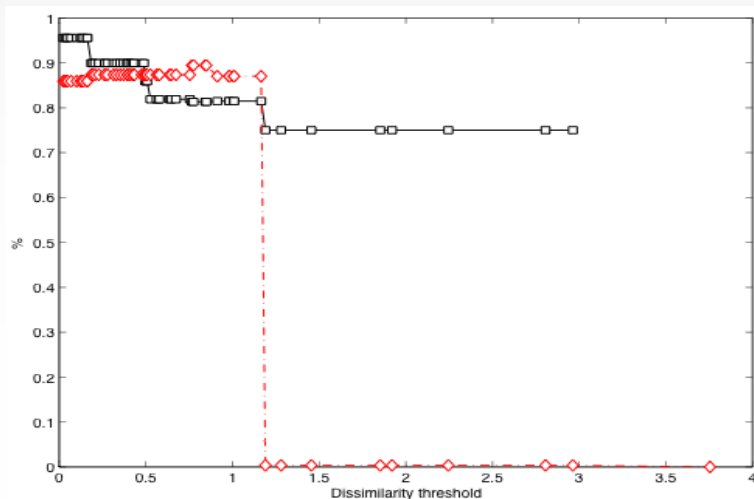
An example of “tuning”

Choice of the clustering

NSR



Efficiency and completeness



e and c estimation

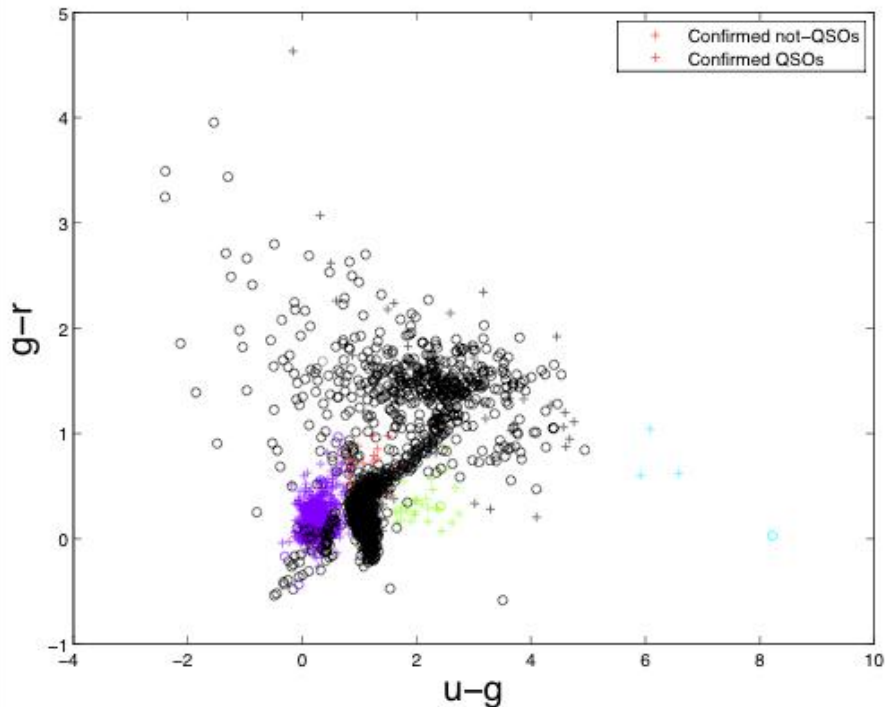
To assess the reliability of the algorithm, the same objects used for the “training” phase have been re-processed using photometric informations only. Results have been compared to the BoK.

algorithm \ labels	QSOs	not QSOs
QSOs	759	72
not QSOs	83	1327

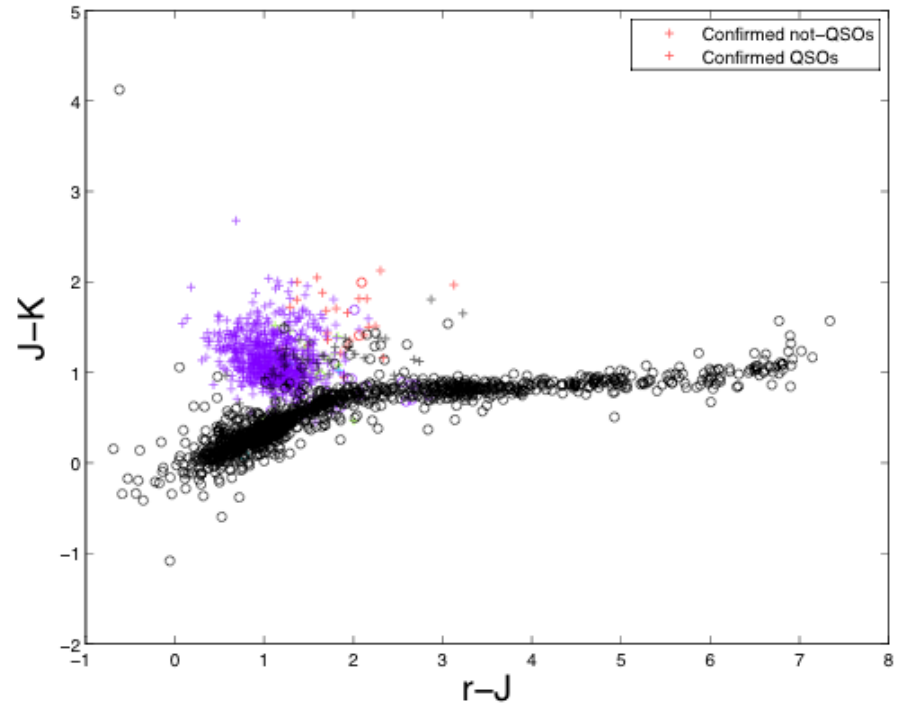
$e = 83.4 \%$ $c = 89.6 \%$

Confusion matrix

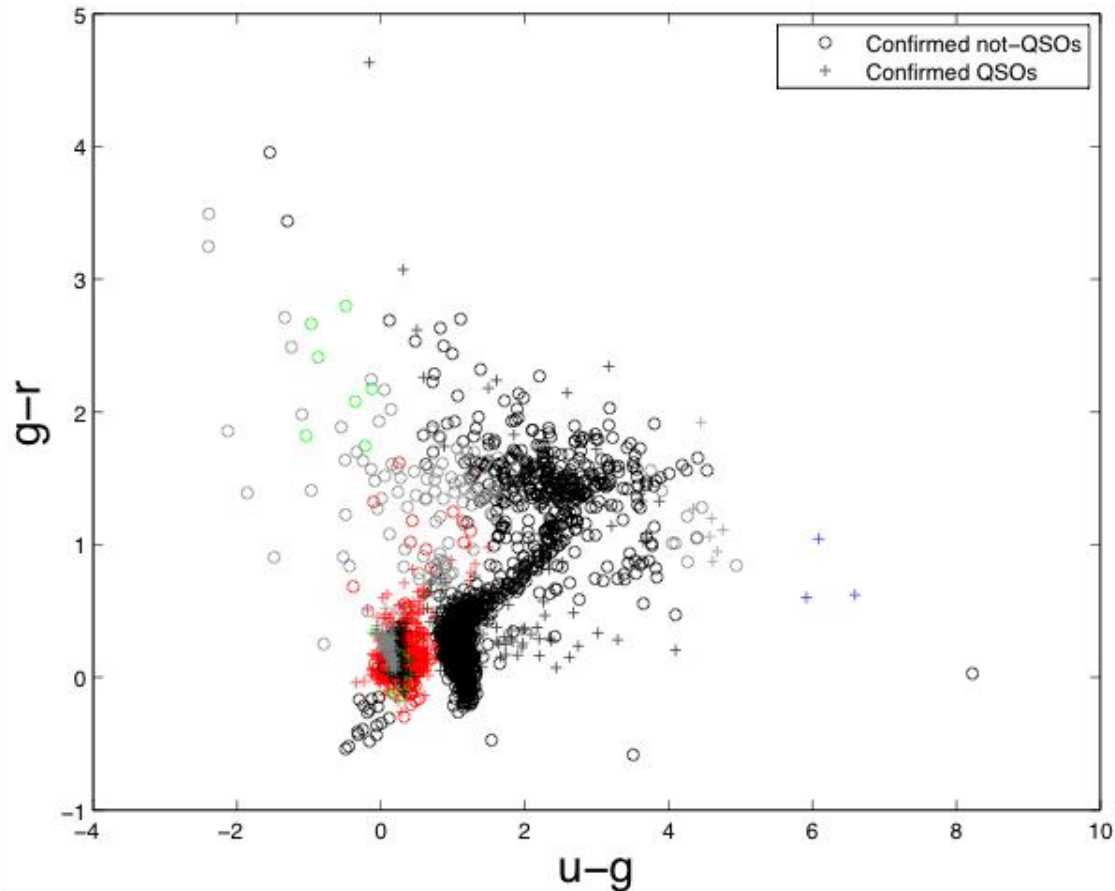
u - g vs g - r



r - J vs J - K

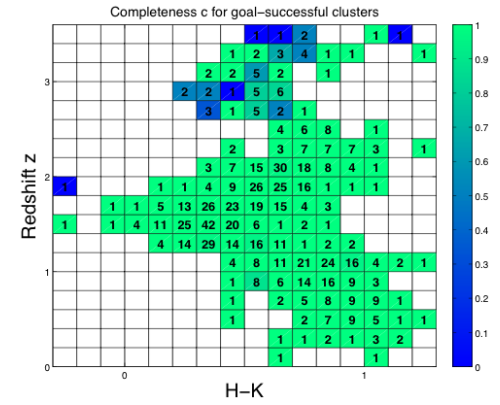
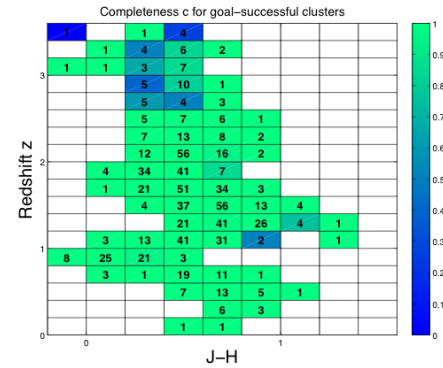
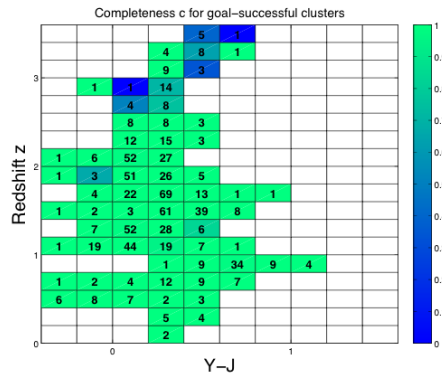
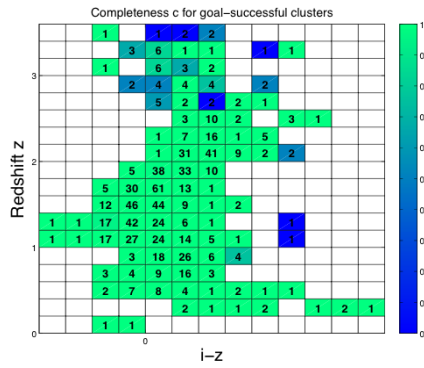
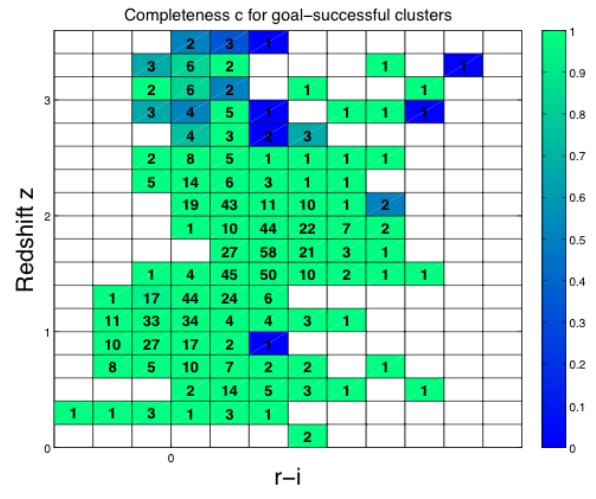
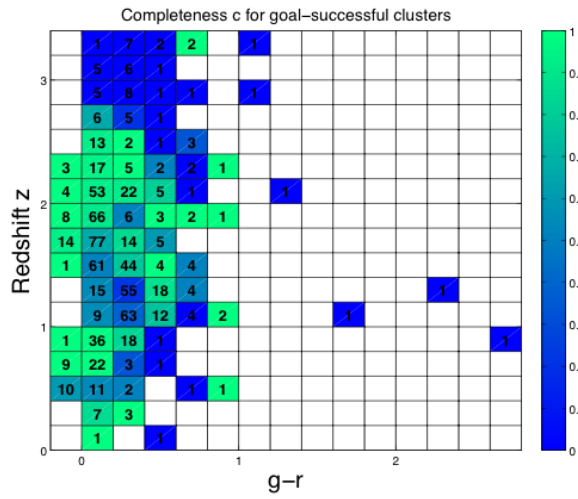
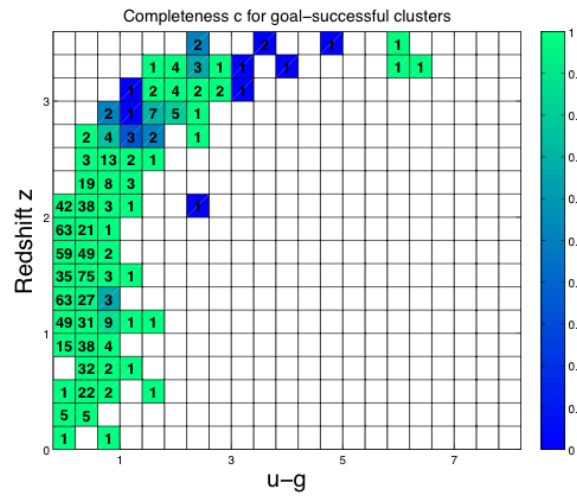


Only a fraction (43%) of these objects have been selected as candidate QSO's by SDSS targeting algorithm in first instance: the remaining sources have been included in the spectroscopic program because they have been selected in other spectroscopic programmes (mainly stars).

u - g vs g - r

In this experiment the clustering has been performed on the same sample of the previous experiment, using only optical colours.

Experiment 2: local values of c

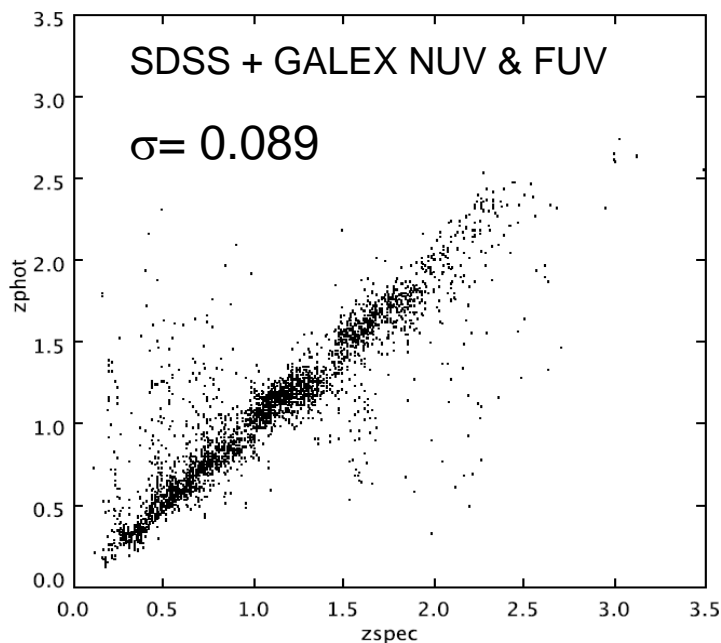
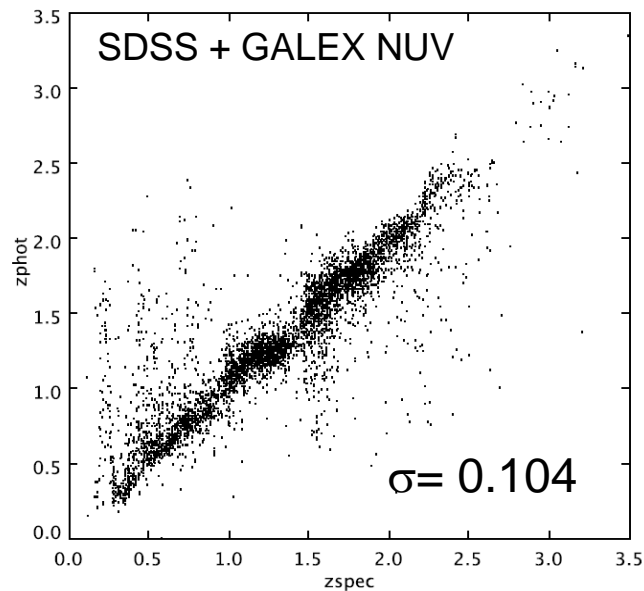
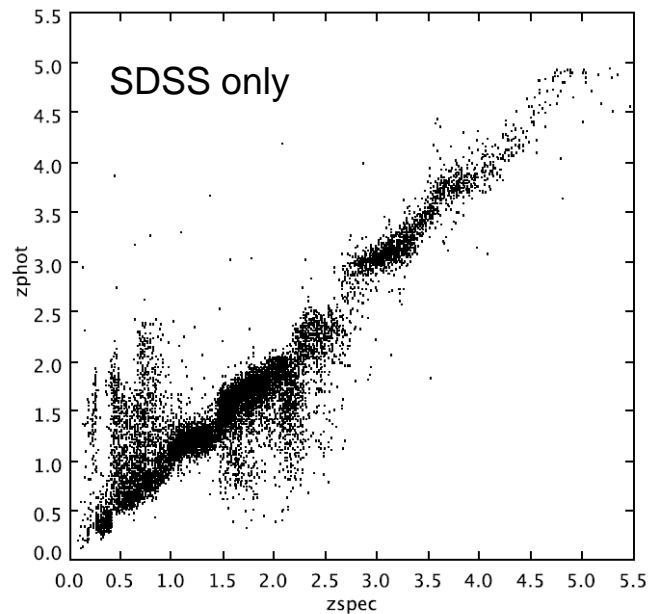


<u>Sample</u>	<u>Parameters</u>	<u>Labels</u>	<u>ϵ_{tot}</u>	<u>C_{tot}</u>	<u>n_{gen}</u>	<u>n_{suc_clus}</u>
Optical QSO candidates (1)	SDSS colours	'specClass'	83.4 % (± 0.3 %)	89.6 % (± 0.6 %)	2	(3,0)
Optical + NIR star-like objects (2)	SDSS colours + UKIDSS colours	'specClass'	91.3 % (± 0.5 %)	90.8 % (± 0.5 %)	3	(3,1,0)
Optical + NIR star-like objects (3)	SDSS colours	'specClass'	92.6 % (± 0.4 %)	91.4 % (± 0.6 %)	3	(3,0,1)

The catalogue of candidate quasars is publicly available at the URL:

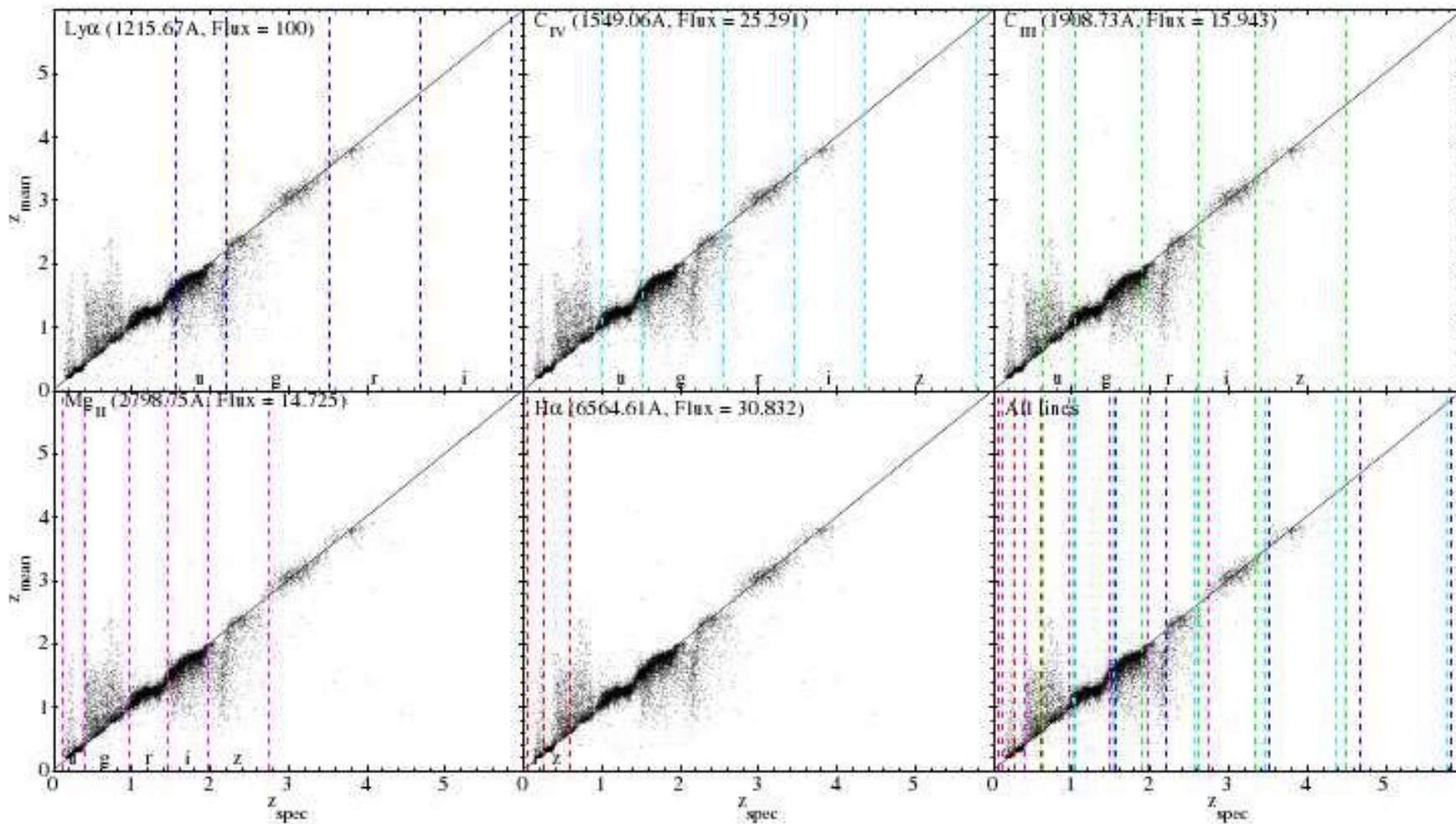
http://voneural.na.infn.it/catalogues_qsos.html

BUT ... LET'S GO BACK TO PHOT-Z



No need for fine tuning !!!

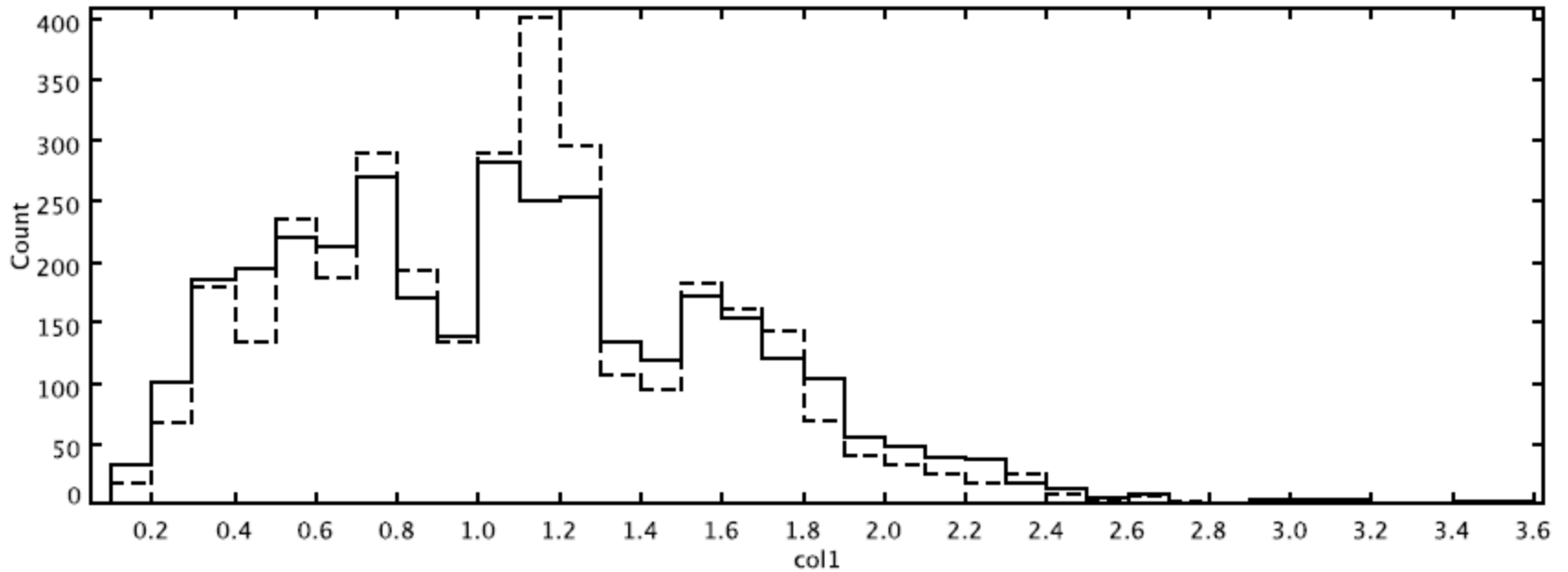
Only New BoK !!!

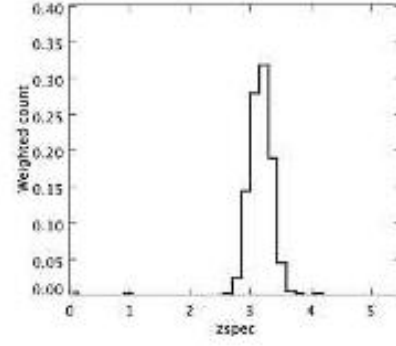
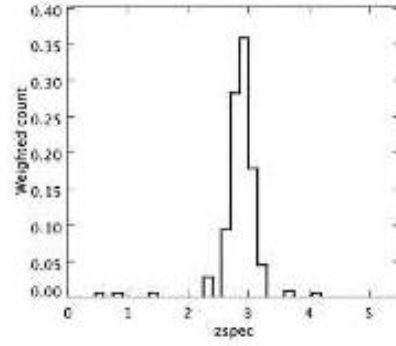
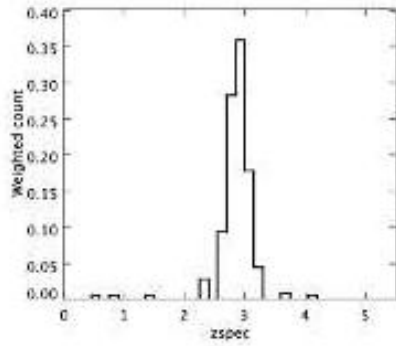
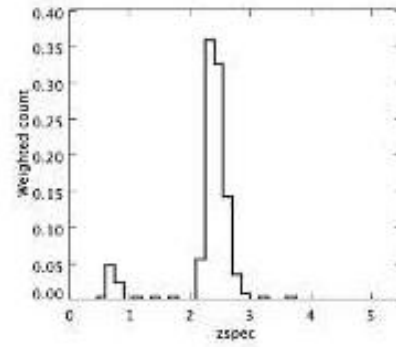
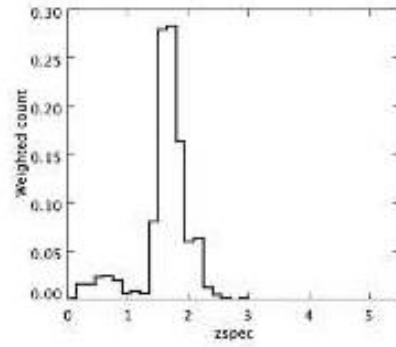
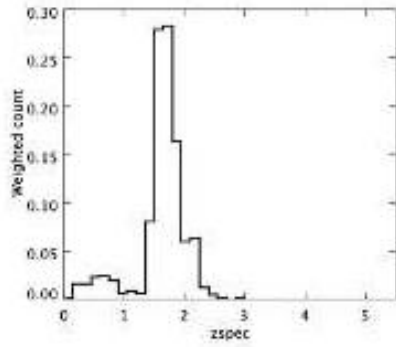
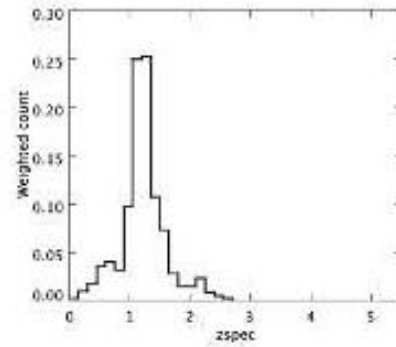
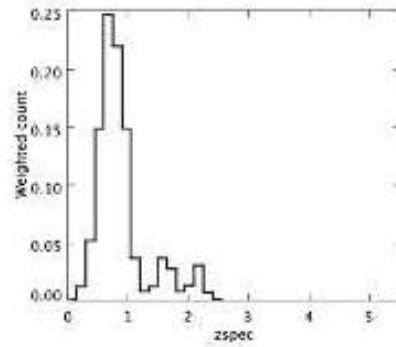
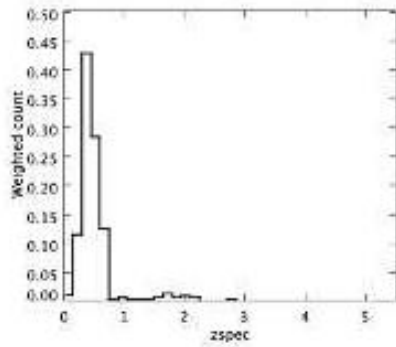


Degeneracy induced by lines exiting photometric bands



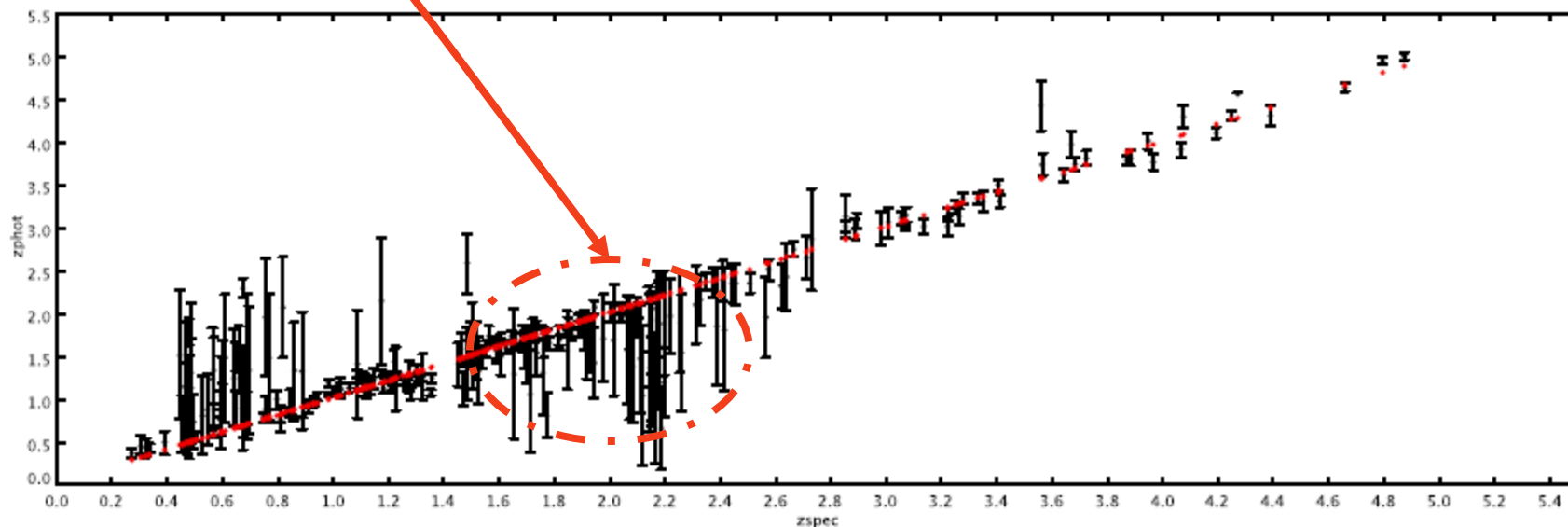
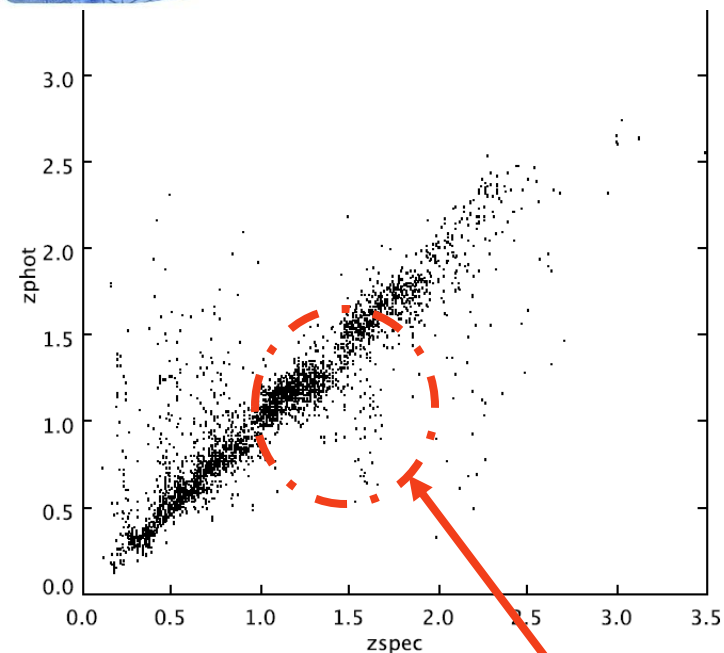
Distribution of Z_spec (solid) and Z_phot (dashed) for test set !!!!

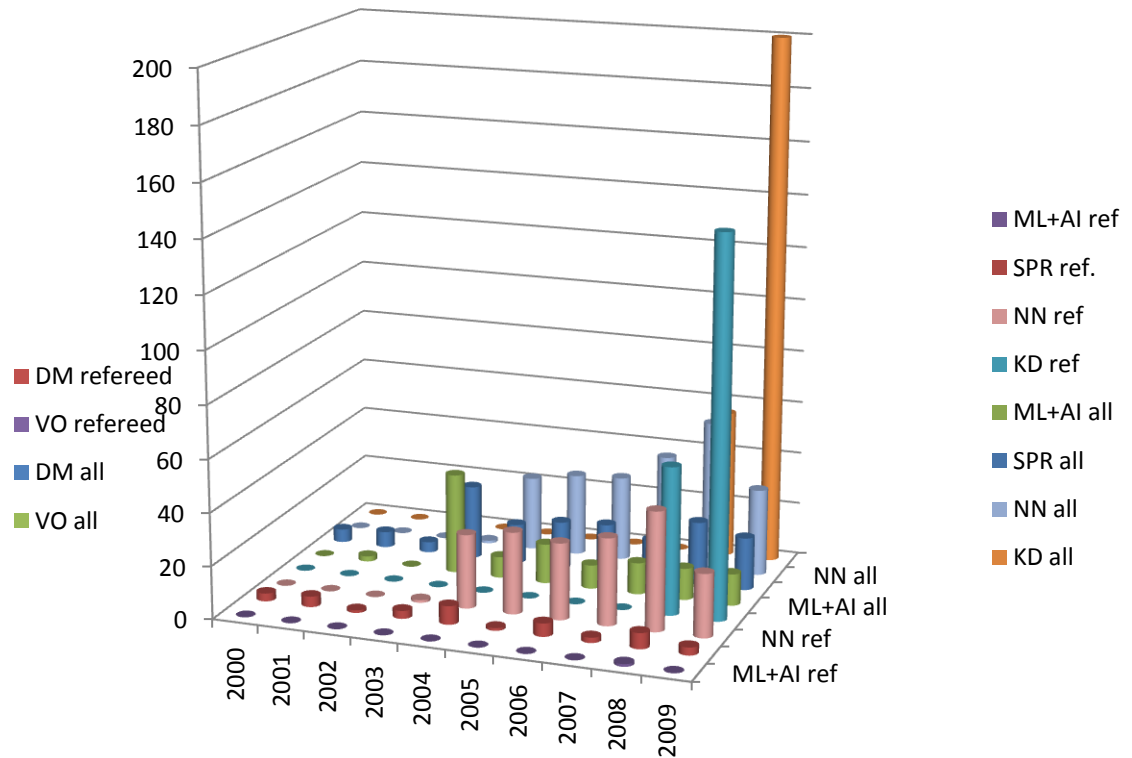
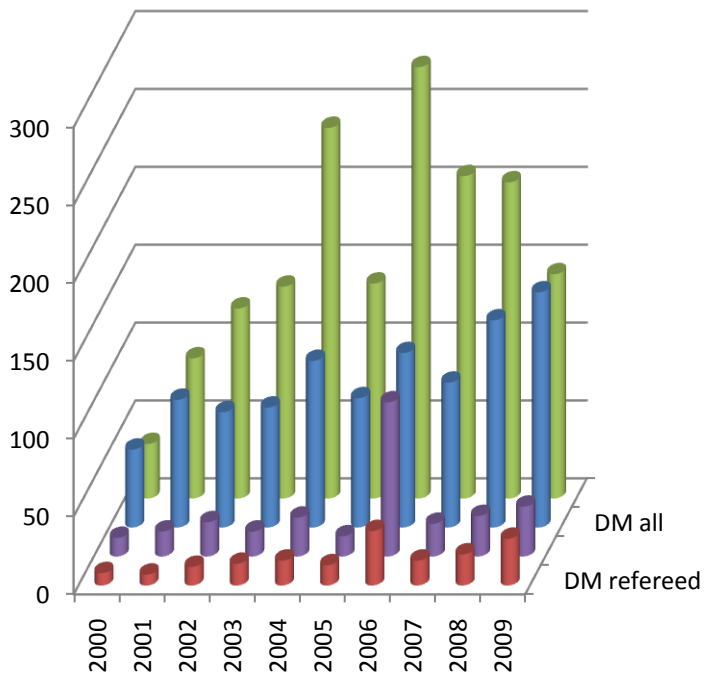




Errors:

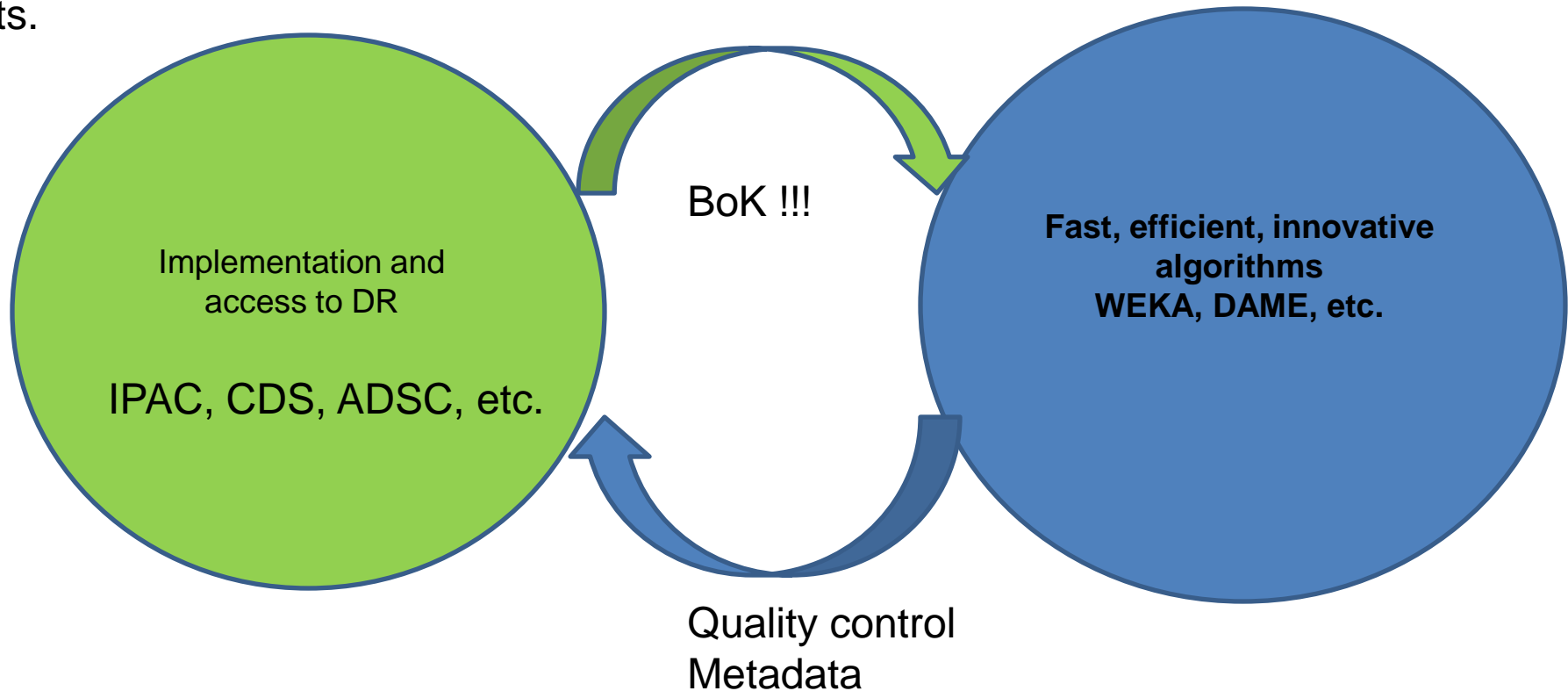
- **Input noise:** error propagation on the input parameter (Ball et al. 2008)
- **Model variance:** different models make differing predictions (Collister & Lahav 2004)
- **Model bias:** different models may be affected by different biases.
- **Target noise:** in some regions of the parameter space, data may represent poorly the relation between featured and targets (*Laurino 2009*).





1. Number of technical/algorithmic papers increases with new funding opportunities. Number of refereed papers remains constant.
2. Most of the work, so far, remains at the implementation stage (computer Science and algorithm development) and does not enter the “science production” stage...
3. Out of one thousand papers checked (galaxies, observational cosmology, survey) over the last two years: DM could be applied or involved in at least 30% of them leading to better results

Machine Learning based Data Mining is unavoidable when working on huge data sets.



Accuracy of results depends on accuracy of BoK !!!!

The extraction of BoK's offers challenges to good data repositories and data archives.

Reliability and completeness of information

(no data is better than bad data)

Compliance with ontologies

Advanced queries in natural language

Max Brescia

