# The VO-Neural/DAME infrastructure: an integrated data mining system support for the e-science community

**OACN**

Massimo Brescia  -  Giuseppe Longo

&

Project Team

*INAF – Osservatorio Astronomico di Capodimonte*

*Dipartimento di Fisica – Università degli Studi di Napoli Federico II*

*California Institute of Technology*

# Project Highlights

Originally named VO-Neural, recently the project is evolving to DAME (DAta Mining & Exploration), but the final name and logo is still under design.

The project, an evolution of the former <u>AstroNeural</u> Collaboration, is financed through:
E.U. grant VOTECH and VO-AIDA
Italian Ministry of Research in the framework of the PON-S.Co.P.E.
 Italian Ministry of Foreign Affairs through a great relevance bilateral project Italy-USA
**VO-Neural/DAME main goal is the design and development of scientific data mining tools, based on Information Technology instruments.**

*Partnership:*
Dipartimento di Fisica (sez. di Astrofisica) - Università degli Studi di Napoli Federico II
INAF - Osservatorio Astronomico di Capodimonte
California Institute of Technology, Pasadena - USA

*Collaborations:*
S.Co.P.E. (high Performance distributed Cooperative System for scientific Experiment)
INAF - Osservatorio Astronomico di Trieste
Dipartimento di Informatica - Università degli Studi di Napoli Federico II
Dipartimento di Ingegneria Informatica -  Università degli Studi di Napoli Federico II
EURO-VO The European Virtual Observatory
IVOA (International Virtual Observatory Alliance)

# Trend of Information Technology

## *Cloud / GRID computing*

Cloud computing is Internet based development and use of computer technology. The cloud is a metaphor for the Internet and is an abstraction for the complex infrastructure it conceals.
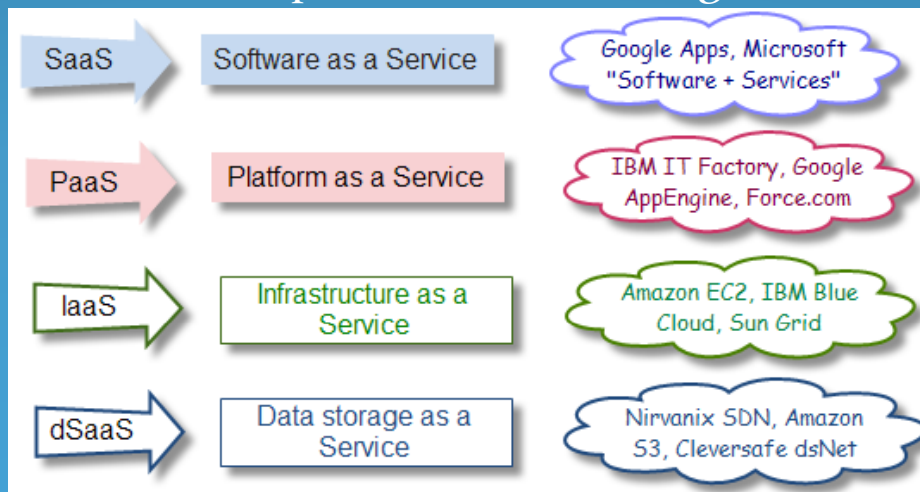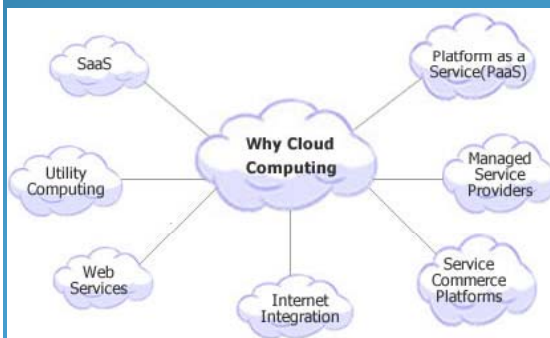
It is a style of computing, provided "*as a service*", to access enabled services from the Internet without knowledge of, expertise with, or control over the technology infrastructure that supports them.

So far, Cloud computing can be considered to implement the following ideas :

**Utility computing** - which was first suggested by *John McCarthy* in 1961, where computing is viewed as a public utility;

**Cluster computing** - which views a group of linked computers as a single virtual computer for high-performance computing (HPC);

**Grid computing** - where the linked computers tend to be organized "*as resources*" to solve a common problem;

Cloud computing landscape

## Scientific community requirements

*Why to land on an "open" distributed infrastructure*

**EGEE GRID**   **PRIVATE GRID**

*Enabling GRID for E-SciencE*

| Target Group | Scientific community | Business |
|---|---|---|
| Service | short-lived batch-style processing (job execution) | long-lived services based on hardware virtualization |
| SLA | Local (between the EGEE project and the resource providers) | Global (between Amazon and users) |
| User Interface | High-level interfaces | HTTP(S), REST, SOAP, Java API, BitTorrent |
| Resource-side middleware | Open Source (Apache 2.0) | Proprietary |
| Ease of Use | Heavy | Light |
| Ease of Deployment | Heavy | Unknown |
| Resource Management | probably similar | |
| Funding Model | Publicly funded | Commercial |

A universal research infrastructure:

"*Un ambiente dove le risorse di ricerca (HW, SW e DATI) possano essere condivise rapidamente e a cui si possa accedere da ovunque sia necessario promuovere una ricerca migliore e più efficace*"

# Virtualized service/resource oriented Infrastructure

Virtualization brings new standardized capabilities to data centers

## Virtualized Data Center



### GRID Computing

Resource oriented

**Grid computing** solutions enable parallel processing of computational tasks, often using idle vs. dedicated capacity.

**Goals**: "Accelerate throughput from decades to days or from months to minutes."

"Enable deep computations that are otherwise intractable."

**Characteristics**: Large numbers of work requests run for short periods of time (minutes / hours).

### CLOUD Computing

Service oriented

**Cloud computing** enables self-service provisioning of virtual machines.

**Goal**: "Simplify deployment of Operating Systems and app servers."

**Characteristics**: Small numbers of VM allocations held for long periods of time (days/months).

*Join data virtualization, resources and services brings to*

## Virtual Organization of data
## HW resource oriented
## SW service oriented
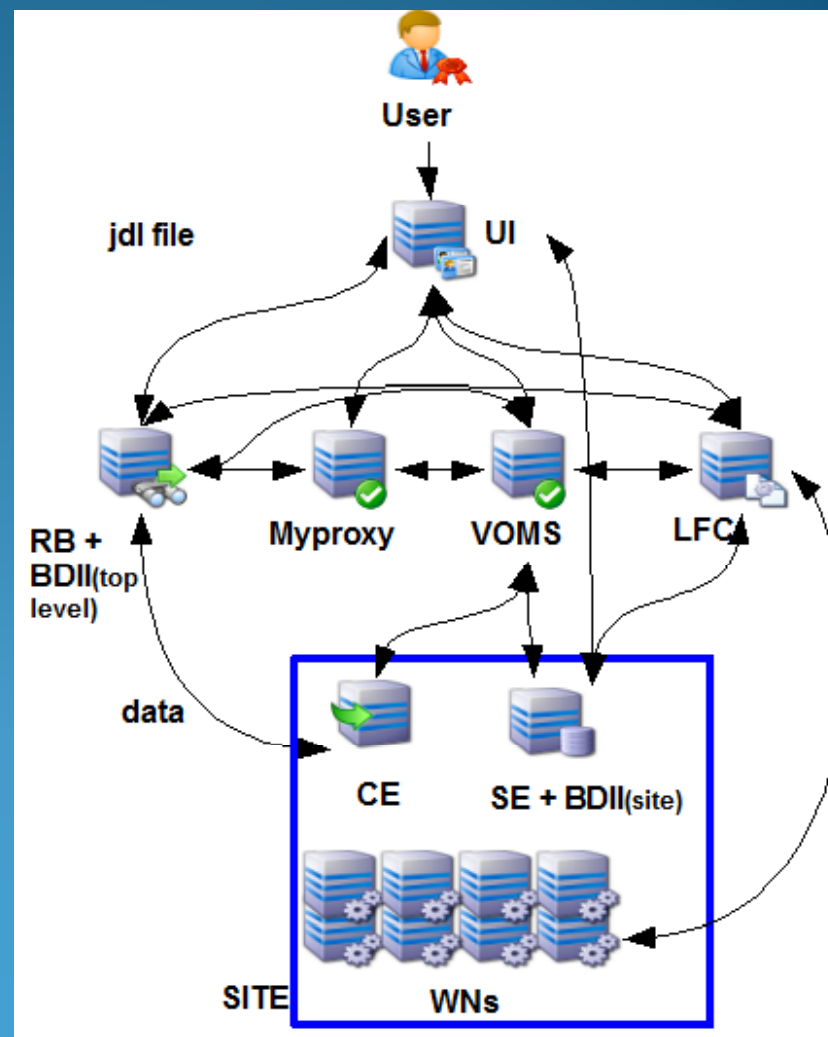
# What is S.Co.P.E. GRID



**INFRASTRUCTURE INTEROPERABILITY**

**APPLICATION INTEROPERABILITY**

Strategic goal:

**ITALIAN e-INFRASTRUCTURE OPEN TO THE COLLABORATION & INTERACTION BETWEEN RESEARCH AND INDUSTRY**

**Resources:**

➢ Computing power of some thousands of cores per project

➢ hundreds of TB per project

➢ distributed resources for massive computing

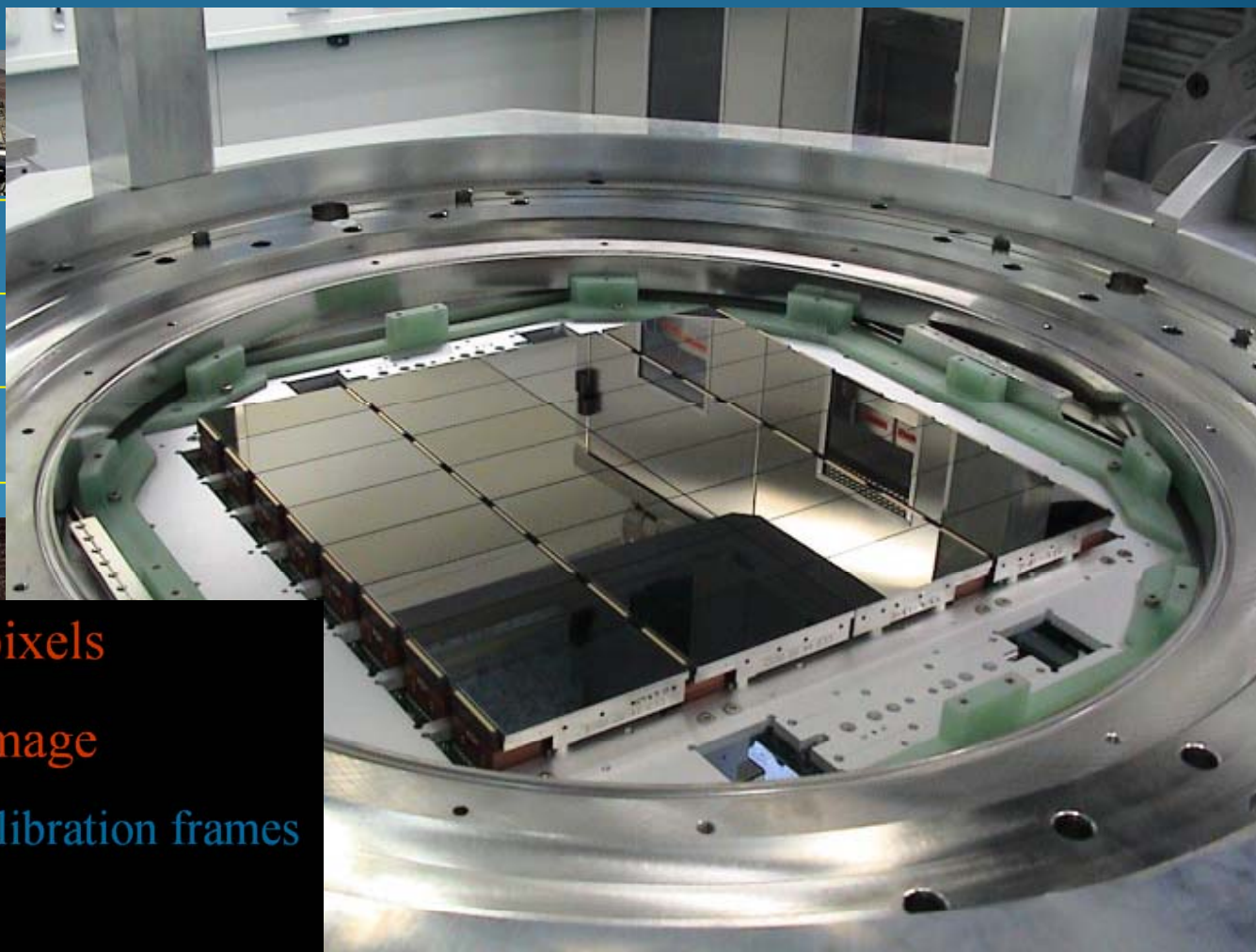# The astrophysical problem

## *Astronomical data rate*

10000000

1000000

100000

10000

1000

100

268,435,456 pixels

0.5 Gbyte x image

50 science frames + 50 calibration frames

⇩

50 Gbyte / Night

## The Astrophysical Application

### Considerations on the next breakthroughs

- We have reached the physical limit of observations (single photon counting) at almost all wavelenght...
- Detectors are linear
- All electromagnetic bands have been opened

### Hence

Our capability to gain new insights on the universe will depend mainly on:

- Capability to recognize patterns or trends in the parameter space (i.e. physical laws) which are not limited to the human 3-D visualization
- Capability to extract patterns from very large multiwavelenght, multiepoch, multi-technique parameter spaces

We need:
data archives organized in a unified Virtual Observatory for wide band cross-correlation;
Data mining software tools based on machine learning and self-adaptive mechanisms;
Distributed high performance computing infrastructure able to work on massive datasets;

# Knowledge Discoveries in Databases (KDD) is in practice still unknown to most astronomers
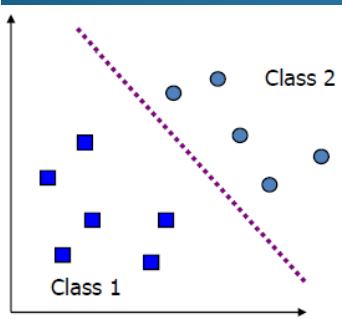
**To implement KDD tools is expensive** (time, computing, need for specialists), requires **coordinated efforts** between astronomers and computer scientists and is aimed to fulfill the needs of **large projects**

## Learning problems as "function approximation"

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

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

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

| variable | characteristics | Type | operation |
|---|---|---|---|
| Quantitative | Numerical with ordering relationship and possibility to define a metric | Actual measurement | regression |
| Categorical (non ordered) | Membership into a finite umber of classes. No ordering relationship. | Numerical codes (targets) arbitrarily orderd | Classification |
| Ordered categorical | Classes orderd by a relationship but there is no metric | Numerical codes n on arbitrarily orderd | Classification |



Class 2

Class 1



Class 2

Class 1



Class 2

Class 1

# Data Mining Exploration with VO-Neural/DAME

## Machine learning methods can be broadly grouped in:

## Supervised methods

They learn how to partition the parameter space by means of a training phase based on examples.

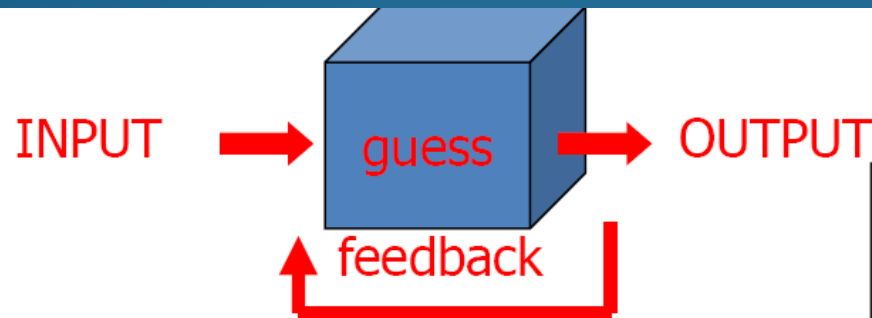Neural Networks such as the Multi Layer Perceptron (MLP), Support Vector Machines (SVM), etc.

## Pro's & Con's

- They are good for interpolation of data, very bad for extrapolations
- They need extensive bases of knowledge (i.e. uniformously sampling the parameter space) which are difficult to obtain;
- Errors are easy to evaluate
- Relatively easy to use

# Supervised Models: Multi Layer Perceptron

INPUT ➡ guess ➡ OUTPUT

↑ feedback

- 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.

## Supervised Models: Multi Layer Perceptron

$$N(U, e, P(e)) = \frac{1}{2k} \sum_j (Y_j - P(e))^2$$

$$\left| Y - P(e) \right| < \varepsilon$$

**backward** ↓

$$net\_input_j = \sum_h w_{jh} O_j$$

$$\delta = f'(o)(Y - P(e))$$

**forward** ↑

$$\delta'' = f'(o) \sum w \delta$$

$$net\_input_h = \sum_i w_{hi} O_i$$

**Back Propagation learning algorithm**

$$w_{ji}(new) = w_{ji}(old) + \eta \delta_i o_j + \alpha \Delta w_{ji}(old)$$

$$f(o) = \frac{1}{1 + \exp(-o)}$$

Output Layer

Hidden Layer

Hidden Layer

# Supervised Models: Support Vector Machine

given a training set formed by pairs [features-label]: $(x_i, y_i)$, i = 1...l
where $x_i \in R^n$ e $y_i \in \{1, -1\}^l$.
Support Vector Machines (SVM) try to solve the following optimization problem:

$$\min_{\omega, b, \xi} \frac{1}{2}\omega^T \omega + C\sum_{i=1}^{l} \xi_i$$

With the condition:

$$y_i(\omega^T \phi(x_i) + b) \geq 1 - \xi_i$$

Vectors $x_i$ are mapped into an higher dimensionality space where the SVM identify an hyper plane which maximizes the distances from the two classes

C > 0 is a classification error correction term

$$K(x_i, x_j) = \phi(x_i)^T(x_j)$$

Is the so called Kernel function

$$K(x_i, x_j) = \exp(-\gamma \|x_i - x_j\|^2), \qquad \gamma > 0$$

- linear: $K(x_i, x_j) = x_i^T x_j$.

- polynomial: $K(x_i, x_j) = (\gamma x_i^T x_j + r)^d$, $\gamma > 0$.

- radial basis function (RBF): $K(x_i, x_j) = \exp(-\gamma \|x_i - x_j\|^2)$, $\gamma > 0$.

- sigmoid: $K(x_i, x_j) = \tanh(\gamma x_i^T x_j + r)$.

Input space

$\phi(.)$

Feature space

– We should maximize the margin, $m$

$$m = \frac{2}{\|\mathbf{w}\|}$$

Class 2

Class 1

$\mathbf{w}^T\mathbf{x} + b = 1$

$\mathbf{w}^T\mathbf{x} + b = -1$

$\mathbf{w}^T\mathbf{x} + b = 0$
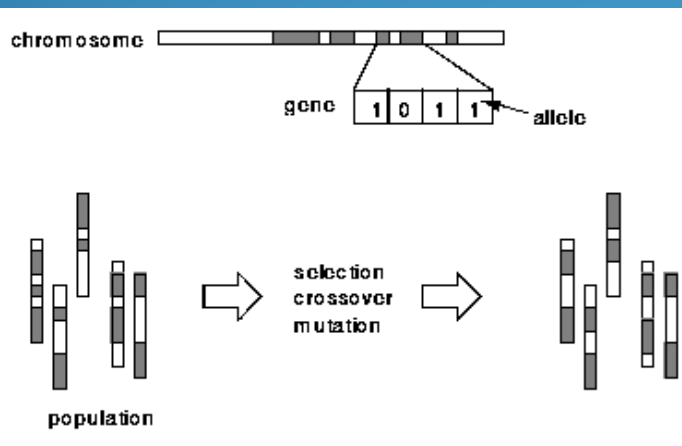
# Supervised Models: MLP & Genetic Algorithms

**Genetic algorithms** are a part of **evolutionary computing**, which is a rapidly growing area of artificial intelligence. As you can guess, genetic algorithms are inspired by Darwin's theory about evolution. Simply said, solution to a problem solved by genetic algorithms is evolved.

If we are solving some problem, we are usually looking for some solution, which will be the best among others. The space of all feasible solutions is called **search space.** Each point in the search space represent one feasible solution. Each feasible solution can be "marked" by its value or fitness for the problem. We are looking for our solution, which is one point (or more) among feasible solutions - that is one point in the search space. The looking for a solution is then equal to a looking for some extreme (minimum or maximum) in the search space. The search space can be whole known by the time of solving a problem, but usually we know only a few points from it and we are generating other points as the process of finding solution continues.



Chromosomes are strings of <u>DNA</u> and serves as a model for the whole organism. A chromosome consist of **genes**, blocks of DNA. Each gene encodes a **trait**, for example color of eyes. Possible settings for a trait (e.g. blue, brown) are called **alleles**. Each gene has its own position in the chromosome. This position is called **locus**.

Complete set of chromosomes is called **genome**.

# Supervised Models: MLP & Genetic Algorithms

$$N(U, e, P(e)) = \frac{1}{2k} \sum_j (Y_j - P(e))^2$$
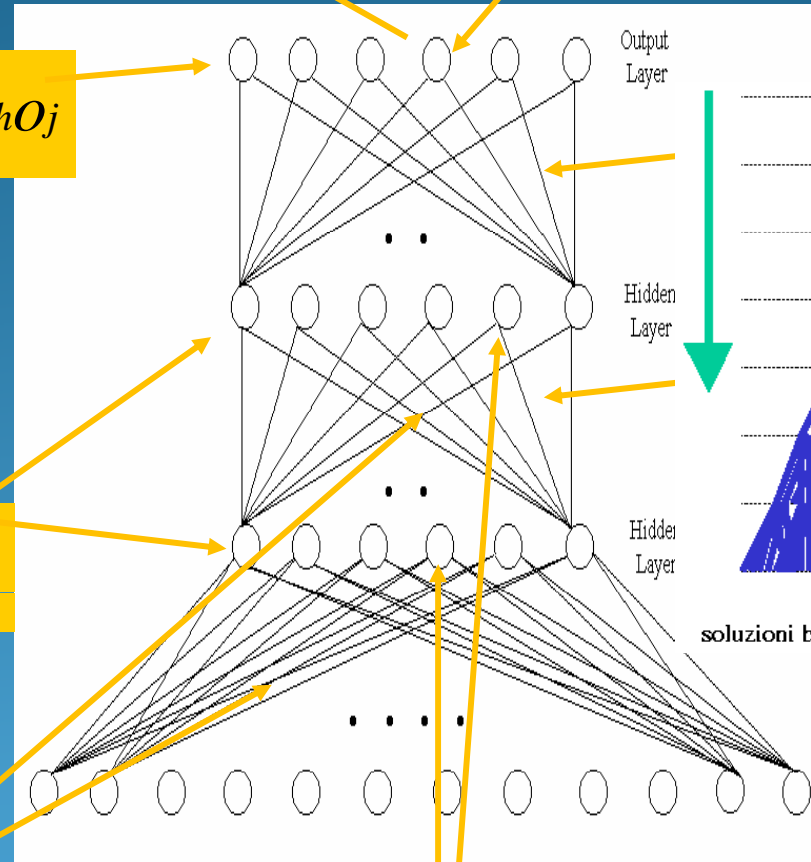
$$|Y - P(e)| < \varepsilon$$

**backward**

$$net\_input_j = \sum_h w_{jh} O_j$$

**forward**

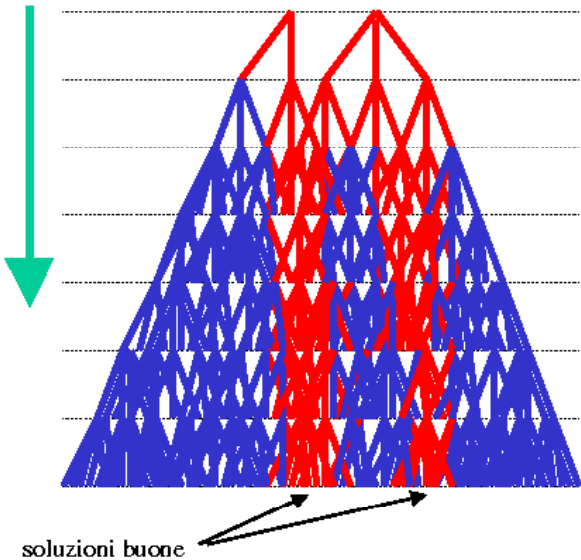$$net\_input_h = \sum_i w_{hi} O_i$$

Output Layer

Hidden Layer

Hidden Layer

soluzioni buone

$$w_{ji}(new) = w_{ji}(old) + \eta \delta_i o_j + \alpha \Delta w_{ji}(old)$$

$$f(o) = \frac{1}{1 + \exp(-o)}$$

Due to the complexity and quantity of source code (different languages, input data formats, multi-platform handling, information flow etc.), internal design standard and protocols became fundamental constraints.

Example of source code design standardization:

Supervised Models: MLP & Genetic Algorithms
UML & OOP approach

http://www.na.astro.it/~brescia/mlpga/html/index.html

# Data Mining Exploration with VO-Neural/DAME

## Unsupervised (clustering) methods

They cluster the data relying on their statistical properties only
Understanding takes place through labeling (very limited BoK).

**Generative Topographic Mapping (GTM), Self Organizing Maps (SOM), Probabilistic Principal Surfaces (PPS), Support Vector Machines (SVM), etc.**

## Pro's & Con's

- In theory they need little or none knowledge a-priori
- Do not reproduce biases present in the BoK

- Evaluation of errors more complex (through complex statistics)
- They are computationally intensive
- They are not user friendly (… more an art than a science; i.e. lot of experience required)

# Unsupervised Models: Self Organizing Maps

The SOM is an algorithm used to visualize and interpret large high-dimensional data sets

The map consists of a regular grid of processing units, "neurons". A vector consisting of features, is associated with each unit. The map attempts to represent all the available observations with optimal accuracy. At the same time vectors become ordered on the grid so that similar vectors are close to each other and dissimilar vectors far from each other.

Fitting of the model vectors is usually carried out by a sequential regression process, where $t = 1,2,...$ is the step index:
For each sample $\mathbf{x}(t)$, first the winner index $c$ (best match) is identified by the condition

$$\forall i, \|\mathbf{x}(t) - \mathbf{m}_c(t)\| \leq \|\mathbf{x}(t) - \mathbf{m}_i(t)\|.$$

After that, all model vectors or a subset of them that belong to nodes centered around node $c = c(\mathbf{x})$ are updated as

$$\mathbf{m}_i(t + 1) = \mathbf{m}_i(t) + h_{c(\mathbf{x}),i}(\mathbf{x}(t) - \mathbf{m}_i(t)).$$

$h_{c(\mathbf{x}),i}$ is the ``neighborhood function'', a decreasing function of the distance between the $i^{\text{th}}$ and $c^{\text{th}}$ nodes on the map grid. This regression is usually reiterated over the available samples.

feature map

weight matrix

Input layer

Input values

(a) Manifold in latent space R³ — ○ x
(b) Manifold in feature space Rᴰ — × t, — y(x)
(c) t projected onto manifold in latent space R³ — × E[x|t]

## NEC: a matter of Gaussians

Clustering method based on the "neg-entropy" NegE, a measure of non gaussianity of a variable. If $A$ is gaussian, then NegE($A$) = 0. Given a threshold $d$:
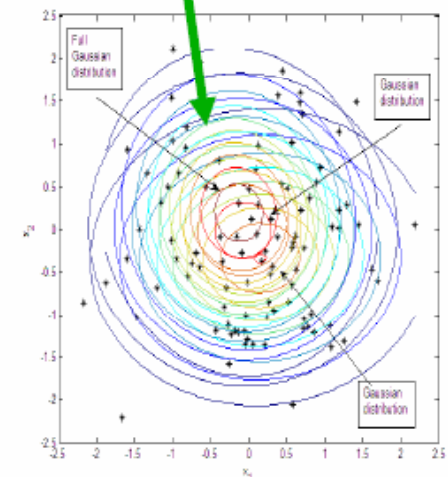
If NegE($A \cup B$) < $d$, then clusters $A$ and $B$ are replaced by cluster $A \cup B$

Not replaced!          Replaced!



NegE=750          NegE=4

## PPS: the Beauty of Spheres

The original $m$-dimensional data space is mapped in a lower $n$-dimensional space, called "latent space". Visualization ease as a spherical manifold is fitted to the data, then projected into the manifold in R³ and plotted as points on the sphere surface. Each latent variable on the sphere is responsible for a number of projected points, which form a "cluster".

# Project Target

Data Gathering (e.g., from sensor networks, telescopes…)

Data Farming:
  Storage/Archiving
  Indexing, Search ability
  Data Fusion, Interoperability

**Database technologies**

Data Mining (or Knowledge Discovery in Databases):
  Pattern or correlation search
  Clustering analysis, automated classification
  Outlier / anomaly searches
  Hyper-dimensional visualization

**Key mathematical issues**

Data understanding
  Computer aided understanding
  KDD
  Etc.

**Ongoing research**

New Knowledge

# Project Kick-off



*In 2007, a group of astronomers, computer scientists, engineers and physicians started to explore possible joined effort to create a data mining toolset, based on a distributed infrastructure, for worldwide users who want to share data, methods and discoveries.*

- **astronomy**: problems, data, understanding of the data structure and biases
- **statistics**: evaluation of the data, falsification/validation of theories/models, etc.
- **computer science**: implementation of infrastructures, databases, middleware, scalable tools, etc.

# Project Team (now)

# Project team (history)

| NOME | RUOLO | TOT MESI | DA | A | POSIZIONE | BACKGROUND |
|------|-------|----------|-----|-----|-----------|------------|
| Longo | P.I. | 28 | | | staff | astrofisica |
| Brescia | P.M. | 28 | | | staff | informatica/astrofisica |
| Djorgovski | SCIENCE+ICT | 28 | | | staff | astrofisica |
| Corazza | SCIENCE+EDU | 22 | 13/06/07 | | staff | informatica |
| Donalek | AI+SW | 28 | | | staff | informatica |
| D'abrusco | SCIENCE | 28 | | | post-doc | fisica |
| Laurino | P.E. | 28 | | | tesi laurea | fisica |
| Garofalo | SW ENG | 28 | | | tesi laurea | ingegneria informatica |
| Nocella | SW ENG | 28 | | | tesi laurea | ingegneria informatica |
| Cavuoti | SCIENCE+SW | 28 | | | contratto | fisica |
| d'Angelo | WEB+DOC+SW | 28 | | | contratto | fisica |
| Manna | SW ENG | 7 | 24/09/08 | | tirocinio | informatica |
| Fiore | SW ENG | 6 | 22/10/08 | | tirocinio | informatica |
| Di Guido | SW ENG | 2 | 02/03/09 | | tirocinio | informatica |
| Deniskina | SW ENG | 12 | 06/09/07 | 04/09/08 | contratto | informatica |
| Skordovski | SW | 19 | 13/06/07 | 31/12/08 | tesi laurea | informatica |
| Russo | SW ENG | 10 | 13/06/07 | 15/04/08 | tesi laurea | informatica |
| Vaccaro | SW ENG | 4 | 13/06/07 | 25/09/07 | tesi laurea | informatica |
| Formicola | SW | 6 | 20/02/08 | 07/07/08 | tesi laurea | informatica |

# Project Management Highlights

## U.M.L. Tools for Suite Functionalities and internal code design



## XP – eXtreme Programming as project life cycle



Copyright 2000 J. Donvan Wells

| PROJECT DOCUMENTATION PRODUCTION | | | | | | |
|---|---|---|---|---|---|---|
| SECTION | CODE | 2007 | 2008 | 2009 | 2010 | MEANING |
| **MANAGEMENT** | SOW | 2 | | | | Statement Of Work |
| | PLA | 1 | 1 | | | Project Plan & Work Breakdown Structure |
| | MIN | 12 | 22 | 8 | | Minutes of Meeting |
| | SCH | 1 | 1 | | | Project schedule |
| | LIS | 12 | 22 | 8 | | Action list |
| | DOC | 8 | 5 | 2 | | Documentation & Website status |
| **DESIGN** | SPE | 8 | 1 | | | Technical Specifications |
| | PDD | | 1 | | | Project Description Document |
| | SRS | | 5 | | | Software Requirement specification |
| | SDD | | | 5 | | Software Design Description |
| **DEVELOPMENT & TEST** | TRE | 5 | 7 | | | Technical Report |
| | PRO | 3 | 1 | 5 | | Test Procedure |
| | VER | 2 | | | | Test Verification & report |
| **RELEASES** | MAN | | 2 | | | User Manual |
| | IDM | | | | | Installation & Deployment Manual |
| | REN | | | | | Release Notes |
| **PUBLICATIONS** | TECH | | 3 | | | Technical papers |
| | SCIENCE | | 5 | 2 | | Scientific papers |
| **TOTALS** | | 54 | 76 | 30 | | 160 |
| **PRESENTATIONS** | PRE | 3 | 2 | 4 | | Meeting talk/poster |

## Project Milestones

| PROJECT MILESTONES | 2007 | 2008 | 2009 | 2010 |
|---|---|---|---|---|
| Statement of work | | | | |
| WBS & Project Plan | | | | |
| Project Design Description | | | | |
| Software Requirement Specifications | | | | |
| DM models Implementation | MLP | SVM | MLPGA NEXT | SOM PPS NEC |
| DM models scientific validation | | | | |
| Software Design Description | | | | |
| Implementation & Test Procedures | | | | |
| Technical Reports | | | | |
| Test Reports | | | | |
| Beta Release Deployment | | | | |
| User & Maintenance Manuals | | | | |
| Commissioning of beta release | | | | |
| Official release 1.0 | | | | |
| New DM models Implementation | | | | |
| New Functionalities Implementation | | | | |

**FRONT- END DEPLOYMENT ENVIRONMENT**

**FRONT- END**
Web Application

VirtualFileStore

**FRAMEWORK DEPLOYMENT ENVIRONMENT**

**GRID** *User Interface*

Protocol

(XML)

**FRAMEWORK**
(restful Web service)

DM Plugins

**DR**
Driver Management System for:
Storage and Execution
Infrastructure Abstraction
Data File Conversion

**REDB**
DBMS for:
Registration
Authentication
Session

**EXECUTION ENVIRONMENT**

**GRID** *Infrastructure*

**DMM**
Data Mining Models
MLP
SVM

**Worker Nodes**

FileStore (Storage Element)

DataBase

**Object Oriented Programming & UML**
**Internal standards and protocols (XML)**
**Java language (generic for DMM)**
**User/Session Registry DB (MySQL)**
**Web-based User I/O**
**Web Application and Web Service Technology**
**Plugin Modularity (easy to be integrated/modified)**
**Hardware independent through platform driver**
**Data conversion and manipulation support**

Logical Flux

Data Trasmission

Protocol Trasmission

*LEGENDA*

# FRONT END Component



**Architecture**:
- Client-server AJAX (Asynchronous JAva-Xml) based;

**Technology**:
- GWT-EXT;

**Features**:
- User GUI deployment and I/O management;
- interaction with internal components through standard protocol (XML);
- Local User/Session data virtualization through Virtual File Store;

## GWT

The **Google Web Toolkit** is an open source toolkit to create client-side applications in Java. GWT compiler translates a Java application into equivalent JavaScript that manipulates a web browser. GWT emphasizes reusable, efficient solutions to asynchronous remote procedure calls, history management and cross-browser portability.

## EXT

Ext is an open-source JavaScript library, for building richly interactive web applications using techniques such as AJAX scripting. Ext JS is an excellent framework for building web applications that have desktop-like functionality in a web browser.

## GWT-EXT

GWT-EXT is a library integrating GWT and EXT. One of the primary goals is to make the GWT-Ext widgets and API's work seamlessly with the core GWT infrastructure and its API's

# DR Component

**Architecture**:

• It depends on the environment choice;

• In S.Co.P.E. DR is a component running on the GRID UI;

**Technology (in S.Co.P.E.)**:

• GRID Software (middleware gLite);

**Features**:

• Storage Device(s) + Execution Environment = Deployment Environment;

• Different Deployment Environments can be more suited for a specific task (e.g. an MLP TEST is unlikely to be a computing intensive task, so GRID latency times are not needed);

• Dynamic Driver Loading => Driver Plugins;

• Drivers are available to the Framework WS and to the Plugins;

• Also used to convert files formats (standard or DMM dependent);



DR translate schema

FITS → DMM format → FITS
VOTable → ASCII
ASCII → CSV
CSV



**Computing Environment**

SA
• Process Environment (CPU)
• FS (HD)

GRID
• Process Environment (WN)
• FS (SE)

OS

• File Sistem interaction
• Execute job
• Monitoring job
• Retrive results

Glite Middleware (CML - API)

• SE interaction
• Submit job
• Monitoring job
• Retrieve results

**DRMS** (Library)

Device
• Processing
• FS

Translator

implements    use

Interface
• processing
• FS
• translating

• Multiple deployments
• Computational methods
• FS methods
• Translating methods

• Manage process environment
• Use different platform
• Access process resources
• Use translating

FW

*Legend*

**Green arrow:** indicates the up-down connections
**Blue arrow:** indicates the bottom-up connections

*Acronyms*

**CPU:** Central Processing Unit
**DRMS:** Driver Management System component
**FS:** File Store
**FW:** Framework component
**HD:** Hard Disk
**SE:** Storage Element
**WN:** Worker Node
**CML:** Command Line
**API:** Application Programming Interface

# DMM Component

## DMPlugin

### DM Functionalities
Classification, Regression, ...

### DM Models
SVM, MLP, PPS, ...

### DM Library wrappers
JNI, SWIG, ...

### DM Libraries
libfann, libsvm, ...

### Low Level Libraries
blas, lapack, gsl, ...

**Architecture**:
- data mining functionality class hierachy;

**Technology**:
- available model packages and libraries;
- custom ad hoc model design and development;
- custom wrappers for internal standardization;

**Features**:
- modularity;
- fast third part application integration;
- functionality specialization;
- multi-language programming support;

# REDB Component

**Architecture**:
• JDBC;

**Technology (in S.Co.P.E.)**:
• MySQL and JDBC API;

**Features**:
• management of user (registration, authentication, working sessions, experiments and files) information and their relationships;
• store and manage information about three different file's categories: "supported", "exotic" and "custom" (datasets, model configuration and intermediate data);

# FRAMEWORK Component

## Architecture:
• Restful Web Service (client-server apps with resource addressable with HTTP methods);
• DM models control interface through Plugin SDK;

## Technology:
• Web container SUN Apache Tomcat;
• Java Servlet for web service;

## Features:
• Internal resource representation through "contextual" VOTables;
• Experiment configuration and execution;
• user authentication and working session management;
• experiment data & working flow trigger and supervision;
• XML based internal communication protocol

## User/Developer Perspective

**A simple user can upload and build his datasets, configure the data mining models available, execute different experiments in service mode, load graphical views of partial/final results.**

**You are not considering yourself as a simple user? Ok, so you think to be a developer. Or at least a scientist who wants to upload and use his application (and possibly to share it with others).**

**Be honest, you don't trust someone else's application.**

**So You want to extend our framework?**

*YES, WE CAN!*

**DM Models Development**
Download our DM Models library;
Add new low level/DM shared libraries and related new wrapper;
Extend the DM class hierarchy;
**Plugin Development**
Download our SDK;
Implement and test the DMPlugin abstract class;
Provide a method to produce the plugin description and Submit for Registration;
The same if you want to develop a new driver for a specific environment or storage system. Just implement the Driver Plugin Interface and register it;

# Application Prototype

**voneural.na.infn.it**

## Application Prototype

**New user registration**

### DAME - DAta Mining and Exploration

- Home
- Sign Up!
- Help & Tutorials
- The Team

## Create an account

First name:

Last name:

Username:

Email address:

Password:

Password again

Click when finished: Register →

Fill out the form to the left (all fields are required), and your account will be created; you'll be sent an email with instructions on how to finish your registration.

We'll only use your email to send you signup instructions. We hate spam as much as you do.

This account will let you subscribe to event streams for future notifications.

# Application Prototype

**logging in**

## DAME - DAta Mining and Exploration

Username
brescia
Password
••••••••
Log in

Home

Sign Up!

Help & Tutorials

The Team

### What is DAME

**DAME** is a web application to perform data mining on massive data sets. In order to ensure scalability it allows the user to access distributed computing facilities provided by the Center for Advanced Research in Computing at Caltech and by the **S.Co.P.E.** project at the University of Napoli Federico II. DAME is derived from the **VO-Neural** project.

As a function of the size and complexity of your task, your computation will be re-directed to larger computing facility.

**DAME** is an evolving platform. Therefore please provide us with your comments and feedbacks.

Start signing up for a new account. Signing up will provide you with a **persistant filestore** on our servers, so that you won't need to upload your datasets each time you want to perform a new calculation.

Your filestore will also contain all the **output files from the experiments you launch**, so that you can visualize or download them when the experiment is done.

During an experiment you can **visualize the log file** showing the status of the experiment and visualize output files. You can also **abort a calculation**.

You can even **download an entire directory** in a compressed zip archive on your hard disk. Output files can be used as inputs for other experiments, and so on...

In the "Help & Tutorials" section you will find **documentation, examples and tutorials**. The first time you login, your filestore will contain some datasets you can use following the tutorials.

Show input
dataset

# Application Prototype

Dame - DAta Mining and Exploration - Mozilla Firefox

File  Modifica  Visualizza  Cronologia  Segnalibri  Yahoo!  Strumenti  ?

http://pcdevauc.na.infn.it:9000/filestore/download/brescia/Samples/provapiccolo.csv/

Google

Più visitati  Come iniziare  Ultime notizie

Yahoo!  Search Web  Mail  Shopping  Personals  My Yahoo!  News  Games  Travel  Finance  Answers  Sports  Sign In

sitemap page  |  VO-Neural Home Page  |  Caricamento in corso...  |  https://imap-ac.n.../src/webmail.php

```
1.692768,0.989927,0.401751,0.315475,0.113712
1.705183,0.972942,0.462513,0.299402,0.104813
1.81883,0.832718,0.387196,0.282825,0.076243
1.754013,0.897587,0.411005,0.295853,0.083096
1.750496,0.987221,0.452448,0.28212,0.135608
2.016077,0.979092,0.447763,0.334382,0.105586
1.725897,1.089231,0.448336,0.378626,0.124194
1.946756,0.970669,0.391537,0.341539,0.105292
1.886108,0.828356,0.398804,0.277565,0.116535
1.963732,0.862139,0.435151,0.337962,0.049855
1.934912,0.979206,0.440548,0.316557,0.114727
1.664152,1.160151,0.517653,0.369753,0.190644
1.809011,1.250761,0.46133,0.358436,0.193634
1.668617,0.915407,0.381332,0.283329,0.143155
1.6555,1.079199,0.450472,0.372208,0.146099
2.03433,0.889687,0.475706,0.338007,0.02892
1.976236,1.048838,0.512547,0.389719,0.12446
2.004997,0.926264,0.402496,0.328529,0.092687
1.87063,0.944035,0.406311,0.32659,0.117726
1.85619,0.996578,0.398401,0.324921,0.117244
2.042259,0.889086,0.421186,0.305461,0.075189
1.971815,0.834901,0.380198,0.290215,0.029776
1.991796,0.981928,0.44222,0.350865,0.103254
1.818487,1.129131,0.449316,0.311396,0.146972
1.797421,0.867332,0.405901,0.326082,0.068398
1.838831,0.746513,0.384384,0.292418,0.068517
1.71731,1.013191,0.414141,0.355978,0.136164
1.750128,0.869583,0.403996,0.245794,0.082915
1.913464,1.039413,0.480566,0.416491,0.118118
2.044985,1.069139,0.447777,0.330549,0.148072
1.825377,0.896421,0.367558,0.299021,0.092071
1.959181,0.947929,0.4189,0.325193,0.09667
1.887531,1.15271,0.464077,0.334187,0.166015
1.66408,0.889212,0.481047,0.372786,0.110893
2.00355,0.890563,0.418589,0.293733,0.062144
1.677416,1.177158,0.472656,0.304817,0.209245
```

**Beta release feature:**

**integrated dataset editor
and builder**

Trasferimento dati da pcdevauc.na.infn.it...

catania_febbraio2009  |  Dame - DAta Minin...  |  2009_tiruvalla_India....  |  voneural-dame   IT  14.07

# Application Prototype

**Check /Edit past experiments**

## DAME - DAta Mining and Exploration

**Massimo Brescia**
Last Login
Thu 02 Apr 2009 11:01AM GMT

Home
MyFilestore
MyExperiments
Logout
Help & Tutorials
The Team

Launch Experiments
New MLP
New SVM
New PhotoZ

### My Experiments

**Experiments List**

| Name | Science case | Mode | Status | Actions |
|------|-------------|------|--------|---------|
| myExperiment | mlpclassification | mlptrain | finished | Remove |

### My Filestore

Sfoglia...

Click here to upload the file

| Dirs | Files | Actions |
|------|-------|---------|
| /brescia | | |
| | No data in this directory! | |
| /brescia/Samples | | Delete Download |
| | iris.dat | Delete |
| | provapiccolo.csv | Delete |
| /brescia/myExperiment | | Delete Download |
| | provapiccolo.csv.fits | Delete |
| | myExperiment.csv | Delete |
| | myExperiment.log | Delete |
| | myExperiment.ERROR | Delete |

# Application Prototype

## Launch new experiments

## Application Prototype

### DAME - DAta Mining and Exploration

**Massimo Brescia**
Last Login
Thu 02 Apr 2009 11:01AM
GMT

- Home
- MyFilestore
- MyExperiments
- Logout
- Help & Tutorials
- The Team

Launch Experiments

New MLP

**Experiment Configuration**

Science case: you can choose whether to classify labeled patterns or to find a regression mapping from examples. You can find documentation here.

Science Case: Regression
- Classification
- Regression

Mode: choose the mode you want to run. If you don't know what these modes mean and how they work, refer to this link.

Mode: Train a New Net

Go!

**Select functionality**

**Select use case and launch the job**

### DAME - DAta Mining and Exploration

**Massimo Brescia**
Last Login
Thu 02 Apr 2009 11:01AM
GMT

- Home
- MyFilestore
- MyExperiments
- Logout
- Help & Tutorials
- The Team

Launch Experiments

New MLP

**Experiment Configuration**

Science case: you can choose whether to classify labeled patterns or to find a regression mapping from examples. You can find documentation here.

Science Case: Regression

Mode: choose the mode you want to run. If you don't know what these modes mean and how they work, refer to this link.

Mode: Train a New Net
- Train a New Net
- Test a Trained Net
- Run
- Full (Train + Test)

# Application Prototype

**Edit and submit model and experiment parameters**

## DAME - DAta Mining and Exploration

**Massimo Brescia**
Last Login
Thu 02 Apr 2009 11:01AM
GMT

- Home
- MyFilestore
- MyExperiments
- Logout
- Help & Tutorials
- The Team

Launch Experiments

- New MLP
- New SVM
- New PhotoZ

### Experiment Configuration

Name. This is the name that will be associated with the experiment. Be sure the name is meaningful to you. When the experiment is done, you will find your files in a directory in your filestore named after your experiment.

| Experiment name: | iris_exp |

Input Nodes. It is the number of input features. If N is the number of input features and M the number of target components, then the training set must have exactly N+M columns.

| Input nodes: | 4 |

**Hidden Nodes Help**

| Hidden nodes: | 3 |

**Output Nodes Help**

| Output nodes: | 3 |
| Max epochs: | 40000 |
| Tolerance: | 1e-05 |
| Training algorithm: | MSE - BATCH |
| Resume training: | ☐ |
| Network: | /Samples/iris.dat |
| Training set: | /Samples/iris.dat |
| Do validation: | ☐ |
| Validation set: | /Samples/iris.dat |

Go!

# Application Prototype

## DAME - DAta Mining and Exploration

**Massimo Brescia**
Last Login
Thu 02 Apr 2009 11:01AM
GMT

**Home**

**MyFilestore**

**MyExperiments**

**Logout**

**Help & Tutorials**

**The Team**

Launch Experiments

New MLP

New SVM

New PhotoZ

### Experiment Details

**Experiment Name: iris_exp**

Started

| Parameter | Value |
|---|---|
| Input Nodes | 4 |
| Hidden Nodes | 3 |
| Output Nodes | 3 |
| Max Epochs | 40000 |
| Tolerance | 1e-05 |
| Training Algorithm | mseBatch |
| Training Set | /brescia/Samples/iris.dat |

| Dirs | Files | Actions | |
|---|---|---|---|
| /brescia/iris_exp | | Download | |
| | iris_exp.ERROR | Delete | |
| | iris_exp_netTrain.mlp | Delete | |
| | iris.dat.fits | Delete | |
| | iris_exp_netTmp.mlp | Delete | |
| | iris_exp.csv | Delete | |
| | iris_exp.log | Delete | |
| | iris_exp.tra | Delete | |

### Experiment Log

```
MLP

Executing option: TRAIN
Input nodes: 4
Output nodes: 3
Nodes in hidden layer: 3
Maximum epochs: 40000
Problem case: Regression
Training algorithm: Batch
Error: MSE
Error tolerance: 1e-05
Input network name: empty
Training dataset: iris_exp/iris.dat.fits
Validation dataset: empty
```

### Plots

**Status during execution**

**Beta release feature:**

**Interactive session optional during execution**

## Our first scientific use cases

**First example**
evaluation of SDSS redshift using supervised NN (MLP)

*Mining the SDSS Archive I. Photometric redshifts in the nearby Universe*, R. D'Abrusco et al. (The Astrophysical Journal, 663: 752-764, 2007 July 10.

**Second example**
Searching for candidate quasars in the SDSS archive

astro-ph/0805.0156v1; to appear soon in MNRAS (R. D'Abrusco et al.)
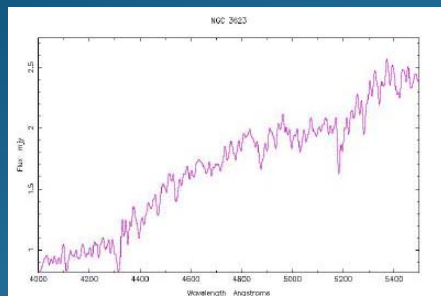
**Third example**
Classifying AGN in SDSS with SVM

Cavuoti 2008, Thesis (VONeural website, voneural.na.infn.it)
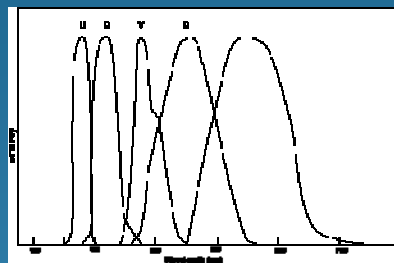
## More infos on WEB site documentation page
### http://voneural.na.infn.it/documents.html

# Science case: Mining the SDSS archive



**Galaxy spectrum - F($\lambda$)**

X



**Photometric system - S$_i$($\lambda$)**

=

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

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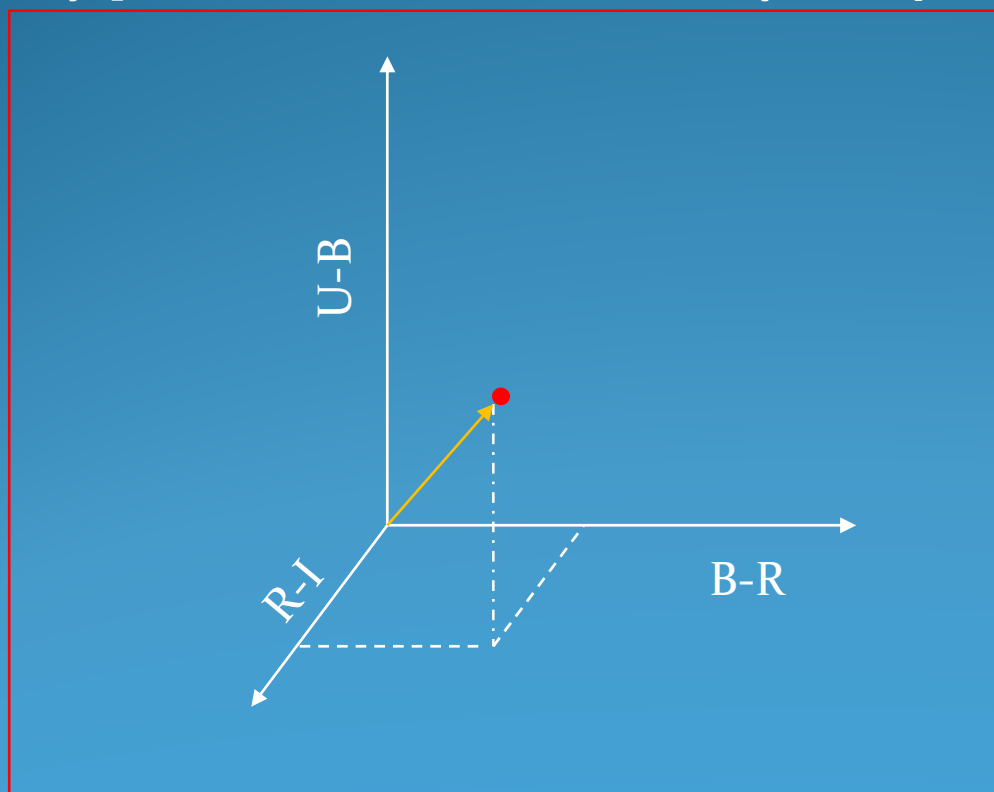
Etc...

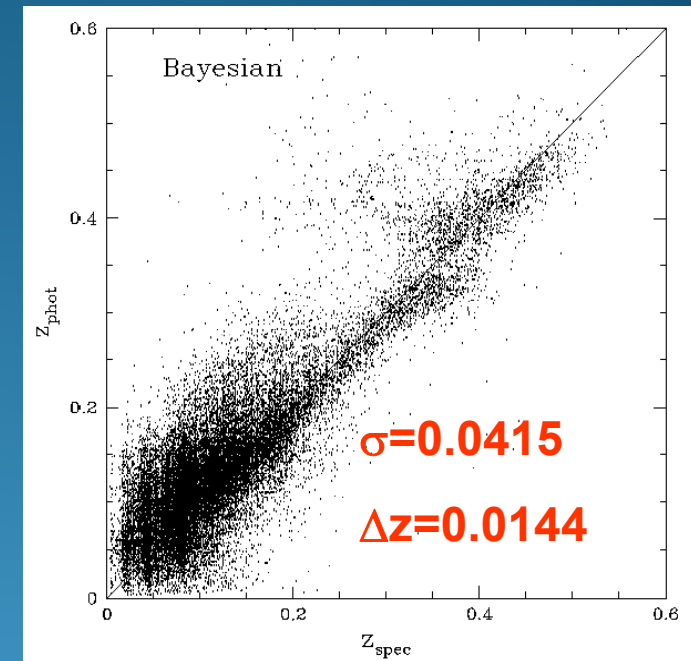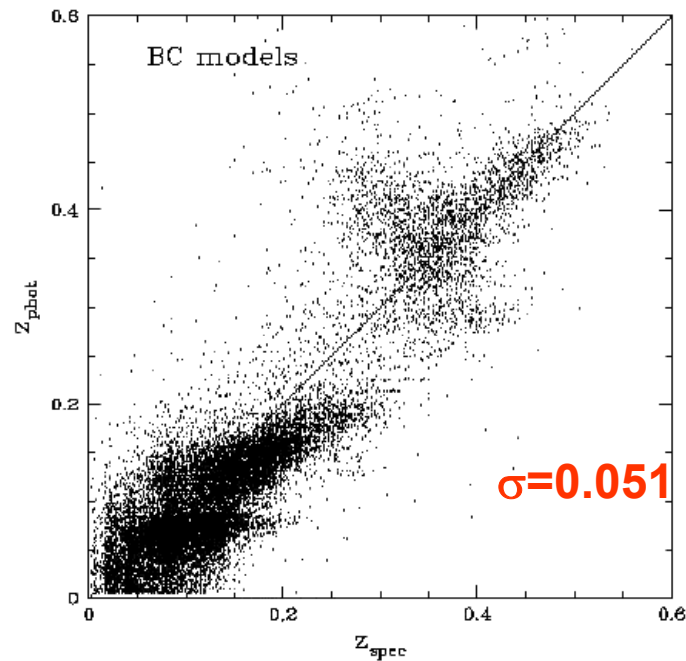Color indexes

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

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

$$etc.$$



U-B

R-I

B-R

- **SED template fitting methods**
- Interpolative methods

# Photometric Redshifts



BC models

$\sigma$=0.051



Bayesian

$\sigma$=0.0415

$\Delta$z=0.0144

| type | method | data | $\Delta z_{rms}$ | Notes | Reference |
|------|--------|------|------------------|-------|-----------|
| | CWW | EDR | 0.0666 | | (Csabai et al. 2003) |
| SEDF | 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) |
| | ANNz | EDR | 0.0229 | | (Collister & Lahav 2004) |
| ML | SVM | EDR | 0.027 | | (Wadadekar 2004) |
| ML | MLP-feed forward | SDSS-DR1 | xx.xxx | yes | (Vanzella et al. 2003) |
| | | SDSS-RLG | | | |

- **SED template fitting methods**
- Interpolative methods

# Photometric Redshifts

- 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 fit a second order polynomial
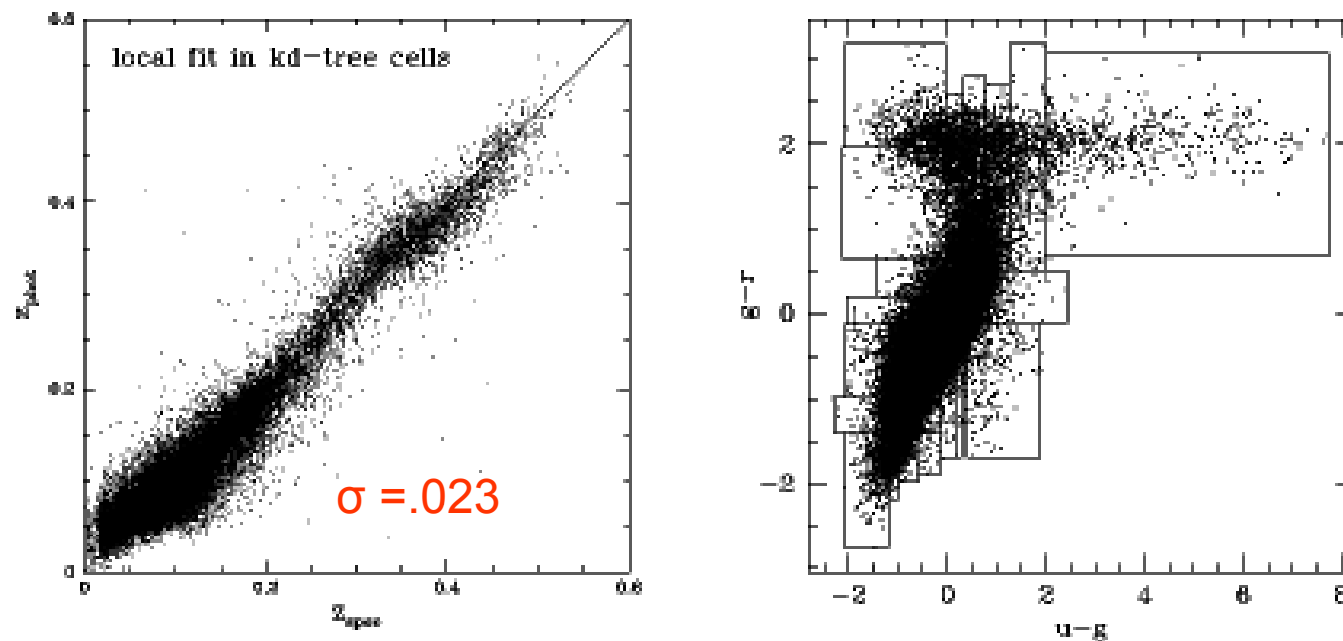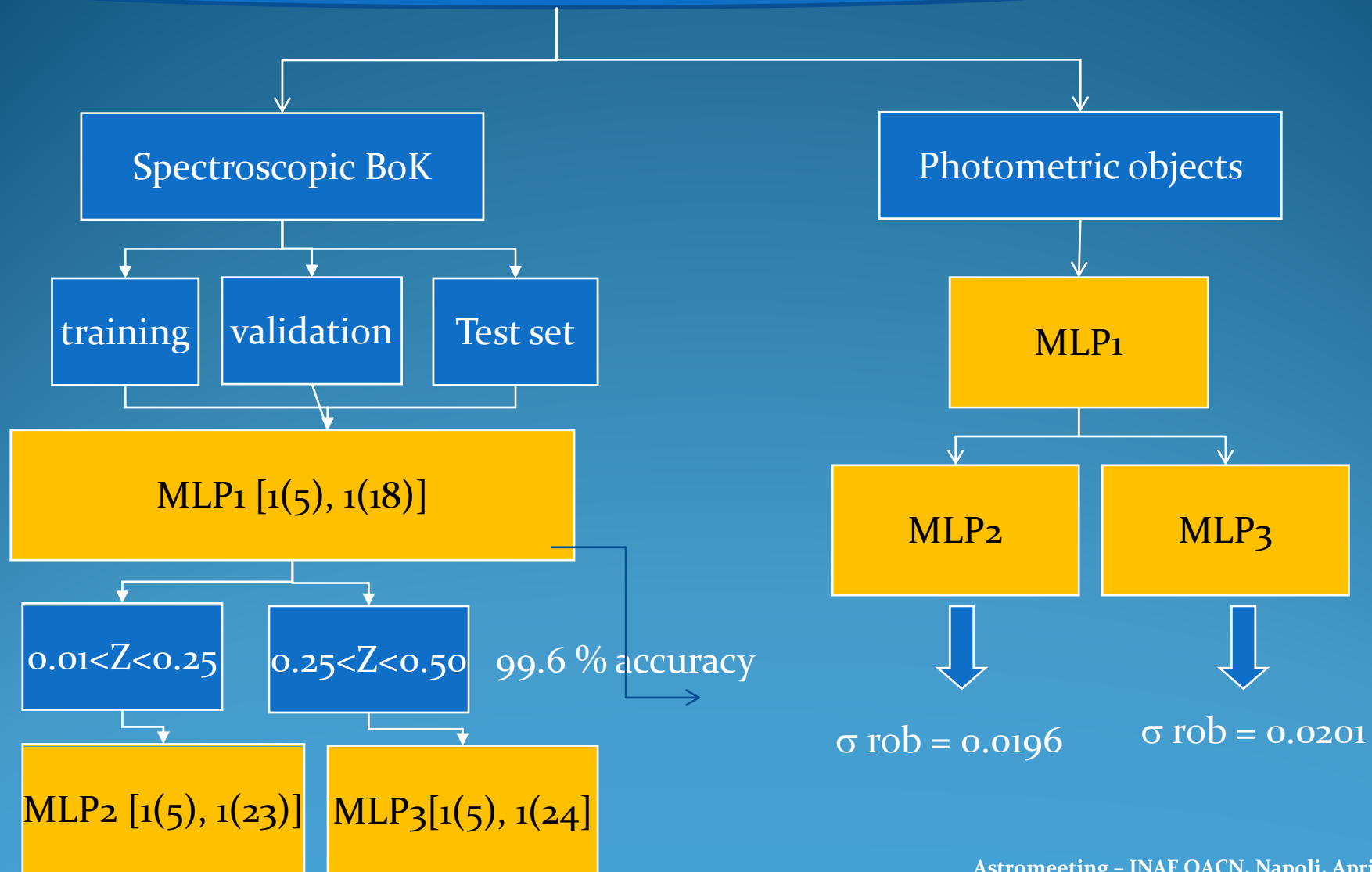


local fit in kd-tree cells

$\sigma = .023$

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.
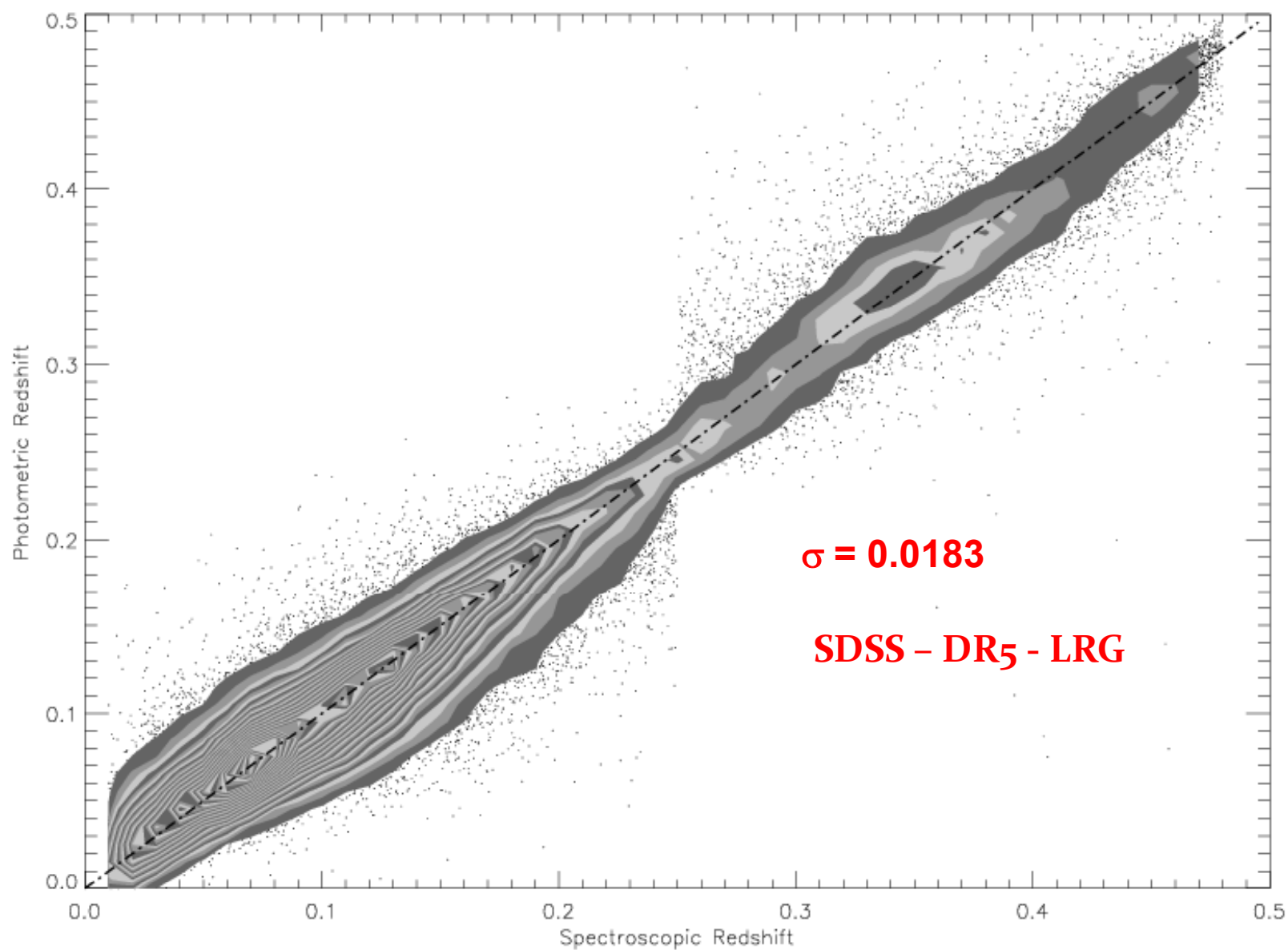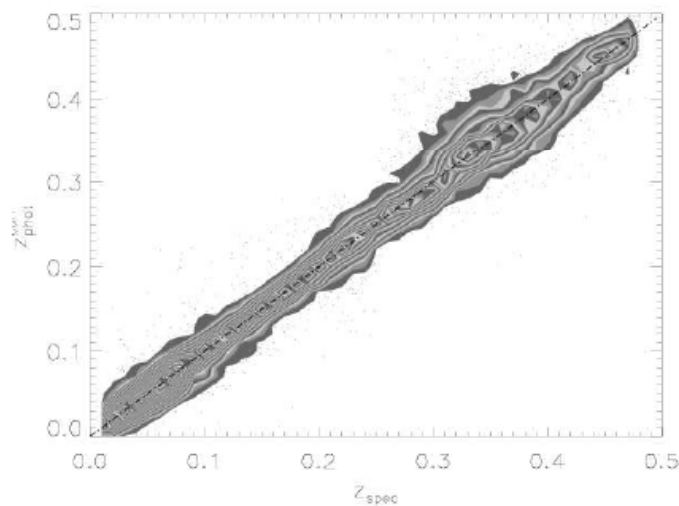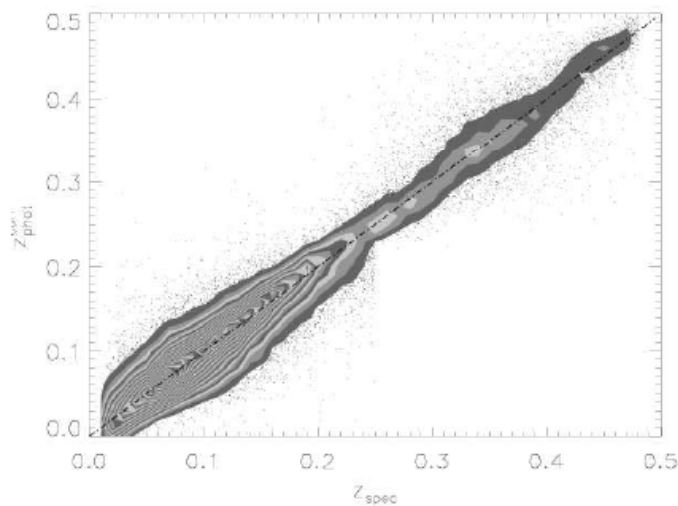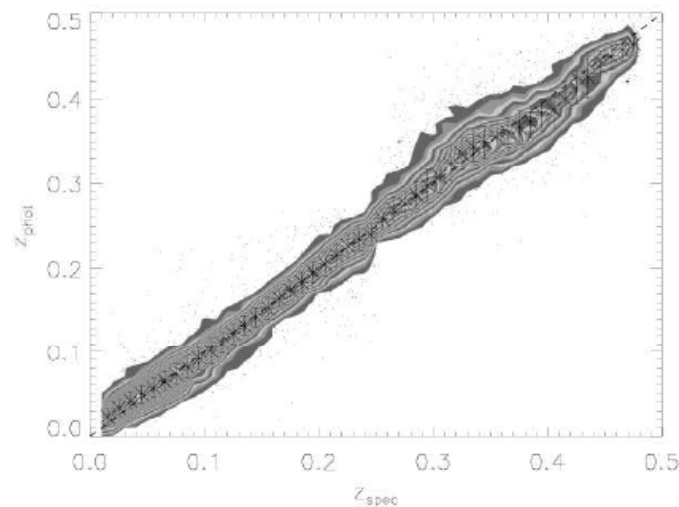
# Science case: Photometric Redshifts



SDSS – DR5

Spectroscopic BoK

training | validation | Test set

MLP1 [1(5), 1(18)]

0.01<Z<0.25 | 0.25<Z<0.50 | 99.6 % accuracy

MLP2 [1(5), 1(23)] | MLP3[1(5), 1(24)]

Photometric objects

MLP1

MLP2 | MLP3

σ rob = 0.0196 | σ rob = 0.0201

# Science case: Photometric Redshifts



$\sigma = 0.0183$

SDSS – DR5 - LRG

# Science case: Photometric Redshifts

General galaxy sample                              LRG sample



Non LRG only

$\sigma = 0.0363$
$\Delta z = -0.0030$

$\sigma = 0.0208$
$\Delta z = -0.0029$

$\sigma = 0.0178$
$\Delta z = -0.0011$

# Science case: Photometric Redshifts
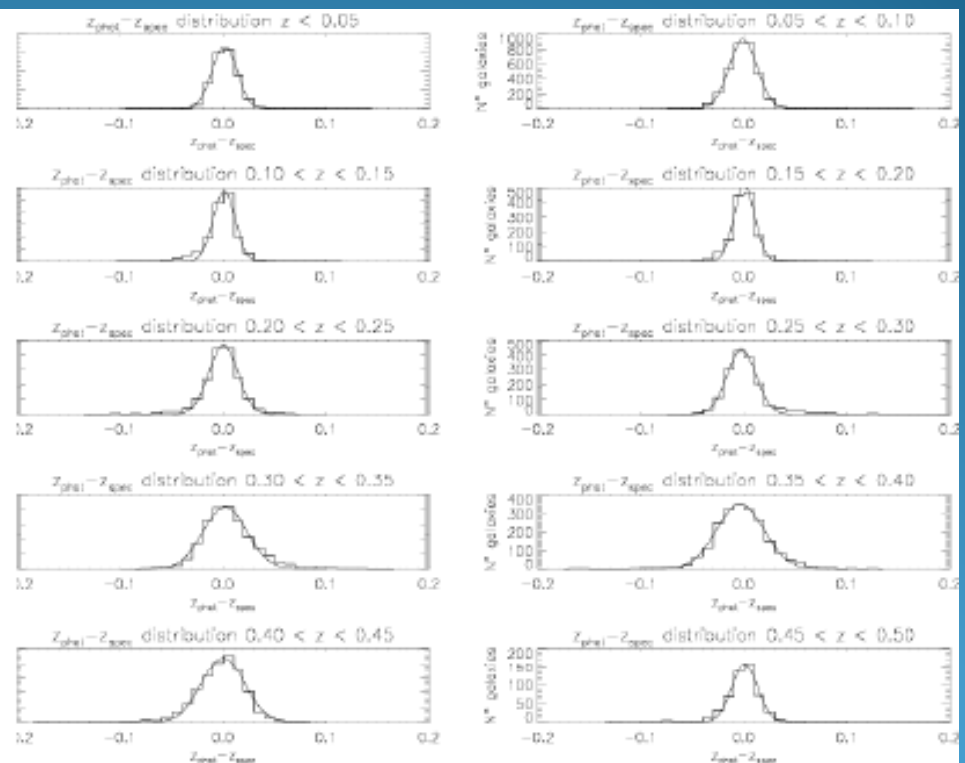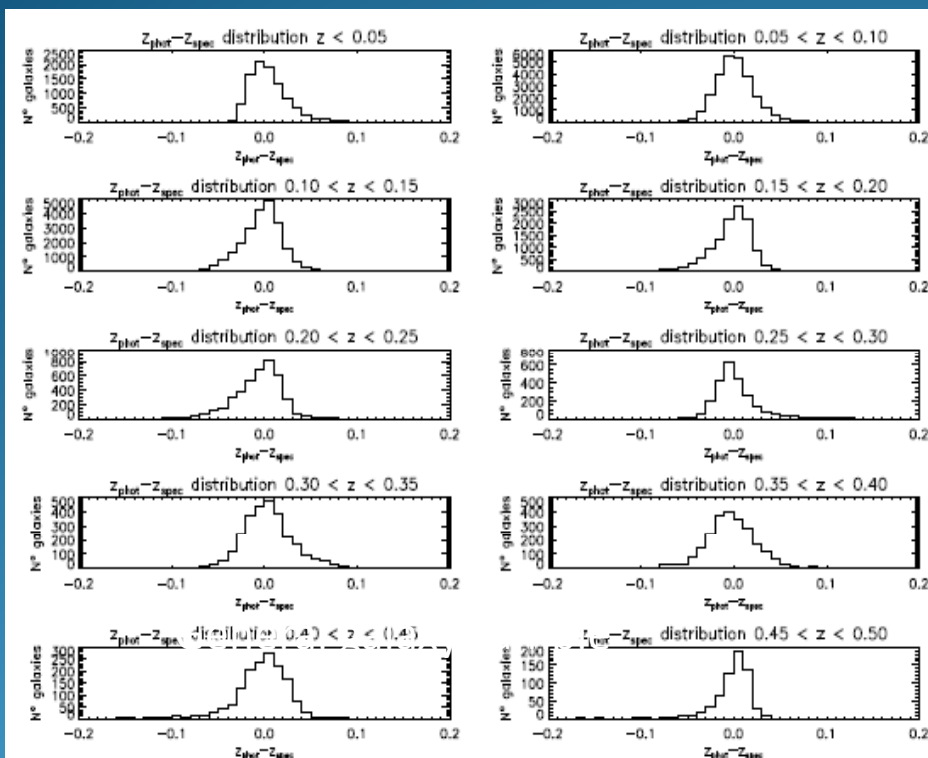


$\sigma = 0.0170$

# Science case: Photometric Redshifts



FIG. 9.— Same as in previous figure but for the LRG sample.

# SDSS galaxies $z_{phot}$

- Generalization of the approach described in the previous paper (D'Abrusco et al. 2007) will be presented at AAS 2009

- K-means algorithm for clustering in the photometric parameter space is applied with an optimal number of clusters $n_{opt}$

- $n_{opt}$ is chosen so that the "weighted accuracy" $\sigma_w$ of the $z_{phot}$ is maximum. Given N clusters with $M_i$ elements each and rms of the $(z_{phot}-z_{spec})$ variable $\sigma_i$ :
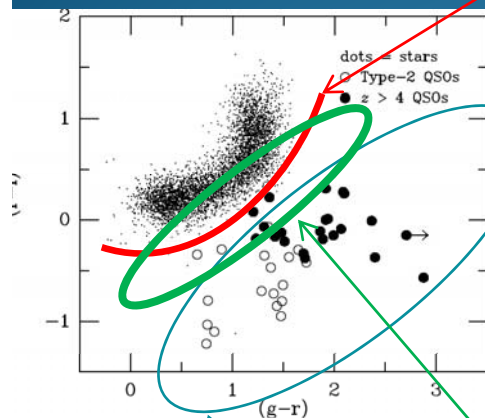
where

$$\sigma_W = M_{tot}^2 \sum_{i=1}^{N} \frac{\sigma_i}{M_i^2}$$

$$M_{tot} = \sum_{i=1}^{N} M_i$$

# Searching for candidate quasars in the SDSS
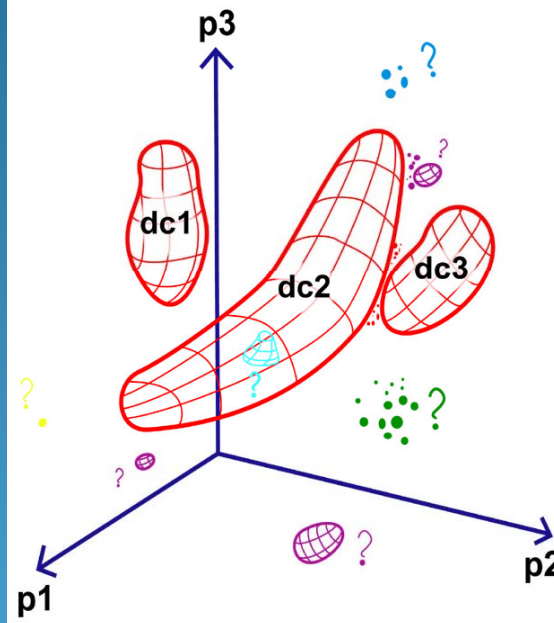
Traditional way to look for candidate QSO in 3 band survey

In 4 bands degeneracy is partially removed

Cutoff line



Candidate QSOs for spectroscopic follow-up's

**Ambiguity zone**

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



R. D'Abrusco

**astro-ph/0805.0156v1**

More are the bands the lower is the degeneracy

How to find the interesting regions (clusters)?
- Data Mining is the answer

How to visualize them ?
- Dimensionality reduction

# Searching for candidate quasars in the SDSS

Several algorithms for "general purpose" photometric identification of candidate QSOs select sources according to different techniques exist.

• Optical surveys: looking for counterparts of strong radio sources (but only ~ 10% of QSO are radio-loud).

• Ultraviolet and optical surveys: looking for star-like sources bluer than stars.

• Multi-colour surveys: looking for star-like objects in colour parameter space lying outside compact regions ("star locus") occupied by stars.
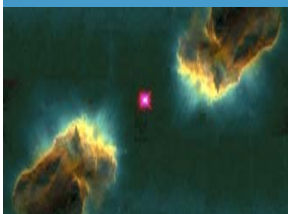
**Overall performances of a generic targeting algorithm are usually expressed by two parameters:**

**Completeness**

$$c = \frac{\text{candidate quasars identified by the algorithm}}{\text{a priori known quasars}}$$
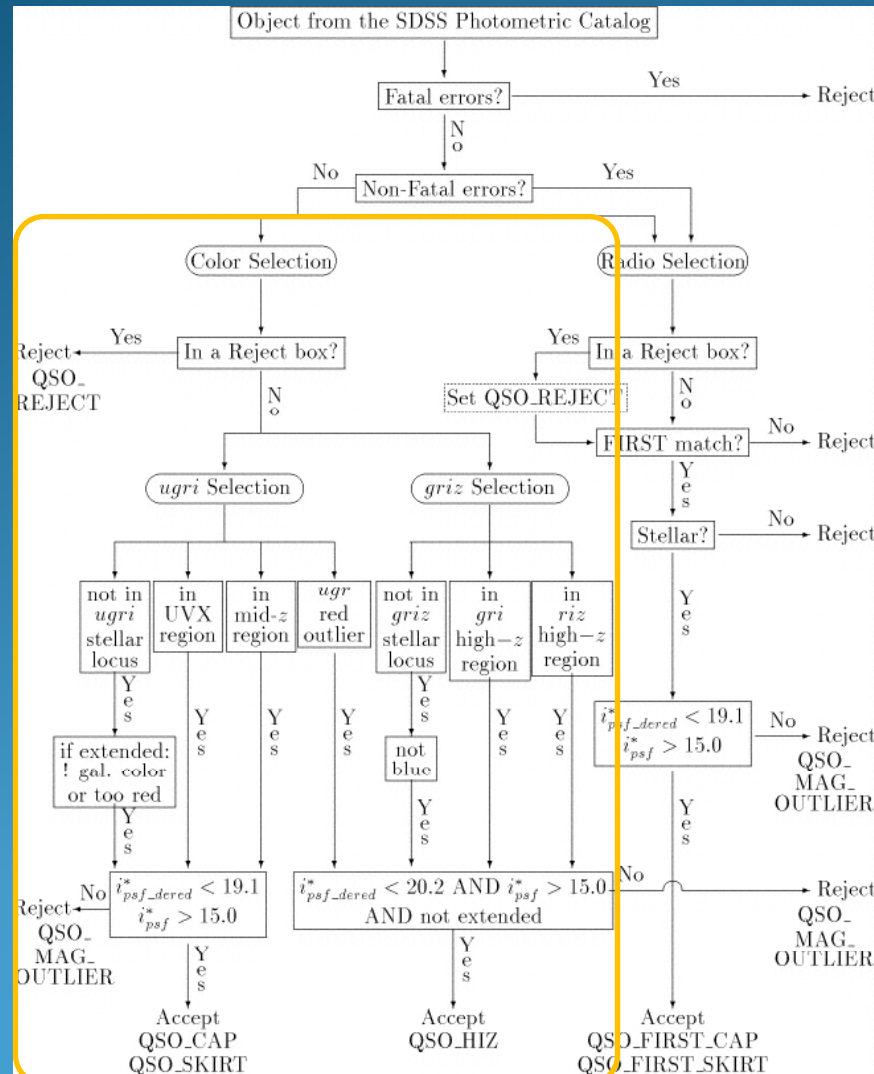
**Efficiency**

$$e = \frac{\text{confirmed quasars identified by the algorithm}}{\text{candidate quasars selected by the algorithm}}$$

# SDSS QSOs targeting algorithm (I)

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:
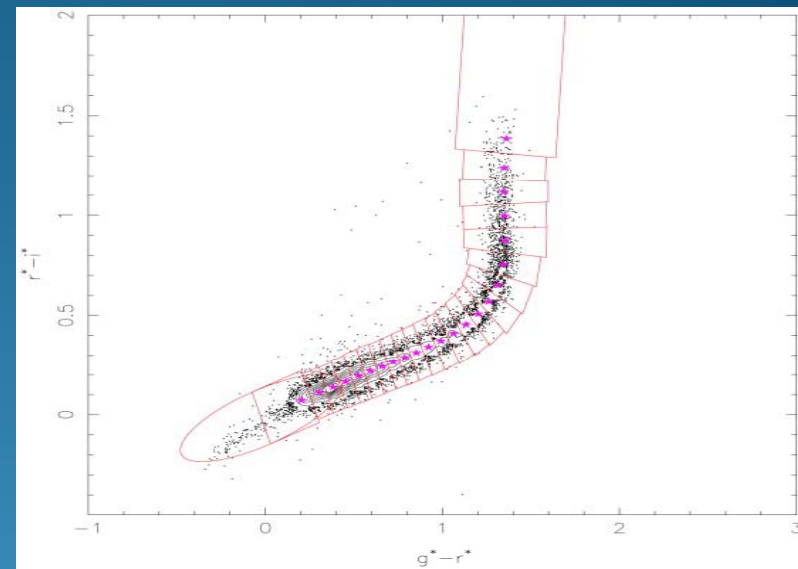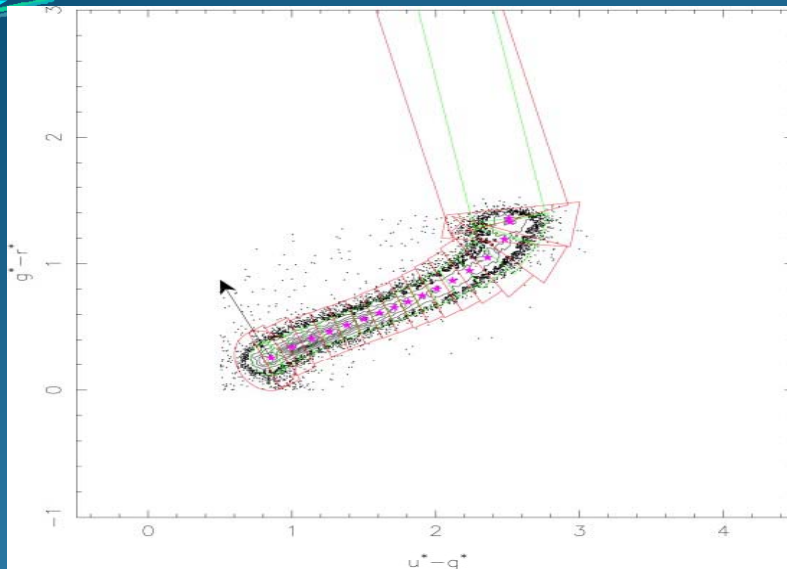


‣ **QSOs are supposed to be placed >4σ far from a cylindrical region containing the "stellar locus"** (S.L.), where σ depends on photometric errors.

**OR**

‣ **QSOs are supposed to be placed inside the inclusion regions**, even if not meeting the previous requirement.
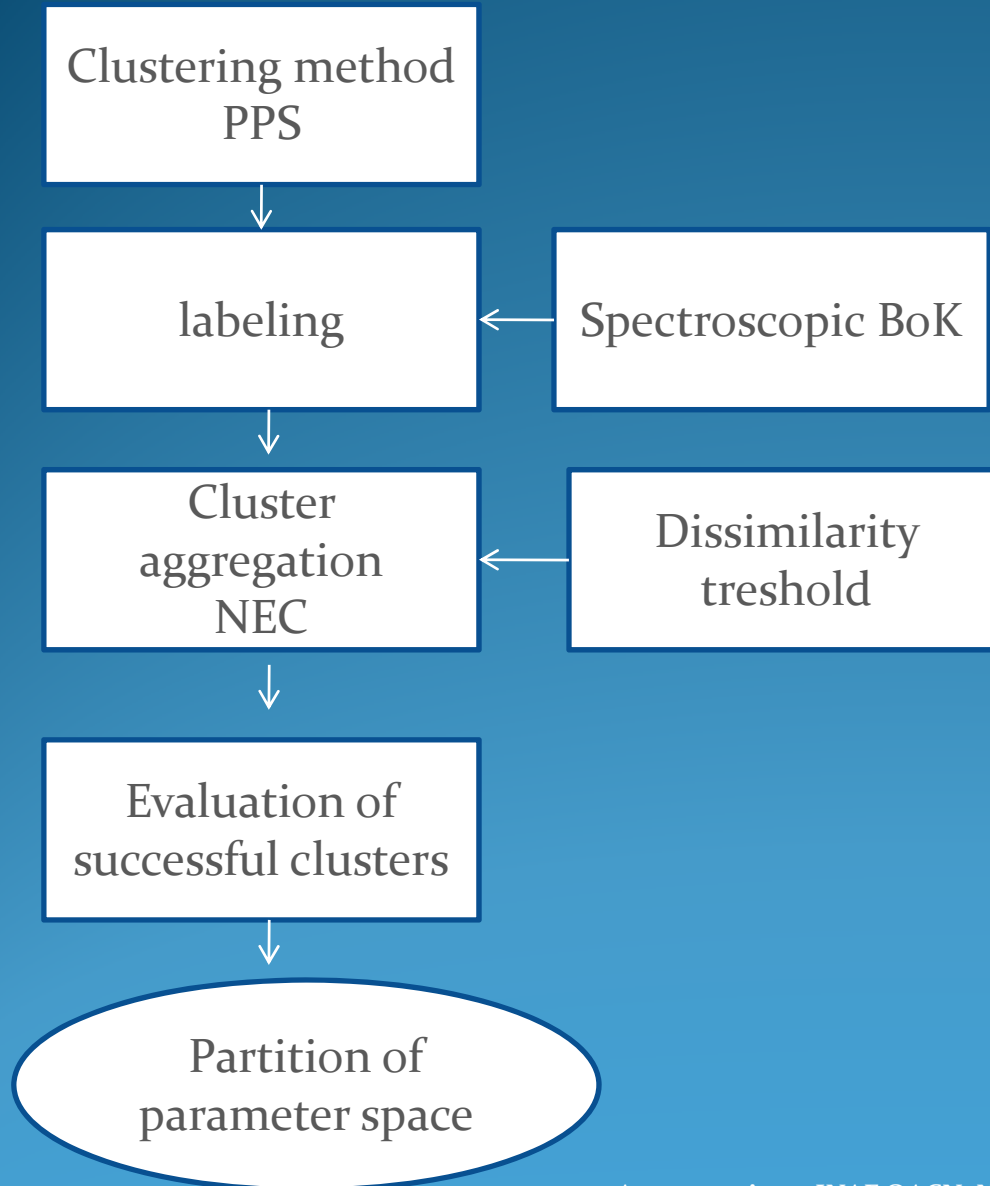
**c = 95%, e = 65%
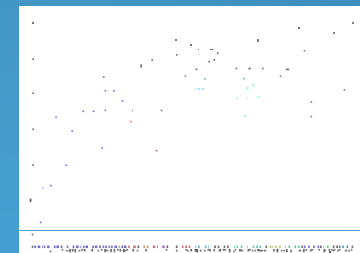locally less**

# SDSS QSOs targeting algorithm (II)



1. **inclusion regions** are regions where S.L. meets QSO's area (due to absorption from Lyα 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.**

# Unsupervised clustering based on latent variable methods

Clustering method
PPS

↓

labeling ← Spectroscopic BoK

↓

Cluster aggregation
NEC ← Dissimilarity treshold

↓

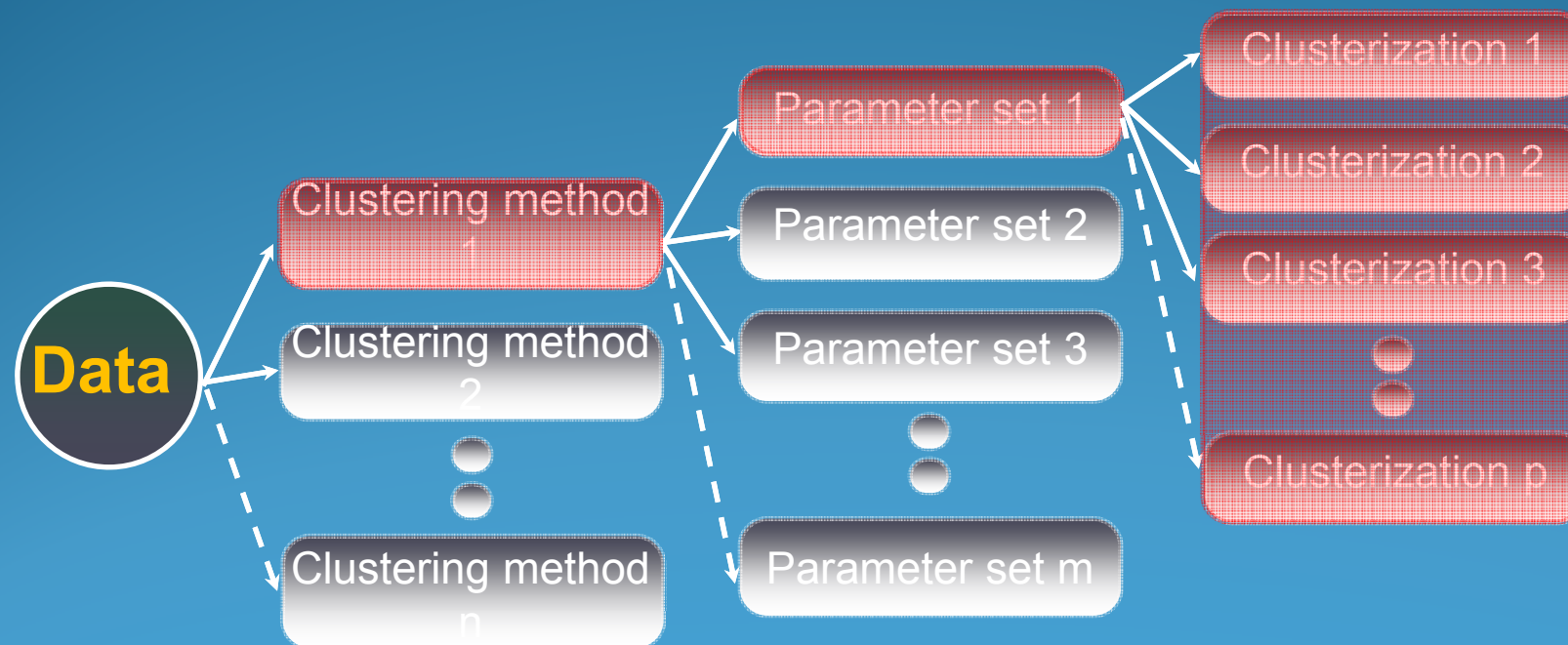Evaluation of
successful clusters

↓

Partition of
parameter space

1. **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.

2. **Dendrogram analysis**: the stability threshold(s) $D_{th}$ can be determined observing the number of branches at different levels of the graph.

# Many experiments are required

1. **Pre-clustering algorithm:** this phase can be accomplished performing a reduction of dimension of the feature space; this reduction via feature extraction/selection can be supervised or unsupervised (our choice in unsupervised).

2. **Agglomerative clustering**: both distance definition and a linkage model (simple, average, complete, Wards, etc.) need to be provided to perform clustering.

**Data**

Clustering method 1
Clustering method 2
Clustering method n

Parameter set 1
Parameter set 2
Parameter set 3
Parameter set m

Clusterization 1
Clusterization 2
Clusterization 3
Clusterization p

# Tuning succesfull clusters

**Once partition of colours space is completed** (as a function of $D_{th}$), **clusters mainly populated by QSO** (according the knowledge-base at our disposal) **are selected and information about these clusters are exploited for the candidate QSO selection.**

To determine the critical dissimilarity $D_{th}$ threshold we rely not only on a stability requirement. Given the following definition:

$$\left[\ \textbf{cluster is "successfull"}\ \right] \xleftrightarrow{\text{Def}} \left[\ \begin{array}{c}\textbf{its fraction of confirmed QSO} \\ \textbf{is higher then a fixed value}\end{array}\ \right]$$

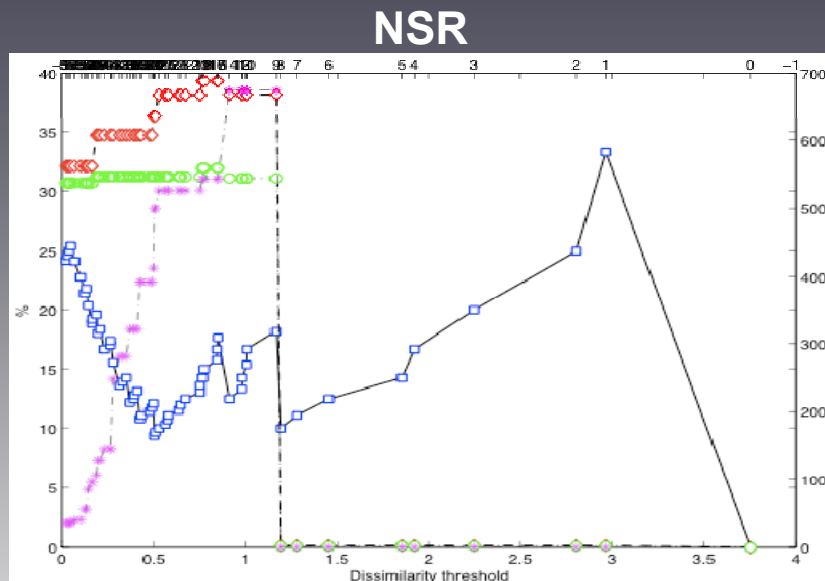we ask $D_{th}$ to maximize the **Normalized Success Ratio** (NSR):

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

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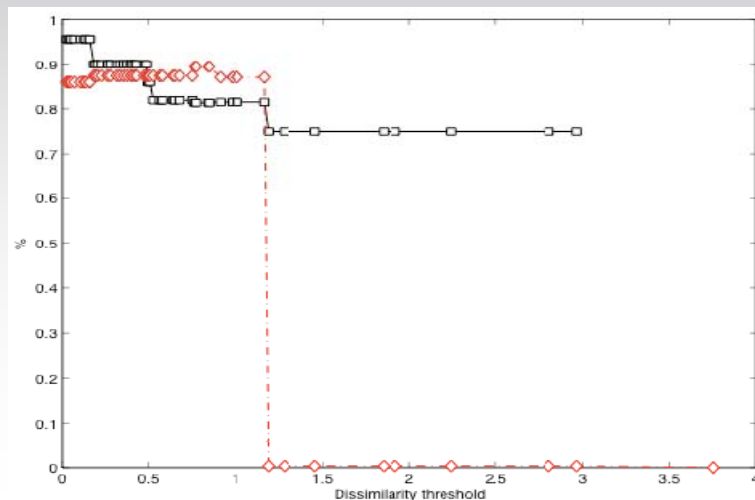
$$e_{tot} = \sum_{i=1}^{n} e_i$$

# 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.

| labels algorithm | QSOs | not QSOs |
|---|---|---|
| QSOs | 759 | 72 |
| not QSOs | 83 | 1327 |

# Data and experiments

## Data samples:

1. **Optical**: sample derived from SDSS database table "Target" queried for QSO candidates, containing ~ $1.11 \cdot 10^5$ records and ~ $5.8 \cdot 10^4$ confirmed QSO ('specClass == 3 OR specClass == 4').

2. **Optical + NIR**: sample derived from positional matching ('best') between SDSS-DR3 database view "Star" queried for all objects with spectroscopic follow-up available and detection in all 5 bands (u,g,r,i,z) with high reliability for redshift estimation and line-fitting classification ('specClass') and high S/N photometry, and UKIDSS-DR1 star-like ('mergedClass == -1') objects fully detected in each of the four Survey bands (Y,J,H,K) and clean photometry

## Experiments:

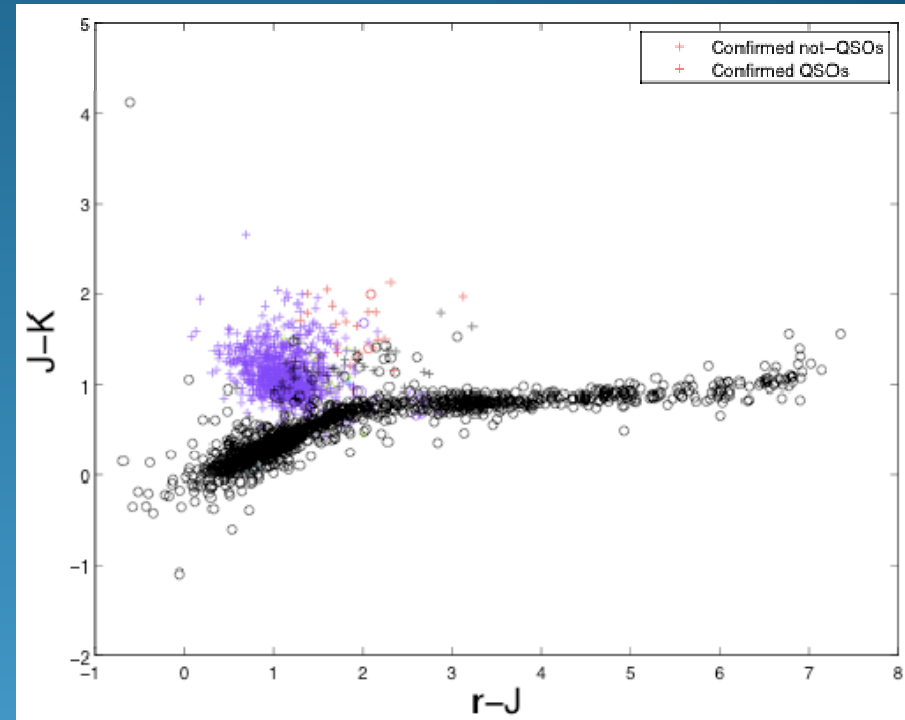| Optical (1) | Optical+NIR (2) | Optical (3) |
|---|---|---|
| candidate QSO | star-like objects | star-like objects |
| 4 colours | 4 + 3 colours | 4 colours |

# Experiment 2: SDSS ∩ UKIDSS
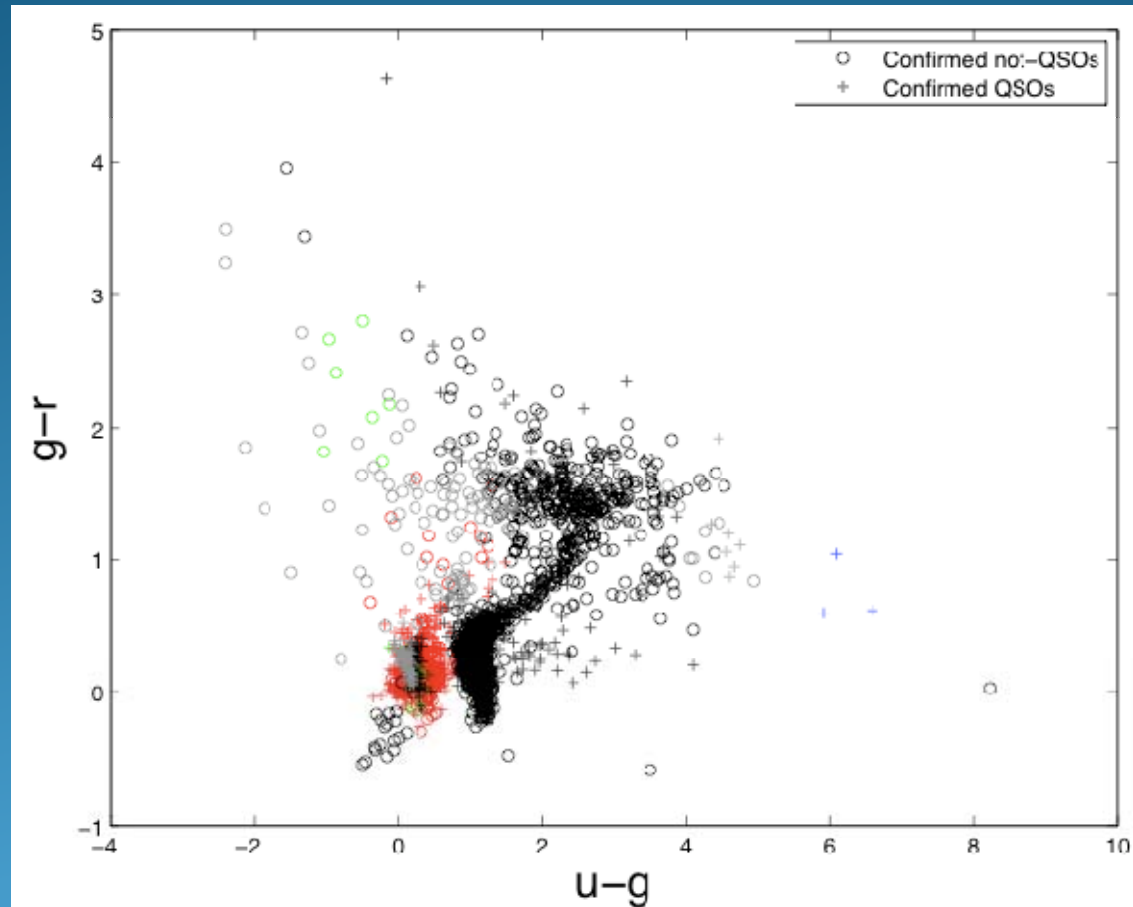
**u - g** vs **g - r**

**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).
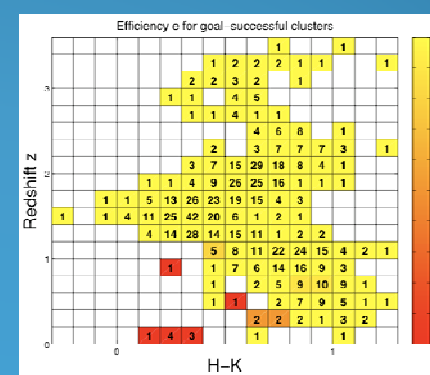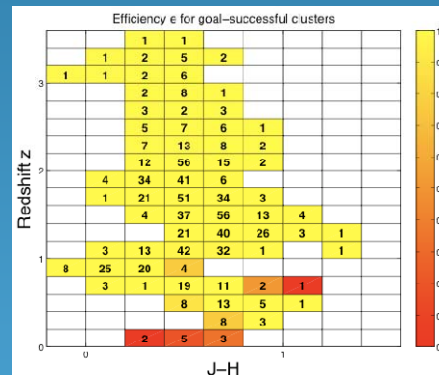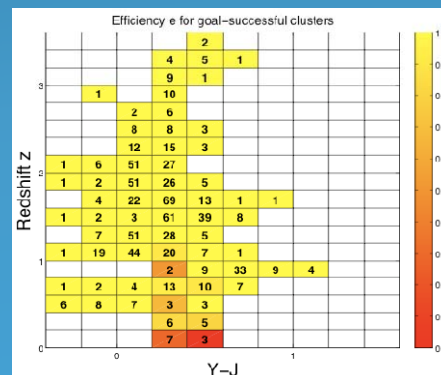
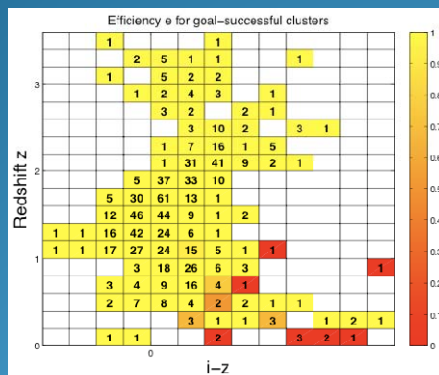# Experiment 3: optical colours

## 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 *e*

# Experiment 2: local values of *c*

# Results (I)

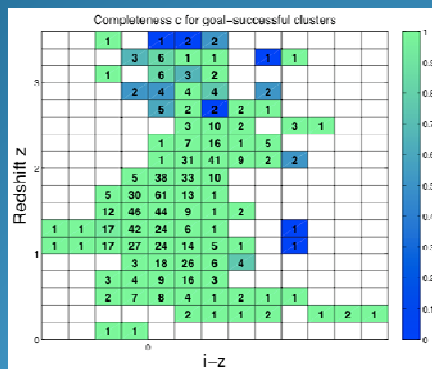| Sample | Parameters | Labels | $e_{tot}$ | $c_{tot}$ | $n_{gen}$ | $n_{suc\_clus}$ |
|---|---|---|---|---|---|---|
| **Optical** QSO candidates (1) | SDSS colours | 'specClass' | 83.4 % (± 0.3 %) | 89.6 % (± 0.6 %) | 2 | (3,0) |
| **Optical + NIR** star-like objects (2) | SDSS colours + UKIDSS colours | '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) |

Won one of the 2 prizes
Scientific application within Vobs

Talk at AAS Meeting - 2009

# Photometric redshifts estimation for QSOs using Neural Networks

G. Barentsen, R. D'Abrusco, O. Laurino, P. Nayak

NVO Summer School 2008, Santa Fe

# Pipeline for Photometric redshift estimation

# A unified vision

Galaxy and QSOs photometric redshifts differences depend only on the different sparseness of the data (BoK).

Few points in a high dimensionality space (i.e. spectroscopic **QSOs**).

High sparseness

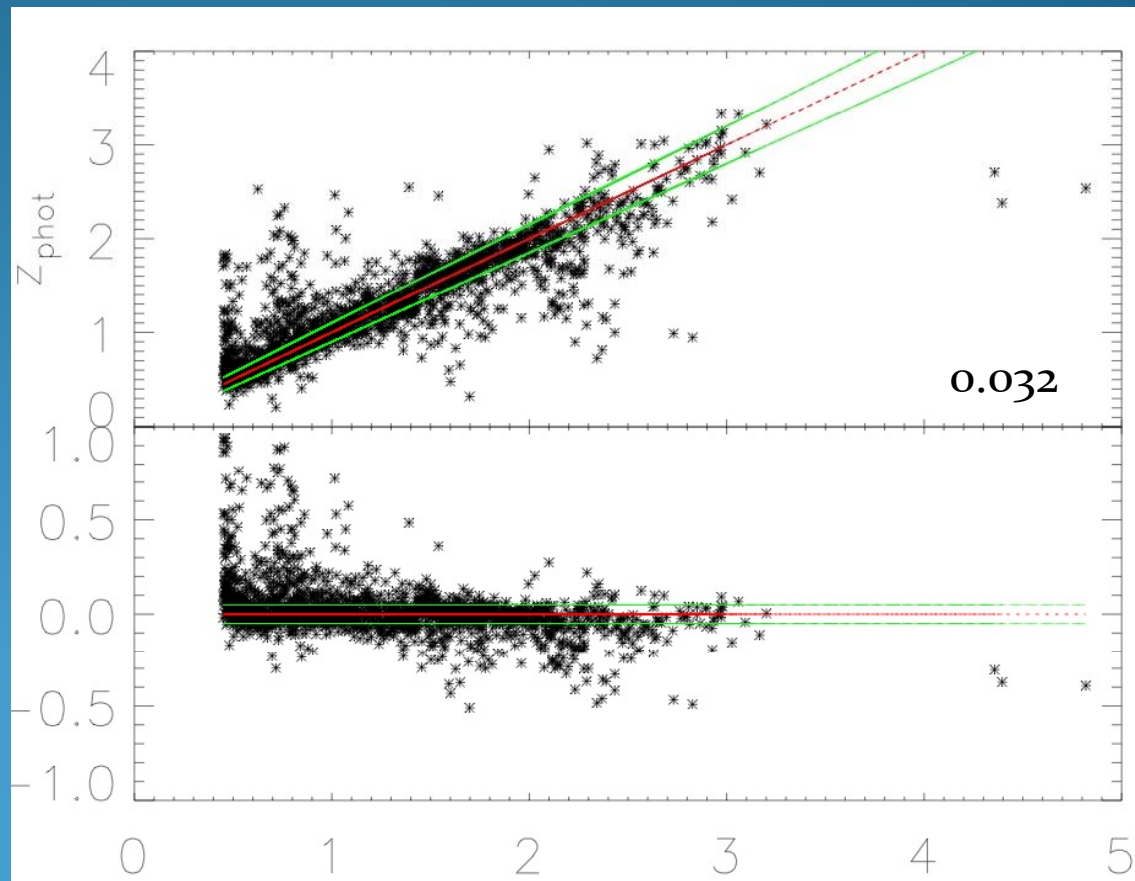Many points in a high dimensionality space (i.e. spectroscopic **galaxies**)

Low sparseness

# Clustering

**Low sparseness**

Crispy clustering, disjoint clusters, no redundancy.

**High sparseness**

Fuzzy clustering, overlapping clusters, redundancy.

Photo sources

N clusters

Disjoint

$$C_i \cap C_j = \emptyset$$
$$\forall i \neq j$$

$$p = \begin{cases} 1 \text{ for one clust.} \\ 0 \text{ for others} \end{cases}$$

Overlapping

$$\exists i \neq j:$$
$$C_i \cap C_j \neq \emptyset$$

$$p = (p_i) \text{ where}$$
$$p_i \neq 0 \text{ and } \sum_{i=1} p_i = 1$$

**Different sparseness of datasets can be taken into account when deriving photometric redshifts, by exploiting <u>redundance</u> between different clusters.**

"For usual (crispy) clustering, assigning a photometric source to one of the closest cluster is straightforward (given a distance definition).

For a fuzzy clustering the probabilistic nature of assignment needs to be taken into consideration. This is the reason why the methods for galaxies and QSOs $z_{phot}$ diverge."

# Recipes: an outlook

**(Low sparseness - galaxies)**

- Each photometric source is assigned to one single cluster.

- The $z_{phot}$ is calculated applying the NN trained on the members of that cluster.

- A unique value of $z_{phot}$ with a unique accuracy and likelihood is produced.

**(High sparseness - QSOs)**

- Each photometric source can have a non-zero probability to belong to every clusters.

- For each source, an estimate of $z_{phot}$ for each cluster is produced.

- A "committee" of NNs is used to determine the most reliable estimate of $z_{phot}$ and the accuracy of the estimate.

# Meta Institute for Computational Astrophysics



**Conference room**



**meeting room**

## Strategy

- To exploit **new communication and interaction tools** (social networks, second life, etc) for teaching and dissemination activities.
- To extend and deepen collaboration with Caltech (and organize a school on e-science (2010) in collaboration with Caltech)
- To extend and deepen collaboration with IUCAA (Inter University Center for Astronomy and Astrophysics, Poona-India)
- To propose a Master in Data Mining and Exploration as joint activity among faculties (Economics, Science and Sociology) and Universities (Federico II, Sannio and Second University)
- To open the use of DAME to new communities (Bioinformatics, economics

# Conclusions ?

- Methods are general and have been widely applied also outside of astronomy.

- In order to produce reliable results a large number of experiments is needed (as well as a good understanding of the tools). SUCCESSFUL SCIENCE CASES ARE A MUST

- Fast, optimized algorithms are required. They allow fast processing, with potentially better accuracy and a more detailed tracing of the process (the whole DR6 Galaxy photometric redshift catalogue went from 11 hs to 2.5 min)

- VO (or just the VO tools?) is not yet ready for data mining. But all that is needed is available. Visualization is still an issue

- BoK are the crucial issue for the future (need to bridge onthologies with intelligent BoK engines)

## Conclusions ?

# Funding

- Italy-USA "great relevance project" financed by MAE has been acknowledged by MAE as best project for 2008 and renewed for 2009.
- Funding pending from MIUR (PRIN) and from EU
- Funding foreseen in the framework of extension of PON-SCOPE

# Technical steps

- To add new data mining models (e.g. SOM, PCA and ICA, Bayesian networks, etc)
- To add web applications for specific applications (e.g. NEXT-II + 2D-Phot)
- To integrate advanced visualization capabilities (STILTS, + VO-PLOT)
- To do a feasibility sudy for the automatic extraction of knowledge from VO archives (spectroscopic knowledge in coll. with Padua University and Padua INAF)
- Time series analysis and classification tools

# Future science cases

- To integrate radioastronomy and optical data for WIMPS candidates to DM
- To improve on available Star/Galaxy classification using priors
- To improve AGN search and classification using supervised methods and improved spectroscopic base of knowledge.
- To study photometric transient classification and apply it to VST surveys.