Astroinformatics - Extracting New Knowledge about Universe

in the Epoch of Petabyte-Scaled Archives

Petr Škoda

Astronomical Institute of the Czech Academy of Sciences

Supported by grant COST LD-15113 of the Czech Ministry of Education Youth and Sports

Python BootCamp IAG, University Sao Paulo , Brazil 16th February 2017

Credits

- The presentation is based on many different sources – mainly the on-line published slides from IVOA meetings, slides from Astroinformatics workshops or pictures found on Internet.
- We acknowledge namely materials of E. Solano, E. Hatsiminaoglu, B.Hanish, G. Djorgovski, G. Longo, T. Hey, F. Le Petit, M. Breddels and presentations from AI2016 in Sorrento

Outline of the Talk

- Data Avalanche in astronomy
- Virtual Observatory
- Astroinformatics
- Visualizations
- Transfer of technology
- Citizen Science

• Examples of our projects

Data Avalanche

















Data Avalanche

Moore law for chips –doubling 1.5 year Data in astronomy – doubling < 1 yr ! (1000/10 yr)







600 000 CD = 372 TB (CD 650MB) 600 000 DVD = 2.5 PB (DVD=4.5GB) Bruce Monro Kilmington UK

Dark Energy Survey Camera

Dark Energy Camera (DECam)





Large Synoptic Survey Telescope



201 CCD 4kx4k, 3.2 Gpix every 20 sec 3.5 deg FOV (64cm) 20 TB/day=6 PB/yr RAW 1.5 PB catalogue !!! detection of changes 60s!

38 billion objects x 1000 32 tril. meas. -5 PB table Cerro Pachón – Future site of the LSST









- One 6.4-gigabyte image every ~17 seconds
- ~1000 visits (two back-to-back images), per night
- 15 terabytes of raw scientific image data / night
- 8.4 terapixel image (movie) of the sky to ~27.5 mag in 6 bands
- A catalog of ~38 billion observed objects (24B galaxies, 14B stars)
- A catalog of ~32 trillion photometric measurements
- ~2000 events per observation (includes variables+asteroids)
- ~2 million events per night, for 10 years
- Requirement: Process & transmit alerts within 60 seconds

Juric 2013

Project EUCLID

EUCLID

CONSORTIUM

The Euclid mission main goal



• What is the Nature of the Dark Matter and Energy?

Dubath 2016

EUCLID principles

DEFLECTION OF LIGHT RAYS CROSSING THE UNIVERSE, EMITTED BY DISTANT GALAXIES



SIMULATION: COURTESY NIC GROUP. S. COLOMBI. IAP

Dubath 2016

Euclid Data Archive



	2021	2022	2023	2024	2025	2026	2027
Storage (PB)	15	30	50	60	75	90	90
Computing (kilo cores / year)	2.5	5	8.5	12	16	20	21

Numbers from Christophe Dabin @ tk1

Atacama Large Milimeter Array ALMA

64 antennas 12m Chajnator 5000m Chile 2008-2013

- it is spectrograph as well as ...
- 0.5-2 PB/yr RAW



LOFAR network



SKA











Square Kilometer Array SKA



Cyber SKA

SQUARE KILOMETRE ARRA

Survey Raw Data Rates out of Correlator



R. Taylor 2013

SKA Data Challenge



SKA Archive Volumes

- ~0.5 10 PB/day of image data
- Source count ~10⁶ sources per square degree
- ~10¹⁰ sources in the accessible SKA sky, 10⁴ numbers/record
- ~1 PB for the catalogued data

100 Pbytes – 3 EBytes / year of fully processed data

Cherenkov Telescope Array

Cherenkov Astronomy and CTA



- Two arrays of 100 (South) et 20 (North) telescopes
- July 2015: sites selection, Chile (ESO) and La Palma
- 2016: pre-production phase
- 2018-2013: production phase
- Observatory open to the community





Millenium Run 10^10 particles Several Gpc to 10 kpc Cube 2 billion ly **One month MPSSC** 25 TB **Evolution of 20 mil** galaxies **Evolution merger tree**

Simulation of the Universe

World's fifth fastest supercomputer

- SPARC64[™] VIIIfx, 2.0GHz octcore (128Gflops / CPU)
 - Total 82944 nodes (663552 CPU core), 10.6 Pflops peak spped
- 16 GB memory / core, Total 1.3PB memory
- 6D torus network

K computer



Simulation of Universe



Problem of 1PB Data Transfer

Data transfer

- If 100 Mb/s network is available
 - ~10TB / day
 - ~100 days / 1PB
- Typically, effective speed is less than 10Mb/s
 - < 1TB / day
 - > 3 years / 1PB ······
- Delivery by car
 - 3 days / 1PB





 From Kobe to Chiba (from Kyoto to Tokyo + 100km , ~600km journey)



Data analysis at storage place Move processing = not data !

D'Abrusco 2010

A growing parameter space



D'Abrusco 2010

Virtual Observatory : Key Definitions

- "The Virtual Observatory will be a system that allows astronomers to interrogate multiple data centers in a seamless and transparent way, which provides new powerful analysis and visualization tools within that system, and which gives data centers a standard framework for publishing and delivering services using their data".
- Standardization of data and metadata, and of data exchange methods.
- Registry, listing available services and what can be done with them.

R.J.Hanisch, P.J.Quinn, in "IVOA – Guidelines for participation"





Ontologies in Astronomy



SKOS, RDF standards, search with understanding (not return QSO as binary star)

From Graham, M. AI2010

Ontologies



Technology of VO

Unified data format– VOTable, UCD (Vizier) Transparent transport (unit conversion) Web services (WS) e-commerce, B2B, J2EE, .Net VOregistry (DNS like) Google for data+WS protocols

- ConeSearch (searching in circle on sky)
- SIAP (Simple Image Access Protocol)
- SSAP(Simple Spectral Access Protocol)
- SLAP(Simple Line Access Protocol)
- TAP (Table Access Protocol)
- VOEVENT (transients, robotic telescopes,Sun)
- more datacubes, on-the-fly data generation

Technology of VO

ADQL (Astronomical Data Query Language)
XMATCH, REGION (2 catalogues - shifted)
Application interoperability – (PLASTIC), SAMP
Allows develop applications as bricks
sending VOTABLES (catalogue-spectra-images)

Commercial interest (GoogleSky, MS WWT)

Workflows - Astrogrid

Running remote services – e.g. Sextractor, CASJobs, AstroNeural MLP....



IVOA Universal Worker Service (UWS)


Ecosystem of VO - level 0



Ecosystem of VO - level 1



Ecosystem of VO - level 2



FITS standard

>30 years, separation of metadata (human readable and data)

```
SIMPLE =
                            T / file does conform to FITS standard
                            16 / number of bits per data pixel
BITPIX =
NAXIS
                             2 / number of data axes
        =
NAXIS1
                          2048 / length of data axis 1
       =
NAXIS2 =
                          2048 / length of data axis 2
                             T / FITS dataset may contain extensions
EXTEND =
          FITS (Flexible Image Transport System) format is defined in 'Astronomy
COMMENT
          and Astrophysics', volume 376, page 359; bibcode: 2001A&A...376..359H
COMMENT
BZERO
                         32768
       =
BSCALE =
                             1 / REAL=TAPE*BSCALE+BZER0
ORIGIN = 'PESO
                               / AsU AV CR Ondrejov
                               / Name of observatory (IRAF style)
OBSERVAT= 'ONDREJOV'
                      49.91056 / Telescope latitude (degrees), +49:54:38.0
LATITUDE=
                      14.78361 / Telescope longitud (degrees), +14:47:01.0
LONGITUD=
                           528 / Height above sea level [m].
HEIGHT =
TELESCOP= 'ZEISS-2m'
                               / 2m Ondrejov observatory telescope
GAIN
                             2 / Electrons per ADU
        =
READNOIS=
                            10 / Readout noise in electrons per pix
TELSYST = 'COUDE
                               / Telescope setup - COUDE or CASSegrain
INSTRUME= '0ES
                               / Coude echelle spectrograph
                               / Camera head name
CAMERA = 'VERSARRAY 2048B'
DETECTOR= 'EEV 2048x2048'
                               / Name of the detector
CHIPID = 'EEV 42-40-1-368'
                               / Name of CCD chip
```

VOTable

```
<TABLE name="SpectroLog">
<FIELD name="Target" ucd="meta.id" datatype="char" arraysize="30*"/>
<FIELD name="Instr" ucd="instr.setup" datatype="char" arraysize="5*"/>
<FIELD name="Dur" ucd="time.expo" datatype="int" width="5" unit="s"/>
<FIELD name="Spectrum" ucd="meta.ref.url" datatype="float" arraysize="*"
    unit="mW/m2/nm" type="location">
<DESCRIPTION>Spectrum absolutely calibrated</DESCRIPTION>
<LINK type="location"
    href="http://ivoa.spectr/server?obsno="/>
</FIELD>
<DATA><TABLEDATA>
<TR><TD>NGC6543</TD><TD>SWS06</TD><TD>2028</TD><TD>01301903</
TD></TR>
<TR><TD>NGC6543</TD><TD>SWS07</TD><TD>2544</TD><TD>01302004</
TD > </TR >
</TABLEDATA></DATA>
</TABLE>
```

Serialization (metadata first, end of data unknown, tree structure)

VOTable Serialization





Universal Content Descriptors

0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0	em.IR em.IR.J em.IR.H em.IR.K em.IR.3-4um em.IR.4-8um em.IR.8-15um em.IR.15-30um em.IR.30-60um em.IR.60-100um	Infrared pa Infrared be Infrared be	art of the spectrum etween 1.0 and 1.5 micron etween 1.5 and 2 micron etween 2 and 3 micron etween 3 and 4 micron etween 4 and 8 micron etween 8 and 15 micron etween 15 and 30 micron etween 30 and 60 micron etween 60 and 100 micron	
	S pos.eq Q pos.eq.dec Q pos.eq.ha Q pos.eq.ra Q pos.eq.spd S pos.errorEllipse Q pos.frame S pos.galactic Q pos.galactic.lat Q pos.galactic.lon	E0 Di Hi Ri Si Ri G Li Li	equatorial coordinates Declination in equatorial coordin lour-angle Right ascension in equatorial co South polar distance in equator Positional error ellipse Reference frame used for posit Galactic coordinates atitude in galactic coordinates ongitude in galactic coordinates	nates bordinates ial coordinates ions (FK5, ICRS,)
		stat.stdev stat.uncalib stat.value stat.variance stat.veight time time.age time.creation time.crossing time.duration	Standar Qualifier Miscella Variance Statistic Time, ge Age Creation Crossin Interval phenome End time	d deviation of a generic incalibrated quantity neous statistical value al weight eneric quantity in units of time or date n time/date (of dataset, file, catalogue,) g time of time describing the duration of a generic event or non e/date of a generic event

Characterization

Curation – long time preservation issues (digital libraries)

Provenance (how was processed, links to other products)

Characterization level 1 (spatial, spectral, temporal, polarization, location, coverage, porosity – SUB-CUBE)

Characterization level 2 (distorsion in images, spectra with nonlinear resolution)

Space-Time-Coordinate Data Model





Cherenkov Telescope Array Data Model



Simple Spectra Access Protocol Spectral Data Model

Simple Spectral Access Protocol V1.04



International Virtual

Observatorv

Alliance

Simple Spectral Access Protocol

Version 1.04 IVOA Recommendation Feb 01, 2008

This version: http://www.ivoa.net/Documents/REC/DAL/SSA-20080201.html Latest version: http://www.ivoa.net/Documents/latest/SSA.html Previous version(s): Version 1.03, December 2007 Version 1.02, September 2007 Version 1.01, June 2007 Version 1.00, May 2007 Version 0.97, November 2006 Version 0.96, September 2006 Version 0.95 May 2006 Version 0.91 October 2005 Version 0.90 May 2005 Editors: D.Tody, M. Dolensky Authors:

D.Tody, M. Dolensky, J. McDowell, F. Bonnarel, T.Budavari, I.Busko, A. Micol, P.Osuna, J.Salgado, P.Skoda, R.Thompson, F.Valdes, and the data access layer working group.



International Virtual Observatory

Observa

Alliance

IVOA Spectral Data Model Version 1.03 IVOA Recommendation 2007-10-29

This version (Recommendation Rev 1)

http://www.ivoa.net/Documents/REC/DM/SpectrumDM-20071029.pdf Latest version: http://www.ivoa.net/Documents/latest/SpectrumDM.html Previous versions:

http://www.ivoa.net/Documents/PR/DM/SpectrumDM-20070913.html

Editors:

Jonathan McDowell, Doug Tody Contributors:

Jonathan McDowell, Doug Tody, Tamas Budavari, Markus Dolensky, Inga Kamp, Kelly McCusker, Pavlos Protopapas, Arnold Rots, Randy Thompson, Frank Valdes, Petr Skoda, and the IVOA Data Access Layer and Data Model Working Groups.

SSAP Parameters

4.1.1 Mandatory Query Parameters

The following parameters must be implemented by a compliant service:

Parameter	Sample value	Physical unit	Datatype	
POS	52,-27.8	degrees; defaults to ICRS	string	
SIZE	0.05	degrees	double	
BAND	2.7E-7/0.13	meters	string	
TIME	1998-05-21/1999	ISO 8601 UTC	string	
FORMAT	votable	-	string	

4.1.2 Recommended and Optional Query Parameters

Parameter	Sample value	Unit	Req	Datatype	
APERTURE	0.00028 (=1")	degrees	OPT	double	
SPECRP	2000	$\lambda/d\lambda$	REC	double	
SPATRES	0.05	degrees	REC	double	
TIMERES	31536000 (=1yr)	seconds	OPT	double	
SNR	5.0	dimensionless	OPT	double	
REDSHIFT	1.3/3.0	dimensionless	OPT	string	
VARAMPL	0.77	dimensionless	OPT	string	
TARGETNAME	mars		OPT	string	
TARGETCLASS	star		OPT	string	
FLUXCALIB	relative		OPT	string	
WAVECALIB	absolute		OPT	string	
PUBDID	ADS/col#R5983		REC	string	
CREATORDID	ivo://auth/col\$R1234		REC	string	
COLLECTION	SDSS-DR5		REC	string	
TOP	20	dimensionless	REC	int	
MAXREC	5000		REC	string	
MTIME	2005-01-01/2006-01-01	ISO 8601	REC	string	
COMPRESS	true		REC	boolean	
RUNID			REC	string	

Big Data handling

- VO Space Moving big tables across (load only results)
- SSO Authentication, authorization, groups and consortia
- UWS Universal worker service (job synch, asynch)
- PDL Parameter Description Language
- SIM-DB Simulations, theory data

SPLAT-VO (Starlink, JAC)

🗙 Starlink SPLAT-VO: <plot0></plot0>								
File Analysis Edit Options Graphics Help								
	JPEG ⇔ ‡		∽ 🖓 🏫 :					
Displaying:	D:\SPEFO\L	A280060.RUI			▼ Y lin	nits 🕼 automatic 🔻		:V-hair
LAMBDA :	6528.643	🗌 :log		D:\SPEFO\LA280060	0.RUI: 0.98	82333	🗌 :log 🗌	:Track free
X scale:	1.0 💌	+ -		۲s	scale: 0.5		+	-
(uм			2-d compou	nd coordinate syste	m			
D:/SPEFO/LA280060.RUI (unkno 0.9 0.7 0.7 0.7 0.7 0.7 0.7 0.7 0.7 0.7 0.65	6300 6	350 6400	6450 65	500 6550 LAMBDA	6600	6650 6700	6750	· · · · · · · · · · · · · · · · · · ·
4								

VOspec (ESAC)



Colour-magnitude diagram



CIELO VO – line catalogue SLAP



(IVOA Line Data Model: Dubernet, Osuna et al., in preparation) (Simple Line Access Protocol: Salgado et al., in preparation)

ISM platform

Interstellar Medium Platform

Bring together expertise in modeling / simulation of the ISM

Provide theoretical services about ISM



Codes - Databases - Tools & services

Complex join of TVO bricks



Data-Knowledge-Wisdom Pyramid



Emergence of a Fourth Research Paradigm

- 1. Thousand years ago Experimental Science
 - Description of natural phenomena
- 2. Last few hundred years Theoretical Science
 - Newton's Laws, Maxwell's Equations...
- 3. Last few decades Computational Science
 - Simulation of complex phenomena
- 4. Today Data-Intensive Science
 - Scientists overwhelmed with data sets
 - from many different sources
 - Data captured by instruments
 - Data generated by simulations
 - Data generated by sensor networks
 - eScience is the set of tools and technologies
 - to support data federation and collaboration
 - For analysis and data mining
 - For data visualization and exploration
 - For scholarly communication and dissemination

(With thanks to Jim Gray)







From T. Hey, Al2010

X-informatics



FOURTH PARADIGM

DATA-INTENSIVE SCIENTIFIC DISCOVERY

EDITED BY TONY HEY, STEWART TANSLEY, AND KRISTIN TOLLS

Downloadable at Microsoft Research site

Changing methodology of the Science

Synergy between different worlds

Sociological aspects (net-based research communities)

Experimental astronomy has become a three players game



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

- Astroinformatics: AAS n. 215, Washington, December 2009, chairperson: K. Borne
- Astroinformatics 2010: Caltech (USA) June 16-19 2010; co-chairpersons: S.G. Djorgovski, G. Longo
- Astroinformatics 2011: UNINA Sorrento, co-chairpersons: S.G. Djorgovski, G. Longo

Longo 2010

Astroinformatics

- Analogy Bioinformatics (Genome analysis with GRIDS, ATB)
- e-Science in Astronomy
- Data mining, Knowledge discovery VO-NEURAL, DAME
- Examples
 - Photometric RedShift
 - Searching for QSO (light curves, MOS)
 - Automatic Light curves classification (GAIA, LSST)
- New ways of scholar communication (VR, 2nd Life, U-Science)
- BIG data problems, GPUs, NoSQL DB, visualization,
- Very NEW emerging discipline



2010 – Pasadena

2011 – Napoli

2012 – Redmond (Microsoft)

2013 – South Africa

Join us on Astroinformatics page on Facebook

IVOA - IG on KDD WIKI

Longo 2013

Data Driven Science

What is Fundamentally New Here?

- The *information volumes and rates* grow exponentially
- Most data will never be seen by humans



- A great increase in the data *information content*
- → Data driven vs. hypothesis driven science
- A great increase in the information complexity
- There are patterns in the data that cannot be comprehended by humans directly



Hidden Patterns in Data

Pattern or structure (Correlations, Clustering, Outliers, etc.) Discovery in High-Dimensional Parameter Spaces



D >> 3 parameter space hypercube

> High-D data cloud: mostly noise, of an arbitrary distribution

But in some corner of some sub-D projection of this data space, there is *something ≠ noise*

Visualization in Machine Learning

A Key Challenge: Visualisating Multidimensional Data Spaces

- Hyperdimensional structures (clusters, correlations, etc.) may be present in many complex data sets, whose dimensionality may be D ~ 10² – 10⁴, or higher
- It is a matter of *data understanding*, choosing the right data mining algorithms, and interpreting the results
- We are biologically limited to perceiving up to ~ 3 - 12(?) dimensions

What good are the data if we cannot effectively extract knowledge from them?

Scientific Communities

"The co-authorship network of scientists represents a prototype of complex evolving networks. In addition, it offers one of the most extensive database to date on social networks."^a

^aBarabàsi et al., "Evolution of the social network of scientific collaborations"

"Social scientists have long recognized the importance of boundaryspanning individuals in diffusing knowledge (Allen 1977; Tushman 1977), and recently, several papers have rigorously demonstrated that technological knowledge diffuses primarily through social relations, not through publications."^a

^aSorenson, and Singh, "Science, Social Networks and Spillovers"

From O. Laurino - AI2010

Motivations of a social networking IT platform for science

- The importance of boundaryspanning individuals in social networks might be what X-informatics is all about;
- we break scientific *cliques* and create new, unexpected, effective links across the science community's network;
- an effective scientific social network platform may be an effective step towards *seamless astronomy*. Seamless not only in terms of data and applications access, but also in terms of social interactions between people in the scientific network.



Virtual Worlds (2nd Life for Science)

http://slurl.com/secondlife/StellaNova/127/129/32

A part of the SciLands virtual continent: http://www.scilands.org/

Now migrating to the *OpenSim*-based VWs, e.g., Intel's *ScienceSim*



Virtual Conferences

Virtual conferences at zero cost

- Problem with time zones
- Outreach, education

Collaboration meetings Public outreach Professional seminars electromagner (Gauge fields mo **Compare With Observations** Tegmark et al. 2004 -0.4 **Big difference:** 3 2 negative charges and repulsion. 6-1.8 Clusters -1.4 + CMB Could DM be cha 0.2 0.4 Matter density 0. Note: w(z) is not constant in our model: $w(0) \sim -0.8$ and Based on 9 minute $w(1100) \sim -0.6$, but these plots were made assuming w = constPresented at American Society, Januar **Nobel laureate** John Mather

From Djorgovski - AI2010

Immersive VR Experiments



Astronomy and data parameter spaces



Scientists immersed in, and interacting with, numerical simulations of star clusters

Citizen Science - Galaxy ZOO



> 20 Science papers published so far

Examples ZOOniverse



Expert vs Non-expert Classifier



Lucy Fortson

AstroInformatics 2010

June 18, 2010

۲
Citizen Science x Expert Science

Verified by human – training sets Independent answers=estimate of error

Serendipitious discovery

J102210.25+311713.9	J123453.39+332430.3	06//3900006/41130 J113857.4+311846.6	J123126.52+405711.5	J110120.35+402242.3	
	٠	-	e	•	Galactic Peas
588017978351616137 J112615.25+385817.4	588017977278988558 J113948.93+382225.9	587739098597687419 J115135.32+375603.6	587739408393044155 J122245.71+360218.4	587739506616631548 J121139.18+330804.5	
	٠	•	•	۰	

Scale - complexity

Knowledge Discovery in U-Science



Known knowns :

Primary task. Data reduction by science team.

Known unknowns :

Related to primary task. Results funneled to specific researchers.

Unknown unknowns :

Serendipity. Currently rely on forum moderators to filter.

Hanny van Arkel - Voorwerp

Light echo of quasar?

Visualization of 1 B points - Gaia DR1



Visualization of Big Data



Visualization of Big Data



Visualization of Radio Data Cubes



3D Slicer provides full linked views, not just slices



Star Forming Regions in Galaxy

the Herschel infrared Galactic Plane Survey

Po-160-250µm composite

from cold starless clumps to hot HII Regions

Sergio Molinari, INAF-IAPS Credits: Gianluca Li Causi (INAF-IAPS)

IAU Astroinformatics 2016, Sorrento

Molinari et al. 2016

CAVE2 Monash University AU



8m diameter, 330 deg FOV , 80x LCD 46" 1366x768 Stereo + head tracking

From Astronomy to Earth Sciences



Big Data Era in Sky and Earth Observation – TD 1403 COST action

Finding Galaxies by Shape NASA



Description: Detecting objects from astronomical measurements by evaluating light measurements in pixels using intelligent software algorithms.

Image Credit: Catalina Sky Survey (CSS), of the Lunar and Planetary Laboratory, University of Arizona, and Catalina Realtime Transient Survey (CRTS), Center for Data-Driven Discovery, Caltech.

Finding Cancer Signatures NASA



Description: Detecting objects from oncology images using intelligent software algorithms transferred to and from space science. Image Credit: EDRN Lung Specimen Pathology image example, University of Colorado

Challenges of (Astro)informatics

- Big Data 3(5)xV
- Complex
- Missing values
- Censoring
- •Upper limits
- •Parallelization (Massive GPU new algorithms)
- •Queries in PB table
- Visualization of many dimensions
- Stream processing
- •Non- Gaussian Statistics, PDF

Ondřejov observatory



Ondřejov 2m Perek Telescope (1967)



Machine Learning of Spectra

Use case: ML of spectra profile of Halpha line (Be stars)



Be Stars : Emission in absorption





LAMOST (Guoshoujing)

Xinglong- China 4m mirror (30 deg meridian) 4000 fibers 10 mil spectra / 5 yr Automatic RV-z





LAMOST Spectral Surveys

DR1 (end 2013) 2 204 860 spectra 1 085 404 stars

DR3 (half 2015) **5 755 126** spectra DR4 (Feb 2016) **+ 741 522**

Each Fiber – 2 motors double arm 33mm circle

Fibre collects light from 3.3 arcsec circle on sky



Hobby Eberly Telescope (HET)



Mc Donald Observatory Texas

Equiv diameter 9.5m (11m)

Fixed in position during observation - only primary tracker



HETDEX Survey

In theory 34944 spectra every 20min !



VIRUS 78 IFU = 156 spectrographs
IFU= 448 fibers
34944 fibers , FOV 22 arcmin, 3500-5500 A, R=800
1 million spectra of galaxies (only part - statistic hits)

Resolution Degradation



OND R=13000

OND R=1800

LAMOST

LAMOST TSNE Structure



Semi-Supervised Training

Not supervised (even if not Domain Adaptation)

- sample of labelled data about 1600
- unlabelled (LAMOST) HDFS limit to 1,048,576 (2^20)

Graph methods:

Label spreading Label propagation



Spark on HDFS - National cloud MetaCentrum

Be Candidates - Semi Supervised ML



Palička 2016













CCD700 Outliers

Unsupervised learning – Local Outlier Factor - LOF



Shakurova 2016

LAMOST Outliers



LAMOST Be star - outlier



Concept of scientific "CLOUD"

ITERATIVE REPEATING of SAME computation (workflow)

Global non-linear optimization (Korel) Synthetic spectra (various elements, wavelength-ranges) Machine Learning (almost all methods)

LARGE stable INPUT data + small changing PARAMS Many runs on SAME data (tuning required)

Graphics visualization from postprocessed output (text) files Using WWW browser - supercomputing in PDA/mobil

Machine Learning of BIG Archive



Principles of SOM

Self-Organizing = Kohonen map



Association (activation) map

How many vectors activate every neuron

Unified Distance Matrix (U-matrix)

Every neuron= sum of distances to neighbours

The higher = more unique (outlier)

Machine Learning of Spectra SW view

ML does not produce new data – same spectra in groups Results the same size as input (+ small overhead)

Tracing visual shape from ML results Solf-Organizing maps – finding outliers Easy trace shape from neuron - clickable maps Visualisation of many spectra in web – dygraph (JS)
Virtual Observatory inside

- OND 2m archive on SSAP protocol (spectra access)
- LAMOST DR1 on SSAP (using DaCHS)
- Preprocessing (rectify, cutout) DataLink on server
- SAMP (send spectra to SPLAT-VO view details)
- Visualization of results
- VO-CLOUD cloud engine based on UWS REST jobs
- Cross-matching (ADQL, TAP, TOPCAT, TAPhandle, pyVO, Vizier)

Conclusions

- Machine learning on big spectra archives may identify new interesting objects yet unknown
- Crucial is interactive visualization of candidates
- VO technology helps in every step
- Future astronomy will be multidisciplinary
- Wide collaboration of experts and informaticians

Symposium S14

29 – 30 June 2017



http://eas.unige.ch/EWASS2017/session.jsp?id=S14

DEMO - create job

3	vo-cloud C	reate new SOM job - Iceweasel		
vo-cloud Create new	v S0 ★ 🕂			
	su. cas.cz /vocloud/jobs/index.xhtml		▼ C Soogle	🔍 🏠 自 🖊
 Most Visited▼ □G	etting Started Connecting			
VO-CLOUD	CREATE NEW SOM JOB			
Home Chulobs	局Create ▼ #Settings Admin ▼ ? Help ×Logout (skoda)			
Project label:	spectra4			
	SQM on spectra labeled in 4 classes			
Description:				
Email me result	s			
Edit config.json				
{ "Name	a":"Stellar spectra".	-		
"Algo	prithm":			
	"Bmu": "normal", "Threads": 1			
}, "Data	a":			
("Path": ["spectra.1863.4"],			
	"File_type": "csv",	_		
Upload paramete	rs			
Please attach	n data with config.json file.			
+ Choose				
Save and run	▼ Cancel			
				(c) mrq 2014 - feedback

DEMO – Job is running

_		vo-cloud Ja	bs - Iceweasel	_	_	_	_		
-cloud	lobs	× +							
🛞 vocl	oud-dev.	asu. cas.cz /vocloud/jobs/index.xhtml		▼ C ⁴	<mark>8</mark> ▼ Google			₫ 🗘	Ê
ost Visit	ed 🕶 🗌	Getting Started Connecting							
VO- C	CLOUI	D JOBS							
Home	🗅 Jobs	E Create ▼ Settings Admin ▼ ? Help × Logout (skoda)							
	_			_		_	_	_	
Туре	Id	Name	Created	Duration	Phase	Action	Delete	Details	
SOM	8603	spectra4	10/8/14	17 sec	EXECUTING	abort	×	C	
SOM	8555	spectra4 (copy)	10/8/14	14 sec	COMPLETED		×	0	
SOM	8550	spectra5	10/7/14	119 sec	COMPLETED		×	C	
SOM	8549	spectra4	10/7/14	62 sec	COMPLETED		×	0	
SOM	8548	iris	10/7/14	0 sec	COMPLETED		×	0	
SOM	8547	ecoli	10/7/14	3 sec	COMPLETED		×	0	
SOM	8537	spectra4_unspec	10/2/14	89 sec	COMPLETED		×	0	
SOM	8536	spectra4	10/2/14	108 sec	COMPLETED		×	C	
SOM	8534	new test of spectra (copy) (copy)	9/26/14	10 sec	COMPLETED		×	0	
SOM	8533	new test of spectra (copy)	9/26/14	0 sec	PENDING	start	×	0	
Korel	8530	testkorel (copy) (copy)	9/26/14	0 sec	COMPLETED		×	6	
SOM	8520	new test of spectra	9/26/14	12 sec	COMPLETED		×	0	
	7602	big job with map (copy)	4/14/13	37 sec	COMPLETED		×		
Korel								4 III III III III III III III III III I	

(c) mrq 2014 - feedback



8			vo-cloud Details of job 8	503 - Iceweasel		_		+ _	□ ×
vo-clo	ud Details of job 8 🗙 🚭								
(€) ® v	/ocloud-dev.asu. cas.cz /vocloud/jobs/index	.xhtml			▼ C B Google	公式	≜ ₽	⋒	≡
🛅 Most V	/isited▼ □ Getting Started □ Connecting	g							
SON	M 8603 COMPLETED	som local	10/8/14 2:31:27 PM	10/8/14 2:31:31 PM	10/8/14 2:32:45 PM	(3 Sec			
2 ri	un again x delete								
Prev	iew								
ind	lex.html - fullscreen								
	1.0				Neuron x: 15 v: 12	A			
	1.0								
	1.8				All Associated Spectra				
	1.7				X+ X- Y+ Y- HOME				
	1.6			1					
	1.5				Display reference vector				
	1.4				Display all spectra				
	1.3								=
	1.2				1. 2_betcmi_na230034 class: 2				
	1.1				2. 2_betcmi_qa070037 class: 2	- 11 1			
		~~~~~		Seener and					
	0.9		AAA MA SA						
	0.8			, b Ai					
	0.7								
	0.0								
	0.6								
	0.5								
	0.4								
	0.3								
	0.2								
	0.1								
	0					-			