Automated classification of Gaia sources

Kester Smith, Gaia AP team

- Introduction to Gaia
- Introduction to Machine Learning
- Application of ML to Gaia data analysis
- Detailed view of the Discrete Source Classifier

- Introduction to Gaia
- Introduction to Machine Learning
- Application of ML to Gaia data analysis
- Detailed view of the Discrete Source Classifier

- Introduction to Gaia
- Introduction to Machine Learning
- Application of ML to Gaia data analysis
- Detailed view of the Discrete Source Classifier

- Introduction to Gaia
- Introduction to Machine Learning
- Application of ML to Gaia data analysis
- Detailed view of the Discrete Source Classifier

Mission characteristics.

- The satellite rotates to scan the whole sky
- 70-90 visits per field
- Limiting magnitude G=20
- Only point sources will be observed.
- \checkmark ~ 1 billion sources expected



Photometric response



















Machine learning

Machine learning: Definition is a little problematic.

Can be thought of as making the computer perform a task for which the instructions have been derived from learning rules + examples, and are not provided directly by the programmer.

Machine learning

Machine learning: Definition is a little problematic.

Can be thought of as making the computer perform a task for which the instructions have been derived from learning rules + examples, and are not provided directly by the programmer.

- **Data space:** The data you actually observe
- **Feature space:** Derived from your data

Can divide ML into

- Supervised learning
- Unsupervised learning
- Semi-supervised learning

Can divide ML into

- Supervised learning DSC
- Unsupervised learning OCA
- Semi-supervised learning













Automated classification of Gaia sources - p. 10/54









Can also divide problems into

- Classification discrete classes, e.g. Star, Galaxy etc.
- Segression continuous variable, e.g. T_{eff} , redshift etc.

For classifiers, it's nice if output is *probabilistic*

Warning: overfitting and regularization

Some non-linear algorithms such as SVM and Neural networks, can fit anything perfectly.

Need to *regularize* the model and control by e.g. *cross* validation

Gaia photometry problem

Low resolution spectrum - 2 \times 180 'pixels'

(although intrinsic resolution \sim 15)

Suffer from the *curse of dimensionality*

Curse



Classifier specification

Need either;

- classifier that works well with a high number of data dimensions,

or

- reduce dimensions (PCA or similar).

also need to combine at least spectra and astrometry, and

possibly other information.

 \rightarrow probabilistic subclassifiers





Classifier specification

fundamental choice;

- supervised
- unsupervised

We use both approaches.

supervisedunsupervisedAll sourcesSupport vector machinek-means based methodOutliersself-organizing map

Map nodes organized in a grid in map space Each node has a random starting vector *W*

Map nodes organized in a grid in map space

Each node has a random starting vector \boldsymbol{W}

For first data point, choose the 'nearest' node

Map nodes organized in a grid in map space

Each node has a random starting vector W

For first data point, choose the 'nearest' node

Map nodes organized in a grid in map space

Each node has a random starting vector W



Map nodes organized in a grid in map space

Each node has a random starting vector W


Self organizing maps

Map nodes organized in a grid in map space

Each node has a random starting vector W





Mapping plot



Classification quality

Completeness_j =
$$\frac{n_{i=j,j}}{N_i}$$

Contamination_j =
$$\frac{\sum_{i \neq j} n_{i,j}}{\sum_{i} n_{i,j}}$$









Parameterization

Uses a variety of methods.

- Support Vector Machines
- Neural Network
- Grid interpolation
- Matisse algorithm



Regression - GSP-Phot



Regression - GSP-Phot



Neural network



Discrete Source Classifier (DSC)

- **•** Task of DSC is classification of all 10^9 Gaia sources.
- Supervised approach
- Available data;
 - BPRP 'photometry'
 - Solution Astrometry: π , proper motions
 - Variability
 - RVS spectrum
 - Possible morphology
 - Position, magnitude

Discrete Source Classifier (DSC)

- **•** Task of DSC is classification of all 10^9 Gaia sources.
- Supervised approach
- Available data;
 - BPRP photometry
 - Solution Astrometry: π , proper motions
 - Variability
 - RVS spectrum
 - Possible morphology
 - Position, magnitude

Modular design with subclassifiers

Modular design with subclassifiers



- Modular design with subclassifiers
- 'Hierarchical' photometric classifier based on Support Vector Machines (SVM)

- Modular design with subclassifiers
- 'Hierarchical' photometric classifier based on Support Vector Machines (SVM)
- Astrometric subclassifier based on Gaussian mixture model

- Modular design with subclassifiers
- 'Hierarchical' photometric classifier based on Support Vector Machines (SVM)
- Astrometric subclassifier based on Gaussian mixture model
- Position G Mag subclassifier based on Kernel Density Estimator

- Modular design with subclassifiers
- 'Hierarchical' photometric classifier based on Support Vector Machines (SVM)
- Astrometric subclassifier based on Gaussian mixture model
- Position G Mag subclassifier based on Kernel Density Estimator

$$P(C|D_1...D_N) = \frac{\prod_{n=1}^{n=N} P(C|D_n)}{P(C)^{N-1}}$$

[Bailer Jones & Smith, GAIA-C8-TN-MPIA-CBJ-053]

Position, G magnitude

Sky density of stellar sources in galactic coordinates LOCLog, of the density of objects, i.e. Log, of the number of objects per square degree). Samples: 5894254, Samples Out: 0



Position, G magnitude

Sky density of goos in galactic coordinates LTotal(Log. of the density of objects, i.e. Log. of the number of objects per square degree). Samples: 104484. Samples Out: 0



KDE

Estimate the underlying probability density fn. from which the data instances were presumably drawn.

$$f(x_0) = \frac{1}{Nh} \sum K_{\lambda}(\frac{x_i - x_0}{\lambda})$$

Estimate probability by

$$Pr(class = j) = \frac{f_j}{\sum_{i=1}^k f_i}$$
 for k classes.

Astrometric Classification

Based on mixture models

$$f(x_i) = \frac{A}{(2\pi)^{3/2} |\Sigma|^{1/2}} \exp(-\frac{1}{2} \times (x_i - \mu_i)^T \Sigma^{-1} (x_i - \mu_i))$$

Effect of convolution;

$$\Sigma \to \Sigma + \Sigma_{noise}$$

Astrometric classifier



PM RA

Astrometric classifier



PM RA

Model fitting

is done by Expectation maximization.

- 1. Determine the probabilities of the points for each existing model
- 2. Maximize the likelihood by adjusting model parameters
- 3. Repeat....









Photometric classifier: Svm



Photometric classifier: Svm



Photometric classifier: Svm



Svm concept: Convex hull



Svm concept: Convex hull










Applying the classifier

For a new point $\mathbf{s},$ the decision boundary is;

$$f = \sum_{i} \alpha_{i} y_{i} \mathbf{x}_{i} \cdot \mathbf{s} + b,$$

$$f > 0. \quad |\mathbf{s} \to class1$$

$$f < 0. \quad |\mathbf{s} \to class2.$$







$$u, v \to u^2, v^2$$



$$u, v \to u^2, v^2$$



Kernel trick

$$f = \sum_{i} \alpha_{i} y_{i} \mathbf{x}_{i} \cdot \mathbf{s} + b,$$

replace dot products with Kernel function. Most commonly radial basis function

$$\exp\left(-\gamma \|x_i - s\|^2\right)$$

need to tune for γ .

Svm: Non separable cases



Svm: Non separable cases



Svm: Non separable cases



Internal structure of DSC



SVM - One Class



SVM - One Class



Probabilities



 ~ 1

Probabilities: pairwise coupling

- classify pairwise (star vs. binary, star vs. quasar etc.)
- Can combine all pairwise probabilities to unique set of $P(C_m)$

(Wu et al. 2004, Journal of Machine learning research)

Training data

Input spectra from codes or from empirical libraries+codes (semi-empirical data) Semi-empirical = SDSS (usually)







Grid	star	wd	bin	qso	gal	???
APec	94.44	_	_	5.55	0.00	0.00
Fastrot	98.26	_	_	0.69	0.00	1.04
Phoenix N	95.73	_	_	0.02	0.54	3.69
Phoenix R	98.67	_	_	0.15	0.85	0.33
SDSS Stars	99.58	_	—	0.20	0.21	0.00
UCD Cond	70.52	_	_	0.45	0.76	28.27
UCD Dust	98.30	_	—	0.00	0.00	1.70
WR	76.74	_	_	9.30	0.00	13.95
SDSS QSOs	0.20	_	—	95.85	2.33	1.60
SDSS galaxies	0.15	—	_	0.50	99.01	0.33
MARCS	82.05	2.80	5.75	4.05	1.90	3.45
0	78.80	11.00	0.40	5.80	1.20	2.80
В	86.20	5.80	1.00	4.40	0.80	1.80
А	89.10	3.20	0.90	3.60	1.00	2.20
Be	82.75	1.72	0.00	10.92	4.59	0.00
C stars	89.48	0.23	1.40	7.24	1.63	0.00
physBins	29.27	1.60	57.40	4.67	2.32	4.72
WDA	18.20	58.42	0.52	16.10	1.95	4.80
WDB	15.67	58.20	0.55	19.37	1.95	4.25

Grid	star	wd	bin	qso	gal	???
APec	94.44	_	_	5.55	0.00	0.00
Fastrot	98.26	_	_	0.69	0.00	1.04
Phoenix N	95.73	_	_	0.02	0.54	3.69
Phoenix R	98.67	_	_	0.15	0.85	0.33
SDSS Stars	99.58	_	_	0.20	0.21	0.00
UCD Cond	70.52	_	_	0.45	0.76	28.27
UCD Dust	98.30	_	—	0.00	0.00	1.70
WR	76.74	_	_	9.30	0.00	13.95
SDSS QSOs	0.20	_	_	95.85	2.33	1.60
SDSS galaxies	0.15	_	_	0.50	99.01	0.33
MARCS	82.05	2.80	5.75	4.05	1.90	3.45
0	78.80	11.00	0.40	5.80	1.20	2.80
В	86.20	5.80	1.00	4.40	0.80	1.80
А	89.10	3.20	0.90	3.60	1.00	2.20
Be	82.75	1.72	0.00	10.92	4.59	0.00
C stars	89.48	0.23	1.40	7.24	1.63	0.00
Binaries	29.27	1.60	57.40	4.67	2.32	4.72
WDA	18.20	58.42	0.52	16.10	1.95	4.80
WDB	15.67	58.20	0.55	19.37	1.95	4.25

Grid	star	wd	bin	qso	gal	???
APec	94.44	_	_	5.55	0.00	0.00
Fastrot	98.26	_	—	0.69	0.00	1.04
Phoenix N	95.73	_	—	0.02	0.54	3.69
Phoenix R	98.67	_	—	0.15	0.85	0.33
SDSS Stars	99.58	—	—	0.20	0.21	0.00
UCD Cond	70.52	_	—	0.45	0.76	28.27
UCD Dust	98.30	—	—	0.00	0.00	1.70
WR	76.74	_	_	9.30	0.00	13.95
SDSS QSOs	0.20	—	—	95.85	2.33	1.60
SDSS galaxies	0.15	_	_	0.50	99.01	0.33
MARCS	82.05	2.80	5.75	4.05	1.90	3.45
0	78.80	11.00	0.40	5.80	1.20	2.80
В	86.20	5.80	1.00	4.40	0.80	1.80
А	89.10	3.20	0.90	3.60	1.00	2.20
Be	82.75	1.72	0.00	10.92	4.59	0.00
C stars	89.48	0.23	1.40	7.24	1.63	0.00
physBins	29.27	1.60	57.40	4.67	2.32	4.72
WDA	18.20	58.42	0.52	16.10	1.95	4.80
WDB	15.67	58.20	0.55	19.37	1.95	4.25

Grid	star	wd	bin	qso	gal	???
APec	94.44	_	_	5.55	0.00	0.00
Fastrot	98.26	_	_	0.69	0.00	1.04
Phoenix N	95.73	_	_	0.02	0.54	3.69
Phoenix R	98.67	_	_	0.15	0.85	0.33
SDSS Stars	99.58	_	_	0.20	0.21	0.00
UCD Cond	70.52	_	_	0.45	0.76	28.27
UCD Dust	98.30	_	_	0.00	0.00	1.70
WR	76.74	_	_	9.30	0.00	13.95
SDSS QSOs	0.20	_	_	95.85	2.33	1.60
SDSS galaxies	0.15	_	_	0.50	99.01	0.33
MARCS	82.05	2.80	5.75	4.05	1.90	3.45
0	78.80	11.00	0.40	5.80	1.20	2.80
В	86.20	5.80	1.00	4.40	0.80	1.80
А	89.10	3.20	0.90	3.60	1.00	2.20
Be	82.75	1.72	0.00	10.92	4.59	0.00
C stars	89.48	0.23	1.40	7.24	1.63	0.00
physBins	29.27	1.60	57.40	4.67	2.32	4.72
WDA	18.20	58.42	0.52	16.10	1.95	4.80
WDB	15.67	58.20	0.55	19.37	1.95	4.25

Grid	star	wd	bin	qso	gal	???
APec	94.44	_	_	5.55	0.00	0.00
Fastrot	98.26	_	—	0.69	0.00	1.04
Phoenix N	95.73	_	—	0.02	0.54	3.69
Phoenix R	98.67	_	—	0.15	0.85	0.33
SDSS Stars	99.58	_	_	0.20	0.21	0.00
UCD Cond	70.52	_	_	0.45	0.76	28.27
UCD Dust	98.30	_	_	0.00	0.00	1.70
WR	76.74	_	_	9.30	0.00	13.95
SDSS QSOs	0.20	_	_	95.85	2.33	1.60
SDSS galaxies	0.15	_	_	0.50	99.01	0.33
MARCS	82.05	2.80	5.75	4.05	1.90	3.45
0	78.80	11.00	0.40	5.80	1.20	2.80
В	86.20	5.80	1.00	4.40	0.80	1.80
А	89.10	3.20	0.90	3.60	1.00	2.20
Be	82.75	1.72	0.00	10.92	4.59	0.00
C stars	89.48	0.23	1.40	7.24	1.63	0.00
physBins	29.27	1.60	57.40	4.67	2.32	4.72
WDA	18.20	58.42	0.52	16.10	1.95	4.80
WDB	15.67	58.20	0.55	19.37	1.95	4.25

Dsc results
















Colours



Colours



Colours



Phoenix results



20







Qso z distribution

Histogram of qsoParamsFaint[, 6]



Ultra-cool dwarfs



Summary

- Gaia data analysis relies heavily on Machine learning algorithms
- These can provide powerful tools for analyzing large data sets (data mining)
- Increasingly used in astronomy, with larger surveys
- But still probably under-utilized

Exercise

- simulated Gaia data set with several classes
- Use R package, which implements many ML algorithms
- Exercise 1: Choose a ML algorithm and classify the sources
- Exercise 2: Choose a ML algorithm and parameterize a class
- Exercise 3: Choose an unsupervised algorithm to cluster the data

Some simple R scripts are provided as a starting point, including a knn classifier and parameterizer.

Algorithms

Algorithm	Classification	Regression	Supervised	Unsupervised
kNN	\checkmark	\checkmark	\checkmark	X
SVM	\checkmark	\checkmark	\checkmark	X
SoM	\checkmark	\checkmark	\checkmark	\checkmark
kMeans	\checkmark	\checkmark	\checkmark	\checkmark
Neural Network	\checkmark	\checkmark	\checkmark	X
PCA+KDE	\checkmark	X	\checkmark	X
Tree methods	\checkmark	\checkmark	\checkmark	\checkmark
Boosting	\checkmark	\checkmark	\checkmark	\checkmark

. . .