



Unsupervised learning in a Nutshell

A probabilistic approach for galaxy emission-line classification

Rafael S. de Souza

EWASS-Prague, 2017 Astroinformatics







The first step in wisdom is to know the things

themselves, Carl Linnaeus

Classification and name-giving will be the foundation of our science.

Hierarchical classification of organisms based on similarities and differences



Classification of objects is a major driver in natural sciences

Morphological feature space



HR diagram: 2D-Stellar classification

Subspace of features

Ad hoc dimensionality reduction



Taxonomical classification





Literature --Name given

Data-driven groups

Pre-Processing

Feature extraction/dimensionality



Feature Selection

LASSO

Ridge







Scientific Case

Test a GMM to discriminate galaxies based on their ionization sources



Rest Wavelength (Å)



Make it simple, but solid: Exploit model-based GMM

- Soft classification
- Generative model
- Stability
- Other Mixture models are possible: Gamma mixture, Poisson mixture, ...





A probabilistic approach to emission-line galaxy classification

R. S. de Souza, M. L. L. Dantas, M. V. Costa-Duarte, E. D. Feigelson, M. Killedar, P.-Y. Lablanche, R. Vilalta, A. Krone-Martins, R. Beck, F. Gieseke



Full solution





LogNII_Ha

How to decide number of clusters?

Few internal validation methods



Explore sampling properties of GMMs to simulate synthetic data



Model-based clustering allows residual analysis



How to measure similarity among groups?

Linear Discriminant Analysis

Project groups into a lower dimensional subspace, while preserving their distances

Kullback-Leibler divergence-relativeentropy1.0
0.8

$$D(f_i||f_j) = \int_{\mathbf{x}} f_i(\mathbf{x}) \ln \frac{f_i(\mathbf{x})}{f_j(\mathbf{x})} d\mathbf{x}.$$

Measures how one probability distribution diverges from a second expected probability distribution

Genomics

COMPARING CHROMOSOMES 1 Outer band represents each species' first Bar charts tell how many base chromosome. Numbers represent millions pairs, 0 to 1 million, match part Chimp of base pairs on the chromosome of the human chromosome. 3555558 1100 Line charts show what percent-ASE age of the human chromosome is similar to each of the other five DNA genomes. Anens Nonkey BRC Chicken 120 Species Regions of highest 30 similarity tend to 240 bundle together. 150 160 Gaps represent areas 170 30 that either haven't or 20. 180 can't be sequenced. 190 180 10 170 160 150 . 40 130 Mous Lines join the 200 regions on each chromosome that are most similar to a human's (based on number of matching base pairs). If a region of the human **Thicker lines** genome, there is reaso represent more basic functions that are similarity.

Visualize multiple associations

Chord diagrams--borrowed from Bioinformatics

Statistical vs Astrophysical classes

An extra group spots the dichotomy Seyfert/LINER

Remarks

- Soft classification provides more realistic predictions, account for boundaries uncertainties;
- Generative models combined with spatial analysis provides a solid recipe to define groups;
- External clustering validation adds semantics to groups via previous domain knowledge.

GMM_Catalogue

Catalogue with probabilistic classification of galaxies based on their ionization source

View the Project on GitHub

This project is maintained by COINtoolbox

Hosted on GitHub Pages — Theme by orderedlist

Galaxy Classification via Gaussian Mixture Models

Catalogue with probabilistic classification of galaxies based on their ionization source using Gaussian Mixture Models

This is one of the products of the third edition of the COIN Residence Program, which took place in August/2016 in Budapest (Hungary).

The catalogue is given bellow. Check the individual folders for detailed information on the files presented here.

GMM Catalogue

This catalogue was designed to provide soft/probabilistic classifications for galaxies based on their ionization source. The data consists of a combination of the BPT and WHAN diagrams.

The probabilistic classification of galaxies via GMM is constructed using data from the SDSS DR7 and the SEAGal/STARLIGHT catalogues and is available at Catalogue.

Short tutorial for the use of GMM in Python

Created by Pierre-Yves and Madura Killedar and available at GMM Python tutorial

BAYESIAN MODELS for Astrophysical Data

Using R, JAGS, Python and Stan

Joseph M. Hilbe, Rafael S. de Souza and Emille E. O. Ishida

