

Overview

Big Data

Definition of BIG DATA

Astronomical Surveys

Supervised ML

Overview

Support Vector Machine

Neural Network

Unsupervised ML

Overview

Agglomerative clustering

LOF

PCA

ML limitations

References

Machine Learning in Astronomy

Reza Monadi

UC Riverside

May 14, 2020

► Astronomy in **Big Data** era

Overview

Big Data

Definition of BIG DATA

Astronomical Surveys

Supervised ML

Overview

Support Vector Machine

Neural Network

Unsupervised ML

Overview

Agglomerative clustering

LOF

PCA

ML limitations

References

Overview

Big Data

Definition of BIG DATA

Astronomical Surveys

Supervised ML

Overview

Support Vector Machine

Neural Network

Unsupervised ML

Overview

Agglomerative clustering

LOF

PCA

ML limitations

References

Overview

Big Data

Definition of BIG DATA

Astronomical Surveys

Supervised ML

Overview

Support Vector Machine

Neural Network

Unsupervised ML

Overview

Agglomerative clustering

LOF

PCA

ML limitations

References

Overview

Big Data

Definition of BIG DATA

Astronomical Surveys

Supervised ML

Overview

Support Vector Machine

Neural Network

Unsupervised ML

Overview

Agglomerative clustering

LOF

PCA

ML limitations

References

3Vs in astronomy

Machine Learning
in Astronomy

Reza Monadi

Overview

Big Data

Definition of BIG DATA

Astronomical Surveys

Supervised ML

Overview

Support Vector Machine

Neural Network

Unsupervised ML

Overview

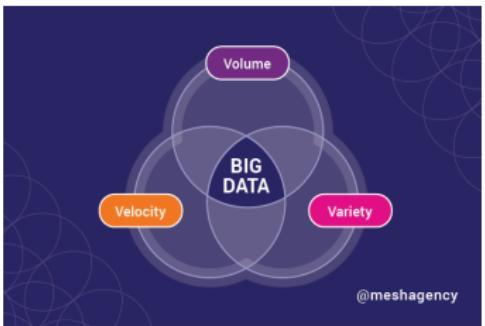
Agglomerative clustering

LOF

PCA

ML limitations

References



- ▶ Volume: larger quantities of data by the advent of better telescopes and surveys

3Vs in astronomy

Machine Learning
in Astronomy

Reza Monadi

Overview

Big Data

Definition of BIG DATA

Astronomical Surveys

Supervised ML

Overview

Support Vector Machine

Neural Network

Unsupervised ML

Overview

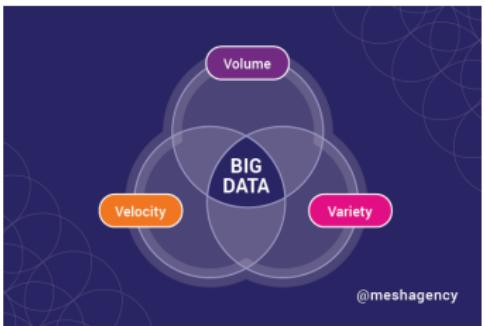
Agglomerative clustering

LOF

PCA

ML limitations

References



- ▶ Volume: larger quantities of data by the advent of better telescopes and surveys
- ▶ Velocity: Higher speed of incoming observational data

3Vs in astronomy

Machine Learning
in Astronomy

Reza Monadi

Overview

Big Data

Definition of BIG DATA

Astronomical Surveys

Supervised ML

Overview

Support Vector Machine

Neural Network

Unsupervised ML

Overview

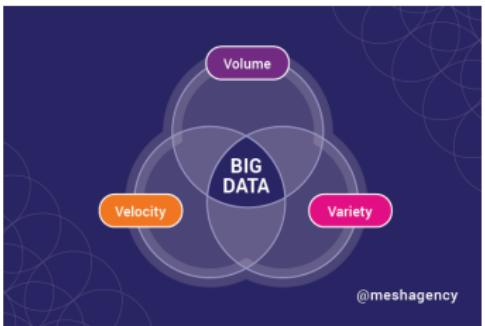
Agglomerative clustering

LOF

PCA

ML limitations

References



- ▶ Volume: larger quantities of data by the advent of better telescopes and surveys
- ▶ Velocity: Higher speed of incoming observational data
- ▶ Variety: multi-wavelength spectroscopic and photometric data from versatile astronomical objects

Astronomical surveys are exponentially growing

Sky Server Projects	Data Volume
DPOSS (The Palomar Digital Sky Survey)	3 TB
2MASS (The Two Micron All-Sky Survey)	10 TB
GBT (Green Bank Telescope)	20 PB
GALEX (The Galaxy Evolution Explorer)	30 TB
SDSS (The Sloan Digital Sky Survey)	40 TB
SkyMapper Southern Sky Survey	500 TB
PanSTARRS (The Panoramic Survey Telescope and Rapid Response System)	40 PB
LSST (The Large Synoptic Survey Telescope)	200 PB
SKA (The Square Kilometer Array)	4.6 EB expected

Zhang et. al 2015



Rubin's Observatory



Large surveys = Big Data
Big Data + ML = Astronomical Knowledge

Overview

Big Data

Definition of BIG DATA

Astronomical Surveys

Supervised ML

Overview

Support Vector Machine

Neural Network

Unsupervised ML

Overview

Agglomerative clustering

LOF

PCA

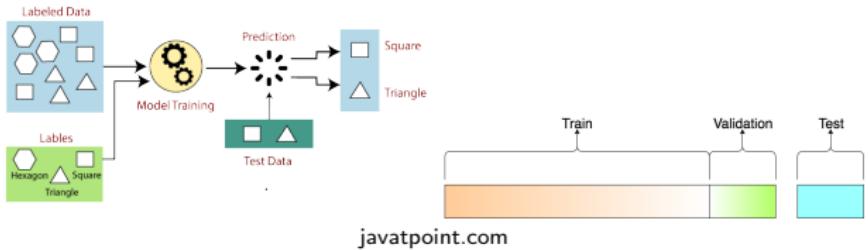
ML limitations

References

Stages of Supervised Learning

Machine Learning
in Astronomy

Reza Monadi



► Training:

Overview

Big Data

Definition of BIG DATA

Astronomical Surveys

Supervised ML

Overview

Support Vector Machine

Neural Network

Unsupervised ML

Overview

Agglomerative clustering

LOF

PCA

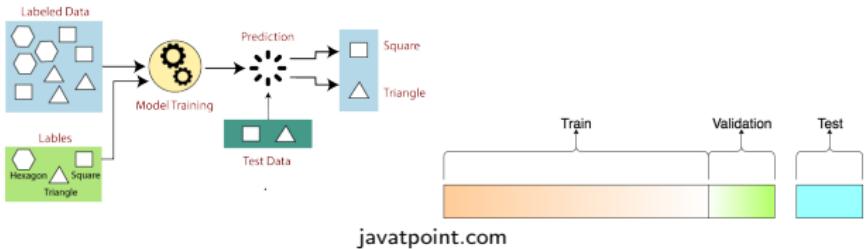
ML limitations

References

Stages of Supervised Learning

Machine Learning
in Astronomy

Reza Monadi



- ▶ Training:
 1. Select a model

Machine Learning
in Astronomy

Overview

Big Data

Definition of BIG DATA

Astronomical Surveys

Supervised ML

Overview

Support Vector Machine

Neural Network

Unsupervised ML

Overview

Agglomerative clustering

LOF

PCA

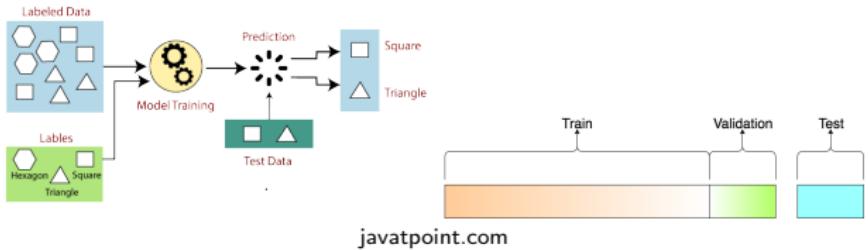
ML limitations

References

Stages of Supervised Learning

Machine Learning
in Astronomy

Reza Monadi



- ▶ **Training:**
 1. Select a model
 2. Train the model by the training set

Supervised ML

Overview
Support Vector Machine
Neural Network

Unsupervised ML

Overview
Agglomerative clustering
LOF
PCA

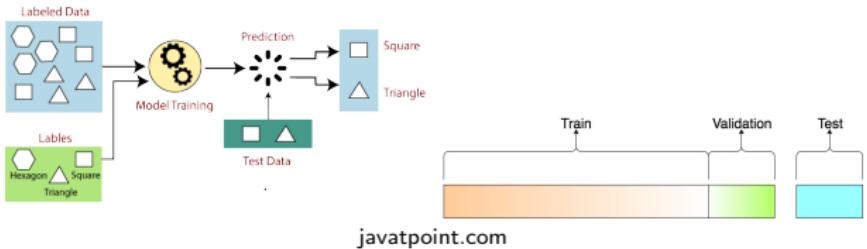
ML limitations

References

Stages of Supervised Learning

Machine Learning
in Astronomy

Reza Monadi



► Training:

1. Select a model
2. Train the model by the training set
3. Validate the model by the validation set

Machine Learning
in Astronomy

Overview

Big Data

Definition of BIG DATA

Astronomical Surveys

Supervised ML

Overview

Support Vector Machine

Neural Network

Unsupervised ML

Overview

Agglomerative clustering

LOF

PCA

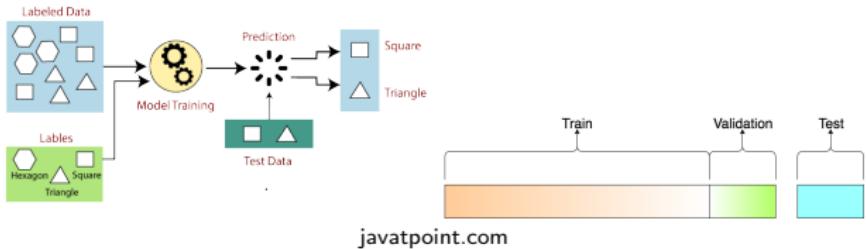
ML limitations

References

Stages of Supervised Learning

Machine Learning
in Astronomy

Reza Monadi



► Training:

1. Select a model
2. Train the model by the training set
3. Validate the model by the validation set
4. Select the optimum model

Machine Learning
in Astronomy

Overview

Big Data

Definition of BIG DATA

Astronomical Surveys

Supervised ML

Overview

Support Vector Machine

Neural Network

Unsupervised ML

Overview

Agglomerative clustering

LOF

PCA

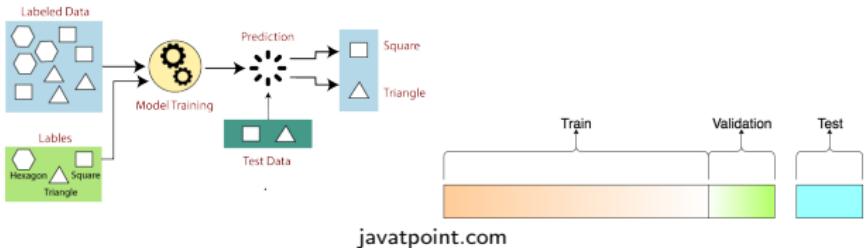
ML limitations

References

Stages of Supervised Learning

Machine Learning
in Astronomy

Reza Monadi



► Training:

1. Select a model
2. Train the model by the training set
3. Validate the model by the validation set
4. Select the optimum model

► Testing:

Overview

Big Data

Definition of BIG DATA

Astronomical Surveys

Supervised ML

Overview

Support Vector Machine

Neural Network

Unsupervised ML

Overview

Agglomerative clustering

LOF

PCA

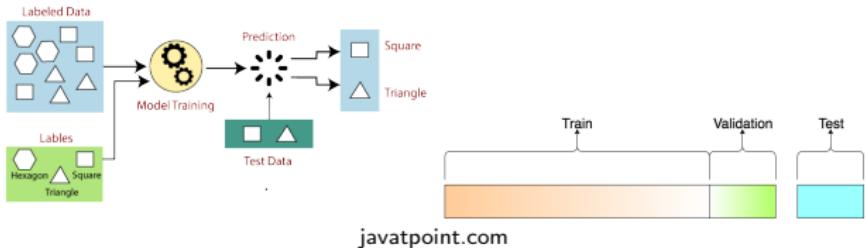
ML limitations

References

Stages of Supervised Learning

Machine Learning
in Astronomy

Reza Monadi



► Training:

1. Select a model
2. Train the model by the training set
3. Validate the model by the validation set
4. Select the optimum model

► Testing:

1. Test learned model by an unseen part of the data-set.

Unsupervised ML

Overview
Agglomerative clustering
LOF
PCA

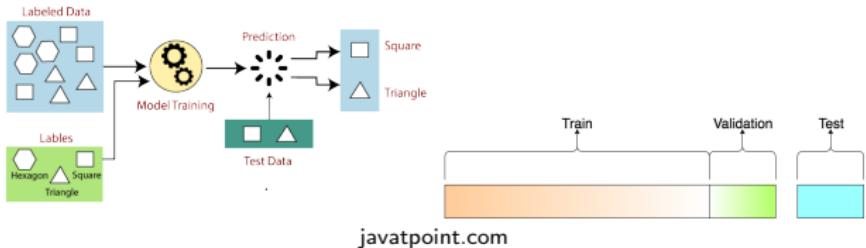
ML limitations

References

Stages of Supervised Learning

Machine Learning
in Astronomy

Reza Monadi



► Training:

1. Select a model
2. Train the model by the training set
3. Validate the model by the validation set
4. Select the optimum model

► Testing:

1. Test learned model by an unseen part of the data-set.
2. Use the best model for predictions about unseen observations.

Supervised ML

Overview
Support Vector Machine
Neural Network

Unsupervised ML

Overview
Agglomerative clustering
LOF
PCA

ML limitations

References

How Support Vector Machine Works?

Machine Learning
in Astronomy

Reza Monadi

Overview

Big Data

Definition of BIG DATA

Astronomical Surveys

Supervised ML

Overview

Support Vector Machine

Neural Network

Unsupervised ML

Overview

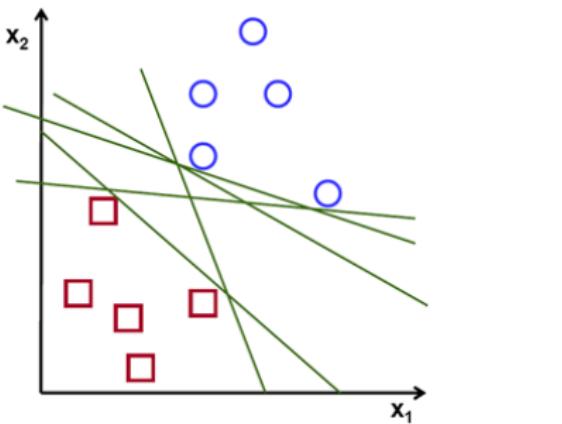
Agglomerative clustering

LOF

PCA

ML limitations

References



How Support Vector Machine Works?

Machine Learning
in Astronomy

Reza Monadi

Overview

Big Data

Definition of BIG DATA

Astronomical Surveys

Supervised ML

Overview

Support Vector Machine

Neural Network

Unsupervised ML

Overview

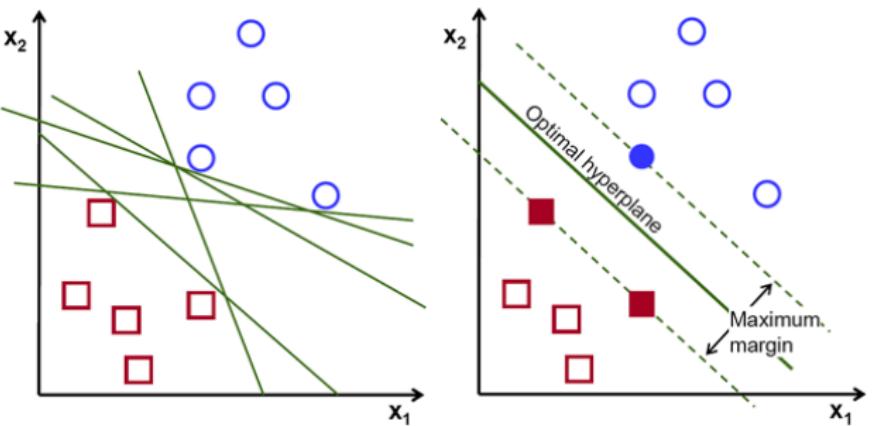
Agglomerative clustering

LOF

PCA

ML limitations

References



towardsdatascience.com

Classifying Pre-Main-Sequence Stars using SVM

Machine Learning
in Astronomy

Reza Monadi

Hubble Tarantula Treasury Project – VI. Identification of Pre-Main-Sequence Stars using Machine Learning techniques

Victor F. Ksoll,^{1,2,*} Dimitrios A. Gouliermis,^{1,3,†} Ralf S. Klessen,¹ Eva K. Grebel,⁴ Elena Sabbi,⁵ Jay Anderson,⁵ Daniel J. Lennon,⁶ Michele Cignoni,⁷ Guido de Marchi,⁸ Linda J. Smith,⁹ Monica Tosi,¹⁰ and Roeland P. van der Marel⁵

¹Institut für Theoretische Astrophysik, Zentrum für Astronomie der Universität Heidelberg, Albert-Ueberle-Str. 2, 69120 Heidelberg, Germany

²Interdisciplinary Center for Scientific Computing, University of Heidelberg, Mathematikon, Im Neuenheimer Feld 205, 69120 Heidelberg, Germany

³Max Planck Institute for Astronomy, Königstuhl 17, 69117 Heidelberg, Germany

⁴Astronomisches Rechen-Institut, Zentrum für Astronomie der Universität Heidelberg, Mönchhofstr. 12-14, 69120 Heidelberg, Germany

⁵Space Telescope Science Institute, 3700 San Martin Drive, Baltimore, MD 21218, USA

⁶ESA – European Space Astronomy Center, Apdo. de Correos 78, E-28691 Associate Villanueva de la Cañada, Madrid, Spain

⁷Department of Physics, University of Pisa, Largo Pontecorvo 3, I-56127 Pisa, Italy

⁸European Space Research and Technology Centre, Keplerlaan 1, 2200 AG Noordwijk, Netherlands

⁹European Space Agency and Space Telescope Science Institute, 3700 San Martin Drive, Baltimore, MD 21218, USA

¹⁰INAF–Osservatorio Astronomico di Bologna, Via Ranzani 1, I-40127 Bologna, Italy

Draft version 21 May 2018

ABSTRACT

The Hubble Tarantula Treasury Project (HTTP) has provided an unprecedented photometric coverage of the entire star-burst region of 30 Doradus down to the half Solar mass limit. We use the deep stellar catalogue of HTTP to identify all the pre-main-sequence (PMS) stars of the region, i.e., stars that have not started their lives on the main-sequence yet. The photometric distinction of these stars from the more evolved populations is not a trivial task due to several factors that alter their colour-magnitude diagram positions. The identification of PMS stars requires, thus, sophisticated statistical methods. We employ Machine Learning Classification techniques on the HTTP survey of more than 800,000 sources to identify the PMS stellar content of the observed field. Our methodology consists of 1) carefully selecting the most probable low-mass PMS stellar population of the star-forming cluster NGC 2070, 2) using this sample to train classification algorithms to build a predictive model for PMS stars, and 3) applying this model in order to identify the most probable PMS content across the entire Tarantula Nebula. We employ Decision Tree, Random Forest and Support Vector Machine classifiers to categorise the stars as PMS and Non-PMS. The Random Forest and Support Vector Machine provided the most accurate models, predicting about 20,000 sources with a candidateness probability higher than 50 percent, and almost 10,000 PMS candidates with a probability higher than 95 percent. This is the richest and most accurate photometric catalogue of extragalactic PMS candidates across the extent of a whole star-forming complex.

Overview

Big Data

Definition of BIG DATA

Astronomical Surveys

Supervised ML

Overview

Support Vector Machine

Neural Network

Unsupervised ML

Overview

Agglomerative clustering

LOF

PCA

ML limitations

References

Classifying Pre-Main-Sequence Stars using SVM

Overview

Big Data

Definition of BIG DATA

Astronomical Surveys

Supervised ML

Overview

Support Vector Machine

Neural Network

Unsupervised ML

Overview

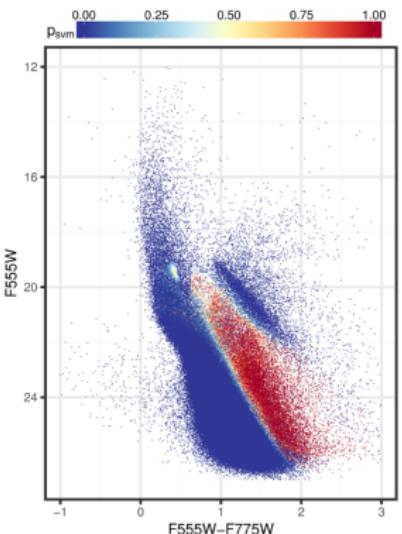
Agglomerative clustering

LOF

PCA

ML limitations

References



Credit: Ksoll, Victor F., et al. 2018

1. Select the most certain PMSs as the training set

Classifying Pre-Main-Sequence Stars using SVM

Overview

Big Data

Definition of BIG DATA

Astronomical Surveys

Supervised ML

Overview

Support Vector Machine

Neural Network

Unsupervised ML

Overview

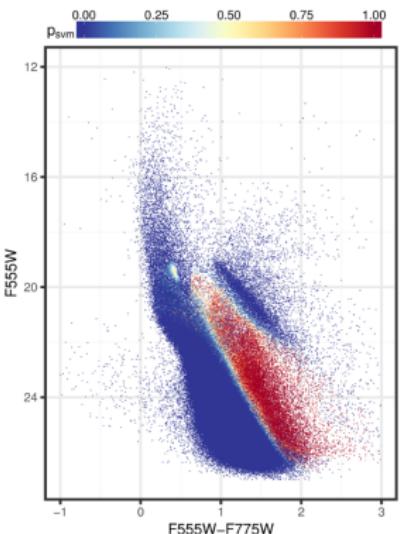
Agglomerative clustering

LOF

PCA

ML limitations

References



Credit: Ksoll, Victor F., et al. 2018

1. Select the most certain PMSs as the training set
2. Train the SVM classifier

Classifying Pre-Main-Sequence Stars using SVM

Overview

Big Data

Definition of BIG DATA

Astronomical Surveys

Supervised ML

Overview

Support Vector Machine

Neural Network

Unsupervised ML

Overview

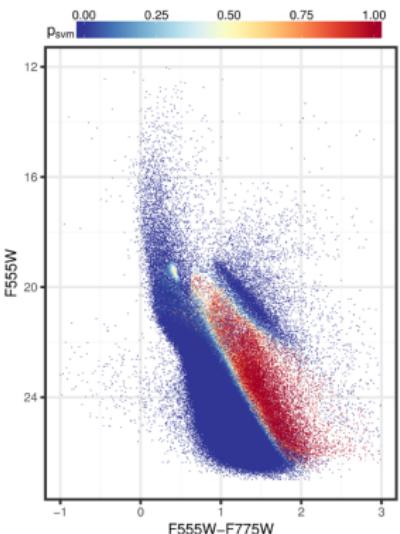
Agglomerative clustering

LOF

PCA

ML limitations

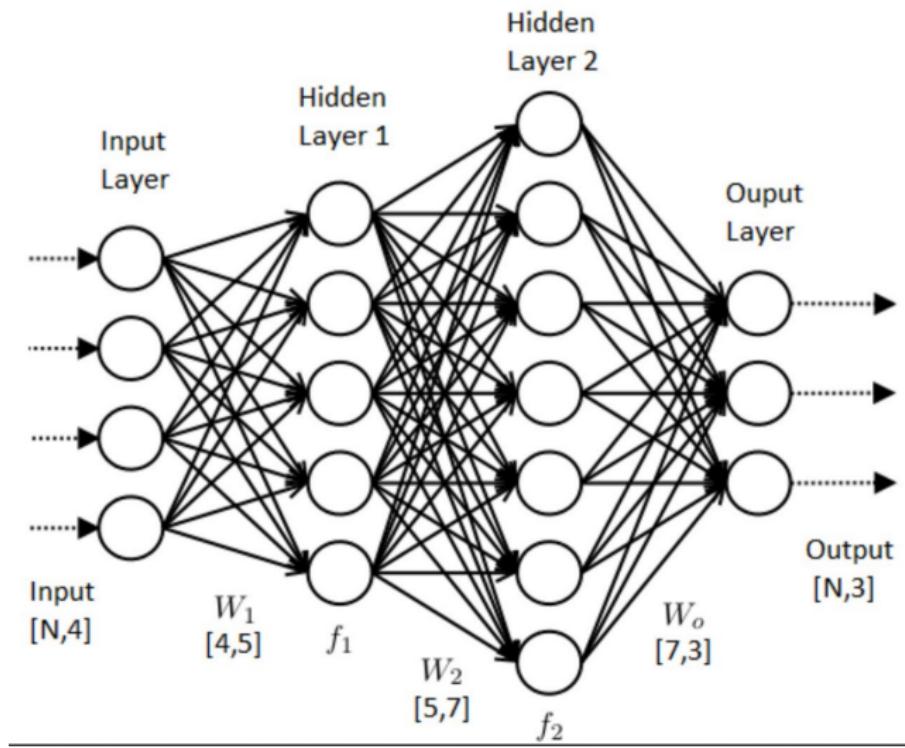
References



Credit: Ksoll, Victor F., et al. 2018

1. Select the most certain PMSs as the training set
2. Train the SVM classifier
3. Use the trained classifier for the entire Tarantula Nebula

How Neural Network works?



Credit: Dalya 2019

Convolutional Neural Network and DLAs

Machine Learning
in Astronomy

Reza Monadi

Deep Learning of Quasar Spectra to Discover and Characterize Damped Ly α Systems

David Parks,¹ J. Xavier Prochaska,² Shawfeng Dong,³ Zheng Cai,^{2,4}

¹Computer Science, UC Santa Cruz, 1156 High St., Santa Cruz, CA 95064 USA — dfparks@ucsc.edu

²Astronomy & Astrophysics, UC Santa Cruz, 1156 High St., Santa Cruz, CA 95064 USA

³Applied Mathematics & Statistics, UC Santa Cruz, 1156 High St., Santa Cruz, CA 95064 USA

⁴Hubble Fellow

Accepted XXX. Received YYY; in original form ZZZ

ABSTRACT

We have designed, developed, and applied a convolutional neural network (CNN) architecture using multi-task learning to search for and characterize strong HI Ly α absorption in quasar spectra. Without any explicit modeling of the quasar continuum nor application of the predicted line-profile for Ly α from quantum mechanics, our algorithm predicts the presence of strong HI absorption and estimates the corresponding redshift z_{abs} and HI column density N_{HI} , with emphasis on damped Ly α systems (DLAs, absorbers with $N_{\text{HI}} \geq 2 \times 10^{20} \text{ cm}^{-2}$). We tuned the CNN model using a custom training set of DLAs injected into DLA-free quasar spectra from the Sloan Digital Sky Survey (SDSS), data release 5 (DR5). Testing on a held-back validation set demonstrates a high incidence of DLAs recovered by the algorithm (97.4% as DLAs and 99% as an HI absorber with $N_{\text{HI}} > 10^{19.5} \text{ cm}^{-2}$) and excellent estimates for z_{abs} and N_{HI} . Similar results are obtained against a human-generated survey of the SDSS DR5 dataset. The algorithm yields a low incidence of false positives and negatives but is challenged by overlapping DLAs and/or very high N_{HI} systems. We have applied this CNN model to the quasar spectra of SDSS-DR7 and the Baryonic Oscillation Spectroscopic Survey (BOSS, data release 12) and provide catalogs of 4,913 and 50,969 DLAs respectively (including 1,659 and 9,230 high-confidence DLAs that were previously unpublished). This work validates the application of deep learning techniques to astronomical spectra for both classification and quantitative measurements.

Overview

Big Data

Definition of BIG DATA

Astronomical Surveys

Supervised ML

Overview

Support Vector Machine

Neural Network

Unsupervised ML

Overview

Agglomerative clustering

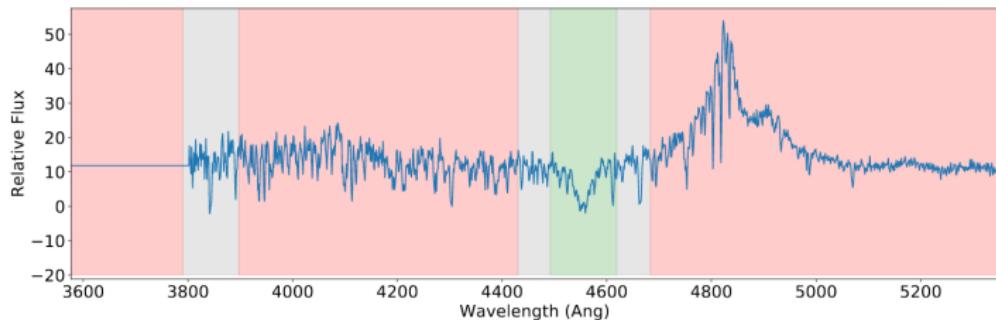
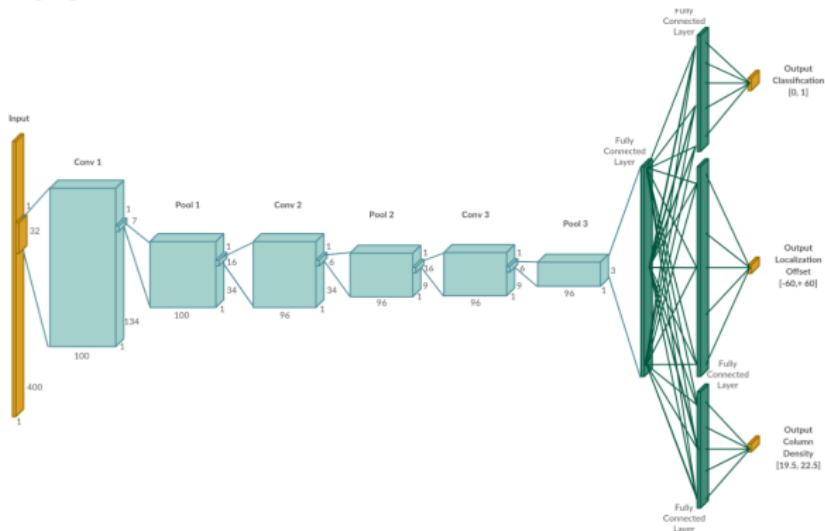
LOF

PCA

ML limitations

References

Training DLA/noDLA models: 1D image recognition



Overview

Big Data

Definition of BIG DATA

Astronomical Surveys

Supervised ML

Overview

Support Vector Machine

Neural Network

Unsupervised ML

Overview

Agglomerative clustering

LOF

PCA

ML limitations

References

Multitask CNN: Location, Density and Class

Machine Learning
in Astronomy

Reza Monadi

Overview

Big Data

Definition of BIG DATA

Astronomical Surveys

Supervised ML

Overview

Support Vector Machine

Neural Network

Unsupervised ML

Overview

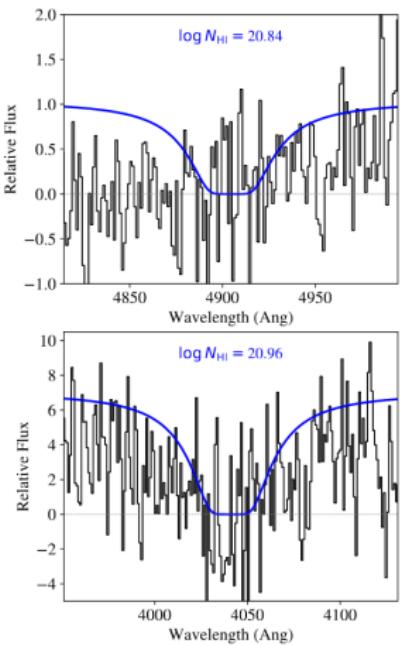
Agglomerative clustering

LOF

PCA

ML limitations

References



Credit: Parks et. al. 2017

How unsupervised learning works?

Machine Learning
in Astronomy

Reza Monadi

Overview

Big Data

Definition of BIG DATA

Astronomical Surveys

Supervised ML

Overview

Support Vector Machine

Neural Network

Unsupervised ML

Overview

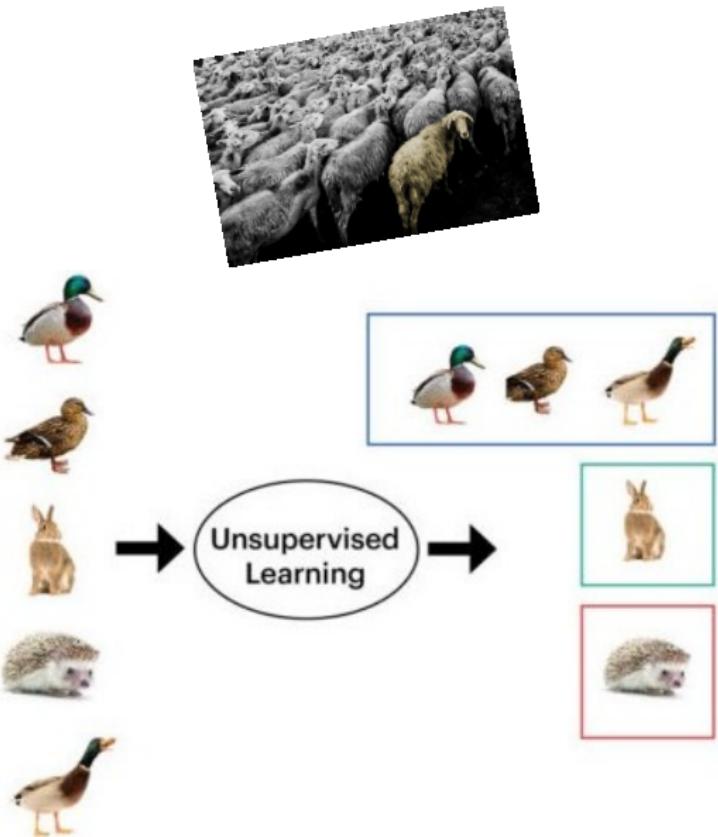
Agglomerative clustering

LOF

PCA

ML limitations

References



Unsupervised Learning and knowledge discovery

Machine Learning
in Astronomy

Reza Monadi

Overview

Big Data

Definition of BIG DATA

Astronomical Surveys

Supervised ML

Overview

Support Vector Machine

Neural Network

Unsupervised ML

Overview

Agglomerative clustering

LOF

PCA

ML limitations

References

Unsupervised Learning and knowledge discovery

Machine Learning
in Astronomy

Reza Monadi

Overview

Big Data

Definition of BIG DATA

Astronomical Surveys

Supervised ML

Overview

Support Vector Machine

Neural Network

Unsupervised ML

Overview

Agglomerative clustering

LOF

PCA

ML limitations

References

Unsupervised Learning and knowledge discovery

Machine Learning
in Astronomy

Reza Monadi

Overview

Big Data

Definition of BIG DATA

Astronomical Surveys

Supervised ML

Overview

Support Vector Machine

Neural Network

Unsupervised ML

Overview

Agglomerative clustering

LOF

PCA

ML limitations

References

How Agglomerative clustering works?

Machine Learning
in Astronomy

Reza Monadi

Overview

Big Data

Definition of BIG DATA

Astronomical Surveys

Supervised ML

Overview

Support Vector Machine

Neural Network

Unsupervised ML

Overview

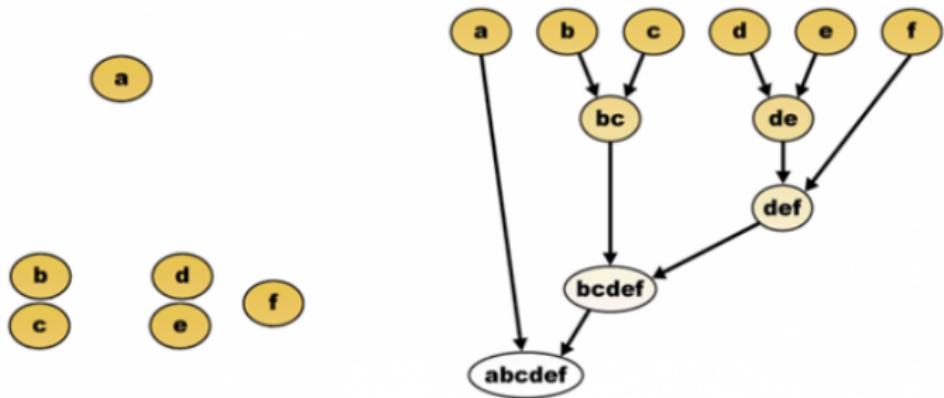
Agglomerative clustering

LOF

PCA

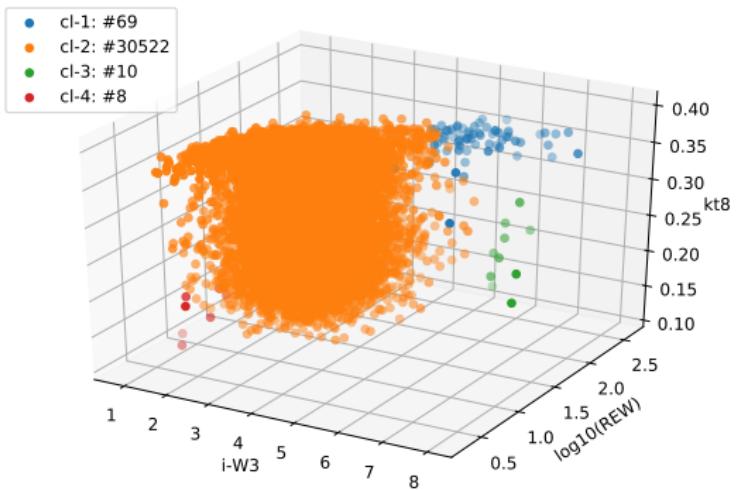
ML limitations

References



Credit: blog.minitab.com

Agglomerative clustering finds weird quasars?



Machine Learning
in Astronomy

Reza Monadi

Overview

Big Data

Definition of BIG DATA

Astronomical Surveys

Supervised ML

Overview

Support Vector Machine

Neural Network

Unsupervised ML

Overview

Agglomerative clustering

LOF

PCA

ML limitations

References

Overview

Big Data

Definition of BIG DATA

Astronomical Surveys

Supervised ML

Overview

Support Vector Machine

Neural Network

Unsupervised ML

Overview

Agglomerative clustering

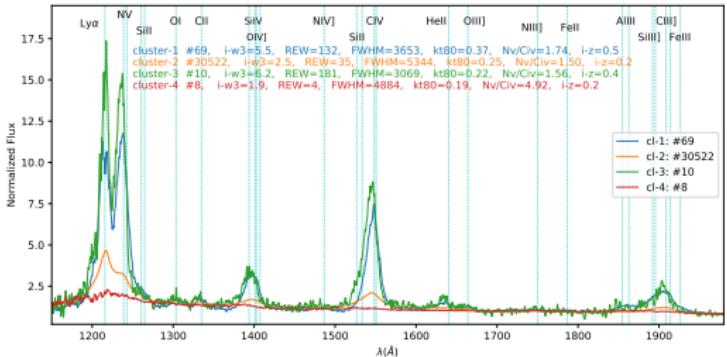
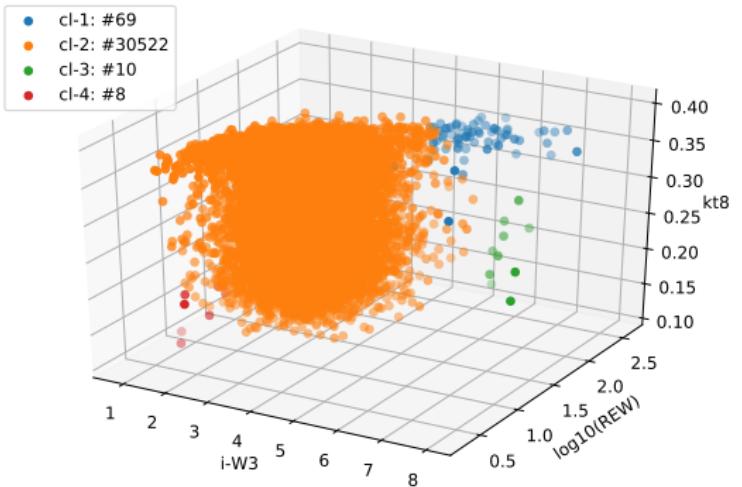
LOF

PCA

ML limitations

References

Agglomerative clustering finds weird quasars?



Local Outlier Factor finds outliers

Overview

Big Data

Definition of BIG DATA

Astronomical Surveys

Supervised ML

Overview

Support Vector Machine

Neural Network

Unsupervised ML

Overview

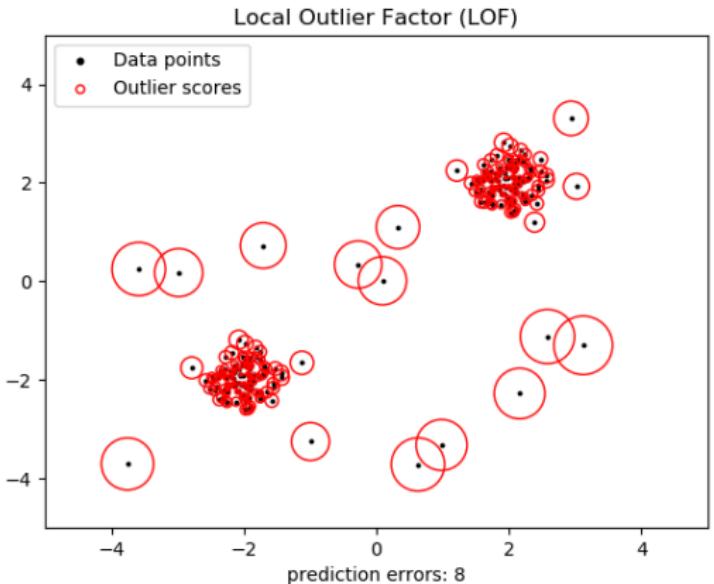
Agglomerative clustering

LOF

PCA

ML limitations

References



Overview

Big Data

Definition of BIG DATA

Astronomical Surveys

Supervised ML

Overview

Support Vector Machine

Neural Network

Unsupervised ML

Overview

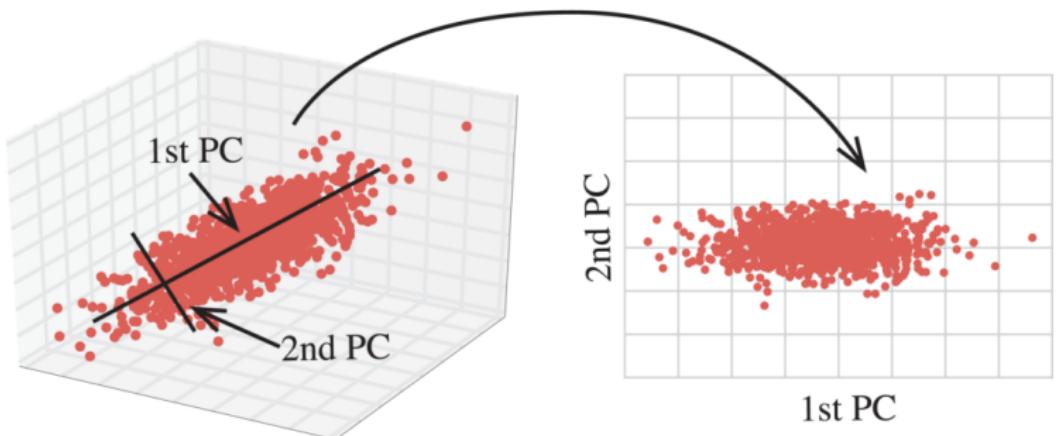
Agglomerative clustering

LOF

PCA

ML limitations

References



Credit: medium.com

Limitation of ML in astronomy

Machine Learning
in Astronomy

Reza Monadi

- ▶ No pre-processing = misleading results

Overview

Big Data

Definition of BIG DATA

Astronomical Surveys

Supervised ML

Overview

Support Vector Machine

Neural Network

Unsupervised ML

Overview

Agglomerative clustering

LOF

PCA

ML limitations

References

Limitation of ML in astronomy

Machine Learning
in Astronomy

Reza Monadi

Overview

Big Data

Definition of BIG DATA

Astronomical Surveys

Supervised ML

Overview

Support Vector Machine

Neural Network

Unsupervised ML

Overview

Agglomerative clustering

LOF

PCA

ML limitations

References

Limitation of ML in astronomy

Machine Learning
in Astronomy

Reza Monadi

Overview

Big Data

Definition of BIG DATA

Astronomical Surveys

Supervised ML

Overview

Support Vector Machine

Neural Network

Unsupervised ML

Overview

Agglomerative clustering

LOF

PCA

ML limitations

References

Limitation of ML in astronomy

Machine Learning
in Astronomy

Reza Monadi

Overview

Big Data

Definition of BIG DATA

Astronomical Surveys

Supervised ML

Overview

Support Vector Machine

Neural Network

Unsupervised ML

Overview

Agglomerative clustering

LOF

PCA

ML limitations

References

Overview

Big Data

Definition of BIG DATA

Astronomical Surveys

Supervised ML

Overview

Support Vector Machine

Neural Network

Unsupervised ML

Overview

Agglomerative clustering

LOF

PCA

ML limitations

References

Limitation of ML in astronomy

- ▶ No pre-processing = misleading results
- ▶ Training on a biased sample = wrong predictions
- ▶ ML is not plug and play, it needs **expertise**
- ▶ There are a variety of methods for each task:
pre-knowledge is needed
- ▶ Some algorithms sacrifice interpretation for functionality:
black boxes

Overview

Big Data

Definition of BIG DATA

Astronomical Surveys

Supervised ML

Overview

Support Vector Machine

Neural Network

Unsupervised ML

Overview

Agglomerative clustering

LOF

PCA

ML limitations

References

Limitation of ML in astronomy

- ▶ No pre-processing = misleading results
- ▶ Training on a biased sample = wrong predictions
- ▶ ML is not plug and play, it needs **expertise**
- ▶ There are a variety of methods for each task:
pre-knowledge is needed
- ▶ Some algorithms sacrifice interpretation for functionality:
black boxes
- ▶ ML algorithms are mostly designed for industry (e.g.
Google, Walmart, Uber, ...) not for astronomy.

[Overview](#)[Big Data](#)[Definition of BIG DATA](#)[Astronomical Surveys](#)[Supervised ML](#)[Overview](#)[Support Vector Machine](#)[Neural Network](#)[Unsupervised ML](#)[Overview](#)[Agglomerative clustering](#)[LOF](#)[PCA](#)[ML limitations](#)[References](#)

Limitation of ML in astronomy

- ▶ No pre-processing = misleading results
- ▶ Training on a biased sample = wrong predictions
- ▶ ML is not plug and play, it needs **expertise**
- ▶ There are a variety of methods for each task:
pre-knowledge is needed
- ▶ Some algorithms sacrifice interpretation for functionality:
black boxes
- ▶ ML algorithms are mostly designed for industry (e.g.
Google, Walmart, Uber, ...) not for astronomy.
- ▶ Collaboration of astronomers with computer scientists
needs to be much stronger

References

1. **ML, General:** Baron, Dalya. "Machine learning in astronomy: A practical overview." arXiv preprint arXiv:1904.07248 (2019).
2. **Big Data:** Zhang, Yanxia, and Yongheng Zhao. "Astronomy in the big data era." Data Science Journal 14 (2015).
3. **SVM:** Ksoll, Victor F., et al. "Hubble Tarantula Treasury ProjectVI. Identification of pre-main-sequence stars using machine-learning techniques." Monthly Notices of the Royal Astronomical Society 479.2 (2018): 2389-2414.
4. **Decision Tree:** Vasconcellos, E. C., et al. "Decision tree classifiers for star/galaxy separation." The Astronomical Journal 141.6 (2011): 189.
5. **Neural Network:** Parks, David, et al. "Deep learning of quasar spectra to discover and characterize damped Ly systems." Monthly Notices of the Royal Astronomical Society 476.1 (2018): 1151-1168.
6. **Agglomerative clustering:** R. Monadi, F. Hamann, S. Bird, Precise Selection of Extremely Red Quasars, (in preparation)

Overview

Big Data

Definition of BIG DATA

Astronomical Surveys

Supervised ML

Overview

Support Vector Machine

Neural Network

Unsupervised ML

Overview

Agglomerative clustering

LOF

PCA

ML limitations

References

