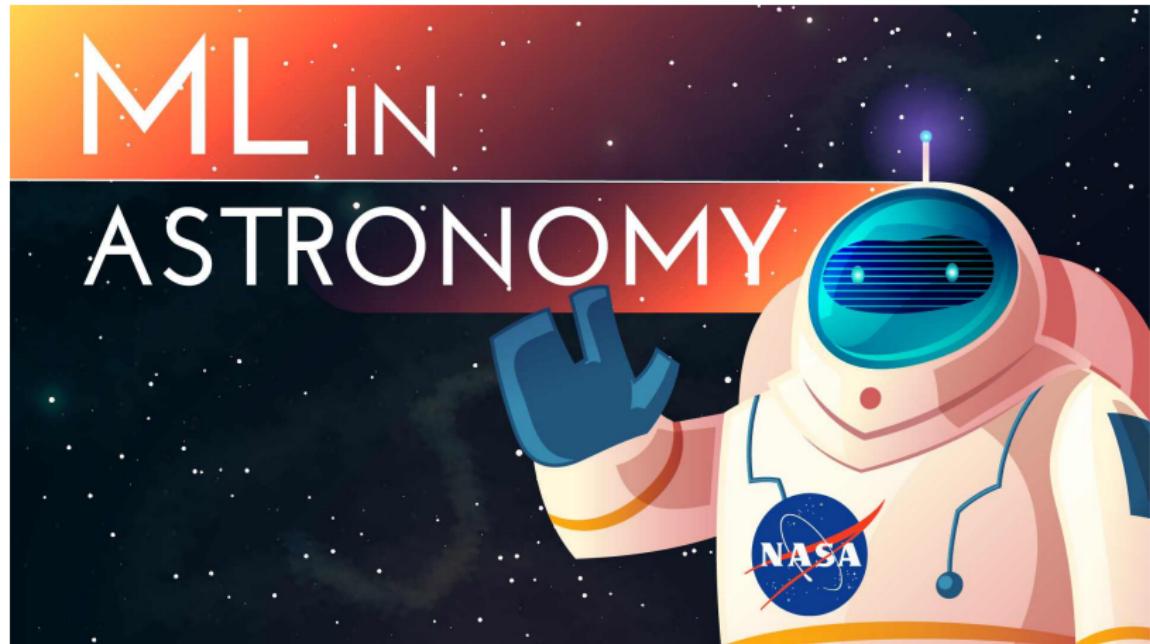


# Machine Learning in Astronomy

Reza Monadi

UC Riverside

May 14, 2020



credit: 365datascience.com

- How astronomy is tied to **BIG DATA?**

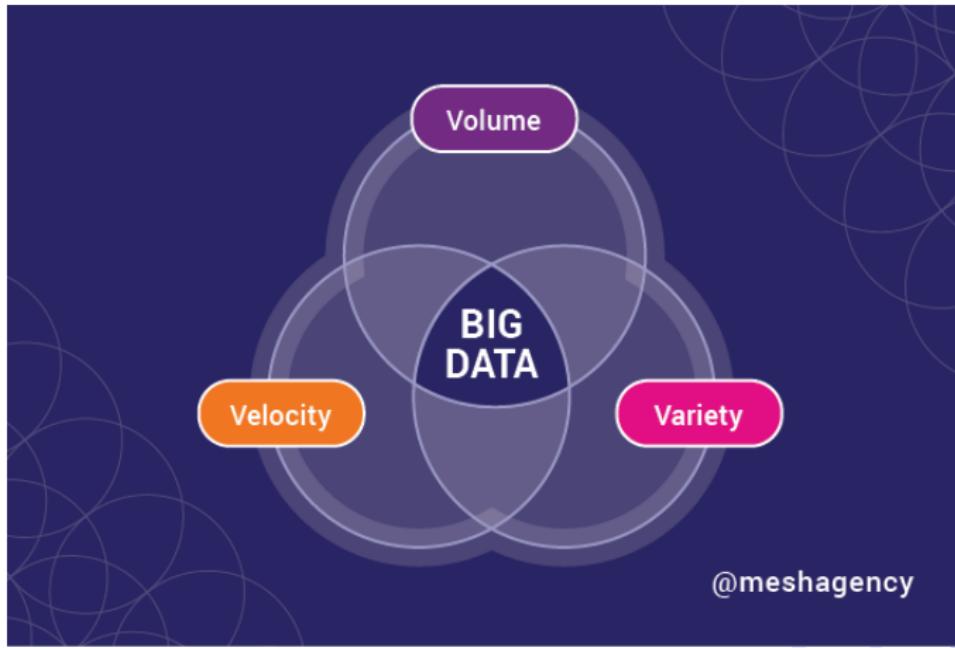
- How astronomy is tied to **BIG DATA?**
- What is supervised and unsupervised **ML?**

- How astronomy is tied to **BIG DATA**?
- What is supervised and unsupervised **ML**?
- How to implement **ML** algorithms in astronomy?

- How astronomy is tied to **BIG DATA**?
- What is supervised and unsupervised **ML**?
- How to implement **ML** algorithms in astronomy?
- How **ML** helps **SKA**?

- How astronomy is tied to **BIG DATA**?
- What is supervised and unsupervised **ML**?
- How to implement **ML** algorithms in astronomy?
- How **ML** helps **SKA**?
- What are the pitfalls of **ML** in astronomy?

# What is BIG DATA?



# VVV in astronomy

- Volume: larger quantities of data by better facilitates

# VVV in astronomy

- Volume: larger quantities of data by better facilitates
- Velocity: Higher speed of getting data

# VVV in astronomy

- Volume: larger quantities of data by better facilitates
- Velocity: Higher speed of getting data
- Verity: More complex structures of data

# Astronomical surveys are really Astronomical.

# Sloan Digital Sky Server $\Rightarrow$ 40 TB



# Sloan Digital Sky Server $\Rightarrow$ 40 TB

## Imaging statistics

<b>Total unique area covered</b>	14,555 square degrees
<b>Total area of imaging (including overlaps)</b>	31,637 square degrees (excluding supernova runs)
<b>Individual image field size</b>	1361x2048 pixels (0.0337 square degrees)
<b>Number of individual fields</b>	938,046 (excluding supernova runs)
<b>Number of catalog objects</b>	1,231,051,050
<b>Number of unique detections</b>	932,891,133
<b>Median PSF FWHM, <i>r</i>-band</b>	1.3 arcsec
<b>Pixel scale</b>	0.396 arcsec
<b>Exposure time per band</b>	53.9 sec
<b>Time difference between observations of each band</b>	71.72 sec (in <i>riuzg</i> order)
<b>Global astrometric precision</b>	0.1 arcsec rms (absolute)

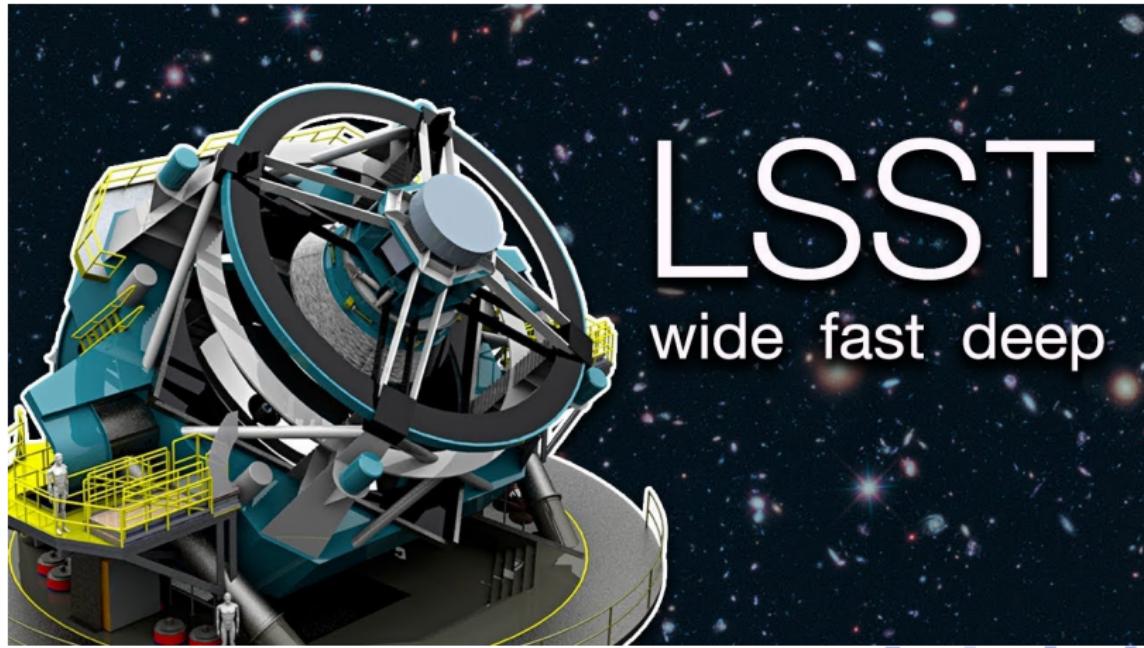
# Sloan Digital Sky Server $\Rightarrow$ 40 TB

## Optical spectroscopy data statistics

### All programs combined

<b>Total spectra</b>	5,789,200
<b>Useful spectra</b>	4,846,156
<b>Galaxies</b>	2,863,635
<b>Quasars</b>	960,678
<b>Stars</b>	1,021,843
<b>Sky</b>	475,531
<b>Standards</b>	108,603
<b>Unknown</b>	352,320

## Large Synaptic Survey Telescope $\Rightarrow$ 200 PB



# Large Synaptic Survey Telescope $\Rightarrow$ 200 PB

- 800+ panoramic images each night

# Large Synaptic Survey Telescope $\Rightarrow$ 200 PB

- 800+ panoramic images each night
- 3.2 billion-pixel camera

# Large Synoptic Survey Telescope $\Rightarrow$ 200 PB

- 800+ panoramic images each night
- 3.2 billion-pixel camera
- Recording the entire visible sky twice each week

# Large Synoptic Survey Telescope $\Rightarrow$ 200 PB

- 800+ panoramic images each night
- 3.2 billion-pixel camera
- Recording the entire visible sky twice each week
- Each patch of sky will be visited 1000 times

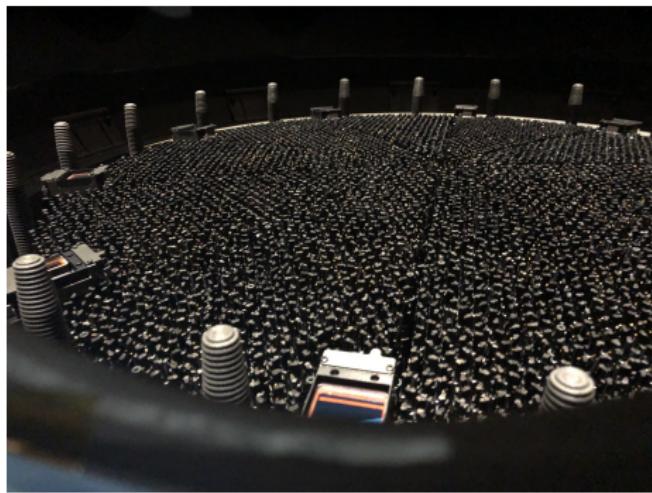
# Zwicky Transient Facility



# Zwicky Transient Facility

Filter(s)	#PSFcat- <i>sci</i> sources	#Aperturecat- <i>sci</i> sources	#PSFcat- <i>ref</i> sources	#Aperturecat- <i>ref</i> sources
<i>g</i>	34,799,157,939	22,008,868,967	2,053,507,697	656,772,487
<i>r</i>	64,698,158,346	40,249,465,415	2,954,247,127	1,013,085,916
<i>i</i>	868,042,655	495,178,759	557,831,277	163,240,538
<i>g + r + i</i>	100,365,358,940	62,753,513,141	5,565,586,101	1,833,098,941

# Dark Energy Spectroscopic Instrument



 DARK ENERGY  
SPECTROSCOPIC  
INSTRUMENT

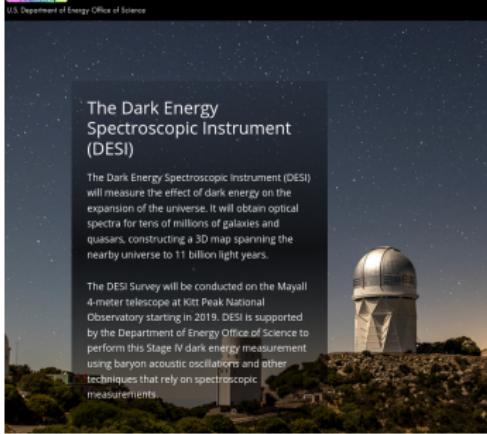
U.S. Department of Energy Office of Science

/ SCIENCE / / INSTRUMENT / / COLLABORATION / / FUNDING

The Dark Energy Spectroscopic Instrument (DESI)

The Dark Energy Spectroscopic Instrument (DESI) will measure the effect of dark energy on the expansion of the universe. It will obtain optical spectra for tens of millions of galaxies and quasars, constructing a 3D map spanning the nearby universe to 11 billion light years.

The DESI Survey will be conducted on the Mayall 4-meter telescope at Kitt Peak National Observatory starting in 2019. DESI is supported by the Department of Energy Office of Science to perform this Stage IV dark energy measurement using baryon acoustic oscillations and other techniques that rely on spectroscopic measurements.



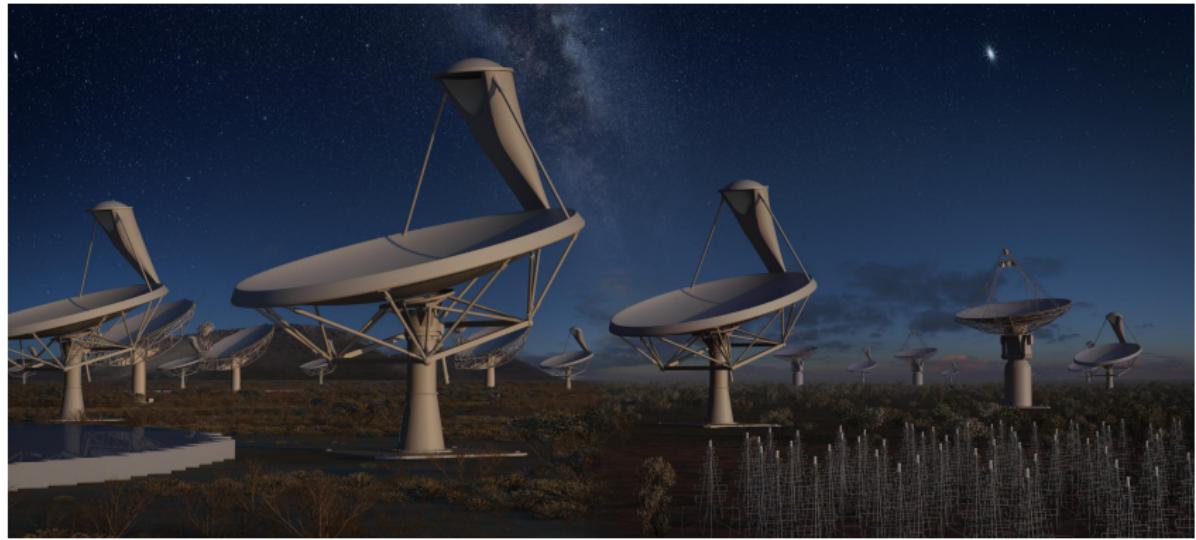
# Dark Energy Spectroscopic Instrument

- Spectra of 25 M galaxies, quasars, stars.

# Dark Energy Spectroscopic Instrument

- Spectra of 25 M galaxies, quasars, stars.
- 5000 spectra per exposure

# Square Kilometer Array $\Rightarrow$ 4.6 EB



# Large surveys

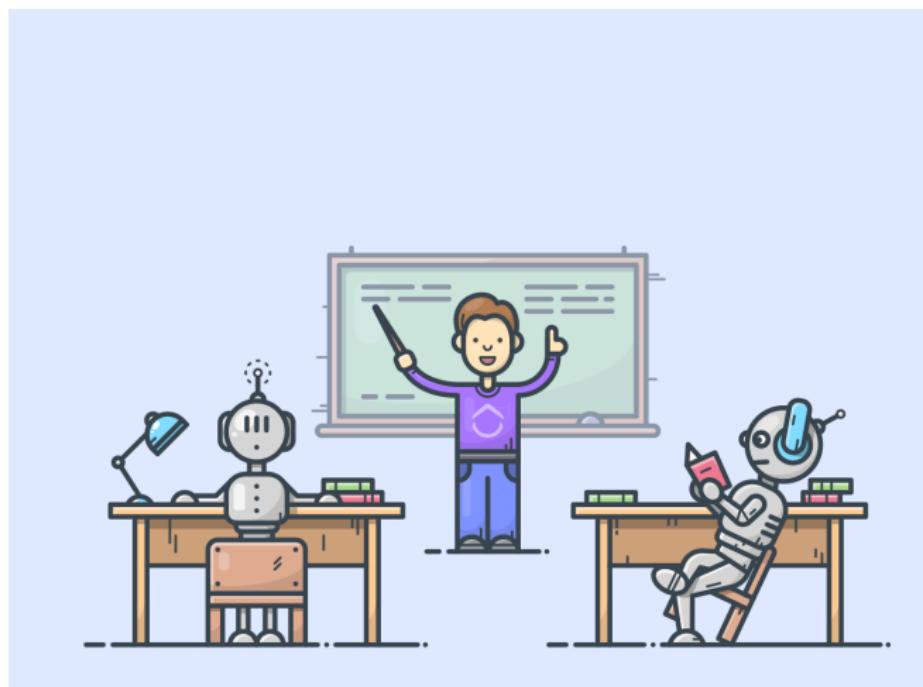
Large surveys ⇒

# Large surveys ⇒ Big Data

Large surveys  $\Rightarrow$  Big Data  $\xrightarrow{ML}$

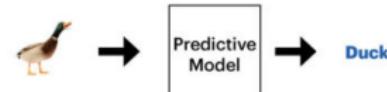
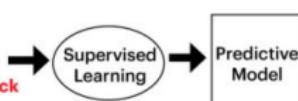
Large surveys  $\Rightarrow$  Big Data  $\xrightarrow{ML}$   
Astronomy Knowledge

# Supervised ML vs. Unsupervised ML

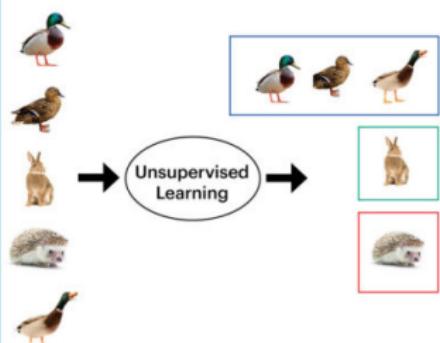


# Supervised ML vs. Unsupervised ML

## Supervised Learning (Classification Algorithm)



## Unsupervised Learning (Clustering Algorithm)



# Stages of Supervised Learning

- Training:

# Stages of Supervised Learning

- Training:
  - ① Select a model

# Stages of Supervised Learning

- Training:
  - 1 Select a model
  - 2 Set up hyper-parameters of model

# Stages of Supervised Learning

- Training:
  - ① Select a model
  - ② Set up hyper-parameters of model
  - ③ Teach the machine by training set

# Stages of Supervised Learning

- Training:
  - ① Select a model
  - ② Set up hyper-parameters of model
  - ③ Teach the machine by training set
- Validation:

# Stages of Supervised Learning

- Training:
  - ① Select a model
  - ② Set up hyper-parameters of model
  - ③ Teach the machine by training set
- Validation:
  - ① Change the hyper-parameters

# Stages of Supervised Learning

- Training:
  - ① Select a model
  - ② Set up hyper-parameters of model
  - ③ Teach the machine by training set
- Validation:
  - ① Change the hyper-parameters
  - ② Select the optimum hyper-parameters

# Stages of Supervised Learning

- Training:
  - ① Select a model
  - ② Set up hyper-parameters of model
  - ③ Teach the machine by training set
- Validation:
  - ① Change the hyper-parameters
  - ② Select the optimum hyper-parameters
- Testing:

# Stages of Supervised Learning

- Training:
  - ① Select a model
  - ② Set up hyper-parameters of model
  - ③ Teach the machine by training set
- Validation:
  - ① Change the hyper-parameters
  - ② Select the optimum hyper-parameters
- Testing:
  - ① Test learned model by an unseen part of the data-set.

# Stages of Supervised Learning

- Training:
  - ① Select a model
  - ② Set up hyper-parameters of model
  - ③ Teach the machine by training set
- Validation:
  - ① Change the hyper-parameters
  - ② Select the optimum hyper-parameters
- Testing:
  - ① Test learned model by an unseen part of the data-set.
  - ② Select the best model and use it for predictions.

# Supervised Learning vs Traditional Model Fitting ?



## Similarities

I

# Supervised Learning vs Traditional Model Fitting ?



## Similarities

- I • Both need a set of labeled measurements and a model

# Supervised Learning vs Traditional Model Fitting ?



## Similarities

- I • Both need a set of labeled measurements and a model



## Differences

# Supervised Learning vs Traditional Model Fitting ?



## Similarities

- I • Both need a set of labeled measurements and a model



## Differences

- I • SML: model is constructed based on input data

# Supervised Learning vs Traditional Model Fitting ?



## Similarities

- I • Both need a set of labeled measurements and a model



## Differences

- SML: model is constructed based on input data
- SML: model can be very nonlinear/complex

# Supervised Learning vs Traditional Model Fitting ?



## Similarities

- I • Both need a set of labeled measurements and a model



## Differences

- SML: model is constructed based on input data
- SML: model can be very nonlinear/complex
- TMF: model is predefined and has limited adaptivity

# Supervised Learning vs Traditional Model Fitting ?



## Similarities

- I • Both need a set of labeled measurements and a model



## Differences

- SML: model is constructed based on input data
- SML: model can be very nonlinear/complex
- TMF: model is predefined and has limited adaptivity
- SML: designed for predicting unseen data

# Supervised Learning vs Traditional Model Fitting ?



## Similarities

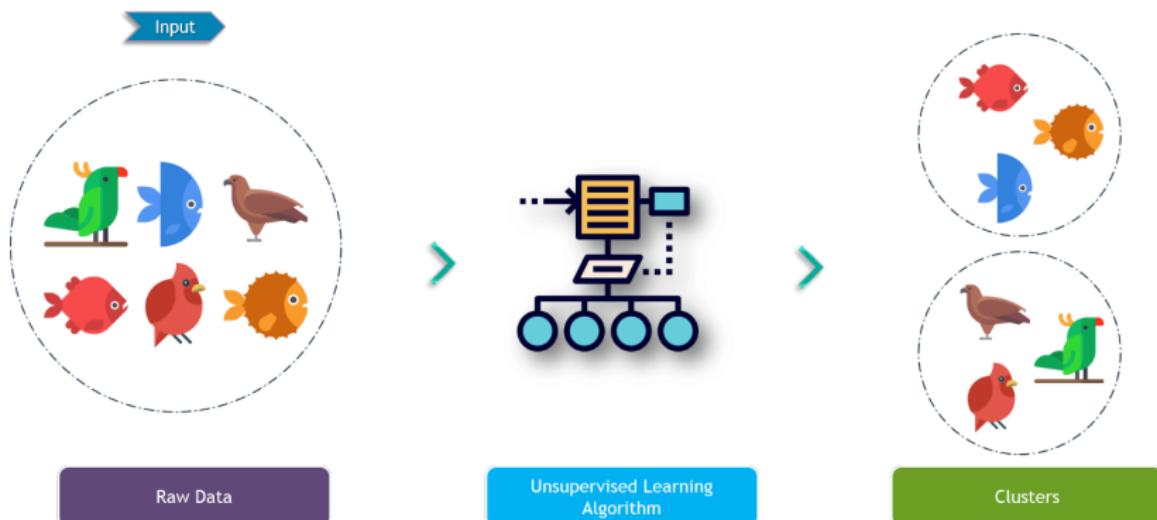
- I • Both need a set of labeled measurements and a model



## Differences

- SML: model is constructed based on input data
- SML: model can be very nonlinear/complex
- TMF: model is predefined and has limited adaptivity
- SML: designed for predicting unseen data
- TMF: infers relationships between features

# How unsupervised learning works?



# Unsupervised Learning and knowledge discovery

- Finding unknown patterns in the data set

# Unsupervised Learning and knowledge discovery

- Finding unknown patterns in the data set
- Distinguishing similar and dissimilar objects

# Unsupervised Learning and knowledge discovery

- Finding unknown patterns in the data set
- Distinguishing similar and dissimilar objects
- Helping for visualization in lower dimensions

# Unsupervised Learning and knowledge discovery

- Finding unknown patterns in the data set
- Distinguishing similar and dissimilar objects
- Helping for visualization in lower dimensions
- Recognizing odd objects (unknown unknowns)

# Unsupervised Learning and knowledge discovery

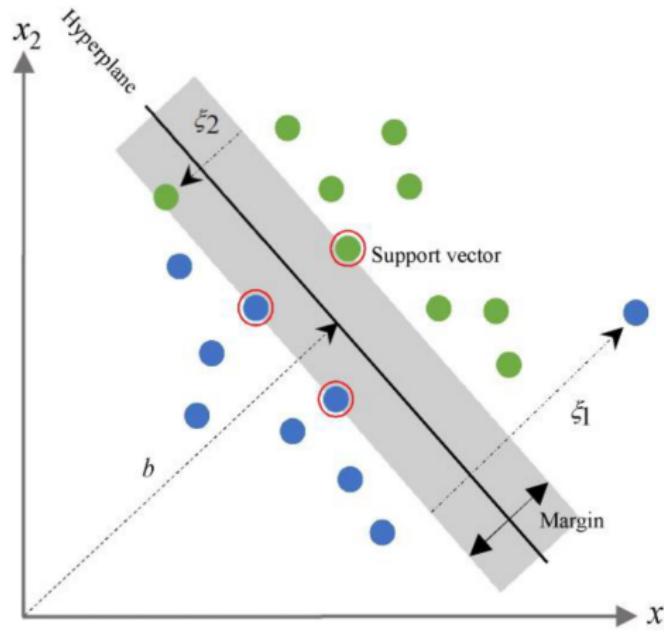
- Finding unknown patterns in the data set
- Distinguishing similar and dissimilar objects
- Helping for visualization in lower dimensions
- Recognizing odd objects (unknown unknowns)

DATA  $\xrightarrow{ML}$  Knowledge

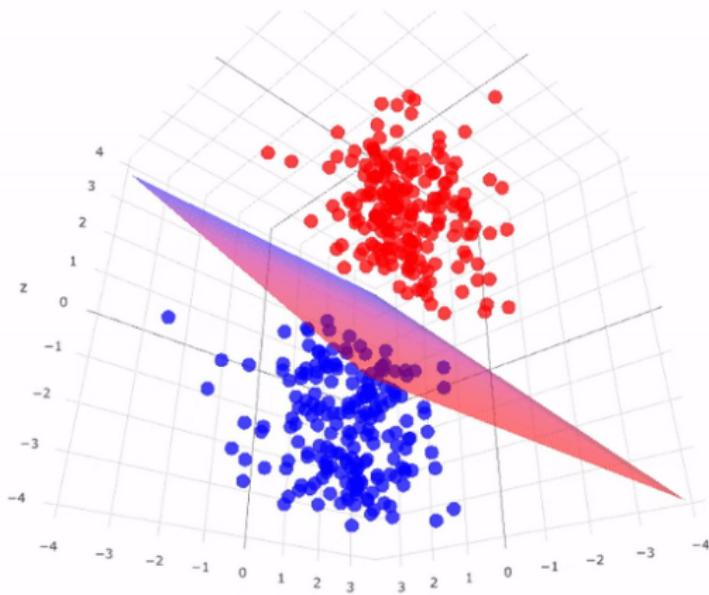
# • Classification

- Classification
- Regression

# How Support Vector Machine Works?



# Hyper-plane in Support Vector Machine



# Classifying Pre-Main-Sequence Stars using SVM

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

Victor F. Ksoll,<sup>1,2,\*</sup> Dimitrios A. Gouliermis,<sup>1,3,4</sup> Ralf S. Klessen,<sup>1</sup> Eva K. Grebel,<sup>4</sup> Elena Sabbi,<sup>5</sup> Jay Anderson,<sup>5</sup> Daniel J. Lennon,<sup>6</sup> Michele Cignoni,<sup>7</sup> Guido de Marchi,<sup>8</sup> Linda J. Smith,<sup>9</sup> Monica Tosi,<sup>10</sup> and Roeland P. van der Marel<sup>11</sup>

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

<sup>2</sup>Institute für Didaktik der Computerwissenschaften, University of Heidelberg, Mathematikon, Im Neuenheimer Feld 205, 69120 Heidelberg, Germany

<sup>3</sup>Max-Planck-Institut für Astronomie, Königstuhl 17, 69117 Heidelberg, Germany

<sup>4</sup>Astronoticias Research Institute, Zentrum für Astronomie der Universität Heidelberg, Mönchhofstr. 12-14, 69120 Heidelberg, Germany

<sup>5</sup>Space Telescope Science Institute, 3700 San Martin Drive, Baltimore, MD 21218, USA

<sup>6</sup>Departamento de Astrofísica, Universidad de La Rioja, Alfonso de la Castaño, 1, 26003 Logroño, Spain

<sup>7</sup>Departament d'Estadística i Investigació Operativa, Facultat d'Economia, Universitat de València, Madrid, Spain

<sup>8</sup>Departament de Física, Universitat de València, Avda. Marcelino Salado s/n, 46130 Paterna, Valencia, Spain

<sup>9</sup>European Space Research and Technology Centre, Keplerlaan 1, 2200 AG Noordwijk, Netherlands

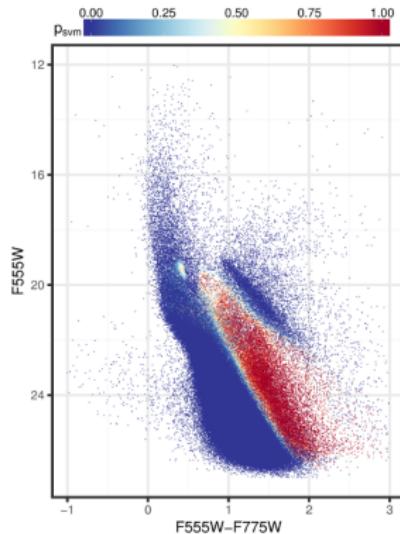
<sup>10</sup>European Space Agency and Space Telescope Science Institute, 3700 San Martin Drive, Baltimore, MD 21218, USA

<sup>11</sup>DINAF-Osservatorio Astronomico di Bologna, Via Ranzani 1, I-40127 Bologna, Italy

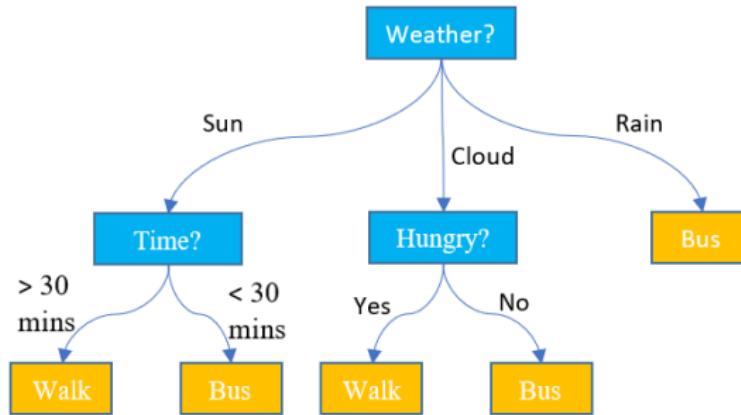
Draft version 21 May 2018

**ABSTRACT**  
The Hubble Tarantula Treasury Project (HTTP) has provided an unprecedented photometric coverage of the entire starburst 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 stars. We train three different machine learning models to identify the most probable the most probable low-mass PMS stellar population of the star-forming cluster NGC2070, 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 predicted the most probable models, predicting about 300,000 sources with a probability higher than 90 percent, among which 100,000 PMS stars 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.

# Classifying Pre-Main-Sequence Stars using SVM



# How Decision Tree works?



# Using DT for classifying galaxies and stars in SDSS?

THE ASTRONOMICAL JOURNAL, 141:189 (12pp), 2011 June  
© 2011. The American Astronomical Society. All rights reserved. Printed in the U.S.A.

doi:10.1088/0004-6256/141/6/189

## DECISION TREE CLASSIFIERS FOR STAR/GALAXY SEPARATION

E. C. VASCONCELLOS<sup>1</sup>, R. R. DE CARVALHO<sup>2</sup>, R. R. GAL<sup>3</sup>, F. L. LABARBERA<sup>4</sup>,

H. V. CAPELATO<sup>5</sup>, H. FRAGO CAMPOS VELHO<sup>6</sup>, M. TREVISAN<sup>6</sup>, AND R. S. R. RUIZ<sup>1</sup>

<sup>1</sup>CAP, National Institute of Space Research, Av. dos Astronautas 1758, São José dos Campos 12227-010, Brazil

<sup>2</sup>DAS, National Institute of Space Research, Av. dos Astronautas 1758, São José dos Campos 12227-010, Brazil

<sup>3</sup>Institute for Astronomy, University of Hawaii, 2680 Woodlawn Dr., Honolulu, HI 96822, USA

<sup>4</sup>INAF-Osservatorio Astronomico di Capodimonte, via Moiariello 16, Napoli 80131, Italy

<sup>5</sup>LAC, National Institute of Space Research, Av. dos Astronautas 1758, São José dos Campos 12227-010, Brazil

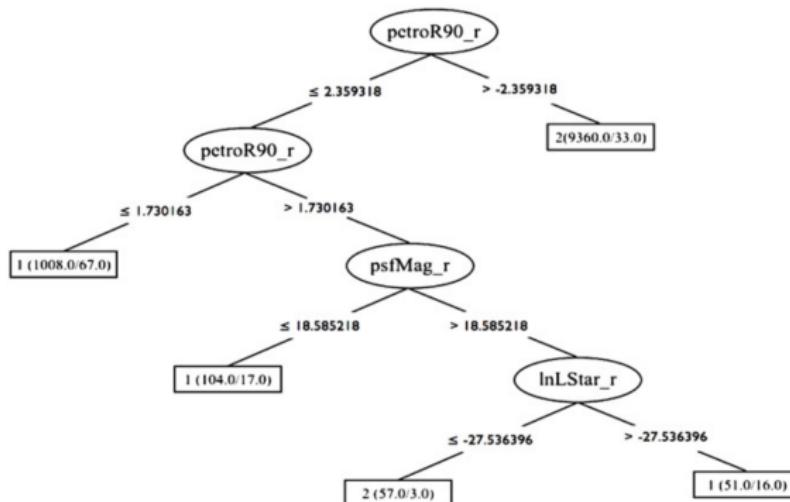
<sup>6</sup>IAG, University of São Paulo, Rua do Matão 1226, São Paulo 05508-090, Brazil

Received 2010 October 27; accepted 2011 February 11; published 2011 May 9

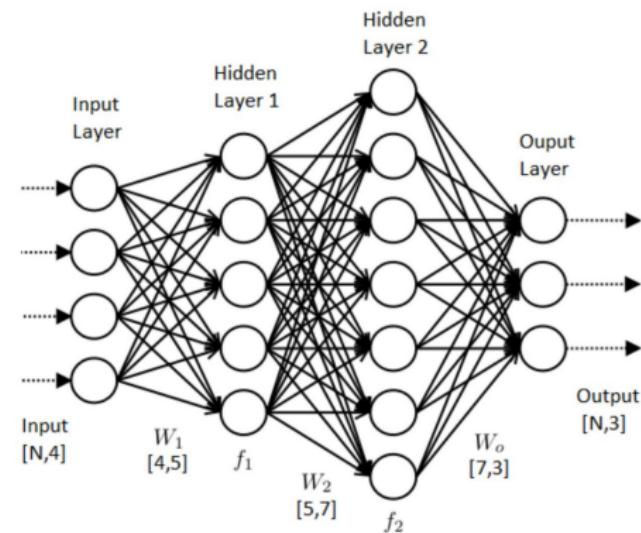
## ABSTRACT

We study the star/galaxy classification efficiency of 13 different decision tree algorithms applied to photometric objects in the Sloan Digital Sky Survey Data Release Seven (SDSS-DR7). Each algorithm is defined by a set of parameters which, when varied, produce different final classification trees. We extensively explore the parameter space of each algorithm, using the set of 884,126 SDSS objects with spectroscopic data as the training set. The efficiency of star-galaxy separation is measured using the completeness function. We find that the Functional Tree algorithm (FT) yields the best results as measured by the mean completeness in two magnitude intervals:  $14 \leq r \leq 21$  (85.2%) and  $r \geq 19$  (82.1%). We compare the performance of the tree generated with the optimal FT configuration to the classifications provided by the SDSS parametric classifier, 2DPHOT, and Ball et al. We find that our FT classifier is comparable to or better in completeness over the full magnitude range  $15 \leq r \leq 21$ , with much lower contamination than all but the Ball et al. classifier. At the faintest magnitudes ( $r > 19$ ), our classifier is the only one that maintains high completeness (>80%) while simultaneously achieving low contamination (~2.5%). We also examine the SDSS parametric classifier (psfMag – modelMag) to see if the dividing line between stars and galaxies can be adjusted to improve the classifier. We find that currently stars in close pairs are often misclassified as galaxies, and suggest a new cut to improve the classifier. Finally, we apply our FT classifier to separate stars from galaxies in the full set of 69,545,326 SDSS photometric objects in the magnitude range  $14 \leq r \leq 21$ .

# Using DT for classifying galaxies and stars in SDSS?



# How Neural Network works?



# Convolutional Neural Network and DLAs

## Deep Learning of Quasar Spectra to Discover and Characterize Damped Ly $\alpha$ Systems

David Parks,<sup>1</sup> J. Xavier Prochaska,<sup>2</sup> Shawfeng Dong,<sup>3</sup> Zheng Cai,<sup>2,4</sup>

<sup>1</sup>Computer Science, UC Santa Cruz, 1156 High St., Santa Cruz, CA 95064 USA — [dparks@ucsc.edu](mailto:dparks@ucsc.edu)

<sup>2</sup>Astronomy & Astrophysics, UC Santa Cruz, 1156 High St., Santa Cruz, CA 95064 USA

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

<sup>4</sup>Hubble Fellow

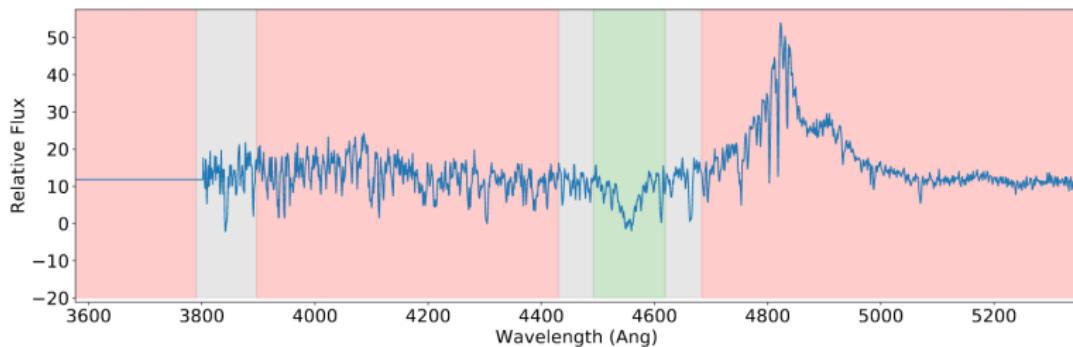
Accepted XXX. Received YYY; in original form ZZZ

### ABSTRACT

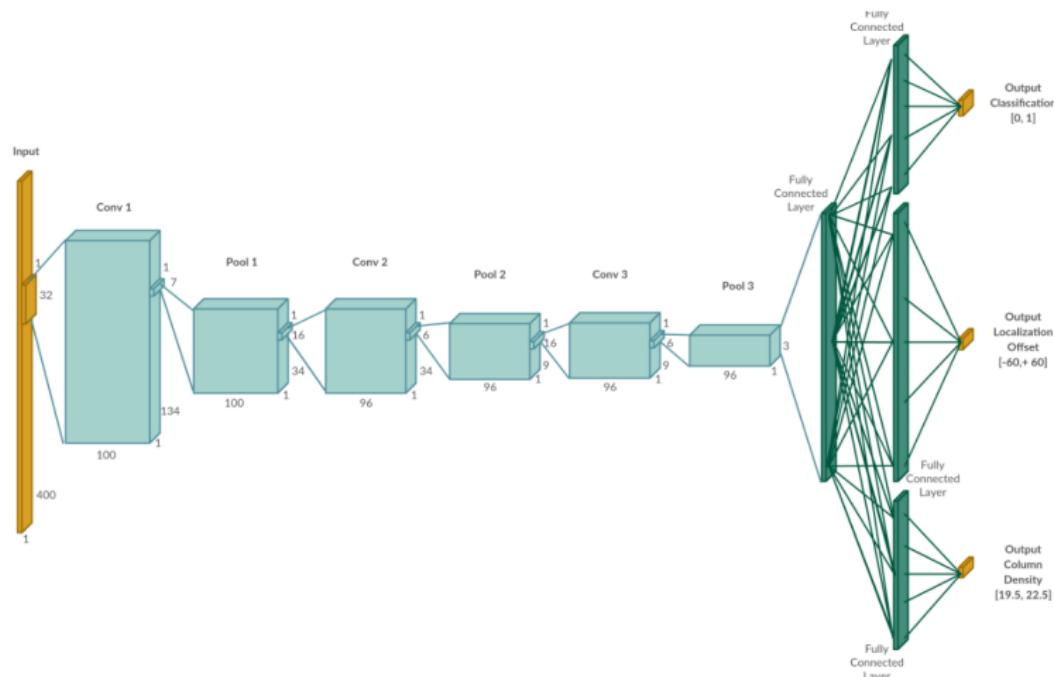
We have designed, developed, and applied a convolutional neural network (CNN) architecture using multi-task learning to search for and characterize strong HI Ly $\alpha$  absorption in quasar spectra. Without any explicit modeling of the quasar continuum nor application of the predicted line-profile for Ly $\alpha$  from quantum mechanics, our algorithm predicts the presence of strong HI absorption and estimates the corresponding redshift  $z_{\text{abs}}$  and HI column density  $N_{\text{HI}}$ , with emphasis on damped Ly $\alpha$  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  $z_{\text{abs}}$  and  $N_{\text{HI}}$ . 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  $N_{\text{HI}}$  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.

# Training DLAs/no DLAs

8 *Parks et al.*



# Treating DLA as a 1D image



Outline  
Big Data  
ML Algorithms  
**Supervised ML in astronomy**  
**Unsupervised ML in astronomy**  
ML and SKA  
ML limitations

**Clustering**  
outlier detection  
dimensionality reduction  
Neural network

Outline  
Big Data  
ML Algorithms  
**Supervised ML in astronomy**  
**Unsupervised ML in astronomy**  
ML and SKA  
ML limitations

Clustering  
**outlier detection**  
dimensionality reduction  
Neural network

Outline  
Big Data  
ML Algorithms  
Supervised ML in astronomy  
Unsupervised ML in astronomy  
ML and SKA  
ML limitations

Clustering  
outlier detection  
**dimensionality reduction**  
Neural network

Outline  
Big Data  
ML Algorithms  
Supervised ML in astronomy  
Unsupervised ML in astronomy  
ML and SKA  
ML limitations

Clustering  
outlier detection  
dimensionality reduction  
Neural network

text

Outline  
Big Data  
ML Algorithms  
Supervised ML in astronomy  
Unsupervised ML in astronomy  
**ML and SKA**  
ML limitations

- ML algorithms are not designed for astronomy but mostly for business.

- ML algorithms are not designed for astronomy but mostly for business.
- Training on a biased sample is very dangerous

- ML algorithms are not designed for astronomy but mostly for business.
- Training on a biased sample is very dangerous
- ML is the best for how only cares about results