Machine Learning on Galaxy Spectra

Graziano Ucci

Kapteyn Astronomical Institute, Groningen



rijksuniversiteit groningen

Suggested Readings



Python Machine Learning

[PACKT]



Hands-On Machine Learning with Scikit-Learn and TensorFlow (Concepts, Tools, and Techniques to Build Intelligent Systems) O'Reilly Media

Ian H. Witten, Eibe Frank, Mark A. Hall **Data Mining: Practical Machine Learning Tools and Techniques** Morgan Kaufmann Series

> Sebastian Raschka Python Machine Learning

> > Packt Publishing

Aurélien Géron

Wes McKinney **Python for Data Analysis** O'Reilly Media



The process of extracting useful insights from (raw) data.



The process of extracting useful insights from (raw) data.

Same as **Data Mining:** the process of discovering valuable information from (large) databases using algorithms able to find hidden patterns in data (Feldman & Dagan 1995).



The process of extracting useful insights from (raw) data.

Same as **Data Mining:** the process of discovering valuable information from (large) databases using algorithms able to find hidden patterns in data (Feldman & Dagan 1995).



OECD 2015 report: countries could be getting much more out of **Data Analytics** in terms of economic and social gains.



The process of extracting useful insights from (raw) data.

Same as **Data Mining:** the process of discovering valuable information from (large) databases using algorithms able to find hidden patterns in data (Feldman & Dagan 1995).





OECD 2015 report: countries could be getting much more out of **Data Analytics** in terms of economic and social gains.

We need a "Data-drive Innovation" where the usage of Data Analytics improves or foster new products, methods and markets.





Image credit: Springel et al. (2005)

- cosmological simulations;
- data from telescopes;



- cosmological simulations;
- data from telescopes;
- genome sequencing;



Source: James King-Holmes / Science Photo Library

- cosmological simulations;
- data from telescopes;
- genome sequencing;
- extreme phenomena in particle physics.



What "Big Data" really is?

The problem arose in the late 1990s within the META Group. For analysts it was becoming evident that their clients were increasingly encumbered by their data assets.

Big Data: a first definition (Gartner Group IT Glossary; Laney 2001)

"Big data is high-volume, -velocity and -variety information assets that demand cost-effective, innovative forms of information processing for enhanced insight and decision making"

Gartner 3 Vs (Laney 2001)



VOLUME: the amount of data that could be generated (**TB up to EB** of data to process, records, transactions, tables, dataframes)



VARIETY: the different types of data we could have in our data sets (structured, semi-structured, unstructured, mixed, multimedia)



VELOCITY: the speed at which new data is generated, collected and analyzed

IBM 4 Vs (2012)



VOLUME: the amount of data that could be generated (**TB up to EB** of data to process, records, transactions, tables, dataframes)



VARIETY: the different types of data we could have in our data sets (structured, semi-structured, unstructured, mixed, multimedia)



VELOCITY: the speed at which new data is generated, collected and analyzed



VERACITY: the messiness or trustworthiness of the data (missing values, authenticity, origin/reputation, uncertainties due to data inconsistency and incompleteness, ambiguities)

5 Vs (2015)



VOLUME: the amount of data that could be generated (**TB up to EB** of data to process, records, transactions, tables, dataframes)



VARIETY: the different types of data we could have in our data sets (structured, semi-structured, unstructured, mixed, multimedia)



VELOCITY: the speed at which new data is generated, collected and analyzed



VERACITY: the messiness or trustworthiness of the data (missing values, authenticity, origin/reputation, uncertainties due to data inconsistency and incompleteness, ambiguities)



VALUE: The ability to turn data into value/money, scientific/business models can be associated to the data (statistics/events, correlations, hypothesis)

The rapid growth in data is leading to a sort of "eScience": we now download the world trying to massively acquire in support of many hypotheses.

The rapid growth in data is leading to a sort of **"eScience"**: we now **download the world** trying to massively acquire in support of many hypotheses.

Science in most cases is driven by data more than by computation: because the cost of storing and acquisition has dropped precipitously, dealing with the scientific stored data requires several levels of analysis and various processing/algorithms to obtain complex models.

The rapid growth in data is leading to a sort of **"eScience"**: we now **download the world** trying to massively acquire in support of many hypotheses.

Science in most cases is driven by data more than by computation: because the cost of storing and acquisition has dropped precipitously, dealing with the scientific stored data requires several levels of analysis and various processing/algorithms to obtain complex models.

In short words, we will not face in science (especially in Astrophysics) an extreme computational complexity, we are already facing it.

Moore's Law has held for decades (the processing power roughly doubles every two years).



Moore's Law has held for decades (the processing power roughly doubles every two years).

However, processing speed is no longer the (only) problem: getting the data to the processor is the bottleneck.



Moore's Law has held for decades (the processing power roughly doubles every two years).

However, processing speed is no longer the (only) problem: getting the data to the processor is the bottleneck.

Example: at a disk transfer rate of 75 MB/sec, the time taken to transfer 100 GB of data is 22 mins (higher if servers have less than 100 GB of RAM).



Moore's Law has held for decades (the processing power roughly doubles every two years).

However, processing speed is no longer the (only) problem: getting the data to the processor is the bottleneck.

Example: at a disk transfer rate of 75 MB/sec, the time taken to transfer 100 GB of data is 22 mins (higher if servers have less than 100 GB of RAM).



This can be considered as the **end of the Moore's Law as we know it**. Increasing the performances cannot be achieved just through increasing hardware speed. We need a new approach: distributed computation must be exploited.

One of the fundamental families of algorithms used to extract information from data is represented by Machine Learning (ML). ML is a process able to map inputs to the output.

One of the fundamental families of algorithms used to extract information from data is represented by **Machine Learning (ML).** ML is a process able to map inputs to the output.

Statistics models the process that gave rise to data

ML tries to make an accurate prediction, given the data

One of the fundamental families of algorithms used to extract information from data is represented by **Machine Learning (ML).** ML is a process able to map inputs to the output.

Statistics models the process that gave rise to data

Traditional techniques rely on relatively small samples combined with heavy assumptions about data and its distributions. ML tries to make an accurate prediction, given the data

ML automatically discovers regularities in data using computational models that generalize the patterns found into new but similar data.

One of the fundamental families of algorithms used to extract information from data is represented by **Machine Learning (ML).** ML is a process able to map inputs to the output.

ML tends to make no pre-assumptions. The usual approach is, in most cases, empirical: the accuracy or even the applicability of the model is checked *a posteriori*.

"Machine Learning is the designing of computational models able to learn from data without being explicitly programmed"

as subfield of the Artificial Intelligence field, ML is mainly concerned with using computers for learning

"Machine Learning is the designing of **computational models** able to learn from data without being explicitly programmed"

as subfield of the Artificial Intelligence field, ML is mainly concerned with using computers for learning

"Machine Learning is the designing of **computational models** able to **learn from data without being explicitly programmed**"

with little to no human involvement, learning means any process whereby the system improves its performance based on experience



In the last decade ML has spread rapidly throughout **computer science** and beyond

web search, spam filters, recommender systems, ads placement, credit scoring, fraud detection, trading, drug design...

Credit Shantanu N. Zagade for Packt Publishing



In the last decade ML has spread rapidly throughout computer science and **beyond**

galaxy morphology classification, photometric redshift determination, cosmological simulations...

Adapted from Fig. 13 of Dieleman et al. (2015)



Case study: AlphaGo

October 2015: AlphaGo defeats Fan Hui ("It's not a human move. I've never seen a human play this move [...] So beautiful. So beautiful.")

March 2016: AlphaGo defeats Lee Sedol





Case study: Image recognition

ML automatically understands the content of images and associate relevant keywords. Deep Convolutional Neural Networks spot faces even if they are partially hidden/upside down.


Case study: GAME

"GAlaxy Machine learning for Emission lines": a code used for astrophysical applications able to infer the physical properties of galaxies from spectra (Ucci et al. 2017, 2018, 2019).

online at: game.sns.it

Case study: Amazon Go

Amazon Go is a store with no checkout required. An experience made possible by computer vision and sensor fusion, combined with deep learning technologies.

The "just walk out technology" automatically detects when products are taken from or returned to the shelves and keeps track of them in a virtual cart.



ML is useful when we do not have a simple and clear algorithmic and/or an analytic representation of our problem

ML is useful when we do not have a simple and clear algorithmic and/or an analytic representation of our problem



luminosity function representation: Schechter function

ML is useful when we do not have a simple and clear algorithmic and/or an analytic representation of our problem





luminosity function representation: Schechter function

handwritten digits representation: ???

Supervised Learning

ML algorithms can figure out how to perform tasks by generalizing from examples. The main goal in Supervised Learning is to learn a model from **labeled data** that allows us to make predictions about unseen or future data.





classification for predicting class labels

regression for predicting continuous outcomes

Supervised Learning











Decision Trees find final decision boundaries automatically based on the data.



PROs

valid for a large range of applications extremely fast and simple easy implementation

CONs

they are prone to **overfitting**

Ensemble methods

Many base learning algorithms (Trees) can be combined into an "ensemble learner" (Forest) which can achieve better results





Ensemble method: AdaBoost

AdaBoost method adds decision trees sequentially to generate a forest



In each iteration we apply the base learners on the training set with updated weights The **final model** is the weighted sum of the n learners.

Deep learning could be defined as a set of algorithms that were developed to train in the most efficient way artificial neural networks with many layers.

hidden layer input layer OUTPUT layer OUTPUT layer OUTPUT layer OUTPUT layer

Neural network

Adapted from: http://neuralnetworksanddeeplearning.com/chap5.html

"Simple" neural networks (i.e. with one hidden layer) can easily classify handwritten digits with an accuracy better than 98%.



Neural network

Adapted from: http://neuralnetworksanddeeplearning.com/chap5.html

Nonetheless, networks with more hidden layers can be even more powerful.



Adapted from: http://neuralnetworksanddeeplearning.com/chap5.html

Deep networks use the intermediate layers to **build up multiple levels of abstraction.**

For example, in visual pattern recognition, the first layer might learn to recognize edges, the neurons in the second layer could learn to recognize more complex shapes, the third layer would then recognize still more complex shapes, and so on.

Deep neural network



Adapted from: http://neuralnetworksanddeeplearning.com/chap5.html

A simple Deep Learning Net for image recognition



Adapted from: https://www.mathworks.com/



Images are matrices of pixel values, eventually one for each channel (i.e. RGB)





Image



Convolved Feature







The purpose of RELU is to introduce non-linearity (given that convolution is a linear operation, i.e. matrix multiplication)





Pooling (also called subsampling or downsampling) downsample, i.e. reduces the dimensionality of each feature map.



Putting all together with a fully connected neural net at the end.

A "not-so-simple" Deep Learning Net for Cosmological applications



Illustration: Nik Spencer; sources: NASA/WMAP Science Team; R. Ellis (Caltech)



Deep learning 21-cm images of the Cosmic Dawn (Gillet et al. 2019)



Most common libraries, programming languages



Most common libraries, programming languages



GAME "GAlaxy Machine learning for Emission lines"

Interstellar Medium (ISM)

Main (baryonic) ingredients of galaxies:

• stars


Interstellar Medium (ISM)

Main (baryonic) ingredients of galaxies:

- stars
- gas



Interstellar Medium (ISM)

Main (baryonic) ingredients of galaxies:

- stars
- gas
- dust



Interstellar Medium (ISM)

Main (baryonic) ingredients of galaxies:

- stars
- gas
- dust
- cosmic rays
- e.m. radiation
- magnetic field







density of an HII region: 100 cm⁻³

size: 10 pc ~ 3×10^{19} cm

density of an HII region: 100 cm⁻³ density of the air (sea level): 10¹⁹ cm⁻³



density of an HII region: 100 cm⁻³ density of the air (sea level): 10¹⁹ cm⁻³



density of an HII region: 100 cm⁻³ density of the air (sea level): 10¹⁹ cm⁻³



density of an HII region: 100 cm⁻³ density of the air (sea level): 10¹⁹ cm⁻³



We detect galaxies at high redshift, but we still do not have a clear understanding of their Interstellar Medium. The main aim is to study the composition and structure of the ISM

Physical properties (SFR, metallicity, etc.) of ISM in galaxies are usually inferred from **spectroscopic information**

Direct method (electron temperature)

(auroral to nebular line ratios, Pérez-Montero 2017)

 $R_{O3} = ([OIII] \lambda 4959 + [OIII] \lambda 5007) / [OIII] \lambda 4363$

 $R_{O2} = ([OII] \lambda 3726 + [OII] \lambda 3729) / ([OII] \lambda 7319 + [OII] \lambda 7330)$



Direct method (electron temperature)

(auroral to nebular line ratios, Pérez-Montero 2017)

 $R_{O3} = ([OIII] \lambda 4959 + [OIII] \lambda 5007) / [OIII] \lambda 4363$

 $\mathsf{R}_{O2} = ([OII] \ \lambda 3726 + [OII] \ \lambda 3729) \ / \ ([OII] \ \lambda 7319 + [OII] \ \lambda 7330)$

relations between the nebular-to-auroral line ratios and the electron temperature as a function of the electron density



Empirical calibrations

(Pagel et al. 1979, Vilchez & Esteban 1996, Pettini & Pagel 2004, Maiolino et al. 2008, Nagao et al. 2011, Marino et al. 2013, Curti et al. 2017)

 $\mathsf{R}_{23} = ([\mathsf{OII}] \ \lambda 3727 + [\mathsf{OIII}] \ \lambda 4959 + [\mathsf{OIII}] \ \lambda 5007) \ / \ \mathsf{H}\beta$

 $N_2 = [NII] \lambda 6583 / Ha$

R = ([OIII] λ 51.80 µm + [OIII] λ 88.33 µm) / [NIII] λ 57.21 µm







Comparison of theoretical spectra from a grid of photoionization models

(McGaugh 1991, Zaritsky et al. 1994, Kewley & Dopita 2002, Kobulnicky & Kewley 2004, Tremonti et al. 2004, Kewley & Ellison 2008, Dopita et al. 2016)

Numerical codes

(*IZI* Blanc et al. 2015, *pyqz* Dopita et al. 2013, *HII-CHI-mistry* Pérez-Montero 2014, *BOND* Vale Asari et al. 2016)

| 1月二下にいるでもロージ走り、手持のなきのり上の作用。 こまきことせきハマる日本に下正式することに思いるり、一下一手、各の子本を以り | |
|--|---|
| ● 10 中位 + 1 - 10 山市山 □ + 1 G N 基本书 O | |
| | 170K+7 |
| いまた、 日本では、日本では、日本では、日本では、日本では、日本では、日本では、日本では、 | |
| | |
| | の中国は行動に行っていり見れる |
| | - day official and an |
| こうしょう るいろう デント・ション 小学校 ひょうかい うまい 中心 ゆうしょう | 1 I OT 6 LA |
| シャンション にどうどう スリー ○中心シャキー キャメロ 小田・中心や日 アルロヨー かいき キャンドン にどうどう | 4日:日、6キャッシャー・キロのキャックのたました |
| | |
| ハー・シスローヨーチャレテ州市山市市の自宅に海路への時一で本心り舟り手れるとして、 ・ ・ トルトロンホスタッリカムまとあるためでしたした子りで、 き | \$ IN STREET |
| | |
| - そくじょそじ過ぎるたちでなっていてきたか。 いきサキキロい・ススになる。このはありい内でする | |
| | 1 ~ 2 0 ~ 2 |
| キャンロット ひつかけ オオ | |
| たいたまたただまたまであがらし、「「「「」」」」」、「」」」、「」」」、「」」」、「」」、「」」」、「」」 | なしてきました。このであるというなななない |
| | |
| | |
| こうにしいない いいいいい いいいいい いいいいい しょうしょう | |
| | |
| | |
| | |
| REPAIRS DE LE ALLE REPAIRS D'ALLE REPAIRS DE LE ALLE REPAIRS DE LE ALL | しょうをすいてきかいました。ここのとのないないできた。 |
| | 1 - 2 - 1 - 1 - 1 - 2 |
| 100-101の11月間を見たるという。2011-2011年11-2014年11-401-401-401-401-401-401-401-401-401-4 | |
| クレートないないないに、4月(東北部ングログルだんだん、小小小・ホテル、キャネル | |
| | |

GAME (**GA**laxy **Machine learning for Emission lines**, Ucci et al. 2017, 2018) is a new, fast code able to reconstruct the physical properties of the ISM in (distant) galaxies





GAME (GAlaxy Machine learning for Emission lines, Ucci et al. 2017, 2018) is a new, fast code able to reconstruct the physical properties of the ISM in (distant) galaxies by using all the available information encrypted in spectra





Inspiration...

Inspiration... Shazam!



you are in a shop and you like the music you are hearing: **tap the button** it there is a **match** you are then given the **metadata of the audio** (title, lyrics, video, artist...)

shazam analyzes the captured sound by creating the fingerprint of the audio, **it starts the search for matches in the database** (Wang 2003)

A "Shazam for galaxies"







you have an observational/synthetic **spectrum of a galaxy**

you are then given the **"metadata" of the galaxy** (gas volume density, ionization state, metallicity...)

the algorithm works by **analyzing the** emission line intensities, it works on the models in the database

A "Shazam for galaxies"







you have an observational/synthetic **spectrum of a galaxy**

you are then given the **"metadata" of the galaxy** (gas volume density, ionization state, metallicity...)

the algorithm works by **analyzing the** emission line intensities, it works on the models in the database



Method (Ucci et al. 2017, 2018)





grid of models: metallicity, column density, ionization parameter, density

| Parameter | minimum | maximum |
|------------------------------|---------|---------|
| $\log(Z/Z_{\odot})$ | -3.0 | 0.5 |
| $\log(n/\mathrm{cm}^{-3})$ | -3.0 | 5.0 |
| $\log(U)$ | -4.0 | 3.0 |
| $\log(N_H/\mathrm{cm}^{-2})$ | 17.0 | 23.0 |



- metallicity (Z / Z_o)
- ionization state (ionization parameter U or FUV flux in the Habing band 6 - 13.6 eV)

$$U = \frac{1}{4\pi R_s^2 nc} \int_{\nu_e}^{\infty} \frac{L_{\nu}}{h\nu} d\nu = \frac{F(H)}{4\pi R_s^2 nc}$$

- density (n / cm⁻³)
- column density (N_H / cm⁻²)

$$N_H = \int_0^R n ds$$



The Rosette nebula (NGC 2237) (Credit WIYN and NOAO/AURA/NSF)

| arid of | arid of modoles motallicity, column donsity | | minimum | maximum | | |
|--|---|------------------------------|---------|---------|--|--|
| grid of | gna of models: metallicity, column density, | | -3.0 | 0.5 | | |
| | ionization parameter, density | $\log(n/\mathrm{cm}^{-3})$ | -3.0 | 5.0 | | |
| | | $\log(U)$ | -4.0 | 3.0 | | |
| , | | $\log(N_H/\mathrm{cm}^{-2})$ | 17.0 | 23.0 | | |
| | | | | | | |
| CLOUDY v13.03 | | | | | | |
| | | | | | | |
| (2) large library of emission line intensities (50,000 models) | | | | | | |









| arid of modeles motallicity, column density | Parameter | minimum | maximum |
|--|------------------------------|--------------|------------|
| grid of models. metallicity, column density, | $\log(Z/Z_{\odot})$ | -3.0 | 0.5 |
| ionization parameter, density | $log(n/cm^{-3})$ log(U) | -3.0 -4.0 | 5.0 3.0 |
| | $\log(N_H/\mathrm{cm}^{-2})$ | 17.0 | 23.0 |
| | | | |

large library of emission line intensities (50,000 models!)

CLOUDY v13.03

GAME

Supervised Machine Learning (SML): physical properties are inferred from an input spectrum using emission line intensities as input features



library of emission line intensities





Integral Field Spectroscopy (IFS)



pixel: a data point on a CCDspaxel: a spectrum in a data cubevoxel: a data point in a data cube

combines imaging and spectroscopy: "3D spectroscopy"

output: data cube



Results: application to IFU observations (Ucci et al. 2018, 2019)



MUSE observations of Henize 2-10 (Cresci et al. 2017)

prototype of starburst HII galaxies (Allen et al. 1976)

$$D = 8.23 \text{ Mpc} (1 \text{ arcsec} = 40 \text{ pc})$$

$$M_{\star} \sim 3.7 \text{ x } 10^9 \text{ M}_{\odot}$$

SFR ~ 1.9
$$M_{\odot}$$
 yr⁻¹ (Reines et al. 2011)

Results: application to IFU observations (Ucci et al. 2018, 2019)





Results: radially-averaged profiles (Ucci et al. 2018, 2019)



GAME Predictive performances (k-fold cross validation)



the library is split into k folds

the code **trains the ML algorithm on k-1** and then tests on the left-out

GAME computes the score k consecutive times, and gives back the mean of these scores for each physical property

GAME Predictive performances (results on 1 fold)





Hands-on session
> cd /media/...

> cd astrophysics_classes/

> cd ucci_machine_learning_on_galaxy_spectra/

> module load python/3.7.2

> jupyter-notebook &