Fifth Workshop on Millimetre Astronomy in Italy **Bologna, June 12-14 2023** 0000100 **UMLAUT: Data-Driven Machine Learning Approach with Automated**¹ **Parameter Tuning and Outlier**¹ Detection Ivano Baronchelli **INAF - IRA** ARC – Italian node

THE ASTROPHYSICAL JOURNAL SUPPLEMENT SERIES, 249:12 (27pp), 2020 July © 2020. The American Astronomical Society. All rights reserved.

https://doi.org/10.3847/1538-4365/ab9a3a



Identification of Single Spectral Lines through Supervised Machine Learning in a Large HST Survey (WISP): A Pilot Study for Euclid and WFIRST

I. Baronchelli¹⁽⁰⁾, C. M. Scarlata²⁽⁰⁾, G. Rodighiero¹⁽⁰⁾, L. Rodríguez-Muñoz¹⁽⁰⁾, M. Bonato^{3,4}⁽⁰⁾, M. Bagley⁵⁽⁰⁾, A. Henry⁶⁽⁰⁾, M. Rafelski^{6,7}⁽⁰⁾, M. Malkan⁸⁽⁰⁾, J. Colbert⁹, Y. S. Dai (戴昱)¹⁰⁽⁰⁾, H. Dickinson^{2,11}⁽⁰⁾, C. Mancini¹⁽⁰⁾, V. Mehta²⁽⁰⁾, L. Morselli¹, and H. I. Teplitz⁹⁽⁰⁾

¹ Dipartimento di Fisica e Astronomia, Università di Padova, vicolo Osservatorio, 3, I-35122 Padova, Italy; ivano.baronchelli@unipd.it ² MN Institute for Astrophysics, University of Minnesota, 116 Church Street SE, Minneapolis, MN 55455, USA ³ INAF-Istituto di Radioastronomia and Italian ALMA Regional Centre, Via Gobetti 101, I-40129, Bologna, Italy ⁴ INAF-Osservatorio Astronomico di Padova, Vicolo dell'Osservatorio 5, I-35122, Padova, Italy ⁵ College of Natural Sciences, The University of Texas at Austin, 2515 Speedway, Austin, TX 78712, USA ⁶ Space Telescope Science Institute, 3700 San Martin Drive, Baltimore, MD 21218, USA ⁷ Department of Physics and Astronomy, Johns Hopkins University, Baltimore, MD 21218, USA ⁸ Department of Physics and Astronomy, UCLA, Physics and Astronomy Building, 3-714, LA CA 90095-1547, USA ⁹ IPAC, Mail Code 314-6, Caltech, 1200 E. California Boulevard, Pasadena, CA 91125, USA ¹⁰ Chinese Academy of Sciences South America Center for Astronomy (CASSACA)/NAOC, 20A Datun Road, Beijing 100101, People's Republic of China ¹¹ School of Physical Sciences, The Open University, Walton Hall, Milton Keynes, MK7 6AA, UK *Received 2019 November 28; revised 2020 May 25; accepted 2020 June 5; published 2020 July 13*

Abstract

Future surveys focusing on understanding the nature of dark energy (e.g., Euclid and WFIRST) will cover large fractions of the extragalactic sky in near-IR slitless spectroscopy. These surveys will detect a large number of galaxies that will have only one emission line in the covered spectral range. In order to maximize the scientific return of these missions, it is imperative that single emission lines are correctly identified. Using a supervised machine-learning approach, we classified a sample of single emission lines extracted from the WFC3 IR Spectroscopic Parallel survey, one of the closest existing analogs to future slitless surveys. Our automatic software integrates a spectral energy distribution (SED)-fitting strategy with additional independent sources of information. We calibrated it and tested it on a "gold" sample of securely identified objects with multiple lines detected. The algorithm correctly classifies real emission lines with an accuracy of 82.6%, whereas the accuracy of the SED-fitting technique alone is low (~50%) due to the limited amount of photometric data available (≤ 6 bands). While not specifically designed for the Euclid and WFIRST surveys, the algorithm represents an important precursor of similar algorithms to be used in these future missions.

THE ASTROPHYSICAL JOURNAL SUPPLEMENT SERIES, 249:12 (27pp), 2020 July © 2020. The American Astronomical Society. All rights reserved.

https://doi.org/10.3847/1538-4365/ab9a3a



Identification of Single Spectral Lines through Supervised Machine Learning in a Large HST Survey (WISP): A Pilot Study for Euclid and WFIRST



unsupervised approach with a supervised approach we previously developed.

THE ASTROPHYSICAL JOURNAL SUPPLEMENT SERIES, 249:12 (27pd), 2020 July © 2020. The American Astronomical Society. All rights reserved.

https://doi.org/10.3847/1538-4365/ab9a3a



Identification of Single Spectral Lines through Supervised Machine Learning in a Large HST Survey (WISP): A Pilot Study for Euclid and WFIRST

correctly identifies real lines in 83.

unsupervised approach with a supe



Autonomous University of Barcelona. Spain

Reviewed by: Kleanthes K. Grohmann, far. More specifically, and for the first time, we discovered that there is a micro-variation resulting from diastratic (age) variation within the Roveda variety, which represents the THE ASTROPHYSICAL JOURNAL SUPPLEMENT SERIES, 249:12 (27pp), 2020 July © 2020. The American Astronomical Society. All rights reserved.

https://doi.org/10.3847/1538-4365/ab9a3a



Identification of Single Spectral Lines through Supervised Machine Learning in a Large HST Survey (WISP): A Pilot Study for Euclid and WFIRST



sky in near-IR slitless spectroscopy correctly identifies real lines in 83. unsupervised approach with a supe

Autonomous University of Barcelona, Spain **Reviewed by:** Kleanthes K, Grohmann,

light on a series of phenomena that have been neglected or poorly understood thus far. More specifically, and for the first time, we discovered that there is a micro-variation resulting from diastratic (age) variation within the Roveda variety, which represents the THE ASTROPHYSICAL JOURNAL SUPPLEMENT SERIES, 249:12 (27pp), 2020 July © 2020. The American Astronomical Society. All rights reserved.

https://doi.org/10.3847/1538-4365/ab9a3a



Identification of Single Spectral Lines through Supervised Machine Learning in a Large HST Survey (WISP): A Pilot Study for Euclid and WFIRST

| I. Baronchelli ¹ ⁰ , C. M. Scarlata ² ⁰ , G. Rodighiero | ¹ ⁽¹⁾ , L. Rodríguez-Muñoz ¹ ⁽¹⁾ , M. Bonato ^{3,4} ⁽¹⁾ | , M. Bagley ⁵ [®] , A. Henry ⁶ [®] , | | |
|--|---|--|---|---|
| M. Rafelski ^{0,7} ⁶⁰ , N | | | | (6) |
| ZECCCO Search | Q Upload | Communities | | → Log in 🕼 Sign up |
| Product -> Solutions -> Open Sou | rce 🗸 Pricing | GitHub | Search | Sign in Sign up |
| Here Bereneballi | Overview Repositories Popular repositories BESTFIT SED fitting code for stellar templates • sed python_library python functions and classes • Python | 4 Projects Package | es ☆ Stars IDL_library Iibrary of IDL functions IDL UMLAUT Unsupervised Machine Learning Algoritt IDL | Public Public hm based on Unbiased Topology |
| Ibaronch | | | | |
| > UMLAUT estimates the value of the (N+1)-th closest data-points of the training set, in a ranke After finding the closest data points, the unknow values assumed by the closest reference data p | parameter for the analysis data point. To d N-dimensional space "associated" with vn parameter is obtained as the combina points along the (N+1)-th dimension. | o this purpose, UMLAUT finds the n the physical parameter space. Ition (ex. from the average) of the | | |
| sky in near-IR slitless spectroscopy correctly identifies real lines in 83. unsupervised approach with a supe | | Autonomous University of Barcelona, Spain Reviewed by: Kleanthes K. Grohmann. | n a series of phenomena that have been re specifically, and for the first time, we a ing from diastratic (age) variation within the | en neglected or poorly understood thus discovered that there is a micro-variation he Roveda variety, which represents the |

Kleanthes K. Grohmann,

Identifying emission lines is quite easy... ... when multiple lines are detected



WISP data (Atek+2010)











WSP PI: M. A. Malkan (University of California - Los Angeles) The WFC3 Infrared Spectroscopic Parallel survey H. Atek et al. (2010)

HST pure parallel large program (~2000 orbits, >1500 arcmin²)

- slitless spectroscopy $\begin{cases} G102: 0.8 1.1 \ \mu m \ (R~210) \\ G141: 1.07 1.7 \ \mu m \ (R~130) \end{cases}$
- direct imaging {
 J band (F110W), H band (F140W/F160W) Ancillary UVIS + IRAC + ground based optical
- One of the most importan proxies For the future surveys of <u>Euclid</u> and <u>Roman</u>
- Now collecting data with <u>JWST</u> (PASSAGE survey, P.I. M. Malkan, proposal ID: GO-1571 - Cycle 1)



An artificial intelligence is not an artificial GOD 0101 $\cap \cap$ \cap $1 \cap 1$ 00100100 010001010 1001 A learner can help exploiting hidden information, but it does not create information from nothing.

Additional information that can help identifying a spectral line: - Apparent Magnitude (F110W); - Apparent Size (F110W); - Color index (J-H); - Line equivalent width; ... etc ...



UMLAUT basics Unbiased Machine Learning Algoirthm Using Topology

Empirical principle:

Objects that are close to each other in N dimensions tend to be close also in an indipendent N+1 dimension, when N>>

UMLAUT basics Unbiased Machine Learning Algoirthm Using Topology

Empirical principle:

Objects that are close to each other in N dimensions tend to be close also in an indipendent N+1 dimension, when N>>



Toucans that are similar to each other in 2 dimensions, tend to be similar also in 3 dimensions...

UMLAUT basics Unbiased Machine Learning Algoirthm Using Topology

Empirical principle:

Objects that are close to each other in N dimensions tend to be close also in an indipendent N+1 dimension, when N>>





Toucans that are similar to each other in 2 dimensions, tend to be similar also in 3 dimensions...

...unless one of the pictures is depicting the picture of a toucan.











Analysis data point • • Х















Among the UMLAUT Advantages:

Single dimensions could retain useful information about the output without showing any correlations at all with the output itself. $\frac{z}{2}$



 \rightarrow UMLAUT takes into account all the input dimensions at the same time!

Among the UMLAUT Advantages:

Dimensionality reduction strategies (such as PCA) often assume orthogonal native axes and/or Gaussian distributions



→ UMLAUT does not require/assume orthogonal native axes or assumptions on the data distribution.

Dimensionality reduction strategy

N-dimensional space of expansion of the N input parameters (N=2 in this example)



Strategy (simplified):

1) Measure the dispersion on the output parameter in positions 1, A and B

2) Move into the direction of decreasing gradient (position 2)

3) Iterate until a local minimum is found



| Naïve Bayes alone | | | | |
|----------------------|---------------|---------------|---------------|--|
| | Completeness | Contamination | Purity | |
| Ha | 88.3 % | 15.2 % | 84.8 % | |
| [OIII] | 77.0 % | 17.9 % | 82.1 % | |
| [011] | 50.0 % | 73.3 % | 26.7 % | |
| UMLAUT alone | | | | |
| | Completeness | Contamination | Purity | |
| Ha | 89.6 % | 16.0 % | 84.0 % | |
| [0111] | 77.8 % | 18.0 % | 82.0 % | |
| [OII] | 12.5 % | 33.3 % | 66.7 % | |
| Naïve Bayes + UMLAUT | | | | |
| | Completeness | Contamination | Purity | |
| Ha | 89.3 % | 12.8 % | 87.2 % | |
| [0111] | 79.6 % | 16.2 % | 83.8 % | |
| [011] | 66.7 % | 70.9 % | 29.1 % | |



Image of a band 3 calibrator

- Are there any other serendipitous detections around?

Compute number counts !



Image of a band 3 calibrator



~ 500 valid calibrator images ~ 700 arcmin² (at 1mJy)

- Are there any other serendipitous detections around?

Compute number counts !









Image of a band 3 calibrator

- Are there any other serendipitous detections around?

Compute number counts !



Image of a band 3 calibrator



Same calibrator (after changing scale)





Same calibrator (after changing scale)





Image of a band 3 calibrator



Same calibrator (after changing scale)



Another calibrator



Same calibrator (after changing scale)



Another calibrator

Same calibrator (after changing scale)





Another calibrator

Same calibrator (after changing scale)

Extensions of the central calibrator are counted as (multiple) serendipitous sources \rightarrow Down to a few sigma, it becomes more and more difficult to identify these contaminants by eye



ALMA Band 3 number counts - simulations



Original image

Simulated image

ALMA Band 3 number counts - simulations



Original image

Simulated image

ALMA Band 3 number counts (5σ)



ALMA Band 3 number counts (5σ)















UMLAUT

Easily adaptable to a wide variety of problems



ORIGINAL RESEARCH published: 11 July 2019 doi: 10.3389/fpsyg.2019.01528



Inter- vs. Intra-Speaker Variation in Mixed Heritage Syntax: A Statistical Analysis

Federica Cognola^{1*}, Ivano Baronchelli^{2*} and Evelina Molinari³

¹ Dipartimento di Studi Europei, Americani e Interculturali, Sapienza Università di Roma, Rome, Italy, ² Dipartimento di Fisica e Astronomia, Università di Padova, Padua, Italy, ³ Istituto Culturale Mòcheno, Palù del Fersina, Italy

Based on the novel data pertaining to five syntactic phenomena (the position of the finite verb in embedded clauses, in sentences with a modal verb, negative concord, the position of focused light/heavy objects in main clauses with a complex tense and scrambling) in the heritage language Mocheno collected via original fieldwork, we show that there are two populations – one exhibiting intra-speaker variation between German and Italian word orders, and one lacking it; and these two populations are the result of diatopic variation and, to a lesser extent, of diastratic variation. The results achieved using quantitative statistical analysis are partially convergent with those arrived at via the traditional theoretical syntax for Mocheno, but our analysis has allowed us to shed new light on a series of phenomena that have been neglected or poorly understood thus far. More specifically, and for the first time, we discovered that there is a micro-variation resulting from diastratic (age) variation within the Roveda variety, which represents the

OPEN ACCESS

Edited by:

Ángel J. Gallego, Autonomous University of Barcelona, Spain

> Reviewed by: Kleanthes K. Grohmann,



Demographics







1811 – Napoleon wants a census of all the dialects spoken in the Italian kindom

→ First important official document written in Mocheno language

| Control | and again and a |
|--|--|
| S Juce Franchite | New right and the |
| East who answered is | 2 and by Benilflain |
| dian is public is entered | A. A. where and he Free whe |
| ter V en normalist Plantes & South | terihi Saran Maran biya Yan pin batu ya yi ngi pit an |
| rates t aria ka geoneticas majat mandade con ella | a water a water handed |
| s sent at the peak to when | 3.65 Sec |
| . 200 wichs bear go leht | 9 to a main side milling |
| esten ene le se pert dest ann est le mais diaets annitél arai | an parte à remark come alle |
| dont buck of the goo pe | the flat on the big time the product |
| nal, dana manané nan- A-Mana samuné nan mejanék | . r. no of a mar light |
| is herean sur pailer. I straine Same court and | 6 and and apple |
| levens, dens ille Cogelaturie | na para para na Sala ing |
| entral. | phase generation and |
| and some an polithe partet. | and the second s |
| na specific gran began reconspirit mena with main and characterizations | an he taise daga a sa ar an ar anna an anna fara da an falainn |
| . Het geveniker helen derdene | V and his in the gran of |
| an words been to world at | a le fifiai - git gaage vergen. In geher - he geher vertagerig |
| nam de second delegendor dans ar. And t | a ganger and a series |
| | 1 |
| | |

Who wrote it?

1811 – Napoleon wants a census of all the dialects spoken in the Italian kindom

→ First important official document written in Mocheno language





Ongoing development:

- Code optimization (same algorithm faster run);
- Translate UMLAUT in different languages (Python)
- Parallelization of the code (NVIDIA-CUDA);
- Improvement of the current dimensionality reduction strategy
- <u>example1</u>: using a PCA-like strategy
- <u>example2</u>: by replacing each of the input dimensions with random ranks and comparing the output variance.

Conclusions:

- New facilities producing unprecedented amount of data are going to be available in the next few years: <u>every</u> <u>aspect of the treatment of such data must be automated;</u>

UMLAUT is:

- a ML algorithm based on the analogy paradigm;
- extremely (and easily) adaptable to very different contexts;
- freely available online

3 applications were presented:

- Identification of single spectral lines
- Recognition of ancient linguistic variants
- Identification of contaminants in sub-mm number counts