Pre-seminar What is ML?

Geneva University, Jan, 2020

- What is Machine learning?
- What is Dimensionality reduction?
- What is Clustering Analysis?
- What am I doing with ML?

What is Machine learning?



y = ax + b









Classification vs. regression

Classification







Supervised vs Unsupervised



What is Dimensionality reduction?







What is Clustering Analysis?



What am I doing with ML?



RFI mitigation



Inpainting CMB







masked









Inpainting CMB



See Co

Inpainting CMB



Medical physics



Original image







Registered image



Mining cosmic datasets +some cool DS stuff

Alireza Vafaei Sadr IPM, Tehran





Outline

Cosmology and BIG data
A Quick Review of Applications
Anomaly detection
DRAMA
Future directions



Cosmology/Astrophysics



E. Siegel, with images derived from ESA/Planck and the DoE/NASA/ NSF interagency task force on CMB research. From his book, Beyond The Galaxy.

Cosmology and Big data



Big Universe, Big Data: Machine Learning and Image Analysis for Astronomy

Jan Kremer, Kristoffer Stensbo-Smidt, Fabian Gieseke, Kim Steenstrup Pedersen, and Christian Igel









xmm-newton

It is getting hotter!



Number of physics submitted manuscripts that include "machine learning" in their abstracts.

A quick review on what people have done



Classification





http://cdn.spacetelescope.org/archives/images/screen/heic99020.jpg

Galaxy zoo challenge



Figure 1. Flowchart of the classification tasks for GZ2, beginning at the top centre. Tasks are colour-coded by their relative depths in the decision tree. Tasks outlined in brown are asked of every galaxy. Tasks outlined in green, blue, and purple are (respectively) one, two or three steps below branching points in the decision tree. Table describes the responses that correspond to the icons in this diagram.

https://www.galaxyzoo.org/

Classifying the Large Scale Structure of the Universe with Deep Neural Networks

M.A. Aragon-Calvo¹ * ¹Instituto de Astronomía, UNAM, Apdo. Postal 106, Ensenada 22800, B.C., México



Detection





RYAN HAUSES1 AND BRANT E. ROBERTSON2.1

A. Vafaei Sadr, ^{1,2,3,4} Etienne, E. Vos, ^{2,4,5}† Bruce A. Bausett, ^{2,15,6}‡ Zaffirah Hosenie,^{25,3} N. Osseer,²³ and Michelle Lochner²³

Data cleansing



Radio frequency interference mitigation using deep convolutional neural networks

Joël Akeret^{a,*}, Chihway Chang^a, Aurelien Lucchi^b, Alexandre Refregier^a



Cleaning our own Dust: Simulating and Separating Galactic Dust Foregrounds with Neural Networks K. Aylor,¹ M. Haq,² L. KNOX,¹ Y. HEZAVEH,^{3,4} AND L. PERREAULT-LEVASSEUR^{3,5,4} Gravitational Wave Denoising of Binary Black Hole Mergers with Deep Learning

Wei Wei^{a,b}, E. A. Huerta^{a,c}

Denoising Weak Lensing Mass Maps with Deep Learning

Masato Shirasaki,¹ Naoki Yoshida,^{2,3,4} and Shiro Ikeda^{5,6}

DENOISING GRAVITATIONAL WAVES WITH ENHANCED DEEP RECURRENT DENOISING AUTO-ENCODERS

Hongyu Shen¹ Daniel George² Eliu. A. Huerta^{2,3} Zhizhen Zhao^{1,3}

Separating the EoR signal with a convolutional denoising autoencoder: a deep-learning-based method

Weitian Li,^{1*} Haiguang Xu,^{12*} Zhixian Ma,³ Ruimin Zhu,⁴ Dan Hu,¹ Zhenghao Zhu,¹ Junhua Gu,⁵ Chenxi Shan,¹ Jie Zhu³ and Xiang-Ping Wu⁵

Solar image denoising with convolutional neural networks

C. J. Díaz Baso¹, J. de la Cruz Rodríguez¹, and S. Danilovic¹

Simulation



CAN: Creative Adversarial Networks Generating "Art" by Learning About Styles and Deviating from Style Norms*

Ahmed Elgammal¹¹ Bingchen Liu¹ Mohamed Elhoseiny² Marian Mazzone³


PROGRESSIVE GROWING OF GANS FOR IMPROVED QUALITY, STABILITY, AND VARIATION

 Tero Karras
 Timo Aila
 Samuli Laine
 Jaakko Lehtinen

 NVIDIA
 NVIDIA
 NVIDIA
 NVIDIA
 NVIDIA

https://towardsdatascience.com/do-gans-really-model-the-true-data-distribution-or-are-they-just-cleverly-fooling-us-d08df69f25eb

Cosmological N-body simulations: a challenge for scalable generative models

Nathanaël Perraudin^{1*}, Ankit Srivastava¹, Aurelien Lucchi², Tomasz Kacprzak³, Thomas Hofmann² and Alexandre Réfrégier³





Mass histogram

Peak histogram



Power spectral density

Real

Faile.

120



From Dark Matter to Galaxies with Convolutional Neural Networks

Jacky H. T. Yip' Xinyue Zhang, Yanfang Wang, Wei Zhang, Yueqiu Sun* Gabriella Contardo, Francisco Villaescusa-Navarro, Siyu He, Shy Genel, Shirley Ho

Deep learning dark matter map reconstructions from DES SV weak lensing data

Niall Jeffrey,¹ • François Lanusse², Ofer Lahav¹, Jean-Luc Starck³

Learning to Predict the Cosmological Structure Formation

Siyu Hender, Yin Line, Yu Fengin, Shirley House Andreas, Slamak Ravanbakhshi, Wel Chen, and Barnabás Póczos

CMB-GAN: Fast Simulations of Cosmic Microwave background Anisotropy maps using Deep Learning

Amit Mishra", Pranath Reddy", Rahul Nigam"



SC-FEGAN: Face Editing Generative Adversarial Network with User's Sketch and Color

Youngjoo Jo Jongyoul Park





https://arxiv.org/pdf/1905.08233v1.pdf

Anomaly detection

The light from quasar pairs reach Earth, although some were absorbed by the gas in the cosmic web, Springel et al. (2005) (cosmic web) / J. Neidel, MPIA











10 FIND THE ODD ONE IN 60 SECONDS





Figure 1 A plot of recent major astronomical discoveries, taken from (Ekers 2009), of which seven were "known-unknowns" (i.e. discoveries made by testing a prediction) and ten were "unknownunknowns" (ie. a serendipitous result found by chance while performing an experiment with different goals). The data in this plot are taken from Wilkinson et al. (2004).

Norris, R. P. (2017). Discovering the unexpected in astronomical survey data. Publications of the Astronomical Society of Australia, 34.

Table 1 Major discoveries made by the Hubble Space Telescope (*HST*). Of the *HST*'s "top ten" discoveries (as ranked by National Geographic magazine), only one was a key project used in the *HST* funding proposal (Lallo 2012). A further four projects were planned in advance by individual scientists but not listed as key projects in the *HST* proposal. Half the "top ten" *HST* discoveries were unplanned, including two of the three most cited discoveries, and including the only *HST* discovery (Dark Energy) to win a Nobel prize. This Table was previously published by Norris et al. (2015).

| Project | Key | Planned? | Nat Geo | Highly | Nobel |
|---|----------|----------------------|--------------|--------------|--------------|
| 0.074 | Project? | | top ten? | cited? | Prize? |
| Use cepheids to improve value of H_0 | 1 | | 1 | 1 | |
| UV spectroscopy of ig medium | 1 | ~ | | | |
| Medium-deep survey | 1 | 1 | | | |
| Image quasar host galaxies | | 1 | 1 | | |
| Measure SMBH masses | | 1 | 1 | | |
| Exoplanet atmospheres | | 1 | 1 | | |
| Planetary Nebulae | | 1 | 1 | | |
| Discover Dark Energy | | | 1 | 1 | \checkmark |
| Comet Shoemaker-Levy | | | 1 | | |
| Deep fields (HDF, HDFS, GOODS, FF, etc) | | | 1 | \checkmark | |
| Proplyds in Orion | | | 1 | | |
| GRB Hosts | | | \checkmark | | 1 |

Neural network-based anomaly detection for high-resolution X-ray spectroscopy

Y. Ichmohe 12. and S. Yamada,³

Search for unusual objects in the WISE Survey

Aleksandra Solarz¹, Mariej Itilirki^{2 † 8} and Agniezzka Pollo^{1 4}

SELF-SUPERVISED ANOMALY DETECTION FOR NARROWBAND SETI

Youfan Gerry Zhang[†], Ki Hyun Woo^{*}, Seung Wool Son[†], Andrew Stemton^{1,3,1,4}, Steve Cruft[†]

Anomaly detection for machine learning redshifts applied to SDSS galaxies

Ben Hoyle^{1,2}, Markus Michael Rau^{1,4}, Kerstin Paech^{1,2}, Christopher Bonnett^a Stella Seitz^{1,4}, Jochen Weller^{1,2,4}

Active Anomaly Detection for time-domain discoveries

E. E. O. Ishida¹^{*}, M. V. Korniko²³^{*}, K. L. Malanchev²³, M. V. Pruzhinskaya², A. A. Volnova⁴, V. S. Korolev⁵⁶, F. Mondon¹, S. Sreejith¹, A. Malancheva² and S. Das⁸

Anomaly Detection in the Open Supernova Catalog

M. V. Pruzhinskaya,¹* K. L. Malanchev,^{1,2}† M. V. Kornilov,^{1,2} E. E. O. Ishida,³ F. Mondon,³ A. A. Volnova⁴ and V. S. Korolev^{5,6}

Current projects in SKA (MeerKAT) DS team:

- Source detection (Mightee, SKAch-I)
- Anomaly detection WTF, PLAsTiCC
- RFI simulation/mitigation
- Extended source simulation
- Observation strategy



A Flexible Framework for Anomaly Detection via Dimensionality Reduction

A. Vafaei Sadr, B. Bassett, M. Kunz ISCMI-2019

No Free Lunch Theorems

- Any two optimization algorithms are equivalent when their performance is averaged across all possible problems
- No anomaly detection algorithm works for all anomalies
- No anomaly algorithm is "best" on average.
- Different algorithms work for different anomalies
- So lets consider families of anomaly algorithms

Dimensionality Reduction Anomaly Meta Algorithm



https://github.com/vafaei-ar/drama

Dimensionality reduction



Clustering



As an example: MNIST





DRAMA (Dimensionality Reduction Anomaly Meta-Algorithm):





Dimensionality Reduction Technique

- Autoencoder
- Variational autoencoder
- principal component analysis
- independent component analysis
- non negative matrix factorization

Also newly added:

- Convolutional (V)AE (1D)
- Convolutional (V)AE (2D)
- Convolutional UMAP

Clustering



Clustering



Distance metrics

| Metric | definition |
|-------------|--|
| LI | $\sum_i u_i - v_i $ |
| L2 | $ u - v _2$ |
| L4 | $ u - v _4$ |
| wL2 | $\left\ \frac{u-v}{\sigma} \right\ _2$ |
| wL4 | $\left\ \frac{u-v}{\sigma}\right\ _4$ |
| Bray-Curtis | $\sum u_i - v_i / \sum u_i + v_i $ |
| Chebyshev | $\max_i u_i - v_i $ |
| Canberra | $\sum_i rac{ u_i - v_i }{ u_i + v_i }$ |
| correlation | $1 - \frac{(u-\bar{u})\cdot(v-\bar{v})}{ (u-\bar{u}) _2 (v-\bar{v}) _2}$ |
| Mahalanobis | $\sqrt{(u-v)C^{-1}(u-v)^T}$ |

Comparison with LOF and i-forest:



| Data set | # point | # outliers | |
|-------------|---------|------------|------|
| lympho | 148 | 18 | 6 |
| breastw | 683 | 9 | 239 |
| wine | 129 | 13 | 10 |
| vertebral | 240 | 6 | 30 |
| glass | 214 | 9 | 9 |
| pima | 768 | 8 | 268 |
| letter | 1600 | 32 | 100 |
| thyroid | 3772 | 6 | 93 |
| ionosphere | 351 | 33 | 126 |
| cardio | 1831 | 21 | 176 |
| wbc | 378 | 30 | 21 |
| annthyroid | 7200 | 6 | 534 |
| arrhythmia | 452 | 274 | 66 |
| vowels | 1456 | 12 | 50 |
| satellite | 6435 | 36 | 2036 |
| satimage-2 | 5803 | 36 | 71 |
| optdigits | 5216 | 64 | 150 |
| mammography | 11183 | 6 | 260 |
| shuttle | 49097 | 9 | 3511 |
| mnist | 7603 | 100 | 700 |
| pendigits | 6870 | 16 | 156 |
| musk | 3062 | 166 | 97 |
| sutp | 95156 | 3 | 30 |
| http | 567498 | 3 | 2211 |
| COVEL | 286048 | 10 | 2747 |
| speech | 3686 | 400 | 61 |

Simulated

Real

Benchmark metrics:



$$MCC = \frac{(TP \times TN) - (FP \times FN)}{\sqrt{(TP + FP)(TP + FN)(TN + FP)(TN + FN)}}$$

$$\text{RWS} = \frac{1}{N(N+1)} \sum_{i=1}^{N} w_i I_i$$

Table 6: Real dataset scores (all $\times 100$). You can see ensemble method results for each DRTs as first five columns (best is yellow). The next four columns are the recommended DRT, splitting level and metrics, LOF results (best for number of neighbors equals 5.10 and 35) and isolated forest, respectively (best is blue). Errors are 1 standard deviation and are only shown if larger than unity.

| Dataset | DRT-metric | AUC(%) | | | MCC(%) | | | RWS(%) | | |
|------------|------------------|-----------|-------------|------------|------------|-------------|-------------|---------------------|-------------|-------------|
| | | DRAMA | LOF | i-Forest | DRAMA | LOF | i-Forest | DRAMA | LOF | i-Forest |
| arrhythmia | NMF-L2 | 83 | 80 | 80 | 41 | 40 ± 3 | 40 ± 3 | 37 | 42 ± 3 | 42 ± 3 |
| shuttle | NMF-L2 | 99 | 95 ± 13 | 100 | 92 | 87 ± 27 | 96 ± 1 | 99 | 85 ± 25 | 93 ± 2 |
| smtp | NMF-L2 | 86 | 91 | 91 | 17 | 0 | 0 | 20 | 0 | 0 |
| lympho | NMF-wL2 | 99 | 100 | 100 | 83 | 86 ± 10 | 88±8 | 211 | 79 ± 14 | 50 ± 13 |
| thyroid | NMF-wL2 | 96 | 97 ± 2 | 198 | 43 | 52 ± 11 | 55 ± 5 | 42 | 19±10 | 52:#4 |
| anuthyroid | NMF-wL2 | 69 | 81±3 | 82 ± 2 | 18 | 26 ± 2 | 26 ± 1 | 18 | 24±1 | 24 ± 1 |
| musk | NMF-wL2 | TON: | 95 ± 14 | 100 | 0.001 | 88 ± 25 | 97 ± 3 | 100 | 80 ± 27 | 94 ± 5 |
| letter | NMF-Mah. | 88 | 65 ± 9 | 62 ± 1 | 42 | 7 ± 14 | 3 ± 2 | 33 | 12±8 | 10 ± 3 |
| vowels | NMF-Mah. | 95 | 77 ± 7 | 75±3 | 36 | 19 ± 7 | 17 ± 6 | 墨 | 18 ± 7 | 15 ± 5 |
| wine | ICA-L2 | 100 | 82 ± 6 | 79 ± 2 | 80 | 17 ± 25 | 8 ± 7 | 12 | 28 ± 22 | 21 ± 13 |
| cardio | ICA-L2 | 0.3 | 90 ± 10 | 93 ± 1 | 50 | 44 ± 11 | 48 ± 3 | 20-1 1 | 43 ± 10 | 46 ± 4 |
| optdigits | ICA-corr. | 81 | 71±5 | 71±5 | Ú. | 012 | 0±2 | -Ű | 4±3 | 3 ± 3 |
| mnist | ICA-corr. | 98 | 79 ± 3 | 80 ± 2 | 52 | 22 ± 5 | 23 ± 5 | 444 | 26 ± 3 | 27 ± 3 |
| glass | ICA-canb. | 03 | 70 ± 3 | 69 ± 1 | 30 | 8 ± 3 | 7 | 11 | 10 ± 4 | 9 ± 4 |
| http: | AE-L2 | 100 | 94 ± 19 | 100 | 89 | 87 ± 29 | 97 ± 2 | 998 | 88 ± 30 | 98 |
| whe | AE-L2 | Direction | 94 | 94 | 012.9 | 51 ± 3 | 51 ± 2 | $53 \pm 10^{\circ}$ | 41±5 | 41±5 |
| maninog. | AE-L2 | 88.+.2 | 85.±4 | 86 | 30:±:7 | 20 ± 3 | 20 ± 3 | 26 ± 8 | 15 ± 3 | 15 ± 2 |
| breastw | AE-L2 | 99 | 94 ± 16 | 99 | 95 | 78 ± 32 | 89 | 2.81 | 84 ± 16 | 89 |
| pima | AE-wL2 | 76.a.l. | 67 ± 3 | 68 | 30=6 | 25 ± 4 | 26 ± 2 | 183 ± 4 | 49±2 | 50 ± 2 |
| vertebral | VAE-cheb. | 79 - 6 | 37 ± 6 | 35 ± 1 | 20 ± 4 | -10 ± 2 | -11 ± 2 | 35±6 | 5 ± 3 | 4 ± 3 |
| ionusphere | VAE-Mah. | 503-11 | 86 ± 2 | 85 | 10+2 | 50 ± 8 | 47 ± 2 | To + 3 | 53 ± 8 | 50 ± 2 |
| pendigits | none-L4 | 981 | 91 ± 12 | 95 | 46 | 31 ± 9 | 34 ± 3 | 53 | 28 ± 8 | 31 ± 3 |
| satimage-2 | none-L4 | 200 | 96 ± 11 | 99 | 91. | 80 ± 23 | 87 ± 2 | 5 | 70 ± 22 | 77 ± 4 |
| COVEL | none-wL4 | 96 | 86 ± 10 | 89 ± 3 | 6 | 8±3 | 9 ± 2 | 50 | 7 ± 2 | 8 ± 2 |
| satellite | none-canb. | 77 | 69 ± 4 | 70 ± 1 | 38 | 33 ± 8 | 36 ± 2 | 39 | 39 ± 2 | 40 ± 1 |
| speech | none-canb. | 501 | 49 ± 5 | 48 ± 1 | 3 | 1 ± 3 | 0 ± 1 | Ω. | 2 ± 2 | 1±1 |

AE, VAE, PCA, ICA, NMF, ...? L1, L2, L4, Chebyshev, Canbera, ...?

Semi-supervised or active learning DRAMA vs. LOF vs. iforest



Local anomaly in low dimension



n_f=100

New class anomaly in low dimension



n_f=100

Local anomaly in high dimension



 $n_{f} = 3000$

New class anomaly in high dimension



 $n_{f} = 3000$

Averaged on real data



Widefield ouTlier Finder (WTF)

Includes 337 boring objects and 17 interesting

AUC: 87 MCC: 26 RWS: 31



boring



interesting

Future directions

- 2D convolutions
- PLAsTiCC (Kaggle competitions)
- Deep learning based clustering
- Graphical user interface
- Active learning mode
- Detect and classify
Thanks for your attention.

You can find DRAMA here https://github.com/vafaei-ar/drama

