# *Knowledge Extraction from Machine Learning*

**Luisa Lucie-Smith**
*University College London*

with H.V. Peiris (UCL, Stockholm), A. Pontzen (UCL),
M. Lochner (AIMS), B. Nord (Fermilab)

*Lucie-Smith, Peiris, Pontzen (2019), arXiv:1906.06339*
*Lucie-Smith, Peiris, Pontzen, Lochner (2018), arXiv:1802.04271*

# *Machine learning in Astronomy*

***Big-data era: many traditional applications of ML.***

- *data mining*
- *classification*
- *data compression*
- *data and model emulation*
- *regression*
- *clustering analyses*
- *outlier detection*

**Can ML enable knowledge extraction?**

- *can we extract new physical insights by studying the learning of ML algorithms?*

# The genetic modification method



Redshift 45.7
0.05 Gyr
Step 0

Suppressed merger        Reference        Enhanced merger

# *"Genetically-modified" galaxies*



$M_\star = 2.7 \times 10^{10}\ M_\odot$
$R_{1/2} = 2.9$ kpc

$M_\star = 2.8 \times 10^{10}\ M_\odot$
$R_{1/2} = 1.2$ kpc

$M_\star = 1.7 \times 10^{10}\ M_\odot$
$R_{1/2} = 2.0$ kpc

Suppressed merger          Reference          Enhanced merger

# Non-linear dark matter halo formation

$\rho$

*Gaussian random field*

x

$\rho$

x

*Dark matter halos*

# N-body simulations



Difficult **physical** interpretation from numerical studies alone

# Insights into dark matter halo collapse from ML?

***Approach***: Train ML algorithm to learn mapping between initial conditions and dark matter halos from N-body simulations



***Aim***: *gain new physical insights into the process of dark matter halo formation*

# ML algorithm: gradient boosted trees

GBTs add new trees to correct the mistakes of the previous ones



$$f_1(x) \quad + \quad f_2(x) \quad + \quad f_3(x) \quad + \quad \ldots$$

*Decision Tree*

# Feature Importance

$$\text{Imp}_t\left(X_j\right) = \sum_n \frac{N_n}{N_t}\left[p - \frac{N_{n_R}}{N_n}p_R - \frac{N_{n_L}}{N_n}p_L\right]$$

fraction of samples

impurity (MSE)

X₂ <= 0.8564

X₃ <= 1.2748

X₁ <= 1.0881

False

True

False

True

False

True

# *ML regression model of halo formation*

*Initial conditions (z=99)*                                    *Final halos (z=0)*

**Features**                          ML algorithm          **Output**
                                         (GBTs)
Properties of the                                      ***Mass of the halo***
***local environment***                                to which each DM
around                                                 particle will belong
DM particles                                           at z=0

***Our choice of features is motivated by existing analytic
approximations of halo collapse***

# *Features based on analytic theories of halo collapse*

1. *Density contrast: motivated by* **extended Press-Schechter theory**



$\delta > \delta c$

R

M(R)

Density contrast above threshold $\delta_c$

Dark matter halo of mass M(R)

2. Tidal shear field (ellipticity/prolateness: motivated by **Sheth-Tormen theory**



Tidal shear forces distort spheres into ellipses

Final halo mass M(R) depends on tidal shear field

Compute features in spheres of 50 different mass scales

# Which features were most informative?



Importance

$M_{\mathrm{smoothing}}/\mathrm{M}_\odot$

**Smoothing mass scale**

# *Machine learning model comparison: Kullback-Leibler (KL) divergence*

*1. Smooth distributions with KDE*

*2. KL divergence prediction vs truth*



**Legend (left plot):**
- Histogram
- KDE

**Legend (right plot):**
- Truth $(t)$
- Density $(d)$, $D_{\mathrm{KL}}(t||d) = 0.028$
- Density+Shear $(d+s)$, $D_{\mathrm{KL}}(t||d+s) = 0.025$

Left plot axes: $n_{\mathrm{particles}}$ vs $\log(M_{\mathrm{true}}/\mathrm{M}_\odot)$

Right plot axes: $n_{\mathrm{particles}}$ vs $\log(M/\mathrm{M}_\odot)$

***Addition of tidal shear information does not yield an improved halo collapse model in contrast to standard interpretations of Sheth-Tormen theory***

# Do the results generalise to independent simulations?

One training simulation

Four independent test simulations



*ML algorithm learnt* **physical connection** *between initial conditions and halo masses*

# *What we have learnt so far...*



**Local overdensity**

$M_{\rm halo}$

**Tidal shear field**

Addition of tidal shear information does not
improve halo collapse model

**How can we go beyond testing current interpretations of halo
collapse?**

# A deep learning approach to halo formation

## Advantages:

- *do not require featurization!*
- *provide as input the "raw data", i.e. the initial density field*

| Input | | Output |
|---|---|---|
| **Initial density field** realization of the simulation | DL algorithm → | **Mass of the halo** to which each DM particle will belong at z=0 |

## Disadvantages:

- *how do we extract physical knowledge from the DL algorithm?*

# Requirements for knowledge extraction from DL

- **interpretability**: *How did the DL model reach its predictions? Produce outputs that help us understand inner workings DL model.*



Edges (layer conv2d0)   Textures (layer mixed3a)   Patterns (layer mixed4a)   Parts (layer mixed4b,c)   Objects (layer mixed4d,e)

- **explainability**: *mapping interpretability onto existing knowledge in the relevant science domain.*

# Learning physical representations

## Human learning



observations → encoding → representation ($x_0 = \_\_$, $v = \_\_$) → decoding → answer

t' question

$x(t') = x_0 + v\, t'$ answer

## SciNet model



question $q$

observation $o$

answer $a$

encoder $E$   $r$   decoder $D$

latent representation

*SciNet learns two relevant physical parameters of damped pendulum problem*



Latent activation 1

Latent activation 2

$b\ [kg/s]$   $\kappa\ [kg/s^2]$

Iten et al. (2018; arXiv:1807.10300)

# Deep learning for knowledge extraction

**Supervised variational encoder**

*Initial conditions density field*

*Latent representation*

$\mathcal{N}(\mu, \sigma)$

*Mass of dark matter halo*

CONV. + POOL

$\mu_1$
$\mu_2$
$\mu_{i-1}$
$\mu_i$

$\sigma_1$
$\sigma_2$
$\sigma_{i-1}$
$\sigma_i$



*Latent variables encode most relevant aspects of initial conditions about final halo masses*

# Deep learning for knowledge extraction



- **Explainability**: *What physical information is compressed by neural network learning? Correlated with overdensities?*

- *Can provide different fields (e.g. density field and tidal shear field) as different 'channels' (like RGB channels for images)*

**Work in progress…**

# *Conclusions*

- ML enabled new, surprising and generalisable insights into halo collapse

- **Work in progress:** interpretable deep learning networks (no featurization) to extract new physical knowledge about cosmological structure formation