# Training dark matter classifiers with machine learning using ATLAS Monte Carlo data

## Introduction

Dark matter is generally considered a large scale feature of physics however it is possible to look for it within elementary particle interactions, more precisely the theoretical dark matter particle called a "weakly interacting massive particle" (WIMP). This project is to perform searches for the dark matter + mono-Z process $(pp \rightarrow \chi\bar{\chi} + Z)$ and is identified by missing energy (due to the dark matter particles) and two opposite charged leptons (See Fig. 1), which is detected by the ATLAS detector in the LHC.

The goal of this project was to develop machine learning models called binary classifiers, which are trained to identify dark matter events from the background events. This was then used in an interactive Jupyter notebook designed to teach six form students and above the fundamental concepts of machine learning and how it can be used in high energy physics, worked on by me and the team at the University of Sussex (See Fig. 2). The models where developed in python using sciKit learn.
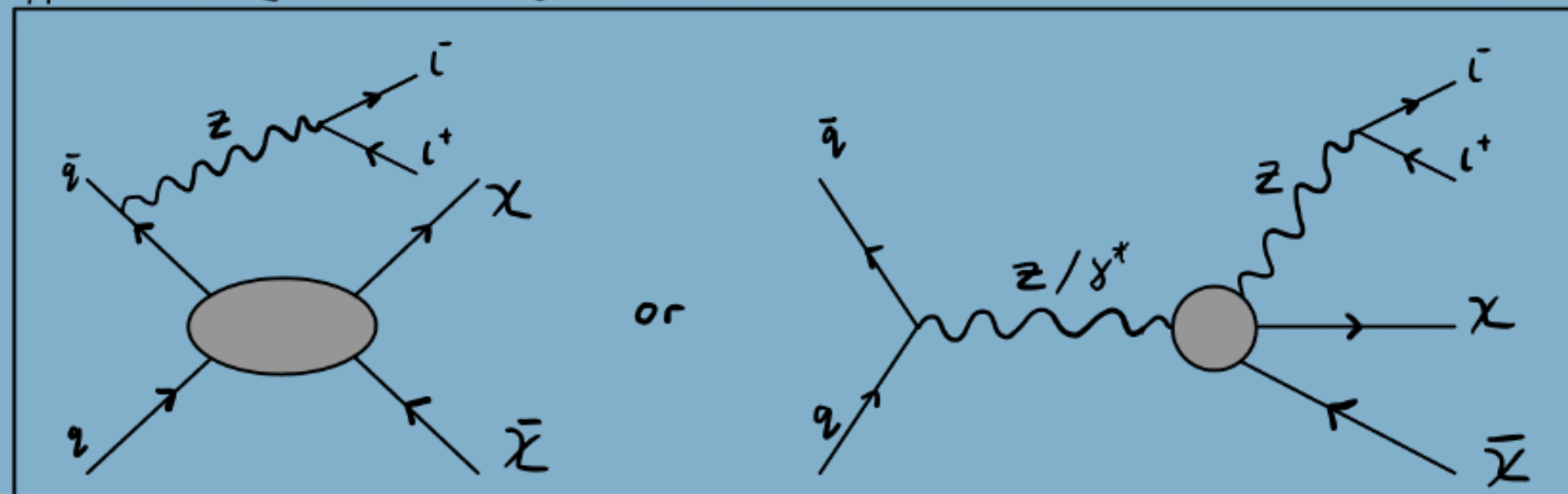


Fig 1. Possible Feynman diagrams for the dark matter + mono-Z process
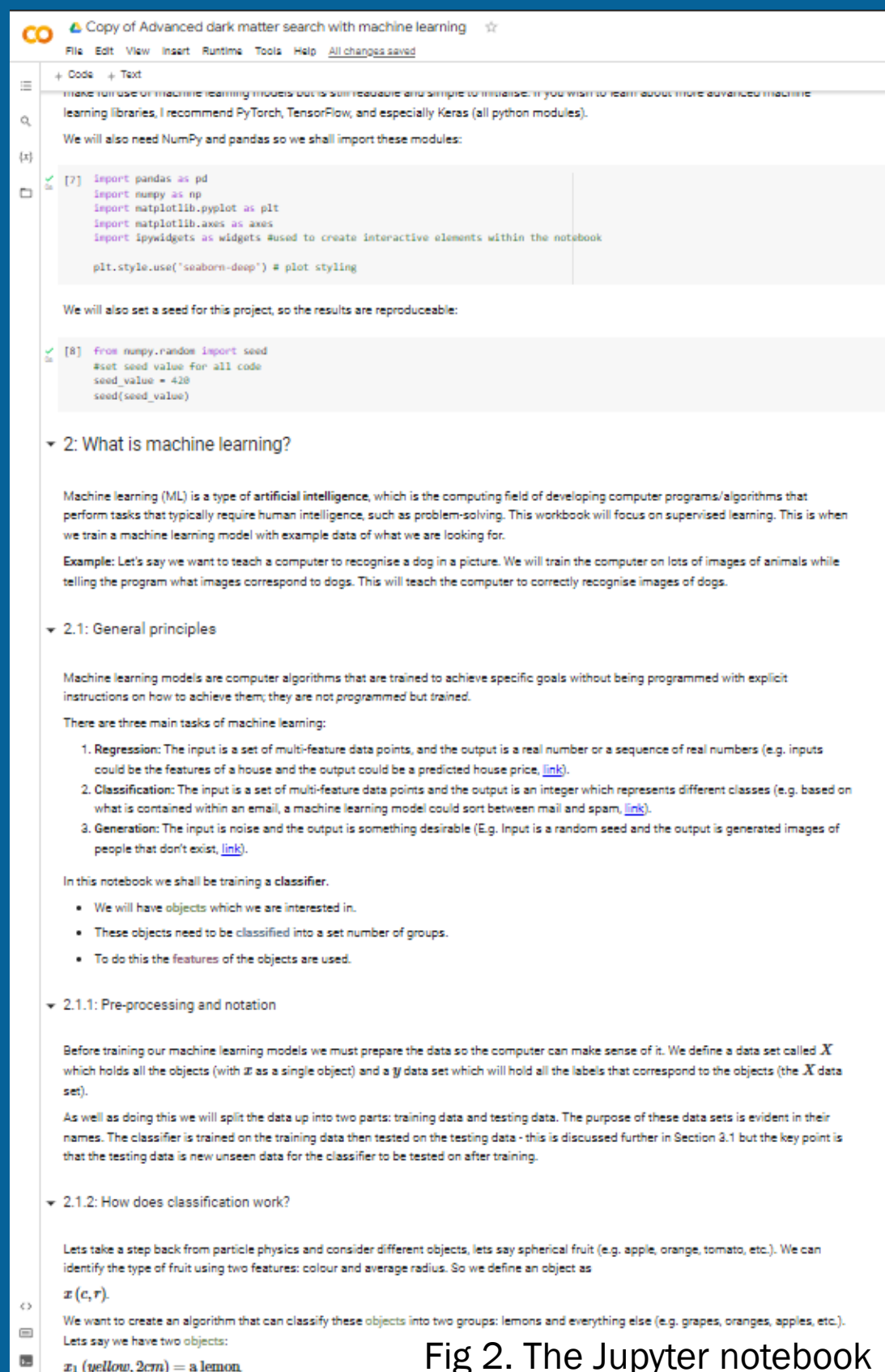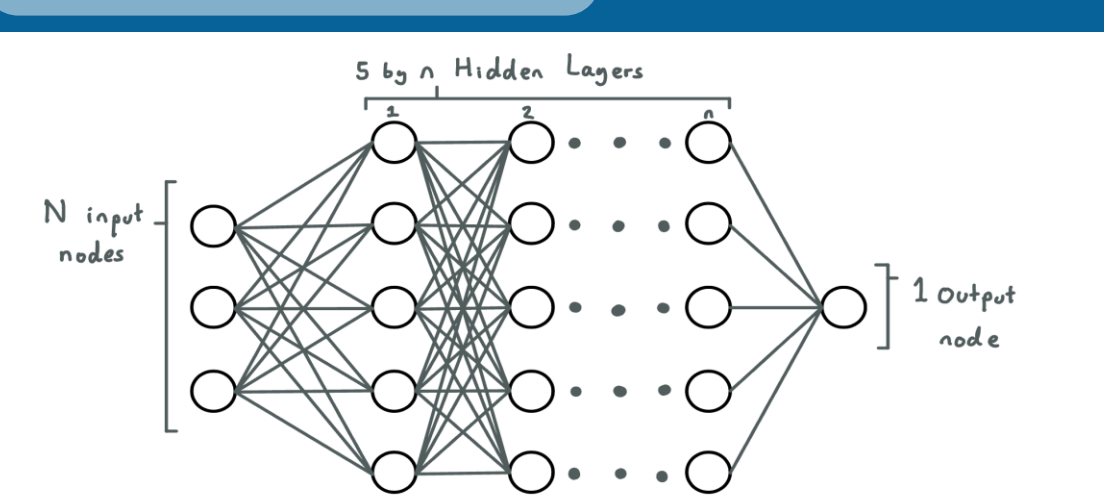


Fig 2. The Jupyter notebook

### ATLAS open data:

ATLAS Open Data provides data from Monte Carlo simulation of the proton-proton collision data collected by the ATLAS experiment for public use, this contains a huge number of labelled collision events including dark matter used in this training.
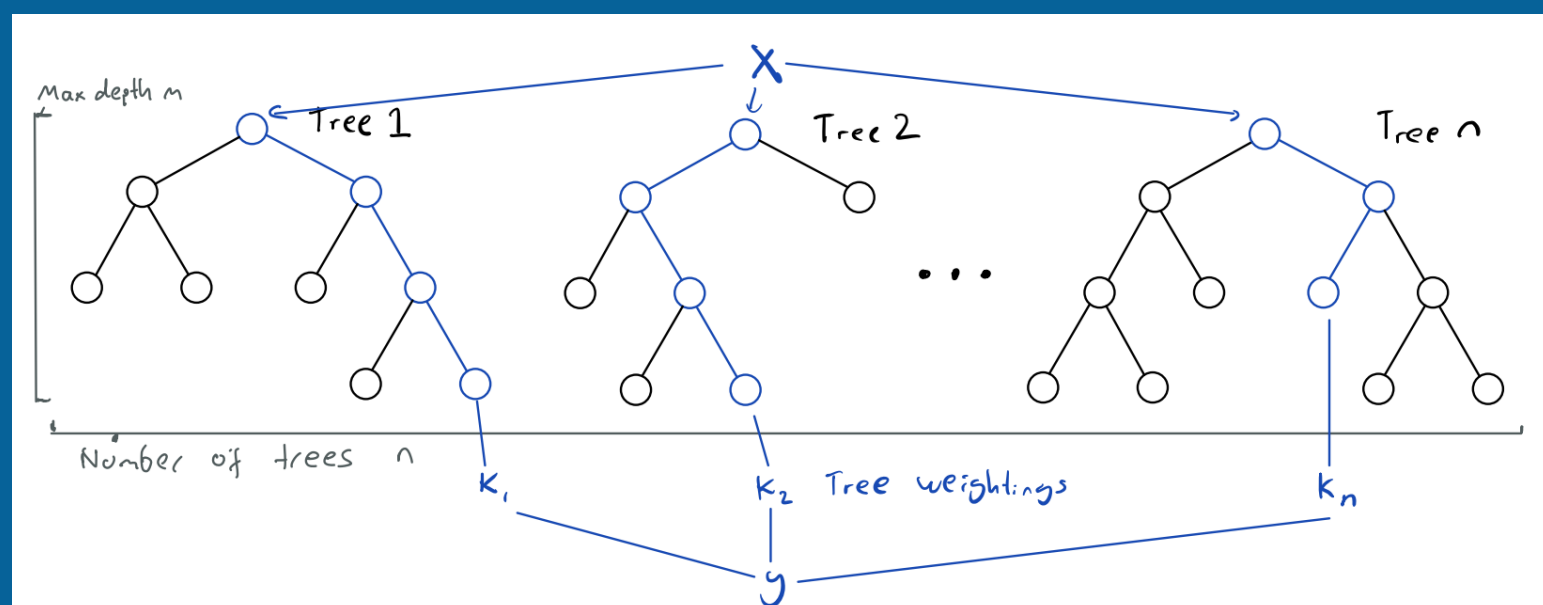
### Why?

- Experimental particle physics involves a lot of data analysis. The LHC produces up to **1 billion proton-proton collisions** per second.
  - Around one petabyte per day, or **106 gigabytes**
- The efficiency and speed of machine learning methods also results in a massive decrease in computing time compared with manual methods

### Neural Network Diagram



Random Forests and neural networks where the machine learning models used in this project.
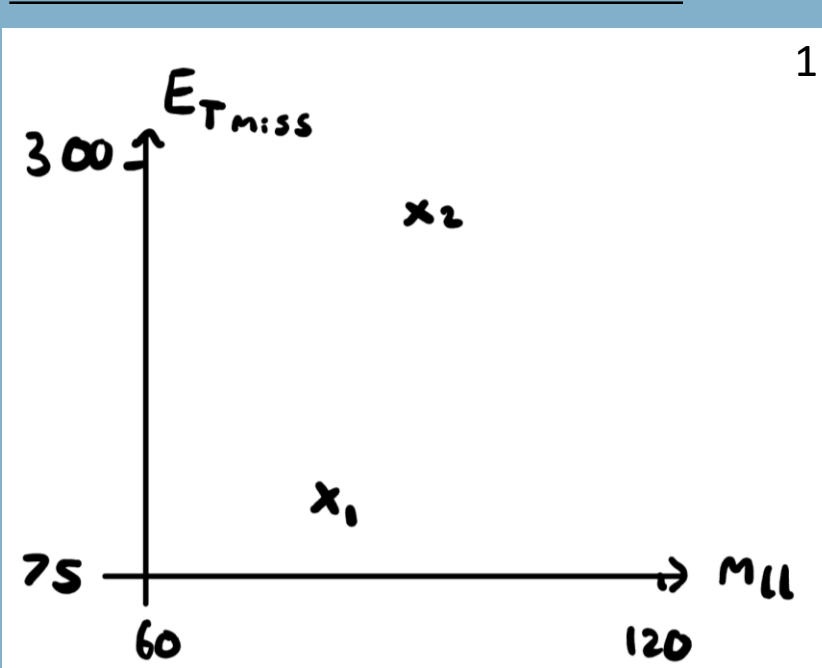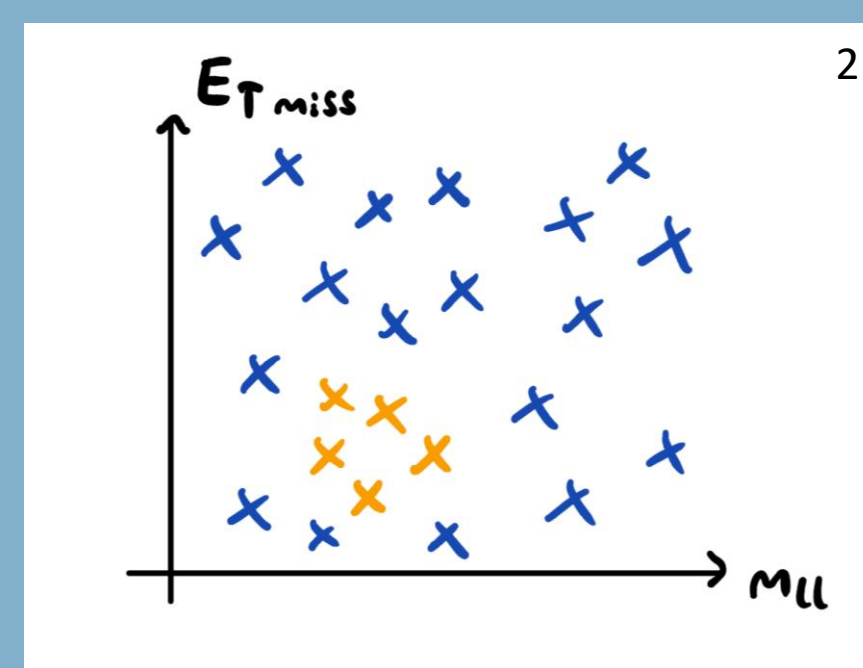
### Random Forest Diagram



### Significance Graphs →

A range of significance is found by performing a scan cut on the classifier response (Fig. 3&4). This measures the significance of the dark matter signal over the background.
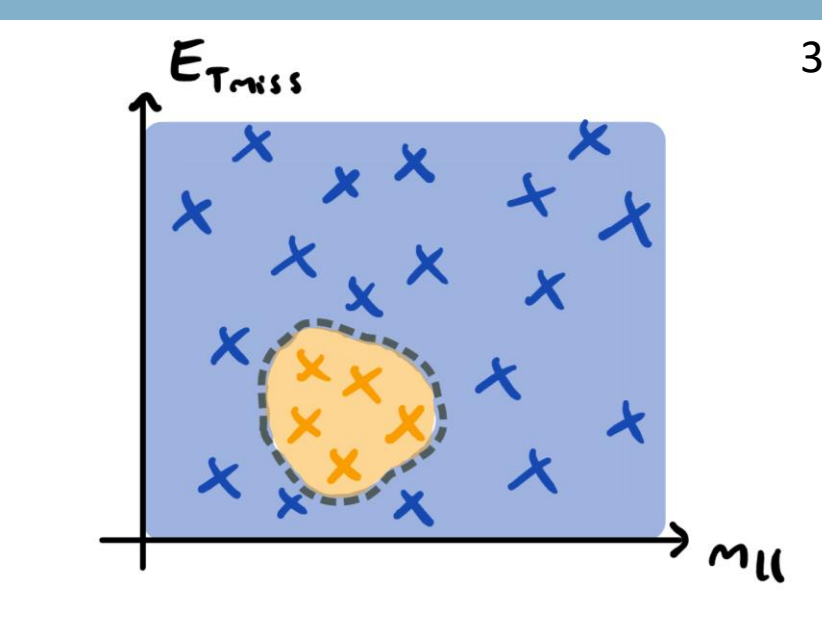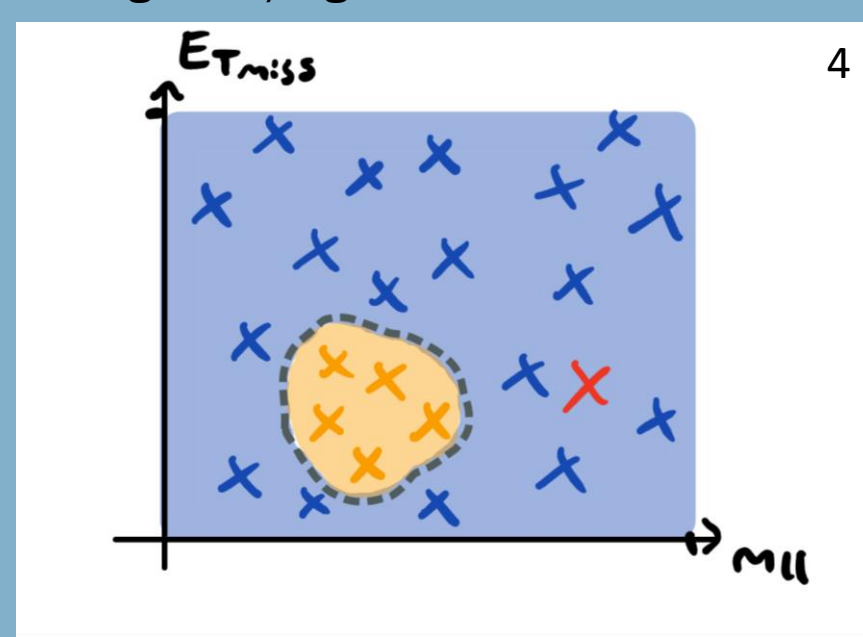
### How does classification work?



We can plot collision events by their features, in an axis of event proprieties.

In this case we have a dataset of labelled data with blue/orange is for background/signal.
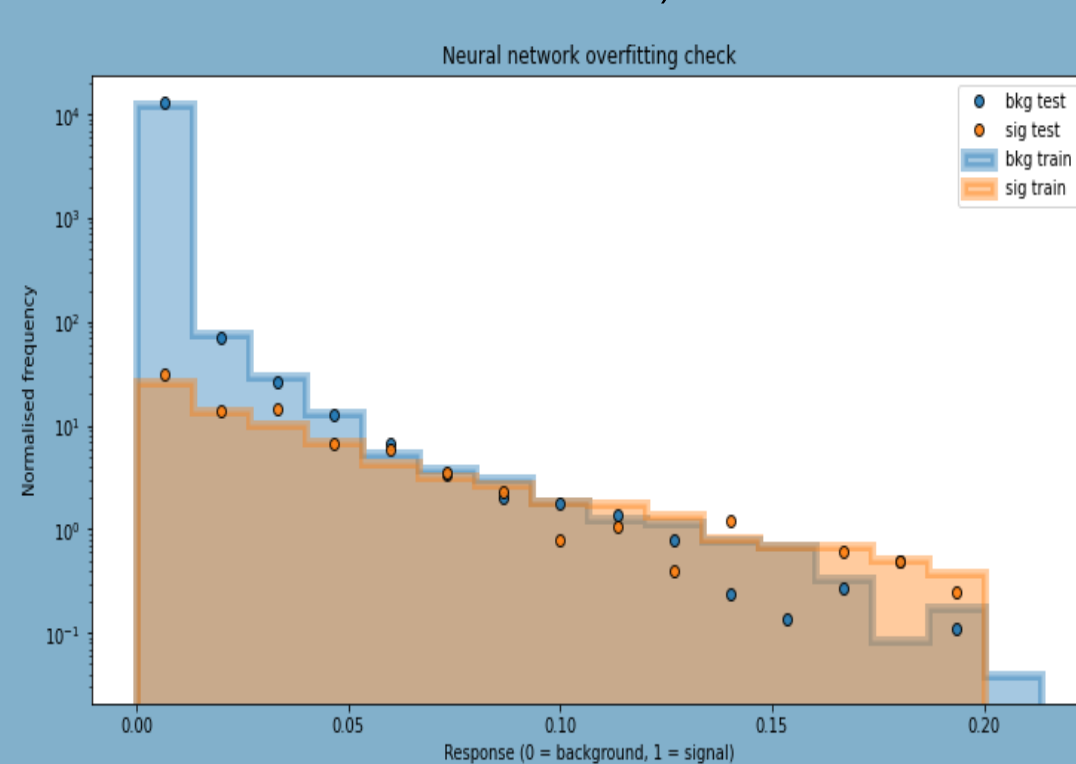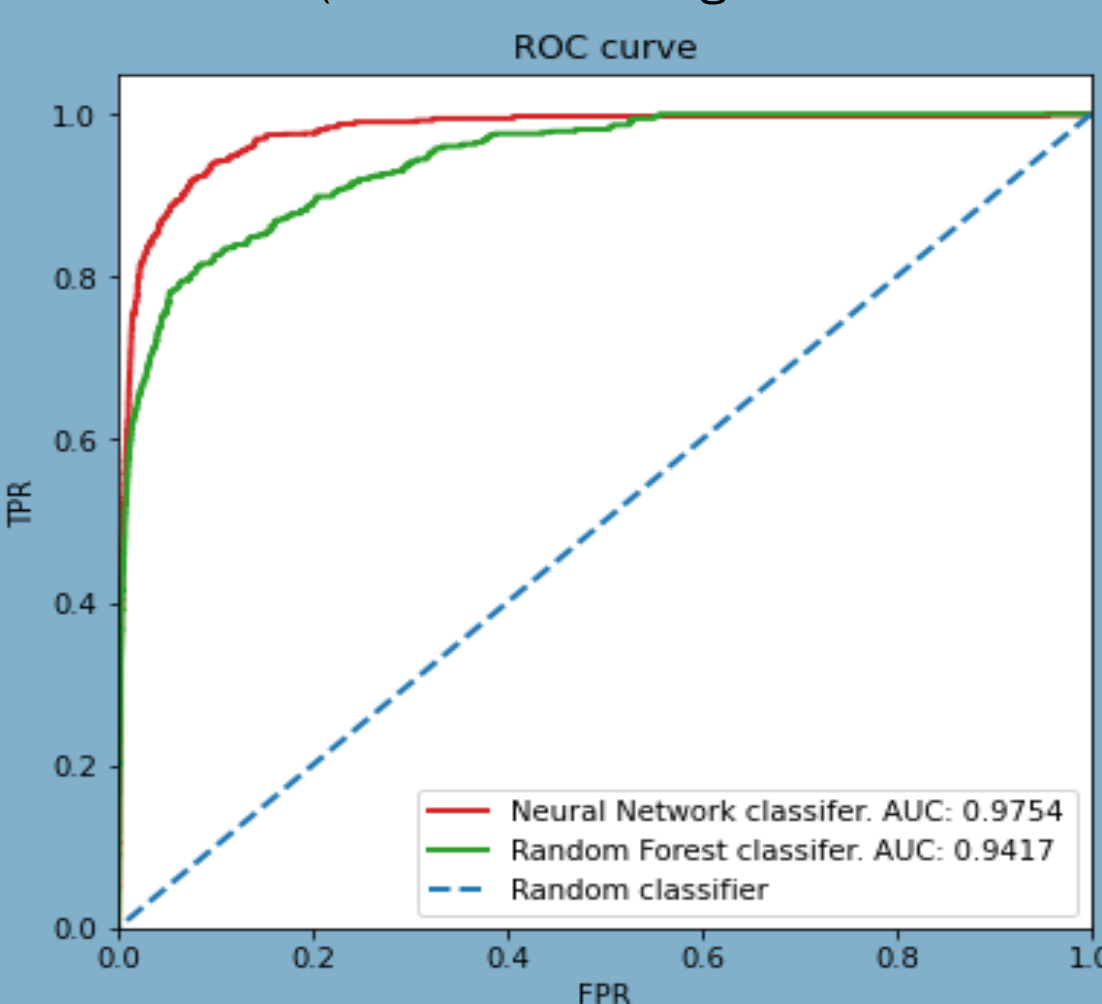
The machine learning model is trained to produce boundaries classifying the points.

If a new point was introduced the model could then predict what it was, here it would predict background.



Fig 3.



Fig 4.

## Model Analysis:

The machine learning models are analysed by plotting model signal and background responses, generating a ROC curve and checking overfitting including a Kolmogorov–Smirnov test (metric measuring match or mis-match between two distributions).



ROC: Receiver operating characteristic
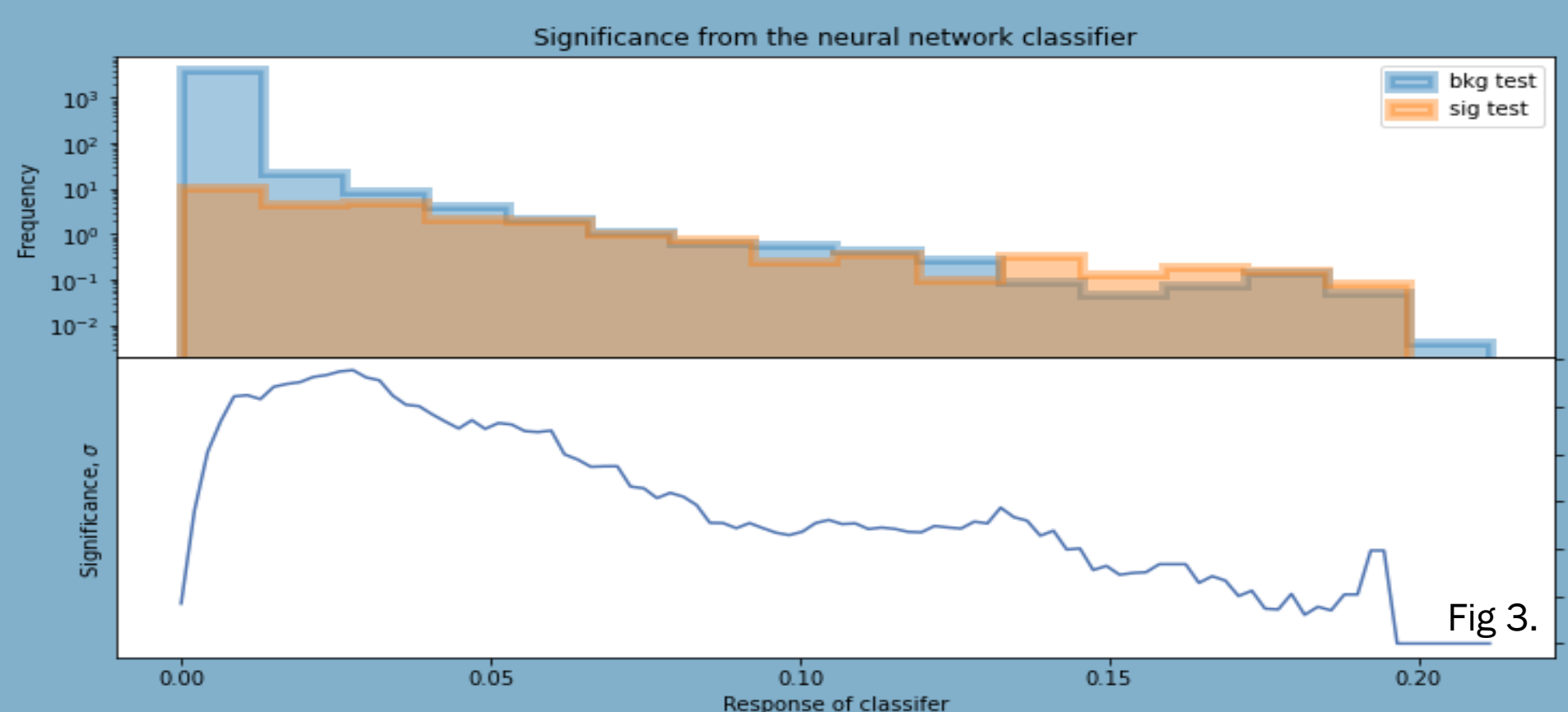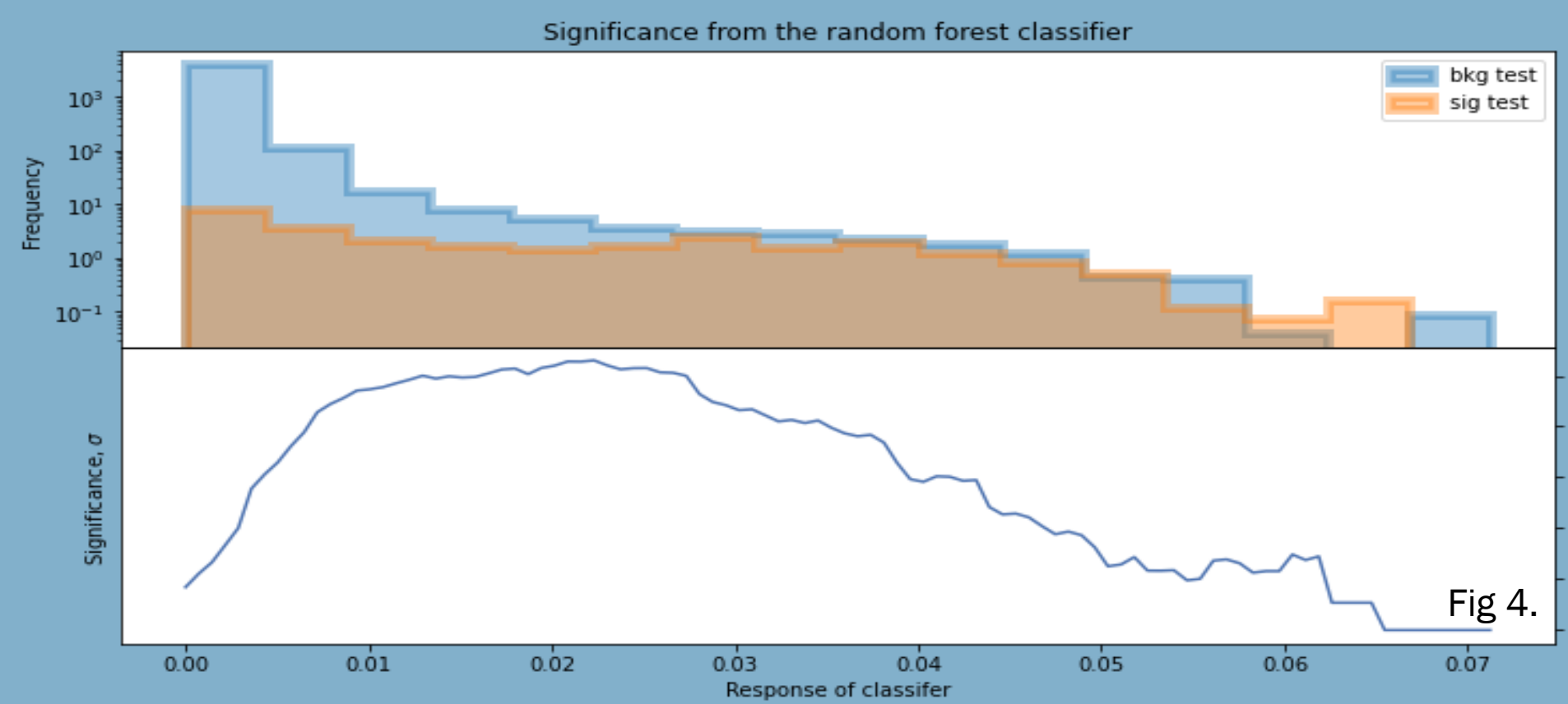
AUC : Area under curve (of ROC)

## Conclusions:

1. Machine learning can be taught in an intuitive way using real scientific applications.
2. Neural networks were found to be the most effective in signal classification.
3. Simplistic CPU bound machine learning models can show the principles of AI driven data analysis and statical significance.
4. Using neural networks and random forests dark matter signal significances of 2.37σ and 2.26σ were found respectively, with limited overfitting and acceptable ROC curves with AUC of 0.9754 and 0.9417.

Oscar Jackson - Physics with Space Science (MPhys) – 4th year
okcj1g19@soton.ac.uk