Maximising AMS by classifying Higgs Boson Data

Kieran Morris, Cecina Babich Morrow and Daniella Montgomery

of the statistical results and R code only.

https://github.com/babichmorrowc/sm1_large_hadron_collider

This RMarkdown document is a companion piece to a Latex document submitted detailing our work on this project. For discussion including references to other works and justification for the techniques and methods used, please refer to that document. This will be an indepth exploration The source code for this Rmarkdown and all annotated R code can be found in the following GitHub repository:

Additionally this project comes with a package which contains all necessary functions which can be downloaded via the following command in the

devtools::install_github("babichmorrowc/higgiesmalls") **Preliminary Observations**

Our data consists of events detected at the Large Hadron Collider which were classified as either Background (b) or Signal (s) - Signal being a

detection of a higgs boson.

We have a few types of data to consider:

2023-11-07

R console:

 All other variables are either direct or indirect measurements and are continuous. Missing Data

What is missing and why?

Show 11 ∨ entries skim_variable

observation.

1 DER mass MMC

n_missing complete_rate 0.8477191232868676 124602

Below are the variables which contain undefined values. By convention it was provided to use with values -999 which is out of range for every

Search:

Search:

21.4252

Difference

18.1382

13.6895

15.7553

15.7553

Background

4	DER_prodeta_jet_jet	580253	0.2908505838154669
5	DER_lep_eta_centrality	580253	0.2908505838154669
6	PRI_jet_leading_pt	327371	0.5999073619167039
7	PRI_jet_leading_eta	327371	0.5999073619167039
8	PRI_jet_leading_phi	327371	0.5999073619167039
9	PRI_jet_subleading_pt	580253	0.2908505838154669
10	PRI_jet_subleading_eta	580253	0.2908505838154669
11	PRI_jet_subleading_phi	580253	0.2908505838154669
Showing 1 to 11 of 11 entries Previous 1 Next Notice that every column besides DER_mass_MMC is a jet variable, and in those, we only have two values for completion_rate . In fact these correspond to different values of PRI_jet_num:			
0 jets: All jet variables were missing data.			
• 1 jet: only PRI_leading_pt , PRI_leading_eta and PRI_leading_phi have data.			
2 or more jets: All jet variables have data.			
The above result can be found in the handbook for the variables provided with the Kaggle challenge. Unfortunately DER_mass_MMC does not have such an explanation and may be a result of some event during measurement. However it is still assigned the same -999 value as the other missing data.			

Despite understanding the cause of (most of) our missing data, we still don't know the impact of it, the following section is dedicated to understanding the correlation between the missing data from each variable and its classification. Below we compute the percentage of NA data in Background and Signal respectively. If this missing data is distributed uniformly across

- Background and Signal then the percentages should be very close.
- Show 11 ∨ entries % of Missing Background and Signal Data

DER_deltaeta_jet_jet 61.9026 75.5921 13.6895 61.9026 75.5921 13.6895 DER_mass_jet_jet DER prodeta jet jet 61.9026 75.5921 13.6895

Signal |

3.287

15.7553 PRI_jet_subleading_pt 61.9026 75.5921 13.6895 PRI_jet_subleading_eta 61.9026 75.5921 13.6895 61.9026 PRI_jet_subleading_phi 75.5921 13.6895 Showing 1 to 11 of 11 entries Previous 1 Next Notice that Signal consistently has 13-18% less missing data than Background . Considering the size of our dataset (\$ ^5) this is statistically significant. This means that in both the DER_mass_MMC and jet-variables cases, the amount of missing data is indicative of a classification. Meaning that including -999 can help with classification and will not be removed. Variable Selection Co-Dependency of Data We have the following algebraic codependencies: • PRI_jet_all_pt = PRI_jet_leading_pt + PRI_jet_subleading_pt • PRI_lep_pt = DER_pt_ratio_lep_tau*PRI_tau_pt We can either remove $PRI_jet_all_pt$ or remove both $PRI_jet_leading_pt$ and $PRI_jet_subleading_pt$ to reduce the dimension by 1 or 2. This same idea can be applied to to the other (multiplicative) dependence. Taking us down to possibly 26 dimensions. Another 'dependence' is the fact that PRI_jet_num characterises the missing data in jet-variables. Meaning we could remove all of these in favour of keeping PRI_jet_num however this is very extreme and may produce detrimental effects. **Density Plots** To get an understanding of the behavior of variable given their classification, we condition each variable on whether y = b or y = s respectively.

Label b s 0.025 -0.15

0.05 -

ί PRI_tau_phi

0.000 -DER_mass_transverse_met_lep

PRI_met_phi PRI_lep_phi PRI_tau_phi PRI_tau_eta DER_pt_tot

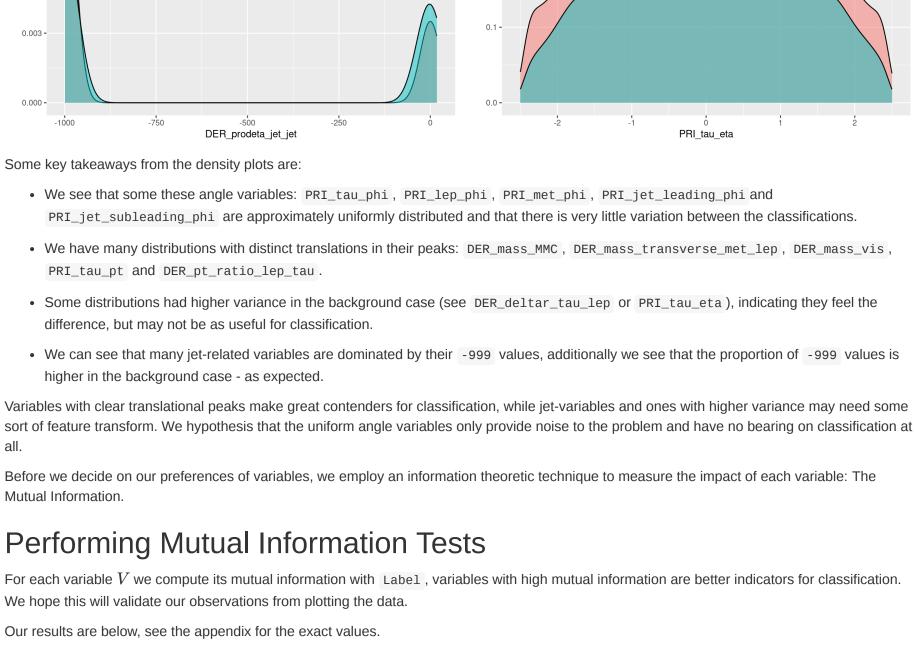
PRI_lep_pt

PRI_jet_subleading_phi

500

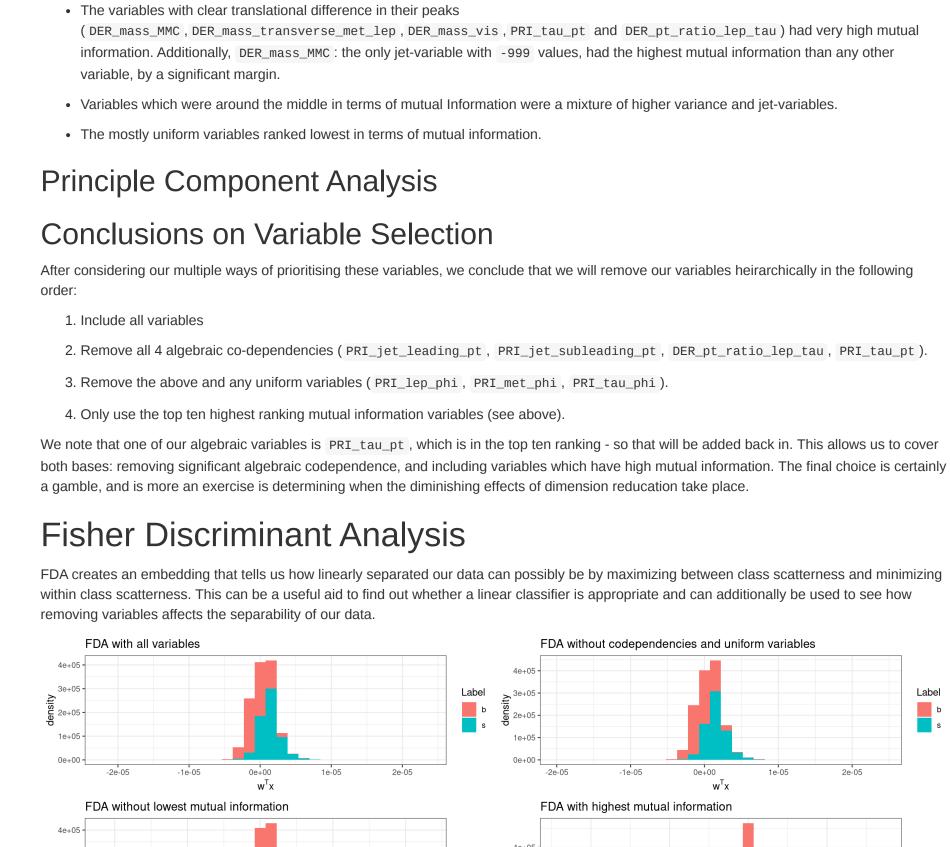
0.015

0.005



PRI_jet_subleading_pt PRI lep eta PRI jet_leading_phi DER_deltar_tau_lep PRI jet num PRI_jet_subleading_eta DER_lep_eta_centrality PRI_jet_leading_eta

Mutual Information Values



For these different FDA embeddings, we can compare the ratio between between class scatterness and within class scatterness:

Label

Scatter ratio

0.3251629060327355

0.3235810938370512

0.3248675982049276

0.2540494515822303

Next

Previous

Search:

Search:

Cost

1

2

2

2

500

DER_pt_h

DER_mass_jet_jet

PRI_tau_pt

PRI_tau_phi

PRI_jet_leading_phi

-1

2000

Previous

1

Gamma

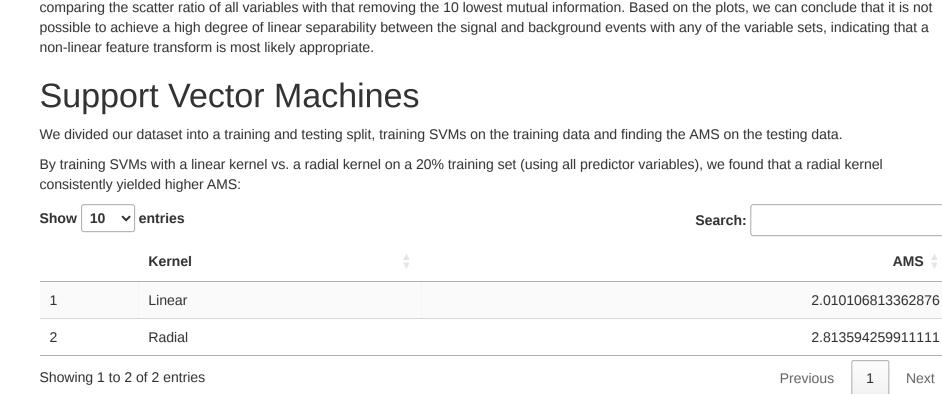
0.1

0.1

0.1

0.5

Next



Based on these results, along with the results of FDA indicating that the dataset is not linearly separable, we moved forward using a radial kernel. We performed 5-fold cross-validation over the 20% training set to tune the two hyperparameters using a grid search: cost and γ (bandwidth). We

dependencies and uniformly distributed variables; 3) Removing the 10 variables with the lowest mutual information; 4) Using only the 10 variables

allowed cost to take on the values 0.5, 1, and 2 and γ to take on the values 0.01, 0.1, and 0.5. Tuning was performed using our function ams_tune_parallel, which modifies the e1071::tune function to maximise AMS and to conduct cross-validation in parallel, dramatically

Below are the results of the hyperparameter tuning using four sets of variables: 1) Including all variables; 2) Removing the algebraic co-

We find that removing variables slightly lowers the linear separability of the data, but in general the values are comparable, particularly when

Variables included **Testing AMS** 1 All variables 2.832435958468929 2 Without co-dependencies and uniform variables 2.856642255556614 3 Without lowest mutual information 2.866093570197652 Highest mutual information 2.812699039205425 Showing 1 to 4 of 4 entries 1 Previous Next We see in the above data that removing the ten lowest mutual information produced the highest AMS when using an RBF kernel, followed by removing algebraic codependencies and the uniform variables. Using only the top ten highest mutual information actually yielded the lowest AMS, highlighting the trade-off between prioritizing high mutual information versus removing too many variables. **Appendix** Density plots of each variable

→ 0.015-

0.005 -

0.015

0.005

0.000

0.008 -0.006 -

0.004 -

0.002 -

-500 DER_prodeta_jet_jet -750 -250 DER_deltar_tau_lep 0.012 -0.08 -0.06 -0.009 -0.006 -0.04 -0.02 -0.003 -0.00 -0.000 -1000 DER_sum_pt 1500 DER_pt_tot 0.8 -0.6 -1.5 -0.4 -1.0 -0.2 -0.5 -0.0 -0.0 --0.5 0.0 DER_met_phi_centrality DER_pt_ratio_lep_tau 0.06 -0.009 -

0.00 -

0.15 -

0.05 -

0.00 -

0.3 -

0.2 -

0.15 -0.02 -0.01 -0.05 -0.00 -0.00 -PRI_lep_phi PRI_met 0.15 -0.003 -0.002 -0.05 -0.001 -PRI_met_sumet 0.005 -0.004 -0.003 -2-0.002 -0.001 -0.000 --1000 -500 PRI_jet_num PRI_jet_leading_pt

0.000 -1000 PRI_jet_subleading_eta PRI_jet_subleading_pt PRI_jet_subleading_phi Mutual Information Results Table Search: **Mutual Information** DER_mass_MMC 0.157847348 0.099242021 DER_mass_transverse_met_lep DER_mass_vis 0.0909557

DER_pt_h DER_deltaeta_jet_jet DER_mass_jet_jet

Previous

Types of Data • Variables KaggleSet and KaggleWeight can be ignored as they denote which data points were in the provided Kaggle challenge and their relative weights, which is irrelevant to us. ullet The discrete classification variable Label which takes values in $\{b,s\}$ • The continuous variable Weight, which will be used to compute the AMS, and will not be used in classification. • The discrete variable PRI_jet_num denotes the amount of jets from each event and takes values in $\{0,1,2,3\}$. We have multiple variables which are the angle of detection for some observation.

2 DER_deltaeta_jet_jet 580253 0.2908505838154669 3 DER mass jet jet 580253 0.2908505838154669 Impact of the Missing Data

DER_mass_MMC

61.9026 75.5921 DER_lep_eta_centrality 29.6369 PRI_jet_leading_pt 45.3922 29.6369 45.3922 PRI_jet_leading_eta 29.6369 PRI_jet_leading_phi 45.3922

The entire list can be found in the appendix but here are a few of note:

PRI_jet_leading_pt PRI_met PRI met_sumet DER_prodeta_jet_jet DER mass jet jet DER_deltaeta_jet_jet DER pt h DER_sum_pt DER_pt_ratio_lep_tau DER_met_phi_centrality PRI tau pt DER_mass_vis R mass transverse met lep DER_mass_MMC 0.05 Mutual Information

Thankfully these results justify many of our direct observations from plotting the densities:

improving the performance of tuning compared to the original package.

Without co-dependencies and uniform variables

Without lowest mutual information

Highest mutual information

with the highest mutual information.

All variables

Variables included

Show 10 ∨ entries

Showing 1 to 4 of 4 entries

-1000

DER_mass_vis

DER_deltaeta_jet_jet

DER_lep_eta_centrality

PRI_tau_eta

200

-750

PRI_jet_leading_eta

0.020

0.015 0.010

0.005 0.000 -

0.006 -

0.003

0.003

0.000

0.000 -

0.2 -

0.1 -

0.0 -

0.04 -

0.03 -

0.02 -0.01 -

0.004

0.002

0.000 -

0.0100 -0.0075 -

density 0.0050 -

0.0025 0.0000

0.009

0.006 -

0.003

dt

1

2

3

14

12

11

10

4

5

6

Show 10 ∨ entries

Variable

PRI_tau_pt

DER_sum_pt

Showing 1 to 10 of 29 entries

DER_pt_ratio_lep_tau

-1000

-1000

5

6

61

9

Show 10 ∨ entries

Showing 1 to 4 of 4 entries

4

Variables included

Highest mutual information

We then used the model using the best hyperparameter values from tuning to find the AMS value on the 80% testing dataset. Show 10 ∨ entries Search: The following are the density plots for our variables, conditioned on whether they are a detection of a higgs boson or not. Note that in the jetvariables and DER_mass_MMC we have peaks at -999 which skew many of the variable plots.

density 90000 0.04 -0.003 -0.02 -

0.004 -

0.002

0.006 -

0.003 -

0.064883381 DER_met_phi_centrality 0.048218121

> 0.033387741 0.032628663 3

0.04355779

0.039696267

0.03402485

Next