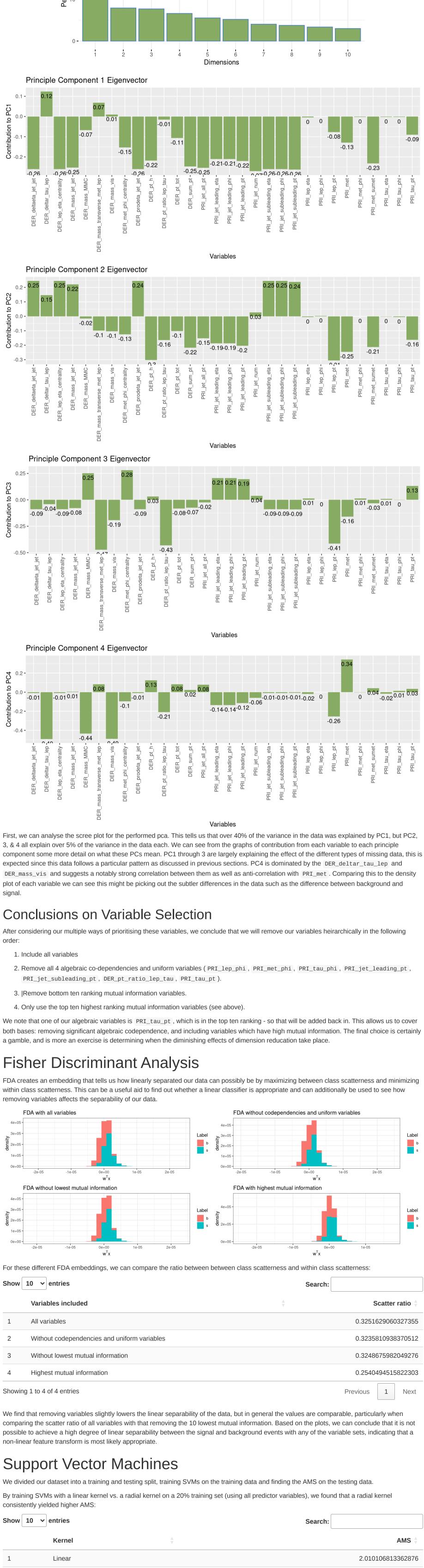
Maximising AMS by classifying Higgs Boson Data Kieran Morris, Cecina Babich Morrow and Daniella Montgomery 2023-11-07 This RMarkdown document is a companion piece to a Latex document submitted detailing our work on this project. For discussion including references to other works and justification for the techniques and methods used, please refer to that document. This will be an indepth exploration of the statistical results and R code only. The source code for this Rmarkdown and all annotated R code can be found in the following GitHub repository: https://github.com/babichmorrowc/sm1_large_hadron_collider Additionally this project comes with a package which contains all necessary functions which can be downloaded via the following command in the R console: 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. Types of Data We have a few types of data to consider: • 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. • 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. All other variables are either direct or indirect measurements and are continuous. Missing Data What is missing and why? 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 observation. Show 11 ∨ entries Search: skim variable n_missing complete_rate DER mass MMC 124602 0.8477191232868676 1 2 0.2908505838154669 DER_deltaeta_jet_jet 580253 3 DER_mass_jet_jet 580253 0.2908505838154669 DER_prodeta_jet_jet 580253 0.2908505838154669 DER lep eta centrality 580253 0.2908505838154669 6 0.5999073619167039 PRI jet leading pt 327371 0.5999073619167039 PRI_jet_leading_eta 327371 8 PRI_jet_leading_phi 327371 0.5999073619167039 9 580253 0.2908505838154669 PRI jet subleading pt PRI_jet_subleading_eta 10 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. Impact of the 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 v entries Search: % of Missing Background and Signal Data Signal Background Difference 3.287 21.4252 18.1382 DER mass MMC DER deltaeta jet jet 61.9026 75.5921 13.6895 61.9026 DER_mass_jet_jet 75.5921 13.6895 61.9026 75.5921 13.6895 DER_prodeta_jet_jet DER_lep_eta_centrality 61.9026 13.6895 75.5921 29.6369 15.7553 PRI_jet_leading_pt 45.3922 PRI_jet_leading_eta 29.6369 45.3922 15.7553 29.6369 45.3922 15.7553 PRI_jet_leading_phi PRI jet subleading pt 61.9026 75.5921 13.6895 PRI jet subleading eta 61.9026 75.5921 13.6895 PRI jet subleading phi 61.9026 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. The entire list can be found in the appendix but here are a few of note: Label b s 0.025 -0.020 0.05 -0.005 0.000 ύ PRI_tau_phi DER_mass_transverse_met_lep 0.003 DER_prodeta_jet_jet PRI_tau_eta Some key takeaways from the density plots are: • We see that some these angle variables: PRI_tau_phi, PRI_lep_phi, PRI_met_phi, PRI_jet_leading_phi and PRI_jet_subleading_phi are approximately uniformly distributed and that there is very little variation between the classifications. • We have many distributions with distinct translations in their peaks: DER_mass_MMC, DER_mass_transverse_met_lep, DER_mass_vis, PRI_tau_pt and DER_pt_ratio_lep_tau. • Some distributions had higher variance in the background case (see DER_deltar_tau_lep or PRI_tau_eta), indicating they feel the difference, but may not be as useful for classification. • We can see that many jet-related variables are dominated by their -999 values, additionally we see that the proportion of -999 values is higher in the background case - as expected. Variables with clear translational peaks make great contenders for classification, while jet-variables and ones with higher variance may need some sort of feature transform. We hypothesis that the uniform angle variables only provide noise to the problem and have no bearing on classification at Before we decide on our preferences of variables, we employ an information theoretic technique to measure the impact of each variable: The Mutual Information. **Performing Mutual Information Tests** For each variable V we compute its mutual information with Label, variables with high mutual information are better indicators for classification. We hope this will validate our observations from plotting the data. Our results are below, see the appendix for the exact values. **Mutual Information Values** PRI_met_phi PRI lep phi PRI tau phi PRI tau eta DER_pt_tot PRI_jet_subleading_phi PRI_lep_pt 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 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 Mutual Information Thankfully these results justify many of our direct observations from plotting the densities: • The variables with clear translational difference in their peaks (DER_mass_MMC , DER_mass_transverse_met_lep , DER_mass_vis , PRI_tau_pt and DER_pt_ratio_lep_tau) had very high mutual information. Additionally, DER_mass_MMC: the only jet-variable with -999 values, had the highest mutual information than any other variable, by a significant margin. • Variables which were around the middle in terms of mutual Information were a mixture of higher variance and jet-variables. • The mostly uniform variables ranked lowest in terms of mutual information. Principle Component Analysis Here we use principle component analysis as a tool for visualization to build an intuition of the predictive power of our dataset when we include the missing data with the values left as -999 and don't include the variables with any data missing, respectively. PC1 vs PC2 PC2 vs PC3 PC2 PC1 PC1 vs PC2 PC2 vs PC3 PC1 vs PC3 -10 PC1 PC1 From these plots we can see that using the missing data gives us improved separability. Scree plot Percentage of explained variances Principle Component 1 Eigenvector Contribution to PC1 0.01 -0.01 -0.07 -0.08 -0.09 -0.11 -0.13 -0.15 -0.21-0.21_{-0.22} DER_sum_pt-0.25 DER_lep_eta_centrality 50 PRI jet_num -2 PRI jet_subleading_eta -2 PRI jet_subleading_phi -2 PRI jet_subleading_pt -2 PRI jet_subleading_pt -2 DER_prodeta_jet_jet_5 PRI_jet_leading_phi -DER_mass_vis-DER_pt_h-PRI_jet_leading_eta -PRI_lep_eta-PRI_lep_phi PRI_met_phi PRI_lep_pt PRI_tau_pt PRI_met DER_met_phi_centrality DER mass MMC Variables Principle Component 2 Eigenvector Contribution to PC2 0 0 -0.1 -0.1 -0.13 -0.1 -0.15 -0.16 DER_pt_h PRI_lep_pt -PRI_tau_pt ading_phi -PRI_lep_eta-PRI_lep_phi PRI_met DER_met_phi_centrality PRI_jet_leading_eta PRI_jet_leading_phi PRI_jet_num PRI_jet_all_pt DER_mass_jet_jet DER mass MMC PRI_jet_subl Principle Component 3 Eigenvector Contribution to PC3 -0.09-0.08 -0.08-0.07 -0.09 -0.09 -0.09-0.09-0.09 -0.16 -0.19 -0.41 DER_pt_tot -DER pt h PRI_lep_pt DER_mass_jet_jet DER_mass_vis PRI_jet_all_pt **Variables** Principle Component 4 Eigenvector Contribution to PC4 -0.01 0.01 -0.01-0.01-0.01_{-0.02} 0 -0.01 -0.14-0.14-0.12 -0.21 -0.26 DER_pt_h PRI_jet_all_pt PRI_m



2

5

6

61

9

1

2

3

dt

1

2

3

14

12

11

10

4

5

6

0.02 -

0.00 -

0.020 -

0.010 -

0.000 -

0.009 -

0.006 -

0.003 -

0.000 -

0.006 -

0.003 -

0.000 -

0.08 -

0.02 -

0.00 -

0.6 -

0.009 -

0.003 -

0.000 -

0.04 -0.03 -0.02 -0.01 -

0.15 -

0.05 -

0.15 -

0.05 -

0.00 -

density

0.002 -

0.000 -

0.0100 -

0.0075 -

0.0025 -

0.009 -

0.003 -

0.000 -

sity 0.006 -

-1000

-2

sity 0.006 -

Radial

improving the performance of tuning compared to the original package.

Without co-dependencies and uniform variables

Without co-dependencies and uniform variables

Without lowest mutual information

Mutual Information Results Table

Highest mutual information

Without lowest mutual information

Highest mutual information

Showing 1 to 2 of 2 entries

with the highest mutual information.

All variables

Variables included

Show 10 ∨ entries

Showing 1 to 4 of 4 entries

Show 10 ∨ entries

Variables included

All variables

Showing 1 to 4 of 4 entries

Appendix

Show 10 ∨ entries

Variable

DER_mass_MMC

DER_mass_vis

DER_met_phi_centrality

DER_pt_ratio_lep_tau

DER_deltaeta_jet_jet

DER_mass_jet_jet

Density plots of each variable

variables and DER_mass_MMC we have peaks at -999 which skew many of the variable plots.

DER_mass_MMC

DER_mass_vis

DER_deltaeta_jet_jet

DER_prodeta_jet_jet

DER_pt_tot

DER_pt_ratio_lep_tau

-500

DER_lep_eta_centrality

PRI_lep_pt

PRI_lep_phi

PRI_met_phi

-500 PRI_jet_leading_eta

PRI_jet_subleading_pt

PRI_jet_subleading_phi

PRI_tau_pt

DER_sum_pt

DER_pt_h

Showing 1 to 10 of 29 entries

DER mass transverse met lep

2.813594259911111

1

Next

Gamma

1

Testing AMS

2.832435958468929

2.856642255556614

2.866093570197652

2.812699039205425

Mutual Information

0.157847348

0.099242021

0.0909557

0.064883381

0.048218121

0.04355779

0.039696267

0.03402485

0.033387741

0.032628663

Next

3

Next

0.1

0.1

0.1

0.5

Next

Previous

Search:

Search:

Search:

Previous

DER_mass_transverse_met_lep

2000 DER_mass_jet_jet

DER_deltar_tau_lep

1000 DER_sum_pt

0.0 DER_met_phi_centrality

PRI_tau_pt

PRI_lep_eta

1000

PRI_met_sumet

PRI_jet_leading_pt

PRI_jet_leading_phi

PRI_jet_subleading_eta

1500

600

Cost

2

2

2

Previous

Previous

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 then used the model using the best hyperparameter values from tuning to find the AMS value on the 80% testing dataset.

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,

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 jet-

Label b s

0.020 -

0.005 -

0.000 -

0.020 -

0.015 -

0.005 -

0.000 -

0.008 -

0.006 -

0.004 -

0.75 **-**

0.25 -

0.009 -

0.003 -

0.000 -

0.5 -

0.04 -

0.02 -

0.00 -

0.05 -

0.02

0.01 -

0.004 -0.003 -0.002 -

0.001 -

0.000 -

0.005 0.004 0.003 0.002 0.001 0.000 -

0.006 -

0.002 -

0.000 -

0.009 -

0.003 -

0.000 -

sity - 900.00

> 0.015 -

highlighting the trade-off between prioritizing high mutual information versus removing too many variables.