

# Project 1 on Machine Learning team Yoor

Sergei Volodin<sup>1</sup>, Baran Nama<sup>1</sup>, and Name here 3<sup>1</sup>

<sup>1</sup>EPFL

{sergei.volodin,baran.nama,email3}@epfl.ch

**Abstract**—A classification dataset from the Large Hadron Collider simulations is being studied. First, the data is thoroughly explored using visual aids. After that, several basic Machine Learning methods are applied on preprocessed data. Results are evaluated using cross-validation. Model overview is given for each considered algorithm and the best model is chosen.

## I. INTRODUCTION

Claim: it is possible to use simple methods for this dataset

- 1) What is the data (Simulation from LHC, details from physics)
- 2) What are we trying to do? Get the best classification score
- 3) Overview of data, diagrams of features, feature selection, feature augmentation
- 4) Methods and their choice (Linear regression, logistic regression, ridge regression) because of simplicity

## II. MODELS AND METHODS

- 1) Least squares. Problem: missing data, overfit
- 2) Mean imputation. Problem: overfit, meaninglessness for some features
- 3) Feature binarization, add new feature 'feature missing', add squares for features for mass
- 4) Ridge regression using k-fold. Problem: low accuracy (?)
- 5) Logistic regression
- 6) Nearest neighbors?

## III. RESULTS

Shows that accuracy is good enough meaning that model selection was good

## IV. DISCUSSION

State that we can improve the accuracy by using non-linear classifiers?

## V. SUMMARY

We have shown that it is possible to detect the Higgs boson using linear methods and feature augmentation.