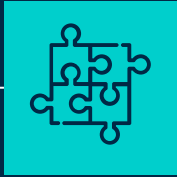# HPV preliminary diagnosis through complete blood sample test, a machine learning approach

Alejandro López Vázquez
Tecnológico de Monterrey
School of Engineering and Science

# ABSTRACT

1. Cervical Cancer is a deadly desease which causes millions of deads yearly world wide. The main goal for this project was to develop a model which can pre-diagnose Human Papiloma Virus (HPV) with certain accuracy bases on a complete blood laboratory exam.
2. Key words:
   - Machine learning.
   - HPV.
   - Early detection.
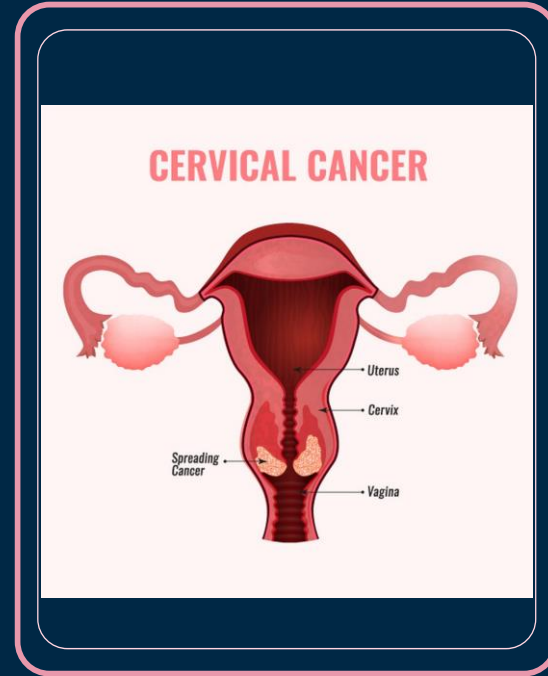   - PCR.
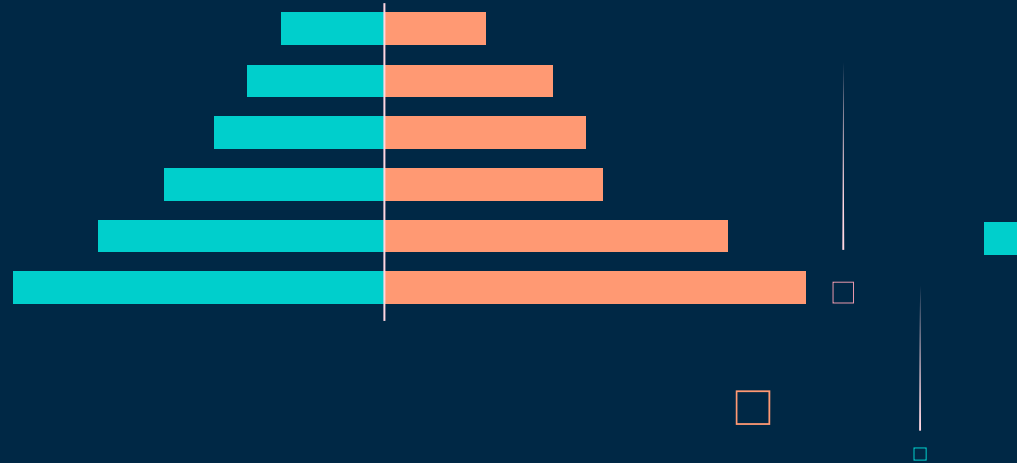   - Neural Network

# TABLE OF CONTENTS

# TABLE OF CONTENTS

# INTRODUCTION

Cervical cancer is the third main cause of dead for women in Mexico. Cervical cancer is caused by some of the many HPV existing types and each can be more or less aggressive, types 16 and 18 being the most dangerous.
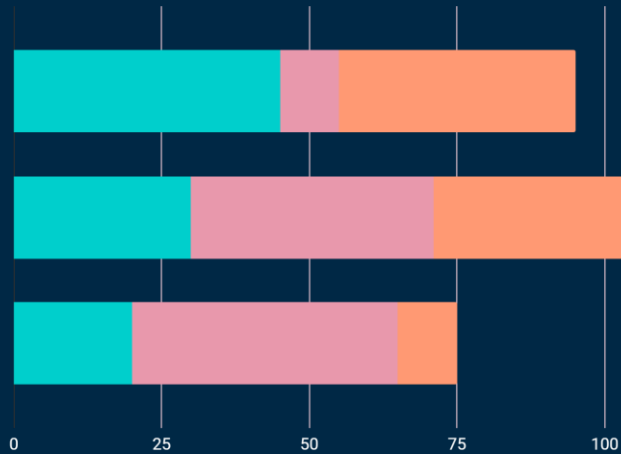


CERVICAL CANCER

# The research questions

1.- Can HPV be dectected through the use of machine learnig with data previously obtained of the patient?

# The research questions



1.- Can HPV be pre-diagnosed based on a complete blood sample test?

# The Data

## Source

The data was obtained from the UU.SS CDC from a the page of the National Health and Nutrition Examination Survey (NHANES).  And we obtained the data specifically of labortory data from the Complete Blood Count with 5-part Differential - Whole Blood and Human Papillomavirus (HPV) DNA - Vaginal Swab from year 2003 to 2016

# Complete Blood Count with 5-part Differential – Whole Blood

- SEQN – Respondent sequence number
- LBXWBCSI – White blood cell count (1000 cells/uL)
- LBXLYPCT – Lymphocyte percent (%)
- LBXMOPCT – Monocyte percent (%)
- LBXNEPCT – Segmented neutrophils percent (%)
- LBXEOPCT – Eosinophils percent (%)
- LBXBAPCT – Basophils percent (%)
- LBDLYMNO – Lymphocyte number (1000 cells/uL)
- LBDMONO – Monocyte number (1000 cells/uL)
- LBDNENO – Segmented neutrophils num (1000 cell/uL)

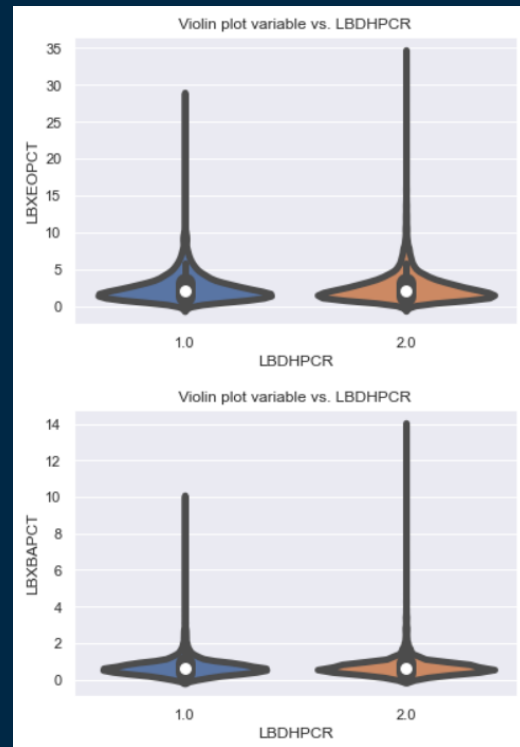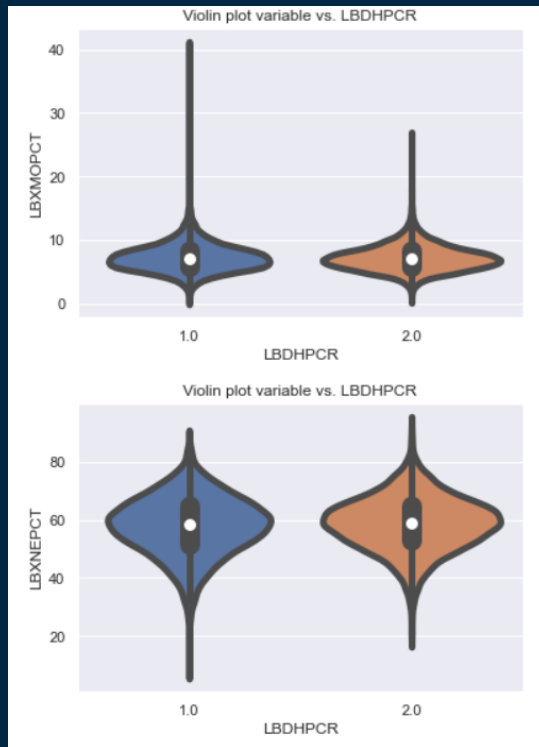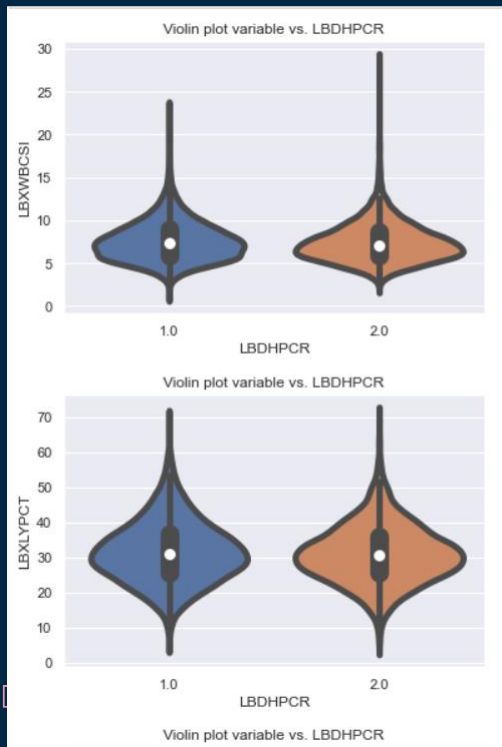# Complete Blood Count with 5-part Differential – Whole Blood

- LBDEONO – Eosinophils number (1000 cells/uL)
- LBDBANO – Basophils number (1000 cells/uL)
- LBXRBCSI – Red blood cell count (million cells/uL)
- LBXHGB – Hemoglobin (g/dL)
- LBXHCT – Hematocrit (%)
- LBXMCVSI – Mean cell volume (fL)
- LBXMCHSI – Mean cell hemoglobin (pg)
- LBXMC – Mean Cell Hgb Conc. (g/dL)
- LBXRDW – Red cell distribution width (%)
- LBXPLTSI – Platelet count (1000 cells/uL)
- LBXMPSI – Mean platelet volume (fL)
- LBXNRBC – Nucleated red blood cells

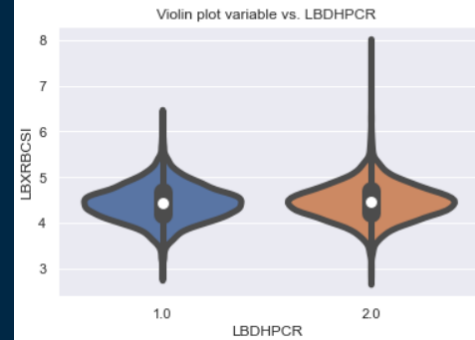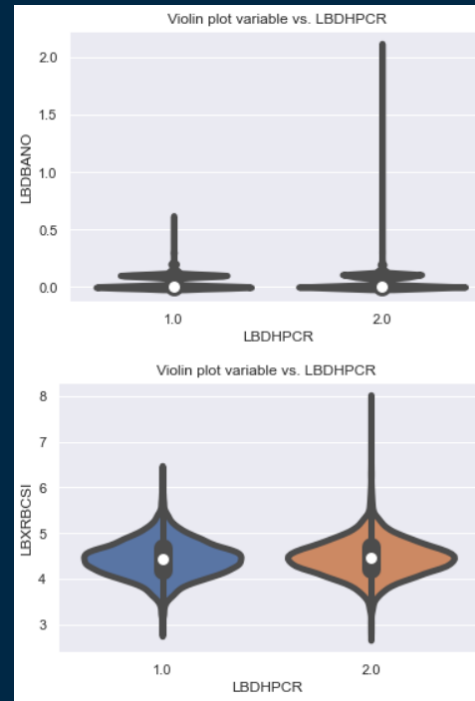# Human Papillomavirus (HPV) DNA – Vaginal Swab: **Roche Linear Array**
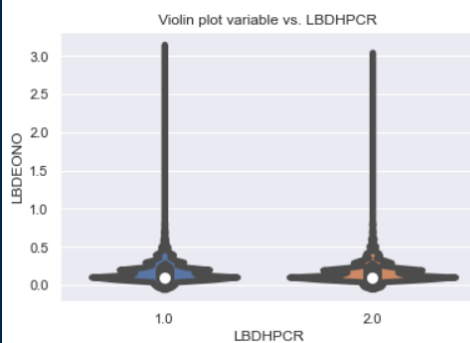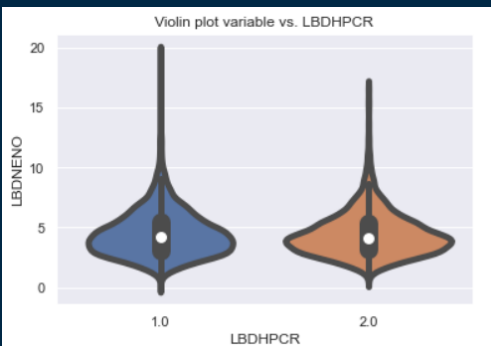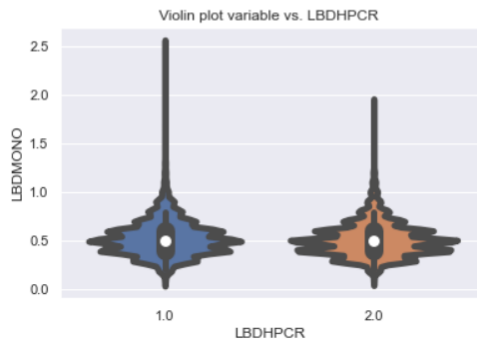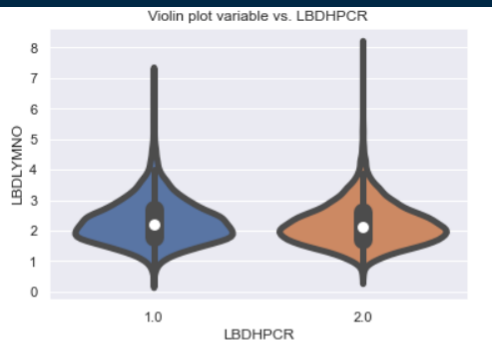
- <u>SEQN – Respondent sequence number</u>

- <u>LBDRPCR - Roche HPV linear array summary result</u>

LBDRPCR  - The HPV PCR Summary variable indicates if at least one type is positive (LBDRPCR=1); the sample is negative (LBDRPCR=2); the sample is inadequate (LBDRPCR=3); or the sample is missing (LBDRPCR=.)
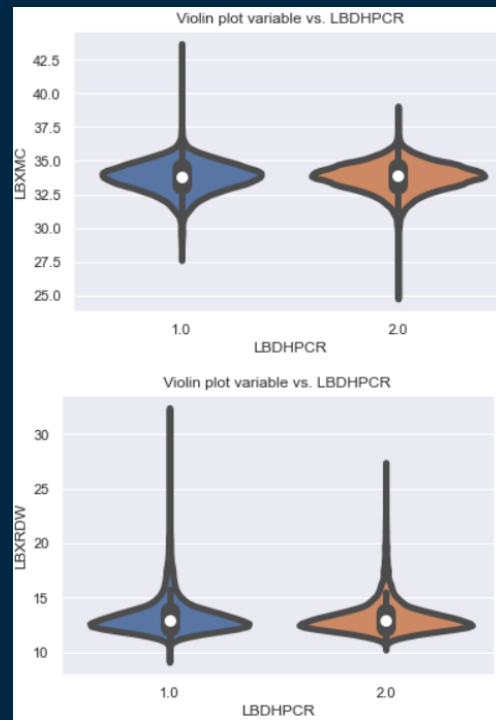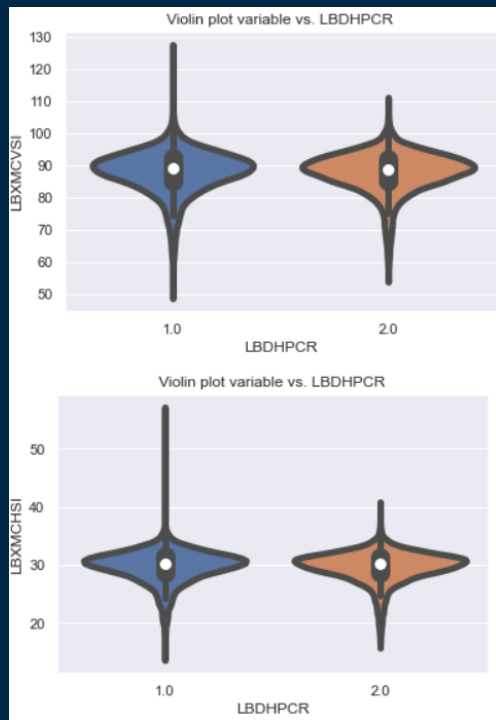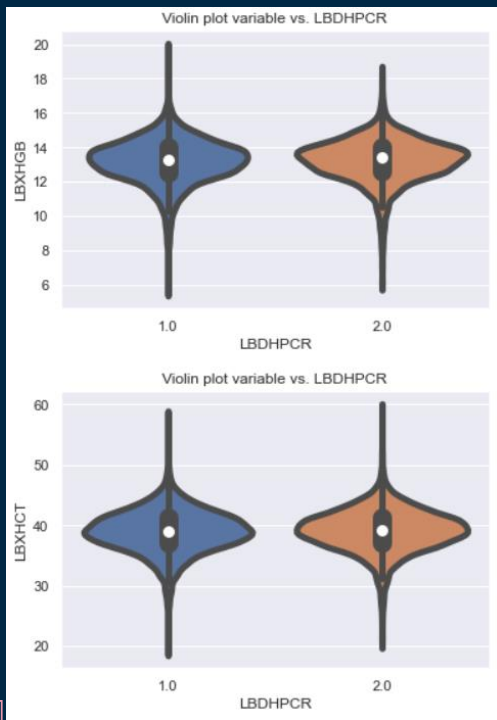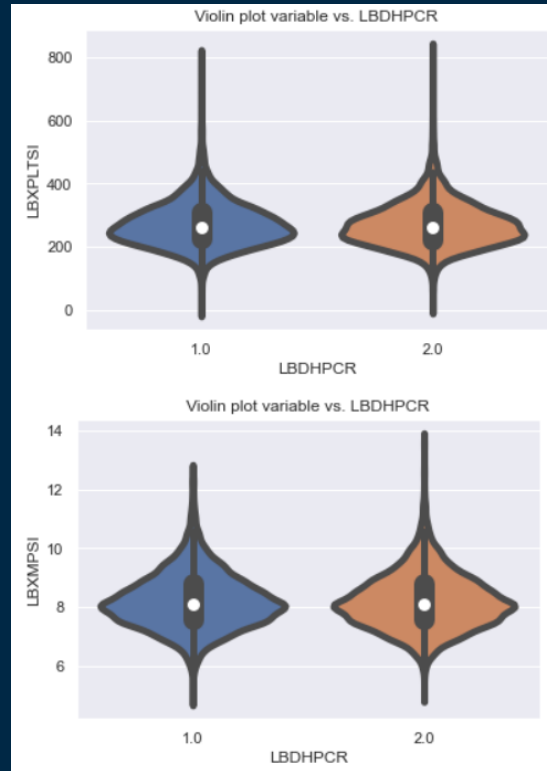
# Data Understanding
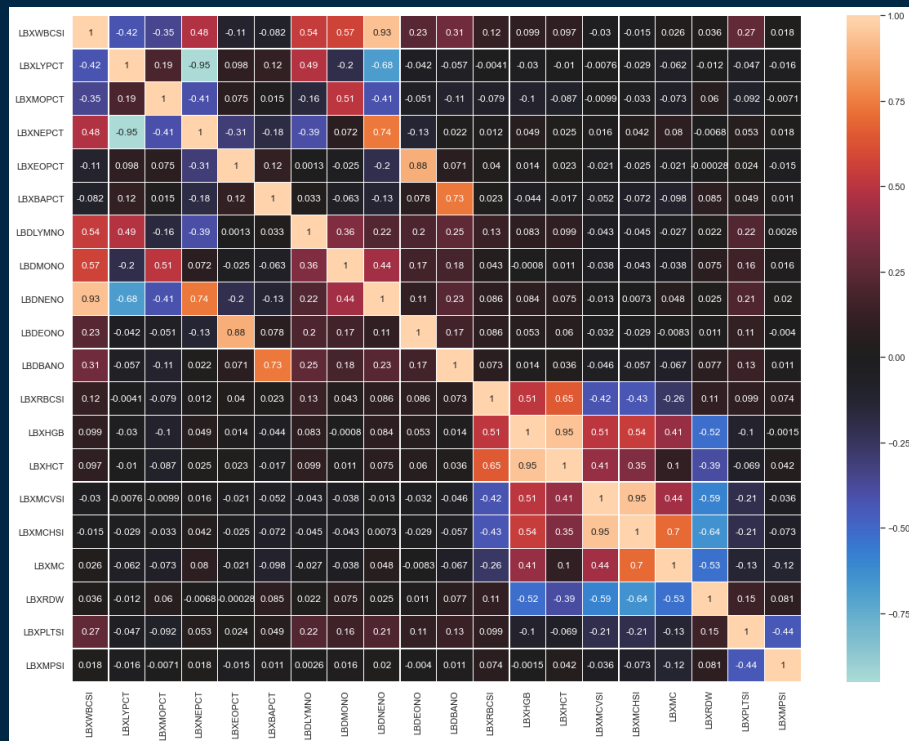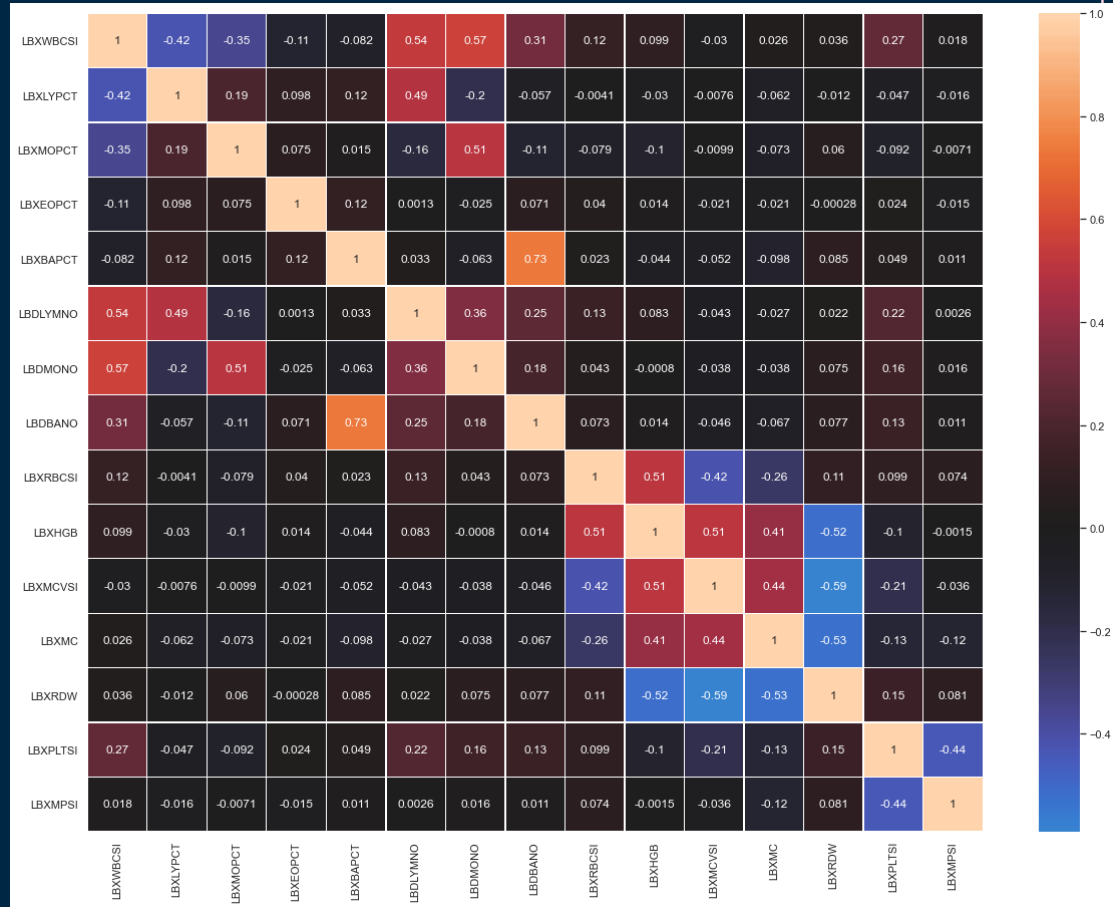
# Data Understanding

# Data Understanding

# Data Understanding

# Correlation Matrix Removed Variables

- LBDEONO
- LBDNENO
- LBXNEPCT
- LBXHCT
- LBXMCHSI

# Data preparation

## Variable normalization

All the variables are continuos so in order to prepare the data a normalization conversion was performed in oorder to place them between the values of 0 and 1

## Train and test splitting

The data was divided into target and predictors, and for both of them a division for a training an a test set was performed with a 75% and a 25% of the data respectively
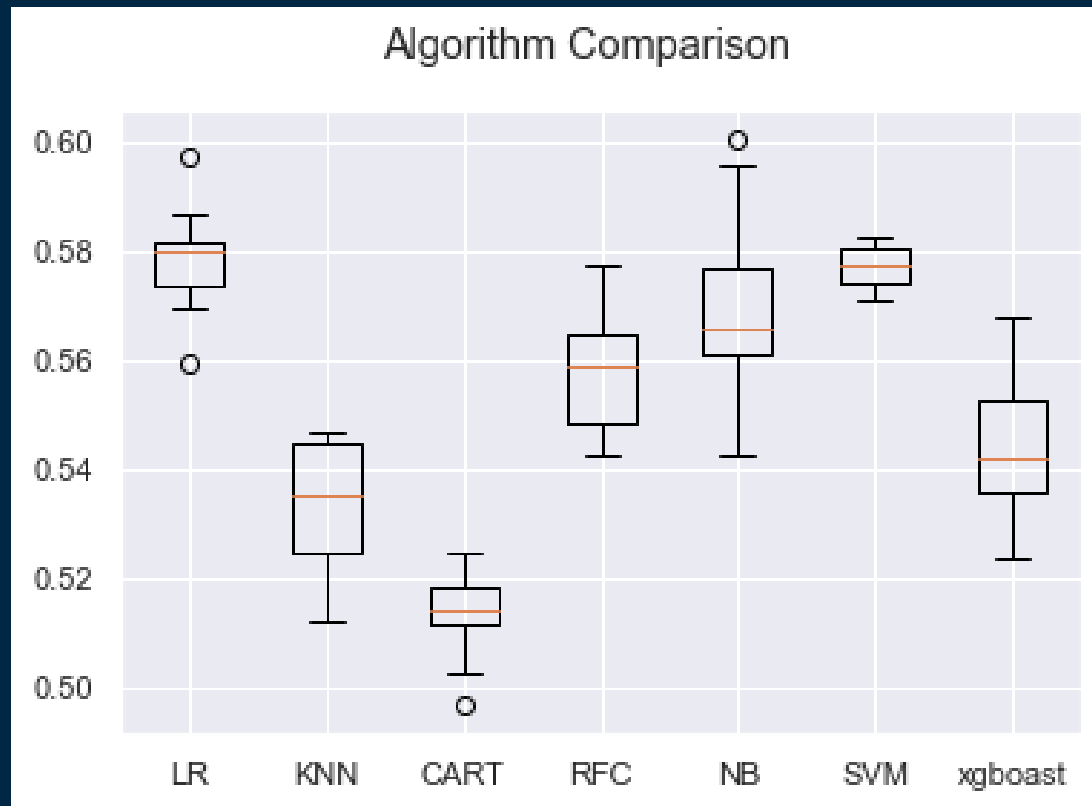
# Models

- Decision Tree
- Logistic Regression
- Naive Bayes
- K neares Neighbors
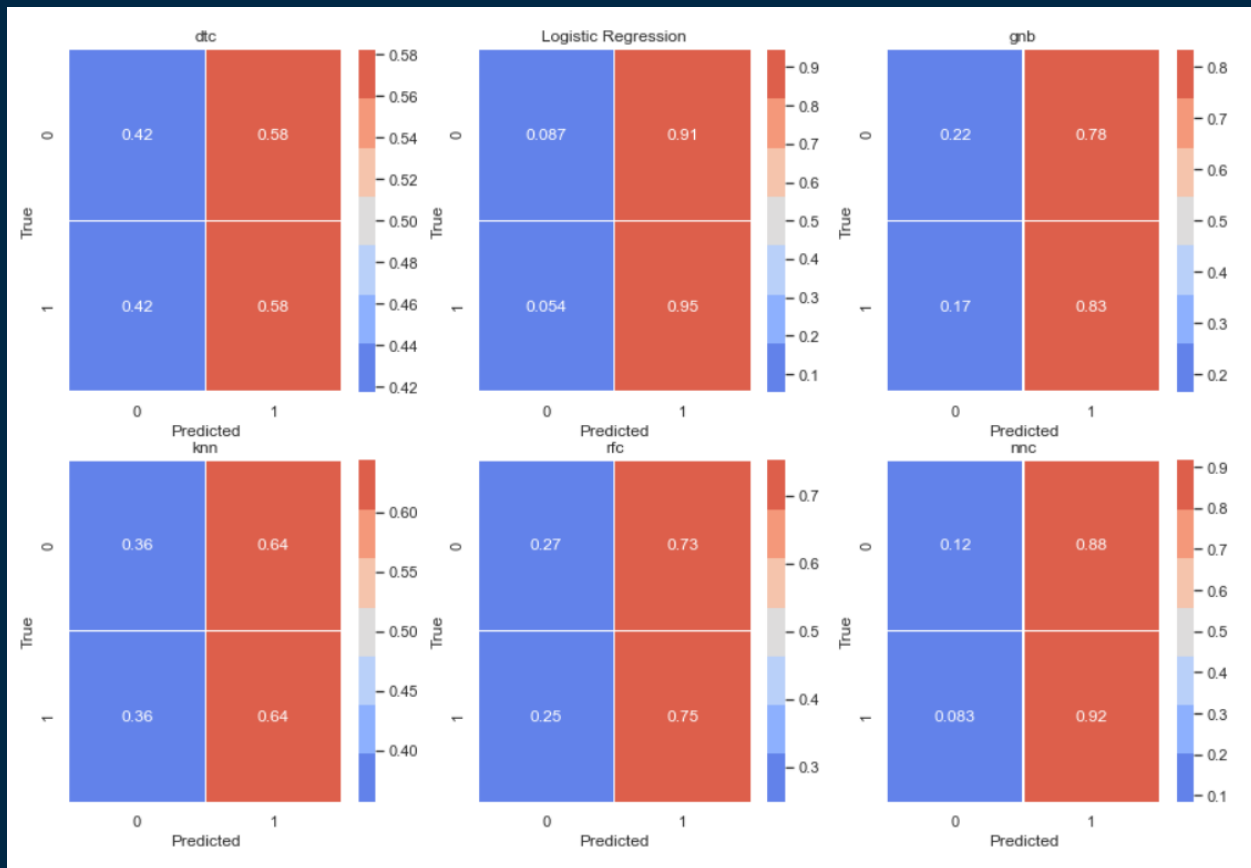- Random Forest
- Multi-Layer Perception

# Results

| | Algorithms | Scores |
|---|---|---|
| 0 | Decision Tree | 0.516331 |
| 1 | Logistic Regression | 0.590234 |
| 2 | Naive Bayes | 0.579677 |
| 3 | K Neighbors Classifier | 0.525239 |
| 4 | Random Ferest | 0.558891 |
| 5 | Neural Network | 0.592874 |
| 6 | Xgboost Classifier | 0.547344 |

| | Algorithm | Accuracy Mean | Accuracy |
|---|---|---|---|
| 0 | LR | 0.578501 | 0.009656 |
| 1 | KNN | 0.532959 | 0.012432 |
| 2 | CART | 0.513737 | 0.008305 |
| 3 | RFC | 0.557794 | 0.010380 |
| 4 | NB | 0.569673 | 0.017123 |
| 5 | SVM | 0.577262 | 0.003621 |
| 6 | xgboast | 0.543932 | 0.013731 |

# Comparisson Boxplot
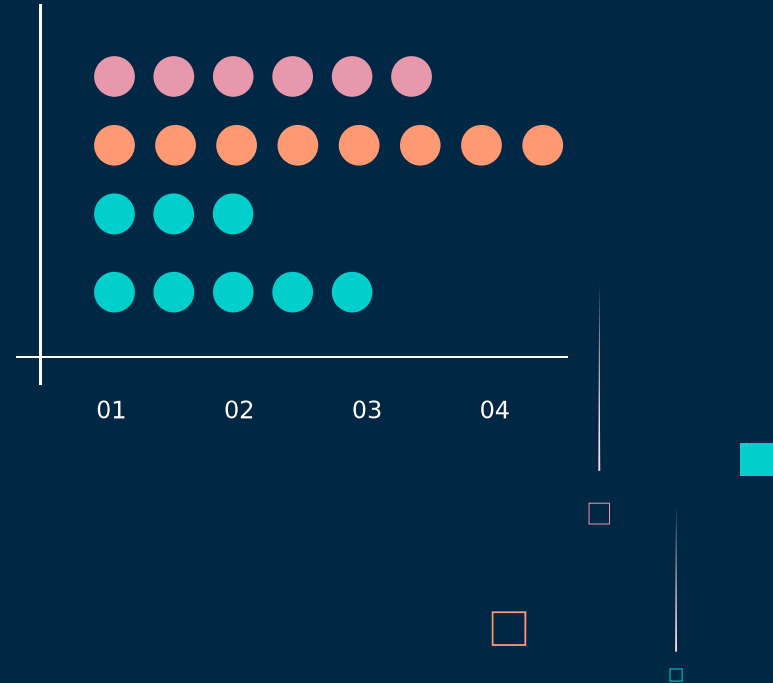


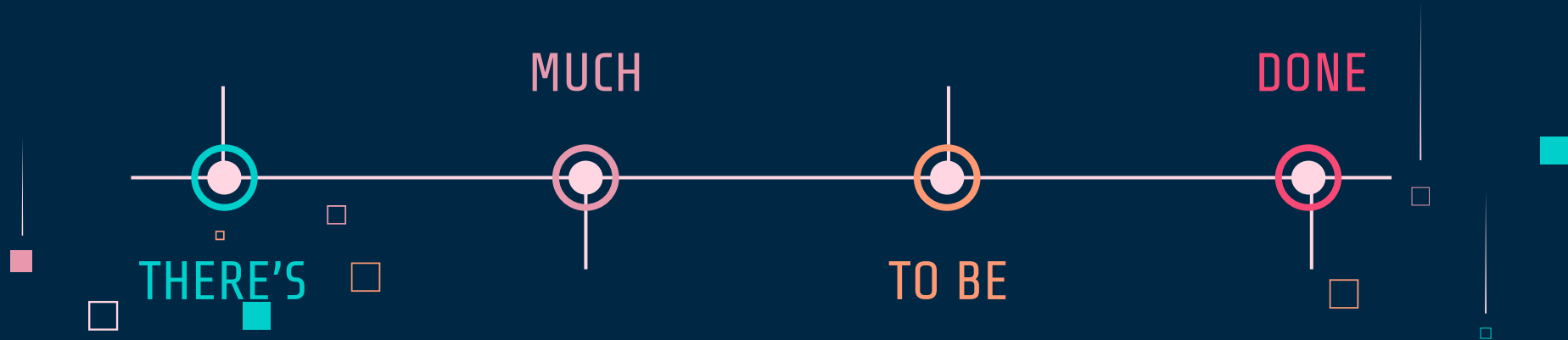Algorithm Comparison

# Correlation matrixes

# Discussion

- The propossed models cannot preddict with Good accuaricy the presence of HPV with the complete blood laboratory results.

# Conclusion

- For future work:
  - There is still a lot of information in the NHANES page, including different laboratories or even quesionares filed by the patients
  - This model can be used also for other diseases uncluded in the same data set.

MUCH

DONE

THERE'S

TO BE

# ANY QUESTIONS?