

TABLE OF CONTENTS

4444

Our team

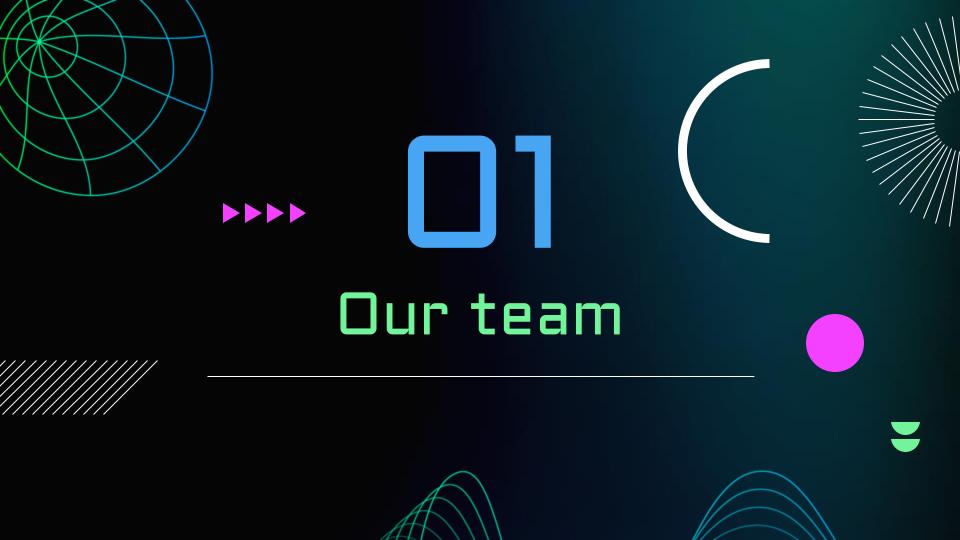
Detailed analysis

O2 Introduction

05 Prediction

O3 Exploratory analysis

06 Conclusion



DUR TEAM













Le Thi Truc Linh



Qinghua Ye



Julian Oliveros Forero











About the Data set

This dataset is a collection of basic health biological signal data.

The goal is to determine the presence or absence of smoking through bio-signals. It will give us an alert for the condition of our body.







Dataset







Shape 55,692 rows, 27 columns





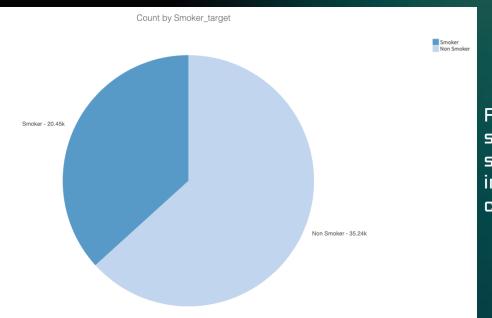
Target Variable
Smoking





Source of DatasetKaggle

Target Variable Analysis



From the pie chart, we can get the statistics that there are 37% of smokers and 63% of non-smokers in the dataset according to the calculation.





Features in the data set

- gender
- age: 5-years gap
- height(cm)
- weight(kg)
- waist(cm): Waist circumference length
- eyesight(left)
- eyesight(right)
- systolic : Blood pressurerelaxation : Blood pressure
- fasting blood sugar
- Cholesterol : total

- Cholesterol: total
- triglyceride
- hemoglobin
- Urine protein
- serum creatinine
- AST : glutamic oxaloacetic transaminase type
- ALT : glutamic oxaloacetic transaminase type
- Gtp : γ-GTP
- dental caries
- tartar: tartar status



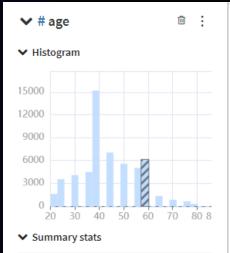




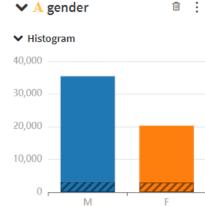
Analysis Feature

Let's look at some of the private data in the dataset.

We can observe that 64% of the data set is made up of men, whereas just 34% of it is made up of women. They may want to boost the male sample size because it's probable that the majority of males smoke frequently. Additionally, we can observe that our dataset's average age is 44 years old, with the oldest person being 85 years old and the youngest person being 20.



N values	55692
N distinct	14
N finite	55692
Mean	44.182916756
Median	40
Std Dev	12.071417567
Min	20
Max	85



Summary stats

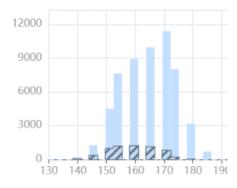
N values	55692
N distinct	2
Mode	М
N empty	0

Frequency table

М		64%	35401
F		36%	20291
	N distinct		2



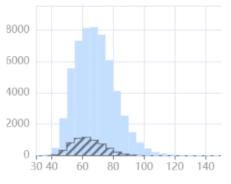




Summary stats

N values	55692
N distinct	13
N finite	55692
Mean	164.64932127
Median	165
Std Dev	9.1945969126
Min	130
Max	190

ŵ



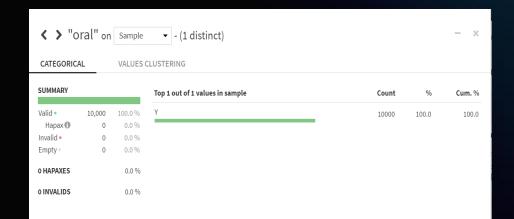
Summary stats

N values	55692
N distinct	22
N finite	55692
Mean	65.864935718
Median	65
Std Dev	12.820305713
Min	30
Max	135

Analysis Feature

Based on the data from the chart below, we can see that the people surveyed in the data set have an average height mainly between 1m60 and 1m70, with an average height of 1m65. Most people's weight is between 60 and 80 kg, and the average weight is 65 kg.





Analysis Feature

After exploring all features of the dataset we see that The Oral column would not be useful for our model because it just has only one value in the sample. This value can't distribute more useful information to the algorithm. In other words, this feature hasn't predictive power. Hence, We decided to drop it out of the dataset.



Correlation matrix

-8.8																								
waist(cm)	1.000	0.042	0.461	0.469	0.287	0.252	0.034	0.259	0.289	0.808	0.041	0.037	0.074	0.388	-0.040	0.099	0.403	0.003	0.319	0.401	-0.396	0.027	0.021	0.228
dental caries	0.042	1.000	0.034	0.061	0.038	0.004	0.002	0.012	0.030	0.072	0.025	0.017	-0.003	0.081	-0.122	-0.004	0.074	0.001	0.029	0.031	-0.030	-0.016	-0.016	0.104
ALT	0.461	0.034	1.000	0.622	0.233	0.731	0.032	0.196	0.211	0.451	0.059	0.054	0.108	0.267	-0.073	0.089	0.420	0.001	0.208	0.357	-0.269	-0.011	-0.015	0.206
Gtp	0.469	0.061	0.622	1.000	0.288	0.472	0.040	0.279	0.280	0.439	0.044	0.041	0.141	0.297	-0.022	0.064	0.449	0.008	0.264	0.461	-0.226	0.009	0.011	0.372
serum creat	0.287	0.038	0.233	0.288	1.000	0.153	0.018	0.074	0.097	0.421	0.107	0.106	0.014	0.477	-0.177	0.044	0.491	0.002	0.078	0.160	-0.212	-0.000	0.003	0.268
AST	0.252	0.004	0.731	0.472	0.153	1.000	0.035	0.118	0.163	0.198	-0.026	-0.028	0.103	0.076	0.114	0.061	0.231	-0.001	0.173	0.196	-0.085	0.037	0.038	0.100
Urine protein	0.034	0.002	0.032	0.040	0.018	0.035	1.000	0.040	0.031	0.026	-0.014	-0.014	-0.006	0.010	0.012	-0.008	0.031	0.001	0.027	0.018	-0.019	0.011	0.016	0.011
fasting bloo	0.259	0.012	0.196	0.279	0.074	0.118	0.040	1.000	0.194	0.180	-0.061	-0.061	0.054	0.033	0.210	0.011	0.113	0.001	0.224	0.269	-0.136	0.042	0.045	0.104
relaxation	0.289	0.030	0.211	0.280	0.097	0.163	0.031	0.194	1.000	0.268	0.007	0.005	0.095	0.120	0.052	0.057	0.234	0.006	0.741	0.231	-0.102	0.006	0.001	0.110
weight(kg)	0.808	0.072	0.451	0.439	0.421	0.198	0.026	0.180	0.268	1.000	0.179	0.176	0.027	0.697	-0.341	0.059	0.542	0.005	0.270	0.351	-0.381	-0.049	-0.053	0.317
eyesight(rig	0.041	0.025	0.059	0.044	0.107	-0.026	-0.014	-0.061	0.007	0.179	1.000	0.695	-0.005	0.247	-0.333	-0.003	0.164	0.003	-0.035	0.021	-0.025	-0.090	-0.100	0.102
eyesight(left)	0.037	0.017	0.054	0.041	0.106	-0.028	-0.014	-0.061	0.005	0.176	0.695	1.000	-0.005	0.242	-0.337	-0.004	0.160	0.004	-0.039	0.021	-0.019	-0.091	-0.098	0.094
Cholesterol	0.074	-0.003	0.108	0.141	0.014	0.103	-0.006	0.054	0.095	0.027	-0.005	-0.005	1.000	-0.078	0.072	0.890	0.050	-0.002	0.057	0.253	0.156	-0.023	-0.020	-0.026 ⁰
height(cm)	0.388	0.081	0.267	0.297	0.477	0.076	0.010	0.033	0.120	0.697	0.247	0.242	-0.078	1.000	-0.499	-0.051	0.584	0.006	0.096	0.169	-0.230	-0.074	-0.074	0.403
age	-0.040	-0.122	-0.073	-0.022	-0.177	0.114	0.012	0.210	0.052	-0.341	-0,333	-0.337	0.072	-0.499	1.000	0.061	-0.316	-0.001	0.113	0.027	0.019	0.176	0.180	-0.164 .
LDL	0.099	-0.004	0.089	0.064	0.044	0.061	-0.008	0.011	0.057	0.059	-0.003	-0.004	0.890	-0.051	0.061	1.000	0.057	-0.001	0.023	0.093	-0.058	-0.017	-0.015	-0.052
hemoglobin	0.403	0.074	0.420	0.449	0.491	0.231	0.031	0.113	0.234	0.542	0.164	0.160	0.050	0.584	-0.316	0.057	1.000	0.008	0.193	0.298	-0.275	-0.033	-0.039	0.417
ID	0.003	0.001	0.001	0.008	0.002	-0.001	0.001	0.001	0.006	0.005	0.003	0.004	-0.002	0.006	-0.001	-0.001	0.008	1.000	0.003	0.001	-0.005	0.003	-0.005	0.011
systolic	0.319	0.029	0.208	0.264	0.078	0.173	0.027	0.224	0.741	0.270	-0.035	-0.039	0.057	0.096	0.113	0.023	0.193	0.003	1.000	0.218	-0.102	0.052	0.045	0.077
triglyceride	0.401	0.031	0.357	0.461	0.160	0.196	0.018	0.269	0.231	0.351	0.021	0.021	0.253	0.169	0.027	0.093	0.298	0.001	0.218	1.000	-0.470	0.006	0.002	0.257
HDL	-0.396	-0.030	-0.269	-0.226	-0.212	-0.085	-0.019	-0.136	-0.102	-0.381	-0.025	-0.019	0.156	-0.230	0.019	-0.058	-0.275	-0.005	-0.102	-0.470	1.000	-0.020	-0.017	-0.195
hearing(left)	0.027	-0.016	-0.011	0.009	-0.000	0.037	0.011	0.042	0.006	-0.049	-0.090	-0.091	-0.023	-0.074	0.176	-0.017	-0.033	0.003	0.052	0.006	-0.020	1.000	0.510	-0.023
hearing(right)	0.021	-0.016	-0.015	0.011	0.003	0.038	0.016	0.045	0.001	-0.053	-0.100	-0.098	-0.020	-0.074	0.180	-0.015	-0.039	-0.005	0.045	0.002	-0.017	0.510	1.000	-0.019
smoking	0.228	0.104	0.206	0.372	0.268	0.100	0.011	0.104	0.110	0.317	0.102	0.094	-0.026	0.403	-0.164	-0.052	0.417	0.011	0.077	0.257	-0.195	-0.023	-0.019	1.000
	waist(cm) d	ental c	ALT	Gtp se	erum cr	AST U	rine pr fa:	sting b re	elaxation w	reight(kg) ey	esight ey	esight Ch	noleste he	eight(cm)	age	LDL he	emoglo	ID	systolic tr	iglyceri	HDL h	earing(l he	aring(s	smoking

Correlation Matrix Conclusion

Based on Correlation matrix we can see that:

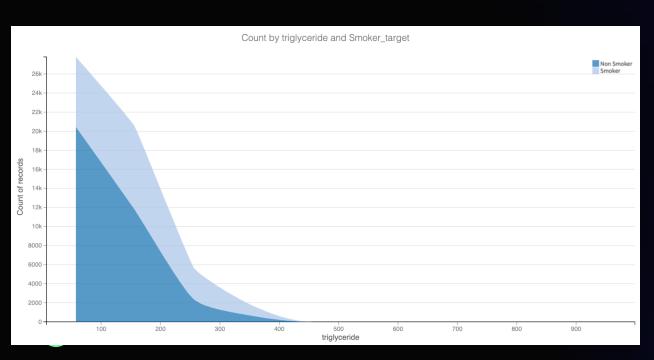
Some features that are highly relevant to the target (Smoking) are respectively hemoglobin (0.417), height(0.403), Gtp(0.372), weight(0.317), serum creatinine(0.268), triglyceride(0.257), waist(0.228) and ALT(0.206). It means that this feature might have big predictive power to the target.

Some features that unrelated to the target (Smoking) are respectively Cholesterol (-0.026), Age(-0.164), LDL(-0.052), HDL(-0.195), hearing left (-0.023), hearing right (-0.019).

The correlation matrix works very well with the continuous variable, but the value type of our target is binary value. It's one of the categorical variables, so maybe the correlation matrix may not really accurately reflects the correlations of the predictors and target. So let's create some graphs between the predictors and the target in the detailed analysis step and see their right relationship.



Bar Chart Gender Analysis

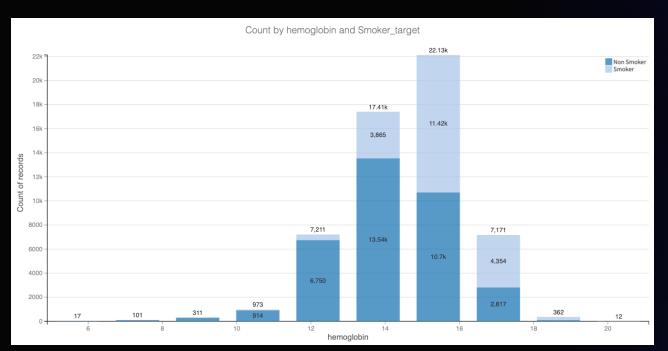


As per the graph majority of people who smoke has low triglyceride level than the people who does not smoke.

So, if we find low triglyceride level we can assume that he might be smoking.



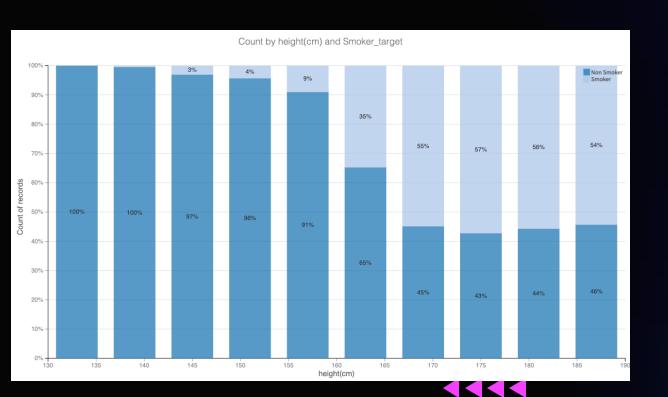
Bar Chart Hemoglobin Analysis



As per the given chart majority of people who smoke their hemoglobin ranges between 15-16 compared to non-smoker hemoglobin ranges between 13-15. So we can assume that if person's hemoglobin level is between 15-16, he might be smoking.



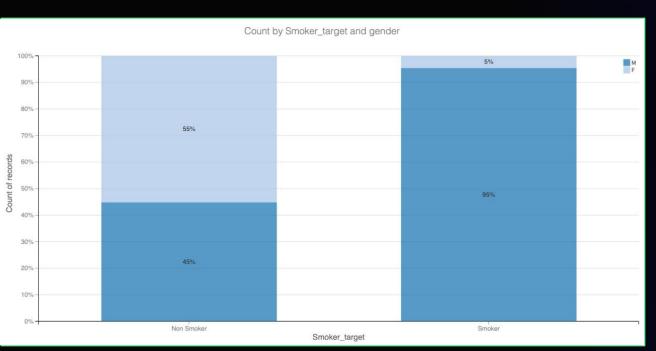
Bar Chart Gender Analysis



Majority of people who has height greater than 165 has high high chances that they smoke and majority of people having height lesser than 160 are non smokers.

So we can assume than person whose height is greater than 165 might be a smoker.

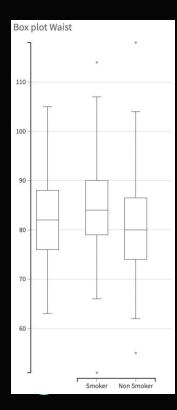
Bar Chart Gender Analysis

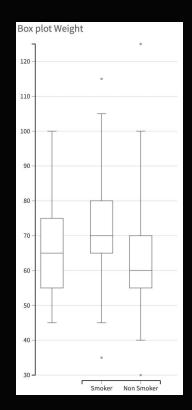


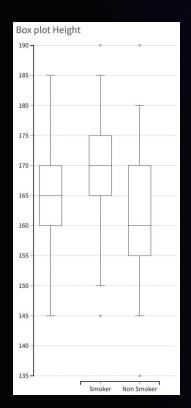
According to the graph, we can see the percentage of the male who doesn't smoke are almost same with that of the female. there are 90% more male smoker than female smoker. Hence, we can say gender has a big influence on people smoking or not.



BOX PLOT Waist, Weight, Height



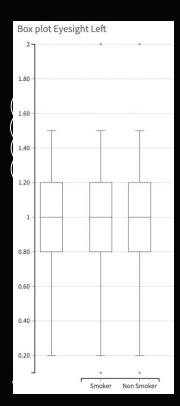


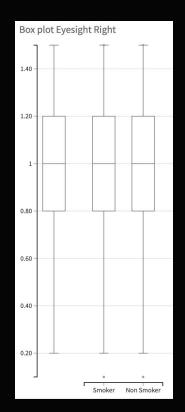


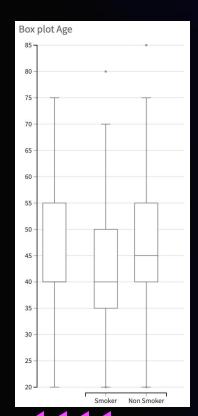
From these box plots, we can see that the variables have high importance in predicting if a person is a smoker since the difference between the values of Waist, Weight, and Height has different average values for smokers and nonsmokers.



BOX PLOT Eyesight Left and Right, Age

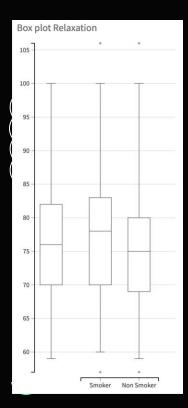


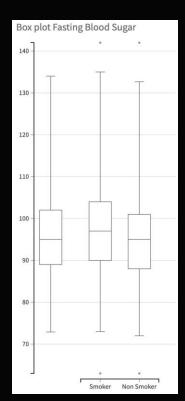


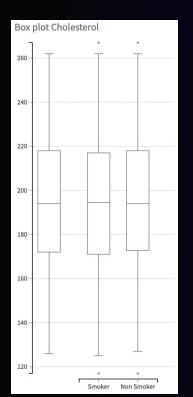


From these box plots, we can see that the variable age has small importance in predicting if a person is a smoker since the average age for smokers is 40 and for non-smokers is 45. On the other hand, for Eyesight left and right we see that these variables have non-relevance in knowing if a person is a smoker since the values for these variables are the same.

BOX PLOT Relaxation, Blood Sugar, Cholesterol



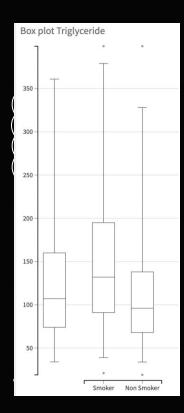


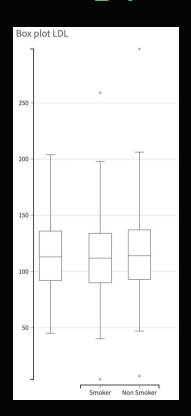


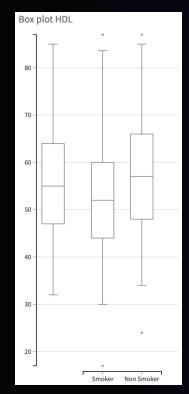
From these box plots, we can see that the variables relaxation and Fasting Blood Sugar have small importance in predicting if a person is a smoker since the average values have a very small difference between smakers and non-smokers. On the other hand, for Cholesterol we see that these variables have non-relevance in knowing if a person is a smoker since the values for these variables are the same for smokers and nonsmokers.



BOX PLOT Triglyceride, LDL, HDL



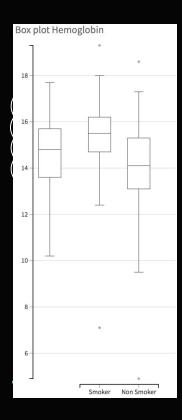


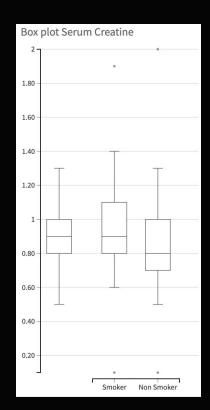


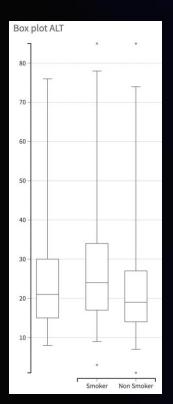
From these box plots, we can see that the variables Triglyceride and HDL have importance in predicting if a person is a smoker since the average values have a small difference between smokers and non-smokers. On the other hand, for LDL we see that these variables have nonrelevance in knowing if a person is a smoker since the values for these variables are the same for smokers and non-smokers.



BOX PLOT Hemoglobin, Serum Creatine, ALT



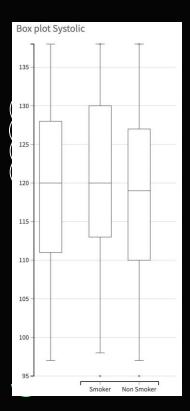


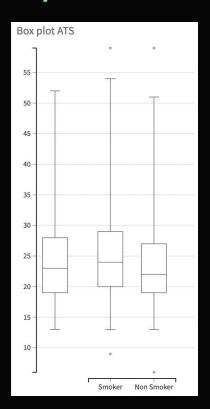


From these box plots, we can see that the variables Hemoglobin have importance in predicting if a person is a smoker since the average values have a big difference between smokers and nonsmokers. On the other hand, for Serum Creatine and ALT we see that these variables have a small difference in knowing if a person is a smoker since the values for these variables only have a little difference for smokers and non-smokers.



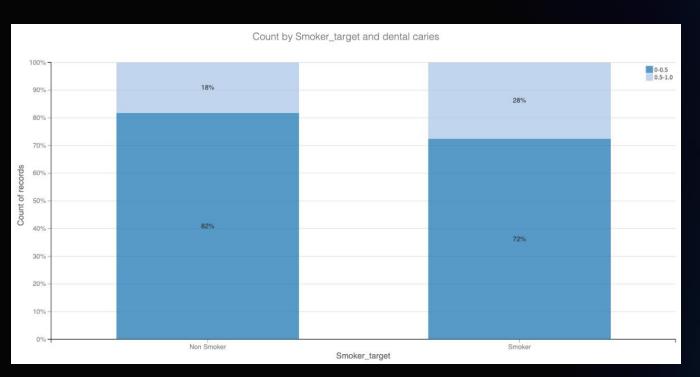
BOX PLOT Systolic, ATS





From the two box plots, we can see there is only a little difference for blood pressure and glutamic oxaloacetic transaminase type among people who smoke or not. So we can say smoking barely influence people's blood pressure and glutamic oxaloacetic transaminase(AST).

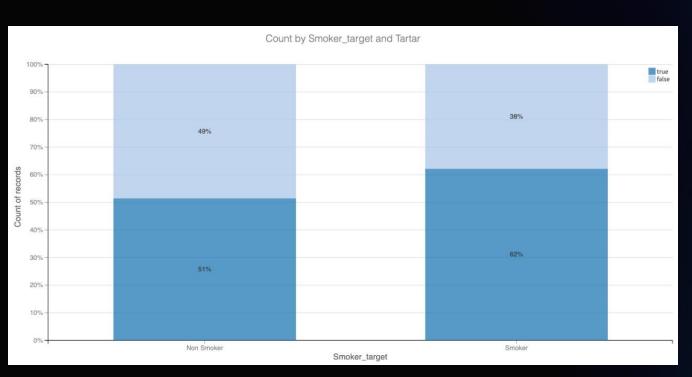
Bar Chart Dental Caries



According to the graph, we can see there is 18% of nonsmokers instead of 28% of smokers who have more dental caries. So, we can say smoking decides slightly if having dental caries or not.



Bar Chart Tartar



According to the graph, we can see that 51% of non-smokers instead of 62% of smokers have tartar on the teeth. Hence, we can say smoking won't cause more tartar issues.











Categorical

- We see that gender has high relevance in knowing if a person is a smoker.
- Dental caries and tartar have small importance in knowing if a person is a smoker.



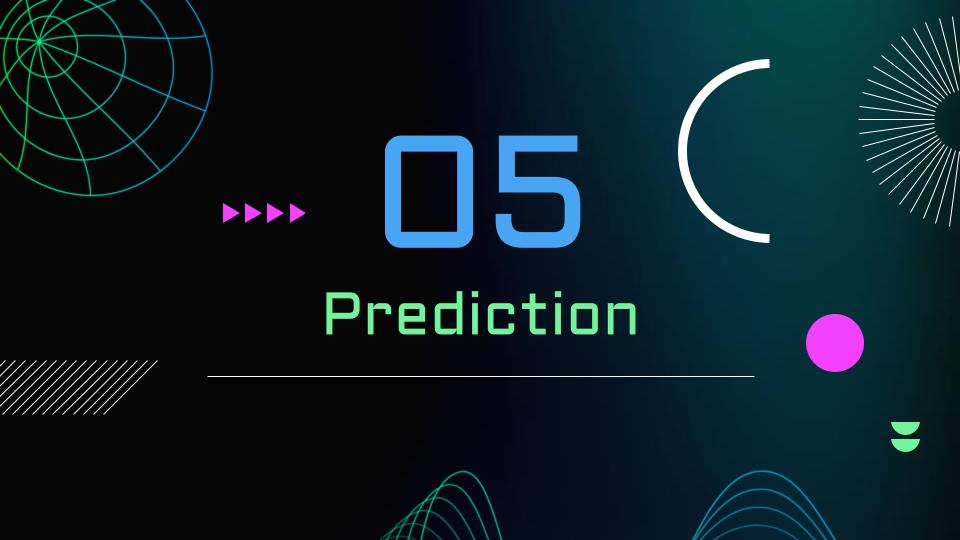
Numerical

The following variables have a small relevance in knowing if a person is a smoker or not.

- Eyesight (left)
- Eyesight (Right)
- Hearing(left)
- Hearing (Right)
- Cholesterol
- LDL







DESING





Algorithm

Since our target variable is binary, we decided to use random Forrest for our prediction model.



Sampling & Splitting

Since we have more people for non-smokers than nonsmokers, we did a rebalance on the data set so that the target values can have a similar amount of values



DESING



Features Handling

After the feature analysis, we decided to train the model with the following features.

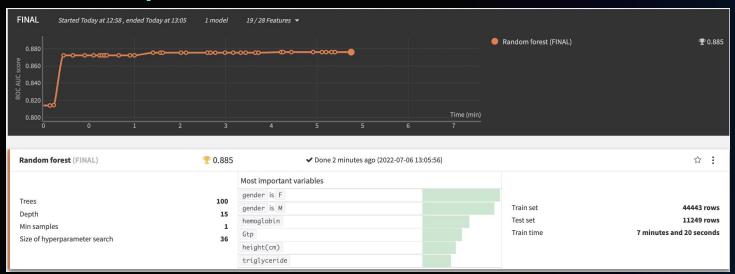
- gender
- age
- height(cm)
- weight(kg)
- waist(cm)

- systolic
- relaxation
- fasting blood sugar
- triglyceride
- hemoglobin

- Urine protein
- serum creatinine
 - AST
- ALT
- Gtp

- dental caries
- tartar

Model performance

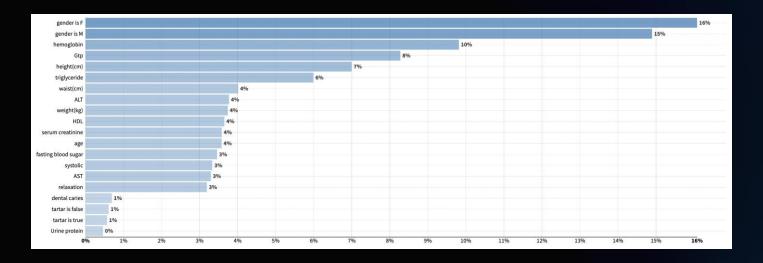


We have achieved 88.5% of ROC with most important features like gender, hemoglobin, height and triglyceride.





Variable Importance



The most important feature which has high impact on model training and its accuracy are gender, hemoglobin, GTP, height and Triglyceride.





Performance metrics

Performance metrics			Thre
Detailed metrics			
Threshold independent		Threshold dependent	
Log loss 🔞	0.4212	Accuracy ②	0.7970
ROC - AUC Score ②	0.8854	Precision ②	0.6871
Calibration loss 🔞	0.0633	Recall ②	0.8266
		F1 Score 🔞	0.7504
		Hamming loss	0.2030
		Cost matrix gain 🔞	0.2634
		Matthews Correlation Coefficient 🔞	0.5889





We have achieved overall accuracy of 79.7%. We majorly focused on recall to improve, because as per the use case if our model predicts that user is non-smoker but actual he was smoker it might cause a problem during medical treatments.

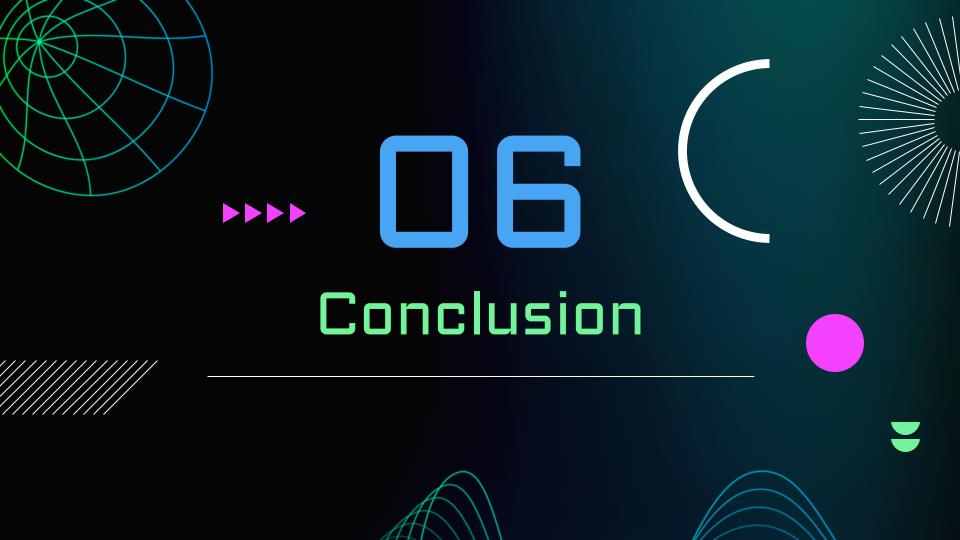
Algorithm Details

Algorithm details						
Algorithm	Randon	n forest cla	ssification	Split quality criterion		Gini
Number of trees	100			Use bootstrap		Yes
Max trees depth	15			Feature sampling strate	egy	auto
Min samples per leaf	1					
Min samples to split Training data	3					
Rows (before preprocess	ing)	44443	Rows (after	r preprocessing)	44443	
Columns (before preprod	essing)	28	Columns (a	after preprocessing)	24	
Matrix type		dense				
Estimated memory usag	е	8.14 MB				

We have used Random Forest Classifier for our model with 100 trees, max depth of 15 and at least 1 sample per leaf, split criterion is Gini and bootstrap enabled.







Conclusion on Analysis and prediction

- Our detailed analysis and model prediction are converging.
- We recommend that if Male smoker whose height is greater than 165 and hemoglobin level is 15-16 can be a sumed as smoker.
- Later using non-body signals and our model we can predict that if person is a smoker or not.
- We majorly concentrated on recall score because our model major use case can be in medical treatment and if my model predict user is non smoker but actually he is a smoker then it might cause problems in medical treatment which even might lead to death.

