

# Exploring the Oncogenic Impact of the **TP53** Gene Variants and Classifying the Pathogenicity



# Quick Recap: What Is the TP53 Gene?

A: Part of the DNA responsible for regulating cell division

B: Contains nucleotides A, T, G, C

C: Present on the 17<sup>th</sup> chromosome

D: None of the above

E: All of the above except D

# Quick Recap: What Is the TP53 Gene?

A: Part of the DNA

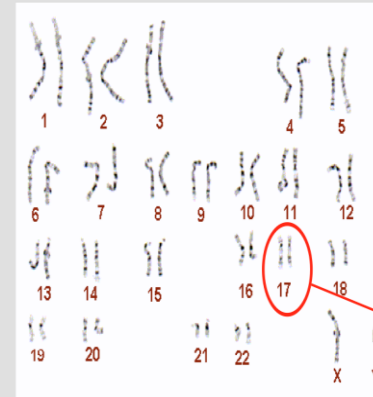
B: Contains nucleotides A, T, G, C

C: Present on the 17<sup>th</sup> chromosome

D: None of the above

E: All of the above except D

## LOCALIZATION OF THE TP53 GENE



17p13.1

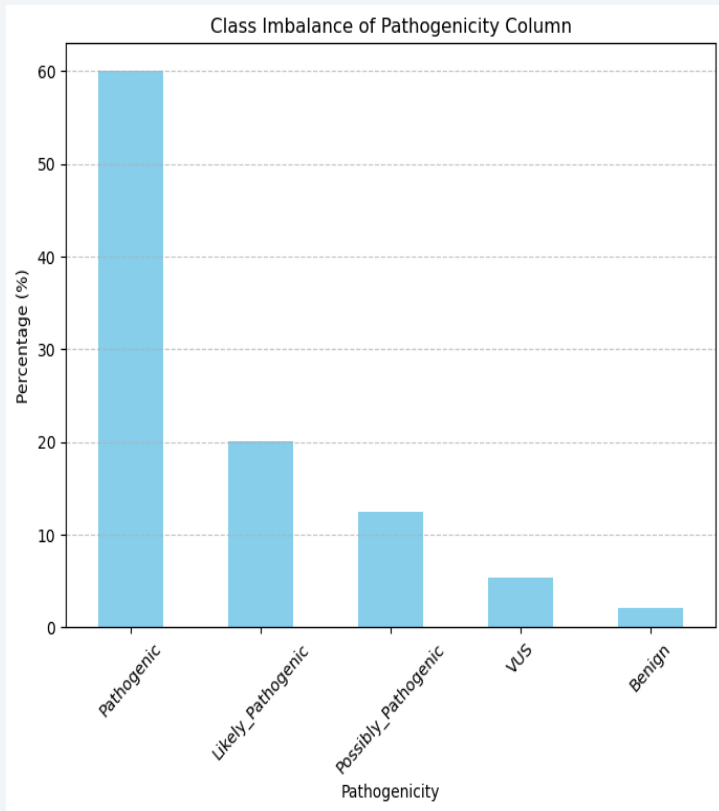
# 01 Data Dictionary

OG shape: 80406, 133

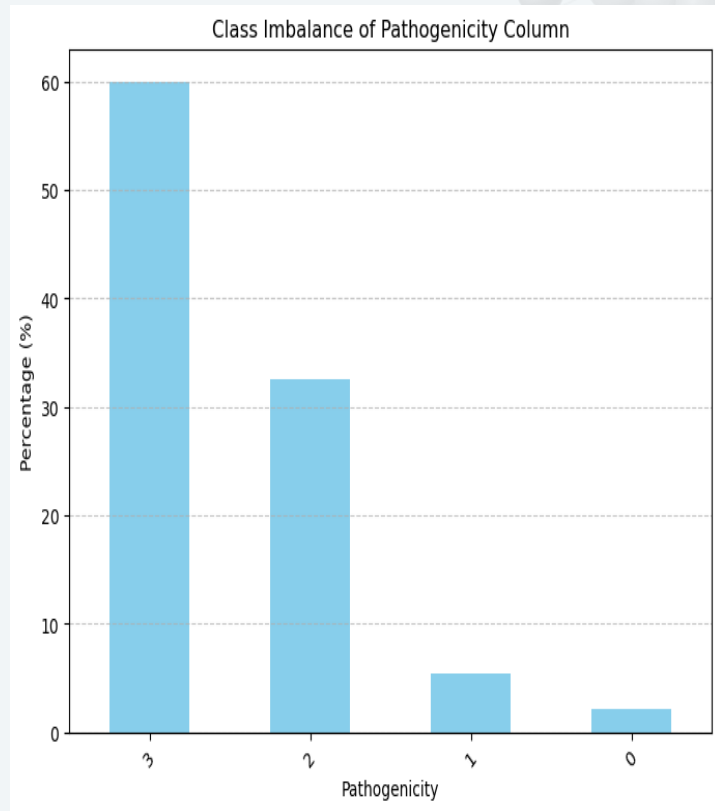
Features	Details
HG38_Start	Mutation start coordinates against the ref genome (int)
Mutant_Codon	<b><u>NNN</u></b> : Sequence of the mutated codon (69 types)
Disease	Name of the disease (11 including unknown/others)
Variant_Type	SNP, DNP, TNP, ONP, INS, DEL
Somatic_Stat	Frequency of the variant (cDNA_nomenclature) found as a somatic event
Target	Details
Pathogenicity	Benign, Unknown, Likely Pathogenic, Pathogenic

AfterEDA: 80337, 78

# EDA

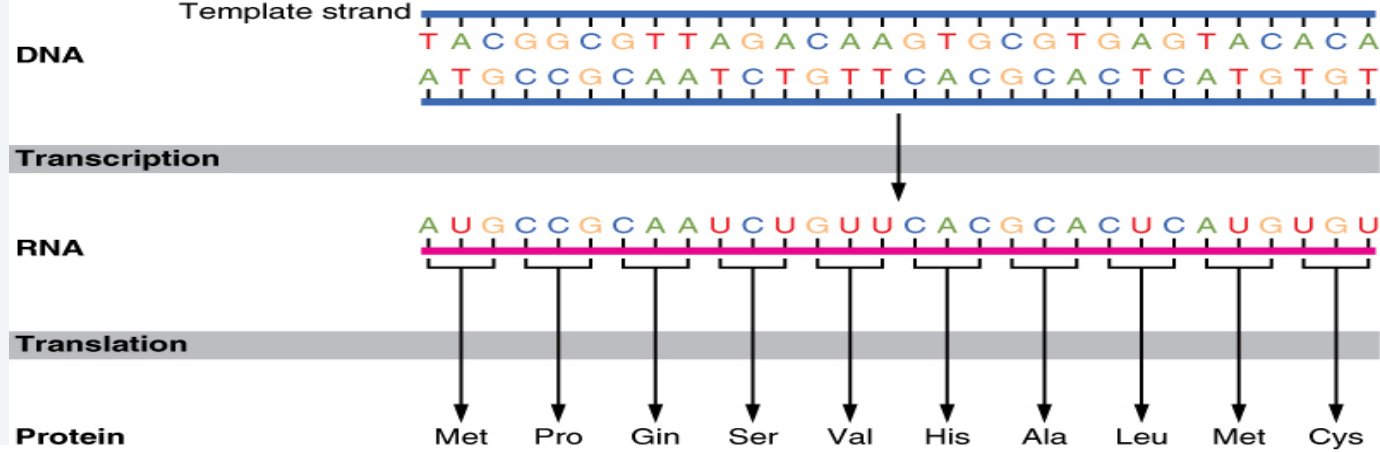


Previously

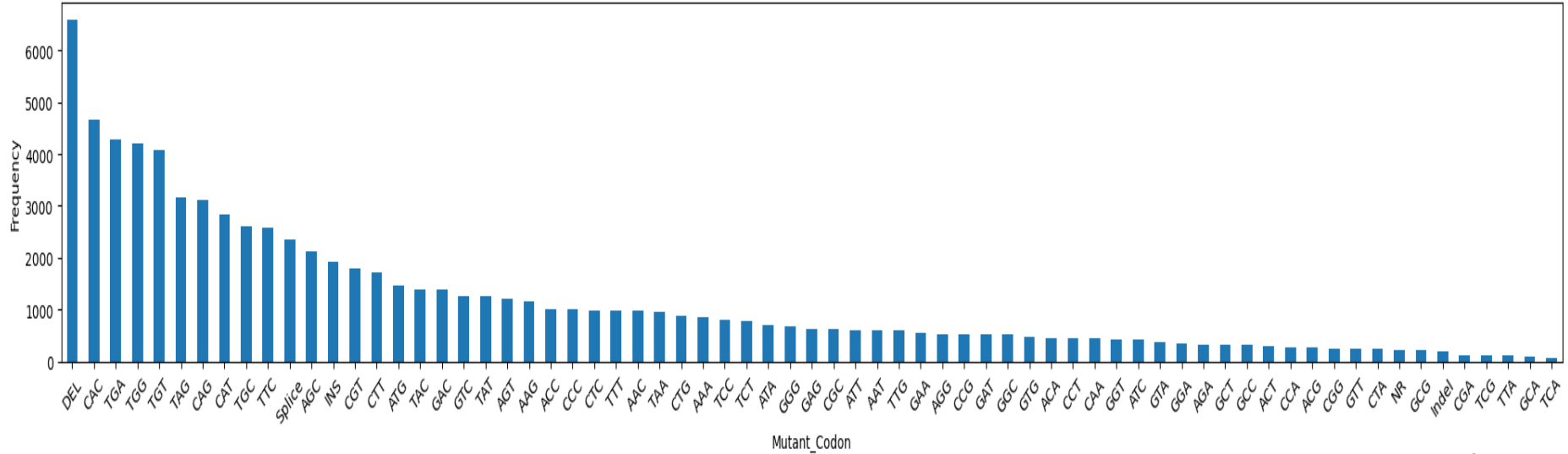


Now

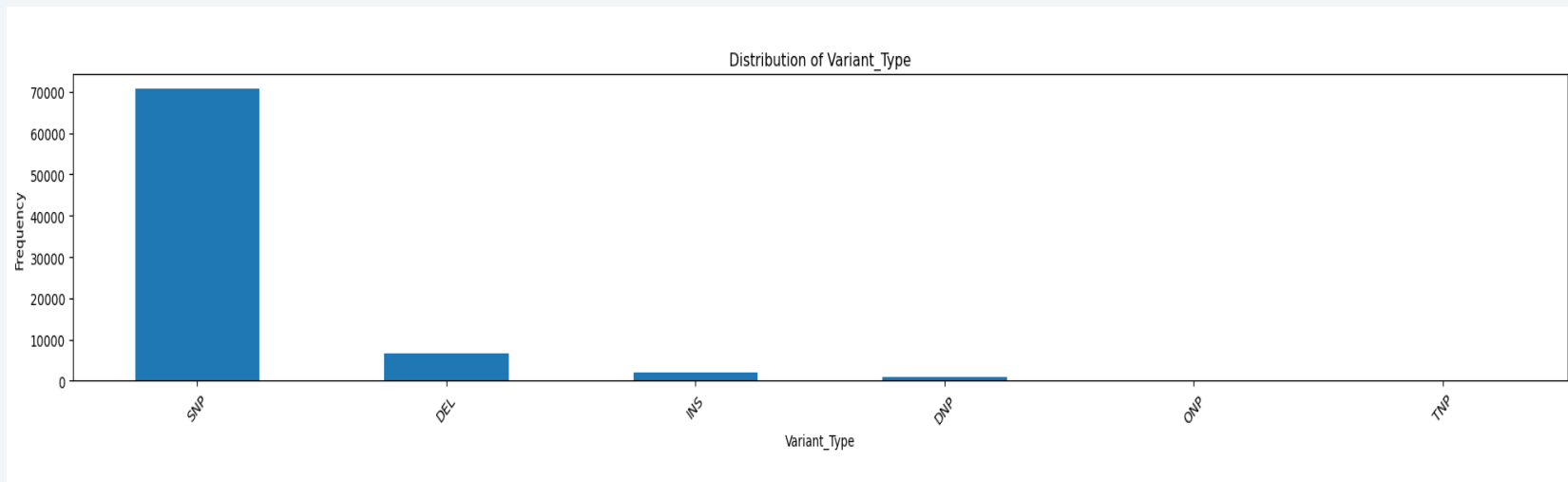
# EDA

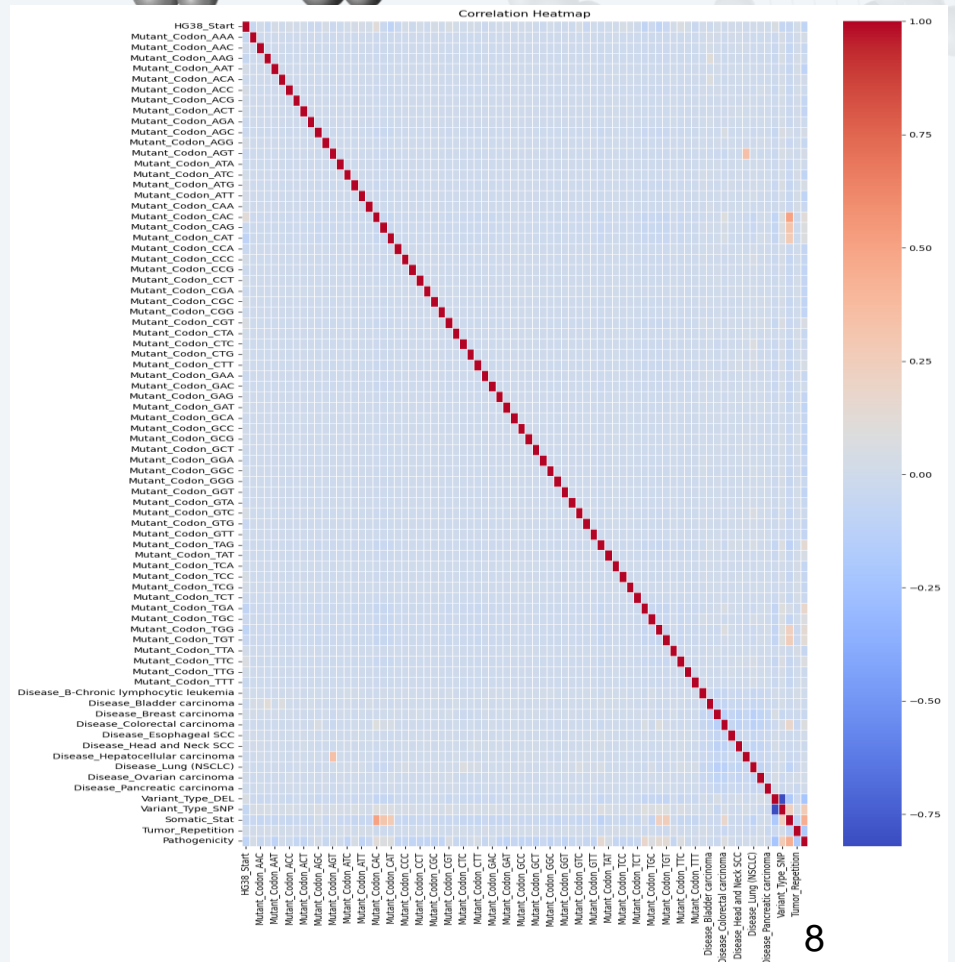
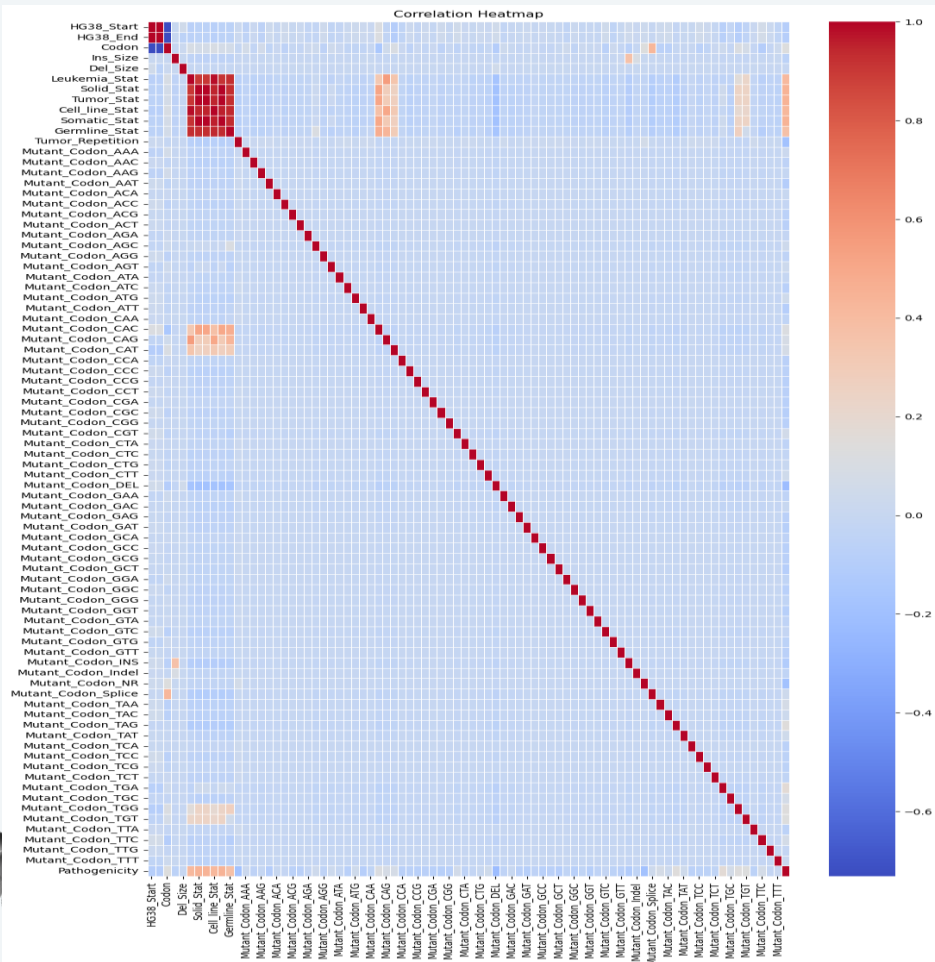


Distribution of Mutant\_Codon



# EDA







# How is my modeling going?

## Baseline Modeling

Logistic  
Regression

Test Accuracy =  
72.08%

Random  
Forest

Test Accuracy =  
99.39%

Random Forest Test Accuracy: 0.88

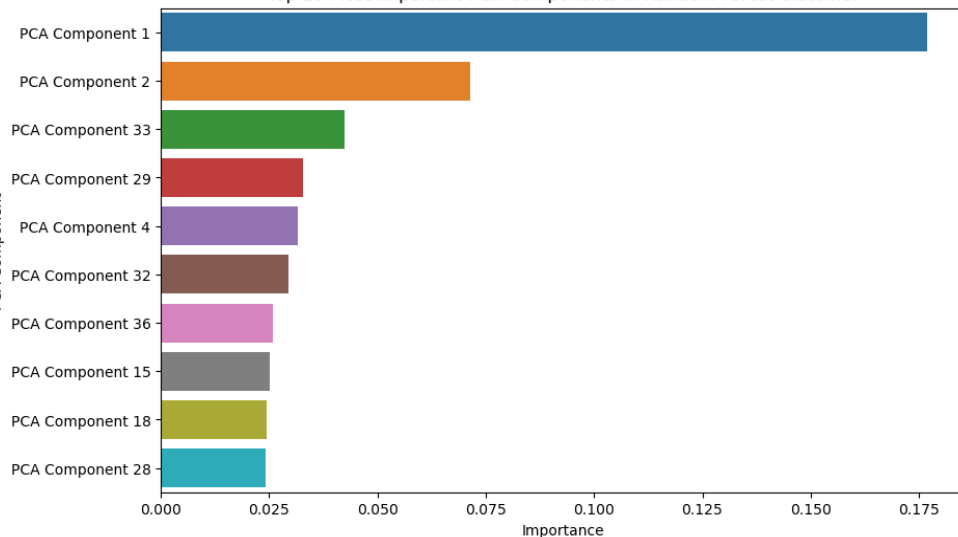
Classification Report for Test Data:

	precision	recall	f1-score	support
0	0.50	0.21	0.29	338
1	0.54	0.36	0.43	860
2	0.81	0.88	0.84	5221
3	0.94	0.94	0.94	9649
accuracy			0.88	16068
macro avg	0.70	0.60	0.63	16068
weighted avg	0.87	0.88	0.87	16068

# Feature contributions



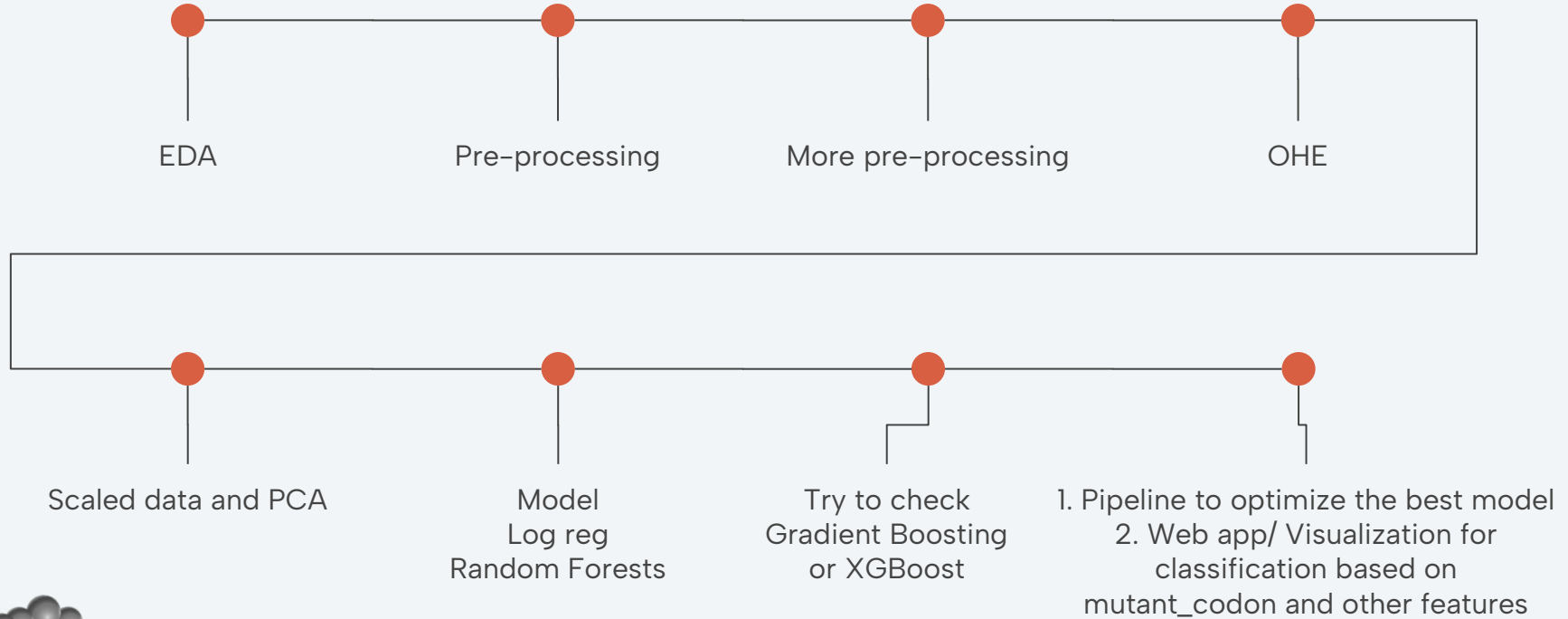
Top 10 Most Important PCA Components in Random Forest Classifier



Top contributing original features:

Variant_Type_DEL	0.052650
Tumor_Repetition	0.038162
Mutant_Codon_TCA	0.034861
Mutant_Codon_CAA	0.032330
Mutant_Codon_CGA	0.031971
HG38_Start	0.031281
Mutant_Codon_CCC	0.029289
Mutant_Codon_AAT	0.029276
Disease_Hepatocellular carcinoma	0.028071
Mutant_Codon_TTA	0.027566

# Milestones reached





# Thanks!

Do you have any questions?

