# Feature Engineering Report
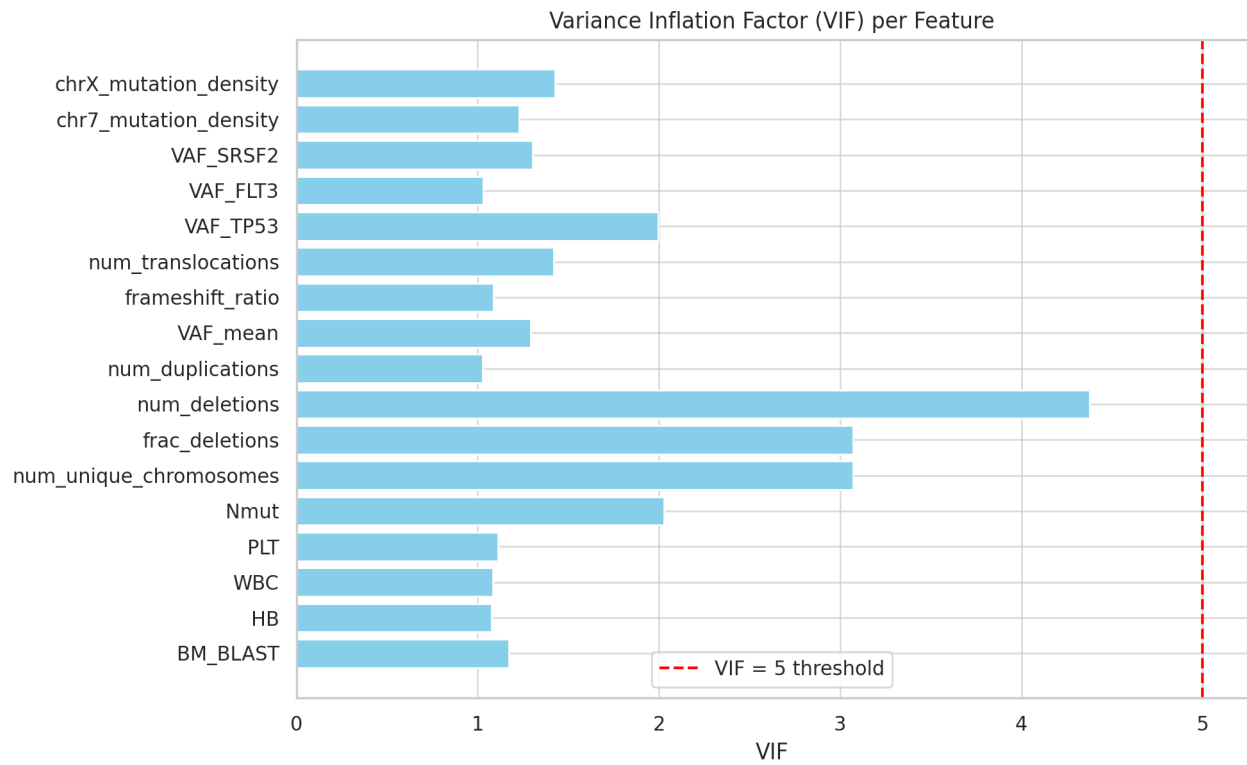# Survival Prediction for Adult Myeloid Leukemia

This report presents an in-depth exploratory analysis and feature engineering process conducted as part of the 2025 QRT Data Challenge in collaboration with Institut Gustave Roussy. The challenge aimed to build predictive models for overall survival (OS) in patients diagnosed with adult myeloid leukemia, using a dataset composed of clinical and molecular data. The dataset includes: - Clinical variables (e.g., blood counts, bone marrow blasts, karyotype abnormalities) - Molecular data from somatic mutations (e.g., affected genes, variant allele frequency) From these raw data, several informative features were derived to capture disease severity and progression patterns. Notable features include: - Mutation burden: number of mutations, unique genes, average VAF - Key gene-specific VAFs (TP53, FLT3, SRSF2) based on their relevance in hematologic malignancies - Chromosome-specific mutation densities (chr7 and chrX) - Cytogenetic abnormalities: deletions, duplications, translocations, and affected chromosomes The goal of this analysis is to visualize these features, assess their predictive potential, and ensure they are suitable for use in survival models.

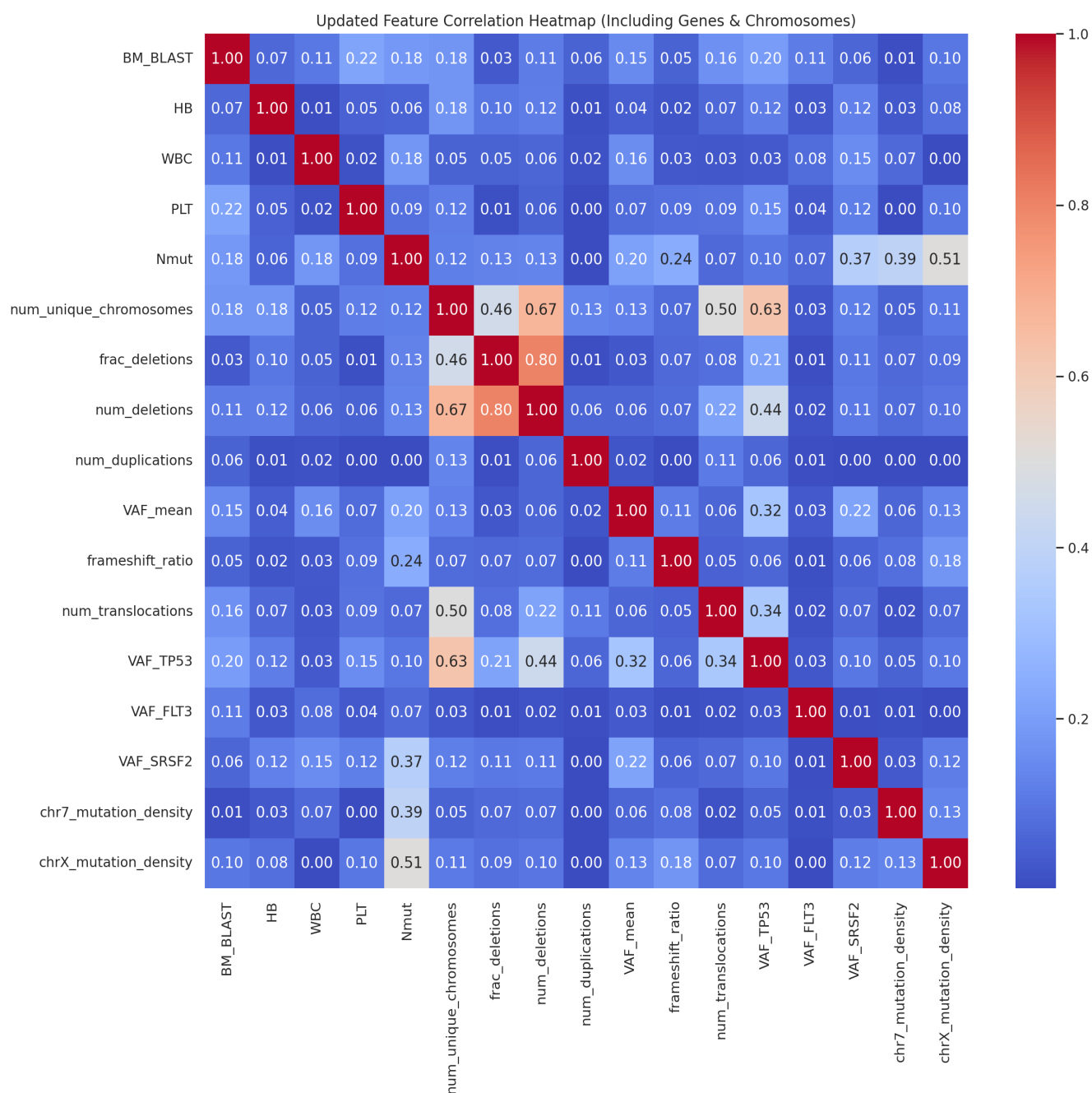## Variance Inflation Factor (VIF) Analysis

| Feature | VIF |
|---|---|
| BM_BLAST | 1.17 |
| HB | 1.08 |
| WBC | 1.08 |
| PLT | 1.11 |
| Nmut | 2.03 |
| num_unique_chromosomes | 3.07 |
| frac_deletions | 3.07 |
| num_deletions | 4.38 |
| num_duplications | 1.03 |
| VAF_mean | 1.29 |
| frameshift_ratio | 1.09 |
| num_translocations | 1.42 |
| VAF_TP53 | 1.99 |
| VAF_FLT3 | 1.03 |
| VAF_SRSF2 | 1.3 |
| chr7_mutation_density | 1.23 |
| chrX_mutation_density | 1.43 |

*Features with VIF < 5 are considered acceptable, indicating low multicollinearity.*

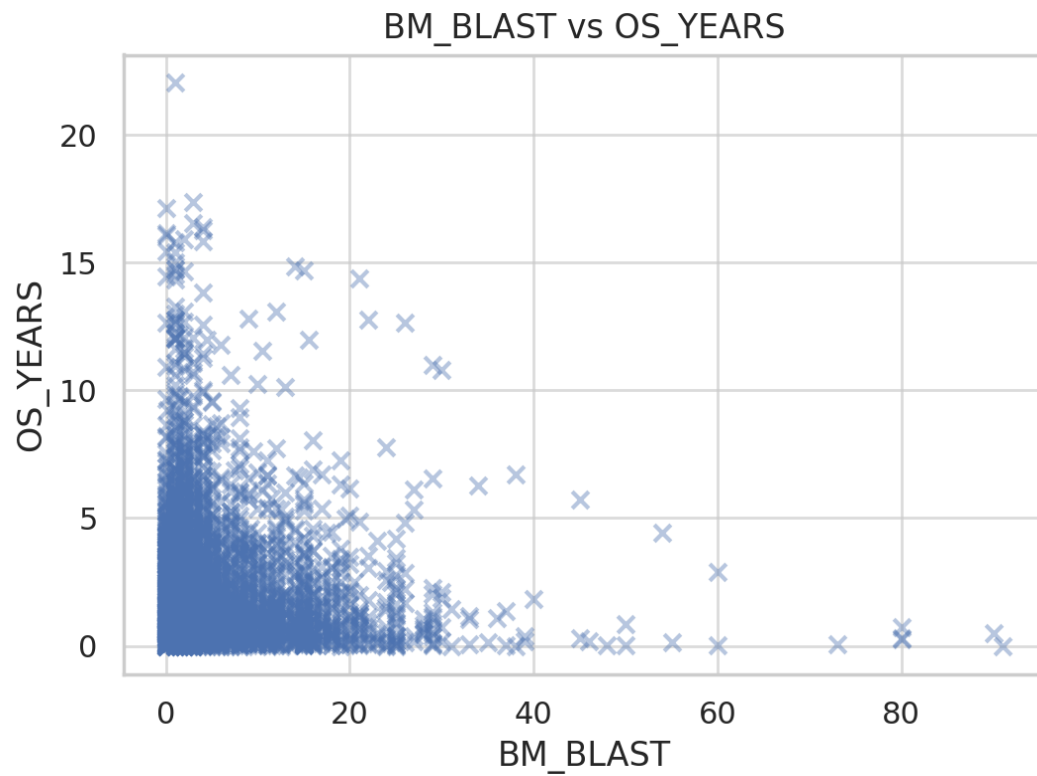Variance Inflation Factor (VIF) per Feature

This bar chart visualizes the Variance Inflation Factor for each feature. The red dashed line indicates a VIF of 5, which is commonly used as a threshold for multicollinearity concerns.
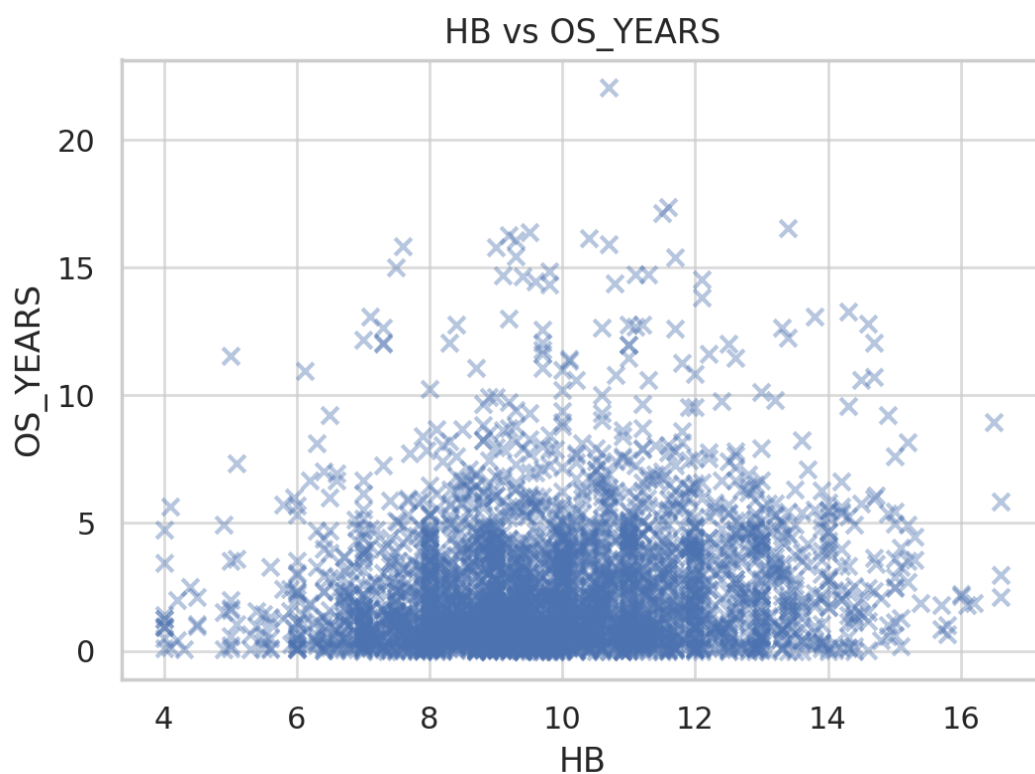
# Feature Correlation Heatmap



Updated Feature Correlation Heatmap (Including Genes & Chromosomes)

This heatmap shows pairwise correlations between engineered features. A high correlation (> 0.8) was detected between 'num_deletions' and 'frac_deletions', which may affect model stability. All other features show acceptable independence.
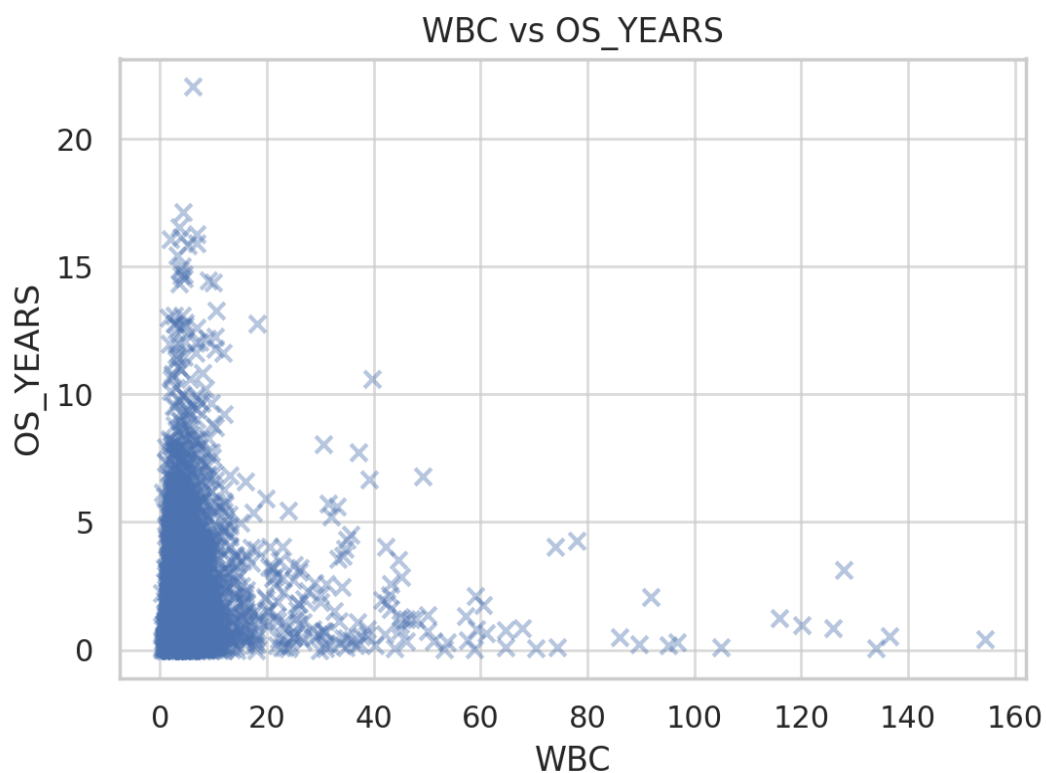
## BM_BLAST vs OS_YEARS



The scatter plot above shows the relationship between BM_BLAST and overall survival (OS_YEARS). This helps assess whether BM_BLAST might influence survival duration.
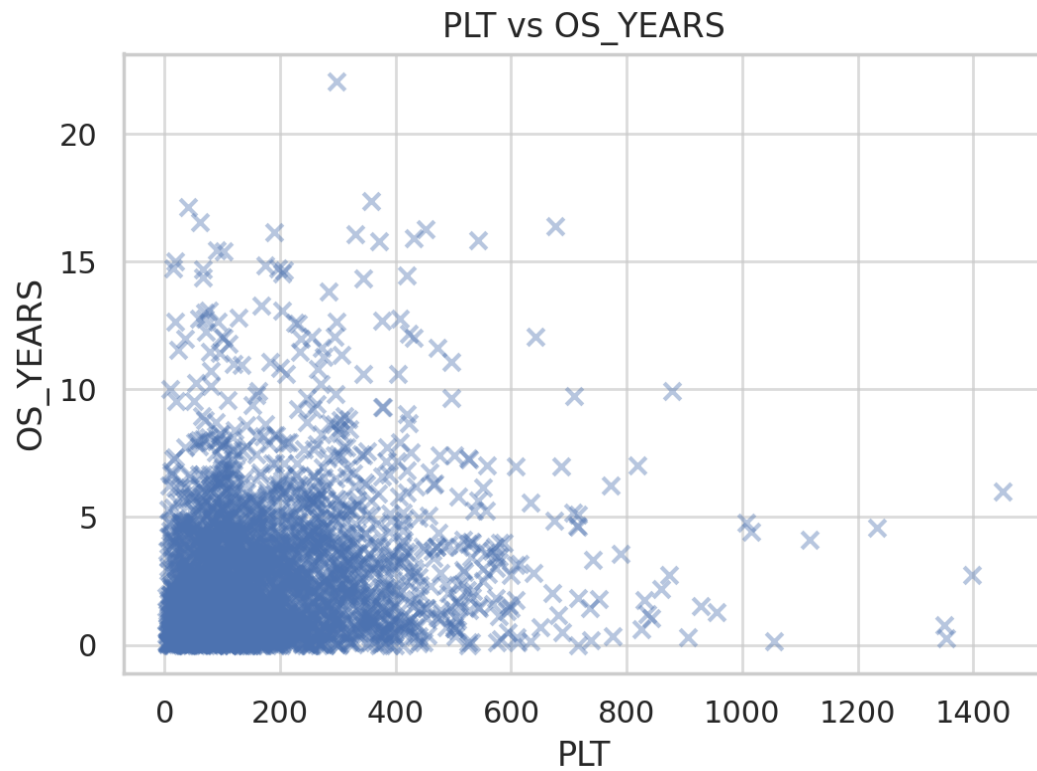
## HB vs OS_YEARS

The scatter plot above shows the relationship between HB and overall survival (OS_YEARS). This helps assess whether HB might influence survival duration.
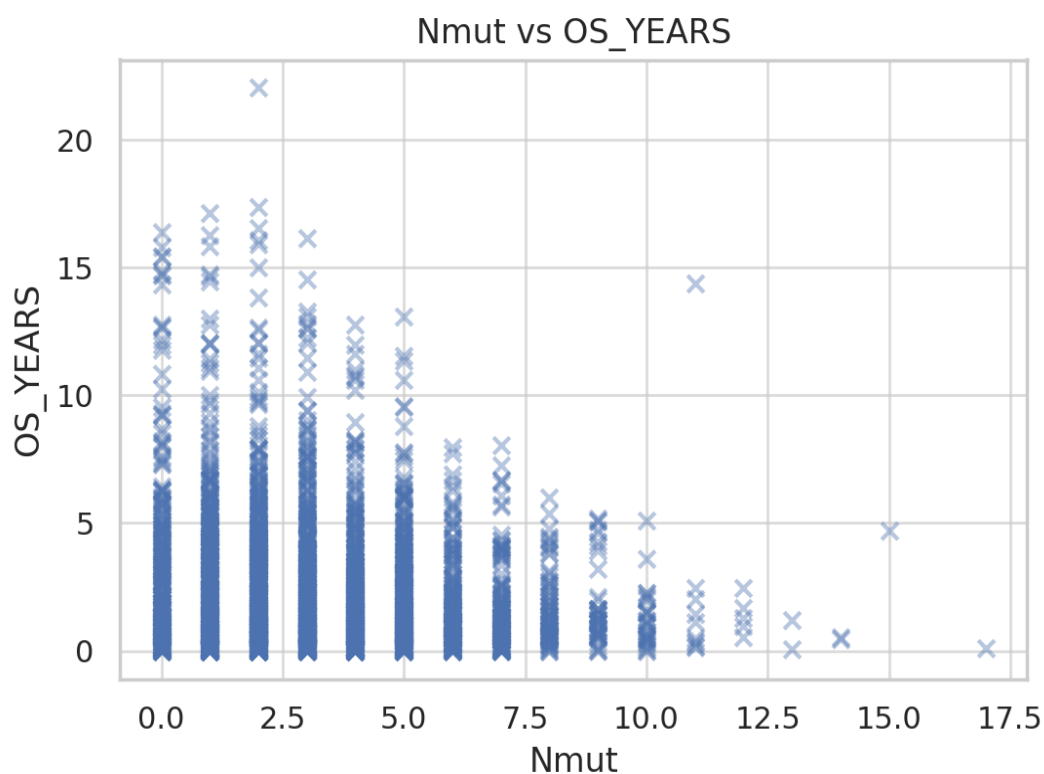

## *WBC vs OS_YEARS*

The scatter plot above shows the relationship between WBC and overall survival (OS_YEARS). This helps assess whether WBC might influence survival duration.
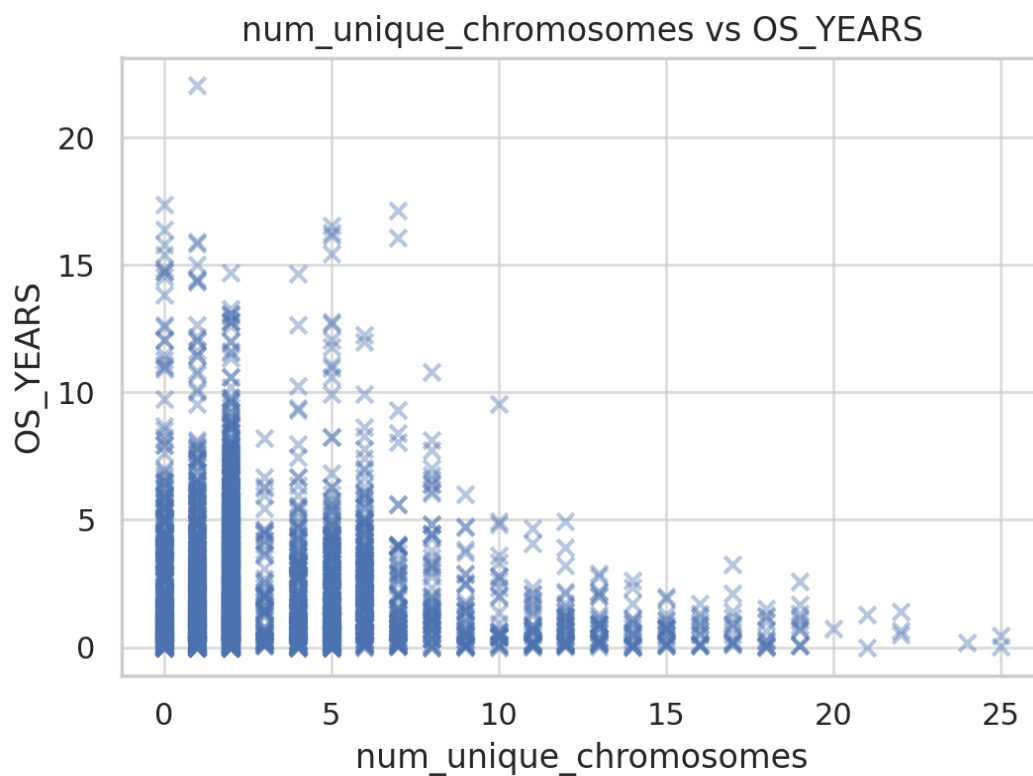
## PLT vs OS_YEARS



The scatter plot above shows the relationship between PLT and overall survival (OS_YEARS). This helps assess whether PLT might influence survival duration.

## Nmut vs OS_YEARS

## Nmut vs OS_YEARS



The scatter plot above shows the relationship between Nmut and overall survival (OS_YEARS). This helps assess whether Nmut might influence survival duration.
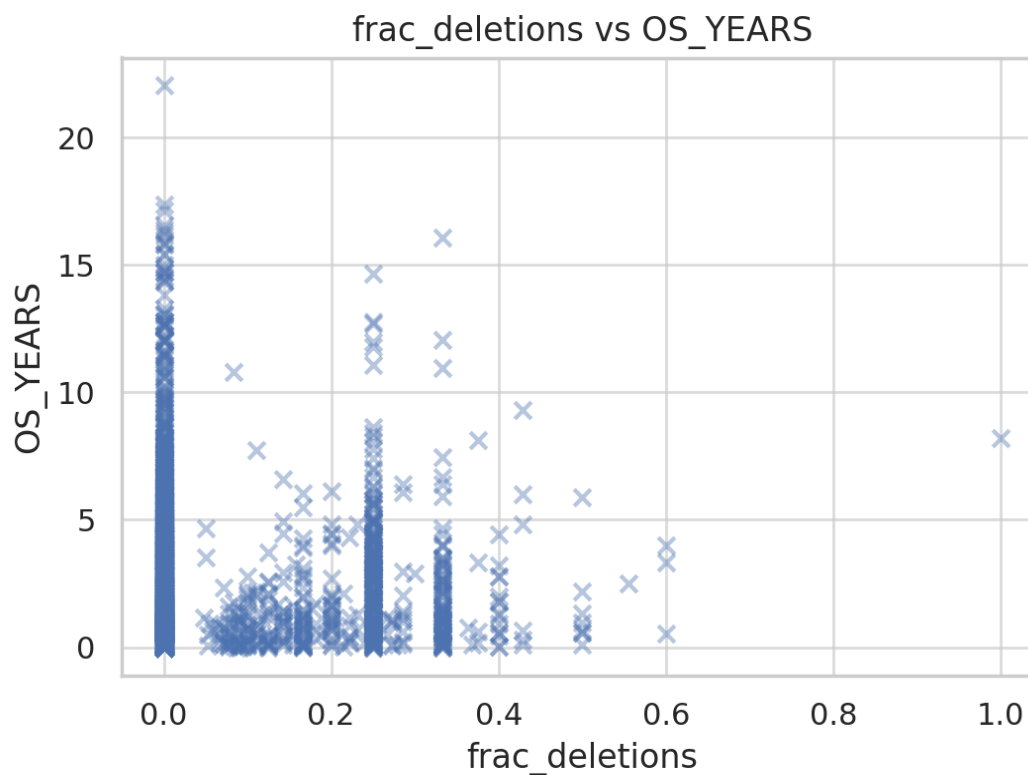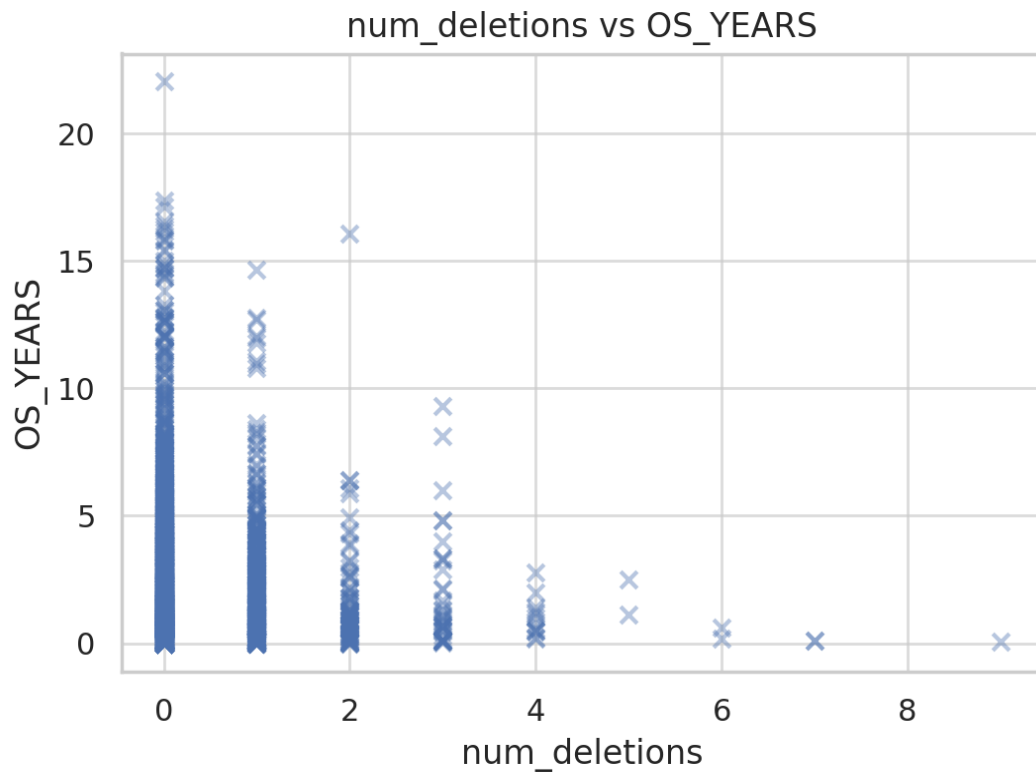
## num_unique_chromosomes vs OS_YEARS

The scatter plot above shows the relationship between num_unique_chromosomes and overall survival (OS_YEARS). This helps assess whether num_unique_chromosomes might influence survival duration.

### *frac_deletions vs OS_YEARS*



frac_deletions vs OS_YEARS
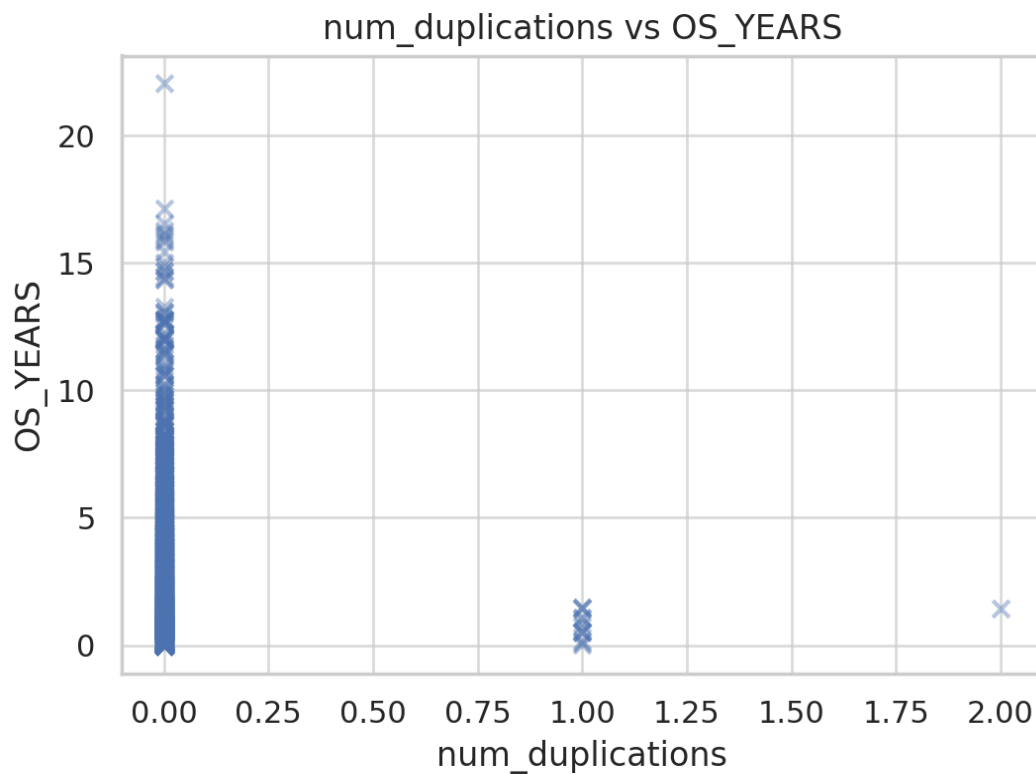
The scatter plot above shows the relationship between frac_deletions and overall survival (OS_YEARS). This helps assess whether frac_deletions might influence survival duration.
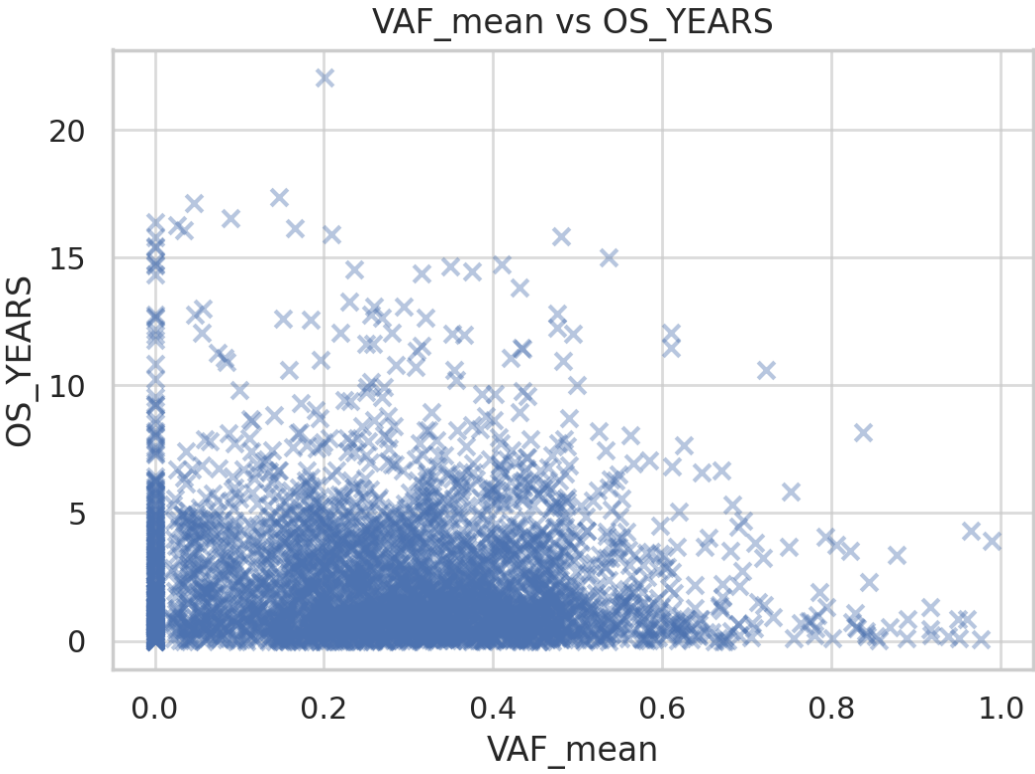
### *num_deletions vs OS_YEARS*

The scatter plot above shows the relationship between num_deletions and overall survival (OS_YEARS). This helps assess whether num_deletions might influence survival duration.
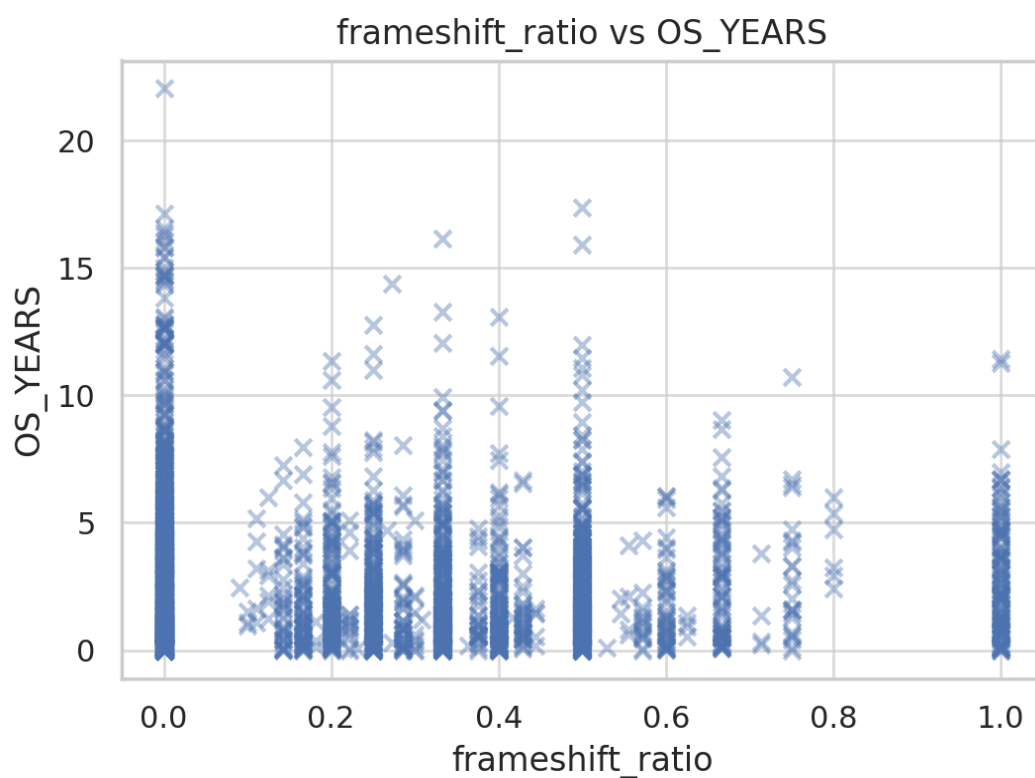
### num_duplications vs OS_YEARS

The scatter plot above shows the relationship between num_duplications and overall survival (OS_YEARS). This helps assess whether num_duplications might influence survival duration.
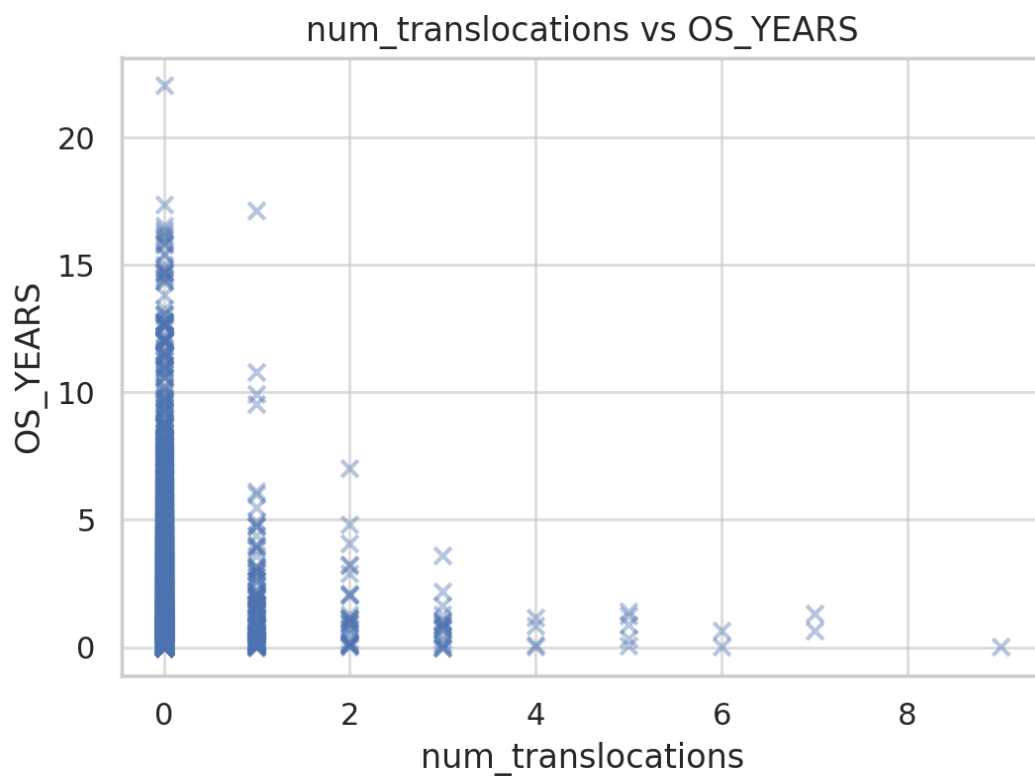
## *VAF_mean vs OS_YEARS*



VAF_mean vs OS_YEARS

The scatter plot above shows the relationship between VAF_mean and overall survival (OS_YEARS). This helps assess whether VAF_mean might influence survival duration.
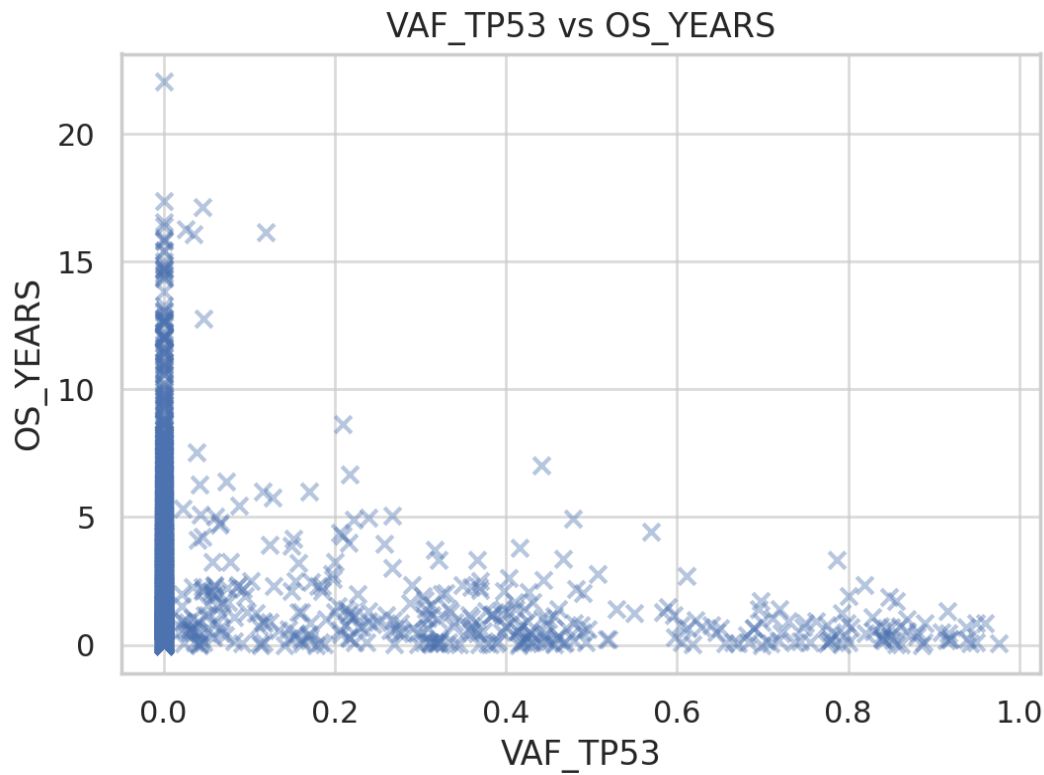
## *frameshift_ratio vs OS_YEARS*

The scatter plot above shows the relationship between frameshift_ratio and overall survival (OS_YEARS). This helps assess whether frameshift_ratio might influence survival duration.
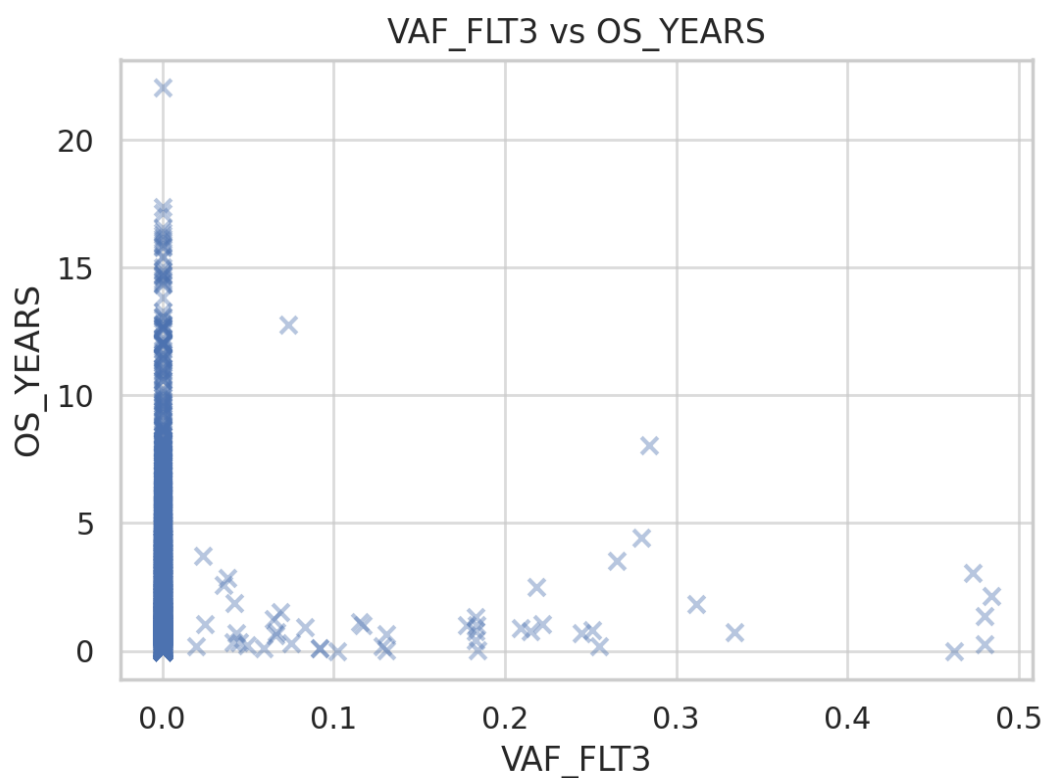
### num_translocations vs OS_YEARS

The scatter plot above shows the relationship between num_translocations and overall survival (OS_YEARS). This helps assess whether num_translocations might influence survival duration.

## *VAF_TP53 vs OS_YEARS*


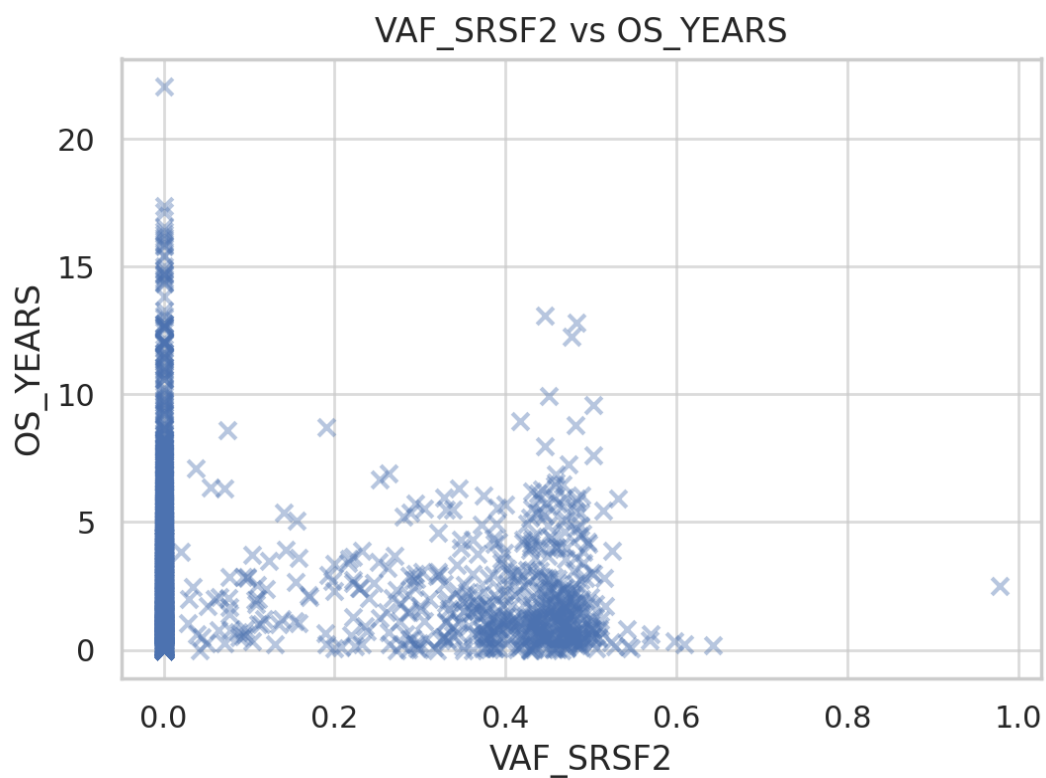
VAF_TP53 vs OS_YEARS

The scatter plot above shows the relationship between VAF_TP53 and overall survival (OS_YEARS). This helps assess whether VAF_TP53 might influence survival duration.
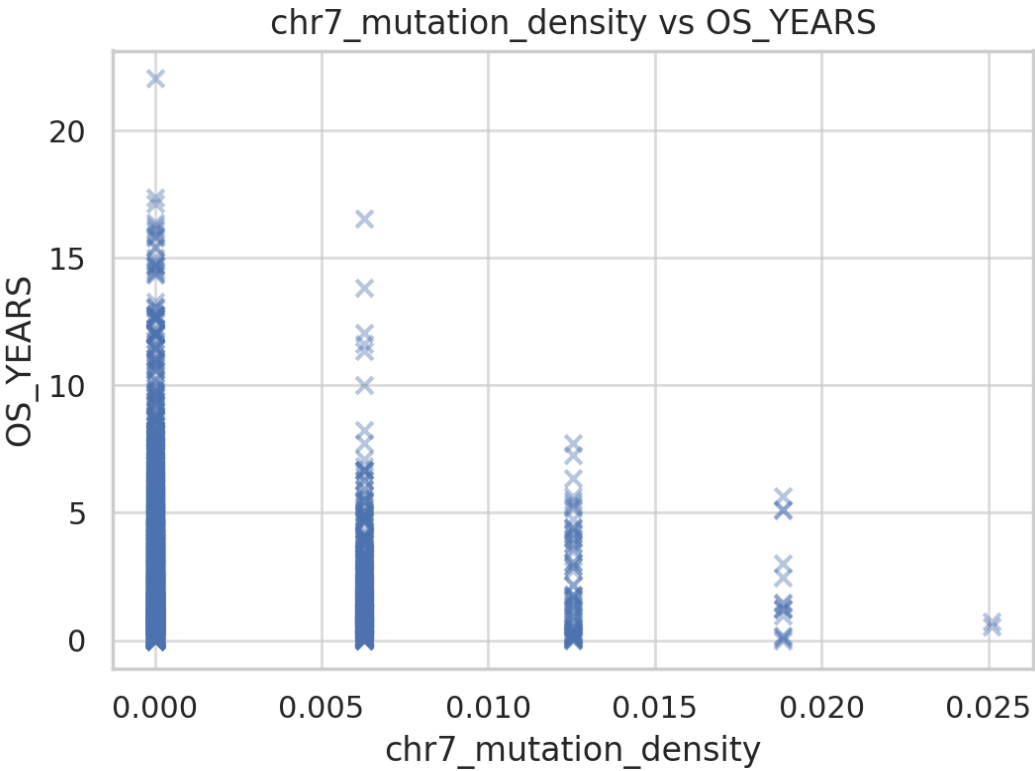
## *VAF_FLT3 vs OS_YEARS*

## VAF_FLT3 vs OS_YEARS



The scatter plot above shows the relationship between VAF_FLT3 and overall survival (OS_YEARS). This helps assess whether VAF_FLT3 might influence survival duration.
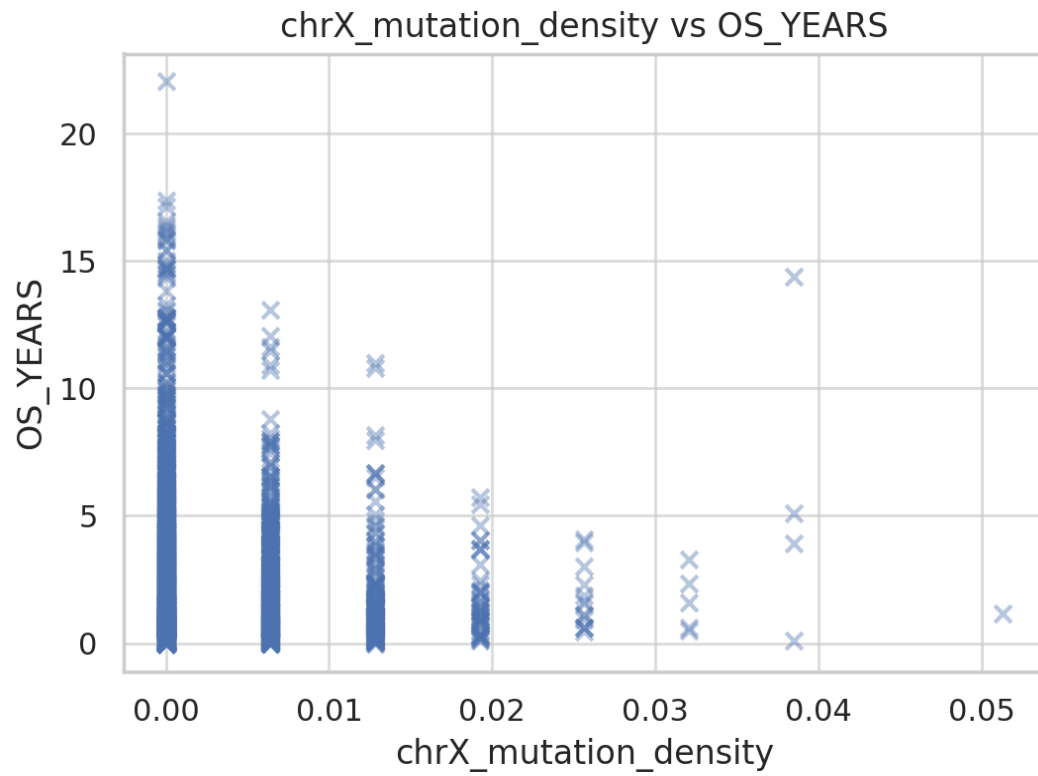
## VAF_SRSF2 vs OS_YEARS

The scatter plot above shows the relationship between VAF_SRSF2 and overall survival (OS_YEARS). This helps assess whether VAF_SRSF2 might influence survival duration.

### *chr7_mutation_density vs OS_YEARS*



chr7_mutation_density vs OS_YEARS

The scatter plot above shows the relationship between chr7_mutation_density and overall survival (OS_YEARS). This helps assess whether chr7_mutation_density might influence survival duration.

### *chrX_mutation_density vs OS_YEARS*

chrX_mutation_density vs OS_YEARS

The scatter plot above shows the relationship between chrX_mutation_density and overall survival (OS_YEARS). This helps assess whether chrX_mutation_density might influence survival duration.

## Conclusion

In conclusion, this exploratory analysis has highlighted several engineered features with strong clinical relevance and potential predictive value for survival in adult myeloid leukemia. The inclusion of gene-specific VAFs and chromosome-specific mutation densities provides a genomic dimension that complements clinical observations. The derived features have demonstrated acceptable multicollinearity and largely low inter-feature correlation, making them well-suited for integration into survival models such as Extra Survival Trees or Cox Proportional Hazards. These findings will guide the subsequent model development phase to improve personalized prognosis and treatment planning in hematologic oncology.