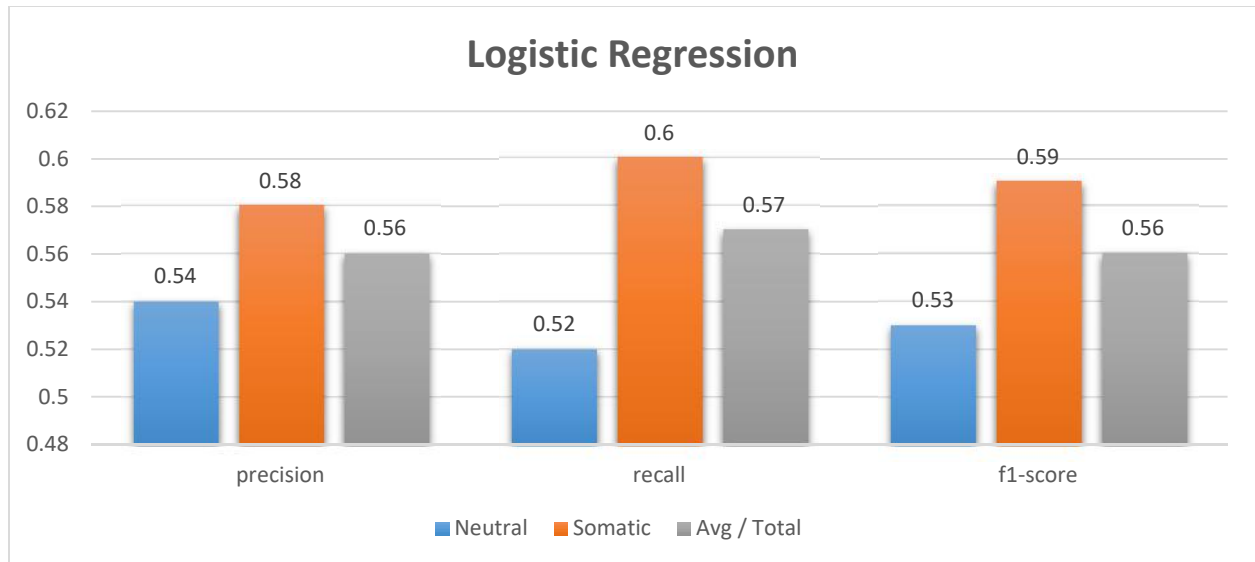**Inference from our project**
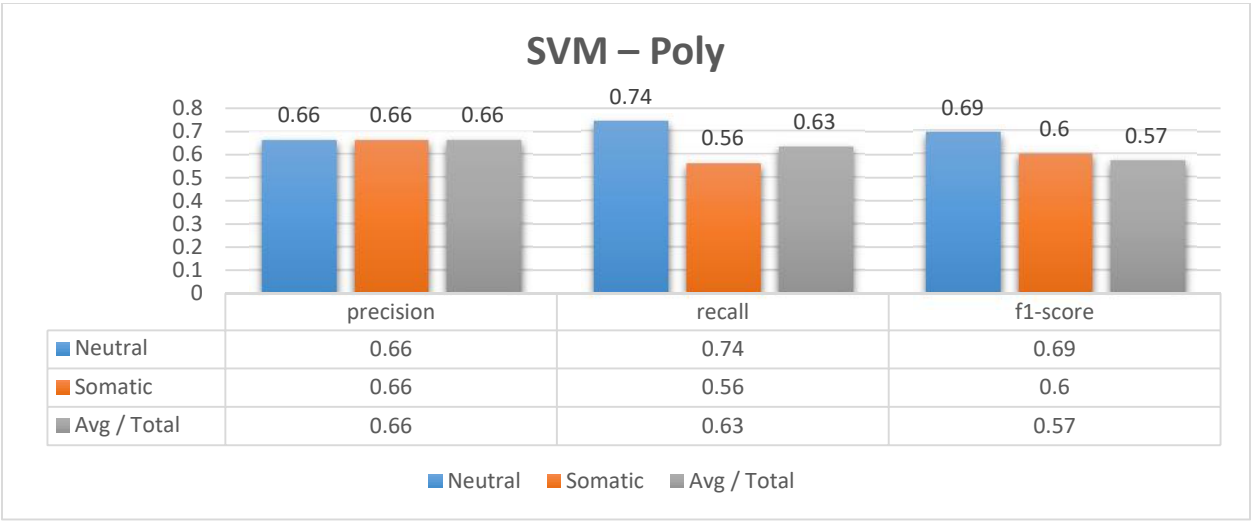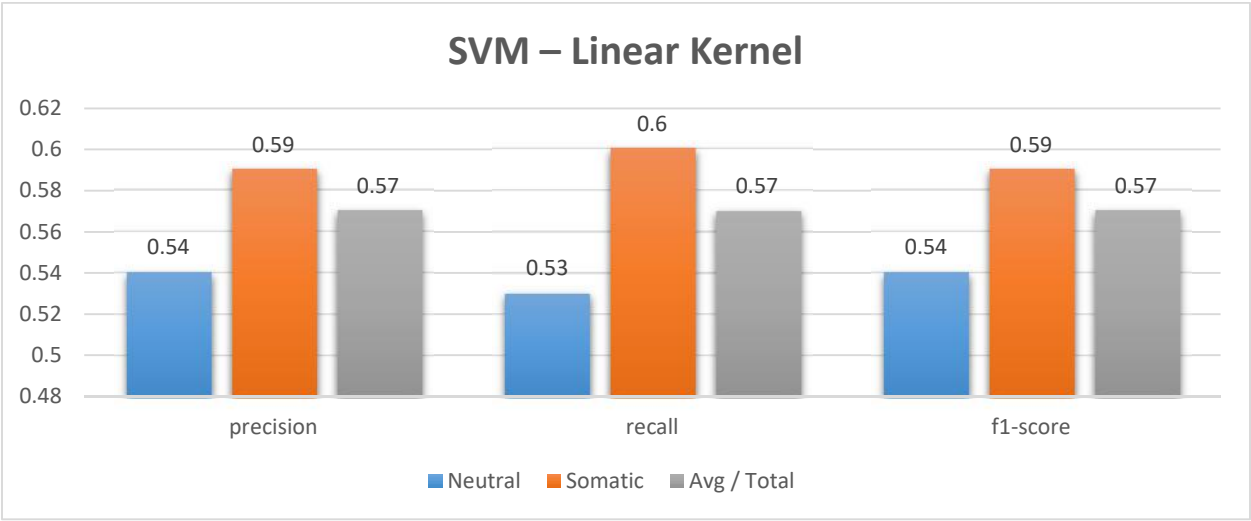
**Logistic Regression**

## Logistic Regression

| | Neutral | Somatic | Avg / Total |
|---|---|---|---|
| precision | 0.54 | 0.58 | 0.56 |
| recall | 0.52 | 0.6 | 0.57 |
| f1-score | 0.53 | 0.59 | 0.56 |

**Random Forest**

## Random Forest

| | Neutral | Somatic | avg / total |
|---|---|---|---|
| precision | 0.57 | 0.67 | 0.62 |
| recall | 0.72 | 0.52 | 0.61 |
| f1-score | 0.64 | 0.59 | 0.61 |

**SVM –Linear Kernel**

## SVM – Linear Kernel

| | precision | recall | f1-score |
|---|---|---|---|
| Neutral | 0.54 | 0.53 | 0.54 |
| Somatic | 0.59 | 0.6 | 0.59 |
| Avg / Total | 0.57 | 0.57 | 0.57 |

Legend: ■ Neutral ■ Somatic ■ Avg / Total

## SVM – Poly

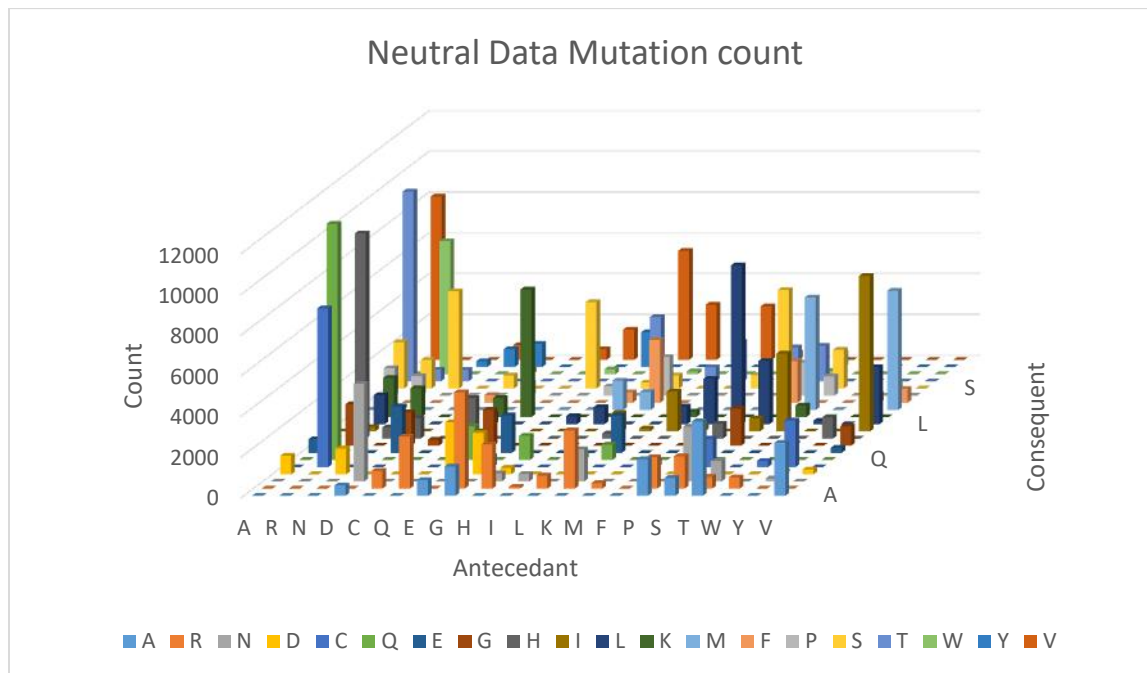| | precision | recall | f1-score |
|---|---|---|---|
| Neutral | 0.66 | 0.74 | 0.69 |
| Somatic | 0.66 | 0.56 | 0.6 |
| Avg / Total | 0.66 | 0.63 | 0.57 |

Legend: ■ Neutral ■ Somatic ■ Avg / Total
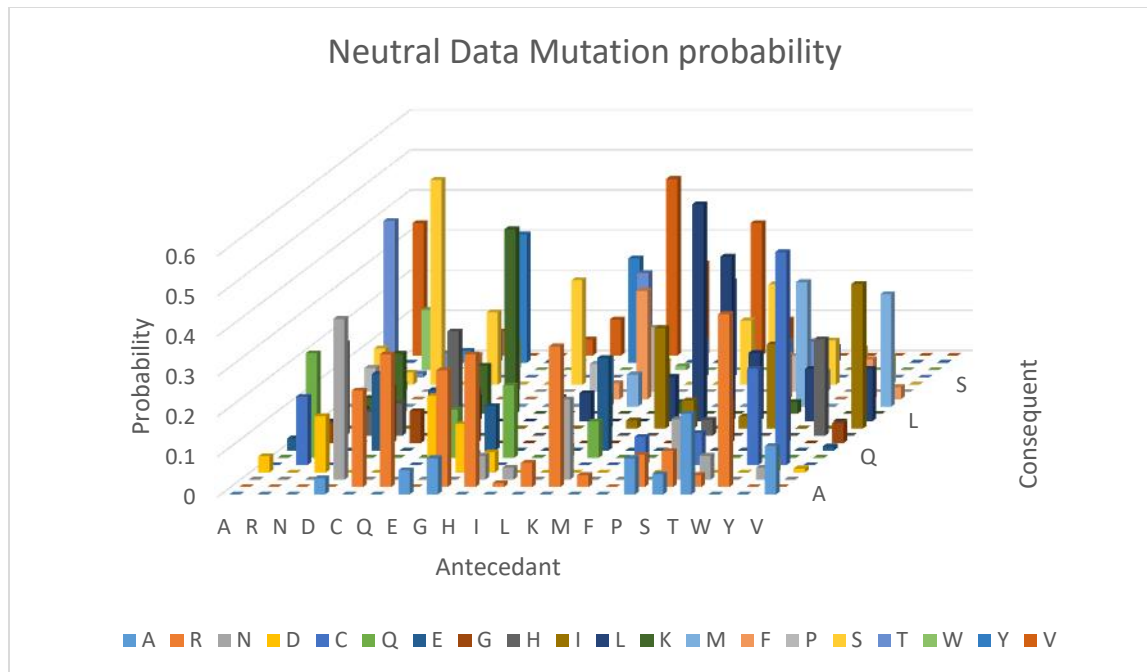
**Accuracy**



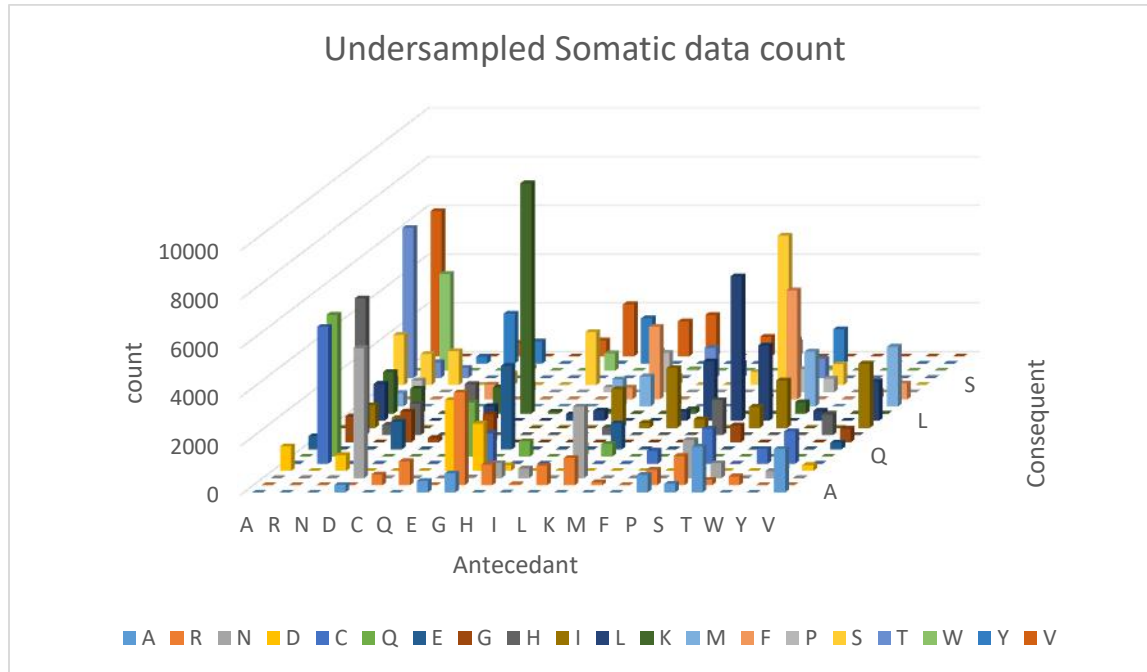**Neutral Data**



**Inference**

- The most frequently occurring mutation is R to Q and they occur 11601 times in the neutral dataset.
- The least occurring mutations are T to Q , S to D, P to D , D to S  and Y to R and they occur only once in the neutral dataset

Neutral Data Mutation probability

**Inference**

- The mutation F to L has the highest probability of occurrence 0.54.
- The mutations V to D , V to E , S to W, R to T, G to W, I to R, L to W, S to W have the lowest probability of occurrence of 0.01
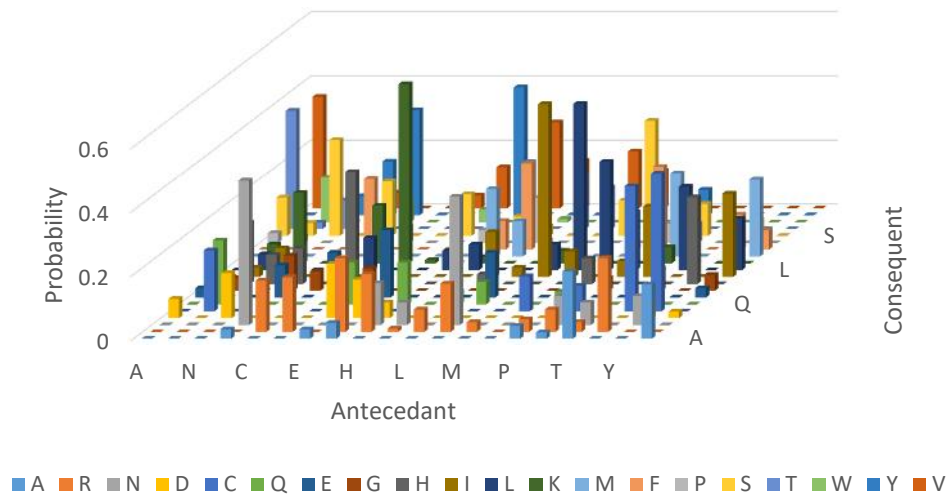
**Under Sampled data**



Undersampled Somatic data count

**Inference**

- The mutation E to K was the most frequently occurring mutation with an occurrence count of 9446
- There were numerous mutations with an occurrence count of one which indicated that there were lots of sparsely occurring mutations in the dataset
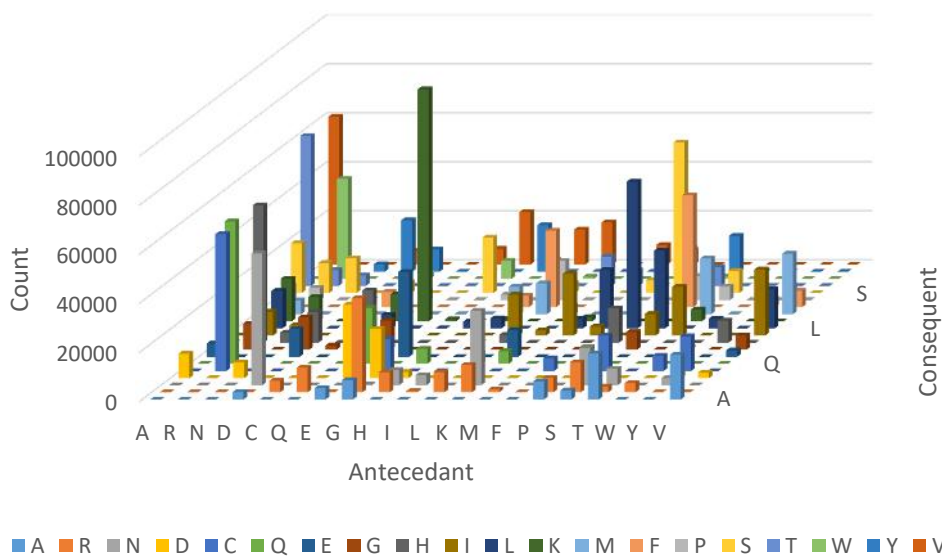
Undersampled Somatic data probability

**Inference**

- The mutation E to K has the maximum probability of occurrence with a probability of occurrence of 0.56
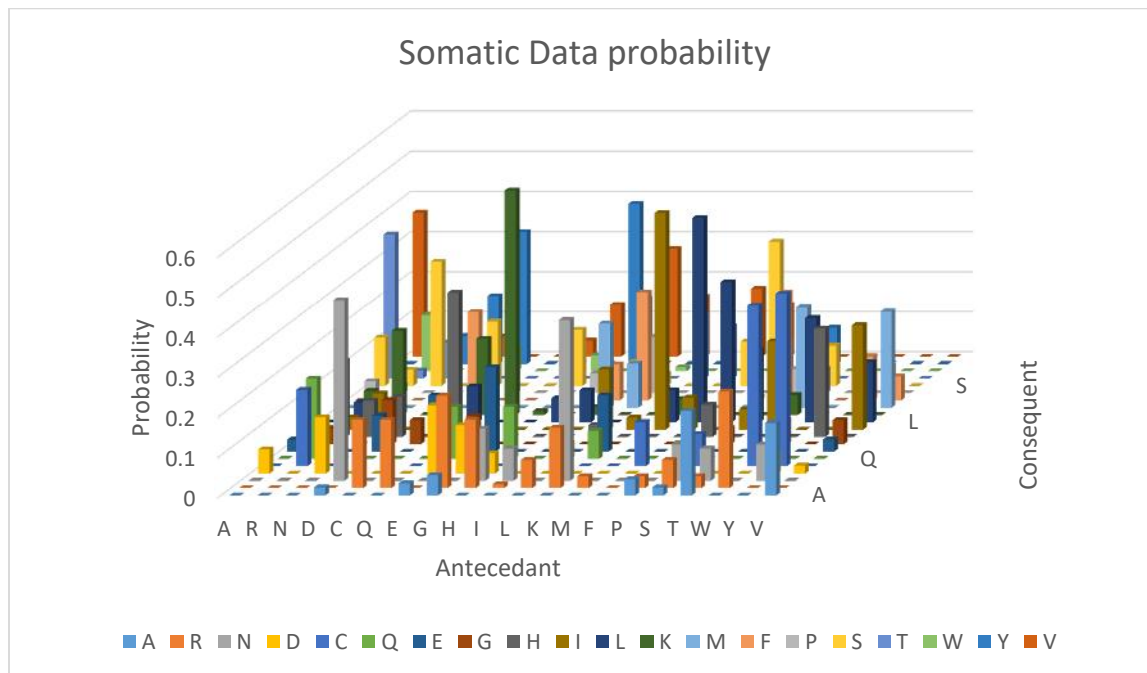- There were numerous mutations with the lowest probability of 0.01

**Complete data**



Somatic data count

**Inference**

- The mutation E to K has the maximum count of occurrence with a count of 94674 times
- There were numerous mutations with an occurrence count of one which indicated that there were lots of sparsely occurring mutations in the dataset



Somatic Data probability

**Inference**

- The mutation E to K has the maximum probability of occurrence with a probability of occurrence of 0.56
- There were numerous mutations with the lowest probability of 0.01.