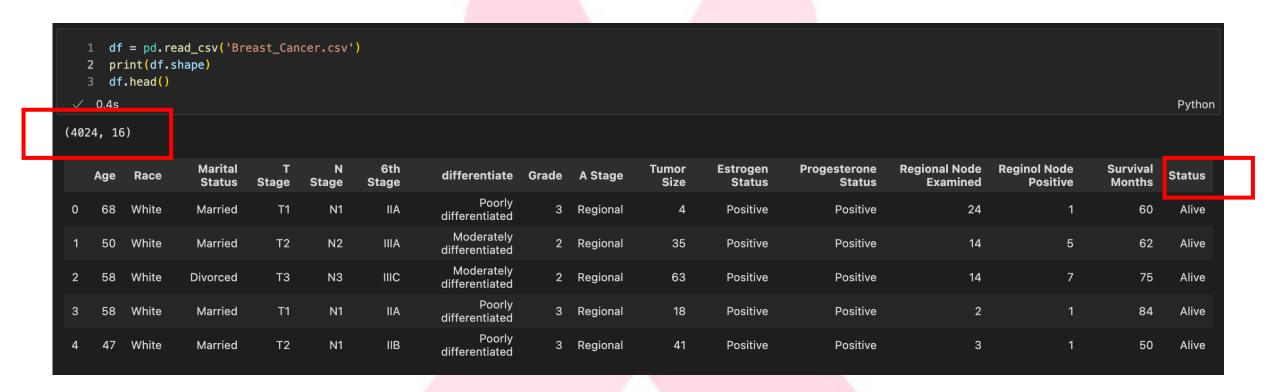
# PROJECT PROPOSAL: BREAST CANCER – MACHINE LEARNING

Source: https://www.kaggle.com/datasets/reihanenamdari/breast-cancer

# Goal:

- To achieve the best accuracy in predicting Status

#### The Dataset:



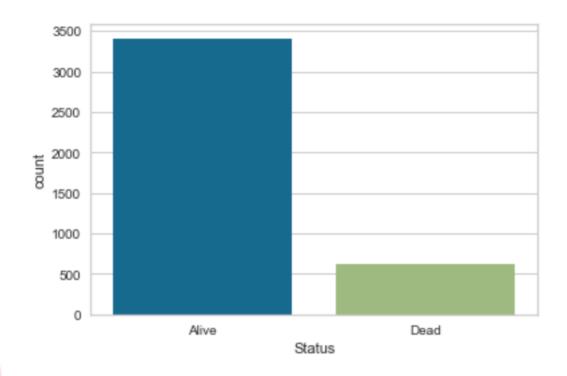
- Classification - Binary

CI_	
	TIIS
<b>349</b>	1440

Alive Dead

Alive 3408 Dead 616

Name: Status, dtype: int64



# T Stage N Stage

# Tumor Size Lymph Node Status

```
1 df['T Stage '].value_counts()

7 0.2s

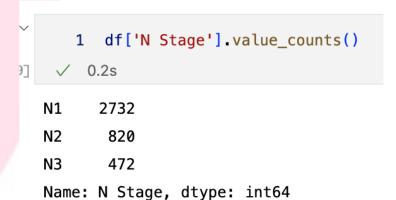
T2 1786

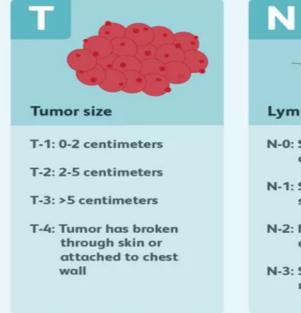
T1 1603

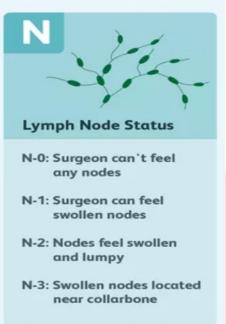
T3 533

T4 102

Name: T Stage , dtype: int64
```



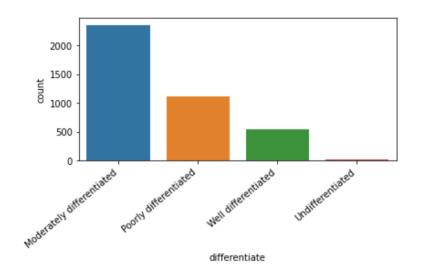


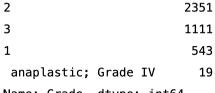


**Differentiate** Grade

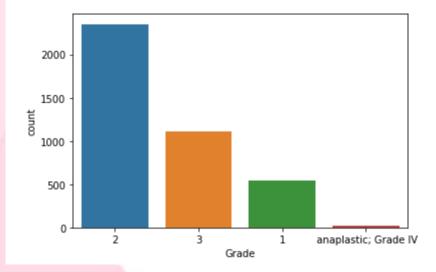
#### How Active is the Cancer

Moderately differentiated 2351
Poorly differentiated 1111
Well differentiated 543
Undifferentiated 19
Name: differentiate, dtype: int64





Name: Grade, dtype: int64



#### **Grading invasive breast cancer cells**

Three features of the invasive breast cancer cell are studied and each is given a score. The scores are then added to get a number between 3 and 9 that is used to get a grade of 1, 2, or 3, which is noted on your pathology report. Sometimes the terms well differentiated, moderately differentiated, and poorly differentiated are used to describe the grade instead of numbers:

- Grade 1 or well differentiated (score 3, 4, or 5). The cells are slower-growing, and look
  more like normal breast cells.
- Grade 2 or moderately differentiated (score 6, 7). The cells are growing at a speed of and look like cells somewhere between grades 1 and 3.
- **Grade 3 or poorly differentiated** (score 8, 9). The cancer cells look very different from normal cells and will probably grow and spread faster.

#### A Stage

#### 6<sup>th</sup> Stage

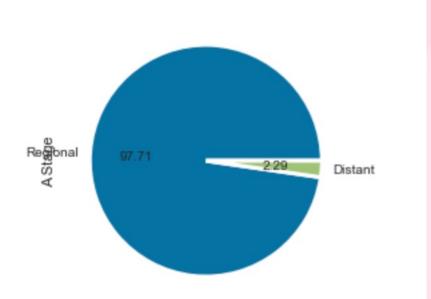
## State of Cancer [ 0 = Regional / 1 = Distant]

## Stage of Cancer

Regional = Cancer has spread to nearby lymph nodes, tissues

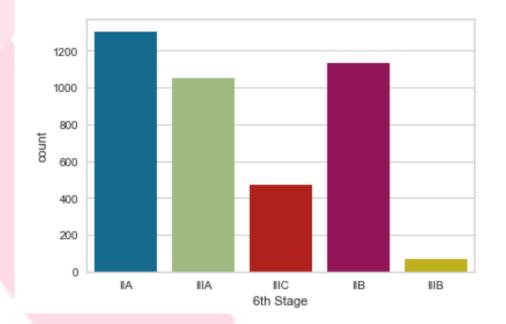
Distant = Cancer has spread to distant parts of the body such as the lungs, liver or bones

https://www.cancer.net/cancer-types/breast-cancer/stages





<AxesSubplot:xlabel='6th Stage', ylabel='count'>



If both positive, cancer is less aggressive

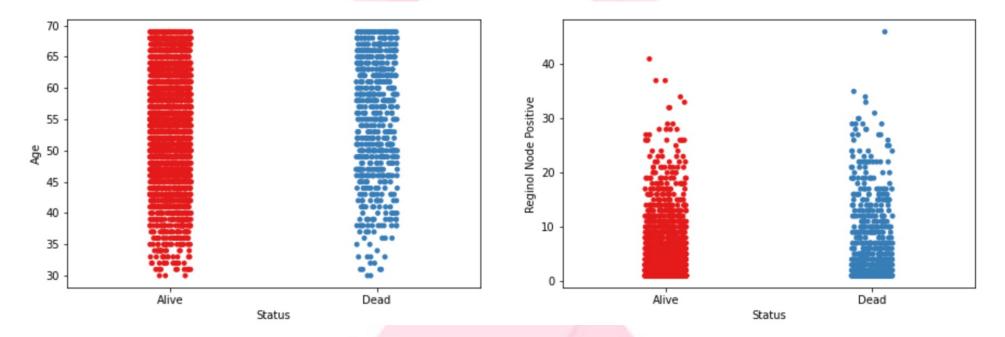
# No empty data cells:

Age	0
Race	0
Marital Status	0
T Stage	0
N Stage	0
6th Stage	0
differentiate	0
Grade	0
A Stage	0
Tumor Size	0
Estrogen Status	0
Progesterone Status	0
Regional Node Examined	0
Reginol Node Positive	0
Survival Months	0
Status	0
dtype: int64	

#### Correlation at first Glance:



#### Numerical Variables:



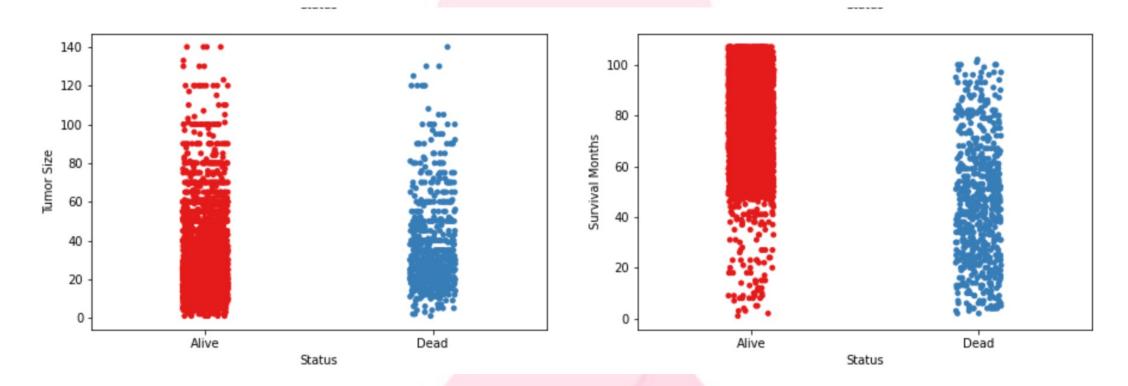
- Here we are checking if Age & Regional Node Positive affects/influence the Status.

#### Age:

It seems that there is no specific range of Age affecting the Status.

### **Regional Node Positive:**

Higher Node positive tend to die



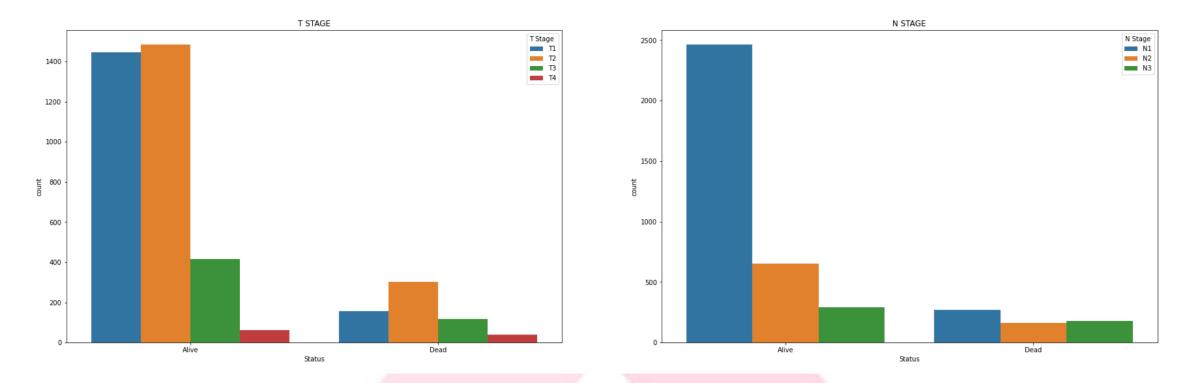
- Here we are checking if Tumor Size & Survival Months affects/influence the Status.

#### **Tumor Size:**

Larger tumor Size tend to die , but seems unclear from this plot

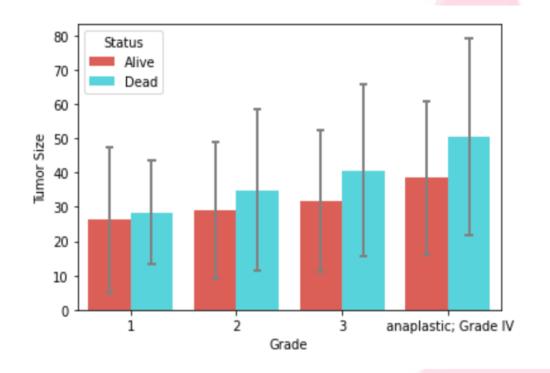
#### **Survival Months:**

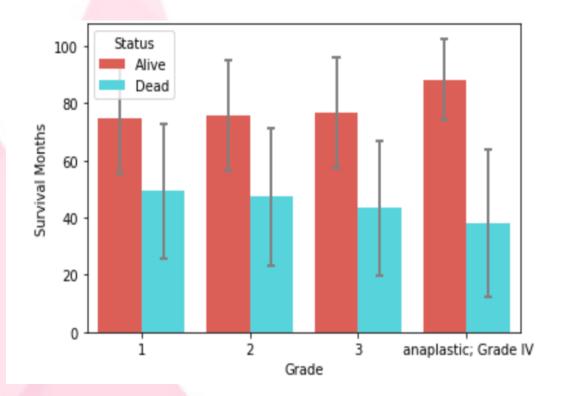
Higher Survival Months tend to survive



- From The T Stage , we can deduce that early stage has higher chance of Survival

#### Comparing Against Grade:





- The higher the grade the larger the tumor size
- Survival Months and Grade has little correlation; Alive goes higher similarly Dead is decreasing

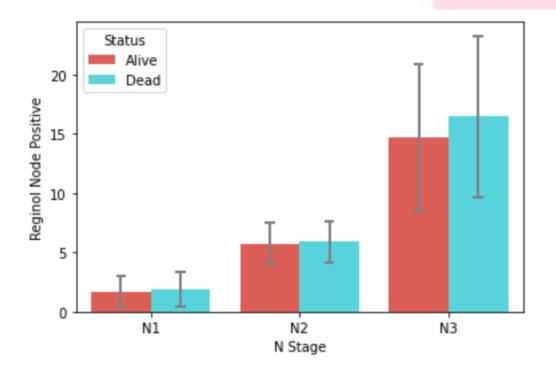
# **Label Encoding:**

	Age	Race	Marital Status	T Stage	N Stage	6th Stage	differentiate	Grade	A Stage	Tumor Size	Estrogen Status	Progesterone Status	Regional Node Examined	Reginol Node Positive	Survival Months	Status
0	38	2	1	0	0	0	1	3	1	3	1	1	23	0	59	0
1	20	2	1	1	1	2	0	2	1	34	1	1	13	4	61	0
2	28	2	0	2	2	4	0	2	1	62	1	1	13	6	74	0
3	28	2	1	0	0	0	1	3	1	17	1	1	1	0	83	0
4	17	2	1	1	0	1	1	3	1	40	1	1	2	0	49	0
4019	32	1	1	0	0	0	0	2	1	8	1	1	0	0	48	0
4020	26	2	0	1	1	2	0	2	1	45	1	1	13	7	68	0
4021	38	2	1	1	0	1	0	2	1	21	1	0	10	2	68	0
4022	28	0	0	1	0	1	0	2	1	43	1	1	10	0	71	0
4023	16	2	1	1	0	1	0	2	1	29	1	1	6	1	99	0

# **Correlation:**

Age -	1	0.08	0.051	-0.067	0.0029	-0.019	0.016	-0.093	0.021	-0.078	0.06	-0.021	-0.034	0.013	-0.0094	0.056
Race -	0.08	1	-0.11	0.0011	-0.032	-0.025	-0.019	-0.057	0.007	-0.0037	0.058	0.034	-0.0044	-0.015	0.041	-0.05
Marital Status -	0.051	-0.11	1	0.0093	0.013	0.003	0.032	-0.0043	0.0049	0.01	-0.018	-0.013	-0.0042	0.0058	-0.026	0.033
T Stage -	-0.067	0.0011	0.0093	1	0.28	0.61	-0.031	0.11	-0.22	0.82	-0.061	-0.058	0.11	0.24	-0.086	0.15
N Stage -	0.0029	-0.032	0.013	0.28	1	0.88	-0.036	0.15	-0.26	0.28	-0.1	-0.094	0.33	0.84	-0.14	0.26
6th Stage -	-0.019	-0.025	0.003	0.61	0.88	1	-0.042	0.17	-0.29	0.52	-0.11	-0.1	0.32	0.78	-0.14	0.26
differentiate -	0.016	-0.019	0.032	-0.031	-0.036	-0.042	1	-0.37	-0.009	-0.031	-0.022	0.009	-0.057	-0.029	-0.0045	-0.019
Grade -	-0.093	-0.057	-0.0043	0.11	0.15	0.17	-0.37	1	-0.044	0.1	-0.19	-0.18	0.083	0.12	-0.058	0.13
A Stage -	0.021	0.007	0.0049	-0.22	-0.26	-0.29	-0.009	-0.044	1	-0.12	0.066	0.027	-0.069	-0.23	0.07	-0.097
Tumor Size -	-0.078	-0.0037	0.01	0.82	0.28	0.52	-0.031	0.1	-0.12	1	-0.059	-0.072	0.11	0.25	-0.087	0.14
Estrogen Status -	0.06	0.058	-0.018	-0.061	-0.1	-0.11	-0.022	-0.19	0.066	-0.059	1	0.51	-0.045	-0.085	0.13	-0.18
Progesterone Status -	-0.021	0.034	-0.013	-0.058	-0.094	-0.1	0.009	-0.18	0.027	-0.072	0.51	1	-0.018	-0.078	0.096	-0.18
Regional Node Examined -	-0.034	-0.0044	-0.0042	0.11	0.33	0.32	-0.057	0.083	-0.069	0.11	-0.045	-0.018	1	0.41	-0.022	0.035
Reginol Node Positive -	0.013	-0.015	0.0058	0.24	0.84	0.78	-0.029	0.12	-0.23	0.25	-0.085	-0.078	0.41	1	-0.14	0.26
Survival Months -	-0.0094	0.041	-0.026	-0.086	-0.14	-0.14	-0.0045	-0.058	0.07	-0.087	0.13	0.096	-0.022	-0.14	1	-0.48
Status -	0.056	-0.05	0.033	0.15	0.26	0.26	-0.019	0.13	-0.097	0.14	-0.18	-0.18	0.035	0.26	-0.48	1
	Age -	Race -	Marital Status -	T Stage -	N Stage -	6th Stage -	differentiate -	Grade -	A Stage -	Tumor Size -	Estrogen Status -	Progesterone Status -	gional Node Examined -	Reginol Node Positive -	Survival Months -	Status -

- 0.8 - 0.6 - 0.2 - 0.0 - -0.4



 Higher N stage corresponds to higher number of regional nodes positive

#### **Handling imbalance Dataset: Using SMOTE**

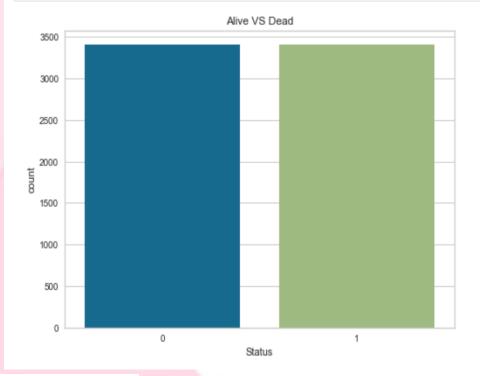
1000

500

```
1 smote = SMOTE(random_state=25)
 2 smote_X_train, smote_Y_train = smote.fit_resample(X,y)
 1 sns.countplot(x=y, data=df)
 3 plt.show()
✓ 0.1s
                           Alive VS Dead
 3500
 3000
 2500
 2000
 1500
```

Status





#### **Train Test Split:**

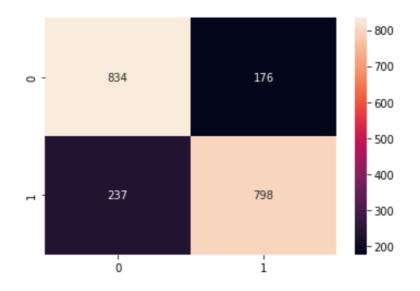
# **Standard Scaler:**

## Algorithm used:

- Logistic Regression
- Naïve Bayes
- KNN
- SVC

70 - 30

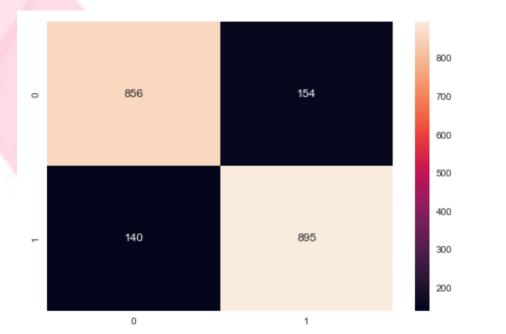
#### **LOGISTIC REGRESSION:**



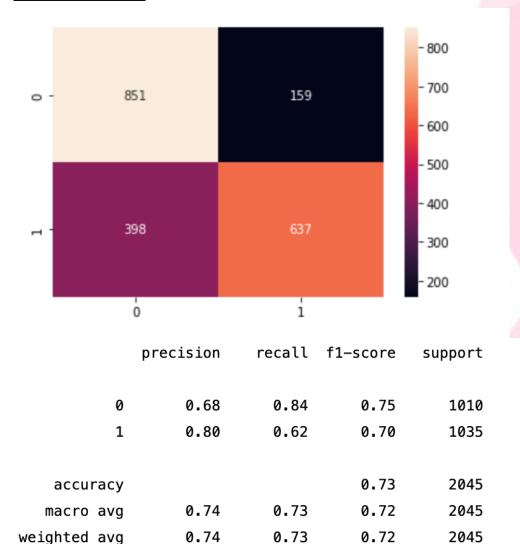
	precision	recall	f1-score	support	
0	0.78	0.83	0.80	1010	
1	0.82	0.77	0.79	1035	
			0.00	2045	
accuracy			0.80	2045	
macro avg	0.80	0.80	0.80	2045	
weighted avg	0.80	0.80	0.80	2045	

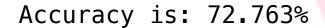
Accuracy is: 79.804%

Accuracy is: 80.507%



#### **NAIVE BAYES:**







Fitting 10 folds for each of 100 candidates, totalling 1000 fits
Tuned Hyperparameters : {'var\_smoothing': 0.006579332246575682}
Accuracy is: 73.821%

	precision	recall	f1-score	support
0	0.68	0.84	0.75	1010
1	0.80	0.61	0.69	1035
26645264			0.73	2045
accuracy			0.73	2045
macro avg	0.74	0.73	0.72	2045
weighted avg	0.74	0.73	0.72	2045

#### KNN:

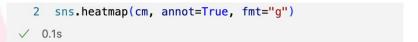
Accuracy is: 83.28%

#### <AxesSubplot:>

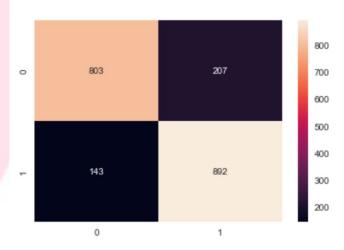


/	1	<pre>print(classification_report(y_test, y_pred))</pre>
]	✓ 0	.2s

support	f1-score	recall	precision	
1010	0.82	0.80	0.85	0
1035	0.84	0.86	0.81	1
2045	0.83			accuracy
2045	0.83	0.83	0.83	macro avg
2045	0.83	0.83	0.83	weighted avg



#### <AxesSubplot:>



#### 

	precision	recall	f1-score	support
0	0.85	0.80	0.82	1010
1	0.81	0.86	0.84	1035
accuracy			0.83	2045
macro avg	0.83	0.83	0.83	2045
weighted avg	0.83	0.83	0.83	2045

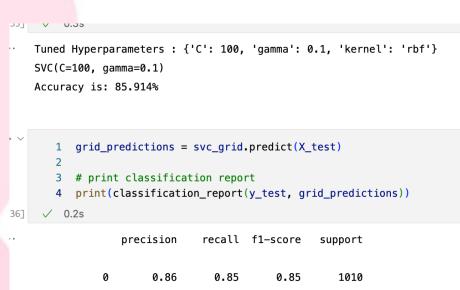
Accuracy for our training dataset with tuning is: 86.59%

# SVC:



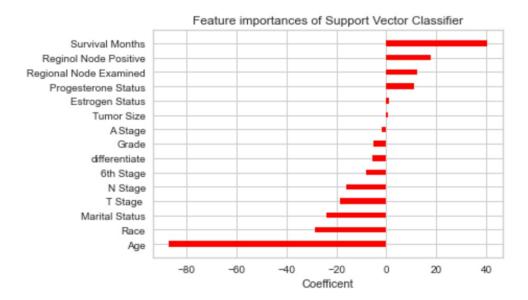
<pre>1 print(classification_report(y_test, y_pred))</pre>										
	precision	recall	f1-score	support						
0	0.80	0.87	0.84	1010						
1	0.87	0.79	0.83	1035						
accuracy			0.83	2045						
macro avg	0.84	0.83	0.83	2045						
weighted avg	0.84	0.83	0.83	2045						

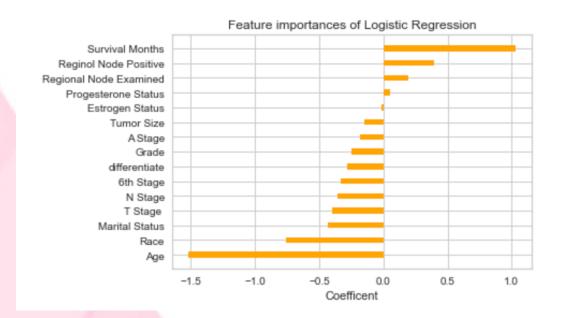
Accuracy is: 83.276%



•		precision	recall	f1–score	support
	0	0.86	0.85	0.85	1010
	1	0.85	0.86	0.86	1035
	accuracy			0.86	2045
	macro avg	0.86	0.86	0.86	2045
	weighted avg	0.86	0.86	0.86	2045

# Feature Importance:





## **Conclusion:**

- Logistic Regression , KNN & SVC are the algorithms to consider
- No pattern detected during EDA
- Our Data might not be a high Quality
  - Need a larger Data