



# Thoracic Surgery

For Lung Cancer Patients

By: Seung Chi



# Problem

- Patients who receive thoracic surgery for lung cancer do so with the expectation that their lives will be prolonged for a sufficient amount of time afterwards.
- The problem to solve is whether there is a way to determine postoperative 1 year survival of lung cancer patients utilizing the patient attributes in the data set.

# Who benefits from answering this problem?

- Patients
- Families of Patients
- Physicians
- Hospitals
- Healthcare Organizations



# Data Set



- Original from UCI Machine Learning Repository
  - Collected retrospectively at Wroclaw Thoracic Surgery Centre for patients who underwent major lung resections for primary lung cancer in the years 2007-2011
  - 470 instances and no missing values
- This report consists of 454 patient data.
  - Excluding outliers from FEV1 and Age columns



# Descriptions of Attributes (1)

| Attribute          | Description   |
|--------------------|---|
| <b>Diagnosis</b>   | ICD-10 codes for primary and secondary as well multiple tumors if any                               |
| <b>FVC</b>         | Amount of air which can be forcibly exhaled from the lungs after taking the deepest breath possible |
| <b>FEV1</b>        | Volume that has been exhaled at the end of the first second of forced expiration                    |
| <b>Performance</b> | Performance status on Zubrod scale, Good (0) to Poor (2)  |
| <b>Pain</b>        | Pain, prior to surgery (T = 1, F = 0)   |
| <b>Haemoptysis</b> | Coughing up blood, prior to surgery (T = 1, F = 0)  |
| <b>Dyspnoea</b>    | Difficult or labored breathing, prior to surgery (T = 1, F = 0)                                     |
| <b>Cough</b>       | Cough, prior to surgery (T = 1, F = 0)  |

# Descriptions of Attributes (2)

| Attribute                | Description   |
|--------------------------|---|
| <b>Weakness</b>          | Weakness, prior to surgery (T = 1, F = 0)                                   |
| <b>Tumor_Size</b>        | T in clinical TNM - size of the original tumor, 1 (smallest) to 4 (largest) |
| <b>Diabetes_Mellitus</b> | Type 2 diabetes mellitus (T = 1, F = 0)                                     |
| <b>MI_6mo</b>            | Myocardial Infarction (Heart Attack) up to 6 months prior (T = 1, F = 0)    |
| <b>PAD</b>               | Peripheral arterial diseases (T = 1, F = 0)                                 |
| <b>Smoking</b>           | Smoking (T = 1, F = 0)  |
| <b>Asthma</b>            | Asthma (T = 1, F = 0)   |
| <b>Age</b>               | Age at surgery  |
| <b>Death_1yr</b>         | 1 year survival period - (T) value if died (T = 1, F = 0)                   |

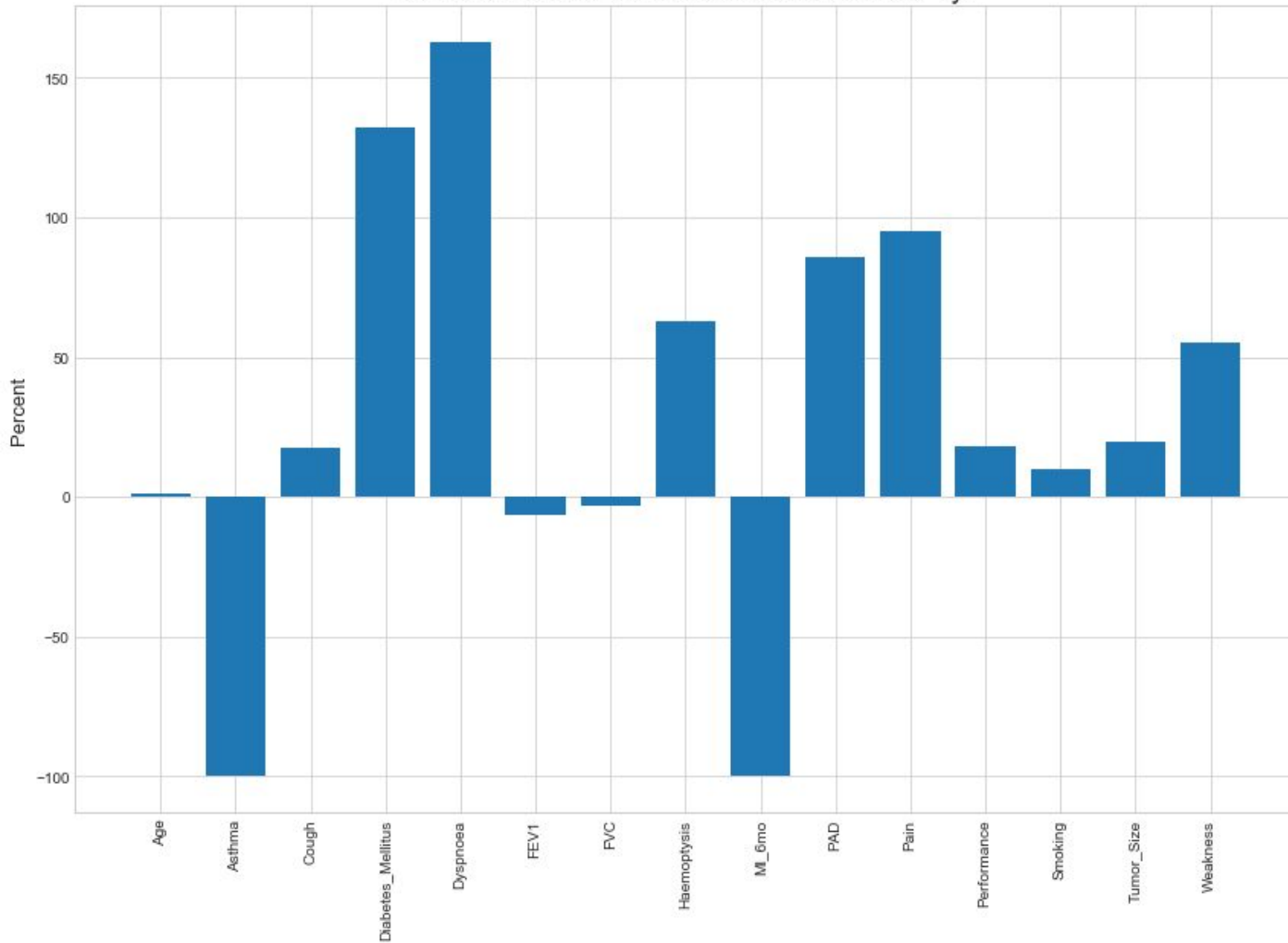
# Difference between 1 year death and live patients

- 69 death out of 454; 15.20% death rate in 1 year post-op.

| Attribute   | Death in 1 year<br>(Mean) | Live 1 year<br>(Mean) |
|-------------|---------------------------|-----------------------|
| FVC         | 3.195072                  | <b>3.304597</b>       |
| FEV1        | 2.383188                  | <b>2.540805</b>       |
| Performance | <b>0.913043</b>           | 0.774026              |
| Pain        | <b>0.101449</b>           | 0.051948              |
| Haemoptysis | <b>0.202899</b>           | 0.124675              |
| Dyspnoea    | <b>0.115942</b>           | 0.044156              |
| Cough       | <b>0.797101</b>           | 0.677922              |
| Weakness    | <b>0.246377</b>           | 0.158442              |

| Attribute         | Death in 1 year<br>(Mean) | Live 1 year<br>(Mean) |
|-------------------|---------------------------|-----------------------|
| Tumor_Size        | <b>2.014493</b>           | 1.683117              |
| Diabetes_Mellitus | <b>0.144928</b>           | 0.062338              |
| MI_6mo            | 0.000000                  | <b>0.005195</b>       |
| PAD               | <b>0.028986</b>           | 0.015584              |
| Smoking           | <b>0.898551</b>           | 0.815584              |
| Asthma            | 0.000000                  | <b>0.005195</b>       |

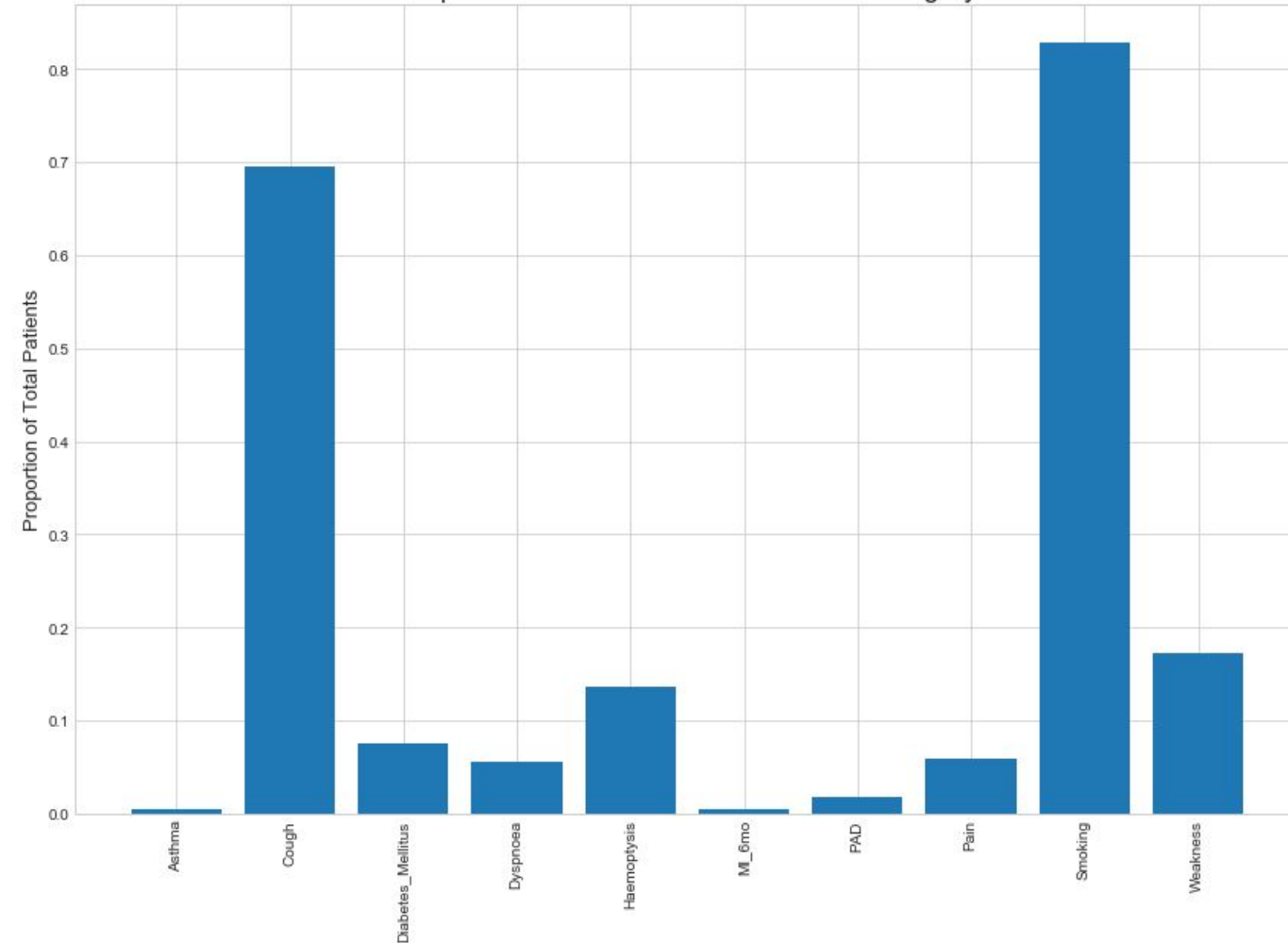
Mean Difference % between Dead and Live 1yr



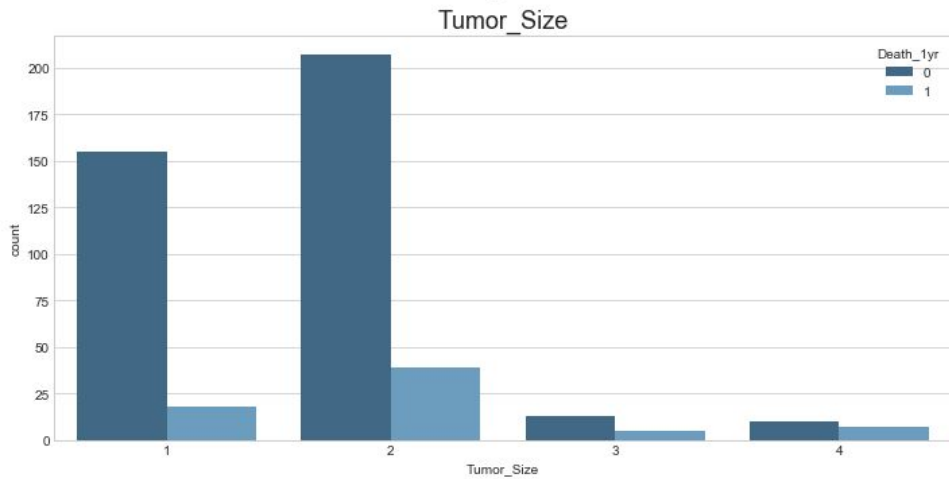
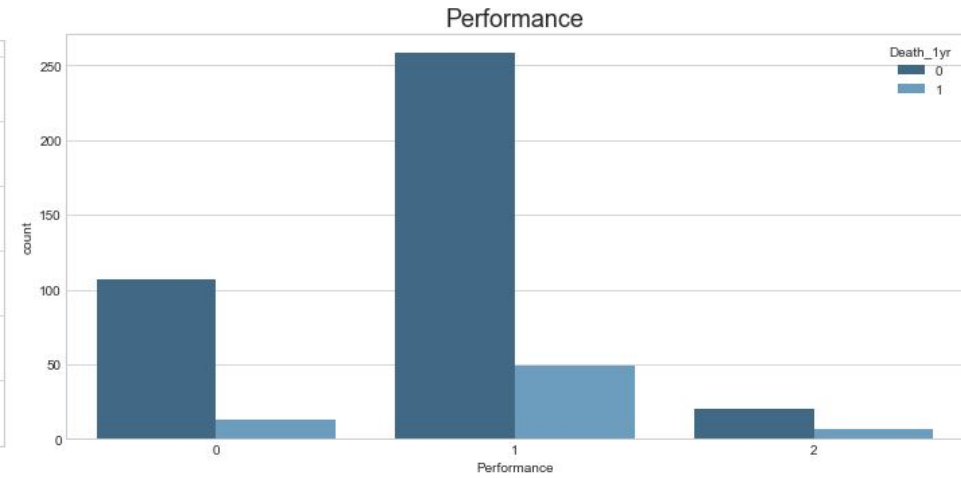
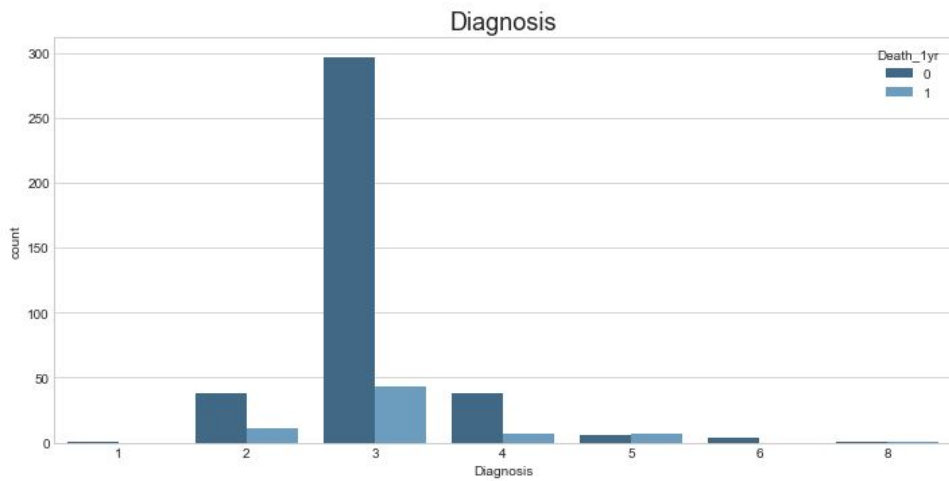
Mean  
Difference %  
of Dead and  
Live (1 yr)



Proportion of Patient Conditions before Surgery



Proportion  
Of Patient  
Conditions  
before  
Surgery



# Categorical Data

# Hypothesis Testing

- Null Hypothesis: The 1 year death and live patients have the same mean, tested for each attribute.
- Test Statistic: Mean difference between death and live patients
- Significance level: 0.05



# Results of Hypothesis Testing

| Attribute          | P value       |
|--------------------|---------------|
| FVC                | 0.1706        |
| FEV1               | 0.0588        |
| <b>Performance</b> | <b>0.0300</b> |
| Pain               | 0.0964        |
| Haemoptysis        | 0.0623        |
| <b>Dyspnoea</b>    | <b>0.0242</b> |
| <b>Cough</b>       | <b>0.0320</b> |

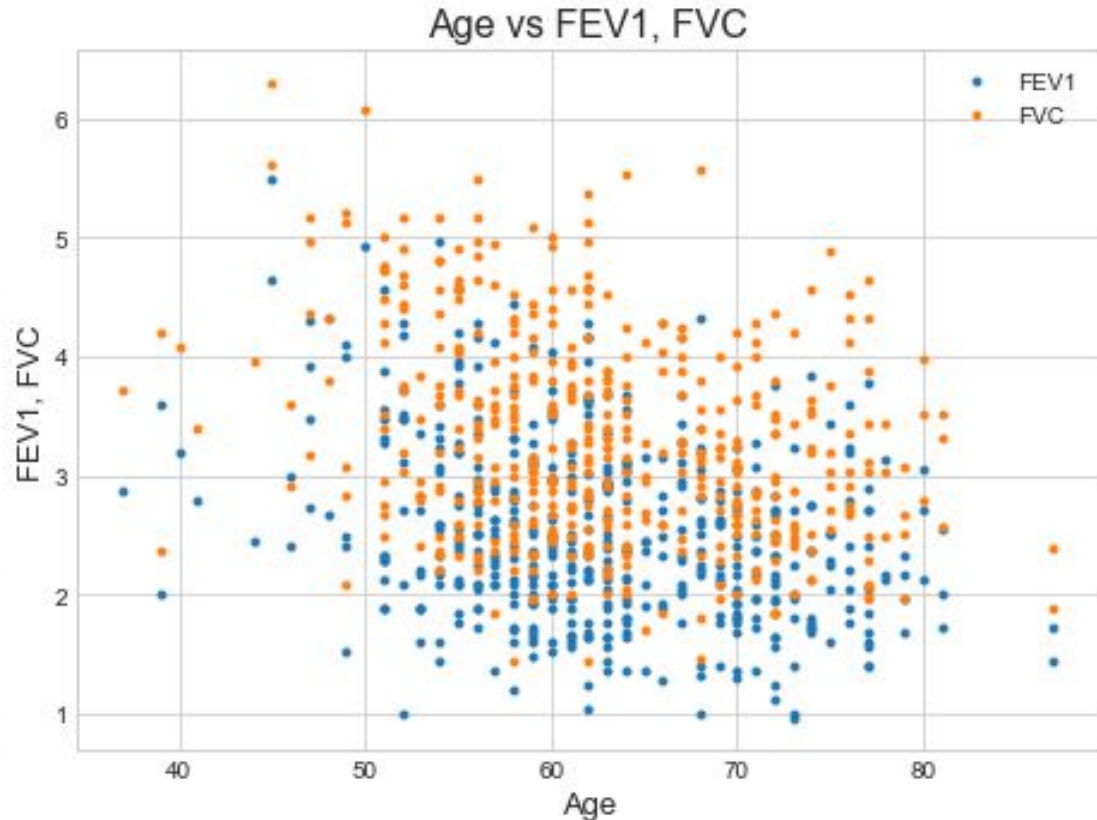
|                          |               |
|--------------------------|---------------|
| Weakness                 | 0.0606        |
| <b>Tumor_Size</b>        | <b>0.0003</b> |
| <b>Diabetes_Mellitus</b> | <b>0.0209</b> |
| MI_6mo                   | 0.7264        |
| PAD                      | 0.3498        |
| Smoking                  | 0.0581        |
| Asthma                   | 0.7178        |
| Age                      | 0.2714        |

# Mean difference % for Attributes of Significance

- Performance = 17.96%
- Dyspnoea = 162.57%
- Cough = 17.58%
- Tumor\_Size = 19.69%
- Diabetes\_Mellitus = 132.49%



# Correlations of Numerical (Age, FVC, FEV1) Data



Correlation Coefficients:

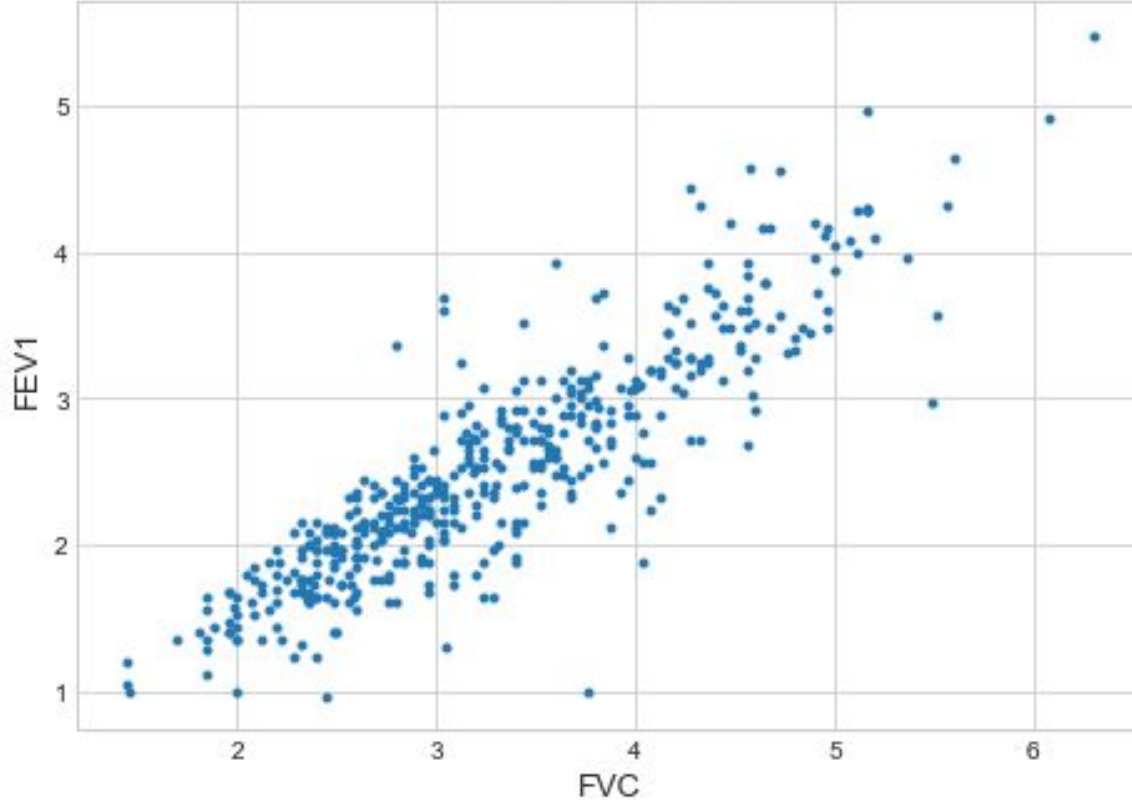
➤ Age & FEV1

○ -0.2994

➤ Age & FVC

○ -0.3096

FVC vs FEV1



# Correlation of FVC and FEV1

Correlation Coefficient:

➤ 0.8875

FEV1/FVC Ratio:

➤ Used in diagnosis of  
obstructive and  
restrictive lung disease

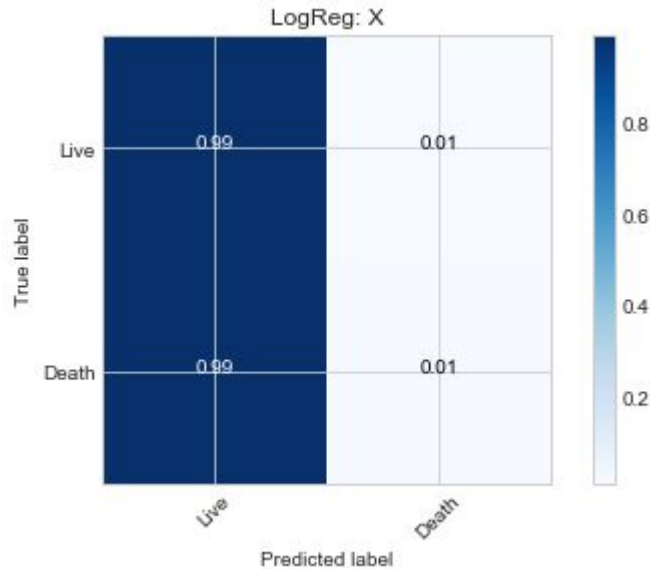
# Machine Learning - Supervised Classification

- Target Variable: Death\_1yr
- X
  - Drops target variable, MI\_6mo, Asthma
- X2
  - Attributes of significance from Hypothesis testing
  - Performance, Dyspnoea, Cough, Tumor\_Size, Diabetes\_Mellitus

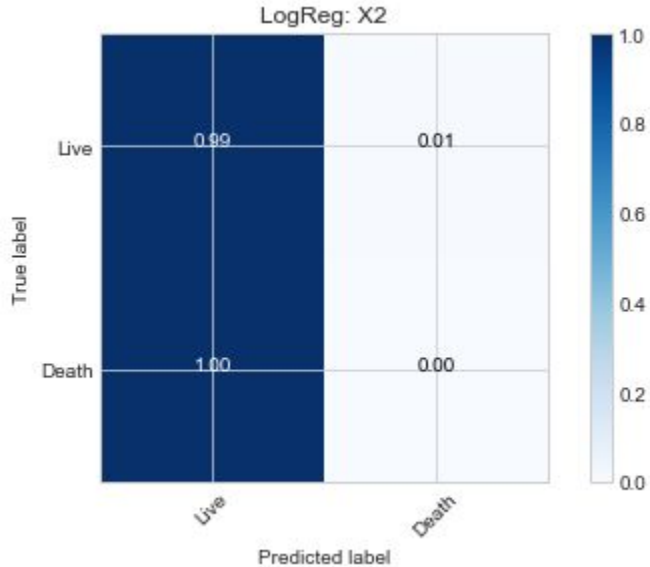




# Logistic Regression

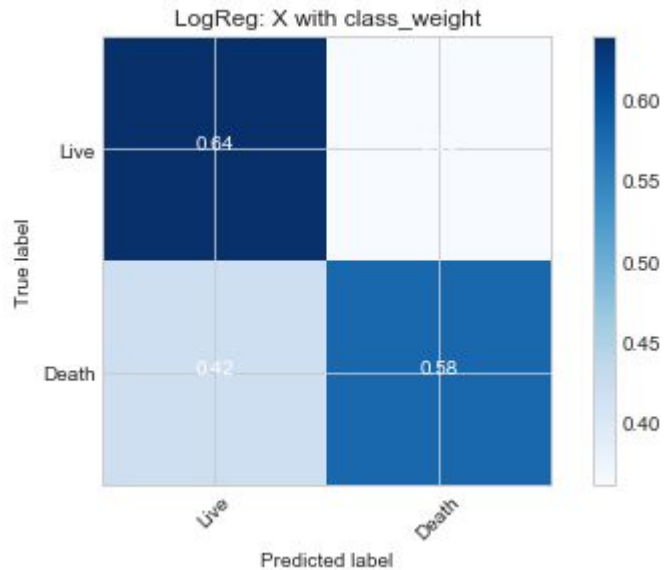


Accuracy score: 84.36%

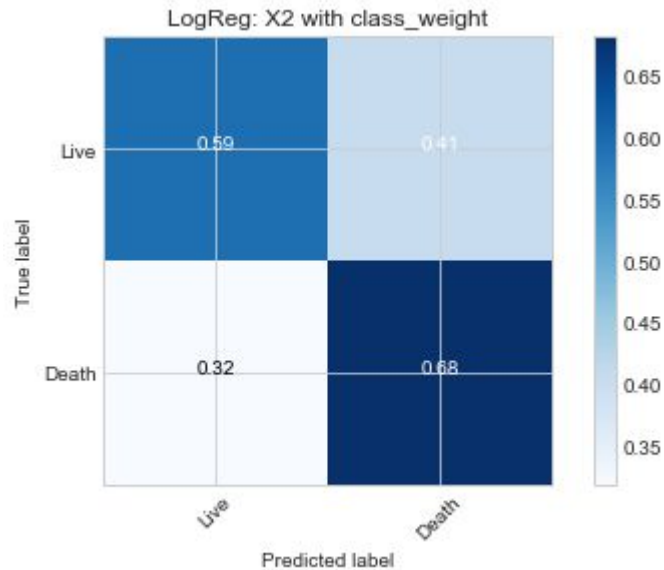


Accuracy score: 84.36%

# Logistic Regression w/ Class Weight

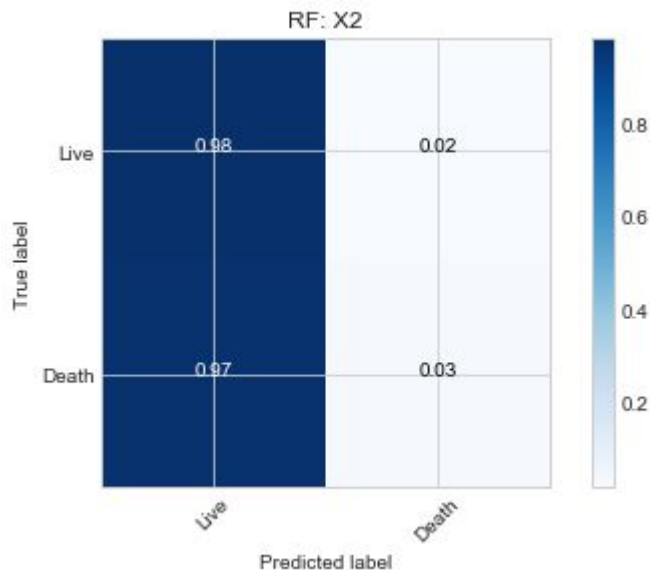


Accuracy score: 63.00%

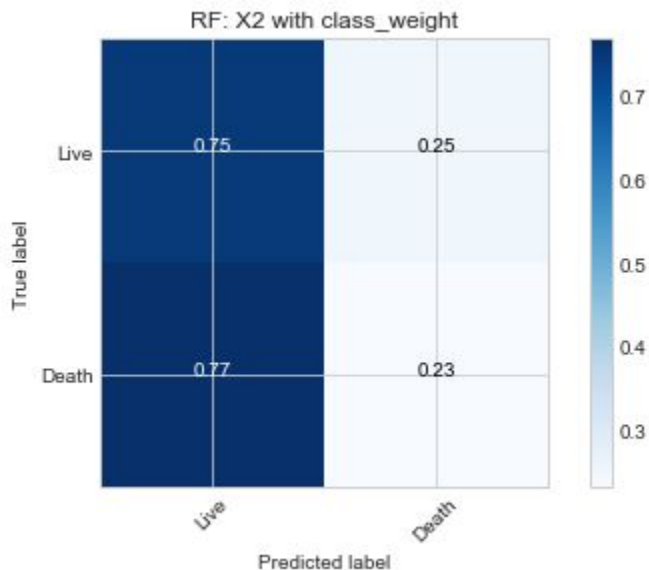


Accuracy score: 60.79%

# Random Forest Classifier (X2 Data)



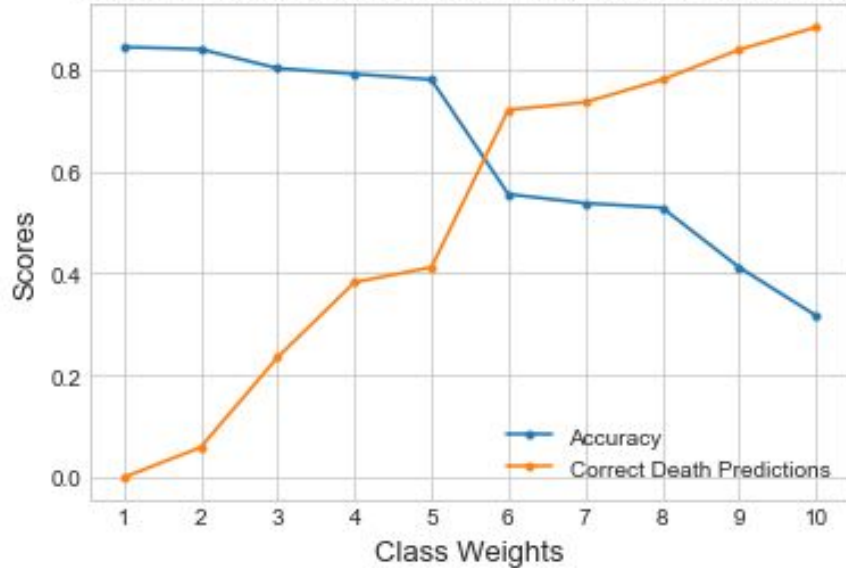
Accuracy score: 83.70%



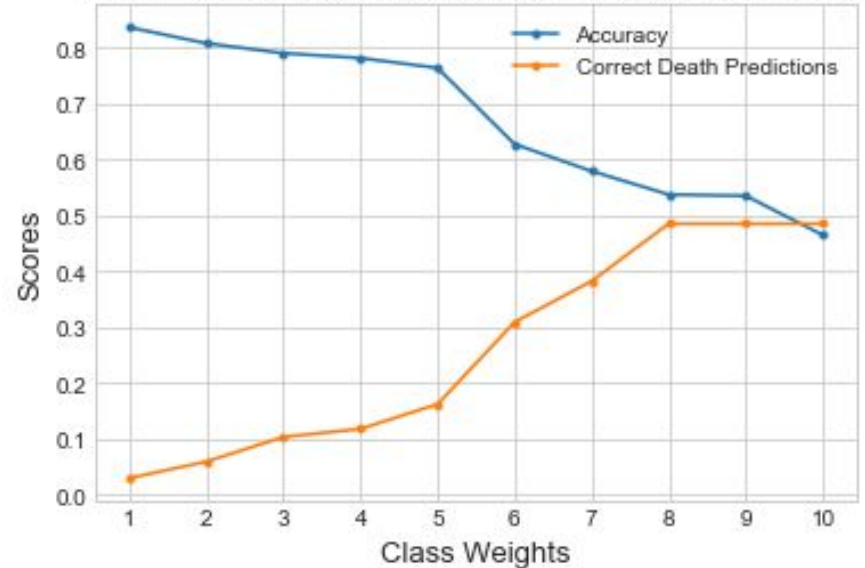
Accuracy score: 66.96%

# Accuracy vs Correct Death Predictions

Class Weights Influence on Log Regression of X2



Class Weights Influence on Random Forest of X2



# Beyond this project...

- Decision on desired outcome  
considering false prediction costs
- More data
- Hyperparameter tuning
- Ensemble method with preferred  
model

