

Breast Cancer Jump Scare Team 5

Cadeem Musgrove, Lucas Perez,
Jaskirat Singh & Aditya Singh



Agenda

This study aims to create models to determine overall survivability and predict clinical data for clinicians

About the Data	Data sources, time range, size, and end to end data flow from source to consumption with components and connectors describing how data evolves through the flow to generate end visuals and for modelling
Exploratory Analysis 1	Modeling Experiment, Evaluation
Exploratory Analysis 2	Modeling Experiment, Evaluation
Exploratory Analysis 3	Modeling Experiment, Evaluation
Considerations & Key Learning	Discussing our key learnings, future implementations, and our key learnings

About the Data

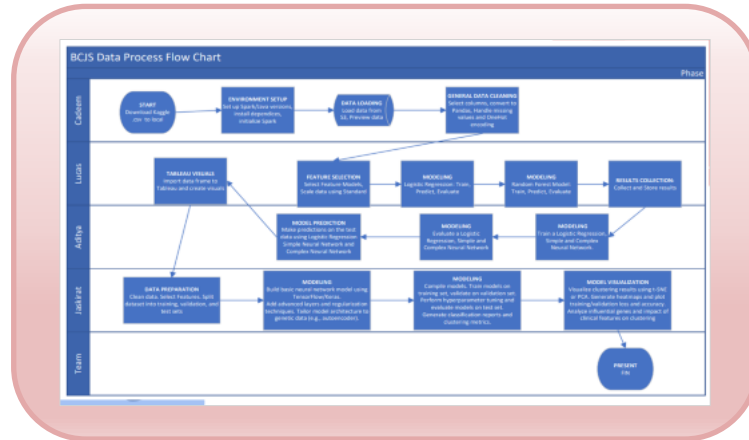


Data Description

Sourced from Kaggle. Data is 1904 Rows x 676 Columns. Data is taken between 2016 to 2020



Data Process



BCJS Data Process Flow Chart

Phase

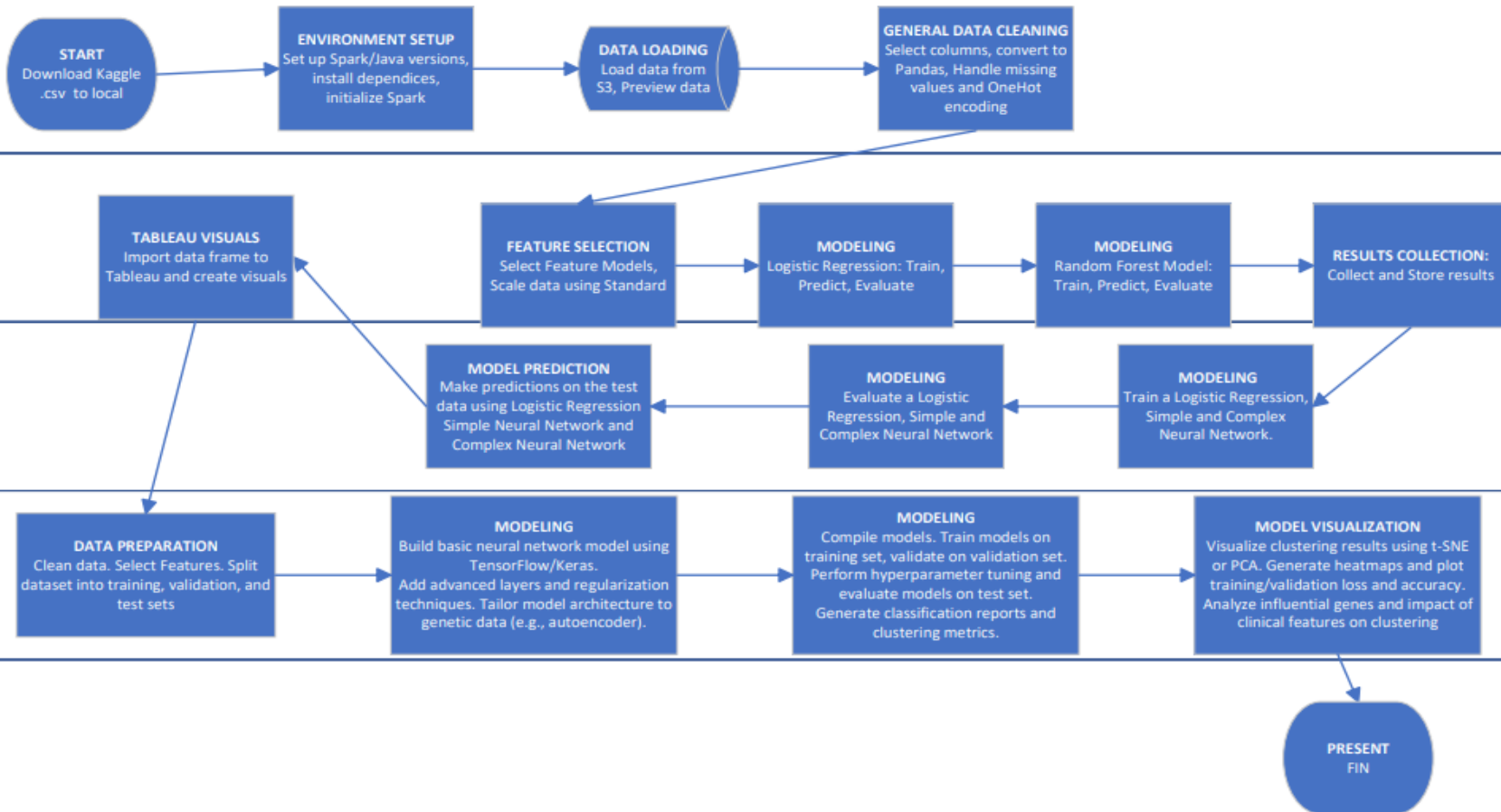
Cadeem

Lucas

Aditya

Jaskirat

Team



Exploratory Analysis 1

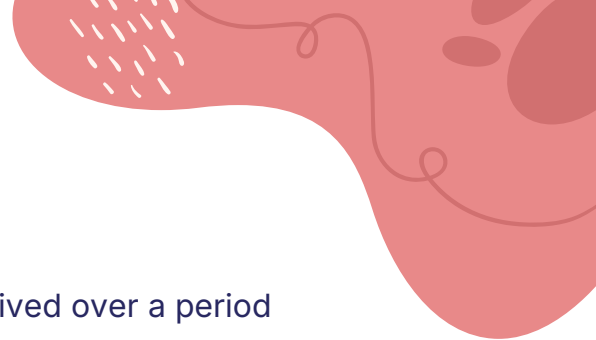
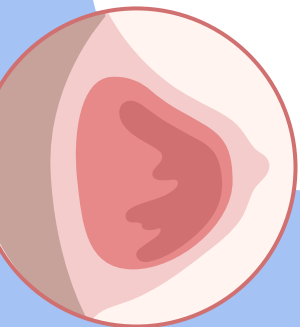
Key Insights on Target Variables

Target variable 'overall_survival' is binary. Indicates whether a patient survived over a period of time.

Modeling Experiment Design

Baseline model: Logistic Regression Model.

Random Survival Forest (RSF) Model employed to handle survivability over time (time-to-event data)



Exploratory Analysis 1

Model Evaluation

Metrics Used

- Accuracy: Primary metric in Logistic Regression
- Survival Probability Metric: RSF Model different time points (ex. 365 days)

Model Comparison

The RSF model of survival data shows better generalization and performance compared to the baseline logistic regression model.

Logistic Regression Model Accuracy

```
print("Classification Report")
print(classification_report(y_test, y_pred))
```

	precision	recall	f1-score	support
0.0	0.57	0.45	0.50	134
1.0	0.56	0.58	0.57	95
accuracy	0.51	0.51	0.52	229
macro avg	0.51	0.51	0.51	229
weighted avg	0.51	0.52	0.52	229

Snippets of Code

Random Survival Forest Predictions

```
1) patients_survival_info.head()
```

	patient_id	survival probability at 300 days	survival probability at 365 days	survival probability at 545 days	survival probability at 750 days
0	0	0.988571	0.878471	0.216091	0.088695
1	1	0.987486	0.921258	0.811910	0.128809
2	2	0.998470	0.987444	0.764962	0.325275
3	3	1.000000	0.983698	0.909005	0.483045
4	4	0.987808	0.828850	0.488934	0.141528

Logistic Regression Model Accuracy

```
print("Classification Report")
print(classification_report(y_test, y_pred))
```

Classification Report

	precision	recall	f1-score	support
0.0	0.67	0.65	0.66	124
1.0	0.56	0.58	0.57	95
accuracy			0.62	219
macro avg	0.61	0.61	0.61	219
weighted avg	0.62	0.62	0.62	219

Random Survival Forest Predictions

```
71] patients_survive_info_df.head()
```

	patient_id	survival probability at 180 days	survival probability at 365 days	survival probability at 545 days	survival probability at 730 days
0	0	0.998571	0.879471	0.216091	0.085850
1	1	0.987496	0.921208	0.811910	0.128509
2	2	0.999470	0.987444	0.764962	0.325275
3	3	1.000000	0.983698	0.939805	0.483045
4	4	0.967509	0.829850	0.495934	0.141528

Next steps: [Generate code with patients_survive_info_df](#)

[View recommended plots](#)

[New interactive sheet](#)

Exploratory Analysis 2

Key Insights

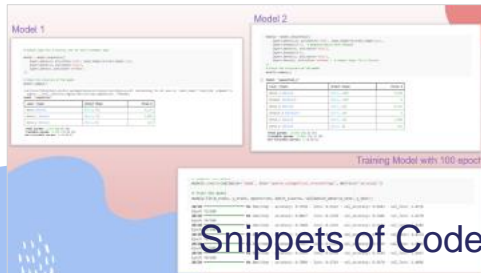
Type of cancer and type of treatment, which treatment is more effective against certain types of cancers.

Model Design

- Created Dummies for testing predictions.
- Ran Logistics Regression as the measurement is categorical

Model evaluation

- Ran with different epochs, neurons, and layers.
- To achieve 56% accuracy.



Model 1

```
Model 1: Logistic Regression with 100 epochs
Accuracy: 0.56
```

Model 2

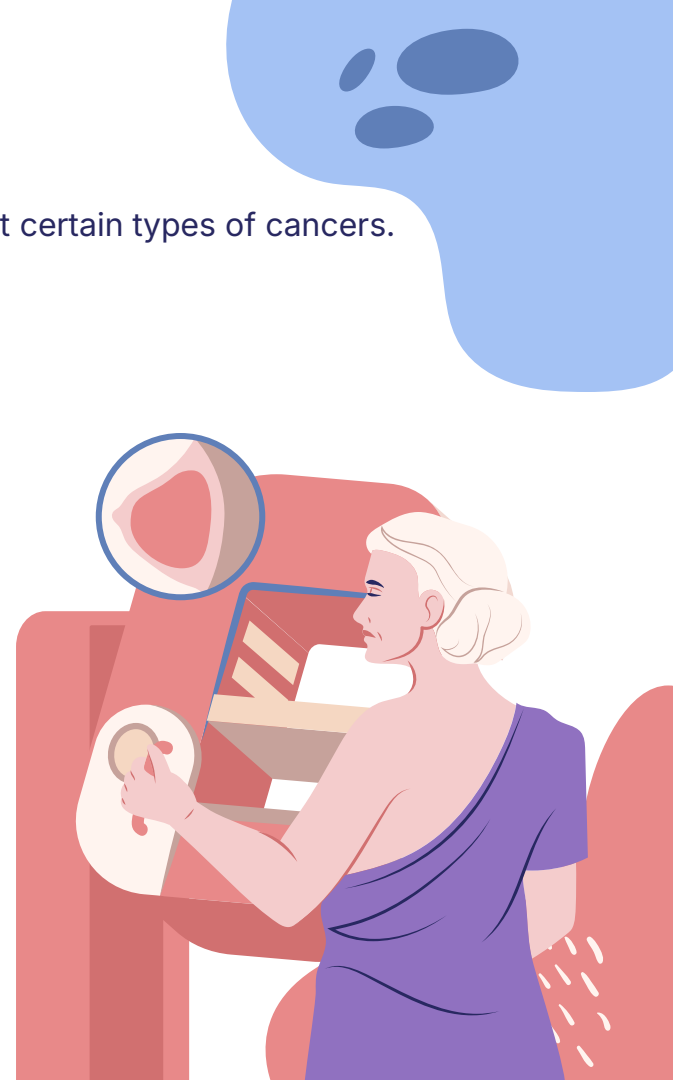
```
Model 2: Logistic Regression with 100 epochs
Accuracy: 0.56
```

Training Model with 100 epochs

```
Epoch 1: Loss: 0.693, Accuracy: 0.500
Epoch 2: Loss: 0.693, Accuracy: 0.500
Epoch 3: Loss: 0.693, Accuracy: 0.500
Epoch 4: Loss: 0.693, Accuracy: 0.500
Epoch 5: Loss: 0.693, Accuracy: 0.500
Epoch 6: Loss: 0.693, Accuracy: 0.500
Epoch 7: Loss: 0.693, Accuracy: 0.500
Epoch 8: Loss: 0.693, Accuracy: 0.500
Epoch 9: Loss: 0.693, Accuracy: 0.500
Epoch 10: Loss: 0.693, Accuracy: 0.500
Epoch 11: Loss: 0.693, Accuracy: 0.500
Epoch 12: Loss: 0.693, Accuracy: 0.500
Epoch 13: Loss: 0.693, Accuracy: 0.500
Epoch 14: Loss: 0.693, Accuracy: 0.500
Epoch 15: Loss: 0.693, Accuracy: 0.500
Epoch 16: Loss: 0.693, Accuracy: 0.500
Epoch 17: Loss: 0.693, Accuracy: 0.500
Epoch 18: Loss: 0.693, Accuracy: 0.500
Epoch 19: Loss: 0.693, Accuracy: 0.500
Epoch 20: Loss: 0.693, Accuracy: 0.500
Epoch 21: Loss: 0.693, Accuracy: 0.500
Epoch 22: Loss: 0.693, Accuracy: 0.500
Epoch 23: Loss: 0.693, Accuracy: 0.500
Epoch 24: Loss: 0.693, Accuracy: 0.500
Epoch 25: Loss: 0.693, Accuracy: 0.500
Epoch 26: Loss: 0.693, Accuracy: 0.500
Epoch 27: Loss: 0.693, Accuracy: 0.500
Epoch 28: Loss: 0.693, Accuracy: 0.500
Epoch 29: Loss: 0.693, Accuracy: 0.500
Epoch 30: Loss: 0.693, Accuracy: 0.500
Epoch 31: Loss: 0.693, Accuracy: 0.500
Epoch 32: Loss: 0.693, Accuracy: 0.500
Epoch 33: Loss: 0.693, Accuracy: 0.500
Epoch 34: Loss: 0.693, Accuracy: 0.500
Epoch 35: Loss: 0.693, Accuracy: 0.500
Epoch 36: Loss: 0.693, Accuracy: 0.500
Epoch 37: Loss: 0.693, Accuracy: 0.500
Epoch 38: Loss: 0.693, Accuracy: 0.500
Epoch 39: Loss: 0.693, Accuracy: 0.500
Epoch 40: Loss: 0.693, Accuracy: 0.500
Epoch 41: Loss: 0.693, Accuracy: 0.500
Epoch 42: Loss: 0.693, Accuracy: 0.500
Epoch 43: Loss: 0.693, Accuracy: 0.500
Epoch 44: Loss: 0.693, Accuracy: 0.500
Epoch 45: Loss: 0.693, Accuracy: 0.500
Epoch 46: Loss: 0.693, Accuracy: 0.500
Epoch 47: Loss: 0.693, Accuracy: 0.500
Epoch 48: Loss: 0.693, Accuracy: 0.500
Epoch 49: Loss: 0.693, Accuracy: 0.500
Epoch 50: Loss: 0.693, Accuracy: 0.500
Epoch 51: Loss: 0.693, Accuracy: 0.500
Epoch 52: Loss: 0.693, Accuracy: 0.500
Epoch 53: Loss: 0.693, Accuracy: 0.500
Epoch 54: Loss: 0.693, Accuracy: 0.500
Epoch 55: Loss: 0.693, Accuracy: 0.500
Epoch 56: Loss: 0.693, Accuracy: 0.500
Epoch 57: Loss: 0.693, Accuracy: 0.500
Epoch 58: Loss: 0.693, Accuracy: 0.500
Epoch 59: Loss: 0.693, Accuracy: 0.500
Epoch 60: Loss: 0.693, Accuracy: 0.500
Epoch 61: Loss: 0.693, Accuracy: 0.500
Epoch 62: Loss: 0.693, Accuracy: 0.500
Epoch 63: Loss: 0.693, Accuracy: 0.500
Epoch 64: Loss: 0.693, Accuracy: 0.500
Epoch 65: Loss: 0.693, Accuracy: 0.500
Epoch 66: Loss: 0.693, Accuracy: 0.500
Epoch 67: Loss: 0.693, Accuracy: 0.500
Epoch 68: Loss: 0.693, Accuracy: 0.500
Epoch 69: Loss: 0.693, Accuracy: 0.500
Epoch 70: Loss: 0.693, Accuracy: 0.500
Epoch 71: Loss: 0.693, Accuracy: 0.500
Epoch 72: Loss: 0.693, Accuracy: 0.500
Epoch 73: Loss: 0.693, Accuracy: 0.500
Epoch 74: Loss: 0.693, Accuracy: 0.500
Epoch 75: Loss: 0.693, Accuracy: 0.500
Epoch 76: Loss: 0.693, Accuracy: 0.500
Epoch 77: Loss: 0.693, Accuracy: 0.500
Epoch 78: Loss: 0.693, Accuracy: 0.500
Epoch 79: Loss: 0.693, Accuracy: 0.500
Epoch 80: Loss: 0.693, Accuracy: 0.500
Epoch 81: Loss: 0.693, Accuracy: 0.500
Epoch 82: Loss: 0.693, Accuracy: 0.500
Epoch 83: Loss: 0.693, Accuracy: 0.500
Epoch 84: Loss: 0.693, Accuracy: 0.500
Epoch 85: Loss: 0.693, Accuracy: 0.500
Epoch 86: Loss: 0.693, Accuracy: 0.500
Epoch 87: Loss: 0.693, Accuracy: 0.500
Epoch 88: Loss: 0.693, Accuracy: 0.500
Epoch 89: Loss: 0.693, Accuracy: 0.500
Epoch 90: Loss: 0.693, Accuracy: 0.500
Epoch 91: Loss: 0.693, Accuracy: 0.500
Epoch 92: Loss: 0.693, Accuracy: 0.500
Epoch 93: Loss: 0.693, Accuracy: 0.500
Epoch 94: Loss: 0.693, Accuracy: 0.500
Epoch 95: Loss: 0.693, Accuracy: 0.500
Epoch 96: Loss: 0.693, Accuracy: 0.500
Epoch 97: Loss: 0.693, Accuracy: 0.500
Epoch 98: Loss: 0.693, Accuracy: 0.500
Epoch 99: Loss: 0.693, Accuracy: 0.500
Epoch 100: Loss: 0.693, Accuracy: 0.500
```

Snippets of Code

Tableau Visuals



Model 1

```
# Output layer has 8 neurons, one for each treatment type

model1 = models.Sequential([
    layers.Dense(64, activation='relu', input_shape=X_train.shape[1:]),
    layers.Dense(32, activation='relu'),
    layers.Dense(8, activation='softmax')
])

# Check the structure of the model
model1.summary()
```

`/usr/local/lib/python3.10/dist-packages/keras/src/layers/core/dense.py:87: UserWarning: Do not pass an 'input_shape'/'input_dim' argument to super().__init__(activity_regularizer=activity_regularizer, **kwargs)`

Model: "sequential"

Layer (type)	Output Shape	Param #
dense_3 (Dense)	(None, 64)	4,224
dense_1 (Dense)	(None, 32)	2,080
dense_2 (Dense)	(None, 8)	264

Total params: 6,568 (25.66 KB)
Trainable params: 6,568 (25.66 KB)
Non-trainable params: 0 (0.00 B)

Model 2

```
model2 = models.Sequential([
    layers.Dense(128, activation='relu', input_shape=X_train.shape[1:]),
    layers.Dropout(0.5), # Regularization with Dropout
    layers.Dense(64, activation='relu'),
    layers.Dropout(0.5),
    layers.Dense(32, activation='relu'),
    layers.Dense(8, activation='softmax') # Output layer for 8 classes
])

# Check the structure of the model
model2.summary()
```

Model: "sequential_1"

Layer (type)	Output Shape	Param #
dense_3 (Dense)	(None, 128)	8,448
dropout (Dropout)	(None, 128)	0
dense_4 (Dense)	(None, 64)	8,256
dropout_1 (Dropout)	(None, 64)	0
dense_5 (Dense)	(None, 32)	2,080
dense_6 (Dense)	(None, 8)	264

Total params: 19,048 (74.41 KB)
Trainable params: 19,048 (74.41 KB)
Non-trainable params: 0 (0.00 B)

Training Model with 100 epoch

```
# COMPILING THE MODEL
model2.compile(optimizer='adam', loss='sparse_categorical_crossentropy', metrics=['accuracy'])

# Train the model
model2.fit(X_train, y_train, epochs=100, batch_size=64, validation_data=(X_test, y_test))
```

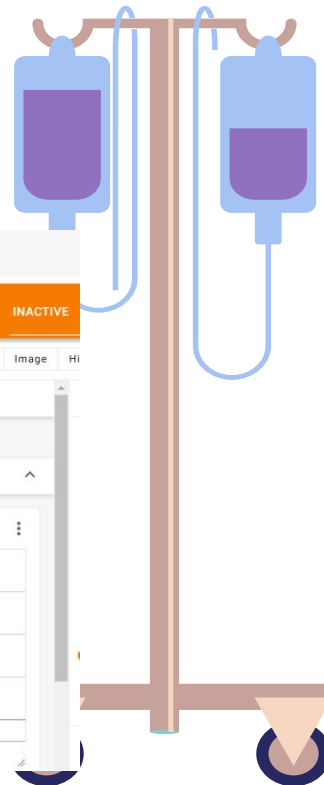
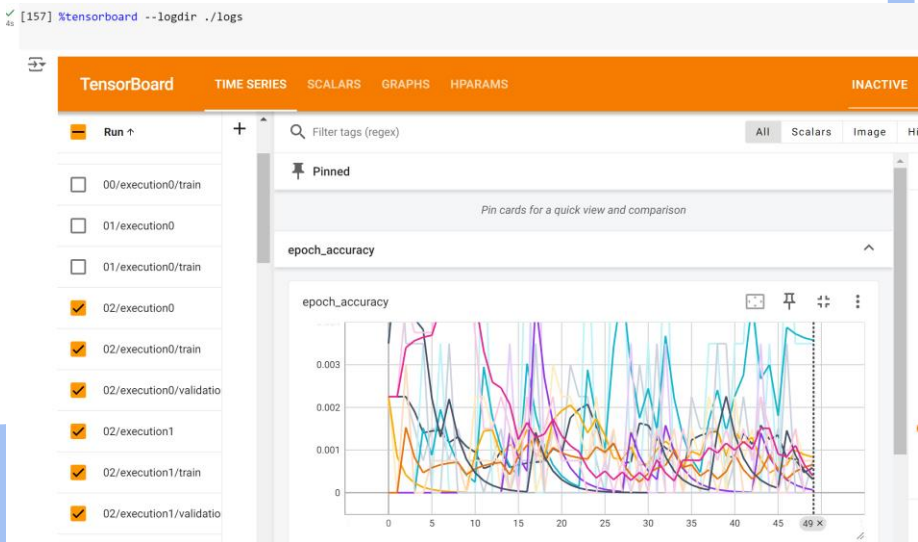
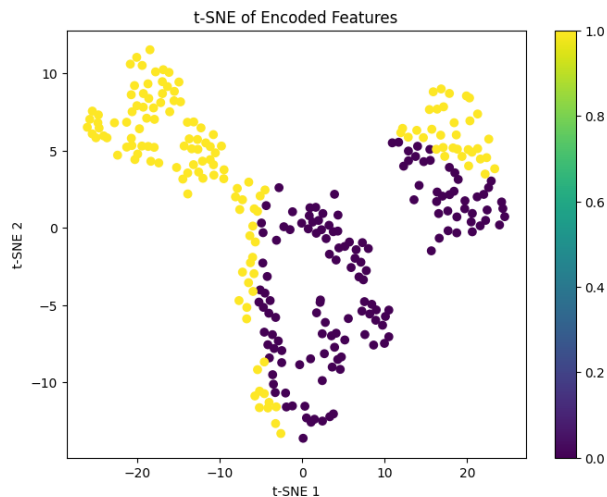
20/20 ————— 0s 5ms/step - accuracy: 0.7930 - loss: 0.5112 - val_accuracy: 0.5643 - val_loss: 1.4376
Epoch 73/100
20/20 ————— 0s 5ms/step - accuracy: 0.8047 - loss: 0.5278 - val_accuracy: 0.5486 - val_loss: 1.4270
Epoch 74/100
20/20 ————— 0s 5ms/step - accuracy: 0.7848 - loss: 0.5254 - val_accuracy: 0.5643 - val_loss: 1.4384
Epoch 75/100
20/20 ————— 0s 4ms/step - accuracy: 0.7992 - loss: 0.4972 - val_accuracy: 0.5580 - val_loss: 1.4372
Epoch 76/100
20/20 ————— 0s 5ms/step - accuracy: 0.7997 - loss: 0.5350 - val_accuracy: 0.5799 - val_loss: 1.4075
Epoch 77/100
20/20 ————— 0s 9ms/step - accuracy: 0.7959 - loss: 0.5370 - val_accuracy: 0.5549 - val_loss: 1.4267
Epoch 78/100
20/20 ————— 0s 8ms/step - accuracy: 0.7809 - loss: 0.5714 - val_accuracy: 0.5674 - val_loss: 1.4606

Exploratory Analysis 3

Key Insights

Metrics Used: Loss and Accuracy for models and clustering metrics (silhouette score and Davies Bouldin score)

Use of TensorBoard for following Model training.



Considerations & Key Lessons

Time

Horizon of the study

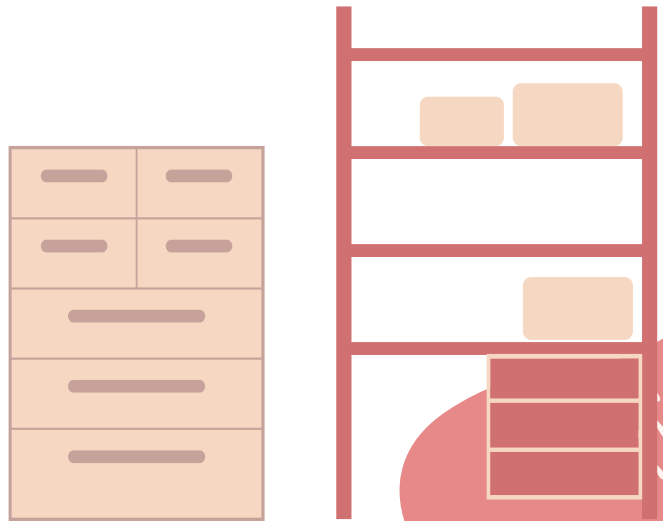
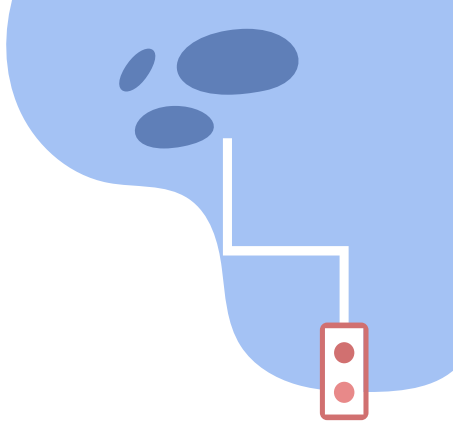
Time constrain on the project time line

Scope

11 Models in 2 weeks

Domain Knowledge

Domain expertise in the medical industry. This model can help clinicians deal with better Prognosis for patients



Questions?

THANK YOU

