

# IBM MACHINE LEARNING

## DEEP AND REINFORCEMENT LEARNING MODELS FOR DIABETES PREDICTION



**Fatima, Sayeda**

**1/13/2022**

## Table of Contents

1) Project Overview and Data Description .....	2
1a) Problem Overview.....	2
1b) About the Dataset.....	2
1b-i) Brief Description of Chosen Data Set .....	2
1b-ii) Summary of Data Attributes .....	2
1c) Data Exploration, Data Cleansing and Features Engineering .....	3
1c-i) Data Exploration .....	3
1c-ii) Data Cleansing Actions & Features Engineering .....	6
2) Main Objectives of the Analysis .....	11
2a) Primary Objective.....	11
2a) Secondary Objectives.....	11
3) Summary of Training Different Deep Learning Models .....	12
3a) Keras Hyperparameter Tuning using Grid Search .....	12
3b) Automated Model Building .....	13
3c) Summarizing Employed Models .....	14
1) Keras Algorithm with Optimizer RMSprop:.....	14
2) Keras Algorithm with Optimizer SGD:.....	16
3) Keras Algorithm with Optimizer Nadam:.....	18
4) Keras Algorithm with Optimizer Adamax: .....	20
5) Keras Algorithm with Optimizer Adadelta: .....	22
6) Keras Algorithm with Optimizer Adam: .....	24
7) Keras Algorithm with Optimizer Adagrad:.....	26
8) Base Model: Random Forest:.....	28
3d) Model Choice.....	29
3d-i) Result Summary .....	29
3d-ii) Model Ranking, Choice and Justification .....	29
4) Summary Key Findings and Insights .....	30
4a) Feature Weights within Models .....	30
4b Overall Feature Prominence .....	31
5) Model Shortcomings and Future Directions.....	32
5a) Model Shortcomings .....	32
5b) Plan of Action to Revisit Analysis .....	32
6 Useful Links .....	33
6a) Guide to Keras.....	33
6b) Keras Hyperparameter Tuning Approaches .....	33
6c) Model Checkpoints.....	33
6d) Get Feature Importance in Keras .....	33
6e) Shap Plots: Some Good Examples.....	33
6f) Matplot Lib .....	33
6f) Other Useful Models.....	33
7) Github Links .....	34
7a) Link to Main Folder.....	34
7b) Link to Assignment Notebook .....	34
8) References .....	35

## 1) Project Overview and Data Description

### 1a) Problem Overview

The drastic and often fatal impacts of diabetes not only are not only worrying but its steady incline expected to reach at 629 million by 2045 has also turned it into global threat (Naz and Ahuja, 2020). Despite its alarming increase, diabetes is principally a preventable disease which can be prevented by adopting healthier lifestyle changes which may also decrease probability of developing other diseases like cancer or heart problems. Hence, early detection of diabetes through a reliable prognosis tool is crucial to either prevent disease onset or stop its further progression.

### 1b) About the Dataset

#### 1b-i) Brief Description of Chosen Data Set

This project uses a hypothetical dataset ‘UCI Pima Diabetes Dataset’ which has been acquired for identifying risk of diabetes and was downloaded from the following link:

<https://www.kaggle.com/datasets/uciml/pima-indians-diabetes-database>

#### 1b-ii) Summary of Data Attributes

The PIMA dataset exhibits 768 data points (rows) and 9 features (columns) reflecting on patients’ characteristics where, based on various factors, each patient has been assigned a Diabetes Score.

Of these, the main features are:

1. times\_pregnant,
2. glucose\_tolerance\_test,
3. blood\_pressure,
4. skin\_thickness,
5. insulin,
6. bmi,
7. pedigree\_function
8. age

## 1c) Data Exploration, Data Cleansing and Features Engineering

Since the quality of any machine learning model highly depends on quality of data, hence, this stage is not only most important but is also time consuming. Hence, it was conducted in a step-by-step process.

### 1c-i) Data Exploration

- Data was first loaded into pandas dataframe

#### Read Dataset

```
1 ## Load in the data set
2 columns = ["times_pregnant", "glucose_tolerance_test", "blood_pressure", "skin_thickness", "insulin",
3           "bmi", "pedigree_function", "age", "has_diabetes"]
4 features = ["times_pregnant", "glucose_tolerance_test", "blood_pressure", "skin_thickness", "insulin",
5           "bmi", "pedigree_function", "age"]
6
7 diabetes_df = pd.read_csv(input_file_path, names=columns, header=0)
```

- Column types were then explored

Data Rows & Columns: (768, 9)

#### Data Types:

	Type
times_pregnant	int64
glucose_tolerance_test	int64
blood_pressure	int64
skin_thickness	int64
insulin	int64
bmi	float64
pedigree_function	float64
age	int64
has_diabetes	int64

#### Data Display:

	times_pregnant	glucose_tolerance_test	blood_pressure	skin_thickness	insulin	bmi	pedigree_function	age	has_diabetes
705	6	80	80	36	0	39.8	0.177	28	0
457	5	86	68	28	71	30.2	0.364	24	0
449	0	120	74	18	63	30.5	0.285	26	0

- Automated Exploratory Data Analysis was performed using Sweetviz

```
1 display_alert("Initial Exploratory Data Analysis", "success")
2 display_alert_color("Perform Quick EDA: To see dataset's distribution and its dispersion.", "teal", "warning")
3 df_eda = sv.analyze(customer_data) # Use Sweetviz for Automated EDA
4 #df_eda.show_html() # Uncomment to Generate Online Report
5 display_alert("Method to Generate Preliminary Exploratory Data Analysis ", "warning")
```

#### Initial Exploratory Data Analysis

Perform Quick EDA: To see dataset's distribution and its dispersion.

Done! Use 'show' commands to display/save. [100%] 00:00 -> (00:00 left)

Report SWEETVIZ\_REPORT.html was generated! NOTEBOOK/COLAB USERS: the web browser MAY not pop up, regardless, the report IS saved in your notebook/colab files.



- Additional descriptive statistics were computed to summarize shape of a dataset's distribution, its dispersion and central tendency

### Summary Statistics

	count	mean	std	min	25%	50%	75%	max
<b>times_pregnant</b>	768.0	3.845052	3.369578	0.000	1.00000	3.0000	6.00000	17.00
<b>glucose_tolerance_test</b>	768.0	120.894531	31.972618	0.000	99.00000	117.0000	140.25000	199.00
<b>blood_pressure</b>	768.0	69.105469	19.355807	0.000	62.00000	72.0000	80.00000	122.00
<b>skin_thickness</b>	768.0	20.536458	15.952218	0.000	0.00000	23.0000	32.00000	99.00
<b>insulin</b>	768.0	79.799479	115.244002	0.000	0.00000	30.5000	127.25000	846.00
<b>bmi</b>	768.0	31.992578	7.884160	0.000	27.30000	32.0000	36.60000	67.10
<b>pedigree_function</b>	768.0	0.471876	0.331329	0.078	0.24375	0.3725	0.62625	2.42
<b>age</b>	768.0	33.240885	11.760232	21.000	24.00000	29.0000	41.00000	81.00
<b>has_diabetes</b>	768.0	0.348958	0.476951	0.000	0.00000	0.0000	1.00000	1.00

### 1c-ii) Data Cleansing Actions & Features Engineering

In machine learning, feature selection is the method to reduce the number of input variables during developing predictive modelling. This reduction in input variables is necessary not only to minimize computational cost of modeling but also to achieve performance improvement of the model.

Among widely practices feature selection approaches include statistical-based feature selection methods which use statistical measures to evaluate relationship between each input variable and the target variable and then select those exhibiting strongest relationship with the latter. While these methods can be both speedy and effective, however, the ultimate choice of statistical measure is largely dependent on data types of both of these variables.

Irrespective of the statistical measure being employed, two dominant feature selection techniques, that is supervised and unsupervised, exist where the former can be further categorized into wrapper, filter and intrinsic techniques. Filter-based feature selection methods employs statistical measures to evaluate correlation between input and output variables so that those exhibiting highest correlations are selected. Statistical measures employed in filter-based feature selection are normally univariate in nature since they evaluate relationship of single input variables one by one with target variable, disregarding their interaction with each other.

Consequently, adopting filter-based feature selection methods, the project approached filter engineering in the following steps.

An automated data cleansing method was created to do the following:

- Drop Columns with Unique Values Less than threshold of 2
- Drop Highly Skewed & Low Correlation Columns with target
- Drop Columns with High Nan Values



## Method to Drop Columns

- 1) With Distinct < 2
- 2) With High Skewness and Low Correlation to Target
- 3) Drop Columns With High Nan Values

```

1 def drop_cols(*name):
2     n = name # Extract Dataframe by Name...this will create a 3d tuple
3     n = (n[0]) # Convert Tuple to DataFrame
4     df_name = [x for x in globals() if globals()[x] is n][0] # Extract Name of Imported Dataframe
5
6     n = n.fillna(0) #Fill Remaining Missing Values with Zero
7     # Find Mean of Null, Nan and Zero Values Before Any Drops
8     m0 = n.isin([' ', 'NULL', 'NaN', 0]).mean().sort_values(axis=0, ascending=False, inplace=False, kind='quicksort', na_posit
9
10
11     # 1) Drop Columns with Unique Values Less than threshold
12     unique_counts = pd.DataFrame.from_records([(col, n[col].nunique()) for col in n.columns], # get unique counts
13         columns=['Column_Name', 'Unique']).sort_values(by=['Unique'])
14     unique = unique_counts[(unique_counts['Unique'] < 2)] #If threshold is Less than 2 then
15     drop1 = (unique['Column_Name']).tolist() # First List of columns to drop
16
17     print(colored("\nDrop 1: ", 'blue', attrs=['bold']))
18     +colored(drop1, 'magenta', attrs=['bold']))
19
20     # 2) Drop Highly Skewed & Low Correlation Columns with target
21     drop2 = analyze[(analyze[column] != 1.000000) & ((analyze['skewness'] > 1) & (analyze[column] < 0))]
22     print(colored("\nDrop 2: \n ", 'blue', attrs=['bold']))
23     +colored(drop2, 'magenta', attrs=['bold']))
24
25     drop2 = drop2.sort_values(by='Columns', ascending=True)
26     drop2 = drop2['Columns'].tolist() # Second List of columns to drop
27
28     drop = drop1 + drop2 + delete_features # Final List of columns to drop
29     print(colored("\nFinal Column Drop List: ", 'blue', attrs=['bold']))
30     +colored(drop, 'magenta', attrs=['bold']))
31
32     if target in drop: # Remove Target from List
33         drop = drop.remove(target)
34     else:
35         print(target, " Not Found.")
36
37     n = n.drop(columns=[col for col in n if col in drop]) # Drop Dataframe Columns if in List
38
39     # Find Mean of Null, Nan and Zero Values Before Dropping
40     m1 = n.isin([' ', 'NULL', 'NaN', 0]).mean().sort_values(axis=0, ascending=False, inplace=False, kind='quicksort', na_posit
41
42     # 3) Drop Columns With High Nan Values
43     #drop_thresh = .90 # Identify Drop Threshold
44     #n = n.loc[:, df.isin([' ', 'NULL', 'NaN', 0]).mean() > drop_thresh] # drop columns if Mean is > 0.90
45
46     n = n.fillna(0) #Fill Remaining Missing Values with Zero
47     #n = n.replace(["NaN"], 0).sort_values(by=target, ascending=False) # Replace all Nan Values with Zero
48
49     # Find Mean of Null, Nan and Zero Values After Dropping
50     m2 = n.isin([' ', 'NULL', 'NaN']).mean().sort_values(axis=0, ascending=False, inplace=False, kind='quicksort', na_positio
51
52     #Print Results
53     print(colored("\nDataframe Average Null Values Before Any Drops\n ", 'blue', attrs=['bold']))
54     +colored(m0, 'magenta', attrs=['bold']))
55     +colored("\n\n Low Distict Columns to Drop: ", 'green', attrs=['bold']))
56     + colored(drop1, 'red', attrs=['bold']))
57     +colored("\n\nDataframe Average Null Values After Dropping Highly Skewed Columns\n ", 'green', attrs=['bold']))
58     +colored(m1, 'red', attrs=['bold']))
59     +colored("\n\n Drop Columns if Mean is > 0.90 \n", 'green', attrs=['bold']))
60     + colored("\n\nDataframe Average Null Values After Drop and 'Nan' Replacement\n", 'blue', attrs=['bold']))
61     +colored(m2, 'magenta', attrs=['bold']))
62     +colored(type(m2), 'magenta', attrs=['bold'))
63
64     return n
65
66 # Return Function
67 n = drop_cols(df)
68 df = n.copy()
69 df.info()

```



Null values were summed and were automatically managed by the above function.

#### Dataframe Average Null Values Before Any Drops

```
has_diabetes      0.651042
insulin           0.486979
skin_thickness    0.295573
times_pregnant    0.144531
blood_pressure    0.045573
bmi               0.014323
glucose_tolerance_test 0.006510
age               0.000000
pedigree_function 0.000000
dtype: float64
```

Low Distict Columns to Drop: []

#### Dataframe Average Null Values After Dropping Highly Skewed Columns

```
has_diabetes      0.651042
insulin           0.486979
skin_thickness    0.295573
times_pregnant    0.144531
blood_pressure    0.045573
bmi               0.014323
glucose_tolerance_test 0.006510
age               0.000000
pedigree_function 0.000000
dtype: float64
```

Drop Columns if Mean is > 0.90

#### Dataframe Average Null Values After Drop and 'Nan' Replacement

```
has_diabetes      0.0
age               0.0
pedigree_function 0.0
bmi               0.0
insulin           0.0
skin_thickness    0.0
blood_pressure    0.0
glucose_tolerance_test 0.0
times_pregnant    0.0
dtype: float64<class 'pandas.core.series.Series'>
```

**Outlier Treatment:** Like Supervised learning, deep learning models are also sensitive to outliers. Hence, an automated method was created to replace outliers with “Mode”, that is the most common value.

Method to Explore and Adjust Outliers:

Replace Outlier Values with Mode (Most Frequent Value)

```

1 def outliers(*name):
2     n = name # Extract Dataframe by Name...this will create a 3d tuple
3     n = (n[0]) # Convert Tuple to To Dataframe
4     df_name = [x for x in globals() if globals()[x] is n][0] # Extract Name of Imported Dataframe
5
6     cols = n.columns # All Columns
7
8     # Numeric Columns
9     numerics = ['int16', 'int32', 'int64', 'float16', 'float32', 'float64']
10    numeric_cols = n.select_dtypes(include=numerics)
11    numeric_cols = numeric_cols.columns.tolist()
12
13    # Object Columns
14    categorical_cols = list(set(cols) - set(numeric_cols))
15
16    # Skewed Columns
17    skewed_cols = analyze[(analyze['skewness'] > 0) | (analyze['skewness'] < 0)]
18    skewed_cols = skewed_cols['Columns'].tolist()
19    skewed_cols.remove(target) # Remove target column
20
21    # Replace Outliers
22    for col in n.columns:
23        if col in skewed_cols:
24            print(colored(col, 'magenta', attrs=['bold']))
25            + colored(" is Skewed...", 'blue', attrs=['bold'])
26            )
27
28            mode = n[col].mode()
29            mode = mode[0]
30
31            if col in numeric_cols:
32                print(colored(col, 'magenta', attrs=['bold']))
33                + colored(" Column Type is: ", 'blue', attrs=['bold'])
34                + colored(n[col].dtypes, 'red', attrs=['bold'])
35                )
36
37                #Calculate quantiles and IQR
38                Q1 = n[col].quantile(0.25) # Same as np.percentile but maps (0,1) and not (0,100)
39                Q3 = n[col].quantile(0.75)
40                IQR = Q3 - Q1
41                # Replace with Mode
42                n[col] = np.where((n[col] < (Q1 - 1.5 * IQR)) | (n[col] > (Q3 + 1.5 * IQR)), mode, n[col])
43
44                print(colored("\nReplaced ", 'blue', attrs=['bold']))
45                + colored(col, 'magenta', attrs=['bold'])
46                + colored(" Skewed Values by Mode: ", 'blue', attrs=['bold'])
47                + colored(mode, 'red', attrs=['bold'])
48                + colored("\n", 'magenta', attrs=['bold'])
49                + colored((n[col]), 'green', attrs=['bold'])
50                )
51
52            else:
53                print("")
54    df_transformed = n.copy()
55    return df_transformed

```

Apply 'outliers' Method to New Dataframe

times\_pregnant is Skewed...

times\_pregnant Column Type is: int64

Replaced times\_pregnant Skewed Values by Mode: 1

```

0      6
1      1
2      8
3      1
4      0
..
763    10
764     2
765     5
766     1
767     1

```

**Data Splitting & Normalization:** To speed up algorithms' learning, data normalization was carried out.

### Provide Method to Split Training and Testing Dataset

```

1 def split_data(*name):
2     n = name # Extract Dataframe by Name...this will create a 3d tuple
3     n = (n[0]) # Convert Tuple to To Dataframe
4     df_name = [x for x in globals() if globals()[x] is n][0] # Extract Name of Imported Dataframe to print Later
5     X = n.loc[ : , n.columns != target] # Remove Target Column
6     y = n[target].astype('int')
7     return train_test_split(X, y, test_size=0.25, stratify=y, random_state = rs)
8
9 display_alert_color_4("Split Data into Test & Train Set: ", "darkgreen", "warning")

```

### Split Data into Test & Train Set:

```

1 X_train, X_test, y_train, y_test = split_data(df_transformed) # Call the Method to Split Training and Testing Dataset
2
3 # Normalize Data
4 normalizer = StandardScaler()
5 X_train_norm = normalizer.fit_transform(X_train)
6 X_test_norm = normalizer.transform(X_test)
7
8 display_alert_color_4("Calculate Mean to Check Proportion of Positive Values: ", "darkbrown", "warning")

```

## 2) Main Objectives of the Analysis

Even though healthcare organizations normally collect big data including electronic health records and images, nevertheless, its robust analysis to acquire meaningful and reliable insights remains a key challenge. Given promising results evidenced by deep and other machine learning methodologies, such medical data can be automatically analysed to discover hidden patterns and factors which may aid in diabetes diagnosis at an early stage.

### 2a) Primary Objective

Hence, the main objective of this project is to present an automated methodology for employing different variations of Keras deep learning model for diabetes prediction using the PIMA dataset.

### 2a) Secondary Objectives

Secondary objective are as follows:

1. Build a baseline performance using Random Forest Model (RFM) for comparison
2. Develop automated grid search method for selecting best hyperparameters for Keras model development
3. Evaluate Keras results and compare with those of RFM to validate model robustness

### 3) Summary of Training Different Deep Learning Models

#### 3a) Keras Hyperparameter Tuning using Grid Search

An automated method was created to find optimal parameters for various Keras Optimizers:

#### GRID Search for Hyperparameter Tuning

##### Get Script Start Time

Script Start Time: 2023-01-09 18:47:23.204268

```

1  # Function to create model, required for KerasClassifier
2  def create_model():
3      # create model
4      model = Sequential()
5      model.add(Dense(6, input_shape=(int(len(X.columns)),), activation='relu'))
6      model.add(Dense(6, input_shape=(int(len(X.columns)),), activation='relu'))
7      model.add(Dense(1, activation='sigmoid'))
8      # Compile model
9      model.compile(loss='binary_crossentropy', optimizer='adam', metrics=['accuracy'])
10     return model
11 # fix random seed for reproducibility
12 #seed = 7
13 #tf.random.set_seed(seed)
14
15 # create model
16 model = KerasClassifier(model=create_model, verbose=0)
17
18 # define the grid search parameters
19 batch_size = [10, 20, 40]
20 epochs = [500, 800, 900]
21 optimizer = ['SGD', 'RMSprop', 'Adagrad', 'Adadelta', 'Adam', 'Adamax', 'Nadam']
22
23 # Search Best Parameters
24 param_grid = dict(batch_size=batch_size, epochs=epochs, optimizer=optimizer)
25 grid = GridSearchCV(estimator=model, param_grid=param_grid, n_jobs=-1, cv=3)
26
27 # Create Model
28 model = KerasClassifier(model=create_model, verbose=0)
29
30 # Fit Model
31 grid_result = grid.fit(X, y)

```

## 3b) Automated Model Building

```

1 optimizer_list = result_df['optimizer'].unique()
2 #optimizer_list = optimizer_list[:1] # Get the number of models from list...use 1 if only the top performing model is required, however, it may not always give best results
3 print(colored("\n Building Keras for Optimizers: " + optimizer_list, "blue", attr=['bold']))
4
5 for i in optimizer_list:
6     opt = result_df.loc[result_df['optimizer'] == i]
7     b_opt = int(opt['batch_size'].iloc[0])
8     e_opt = int(opt['epochs'].iloc[0])
9
10    # Create model
11    print(colored("\n" + str(i) + " Model with Batch Size: " + str(b_opt) + " and Epochs: " + str(e_opt)), 'magenta', attr=['bold'])
12    +colored("\nCreating Model", 'green', attr=['bold'])
13    model = Sequential()
14    model.add(Dense(6, input_shape=(c,), activation='relu'))
15    model.add(Dense(6, activation='relu'))
16    model.add(Dense(1, activation='sigmoid'))
17
18    # Compile model
19    model.compile(optimizer=i, loss='binary_crossentropy', metrics=['accuracy'])
20
21    # Set Callbacks
22    # Simple Early Stopping to stop training As soon as the loss of the model begins to increase on the test dataset
23    es = EarlyStopping(monitor='val_loss', mode='min', patience=(e_opt/4), verbose=0)# Add patience=(e_opt/4) to not stop
24    # ModelCheckpoint: The mode parameter controls whether the ModelCheckpoint should be looking for values that minimize our metric or maximize it.
25    # Since we are working with loss, lower is better, so we set mode="min". If we were instead working with val_acc, we would set mode="max" (since higher accuracy is better).
26    mc = ModelCheckpoint("best_model.h5", monitor='val_loss', mode='min', verbose=0, save_best_only=True)
27
28    print(colored('Fitting Model', 'blue', attr=['bold']))
29    # the fit function returns the run history and is very convenient, as it contains information about the model fit, iterations etc.
30    #By setting verbose 0, 1 or 2 you just say how do you want to "see" the training progress for each epoch.
31    #verbose=0 will show you nothing (silent)
32    #verbose=1 will show you an animated progress bar #####
33    #verbose=2 will just mention the number of epoch [Epoch 1/10]
34
35    run_hist = model.fit(X_train_norm, y_train, validation_data=(X_test_norm, y_test), epochs=e_opt, batch_size=b_opt, verbose=0, callbacks=[es, mc])
36
37    # Let's look at the run_hist object that was created, specifically its history attribute.
38    #run_hist.history.keys()
39
40    print(colored('Making Model Predictions', 'magenta', attr=['bold']))
41    y_pred_prob_nn = model.predict(X_test_norm)#.ravel()
42    y_pred_class_nn = np.argmax(y_pred_prob_nn,axis=0)
43
44
45    print(colored('Calculating AUC Score', 'magenta', attr=['bold']))
46    fpr, tpr, thresholds = roc_curve(y_test, y_pred_prob_nn)
47    auc_score = auc(fpr, tpr)
48    print(colored(str(i)+ " Model auc_score :"+ str(auc_score)), 'magenta', attr=['bold']))
49
50    # Print model validation loss performance
51    print(colored('14: Computing Model validation loss performance', 'magenta', attr=['bold']))
52    n = ((run_hist.history))
53    n = pd.DataFrame.from_dict(n, orient='columns')
54    print(colored(n.tail(1), 'magenta', attr=['bold']))
55    # Get val_loss value for head and tail
56    vlh = ((n['val_loss']).head(1).iloc[0])
57    vlt = ((n['val_loss']).tail(1).iloc[0])
58    vah = ((n['val_accuracy']).head(1).iloc[0])
59    vat = ((n['val_accuracy']).tail(1).iloc[0])
60    e_act = (len(run_hist.history['loss']))
61
62    # Initialize data of lists.
63    data = {'Optimizer': [i],
64            'AUC Score': [auc_score],
65            'Input Val Loss': [vlh],
66            'Output Val Loss': [vlt],
67            'Input Val Accuracy': [vah],
68            'Output Val Accuracy': [vat],
69            'Epochs': [e_opt],
70            'Actual Epochs': [e_act]}
71
72    # Create DataFrame
73    n = pd.DataFrame(data)
74    res = res.append(n, ignore_index=True)
75    print(colored(res, 'blue', attr=['bold']))
76
77    # Get Feature Importance
78    e = shap.KernelExplainer(model, X_train_norm)
79    shap_values = e.shap_values(X_test_norm)
80
81    # Get Feature Importance in DataFrame
82    shap_mean = np.abs(shap_values).mean(axis=0)
83    shap_mean = pd.DataFrame(shap_mean, columns = features)
84    shap_mean = shap_mean.mean()
85    shap_df = shap_mean.to_frame().reset_index()
86    shap_df = shap_df.rename({'index': 'Features', 0: 'Importance'}, axis=1)#.sort_values(by='Importance', ascending=True)
87    shap_df['Importance'] = abs(shap_df['Importance'])
88    shap_df['Model'] = [i]
89    shap_df = shap_df.sort_values(by='Importance', ascending=False)
90    model_features = model.features.append(shap_df, ignore_index=True)
91
92    print(colored('14: List Feature Importance\n', 'magenta', attr=['bold']))
93    +colored(shap_df, 'blue', attr=['bold']))
94
95    # Plot the roc curve
96    print(colored("\n14: ROC Curve:", 'magenta', attr=['bold']))
97    plot_roc(y_test, y_pred_prob_nn, 'NW')
98
99    # Let's plot the training loss and the validation loss over the different epochs and see how it looks.
100    fig, (ax1, ax2) = plt.subplots(2, 1)
101    ax1.plot(run_hist.history['loss'], "r", markers='.', label="Train Loss")
102    ax1.plot(run_hist.history['val_loss'], "b", markers='.', label="Validation Loss")
103    ax1.set(title= (i + ": Displaying Model Validation Loss").format(i))
104    ax1.legend()
105    ax2.set(title= (i + ": Displaying SHAP Feature Importance").format(i))
106
107    #We can use the shapely values to interpret our model. 'force_plot' showing how each feature influences the output.
108    x1 = shap.summary_plot(shap_values[0], X_test, feature_names=features)
109    x1 = shap.summary_plot(shap_values, X_test, feature_names=features, plot_type = 'bar', title=(i + ": Displaying SHAP Feature Importance with Bars"))
110
111    # We can use the shapely values to interpret our model.
112    shap.initjs()
113    # visualize the first prediction's explanation with a force plot
114    shap.force_plot(e.expected_value[0], shap_values[0][0], features = features)
115    shap.plots.force(e.expected_value[0], shap_values[0][0], features = features, matplotlib=True, show=False)
116    plt.title(i + " Model Feature Importance", y=1.75)
117    plt.show()

```

### 3c) Summarizing Employed Models

Following is a summary of all Keras Model Variations which have been used to predict onset of diabetes.

#### 1) Keras Algorithm with Optimizer RMSprop:

RMSprop Model with Batch Size: 10 and Epochs: 900

Creating Model

Fitting Model

Using 576 background data samples could cause slower run times. Consider

Making Model Predictions

Calculating AUC Score

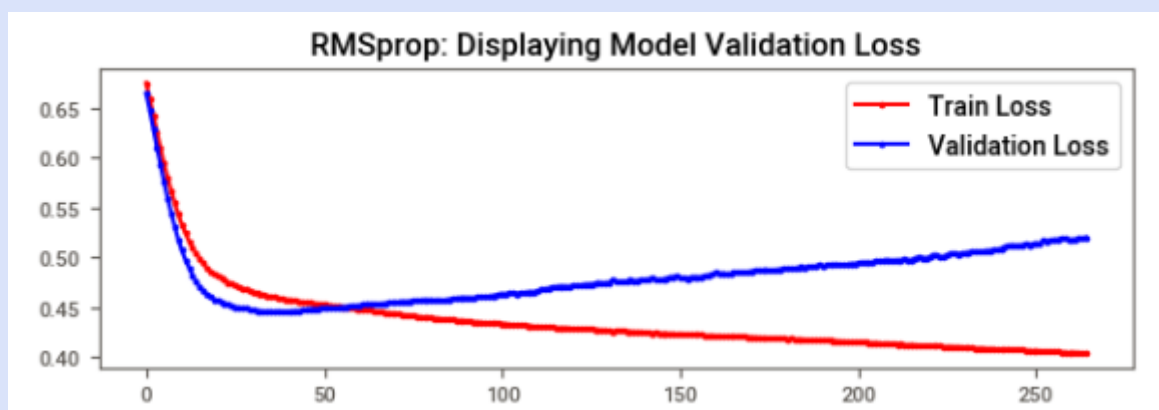
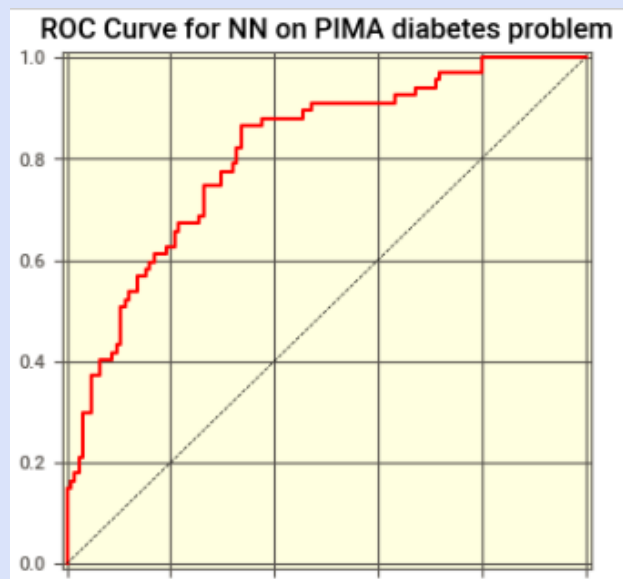
RMSprop Model auc\_score :0.8108656716417909

RMSprop: Computing Model validation loss performance

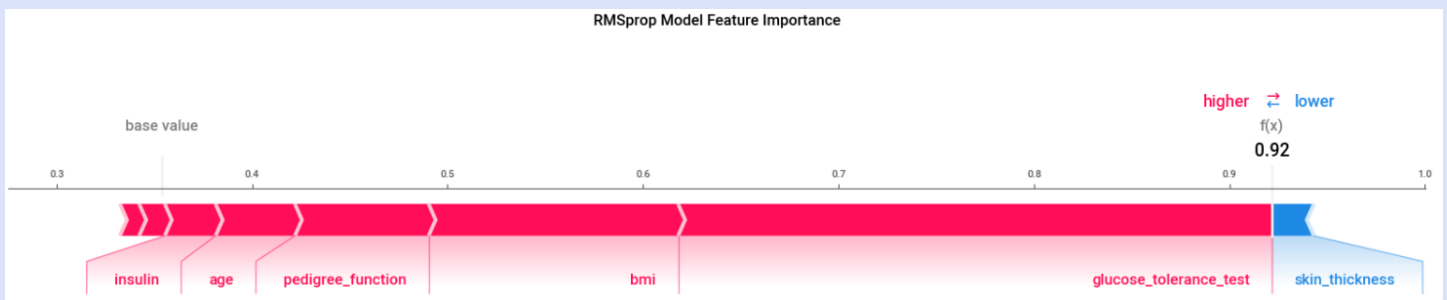
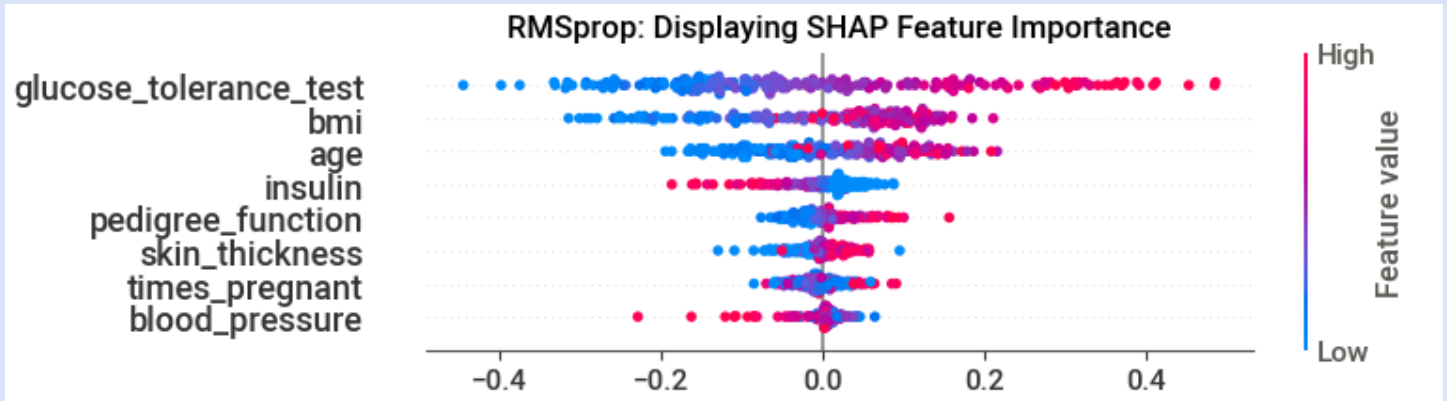
	loss	accuracy	val_loss	val_accuracy
264	0.403722	0.814236	0.519842	0.744792

	Optimizer	Input	Val Loss	Output	Val Loss	Input	Val Accuracy
0		0	0.000000		0.000000		0.000000
1	RMSprop		0.664082		0.519842		0.697917

	Output	Val Accuracy	AUC Score	Epochs	Actual Epochs
0		0.000000	NaN	NaN	NaN
1		0.744792	0.810866	900.0	265.0







### RMSprop: List Feature Importance

	Features	Importance	Model
1	glucose_tolerance_test	0.168594	RMSprop
5	bmi	0.098901	RMSprop
7	age	0.075343	RMSprop
4	insulin	0.035594	RMSprop
6	pedigree_function	0.024697	RMSprop
3	skin_thickness	0.022470	RMSprop
0	times_pregnant	0.021662	RMSprop
2	blood_pressure	0.017175	RMSprop

## 2) Keras Algorithm with Optimizer SGD:

SGD Model with Batch Size: 20 and Epochs: 800

Creating Model

Fitting Model

Using 576 background data samples could cause slower run times. Consider

Making Model Predictions

Calculating AUC Score

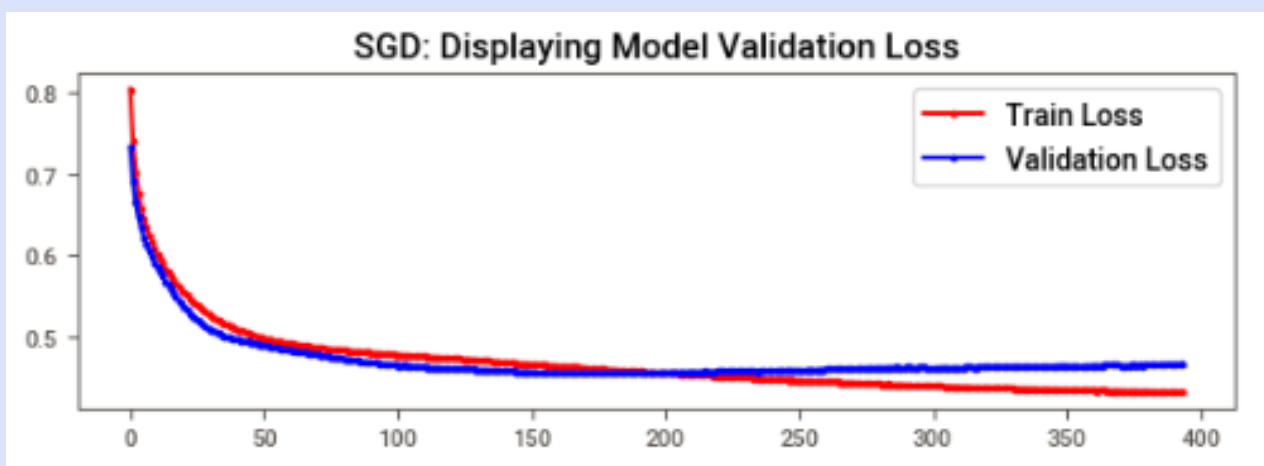
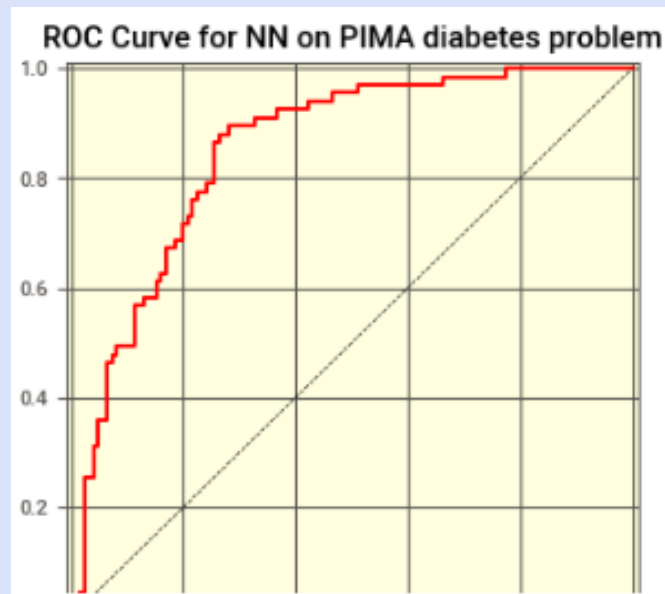
SGD Model auc\_score :0.8511044776119403

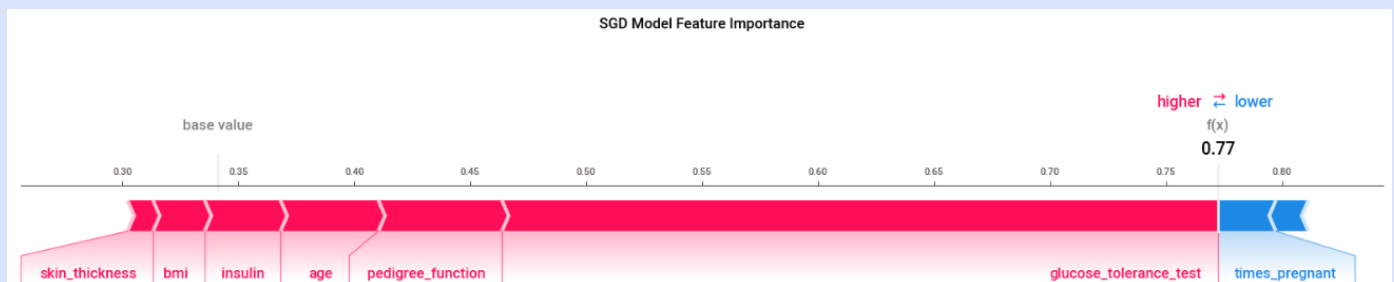
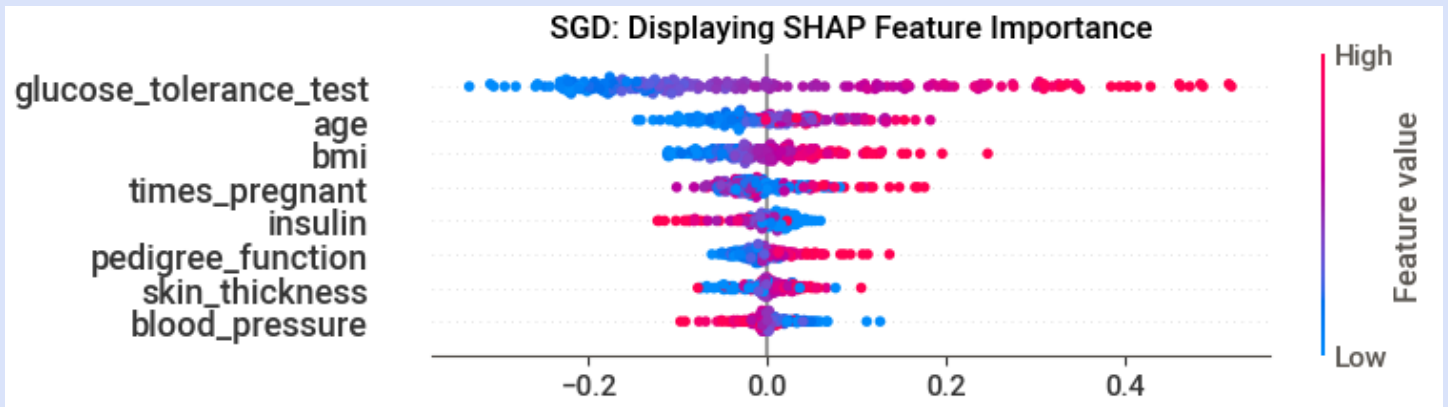
SGD: Computing Model validation loss performance

	loss	accuracy	val_loss	val_accuracy
393	0.432101	0.814236	0.46756	0.765625

	Optimizer	Input Val Loss	Output Val Loss	Input Val Accuracy \
0	0	0.000000	0.000000	0.000000
1	RMSprop	0.664082	0.519842	0.697917
2	SGD	0.732052	0.467560	0.572917

	Output Val Accuracy	AUC Score	Epochs	Actual Epochs
0	0.000000	NaN	NaN	NaN
1	0.744792	0.810866	900.0	265.0





## SGD: List Feature Importance

	Features	Importance	Model
1	glucose_tolerance_test	0.175242	SGD
7	age	0.050878	SGD
5	bmi	0.045077	SGD
0	times_pregnant	0.034531	SGD
4	insulin	0.023013	SGD
6	pedigree_function	0.021368	SGD
3	skin_thickness	0.020594	SGD
2	blood_pressure	0.016589	SGD

### 3) Keras Algorithm with Optimizer Nadam:

Nadam Model with Batch Size: 10 and Epochs: 800

Creating Model

Fitting Model

Using 576 background data samples could cause slower run times. Co

Making Model Predictions

Calculating AUC Score

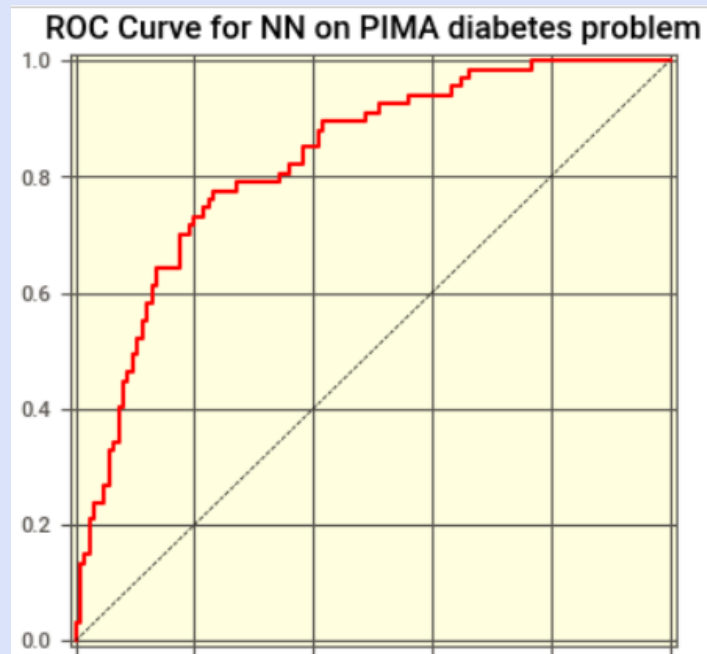
Nadam Model auc\_score :0.8262686567164179

Nadam: Computing Model validation loss performance

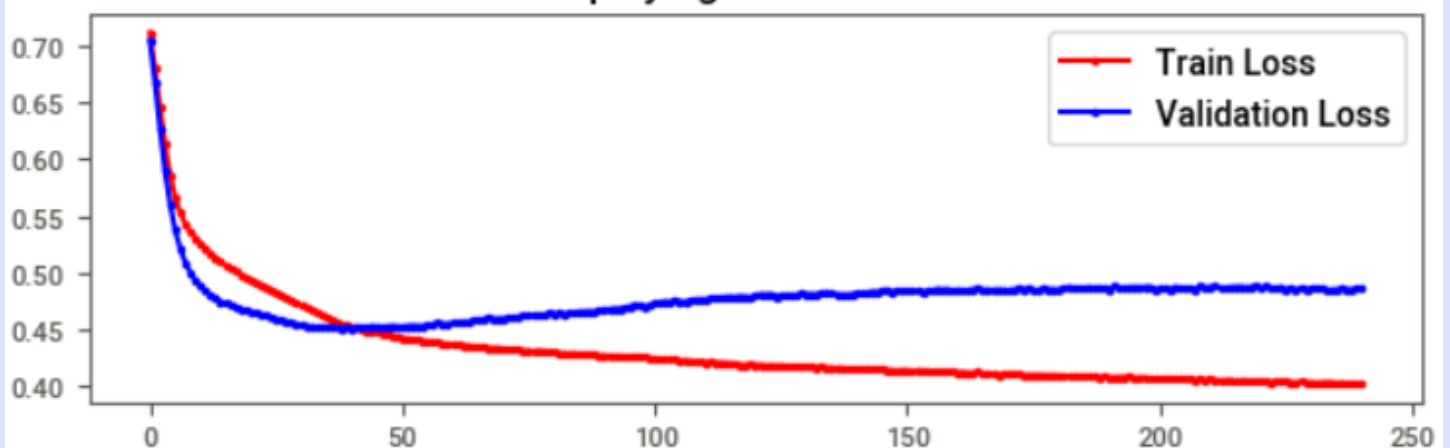
	loss	accuracy	val_loss	val_accuracy
240	0.402151	0.810764	0.486368	0.776042

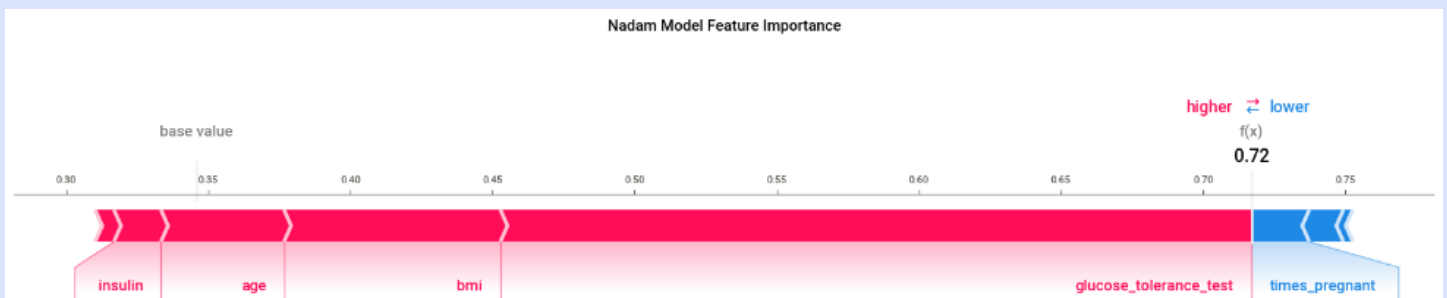
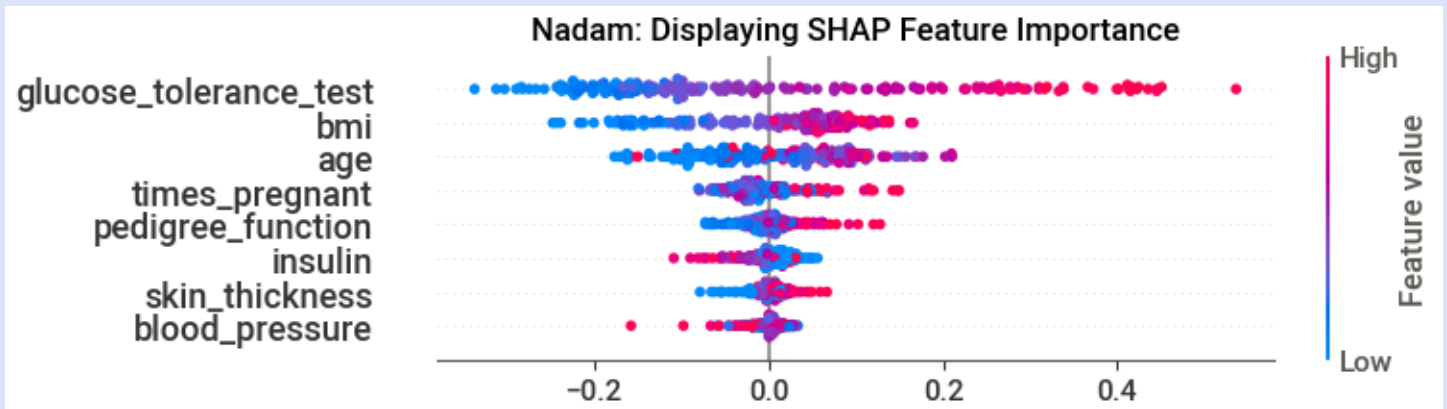
  

Optimizer	Input	Val Loss	Output	Val Loss	Input	Val Accuracy
0	0	0.000000		0.000000		0.000000
1	RMSprop	0.664082		0.519842		0.697917
2	SGD	0.732052		0.467560		0.572917
3	Nadam	0.704898		0.486368		0.651042



Nadam: Displaying Model Validation Loss





## Nadam: List Feature Importance

	Features	Importance	Model
1	glucose_tolerance_test	0.179169	Nadam
5	bmi	0.079997	Nadam
7	age	0.071661	Nadam
0	times_pregnant	0.030211	Nadam
6	pedigree_function	0.022472	Nadam
4	insulin	0.017829	Nadam
3	skin_thickness	0.013714	Nadam
2	blood_pressure	0.011544	Nadam

#### 4) Keras Algorithm with Optimizer Adamax:

Adamax Model with Batch Size: 10 and Epochs: 900

Creating Model

Fitting Model

Using 576 background data samples could cause slower run times. Consider

Making Model Predictions

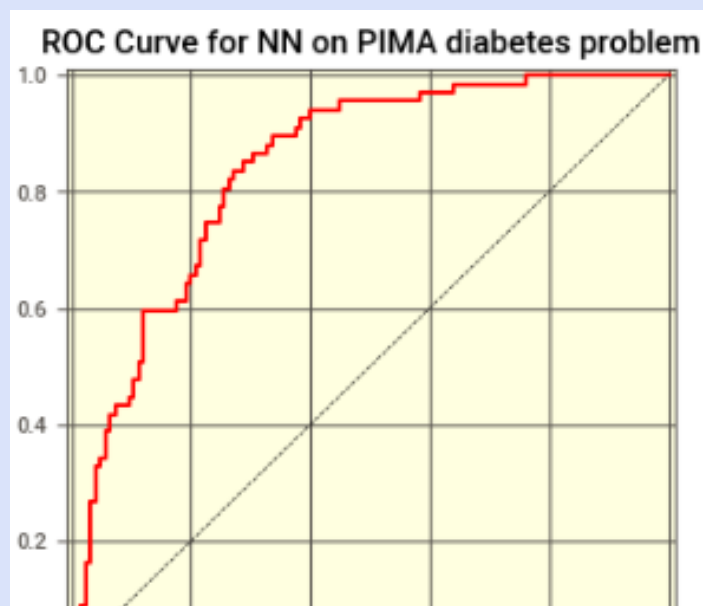
Calculating AUC Score

Adamax Model auc\_score :0.8411940298507463

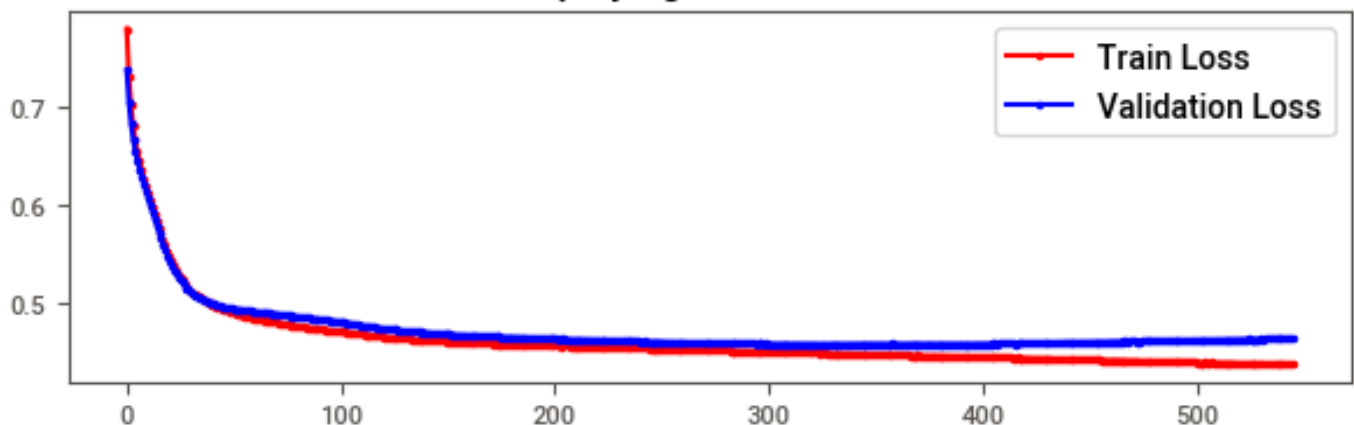
Adamax: Computing Model validation loss performance

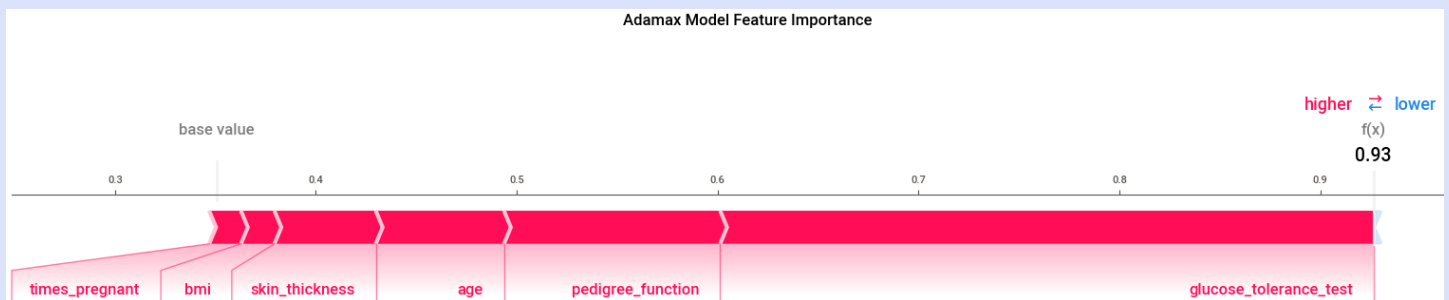
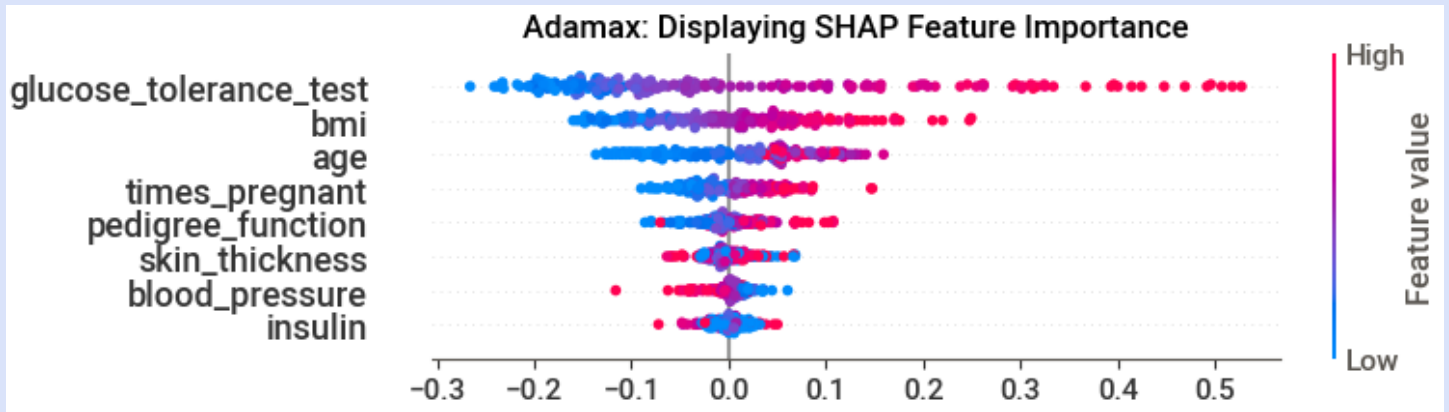
	loss	accuracy	val_loss	val_accuracy
545	0.438843	0.788194	0.46459	0.776042

Optimizer	Input	Val Loss	Output	Val Loss	Input	Val Accuracy	\
0	0	0.000000		0.000000		0.000000	
1	RMSprop	0.664082		0.519842		0.697917	
2	SGD	0.732052		0.467560		0.572917	
3	Nadam	0.704898		0.486368		0.651042	
4	Adamax	0.736525		0.464590		0.432292	



Adamax: Displaying Model Validation Loss





## Adamax: List Feature Importance

	Features	Importance	Model
1	glucose_tolerance_test	0.158482	Adamax
5	bmi	0.069657	Adamax
7	age	0.061823	Adamax
0	times_pregnant	0.032609	Adamax
6	pedigree_function	0.023062	Adamax
3	skin_thickness	0.015775	Adamax
2	blood_pressure	0.013057	Adamax
4	insulin	0.012496	Adamax



## 5) Keras Algorithm with Optimizer Adadelata:

Adadelata Model with Batch Size: 20 and Epochs: 800

Creating Model

Fitting Model

Using 576 background data samples could cause slower run times. Consider using a smaller batch size.

Making Model Predictions

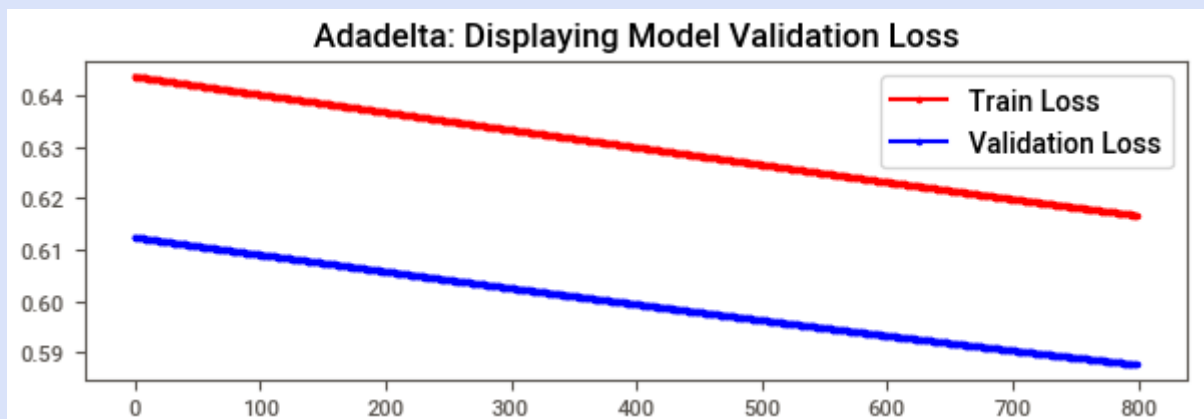
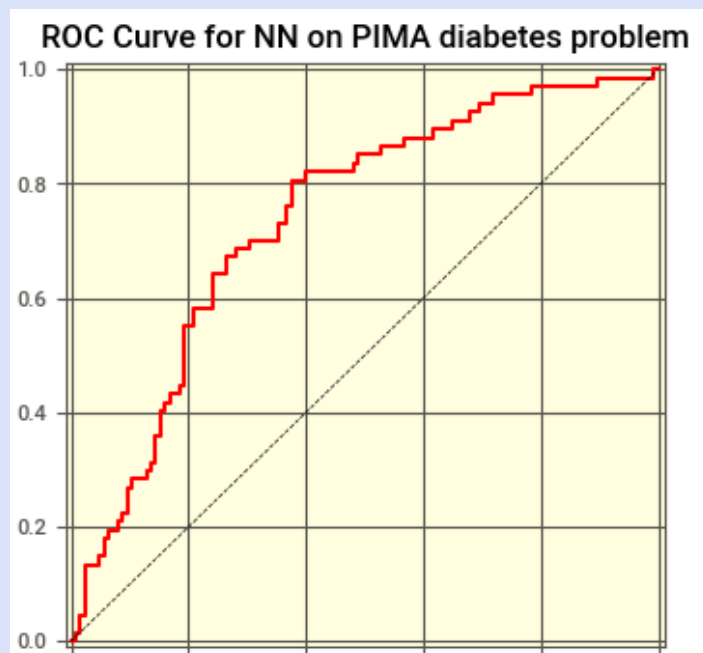
Calculating AUC Score

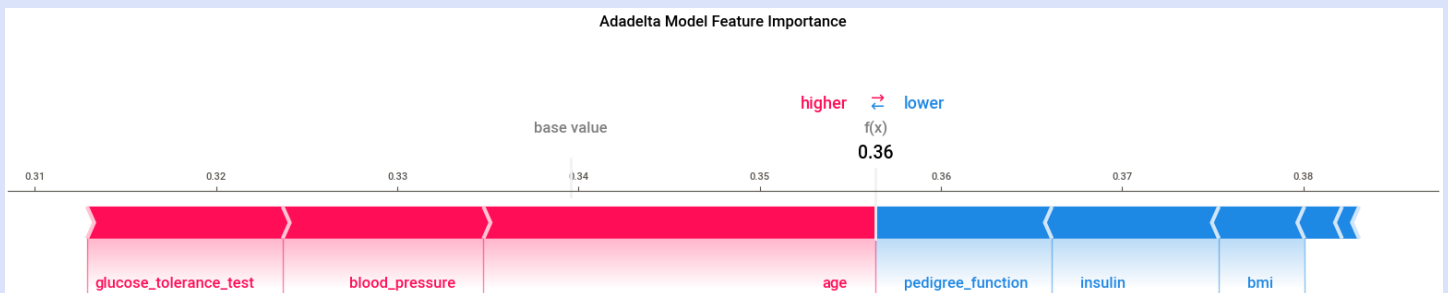
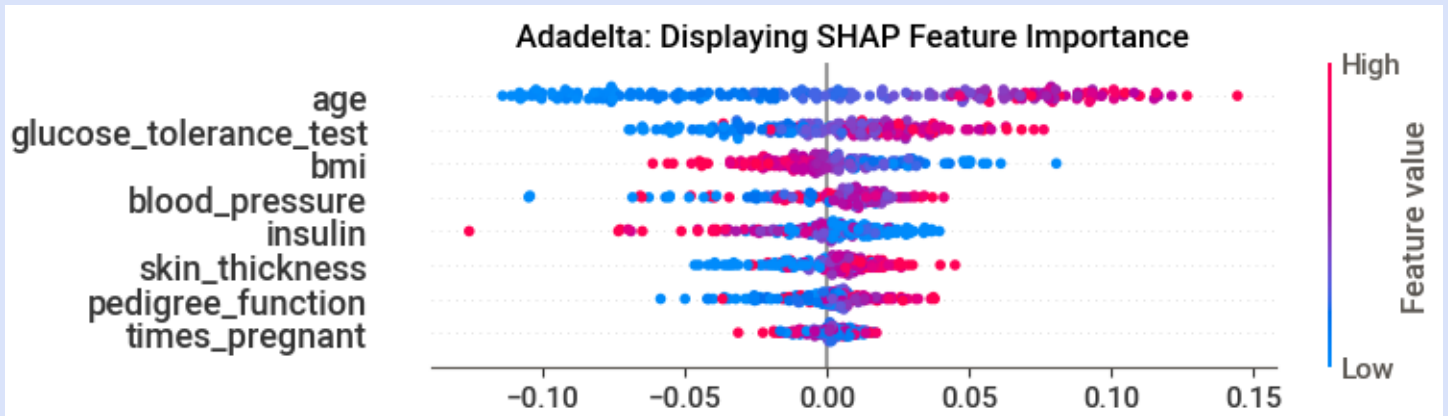
Adadelata Model auc\_score :0.7383880597014926

Adadelata: Computing Model validation loss performance

	loss	accuracy	val_loss	val_accuracy
799	0.616719	0.651042	0.587543	0.651042

	Optimizer	Input	Val Loss	Output	Val Loss	Input	Val Accuracy
0	0		0.000000		0.000000		0.000000
1	RMSprop		0.664082		0.519842		0.697917
2	SGD		0.732052		0.467560		0.572917
3	Nadam		0.704898		0.486368		0.651042
4	Adamax		0.736525		0.464590		0.432292
5	Adadelata		0.612426		0.587543		0.651042





## Adadelata: List Feature Importance

	Features	Importance	Model
7	age	0.060745	Adadelata
1	glucose_tolerance_test	0.025040	Adadelata
5	bmi	0.018069	Adadelata
2	blood_pressure	0.016655	Adadelata
4	insulin	0.016061	Adadelata
3	skin_thickness	0.012444	Adadelata
6	pedigree_function	0.011359	Adadelata
0	times_pregnant	0.006161	Adadelata

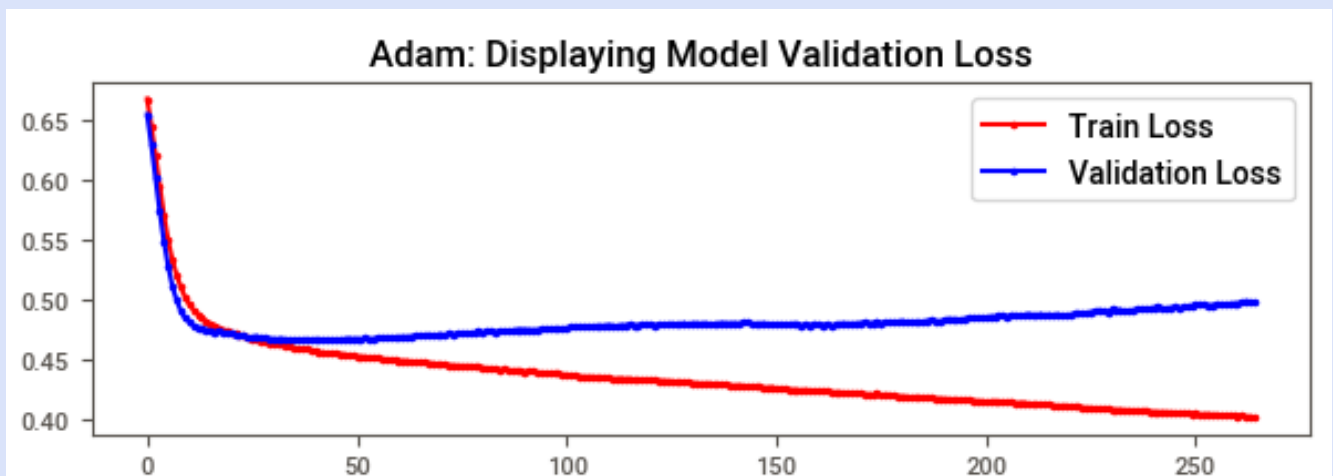
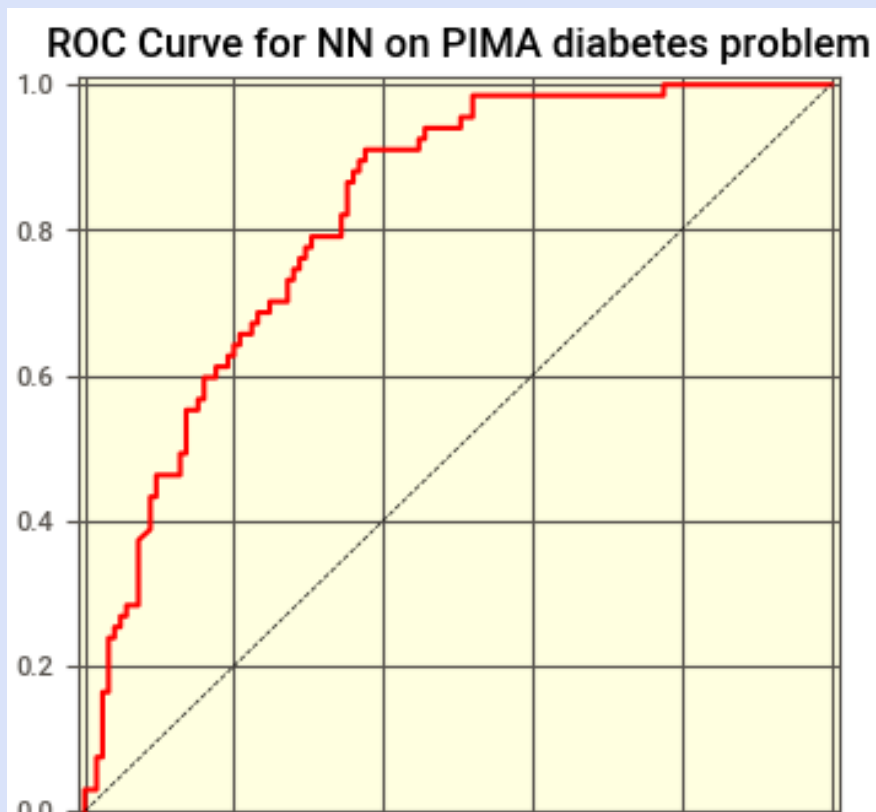
## 6) Keras Algorithm with Optimizer Adam:

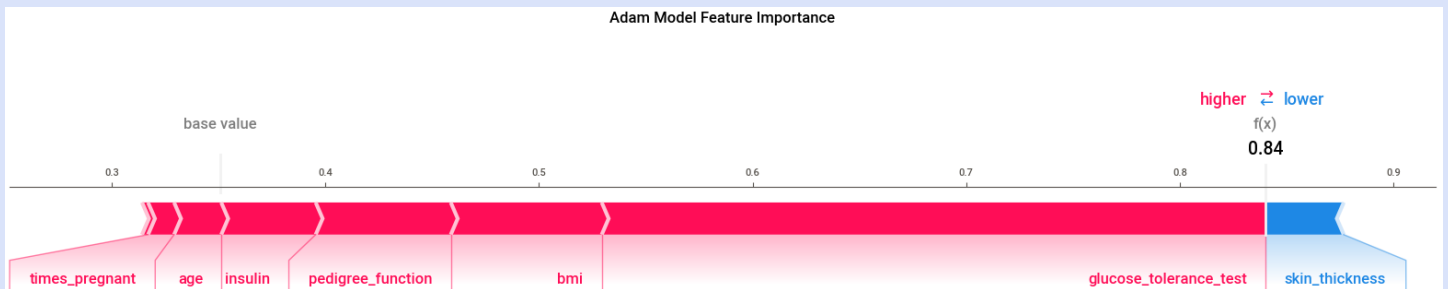
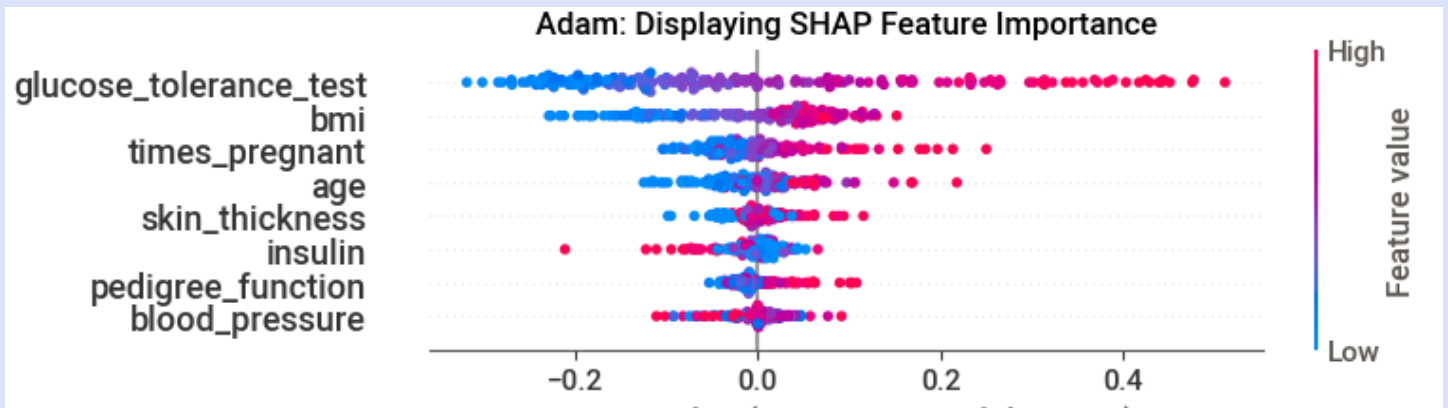
Adam Model auc\_score :0.8208955223880596

Adam: Computing Model validation loss performance

	loss	accuracy	val_loss	val_accuracy
264	0.401989	0.815972	0.498047	0.744792

	Optimizer	Input Val Loss	Output Val Loss	Input Val Accuracy \
0	0	0.000000	0.000000	0.000000
1	RMSprop	0.664082	0.519842	0.697917
2	SGD	0.732052	0.467560	0.572917
3	Nadam	0.704898	0.486368	0.651042
4	Adamax	0.736525	0.464590	0.432292
5	Adadelta	0.612426	0.587543	0.651042
6	Adam	0.655232	0.498047	0.671875





## Adam: List Feature Importance

	Features	Importance	Model
1	glucose_tolerance_test	0.170184	Adam
5	bmi	0.072324	Adam
0	times_pregnant	0.040196	Adam
7	age	0.035781	Adam
3	skin_thickness	0.020557	Adam
4	insulin	0.019800	Adam
6	pedigree_function	0.019177	Adam
2	blood_pressure	0.018084	Adam

## 7) Keras Algorithm with Optimizer Adagrad:

Making Model Predictions

Calculating AUC Score

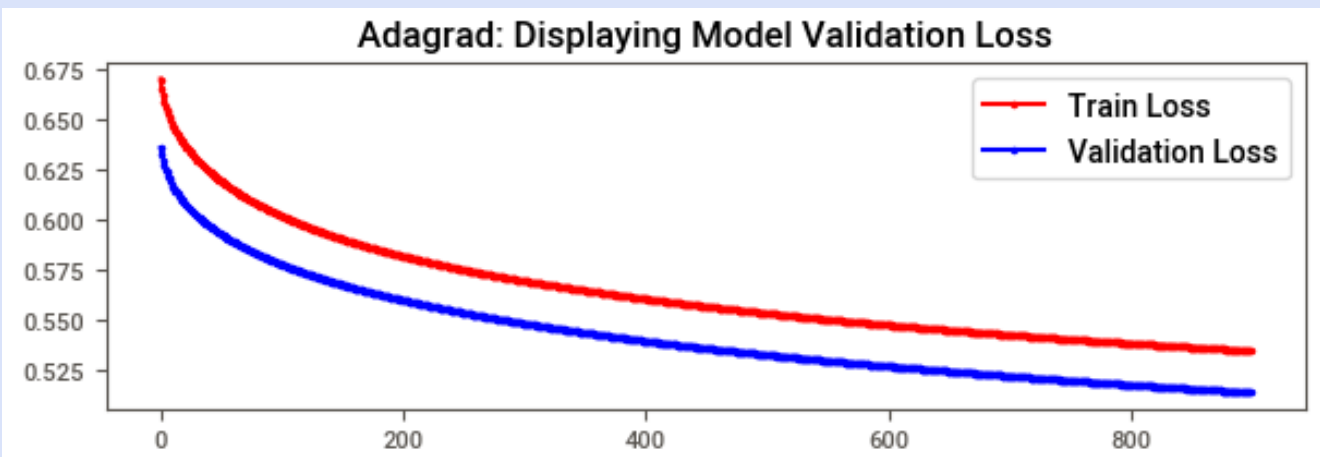
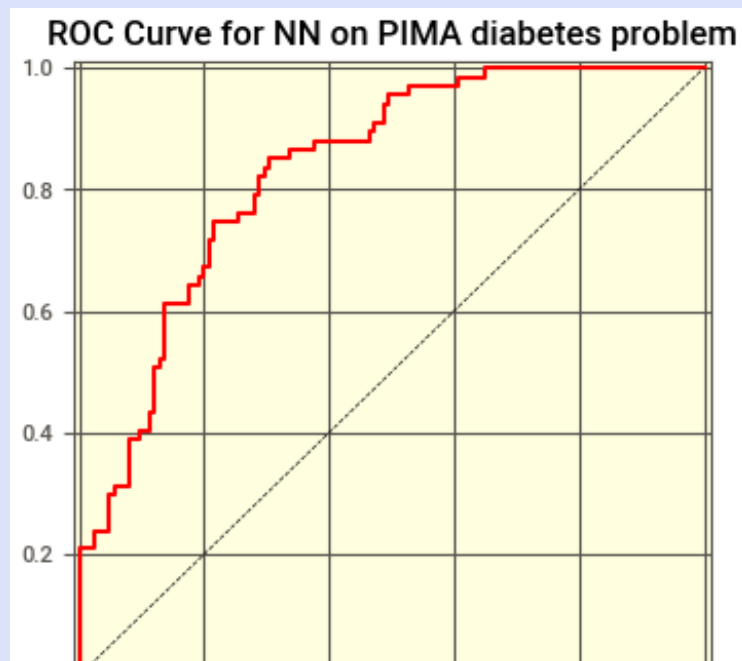
Adagrad Model auc\_score :0.8320000000000001

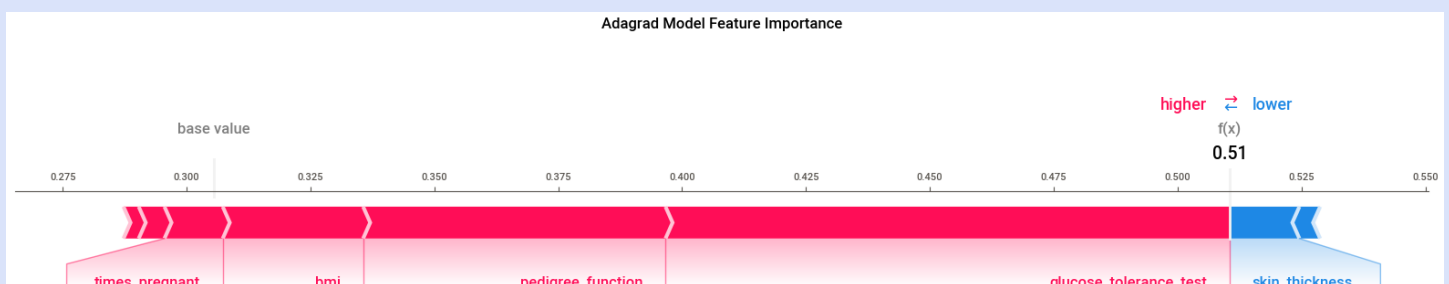
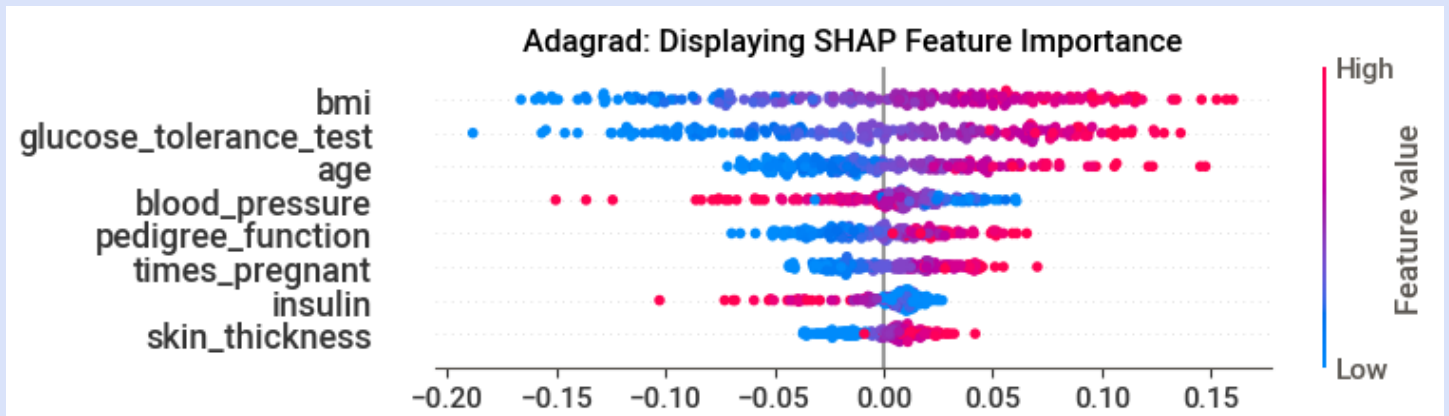
Adagrad: Computing Model validation loss performance

loss accuracy val\_loss val\_accuracy

899 0.534858 0.682292 0.514088 0.71875

	Optimizer	Input	Val Loss	Output	Val Loss	Input	Val Accuracy	\
0	0		0.000000		0.000000		0.000000	
1	RMSprop		0.664082		0.519842		0.697917	
2	SGD		0.732052		0.467560		0.572917	
3	Nadam		0.704898		0.486368		0.651042	
4	Adamax		0.736525		0.464590		0.432292	
5	Adadelat		0.612426		0.587543		0.651042	
6	Adam		0.655232		0.498047		0.671875	
7	Adagrad		0.635958		0.514088		0.645833	





## Adagrad: List Feature Importance

	Features	Importance	Model
5	bmi	0.064633	Adagrad
1	glucose_tolerance_test	0.059301	Adagrad
7	age	0.036672	Adagrad
2	blood_pressure	0.023516	Adagrad
6	pedigree_function	0.022678	Adagrad
0	times_pregnant	0.021181	Adagrad
4	insulin	0.013991	Adagrad
3	skin_thickness	0.013115	Adagrad

## 8) Base Model: Random Forest:

### RF Model WITH Auto-Hyptuned Parameters

```

1 #Running the Model...Call Method to 'Build Simple Random Forest With Class Weights'
2 op_preds, op_model = build_op_rf(X_train, y_train, X_test, best_params=optimal_rf_params) # Call Method to 'Build Simple Random Forest Without Class Weights'
3
4 #Collect Results
5 original_results = evaluate(y_test, op_preds, eval_type="RForest")
6 original_results["Factors"] = 'RF1_features_df'
7 rf_acc = original_results['auc'] # Extract Roc_acc for comparison with deep Learning models
8 print(original_results, "\nRF_acc: ", rf_acc)
9
10 # Get Factors Contributing to Diabetes
11 RF1_features_df = feature_importance(X, op_model)[:40]
12 RF1_features_df = RF1_features_df[RF1_features_df['Importance'] > 0] # Filter Features with Negative Correlation
13 RF1_features_df = RF1_features_df.sort_values(by=['Importance'], ascending=True) # Attach Importance to data before split
14
15 # Capture Model Drivers in Dataframe
16 RF1_drivers = [['Uses train test set to achieve best hyptuning to resample data from train set to achieve balanced class.'],
17               ['No optimal parameters were used.'],
18               ['Runs random forest model with above drivers to get top contributory features.']] # initialize List of Lists
19 RF1_drivers = pd.DataFrame(RF1_drivers, columns=['Main Drivers behind RF1 model are as follows:']) # Create the pandas DataFrame
20 RF1_drivers.index = np.arange(1, len(RF1_drivers)+1)
21
22 # Plot Features
23 RF1_features_df.set_index('Factors', inplace=True)
24 RF1_features_df.plot(kind='barh', figsize=(4, 4))
25 plt.title('Feature Importance with RF Model 1 With `Optimal Parameters`')
26
27 # Merge Description to Encoded Features/Factors
28 RF1_features_df = RF1_features_df.rename_axis('Factors').reset_index() # Reset index before copying it to column
29 #RF1_features_df = pd.merge(RF1_features_df, factor_list, on='Factors', how='left', indicator=True).replace(np.nan, "") # Merge Data
30 #RF1_features_df['Description'] = np.where(RF1_features_df['Description'] == '/', RF1_features_df['Factors'], RF1_features_df['Description'])
31 #RF1_features_df['Description'] = RF1_features_df['Description'].replace(r'_+', ' ', regex=True) # Replace Underscore with single space
32 final_columns = ['Factors', 'Importance', 'Description'] # List of Columns to keep
33 RF1_features_df = RF1_features_df.drop(columns=[col for col in RF1_features_df if col not in final_columns]) # Drop Columns if not in List
34 #print(op_preds)

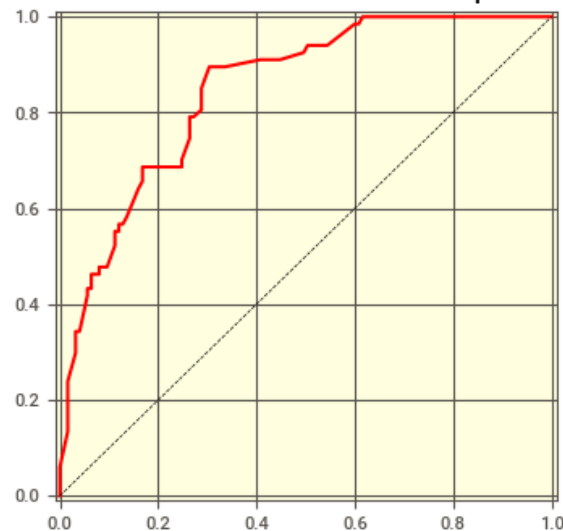
```

```

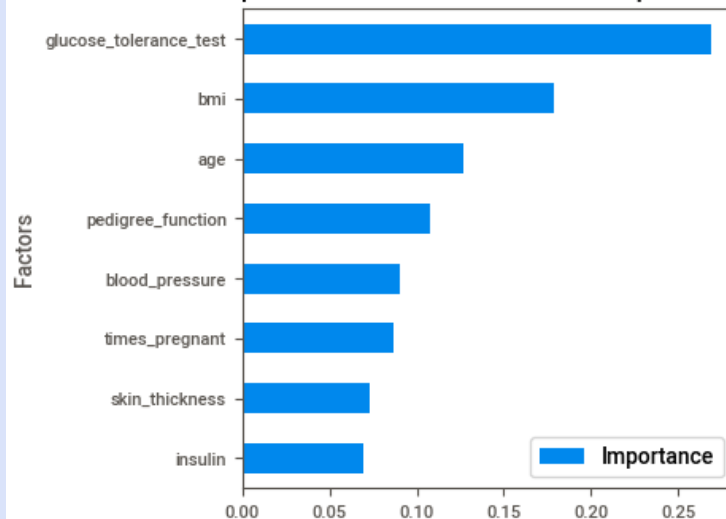
{'type': 'RForest', 'accuracy': 0.7760416666666666, 'recall': 0.6865671641791045, 'precision': 0.6764705882352942, 'fscore': 0.6861732644865174, 'auc': 0.7552835820895524, 'Factors': 'RF1_features_df'}
RF_acc: 0.7552835820895524

```

ROC Curve for RF on PIMA diabetes problem



Feature Importance with RF Model 1 With `Optimal Parameters`





### 3d) Model Choice

#### 3d-i) Result Summary

Although, the assignment required construction of Deep Learning Models (DLM), however, Random Forest was also included as a baseline measure to compare performance of DLM against other unsupervised classification algorithms. As evident from the Result Summary Table below, based on ROC AUC Score, all variations of DLM outperformed Random Forest Model, indicating it to be a better choice for diabetes prediction.

Result Summary Table:

	Optimizer	Input Val Loss	Output Val Loss	Input Val Accuracy	Output Val Accuracy	AUC Score	Epochs	Actual Epochs	val_loss_diff	val_acc_diff	per_epochs_diff	epochs_diff
0	Adagrad	0.636	0.514	0.646	0.719	0.832	900	900	-0.122	0.073	100	0
1	RF	0.000	0.000	0.000	0.000	0.755	0	0	0.000	0.000	0	0
2	Adadelta	0.612	0.588	0.651	0.651	0.738	800	800	-0.025	0.000	100	0
3	Adamax	0.737	0.465	0.432	0.776	0.841	900	546	-0.272	0.344	60	354
4	SGD	0.732	0.468	0.573	0.766	0.851	800	394	-0.264	0.193	49	406
5	Nadam	0.705	0.486	0.651	0.776	0.826	800	241	-0.219	0.125	30	559
6	Adam	0.655	0.498	0.672	0.745	0.821	900	265	-0.157	0.073	29	635
7	RMSprop	0.664	0.520	0.698	0.745	0.811	900	265	-0.144	0.047	29	635

#### 3d-ii) Model Ranking, Choice and Justification

	Optimizer	Input Val Loss	Output Val Loss	Input Val Accuracy	Output Val Accuracy	AUC Score	Epochs	Actual Epochs	val_loss_diff	val_acc_diff	per_epochs_diff	epochs_diff	Rank
0	Adagrad	0.636	0.514	0.646	0.719	0.832	900	900	-0.122	0.073	100	0	1
1	SGD	0.732	0.468	0.573	0.766	0.851	800	394	-0.264	0.193	49	406	2
2	Adamax	0.737	0.465	0.432	0.776	0.841	900	546	-0.272	0.344	60	354	2
3	Nadam	0.705	0.486	0.651	0.776	0.826	800	241	-0.219	0.125	30	559	2
4	Adam	0.655	0.498	0.672	0.745	0.821	900	265	-0.157	0.073	29	635	2
5	RMSprop	0.664	0.520	0.698	0.745	0.811	900	265	-0.144	0.047	29	635	2
6	RF	0.000	0.000	0.000	0.000	0.755	0	0	0.000	0.000	0	0	3
7	Adadelta	0.612	0.588	0.651	0.651	0.738	800	800	-0.025	0.000	100	0	4

#### Recommended Model and Justification:

Since performance measurement is an essential task in machine learning, hence, for classification problems, AUC-ROC Curve is a reliable measure. As such, for visualizing performance of a multi-class diabetes prediction problem, the AUC (Area Under The Curve) ROC (Receiver Operating Characteristics) curve has been utilized since it is one of the most important evaluation metrics for checking any classification model's performance.

Apart from AUC Score, validation loss was also taken into account which allowed final model selection on the premise that validation loss continuously declined to ensure no overfitting has occurred.

Based on the above criterion, Best Model was found to be **Adagrad** with an **roc\_auc** of **0.832**. which successfully satisfied both primary and secondary objectives.

Hence, we recommend this model for Diabetes Prediction.

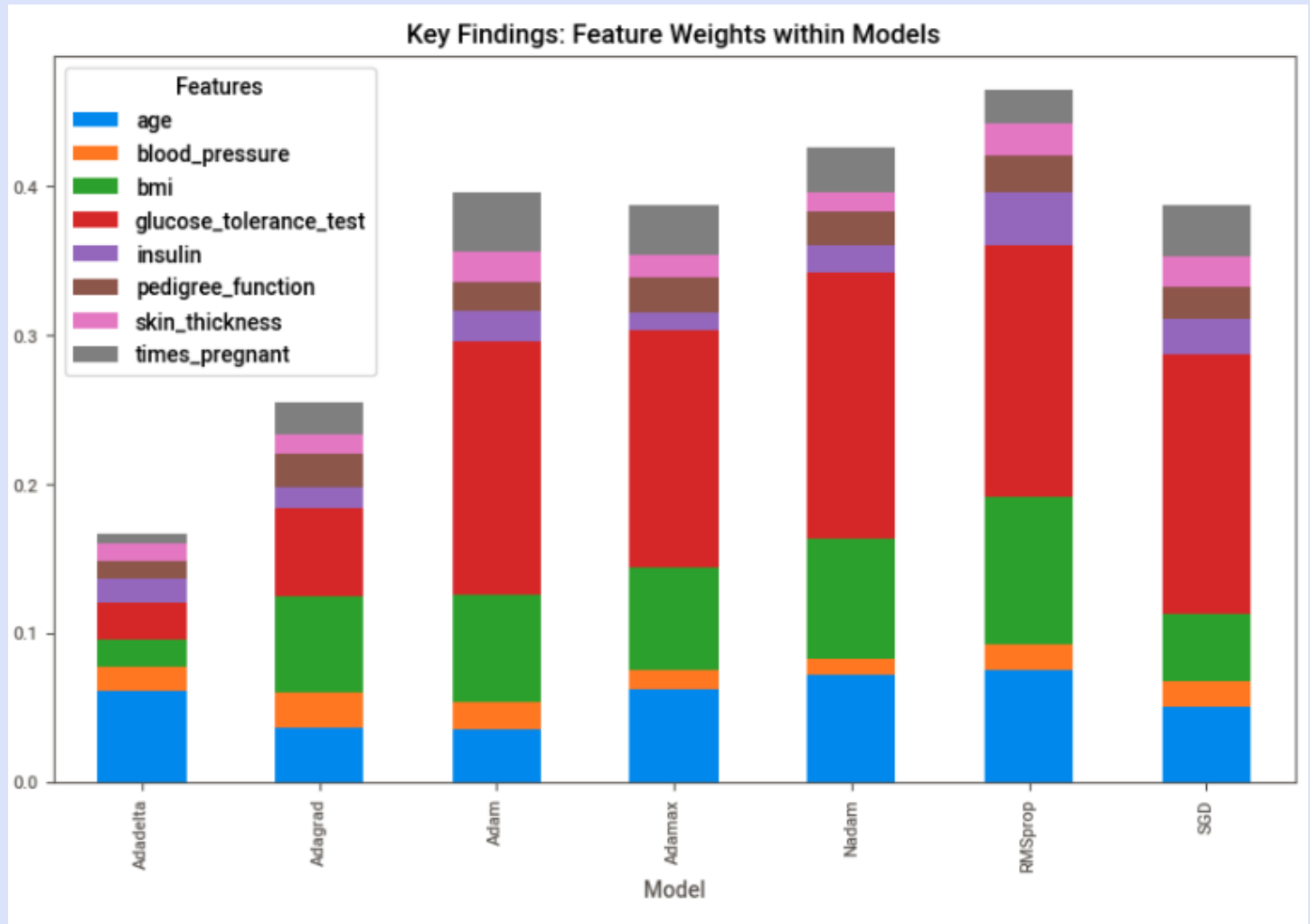
**Adagrad** Model Features in terms of importance are given below:

	Features	Importance
0	bmi	0.064633
1	glucose_tolerance_test	0.059301
2	age	0.036672
3	blood_pressure	0.023516
4	pedigree_function	0.022678
5	times_pregnant	0.021181
6	insulin	0.013991
7	skin_thickness	0.013115

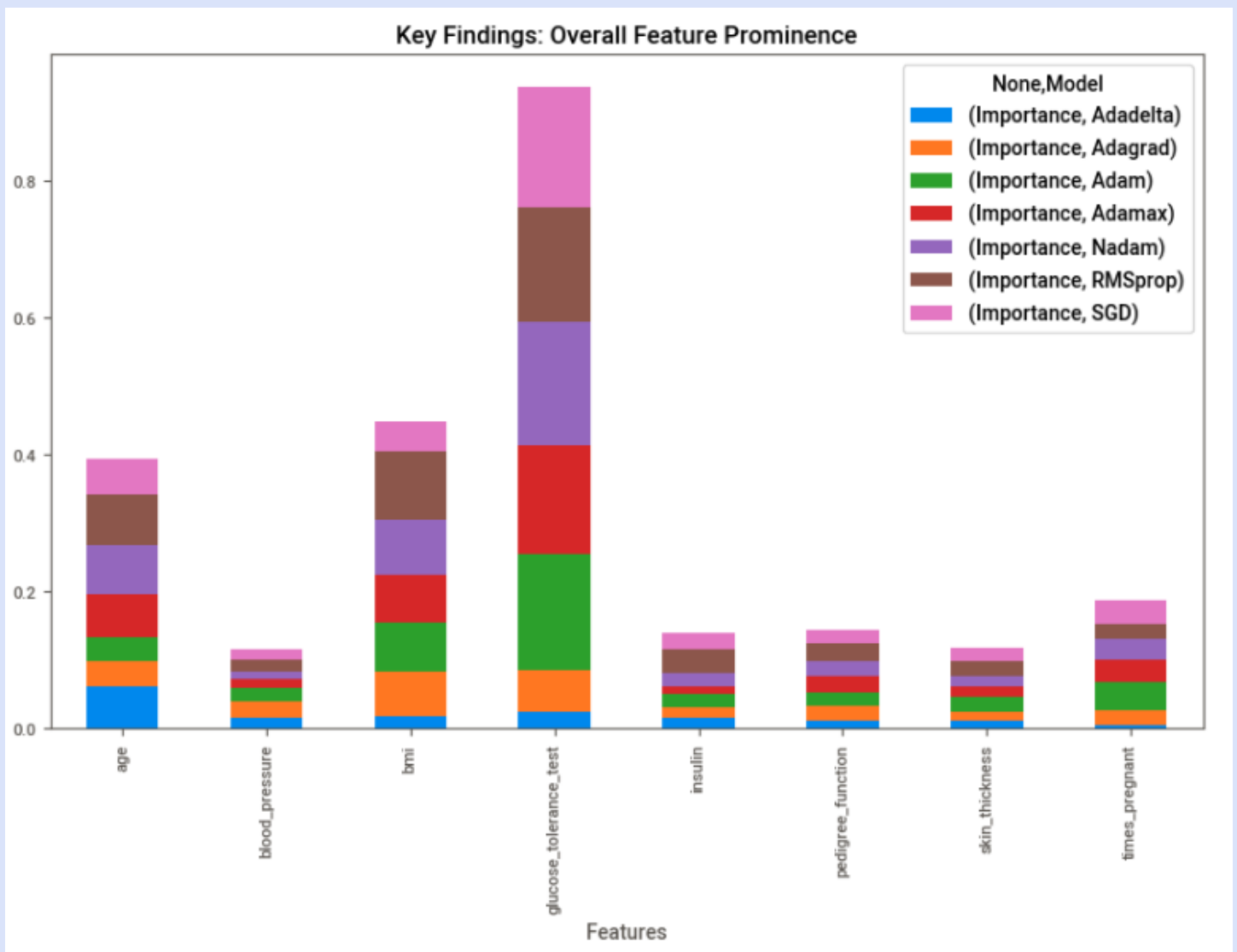
#### 4) Summary Key Findings and Insights

Model output clearly shows glucose tolerance test, bmi and age as most prominent features in most models irrespective of operator variation or auc score. Consequently, feature weights and their individual prominence is depicted as follows:

##### 4a) Feature Weights within Models



## 4b Overall Feature Prominence



## 5) Model Shortcomings and Future Directions

### 5a) Model Shortcomings

Despite its novel approach to automatically pick best model with best parameters, the project is not without its shortcomings.

To begin with, it only utilized a very small dataset with merely nine features.

Hence, it may have discounted other important factors that may impact on set of diabetes.

Furthermore, the project could have employed other deep learning models to ensure even more improved results.

Then, employed Grid Search Hyperparameter Tuning proved to be computationally intensive and could have been replaced by more efficient methods. This also led to a basic model architecture which discounted incorporating other parameters like different learning rates, momentum, etc.

### 5b) Plan of Action to Revisit Analysis

Future recommendations are, hence, as follows:

- 1) Use a larger dataset with more features.
- 2) Employ more unsupervised models like LSTM etc.
- 3) Try other Hyperparameter Tuning methods like Keras Tuner or Bayesian optimization (See, Section 5).
- 4) Include additional parameter search like learning rate, momentum, neuron activation function, etc. to make a more robust model.

## 6 Useful Links

### 6a) Guide to Keras

<https://machinelearningmastery.com/use-keras-deep-learning-models-scikit-learn-python/>

### 6b) Keras Hyperparameter Tuning Approaches

<https://machinelearningmastery.com/grid-search-hyperparameters-deep-learning-models-python-keras/>

<https://medium.datadriveninvestor.com/hyperparameter-tuning-with-deep-learning-grid-search-8630aa45b2da>

<https://towardsai.net/p/l/stop-using-grid-search-the-complete-practical-tutorial-on-keras-tuner>

### 6c) Model Checkpoints

<https://machinelearningmastery.com/how-to-stop-training-deep-neural-networks-at-the-right-time-using-early-stopping/>

<https://pyimagesearch.com/2021/06/30/how-to-use-the-modelcheckpoint-callback-with-keras-and-tensorflow/>

### 6d) Get Feature Importance in Keras

<https://stackoverflow.com/questions/45361559/feature-importance-chart-in-neural-network-using-keras-in-python>

<https://www.kdnuggets.com/2020/01/explaining-black-box-models-ensemble-deep-learning-lime-shap.html>

### 6e) Shap Plots: Some Good Examples

[https://shap.readthedocs.io/en/latest/example\\_notebooks/api\\_examples/plots/decision\\_plot.html](https://shap.readthedocs.io/en/latest/example_notebooks/api_examples/plots/decision_plot.html)

### 6f) Matplot Lib

[https://matplotlib.org/stable/gallery/subplots\\_axes\\_and\\_figures/subplot.html](https://matplotlib.org/stable/gallery/subplots_axes_and_figures/subplot.html)

<https://medium.com/mlearning-ai/shap-force-plots-for-classification-d30be430e195>

### 6f) Other Useful Models

<https://github.com/AI-MOO/IBM-Machine-Learning-Professional-Certificate>

<https://github.com/topics/ibm-machine-learning>

<https://samyzaf.com/ML/pima/pima.html>

## 7) Github Links

### 7a) Link to Main Folder

<https://github.com/FATIMASP/IBM-MACHINE-LEARNING-CERTIFICATION/tree/main/Deep%20Learning%20and%20Reinforcement%20Learning>

### 7b) Link to Assignment Notebook

<https://github.com/FATIMASP/IBM-MACHINE-LEARNING-CERTIFICATION/blob/main/Deep%20Learning%20and%20Reinforcement%20Learning/DEEP%20LEARNING%20%26%20REINFORCEMENT%20LEARNING%20FINAL%20MODEL.ipynb>

## 8) References

Naz, H. & Ahuja, S., 2020. Deep learning approach for diabetes prediction using PIMA Indian dataset. *Journal of Diabetes and Metabolic Disorders*, 19(1), p. 391–403.



The grading will center around 5 main points:

1. Does the report include a section describing the data?
2. Does the report include a paragraph detailing the main objective(s) of this analysis?
3. Does the report include a section with variations of a Deep Learning model and specifies which one is the model that best suits the main objective(s) of this analysis?
4. Does the report include a clear and well presented section with key findings related to the main objective(s) of the analysis?
5. Does the report highlight possible flaws in the model and a plan of action to revisit this analysis with additional data or different modeling techniques?