```python
In [1]:  import pandas as pd
         import matplotlib.pyplot as plt
         import seaborn as sns
         import numpy as np
         sns.set_theme(color_codes=True)
```

```python
In [2]:  df = pd.read_csv('diabetes_prediction_dataset.csv')
         df.head()
```

Out[2]:

|   | gender | age | hypertension | heart_disease | smoking_history | bmi | HbA1c_level | blood_glucose_level |
|---|--------|-----|--------------|---------------|-----------------|-----|-------------|---------------------|
| 0 | Female | 80.0 | 0 | 1 | never | 25.19 | 6.6 | 140 |
| 1 | Female | 54.0 | 0 | 0 | No Info | 27.32 | 6.6 | 80 |
| 2 | Male | 28.0 | 0 | 0 | never | 27.32 | 5.7 | 158 |
| 3 | Female | 36.0 | 0 | 0 | current | 23.45 | 5.0 | 155 |
| 4 | Male | 76.0 | 1 | 1 | current | 20.14 | 4.8 | 155 |

```python
In [3]:  #Cheeck Object data types unique value
         df.select_dtypes(include='object').nunique()
```

```
Out[3]:  gender            3
         smoking_history   6
         dtype: int64
```

# Exploratory Data Analysis

# Categorial Data

In [4]:
```python
# list of categorical variables to plot
cat_vars = ['gender', 'smoking_history', 'hypertension', 'heart_disease']

# create figure with subplots
fig, axs = plt.subplots(nrows=2, ncols=2, figsize=(15, 10))
axs = axs.flatten()

# create barplot for each categorical variable
for i, var in enumerate(cat_vars):
    sns.countplot(x=var, hue='diabetes', data=df, ax=axs[i])
    axs[i].set_xticklabels(axs[i].get_xticklabels(), rotation=90)

# adjust spacing between subplots
fig.tight_layout()

# show plot
plt.show()
```
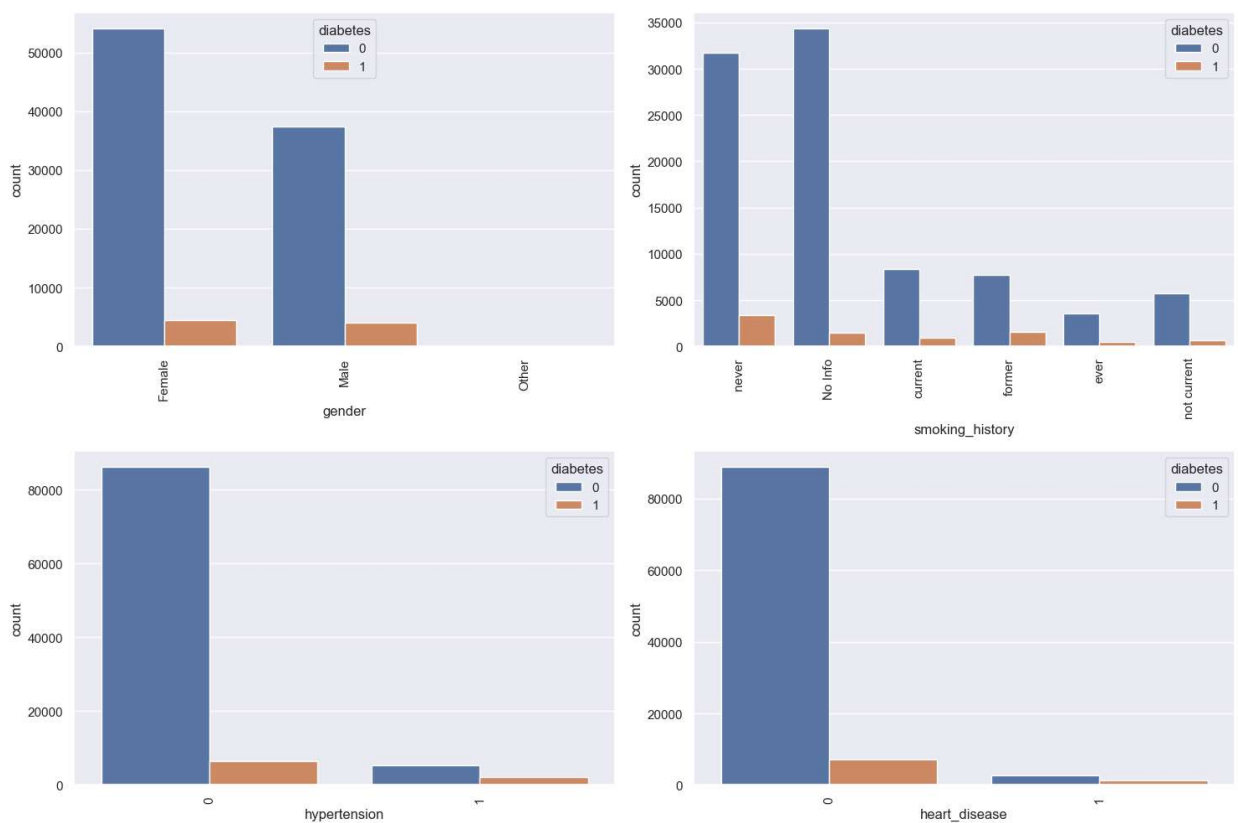
```
In [5]:  import warnings
         warnings.filterwarnings("ignore")
         # get list of categorical variables
         cat_vars = ['gender', 'smoking_history', 'hypertension', 'heart_disease']

         # create figure with subplots
         fig, axs = plt.subplots(nrows=2, ncols=2, figsize=(15, 10))
         axs = axs.flatten()

         # create histplot for each categorical variable
         for i, var in enumerate(cat_vars):
             sns.histplot(x=var, hue='diabetes', data=df, ax=axs[i], multiple="fill", kde=Fals
             axs[i].set_xticklabels(df[var].unique(), rotation=90)
             axs[i].set_xlabel(var)

         # adjust spacing between subplots
         fig.tight_layout()

         # show plot
         plt.show()
```
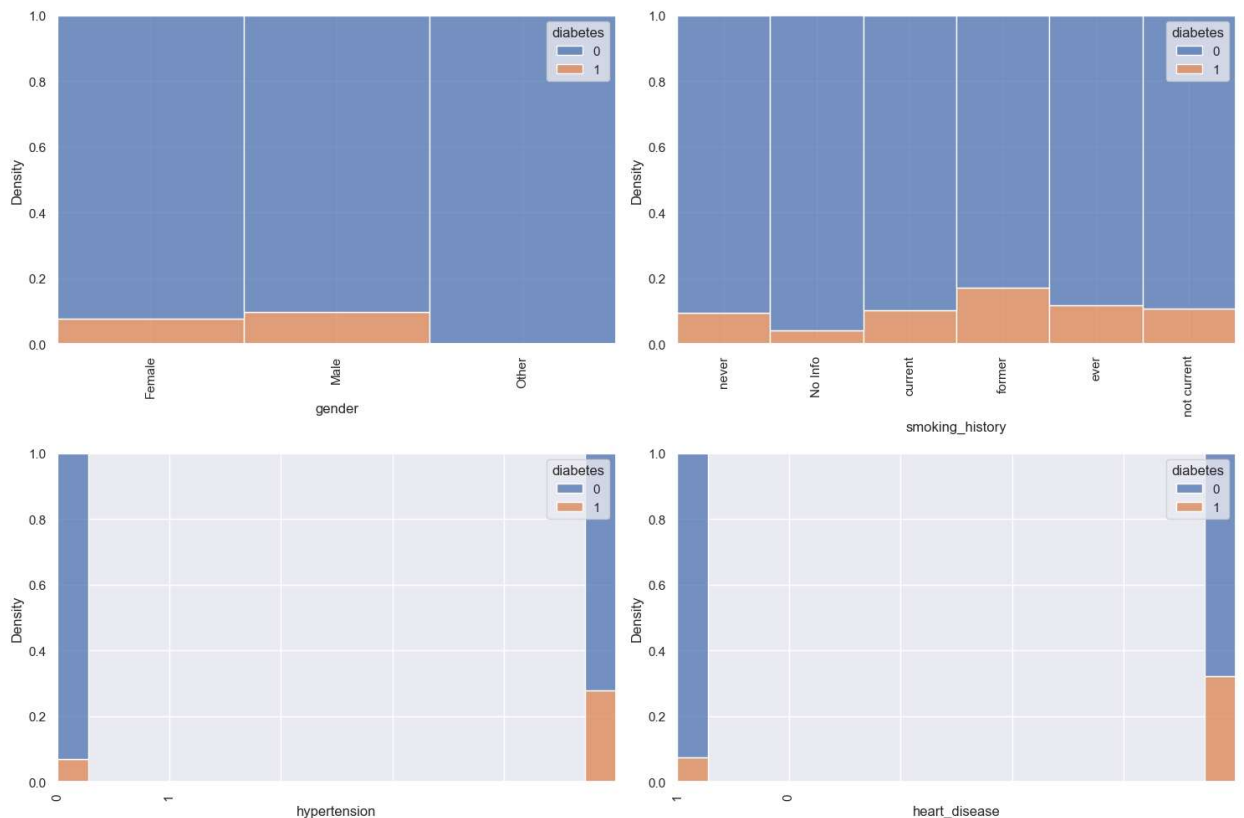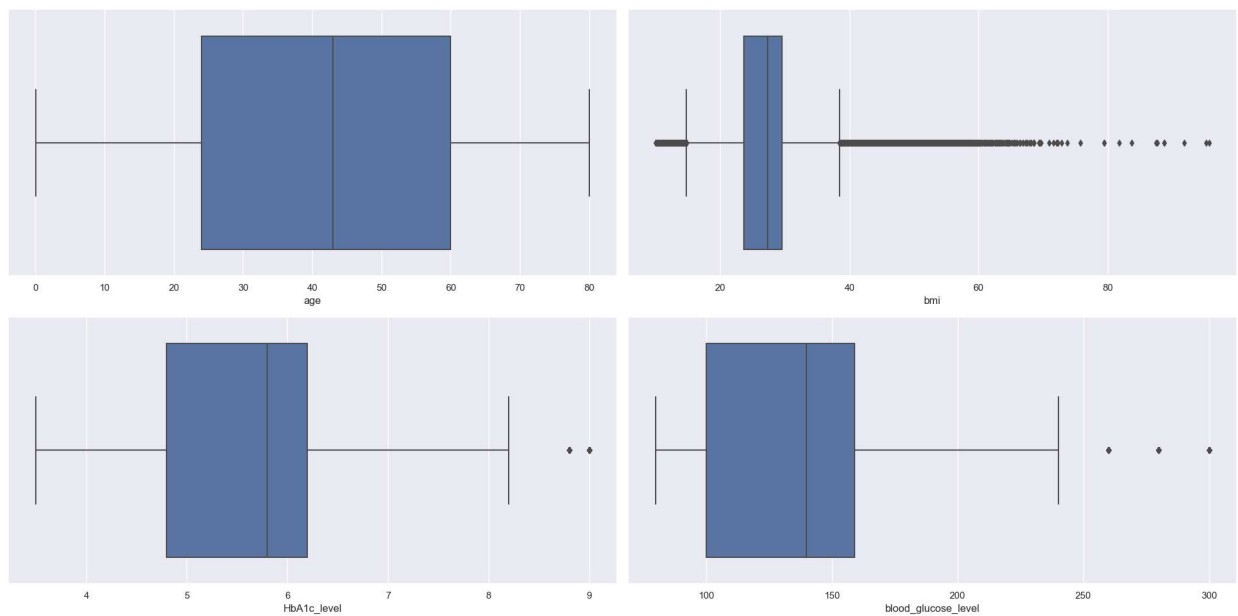


# Numerical EDA

In [6]:
```python
num_vars = ['age', 'bmi', 'HbA1c_level', 'blood_glucose_level']

fig, axs = plt.subplots(nrows=2, ncols=2, figsize=(20, 10))
axs = axs.flatten()

for i, var in enumerate(num_vars):
    sns.boxplot(x=var, data=df, ax=axs[i])

fig.tight_layout()

plt.show()
```
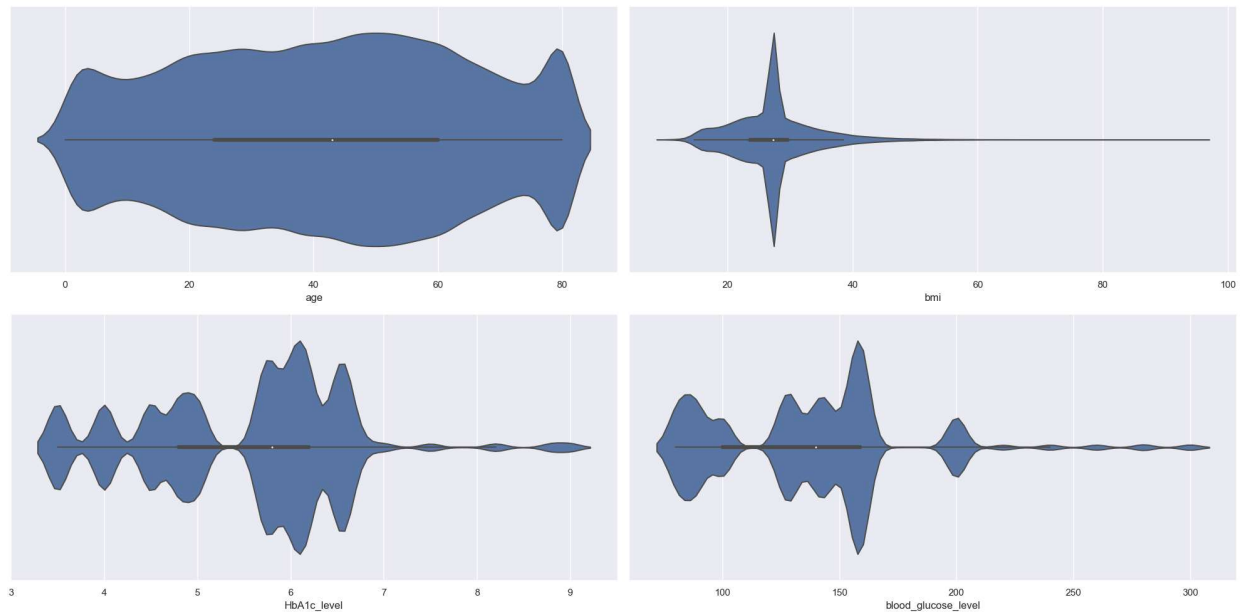
In [7]:
```python
num_vars = ['age', 'bmi', 'HbA1c_level', 'blood_glucose_level']

fig, axs = plt.subplots(nrows=2, ncols=2, figsize=(20, 10))
axs = axs.flatten()

for i, var in enumerate(num_vars):
    sns.violinplot(x=var, data=df, ax=axs[i])

fig.tight_layout()
plt.show()
```
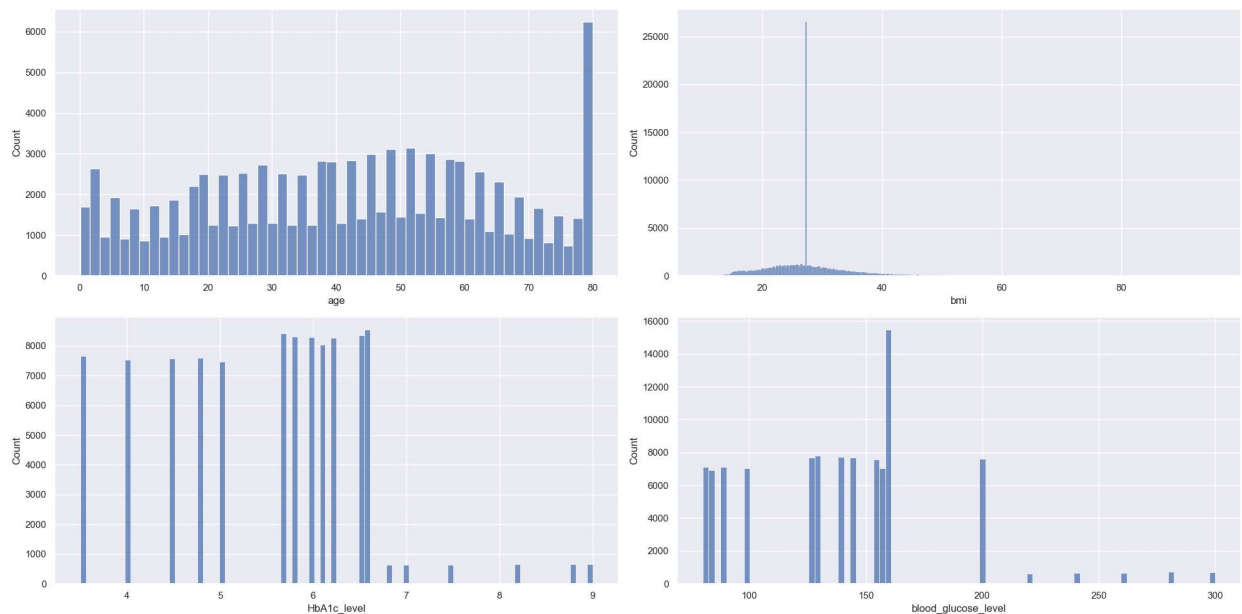
```python
In [8]:  num_vars = ['age', 'bmi', 'HbA1c_level', 'blood_glucose_level']

         fig, axs = plt.subplots(nrows=2, ncols=2, figsize=(20, 10))
         axs = axs.flatten()

         for i, var in enumerate(num_vars):
             sns.histplot(x=var, data=df, ax=axs[i])

         fig.tight_layout()

         plt.show()
```



# Data Preprocessing

```python
In [9]:  df.shape
```

```
Out[9]:  (100000, 9)
```

```python
In [10]:  #Check missing value
          check_missing = df.isnull().sum() * 100 / df.shape[0]
          check_missing[check_missing > 0].sort_values(ascending=False)
```

```
Out[10]:  Series([], dtype: float64)
```

# Label encoding each categorial column

```python
In [11]:  # Loop over each column in the DataFrame where dtype is 'object'
          for col in df.select_dtypes(include=['object']).columns:

              # Print the column name and the unique values
              print(f"{col}: {df[col].unique()}")
```

```
gender: ['Female' 'Male' 'Other']
smoking_history: ['never' 'No Info' 'current' 'former' 'ever' 'not current']
```

In [12]:
```python
from sklearn import preprocessing

# Loop over each column in the DataFrame where dtype is 'object'
for col in df.select_dtypes(include=['object']).columns:

    # Initialize a LabelEncoder object
    label_encoder = preprocessing.LabelEncoder()

    # Fit the encoder to the unique values in the column
    label_encoder.fit(df[col].unique())

    # Transform the column using the encoder
    df[col] = label_encoder.transform(df[col])

    # Print the column name and the unique encoded values
    print(f"{col}: {df[col].unique()}")
```
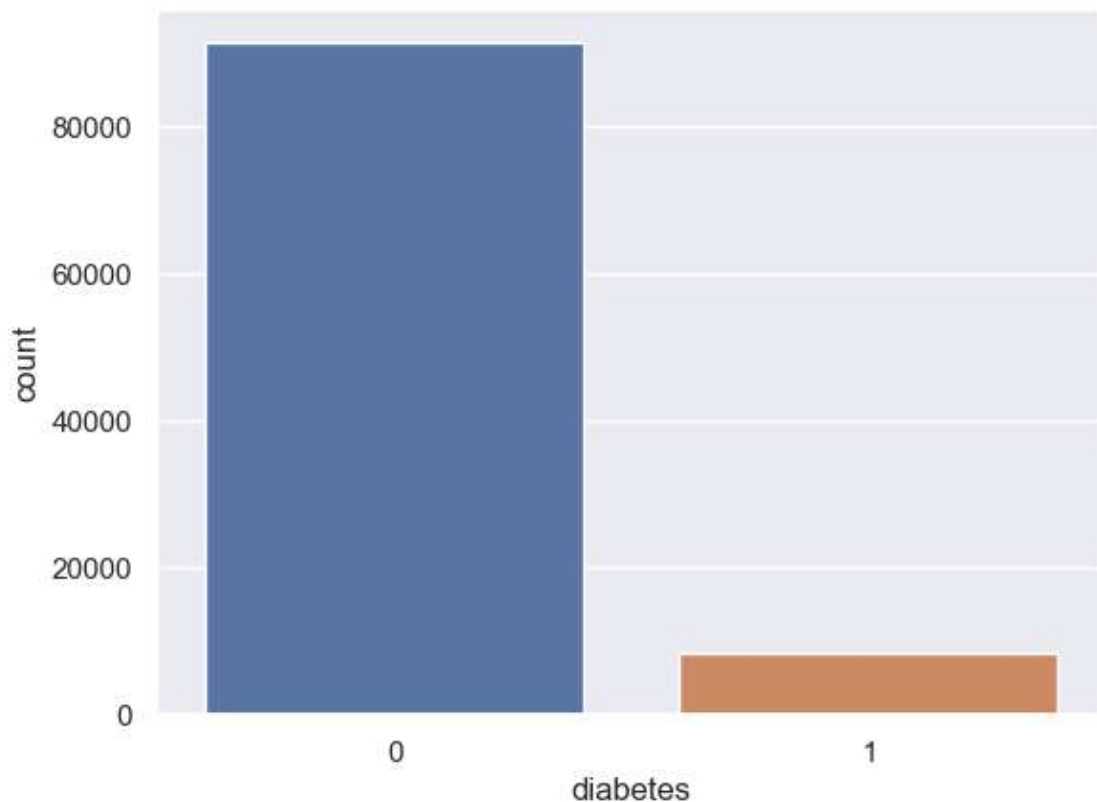
```
gender: [0 1 2]
smoking_history: [4 0 1 3 2 5]
```

# Check the Label value 'Diabetes' if its balanced or not

In [13]:
```python
sns.countplot(df['diabetes'])
df['diabetes'].value_counts()
```

Out[13]:
```
0    91500
1     8500
Name: diabetes, dtype: int64
```

In [14]:
```python
# Undersampling majority class
from imblearn.under_sampling import RandomUnderSampler

X = df.drop('diabetes', axis=1)
y = df['diabetes']

rus = RandomUnderSampler(random_state=42)
X_resampled, y_resampled = rus.fit_resample(X, y)

# create new DataFrame with undersampled data
df_resampled = pd.concat([X_resampled, y_resampled], axis=1)
```
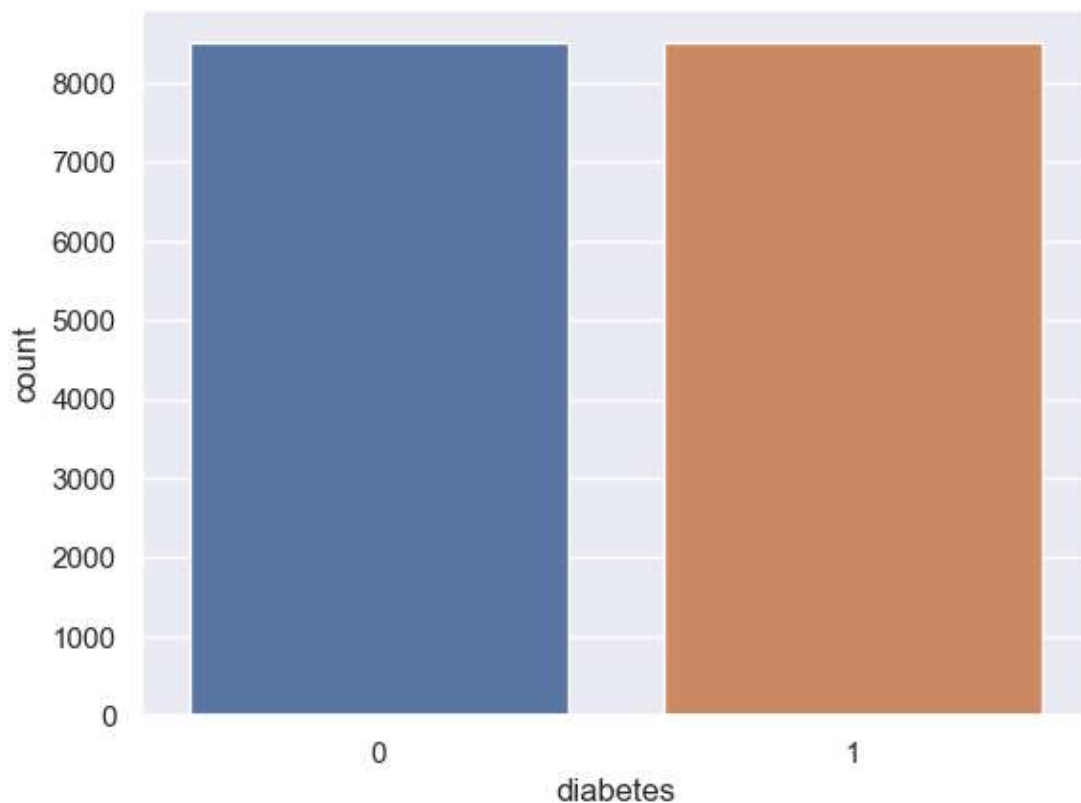
In [15]:
```python
sns.countplot(df_resampled['diabetes'])
df_resampled['diabetes'].value_counts()
```

Out[15]:
```
0    8500
1    8500
Name: diabetes, dtype: int64
```



# Check the Outliers using Z-Score

4/26/23, 7:10 PM Diabetes Prediction - Jupyter Notebook

```python
from scipy import stats

# define a function to remove outliers using z-score for only selected numerical colum
def remove_outliers(df_resampled, cols, threshold=3):
    # loop over each selected column
    for col in cols:
        # calculate z-score for each data point in selected column
        z = np.abs(stats.zscore(df_resampled[col]))
        # remove rows with z-score greater than threshold in selected column
        df_resampled = df_resampled[(z < threshold) | (df_resampled[col].isnull())]
    return df_resampled
```

In [17]:
```python
selected_cols = ['bmi', 'HbA1c_level', 'blood_glucose_level']
df_clean = remove_outliers(df_resampled, selected_cols)
df_clean.shape
```

Out[17]: (16786, 9)

# Heatmap Correlattion

In [18]:
```python
#Correlation Heatmap
plt.figure(figsize=(20, 16))
sns.heatmap(df_clean.corr(), fmt='.2g', annot=True)
```

Out[18]: <AxesSubplot:>



## Train Test Split

In [19]:
```python
X = df_clean.drop('diabetes', axis=1)
y = df_clean['diabetes']
```

In [20]:
```python
#test size 20% and train size 80%
from sklearn.model_selection import train_test_split
from sklearn.metrics import accuracy_score
X_train, X_test, y_train, y_test = train_test_split(X,y, test_size=0.2,random_state=0
```

## Decision Tree

In [21]:
```python
from sklearn.tree import DecisionTreeClassifier
from sklearn.model_selection import GridSearchCV
dtree = DecisionTreeClassifier()
param_grid = {
    'max_depth': [3, 4, 5, 6, 7, 8],
    'min_samples_split': [2, 3, 4],
    'min_samples_leaf': [1, 2, 3, 4]
}

# Perform a grid search with cross-validation to find the best hyperparameters
grid_search = GridSearchCV(dtree, param_grid, cv=5)
grid_search.fit(X_train, y_train)

# Print the best hyperparameters
print(grid_search.best_params_)
```

```
{'max_depth': 8, 'min_samples_leaf': 1, 'min_samples_split': 2}
```

In [22]:
```python
from sklearn.tree import DecisionTreeClassifier
dtree = DecisionTreeClassifier(random_state=0, max_depth=8, min_samples_leaf=1, min_s
dtree.fit(X_train, y_train)
```

Out[22]: DecisionTreeClassifier(max_depth=8, random_state=0)

In [23]:
```python
y_pred = dtree.predict(X_test)
print("Accuracy Score :", round(accuracy_score(y_test, y_pred)*100 ,2), "%")
```
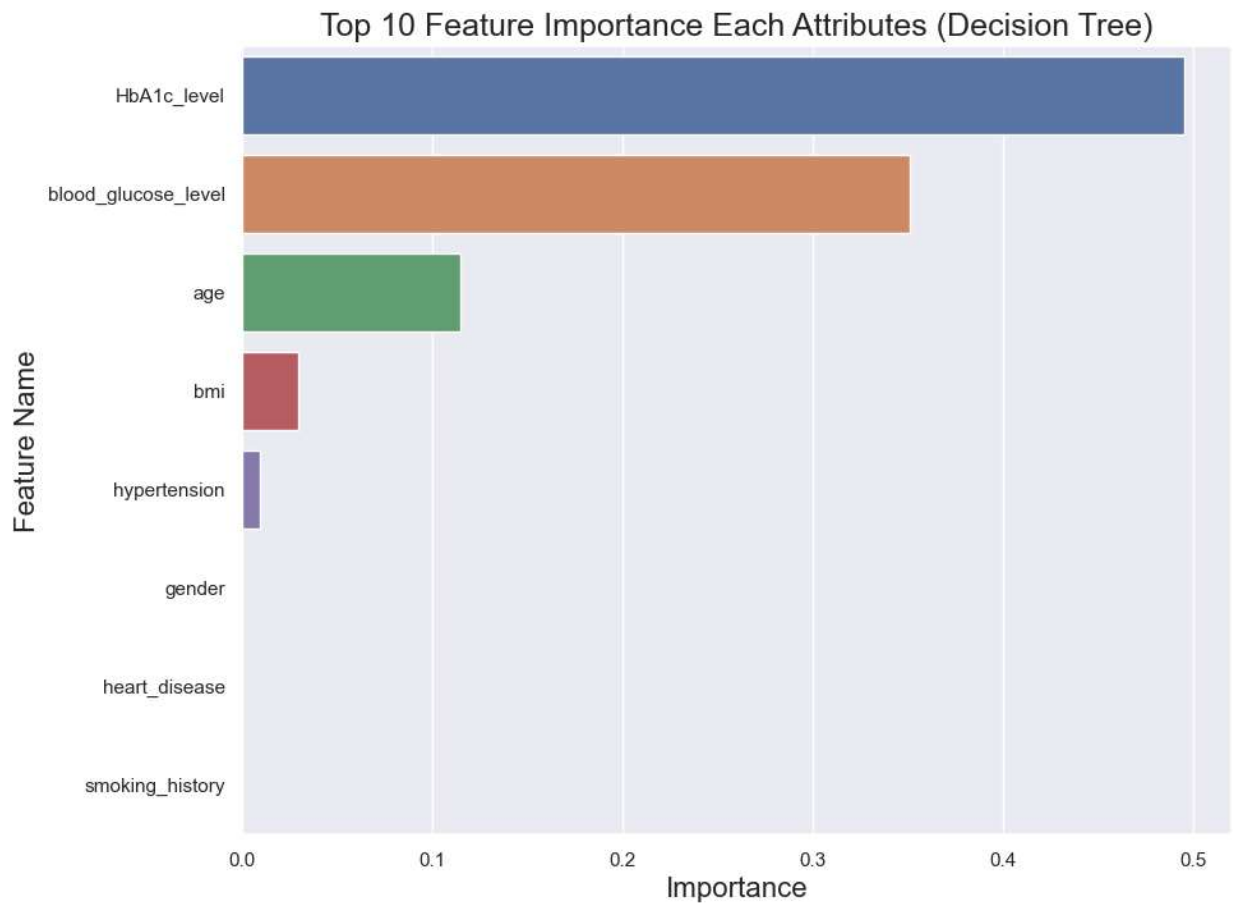
```
Accuracy Score : 89.79 %
```

In [24]:
```python
from sklearn.metrics import accuracy_score, f1_score, precision_score, recall_score,
print('F-1 Score : ',(f1_score(y_test, y_pred, average='micro')))
print('Precision Score : ',(precision_score(y_test, y_pred, average='micro')))
print('Recall Score : ',(recall_score(y_test, y_pred, average='micro')))
print('Jaccard Score : ',(jaccard_score(y_test, y_pred, average='micro')))
print('Log Loss : ',(log_loss(y_test, y_pred)))
```
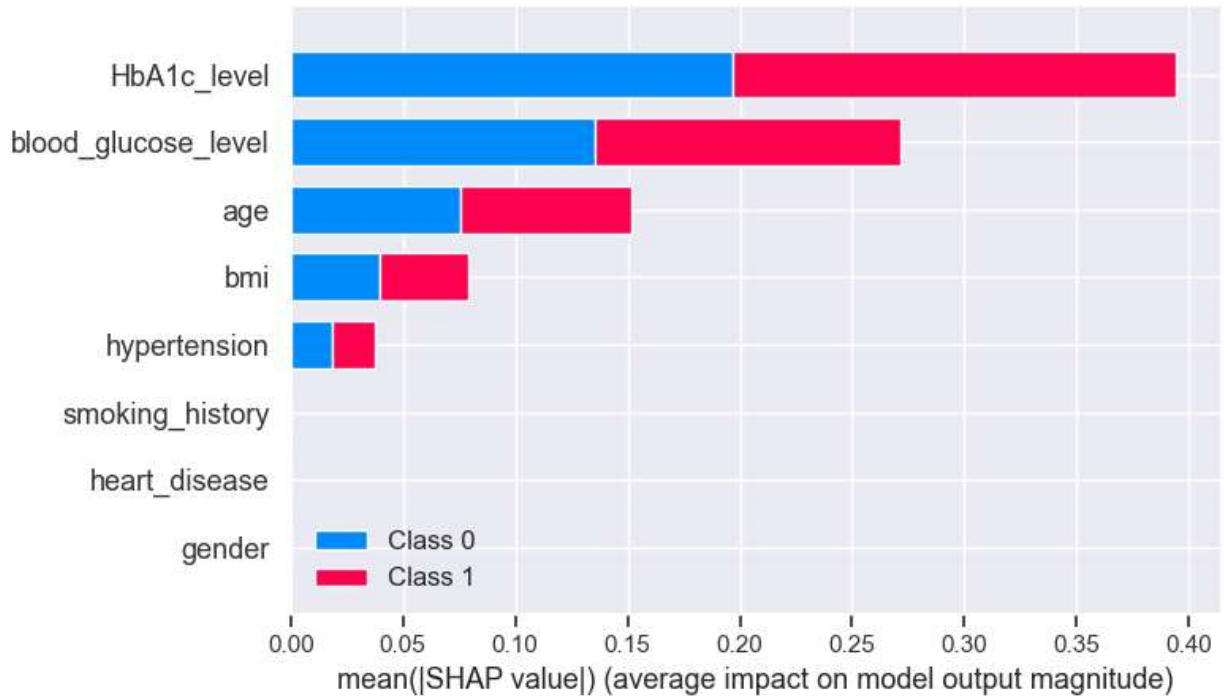
```
F-1 Score :   0.8978558665872544
Precision Score :   0.8978558665872544
Recall Score :   0.8978558665872544
Jaccard Score :   0.8146446906241557
Log Loss :   3.5279805311928363
```

In [25]:
```python
imp_df = pd.DataFrame({
    "Feature Name": X_train.columns,
    "Importance": dtree.feature_importances_
})
fi = imp_df.sort_values(by="Importance", ascending=False)

fi2 = fi.head(10)
plt.figure(figsize=(10,8))
sns.barplot(data=fi2, x='Importance', y='Feature Name')
plt.title('Top 10 Feature Importance Each Attributes (Decision Tree)', fontsize=18)
plt.xlabel ('Importance', fontsize=16)
plt.ylabel ('Feature Name', fontsize=16)
plt.show()
```
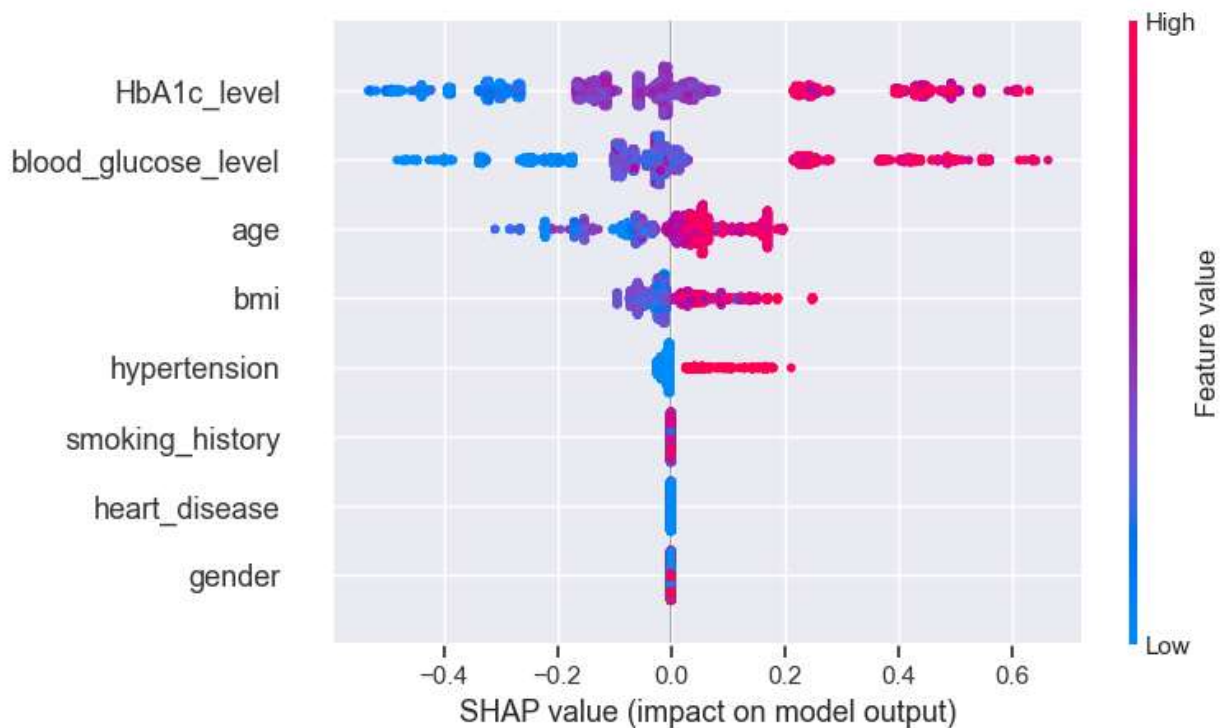
In [26]:
```python
import shap
explainer = shap.TreeExplainer(dtree)
shap_values = explainer.shap_values(X_test)
shap.summary_plot(shap_values, X_test)
```
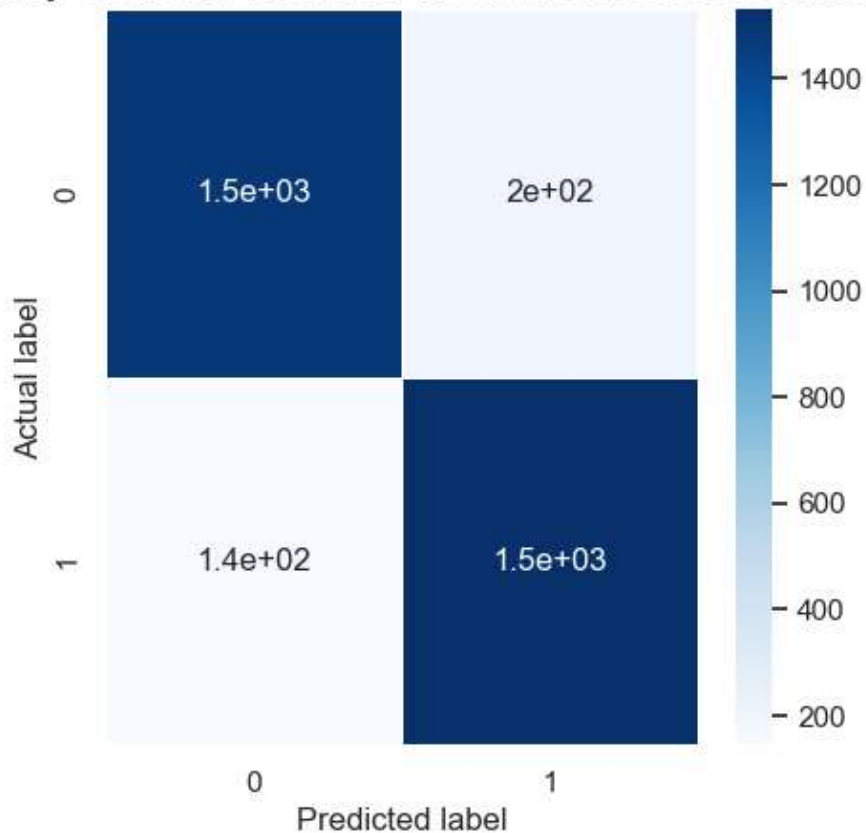


In [27]:
```python
# compute SHAP values
explainer = shap.TreeExplainer(dtree)
shap_values = explainer.shap_values(X_test)
shap.summary_plot(shap_values[1], X_test.values, feature_names = X_test.columns)
```

In [28]:
```python
from sklearn.metrics import confusion_matrix
cm = confusion_matrix(y_test, y_pred)
plt.figure(figsize=(5,5))
sns.heatmap(data=cm,linewidths=.5, annot=True,  cmap = 'Blues')
plt.ylabel('Actual label')
plt.xlabel('Predicted label')
all_sample_title = 'Accuracy Score for Decision Tree: {0}'.format(dtree.score(X_test,
plt.title(all_sample_title, size = 15)
```

Out[28]: Text(0.5, 1.0, 'Accuracy Score for Decision Tree: 0.8978558665872544')

In [29]:
```python
from sklearn.metrics import roc_curve, roc_auc_score
y_pred_proba = dtree.predict_proba(X_test)[:][:,1]

df_actual_predicted = pd.concat([pd.DataFrame(np.array(y_test), columns=['y_actual'])
df_actual_predicted.index = y_test.index

fpr, tpr, tr = roc_curve(df_actual_predicted['y_actual'], df_actual_predicted['y_pred
auc = roc_auc_score(df_actual_predicted['y_actual'], df_actual_predicted['y_pred_prob

plt.plot(fpr, tpr, label='AUC = %0.4f' %auc)
plt.plot(fpr, fpr, linestyle = '--', color='k')
plt.xlabel('False Positive Rate')
plt.ylabel('True Positive Rate')
plt.title('ROC Curve', size = 15)
plt.legend()
```
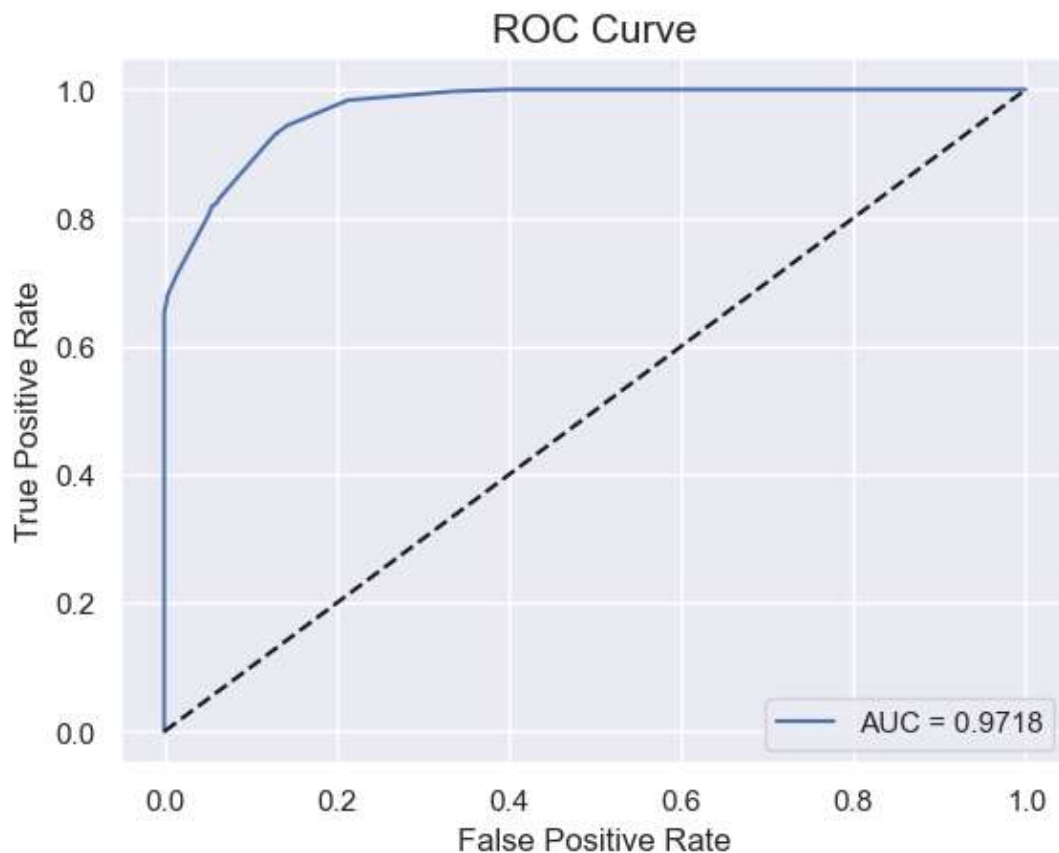
Out[29]:  <matplotlib.legend.Legend at 0x2578c0abd00>



# Random Forest

In [30]:
```python
from sklearn.ensemble import RandomForestClassifier
from sklearn.model_selection import GridSearchCV
rfc = RandomForestClassifier()
param_grid = {
    'n_estimators': [100, 200],
    'max_depth': [None, 5, 10],
    'max_features': ['sqrt', 'log2', None]
}

# Perform a grid search with cross-validation to find the best hyperparameters
grid_search = GridSearchCV(rfc, param_grid, cv=5)
grid_search.fit(X_train, y_train)

# Print the best hyperparameters
print(grid_search.best_params_)
```

{'max_depth': 10, 'max_features': 'sqrt', 'n_estimators': 100}

In [31]:
```python
from sklearn.ensemble import RandomForestClassifier
rfc = RandomForestClassifier(random_state=0, max_features='sqrt', n_estimators=100, ma
rfc.fit(X_train, y_train)
```

Out[31]: RandomForestClassifier(max_depth=10, max_features='sqrt', random_state=0)

In [32]:
```python
y_pred = rfc.predict(X_test)
print("Accuracy Score :", round(accuracy_score(y_test, y_pred)*100 ,2), "%")
```

Accuracy Score : 90.41 %

In [33]:
```python
from sklearn.metrics import accuracy_score, f1_score, precision_score, recall_score, 
print('F-1 Score : ',(f1_score(y_test, y_pred, average='micro')))
print('Precision Score : ',(precision_score(y_test, y_pred, average='micro')))
print('Recall Score : ',(recall_score(y_test, y_pred, average='micro')))
print('Jaccard Score : ',(jaccard_score(y_test, y_pred, average='micro')))
print('Log Loss : ',(log_loss(y_test, y_pred)))
```
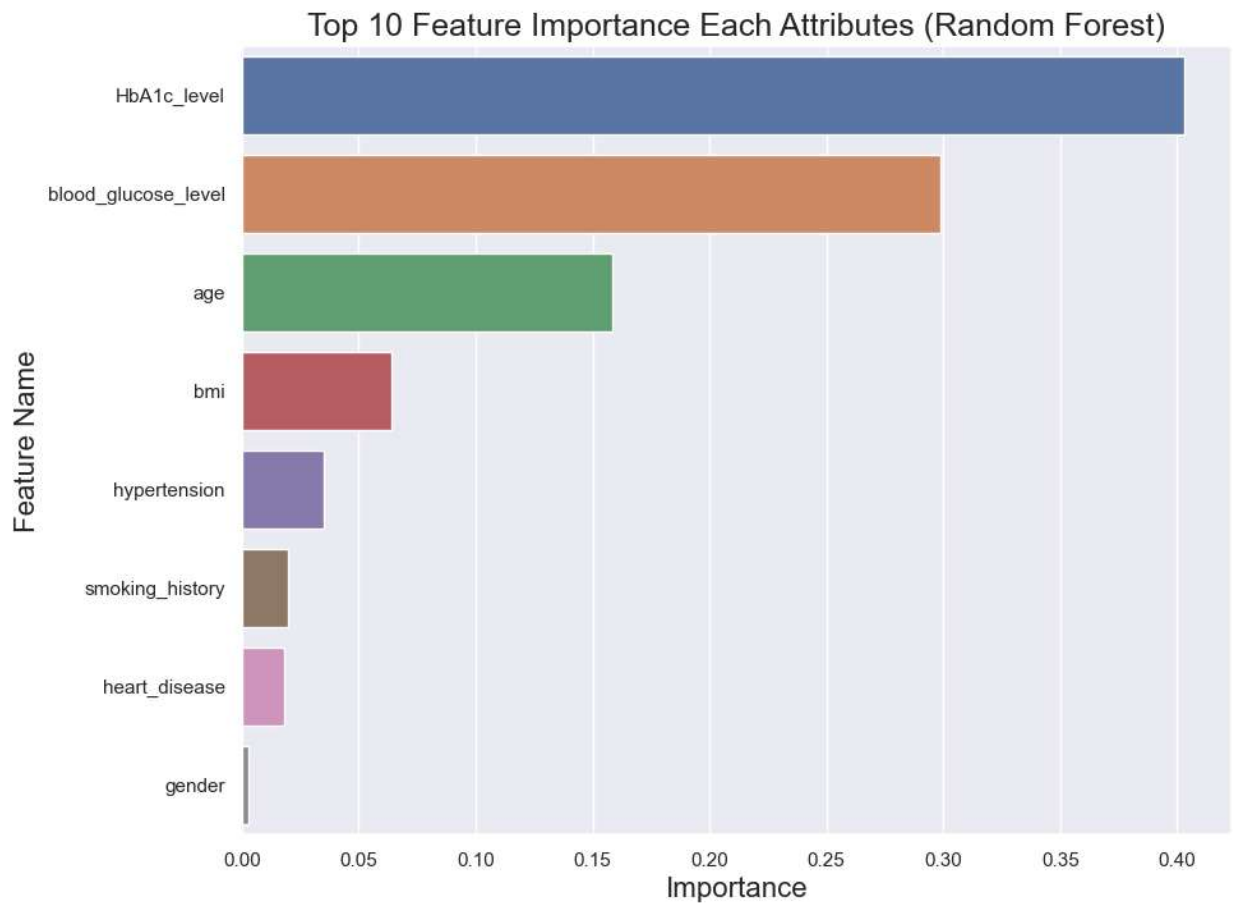
```
F-1 Score :  0.9041095890410958
Precision Score :  0.9041095890410958
Recall Score :  0.9041095890410958
Jaccard Score :  0.825
Log Loss :  3.3119765137402606
```
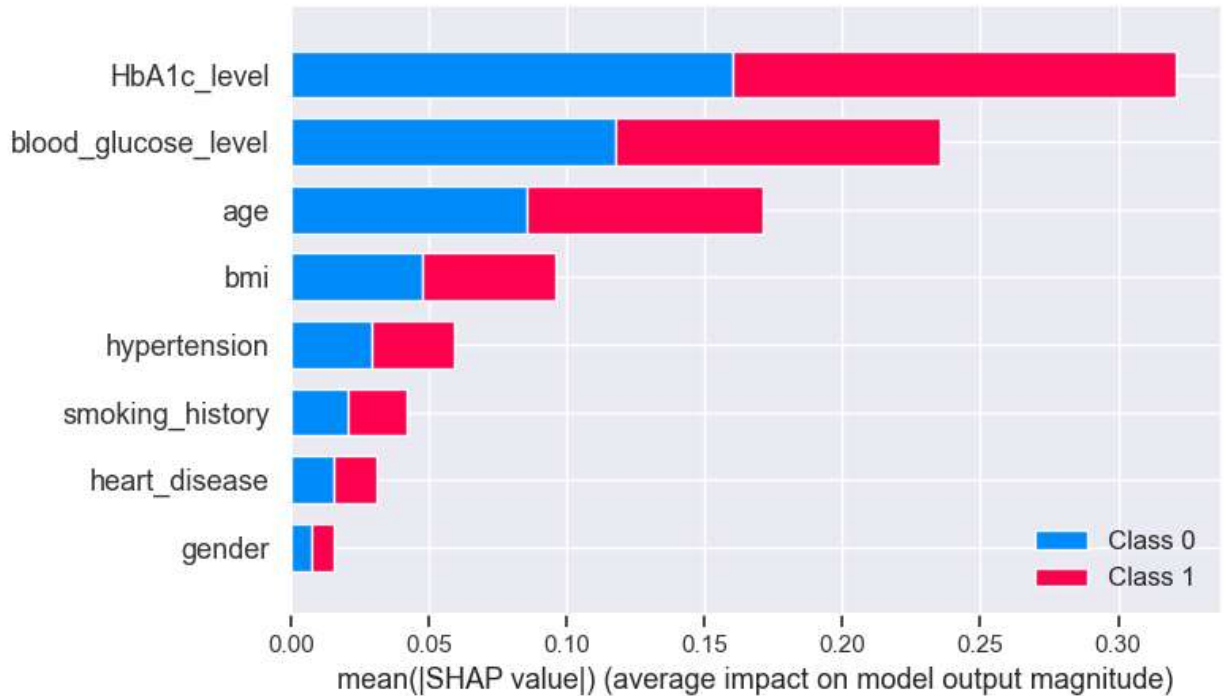
In [34]:
```python
imp_df = pd.DataFrame({
    "Feature Name": X_train.columns,
    "Importance": rfc.feature_importances_
})
fi = imp_df.sort_values(by="Importance", ascending=False)

fi2 = fi.head(10)
plt.figure(figsize=(10,8))
sns.barplot(data=fi2, x='Importance', y='Feature Name')
plt.title('Top 10 Feature Importance Each Attributes (Random Forest)', fontsize=18)
plt.xlabel ('Importance', fontsize=16)
plt.ylabel ('Feature Name', fontsize=16)
plt.show()
```
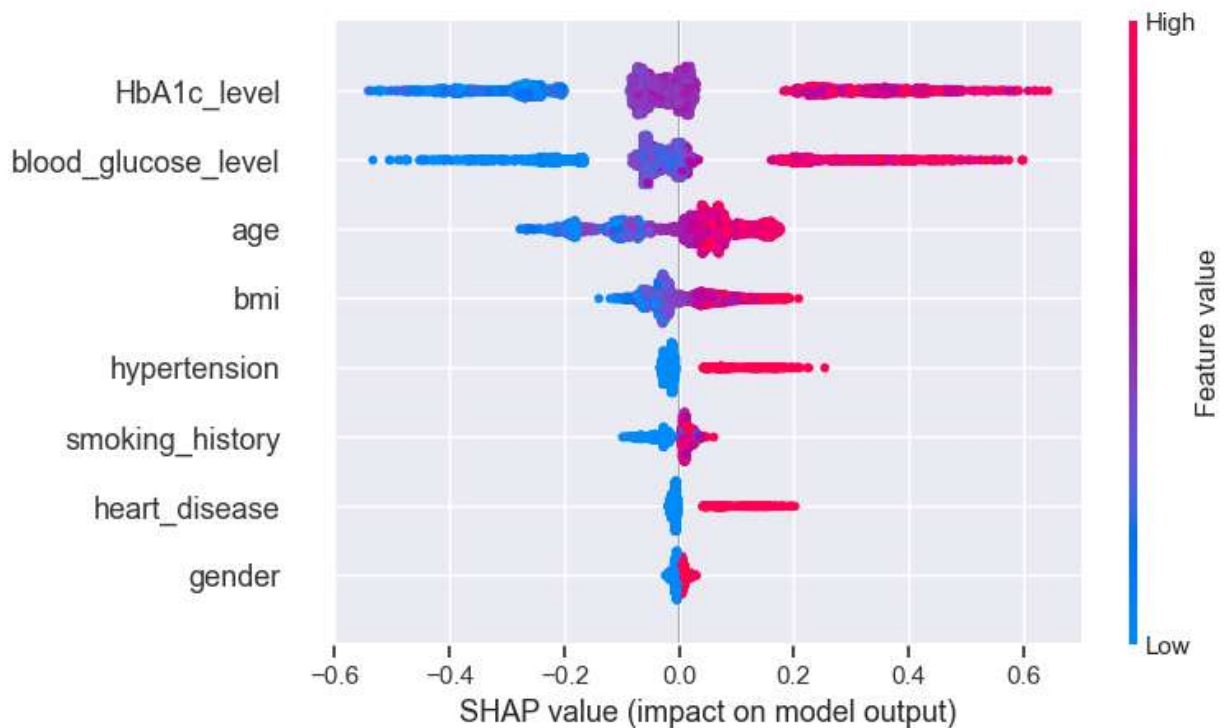


Top 10 Feature Importance Each Attributes (Random Forest)

In [35]:
```python
import shap
explainer = shap.TreeExplainer(rfc)
shap_values = explainer.shap_values(X_test)
shap.summary_plot(shap_values, X_test)
```



In [36]:
```python
# compute SHAP values
explainer = shap.TreeExplainer(rfc)
shap_values = explainer.shap_values(X_test)
shap.summary_plot(shap_values[1], X_test.values, feature_names = X_test.columns)
```
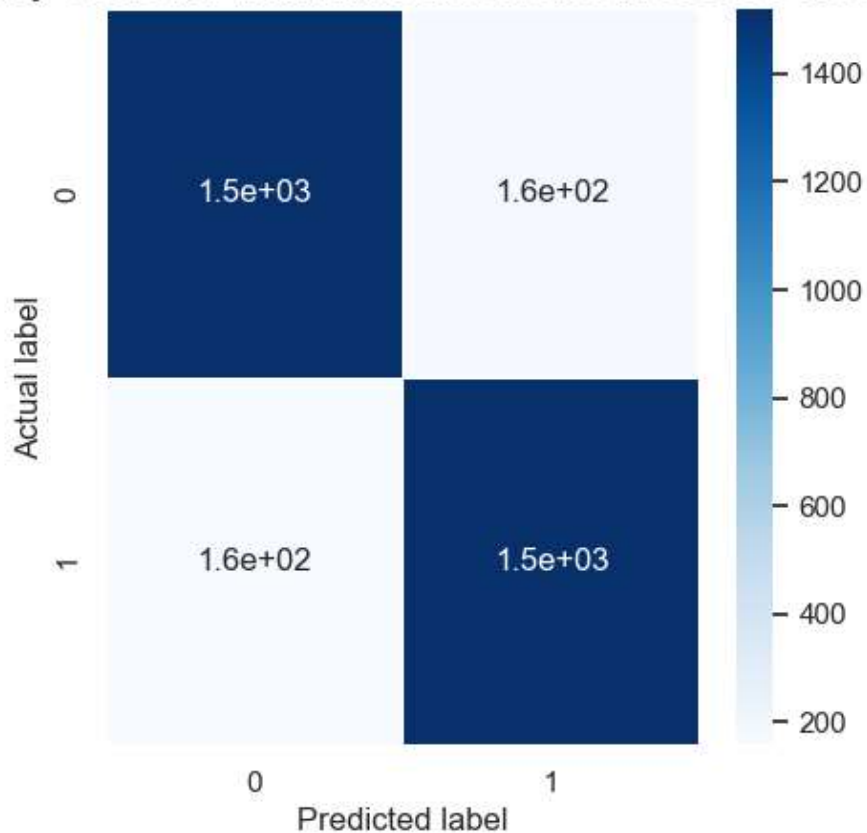
In [37]: 
```python
from sklearn.metrics import confusion_matrix
cm = confusion_matrix(y_test, y_pred)
plt.figure(figsize=(5,5))
sns.heatmap(data=cm,linewidths=.5, annot=True,  cmap = 'Blues')
plt.ylabel('Actual label')
plt.xlabel('Predicted label')
all_sample_title = 'Accuracy Score for Random Forest: {0}'.format(rfc.score(X_test, y_
plt.title(all_sample_title, size = 15)
```

Out[37]: Text(0.5, 1.0, 'Accuracy Score for Random Forest: 0.9041095890410958')

```
In [38]: from sklearn.metrics import roc_curve, roc_auc_score
         y_pred_proba = rfc.predict_proba(X_test)[:][:,1]

         df_actual_predicted = pd.concat([pd.DataFrame(np.array(y_test), columns=['y_actual'])
         df_actual_predicted.index = y_test.index

         fpr, tpr, tr = roc_curve(df_actual_predicted['y_actual'], df_actual_predicted['y_pred
         auc = roc_auc_score(df_actual_predicted['y_actual'], df_actual_predicted['y_pred_proba

         plt.plot(fpr, tpr, label='AUC = %0.4f' %auc)
         plt.plot(fpr, fpr, linestyle = '--', color='k')
         plt.xlabel('False Positive Rate')
         plt.ylabel('True Positive Rate')
         plt.title('ROC Curve', size = 15)
         plt.legend()
```

Out[38]: <matplotlib.legend.Legend at 0x2578e5cfd60>