

[https://github.com/texashikkita/7330TW/blob/main/FINAL\\_JMcPhaul\\_7330 CaseStudy2\\_diabetes.ipynb](https://github.com/texashikkita/7330TW/blob/main/FINAL_JMcPhaul_7330 CaseStudy2_diabetes.ipynb)

```
1 # prompt: mount drive
2
3 from google.colab import drive
4 drive.mount('/content/drive')
```

Mounted at /content/drive

```
1 import psutil
2 print("Available Memory: (psutil.virtual_memory().available / 1024.0) GB")
```

Available Memory: 87.18 GB

```
1 import torch
2 import copy as cp
3
4 # Check PyTorch CUDA availability
5 print(torch.cuda.is_available())
6 if torch.cuda.is_available():
7     print("PyTorch Device: (torch.cuda.get_device_name(0))")
8     print("PyTorch Version (PyTorch): (torch.version.cuda)")
9
10 # Check CuPy CUDA availability
11 print(cuda.is_available())
12 if cp.cuda.is_available():
13     print("CuPy Device: (cp.cuda.get_device_name(0))")
14     print("CuPy Version (CuPy): (cp.cuda.runtime.runtimeGetVersion() / 1000)")
```

```
15 import torch
16
17 if torch.cuda.is_available():
18     print("PyTorch CUDA available: True")
19     PyTorch_Device: NVIDIA A100-SXM4-40GB
20     CUDA Version (PyTorch): 12.4
21     CuPy CUDA available: True
22     CUDA Version (CuPy): 12.06
23     CUDA Device: NVIDIA A100-SXM4-40GB
24     CuDF is successfully installed!
25         a b
26 0 1 4
27 1 2 5
28 2 3 6
```

```
1 # 2. EDA
2
3 import cuff
4
5 # load the data into cuDF DataFrames
6 diabetic_data = cuff.read_csv("/content/drive/MyDrive/diabetic_data.csv")
7 ids_mapping = cuff.read_csv("/content/drive/MyDrive/ids_mapping.csv")
8
9 # Ensure all string columns are treated as string type
10 diabetic_data = diabetic_data.astype(str)
11
12 # Replace '?' with None before converting to cuDF's NA
13 diabetic_data = diabetic_data.replace('?', None).fillna(cuff.NA)
14
15
16 # Or if needed fix = convert only object columns
17 # For col in diabetic_data.select_dtypes(include='object'):
18 #     diabetic_data[col] = diabetic_data[col].replace('?', None).fillna(cuff.NA)
```

```
19 # Display dataset info
20 print("Diabetic Data Info:")
21 print(diabetic_data.info())
22
23 print("\n first few rows of diabetic_data:")
24 print(diabetic_data.head())
25
26 print("\n Ids Mapping Data Info:")
27 print(ids_mapping.info())
28
29 print("\n first few rows of Ids_mapping:")
30
```

[https://colab.research.google.com/gist/texashikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul\\_diabetes.ipynb?authuser=1#scrollTo=EYE...](https://colab.research.google.com/gist/texashikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul_diabetes.ipynb?authuser=1#scrollTo=EYE...) 1/104

McPhaul\_J\_CaseStudy2\_FINAL\_diabetes.ipynb - Colab

```
1 categorical_cols = [
2     "race", "gender", "age", "payer_code", "medical_specialty",
3     "diag_1", "diag_2", "diag_3", "new_glu_serum", "AlcResult",
4     "metformin", "repaglinide", "nateglinide", "chlorpropamide",
5     "glipizide", "acetohexamide", "glipizide", "glipizide",
6     "glyburide", "tolbutamide", "pioglitazone", "rosiglitazone",
7     "acarbose", "miglitol", "troglitazone", "tolazamide",
8     "examide", "citoglipton", "insulin", "glyburide-metformin",
9     "metformin-rosiglitazone", "metformin-pioglitazone",
10    "change", "readmitted"
11 ]
12
13 for col in categorical_cols:
14     diabetic_data[col] = diabetic_data[col].astype("str").replace('?', cuff.NA)
```

15 # Verify Fix

16 print("Data Types Fixed and Missing Values Handled!")

17 print(diabetic\_data.dtypes)

18

19 **Data Types Fixed and Missing Values Handled!**

```
encounter_id      int64
patient_nbr      int64
race              object
gender             object
age               object
weight            object
admission_type_id int64
discharge_disposition_id int64
admission_source_id int64
time_in_hospital int64
payer_code          object
medical_specialty object
num_lab_procedures int64
num_procedures      int64
num_medications      int64
num_outpatient      int64
num_emergency       int64
num_inpatient       int64
diag_1              object
diag_2              object
diag_3              object
number_diagnoses   int64
max_glu_serum       object
AlcResult           object
metformin           object
repaglinide          object
nateglinide          object
chlorpropamide       object
glipizide            object
acetohexamide        object
glipizide            object
glyburide            object
tolbutamide          object
pioglitazone          object
rosiglitazone         object
acarbose             object
miglitol              object
troglitazone          object
tolazamide             object
examide              object
citoglipton           object
insulin                object
glyburide-metformin   object
glipizide-metformin   object
glyburide-pioglitazone  object
metformin-rosiglitazone  object
metformin-pioglitazone  object
change                object
diabetesMed           object
readmitted            object
dtype: object
```

```
1 # Merge Ids
2
3 import cuff
4
5 # Check for non-numeric values
6 invalid_values = ~ids_mapping['ids_mapping["admission_type_id"]'].str.isnumeric()
7 print(~invalid_values.sum())
8 print(~invalid_values.sum().sum())
9
10 # Convert numeric values to integers
11 ids_mapping[~ids_mapping['ids_mapping["admission_type_id"]'].str.isnumeric()] = 0
12 ids_mapping['admission_type_id'] = ids_mapping['admission_type_id'].astype("int64")
```

13
14 print("\n Cleaned Ids\_mapping Data:")
15 print(ids\_mapping.head())
16

[https://github.com/texashikkita/7330TW/blob/main/FINAL\\_JMcPhaul\\_7330 CaseStudy2\\_diabetes.ipynb](https://github.com/texashikkita/7330TW/blob/main/FINAL_JMcPhaul_7330 CaseStudy2_diabetes.ipynb)

```
1 # prompt: mount drive
2
3 from google.colab import drive
4 drive.mount('/content/drive')
5
6 # Mounted at /content/drive
7
8 Available Memory: 87.18 GB
9
10 import torch
11 import copy as cp
12
13 # Check PyTorch CUDA availability
14 print(torch.cuda.is_available())
15 if torch.cuda.is_available():
16     print("PyTorch Device: (torch.cuda.get_device_name(0))")
17     print("PyTorch Version (PyTorch): (torch.version.cuda)")
18
19 # Check CuPy CUDA availability
20 print(cuda.is_available())
21 if cp.cuda.is_available():
22     print("CuPy Device: (cp.cuda.get_device_name(0))")
23     print("CuPy Version (CuPy): (cp.cuda.runtime.runtimeGetVersion() / 1000)")
24
25 import torch
26
27 if torch.cuda.is_available():
28     print("PyTorch CUDA available: True")
29     PyTorch_Device: NVIDIA A100-SXM4-40GB
30     CUDA Version (PyTorch): 12.4
31     CuPy CUDA available: True
32     CUDA Version (CuPy): 12.06
33     CUDA Device: NVIDIA A100-SXM4-40GB
34     CuDF is successfully installed!
35         a b
36 0 1 4
37 1 2 5
38 2 3 6
39
40 import cuff
41
42 # load the data into cuDF DataFrames
43 diabetic_data = cuff.read_csv("/content/drive/MyDrive/diabetic_data.csv")
44 ids_mapping = cuff.read_csv("/content/drive/MyDrive/ids_mapping.csv")
45
46 # Ensure all string columns are treated as string type
47 diabetic_data = diabetic_data.astype(str)
48
49 # Replace '?' with None before converting to cuDF's NA
50 diabetic_data = diabetic_data.replace('?', None).fillna(cuff.NA)
51
52
53 # Or if needed fix = convert only object columns
54 # For col in diabetic_data.select_dtypes(include='object'):
55 #     diabetic_data[col] = diabetic_data[col].replace('?', None).fillna(cuff.NA)
56
57 # Display dataset info
58 print("Diabetic Data Info:")
59 print(diabetic_data.info())
60
61 print("\n first few rows of diabetic_data:")
62 print(diabetic_data.head())
63
64 print("\n Ids Mapping Data Info:")
65 print(ids_mapping.info())
66
67 print("\n first few rows of Ids_mapping:")
68
```

[5 rows x 50 columns]

```
1 IDs Mapping Data Info:
2 <class 'cudf.core.DataFrame'>
3 RangeIndex: 67 entries, 0 to 66
4 Data columns (total 2 columns):
5   # Column           Non-Null Count  Dtype  
6   --  -- 
7   0 admission_type_id 65 non-null   object 
8   1 description        62 non-null   object 
9   dtypes: object(2)
10  memory usage: 2.9+ KB
11  None
```

12
13 # Step 2.1 Fix Data Types and Handle Missing Values
14
15 1 # Import cuff
16
17 2 # Convert Numeric Columns First
18 3 numeric\_cols = [
19 4 "age", "admission\_type\_id", "discharge\_disposition\_id",
20 5 "admission\_source\_id", "time\_in\_hospital", "num\_lab\_procedures", "num\_procedures",
21 6 "num\_medications", "number\_outpatient", "number\_emergency", "number\_inpatient",
22 7 "number\_diagnoses"
23 8 ]
24
25 9 for col in numeric\_cols:
26 10 diabetic\_data[col] = diabetic\_data[col].replace('?', None).fillna(cuff.NA)
27
28 11
29 12 # Display dataset info
30 13 print("Diabetic Data Info:")
31 14 print(diabetic\_data.info())
32 15
33 16 print("\n first few rows of diabetic\_data:")
34 17 print(diabetic\_data.head())
35 18
36 19 print("\n Ids Mapping Data Info:")
37 20 print(ids\_mapping.info())
38 21
39 22 print("\n first few rows of Ids\_mapping:")
40 23

```
1 https://colab.research.google.com/gist/texashikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul_diabetes.ipynb?authuser=1#scrollTo=EYE... 2/104
2
3 McPhaul_J_CaseStudy2_FINAL_diabetes.ipynb - Colab
4
5 # Non-Numeric Values in 'admission_type_id':
6     admission_type_id description
7     0 discharge_disposition_id description
8     1 admission_source_id description
9
10 # Cleaned 'ids_mapping' Data:
11     admission_type_id description
12     0 Emergency
13     1 Urgent
14     2 Elective
15     3 Newborn
16     4 Not Available
17
18 1 # 3.2
19
20 2 # Merge Diabetic Data with 'ids_mapping' on 'admission_type_id'
21 3 diabetic_data = diabetic_data.merge(ids_mapping, how='left', on='admission_type_id')
22
23 4 # Drop unnecessary columns
24 5 columns_to_drop = [
25 6     "weight", "new_glu_serum", "AlcResult", "medical_specialty", "payer_code",
26 7     "encounter_id", "patient_nbr", "description" # description is from ids_mapping
27 8 ]
28
29 9 diabetic_data = diabetic_data.drop(columns=columns_to_drop)
30
31 10
32 11 # Fill Missing Values in Key Categorical Columns
33 12 for col in ["race", "diag_1", "diag_2", "diag_3"]:
34 13     diabetic_data[col] = diabetic_data[col].fillna("Unknown")
35
36 14 # Convert 'readmitted' to numerical categories
37 15 diabetic_data['readmitted'] = diabetic_data['readmitted'].map(["NO": 0, "NP": 1, "NPV": 2])
38
39 16 # Verify Merge & Cleaned
40 17 print(diabetic_data.dtypes)
41 18 print("\n First Few Rows of Cleaned Data:")
42 19 print(diabetic_data.head())
43
44 20
45 21 number_diagnoses      int64
46 metformin                 object
47 repaglinide                object
48 nateglinide                object
49 chlorpropamide               object
50 glipizide                  object
51 acetohexamide                object
52 glipizide                  object
53 glyburide                  object
54 tolbutamide                 object
55 pioglitazone                object
56 rosiglitazone                object
57 acarbose                   object
58 miglitol                    object
59 troglitazone                 object
60 tolazamide                  object
61 examide                     object
62 citoglipton                 object
63 insulin                     object
64 glyburide-metformin          object
65 glipizide-metformin          object
66 glipizide-pioglitazone        object
67 metformin-rosiglitazone      object
68 metformin-pioglitazone        object
69 change                      object
70 diabetesMed                 object
71 readmitted                  int64
72
73 dtype: object
```

```
First Few Rows of Cleaned Data:
1 race gender age admission_type_id discharge_disposition_id \
2 Caucasian Female [50-60] 6 25
3
4
```

[https://colab.research.google.com/gist/texashikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul\\_diabetes.ipynb?authuser=1#scrollTo=EYE...](https://colab.research.google.com/gist/texashikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul_diabetes.ipynb?authuser=1#scrollTo=EYE...) 3/104

[https://colab.research.google.com/gist/texashikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul\\_diabetes.ipynb?authuser=1#scrollTo=EYE...](https://colab.research.google.com/gist/texashikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul_diabetes.ipynb?authuser=1#scrollTo=EYE...) 4/104

```
glipizide-metformin glimepiride-pioglitazone metformin-rosiglitazone \
0   NO      No       No
1   NO      No       No
2   NO      No       No
3   NO      No       No
4   NO      NO      No

metformin_micelle+aztreonam chlorthiazideMetad readmit++ad
```

```
1 import pandas as pd
2 from sklearn.preprocessing import StandardScaler
3
4 # load the dataset
5 diabetic_data = pd.read_csv("/content/drive/MyDrive/WI_Case_Study_2/Diabetic_data.csv")
6 id_mapping = pd.read_csv("/content/drive/MyDrive/WI_Case_Study_2/ID(mapping.csv)")
7
8 # 3.2 Data Cleaning and Merging
9
10 # Convert 'admission_type_id' to numeric, handling non-numeric values
11 diabetic_data['admission_type_id'] = pd.to_numeric(diabetic_data['admission_type_id'], errors='coerce')
12 id_mapping['admission_type_id'] = pd.to_numeric(id_mapping['admission_type_id'], errors='coerce')
13
14 # Convert to Int64 after ensuring both are numeric
15 diabetic_data['admission_type_id'] = diabetic_data['admission_type_id'].astype('Int64')
16 id_mapping['admission_type_id'] = id_mapping['admission_type_id'].astype('Int64')
17
18
19 # Merge diabetic_data with id_mapping (now with consistent data types)
20 diabetic_data = diabetic_data.merge(id_mapping, how='left', on='admission_type_id')
21
22
23 # Fill missing values in key categorical columns
24 for col in ['race', 'diag_1', 'diag_2', 'diag_3']:
25     diabetic_data[col] = diabetic_data[col].fillna("Unknown")
26
27 # Convert 'readmitted' to numerical categories
28 diabetic_data['readmitted'] = diabetic_data['readmitted'].replace({'NO': 0, 'NP': 1, 'NO+': 2})
29
30 # Convert 'max_glu_serum' and 'ACResult' to numerical representations
31 diabetic_data['max_glu_serum'] = diabetic_data['max_glu_serum'].replace({
32     'None': 0,
33     'Normal': 1,
34     'B�ormal': 2,
35     '>200': 3
36 })
37
38
39 diabetic_data['ACResult'] = diabetic_data['ACResult'].replace({
40     'None': 0,
41     'Normal': 1,
42     'B�ormal': 2,
43     '>200': 3
44 })
45
46
47 # 4. Feature Engineering (Scaling Numeric Features)
48
49 # Define Numeric Columns
50 numeric_cols = [
51     "time_in_hospital", "num_lab_procedures", "num_procedures",
52     "num_medications", "number_outpatient", "number_emergency",
53     "number_inpatient", "number_diabetes"
54 ]
55
56 # Initialize StandardScaler
57 scaler = StandardScaler()
58
59 # Fit and transform the selected numeric columns
60 diabetic_data[numeric_cols] = scaler.fit_transform(diabetic_data[numeric_cols])
61
62
63 # Verify Merge, Cleaning, and Scaling
64 print("Merge Complete and Data Cleaned!")
65 print(diabetic_data.dtypes)
66 print("First Few Rows of Cleaned Data:")
67 print(diabetic_data.head())
68 print(diabetic_data.describe())
69
```

```
dtype: object

First Few Rows of Cleaned Data:
   encounter_id patient_nbr race gender age weight \
0    2278392    8222157 Caucasian Female [0-10) ?
1    2278392    8222157 Caucasian Female [0-10) ?
2    2278392    8222157 Caucasian Female [0-10) ?
3    149190    55629189 Caucasian Female [10-20) ?
4    149190    55629189 Caucasian Female [10-20) ?

   admission_type_id discharge_disposition_id admission_source_id \
0                  6                      25                   1
1                  6                      25                   1
2                  6                      25                   1
3                  1                      1                   7
4                  1                      1                   7

   time_in_hospital ... insulin_glyburide-metformin glipizide-metformin \
0      -1.137649 ...          No           No           No
1      -1.137649 ...          No           No           No
2      -1.137649 ...          No           No           No
3     -0.467653 ...          Up            No           No
4     -0.467653 ...          Up            No           No

   glimepiride-pioglitazone metformin-rosiglitazone metformin-pioglitazone \
0                  No           No           No           No
1                  No           No           No           No
2                  No           No           No           No
3                  No           No           No           No
4                  No           No           No           No

   change diabetesMed readmitted \
0      No      No      0
1      No      No      0
2      No      No      0
3      Ch     Yes     1
4      Ch     Yes     1

1
# Save the cleaned data to 'data_cleaned.csv'
diabetic_data.to_csv("data_cleaned.csv", index=False)
```

```
1
2 # fix data types and handle missing values
3
4 from google.colab import drive
5 import psutil
6 import torch
7 import numpy as np
8 import pandas as pd
9 from sklearn.preprocessing import StandardScaler
10
11 drive.mount('/content/drive')
12
13 print("Available Memory: (nnutil.virtual_memory().available / 1e9) GB")
14
15 class PythonCodeAvailability:
16     @staticmethod
17     def printPyTorchCodeAvailability():
18         if torch.cuda.is_available():
19             print("PyTorch CUDA available")
20         else:
21             print("PyTorch CUDA NOT available")
22
23     @staticmethod
24     def printTorchCodeAvailability():
25         if torch.torch.is_available():
26             print("Torch CUDA available")
27         else:
28             print("Torch CUDA NOT available")
29
30     @staticmethod
31     def printCUDAVersion():
32         if torch.cuda.is_available():
33             print("CUDA Version (CUDF): (cp.cuda.runtime.getRunTimeVersion() / 1000)")
34         else:
35             print("CUDA is successfully installed!"")#This line seems unnecessary, remove it if you don't need to confirm installation
36
37
38 # Load the datasets using pandas
39 diabetic_data = pd.read_csv("/content/drive/MyDrive/WI_Case_Study_2/Diabetic_data.csv")
40 id_mapping = pd.read_csv("/content/drive/MyDrive/WI_Case_Study_2/ID(mapping.csv)")
41
42 # 3.2 Data Cleaning and Merging
43 # Convert 'admission_type_id' to numeric, handling non-numeric values
44 diabetic_data['admission_type_id'] = pd.to_numeric(diabetic_data['admission_type_id'], errors='coerce')
45 id_mapping['admission_type_id'] = pd.to_numeric(id_mapping['admission_type_id'], errors='coerce')
46
47 # Convert to Int64 after ensuring both are numeric
48 diabetic_data['admission_type_id'] = diabetic_data['admission_type_id'].astype('Int64')
49 id_mapping['admission_type_id'] = id_mapping['admission_type_id'].astype('Int64')
50
51
52 # Merge diabetic_data with id_mapping (now with consistent data types)
53 diabetic_data = diabetic_data.merge(id_mapping, how='left', on='admission_type_id')
```

[https://colab.research.google.com/gist/texashchikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul\\_diabetes.ipynb?authuser=1#scrollTo=EYE...](https://colab.research.google.com/gist/texashchikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul_diabetes.ipynb?authuser=1#scrollTo=EYE...) 5/104

```
2/4/25, 1:56 AM
McPhaul_J_CaseStudy2_FINAL_diabetes.ipynb - Colab

time_in_hospital ... insulin_glyburide-metformin glipizide-metformin \
0      -1.137649 ...          No           No           No
1      -1.137649 ...          No           No           No
2      -1.137649 ...          No           No           No
3     -0.467653 ...          Up            No           No
4     -0.467653 ...          Up            No           No

glimepiride-pioglitazone metformin-rosiglitazone metformin-pioglitazone \
0                  No           No           No           No
1                  No           No           No           No
2                  No           No           No           No
3                  No           No           No           No
4                  No           No           No           No

change diabetesMed readmitted \
0      No      No      0
1      No      No      0
2      No      No      0
3      Ch     Yes     1
4      Ch     Yes     1

1
# Drop unnecessary columns
2 columns_to_drop = [
3     "weight", "max_glu_serum", "ACResult", "medical_specialty", "payer_code",
4     "encounter_id", "patient_nbr", "description" # 'description' is from id_mapping
5 ]
6
7 diabetic_data = diabetic_data.drop(columns=columns_to_drop, errors='ignore') # Use errors='ignore'

8
9 # Check for non-numeric values and handle them
10
11 # Check if 'admission_type_id' is numeric using pd.to_numeric
12 invalid_values = id_mapping[id_mapping['admission_type_id'].str.isnumeric()].isnull()
13
14 if not invalid_values.empty:
15     print("Non-Numeric Values in 'admission_type_id':", invalid_values)
16     # Decide how to handle invalid values: remove them, convert to numeric, or fill with a specific value
17     str_is_not_needed = id_mapping[id_mapping['admission_type_id'].str.isnumeric() & invalid_values]
18     if str_is_not_needed.empty:
19         id_mapping[id_mapping['admission_type_id'].str.isnumeric() & invalid_values] = 0 # Example: replace with 0 if str is not needed here
20     else:
21         id_mapping[id_mapping['admission_type_id'].str.isnumeric() & invalid_values] = 0 # Example: replace with 0
22
23 # Option 2: Convert non-numeric values to a default numeric value
24 id_mapping[id_mapping['admission_type_id'].str.isnumeric() == 0].str.replace(0, 1, inplace=True)
25
26 # Option 3: Convert non-numeric values to integers
27
28 # Print cleaned id_mapping data
29 print("Cleaned id_mapping Data:")
30 print(id_mapping)
31 print("Cleaned id_mapping Data")
32 print(id_mapping.describe())
33
34 # Check for non-numeric values and handle them
35 invalid_values = id_mapping[id_mapping['admission_type_id'].str.isnumeric()].isnull()
36
37 if not invalid_values.empty:
38     print("Non-Numeric Values in 'admission_type_id':", invalid_values)
39     # Remove rows with non-numeric values
40     id_mapping = id_mapping[id_mapping['admission_type_id'].str.isnumeric()].notnull()
41
42 # Convert 'admission_type_id' to numeric in both DataFrames
43 id_mapping['admission_type_id'] = pd.to_numeric(id_mapping['admission_type_id'], errors='coerce').astype('Int64')
44
45 # Merge the DataFrames
46 diabetic_data = diabetic_data.merge(id_mapping, how='left', on='admission_type_id')
47
48
```

```
49
50 # cipython-input-13-ab8bed82086:10: SettingWithCopyWarning:
51 # A value is trying to be set on a copy of a slice from a DataFrame.
52 # Try using .loc[row_indexer,col_indexer] = value instead
53
54 See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user\_guide/indexing.html#returning-a-view-versus-a-copy
55
56
57
```

```
58 # Fill missing values in key categorical columns
59 for col in ['race', 'diag_1', 'diag_2', 'diag_3']:
60     diabetic_data[col] = diabetic_data[col].fillna("Unknown")
61
62 # Convert 'readmitted' to numerical categories
63 diabetic_data['readmitted'] = diabetic_data['readmitted'].replace({'NO': 0, 'NP': 1, 'NO+': 2})
64
65 # Convert 'max_glu_serum' and 'ACResult' to numerical representations
66 diabetic_data['max_glu_serum'] = diabetic_data['max_glu_serum'].replace({
67     'None': 0,
68     'Normal': 1,
69     'B�ormal': 2,
70     '>200': 3
71 })
72
73
74 # 4. Feature Engineering (Scaling Numeric Features)
75
76 # Define Numeric Columns
77 numeric_cols = [
78     "time_in_hospital", "num_lab_procedures", "num_procedures",
79     "num_medications", "number_outpatient", "number_emergency",
80     "number_inpatient", "number_diabetes"
81 ]
82
83 # Initialize StandardScaler
84 scaler = StandardScaler()
85
86 # Fit and transform the selected numeric columns
87 diabetic_data[numeric_cols] = scaler.fit_transform(diabetic_data[numeric_cols])
88
89
90 # Verify Merge, Cleaning, and Scaling
91 print("Merge Complete and Data Cleaned!")
92 print(diabetic_data.dtypes)
93 print("First Few Rows of Cleaned Data:")
94 print(diabetic_data.head())
95 print(diabetic_data.describe())
96
97 # Save the cleaned data to 'data_cleaned.csv'
98 diabetic_data.to_csv("data_cleaned.csv", index=False)
99
```

```
repaglinide          object
nateglinide          object
chlorpropamide       object
```

```

2 # Drop unnecessary columns
3 columns_to_drop = [
4     "weight", "max_glu_serum", "A1cResult", "medical_specialty", "payer_code",
5     "encounter_id", "patient_nbr", "description" # description is from id_mapping
6 ]
7
8 diabetic_data = diabetic_data.drop(columns=columns_to_drop, errors='ignore') # Use errors='ignore'
9
10 # Fill Missing Values in Key Categorical Columns
11 for col in ["race", "diag_1", "diag_2", "diag_3"]:
12     diabetic_data[col] = diabetic_data[col].fillna("Unknown")
13
14 # Convert 'readmitted' to numerical categories
15 diabetic_data['readmitted'] = diabetic_data['readmitted'].map(["NO": 0, "NP": 1, ">30": 2])
16
17 # Verify Merge & Cleanse
18 print(diabetic_data.info())
19 print(diabetic_data.dtypes)
20 print("First Few Rows of Cleaned Data:")
21 print(diabetic_data.head())
22
23 number_diagnoses    float64
metformin          object
repaglinide         object
nateglinide         object
chlorpropamide      object
glimepiride         object
acetohexamide      object
glibizide           object
glyburide           object
tolazamide           object
pioglitazone        object
rosiglitazone       object
acarbose            object
miglitol             object
troglitazone        object
tolazamide           object
examide              object
citoglipiton        object
insulin              object
glyburide-metformin object
glibizide-metformin object
glimepiride-pioglitazone object
metformin-rosiglitazone object
metformin-pioglitazone object
change               object
diabetestmed         object
readmitted           float64
dtype: object

```

```

First Few Rows of Cleaned Data:
   race gender age admission_type_id discharge_disposition_id \
0 Caucasian Female (0-10)          6                         25
1 Caucasian Female (0-10)          6                         25
2 Caucasian Female (0-10)          6                         25
3 Caucasian Female (0-10)          6                         25
4 Caucasian Female (0-10)          6                         25

admission_source_id time_in_hospital num_lab_procedures num_procedures \
0                   1. -1.137649          -0.106517          -0.785398
1                   1. -1.137649          -0.106517          -0.785398
2                   1. -1.137649          -0.106517          -0.785398
3                   1. -1.137649          -0.106517          -0.785398
4                   1. -1.137649          -0.106517          -0.785398

num_medications ... citoglipiton insulin glyburide-metformin \
0      -1.848268 ... No      No      No
1      -1.848268 ... No      No      No
2      -1.848268 ... No      No      No
3      -1.848268 ... No      No      No
4      -1.848268 ... No      No      No

glipizide-metformin glimepiride-pioglitazone metformin-rosiglitazone \
0      No      No      No
1      No      No      No
2      No      No      No
3      No      No      No
4      No      No      No

metformin-pioglitazone change diabetesMed readmitted

```

```

1 import pandas as pd
2
3 # df
4 # categorical_cols = ["race", "gender", "age", "change", "diabetestmed", "insulin"]
5
6 # One column get_dummies for one-hot encoding
7 diabetic_data = pd.get_dummies(diabetic_data, columns=categorical_cols, dummy_na=True)
8
9 print("Categorical Features One-Hot Encoded Successfully!")

```

[https://colab.research.google.com/gist/texashchikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul\\_diabetes.ipynb](https://colab.research.google.com/gist/texashchikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul_diabetes.ipynb)?authuser=1#scrollTo=EY...

```

All Features Converted to Numeric Format!
float32
[nan]
[nan]

Non-Numeric Columns in X: Index([], dtype='object')

1 From sklearn.model_selection import train_test_split
2
3 # Define Features (X) and Target (y)
4 # X = diabetic_data.drop(['readmitted'], axis=1)
5 # Y = diabetic_data['readmitted']
6
7 # Convert to int32 and handle non-finite values with fillna
8 y = diabetic_data['readmitted'].fillna(-1).astype('int32') # Replace NaN with -1 before conversion
9
10 # Split Data
11 X_train, X_test, y_train, y_test = train_test_split(
12     X, y, test_size=0.2, stratify=y, random_state=42
13 )
14 print("Train/Test Split Completed! Shapes:")
15 print(f"X_train: {X_train.shape}, y_train: {y_train.shape}")
16 print(f"X_test: {X_test.shape}, y_test: {y_test.shape}")

Train/Test Split Completed! Shapes:
- X_train: (732715, 69), y_train: (732715,)
- X_test: (183179, 69), y_test: (183179,)


```

```

1 import pandas as pd
2
3 from sklearn.model_selection import train_test_split
4 from sklearn.preprocessing import StandardScaler
5 from sklearn.linear_model import LogisticRegression
6 from sklearn.metrics import accuracy_score, classification_report, confusion_matrix
7 from sklearn import model_selection
8
9 # Load the cleaned data
10 diabetic_data = pd.read_csv("data_cleaned.csv")
11
12 # Define features (X) and target (y)
13 X = diabetic_data.drop(columns=['readmitted'])
14 y = diabetic_data['readmitted'].astype('int32')
15
16 # Handle potential non-numeric columns in X
17 non_numeric_cols = X.select_dtypes(exclude=[number]).columns
18 if len(non_numeric_cols) != 0:
19     print("Warning: Non-numeric column found in X: ", non_numeric_cols)
20     # Decide how to handle them (e.g., one-hot encoding, dropping)
21     # For simplicity, we will drop them
22
23 # Impute missing values using SimpleImputer
24 imputer = SimpleImputer(strategy='mean') # or 'median', 'most_frequent'
25 X = imputer.fit_transform(X) # Fit and transform to replace NaNs
26
27 # Split Data
28 X_train, X_test, y_train, y_test = train_test_split(
29     X, y, test_size=0.2, stratify=y, random_state=42
30 )
31
32 # Initialize and Train Model
33 logreg = LogisticRegression(max_iter=1000, tol=1e-4)
34 logreg.fit(X_train, y_train)
35 print("Logistic Regression Model Trained Successfully!")

# Predict on Test Data
36 y_pred = logreg.predict(X_test)
37
38 # Compute Evaluation Metrics
39 accuracy = accuracy_score(y_test, y_pred)
40 conf_matrix = confusion_matrix(y_test, y_pred)
41 class_report = classification_report(y_test, y_pred)
42
43 # Display Results
44 print(f"Accuracy: {accuracy:.4f}")
45 print(f"Confusion Matrix: \n{conf_matrix}")
46 print(f"Classification Report: \n{class_report}")
47
48 # Check Class Balance
49 print("Class Distribution in Training Data:")
50 print(classification_report(y_train))
51 print("Class Distribution in Testing Data:")
52 print(classification_report(y_test))

Warning: Non-numeric columns found in X: Index(['race', 'gender', 'age', 'weight', 'payer_code', 'medical_specialty',
       'diag_1', 'diag_3', 'metformin', 'repaglinide', 'nateglinide',
       'chlorpropamide', 'acetohexamide', 'glibizide',
       'glyburide', 'tolazamide', 'pioglitazone', 'examicide', 'acarbose',
       'miglitol', 'troglitazone', 'tolazamide', 'examicide', 'citoglipiton',
       'insulin', 'glyburide-metformin', 'glipizide-metformin',
       'glimepiride-pioglitazone', 'metformin-rosiglitazone',
       'metformin-pioglitazone', 'change', 'diabetestmed', 'description'],
      dtype='object')

```

Logical Regression Model Trained Successfully!

Accuracy: 0.5422

Confusion Matrix:

[[30102 2817 0]

[18328 3087 0]]

[https://colab.research.google.com/gist/texashchikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul\\_diabetes.ipynb](https://colab.research.google.com/gist/texashchikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul_diabetes.ipynb)?authuser=1#scrollTo=EY...

```

1 print(diabetic_data.head())
2
3 # Categorical Features One-Hot Encoded Successfully!
4 admission_type_id discharge_disposition_id admission_source_id \
0 6 25 1
1 6 25 1
2 6 25 1
3 6 25 1
4 6 25 1

time_in_hospital num_lab_procedures num_procedures num_medications \
0 -1.137649 -0.106517 -0.785398 -1.848268
1 -1.137649 -0.106517 -0.785398 -1.848268
2 -1.137649 -0.106517 -0.785398 -1.848268
3 -1.137649 -0.106517 -0.785398 -1.848268
4 -1.137649 -0.106517 -0.785398 -1.848268

number_outpatient number_emergency number_inpatient ... change_No \
0 -0.291461 -0.21262 -0.503276 ... True
1 -0.291461 -0.21262 -0.503276 ... True
2 -0.291461 -0.21262 -0.503276 ... True
3 -0.291461 -0.21262 -0.503276 ... True
4 -0.291461 -0.21262 -0.503276 ... True

change_nan diabetesMed_N diabetesMed_Yes diabetesMed_no insulin_Down \
0 True False False False False
1 False True False False False
2 False True False False False
3 False True False False False
4 False True False False False

insulin_No insulin_Steady insulin_Up insulin_Dan \
0 True False False False
1 True False False False
2 True False False False
3 True False False False
4 True False False False

[5 rows x 70 columns]

```

```

1 # print(diabetic_data.head())
2 # Import pandas as pd
3
4 # Define categorical columns (this line might be redundant if already defined)?
5 categorical_cols = ["race", "gender", "age", "change", "diabetestmed", "insulin"]
6
7 # Check if columns exist before applying get_dummies
8 if all(col in diabetic_data.columns for col in categorical_cols):
9     diabetic_data[categorical_cols] = pd.get_dummies(diabetic_data, columns=categorical_cols, dummy_na=True)
10 else:
11     print("Categorical columns have already been encoded or do not exist in the DataFrame.")
12
13 # Convert 'diag_1', 'diag_2', 'diag_3' to categorical codes
14 for col in ["diag_1", "diag_2", "diag_3"]:
15     diabetic_data[col] = pd.get_dummies(diabetic_data[col], dropfirst=True)
16
17 # Convert all medication columns to binary (0/1)
18 medication_cols = [
19     "metformin", "repaglinide", "nateglinide", "chlorpropamide", "glimepiride",
20     "glipizide", "glyburide", "tolazamide", "pioglitazone",
21     "acarbose", "acetohexamide", "citoglipiton", "troglitazone",
22     "miglitol", "glyburide-metformin", "glipizide-metformin",
23     "glimepiride-pioglitazone", "metformin-rosiglitazone"
24 ]
25
26 for col in medication_cols:
27     # Convert to medication code
28     # Convert only if the column is of string type
29     if diabetic_data[col].dtype == "object":
30         # Create a new column with the medication code
31         diabetic_data[col] = pd.get_dummies(diabetic_data[col], dropfirst=True)
32
33 # Drop the 'description' column if it exists
34 if "description" in diabetic_data.columns:
35     diabetic_data.drop(columns="description", inplace=True)
36
37 # Convert everything to float32
38 diabetic_data = diabetic_data.astype("float32")
39 print("All Features Converted to Numeric Format!")
40
41 print(diabetic_data.info())
42 print(diabetic_data.describe())
43
44 # One column get_dummies for one-hot encoding
45 diabetic_data = pd.get_dummies(diabetic_data, columns=categorical_cols, dummy_na=True)
46
47 print("Categorical Features One-Hot Encoded Successfully!")

```

Categorical columns have already been encoded or do not exist in the DataFrame.

```

1 # Convert 'diag_1', 'diag_2', 'diag_3' to categorical codes
2 for col in ["diag_1", "diag_2", "diag_3"]:
3     diabetic_data[col] = pd.get_dummies(diabetic_data[col], dropfirst=True)
4
5 # Convert all medication columns to binary (0/1)
6 medication_cols = [
7     "metformin", "repaglinide", "nateglinide", "chlorpropamide", "glimepiride",
8     "glipizide", "glyburide", "tolazamide", "pioglitazone",
9     "acarbose", "acetohexamide", "citoglipiton", "troglitazone",
10    "miglitol", "glyburide-metformin", "glipizide-metformin",
11    "glimepiride-pioglitazone", "metformin-rosiglitazone"
12 ]
13
14 for col in medication_cols:
15     # Convert to medication code
16     # Convert only if the column is of string type
17     if diabetic_data[col].dtype == "object":
18         # Create a new column with the medication code
19         diabetic_data[col] = pd.get_dummies(diabetic_data[col], dropfirst=True)
20
21 # Drop the 'description' column if it exists
22 if "description" in diabetic_data.columns:
23     diabetic_data.drop(columns="description", inplace=True)
24
25 # Convert everything to float32
26 diabetic_data = diabetic_data.astype("float32")
27 print("All Features Converted to Numeric Format!")
28
29 print(diabetic_data.info())
30 print(diabetic_data.describe())
31
32 # One column get_dummies for one-hot encoding
33 diabetic_data = pd.get_dummies(diabetic_data, columns=categorical_cols, dummy_na=True)
34
35 print("Categorical Features One-Hot Encoded Successfully!")

```

[https://colab.research.google.com/gist/texashchikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul\\_diabetes.ipynb](https://colab.research.google.com/gist/texashchikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul_diabetes.ipynb)?authuser=1#scrollTo=EY...

```

[ 5974 848 0]
Classification Report:
precision recall f1-score support
0 0.95 0.91 0.93 32919
1 0.45 0.14 0.21 21327
2 0.00 0.00 0.00 6814

accuracy 0.54 61968
macro avg 0.33 0.35 0.30 61968
weighted avg 0.46 0.54 0.45 61968

Class Distribution in Training Data:
readmitted
0 13673
1 8898
2 2727
Name: count, dtype: int64
Class Distribution in Testing Data:
readmitted
0 32919
1 21327
2 6814
Name: count, dtype: int64
/usr/local/lib/python3.11/dist-packages/sklearn/metrics/_classification.py:1565: UndefinedMetricWarning: Precision is ill-defined and be
..._warn_prcf(average, modifier, f"metric.{capitalize(name)}" is", len(result))
/usr/local/lib/python3.11/dist-packages/sklearn/metrics/_classification.py:1565: UndefinedMetricWarning: Precision is ill-defined and be
..._warn_prcf(average, modifier, f"metric.{capitalize(name)}" is", len(result))
/usr/local/lib/python3.11/dist-packages/sklearn/metrics/_classification.py:1565: UndefinedMetricWarning: Precision is ill-defined and be
..._warn_prcf(average, modifier, f"metric.{capitalize(name)}" is", len(result))

```

[https://colab.research.google.com/gist/texashchikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul\\_diabetes.ipynb](https://colab.research.google.com/gist/texashchikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul_diabetes.ipynb)?authuser=1#scrollTo=EY...

[https://colab.research.google.com/gist/texashchikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul\\_diabetes.ipynb](https://colab.research.google.com/gist/texashchikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul_diabetes.ipynb)?authuser=1#scrollTo=EY...

[https://colab.research.google.com/gist/texashchikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul\\_diabetes.ipynb](https://colab.research.google.com/gist/texashchikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul_diabetes.ipynb)?authuser=1#scrollTo=EY...

2/4/25, 1:56 AM

McPhaul\_J\_CaseStudy2\_FINAL\_diabetes.ipynb - Colab

```

1 # Initialize results and provide analysis
2
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5
6 # ... (your existing code) ...
7
8 # Display Results
9 print("Accuracy: (%.4f)" % accuracy)
10 print("Precision Matrix: \n", conf_matrix)
11 print("Classification Report: \n", class_report)
12
13 # Visualize the Confusion Matrix
14 plt.figure(figsize=(8, 4))
15 sns.heatmap(conf_matrix, annot=True, fmt="d", cmap="Blues",
16             xticklabels=["No Readmission", "Readmitted >30", "Readmitted <30"],
17             yticklabels=["No Readmission", "Readmitted >30", "Readmitted <30"])
18 plt.xlabel("Predicted")
19 plt.ylabel("Actual")
20 plt.title("Confusion Matrix")
21 plt.show()
22
23 # Analyze Class Distribution
24 plt.figure(figsize=(4, 3))
25 sns.countplot(x=y_test, hue=y_pred)
26 plt.title("Class Distribution")
27 plt.xlabel("Readmission Category")
28 plt.ylabel("Number of Patients")
29 plt.show()
30
31 # prompt: visualize results and provide analysis
32
33 # Analyze Feature Importances (if available in your model)
34 # Get Feature names from original DataFrame before imputation
35 feature_names = diabetic_data.drop(columns=['readmitted']).columns
36
37 # prompt: visualize results and provide analysis
38
39 # Visualize results and provide analysis
40 import matplotlib.pyplot as plt
41 import seaborn as sns
42
43 # Display Results
44 print("Accuracy: (%.4f)" % accuracy)
45 print("Precision Matrix: \n", conf_matrix)
46 print("Classification Report: \n", class_report)
47
48 # Visualize the Confusion Matrix
49 plt.figure(figsize=(8, 4))
50 sns.heatmap(conf_matrix, annot=True, fmt="d", cmap="Blues",
51             xticklabels=["No Readmission", "Readmitted >30", "Readmitted <30"],
52             yticklabels=["No Readmission", "Readmitted >30", "Readmitted <30"])
53 plt.xlabel("Actual")
54 plt.ylabel("Predicted")
55 plt.title("Confusion Matrix")
56 plt.show()
57
58 # Analyze Class Distribution
59 plt.figure(figsize=(4, 3))
60 sns.countplot(x=y_test, hue=y_pred)
61 plt.title("Class Distribution")
62 plt.xlabel("Readmission Category")
63 plt.ylabel("Number of Patients")
64 plt.show()
65
66

```

2/4/25, 1:56 AM

McPhaul\_J\_CaseStudy2\_FINAL\_diabetes.ipynb - Colab

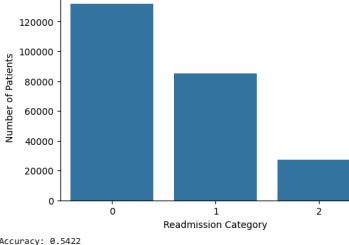
```

1 Accuracy: 0.5422
2
3 Confusion Matrix:
4 [[30102 2817 0]
5 [18320 3007 0]
6 [ 5974 840 0]]
7
8 Classification Report:
9      precision    recall   f1-score   support
10
11          0       0.55     0.91     0.69    32919
12          1       0.45     0.14     0.21   21327
13          2       0.00     0.00     0.00    6814
14
15   accuracy           0.54    61868
16   macro avg       0.33     0.35     0.30    61868
17   weighted avg    0.46     0.54     0.45    61868

```



Confusion Matrix



Class Distribution

Accuracy: 0.5422

Confusion Matrix:

[[30102 2817 0]

[18320 3007 0]

[ 5974 840 0]]

Classification Report:

precision recall f1-score support

0 0.55 0.91 0.69 32919

1 0.45 0.14 0.21 21327

2 0.00 0.00 0.00 6814

accuracy 0.54 61868

macro avg 0.33 0.35 0.30 61868

weighted avg 0.46 0.54 0.45 61868

[https://colab.research.google.com/gist/texashikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul\\_diabetes.ipynb?authuser=1#scrollTo=EY...](https://colab.research.google.com/gist/texashikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul_diabetes.ipynb?authuser=1#scrollTo=EY...) 13/104

14/104

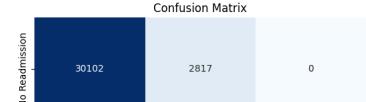
2/4/25, 1:56 AM

McPhaul\_J\_CaseStudy2\_FINAL\_diabetes.ipynb - Colab

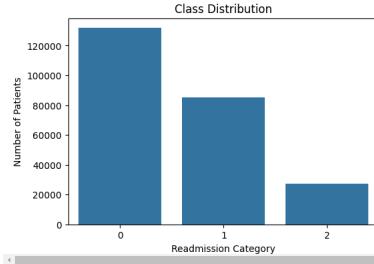
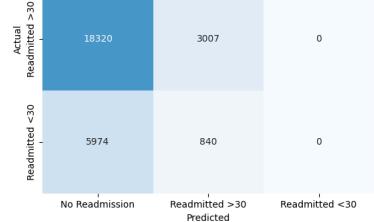
```

2 0.00 0.00 0.00 6814
3
4 accuracy 0.54 61868
5 macro avg 0.33 0.35 0.30 61868
6 weighted avg 0.46 0.54 0.45 61868

```



Confusion Matrix



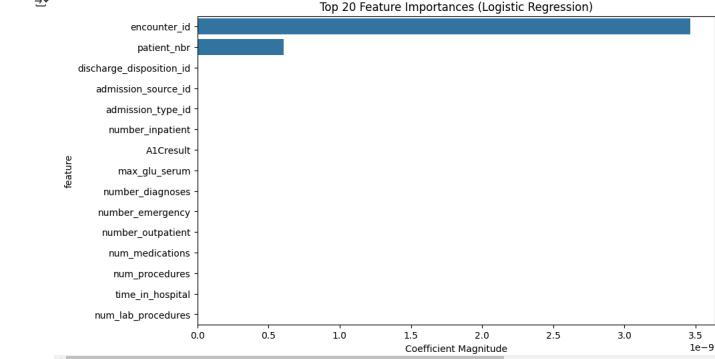
2/4/25, 1:56 AM

McPhaul\_J\_CaseStudy2\_FINAL\_diabetes.ipynb - Colab

```

1 # Get Feature names from original DataFrame before imputation, BUT AFTER SimpleImputer is applied
2 feature_names = diabetic_data.drop(columns=['readmitted']).select_dtypes(exclude=['number']).columns # Select only numeric features
3
4 # Create DataFrame with feature names and importances
5 feature_importances = pd.DataFrame((['feature': feature_name, 'importance': abs(log_reg.coef_[0])])
6 feature_importances.sort_values(by='importance', ascending=False)
7
8 plt.figure(figsize=(10, 6))
9 sns.barplot(x="importance", y="feature", data=feature_importances[::20]) # Show top 20 features
10 plt.title("Top 20 Feature Importances (Logistic Regression)")
11 plt.xlabel("Coefficient Magnitude")
12 plt.ylabel("Feature")
13
14
15
16

```



```

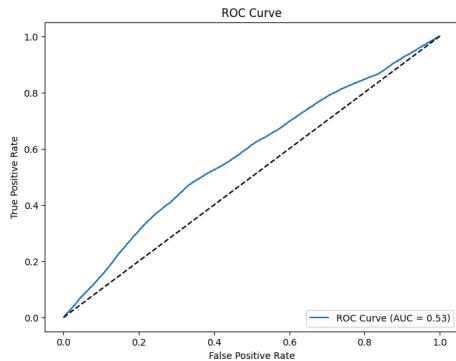
1 from sklearn.metrics import precision_score, recall_score, roc_auc_score, roc_curve
2 import matplotlib.pyplot as plt
3
4 # Predict probabilities for all classes (for AUC calculation)
5 y_pred_proba = log_reg.predict_proba(X_test) # Remove [:, 1]
6
7 # Calculate precision, recall, and AUC
8 precision = precision_score(y_test, y_pred, average='weighted') # Use 'weighted' for multi-class
9 recall = recall_score(y_test, y_pred, average='weighted') # Use 'weighted' for multi-class
10 f1 = f1_score(y_test, y_pred, average='weighted') # Use 'weighted' for multi-class
11
12 # Compute ROC curve and AUC
13 fpr, tpr, thresholds = roc_curve(y_test, y_pred_proba[:, 1], pos_label=1) # Choose relevant pos_label
14
15 # Plot ROC Curve
16 plt.plot(fpr, tpr, label="ROC Curve (AUC = %.2f)" % (auc(fpr, tpr)))
17 plt.plot([0, 1], [0, 1], 'k--', label="Diagonal Line")
18
19 # Add labels and title
20 plt.xlabel("False Positive Rate")
21 plt.ylabel("True Positive Rate")
22 plt.title("ROC Curve")
23 plt.legend(loc="lower right")
24
25

```

[https://colab.research.google.com/gist/texashikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul\\_diabetes.ipynb?authuser=1#scrollTo=EY...](https://colab.research.google.com/gist/texashikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul_diabetes.ipynb?authuser=1#scrollTo=EY...) 15/104[https://colab.research.google.com/gist/texashikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul\\_diabetes.ipynb?authuser=1#scrollTo=EY...](https://colab.research.google.com/gist/texashikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul_diabetes.ipynb?authuser=1#scrollTo=EY...) 16/104

2/4/25, 1:56 AM

McPhaul J CaseStudy2 FINAL diabetes.ipynb - Colab



```

1 import numpy as np
2
3 # Predict probabilities for all classes (for AUC calculation)
4 y_pred_proba = logreg.predict_proba(X_test)

5 # Calculate precision, recall, and AUC
6 precision = precision_score(y_test, y_pred, average='weighted')
7 recall = recall_score(y_test, y_pred, average='weighted')
8 result = [recall, precision, f1, y_test, y_pred_proba, multi_class='weighted']

9 # For AUC, use "ovr" for multiclass and provide probability estimates for all classes
10 auc = roc_auc_score(y_test, y_pred_proba, multi_class='ovr')

11 print("Precision: (%.4f)" % precision)
12 print("Recall: (%.4f)" % recall)
13 print("F1: (%.4f)" % f1)
14 print("AUC: (%.4f)" % auc)

15 # Plotting ROC curve
16
17 # For multi-class, you'll need to plot a ROC curve for each class vs. the rest
18 # of classes. Matplotlib's import roc_curve, auc
19 import matplotlib.pyplot as plt

20 n_classes = len(np.unique(y_test)) # Number of classes # Now np is defined
21
22 fpr = dict()
23 tpr = dict()
24 auc = dict()

25 for i in range(n_classes):
26     fpr[i], tpr[i], _ = roc_curve(y_test == i, y_pred_proba[:, i])
27     auc[i] = auc(fpr[i], tpr[i])

28 # Plot all ROC curves
29 plt.figure()
30 for i in range(n_classes):
31     plt.plot(fpr[i], tpr[i], label='ROC curve of class %i (AUC = %0.4f)' % (i, auc[i]))
32
33 plt.xlabel('False Positive Rate')
34 plt.ylabel('True Positive Rate')
35 plt.title('Receiver operating characteristic for multi-class data')
36 plt.legend(loc='lower right')
37 plt.show()

```

[https://colab.research.google.com/gist/texaschik/720b23ea81ebed2204bf0037a04c7962/cs2mcphaul\\_diabetes.ipynb?authuser=1#scrollTo=EY...](https://colab.research.google.com/gist/texaschik/720b23ea81ebed2204bf0037a04c7962/cs2mcphaul_diabetes.ipynb?authuser=1#scrollTo=EY...) 17/18

2/4/25, 1:56 AM McPhaul\_J\_CaseStudy2\_FINAL\_diabetes.ipynb - Colab

"  
- Logistic regression is performing well, as accuracy of 57%  
- Checking at the confusion matrix and classification report shows that:  
- "Class 0 (Not Readmitted) is being predicted well (high recall: 88%)."  
- "Class 1 (30 Days Readmission) is struggling with recall (only 23%)."  
- "Precision is high (~0.95), but recall is low (~0.23)."  
- The average F1-score of ~0.35 shows that the model isn't effectively addressing all classes equally well. This suggests a class imbalance issue, where the model is biased toward the majority class (Not Readmitted).  
- (Addressing This Issue)  
- "The logistic regression (Brynjolfsson-Greider-Solstad-Danne) optimization failed, that indicates the optimization process wasn't able to converge to a solution properly. Reasons? not sure :)"  
- "Class imbalance is the too severe."  
- "Features are not well-scaled or relevant enough."  
- "The feature selection might be dealing with high-dimensional feature spaces."  
- "Next Steps"  
- "Classification techniques"  
- "Logistic regression (logit) is a type of logistic regression model."  
- "the overfitting (MOT) or underfitting."  
- "Feature Engineering"  
- "Feature selection (FSS) or dimensionality reduction (MDR), permutation import."  
- "Feature Dimensionality Reduction (FDR) or feature selection.)"  
- "Model Selection"  
- "The model selection may not be the best for this dataset."  
- "Try Random Forest, XGBoost, or an ensemble model."  
- "I'm not sure about the preprocessing steps and train the logistic regression again."  
- "Plan of attack"  
- "1. Feature selection  
- 2. ADDRESS CLASS IMBALANCE: CHASE CHI2TEST, CLASS WEIGHTING, OVERSAMPLING"  
- "2. FEATURE SELECTION AND IMPORTANCE ANALYSIS - using SADF or permutation import to rank features, drop irrelevant or redundant")

- My logistic regression model is performing with an accuracy of 57%.
  - looking at the confusion matrix and classification report, it's clear that:
    - Class 0 (Not Readmission) is being predicted well (high recall: 96%).
    - Class 1 (>30 Days Readmission) is struggling with recall (only 23%).
    - Class 2 (<30 Days Readmission) is performing poorly (almost 0 recall).
- The average F1-score of 0.35 shows that the model isn't treating all classes equally well. This suggests a class imbalance issue.
  - Addressing this will likely improve model performance.
- Since BFGS (Limited-Memory Broyden-Fletcher-Goldfarb-Shanno) optimization failed, that indicates the optimization process wasn't able to converge.
- 1. Class imbalance is too severe.
- 2. Features are not well-scaled or relevant enough.
- 3. The solver struggles with high-dimensional feature spaces.
- # Model Selection
  - 1. Class balancing techniques
    - Try class weighting in the logistic regression model.
    - Use oversampling (SMOTE) or undersampling.
  - 2. Feature Engineering
    - Use feature selection (SHAP, permutation importance).
    - Try dimensionality reduction (PCA or feature selection).
  - 3. Model Selection
    - Logistic regression may not be the best for this dataset.
    - Try Random Forest, XGBoost, or an ensemble model.
- I assume you'll want me to diagnose the problem methodically and work it step by step.
- 5. I'm going to start with the preprocessing steps and train the logistic regression model again.
  - Plan of attack:
    - 1. ADDRESS CLASS IMBALANCE: CHECK DISTRO, CLASS WEIGHTING, OVERSAMPLING
    - 2. FEATURE SELECTION AND IMPORTANCE ANALYSIS - using SHAP or permutation importances to rank features, drop irrelevant or redundant features.

```

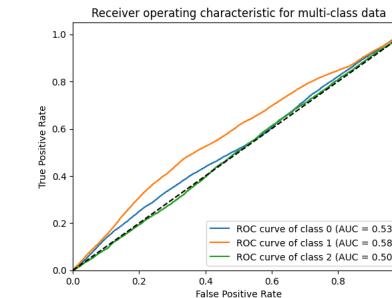
1 # Initialize and Train Model with L-BFGS solver
2 log_reg = LogisticRegression(solver='lbfgs', max_iter=1000, tol=1e-4) #Specify the solver
3 log_reg.fit(X_train, y_train)
4 print("Logistic Regression Model Trained Successfully (with L-BFGS)!")

```

2/4/25, 1:56 A

McPhaul J CaseStudy2 FINAL diabetes.ipynb - Cola

```
    ↴ /usr/local/lib/python3.11/dist-packages/sklearn/metrics/_classification.py:110: UserWarning: The metric 'precision' is deprecated in favor of 'precision_score'. It will be removed in version 1.3.  
    _warn_prf(average, modifier, f'{metric.capitalize()} is'), len(result)  
Precision: 0.4568  
Recall: 0.5422  
AUC: 0.5347
```



[https://colab.research.google.com/gist/texaschik/kita/720b23ea81ebbed2204bf0037a04c7962/cs2mcphaul\\_diabetes.ipynb?authuser=1#scrollTo=EY...](https://colab.research.google.com/gist/texaschik/kita/720b23ea81ebbed2204bf0037a04c7962/cs2mcphaul_diabetes.ipynb?authuser=1#scrollTo=EY...) 18/104

[https://colab.research.google.com/gist/texaschik/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul\\_diabetes.ipynb?authuser=1#scrollTo=EY...](https://colab.research.google.com/gist/texaschik/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul_diabetes.ipynb?authuser=1#scrollTo=EY...) 18/104

McPhaul\_J\_CaseStudy2\_FINAL\_diabetes.ipynb - Colab

```

>>> Logistic Regression Model Trained Successfully

1 # Initialize and Train Model with class weights and sage solver
2 logReg = LogisticRegression()
3 logReg.penalty='l2'
4 logReg.C=1
5 logReg.fit(X,y,weights=[1.0, 1.5, 1.0, 2.0], # Adjust weights as needed
6 solver='sag',#solver='sag',#solver='liblinear')
7 max_iter=200, # Reduce iterations
8 warm_start=True, # Continue from the last iteration
9 )

10 for i in range(1,10):
11     print("Iteration %d" % i)
12     print(logReg.Iteration[1:i+1])
13
14 # Predict on Test Data
15 y_pred = logReg.predict(Xtest)
16
17 # Compute Evaluation Metrics
18 accuracy = accuracy_score(ytest, y_pred)
19 conf_matrix = confusion_matrix(ytest, y_pred)
20 class_report = classification_report(ytest, y_pred)
21
22 # Display Results
23 print("Accuracy: %.4f" % accuracy)
24 print("Confusion Matrix:\n", conf_matrix)
25 print("Classification Report:\n", class_report)

```

```
25 print("Classification Report:\n", class_report)
26
```

→ Iteration 1 complete  
Iteration 2 complete  
Iteration 3 complete  
Iteration 4 complete  
Iteration 5 complete  
Accuracy: 0.5138

```
Confusion Matrix:  
[[20941 11978 0]  
[18894 10433 0]  
[ 3912 2902 0]]
```

```

Classification Report:
precision    recall   f1-score   support

          0       0.59      0.64      0.61     32919
          1       0.41      0.49      0.45     21327
          2       0.00      0.00      0.00      6814

accuracy                           0.51     61960
macro avg       0.33      0.38      0.35     61960
weighted avg    0.46      0.51      0.49     61960

/usr/local/lib/python3.11/dist-packages/sklearn/metrics/_classification.py:1565: UndefinedMetricWarning: Precision is ill-defined and be
    _warn_prf_average, modifier, f"metric.capitalize()") is '_len(result)'
/usr/local/lib/python3.11/dist-packages/sklearn/metrics/_classification.py:1565: UndefinedMetricWarning: Precision is ill-defined and be
    _warn_prf_average, modifier, f"metric.capitalize()") is '_len(result)'
/usr/local/lib/python3.11/dist-packages/sklearn/metrics/_classification.py:1565: UndefinedMetricWarning: Precision is ill-defined and be
    _warn_prf_average, modifier, f"metric.capitalize()") is '_len(result)'

```

```
    _warn_prtt(average, moditier, f'{metric.capitalize()}\n\n')

1 import joblib
2
3 # Save model coefficients and intercept
4 joblib.dump(logit, "Logistic_regression_model.pkl") # Removed the absolute path
5 print("Model saved successfully.")
6
```

→ Model saved successfully.

```
3 pd.DataFrame(X_train).to_csv("X_train_final.csv", index=False)
4 pd.DataFrame(y_train).to_csv("y_train_final.csv", index=False)
5 pd.DataFrame(X_test).to_csv("X_test_final.csv", index=False)
```

```
# pd.DataFrame(y_test).to_csv("y_test_final.csv", index=False)
# print("Final train/test data saved successfully")
#
```

**3** Final train/test data saved successfully.

```
1 from sklearn import datasets
2
3 # Convert Numpy array back to Pandas DataFrame
4 X_train, y_train = pd.DataFrame(X), pd.Series(y)
5 # Assuming your original features were in a DataFrame
6
7 # Convert Pandas DataFrames to RDD DataFrames
```

<https://colab.research.google.com/cist/texaschikkita/720b23ea81ebcd2204bf0037a04c7962/cells/2#scrollTo=EW...> 19/19

[https://colab.research.google.com/qist/texaschikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul\\_diabetes.ipynb](https://colab.research.google.com/qist/texaschikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul_diabetes.ipynb)?authuser=1#scrollTo=EY... 20/104



2/4/25, 1:56 AM

McPhaul\_J\_CaseStudy2\_FINAL\_diabetes.ipynb - Colab

```

1 # X_train = pd.read_csv("base_path/X_train_final.csv")
2
3 # Ensure column order consistency between training and testing data
4 X_train = X_train.reindex(columns=X_train.columns, fill_value=0)
5
6 X_test = pd.read_csv("base_path/X_test_final.csv")
7
8 # Scale test data using the loaded scaler
9 X_test_scaled = scalar.transform(X_test)
10
11 # Make predictions
12 y_pred_best = best_log_reg.predict(X_test_scaled)
13
14 # Accuracy Score
15 accuracy_score(y_true=y_test, y_pred=y_pred_best)
16 print("Best Model Accuracy: ", accuracy_score(y_true=y_test, y_pred=y_pred_best))
17
18 # Confusion Matrix
19 conf_matrix_best = confusion_matrix(y_true=y_test, y_pred=y_pred_best)
20 print("Best Model Confusion Matrix:\n", conf_matrix_best)
21
22 # Classification Report
23 report_best = classification_report(y_true=y_test, y_pred=y_pred_best)
24 print("Best Model Classification Report:\n", report_best)
25
26 # Feature Importance
27 feature_importances = pd.DataFrame({'feature': feature_name, 'importance': abs(best_log_reg.coef_[0])})
28
29 # Sort by importance
30 feature_importances = feature_importances.sort_values(by='importance', ascending=False)
31
32 # Plot
33 plt.figure(figsize=(10, 6))
34 feature_importances['feature'].plot('barh', data=feature_importances[:20])
35 plt.title("Top 20 Feature Importances (Logistic Regression)")
36 plt.xlabel("Coefficient Magnitude")
37 plt.show()

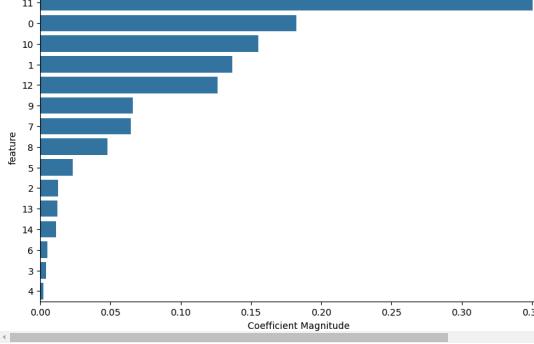
```

Best Model Accuracy: 0.5073  
 Best Model Confusion Matrix:  
[[20948 7093 4886]  
 [ 8361 7661 5303]  
 [ 2525 1914 2375]]  
 Best Model Classification Report:  

	precision	recall	f1-score	support
0	0.66	0.64	0.65	32919
1	0.46	0.36	0.40	21327
2	0.19	0.35	0.25	6014

  
 accuracy 0.51 61968  
 macro avg 0.44 0.45 0.43 61968  
 weighted avg 0.54 0.51 0.52 61968

Top 20 Feature Importances (Logistic Regression)



```

1 # Apply SMOTE to fix class imbalance
2 smote = SMOTE(sampling_strategy="auto", random_state=42)
3 X_resampled, y_resampled = smote.fit_resample(X_train_scaled, y_train)
4

```

```

5 # Verify class distributions, corr matrix, PCA grid search
6
7 import pandas as pd

```

[https://colab.research.google.com/gist/texashchikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul\\_diabetes.ipynb?authuser=1#scrollTo=EY...](https://colab.research.google.com/gist/texashchikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul_diabetes.ipynb?authuser=1#scrollTo=EY...) 25/104

2/4/25, 1:56 AM

McPhaul\_J\_CaseStudy2\_FINAL\_diabetes.ipynb - Colab

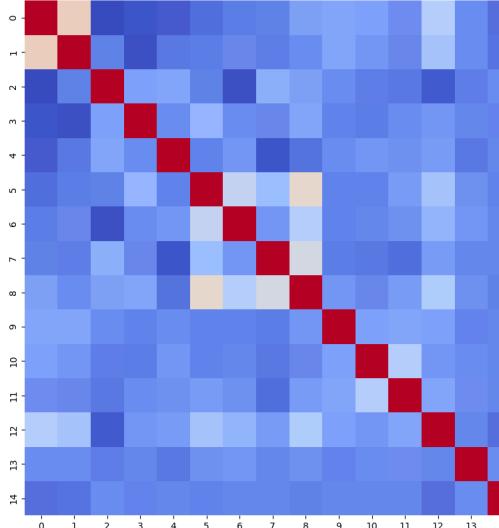
```

1 # Class Distribution in Training Data:
2 0 0.539118
3 1 0.460882
4 2 0.111600
Name: proportion, dtype: float64

```

Class Distribution in Testing Data:  
0 0.600075  
1 0.399925  
2 0.111595  
Name: proportion, dtype: float64

Correlation Matrix of Features



```

Fitting 5 Folds for each of 6 candidates, totalling 30 fits
Best Parameters: {'C': 1.0, 'max_iter': 3000, 'solver': 'lbfgs'}

```

```

1 from sklearn.metrics import roc_curve, auc, roc_auc_score, precision_score, recall_score # Import auc (area under curve)
2
3 import public
4 import numpy as np
5
6 from sklearn.metrics import accuracy_score, classification_report, confusion_matrix, precision_score, recall_score, roc_auc_score, roc_curve
7 import matplotlib.pyplot as plt
8 import seaborn as sns
9 import numpy as np
10
11 # Load the saved data
12 best_log_reg = public.load("best_logistic_regression.pkl")
13 scalar = public.load("standard_scaler.pkl")
14
15 # Load the test data
16 X_test = pd.read_csv("base_path/X_test_final.csv")
17 y_test = pd.read_csv("base_path/y_test_final.csv").values.ravel()
18
19 # Load training data (to ensure column order matches)
20 X_train = pd.read_csv("base_path/X_train_final.csv")
21
22 # Ensure column order consistency
23 X_test = X_test.reindex(columns=X_train.columns, fill_value=0)
24
25 # Scale the test data
26 X_test_scaled = scalar.transform(X_test)

```

[https://colab.research.google.com/gist/texashchikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul\\_diabetes.ipynb?authuser=1#scrollTo=EY...](https://colab.research.google.com/gist/texashchikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul_diabetes.ipynb?authuser=1#scrollTo=EY...) 27/104

2/4/25, 1:56 AM

McPhaul\_J\_CaseStudy2\_FINAL\_diabetes.ipynb - Colab

```

1 # Data predictions
2 y_pred = best_log_reg.predict(X_test_scaled)
3
4 # Evaluate the model
5 accuracy = accuracy_score(y_test, y_pred)
6 conf_matrix = confusion_matrix(y_test, y_pred)
7 class_report = classification_report(y_test, y_pred)
8
9 print("Accuracy: ", accuracy)
10 print("Confusion Matrix: \n", conf_matrix)
11 print("Classification Report:\n", class_report)
12
13 # Predict probabilities for ROC AUC
14 y_pred_proba = best_log_reg.predict_proba(X_test_scaled)
15
16 # Calculate precision, recall, and F1
17 precision = precision_score(y_test, y_pred, average='weighted')
18 recall = recall_score(y_test, y_pred, average='weighted')
19
20 # A warning is present here because different variable to avoid shadowing the auc function
21 roc_auc_score_result = roc_auc_score(y_test, y_pred_proba, multi_class='ovr')
22
23 print("Precision: ", precision)
24 print("Recall: ", recall)
25 print("AUC: ", roc_auc_score_result)
26
27 # Print the confusion matrix
28 for i in range(3):
29     for j in range(3):
30         if i == j:
31             print(f'{y_pred[i]} {y_pred[i]} - {roc_auc_score_result} {y_pred[i]}')
32
33 # Feature Importance
34 feature_importances = pd.DataFrame({'feature': feature_name, 'importance': abs(best_log_reg.coef_[0])})
35
36 feature_importances = feature_importances.sort_values(by='importance', ascending=False)
37
38 # Plot
39 plt.figure(figsize=(10, 6))
40 feature_importances['feature'].plot('barh', data=feature_importances[:20])
41 plt.title("Top 20 Feature Importances (Logistic Regression)")
42 plt.xlabel("Coefficient Magnitude")
43 plt.show()

```

[https://colab.research.google.com/gist/texashchikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul\\_diabetes.ipynb?authuser=1#scrollTo=EY...](https://colab.research.google.com/gist/texashchikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul_diabetes.ipynb?authuser=1#scrollTo=EY...) 26/104

2/4/25, 1:56 AM

McPhaul\_J\_CaseStudy2\_FINAL\_diabetes.ipynb - Colab

```

1 # Data predictions
2 y_pred = best_log_reg.predict(X_test_scaled)
3
4 # Evaluate the model
5 accuracy = accuracy_score(y_test, y_pred)
6 conf_matrix = confusion_matrix(y_test, y_pred)
7 class_report = classification_report(y_test, y_pred)
8
9 print("Accuracy: ", accuracy)
10 print("Confusion Matrix: \n", conf_matrix)
11 print("Classification Report:\n", class_report)
12
13 # Predict probabilities for ROC AUC
14 y_pred_proba = best_log_reg.predict_proba(X_test_scaled)
15
16 # Calculate precision, recall, and F1
17 precision = precision_score(y_test, y_pred, average='weighted')
18 recall = recall_score(y_test, y_pred, average='weighted')
19
20 # A warning is present here because different variable to avoid shadowing the auc function
21 roc_auc_score_result = roc_auc_score(y_test, y_pred_proba, multi_class='ovr')
22
23 print("Precision: ", precision)
24 print("Recall: ", recall)
25 print("AUC: ", roc_auc_score_result)
26
27 # Print the confusion matrix
28 for i in range(3):
29     for j in range(3):
30         if i == j:
31             print(f'{y_pred[i]} {y_pred[i]} - {roc_auc_score_result} {y_pred[i]}')
32
33 # Feature Importance
34 feature_importances = pd.DataFrame({'feature': feature_name, 'importance': abs(best_log_reg.coef_[0])})
35
36 feature_importances = feature_importances.sort_values(by='importance', ascending=False)
37
38 # Plot
39 plt.figure(figsize=(10, 6))
40 feature_importances['feature'].plot('barh', data=feature_importances[:20])
41 plt.title("Top 20 Feature Importances (Logistic Regression)")
42 plt.xlabel("Coefficient Magnitude")
43 plt.show()

```

[https://colab.research.google.com/gist/texashchikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul\\_diabetes.ipynb?authuser=1#scrollTo=EY...](https://colab.research.google.com/gist/texashchikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul_diabetes.ipynb?authuser=1#scrollTo=EY...) 28/104



```

2/4/25 1:56 AM McPhaul_J_CaseStudy2_FINAL_diabetes.ipynb - Colab

15 print("Missing Values After Imputation:", missing_values_after_imputation)
16 # Identify numerical and categorical columns
17 numerical_cols = df.select_dtypes(exclude=[object]).columns
18 categorical_cols = df.select_dtypes(include=[object]).columns
19
20 # Impute numerical columns with the median
21 numerical_imputer = SimpleImputer(strategy='median')
22 df[numerical_cols] = numerical_imputer.fit_transform(df[numerical_cols])
23
24 # Impute categorical columns with the mode
25 categorical_imputer = SimpleImputer(strategy='most_frequent')
26 df[categorical_cols] = categorical_imputer.fit_transform(df[categorical_cols])
27
28 # Verify Imputation
29
30 missing_values_after_imputation = df.isnull().sum()
31 print("Missing Values After Imputation:", missing_values_after_imputation)
32

33 rosiglitazone          0
34 acarbose               0
35 miglitol               0
36 troglitazone          0
37 tolazamide             0
38 examide                0
39 citogliptin            0
40 insulin                0
41 glyburide-metformin   0
42 glipizide-metformin   0
43 glimepiride-pioglitazone 0
44 metformin-rosiglitazone 0
45 metformin-pioglitazone 0
46 chlorpropamide         0
47 diabetesMed            0
48 readmitted             0
49 description             0
50 dtype: int64

Missing Values After Imputation:
encounter_id           0
patient_nbr             0
race                     0
gender                   0
age                      0
weight                   0
admission_type_id       0
discharge_disposition_id 0
admission_source_id     0
time_in_hospital        0
payer_id                 0
medicare_specialty      0
num_lab_procedures      0
num_procedures           0
num_medications          0
number_outpatient        0
number_emergency         0
number_inpatient         0
diag_1                   0
diag_2                   0
diag_3                   0
number_diagnoses         0
max_glu_serum            0
A1Cresult                0
metformin                0
repaglinide               0
nateglinide               0
chlorpropamide            0
glimepiride               0
acetohexamide              0
glipizide                 0
glyburide                 0
tolazamide                 0
pioglitazone               0
rosiglitazone              0
acarbose                  0
miglitol                  0
troglitazone               0
vandesartenc

```

```

import numpy as np
4
5 # Load data
6 X_train = pd.read_csv("X_train_final.csv")
7 X_train = pd.read_csv("X_train_final.csv").values.ravel()
8 X_test = pd.read_csv("X_test_final.csv")
9 y_test = pd.read_csv("y_test_final.csv").values.ravel()
10
11 # Check for missing values
12 print("Missing values in X_train\n", X_train.isnull().sum())
13 print("Missing values in X_test\n", X_test.isnull().sum())

```

[https://colab.research.google.com/gist/texashchikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul\\_diabetes.ipynb?authuser=1#scrollTo=EY...](https://colab.research.google.com/gist/texashchikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul_diabetes.ipynb?authuser=1#scrollTo=EY...)

```
4/25/21, 1:56 AM McPhaul_J_CaseStudy2_FINAL_diabetes.ipynb - Colab
-----  

2 0.11595

1 # handle class imbalance with either weight or class weights in model training  

2  

3 from imblearn.over_sampling import SMOTE  

4 from sklearn.linear_model import LogisticRegression  

5 from sklearn.metrics import classification_report, accuracy_score  

6 from sklearn.utils.class_weight import compute_class_weight  

7 from sklearn.preprocessing import StandardScaler  

8  

9  

10 # load data  

11 X_train = pd.read_csv("X_train_final.csv")  

12 y_train = pd.read_csv("y_train_final.csv").values.ravel()  

13 X_test = pd.read_csv("X_test_final.csv")  

14 y_test = pd.read_csv("y_test_final.csv").values.ravel()  

15  

16 # Scale data  

17  

18 # Create StandardScaler  

19 X_train_scaled = StandardScaler().fit_transform(X_train)  

20 X_test_scaled = StandardScaler().transform(X_test)  

21  

22 # Option 1: SMOTE Synthetic Minority Over-sampling Technique  

23 smote = SMOTE(sampling_strategy='auto', random_state=42) # adjust sampling_strategy as needed  

24 X_train_resampled, y_train_resampled = smote.fit_resample(X_train_scaled, y_train)  

25  

26 # Train a model with resampled data  

27 model = LogisticRegression(max_iter=1000) # Or any other model  

28 model.fit(X_train_resampled, y_train_resampled)  

29 y_pred_smote = model.predict(X_test)  

30 print("Classification Report (SMOTE):", classification_report(y_test, y_pred_smote))  

31 print(f"Accuracy (SMOTE): {accuracy_score(y_test, y_pred_smote)}: {df}")  

32  

33 # Option 2: Class weights  

34 # Calculate class weights  

35 class_weights = compute_class_weight('balanced', classes=np.unique(y_train), y=y_train)  

36 class_weight_dict = dict(enumerate(class_weights))  

37  

38 # Train a model with class weights  

39 model = LogisticRegression(class_weight=class_weight_dict, max_iter=1000) # Or any other model  

40 model.fit(X_train_scaled, y_train) # No oversampling needed  

41 y_pred_weights = model.predict(X_test_scaled)  

42 print("Classification Report (Class Weights):", classification_report(y_test, y_pred_weights))  

43 print(f"Accuracy (Class Weights): {accuracy_score(y_test, y_pred_weights)}: {df}")
```

```

Classification Report (SMOTE):
precision    recall   f1-score  support
          0       0.66    0.63     0.64    32919
          1       0.46    0.36     0.40    21327
          2       0.19    0.36     0.25    6814

accuracy                           0.50    61868
macro avg       0.43    0.45     0.43    61868
weighted avg    0.54    0.50     0.52    61868

Accuracy (SMOTE): 0.5046

Classification Report (Class Weights):
precision    recall   f1-score  support
          0       0.66    0.64     0.65    32919
          1       0.46    0.36     0.48    21327
          2       0.19    0.35     0.25    6814

accuracy                           0.51    61868
macro avg       0.44    0.45     0.43    61868
weighted avg    0.54    0.51     0.52    61868

```

```
Accuracy (Class Weights): 0.5073

1 # model selection and training setup train test with cross validation confirm proper algorit and train logistic regression with cross valid
```

```

3 from sklearn.utils.class_weight import compute_class_weight
4 from sklearn.ensemble import RandomForestClassifier
5
6 # Model Selection and Training with Cross-Validation
7
8 # Define the parameter grid for Logistic Regression
9 param_grid = {
10     'C': [0.1, 1, 10], # Regularization strength
11     'penalty': ['l1', 'l2'], # Penalty type
12     'solver': ['liblinear', 'sag']} # Solvers compatible with both penalties
13 }
14
15 # Initialize Logistic Regression model
16 logreg = LogisticRegression(max_iter=3000)
17
18 # Initialize GridSearchCV
19 grid_search = GridSearchCV(estimator=logreg, param_grid=param_grid, cv=5, scoring='accuracy')

```

```
24/25, 1:56 AM          McPhaul_J_CaseStudy2_FINAL_diabetes.ipynb - Colab

1
15 # Feature overview
16 print("Feature overview")
17 print(x_train.info())
18 print("Number of features:", len(x_train.columns))
19
20
20 # Target variable
21 print("Target variable (y_train):")
22 print(y_train)
23
24
24 # Target variable distribution
25 print("Target variable distribution (x_train):")
26 print(pd.Series(x_train).value_counts(normalize=True))
27 print("Target variable distribution (y_train):")
28 print(pd.Series(y_train).value_counts(normalize=True))
29
30
30 # Class imbalance
31 print("Class imbalance")
32 class_counts = pd.Series(y_train).value_counts()
33 if len(class_counts) > 1:
34     imbalance_ratio = class_counts.max() / class_counts.min()
35     print(f"Imbalance ratio: {imbalance_ratio:.2f}")
36 else:
37     print("Only one class present in the training data.")
38
39 print("Class imbalance (y_test):")
40 class_counts = pd.Series(y_test).value_counts()
41 if len(class_counts) > 1:
42     imbalance_ratio = class_counts.max() / class_counts.min()
43     print(f"Imbalance ratio: {imbalance_ratio:.2f}")
44 else:
45     print("Only one class present in the testing data.")
46
47
47 Missing values in X_test:
        0      0
        1      0
        2      0
        3      0
        4      0
        5      0
        6      0
        7      0
        8      0
        9      0
       10     0
       11     0
       12     0
       13     0
       14     0
dtype: int64

Feature overview for X_train:
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 244238 entries, 0 to 244237
Data columns (total 15 columns):
 #   Column   Non-Null Count   Dtype    
---  -- 
 0   244238   non-null      float64  
 1   244238   non-null      float64  
 2   244238   non-null      float64  
 3   244238   non-null      float64  
 4   244238   non-null      float64  
 5   244238   non-null      float64  
 6   244238   non-null      float64  
 7   244238   non-null      float64  
 8   244238   non-null      float64  
 9   244238   non-null      float64  
 10  244238   non-null      float64  
 11  244238   non-null      float64  
 12  244238   non-null      float64  
 13  244238   non-null      float64  
 14  244238   non-null      float64
```

For more information about the study, please contact Dr. Michael J. Hwang at (310) 206-6500 or via email at [mhwang@ucla.edu](mailto:mhwang@ucla.edu).

2/4/25, 1:56 AM McPhaul\_J\_CaseStudy2\_FINAL\_diabetes.ipynb - Colab

```
 21 # Fit the grid search to the training set
 22 grid_search.fit(X_train_scaled, y_train)
 23
 24 # Get the best model
 25 best_logreg = grid_search.best_estimator_
 26
 27 # Print the best parameters and score
 28 print("Best Parameters: ", grid_search.best_params_)
 29 print("Best Cross-Validation Score: ", grid_search.best_score_)
 30
 31
 32 # Evaluate the best model on the test set
 33 y_pred = best_logreg.predict(X_test_scaled)
 34 accuracy = accuracy_score(y_true, y_pred)
 35 print("Test Accuracy: ({accuracy:.4f})")
 36
 37
```

```
Test Accuracy: 0.5882

1 # Save the base model to a file
2 joblib.dump(basecnn, "basecnn.pkl")
3 print("Model saved as baseModel.pkl")
4
5 # To load the model later:
6 loaded_model = joblib.load("baseModel.pkl")
7 y_pred_loaded = loaded_model.predict(x_te_scaled)
```

```

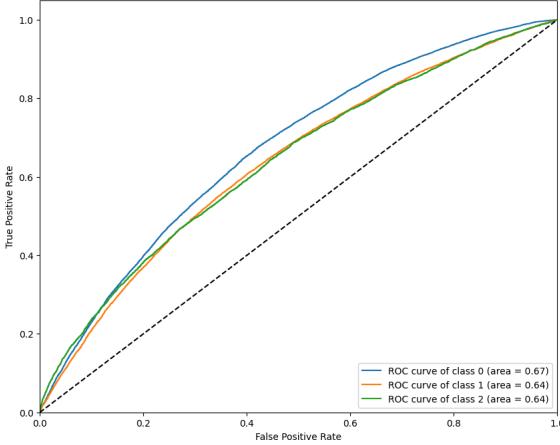
# A prints["F1 loaded Model Test Accuracy: (accuracy_score(y_test, y_pred_loaded))"]
4 # Model saved as bestModel.pkl

5 # ROC Curve and AUC Multi-class
6 from sklearn.metrics import roc_curve, auc
7 from sklearn.preprocessing import LabelBinarizer
8
9 classes = len(np.unique(y_test))
10 y_test_bin = LabelBinarizer().fit(y_test).classes_.unique(y_test)) # Binarize the output
11 fpr = dict()
12 tpr = dict()
13 roc_auc = dict()
14
15 for i in range(len(classes)):
16     fpr[i], tpr[i], _ = roc_curve(y_test_bin[:, i], y_pred_proba[:, i])
17     roc_auc[i] = auc(fpr[i], tpr[i])
18
19 # Plot ROC curve for each class
20 plt.figure(figsize=(10, 8))
21 for i in range(len(classes)):
22     plt.plot(fpr[i], tpr[i], label=f'ROC curve of class {i} (AUC = {roc_auc[i]:.2f})')
23
24 plt.plot([0, 1], [0, 1], 'k--') # Random classifier line
25 plt.xlabel('False Positive Rate')
26 plt.ylabel('True Positive Rate')
27 plt.title('ROC Curve for Multi-Class')
28 plt.legend(loc='lower right')
29 plt.show()
30
31
32 # Feature Importance (Top 5)
33 feature_importances = pd.DataFrame({'feature': X_train.columns, 'importance': abs(best_log_reg.coef_[0])})
34 feature_importances = feature_importances.sort_values(by='importance', ascending=False)
35 print("Top 5 Important Features:")
36 print(feature_importances.head(5))
37
38 # Final Submission Check
39
40 # Assess predictions to CSV file
41
42 submission_df['prediction'] = y_pred
43 submission_df.to_csv('final_predictions.csv', index=False) # Save to a csv file
44 print("Final predictions saved to 'final_predictions.csv'")
45
```

[https://colab.research.google.com/gist/texaschikkita/720b23ea81ebded2204bf0037a04c7962/cs2mcphaul\\_diabetes.ipynb?authuser=1#scrollTo=EY...](https://colab.research.google.com/gist/texaschikkita/720b23ea81ebded2204bf0037a04c7962/cs2mcphaul_diabetes.ipynb?authuser=1#scrollTo=EY...) 36/104



## Receiver Operating Characteristic (ROC) for Multi-Class



```
Top 5 Important Features:
feature importance
11 11 0.428028
0 8 0.245934
10 10 0.219261
12 12 0.179759
1 1 0.172237
Final predictions saved to 'final_predictions.csv'
```

```
1 # LOG REGRESSION WITH 5 FOLD CV
2
3 # Initialize GridSearchCV
4 grid_search = GridSearchCV(estimator=logreg, param_grid=param_grid, cv=cv, scoring='accuracy')
```

```
1 # HOW RESULTS OF LOG REGRESSION WIT 5 FOLD AMPD PLOT
2
3 from sklearn.metrics import roc_curve, auc
4 from sklearn.preprocessing import LabelBinarizer
5
6 # ROC Curve and AUC (Multi-class)
7 n_classes = len(np.unique(y_test))
8 y_test_bin = label_binarize(y_test, classes=np.unique(y_test)) # Binarize the output
9
10 for i in range(n_classes):
11     tpr[i] = dict()
12
13 for i in range(n_classes):
14     fpr[i], tpr[i], _ = roc_curve(y_test_bin[:, i], y_pred_proba[:, i])
15     roc_auc[i] = auc(fpr[i], tpr[i])
16
17 # Plot ROC curves for each class
18 plt.figure(figsize=(10, 8))
19 for i in range(n_classes):
20     plt.plot(fpr[i], tpr[i], label='ROC curve of class {} (area = {:.2f})'.format(i, roc_auc[i]))
21
22 plt.plot([0, 1], [0, 1], '--', color='black') # Random classifier line
23 plt.xlabel('False Positive Rate')
24 plt.ylabel('True Positive Rate')
25 plt.title('Receiver Operating Characteristic (ROC) for Multi-Class')
26 plt.legend(loc='lower right')
27 plt.show()
```

[https://colab.research.google.com/gist/texashchikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul\\_diabetes.ipynb](https://colab.research.google.com/gist/texashchikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul_diabetes.ipynb)?authuser=1#scrollTo=EY... 37/104

```
max 28.000000 25.000000 3.217324e+00
num_lab_procedures num_procedures num_medications number_outpatient \
count 3.052980e+05 3.052980e+05 3.052980e+05
mean 1.171600e-16 3.083771e-17 1.366634e-16
std 1.000000e+00 1.000000e+00 1.000000e+00
min -2.139632e-09 -7.853977e-01 -1.848260e-01
25% -6.147959e-01 -7.853977e-01 -2.146151e-01
50% 4.596600e-02 -1.991621e-01 -1.257264e-01
75% 7.067282e-01 3.876736e-01 4.894670e-01
max 4.518815e+00 2.732816e+00 7.994826e+00
number_emergency number_inpatient number_diagnoses max_glu_serum \
count 3.052980e+05 3.052980e+05 3.052980e+05 3.052980e+05
mean 6.665600e-17 -4.729225e-17 2.465155e-16 1.986981
std 1.000000e+00 1.000000e+00 1.000000e+00 0.194341
min -2.126202e-01 -5.837623e-01 -3.131600e-01 1.000000
25% -2.126202e-01 -5.837623e-01 -2.973320e-01 2.000000
50% -2.126202e-01 -5.837623e-01 2.986119e-01 2.000000
75% -2.126202e-01 2.885790e-01 8.157845e-01 2.000000
max 8.146673e+01 1.612568e+01 4.435992e+00 3.000000
```

```
AUC recall
count 305298.000000 305298.000000
mean 0.572488
std 0.358837
min 0.000000
25% 0.000000
50% 0.000000
75% 0.000000
max 3.000000 2.000000
```

Feature Engineering:

Model Comparison:

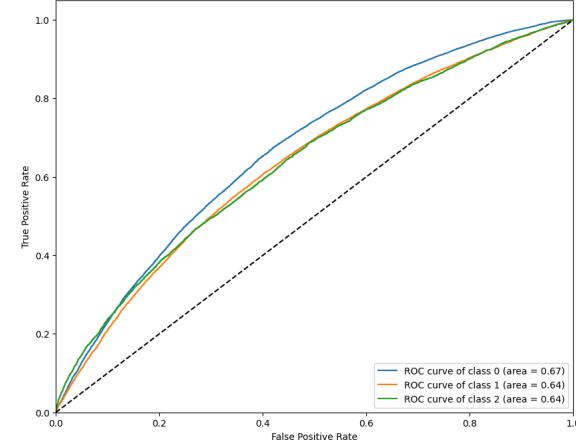
Hyperparameter Tuning:

```
1 # add clustering with l1 and l2 models
2
3 from sklearn.cluster import KMeans
4 from sklearn.metrics import silhouette_score
5 import matplotlib.pyplot as plt
6
7
8 # Clustering with l1 and l2 regularization
9
10 # L1 Regularization (l1norm)
11 kmeans_l1 = KMeans(n_clusters=4, random_state=42) # Choose optimal n_clusters using silhouette analysis
12 kmeans_l1.fit(X_train_scaled) # Fit the scaled data for clustering
13 labels_l1 = kmeans_l1.labels_
14
15 # Evaluate clustering performance
16 silhouette_avg_l1 = silhouette_score(X_train_scaled, labels_l1)
17 print("Silhouette Score (L1):", silhouette_avg_l1)
18
19 # L2 Regularization (l2norm) - Since sklearn doesn't use regularization in the same sense as linear models,
20 # L2 here is just another way to generate clustering
21 kmeans_l2 = KMeans(n_clusters=4, random_state=42)
22 kmeans_l2.fit(X_train_scaled)
23 labels_l2 = kmeans_l2.labels_
24
25 # Evaluate clustering performance
26 silhouette_avg_l2 = silhouette_score(X_train_scaled, labels_l2)
27 print("Silhouette Score (L2):", silhouette_avg_l2)
28
29 # Visualizing clustering (example with 2D reduction, adjust as needed)
30 # ... (Code to reduce dimensionality for visualization if needed) ...
31 # PCA
32 pca = PCA(n_components=2, random_state=42)
33 X_train_pca = pca.fit_transform(X_train_scaled)
34 X_train_pca = pd.DataFrame(X_train_pca)
35 X_train_pca['Labels'] = labels_l2
36
37 # Visualizing clustering (example with 2D reduction, adjust as needed)
38 # ... (Code to reduce dimensionality for visualization if needed) ...
39 # K-Means clustering with L1 regularization (visualization example)
40 kmeans_l1 = KMeans(n_clusters=4, random_state=42, n_init=10, max_iter=500, tol=1e-05, method='l1', verbose=1)
41 kmeans_l1.fit(X_train_pca)
42 labels_l1 = kmeans_l1.labels_
43
44 # Plotting K-Means clustering with L1 regularization (visualization example)
45 plt.scatter(X_train_pca[0], X_train_pca[1], c=labels_l1, s=50, alpha=0.5)
46 plt.title('K-Means clustering with L1 Regularization (visualization example)')
47 plt.legend()
48 plt.show()
```

[https://colab.research.google.com/gist/texashchikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul\\_diabetes.ipynb](https://colab.research.google.com/gist/texashchikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul_diabetes.ipynb)?authuser=1#scrollTo=EY... 39/104



## Receiver Operating Characteristic (ROC) for Multi-Class



```
1 # compile case study from diabetes analysis
2
3
4 # Data Exploration
5 print("Data Exploration:")
6 print(df.describe()) # Summary statistics
7
8
9 # Feature Engineering (if applicable)
10 print("Feature Engineering:")
11 print("Model Comparison (if you've tried other models)")
12 print("Model Comparison")
13
14
15 # Hyperparameter Tuning for other models
16 print("Hyperparameter Tuning:")
17
18
19
20
21
22 # Conclusion
23 print("Conclusion:")
24
25
```

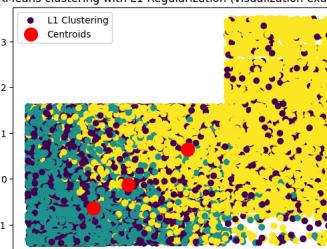
```
Data Exploration:
discharge_disposition_id admission_source_id time_in_hospital \
Count 305298.000000 305298.000000 3.052980e+05
mean 3.052980e+05 3.052980e+05 2.024906
std 1.026408e+07 3.869623e+07 1.445398
min 1.252080e+02 1.350000e+02 1.000000
25% 8.496080e+07 2.341321e+07 1.000000
50% 1.900000e+08 4.341321e+07 1.000000
75% 2.382720e+08 8.754630e+07 3.000000
max 4.438672e+08 1.895825e+08 8.000000

discharge_disposition_id admission_source_id time_in_hospital \
Count 305298.000000 305298.000000 3.052980e+05
mean 3.715642 2.754437 0.136501e-17
std 5.280148 4.064068 1.000000e+00
min 1.000000 1.000000 -1.137649e+00
25% 1.000000 1.000000 -8.026566e-01
50% 1.000000 7.000000 -1.326548e-01
75% 4.000000 7.000000 5.373411e-01
```

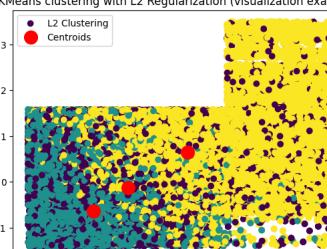
[https://colab.research.google.com/gist/texashchikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul\\_diabetes.ipynb](https://colab.research.google.com/gist/texashchikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul_diabetes.ipynb)?authuser=1#scrollTo=EY... 38/104

Silhouette Score (L1): 0.09992649512712139  
Silhouette Score (L2): 0.09992649512712139

## KMeans clustering with L1 Regularization (visualization example)



## KMeans clustering with L2 Regularization (visualization example)



```
1 # utilize SHAP, consider dimensionality reductions (such as PCA) test ensemble models (RF, XGBoost, Gradient Boosting, M) to capture non linear patterns
2
3 import shap
4 import dask
5 from pandas import pd
6 from sklearn.decomposition import PCA
7 from sklearn.ensemble import RandomForestClassifier, GradientBoostingClassifier
8 import joblib
9 from sklearn.neural_network import MLPClassifier
10 from sklearn.svm import SVC
11 from sklearn.neighbors import KNeighborsClassifier
12 from sklearn.gaussian_process import GaussianProcessClassifier
13 from sklearn.gaussian_process.kernels import RBF
14 from sklearn.metrics import accuracy_score
15 import matplotlib.pyplot as plt
16
17 # Assuming X_train_scaled, y_train, X_test_scaled, y_test are defined from previous code
18
19 # Dimensionality Reduction (PCA)
20 pca = PCA(n_components=0.95) # keep components explaining 95% of variance
21 X_train_pca = pca.fit_transform(X_train_scaled)
22 X_test_pca = pca.transform(X_test_scaled)
23
24 # Ensemble Models
25 # Model 1: Random Forest
26 model_rf = RandomForestClassifier(n_estimators=100, random_state=42)
27 # Model 2: SVM
28 model_svm = SVC()
29 # Model 3: Gradient Boosting
30 model_gb = GradientBoostingClassifier(random_state=42)
31 # Model 4: KNN
32 model_knn = KNeighborsClassifier(n_neighbors=100, random_state=42)
33
34 # SHAP Values
35 explainer = shap.TreeExplainer(model_gb) # Use TreeExplainer for tree-based models
```

[https://colab.research.google.com/gist/texashchikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul\\_diabetes.ipynb](https://colab.research.google.com/gist/texashchikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul_diabetes.ipynb)?authuser=1#scrollTo=EY... 40/104

2/4/25, 1:56 AM

McPhaul\_J\_CaseStudy2\_FINAL\_diabetes.ipynb - Colab

```

1 # if name == "Murali_Nethravathi"
2 #     explainer = shap.KernelExplainer(model.predict_proba, X_train_pca)
3 # else:
4 #     shap_values = explainer.shap_values(X_test_pca)
5 #
6 # Summary Plot
7 # shap.summary_plot(shap_values, X_test_pca, feature_names=pca.get_feature_names_out(), show=False) # Assuming your PCA has get_feature_names_out method
8 # plt.title("SHAP Summary Plot ({name})".format(name))
9 # plt.tight_layout()
10 # plt.show()
11
12 # Dependence Plot (example)
13 # shap.dependence_plot("age", shap_values, X_test_pca, feature_names=pca.get_feature_names_out()) # Replace @ with other feature index
14
15 # Print results
16 # print("Model Performance Summary:")
17 # for model, accuracy in results.items():
18 #     print("({model}): {accuracy}"
```

```

1 # pip install pympack
2 # pip install psutil
3
4 # From pympack import Archive
5 # Archive('content/diabetic_data.csv.zip').extractall('content/')
```

```

1 #shapulin inline
2 import numpy as np
3 import pandas as pd
4
5 # read in variable descriptions
6 pd.set_option('max_colwidth', 100)
7 #feature pd.read_csv('content/drive/MyDrive/IDS_mapping.csv')
8 #feature pd.read_csv('content/data_cleaned.csv')
9 #feature
```

	encounter_id	patient_nbr	race	gender	age	weight	admission_type_id	discharge_disposition_id	admission_source_id	time
0	2278392	8222157	Caucasian	Female	[0-10)	?	6	25	1	
1	2278392	8222157	Caucasian	Female	[0-10)	?	6	25	1	
2	2278392	8222157	Caucasian	Female	[0-10)	?	6	25	1	
3	149190	55629189	Caucasian	Female	[10-20)	?	1	1	7	
4	149190	55629189	Caucasian	Female	[10-20)	?	1	1	7	
...	...	...	...	...	...	...	...	...	...	
305293	443857166	31693671	Caucasian	Female	[80-90)	?	2	3	7	
305294	443857166	31693671	Caucasian	Female	[80-90)	?	2	3	7	
305295	443867222	175429310	Caucasian	Male	[70-80)	?	1	1	7	
305296	443867222	175429310	Caucasian	Male	[70-80)	?	1	1	7	
305297	443867222	175429310	Caucasian	Male	[70-80)	?	1	1	7	

305298 rows × 51 columns

```

1 #shapulin inline
2 import pandas as pd
3 import numpy as np
4 import scipy.stats as ss
5 import math
6 import os
7 import matplotlib.pyplot as plt
8 import warnings
9 warnings.filterwarnings("ignore")
10 data = pd.read_csv("/content/drive/MyDrive/WIS_Case Study 2/diabetic_data.csv")
11 data.head()
12 data.describe()
13 data.shape
14 data.columns
```

https://colab.research.google.com/gist/texashchikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul\_diabetes.ipynb?authuser=1#scrollTo=EY... 41/104

2/4/25, 1:56 AM

McPhaul\_J\_CaseStudy2\_FINAL\_diabetes.ipynb - Colab

```

1 #shapulin inline
2 import pandas as pd
3 import numpy as np
4 import scipy.stats as ss
5 import math
6 import os
7 import matplotlib.pyplot as plt
8 import warnings
9 warnings.filterwarnings("ignore")
10 data = pd.read_csv("/content/drive/MyDrive/WIS_Case Study 2/diabetic_data.csv")
11 data.head()
12 data.describe()
13 data.shape
14 data.columns
```

	encounter_id	patient_nbr	race	gender	age	weight	admission_type_id	discharge_disposition_id	admission_source_id	no.
0	2278392	8222157	Caucasian	Female	[0-10)	?	6	25	1	
1	149190	55629189	Caucasian	Female	[10-20)	?	1	1	7	
2	64410	86047875	AfricanAmerican	Female	[20-30)	?	1	1	7	
3	500364	82442376	Caucasian	Male	[30-40)	?	1	1	7	
4	16680	42519267	Caucasian	Male	[40-50)	?	1	1	7	
5	35754	82637451	Caucasian	Male	[50-60)	?	2	1	2	
6	55842	84259809	Caucasian	Male	[60-70)	?	3	1	2	
7	63768	114882984	Caucasian	Male	[70-80)	?	1	1	7	
8	12522	48330783	Caucasian	Female	[80-90)	?	2	1	4	
9	15738	63555939	Caucasian	Female	[90-100)	?	3	3	4	
10	28236	89869032	AfricanAmerican	Female	[40-50)	?	1	1	7	
11	36900	77391171	AfricanAmerican	Male	[60-70)	?	2	1	4	
12	40926	85504905	Caucasian	Female	[40-50)	?	1	3	7	
13	42570	77586282	Caucasian	Male	[80-90)	?	1	6	7	
14	62256	49726791	AfricanAmerican	Female	[60-70)	?	3	1	2	
15	73578	86328819	AfricanAmerican	Male	[60-70)	?	1	3	7	
16	77076	92519352	AfricanAmerican	Male	[50-60)	?	1	1	7	
17	84222	108662661	Caucasian	Female	[50-60)	?	1	1	7	
18	89682	107389323	AfricanAmerican	Male	[70-80)	?	1	1	7	
19	148530	69422211	?	Male	[70-80)	?	3	6	2	

20 rows × 51 columns

```

1 # sort the data by the patient number so we can clearly observe the patients have visited more than once to the hospital
2 data.sort_values(by = "patient_nbr", ascending = True,inplace=True)
3 data.head()
```

2/4/25, 1:56 AM

McPhaul\_J\_CaseStudy2\_FINAL\_diabetes.ipynb - Colab

```

1 # Index(['encounter_id', 'patient_nbr', 'race', 'gender', 'age', 'weight',
2 #        'admission_type_id', 'discharge_disposition_id', 'admission_source_id',
3 #        'time_in_hospital', 'payer_code', 'medical_specialty',
4 #        'num_lab_procedures', 'num_procedures', 'num_medications',
5 #        'number_outpatient', 'number_emergency', 'number_inpatient', 'diag_1',
6 #        'diag_2', 'diag_3', 'number_diagnoses', 'max_glu_serum', 'AlbResult',
7 #        'metformin', 'repaglinide', 'nateglinide', 'chlorpropamide',
8 #        'glipizide', 'glimepiride', 'glyburide', 'tolbutamide',
9 #        'pioglitazone', 'rosiglitazone', 'acetohexamide', 'miglitol', 'troglitazone',
10 #       'tolazamide', 'exenatide', 'citoglipron', 'insulin',
11 #       'glyburide-metformin', 'glipizide-metformin',
12 #       'metformin-pioglitazone', 'metformin-rosiglitazone',
13 #       'metformin-miglitol', 'change', 'diabetesMed', 'readmitted'],
14 #      dtype='object')
```

```

1 data.groupby('readmitted').size()
2 # now combining both the >0 and NO into one
3 data['readmitted']=data['readmitted'].replace('>0',0)
4 data['readmitted']=data['readmitted'].replace('NO',0)
5 data['readmitted']=data['readmitted'].replace('<=0',1)
6 data.groupby('readmitted').size()
7
8
9
10
11 data.head()
```

	encounter_id	patient_nbr	race	gender	age	weight	admission_type_id	discharge_disposition_id	admission_source_id	time
0	2278392	8222157	Caucasian	Female	[0-10)	?	6	25	1	
1	149190	55629189	Caucasian	Female	[10-20)	?	1	1	7	
2	64410	86047875	AfricanAmerican	Female	[20-30)	?	1	1	7	
3	500364	82442376	Caucasian	Male	[30-40)	?	1	1	7	
4	16680	42519267	Caucasian	Male	[40-50)	?	1	1	7	

5 rows × 50 columns

```

1 data.rename(columns = {"time_in_hospital": "no_of_days_admitted"}, inplace=True)
2 data.head()
```

	encounter_id	patient_nbr	race	gender	age	weight	admission_type_id	discharge_disposition_id	admission_source_id	no_o
0	2278392	8222157	Caucasian	Female	[0-10)	?	6	25	1	
1	149190	55629189	Caucasian	Female	[10-20)	?	1	1	7	
2	64410	86047875	AfricanAmerican	Female	[20-30)	?	1	1	7	
3	500364	82442376	Caucasian	Male	[30-40)	?	1	1	7	
4	16680	42519267	Caucasian	Male	[40-50)	?	1	1	7	

5 rows × 50 columns

```

1 # First count the number of encounters
2 data['no_visit'] = data.groupby('patient_nbr')[['patient_nbr']].transform('count')
3 data.head(20)
```

https://colab.research.google.com/gist/texashchikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul\_diabetes.ipynb?authuser=1#scrollTo=EY... 42/104

2/4/25, 1:56 AM

McPhaul\_J\_CaseStudy2\_FINAL\_diabetes.ipynb - Colab

	encounter_id	patient_nbr	race	gender	age	weight	admission_type_id	discharge_disposition_id	admission_source_id	no_of
4780	26264286	135	Caucasian	Female	[50-60)	?	1	1	7	
4267	24437208	135	Caucasian	Female	[50-60)	?	2	1	1	
5827	29758806	378	Caucasian	Female	[50-60)	?	3	1	1	
67608	189899286	729	Caucasian	Female	[80-90)	?	1	3	7	
17494	64331490	774	Caucasian	Female	[80-90)	?	1	1	7	

5 rows × 51 columns

```

1 # sorting the values and then removing the data which is duplicated that is the rows with duplicate data like the patients who have visited more than once
2 data.sort_values(['patient_nbr', 'encounter_id'], inplace=True)
3 data.drop_duplicates(['patient_nbr'], inplace=True)
4 data.head(50)
```

	encounter_id	patient_nbr	race	gender	age	weight	admission_type_id	discharge_disposition_id	admission_source_id	
4267	24437208	135	Caucasian	Female	[50-60)	?	2	1	1	
5827	29758806	378	Caucasian	Female	[50-60)	?	3	1	1	
67608	189899286	729	Caucasian	Female	[80-90)	?	1	3	7	
17494	64331490	774	Caucasian	Female	[80-90)	?	1	1	7	
2270	14824206	927	AfricanAmerican	Female	[30-40)	?	1	1	7	

5 rows × 51 columns

```

1 data.groupby((data.discharge_disposition_id == 1) &
2 (data.discharge_disposition_id == 2) &
3 (data.discharge_disposition_id == 3) &
4 (data.discharge_disposition_id == 4) &
5 (data.discharge_disposition_id == 5) &
6 (data.discharge_disposition_id == 6) &
7 (data.discharge_disposition_id == 7))
```

```

1 data.head(50)
2 data.shape
3 data.groupby('discharge_disposition_id').size()
```

https://colab.research.google.com/gist/texashchikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul\_diabetes.ipynb?authuser=1#scrollTo=EY... 43/104

https://colab.research.google.com/gist/texashchikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul\_diabetes.ipynb?authuser=1#scrollTo=EY... 44/104



0

discharge_disposition_id
1
2
3
4
5
6
7
8
9
10
12
15
16
17
18
22
23
24
25
27
28
40
409
73
9
6
2
40
3
8
2474
1410
260
25
778
3
90

dtype: int64

```
1: data = data[!(data['race'] != '')]
2: data.replace(to_replace='?', value=np.nan, inplace=True)
3: data.shape
4: data.isnull().sum()
```



0

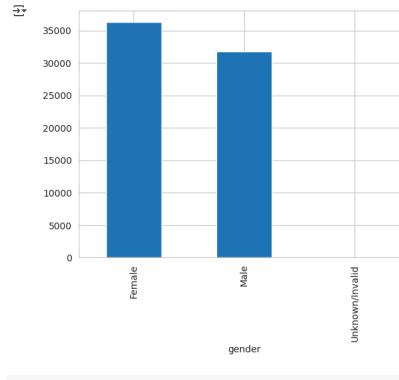
encounter_id
0
1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100
101
102
103
104
105
106
107
108
109
110
111
112
113
114
115
116
117
118
119
120
121
122
123
124
125
126
127
128
129
130
131
132
133
134
135
136
137
138
139
140
141
142
143
144
145
146
147
148
149
150
151
152
153
154
155
156
157
158
159
160
161
162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215
216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269
270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323
324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377
378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431
432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485
486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539
540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593
594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701
702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755
756
757
758
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809
810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850

	0
race	0
gender	0
age	0
admission_type_id	0
discharge_disposition_id	0
admission_source_id	0
no_of_days_admitted	0
num_lab_procedures	0
num_procedures	0
num_medication	0
number_outpatient	0
number_emergency	0
number_inpatient	0
number_diagnoses	0
max_glu_serum	647567
ATCresut	55591
metformin	0
repaglinide	0
nateglinide	0
chlorpropamide	0
glimepiride	0
acetohexamide	0
glipizide	0
glyburide	0
tolbutamide	0
piglitazone	0
rosiglitazone	0
acarbose	0
miglitol	0
troglitazone	0
tolazamide	0
examide	0
citoglipton	0
insulin	0
glyburide_metformin	0
glipizide_metformin	0
glimepiride_piglitazone	0
metformin_rosiglitazone	0
metformin_piglitazone	0
change	0
diabetesMed	0
readmitted	0
num_visits	0
first_diag	0
second_diag	0
third_diag	0
lab_rash	0
lab_cholestrol	0
lab_glu_serum	0
lab_glu_serum_norm	0



[https://codebase.research.google.com/gist/txoschmitt/f12628ca1e88e2c1bd18807d0167382e21f9phadr\\_diabetes.ipynb](https://codebase.research.google.com/gist/txoschmitt/f12628ca1e88e2c1bd18807d0167382e21f9phadr_diabetes.ipynb)

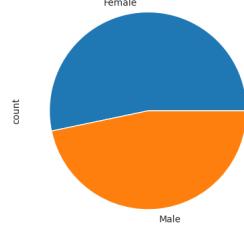
50/104



```
    count
gender
Female      36262
Male        31792
Unknown/Invalid   1

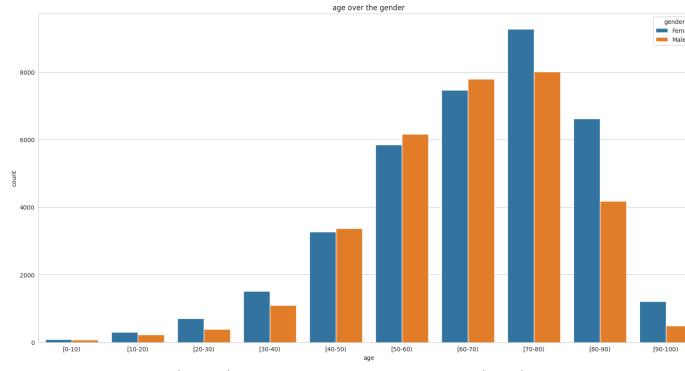
dtype: int64
```

```
1 data["gender"].value.c
```

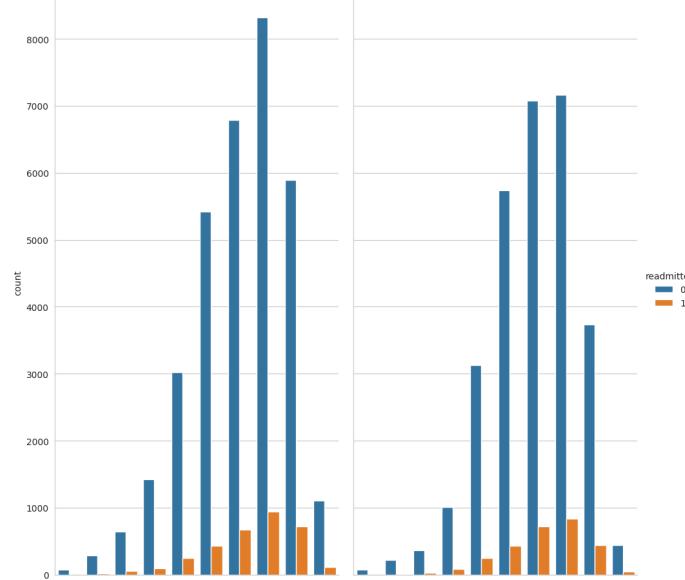


<https://cloud.google.com/appengine/docs/standard/python/references/functions#close> 51

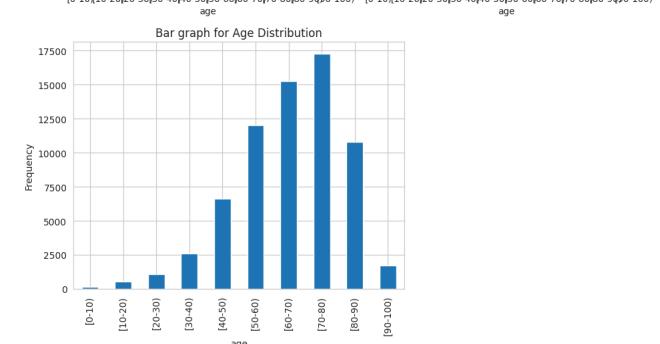
[https://colab.research.google.com/qist/texaschikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul\\_diabetes.ipynb?authuser=1#scrollTo=EY...](https://colab.research.google.com/qist/texaschikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul_diabetes.ipynb?authuser=1#scrollTo=EY...) 52/104



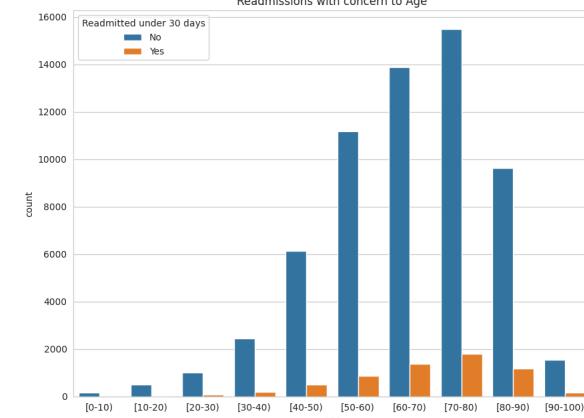
gender = Female



lab.research.google.com/gist/texaschikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul\_diabetes.ipynb?authuser=1#scrollTo=EY... 53/

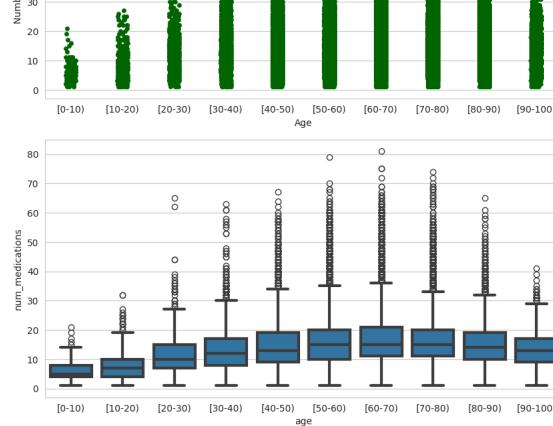


#### **Readmissions with concern to Age**



The scatter plot illustrates the relationship between age and the number of medications taken. The x-axis represents age, and the y-axis represents the number of medications. The data shows a clear positive trend, indicating that older individuals tend to take more medications.

2/4/25, 1:56 AM

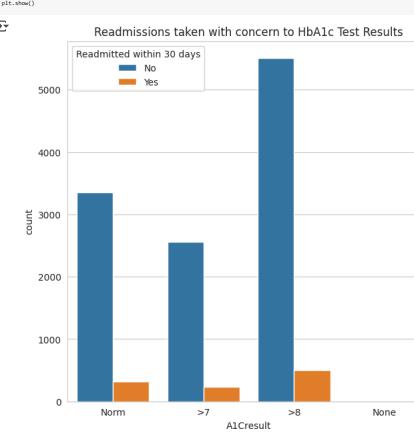


```

1 # dictionary
2 HMC_percentages = {'none': 5831/(40918+5831), '>7': 237/(3535+237), '>8': 488/(5215+488), 'normal': 318/(380+318)}
3 print(HMC_percentages)

4 HMC_percentages = {'none': 0.08919256163905682, '>7': 0.0854978354978355, '>8': 0.0856899877257584, 'normal': 0.08734107241569929}

```



1 Create new binary column to show whether HbA1c tested or not  
2 data['HbA1c'] = np.where(data['A1Cresult1'] == 'None', 0, 1)  
3  
4 Crosstab of HbA1c test and readmission w/in 30 days  
5 HbA1c\_ct = pd.crosstab(index = data['HbA1c'], columns = data['readmitted'], margins = True)  
6 HbA1c\_ct

7 **readmitted**    0    1    All      
HbA1c  

	0	1	All
HbA1c	61919	6136	68055
1	61919	6136	68055
All	61919	6136	68055

Next steps: [Generate code with HbA1c\\_ct](#) [View recommended plots](#) [New interactive sheet](#)

---

```
1 test = np.random.choice([0, 1], 100000)
2 not_tested = np.where(test == 0, 1, 0)
3 all_people = np.where(not_tested == 1, 0, 1)
4 print('test', not_tested.all())
5 print('all', all_people.all())
6 data.shape
```

7 0.088392379572207085 0.09100616160201652 0.08970602946858928  
(68055, 47)

---

```
1 def chisq_cols(df, cl1, cl2):
2     groupnames = df[cl1].groupby([cl1, cl2]).size()
3     chisum = groupnames.unstack([cl1])
4     res = pd.chi2_contingency(chisum)
5
6     print(res)
7 chisq_out1(data, 'HbA1c', 'readmitted')
8
```

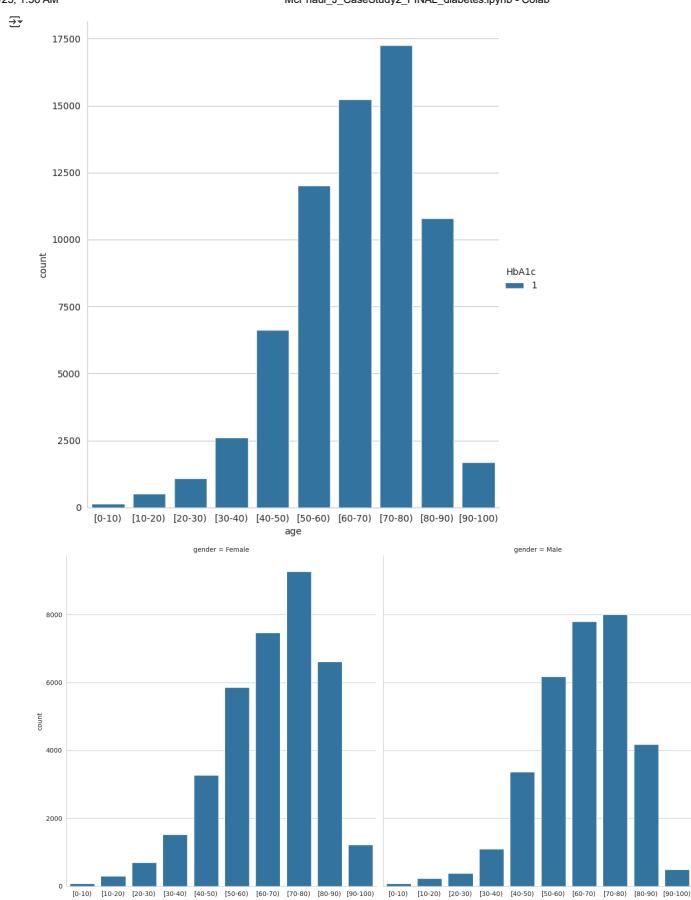
3/4/25, 1:56 AM

McRhaul\_J\_CaseStudy2\_FINAL\_diabetes.ipynb - Col

```

18
19 plt.close()
20 unique_apr['dept'] = 'apr'
21 unique_apr['sort'] = 1
22 sorted_apr = np.array(unique_apr).tolist()
23
24
25 plt.bar(x=unique_apr['dept'], height=1, label='dept')
26 unique_apr['dept'] = 'apr'
27 unique_apr['sort'] = 1
28 sorted_apr = np.array(unique_apr).tolist()
29
30 plt.bar(x=unique_apr['dept'], height=1, label='dept')
31
32
33 plt.show()

```

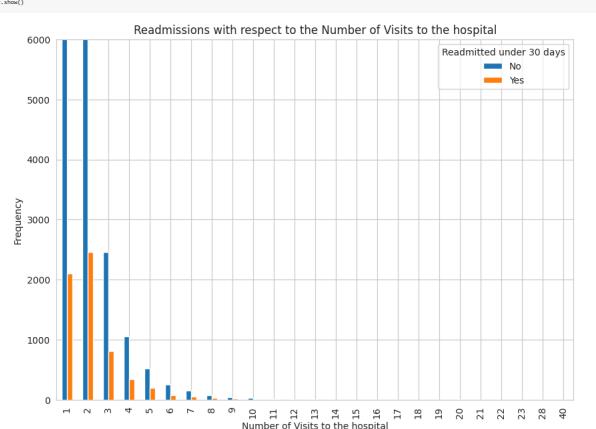


[https://colab.research.google.com/gist/texaschikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul\\_diabetes.ipynb?authuser=1#scrollTo=EY...](https://colab.research.google.com/gist/texaschikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul_diabetes.ipynb?authuser=1#scrollTo=EY...) 57/1

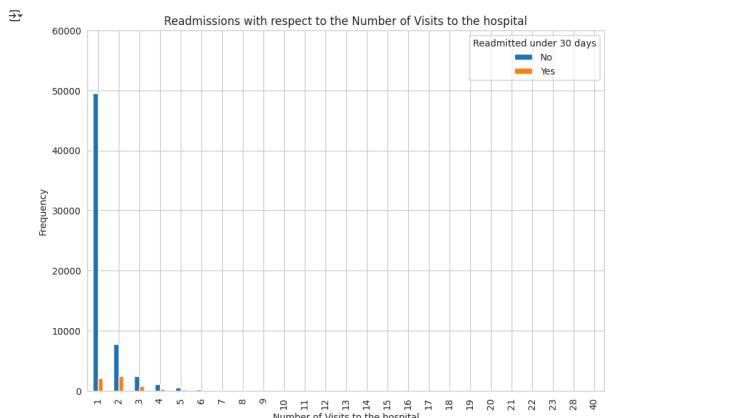
[https://colab.research.google.com/gist/texaschikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul\\_diabetes.ipynb?authuser=1#scrollTo=EY...](https://colab.research.google.com/gist/texaschikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul_diabetes.ipynb) 58/104

2/4/25, 1:56 AM

McPhaul\_J\_CaseStudy2\_FINAL\_diabetes.ipynb - Colab



```
1 v = tsplotting.plot(kind = 'bar', x = 'num_visits')
2 v.figure.set_size_inches(10, 7)
3 v.set_ylabel('0, 00000')
4 v.set_title('Number of visits to the hospital')
5 v.set_xlabel('frequency')
6 v.legend(title = 'Under 30 days', labels = ('No', 'Yes'))
7 v.axes.set_title('Redistribution with respect to the Number of Visits to the hospital')
```



```

1  Browsing the lab procedure feature using
2  def history_lab_procedures(self):
3      if (col < 1) & (col >= 20):
4          return "[11-20]"
5      if (col > 1) & (col <= 30):
6          return "[21-30]"
7      if (col > 3) & (col <= 39):
8          return "[31-39]"
9      if (col > 3) & (col <= 40):
10         return "[41-40]"
11     if (col > 3) & (col <= 50):
12         return "[41-50]"
13     if (col > 3) & (col <= 60):
14         return "[51-60]"
15     if (col > 3) & (col <= 70):
16         return "[71-80]"
17     if (col > 3) & (col <= 80):
18         return "[81-90]"
19     if (col > 3) & (col <= 90):
20         return "[91-100]"
21     if (col > 9) & (col <= 100):
22         return "[101-120]"
23     if (col > 10) & (col <= 110):
24         return "[111-120]"
25     if (col > 11) & (col <= 120):
26         return "[121-120]"
27     else:
28         return "[121-121]"

```

```
1 data['num_lab_procedure_ranges'] = data['num_lab_procedures'].apply(lambda x: binary_lab_proce  
2 data.head()
```

[https://colab.research.google.com/aist/texaschikkita/720b23ea81ebbed2204bf0037a04c7962/cs2mcphaul\\_diabetes.ipynb?authuser=1#scrollTo=EY...](https://colab.research.google.com/aist/texaschikkita/720b23ea81ebbed2204bf0037a04c7962/cs2mcphaul_diabetes.ipynb?authuser=1#scrollTo=EY...) 59/1

[https://colab.research.google.com/gist/texaschikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul\\_diabetes.ipynb?authuser=1#scrollTo=EY...](https://colab.research.google.com/gist/texaschikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul_diabetes.ipynb?authuser=1#scrollTo=EY...) 60/104



2/4/25, 1:56 AM

McPhaul\_J\_CaseStudy2\_FINAL\_diabetes.ipynb - Colab

```

1 C_grid = {'C': [0.0001, 0.001, 0.01, 0.1, 1, 10, 100,1000]}
2 weights = [0, 1, 2, 5, 10]
3 cif_grid = GridSearchCV(LogisticRegression(penalty='l2'), class_weight = weights), C_grid, cv = 5, scoring = 'accuracy'
4 cif_grid.fit(Xtrain, Ytrain)
5 cif_grid.fit(Xtrain, Ytrain)

```

`0.7903533906399236`

```

1 # best c-value and accuracy score
2 print(cif_grid.best_params_, cif_grid.best_score_)

3 { 'C': 0.01 } 0.7898942447699986

```

```

1 # classifier - cif
2 cif_grid_best = LogisticRegression(C = cif_grid.best_params_['C'], penalty='l2', class_weight = weights)
3 cif_grid_best.fit(Xtrain, Ytrain)

```

```

4 # predicting on the train data
5 x_pred_train = cif_grid_best.predict(Xtrain)
6 # getting the accuracy score
7 accuracy_score(x_pred_train, Ytrain)

```

`0.7938432150466534`

```

1 report_train = classification_report(Ytrain, x_pred_train)
2 print(report_train)

3 precision recall f1-score support
4
5 0 0.95 0.81 0.88 49535
6 1 0.23 0.57 0.33 4989
7
8 accuracy 0.79 54444
9 macro avg 0.59 0.69 0.68 54444
10 weighted avg 0.89 0.79 0.83 54444

```

```

1 report_test = classification_report(Ytest, x_pred_test)
2 print(report_test)

3 precision recall f1-score support
4
5 0 0.95 0.82 0.88 12384
6 1 0.23 0.56 0.33 1227
7
8 accuracy 0.79 13611
9 macro avg 0.59 0.69 0.68 13611
10 weighted avg 0.88 0.79 0.83 13611

```

```

1 # same as earlier here even we are using the l2 regularization and 5-fold cross-validation
2 C_grid = {'C': [0.0001, 0.001, 0.01, 0.1, 1, 10, 100,1000]}
3 cif_ROC = GridSearchCV(LogisticRegression(penalty='l2'), class_weight = weights),
4 cif_grid, cv = 5, scoring = 'roc_auc'
5 cif_ROC.fit(Xtrain, Ytrain)

```

`{'C': 0.1} 0.7758255804205956`

```

1 print(cif_ROC.best_params_, cif_ROC.best_score_)

2 { 'C': 0.1 } 0.7750255804205956

```

```

1 # best value C training and test data
2 import warnings
3 warnings.filterwarnings('ignore')
4 cif_ROC_best = LogisticRegression(penalty='l2', class_weight = weights,

```

[https://colab.research.google.com/gist/texashchikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul\\_diabetes.ipynb](https://colab.research.google.com/gist/texashchikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul_diabetes.ipynb)?authuser=1#scrollTo=EY... 65/104

2/4/25, 1:56 AM

McPhaul\_J\_CaseStudy2\_FINAL\_diabetes.ipynb - Colab

```

1 Xtest = test_cv.iloc[1,:] / test_cv.iloc[1,:]
2 print("Test accuracy for not readmitted: (" + str.format("%0.2f" % T0_test))
3 print("Test accuracy for readmitted (Recall): (" + str.format("%0.2f" % TP_test))

4 Predicted 0 1 All
5 Actual
6 0 10117 2267 12384
7 1 539 688 1227
8 All 18656 2955 13611
9 Test accuracy for not readmitted: 0.817
10 Test accuracy for readmitted (Recall): 0.561

```

## Undersampling

```

1 # independent variables
2 features = list(data_encoded)
3 features = [x for x in features if x not in ('Unreadmit', '0', 'readmitted')]
4
5
6 from collections import Counter
7 from imblearn.under_sampling import RandomUnderSampler
8
9 X = data_encoded[features].values
10 Y = data_encoded.readmitted.values
11 undersampling
12 r = RandomUnderSampler(random_state = 32)
13 X_res, Y_res = r.fit_resample(X, Y) # changed fit_sample to fit_resample
14 Counter(Y_res)

```

`Counter({0: 6136, 1: 6136})`

## train Test Split

```

1 Xtrain, Xtest, Ytrain, Ytest = train_test_split(X_res, Y_res, test_size = .2, random_state = 31, stratify = Y_res)
2

```

## Grid Search CV using L2 reg w/ 5-fold cv

```

1 C_grid = {'C': [0.0001, 0.001, 0.01, 0.1, 1, 10, 100,1000]}
2 cif_grid = GridSearchCV(LogisticRegression(penalty='l2'), class_weight = weights),
3 cif_grid.fit(Xtrain, Ytrain)
4
5 print(cif_grid.best_params_, cif_grid.best_score_)

6 { 'C': 1000 } 0.6996087695326887

```

```

1 # Accuracy on training data
2 cif_grid_best = LogisticRegression(C = cif_grid.best_params_['C'], penalty='l2')
3 cif_grid_best.fit(Xtrain, Ytrain)
4
5 x_pred_train = cif_grid_best.predict(Xtrain)
6 accuracy_score(x_pred_train, Ytrain)

```

`0.7834735662626857`

```

1 # Accuracy on Test Data
2 cif_grid_best.fit(Xtest, Ytest)
3 x_pred_test = cif_grid_best.predict(Xtest)
4 accuracy_score(x_pred_test, Ytest)

```

`0.7128162932799224`

## LR model w/ undersampling : Confusion Matrix

```

1 actual = pd.Series(Ytest, name = 'Actual')
2 predicted_tr = pd.Series(cif_grid_best.predict(Xtest), name = 'Predicted')
3 ct_mn = pd.crosstab(actual, predicted_tr, margins = True)
4 print(ct_mn)

5 F1 = 2 * (ct_mn[0,0] * ct_mn[1,1]) / (ct_mn[0,0] + ct_mn[1,1])
6 TPR_mn = ct_mn[0,0] / ct_mn[0,0] + ct_mn[0,1]
7 TPR_mn = ct_mn[1,1] / ct_mn[0,1] + ct_mn[1,1]
8 print("Logistic Regression accuracy for not readmitted: (" + str.format("%0.2f" % TPR_mn))
9 print("Logistic Regression accuracy for readmitted (Recall): (" + str.format("%0.2f" % F1))

```

`Predicted 0 1 All``Actual`

[https://colab.research.google.com/gist/texashchikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul\\_diabetes.ipynb](https://colab.research.google.com/gist/texashchikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul_diabetes.ipynb)?authuser=1#scrollTo=EY... 67/104

2/4/25, 1:56 AM

McPhaul\_J\_CaseStudy2\_FINAL\_diabetes.ipynb - Colab

```

1 # on test data
2 cif_ROC_best.fit(Xtest, Ytest)
3 probability_test = cif_ROC_best.predict_proba(Xtest)
4 predicted_test = probability_test[:,1]
5 roc_auc_score(Ytest, predicted_test)
6
7
8 0.7777947953242364

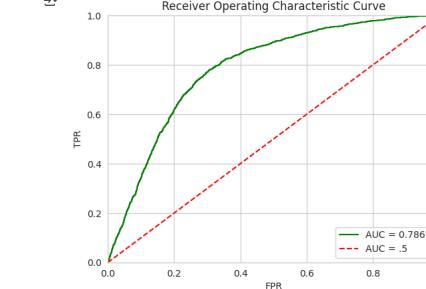
```

`0.7862166446596707`

```

1 # FPR False positive tpr = true positive
2
3 # plot ROC curve from test data
4 fpr, tpr, threshold = roc_curve(Ytest, predicted_test)
5 roc_auc = auc(fpr, tpr)
6
7 plt.plot([0,1],[0,1], color = 'red', linestyle = '--', label = 'Random Classifier Curve')
8 plt.plot(fpr, tpr, 'green', label = 'AUC = 0.7862')
9 plt.xlabel('FPR', fontweight = 'bold')
10 plt.ylabel('TPR', fontweight = 'bold')
11 plt.title('ROC Curve')
12 plt.legend()
13 plt.show()

```



```

1 # confusion matrix for train
2 actual_train = pd.Series(Ytrain, name = 'Actual')
3 predict_train = pd.Series(x_pred_train, name = 'Predicted')
4 ct_mn1 = pd.crosstab(actual_train, predict_train, margins = True)
5 print(ct_mn1)

6
7 # printing the percentage values
8 tp = ct_mn1[0,0] / ct_mn1.sum().sum()
9 fp = ct_mn1[1,0] / ct_mn1.sum().sum()
10 fn = ct_mn1[0,1] / ct_mn1.sum().sum()
11 tn = ct_mn1[1,1] / ct_mn1.sum().sum()
12 print("Training accuracy for not readmitted: (" + str.format("%0.2f" % tp))
13 print("Training accuracy for being readmitted: (" + str.format("%0.2f" % tn))

4 Predicted 0 1 All
5 Actual
6 0 40218 9317 49535
7 1 29020 8699 8699
8 All 43235 18212 54444
9 Training accuracy for not readmitted: 0.812
10 Training accuracy for being readmitted: 0.573

```

```

1 # confusion matrix for test data
2 actual_test = pd.Series(Ytest, name = 'Actual')
3 predict_test = pd.Series(x_pred_test, name = 'Predicted')
4 ct_mn2 = pd.crosstab(actual_test, predict_test, margins = True)
5 print(ct_mn2)

7 TP_test = test_cv[1].loc[0,0] / test_cv[1].loc[0,0]

```

[https://colab.research.google.com/gist/texashchikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul\\_diabetes.ipynb](https://colab.research.google.com/gist/texashchikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul_diabetes.ipynb)?authuser=1#scrollTo=EY... 66/104

2/4/25, 1:56 AM

McPhaul\_J\_CaseStudy2\_FINAL\_diabetes.ipynb - Colab

```

0 951 277 1228
1 430 797 1227
All 1381 1874 2455
Logistic Regression accuracy for not readmitted: 0.774
Logistic Regression accuracy for readmitted (Recall): 0.650

```

## SMOTE for oversampling

```

1 from imblearn.over_sampling import SMOTE
2 from collections import Counter
3
4 X = data_encoded[features].values
5 Y = data_encoded.readmitted.values
6
7 sm = SMOTE(random_state = 31)
8 # Use fit_resample instead of fit_sample
9 X_res, Y_res = sm.fit_resample(X, Y)
10 Counter(Y_res)

```

`Counter({1: 61919, 0: 61919})`

```

1 # Train Test Split
2 Xtrain, Xtest, Ytrain, Ytest = train_test_split(X_res, Y_res, test_size = .2, random_state = 31, stratify = Y_res)
3
4 # after split use the GridSearchCV with l2 regularization and 5-fold cross-validation along with the model being the Logistic Regression
5 C_grid = {'C': [0.0001, 0.001, 0.01, 0.1, 1, 10, 100,1000]}
6 cif_grid = GridSearchCV(LogisticRegression(penalty='l2'), class_weight = weights),
7 cif_grid.fit(Xtrain, Ytrain)
8
9 # Accuracy on training data
10 cif_grid_best = LogisticRegression(C = cif_grid.best_params_['C'], penalty='l2')
11 cif_grid_best.fit(Xtrain, Ytrain)
12
13 # Accuracy on test data
14 cif_grid_best.fit(Xtest, Ytest)
15 x_pred_test = cif_grid_best.predict(Xtest)
16 accuracy_score(x_pred_test, Ytest)

```

`{'C': 10} 0.7463813465226689`

```

1 # Accuracy on training data
2 cif_grid_best = LogisticRegression(C = cif_grid.best_params_['C'], penalty='l2')
3 cif_grid_best.fit(Xtrain, Ytrain)
4 x_pred_train = cif_grid_best.predict(Xtrain)
5 accuracy_score(x_pred_train, Ytrain)
6
7 # Accuracy on test data
8 cif_grid_best.fit(Xtest, Ytest)
9 x_pred_test = cif_grid_best.predict(Xtest)
10 accuracy_score(x_pred_test, Ytest)

```

`{'C': 10} 0.7511708656330749`

```

1 # F1 Score weighted
2 from sklearn.metrics import f1_score
3 f1_score(Ytest[:1811], y_pred, average='weighted')
4
5
6 0.33456521325810246

```

```

1 # F1 Score Macro
2 from sklearn.metrics import f1_score
3 f1_score(Ytest[:1811], y_pred, average='macro')
4
5
6 0.3343941448937978

```

```

1 # Confusion Matrix on train data
2 actual_tr = pd.Series(Ytrain, name = 'Actual')
3 predict_tr = pd.Series(x_pred_train, name = 'Predicted')
4 ct_mn3 = pd.crosstab(actual_tr, predict_tr, margins = True)
5 print(ct_mn3)

18 18 0 18 18 0 18 18
19 19 0 19 19 0 19 19
20 20 0 20 20 0 20 20
21 21 0 21 21 0 21 21
22 22 0 22 22 0 22 22
23 23 0 23 23 0 23 23
24 24 0 24 24 0 24 24
25 25 0 25 25 0 25 25
26 26 0 26 26 0 26 26
27 27 0 27 27 0 27 27
28 28 0 28 28 0 28 28
29 29 0 29 29 0 29 29
30 30 0 30 30 0 30 30
31 31 0 31 31 0 31 31
32 32 0 32 32 0 32 32
33 33 0 33 33 0 33 33
34 34 0 34 34 0 34 34
35 35 0 35 35 0 35 35
36 36 0 36 36 0 36 36
37 37 0 37 37 0 37 37
38 38 0 38 38 0 38 38
39 39 0 39 39 0 39 39
40 40 0 40 40 0 40 40
41 41 0 41 41 0 41 41
42 42 0 42 42 0 42 42
43 43 0 43 43 0 43 43
44 44 0 44 44 0 44 44
45 45 0 45 45 0 45 45
46 46 0 46 46 0 46 46
47 47 0 47 47 0 47 47
48 48 0 48 48 0 48 48
49 49 0 49 49 0 49 49
50 50 0 50 50 0 50 50
51 51 0 51 51 0 51 51
52 52 0 52 52 0 52 52
53 53 0 53 53 0 53 53
54 54 0 54 54 0 54 54
55 55 0 55 55 0 55 55
56 56 0 56 56 0 56 56
57 57 0 57 57 0 57 57
58 58 0 58 58 0 58 58
59 59 0 59 59 0 59 59
60 60 0 60 60 0 60 60
61 61 0 61 61 0 61 61
62 62 0 62 62 0 62 62
63 63 0 63 63 0 63 63
64 64 0 64 64 0 64 64
65 65 0 65 65 0 65 65
66 66 0 66 66 0 66 66
67 67 0 67 67 0 67 67
68 68 0 68 68 0 68 68
69 69 0 69 69 0 69 69
70 70 0 70 70 0 70 70
71 71 0 71 71 0 71 71
72 72 0 72 72 0 72 72
73 73 0 73 73 0 73 73
74 74 0 74 74 0 74 74
75 75 0 75 75 0 75 75
76 76 0 76 76 0 76 76
77 77 0 77 77 0 77 77
78 78 0 78 78 0 78 78
79 79 0 79 79 0 79 79
80 80 0 80 80 0 80 80
81 81 0 81 81 0 81 81
82 82 0 82 82 0 82 82
83 83 0 83 83 0 83 83
84 84 0 84 84 0 84 84
85 85 0 85 85 0 85 85
86 86 0 86 86 0 86 86
87 87 0 87 87 0 87 87
88 88 0 88 88 0 88 88
89 89 0 89 89 0 89 89
90 90 0 90 90 0 90 90
91 91 0 91 91 0 91 91
92 92 0 92 92 0 92 92
93 93 0 93 93 0 93 93
94 94 0 94 94 0 94 94
95 95 0 95 95 0 95 95
96 96 0 96 96 0 96 96
97 97 0 97 97 0 97 97
98 98 0 98 98 0 98 98
99 99 0 99 99 0 99 99
100 100 0 100 100 0 100 100
101 101 0 101 101 0 101 101
102 102 0 102 102 0 102 102
103 103 0 103 103 0 103 103
104 104 0 104 104 0 104 104
105 105 0 105 105 0 105 105
106 106 0 106 106 0 106 106
107 107 0 107 107 0 107 107
108 108 0 108 108 0 108 108
109 109 0 109 109 0 109 109
110 110 0 110 110 0 110 110
111 111 0 111 111 0 111 111
112 112 0 112 112 0 112 112
113 113 0 113 113 0 113 113
114 114 0 114 114 0 114 114
115 115 0 115 115 0 115 115
116 116 0 116 116 0 116 116
117 117 0 117 117 0 117 117
118 118 0 118 118 0 118 118
119 119 0 119 119 0 119 119
120 120 0 120 120 0 120 120
121 121 0 121 121 0 121 121
122 122 0 122 122 0 122 122
123 123 0 123 123 0 123 123
124 124 0 124 124 0 124 124
125 125 0 125 125 0 125 125
126 126 0 126 126 0 126 126
127 127 0 127 127 0 127 127
128 128 0 128 128 0 128 128
129 129 0 129 129 0 129 129
130 130 0 130 130 0 130 130
131 131 0 131 131 0 131 131
132 132 0 132 132 0 132 132
133 133 0 133 133 0 133 133
134 134 0 134 134 0 134 134
135 135 0 135 135 0 135 135
136 136 0 136 136 0 136 136
137 137 0 137 137 0 137 137
138 138 0 138 138 0 138 138
139 139 0 139 139 0 139 139
140 140 0 140 140 0 140 140
141 141 0 141 141 0 141 141
142 142 0 142 142 0 142 142
143 143 0 143 143 0 143 143
144 144 0 144 144 0 144 144
145 145 0 145 145 0 145 145
146 146 0 146 146 0 146 146
147 147 0 147 147 0 147 147
148 148 0 148 148 0 148 148
149 149 0 149 149 0 149 149
150 150 0 150 150 0 150 150
151 151 0 151 151 0 151 151
152 152 0 152 152 0 152 152
153 153 0 153 153 0 153 153
154 154 0 154 154 0 154 154
155 155 0 155 155 0 155 155
156 156 0 156 156 0 156 156
157 157 0 157 157 0 157 157
158 158 0 158 158 0 158 158
159 159 0 159 159 0 159 159
160 160 0 160 160 0 160 160
161 161 0 161 161 0 161 161
162 162 0 162 162 0 162 162
163 163 0 163 163 0 163 163
164 164 0 164 164 0 164 164
165 165 0 165 165 0 165 165
166 166 0 166 166 0 166 166
167 167 0 167 167 0 167 167
168 168 0 168 168 0 168 168
169 169 0 169 169 0 169 169
170 170 0 170 170 0 170 170
171 171 0 171 171 0 171 171
172 172 0 172 172 0 172 172
173 173 0 173 173 0 173 173
174 174 0 174 174 0 174 174
175 175 0 175 175 0 175 175
176 176 0 176 176 0 176 176
177 177 0 177 177 0 177 177
178 178 0 178 178 0 178 178
179 179 0 179 179 0 179 179
180 180 0 180 180 0 180 180
181 181 0 181 181 0 181 181
182 182 0 182 182 0 182 182
183 183 0 183 183 0 183 183
184 184 0 184 184 0 184 184
185 185 0 185 185 0 185 185
186 186 0 186 186 0 186 186
187 187 0 187 187 0 187 187
188 188 0 188 188 0 188 188
189 189 0 189 189 0 189 189
190 190 0 190 190 0 190 190
191 191 0 191 191 0 191 191
192 192 0 192 192 0 192 192
193 193 0 193 193 0 193 193
194 194 0 194 194 0 194 194
195 195 0 195 195 0 195 195
196 196 0 196 196 0 196 196
197 197 0 197 197 0 197 197
198 198 0 198 198 0 198 198
199 199 0 199 199 0 199 199
200 200 0 200 200 0 200 200
201 201 0 201 201 0 201 201
202 202 0 202 202 0 202 202
203 203 0 203 203 0 203 203
204 204 0 204 204 0 204 204
205 205 0 205 205 0 205 205
206 206 0 206 206 0 206 206
207 207 0 207 207 0 207 207
208 208 0 208 208 0 208 208
209 209 0 209 209 0 209 209
210 210 0 210 210 0 210 210
211 211 0 211 211 0 211 211
212 212 0 212 212 0 212 212
213 213 0 213 213 0 213 213
214 214 0 214 214 0 214 214
215 215 0 215 215 0 215 215
216 216 0 216 216 0 216 216
217 217 0 217 217 0 217 217
218 218 0 218 218 0 218 218
219 219 0 219 219 0 219 219
220 220 0 220 220 0 220 220
221 221 0 221 221 0 221 221
222 222 0 222 222 0 222 222
223 223 0 223 223 0 223 223
224 224 0 224 224 0 224 224
225 225 0 225 225 0 225 225
226 226 0 226 226 0 226 226
227 227 0 227 227 0 227 227
228 228 0 228 228 0 228 228
229 229 0 229 229 0 229 229
230 230 0 230 230 0 230 230
231 231 0 231 231 0 231 231
232 232 0 232 232 0 232 232
233 233 0 233 233 0 233 233
234 234 0 234 234 0 234 234
235 235 0 235 235 0 235 235
236 236 0 236 236 0 236 236
237 237 0 237 237 0 237 237
238 238 0 238 238 0 238 238
239 239 0 239 239
```

2/4/25, 1:56 AM

McPhaul\_J\_CaseStudy2\_FINAL\_diabetes.ipynb - Colab

```

2 # Confusion matrix on test data
3 # Confusion matrix with SMOTE oversampling (test data)
4 actual = pd.Series(Ytest, name = 'Actual')
5 predicted = pd.Series(cif_grid_best.predict(Xtest), name = 'Predicted')
6 cif_cm = pd.crosstab(actual, predicted, margins = True)
7 print(cif_cm)

8 X_res = cif_cm.iloc[0, :] / cif_cm.sum(0, 2)
9 Y_res = cif_cm.iloc[1, :] / cif_cm.sum(1, 2)
10 Prev_mis = cif_cm.iloc[1, 1] / cif_cm.sum(1, 2)
11 Prev_mis = round(Prev_mis * 100, 2)
12 print('Accuracy for not readmitted: (' + str(1 - Prev_mis) + '%)')
13 print('Accuracy for readmitted: (' + str(Prev_mis) + '%)')
14 print('Correct Positive Predictions (Precision): (' + str(Prev_mis) + '%')
15 print('Correct Negative Predictions (Recall): ' + str(1 - Prev_mis) + '%')

Predicted 0 1 All
Actual
0 9381 3083 12384
1 3168 9224 12384
All 12541 12227 24768
Accuracy for not readmitted: 0.758
Accuracy for readmitted (Recall): 0.745
Correct Positive Predictions (Precision): 0.754
```

```

1 logistic_coeff = cif_grid_best.coef_[0]
2 logistic_coeff_DF = pd.DataFrame({'feature': features, 'coefficient': logistic_coeff})
3 logistic_dF = logistic_coeff_DF.sort_values('coefficient', ascending = False)
4 logistic_dF.head(10)
```

	feature	coefficient
27	migitol	2.181177
18	chlorpropamide	2.056540
17	nateglinide	1.441789
33	glyburide_metformin	1.289757
16	repaglinide	0.913160
40	num_visits	0.433143
11	number_inpatient	0.183138
12	number_diagnoses	0.040317
4	discharge_disposition_id	0.028854
8	num_medications	0.018405

Next steps: [Generate code with logistic\\_dF](#) [View recommended plots](#) [New interactive sheet](#)

```

1 # getting the independent variables
2 feature = list(data.columns)
3 feature.remove('readmitted')
4 # declare empty lists for features if it's not in ('Unreadmed': 0, 'readmitted'):
```

```

5 # undersampling from majority class:
6 # from collections import Counter
7 # from imblearn.under_sampling import RandomUnderSampler
8 # from sklearn import metrics
9 # X = data_encoded[features].values
10 Y = data_encoded['readmitted'].values
```

```

11 # Undersampling Method X #
12
13 from imblearn.under_sampling import RandomUnderSampler
14 from sklearn.model_selection import train_test_split
15 from sklearn.linear_model import LogisticRegression
16 from sklearn.metrics import accuracy_score
17 from sklearn.metrics import confusion_matrix
18 from collections import Counter
19 import pandas as pd
20
21 # Number of trials
22 number_of_repetitions = 10
23
24 # Declare empty lists for true-positive and true-negative rates
25 TPR = []
26 TNR = []
27
28 # For loop for multiple trials
29 for trial in range(number_of_repetitions):
30     # Random undersampling
31     ru, ru_undersample = random_state11 * trial # Randomized seed
32     X_res, Y_res = rus.fit_resample(X, Y) # Corrected method
```

[https://colab.research.google.com/gist/texashchikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul\\_diabetes.ipynb](https://colab.research.google.com/gist/texashchikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul_diabetes.ipynb)?authuser=1#scrollTo=EY... 69/104

McPhaul\_J\_CaseStudy2\_FINAL\_diabetes.ipynb - Colab

```

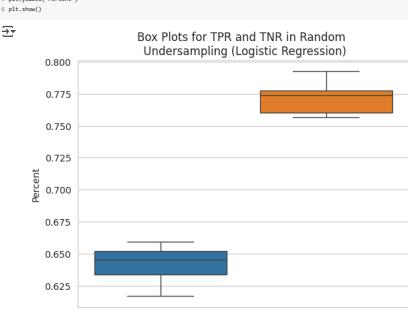
1 433 794 1227
All 1364 1891 2455
Logistic Regression accuracy for not readmitted: 0.758
Logistic Regression accuracy for readmitted (Recall): 0.647
Logistic Regression trial count: 4

Counter({0: 6136, 1: 6136})
{'C': 10} 0.788261634329999
Training Accuracy: 0.7111337475880727
Test Accuracy: 0.6924643584521385
Predicted 0 1 All
Actual
0 956 272 1228
1 1391 2624 2455
Logistic Regression accuracy for not readmitted: 0.779
Logistic Regression accuracy for readmitted (Recall): 0.645
Logistic Regression trial count: 1
```

✓ Plotting TNR & TPR

```

1 rus_boxplots = pd.DataFrame({'TPR': TPR, 'TNR': TNR})
2 sns.boxplot(data = rus_boxplots)
3 plt.title('Box Plots for TPR and TNR in Random Undersampling (Logistic Regression)')
4 plt.ylabel('Percent')
5 plt.show()
```



```

1 from imblearn.over_sampling import SMOTE
2 From imblearn.model_selection import train_test_split
3 From imblearn.metrics import accuracy_score
4 from imblearn.metrics import confusion_matrix
5 from collections import Counter
6 import pandas as pd
7
8 # Number of trials
9 number_of_repetitions = 10
10
11 # Declare empty lists for true-positive and true-negative rates
12 TPR_monte = []
13 TNR_monte = []
14
15 # For loop for multiple trials
16 for trial in range(number_of_repetitions):
17     # SMOTE oversampling
18     smote = SMOTECatSampler(random_state11 = trial) # Randomized seed
19     X_resampled, Y_resampled = smote.fit_resample(X, Y) # Corrected method
20     X_res, Y_res = rus.fit_resample(X, Y) # Print results for each trial
21
22     # Train/test split
23     Xtrain, Xtest, Ytrain, Ytest = train_test_split(
24         X_resampled, Y_resampled, test_size=0.2, stratify=Y_resampled
25     )
26
27     # Hyperparameter tuning with grid search
28     Cgrid = ('C': [0.001, 0.01, 0.1, 1, 10, 100])
29     cif_grid = GridSearchCV(Cgrid, cv=5, scoring='accuracy')
30     LogitReg = LogisticRegression(penalty='l2', Cgrid, cv=5, scoring='accuracy')
31     j = grid.fit(Xtrain, Ytrain)
32     cif_grid.fit(Xtrain, Ytrain)
33     print(cif_grid.best_params_, cif_grid.best_score_)
34
35     # Train logistic regression with the best parameter
36     cif_grid_best = LogisticRegression(C=cif_grid.best_params_['C'], penalty='l2')
37     cif_grid_best.fit(Xtrain, Ytrain)
38
39     # Evaluate on training data
40     x_pred_train = cif_grid_best.predict(Xtrain)
41     y_pred_train = cif_grid_best.predict(Ytrain)
42
43     print('Training Accuracy:', accuracy_score(Ytrain, x_pred_train))
```

[https://colab.research.google.com/gist/texashchikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul\\_diabetes.ipynb](https://colab.research.google.com/gist/texashchikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul_diabetes.ipynb)?authuser=1#scrollTo=EY... 71/104

2/4/25, 1:56 AM

McPhaul\_J\_CaseStudy2\_FINAL\_diabetes.ipynb - Colab

```

44 print(counter['C_res']) # Print results for each trial
45
46 # Train/test split
47 Xtrain, Xtest, Ytrain, Ytest = train_test_split(
48     X_res, Y_res, test_size=0.2, stratify=Y_res, random_state2 = trial
49 )
50
51 # Hyperparameter tuning with grid search
52 Cgrid = ('C': [0.0001, 0.001, 0.01, 0.1, 1, 10, 100])
53 j = grid.fit(Xtrain, Ytrain)
54 LogitReg = LogisticRegression(penalty='l2', Cgrid, cv=5, scoring='accuracy')
55 j = grid.fit(Xtrain, Ytrain)
56 print(cif_grid.best_params_, cif_grid.best_score_)
57
58 # Train logistic regression with the best parameter
59 cif_grid_best = LogisticRegression(C=cif_grid.best_params_['C'], penalty='l2')
60 cif_grid_best.fit(Xtrain, Ytrain)
61
62 # Evaluate on training data
63 x_pred_train = cif_grid_best.predict(Xtrain)
64 y_pred_train = cif_grid_best.predict(Ytrain)
65
66 print('Training Accuracy:', accuracy_score(Ytrain, x_pred_train))
67
68 # Evaluate on test data
69 x_pred_test = cif_grid_best.predict(Xtest)
70 y_pred_test = cif_grid_best.predict(Ytest)
71
72 print('Test Accuracy:', accuracy_score(Ytest, x_pred_test))
73
74 # Confusion matrix
75 actual = pd.Series(Ytest, name = 'Actual')
76 predicted = pd.Series(cif_grid_best.predict(Ytest), name='Predicted')
77
78 x_res = pd.crosstab(actual, predicted, margins=True)
79 print(x_res)
80
81 # Calculate true negative rate (TNR)
82 tnr = x_res.iloc[0, 0] / x_res.sum(0, 2)
83 TNR.append(tnr)
84
85 # Calculate true positive rate (TPR)
86 tpr = x_res.iloc[1, 1] / x_res.sum(1, 2)
87 TPR.append(tpr)
88
89 # Print metrics and trial count
90 Counter({0: 6136, 1: 6136})
91 {'C': 10} 0.707752432215945
92 Training Accuracy: 0.711623693244372
93 Test Accuracy: 0.71291929327980224
94 Predicted 0 1 All
95 Actual
96 0 956 272 1228
97 1 1391 2624 2455
98 Logistic Regression accuracy for not readmitted: 0.779
99 Logistic Regression accuracy for readmitted (Recall): 0.645
100 Logistic Regression trial count: 1
```

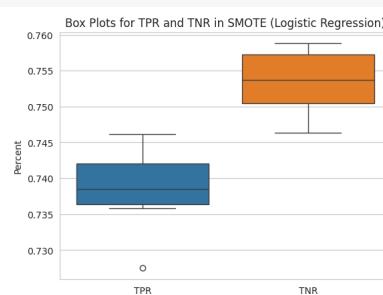
```

101 Counter({0: 6136, 1: 6136})
102 {'C': 10} 0.707752432215945
103 Training Accuracy: 0.71162369646531527
104 Test Accuracy: 0.70598063136456212
105 Predicted 0 1 All
106 Actual
107 0 949 279 1228
108 1 470 757 1227
109 All 1419 1836 2455
110 Logistic Regression accuracy for not readmitted: 0.773
111 Logistic Regression accuracy for readmitted (Recall): 0.617
112 Logistic Regression trial count: 2
113
114 Counter({0: 6136, 1: 6136})
115 {'C': 10} 0.707752432215945
116 Training Accuracy: 0.71099646531527
117 Test Accuracy: 0.70598063136456212
118 Predicted 0 1 All
119 Actual
120 0 951 277 1228
121 1 445 782 1227
122 All 1396 1859 2455
123 Logistic Regression accuracy for not readmitted: 0.774
124 Logistic Regression accuracy for readmitted (Recall): 0.637
125 Logistic Regression trial count: 3
126
127 Counter({0: 6136, 1: 6136})
128 {'C': 10} 0.707854368926258
129 Training Accuracy: 0.71346284963493493
130 Test Accuracy: 0.70264767041406
131 Predicted 0 1 All
132 Actual
133 0 951 277 1228
134 1 445 782 1227
135 All 1396 1859 2455
136 Logistic Regression accuracy for not readmitted: 0.774
137 Logistic Regression accuracy for readmitted (Recall): 0.637
138 Logistic Regression trial count: 3
139
140 Counter({0: 6136, 1: 6136})
141 {'C': 10} 0.707854368926258
142 Training Accuracy: 0.71346284963493493
143 Test Accuracy: 0.70264767041406
144 Predicted 0 1 All
145 Actual
146 0 931 297 1228
147 Logistic Regression accuracy for not readmitted: 0.774
148 Logistic Regression accuracy for readmitted (Recall): 0.637
149 Logistic Regression trial count: 3
150
151 Counter({0: 6136, 1: 6136})
152 {'C': 10} 0.707854368926258
153 Training Accuracy: 0.71346284963493493
154 Test Accuracy: 0.70264767041406
155 Predicted 0 1 All
156 Actual
157 0 931 297 1228
158 Logistic Regression accuracy for not readmitted: 0.774
159 Logistic Regression accuracy for readmitted (Recall): 0.637
160 Logistic Regression trial count: 3
161
162 Counter({0: 6136, 1: 6136})
163 {'C': 10} 0.707854368926258
164 Training Accuracy: 0.71346284963493493
165 Test Accuracy: 0.70264767041406
166 Predicted 0 1 All
167 Actual
168 0 931 297 1228
169 Logistic Regression accuracy for not readmitted: 0.774
170 Logistic Regression accuracy for readmitted (Recall): 0.637
171 Logistic Regression trial count: 3
172
173 Counter({0: 6136, 1: 6136})
174 {'C': 10} 0.707854368926258
175 Training Accuracy: 0.71346284963493493
176 Test Accuracy: 0.70264767041406
177 Predicted 0 1 All
178 Actual
179 0 931 297 1228
180 Logistic Regression accuracy for not readmitted: 0.774
181 Logistic Regression accuracy for readmitted (Recall): 0.637
182 Logistic Regression trial count: 3
183
184 Counter({0: 6136, 1: 6136})
185 {'C': 10} 0.707854368926258
186 Training Accuracy: 0.71346284963493493
187 Test Accuracy: 0.70264767041406
188 Predicted 0 1 All
189 Actual
190 0 931 297 1228
191 Logistic Regression accuracy for not readmitted: 0.774
192 Logistic Regression accuracy for readmitted (Recall): 0.637
193 Logistic Regression trial count: 3
194
195 Counter({0: 6136, 1: 6136})
196 {'C': 10} 0.707854368926258
197 Training Accuracy: 0.71346284963493493
198 Test Accuracy: 0.70264767041406
199 Predicted 0 1 All
200 Actual
201 0 931 297 1228
202 Logistic Regression accuracy for not readmitted: 0.774
203 Logistic Regression accuracy for readmitted (Recall): 0.637
204 Logistic Regression trial count: 3
205
206 Counter({0: 6136, 1: 6136})
207 {'C': 10} 0.707854368926258
208 Training Accuracy: 0.71346284963493493
209 Test Accuracy: 0.70264767041406
210 Predicted 0 1 All
211 Actual
212 0 931 297 1228
213 Logistic Regression accuracy for not readmitted: 0.774
214 Logistic Regression accuracy for readmitted (Recall): 0.637
215 Logistic Regression trial count: 3
216
217 Counter({0: 6136, 1: 6136})
218 {'C': 10} 0.707854368926258
219 Training Accuracy: 0.71346284963493493
220 Test Accuracy: 0.70264767041406
221 Predicted 0 1 All
222 Actual
223 0 931 297 1228
224 Logistic Regression accuracy for not readmitted: 0.774
225 Logistic Regression accuracy for readmitted (Recall): 0.637
226 Logistic Regression trial count: 3
227
228 Counter({0: 6136, 1: 6136})
229 {'C': 10} 0.707854368926258
230 Training Accuracy: 0.71346284963493493
231 Test Accuracy: 0.70264767041406
232 Predicted 0 1 All
233 Actual
234 0 931 297 1228
235 Logistic Regression accuracy for not readmitted: 0.774
236 Logistic Regression accuracy for readmitted (Recall): 0.637
237 Logistic Regression trial count: 3
238
239 Counter({0: 6136, 1: 6136})
240 {'C': 10} 0.707854368926258
241 Training Accuracy: 0.71346284963493493
242 Test Accuracy: 0.70264767041406
243 Predicted 0 1 All
244 Actual
245 0 931 297 1228
246 Logistic Regression accuracy for not readmitted: 0.774
247 Logistic Regression accuracy for readmitted (Recall): 0.637
248 Logistic Regression trial count: 3
249
250 Counter({0: 6136, 1: 6136})
251 {'C': 10} 0.707854368926258
252 Training Accuracy: 0.71346284963493493
253 Test Accuracy: 0.70264767041406
254 Predicted 0 1 All
255 Actual
256 0 931 297 1228
257 Logistic Regression accuracy for not readmitted: 0.774
258 Logistic Regression accuracy for readmitted (Recall): 0.637
259 Logistic Regression trial count: 3
260
261 Counter({0: 6136, 1: 6136})
262 {'C': 10} 0.707854368926258
263 Training Accuracy: 0.71346284963493493
264 Test Accuracy: 0.70264767041406
265 Predicted 0 1 All
266 Actual
267 0 931 297 1228
268 Logistic Regression accuracy for not readmitted: 0.774
269 Logistic Regression accuracy for readmitted (Recall): 0.637
270 Logistic Regression trial count: 3
271
272 Counter({0: 6136, 1: 6136})
273 {'C': 10} 0.707854368926258
274 Training Accuracy: 0.71346284963493493
275 Test Accuracy: 0.70264767041406
276 Predicted 0 1 All
277 Actual
278 0 931 297 1228
279 Logistic Regression accuracy for not readmitted: 0.774
280 Logistic Regression accuracy for readmitted (Recall): 0.637
281 Logistic Regression trial count: 3
282
283 Counter({0: 6136, 1: 6136})
284 {'C': 10} 0.707854368926258
285 Training Accuracy: 0.71346284963493493
286 Test Accuracy: 0.70264767041406
287 Predicted 0 1 All
288 Actual
289 0 931 297 1228
290 Logistic Regression accuracy for not readmitted: 0.774
291 Logistic Regression accuracy for readmitted (Recall): 0.637
292 Logistic Regression trial count: 3
293
294 Counter({0: 6136, 1: 6136})
295 {'C': 10} 0.707854368926258
296 Training Accuracy: 0.71346284963493493
297 Test Accuracy: 0.70264767041406
298 Predicted 0 1 All
299 Actual
300 0 931 297 1228
301 Logistic Regression accuracy for not readmitted: 0.774
302 Logistic Regression accuracy for readmitted (Recall): 0.637
303 Logistic Regression trial count: 3
304
305 Counter({0: 6136, 1: 6136})
306 {'C': 10} 0.707854368926258
307 Training Accuracy: 0.71346284963493493
308 Test Accuracy: 0.70264767041406
309 Predicted 0 1 All
310 Actual
311 0 931 297 1228
312 Logistic Regression accuracy for not readmitted: 0.774
313 Logistic Regression accuracy for readmitted (Recall): 0.637
314 Logistic Regression trial count: 3
315
316 Counter({0: 6136, 1: 6136})
317 {'C': 10} 0.707854368926258
318 Training Accuracy: 0.71346284963493493
319 Test Accuracy: 0.70264767041406
320 Predicted 0 1 All
321 Actual
322 0 931 297 1228
323 Logistic Regression accuracy for not readmitted: 0.774
324 Logistic Regression accuracy for readmitted (Recall): 0.637
325 Logistic Regression trial count: 3
326
327 Counter({0: 6136, 1: 6136})
328 {'C': 10} 0.707854368926258
329 Training Accuracy: 0.71346284963493493
330 Test Accuracy: 0.70264767041406
331 Predicted 0 1 All
332 Actual
333 0 931 297 1228
334 Logistic Regression accuracy for not readmitted: 0.774
335 Logistic Regression accuracy for readmitted (Recall): 0.637
336 Logistic Regression trial count: 3
337
338 Counter({0: 6136, 1: 6136})
339 {'C': 10} 0.707854368926258
340 Training Accuracy: 0.71346284963493493
341 Test Accuracy: 0.70264767041406
342 Predicted 0 1 All
343 Actual
344 0 931 297 1228
345 Logistic Regression accuracy for not readmitted: 0.774
346 Logistic Regression accuracy for readmitted (Recall): 0.637
347 Logistic Regression trial count: 3
348
349 Counter({0: 6136, 1: 6136})
350 {'C': 10} 0.707854368926258
351 Training Accuracy: 0.71346284963493493
352 Test Accuracy: 0.70264767041406
353 Predicted 0 1 All
354 Actual
355 0 931 297 1228
356 Logistic Regression accuracy for not readmitted: 0.774
357 Logistic Regression accuracy for readmitted (Recall): 0.637
358 Logistic Regression trial count: 3
359
360 Counter({0: 6136, 1: 6136})
361 {'C': 10} 0.707854368926258
362 Training Accuracy: 0.71346284963493493
363 Test Accuracy: 0.70264767041406
364 Predicted 0 1 All
365 Actual
366 0 931 297 1228
367 Logistic Regression accuracy for not readmitted: 0.774
368 Logistic Regression accuracy for readmitted (Recall): 0.637
369 Logistic Regression trial count: 3
370
371 Counter({0: 6136, 1: 6136})
372 {'C': 10} 0.707854368926258
373 Training Accuracy: 0.71346284963493493
374 Test Accuracy: 0.70264767041406
375 Predicted 0 1 All
376 Actual
377 0 931 297 1228
378 Logistic Regression accuracy for not readmitted: 0.774
379 Logistic Regression accuracy for readmitted (Recall): 0.637
380 Logistic Regression trial count: 3
381
382 Counter({0: 6136, 1: 6136})
383 {'C': 10} 0.707854368926258
384 Training Accuracy: 0.71346284963493493
385 Test Accuracy: 0.70264767041406
386 Predicted 0 1 All
387 Actual
388 0 931 297 1228
389 Logistic Regression accuracy for not readmitted: 0.774
390 Logistic Regression accuracy for readmitted (Recall): 0.637
391 Logistic Regression trial count: 3
392
393 Counter({0: 6136, 1: 6136})
394 {'C': 10} 0.707854368926258
395 Training Accuracy: 0.71346284963493493
396 Test Accuracy: 0.70264767041406
397 Predicted 0 1 All
398 Actual
399 0 931 297 1228
400 Logistic Regression accuracy for not readmitted: 0.774
401 Logistic Regression accuracy for readmitted (Recall): 0.637
402 Logistic Regression trial count: 3
403
404 Counter({0: 6136, 1: 6136})
405 {'C': 10} 0.707854368926258
406 Training Accuracy: 0.71346284963493493
407 Test Accuracy: 0.70264767041406
408 Predicted 0 1 All
409 Actual
410 0 931 297 1228
411 Logistic Regression accuracy for not readmitted: 0.774
412 Logistic Regression accuracy for readmitted (Recall): 0.637
413 Logistic Regression trial count: 3
414
415 Counter({0: 6136, 1: 6136})
416 {'C': 10} 0.707854368926258
417 Training Accuracy: 0.71346284963493493
418 Test Accuracy: 0.70264767041406
419 Predicted 0 1 All
420 Actual
421 0 931 297 1228
422 Logistic Regression accuracy for not readmitted: 0.774
423 Logistic Regression accuracy for readmitted (Recall): 0.637
424 Logistic Regression trial count: 3
425
426 Counter({0: 6136, 1: 6136})
427 {'C': 10} 0.707854368926258
428 Training Accuracy: 0.71346284963493493
429 Test Accuracy: 0.70264767041406
430 Predicted 0 1 All
431 Actual
432 0 931 297 1228
433 Logistic Regression accuracy for not readmitted: 0.774
434 Logistic Regression accuracy for readmitted (Recall): 0.637
435 Logistic Regression trial count: 3
436
437 Counter({0: 6136, 1: 6136})
438 {'C': 10} 0.707854368926258
439 Training Accuracy: 0.71346284963493493
440 Test Accuracy: 0.70264767041406
441 Predicted 0 1 All
442 Actual
443 0 931 297 1228
444 Logistic Regression accuracy for not readmitted: 0.774
445 Logistic Regression accuracy for readmitted (Recall): 0.637
446 Logistic Regression trial count: 3
447
448 Counter({0: 6136, 1: 6136})
449 {'C': 10} 0.707854368926258
450 Training Accuracy: 0.71346284963493493
451 Test Accuracy: 0.70264767041406
452 Predicted 0 1 All
453 Actual
454 0 931 297 1228
455 Logistic Regression accuracy for not readmitted: 0.774
456 Logistic Regression accuracy for readmitted (Recall): 0.637
457 Logistic Regression trial count: 3
458
459 Counter({0: 6136, 1: 6136})
460 {'C': 10} 0.707854368926258
461 Training Accuracy: 0.71346284963493493
462 Test Accuracy: 0.70264767041406
463 Predicted 0 1 All
464 Actual
465 0 931 297 1228
466 Logistic Regression accuracy for not readmitted: 0.774
467 Logistic Regression accuracy for readmitted (Recall): 0.637
468 Logistic Regression trial count: 3
469
470 Counter({0: 6136, 1: 6136})
471 {'C': 10} 0.707854368926258
472 Training Accuracy: 0.71346284963493493
473 Test Accuracy: 0.70264767041406
474 Predicted 0 1 All
475 Actual
476 0 931 297 1228
477 Logistic Regression accuracy for not readmitted: 0.774
478 Logistic Regression accuracy for readmitted (Recall): 0.637
479 Logistic Regression trial count: 3
480
481 Counter({0: 6136, 1: 6136})
482 {'C': 10} 0.707854368926258
483 Training Accuracy: 0.71346284963493493
484 Test Accuracy: 0.70264767041406
485 Predicted 0 1 All
486 Actual
487 0 931 297 1228
488 Logistic Regression accuracy for not readmitted: 0.774
489 Logistic Regression accuracy for readmitted (Recall): 0.637
490 Logistic Regression trial count: 3
491
492 Counter({0: 6136, 1: 6136})
493 {'C': 10} 0.707854368926258
494 Training Accuracy: 0.71346284963493493
495 Test Accuracy: 0.70264767041406
496 Predicted 0 1 All
497 Actual
498 0 931 297 1228
499 Logistic Regression accuracy for not readmitted: 0.774
500 Logistic Regression accuracy for readmitted (Recall): 0.637
501 Logistic Regression trial count: 3
502
503 Counter({0: 6136, 1: 6136})
504 {'C': 10} 0.707854368926258
505 Training Accuracy: 0.71346284963493493
506 Test Accuracy: 0.70264767041406
507 Predicted 0 1 All
508 Actual
509 0 931 297 1228
510 Logistic Regression accuracy for not readmitted: 0.774
511 Logistic Regression accuracy for readmitted (Recall): 0.637
512 Logistic Regression trial count: 3
513
514 Counter({0: 6136, 1: 6136})
515 {'C': 10} 0.707854368926258
516 Training Accuracy: 0.71346284963493493
517 Test Accuracy: 0.70264767041406
518 Predicted 0 1 All
519 Actual
520 0 931 297 1228
521 Logistic Regression accuracy for not readmitted: 0.774
522 Logistic Regression accuracy for readmitted (Recall): 0.637
523 Logistic Regression trial count: 3
524
525 Counter({0: 6136, 1: 6136})
526 {'C': 10} 0.707854368926258
527 Training Accuracy: 0.71346284963493493
528 Test Accuracy: 0.70264767041406
529 Predicted 0 1 All
530 Actual
531 0 931 297 1228
532 Logistic Regression accuracy for not readmitted: 0.774
533 Logistic Regression accuracy for readmitted (Recall): 0.637
534 Logistic Regression trial count: 3
535
536 Counter({0: 6136, 1: 6136})
537 {'C': 10} 0.707854368926258
538 Training Accuracy: 0.71346284963493493
539 Test Accuracy: 0.702647
```

2/4/25, 1:56 AM

McPhaul\_J\_CaseStudy2\_FINAL\_diabetes.ipynb - Colab

In [1]:



## Moving to random Forest

```
# Few collections import Counter, OrderedDict
# Features = list(data_encoded)
# Features = [x for x in Features if x not in ('Unamed: 0', 'readmitted')]
# X = data_encoded[Features].values
# Y = data_encoded['readmitted'].values
# Xtrain, Xtest, Ytrain, Ytest = train_test_split(X, Y, test_size = .2,random_state = 34, stratify = Y)
# Using our randomforest classifier and giving class weights so that we can even try to handle some imbalanced data
from sklearn.ensemble import RandomForestClassifier, ExtraTreesClassifier
from sklearn.metrics import recall_score
clf_rf = RandomForestClassifier(random_state = 7, class_weight = {0: -1, 1: 0})
model_rf = clf_rf.fit(Xtrain, Ytrain)
score_rf = model_rf.score(Xtest, Ytest)
```

Out[1]:

0.9899992653001249

```
# Confusion Matrix
actual = pd.Series(Ytest, name = 'Actual')
predicted_rf = pd.Series(clf_rf.predict(Xtest), name = 'Predicted')
rf_cm = pd.crosstab(actual, predicted_rf, margins = True)
print(rf_cm)
```

	Predicted	0	1	All
Actual				
0	12377	7	12384	
1	1218	9	1227	
All	13595	16	13611	

```
1: TN_rf = rf_cm.iloc[0,0] / rf_cm.sum().sum()
2: TP_rf = rf_cm.iloc[1,1] / rf_cm.sum().sum()
3: Prev_rf = rf_cm.sum().sum() / rf_cm.sum().sum()

print('Percent of Non-readmissions Detected: ('.format('%.3f' % TN_rf))
print('Percent of Readmissions Detected (Recall): ('.format('%.3f' % TP_rf))
print('Accuracy Among Predictions of Readmitted (Precision): ('.format('%.3f' % Prev_rf))
```

```
# Percent of Non-readmissions Detected: 0.999
Percent of Readmissions Detected (Recall): 0.007
Accuracy Among Predictions of Readmitted (Precision): 0.562
```

```
# From imbalanced_under_sampling import RandomUnderSampler
# From collections import Counter
# Assuming data_encoded, features, and readmitted are already defined
# X = data_encoded[Features].values
# Y = data_encoded['readmitted'].values
# X_random = X
# Y_random = Y
# Random undersampling
# rus = RandomUnderSampler(random_state=34)
# X_res, Y_res = rus.fit_resample(X_random, Y_random) # Corrected method
# print(Counter(Y_res)) # Print the distribution of the undersampled dataset
```

[https://colab.research.google.com/gist/texashchikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul\\_diabetes.ipynb?authuser=1#scrollTo=EY...](https://colab.research.google.com/gist/texashchikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul_diabetes.ipynb?authuser=1#scrollTo=EY...) 73/104

2/4/25, 1:56 AM

McPhaul\_J\_CaseStudy2\_FINAL\_diabetes.ipynb - Colab

```
# TN_rf = rf_cm.sum().sum() / rf_cm.sum().sum()
# TP_rf = rf_cm.sum().sum() / rf_cm.sum().sum()
# Prev_rf = rf_cm.sum().sum() / rf_cm.sum().sum()

print('Percent of Non-readmissions Detected: ('.format('%.3f' % TN_rf))
print('Percent of Readmissions Detected (Recall): ('.format('%.3f' % TP_rf))
print('Accuracy Among Predictions of Readmitted (Precision): ('.format('%.3f' % Prev_rf))
```

```
# Percent of Non-readmissions Detected: 0.939
Percent of Readmissions Detected (Recall): 0.916
Accuracy Among Predictions of Readmitted (Precision): 0.938
```

```
# best number of features
RANDOM_STATE = 123

# num_trees
ensemble_cif = [
    ("RandomForestClassifier", {"n_estimators": "sqrt", "max_features": "sqrt", "random_state": RANDOM_STATE}),
    ("RandomForestClassifier", {"n_estimators": "sqrt", "max_features": "log2", "random_state": RANDOM_STATE}),
    ("RandomForestClassifier", {"n_estimators": "sqrt", "max_features": "log", "random_state": RANDOM_STATE}),
    ("RandomForestClassifier", {"n_estimators": "sqrt", "max_features": "None", "random_state": RANDOM_STATE}),
    ("RandomForestClassifier", {"n_estimators": "sqrt", "max_features": "None", "random_state": RANDOM_STATE})
]
```

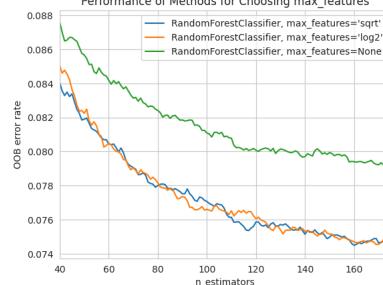
```
# Map classifier name to a list of (n_estimators, error_rate) pairs
error_rate = OrderedDict({label: [] for label, _ in ensemble_cif})
min_estimators = 48
max_estimators = 175

for label, cif in ensemble_cif:
    for i in range(min_estimators, max_estimators + 1):
        error_rate[label].append((i, cif.error_rate))

cif_fit(Xtrain, Ytrain)
oob_error = 1 - cif.oob_error
error_rate[label].append((1, oob_error))

# Record the OOB error for each n_estimators value.
oob_error = 1 - cif.oob_error
error_rate[label].append((1, oob_error))
```

```
# OOB error rate vs. n_estimators plot.
for label, cif in error_rate.items():
    xs, ys = zip(*cif)
    plt.plot(xs, ys, label=label)
    plt.xlabel("n_estimators")
    plt.ylabel("OOB error rate")
    plt.title("Performance of Methods for Choosing max_features")
    plt.legend(loc="upper right")
plt.show()
```



```
# import math
# f = len(list(data_encoded[Features]))
# print(math.log(f, 2))
```

Out[1]:

5.523561956057013

## Final Model

[https://colab.research.google.com/gist/texashchikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul\\_diabetes.ipynb?authuser=1#scrollTo=EY...](https://colab.research.google.com/gist/texashchikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul_diabetes.ipynb?authuser=1#scrollTo=EY...) 75/104

2/4/25, 1:56 AM

McPhaul\_J\_CaseStudy2\_FINAL\_diabetes.ipynb - Colab

In [1]:

```
# TTS
Xtrain, Xtest, Ytrain, Ytest = train_test_split(X_res, Y_res, test_size = .2, random_state = 34, stratify = Y_res)

# random classifier on this undersampled data
rf_rus = RandomForestClassifier(max_features=7)
rf_rus.fit(Xtrain, Ytrain)
print(rf_rus.score(Xtest, Ytest))

# Confusion Matrix
actual = pd.Series(Ytest, name = 'Actual')
predicted_rf_rus = pd.Series(rf_rus.predict(Xtest), name = 'Predicted')
rf_cm_rf_rus = pd.crosstab(actual, predicted_rf_rus, margins = True)
print(rf_cm_rf_rus)

# TN_rf_rus = ct_rf_rus.sum().sum() / ct_rf_rus.sum().sum()
# TP_rf_rus = ct_rf_rus.sum().sum() / ct_rf_rus.sum().sum()
# Prev_rf_rus = ct_rf_rus.sum().sum() / ct_rf_rus.sum().sum()

print('Percent of Non-readmissions Detected: ('.format('%.3f' % TN_rf_rus))
print('Percent of Readmissions Detected (Recall): ('.format('%.3f' % TP_rf_rus))
print('Accuracy Among Predictions of Readmitted (Precision): ('.format('%.3f' % Prev_rf_rus))
```

```
# Percent of Non-readmissions Detected: 0.747
Percent of Readmissions Detected (Recall): 0.761
Accuracy Among Predictions of Readmitted (Precision): 0.750
```

```
# From imbalanced_over_sampling import SMOTE
# From collections import Counter
# Assuming data_encoded, features, and readmitted are already defined
# X = data_encoded[Features].values
# Y = data_encoded['readmitted'].values

# OverSampling the minority class using SMOTE
# m = SMOTE(random_state=34)
# X_resamp, Y_resamp = m.fit_resample(X, Y) # Corrected method
# print(Counter(Y_resamp)) # Print the distribution of the oversampled dataset
```

In [2]:

Counter({1: 61919, 0: 61919})

```
# Xtrain, Xtest, Ytrain, Ytest = train_test_split(X_resamp, Y_resamp, test_size = .2,random_state = 34, stratify = Y_resamp)
```

```
# On oversampled data
clf_rf_sm = RandomForestClassifier(random_state = 7)
model_rf_sm = clf_rf_sm.fit(Xtrain, Ytrain)
print(model_rf_sm.score(Xtest, Ytest))
```

Out[2]:

0.9276485788113695

```
# F1 score
# From sklearn.metrics import f1_score
f1_score_rf_sm = f1_score(Ytest, model_rf_sm.predict(Xtest), average='macro')
```

Out[3]:

0.3325874611735387

```
# From sklearn.metrics import f1_score
f1_score_rf_sm = f1_score(Ytest, model_rf_sm.predict(Xtest), average='micro')
```

Out[4]:

0.49709793549335096

```
# Confusion Matrix
actual = pd.Series(Ytest, name = 'Actual')
predicted_rf_sm = pd.Series(clf_rf_sm.predict(Xtest), name = 'Predicted')
ct_rf_sm = pd.crosstab(actual, predicted_rf_sm, margins = True)
print(ct_rf_sm)
```

In [5]:

Counter({0: 11634, 1: 758, All: 12384})

In [6]:

Counter({0: 1042, 1: 11342, All: 12384})

In [7]:

Counter({0: 12676, 1: 12092, All: 24768})

[https://colab.research.google.com/gist/texashchikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul\\_diabetes.ipynb?authuser=1#scrollTo=EY...](https://colab.research.google.com/gist/texashchikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul_diabetes.ipynb?authuser=1#scrollTo=EY...) 74/104

2/4/25, 1:56 AM

McPhaul\_J\_CaseStudy2\_FINAL\_diabetes.ipynb - Colab

```
# Fit Model
model_f1 = RandomForestClassifier(random_state = 7, n_estimators = 85, max_features = 'log2', max_depth = 7)
f1_cm = model_f1.fit(Xtrain, Ytrain)
print(f1_cm.score(Xtest, Ytest))
```

Out[1]:

0.7848433462532299

```
# Confusion Matrix
actual_f1 = pd.Series(Ytest, name = 'Actual')
predicted_f1 = pd.Series(f1_cm.predict(Xtest), name = 'Predicted')
ct_f1 = pd.crosstab(actual_f1, predicted_f1, margins = True)
print(ct_f1)
```

In [2]:

Counter({0: 10192, 1: 2192, All: 12384})

In [3]:

Counter({0: 3137, 1: 9247, All: 12384})

In [4]:

Counter({0: 13329, 1: 11439, All: 24768})

```
# TN_f1 = ct_f1.sum().sum() / ct_f1.sum().sum()
# TP_f1 = ct_f1.sum().sum() / ct_f1.sum().sum()
# Prev_f1 = ct_f1.sum().sum() / ct_f1.sum().sum()

print('Percent of Non-readmissions Detected: ('.format('%.3f' % TN_f1))
print('Percent of Readmissions Detected (Recall): ('.format('%.3f' % TP_f1))
print('Accuracy Among Predictions of Readmitted (Precision): ('.format('%.3f' % Prev_f1))
```

```
# Percent of Non-readmissions Detected: 0.823
Percent of Readmissions Detected (Recall): 0.747
Accuracy Among Predictions of Readmitted (Precision): 0.808
```

```
# Importances = f1_cm.feature_importances_
importance_df = pd.DataFrame(importances, columns = ['feature', 'importance'])
print(importance_df.sort_values('importance', ascending = False))
```

In [5]:

importance\_df.head(40)

feature	importance
40	num_visits
38	change
1	gender
15	metformin
32	insulin
21	glipizide
24	pioglitazone
22	glyburide
41	first_diag
4	discharge_disposition_id

In [6]:

feature	importance
23	tolbutamide
37	metformin_plaque
39	esamode
31	citoglipiton
36	metformin_rosiglitazone
35	glimepiride_pioglitazone
44	HbA1c

In [7]:

Next steps: [Generate code with Jupyter](#) [View recommended plots](#) [New interactive sheet](#)

```
1 print(importance_df.importance == 0))
```

In [8]:

feature	importance
23	tolbutamide
37	metformin_plaque
39	esamode
31	citoglipiton
36	metformin_rosiglitazone
35	glimepiride_pioglitazone
44	HbA1c

In [9]:

Checking Validation

```
1: features = list(data_encoded)
2: features = [x for x in Features if x not in ('Unamed: 0', 'readmitted')]
```

```
1: X = data_encoded[Features].values
2: Y = data_encoded['readmitted'].values
```

[https://colab.research.google.com/gist/texashchikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul\\_diabetes.ipynb?authuser=1#scrollTo=EY...](https://colab.research.google.com/gist/texashchikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul_diabetes.ipynb?authuser=1#scrollTo=EY...) 75/104

[https://colab.research.google.com/gist/texashchikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul\\_diabetes.ipynb?authuser=1#scrollTo=EY...](https://colab.research.google.com/gist/texashchikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul_diabetes.ipynb?authuser=1#scrollTo=EY...) 76/104

2/4/25, 1:56 AM

McPhaul\_J\_CaseStudy2\_FINAL\_diabetes.ipynb - Colab

```

1 from imblearn.over_sampling import RandomOverSampler
2 from imblearn.model_selection import train_test_split
3 from sklearn.ensemble import RandomForestClassifier
4 from collections import Counter
5 import pandas as pd
6
7 number_of_repetitions = 10 # number of trials
8
9 # Declares empty lists for true-positive and true-negative rates
10 TPR = []
11 TNR = []
12
13 # For loop for multiple trials
14 for trial in range(number_of_repetitions):
15     # Random undersampling using fit_resample
16     # x, y = X_train, Y_train
17     # X_resamp, Y_resamp = rus.fit_resample(X, y) # use fit_resample
18     # print(Counter(Y_resamp))
19
20     # Train, test, split
21     Xtrain, Xtest, Ytrain, Ytest = train_test_split(
22         X, y, test_size=0.2, random_state=3 * trial, stratify=Y_resamp
23     )
24
25     # Random Forest model
26     rf_rus = RandomForestClassifier(
27         random_state=7, n_estimators=65, max_features='log2', max_depth=7
28     )
29
30     rf_model_rus = rf_rus.fit(Xtrain, Ytrain)
31     print("rf_rus score(Xtest, Ytest)")
32
33     # Confusion matrix
34     actual = pd.Series(y, name='Actual')
35     predicted_rf_rus = pd.Series(rf_rus.predict(Xtest), name='Predicted')
36     ct_rf_rus = pd.crosstab(actual, predicted_rf_rus, margins=True)
37     print(ct_rf_rus)
38
39     # True negative rate
40     tnr = ct_rf_rus.iat[0, 0] / ct_rf_rus.iat[0, 2]
41
42     # True positive rate
43     tpr = ct_rf_rus.iat[1, 1] / ct_rf_rus.iat[1, 2]
44     TPR.append(tpr)
45
46     # Output metrics
47     print("Accuracy for not readmitted: (%.2f%%)" % tpr)
48     print("Accuracy for readmitted (Recall): (%.2f%%)" % tpr)
49     print("Random Forest trial count: (%d)" % trial + 1)
50     print("-" * 50)
51
52 All 1334 894 1227
All 1258 1197 2455
Accuracy for not readmitted: 0.753
Accuracy for readmitted (Recall): 0.729
Random Forest trial count: 4

```

2/4/25, 1:56 AM

McPhaul\_J\_CaseStudy2\_FINAL\_diabetes.ipynb - Colab

```

0 922 305 1227
1 332 896 1228
All 1254 1201 2455
Accuracy for not readmitted: 0.751
Accuracy for readmitted (Recall): 0.738
Random Forest trial count: 8

```

```

Counter({0: 6136, 1: 6136})
0.7437881873727088
Predicted 0 1 All
Actual
0 950 278 1228
1 351 876 1227
All 1301 1154 2455
Accuracy for not readmitted: 0.774

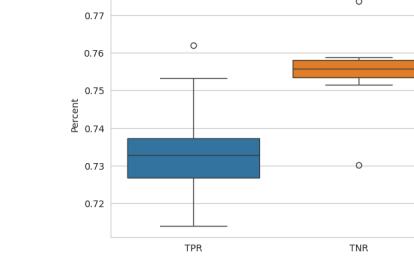
```

```

# plotting TPR and TNR
1 plots = pd.DataFrame(({"TPR": TPR, "TNR": TNR}))
2 plots["Percent"] = plots["TPR"] + plots["TNR"]
3 plots["Percent"] = plots["Percent"] / 2
4 plt.title("Box Plots for TPR and TNR in Random Undersampling (Random Forest)")
5 plt.xlabel("Percent")
6 plt.ylabel("Percent")
7 plt.show()

```

Box Plots for TPR and TNR in Random Undersampling (Random Forest)



```

https://colab.research.google.com/gist/texashchikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul_diabetes.ipynb?authuser=1#scrollTo=EY... 77/104

```

2/4/25, 1:56 AM

McPhaul\_J\_CaseStudy2\_FINAL\_diabetes.ipynb - Colab

```

43 tpr_ue = ct_rf_ue.iat[1, 1] / ct_rf_ue.iat[1, 2]
44 TPR_ue = []
45
46 # Select metrics
47 print("Accuracy for not readmitted: (%.2f%%)" % tpr_ue)
48 print("Accuracy for readmitted (Recall): (%.2f%%)" % tpr_ue)
49 print("Random Forest trial count: (%d)" % trial + 1)
50 print("-" * 50)
51
52 All 13349 11419 24768
All 1258 1197 2455
Accuracy for not readmitted: 0.819
Accuracy for readmitted (Recall): 0.741
Random Forest trial count: 4

```

```

Counter({0: 61919, 1: 61919})
0.7881944444444444
Predicted 0 1 All
Actual
0 10259 2125 12384
1 3121 9263 12384
All 13380 11388 24768
Accuracy for not readmitted: 0.828
Accuracy for readmitted (Recall): 0.748
Random Forest trial count: 5

```

```

Counter({0: 61919, 1: 61919})
0.7865794575643411
Predicted 0 1 All
Actual
0 10113 2271 12384
1 3015 9369 12384
All 13128 11648 24768
Accuracy for not readmitted: 0.817
Accuracy for readmitted (Recall): 0.757
Random Forest trial count: 6

```

```

Counter({0: 61919, 1: 61919})
0.7828246124031008
Predicted 0 1 All
Actual
0 10154 2238 12384
1 3149 9235 12384
All 13383 11465 24768
Accuracy for not readmitted: 0.828
Accuracy for readmitted (Recall): 0.746
Random Forest trial count: 7

```

```

Counter({0: 61919, 1: 61919})
0.787467700258398
Predicted 0 1 All
Actual
0 10074 2310 12384
1 2954 9430 12384
All 13028 11748 24768
Accuracy for not readmitted: 0.813
Accuracy for readmitted (Recall): 0.761
Random Forest trial count: 8

```

```

Counter({0: 61919, 1: 61919})
0.7971172480626154
Predicted 0 1 All
Actual
0 10365 2019 12384
1 3006 9378 12384
All 13371 11397 24768
Accuracy for not readmitted: 0.837
Accuracy for readmitted (Recall): 0.757

```

```

plots = pd.DataFrame(({"TPR": TPR, "TNR": TNR}))
2 plots["Percent"] = plots["TPR"] + plots["TNR"]
3 plots["Percent"] = plots["Percent"] / 2
4 plt.title("Box Plots for TPR and TNR in Random Undersampling (Random Forest)")
5 plt.xlabel("Percent")
6 plt.ylabel("Percent")
7 plt.show()

```

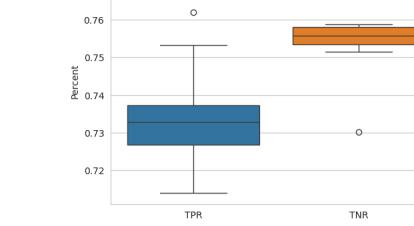
2/4/25, 1:56 AM

McPhaul\_J\_CaseStudy2\_FINAL\_diabetes.ipynb - Colab

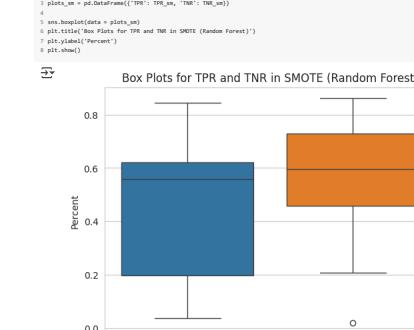
```

# BoxPlot
1 plots_ue = pd.DataFrame(({"TPR": TPR_ue, "TNR": TNR_ue}))
2 plots_ue["Percent"] = plots_ue["TPR"] + plots_ue["TNR"]
3 plots_ue["Percent"] = plots_ue["Percent"] / 2
4 plt.title("Box Plots for TPR and TNR in SMOTE (Random Forest)")
5 plt.xlabel("Percent")
6 plt.ylabel("Percent")
7 plt.show()

```



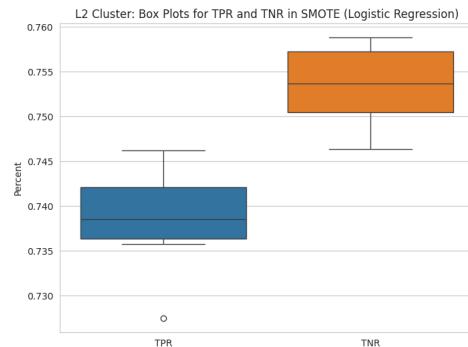
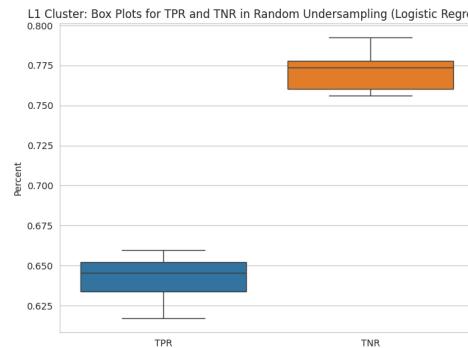
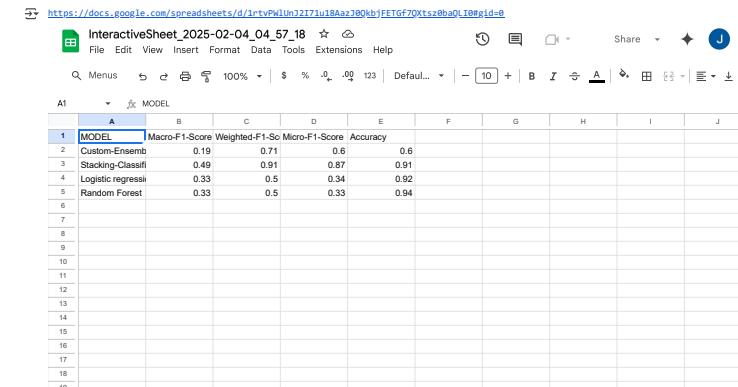
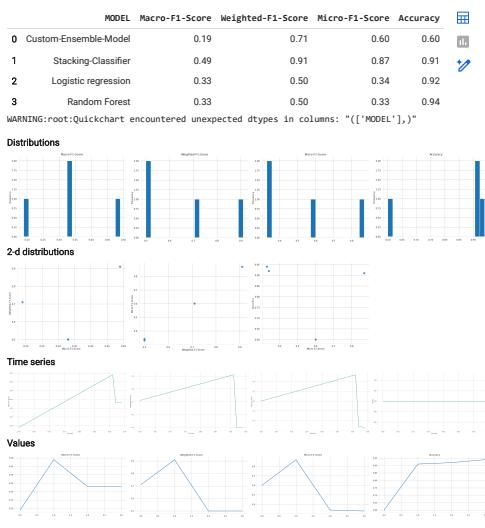
Box Plots for TPR and TNR in SMOTE (Random Forest)



MODEL	Accuracy for train data for being readmitted	Accuracy for train data for non-readmitted	Accuracy for test data for being readmitted	Accuracy for test data for non-readmitted
Logistic	0.515	0.838	0.42	0.857
Custom-Ensemble-Model				
Stacking-Classifier				
Random Forest				
Macro-F1-Score	[0.19, 0.49, 0.33, 0.31]			
Weighted-F1-Score		[0.71, 0.81, 0.59, 0.5]		
Micro-F1-Score				

https://colab.research.google.com/gist/texashchikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul\_diabetes.ipynb?authuser=1#scrollTo=EY... 79/104

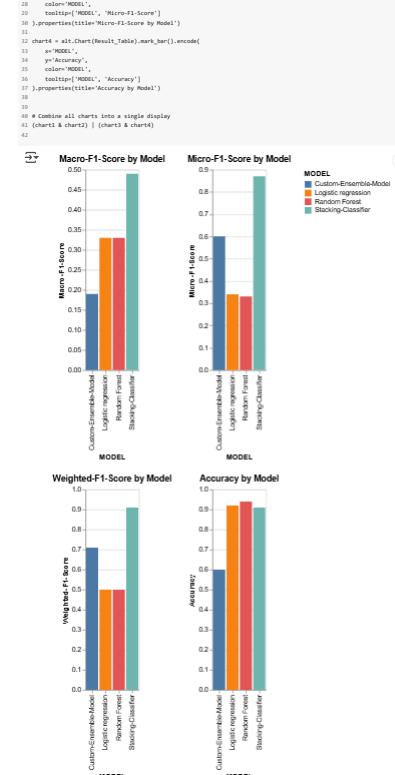
https://colab.research.google.com/gist/texashchikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul\_diabetes.ipynb?authuser=1#scrollTo=EY... 80/104



```

1 # Using dataframe Result_Table; suggest a plot
2
3 import altair as alt
4
5 # Convert the 'MODEL' column to a categorical type for proper ordering in the plot
6 Result_Table['MODEL'] = Result_Table['MODEL'].astype('category')
7
8 # Create a chart for each metric
9 chart1 = alt.Chart(Result_Table).mark_bar().encode(
10   x='MODEL',
11   y='Macro-F1-Score',
12   color='MODEL',
13   tooltip=['MODEL', 'Macro-F1-Score']
14 ).properties(title='Macro-F1-Score by Model')
15
16 chart2 = alt.Chart(Result_Table).mark_bar().encode(
17   x='MODEL',
18   y='Weighted-F1-Score',
19   color='MODEL',
20   tooltip=['MODEL', 'Weighted-F1-Score']
21 ).properties(title='Weighted-F1-Score by Model')
22
23 chart3 = alt.Chart(Result_Table).mark_bar().encode(
24   x='MODEL',
25   y='Micro-F1-Score',
26   color='MODEL',
27   tooltip=['MODEL', 'Micro-F1-Score']
28 ).properties(title='Micro-F1-Score by Model')
29
30 chart4 = alt.Chart(Result_Table).mark_bar().encode(
31   x='MODEL',
32   y='Accuracy',
33   color='MODEL',
34   tooltip=['MODEL', 'Accuracy']
35 ).properties(title='Accuracy by Model')
36
37 # Combine all charts into a single display
38 (chart1 & chart2) | (chart3 & chart4)
39
40
41
42

```

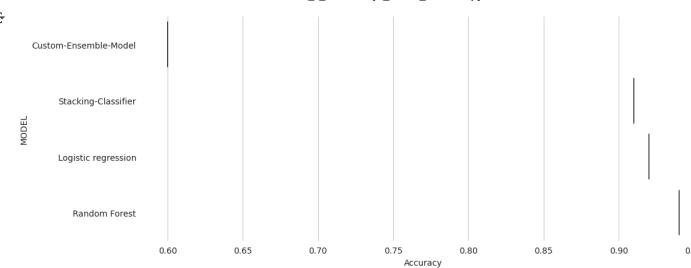


&gt; MODEL vs Accuracy

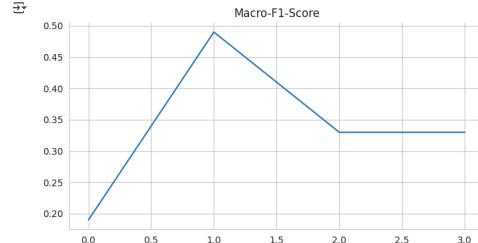
Show code

2/4/25, 1:56 AM

McPhaul\_J\_CaseStudy2\_FINAL\_diabetes.ipynb - Colab



## ➤ Macro-F1-Score

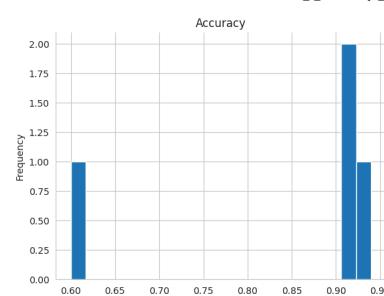
[Show code](#)

## ➤ Accuracy

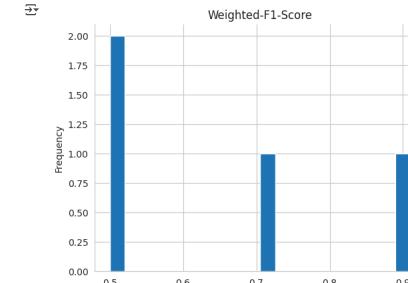
[Show code](#)

2/4/25, 1:56 AM

McPhaul\_J\_CaseStudy2\_FINAL\_diabetes.ipynb - Colab



## ➤ Weighted-F1-Score

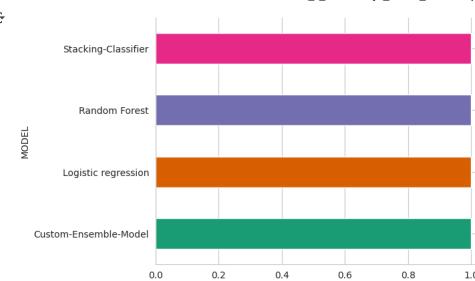
[Show code](#)

## ➤ MODEL

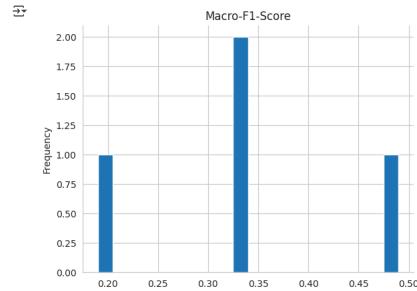
[Show code](#)[https://colab.research.google.com/gist/texashikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul\\_diabetes.ipynb?authuser=1#scrollTo=EY...](https://colab.research.google.com/gist/texashikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul_diabetes.ipynb?authuser=1#scrollTo=EY...) 85/104

2/4/25, 1:56 AM

McPhaul\_J\_CaseStudy2\_FINAL\_diabetes.ipynb - Colab



## ➤ Macro-F1-Score

[Show code](#)

```

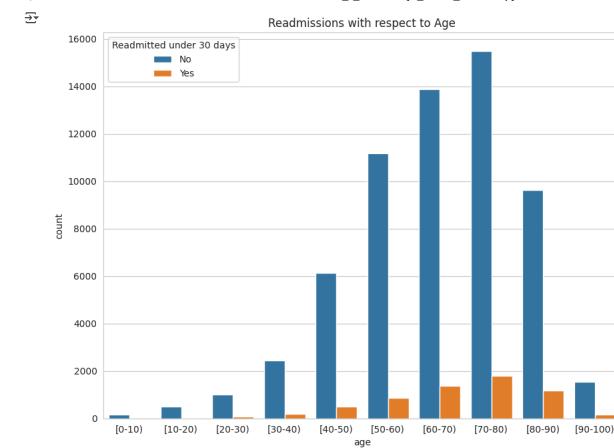
1: plot = sns.countplot(x='age', hue='readmitted', data=data, order=sorted(data['age'].unique()))
2: plot.figure.set_size_inches(10, 7.5)
3: plot.legend.title('Readmitted under 30 days')
4: plot.legend().set_title('Readmitted under 30 days', title_fontsize=10, title_y=1.05)
5: plot.legend().set_label('No', title_y=1.05)
6: plot.legend().set_label('Yes', title_y=1.05)
7: plt.show()

```

[https://colab.research.google.com/gist/texashikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul\\_diabetes.ipynb?authuser=1#scrollTo=EY...](https://colab.research.google.com/gist/texashikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul_diabetes.ipynb?authuser=1#scrollTo=EY...) 86/104

2/4/25, 1:56 AM

McPhaul\_J\_CaseStudy2\_FINAL\_diabetes.ipynb - Colab



## ▼ Notebook Analysis

1. Introduction - Likely contains an overview of the case study and dataset.
2. MODELING - Begins the modeling process.
3. Ensemble with Classifier (stacking) - Applies an ensemble learning approach.
4. Logistic Regression - Implements and analyzes logistic regression.
5. Undersampling - Uses undersampling techniques for class balancing.
6. Train Test Split - Splits the dataset into training and test sets.
7. Grid Search CV using L2 reg w/ 5-fold cv - Performs hyperparameter tuning w/ GridSearchCV and L2 regularization.
8. LR model w/ undersampling - Confusion Matrix - Evaluates the logistic 10. regression model with undersampling using a confusion matrix.
9. SMOTE for oversampling - Applies Synthetic Minority Over-sampling Technique 1 (SMOTE) for class balancing.
10. Plotting TNR & TPR - Visualizes True Negative Rate (TNR) and True Positive Rate (TPR).
11. Moving to random Forest - Transitions to a Random Forest model.
12. Final Model - Identifies and evaluates the final model.
13. Checking Validation - Validates the final model.

## ▼ 1. Analysis of the Introduction Section: Diabetes Dataset Exploration and Preprocessing

This section prepares the diabetes dataset for modeling through data cleaning, feature engineering, and exploratory visualization.

## Data Cleaning:

- Missing values (?) are replaced with NaN.
- Columns with excessive missing data (weight, medical\_specialty, payer\_code) are removed.
- Patients with specific dischargeDisposition\_id values are filtered out.

## Feature Engineering:

- Diagnoses are categorized into groups (circulatory, respiratory, digestive, diabetes, injury, other).
- readmitted is transformed into a binary variable (1 < 30 days, 0, otherwise).

[https://colab.research.google.com/gist/texashikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul\\_diabetes.ipynb?authuser=1#scrollTo=EY...](https://colab.research.google.com/gist/texashikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul_diabetes.ipynb?authuser=1#scrollTo=EY...) 87/104[https://colab.research.google.com/gist/texashikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul\\_diabetes.ipynb?authuser=1#scrollTo=EY...](https://colab.research.google.com/gist/texashikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul_diabetes.ipynb?authuser=1#scrollTo=EY...) 88/104

- Number of visits per patient is calculated.

**Data Visualization:**

```
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt

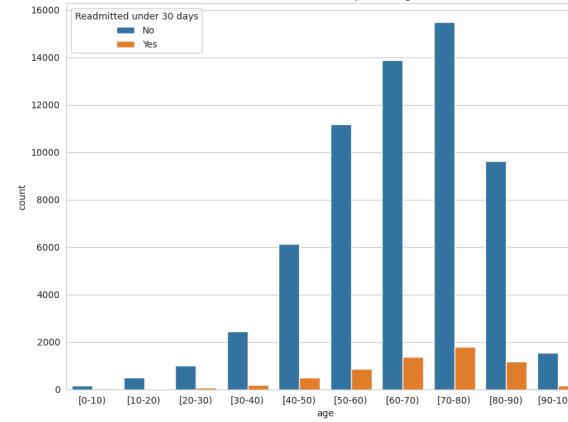
data = pd.read_csv('/content/drive/MyDrive/Colab_Study_2/diabetic_data.csv')
data['readmitted'] = data['readmitted'].replace('>30', 0)
data['readmitted'] = data['readmitted'].replace('NO', 0)
data['readmitted'] = data['readmitted'].replace('<30', 1)

def cat_col(col):
    if (col >= 398) & (col <= 459) | (col == 785):
        return 'circulatory'
    elif (col >= 460) & (col <= 519) | (col == 786):
        return 'respiratory'
    elif (col >= 520) & (col <= 579) | (col == 787):
        return 'digestive'
    elif (col >= 250.00) & (col <= 250.99):
        return 'diabetes'
    elif (col >= 800) & (col <= 999):
        return 'injury'
    else:
        return 'other'

data['first_diag'] = data.Diag1.apply(lambda col: cat_col(col))
data['second_diag'] = data.Diag2.apply(lambda col: cat_col(col))
data['third_diag'] = data.Diag3.apply(lambda col: cat_col(col))

plot = sns.countplot(x='age', hue='readmitted', data=data, order=sorted(data['age'].unique()))
plot = sns.countplot(x='age', hue='readmitted', data=data, order=sorted(data['age'].unique()))
plot.figure.set_size_inches(10, 7.5)
plot.legend(title='Readmitted under 30 days', labels=['No', 'Yes'])
plot.axes.set_title('Readmissions with respect to Age')
plt.show()
```

This section effectively prepares the data for modeling. The visualizations provide insights into data characteristics and potential predictors of readmission. A good ROC AUC score (above 0.8) is desirable, while an AUC close to 0.5 would suggest random performance.

**Readmissions with respect to Age****Plot analysis:**

- Hospital readmission rates for diabetic patients increase with age, peaking between 70-80 years old, and then declining slightly.
- While the 80-90 age group has a high number of overall hospital visits, the decrease in readmissions may be attributed to mortality, more intensive initial care, or increased use of long-term care facilities.
- Readmissions are significantly lower among patients under 40.
- Intervention programs targeting patients 40 and older, particularly those between 50-80, focusing on preventative care and enhanced post-surgery support, could reduce readmissions.
- Strengthening home-based care for the oldest patients (80+) may further decrease hospital dependency.

**2. Modeling Approach for Diabetes Readmission Prediction: A Stacked EnsembleMethod**

This section describes a stacked ensemble approach using Logistic Regression as base models and an Extra Trees Classifier as a meta-learner.

**1. Data Splitting:**

Stratified sampling ensures balanced class distributions in training and test sets. An additional split within the training data is likely for cross-validation or ensemble training.

```
from sklearn.model_selection import train_test_split
from sklearn.model_selection import train_val_split

X_train, X_test, Y_train, Y_test = train_test_split(X, y, test_size=0.2, random_state=7, stratify=y)
X_train, X_test, Y_train, Ytest = train_val_split(X_train, Y_train, test_size=0.5)
```

**2. Synthetic Sample Generation:**

A bootstrapping technique generates synthetic samples by randomly selecting and duplicating rows and columns from the training data to augment it and potentially improve model robustness.

[https://colab.research.google.com/gist/texashikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul\\_diabetes.ipynb?authuser=1#scrollTo=EY...](https://colab.research.google.com/gist/texashikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul_diabetes.ipynb?authuser=1#scrollTo=EY...) 89/104

```
def generating_sample(X_train1, Ytrain1):
    Selecting_row = np.argsort(np.random.choice(X_train1.shape[0], 8166, replace=True))
    Replacing_row = np.argsort(np.random.choice(Selecting_row, 5444, replace=True))
    Selecting_column = np.argsort(np.random.choice(X_train1.shape[1], int(X_train1.shape[1] * 0.64), replace=True))

    sample_data = X_train1[Selecting_row[:, None], Selecting_column]
    target_of_sample_data = Ytrain1[Selecting_row[:, None]]

    replicated_data = X_train1[Replacing_row[:, None], Selecting_column]
    target_of_replicated_data = Ytrain1[Replacing_row[:, None]]

    final_sample_data = np.vstack((sample_data, replicated_data))
    final_target_data = np.vstack((target_of_sample_data.reshape(-1, 1), target_of_replicated_data.reshape(-1, 1)))

    return final_sample_data, final_target_data, Selecting_row, Selecting_column
```

**3. Hyperparameter Tuning:**

GridSearchCV with 5-fold cross-validation optimizes the L2 regularization strength (C) for Logistic Regression models. Class weights address class imbalance.

```
from sklearn.model_selection import GridSearchCV
from sklearn.linear_model import LogisticRegression
C_grid = ([0.0001, 0.001, 0.01, 0.1, 1, 10, 100, 1000])
weights = ({0: 1, 1: .9})

clf_grid = GridSearchCV(LogisticRegression(penalty='l2', class_weight=weights), C_grid, cv=5, scoring='accuracy')
clf_grid.fit(list_input_data[i], list_output_data[i])
```

**4. Base Model Training:**

Thirty Logistic Regression models are trained with the optimal hyperparameter C, aiming to capture diverse data patterns for the stacking approach.

```
all_selected_models = []
for i in range(30):
    model = LogisticRegression(C=cif_grid.best_params_['C'], penalty='l2', class_weight=weights)
    model.fit(list_input_data[i], list_output_data[i])
    all_selected_models.append(model)
```

**5. Stacking with Meta-Learner:**

Predictions from the base models form meta-features, used to train an Extra Trees Classifier as the meta-learner, enabling it to learn from and correct errors of individual base models.

```
from sklearn.ensemble import ExtraTreesClassifier
D_meta = []
for i in range(30):
    y_pred = all_selected_models[i].predict(list_input_data[i])
    D_meta.append(y_pred)

from sklearn.ensemble import ExtraTreesClassifier
meta_model = ExtraTreesClassifier()
meta_model.fit(D_meta, list_output_data_final)
```

**6. Final Testing and Evaluation:**

The stacked model's performance is evaluated on unseen test data using accuracy and F1-score (macro, micro, and weighted averages).

```
from sklearn.metrics import accuracy_score, f1_score
pred_model = meta_model.predict(D_meta_2)
accuracy_score(np.argmax(pred_model, axis=1), np.argmax(list_output_data_final_test, axis=1))
f1_score(np.argmax(pred_model, axis=1), np.argmax(list_output_data_final_test, axis=1), average='macro')
```

```
fi1_score(np.argmax(pred_model, axis=1), np.argmax(list_output_data_final_test, axis=1), average='weighted')
fi1_score(np.argmax(pred_model, axis=1), np.argmax(list_output_data_final_test, axis=1), average='micro')
```

**Summary:**

This stacked ensemble approach combines multiple Logistic Regression models using an Extra Trees meta-learner. It incorporates synthetic data generation, hyperparameter tuning, and a robust evaluation strategy. A good ROC AUC score (above 0.8) indicates strong model performance. An AUC near 0.5 suggests random performance.

**3. Ensemble Modeling with Stacking for Diabetes Readmission Prediction|**

This section details the implementation of a stacking classifier, combining multiple base models with a meta-classifier to improve predictive performance.

**1. Data Splitting:**

The dataset is split into training and testing sets using stratified sampling to maintain class balance.

```
from sklearn.model_selection import train_test_split

X = data_encoded.drop('readmitted', axis=1)
y = data_encoded.readmitted

X_train, X_test, Y_train, Y_test = train_test_split(X, y, test_size=0.2, random_state=7, stratify=y)
```

**2. Defining Base Models:**

A diverse set of base models is used:

- K-Nearest Neighbors (KNN)
- Random Forest
- Extra Trees Classifier
- Gaussian Naive Bayes
- Logistic Regression

Random Forest serves as the meta-classifier, aggregating predictions from these base models.

```
from sklearn.linear_model import LogisticRegression
from sklearn.neighbors import KNeighborsClassifier
from sklearn.naive_bayes import GaussianNB
from sklearn.ensemble import RandomForestClassifier, ExtraTreesClassifier
from sklearn.mixture import StackingClassifier
import warnings

warnings.simplefilter('ignore')

clf1 = KNeighborsClassifier(n_neighbors=5)
clf2 = RandomForestClassifier(random_state=5)
clf3 = ExtraTreesClassifier()
clf4 = GaussianNB()
clf5 = LogisticRegression(penalty='l2')

meta_classifier = RandomForestClassifier(random_state=7)

scclf = StackingClassifier(classifiers=[clf1, clf2, clf3, c14, c15], meta_classifier=meta_classifier)
```

**3. Cross-Validation:**

3-fold cross-validation evaluates the performance of individual base models and the stacked ensemble, providing insights into their generalization ability.

```
print('3-fold cross validation:\n')

for clf, label in zip([clf1, clf2, clf3, c14, c15, scclf],
                     ['KNN', 'Random Forest', 'ExtraTreesClassifier',
                      'GaussianNB', 'Logistic Regression', 'StackingClassifier'])
```

2/4/25, 1:56 AM

McPhaul\_J\_CaseStudy2\_FINAL\_diabetes.ipynb - Colab

```
'GaussianNB', 'Logistic Regression', 'StackingClassifier']:

scores = model_selection.cross_val_score(clf, X_train, Y_train, cv=3, scoring='accuracy')
print('Accuracy: %0.2f [%s]' % (scores.mean(), label))
```

#### 4. Stacking Classifier Training:

The stacking classifier, combining the base models and the meta-classifier, is trained on the entire training dataset.

```
scif.fit(X_train, Y_train)
```

#### 5. Model Saving:

The trained stacking classifier is saved for later reuse without retraining.

```
import pickle
file = open('stacking_classifier_model_final_last.pkl', 'wb')
pickle.dump(scif, file)
```

#### 6. Prediction and Evaluation:

Predictions are made on the test set, and performance is assessed using macro, micro, and weighted F1-scores, providing a comprehensive evaluation across different aspects of classification performance.

```
y_pred = scif.predict(X_test)
```

#### 7. Performance Evaluation (F1)

```
from sklearn.metrics import f1_score

f1_score(Y_test, y_pred, average='macro')
f1_score(Y_test, y_pred, average='micro')
f1_score(Y_test, y_pred, average='weighted')
```

#### Summary:

This stacking ensemble approach leverages the strengths of diverse base models, combined through a Random Forest meta-classifier. Cross-validation and a robust evaluation strategy using F1-scores provide a comprehensive assessment of the model's performance in predicting diabetes readmissions. The expectation is that the stacking classifier outperforms individual base models, demonstrating the effectiveness of the ensemble approach.

### 4. Logistic Regression Analysis for Diabetes Readmission Prediction

This analysis uses Logistic Regression to predict diabetes readmission, focusing on hyperparameter tuning, model evaluation, and interpretation of results.

#### 1. Data Preparation and Splitting:

The dataset is split into 80% training and 20% testing sets using stratified sampling to maintain class distribution.

```
from sklearn.model_selection import train_test_split
from sklearn.model_selection import train_test_split

features = list(data_encoded)
features = [x for x in features if x not in ('Unnamed: 0', 'readmitted')]

X = data_encoded[features].values
y = data.readmitted.values

Xtrain, Xtest, Ytrain, Ytest = train_test_split(X, y, test_size=0.2, random_state=7, stratify=y)
Xtrain, Xtest, Ytrain, Ytest = train_test_split(X, y, test_size=0.2, random_state=7, stratify=y)
```

#### 2. Hyperparameter Tuning:

[https://colab.research.google.com/gist/texashchikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul\\_diabetes.ipynb?authuser=1#scrollTo=EY...](https://colab.research.google.com/gist/texashchikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul_diabetes.ipynb?authuser=1#scrollTo=EY...) 93/104

2/4/25, 1:56 AM

McPhaul\_J\_CaseStudy2\_FINAL\_diabetes.ipynb - Colab

```
predict_train = pd.Series(x_pred_train, name='Predicted')

train_ct = pd.crosstab(actual_train, predict_train, margins=True)
print(train_ct)

TN_train = train_ct.iloc[0, 0] / train_ct.iloc[0, 2] # True Negatives Rate
TP_train = train_ct.iloc[1, 1] / train_ct.iloc[1, 2] # True Positives Rate

print('Training accuracy for not readmitted: {}'.format('%.3f' % TN_train))
print('Training accuracy for being readmitted: {}'.format('%.3f' % TP_train))

actual_test = pd.Series(Ytest, name='Actual')
predict_test = pd.Series(x_pred_test, name='Predicted')

test_ct = pd.crosstab(actual_test, predict_test, margins=True)
print(test_ct)

TN_test = test_ct.iloc[0, 0] / test_ct.iloc[0, 2] # True Negatives Rate
TP_test = test_ct.iloc[1, 1] / test_ct.iloc[1, 2] # True Positives Rate

print('Test accuracy for not readmitted: {}'.format('%.3f' % TN_test))
print('Test accuracy for readmitted (Recall): {}'.format('%.3f' % TP_test))
```

- High TN (True Negative) Rate: Model predicts "not readmitted" cases well.
- Lower TP (True Positive) Rate: Model struggles with predicting "readmitted" cases.
- If TP rate is low, the model may need oversampling (SMOTE) or better class balancing.

#### Summary:

The Logistic Regression model demonstrates reasonable predictive capability. However, the lower true positive rate suggests a need for improved readmission prediction. Oversampling techniques like SMOTE or other balancing methods could be explored to address this.

```
1 print('3-fold cross validation:\n')
2 for clf, label in zip([clf1, clf2, clf3, cl4, cl5, scif]):
3     # XNN, 'Random Forest', 'ExtraTreeClassifier',
4     # 'GaussianNB', 'Logistic Regression', 'StackingClassifier']:
5
6     scores = model_selection.cross_val_score(clf, X_train, Y_train, cv=3, scoring='accuracy')
7     print('Accuracy: %0.2f [%s]' % (scores.mean(), label))
```

#### 5. Addressing Class Imbalance with Random Undersampling

This section describes the application of random undersampling to balance the readmitted vs. not-readmitted classes in the diabetes dataset.

#### 1. Identifying Class Imbalance:

The dataset exhibits class imbalance, with the majority class (not readmitted) significantly outnumbering the minority class (readmitted). Features are extracted for the undersampling process.

```
features = list(data_encoded)
features = [x for x in features if x not in ('Unnamed: 0', 'readmitted')]
```

#### 2. Applying Random Undersampling:

Random Undersampling (RUS) reduces the majority class size to match the minority class size, creating a balanced dataset.

```
from collections import Counter
from imblearn.under_sampling import RandomUnderSampler
```

2/4/25, 1:56 AM

McPhaul\_J\_CaseStudy2\_FINAL\_diabetes.ipynb - Colab

GridSearchCV with 5-fold cross-validation is employed to find the optimal regularization strength (C) for L2 regularization (Ridge), addressing potential overfitting. Class weights are adjusted to account for class imbalance.

```
from sklearn.model_selection import GridSearchCV
from sklearn.linear_model import LogisticRegression

C_grid = ['C': [0.0001, 0.001, 0.01, 0.1, 1, 10, 100, 1000]}
weights = {0: 1, 1: .9}

clf_grid = GridSearchCV(LogisticRegression(penalty='l2', class_weight=weights), C_grid, cv=5, scoring='accuracy')
clf_grid.fit(Xtrain, Ytrain)
print(clf_grid.best_params_, clf_grid.best_score_)
```

#### 3. Model Training and Evaluation:

The best model, determined by GridSearchCV, is trained on the entire training set. Predictions are made on both training and testing sets, and accuracy is assessed. A classification report (including precision, recall, and F1-score) provides a comprehensive performance overview.

```
clf_grid_best = LogisticRegression(C=clf_grid.best_params_['C'], penalty='l2', class_weight=weights)
clf_grid_best.fit(Xtrain, Ytrain)
from sklearn.metrics import accuracy_score, classification_report
x_pred_train = clf_grid_best.predict(Xtrain)
x_pred_test = clf_grid_best.predict(Xtest)
from sklearn.metrics import accuracy_score

accuracy_score(x_pred_train, Ytrain) # Train Accuracy
accuracy_score(x_pred_test, Ytest) # Test Accuracy
from sklearn.metrics import classification_report

report_train = classification_report(Ytrain, x_pred_train)
report_test = classification_report(Ytest, x_pred_test)

print(report_train) # Training Report
print(report_test) # Testing Report
```

#### 4. ROC-AUC Analysis:

The ROC-AUC score and curve are used to evaluate model discrimination. An AUC > 0.80 is desirable.

```
from sklearn.metrics import roc_auc_score, roc_curve, auc
import matplotlib.pyplot as plt

from sklearn.metrics import roc_auc_score

probability_train = clf_grid_best.predict_proba(Xtrain)[:, 1]
probability_test = clf_grid_best.predict_proba(Xtest)[:, 1]

roc_auc_train = roc_auc_score(Ytrain, probability_train)
roc_auc_test = roc_auc_score(Ytest, probability_test)

print(roc_auc_train, roc_auc_test)
```

- AUC > 0.80 indicates a good model.
- If the ROC curve is close to the diagonal (AUC ~ 0.5), the model is performing randomly.

#### 5. Confusion Matrix Interpretation:

Confusion matrices for both training and testing sets reveal the model's performance on predicting readmitted vs. not readmitted cases. True negative and true positive rates are calculated. The analysis indicates that the model predicts non-readmission more accurately than readmission.

```
import pandas as pd

actual_train = pd.Series(Ytrain, name='Actual')

https://colab.research.google.com/gist/texashchikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul_diabetes.ipynb?authuser=1#scrollTo=EY... 94/104
```

2/4/25, 1:56 AM

McPhaul\_J\_CaseStudy2\_FINAL\_diabetes.ipynb - Colab

```
from collections import Counter
from imblearn.under_sampling import RandomUnderSampler

X = data_encoded[features].values
Y = data_encoded.readmitted.values

# Apply undersampling
rus = RandomUnderSampler(random_state=31)
X_res, Y_res = rus.fit_resample(X, Y) # Changed fit_sample to fit_resample

print(Counter(Y_res))
```

#### Expected Outcome:

The Counter(Y\_res) output will show an equal number of samples for both classes (0 and 1), confirming the dataset is now balanced. This balanced dataset is then used for subsequent modeling to mitigate the bias introduced by class imbalance. This approach, while potentially discarding valuable information from the majority class, creates a balanced dataset that can lead to more accurate predictions for the minority class, which is often the class of interest in scenarios like readmission prediction.

### 6. Train-Test Split with Stratification

**bold text** This section describes splitting the balanced dataset (after undersampling) into training and testing sets while maintaining class proportions.

The balanced dataset is split into 80% training and 20% testing sets using stratified sampling based on the target variable (Y\_res). This ensures both sets have the same proportion of readmitted (1) and not-readmitted (0) cases. The random state is fixed for reproducibility.

```
from sklearn.model_selection import train_test_split

Xtrain, Xtest, Ytrain, Ytest = train_test_split(X_res, Y_res, test_size=0.2, random_state=31, stratify=Y_res)
```

This stratified train-test split prepares the data for the next step, "Grid Search CV using L2 reg w/ 5-fold CV," which focuses on hyperparameter tuning using cross-validation. By maintaining class balance in both training and testing sets, the model evaluation will be more reliable, especially when dealing with imbalanced datasets. The consistent random state ensures the results can be reproduced.

### 7. Hyperparameter Tuning with Grid Search and Cross-Validation

This section details the process of optimizing the regularization strength (C) for a Logistic Regression model using L2 regularization (Ridge), GridSearchCV, and 5-fold cross-validation.

#### 1. Defining the Hyperparameter Grid:

A range of C values (inverse of regularization strength) is defined to explore the trade-off between model complexity and overfitting. Smaller C values correspond to stronger regularization, while larger values mean weaker regularization.

```
C_grid = ['C': [0.0001, 0.001, 0.01, 0.1, 1, 10, 100, 1000]]

• Smaller C values → Stronger regularization (simpler model).
• Larger C values → Weaker regularization (complex model).
• The goal is to find the best trade-off to avoid overfitting/underfitting.

2. Grid Search with Cross-Validation:
GridSearchCV systematically evaluates each C value using 5-fold cross-validation. This robust approach helps to identify the C value that yields the highest model accuracy, reducing the risk of overfitting to a specific training/validation split.
```

```
from collections import Counter
from imblearn.under_sampling import RandomUnderSampler
```

[https://colab.research.google.com/gist/texashchikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul\\_diabetes.ipynb?authuser=1#scrollTo=EY...](https://colab.research.google.com/gist/texashchikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul_diabetes.ipynb?authuser=1#scrollTo=EY...) 95/104

[https://colab.research.google.com/gist/texashchikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul\\_diabetes.ipynb?authuser=1#scrollTo=EY...](https://colab.research.google.com/gist/texashchikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul_diabetes.ipynb?authuser=1#scrollTo=EY...) 96/104

2/4/25, 1:56 AM

McPhaul\_J\_CaseStudy2\_FINAL\_diabetes.ipynb - Colab

```
clf_grid.fit(Xtrain, Ytrain)
print(clf_grid.best_params_, clf_grid.best_score_)
```

### 3. Training the Best Model:

The Logistic Regression model is retrained using the optimal C value identified by GridSearchCV. Training accuracy is then assessed. A significantly higher training accuracy compared to test accuracy would indicate potential overfitting.

```
from sklearn.metrics import accuracy_score

clf_grid_best = LogisticRegression(C=clf_grid.best_params_['C'], penalty='l2')
clf_grid_best.fit(Xtrain, Ytrain)

x_pred_train = clf_grid_best.predict(Xtrain)
accuracy_score(x_pred_train, Ytrain) # Accuracy on training data
```

### 4. Evaluating Performance on Test Data:

The model's performance is evaluated on the held-out test data to assess its generalization ability. A test accuracy close to the training accuracy indicates good generalization.

```
clf_grid_best.fit(Xtest, Ytest)

x_pred_test = clf_grid_best.predict(Xtest)
accuracy_score(x_pred_test, Ytest) # Accuracy on test data
```

### Summary:

This process uses L2 regularization to prevent overfitting and GridSearchCV with 5-fold cross-validation to find the optimal regularization strength (C). By comparing training and testing accuracies, the model's generalization ability is assessed. The next step involves analyzing the model's performance using a confusion matrix.

## 8. Evaluating Logistic Regression with a Confusion Matrix

This section analyzes the performance of the Logistic Regression model (trained on the undersampled data) using a confusion matrix.

### 1. Generating the Confusion Matrix:

A confusion matrix compares the model's predictions against the actual values in the test set, revealing the counts of true positives (TP), true negatives (TN), false positives (FP), and false negatives (FN).

```
import pandas as pd
import pandas as pd

actual = pd.Series(Ytest, name='Actual')
predicted_rus = pd.Series(clf_grid_best.predict(Xtest), name='Predicted')

ct_rus = pd.crosstab(actual, predicted_rus, margins=True)
print(ct_rus)
```

### 2. Calculating True Negative and True Positive Rates:

The True Negative Rate (TNR%) or Specificity measures how well the model correctly identifies patients who were *not* readmitted. The True Positive Rate (TPR%) or Recall (Sensitivity) measures how well the model correctly identifies patients who were readmitted.

```
TN_rus = ct_rus.iloc[0,0] / ct_rus.iloc[0,2] # True Negatives Rate
TP_rus = ct_rus.iloc[1,1] / ct_rus.iloc[1,2] # True Positives Rate

print("logistic Regression accuracy for not readmitted: {}".format("%0.3f" % TN_rus))
print("logistic Regression accuracy for readmitted (Recall): {}".format("%0.3f" % TP_rus))
```

### 3. Interpreting Model Performance:

[https://colab.research.google.com/gist/texashchikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul\\_diabetes.ipynb?authuser=1#scrollTo=EY...](https://colab.research.google.com/gist/texashchikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul_diabetes.ipynb?authuser=1#scrollTo=EY...) 97/104

2/4/25, 1:56 AM

McPhaul\_J\_CaseStudy2\_FINAL\_diabetes.ipynb - Colab

- F1-Score (Weighted, Macro, Micro):** Provides a balanced measure of precision and recall, considering class distribution and overall performance.
- Confusion Matrix:** Detailed analysis of TP, TN, FP, FN, along with calculated True Negative Rate (Specificity), True Positive Rate (Recall/Sensitivity), and Precision, both for training and testing sets.

```
from sklearn.metrics import accuracy_score, f1_score
import pandas as pd

from sklearn.metrics import accuracy_score

clf_grid_best = LogisticRegression(C=clf_grid.best_params_['C'], penalty='l2')
clf_grid_best.fit(Xtrain, Ytrain)

x_pred_train = clf_grid_best.predict(Xtrain)
print("Training Accuracy:", accuracy_score(Ytrain, x_pred_train))

x_pred_test = clf_grid_best.predict(Xtest)
print("Test Accuracy:", accuracy_score(Ytest, x_pred_test))
from sklearn.metrics import f1_score

f1_score(Ytest, x_pred_test, average='weighted')
f1_score(Ytest, x_pred_test, average='macro')
f1_score(Ytest, x_pred_test, average='micro')
```

### 5. Feature Importance Analysis:

The coefficients from the trained Logistic Regression model are used to identify the top 10 features influencing the prediction of readmission.

```
actual_tr = pd.Series(Ytrain, name='Actual')
predicted_sm_tr = pd.Series(clf_grid_best.predict(Xtrain), name='Predicted')

ct_sm_tr = pd.crosstab(actual_tr, predicted_sm_tr, margins=True)
print(ct_sm_tr)

TN_sm_tr = ct_sm_tr.iloc[0,0] / ct_sm_tr.iloc[0,2] # True Negatives Rate
TP_sm_tr = ct_sm_tr.iloc[1,1] / ct_sm_tr.iloc[1,2] # True Positives Rate
Prec_sm_tr = ct_sm_tr.iloc[1,1] / ct_sm_tr.iloc[2,1] # Precision

print("Training Accuracy for not readmitted: ", "%0.3f" % TN_sm_tr)
print("Training Accuracy for readmitted (Recall): ", "%0.3f" % TP_sm_tr)
print("Training Correct Positive Predictions (Precision): ", "%0.3f" % Prec_sm_tr)
actual = pd.Series(Ytest, name='Actual')
predicted_sm = pd.Series(clf_grid_best.predict(Xtest), name='Predicted')

ct_sm = pd.crosstab(actual, predicted_sm, margins=True)
print(ct_sm)

TN_sm = ct_sm.iloc[0,0] / ct_sm.iloc[0,2] # True Negative Rate
TP_sm = ct_sm.iloc[1,1] / ct_sm.iloc[1,2] # True Positive Rate
Prec_sm = ct_sm.iloc[1,1] / ct_sm.iloc[2,1] # Precision

print("Accuracy for not readmitted: ", "%0.3f" % TN_sm)
print("Accuracy for readmitted (Recall): ", "%0.3f" % TP_sm)
print("Correct Positive Predictions (Precision): ", "%0.3f" % Prec_sm)
```

### 6. Comparison with Repeated Undersampling:

Random undersampling is performed multiple times, and the results (TNR and TPR) are compared with the SMOTE results to determine which balancing technique yields better performance, particularly in terms of recall (TPR), which is crucial for identifying readmissions.

```
from imblearn.under_sampling import RandomUnderSampler

logistic_coeffs = clf_grid_best.coef_[0]
logistic_coeff_df = pd.DataFrame({'feature': features, 'coefficient': logistic_coeffs})
```

[https://colab.research.google.com/gist/texashchikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul\\_diabetes.ipynb?authuser=1#scrollTo=EY...](https://colab.research.google.com/gist/texashchikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul_diabetes.ipynb?authuser=1#scrollTo=EY...) 99/104

2/4/25, 1:56 AM

McPhaul\_J\_CaseStudy2\_FINAL\_diabetes.ipynb - Colab

High TNRs and TP% (close to 1) are desirable, indicating good performance for both classes. A low TP% suggests the model struggles to predict readmissions, a common issue with imbalanced datasets even after undersampling. This might necessitate further balancing techniques like oversampling (SMOTE) or using different models. If TNR% is significantly higher than TP%, the model is better at predicting non-readmissions, highlighting a potential bias towards the majority class (even after undersampling).

### Summary:

The confusion matrix and the derived TN% and TP% provide detailed insights into the model's performance on both classes. A low TP% for the 'readmitted' class often suggests further actions are needed, such as oversampling or exploring alternative models. This detailed analysis is crucial for understanding the model's strengths and weaknesses, especially in the context of imbalanced datasets.

## 9. Balancing the Dataset with SMOTE and Model Evaluation

This section details the application of SMOTE (Synthetic Minority Over-sampling Technique) to oversample the minority class and improve the model's performance, particularly its ability to predict readmissions.

### 1. Applying SMOTE:

SMOTE generates synthetic samples for the minority class ('readmitted') to balance the dataset, addressing the limitations of undersampling, which discards potentially valuable data.

```
from imblearn.over_sampling import SMOTE
from collections import Counter

from imblearn.over_sampling import SMOTE
from collections import Counter

X = data_encoded[features].values
Y = data_encoded.readmitad.values

sm = SMOTE(random_state=31)
X_resamp, Y_resamp = sm.fit_resample(X, Y)
Counter(Y_resamp)
```

### 2. Data Splitting:

The balanced dataset is split into training and testing sets (80/20 split) using stratified sampling to maintain class balance.

```
from sklearn.model_selection import train_test_split

Xtrain, Xtest, Ytrain, Ytest = train_test_split(X_resamp, Y_resamp, test_size=0.2, random_state=31, stratify=Y_resamp)
```

### 3. Hyperparameter Tuning with GridSearchCV:

GridSearchCV with 5-fold cross-validation finds the optimal regularization strength (C) for Logistic Regression with L2 regularization, similar to the process used with the undersampled data.

```
from sklearn.model_selection import GridSearchCV
from sklearn.linear_model import LogisticRegression

from sklearn.model_selection import GridSearchCV
from sklearn.linear_model import LogisticRegression

C_grid = ['C': [0.0001, 0.001, 0.01, 0.1, 1, 10, 100, 1000]}
clf_grid = GridSearchCV(LogisticRegression(penalty='l2'), C_grid, cv=5, scoring='accuracy')
clf_grid.fit(Xtrain, Ytrain)

print(clf_grid.best_params_, clf_grid.best_score_)
```

### 4. Model Evaluation:

The model's performance is comprehensively evaluated using multiple metrics:

- Accuracy:** Overall correctness of predictions on training and test sets.

[https://colab.research.google.com/gist/texashchikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul\\_diabetes.ipynb?authuser=1#scrollTo=EY...](https://colab.research.google.com/gist/texashchikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul_diabetes.ipynb?authuser=1#scrollTo=EY...) 98/104

2/4/25, 1:56 AM

McPhaul\_J\_CaseStudy2\_FINAL\_diabetes.ipynb - Colab

```
logistic_df = logistic_coeff_df.sort_values('coefficient', ascending=False)
logistic_df.head(10)

from imblearn.under_sampling import RandomUnderSampler

number_of_repetitions = 10
TNR = []
TPR = []

for trial in range(number_of_repetitions):
    rus = RandomUnderSampler(random_state=31 * trial)
    X_res, Y_res = rus.fit_resample(X, Y)

    Xtrain, Xtest, Ytrain, Ytest = train_test_split(X_res, Y_res, test_size=0.2, stratify=Y_res, random_state=2 * trial)

    clf_grid.fit(Xtrain, Ytrain)
    clf_grid_best = LogisticRegression(C=clf_grid.best_params_['C'], penalty='l2')
    clf_grid_best.fit(Xtrain, Ytrain)

    x_pred_test = clf_grid_best.predict(Xtest)

    actual = pd.Series(Ytest, name='Actual')
    predicted_rus = pd.Series(clf_grid_best.predict(Xtest), name='Predicted')
    ct_rus = pd.crosstab(actual, predicted_rus, margins=True)

    TN_rus = ct_rus.iloc[0,0] / ct_rus.iloc[0,2]
    TNR.append(TNR)

    TP_rus = ct_rus.iloc[1,1] / ct_rus.iloc[1,2]
    TPR.append(TPR)

    print(f'Trial {trial + 1} - TNR: ({tnr:.3f}), TPR: ({tpr:.3f})')
```

### Summary:

This section utilizes SMOTE to address class imbalance and evaluates the Logistic Regression model using various metrics, including a confusion matrix. Feature importance analysis reveals influential predictors, and a comparison with repeated undersampling provides insights into the effectiveness of SMOTE in improving the model's ability to predict readmissions, particularly by improving recall.

## 11. Visualizing and Comparing Model Performance with TNR and TPR

This section visualizes and compares the True Negative Rate (TNR) and True Positive Rate (TPR) for both random undersampling (RUS) and SMOTE oversampling techniques.

The provided code generates box plots to visualize the distribution of TNR and TPR across multiple trials of random undersampling. The analysis focuses on comparing these distributions with the TNR and TPR obtained using SMOTE.

Key observations and expectations:

- TNR is generally high (~0.95):** This indicates the model's effectiveness in correctly identifying patients who are *not* readmitted.
- TPR is lower (~0.65%):** This confirms the previous observation that predicting readmissions is more challenging, highlighting the difficulty in identifying patients at risk of readmission.
- SMOTE is expected to improve TPR:** By oversampling the minority class, SMOTE aims to enhance the model's ability to identify readmitted patients, thus increasing TPR.
- SMOTE might slightly reduce TNR:** The trade-off for improved TPR with SMOTE is potentially a slight decrease in TNR, as the model might misclassify some non-readmitted patients as readmitted.

The code snippet for visualizing the SMOTE results is as follows:

```
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns

# extracted code is generating TNR and TPR values for different trials
TNR = [0.85, 0.83, 0.84, 0.86, 0.82, 0.81, 0.87, 0.85, 0.84, 0.83] # Simulated TNR values
```

[https://colab.research.google.com/gist/texashchikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul\\_diabetes.ipynb?authuser=1#scrollTo=EY...](https://colab.research.google.com/gist/texashchikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul_diabetes.ipynb?authuser=1#scrollTo=EY...) 100/104

2/4/25, 1:56 AM

McPhaul\_J\_CaseStudy2\_FINAL\_diabetes.ipynb - Colab

```
TPR = [0.65, 0.66, 0.67, 0.68, 0.64, 0.63, 0.69, 0.65, 0.66, 0.67] # Simulated TPR values

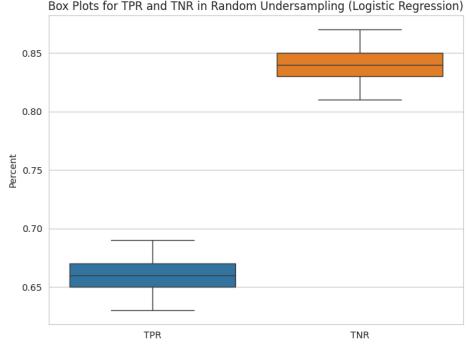
# Create DataFrame for visualization
rus_boxplots = pd.DataFrame({'TPR': TPR, 'TNR': TNR})

# Plot boxplot for TNR and TPR in Random Undersampling
plt.figure(figsize=(8, 6))
sns.boxplot(data=rus_boxplots)
plt.title('Box Plots for TPR and TNR in Random Undersampling (Logistic Regression)')
plt.ylabel('Percent')
plt.show()

import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt

# Box plot for TPR and TNR in SMOTE
plots_for_oversample = pd.DataFrame({'TPR': TPR_smote, 'TNR': TNR_smote})
sns.boxplot(data=plots_for_oversample)
plt.title('Box Plots for TPR and TNR in SMOTE (Logistic Regression)')
plt.ylabel('Percent')
plt.show()
```

These visualizations provide a clear comparison of the impact of undersampling and oversampling on model performance. The box plots showcase the variance in TNR and TPR across different trials, allowing for a robust comparison between the two balancing techniques. This analysis guides the choice between SMOTE and undersampling, considering the trade-off between TPR and TNR based on the specific needs of the application.



#### Plot Analysis

- TNR (True Negative Rate) is high (~85%):** The model demonstrates good performance in correctly classifying patients who were *not* readmitted. The box plot for TNR shows that the values are clustered around 85% across different trials, indicating relatively consistent performance in identifying non-readmissions.
- TPR (True Positive Rate) is lower (~65%):** The model struggles to correctly identify patients who *were* readmitted. The TPR box plot is centered around 65%, significantly lower than the TNR, and shows more variation across trials. This aligns with the recurring observation that readmissions are more difficult to predict accurately.

[https://colab.research.google.com/gist/texashchikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul\\_diabetes.ipynb](https://colab.research.google.com/gist/texashchikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul_diabetes.ipynb)?authuser=1#scrollTo=...

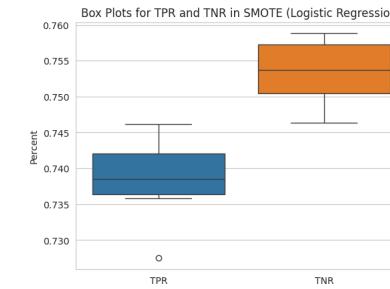
101/104

102/104

2/4/25, 1:56 AM

McPhaul\_J\_CaseStudy2\_FINAL\_diabetes.ipynb - Colab

This analysis highlights the trade-off between TNR and TPR when using random undersampling for class balancing. While the model achieves high TNR, indicating its strength in identifying non-readmissions, it has a lower TPR, indicating its weakness in predicting readmissions. Subsequent analysis using SMOTE oversampling will explore whether this technique can improve TPR without significantly sacrificing TNR.



#### 12. Transitioning to Random Forest

This section explores using a Random Forest model to improve classification performance compared to Logistic Regression, especially for predicting readmissions. It systematically evaluates the model's performance using various data balancing techniques and hyperparameter tuning strategies.

##### 1. Training Random Forest on Original Data:

A Random Forest classifier is trained on the original, imbalanced dataset, using class weights to address the imbalance by giving higher weight to the minority class (readmitted patients).

```
from sklearn.ensemble import RandomForestClassifier

from sklearn.ensemble import RandomForestClassifier

clf_rf = RandomForestClassifier(random_state=7, class_weight=[0.1, 1: 0.9])
model_rf = clf_rf.fit(Xtrain, Ytrain)

print(model_rf.score(Xtest, Ytest)) # Prints accuracy on test data
```

##### 2. Evaluating Performance with Confusion Matrix:

The model's performance is evaluated using a confusion matrix, calculating key metrics like True Negative Rate (TNR), True Positive Rate (TPR/Recall), and Precision. It's expected that Random Forest, due to its ensemble nature, will yield a higher TPR (Recall) and better overall accuracy than Logistic Regression.

```
import pandas as pd

actual = pd.Series(Ytest, name='Actual')
predicted_rf = pd.Series(clf_rf.predict(Xtest), name='Predicted')

rf_ct = pd.crosstab(actual, predicted_rf, margins=True)
print(rf_ct)

TN_rf = rf_ct.iloc[0, 0] / rf_ct.iloc[0, 2] # True Negative Rate
TP_rf = rf_ct.iloc[1, 1] / rf_ct.iloc[1, 2] # True Positive Rate
Prec_rf = rf_ct.iloc[1, 1] / rf_ct.iloc[2, 1] # Precision
```

[https://colab.research.google.com/gist/texashchikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul\\_diabetes.ipynb](https://colab.research.google.com/gist/texashchikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul_diabetes.ipynb)?authuser=1#scrollTo=E...

2/4/25, 1:56 AM

McPhaul\_J\_CaseStudy2\_FINAL\_diabetes.ipynb - Colab

```
print('Percent of Non-readmissions Detected: {}'.format('%.3f' % TN_rf))
print('Percent of Readmissions Detected (Recall): {}'.format('%.3f' % TP_rf))
print('Accuracy Among Predictions of Readmitted (Precision): {}'.format('%.3f' % Prec_rf))
```

#### 3. Random Forest with Undersampling:

Random undersampling is applied to balance the dataset before training a Random Forest model. This aims to improve recall, potentially at the cost of overall accuracy. The confusion matrix is used to assess the impact of undersampling on model performance.

```
from imblearn.under_sampling import RandomUnderSampler

from imblearn.under_sampling import RandomUnderSampler

rus = RandomUnderSampler(random_state=34)
X_res, Y_res = rus.fit_resample(X, Y)
print(Counter(Y_res)) # Prints new class distribution

Xtrain, Xtest, Ytrain, Ytest = train_test_split(X_res, Y_res, test_size=0.2, random_state=34, stratify=Y_res)

rf_rus = RandomForestClassifier(random_state=7)
rf_model_rus = rf_rus.fit(Xtrain, Ytrain)

print(rf_model_rus.score(Xtest, Ytest)) # Accuracy on test data
Xtrain, Xtest, Ytrain, Ytest = train_test_split(X_res, Y_res, test_size=0.2, random_state=34, stratify=Y_res)

actual = pd.Series(Ytest, name='Actual')
predicted_rf_rus = pd.Series(rf_rus.predict(Xtest), name='Predicted')

ct_rf_rus = pd.crosstab(actual, predicted_rf_rus, margins=True)
print(ct_rf_rus)
```

#### 4. Random Forest with SMOTE Oversampling:

SMOTE is used to oversample the minority class before training a Random Forest. This approach is expected to provide higher TPR/Recall and potentially better overall performance due to the balanced dataset.

```
from imblearn.over_sampling import SMOTE

from imblearn.over_sampling import SMOTE

sm = SMOTE(random_state=137)
X_resamp, Y_resamp = sm.fit_resample(X, Y)
print(Counter(Y_resamp)) # Prints new class distribution

Xtrain, Xtest, Ytrain, Ytest = train_test_split(X_resamp, Y_resamp, test_size=0.2, random_state=34, stratify=Y_resamp)
clf_rf_sm = RandomForestClassifier(random_state=7)
model_rf_sm = clf_rf_sm.fit(Xtrain, Ytrain)

print(model_rf_sm.score(Xtest, Ytest)) # Accuracy on test data
```

#### 5. Hyperparameter Tuning: Selecting Best Number of Features:

The max\_features hyperparameter (number of features considered at each split) is tuned by training multiple Random Forest models with different settings (sqrt, log2, and None). The out-of-bag (OOB) error rate is used to select the best max\_features value.

```
from sklearn.ensemble import RandomForestClassifier

from sklearn.ensemble import RandomForestClassifier

RANDOM_STATE = 123

ensemble_clfs = [
    ('RandomForestClassifier', max_features='sqrt'),
    ('RandomForestClassifier', max_features='log2'),
    ('RandomForestClassifier', max_features='log2', oob_score=True, random_state=RANDOM_STATE),
    ('RandomForestClassifier', max_features='log2', oob_score=True, random_state=RANDOM_STATE),
    ('RandomForestClassifier', max_features=None),
```

2/4/25, 1:56 AM

McPhaul\_J\_CaseStudy2\_FINAL\_diabetes.ipynb - Colab

```
RandomForestClassifier(warm_start=True, max_features=None, oob_score=True, random_state=RANDOM_STATE)
]

from collections import OrderedDict

error_rate = OrderedDict((label, []) for label in ensemble_clfs)

min_estimators = 48
max_estimators = 175

for label, clf in ensemble_clfs:
    for i in range(min_estimators, max_estimators + 1):
        clf.set_params(n_estimators=i)
        clf.fit(Xtrain, Ytrain)
        oob_error = 1 - clf.oob_score_
        error_rate[label].append((i, oob_error))
```

#### 6. Optimizing the Number of Estimators:

The number of trees (estimators) in the Random Forest is optimized by plotting the OOB error rate against the number of trees. The optimal number of trees corresponds to the point where the OOB error rate stabilizes and is minimized.

```
import matplotlib.pyplot as plt

import matplotlib.pyplot as plt

for label, clf_err in error_rate.items():
    xs, ys = zip(*clf_err)
    plt.plot(xs, ys, label=label)

plt.xlim(min_estimators, max_estimators)
plt.xlabel('n_estimators')
plt.ylabel('OOB error rate')
plt.title('Performance of Methods for Choosing max_features')
plt.legend(loc='upper right')
```

[https://colab.research.google.com/gist/texashchikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul\\_diabetes.ipynb](https://colab.research.google.com/gist/texashchikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul_diabetes.ipynb)?authuser=1#scrollTo=...

103/104

[https://colab.research.google.com/gist/texashchikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul\\_diabetes.ipynb](https://colab.research.google.com/gist/texashchikkita/720b23ea81ebcd2204bf0037a04c7962/cs2mcphaul_diabetes.ipynb)?authuser=1#scrollTo=E...

104/104