

```

1 # Install required libraries
2 !pip install pandas numpy matplotlib seaborn scikit-learn tensorflow keras imbalanced-learn xgboost
3
4 # Import libraries
5 import pandas as pd
6 import numpy as np
7 import matplotlib.pyplot as plt
8 import seaborn as sns
9 from sklearn.model_selection import train_test_split, cross_val_score, GridSearchCV
10 from sklearn.preprocessing import StandardScaler, OneHotEncoder
11 from sklearn.impute import SimpleImputer
12 from sklearn.compose import ColumnTransformer
13 from sklearn.pipeline import Pipeline
14 from sklearn.metrics import mean_absolute_error, mean_squared_error, r2_score, accuracy_score, precision_score, recall_score, f1_score
15 from sklearn.linear_model import LinearRegression, Ridge, Lasso
16 from sklearn.ensemble import RandomForestRegressor, GradientBoostingRegressor, RandomForestClassifier
17 from sklearn.svm import SVR, SVC
18 from sklearn.tree import DecisionTreeRegressor, DecisionTreeClassifier
19 from imblearn.over_sampling import SMOTE
20 import warnings
21 warnings.filterwarnings('ignore')
22

→ Requirement already satisfied: numpy in /usr/local/lib/python3.11/dist-packages (2.0.2)
Requirement already satisfied: matplotlib in /usr/local/lib/python3.11/dist-packages (3.10.0)
Requirement already satisfied: seaborn in /usr/local/lib/python3.11/dist-packages (0.13.2)
Requirement already satisfied: scikit-learn in /usr/local/lib/python3.11/dist-packages (1.6.1)
Requirement already satisfied: tensorflow in /usr/local/lib/python3.11/dist-packages (2.18.0)
Requirement already satisfied: keras in /usr/local/lib/python3.11/dist-packages (3.8.0)
Requirement already satisfied: imbalanced-learn in /usr/local/lib/python3.11/dist-packages (0.13.0)
Requirement already satisfied: xgboost in /usr/local/lib/python3.11/dist-packages (2.1.4)
Requirement already satisfied: python-dateutil>=2.8.2 in /usr/local/lib/python3.11/dist-packages (from pandas) (2.9.0.post0)
Requirement already satisfied: pytz>=2020.1 in /usr/local/lib/python3.11/dist-packages (from pandas) (2025.2)
Requirement already satisfied: tzdata>=2022.7 in /usr/local/lib/python3.11/dist-packages (from pandas) (2025.2)
Requirement already satisfied: contourpy>=1.0.1 in /usr/local/lib/python3.11/dist-packages (from matplotlib) (1.3.2)
Requirement already satisfied: cycler>=0.10 in /usr/local/lib/python3.11/dist-packages (from matplotlib) (0.12.1)
Requirement already satisfied: fonttools>=4.22.0 in /usr/local/lib/python3.11/dist-packages (from matplotlib) (4.57.0)
Requirement already satisfied: kiwisolver>=1.3.1 in /usr/local/lib/python3.11/dist-packages (from matplotlib) (1.4.8)
Requirement already satisfied: packaging>=20.0 in /usr/local/lib/python3.11/dist-packages (from matplotlib) (24.2)
Requirement already satisfied: pillow>=8 in /usr/local/lib/python3.11/dist-packages (from matplotlib) (11.2.1)
Requirement already satisfied: pyparsing>=2.3.1 in /usr/local/lib/python3.11/dist-packages (from matplotlib) (3.2.3)
Requirement already satisfied: scipy>=1.6.0 in /usr/local/lib/python3.11/dist-packages (from scikit-learn) (1.15.2)
Requirement already satisfied: joblib>=1.2.0 in /usr/local/lib/python3.11/dist-packages (from scikit-learn) (1.5.0)
Requirement already satisfied: threadpoolctl>=3.1.0 in /usr/local/lib/python3.11/dist-packages (from scikit-learn) (3.6.0)
Requirement already satisfied: absl-py>=1.0.0 in /usr/local/lib/python3.11/dist-packages (from tensorflow) (1.4.0)
Requirement already satisfied: astunparse>=1.6.0 in /usr/local/lib/python3.11/dist-packages (from tensorflow) (1.6.3)
Requirement already satisfied: flatbuffers>=24.3.25 in /usr/local/lib/python3.11/dist-packages (from tensorflow) (25.2.10)
Requirement already satisfied: gast!=0.5.0,!0=0.5.1,!0=0.5.2,>0=0.2.1 in /usr/local/lib/python3.11/dist-packages (from tensorflow) (0.5.2)
Requirement already satisfied: google-pasta>=0.1.1 in /usr/local/lib/python3.11/dist-packages (from tensorflow) (0.2.0)
Requirement already satisfied: libclang>=13.0.0 in /usr/local/lib/python3.11/dist-packages (from tensorflow) (18.1.1)
Requirement already satisfied: opt-eins>=2.3.2 in /usr/local/lib/python3.11/dist-packages (from tensorflow) (3.4.0)
Requirement already satisfied: protobuf!=4.21.0,!0=4.21.1,!0=4.21.2,!0=4.21.3,!0=4.21.4,!0=4.21.5,<6.0.0dev,>=3.20.3 in /usr/local/lib/python3.11/dist-packages (from tensorflow) (4.21.5)
Requirement already satisfied: requests<3,>=2.21.0 in /usr/local/lib/python3.11/dist-packages (from tensorflow) (2.32.3)
Requirement already satisfied: setuptools in /usr/local/lib/python3.11/dist-packages (from tensorflow) (75.2.0)
Requirement already satisfied: six>=1.12.0 in /usr/local/lib/python3.11/dist-packages (from tensorflow) (1.17.0)
Requirement already satisfied: termcolor>=1.1.0 in /usr/local/lib/python3.11/dist-packages (from tensorflow) (3.1.0)
Requirement already satisfied: typing-extensions>=3.6.6 in /usr/local/lib/python3.11/dist-packages (from tensorflow) (4.13.2)
Requirement already satisfied: wrapt>=1.11.0 in /usr/local/lib/python3.11/dist-packages (from tensorflow) (1.17.2)
Requirement already satisfied: grpcio<2.0,>=1.24.3 in /usr/local/lib/python3.11/dist-packages (from tensorflow) (1.71.0)
Requirement already satisfied: tensorboard<2.19,>=2.18 in /usr/local/lib/python3.11/dist-packages (from tensorflow) (2.18.0)
Requirement already satisfied: h5py>=3.11.0 in /usr/local/lib/python3.11/dist-packages (from tensorflow) (3.13.0)
Requirement already satisfied: ml-dtypes<0.5.0,>=0.4.0 in /usr/local/lib/python3.11/dist-packages (from tensorflow) (0.4.1)
Requirement already satisfied: tensorflow-io-gcs-filesystem>=0.23.1 in /usr/local/lib/python3.11/dist-packages (from tensorflow) (0.23.1)
Requirement already satisfied: rich in /usr/local/lib/python3.11/dist-packages (from keras) (13.9.4)
Requirement already satisfied: namex in /usr/local/lib/python3.11/dist-packages (from keras) (0.0.9)
Requirement already satisfied: optree in /usr/local/lib/python3.11/dist-packages (from keras) (0.15.0)
Requirement already satisfied: sklearn-compat<1,>=0.1 in /usr/local/lib/python3.11/dist-packages (from imbalanced-learn) (0.1.3)
Requirement already satisfied: nvidia-nccl-cu12 in /usr/local/lib/python3.11/dist-packages (from xgboost) (2.21.5)
Requirement already satisfied: wheel<1.0,>=0.23.0 in /usr/local/lib/python3.11/dist-packages (from astunparse>=1.6.0->tensorflow) (0.23.0)
Requirement already satisfied: charset-normalizer<4,>=2 in /usr/local/lib/python3.11/dist-packages (from requests<3,>=2.21.0->tensorflow) (4.1.0)
Requirement already satisfied: idna<4,>=2.5 in /usr/local/lib/python3.11/dist-packages (from requests<3,>=2.21.0->tensorflow) (3.3.0)
Requirement already satisfied: urllib3<3,>=1.21.1 in /usr/local/lib/python3.11/dist-packages (from requests<3,>=2.21.0->tensorflow) (1.21.1)
Requirement already satisfied: certifi>=2017.4.17 in /usr/local/lib/python3.11/dist-packages (from requests<3,>=2.21.0->tensorflow) (2017.4.17)
Requirement already satisfied: markdown>=2.6.8 in /usr/local/lib/python3.11/dist-packages (from tensorflow<2.19,>=2.18->tensorflow) (2.6.8)
Requirement already satisfied: tensorboard-data-server<0.8.0,>=0.7.0 in /usr/local/lib/python3.11/dist-packages (from tensorflow<2.19,>=2.18->tensorflow) (0.7.0)
Requirement already satisfied: werkzeug>=1.0.1 in /usr/local/lib/python3.11/dist-packages (from tensorflow<2.19,>=2.18->tensorflow) (1.0.1)
Requirement already satisfied: markdown-it-py>=2.2.0 in /usr/local/lib/python3.11/dist-packages (from rich->keras) (3.0.0)
Requirement already satisfied: pygments<3.0.0,>=2.13.0 in /usr/local/lib/python3.11/dist-packages (from rich->keras) (2.19.1)
Requirement already satisfied: mdurl~0.1 in /usr/local/lib/python3.11/dist-packages (from markdown-it-py>=2.2.0->rich->keras) (0.1.0)
Requirement already satisfied: MarkupSafe>=2.1.1 in /usr/local/lib/python3.11/dist-packages (from werkzeug>=1.0.1->tensorflow<2.19,>=2.18->tensorflow) (2.1.1)

```

```

1 import pandas as pd # Make sure to import pandas first
2
3 # Read the Excel file - specify sheet_name as a string or index
4 warfarin_data = pd.read_excel('PS206767-553247439.xls', sheet_name='Subject Data') # Adjust sheet name as needed

```

```

5
6 # Display the first few rows
7 print(warfarin_data.head())

```

	PharmGKB Subject ID	PharmGKB Sample ID	Project Site	Gender	
0	PA135312261	PA135312629	1	male	
1	PA135312262	PA135312630	1	female	
2	PA135312263	PA135312631	1	female	
3	PA135312264	PA135312632	1	male	
4	PA135312265	PA135312633	1	male	

	Race (Reported)	Race (OMB)	Ethnicity (Reported)	Ethnicity (OMB)	
0	White	White	not Hispanic or Latino	not Hispanic or Latino	
1	White	White	not Hispanic or Latino	not Hispanic or Latino	
2	White	White	not Hispanic or Latino	not Hispanic or Latino	
3	White	White	not Hispanic or Latino	not Hispanic or Latino	
4	White	White	not Hispanic or Latino	not Hispanic or Latino	

	Age	Height (cm)	...	
0	60 - 69	193.040	...	
1	50 - 59	176.530	...	
2	40 - 49	162.560	...	
3	60 - 69	182.245	...	
4	50 - 59	167.640	...	

	VKORC1 QC genotype: -4451 C>A (861); Chr16:31018002; rs17880887; A/C			
0				NaN
1				C/C
2				NaN
3				NaN
4				NaN

	CYP2C9 consensus	VKORC1 -1639 consensus	VKORC1 497 consensus	
0	*1/*1	A/G	G/T	
1	*1/*1	A/A	G/T	
2	*1/*1	G/G	T/T	
3	*1/*1	A/G	G/T	
4	*1/*3	A/G	T/T	

	VKORC1 1173 consensus	VKORC1 1542 consensus	VKORC1 3730 consensus	
0	NaN	C/G	A/G	
1	T/T	C/C	G/G	
2	NaN	G/G	A/G	
3	NaN	C/G	G/G	
4	NaN	C/G	A/G	

	VKORC1 2255 consensus	VKORC1 -4451 consensus	
0	NaN	NaN	
1	T/T	C/C	
2	NaN	NaN	
3	NaN	NaN	
4	NaN	NaN	

	Comments regarding Project Site Dataset		
0		Project 1:	
1		Warfarin Therapeutic Dose Definition:	
2	The dose (unchanged for 6 days) that yielded a...		
3		NaN	
4		NaN	

[5 rows x 68 columns]

```
1 warfarin_data.head()
```

[

	PharmGKB Subject ID	PharmGKB Sample ID	Project Site	Gender	Race (Reported)	Race (OMB)	Ethnicity (Reported)	Ethnicity (OMB)	Age	Height (cm)	...	VKORC1 QC genotype: -4451 C>A (861); Chr16:31018002; rs17880887; A/C	CYP2C12 consensus:
0	PA135312261	PA135312629		1 male	White	White	Hispanic or Latino	not Hispanic or Latino	60 - 69	193.040	...	NaN	*1/*
1	PA135312262	PA135312630		1 female	White	White	Hispanic or Latino	not Hispanic or Latino	50 - 59	176.530	...	C/C	*1/*
2	PA135312263	PA135312631		1 female	White	White	Hispanic or Latino	not Hispanic or Latino	40 - 49	162.560	...	NaN	*1/*
3	PA135312264	PA135312632		1 male	White	White	Hispanic or Latino	not Hispanic or Latino	60 - 69	182.245	...	NaN	*1/*
4	PA135312265	PA135312633		1 male	White	White	Hispanic or Latino	not Hispanic or Latino	50 - 59	167.640	...	NaN	*1/*

5 rows × 68 columns

1 warfarin_data.columns

Index(['PharmGKB Subject ID', 'PharmGKB Sample ID', 'Project Site', 'Gender', 'Race (Reported)', 'Race (OMB)', 'Ethnicity (Reported)', 'Ethnicity (OMB)', 'Age', 'Height (cm)', 'Weight (kg)', 'Indication for Warfarin Treatment', 'Comorbidities', 'Diabetes', 'Congestive Heart Failure and/or Cardiomyopathy', 'Valve Replacement', 'Medications', 'Aspirin', 'Acetaminophen or Paracetamol (Tylenol)', 'Was Dose of Acetaminophen or Paracetamol (Tylenol) >1300mg/day', 'Simvastatin (Zocor)', 'Atorvastatin (Lipitor)', 'Fluvastatin (Lescol)', 'Lovastatin (Mevacor)', 'Pravastatin (Pravachol)', 'Rosuvastatin (Crestor)', 'Cerivastatin (Baycol)', 'Amiodarone (Cordarone)', 'Carbamazepine (Tegretol)', 'Phenytoin (Dilantin)', 'Rifampin or Rifampicin', 'Sulfonamide Antibiotics', 'Macrolide Antibiotics', 'Anti-fungal Azoles', 'Herbal Medications, Vitamins, Supplements', 'Target INR', 'Estimated Target INR Range Based on Indication', 'Subject Reached Stable Dose of Warfarin', 'Therapeutic Dose of Warfarin', 'INR on Reported Therapeutic Dose of Warfarin', 'Current Smoker', 'Cyp2C9 genotypes', 'Genotyped QC Cyp2C9*2', 'Genotyped QC Cyp2C9*3', 'Combined QC Cyp2C9', 'VKORC1 genotype: -1639 G>A (3673); chr16:31015190; rs9923231; C/T', 'VKORC1 QC genotype: -1639 G>A (3673); chr16:31015190; rs9923231; C/T', 'VKORC1 genotype: 497T>G (5808); chr16:31013055; rs2884737; A/C', 'VKORC1 QC genotype: 497T>G (5808); chr16:31013055; rs2884737; A/C', 'VKORC1 genotype: 1173 C>T(6484); chr16:31012379; rs9934438; A/G', 'VKORC1 QC genotype: 1173 C>T(6484); chr16:31012379; rs9934438; A/G', 'VKORC1 genotype: 1542G>C (6853); chr16:31012010; rs8050894; C/G', 'VKORC1 QC genotype: 1542G>C (6853); chr16:31012010; rs8050894; C/G', 'VKORC1 genotype: 3730 G>A (9041); chr16:31009822; rs7294; A/G', 'VKORC1 QC genotype: 3730 G>A (9041); chr16:31009822; rs7294; A/G', 'VKORC1 genotype: 2255C>T (7566); chr16:31011297; rs2359612; A/G', 'VKORC1 QC genotype: 2255C>T (7566); chr16:31011297; rs2359612; A/G', 'VKORC1 genotype: -4451 C>A (861); Chr16:31018002; rs17880887; A/C', 'VKORC1 QC genotype: -4451 C>A (861); Chr16:31018002; rs17880887; A/C', 'CYP2C9 consensus', 'VKORC1 -1639 consensus', 'VKORC1 497 consensus', 'VKORC1 1173 consensus', 'VKORC1 1542 consensus', 'VKORC1 3730 consensus', 'VKORC1 2255 consensus', 'VKORC1 -4451 consensus', 'Comments regarding Project Site Dataset'], dtype='object')

1 warfarin_data.describe()

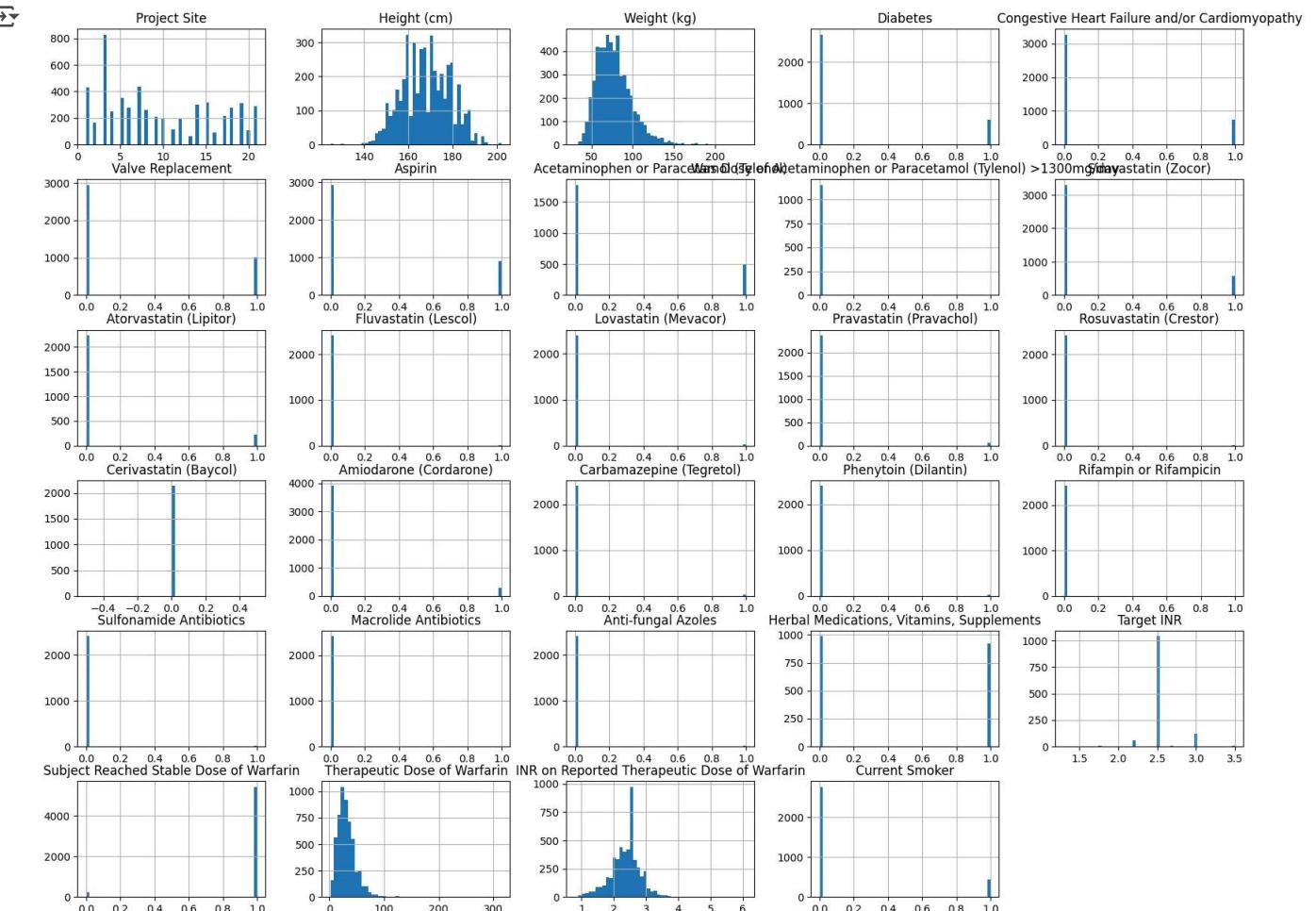
	Project Site	Height (cm)	Weight (kg)	Diabetes	Congestive Heart Failure and/or Cardiomyopathy	Valve Replacement	Aspirin	Acetaminophen or Paracetamol (Tylenol)	Acetaminophen or Paracetamol (Tylenol) >1300mg/day	Sim
count	5700.000000	4554.000000	5413.000000	3283.000000	4012.000000	3958.000000	3840.000000	2274.000000	1159.000000	386
mean	9.555965	168.047778	77.852569	0.187024	0.183699	0.256443	0.238542	0.220317	0.006040	
std	6.340864	10.845992	21.859764	0.389990	0.387287	0.436724	0.426247	0.414551	0.077514	
min	1.000000	124.968000	30.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	
25%	4.000000	160.020000	62.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	
50%	8.000000	167.894000	75.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	
75%	15.000000	176.022000	90.000000	0.000000	0.000000	1.000000	0.000000	0.000000	0.000000	
max	21.000000	202.000000	237.700000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	

8 rows × 29 columns

```

1 %matplotlib inline
2 import matplotlib.pyplot as plt
3 warfarin_data.hist(bins=50, figsize=(20,15))
4 plt.show()

```



```

1 np.random.seed(10)

1 # For illustration only. Sklearn has train_test_split()
2 def split_train_test(data, test_ratio):
3     shuffled_indices = np.random.permutation(len(data))
4     test_set_size = int(len(data) * test_ratio)
5     test_indices = shuffled_indices[:test_set_size]
6     train_indices = shuffled_indices[test_set_size:]
7     return data.iloc[train_indices], data.iloc[test_indices]

```

```

1 # run the function to get the train & test set
2 train_set, test_set = split_train_test(warfarin_data, 0.2)

```

```
1 train_set.info()
```

12	Comorbidities	3028	non-null	object
13	Diabetes	2651	non-null	float64
14	Congestive Heart Failure and/or Cardiomyopathy	3218	non-null	float64
15	Valve Replacement	3167	non-null	float64
16	Medications	3454	non-null	object
17	Aspirin	3094	non-null	float64
18	Acetaminophen or Paracetamol (Tylenol)	1826	non-null	float64
19	Was Dose of Acetaminophen or Paracetamol (Tylenol) >1300mg/day	952	non-null	float64
20	Simvastatin (Zocor)	3112	non-null	float64
21	Atorvastatin (Lipitor)	1988	non-null	float64
22	Fluvastatin (Lescol)	1957	non-null	float64
23	Lovastatin (Mevacor)	1964	non-null	float64
24	Pravastatin (Pravachol)	1965	non-null	float64
25	Rosuvastatin (Crestor)	1957	non-null	float64
26	Cerivastatin (Baycol)	1714	non-null	float64
27	Amiodarone (Cordarone)	3359	non-null	float64
28	Carbamazepine (Tegretol)	1963	non-null	float64
29	Phenytoin (Dilantin)	1958	non-null	float64
30	Rifampin or Rifampicin	1957	non-null	float64
31	Sulfonamide Antibiotics	1955	non-null	float64
32	Macrolide Antibiotics	1957	non-null	float64
33	Anti-fungal Azoles	1960	non-null	float64
34	Herbal Medications, Vitamins, Supplements	1544	non-null	float64
35	Target INR	1005	non-null	float64
36	Estimated Target INR Range Based on Indication	3072	non-null	object
37	Subject Reached Stable Dose of Warfarin	4527	non-null	float64
38	Therapeutic Dose of Warfarin	4425	non-null	float64
39	INR on Reported Therapeutic Dose of Warfarin	3971	non-null	float64
40	Current Smoker	2572	non-null	float64
41	Cyp2C9 genotypes	4449	non-null	object
42	Genotyped QC Cyp2C9*2	388	non-null	object
43	Genotyped QC Cyp2C9*3	388	non-null	object
44	Combined QC CYP2C9	388	non-null	object
45	VKORC1 genotype: -1639 G>A (3673); chr16:31015190; rs9923231; C/T	3214	non-null	object
46	VKORC1 QC genotype: -1639 G>A (3673); chr16:31015190; rs9923231; C/T	388	non-null	object
47	VKORC1 genotype: 497T>G (5808); chr16:31013055; rs2884737; A/C	1277	non-null	object
48	VKORC1 QC genotype: 497T>G (5808); chr16:31013055; rs2884737; A/C	388	non-null	object
49	VKORC1 genotype: 1173 C>T(6484); chr16:31012379; rs9934438; A/G	2687	non-null	object
50	VKORC1 QC genotype: 1173 C>T(6484); chr16:31012379; rs9934438; A/G	388	non-null	object
51	VKORC1 genotype: 1542G>C (6853); chr16:31012010; rs8050894; C/G	2661	non-null	object
52	VKORC1 QC genotype: 1542G>C (6853); chr16:31012010; rs8050894; C/G	388	non-null	object
53	VKORC1 genotype: 3730 G>A (9041); chr16:31009822; rs7294; A/G	2476	non-null	object
54	VKORC1 QC genotype: 3730 G>A (9041); chr16:31009822; rs7294; A/G	388	non-null	object
55	VKORC1 genotype: 2255C>T (7566); chr16:31011297; rs2359612; A/G	1755	non-null	object
56	VKORC1 QC genotype: 2255C>T (7566); chr16:31011297; rs2359612; A/G	388	non-null	object
57	VKORC1 genotype: -4451 C>A (861); Chr16:31018002; rs17880887; A/C	428	non-null	object
58	VKORC1 QC genotype: -4451 C>A (861); Chr16:31018002; rs17880887; A/C	364	non-null	object
59	CYP2C9 consensus	4442	non-null	object
60	VKORC1 -1639 consensus	3341	non-null	object
61	VKORC1 497 consensus	1509	non-null	object
62	VKORC1 1173 consensus	2873	non-null	object
63	VKORC1 1542 consensus	2818	non-null	object
64	VKORC1 3730 consensus	2661	non-null	object
65	VKORC1 2255 consensus	2011	non-null	object
66	VKORC1 -4451 consensus	727	non-null	object
67	Comments regarding Project Site Dataset	129	non-null	object

```
dtypes: float64(28), int64(1), object(39)
```

```
memory usage: 2.4+ MB
```

```
1 # After identifying the correct column names, let's define our features and target
```

```
2 # These need to be adjusted based on the actual column names from the dataset
```

```
3
```

```
4 # First check for the correct therapeutic dose column name
```

```
5 dose_columns = [col for col in warfarin_data.columns if 'dose' in col.lower() or 'therapeutic' in col.lower()]
```

```
6 print("Possible therapeutic dose columns:", dose_columns)
```

```
7
```

```
8 # Now check for other feature columns
```

```
9 demographic_columns = [col for col in warfarin_data.columns if any(term in col.lower() for term in ['gender', 'age', 'race', 'height'])]
```

```
10 print("Demographic columns:", demographic_columns)
```

```
11
```

```
12 medication_columns = [col for col in warfarin_data.columns if any(term in col.lower() for term in ['diabetes', 'simvastatin', 'amiodarone'])]
13 print("Medication and genetic columns:", medication_columns)
14
15 # Now define the features and target correctly
16 # Adjust these based on actual column names
17 target_regression = [col for col in dose_columns if 'therapeutic' in col.lower()][0]
18 print(f"Target regression column: {target_regression}")
19
20 # List of features - adjust based on your actual column names
21 features = demographic_columns + medication_columns
22 print(f"Features: {features}")
23
24 # Now display data types and missing values with correct column names
25 print("\nData types:")
26 warfarin_data[features + [target_regression]].info()
27
28 print("\nMissing values:")
29 print(warfarin_data[features + [target_regression]].isnull().sum())
30
31 # Descriptive statistics
32 print("\nDescriptive statistics:")
33 warfarin_data[features + [target_regression]].describe()
```

Possible therapeutic dose columns: ['Was Dose of Acetaminophen or Paracetamol (Tylenol) >1300mg/day', 'Subject Reached Stable Dos Demographic columns: ['Gender', 'Race (Reported)', 'Race (OMB)', 'Age', 'Height (cm)', 'Weight (kg)', 'Acetaminophen or Paracetam Medication and genetic columns: ['Diabetes', 'Simvastatin (Zocor)', 'Amiodarone (Cordarone)', 'Target INR', 'Estimated Target INR Target regression column: Therapeutic Dose of Warfarin Features: ['Gender', 'Race (Reported)', 'Race (OMB)', 'Age', 'Height (cm)', 'Weight (kg)', 'Acetaminophen or Paracetamol (Tylenol

Data types:

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 5700 entries, 0 to 5699

Data columns (total 41 columns):

#	Column	Non-Null Count	Dtype
0	Gender	5696	object
1	Race (Reported)	5194	object
2	Race (OMB)	5700	object
3	Age	5658	object
4	Height (cm)	4554	float64
5	Weight (kg)	5413	float64
6	Acetaminophen or Paracetamol (Tylenol)	2274	float64
7	Was Dose of Acetaminophen or Paracetamol (Tylenol) >1300mg/day	1159	float64
8	Diabetes	3283	float64
9	Simvastatin (Zocor)	3861	float64
10	Amiodarone (Cordarone)	4182	float64
11	Target INR	1259	float64
12	Estimated Target INR Range Based on Indication	3818	object
13	INR on Reported Therapeutic Dose of Warfarin	4968	float64
14	Cyp2C9 genotypes	5567	object
15	Genotyped QC Cyp2C9*2	478	object
16	Genotyped QC Cyp2C9*3	478	object
17	Combined QC CYP2C9	478	object
18	VKORC1 genotype: -1639 G>A (3673); chr16:31015190; rs9923231; C/T	4046	object
19	VKORC1 QC genotype: -1639 G>A (3673); chr16:31015190; rs9923231; C/T	478	object
20	VKORC1 genotype: 497T>G (5808); chr16:31013055; rs2884737; A/C	1599	object
21	VKORC1 QC genotype: 497T>G (5808); chr16:31013055; rs2884737; A/C	478	object
22	VKORC1 genotype: 1173 C>T(6484); chr16:31012379; rs9934438; A/G	3363	object
23	VKORC1 QC genotype: 1173 C>T(6484); chr16:31012379; rs9934438; A/G	478	object
24	VKORC1 genotype: 1542G>C (6853); chr16:31012010; rs8050894; C/G	3330	object
25	VKORC1 QC genotype: 1542G>C (6853); chr16:31012010; rs8050894; C/G	478	object
26	VKORC1 genotype: 3730 G>A (9041); chr16:31009822; rs7294; A/G	3124	object
27	VKORC1 QC genotype: 3730 G>A (9041); chr16:31009822; rs7294; A/G	478	object
28	VKORC1 genotype: 2255C>T (7566); chr16:31011297; rs2359612; A/G	2223	object
29	VKORC1 QC genotype: 2255C>T (7566); chr16:31011297; rs2359612; A/G	478	object
30	VKORC1 genotype: -4451 C>A (861); Chr16:31018002; rs17880887; A/C	540	object
31	VKORC1 QC genotype: -4451 C>A (861); Chr16:31018002; rs17880887; A/C	451	object
32	CYP2C9 consensus	5556	object
33	VKORC1 -1639 consensus	4201	object
34	VKORC1 497 consensus	1883	object
35	VKORC1 1173 consensus	3591	object
36	VKORC1 1542 consensus	3521	object
37	VKORC1 3730 consensus	3344	object
38	VKORC1 2255 consensus	2536	object
39	VKORC1 -4451 consensus	912	object
40	Therapeutic Dose of Warfarin	5528	float64

dtypes: float64(10), object(31)

memory usage: 1.8+ MB

Missing values:

Gender	4
Race (Reported)	506
Race (OMB)	0
Age	42
Height (cm)	1146
Weight (kg)	287
Acetaminophen or Paracetamol (Tylenol)	3426
Was Dose of Acetaminophen or Paracetamol (Tylenol) >1300mg/day	4541
Diabetes	2417
Simvastatin (Zocor)	1839
Amiodarone (Cordarone)	1518
Target INR	4441
Estimated Target INR Range Based on Indication	1882
INR on Reported Therapeutic Dose of Warfarin	732
Cyp2C9 genotypes	133
Genotyped QC Cyp2C9*2	5222
Genotyped QC Cyp2C9*3	5222
Combined QC CYP2C9	5222
VKORC1 genotype: -1639 G>A (3673); chr16:31015190; rs9923231; C/T	1654
VKORC1 QC genotype: -1639 G>A (3673); chr16:31015190; rs9923231; C/T	5222
VKORC1 genotype: 497T>G (5808); chr16:31013055; rs2884737; A/C	4101
VKORC1 QC genotype: 497T>G (5808); chr16:31013055; rs2884737; A/C	5222
VKORC1 genotype: 1173 C>T(6484); chr16:31012379; rs9934438; A/G	2337
VKORC1 QC genotype: 1173 C>T(6484); chr16:31012379; rs9934438; A/G	5222
VKORC1 genotype: 1542G>C (6853); chr16:31012010; rs8050894; C/G	2370
VKORC1 QC genotype: 1542G>C (6853); chr16:31012010; rs8050894; C/G	5222
VKORC1 genotype: 3730 G>A (9041); chr16:31009822; rs7294; A/G	2576
VKORC1 QC genotype: 3730 G>A (9041); chr16:31009822; rs7294; A/G	5222
VKORC1 genotype: 2255C>T (7566); chr16:31011297; rs2359612; A/G	3477
VKORC1 QC genotype: 2255C>T (7566); chr16:31011297; rs2359612; A/G	5222
VKORC1 genotype: -4451 C>A (861); Chr16:31018002; rs17880887; A/C	5160
VKORC1 QC genotype: -4451 C>A (861); Chr16:31018002; rs17880887; A/C	5249
CYP2C9 consensus	144

```

VKORC1 -1639 consensus          1499
VKORC1 497 consensus           3817
VKORC1 1173 consensus          2109
VKORC1 1542 consensus          2179
VKORC1 3730 consensus          2356
VKORC1 2255 consensus          3164
VKORC1 -4451 consensus          4788
Therapeutic Dose of Warfarin      172
dtype: int64

```

Descriptive statistics:

	Height (cm)	Weight (kg)	Acetaminophen or Paracetamol (Tylenol)	Acetaminophen or Paracetamol (Tylenol) >1300mg/day	Was Diabetes	Simvastatin (Zocor)	Amiodarone (Cordarone)	INR	INR on Reported Therapeutic Dose of Warfarin	The
count	4554.000000	5413.000000	2274.000000	1159.000000	3283.000000	3861.000000	4182.000000	1259.000000	4968.000000	552
mean	168.047778	77.852569	0.220317	0.006040	0.187024	0.146335	0.066236	2.538324	2.364438	3
std	10.845992	21.859764	0.414551	0.077514	0.389990	0.353488	0.248724	0.198140	0.464644	1
min	124.968000	30.000000	0.000000	0.000000	0.000000	0.000000	0.000000	1.300000	0.800000	
25%	160.020000	62.000000	0.000000	0.000000	0.000000	0.000000	0.000000	2.500000	2.100000	1
50%	167.894000	75.000000	0.000000	0.000000	0.000000	0.000000	0.000000	2.500000	2.400000	2
75%	176.022000	90.000000	0.000000	0.000000	0.000000	0.000000	0.000000	2.500000	2.610000	3
max	202.000000	237.700000	1.000000	1.000000	1.000000	1.000000	1.000000	3.500000	6.100000	31

```

1 # Step 3: Data Preprocessing
2
3 # 3.1 Handle Missing Values
4 print("\n--- Handling Missing Values ---")
5
6 # First, replace any non-standard missing values with NaN
7 warfarin_data.replace(['NA', 'N/A', '-99', -99], np.nan, inplace=True)
8
9 # Check updated missing values count
10 print("Missing values after standardization:")
11 print(warfarin_data[features + [target_regression]].isnull().sum())
12
13 # 3.2 Impute missing values
14 # For numeric features
15 numeric_features = warfarin_data[features].select_dtypes(include=['int64', 'float64']).columns
16 for col in numeric_features:
17     median_value = warfarin_data[col].median()
18     warfarin_data[col].fillna(median_value, inplace=True)
19     print(f"Imputed {col} with median: {median_value}")
20
21 # For categorical features
22 categorical_features = warfarin_data[features].select_dtypes(include=['object']).columns
23 for col in categorical_features:
24     mode_value = warfarin_data[col].mode()[0]
25     warfarin_data[col].fillna(mode_value, inplace=True)
26     print(f"Imputed {col} with mode: {mode_value}")
27
28 # 3.3 Create binary classification target
29 # Following the project description, create HRD (>30 mg/wk) and LRD (<=30 mg/wk) classes
30 print("\n--- Creating Classification Target ---")
31 dose_threshold = 30 # mg/wk
32 warfarin_data['dose_class'] = (warfarin_data[target_regression] > dose_threshold).astype(int)
33 warfarin_data['dose_class'] = warfarin_data['dose_class'].map({1: 'HRD', 0: 'LRD'})
34 print(f"Dose class distribution: {warfarin_data['dose_class'].value_counts()}")
35
36 # 3.4 Check for outliers in numeric features
37 print("\n--- Checking for Outliers ---")
38 for col in numeric_features:
39     q1 = warfarin_data[col].quantile(0.25)
40     q3 = warfarin_data[col].quantile(0.75)
41     iqr = q3 - q1
42     lower_bound = q1 - 1.5 * iqr
43     upper_bound = q3 + 1.5 * iqr
44     outliers = warfarin_data[(warfarin_data[col] < lower_bound) | (warfarin_data[col] > upper_bound)][col]
45     if not outliers.empty:

```

```

46     print(f"Outliers in {col}: {len(outliers)} values")
47     # You can decide whether to cap outliers, remove them, or keep them based on domain knowledge
→ 48     dtype: int64
    Imputed Height (cm) with median: 167.89399999981867
    Imputed Weight (kg) with median: 75.0
    Imputed Acetaminophen or Paracetamol (Tylenol) with median: 0.0
    Imputed Was Dose of Acetaminophen or Paracetamol (Tylenol) >1300mg/day with median: 0.0
    Imputed Diabetes with median: 0.0
    Imputed Simvastatin (Zocor) with median: 0.0
    Imputed Amiodarone (Cordarone) with median: 0.0
    Imputed Target INR with median: 2.5
    Imputed INR on Reported Therapeutic Dose of Warfarin with median: 2.4
    Imputed Gender with mode: male
    Imputed Race (Reported) with mode: White
    Imputed Race (OMB) with mode: White
    Imputed Age with mode: 70 - 79
    Imputed Estimated Target INR Range Based on Indication with mode: 2-3
    Imputed Cyp2C9 genotypes with mode: *1/*1
    Imputed Genotyped QC Cyp2C9*2 with mode: *1/*1
    Imputed Genotyped QC Cyp2C9*3 with mode: *1/*1
    Imputed Combined QC CYP2C9 with mode: *1/*1
    Imputed VKORC1 genotype: -1639 G>A (3673); chr16:31015190; rs9923231; C/T with mode: A/A
    Imputed VKORC1 QC genotype: -1639 G>A (3673); chr16:31015190; rs9923231; C/T with mode: G/G
    Imputed VKORC1 genotype: 497T>G (5808); chr16:31013055; rs2884737; A/C with mode: T/T
    Imputed VKORC1 QC genotype: 497T>G (5808); chr16:31013055; rs2884737; A/C with mode: T/T
    Imputed VKORC1 genotype: 1173 C>T(6484); chr16:31012379; rs9934438; A/G with mode: T/T
    Imputed VKORC1 QC genotype: 1173 C>T(6484); chr16:31012379; rs9934438; A/G with mode: C/C
    Imputed VKORC1 genotype: 1542G>C (6853); chr16:31012010; rs8050894; C/G with mode: C/C
    Imputed VKORC1 QC genotype: 1542G>C (6853); chr16:31012010; rs8050894; C/G with mode: C/G
    Imputed VKORC1 genotype: 3730 G>A (9041); chr16:31009822; rs7294; A/G with mode: G/G
    Imputed VKORC1 QC genotype: 3730 G>A (9041); chr16:31009822; rs7294; A/G with mode: G/G
    Imputed VKORC1 genotype: 2255C>T (7566); chr16:31011297; rs2359612; A/G with mode: T/T
    Imputed VKORC1 QC genotype: 2255C>T (7566); chr16:31011297; rs2359612; A/G with mode: C/C
    Imputed VKORC1 genotype: -4451 C>A (861); Chr16:31018002; rs17880887; A/C with mode: C/C
    Imputed VKORC1 QC genotype: -4451 C>A (861); Chr16:31018002; rs17880887; A/C with mode: C/C
    Imputed CYP2C9 consensus with mode: *1/*1
    Imputed VKORC1 -1639 consensus with mode: A/A
    Imputed VKORC1 497 consensus with mode: T/T
    Imputed VKORC1 1173 consensus with mode: T/T
    Imputed VKORC1 1542 consensus with mode: C/C
    Imputed VKORC1 3730 consensus with mode: G/G
    Imputed VKORC1 2255 consensus with mode: T/T
    Imputed VKORC1 -4451 consensus with mode: C/C

    --- Creating Classification Target ---
Dose class distribution: dose_class
    LRD      3296
    HRD      2404
    Name: count, dtype: int64

    --- Checking for Outliers ---
    Outliers in Height (cm): 69 values
    Outliers in Weight (kg): 153 values
    Outliers in Acetaminophen or Paracetamol (Tylenol): 501 values
    Outliers in Was Dose of Acetaminophen or Paracetamol (Tylenol) >1300mg/day: 7 values
    Outliers in Diabetes: 614 values
    Outliers in Simvastatin (Zocor): 565 values
    Outliers in Amiodarone (Cordarone): 277 values
    Outliers in Target INR: 219 values
    Outliers in INR on Reported Therapeutic Dose of Warfarin: 363 values

```

```

1 # Step 4: Data Visualization and Exploratory Analysis
2
3 # 4.1 Visualize the distribution of the target variable
4 plt.figure(figsize=(10, 6))
5 plt.hist(warfarin_data[target_regression], bins=30, color='skyblue', edgecolor='black')
6 plt.title('Distribution of Therapeutic Dose of Warfarin')
7 plt.xlabel('Dose (mg/week)')
8 plt.ylabel('Frequency')
9 plt.grid(True, alpha=0.3)
10 plt.savefig('dose_distribution.png')
11 plt.close()
12
13 # 4.2 Visualize the relationship between key features and the target
14 # Select a few important numeric features
15 key_numeric_features = numeric_features[:4] # Adjust based on your data
16
17 plt.figure(figsize=(15, 10))
18 for i, feature in enumerate(key_numeric_features):
19     plt.subplot(2, 2, i+1)
20     plt.scatter(warfarin_data[feature], warfarin_data[target_regression], alpha=0.5)
21     plt.title(f'{feature} vs Therapeutic Dose')
22     plt.xlabel(feature)
23     plt.ylabel('Therapeutic Dose (mg/week)')
24     plt.grid(True, alpha=0.3)
25 plt.tight_layout()

```

```

26 plt.savefig('feature_vs_dose.png')
27 plt.close()
28
29 # 4.3 Visualize the class distribution
30 plt.figure(figsize=(8, 6))
31 warfarin_data['dose_class'].value_counts().plot(kind='bar', color=['skyblue', 'salmon'])
32 plt.title('Distribution of Dose Classes')
33 plt.xlabel('Dose Class')
34 plt.ylabel('Count')
35 plt.xticks(rotation=0)
36 plt.grid(True, alpha=0.3, axis='y')
37 plt.savefig('dose_class_distribution.png')
38 plt.close()
39
40 # 4.4 Correlation heatmap for numeric features
41 corr_matrix = warfarin_data[list(numeric_features) + [target_regression]].corr()
42 plt.figure(figsize=(12, 10))
43 sns.heatmap(corr_matrix, annot=True, cmap='coolwarm', fmt='.2f', linewidths=0.5)
44 plt.title('Correlation Matrix')
45 plt.tight_layout()
46 plt.savefig('correlation_heatmap.png')
47 plt.close()
48
49 print("Visualizations have been saved as PNG files.")

```

→ Visualizations have been saved as PNG files.

```

1 # 1. Start with cleaned features and proper target
2 X_filtered = warfarin_data.drop(columns=[target_regression, 'dose_class'] +
3                                     [col for col in warfarin_data.columns if
4                                      'dose' in col.lower() or
5                                      'therapeutic' in col.lower() or
6                                      'warfarin' in col.lower() or
7                                      'inr' in col.lower()])
8
9 # Keep only demographics and genetic features as mentioned in the project description
10 valid_features = ['Gender', 'Race (Reported)', 'Age', 'Height (cm)', 'Weight (kg)',
11                    'Diabetes', 'Simvastatin (Zocor)', 'Amiodarone (Cordarone)',
12                    'Cyp2C9 genotypes', 'VKORC1 genotype: -1639 G>A (3673); chr16:31015190; rs9923231; C/T']
13
14 X_filtered = X_filtered[[col for col in X_filtered.columns if any(feat in col for feat in valid_features)]]
15
16 # 2. Handle missing values and preprocess
17 numeric_features = X_filtered.select_dtypes(include=['int64', 'float64']).columns.tolist()
18 categorical_features = X_filtered.select_dtypes(include=['object', 'category']).columns.tolist()
19
20 # Process numeric features
21 numeric_imputer = SimpleImputer(strategy='median')
22 X_numeric = X_filtered[numeric_features].copy() if numeric_features else pd.DataFrame()
23 if not X_numeric.empty:
24     X_numeric_imputed = pd.DataFrame(
25         numeric_imputer.fit_transform(X_numeric),
26         columns=numeric_features,
27         index=X_numeric.index
28     )
29
30 # Process categorical features
31 if categorical_features:
32     cat_imputer = SimpleImputer(strategy='most_frequent')
33     X_cat = X_filtered[categorical_features].copy()
34     X_cat_imputed = pd.DataFrame(
35         cat_imputer.fit_transform(X_cat),
36         columns=categorical_features,
37         index=X_cat.index
38     )
39
40 # One-hot encode categorical features
41 X_cat_encoded = pd.get_dummies(X_cat_imputed, drop_first=True)
42
43 # Combine numeric and categorical features
44 X_processed = pd.concat([X_numeric_imputed, X_cat_encoded], axis=1) if not X_numeric.empty else X_cat_encoded
45 else:
46     X_processed = X_numeric_imputed
47
48 # 3. Define target variables
49 y_reg = warfarin_data[target_regression].copy()
50 y_cls = (y_reg > 30).astype(int) # Binary classification (>30 mg/wk = high dose)
51
52 # 4. Split data
53 X_train_reg, X_test_reg, y_train_reg, y_test_reg = train_test_split(
54     X_processed, y_reg, test_size=0.2, random_state=42)
55

```

```

56 X_train_cls, X_test_cls, y_train_cls, y_test_cls = train_test_split(
57     X_processed, y_cls, test_size=0.2, random_state=42)

1 # Step 6: Split Data for Modeling
2
3 # 6.1 Prepare data for both regression and classification tasks using processed data
4 print("\n--- Splitting Data for Modeling ---")
5
6 # Define the regression target (Update this column name if needed)
7 # You can also use: warfarin_data.columns to confirm the exact name
8 y_reg_imputed = warfarin_data['Therapeutic Dose of Warfarin'] # <-- Replace with actual column if different
9
10 # For regression task
11 X_train_reg, X_test_reg, y_train_reg, y_test_reg = train_test_split(
12     X_processed, y_reg_imputed, test_size=0.2, random_state=42
13 )
14
15 # For classification task - create the binary target first
16 dose_threshold = 30 # mg/week
17 y_cls = (y_reg_imputed > dose_threshold).astype(int)
18 print(f"Dose class distribution: {y_cls.value_counts()}")
19
20 # Split for classification
21 X_train_cls, X_test_cls, y_train_cls, y_test_cls = train_test_split(
22     X_processed, y_cls, test_size=0.2, random_state=42
23 )
24
25 # Print shapes
26 print(f"Training set size: {X_train_reg.shape[0]} samples")
27 print(f"Testing set size: {X_test_reg.shape[0]} samples")
28 print(f"Number of features: {X_train_reg.shape[1]}")
29

```



```

--- Splitting Data for Modeling ---
Dose class distribution: Therapeutic Dose of Warfarin
0    3296
1    2404
Name: count, dtype: int64
Training set size: 4560 samples
Testing set size: 1140 samples
Number of features: 47

```

```

1 # Clean target variable before modeling
2 y_reg_imputed = warfarin_data['Therapeutic Dose of Warfarin'] # Replace with actual column if different
3 non_nan_mask = y_reg_imputed.notna()
4
5 # Filter both features and target
6 X_clean = X_processed[non_nan_mask]
7 y_clean = y_reg_imputed[non_nan_mask]
8
9 # Train-test split
10 from sklearn.model_selection import train_test_split
11 X_train_reg, X_test_reg, y_train_reg, y_test_reg = train_test_split(
12     X_clean, y_clean, test_size=0.2, random_state=42
13 )
14
15
16 # Step 7: Build and Evaluate Regression Models
17
18 from sklearn.linear_model import LinearRegression, Ridge, Lasso
19 from sklearn.ensemble import RandomForestRegressor, GradientBoostingRegressor
20 from sklearn.svm import SVR
21 from sklearn.tree import DecisionTreeRegressor
22 from sklearn.metrics import mean_squared_error, mean_absolute_error, r2_score
23 from sklearn.model_selection import cross_val_score
24 import matplotlib.pyplot as plt
25 import numpy as np
26
27 print("\n--- Building Regression Models ---")
28
29 # 7.1 Define the regression models
30 regression_models = {
31     'Linear Regression': LinearRegression(),
32     'Ridge Regression': Ridge(alpha=1.0),
33     'Lasso Regression': Lasso(alpha=0.1),
34     'Decision Tree': DecisionTreeRegressor(random_state=42),
35     'Random Forest': RandomForestRegressor(n_estimators=100, random_state=42),
36     'Gradient Boosting': GradientBoostingRegressor(n_estimators=100, random_state=42),
37     'SVR': SVR(kernel='rbf')
38 }
39

```

```

40 # 7.2 Train and evaluate each model
41 results_reg = {}
42 best_model_reg = None
43 best_score_reg = float('inf') # Lower RMSE is better
44
45 for name, model in regression_models.items():
46     print(f"\nTraining {name}...")
47
48     # Train the model
49     model.fit(X_train_reg, y_train_reg)
50
51     # Make predictions
52     y_pred_train = model.predict(X_train_reg)
53     y_pred_test = model.predict(X_test_reg)
54
55     # Calculate metrics
56     train_mse = mean_squared_error(y_train_reg, y_pred_train)
57     test_mse = mean_squared_error(y_test_reg, y_pred_test)
58     train_rmse = np.sqrt(train_mse)
59     test_rmse = np.sqrt(test_mse)
60     train_mae = mean_absolute_error(y_train_reg, y_pred_train)
61     test_mae = mean_absolute_error(y_test_reg, y_pred_test)
62     train_r2 = r2_score(y_train_reg, y_pred_train)
63     test_r2 = r2_score(y_test_reg, y_pred_test)
64
65     # Cross-validation score
66     cv_scores = cross_val_score(model, X_train_reg, y_train_reg, cv=5, scoring='neg_mean_squared_error')
67     cv_rmse = np.sqrt(-np.mean(cv_scores))
68
69     # Store results
70     results_reg[name] = {
71         'Train MSE': train_mse,
72         'Test MSE': test_mse,
73         'Train RMSE': train_rmse,
74         'Test RMSE': test_rmse,
75         'Train MAE': train_mae,
76         'Test MAE': test_mae,
77         'Train R22

```

```
122 # Evaluation metrics
123 train_mse_mlp = mean_squared_error(y_train_reg, y_pred_train_mlp)
124 test_mse_mlp = mean_squared_error(y_test_reg, y_pred_test_mlp)
125 train_rmse_mlp = np.sqrt(train_mse_mlp)
126 test_rmse_mlp = np.sqrt(test_mse_mlp)
127 train_mae_mlp = mean_absolute_error(y_train_reg, y_pred_train_mlp)
128 test_mae_mlp = mean_absolute_error(y_test_reg, y_pred_test_mlp)
129 train_r2_mlp = r2_score(y_train_reg, y_pred_train_mlp)
130 test_r2_mlp = r2_score(y_test_reg, y_pred_test_mlp)
131
132 # Print results
133 print("MLP Regressor Results:")
134 print(f" Train RMSE: {train_rmse_mlp:.4f}")
135 print(f" Test RMSE: {test_rmse_mlp:.4f}")
136 print(f" Train MAE: {train_mae_mlp:.4f}")
137 print(f" Test MAE: {test_mae_mlp:.4f}")
138 print(f" Train R2: {train_r2_mlp:.4f}")
139 print(f" Test R2: {test_r2_mlp:.4f}")
140
141 # Plot predicted vs actual for test set
142 plt.figure(figsize=(6, 5))
143 plt.scatter(y_test_reg, y_pred_test_mlp, alpha=0.5)
144 plt.xlabel("Actual Dose")
145 plt.ylabel("Predicted Dose (MLP)")
146 plt.title("MLP Regressor - Actual vs Predicted")
147 plt.grid(True)
148 plt.show()
149
150
151
152
```



--- Building Regression Models ---

Training Linear Regression...

Results for Linear Regression:
Train MSE: 181.9974
Test MSE: 145.7731
Train RMSE: 13.4906
Test RMSE: 12.0737
Train MAE: 9.2498
Test MAE: 8.8517
Train R²: 0.3759
Test R²: 0.4043
CV RMSE: 13.6324

Training Ridge Regression...

Results for Ridge Regression:
Train MSE: 182.1058
Test MSE: 146.0475
Train RMSE: 13.4947
Test RMSE: 12.0850
Train MAE: 9.2565
Test MAE: 8.8588
Train R²: 0.3755
Test R²: 0.4031
CV RMSE: 13.6244

Training Lasso Regression...

Results for Lasso Regression:
Train MSE: 191.7559
Test MSE: 154.8631
Train RMSE: 13.8476
Test RMSE: 12.4444
Train MAE: 9.5593
Test MAE: 9.1486
Train R²: 0.3424
Test R²: 0.3671
CV RMSE: 13.9295

Training Decision Tree...

Results for Decision Tree:
Train MSE: 2.9176
Test MSE: 421.7495
Train RMSE: 1.7081
Test RMSE: 20.5365
Train MAE: 0.2692
Test MAE: 13.7100
Train R²: 0.9900
Test R²: -0.7236
CV RMSE: 19.4312

Training Random Forest...

Results for Random Forest:
Train MSE: 31.0276
Test MSE: 193.8273
Train RMSE: 5.5702
Test RMSE: 13.9222
Train MAE: 3.8813
Test MAE: 10.0460
Train R²: 0.8936
Test R²: 0.2079
CV RMSE: 14.7079

Training Gradient Boosting...

Results for Gradient Boosting:
Train MSE: 158.7664
Test MSE: 153.9739
Train RMSE: 12.6003
Test RMSE: 12.4086
Train MAE: 8.9534
Test MAE: 9.0369
Train R²: 0.4555
Test R²: 0.3707
CV RMSE: 13.9424

Training SVR...

Results for SVR:
Train MSE: 254.4010
Test MSE: 209.2221
Train RMSE: 15.9500
Test RMSE: 14.4645
Train MAE: 11.0463
Test MAE: 10.7077
Train R²: 0.1276
Test R²: 0.1449
CV RMSE: 15.9813

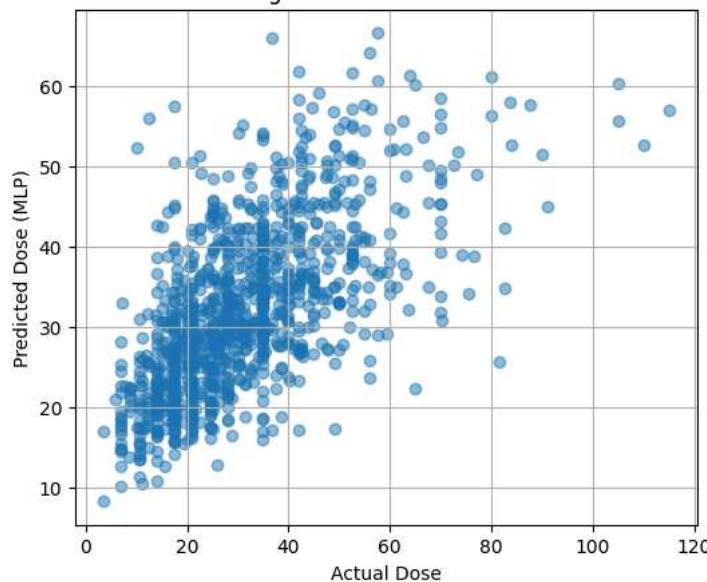
Best regression model: Linear Regression
Best RMSE: 12.0737

```
--- Training MLP Regressor (Deep Learning) ---
```

```
MLP Regressor Results:
```

```
Train RMSE: 13.4102
Test RMSE: 12.1261
Train MAE: 9.3065
Test MAE: 9.0232
Train R2: 0.3833
Test R2: 0.3991
```

MLP Regressor - Actual vs Predicted



```
1 # Step 8: Build and Evaluate Classification Models
2
3 from sklearn.linear_model import LogisticRegression
4 from sklearn.ensemble import RandomForestClassifier, GradientBoostingClassifier
5 from sklearn.svm import SVC
6 from sklearn.tree import DecisionTreeClassifier
7 from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score, roc_auc_score, confusion_matrix
8
9 # 8.1 Define the classification models
10 print("\n--- Building Classification Models ---")
11 classification_models = {
12     'Logistic Regression': LogisticRegression(random_state=42, max_iter=1000),
13     'Decision Tree': DecisionTreeClassifier(random_state=42),
14     'Random Forest': RandomForestClassifier(n_estimators=100, random_state=42),
15     'Gradient Boosting': GradientBoostingClassifier(n_estimators=100, random_state=42),
16     'SVC': SVC(probability=True, random_state=42)
17 }
18
19 # 8.2 Train and evaluate each model
20 results_cls = {}
21 best_model_cls = None
22 best_score_cls = 0 # Higher accuracy is better
23
24 for name, model in classification_models.items():
25     print(f"\nTraining {name}...")
26
27     # Train the model
28     model.fit(X_train_cls, y_train_cls)
29
30     # Make predictions
31     y_pred_train = model.predict(X_train_cls)
32     y_pred_test = model.predict(X_test_cls)
33
34     # Calculate metrics
35     train_acc = accuracy_score(y_train_cls, y_pred_train)
36     test_acc = accuracy_score(y_test_cls, y_pred_test)
37     test_precision = precision_score(y_test_cls, y_pred_test, average='weighted')
38     test_recall = recall_score(y_test_cls, y_pred_test, average='weighted')
39     test_f1 = f1_score(y_test_cls, y_pred_test, average='weighted')
40
41     # ROC AUC for binary classification
42     if hasattr(model, 'predict_proba'):
43         y_pred_proba = model.predict_proba(X_test_cls)[:, 1]
44         test_roc_auc = roc_auc_score(y_test_cls, y_pred_proba)
```

```

45     else:
46         test_roc_auc = None
47
48     # Cross-validation score
49     cv_scores = cross_val_score(model, X_train_cls, y_train_cls, cv=5, scoring='accuracy')
50     cv_acc = np.mean(cv_scores)
51
52     # Store results
53     results_cls[name] = {
54         'Train Accuracy': train_acc,
55         'Test Accuracy': test_acc,
56         'Test Precision': test_precision,
57         'Test Recall': test_recall,
58         'Test F1': test_f1,
59         'Test ROC AUC': test_roc_auc,
60         'CV Accuracy': cv_acc,
61         'model': model # Store the model object
62     }
63
64     # Print results
65     print(f"Results for {name}:")
66     for metric, value in results_cls[name].items():
67         if value is not None and metric != 'model':
68             print(f" {metric}: {value:.4f}")
69         elif metric != 'model':
70             print(f" {metric}: N/A")
71
72     # Confusion matrix
73     cm = confusion_matrix(y_test_cls, y_pred_test)
74     print(f"Confusion Matrix:\n{cm}")
75
76     # Check if this is the best model
77     if test_acc > best_score_cls:
78         best_score_cls = test_acc
79         best_model_cls = model
80         best_model_name_cls = name
81
82 print(f"\nBest classification model: {best_model_name_cls}")
83 print(f"Best accuracy: {best_score_cls:.4f}")
84
85 # 8.3 Visualize model performance comparison
86 model_names = list(results_cls.keys())
87 test_acc_values = [results_cls[name]['Test Accuracy'] for name in model_names]
88 test_f1_values = [results_cls[name]['Test F1'] for name in model_names]
89
90 plt.figure(figsize=(12, 10))
91
92 # Plot Accuracy
93 plt.subplot(2, 1, 1)
94 bars = plt.bar(model_names, test_acc_values, color='lightgreen')
95 plt.title('Test Accuracy by Model (Higher is Better)')
96 plt.ylabel('Accuracy')
97 plt.xticks(rotation=45)
98 plt.grid(axis='y', alpha=0.3)
99 for bar in bars:
100     plt.text(bar.get_x() + bar.get_width()/2, bar.get_height() + 0.02,
101             f'{bar.get_height():.2f}', ha='center')
102
103 # Plot F1 Score
104 plt.subplot(2, 1, 2)
105 bars = plt.bar(model_names, test_f1_values, color='lightblue')
106 plt.title('Test F1 Score by Model (Higher is Better)')
107 plt.ylabel('F1 Score')
108 plt.xticks(rotation=45)
109 plt.grid(axis='y', alpha=0.3)
110 for bar in bars:
111     plt.text(bar.get_x() + bar.get_width()/2, bar.get_height() + 0.02,
112             f'{bar.get_height():.2f}', ha='center')
113
114 plt.tight_layout()
115 plt.savefig('classification_model_comparison.png')
116 plt.close()
117 print("Classification model comparison plot saved as 'classification_model_comparison.png'")
```



```

1 test F1: 0.6410
2 Test ROC AUC: 0.6311
3 CV Accuracy: 0.6397
4 Confusion Matrix:
5 [[440 213]
6 [197 290]]

7 Training Random Forest...
8 Results for Random Forest:
9   Train Accuracy: 0.9917
10  Test Accuracy: 0.6982
11  Test Precision: 0.6982
12  Test Recall: 0.6982
13  Test F1: 0.6982
14  Test ROC AUC: 0.7592
15  CV Accuracy: 0.6928
16 Confusion Matrix:
17 [[481 172]
18 [172 315]]

19 Training Gradient Boosting...
20 Results for Gradient Boosting:
21   Train Accuracy: 0.7520
22   Test Accuracy: 0.7360
23   Test Precision: 0.7344
24   Test Recall: 0.7360
25   Test F1: 0.7332
26   Test ROC AUC: 0.8087
27   CV Accuracy: 0.7274
28 Confusion Matrix:
29 [[532 121]
30 [180 307]]

31 Training SVC...
32 Results for SVC:
33   Train Accuracy: 0.6520
34   Test Accuracy: 0.6561
35   Test Precision: 0.6684
36   Test Recall: 0.6561
37   Test F1: 0.6244
38   Test ROC AUC: 0.7032
39   CV Accuracy: 0.6542
40 Confusion Matrix:
41 [[582 71]
42 [321 166]]

43 Best classification model: Logistic Regression
44 Best accuracy: 0.7412
45 Classification model comparison plot saved as 'classification_model_comparison.png'

```

```

1 # Step 9: Hyperparameter Tuning for Best Models
2
3 from sklearn.model_selection import GridSearchCV
4
5 # 9.1 Tune the best regression model
6 print("\n--- Hyperparameter Tuning for Regression Model ---")
7 print(f'Tuning {best_model_name_reg}...')
8
9 # Define parameter grid based on the best model type
10 if best_model_name_reg == 'Random Forest':
11     param_grid = {
12         'n_estimators': [50, 100, 200],
13         'max_depth': [None, 10, 20, 30],
14         'min_samples_split': [2, 5, 10],
15         'min_samples_leaf': [1, 2, 4]
16     }
17 elif best_model_name_reg == 'Gradient Boosting':
18     param_grid = {
19         'n_estimators': [50, 100, 200],
20         'learning_rate': [0.01, 0.1, 0.2],
21         'max_depth': [3, 5, 7],
22         'min_samples_split': [2, 5]
23     }
24 elif best_model_name_reg in ['Ridge Regression', 'Lasso Regression']:
25     param_grid = {
26         'alpha': [0.001, 0.01, 0.1, 1.0, 10.0, 100.0]
27     }
28 elif best_model_name_reg == 'SVR':
29     param_grid = {
30         'C': [0.1, 1, 10, 100],
31         'gamma': ['scale', 'auto', 0.01, 0.1],
32         'kernel': ['linear', 'rbf']
33     }
34 else:
35     param_grid = {} # Empty grid for models without tunable hyperparameters
36
37 # Only run grid search if there are parameters to tune

```

```

38 if param_grid:
39     grid_search = GridSearchCV(best_model_reg, param_grid, cv=5,
40                             scoring='neg_mean_squared_error', n_jobs=-1)
41     grid_search.fit(X_train_reg, y_train_reg)
42
43     print(f"Best parameters: {grid_search.best_params_}")
44     print(f"Best CV score: {np.sqrt(-grid_search.best_score_):.4f} RMSE")
45
46     # Update the best model
47     best_model_reg_tuned = grid_search.best_estimator_
48
49     # Evaluate the tuned model
50     y_pred_test_tuned = best_model_reg_tuned.predict(X_test_reg)
51     test_rmse_tuned = np.sqrt(mean_squared_error(y_test_reg, y_pred_test_tuned))
52     test_r2_tuned = r2_score(y_test_reg, y_pred_test_tuned)
53
54     print(f"\nTuned model performance:")
55     print(f" Test RMSE: {test_rmse_tuned:.4f}")
56     print(f" Test R2: {test_r2_tuned:.4f}")
57
58     # Compare with the untuned model
59     print(f"\nImprovement over untuned model:")
60     print(f" RMSE improvement: {best_score_reg - test_rmse_tuned:.4f}")
61     print(f" R2 improvement: {test_r2_tuned - results_reg[best_model_name_reg]['Test R2']:.4f}")
62
63     # Update the best model
64     best_model_reg = best_model_reg_tuned
65     best_score_reg = test_rmse_tuned
66 else:
67     print(f"\nNo hyperparameter tuning needed for {best_model_name_reg}")
68
69 # 9.2 Tune the best classification model
70 print("\n--- Hyperparameter Tuning for Classification Model ---")
71 print(f"\nTuning {best_model_name_cls}...")
72
73 # Define parameter grid based on the best model type
74 if best_model_name_cls == 'Random Forest':
75     param_grid = {
76         'n_estimators': [50, 100, 200],
77         'max_depth': [None, 10, 20, 30],
78         'min_samples_split': [2, 5, 10],
79         'min_samples_leaf': [1, 2, 4]
80     }
81 elif best_model_name_cls == 'Gradient Boosting':
82     param_grid = {
83         'n_estimators': [50, 100, 200],
84         'learning_rate': [0.01, 0.1, 0.2],
85         'max_depth': [3, 5, 7],
86         'min_samples_split': [2, 5]
87     }
88 elif best_model_name_cls == 'Logistic Regression':
89     param_grid = {
90         'C': [0.001, 0.01, 0.1, 1, 10, 100],
91         'penalty': ['l2'],
92         'solver': ['lbfgs', 'liblinear']
93     }
94 elif best_model_name_cls == 'SVC':
95     param_grid = {
96         'C': [0.1, 1, 10, 100],
97         'gamma': ['scale', 'auto', 0.01, 0.1],
98         'kernel': ['linear', 'rbf']
99     }
100 else:
101     param_grid = {}
102
103 # Only run grid search if there are parameters to tune
104 if param_grid:
105     grid_search = GridSearchCV(best_model_cls, param_grid, cv=5,
106                             scoring='accuracy', n_jobs=-1)
107     grid_search.fit(X_train_cls, y_train_cls)
108
109     print(f"\nBest parameters: {grid_search.best_params_}")
110     print(f"Best CV score: {grid_search.best_score_:.4f} accuracy")
111
112     # Update the best model
113     best_model_cls_tuned = grid_search.best_estimator_
114
115     # Evaluate the tuned model
116     y_pred_test_tuned = best_model_cls_tuned.predict(X_test_cls)
117     test_acc_tuned = accuracy_score(y_test_cls, y_pred_test_tuned)
118     test_f1_tuned = f1_score(y_test_cls, y_pred_test_tuned, average='weighted')
119

```

```

120     print(f"Test Accuracy: {test_acc_tuned:.4f}")
121     print(f"Test F1 Score: {test_f1_tuned:.4f}")
122
123
124     # Compare with the untuned model
125     print(f"Improvement over untuned model:")
126     print(f"Accuracy improvement: {test_acc_tuned - best_score_cls:.4f}")
127     print(f"F1 Score improvement: {test_f1_tuned - results_cls[best_model_name_cls]['Test F1']:.4f}")
128
129     # Update the best model
130     best_model_cls = best_model_cls_tuned
131     best_score_cls = test_acc_tuned
132 else:
133     print(f"No hyperparameter tuning needed for {best_model_name_cls}")
134
135     # 9.3 Tune the Deep Learning Model (MLP Regressor)
136 print("\n--- Hyperparameter Tuning for MLP Regressor (Deep Learning) ---")
137
138 from sklearn.neural_network import MLPRegressor
139
140 # Define parameter grid
141 mlp_param_grid = {
142     'hidden_layer_sizes': [(50,), (100,), (100, 50), (128, 64)],
143     'activation': ['relu', 'tanh'],
144     'solver': ['adam'],
145     'learning_rate_init': [0.001, 0.01],
146     'max_iter': [500, 1000]
147 }
148
149 mlp_grid_search = GridSearchCV(
150     estimator=MLPRegressor(random_state=42),
151     param_grid=mlp_param_grid,
152     cv=3,
153     scoring='neg_mean_squared_error',
154     n_jobs=-1
155 )
156
157 mlp_grid_search.fit(X_train_reg, y_train_reg)
158
159 print("Best MLP Parameters:", mlp_grid_search.best_params_)
160
161 # Evaluate tuned MLP
162 best_mlp_tuned = mlp_grid_search.best_estimator_
163 y_pred_mlp_tuned = best_mlp_tuned.predict(X_test_reg)
164
165 mlp_rmse_tuned = np.sqrt(mean_squared_error(y_test_reg, y_pred_mlp_tuned))
166 mlp_r2_tuned = r2_score(y_test_reg, y_pred_mlp_tuned)
167
168 print("Tuned MLP Performance:")
169 print(f"Test RMSE: {mlp_rmse_tuned:.4f}")
170 print(f"Test R^2: {mlp_r2_tuned:.4f}")
171

```

--- Hyperparameter Tuning for Regression Model ---
Tuning Linear Regression...
No hyperparameter tuning needed for Linear Regression

--- Hyperparameter Tuning for Classification Model ---
Tuning Logistic Regression...
Best parameters: {'C': 1, 'penalty': 'l2', 'solver': 'lbfgs'}
Best CV score: 0.7357 accuracy
Tuned model performance:
Test Accuracy: 0.7412
Test F1 Score: 0.7381
Improvement over untuned model:
Accuracy improvement: 0.0000
F1 Score improvement: 0.0000

--- Hyperparameter Tuning for MLP Regressor (Deep Learning) ---
Best MLP Parameters: {'activation': 'relu', 'hidden_layer_sizes': (100,), 'learning_rate_init': 0.01, 'max_iter': 500, 'solver': 'ad

```

1 # Step 10: Feature Importance Analysis
2
3 # 10.1 Analyze feature importance for regression model
4 print("\n--- Feature Importance for Regression Model ---")
5
6 if hasattr(best_model_reg, 'feature_importances_'):
7     # For tree-based models

```

```

8 importances = best_model_reg.feature_importances_
9 indices = np.argsort(importances)[::-1]
10
11 # Print feature ranking
12 print("Feature ranking for regression model:")
13 for i, idx in enumerate(indices[:15]): # Print top 15 features
14     if i < len(X_processed.columns):
15         print(f"{i+1}. {X_processed.columns[idx]} ({importances[idx]:.4f})")
16
17 # Plot feature importances
18 plt.figure(figsize=(12, 8))
19 plt.title(f"Feature Importances for {best_model_name_reg}")
20 plt.bar(range(min(15, len(indices))), importances[indices[:15]], align="center")
21 plt.xticks(range(min(15, len(indices))), [X_processed.columns[i] for i in indices[:15]], rotation=90)
22 plt.tight_layout()
23 plt.savefig('regression_feature_importance.png')
24 plt.close()
25 print("Feature importance plot saved as 'regression_feature_importance.png'")
26
27 elif hasattr(best_model_reg, 'coef_'):
28     # For linear models
29     coefficients = best_model_reg.coef_
30     importance = np.abs(coefficients)
31     indices = np.argsort(importance)[::-1]
32
33 # Print feature ranking
34 print("Feature ranking for regression model:")
35 for i, idx in enumerate(indices[:15]): # Print top 15 features
36     if i < len(X_processed.columns):
37         print(f"{i+1}. {X_processed.columns[idx]} ({coefficients[idx]:.4f})")
38
39 # Plot feature coefficients
40 plt.figure(figsize=(12, 8))
41 plt.title(f"Feature Coefficients for {best_model_name_reg}")
42 plt.bar(range(min(15, len(indices))), coefficients[indices[:15]], align="center")
43 plt.xticks(range(min(15, len(indices))), [X_processed.columns[i] for i in indices[:15]], rotation=90)
44 plt.tight_layout()
45 plt.savefig('regression_feature_coefficients.png')
46 plt.close()
47 print("Feature coefficients plot saved as 'regression_feature_coefficients.png'")
48
49 # 10.2 Analyze feature importance for classification model
50 print("\n--- Feature Importance for Classification Model ---")
51
52 if hasattr(best_model_cls, 'feature_importances_'):
53     # For tree-based models
54     importances = best_model_cls.feature_importances_
55     indices = np.argsort(importances)[::-1]
56
57 # Print feature ranking
58 print("Feature ranking for classification model:")
59 for i, idx in enumerate(indices[:15]): # Print top 15 features
60     if i < len(X_processed.columns):
61         print(f"{i+1}. {X_processed.columns[idx]} ({importances[idx]:.4f})")
62
63 # Plot feature importances
64 plt.figure(figsize=(12, 8))
65 plt.title(f"Feature Importances for {best_model_name_cls}")
66 plt.bar(range(min(15, len(indices))), importances[indices[:15]], align="center")
67 plt.xticks(range(min(15, len(indices))), [X_processed.columns[i] for i in indices[:15]], rotation=90)
68 plt.tight_layout()
69 plt.savefig('classification_feature_importance.png')
70 plt.close()
71 print("Feature importance plot saved as 'classification_feature_importance.png'")
72
73 elif hasattr(best_model_cls, 'coef_'):
74     # For linear models
75     coefficients = best_model_cls.coef_[0] if best_model_cls.coef_.ndim > 1 else best_model_cls.coef_
76     importance = np.abs(coefficients)
77     indices = np.argsort(importance)[::-1]
78
79 # Print feature ranking
80 print("Feature ranking for classification model:")
81 for i, idx in enumerate(indices[:15]): # Print top 15 features
82     if i < len(X_processed.columns):
83         print(f"{i+1}. {X_processed.columns[idx]} ({coefficients[idx]:.4f})")
84
85 # Plot feature coefficients
86 plt.figure(figsize=(12, 8))
87 plt.title(f"Feature Coefficients for {best_model_name_cls}")
88 plt.bar(range(min(15, len(indices))), coefficients[indices[:15]], align="center")
89 plt.xticks(range(min(15, len(indices))), [X_processed.columns[i] for i in indices[:15]], rotation=90)

```

```

90     plt.tight_layout()
91     plt.savefig('classification_feature_coefficients.png')
92     plt.close()
93     print("Feature coefficients plot saved as 'classification_feature_coefficients.png''")
94
95     # 10.3 Analyze feature importance for MLP Regressor (Deep Learning)
96
97 from sklearn.inspection import permutation_importance
98
99 print("\n--- Feature Importance for MLP Regressor (Deep Learning) ---")
100
101 # Run permutation importance
102 perm_importance = permutation_importance(
103     estimator=best_mlp_tuned,           # this is the tuned MLP from Step 9.3
104     X=X_test_reg,
105     y=y_test_reg,
106     scoring='neg_mean_squared_error',
107     n_repeats=10,
108     random_state=42,
109     n_jobs=-1
110 )
111
112 # Get importances and sort
113 importances = perm_importance.importances_mean
114 indices = np.argsort(importances)[::-1]
115
116 # Print top 15 features
117 print("Feature ranking for MLP Regressor:")
118 for i, idx in enumerate(indices[:15]):
119     print(f"{i+1}. {X_processed.columns[idx]} ({importances[idx]:.4f})")
120
121 # Plot
122 plt.figure(figsize=(12, 8))
123 plt.title("Permutation Importance for MLP Regressor")
124 plt.bar(range(15), importances[indices[:15]], align="center")
125 plt.xticks(range(15), [X_processed.columns[i] for i in indices[:15]], rotation=90)
126 plt.tight_layout()
127 plt.savefig("mlp_regressor_feature_importance.png")
128 plt.close()
129 print("Permutation importance plot saved as 'mlp_regressor_feature_importance.png'")
130

```

--- Feature Importance for Regression Model ---
 Feature ranking for regression model:
 1. Cyp2C9 genotypes_*3/*3 (-23.1252)
 2. Cyp2C9 genotypes_*2/*3 (-18.4668)
 3. Cyp2C9 genotypes_*1/*14 (-17.8902)
 4. Race (Reported)_Black African (-15.9044)
 5. Cyp2C9 genotypes_*1/*13 (-15.6725)
 6. Age_90+ (-15.2158)
 7. Age_80 - 89 (-12.1793)
 8. Race (Reported)_Hispanic (-11.0533)
 9. Race (Reported)_other (-11.0229)
 10. Cyp2C9 genotypes_*2/*2 (-11.0174)
 11. Race (Reported)_Japanese (-10.7450)
 12. Age_70 - 79 (-10.6297)
 13. Race (Reported)_Han Chinese (-10.3989)
 14. VKORC1 genotype: -1639 G>A (3673); chr16:31015190; rs9923231; C/T_G/G (10.3562)
 15. Cyp2C9 genotypes_*1/*3 (-9.7862)
 Feature coefficients plot saved as 'regression_feature_coefficients.png'.

--- Feature Importance for Classification Model ---
 Feature ranking for classification model:
 1. Cyp2C9 genotypes_*3/*3 (-2.4757)
 2. Cyp2C9 genotypes_*2/*3 (-2.4307)
 3. Race (Reported)_Black other (-1.4596)
 4. Age_90+ (-1.4376)
 5. Cyp2C9 genotypes_*1/*3 (-1.3811)
 6. Cyp2C9 genotypes_*2/*2 (-1.3441)
 7. Race (Reported)_Japanese (-1.3007)
 8. Race (Reported)_Indian (1.2518)
 9. Race (Reported)_Han Chinese (-1.2369)
 10. VKORC1 genotype: -1639 G>A (3673); chr16:31015190; rs9923231; C/T_G/G (1.1361)
 11. Race (Reported)_Asian (-1.1242)
 12. Age_80 - 89 (-1.0865)
 13. Age_70 - 79 (-1.0039)
 14. Cyp2C9 genotypes_*1/*5 (-0.9940)
 15. Amiodarone (Cordarone) (-0.9700)
 Feature coefficients plot saved as 'classification_feature_coefficients.png'.

--- Feature Importance for MLP Regressor (Deep Learning) ---
 Feature ranking for MLP Regressor:
 1. VKORC1 genotype: -1639 G>A (3673); chr16:31015190; rs9923231; C/T_G/G (33.7822)
 2. Weight (kg) (26.9051)
 3. Cyp2C9 genotypes_*1/*3 (16.6198)

```

4. Race (Reported)_Japanese (14.3420)
5. Age_70 - 79 (10.5568)
6. Cyp2C9 genotypes_*2/*3 (10.3220)
7. Cyp2C9 genotypes_*1/*2 (6.5104)
8. Age_80 - 89 (6.2551)
9. Race (Reported)_Han Chinese (5.9399)
10. Amiodarone (Cordarone) (4.1803)
11. Race (Reported)_Caucasian (3.6003)
12. Cyp2C9 genotypes_*3/*3 (3.5078)
13. Age_40 - 49 (3.2009)
14. Race (Reported)_Indian (2.3737)
15. Cyp2C9 genotypes_*2/*2 (2.3593)
Permutation importance plot saved as 'mlp_regressor_feature_importance.png'

1 # Step 11: Model Visualization and Interpretation
2
3 # 11.1 Visualize regression model predictions
4 print("\n--- Visualizing Regression Model Predictions ---")
5
6 # Predict values using the best regression model
7 y_pred_reg = best_model_reg.predict(X_test_reg)
8
9 # Plot predicted vs actual values
10 plt.figure(figsize=(10, 6))
11 plt.scatter(y_test_reg, y_pred_reg, alpha=0.5)
12 plt.plot([y_test_reg.min(), y_test_reg.max()], [y_test_reg.min(), y_test_reg.max()], 'k--', lw=2)
13 plt.xlabel('Actual Therapeutic Dose')
14 plt.ylabel('Predicted Therapeutic Dose')
15 plt.title(f'Actual vs. Predicted Therapeutic Dose: {best_model_name_reg}')
16 plt.grid(True, alpha=0.3)
17 plt.savefig('regression_predictions.png')
18 plt.close()
19 print("Regression predictions plot saved as 'regression_predictions.png'")
20
21 # Plot residuals
22 plt.figure(figsize=(10, 6))
23 residuals = y_test_reg - y_pred_reg
24 plt.scatter(y_pred_reg, residuals, alpha=0.5)
25 plt.axhline(y=0, color='r', linestyle='--')
26 plt.xlabel('Predicted Therapeutic Dose')
27 plt.ylabel('Residuals')
28 plt.title(f'Residual Plot: {best_model_name_reg}')
29 plt.grid(True, alpha=0.3)
30 plt.savefig('regression_residuals.png')
31 plt.close()
32 print("Regression residuals plot saved as 'regression_residuals.png'")
33
34 # Residual histogram
35 plt.figure(figsize=(10, 6))
36 plt.hist(residuals, bins=30, alpha=0.7, color='skyblue', edgecolor='black')
37 plt.axvline(x=0, color='r', linestyle='--')
38 plt.xlabel('Residual Value')
39 plt.ylabel('Frequency')
40 plt.title('Histogram of Residuals')
41 plt.grid(True, alpha=0.3)
42 plt.savefig('regression_residuals_histogram.png')
43 plt.close()
44 print("Residuals histogram saved as 'regression_residuals_histogram.png'")
45
46 # 11.2 Visualize classification model predictions
47 print("\n--- Visualizing Classification Model Predictions ---")
48
49 # Predict values using the best classification model
50 y_pred_cls = best_model_cls.predict(X_test_cls)
51
52 # Create confusion matrix
53 cm = confusion_matrix(y_test_cls, y_pred_cls)
54
55 # Plot confusion matrix
56 plt.figure(figsize=(8, 6))
57 sns.heatmap(cm, annot=True, fmt='d', cmap='Blues', cbar=False)
58 plt.xlabel('Predicted Label')
59 plt.ylabel('True Label')
60 plt.title(f'Confusion Matrix: {best_model_name_cls}')
61 plt.savefig('classification_confusion_matrix.png')
62 plt.close()
63 print("Classification confusion matrix saved as 'classification_confusion_matrix.png'")
64
65 # Plot ROC curve for binary classification
66 if hasattr(best_model_cls, 'predict_proba'):
67     y_pred_proba = best_model_cls.predict_proba(X_test_cls)[:, 1]
68     fpr, tpr, thresholds = roc_curve(y_test_cls, y_pred_proba)
69     roc_auc = roc_auc_score(y_test_cls, y_pred_proba)

```

```

70
71 plt.figure(figsize=(8, 6))
72 plt.plot(fpr, tpr, color='darkorange', lw=2, label=f'ROC curve (area = {roc_auc:.2f})')
73 plt.plot([0, 1], [0, 1], color='navy', lw=2, linestyle='--')
74 plt.xlim([0.0, 1.0])
75 plt.ylim([0.0, 1.05])
76 plt.xlabel('False Positive Rate')
77 plt.ylabel('True Positive Rate')
78 plt.title(f'ROC Curve: {best_model_name_cls}')
79 plt.legend(loc="lower right")
80 plt.grid(True, alpha=0.3)
81 plt.savefig('classification_roc_curve.png')
82 plt.close()
83 print("ROC curve saved as 'classification_roc_curve.png'")

```

→ --- Visualizing Regression Model Predictions ---
 Regression predictions plot saved as 'regression_predictions.png'
 Regression residuals plot saved as 'regression_residuals.png'
 Residuals histogram saved as 'regression_residuals_histogram.png'

--- Visualizing Classification Model Predictions ---
 Classification confusion matrix saved as 'classification_confusion_matrix.png'
 ROC curve saved as 'classification_roc_curve.png'

```

1 # Step 12: Save Final Models
2
3 import pickle
4 import joblib
5
6 # 12.1 Create folder for saving models
7 import os
8 os.makedirs('models', exist_ok=True)
9
10 # 12.2 Save the best regression model
11 print("\n--- Saving Best Models ---")
12 joblib.dump(best_model_reg, 'models/best_regression_model.pkl')
13 print(f"Best regression model ({best_model_name_reg}) saved as 'models/best_regression_model.pkl'")
14
15 # 12.3 Save the best classification model
16 joblib.dump(best_model_cls, 'models/best_classification_model.pkl')
17 print(f"Best classification model ({best_model_name_cls}) saved as 'models/best_classification_model.pkl'")
18
19 # 12.4 Save preprocessing objects
20 joblib.dump(numeric_imputer, 'models/numeric_imputer.pkl')
21 print(f"Numeric imputer saved as 'models/numeric_imputer.pkl'")
22
23 if 'cat_imputer' in locals():
24     joblib.dump(cat_imputer, 'models/categorical_imputer.pkl')
25     print(f"Categorical imputer saved as 'models/categorical_imputer.pkl'")
26
27 # 12.5 Save feature names
28 with open('models/feature_names.pkl', 'wb') as f:
29     pickle.dump(X_processed.columns.tolist(), f)
30 print("Feature names saved as 'models/feature_names.pkl'")

```

→ --- Saving Best Models ---
 Best regression model (Linear Regression) saved as 'models/best_regression_model.pkl'
 Best classification model (Logistic Regression) saved as 'models/best_classification_model.pkl'
 Numeric imputer saved as 'models/numeric_imputer.pkl'
 Categorical imputer saved as 'models/categorical_imputer.pkl'
 Feature names saved as 'models/feature_names.pkl'

```

1 # Step 13: Create Web Application with Gradio
2
3 # 13.1 Generate Gradio application code
4 print("\n--- Creating Gradio Web Application ---")
5
6 gradio_code = """
7 import gradio as gr
8 import pandas as pd
9 import numpy as np
10 import joblib
11 import pickle
12
13 # Load the trained models and preprocessing components
14 best_reg_model = joblib.load('models/best_regression_model.pkl')
15 best_cls_model = joblib.load('models/best_classification_model.pkl')
16 numeric_imputer = joblib.load('models/numeric_imputer.pkl')
17
18 # Load feature names

```

```

19 with open('models/feature_names.pkl', 'rb') as f:
20     feature_names = pickle.load(f)
21
22 def predict_warfarin_dose(gender, race, age, height, weight, diabetes,
23                           simvastatin, amiodarone, cyp2c9, vkorc1):
24     # Create input data dictionary
25     input_data = {
26         'Gender': [gender],
27         'Race (Reported)': [race],
28         'Age': [age],
29         'Height (cm)': [height],
30         'Weight (kg)': [weight],
31         'Diabetes': [1 if diabetes else 0],
32         'Simvastatin (Zocor)': [1 if simvastatin else 0],
33         'Amiodarone (Cordarone)': [1 if amiodarone else 0],
34         'Cyp2C9 genotypes': [cyp2c9],
35         'VKORC1 genotype: -1639 G>A (3673); chr16:31015190; rs9923231; C/T': [vkorc1]
36     }
37
38     # Convert to DataFrame
39     input_df = pd.DataFrame(input_data)
40
41     # Preprocess input
42     # One-hot encode categorical variables
43     cat_cols = [col for col in input_df.columns if input_df[col].dtype == 'object']
44     num_cols = [col for col in input_df.columns if col not in cat_cols]
45
46     # Handle numeric columns
47     if num_cols:
48         input_df[num_cols] = numeric_imputer.transform(input_df[num_cols])
49
50     # Handle categorical columns
51     if cat_cols:
52         input_df_encoded = pd.get_dummies(input_df, columns=cat_cols)
53     else:
54         input_df_encoded = input_df.copy()
55
56     # Align columns with training data
57     final_df = pd.DataFrame(0, index=[0], columns=feature_names)
58     for col in input_df_encoded.columns:
59         if col in feature_names:
60             final_df[col] = input_df_encoded[col].values
61
62     # Make predictions
63     dose_prediction = best_reg_model.predict(final_df)[0]
64     class_prediction = best_cls_model.predict(final_df)[0]
65
66     # Get probability for classification
67     class_proba = best_cls_model.predict_proba(final_df)[0]
68     high_dose_prob = class_proba[1] if len(class_proba) > 1 else 0
69
70     # Format results
71     class_label = "High Dose (>30 mg/week)" if class_prediction == 1 else "Low Dose (≤30 mg/week)"
72
73     # Create detailed results
74     results = f"## Warfarin Dosing Prediction Results\n\n"
75     results += f"### Predicted Weekly Dose\n"
76     results += f"**{dose_prediction:.1f} mg/week**\n\n"
77     results += f"### Dose Classification\n"
78     results += f"**{class_label}** (Confidence: {high_dose_prob*100:.1f}%)"
79     results += f"### Patient Information\n"
80     results += f"- Gender: {gender}\n"
81     results += f"- Race: {race}\n"
82     results += f"- Age: {age} years\n"
83     results += f"- Height: {height} cm\n"
84     results += f"- Weight: {weight} kg\n"
85     results += f"- Diabetes: {'Yes' if diabetes else 'No'}\n"
86     results += f"- Taking Simvastatin: {'Yes' if simvastatin else 'No'}\n"
87     results += f"- Taking Amiodarone: {'Yes' if amiodarone else 'No'}\n"
88     results += f"- CYP2C9 Genotype: {cyp2c9}\n"
89     results += f"- VKORC1 Genotype: {vkorc1}\n\n"
90     results += f"### Clinical Considerations\n"
91     results += f"- This prediction is based on a machine learning model and should be used as a decision support tool only.\n"
92     results += f"- Regular INR monitoring is still necessary regardless of the predicted dose.\n"
93     results += f"- Clinical judgment should always prevail when determining the appropriate warfarin dose.\n"
94
95     return results
96
97 # Create Gradio interface
98 with gr.Blocks(title="Warfarin Dosing Prediction Tool") as demo:
99     gr.Markdown("# Warfarin Dosing Prediction Tool")
100    gr.Markdown("Enter patient information to predict the optimal warfarin dosage")

```

```

101
102     with gr.Row():
103         with gr.Column():
104             # Demographics
105             gr.Markdown("### Patient Demographics")
106             gender = gr.Dropdown(
107                 ["Male", "Female"],
108                 label="Gender"
109             )
110             race = gr.Dropdown(
111                 ["White", "Asian", "Black or African American", "Unknown"],
112                 label="Race"
113             )
114             age = gr.Slider(
115                 minimum=0, maximum=100, step=1, value=60,
116                 label="Age (years)"
117             )
118             height = gr.Number(
119                 value=170,
120                 label="Height (cm)"
121             )
122             weight = gr.Number(
123                 value=70,
124                 label="Weight (kg)"
125             )
126
127             # Medical information
128             gr.Markdown("### Medical Information")
129             diabetes = gr.Checkbox(
130                 label="Diabetes"
131             )
132             simvastatin = gr.Checkbox(
133                 label="Taking Simvastatin (Zocor)"
134             )
135             amiodarone = gr.Checkbox(
136                 label="Taking Amiodarone (Cordarone)"
137             )
138
139             # Genetic information
140             gr.Markdown("### Genetic Information")
141             cyp2c9 = gr.Dropdown(
142                 ["*1/*1", "*1/*2", "*1/*3", "*2/*2", "*2/*3", "*3/*3"],
143                 label="CYP2C9 Genotype"
144             )
145             vkorc1 = gr.Dropdown(
146                 ["G/G", "A/G", "A/A"],
147                 label="VKORC1 Genotype"
148             )
149
150             # Submit button
151             submit_btn = gr.Button("Predict Warfarin Dose", variant="primary")
152
153             with gr.Column():
154                 # Output
155                 output = gr.Markdown()
156
157             # Set default values
158             gender.value = "Male"
159             race.value = "White"
160             cyp2c9.value = "*1/*1"
161             vkorc1.value = "G/G"
162
163             # Connect the function
164             submit_btn.click(
165                 predict_warfarin_dose,
166                 inputs=[gender, race, age, height, weight, diabetes, simvastatin, amiodarone, cyp2c9, vkorc1],
167                 outputs=output
168             )
169
170             gr.Markdown("## About This Model")
171             gr.Markdown('''
172             This tool uses a machine learning model trained on data from the International Warfarin Pharmacogenetics Consortium.
173             The model considers patient demographics, medical conditions, and genetic factors to predict the optimal warfarin dose.
174
175             **Important Disclaimer**: This prediction tool is for educational and research purposes only.
176             It should not replace professional medical advice, diagnosis, or treatment. Always consult with a healthcare provider
177             for medical decisions regarding warfarin dosing.
178             ''')
179
180             # For running the app locally
181             if __name__ == "__main__":
182                 demo.launch()

```