

DTSA 5509 - Introduction to Machine Learning Final Project: Supervised Learning with Alzheimer's Disease Datasets

Israel Johnson
June 25, 2023



University of Colorado **Boulder**

PROBLEM STATEMENT

- There has been ongoing research in the area of developing artificial intelligence (AI) and machine learning (ML) applications for Alzheimer's disease and dementia. These applications range from detecting early stages of Alzheimer's disease (and dementia) to measuring the impact of the disease in the presence of other diseases.
- Researchers are still exploring different ML model approaches for analyzing datasets related to Alzheimer's disease and dementia, seeking which are the best models to apply in given scenarios related to the disease.



University of Colorado **Boulder**

PROJECT OBJECTIVE

- The goal for this project is to apply model analysis for two different scenarios, where both scenarios relate to datasets about Alzheimer's disease and dementia
- The two scenarios are as follows:
 - Determine if Alzheimer's disease is an effective feature with data related to the COVID19 pandemic, when compared to other disease features that are present. A multilinear regression model is used for model analysis in this scenario.
 - Determine if whether or not there is a strong correlation between age and gender features among patients who have Alzheimer's disease or dementia. Also, in either case of correlation conditions, find ways to improve model accuracy and performance, if necessary. For the model approach for this case, a logistic regression model is applied to the corresponding data.



University of Colorado **Boulder**

DATA SOURCES

- 1.) U.S Department of Health and Human Services, Centers for Disease Control and Prevention (2021). Monthly Counts of Deaths By Select Causes (2020-2021) [Data set]. Retrieved from <https://catalog.data.gov/dataset/monthly-counts-of-deaths-by-select-causes-2020-2021-2785a>
- 2.) U.S Department of Health and Human Services, Centers for Disease Control and Prevention (2021). Conditions Contributing to Deaths Involving Coronavirus Disease 2019 (COVID-19) by Age Group [Data set]. Retrieved from <https://catalog.data.gov/dataset/conditions-contributing-to-deaths-involving-coronavirus-disease-2019-covid-19-by-age-group>
- 3.) Open Access Series of Imaging Studies (OASIS) Brains Project. OASIS: Cross-sectional MRI Data in Young, Middle-Aged, Nondemented, and Demented Older Adults (2007) [Data set]. Retrieved from <https://www.oasis-brains.org/#oasis1>
- 4.) Open Access Series of Imaging Studies (OASIS) Brains Project. OASIS: Longitudinal MRI Data in Nondemented and Demented Older Adults (2010) [Data set]. Retrieved from <https://www.oasis-brains.org/#oasis2>
- 5.) Dinner, Baris (2022). Alzheimer Features (2022) [Data set]. Retrieved from <https://www.kaggle.com/datasets/brsdincer/alzheimer-features>

GITHUB REPOSITORY SOURCE

- **Source:** https://github.com/IsraelsLibrary/DTSA_5509_Intro_to_Machine_Learning



University of Colorado **Boulder**

DATA CLEANING PROCESS

- For data in scenario one, various column features were found to contain null values. A four percent null value threshold was established to limit the number of null values that were permitted. Any column features that exceeded the given threshold were removed from the dataset.
- In preparation for scenario two, two ‘helper’ functions were created to drop unrelevant columns from associated data as well as remove any row features that contained null values. Unlike scenario one where data cleaning was done for one data source, scenario two involved four different sources, which required more effort when it came to the data cleaning process. The data cleaning process for scenario two occurs after further data transformations are made (extracting key features from the different data sources and combining them into one dataset).



University of Colorado **Boulder**

DATA CLEANING PROCESS

```
1 threshold = int(0.04 * len(monthly_provisional_data))
2
3 features_to_impute = [col for col in monthly_provisional_data.columns if
4                         monthly_provisional_data[col].isnull().sum() != 0 and
5                         monthly_provisional_data[col].isnull().sum() <= threshold]
6 features_to_throw = [col for col in monthly_provisional_data.columns if
7                         monthly_provisional_data[col].isnull().sum() > threshold]
8
9 old_monthly_provisional_data = monthly_provisional_data
10
11 monthly_provisional_data.drop(columns=features_to_throw, axis=1, inplace=True)
```



University of Colorado **Boulder**

DATA CLEANING PROCESS

```
1 # 'Helper' functions that perform data cleaning for the datasets in scenarios one and two.  
2  
3 def drop_columns(columns_to_drop, dataset):  
4     for col in columns_to_drop:  
5         if col in dataset.columns:  
6             dataset = dataset.drop(columns=col, axis=1)  
7     return dataset  
8  
9 def clean_data(dataset):  
10    dataset = dataset.dropna()  
11    if dataset.isnull().sum().sum() > 0:  
12        clean_data(dataset)
```



University of Colorado **Boulder**

DATA CLEANING PROCESS

```
1 # 'Helper' function that forms the final version of the dataset for scenario 2.
2 def transform_dataset2(datasets, limit):
3     new_dataset = {}
4
5     new_dataset['age'] = []
6     new_dataset['sex'] = []
7
8     for dataset in datasets.values():
9         new_dataset['age'].extend(dataset[col].tolist()[:limit] for col in dataset if col in 'Age')
10        new_dataset['sex'].extend(dataset[col].tolist()[:limit] for col in dataset if col in 'M/F')
11
12    all_ages = []
13    genders = []
14    group = []
15
16    for age_, sex_ in zip(new_dataset['age'], new_dataset['sex']):
17        all_ages += age_
18        genders += sex_
19
20    for ind in range(len(genders)):
21        if genders[ind] == 'M':
22            genders[ind] = 0
23        else:
24            genders[ind] = 1
25
26    new_dataset['age'] = all_ages
27    new_dataset['sex'] = genders
28
29    return pd.DataFrame.from_dict(new_dataset)
30
```



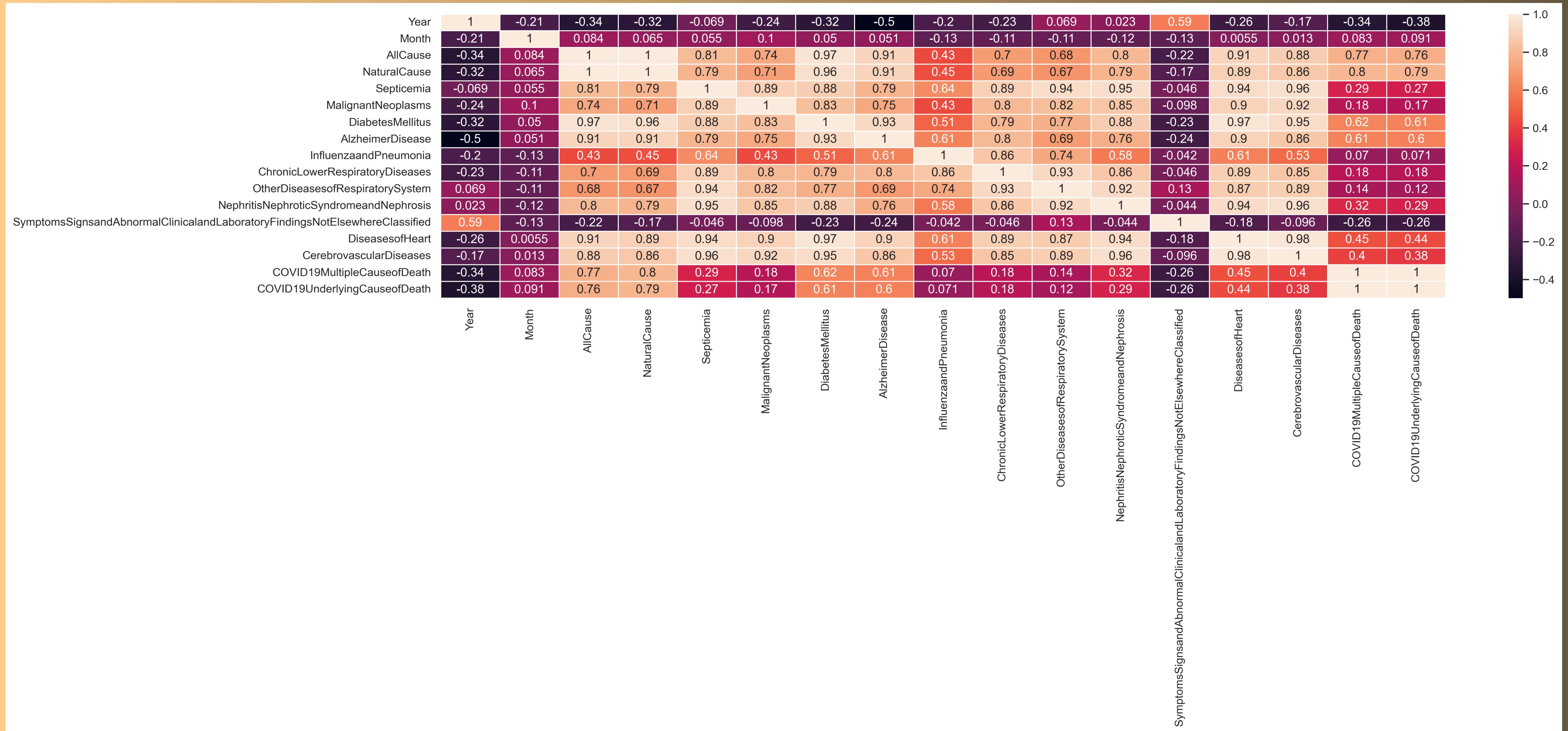
University of Colorado **Boulder**

EXPLORATORY DATA ANALYSIS (SCENARIO ONE)

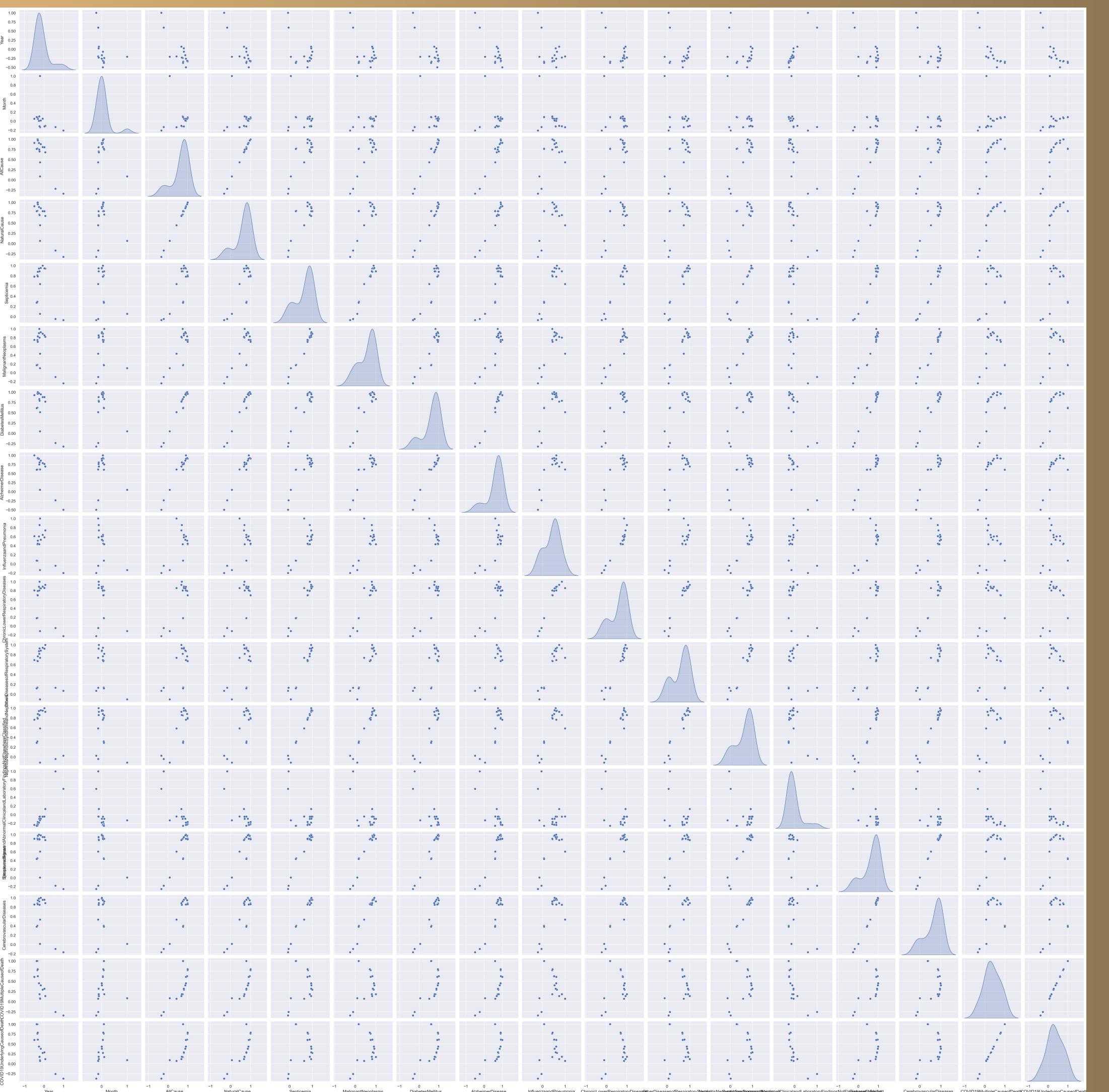
- The first stage of exploratory data analysis (EDA) in scenario one was to generate a correlation matrix to see where strong correlation and multicollinearity can be found.
- Along with a correlation matrix, a series of pair plots were also created to explore correlation and collinearity among features



University of Colorado **Boulder**



University of Colorado **Boulder**



University of Colorado **Boulder**

EDA SUMMARY AND RESULTS (SCENARIO ONE)

- Strong correlation and multicollinearity was found between the features related to types of causes of death and the features related to different disease types. This conclusion serves as the basis for the multilinear regression model approach and further analysis.



University of Colorado **Boulder**

MULTILINEAR REGRESSION MODEL (SCENARIO ONE)

```
1 # Creating the first iteration of the Multilinear Regression model with 'AllCause' as the predicting label.
2 x = monthly_provisional_data[['MalignantNeoplasms', 'DiabetesMellitus','DiseasesofHeart','AlzheimerDisease',
3                               'InfluenzaandPneumonia', 'ChronicLowerRespiratoryDiseases', 'OtherDiseasesofRespirator
4                               'DiseasesofHeart']]
5 y = monthly_provisional_data[['AllCause']]
6 regr = linear_model.LinearRegression()
7 regr.fit(x, y)
8
9 # with statsmodels
10 x = sm.add_constant(x) # adding a constant
11
12 model = sm.OLS(y, x).fit()
13 predictions = model.predict(x)
```



University of Colorado **Boulder**

EVALUATION METRICS (SCENARIO ONE)

Out[71]: OLS Regression Results											
Dep. Variable:		AllCause	R-squared:		0.966						
Model:		OLS	Adj. R-squared:		0.958						
Method:		Least Squares	F-statistic:		136.3						
Date:		Sun, 25 Jun 2023	Prob (F-statistic):		5.63e-23						
Time:		08:26:55	Log-Likelihood:		-439.56						
No. Observations:		42	AIC:		895.1						
Df Residuals:		34	BIC:		909.0						
Df Model:		7									
Covariance Type: nonrobust											
		coef	std err	t	P> t	[0.025	0.975]				
	const	7829.8612	1.41e+04	0.554	0.583	-2.09e+04	3.65e+04				
	MalignantNeoplasms	-0.7712	0.969	-0.796	0.432	-2.740	1.197				
	DiabetesMellitus	49.9709	9.938	5.028	0.000	29.775	70.167				
	DiseasesofHeart	-2.0285	1.151	-1.762	0.087	-4.368	0.311				
	AlzheimerDisease	7.0316	3.323	2.116	0.042	0.279	13.785				
	InfluenzaandPneumonia	-4.7917	3.937	-1.217	0.232	-12.792	3.208				
	ChronicLowerRespiratoryDiseases	0.3633	5.382	0.068	0.947	-10.574	11.301				
	OtherDiseasesofRespiratorySystem	16.7025	9.921	1.683	0.101	-3.460	36.865				
	DiseasesofHeart	-2.0285	1.151	-1.762	0.087	-4.368	0.311				
Omnibus:		1.254	Durbin-Watson:		1.389						
Prob(Omnibus):		0.534	Jarque-Bera (JB):		1.117						
Skew:		0.382	Prob(JB):		0.572						
Kurtosis:		2.764	Cond. No.		1.43e+18						

Out[18]: OLS Regression Results											
Dep. Variable:		COVID19MultipleCauseofDeath	R-squared:		0.883						
Model:		OLS	Adj. R-squared:		0.859						
Method:		Least Squares	F-statistic:		36.65						
Date:		Sun, 25 Jun 2023	Prob (F-statistic):		4.93e-14						
Time:		15:27:55	Log-Likelihood:		-441.40						
No. Observations:		42	AIC:		898.8						
Df Residuals:		34	BIC:		912.7						
Df Model:		7									
Covariance Type: nonrobust											
		coef	std err	t	P> t	[0.025	0.975]				
	const	1.583e+04	1.48e+04	1.073	0.291	-1.42e+04	4.58e+04				
	MalignantNeoplasms	-3.3477	1.012	-3.308	0.002	-5.405	-1.291				
	DiabetesMellitus	48.2608	10.383	4.648	0.000	27.160	69.362				
	DiseasesofHeart	-2.8382	1.203	-2.360	0.024	-5.282	-0.394				
	AlzheimerDisease	7.2680	3.472	2.093	0.044	0.212	14.324				
	InfluenzaandPneumonia	-6.3079	4.113	-1.534	0.134	-14.666	2.050				
	ChronicLowerRespiratoryDiseases	1.1252	5.623	0.200	0.843	-10.302	12.553				
	OtherDiseasesofRespiratorySystem	9.8939	10.366	0.954	0.347	-11.172	30.960				
	DiseasesofHeart	-2.8382	1.203	-2.360	0.024	-5.282	-0.394				
Omnibus:		1.176	Durbin-Watson:		1.306						
Prob(Omnibus):		0.555	Jarque-Bera (JB):		1.190						
Skew:		0.347	Prob(JB):		0.552						
Kurtosis:		2.556	Cond. No.		1.43e+18						



University of Colorado **Boulder**

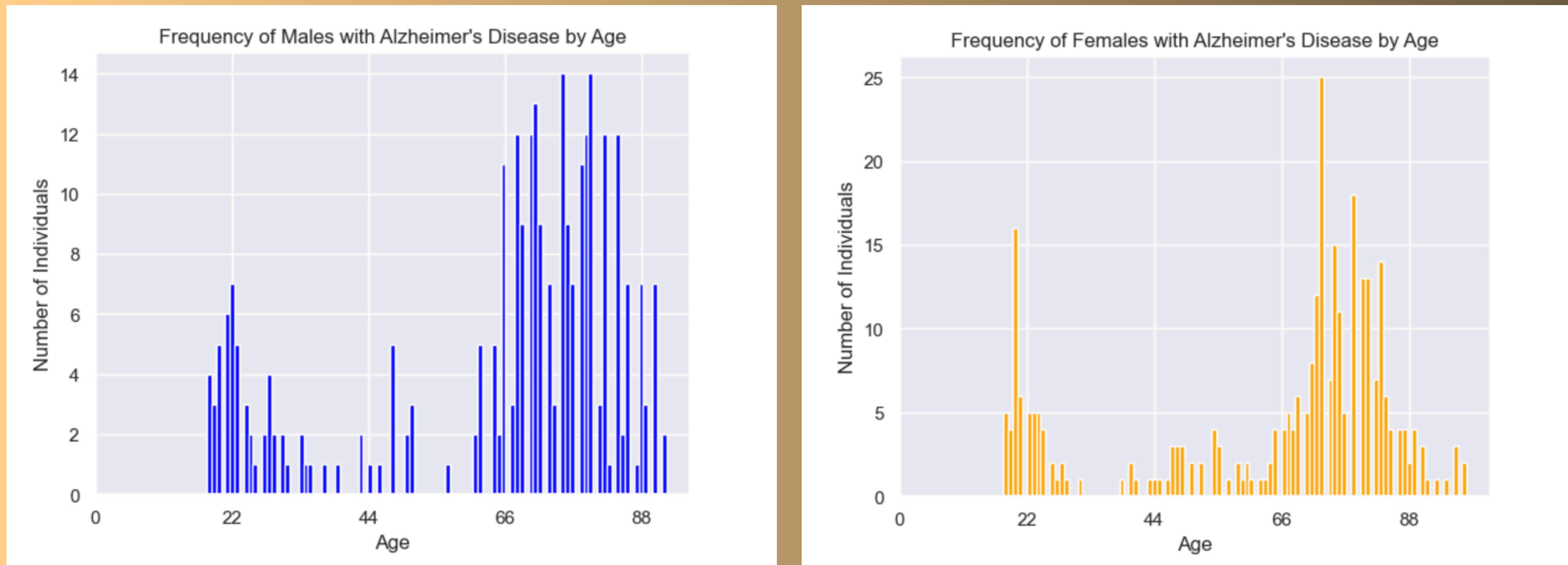
EXPLORATORY DATA ANALYSIS (SCENARIO TWO)

- Selected features for the second scenario of the project are as follows: ‘age’ and ‘gender’.
- Generated histograms show the frequency of male and female patients with Alzheimer’s disease among different ages.



University of Colorado **Boulder**

EXPLORATORY DATA ANALYSIS (SCENARIO TWO)



University of Colorado **Boulder**

LOGISTIC REGRESSION MODEL (SCENARIO TWO)

```
In [48]: 1 # Forms the first iteration of the logistic regression for scenario 2.  
2  
3 X = scenario2_data.iloc[:, :-1].values  
4 y = scenario2_data.iloc[:, -1].values  
5 x_train, x_test, y_train, y_test = model_selection.train_test_split(X,y,test_size=0.2, random_state=5)  
6  
7 LogReg = LogisticRegression(solver='liblinear')  
8 LogReg.fit(x_train, y_train)
```

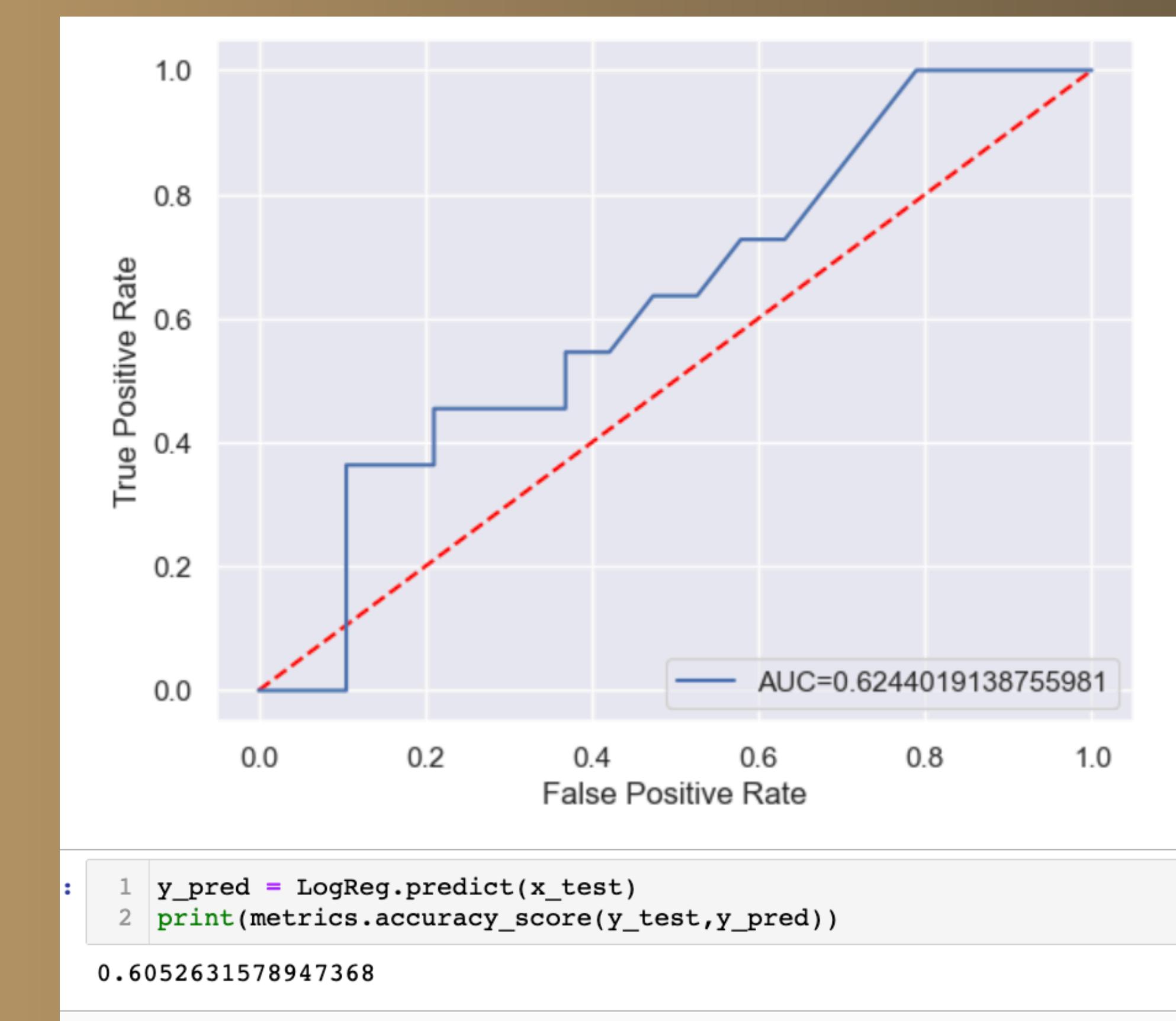
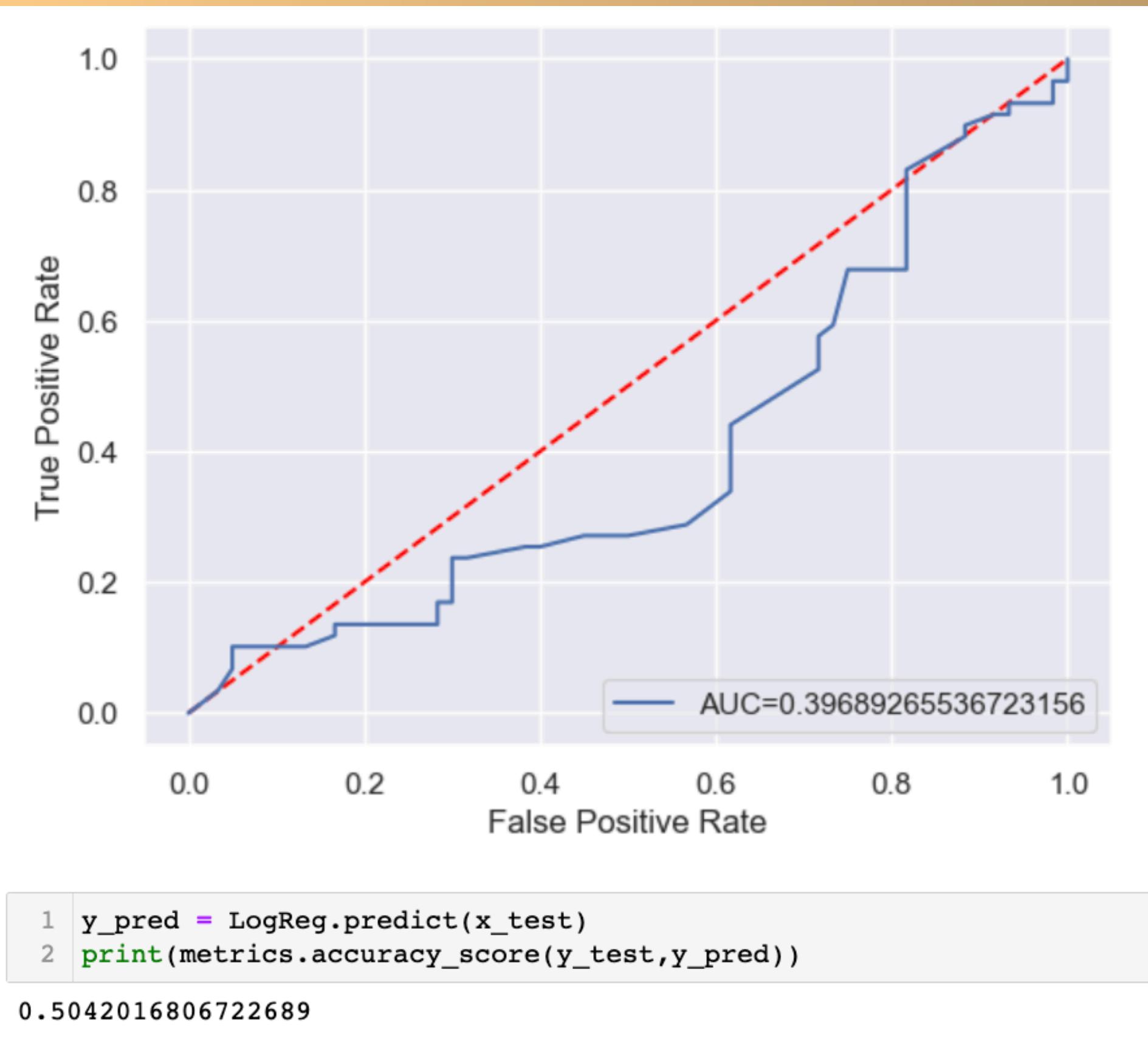
Out[48]:

```
LogisticRegression  
LogisticRegression(solver='liblinear')
```



University of Colorado **Boulder**

EVALUATION METRICS (SCENARIO TWO)



University of Colorado **Boulder**

RESULTS AND DISCUSSION

- For scenario one, the following column features were selected as dependent variables for the multilinear regression model: ‘AllCause’, ‘COVID19MultipleCauseofDeath’, and ‘COVIDUnderlyingCauseofDeath’. As a column feature, Alzheimer’s disease was found to be significant for two of these variables: ‘AllCause’ (category for all causes of death) and ‘COVID19MultipleCauseofDeath’ (category for multiple causes leading to COVID19 deaths). The goal for scenario one was achievable.
- In scenario two, the selected features (‘age’ and ‘gender’) were found to be not well correlated. Despite this, model analysis revealed effective methods for necessary improvements for model performance and accuracy: dataset pruning and increasing random shuffling on input data. As a result, concluding AUC metrics and accuracy increased for this scenario. The initial goals for scenario two were achievable.



University of Colorado **Boulder**

LESSONS LEARNED

- Data cleaning was revealed to be a critical aspect for the project.
- Repeating scenarios with different datasets can reveal different correlation results for analysis.
- Testing the same machine learning models with larger scales of data can determine robustness of the given models.



University of Colorado **Boulder**

CONCLUSION

- In conclusion, multilinear regression models were found to be an effective model approach in identifying Alzheimer's disease as a significant feature in the presence of other diseases. This model approach can be beneficial for other similar scenarios where multiple diseases are involved.
- For scenario two, even though strong correlation was not found, dataset pruning and increased random shuffling of input data were found to be effective methods when improving accuracy of the given logistic regression model. Future tests will need to be conducted, and repeating the scenario with different features would be beneficial in finding stronger correlation and collinearity.



University of Colorado **Boulder**

THANK YOU



University of Colorado **Boulder**