

# CS 536: Final Project - Data Completion and Interpolation

Name:- **Maresh Reddy Annapureddy**

Netid:-**ms1700**

## 1. Dataset and handling

In [1]:

```
import numpy as np
import pandas as pd
import random
import matplotlib.pyplot as plt
```

The dataset is **HCVdata.csv** from UCI machine learning repository. The data set contains laboratory values of blood donors and Hepatitis C patients and demographic values like age. It contains totally 114 attributes and 615 instances. All attributes except for sex and category are numerical. In this project, data interpolation is purely concentrated on continuous values. so, categorical attributes i.e, 'index', 'Category', 'Age', 'sex' are dropped from the dataset. Now, dataset contains only continuous numerical values i.e, 'ALB', 'ALP', 'ALT', 'AST', 'BIL', 'CHE', 'CHOL', 'CREA', 'GGT', 'PROT'. Below is the code to drop categorical columns.

In [2]:

```
data = pd.read_csv("C:/Users/vrma/Documents/2nd semester/Machine Learning/hcvdata.csv")
data=data.drop(columns=['index','Category','Age','Sex'])
print(data)

ALB    ALP    ALT    AST    BIL    CHE    CHOL    CREA    GGT    PROT
0    38.5    52.5    7.7    22.1    7.5    6.93    3.23    106.0    12.1    69.0
1    38.5    70.3    18.0    24.7    3.9    11.17    4.80    74.0    15.6    76.5
2    46.9    74.7    36.2    52.6    6.1    8.84    5.20    86.0    33.2    79.3
3    43.2    52.0    30.6    22.6    18.9    7.33    4.74    80.0    33.8    75.7
4    39.2    74.1    32.6    24.8    9.6    9.15    4.32    76.0    29.9    68.7
..    ...    ...    ...    ...    ...    ...    ...    ...    ...
610    32.0    416.6    5.9    110.3    50.0    5.57    6.30    55.7    650.9    68.5
611    24.0    102.8    2.9    44.4    20.0    1.54    3.02    63.0    35.9    71.3
612    29.0    87.3    3.5    99.0    48.0    1.66    3.63    66.7    64.2    82.0
613    33.0    NaN    39.0    62.0    20.0    3.56    4.20    52.0    50.0    71.0
614    36.0    NaN    100.0    80.0    12.0    9.07    5.30    67.0    34.0    68.0

[615 rows x 10 columns]
```

## Removing values randomly from dataset:-

As the chosen dataset does not have missing values, randomly values are deleted from the dataset. To this, function named **remove\_features()** is defined. This function deletes given percent of data for the given attributes of the dataset. Below is the code for **remove\_features()** function.

In [16]:

```
def remove_features(data,remove_perc,cols):
    data1=data.copy()
    n_samples,n_features = data1.shape
    random_columns=random.sample(range(n_features),cols)
    for col in random_columns:
        no_samples_to_remove = random.sample(range(n_samples),int(remove_perc*n_samples/100))
        data1.iloc[:,col].where(data1.iloc[:,col].isnull() < 0 , inplace = True)
    return data1
```

In [17]:

```
true_data=data.copy()
new_data=remove_features(data,30,5)
print(data)
print(new_data)

ALB    ALP    ALT    AST    BIL    CHE    CHOL    CREA    GGT    PROT
0    38.5    52.5    7.7    22.1    7.5    6.93    3.23    106.0    12.1    69.0
1    38.5    70.3    18.0    24.7    3.9    11.17    4.80    74.0    15.6    76.5
2    46.9    74.7    36.2    52.6    6.1    8.84    5.20    86.0    33.2    79.3
3    43.2    52.0    30.6    22.6    18.9    7.33    4.74    80.0    33.8    75.7
4    39.2    74.1    32.6    24.8    9.6    9.15    4.32    76.0    29.9    68.7
..    ...    ...    ...    ...    ...    ...    ...    ...    ...
610    32.0    416.6    5.9    110.3    50.0    5.57    6.30    55.7    650.9    68.5
611    24.0    102.8    2.9    44.4    20.0    1.54    3.02    63.0    35.9    71.3
612    29.0    87.3    3.5    99.0    48.0    1.66    3.63    66.7    64.2    82.0
613    33.0    NaN    39.0    62.0    20.0    3.56    4.20    52.0    50.0    71.0
614    36.0    NaN    100.0    80.0    12.0    9.07    5.30    67.0    34.0    68.0

[615 rows x 10 columns]
```

To find if the which features are randomly removed find **missing\_features()** function is defined. And find **missing\_columns()** function return the attributes that have missing values and in the descending order of number of missing values. And this is all about the handling of the data

In [8]:

```
def find_missing_features(data):
    data=data.copy()
    n_samples,n_features = data.shape
    missing_column= data.isna()
    missing_columns=missing_column.any()
    missing_columns=[]
    for i in range(n_features):
        if missing_column[i]:
            missing_columns.append(i)
    return missing_columns

def find_missing_columns(data):
    data1=data.copy()
    missing_features = find_missing_features(data1)
    n_samples,n_features = data1.shape
    zero_features_num = new_data.isnull()
    zero_features_num = zero_features_num.sum(axis=0)
    missing_features_column= zero_features_num.nlargest(len(missing_features))
    missing_columns = list(missing_features_column.index.values)
    missing_columns.reverse()
    print(missing_columns)
    return missing_columns
```

In [9]:

```
k=find_missing_features(new_data)
print(k)
k=find_missing_columns(new_data)
```

```
[0, 1, 2, 4, 6, 8, 9]
```

```
['ALB', 'ALP', 'GGT', 'BIL', 'PROT', 'ALT', 'CHOL']
```

## 2. Naive method:-

Here in this method, each missing features columns mean is calculated and replaced the missing values with the mean. For this **mean\_method()** function is defined. It takes the data and returns the total filled data. Here **mean\_squared method** is defined which calculates mean squared error of the predicted data and original data.

In [10]:

```
def mean_squared_error(data_1,data_2,perc):
    k=(data_1-data_2)**2).sum()
    n=len(data_1)*100
    k=k/n
    return k

def mean_method(data):
    data1=data.copy()
    n_samples,n_features = data_1.shape
    missing_features=find_missing_features(data_1)
    for i in missing_features:
        mean = data_1.iloc[:, i].mean()
        mean = round(mean,2)
        data_1.iloc[:, i].fillna(mean, inplace=True)
    return data_1
```

In [18]:

```
data=true_data.copy()
new_data=remove_features(data,30,10)
mean_data=mean_method(new_data)
print(mean_data)
perc=10
mse_all=[]
percent=[]
print(true_data)
for i in range(9):
    percent.append(perc)
    new_data=remove_features(data,perc,5)
    mean_data=mean_method(new_data)
    k=mean_squared_error(true_data,mean_data,perc)
    k=list(k)
    #print(k)
    print("percent of the data is missing:",perc-10)
    mse_all.append(sum(k)/10)
    #print(k)
    print(percent)
    plt.plot(percent,mse_all)
    plt.xlabel("Percent of data removed")
    plt.ylabel("Mean squared error")
plt.title("Mean squared error as a function of percent of data removed over all the data")
```

```
ALB    ALP    ALT    AST    BIL    CHE    CHOL    CREA    GGT    PROT
0    38.5    52.5    7.7    22.1    7.5    6.93    3.23    106.0    12.1    69.0
1    38.5    70.3    18.0    24.7    3.9    11.17    4.80    74.0    15.6    76.5
2    46.9    74.7    36.2    52.6    6.1    8.84    5.20    86.0    33.2    79.3
3    43.2    52.0    30.6    22.6    18.9    7.33    4.74    80.0    33.8    75.7
4    39.2    74.1    32.6    24.8    9.6    9.15    4.32    76.0    29.9    68.7
..    ...    ...    ...    ...    ...    ...    ...    ...    ...
610    32.0    416.6    5.9    110.3    50.0    5.57    6.30    55.7    650.9    68.5
611    24.0    102.8    2.9    44.4    20.0    1.54    3.02    63.0    35.9    71.3
612    29.0    87.3    3.5    99.0    48.0    1.66    3.63    66.7    64.2    82.0
613    33.0    NaN    39.0    62.0    20.0    3.56    4.20    52.0    50.0    71.0
614    36.0    NaN    100.0    80.0    12.0    9.07    5.30    67.0    34.0    68.0

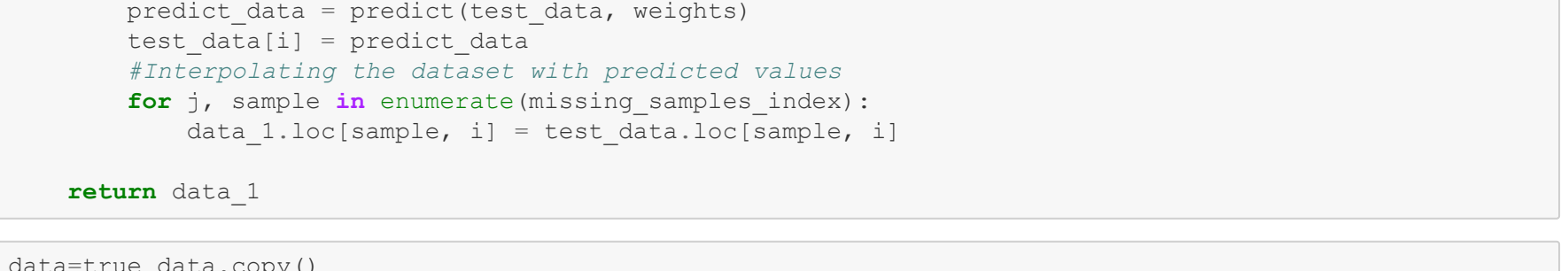
[615 rows x 10 columns]
```

```
ALB    ALP    ALT    AST    BIL    CHE    CHOL    CREA    GGT    PROT
0    38.5    52.5    7.7    22.1    7.5    6.93    3.23    106.0    12.1    69.0
1    38.5    70.3    18.0    24.7    3.9    11.17    4.80    74.0    15.6    76.5
2    46.9    74.7    36.2    52.6    6.1    8.84    5.20    86.0    33.2    79.3
3    43.2    52.0    30.6    22.6    18.9    7.33    4.74    80.0    33.8    75.7
4    39.2    74.1    32.6    24.8    9.6    9.15    4.32    76.0    29.9    68.7
..    ...    ...    ...    ...    ...    ...    ...    ...    ...
610    32.0    416.6    5.9    110.3    50.0    5.57    6.30    55.7    650.9    68.5
611    24.0    102.8    2.9    44.4    20.0    1.54    3.02    63.0    35.9    71.3
612    29.0    87.3    3.5    99.0    48.0    1.66    3.63    66.7    64.2    82.0
613    33.0    NaN    39.0    62.0    20.0    3.56    4.20    52.0    50.0    71.0
614    36.0    NaN    100.0    80.0    12.0    9.07    5.30    67.0    34.0    68.0

[615 rows x 10 columns]
```

```
[615 rows x 10 columns]
percent of the data is missing: 10
percent of the data is missing: 20
percent of the data is missing: 30
percent of the data is missing: 40
percent of the data is missing: 50
percent of the data is missing: 60
percent of the data is missing: 70
percent of the data is missing: 80
percent of the data is missing: 90
[10, 20, 30, 40, 50, 60, 70, 80, 90]
[338.0638038943087, 640.1607543089431, 583.2243667479675, 287.73867247967485, 820.430790569105, 23
3.586823062350577, 577.5590562835913, 61.799361321138214, 379.9559711111111]
```

Out[18]: Text (0.5, 1.0, 'Mean squared error as a function of percent of data removed over all the data')



We can see that from as the percent of the data removed changes the mean squared error changes irregularly as it is random. Here the metric used to evaluate the model is mean squared error.

## 3. Regression Model:-

Now same data interpolation is performing with the regression models i.e. Naive regression,ridge and lasso models. Here, we use the predefined functions (i.e. find\_missing\_features()) to remove the datapoints randomly from the dataset. Other functions to find the columns that has missing values so the method that followed is first we find the which columns have missing values. Now we will find among the missing value features, the least missing feature is considered as a value and remaining features are considered as x-values. we train these data values using the regression models and we predict the missing values.

### Linear Regression:-

In [14]:

```
def linear_regression(x,y):
    w = np.dot(np.dot(np.linalg.inv(np.dot(x.T, x)), x.T), y)
    return w

def predict(data, weights):
    predicted=np.dot(data, weights)
    return predicted

def linear_regression_fit(data):
    data1=data.copy()
    # Finding missing columns name in descending order
    missing_columns=find_missing_columns(data)
    #maintaining a copy of missing columns
    missing_columns_c=missing_columns.copy()

    for i in missing_columns:
        #removing a filled column after each iteration
        missing_columns_c.remove(i)
        all_data = data1[data1.columns.drop(missing_columns_c)]
        test_data = all_data[all_data.isna().any(axis=1)]
        #finding indices of missing samples data
        missing_samples_index = test_data.index
        #splitting data into train data and test data
        train_data = all_data.dropna(how='any', axis=0)
        test_data = test_data.dropna(axis=0)
        #splitting train data into x and y data
        x_data=train_data.drop(['AST', 'ALT', 'BIL', 'CHE', 'CHOL', 'CREA', 'GGT', 'PROT'], axis=1)
        y_data=train_data[['AST', 'ALT', 'BIL', 'CHE', 'CHOL', 'CREA', 'GGT', 'PROT']]

        #Finding the linear regression
        weights = linear_regression(x_data,y_data)
        #predicting weights with trained weights
        predict_data = predict(test_data, weights)
        test_data[1] = predict_data
        #interpolating the dataset with predicted values
        for j, sample in enumerate(missing_samples_index):
            data_1.loc[sample, i] = test_data.loc[sample, i]
    return data_1
```

In [19]:

```
new_data=true_data.copy()
new_data=remove_features(data,30,5)
linear_data=linear_regression_fit(new_data)
print(linear_data)
perc=10
mse_all=[]
percent=[]
print(true_data)
for i in range(9):
    percent.append(perc)
    new_data=true_data.copy()
    new_data=remove_features(data,perc,5)
    mean_data=linear_regression_fit(new_data)
    k=mean_squared_error(true_data,mean_data,perc)
    k=list(k)
    #print(k)
    print("percent of the data is missing:",perc-10)
    mse_all.append(sum(k)/10)
    #print(k)
    print(percent)
    plt.plot(percent,mse_all)
    plt.xlabel("Percent of data removed")
    plt.ylabel("Mean squared error")
plt.title("Mean squared error as a function of percent of data removed over all the data")
```

```
['ALT', 'ALP', 'PROT', 'GGT', 'CHE', 'BIL', 'CHOL']
ALB    ALP    ALT    AST    BIL    CHE    CHOL    CREA    GGT    PROT
0    38.500000    52.500000    7.7    22.1    7.5    6.93    3.230000    106.0    12.1    69.0
1    38.500000    70.300000    18.0    24.7    3.9    11.17    4.80    74.0    15.6    76.5
2    46.900000    74.700000    36.2    52.6    6.1    8.84    5.200000    86.0    33.2    79.3
3    43.200000    52.000000    30.6    22.6    18.9    7.33    4.740000    80.0    33.8    75.7
4    39.15680    74.100000    32.6    24.8    9.6    9.15    4.320000    76.0    29.9    68.7
..    ..    ..    ..    ..    ..    ..    ..    ..    ..
610    32.00000    416.600000    5.9    110.3    50.0    5.57    6.180446    55.7    650.9    68.5
611    24.000000    102.800000    2.9    44.4    20.0    1.54    3.163603    63.0    35.9    71.3
612    29.000000    87.300000    3.5    99.0    48.0    1.66    3.630000    66.7    64.2    82.0
613    29.053833    67.344346    39.0    62.0    20.0    3.56    3.693227    52.0    50.0    71.0
614    36.060000    121.216760    100.0    80.0    12.0    9.07    5.300000    67.0
```

```
GGT    PROT
0    12.100000    50.725068
1    15.600000    76.500000
2    33.200000    79.300000
3    33.800000    75.700000
4    29.900000    64.908339
610    650.900000    68.500000
611    71.634911    71.300000
612    95.532236    82.000000
613    50.000000    71.000000
614    34.000000    68.000000

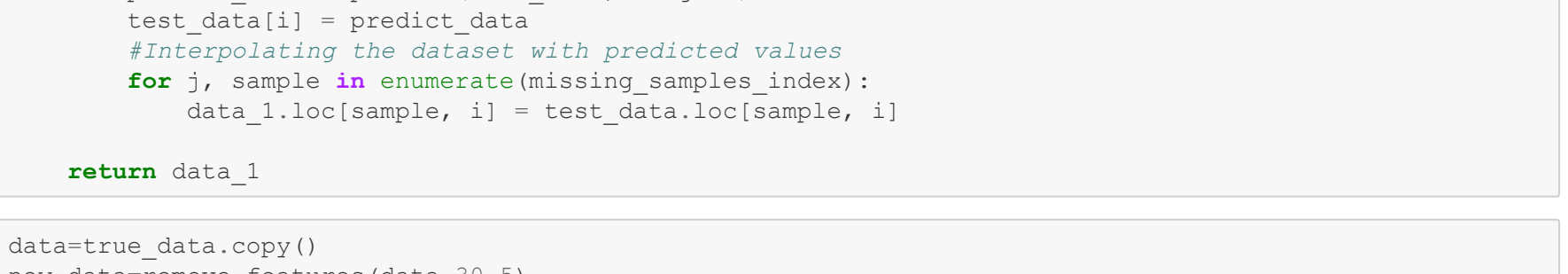
[615 rows x 10 columns]
```

```
ALB    ALP    ALT    AST    BIL    CHE    CHOL    CREA    GGT    PROT
0    38.5    52.5    7.7    22.1    7.5    6.93    3.23    106.0    12.1    69.0
1    38.5    70.3    18.0    24.7    3.9    11.17    4.80    74.0    15.6    76.5
2    46.9    74.7    36.2    52.6    6.1    8.84    5.20    86.0    33.2    79.3
3    43.2    52.0    30.6    22.6    18.9    7.33    4.74    80.0    33.8    75.7
4    39.2    74.1    32.6    24.8    9.6    9.15    4.32    76.0    29.9    68.7
..    ...    ...    ...    ...    ...    ...    ...    ...    ...
610    32.0    416.6    5.9    110.3    50.0    5.57    6.30    55.7    650.9    68.5
611    24.0    102.8    2.9    44.4    20.0    1.54    3.02    63.0    35.9    71.3
612    29.0    87.3    3.5    99.0    48.0    1.66    3.63    66.7    64.2    82.0
613    33.0    NaN    39.0    62.0    20.0    3.56    4.20    52.0    50.0    71.0
614    36.0    NaN    100.0    80.0    12.0    9.07    5.30    67.0    34.0    68.0

[615 rows x 10 columns]
```

```
['ALB', 'PROT', 'ALT', 'CHOL', 'ALP']
percent of the data is missing: 10
['PROT', 'ALT', 'ALP', 'GGT', 'BIL', 'AST', 'CHOL', 'ALP']
percent of the data is missing: 20
['PROT', 'ALB', 'CREA', 'BIL', 'ALT', 'CHOL', 'ALP']
percent of the data is missing: 30
['PROT', 'ALB', 'CREA', 'AST', 'ALT', 'CHOL', 'ALP']
percent of the data is missing: 40
['ALT', 'CHOL', 'ALP', 'PROT', 'GGT', 'CHE', 'BIL', 'ALB']
percent of the data is missing: 50
['ALT', 'BIL', 'ALB', 'PROT', 'CHOL', 'ALP']
percent of the data is missing: 60
['ALT', 'CHOL', 'ALP', 'PROT', 'GGT', 'CHE', 'BIL', 'ALB']
percent of the data is missing: 70
['ALT', 'BIL', 'ALB', 'PROT', 'CHOL', 'ALP']
percent of the data is missing: 80
['PROT', 'ALP', 'CREA', 'AST', 'ALT', 'ALB', 'CHOL']
percent of the data is missing: 90
[10, 20, 30, 40, 50, 60, 70, 80, 90]
[138.52860729705105, 855.8622665307139, 371.84139573004704, 267.7816059912692, 211.1448107400627, 12
7.69946225056888, 593.4844230263159, 526.1250757002268, 692.41718689224]
```

Out[19]: Text (0.5, 1.0, 'Mean squared error as a function of percent of data removed over all the data')



We can see that from as the percent of the data removed changes the mean squared error changes irregularly as it is random. Here the metric used to evaluate the model is mean squared error.

### Ridge Regression:-

In [34]:

```
def ridge_regression(x,y,lamda):
    print(lamda)
    w = np.dot(np.dot(np.linalg.inv(np.dot(x.T, x) + lamda*np.identity(features)), x.T), y)
    return w

def ridge_regression_fit(data):
    data1=data.copy()
    # Finding missing columns name in descending order
    missing_columns=find_missing_columns(data)
    #maintaining a copy of missing columns
    missing_columns_c=missing_columns.copy()

    for i in missing_columns:
        #removing a filled column after each iteration
        missing_columns_c.remove(i)
        all_data = data1[data1.columns.drop(missing_columns_c)]
        test_data = all_data[all_data.isna().any(axis=1)]
        #finding indices of missing samples data
        missing_samples_index = test_data.index
        #splitting data into train data and test data
        train_data = all_data.dropna(how='any', axis=0)
        test_data = test_data.dropna(axis=0)
        #splitting train data into x and y data
        x_data=train_data.drop(['AST', 'ALT', 'BIL', 'CHE', 'CHOL', 'CREA', 'GGT', 'PROT'], axis=1)
        y_data=train_data[['AST', 'ALT', 'BIL', 'CHE', 'CHOL', 'CREA', 'GGT', 'PROT']]

        #Finding the ridge regression
        weights = ridge_regression(x_data,y_data,0.1)
        #predicting weights with trained weights
        predict_data = predict(test_data, weights)
        test_data[1] = predict_data
        #interpolating the dataset with predicted values
        for j, sample in enumerate(missing_samples_index):
            data_1.loc[sample, i] = test_data.loc[sample, i]
    return data_1
```

In [35]:

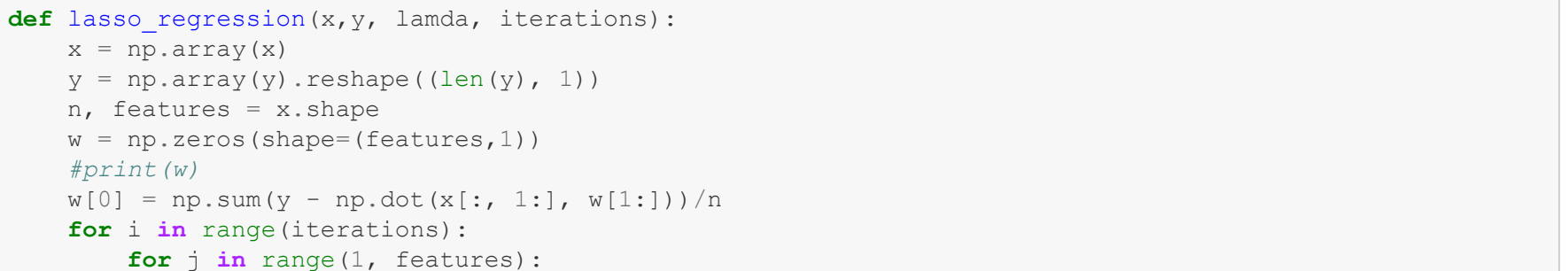
```
data=true_data.copy()
new_data=remove_features(data,30,5)
linear_data=ridge_regression_fit(new_data)
print(linear_data)
perc=10
mse_all=[]
percent=[]
print(true_data)
for i in range(9):
    percent.append(perc)
    new_data=true_data.copy()
    new_data=remove_features(data,perc,5)
    mean_data=ridge_regression_fit(new_data)
    k=mean_squared_error(true_data,mean_data,perc)
    k=list(k)
    #print(k)
    print("percent of the data is missing:",perc-10)
    mse_all.append(sum(k)/10)
    #print(k)
    print(percent)
    plt.plot(percent,mse_all)
    plt.xlabel("Percent of data removed")
    plt.ylabel("Mean squared error")
plt.title("Mean squared error as a function of percent of data removed over all the data")
```

```
['ALB', 'ALP', 'CREA', 'BIL', 'PROT', 'AST', 'CHOL']
ALB    ALP    ALT    AST    BIL    CHE    CHOL    CREA    GGT    PROT
0    38.5    52.500000    7.7    22.1    7.5    6.93    3.230000    106.0    12.1    69.0
1    38.5    70.300000    18.0    24.7    3.9    11.17    4.80    74.0    15.6    76.5
2    46.9    74.700000    36.2    52.6    6.1    8.84    5.20    86.0    33.2    79.3
3    43.2    52.000000    30.6    22.6    18.9    7.33    4.74    80.0    33.8    75.7
4    39.2    74.100000    32.6    24.8    9.6    9.15    4.32    76.0    29.9    68.7
..    ..    ..    ..    ..    ..    ..    ..    ..    ..
610    32.0    416.600000    112.506496    110.3    69.584499    5.67    6.300000
611    24.0    102.800000    2.9    44.4    20.000000    1.54    3.020000    63.0    35.9    71.3
612    29.0    87.300000    3.5    99.0    48.000000    1.66    3.630000    66.7    64.2    82.0
613    33.0    49.942454    39.000000    62.0    18.755176    3.56    4.200000
614    36.0    NaN    100.000000    80.0    12.0    9.07    5.30    67.0    34.0    68.0

[615 rows x 10 columns]
```

```
['ALB', 'PROT', 'ALT', 'CHOL', 'ALP']
percent of the data is missing: 10
['PROT', 'ALB', 'CREA', 'AST', 'ALT', 'CHOL', 'ALP']
percent of the data is missing: 20
['PROT', 'ALB', 'CREA', 'BIL', 'ALT', 'CHOL', 'ALP']
percent of the data is missing: 30
['ALT', 'CHOL', 'ALP', 'PROT', 'GGT', 'CHE', 'BIL', 'ALB']
percent of the data is missing: 40
['ALT', 'BIL', 'ALB', 'PROT', 'CHOL', 'ALP']
percent of the data is missing: 50
['ALT', 'CHOL', 'ALP', 'PROT', 'GGT', 'CHE', 'BIL', 'ALB']
percent of the data is missing: 60
['ALT', 'CHOL', 'ALP', 'PROT', 'GGT', 'CHE', 'BIL', 'ALB']
percent of the data is missing: 70
['ALT', 'CHOL', 'ALP', 'PROT', 'GGT', 'CHE', 'BIL', 'ALB']
percent of the data is missing: 80
['PROT', 'ALP', 'CREA', 'AST', 'ALT', 'ALB', 'CHOL']
percent of the data is missing: 90
[10, 20, 30, 40, 50, 60, 70, 80, 90]
[127447.73520102619, 747630.804077897, 87441.97196282801, 88101.90851031643, 91861.2113678025, 2471.
1418454923577, 31.298241.7436191, 1286.682.5368159, 21.04374.073217697]
```

Out[35]: Text (0.5, 1.0, 'Mean squared error as a function of percent of data removed over all the data')



From the above plot, we can say that as the percent of data removed is increased mean squared error is increasing. And from the plot we can say that ridge regression performed on par with linear model. After experimenting with lambda values, lambda=0.1 is the optimal value for the model.

In [32]:

```
def lasso_regression(x,y,lamda, iterations):
    n, features = x.shape
    y = np.array(y).reshape((len(y), 1))
    w = np.zeros(shape=(features,1))
    #print(w)
    w[0] = np.sum(y - np.dot(x[1:, 1:], w[1:]))/n
    for i in range(iterations):
        w = w.copy()
        r = y - np.dot(x, w)
        B = lamda/2
        A = np.dot(x[1:, 1:], r)
        first=(A+B)/np.dot(x[1:, 1:],T,x[1:, 1])
        second=(A-B)/np.dot(x[1:, 1:],T,x[1:, 1])
        if w[j]>first:
            w[j]=w[j]-first
        elif w[j]<-second:
            w[j]=w[j]+second
        else:
            w[j]=0
    w = np.reshape(w,(len(w),-1))
    return w

def lasso_regression_fit(data):
    data1=data.copy()
    # Finding missing columns name in descending order
    missing_columns=find_missing_columns(data)
    #maintaining a copy of missing columns
    missing_columns_c=missing_columns.copy()

    for i in missing_columns:
        #removing a filled column after each iteration
        missing_columns_c.remove(i)
        all_data = data1[data1.columns.drop(missing_columns_c)]
        test_data = all_data[all_data.isna().any(axis=1)]
        #finding indices of missing samples data
        missing_samples_index = test_data.index
        #splitting data into train data and test data
        train_data = all_data.dropna(how='any', axis=0)
        test_data = test_data.dropna(axis=0)
        #splitting train data into x and y data
        x_data=train_data.drop(['AST', 'ALT', 'BIL', 'CHE', 'CHOL', 'CREA', 'GGT', 'PROT'], axis=1)
        y_data=train_data[['AST', 'ALT', 'BIL', 'CHE', 'CHOL', 'CREA', 'GGT', 'PROT']]

        #Finding the lasso regression
        weights = lasso_regression(x_data,y_data,1.0000)
        #predicting weights with trained weights
        predict_data = predict(test_data, weights)
        test_data[1]
```