```
ECE204 Data Science & Engineering
         This notebook uses X and df for multiple problems, and those variables could be mistakenly used for the wrong problem. To start fresh, restarting the
         kernel in the Kernel menu or with the button above.
         Another useful command "Kernel > Restart and run all" in the menu above. This clears all variables and runs the notebook top to bottom.
         Items 2-5 uses the wine-quality.csv dataset. Each row of this dataset corresponds to data about a different wine sample. The first 11 columns
         correspond to observations made by physiochemical tests on a wine sample (acidity level, density, sugar content, etc.). The 12th column is a rated quality
         score for the wine sample, on a 0-10 scale.
 In [2]: 1 # Necessary imports (run this cell!)
           2 import pandas as pd
           3 import numpy as np
           4 from sklearn.model_selection import train_test_split
         Intel MKL WARNING: Support of Intel(R) Streaming SIMD Extensions 4.2 (Intel(R) SSE4.2) enabled only processors has
          been deprecated. Intel oneAPI Math Kernel Library 2025.0 will require Intel(R) Advanced Vector Extensions (Intel(R)
         AVX) instructions.
         Intel MKL WARNING: Support of Intel(R) Streaming SIMD Extensions 4.2 (Intel(R) SSE4.2) enabled only processors has
         been deprecated. Intel oneAPI Math Kernel Library 2025.0 will require Intel(R) Advanced Vector Extensions (Intel(R)
         AVX) instructions.
         Item 1: Rules & Honor Code
                See Canvas.
          Dataset manipulation and information summary
          Item 3 a)
 In [7]:
           1 ### Starter code
             df = pd.read_csv("wine-quality.csv")
 Out[7]:
               fixed acidity volatile acidity citric acid residual sugar chlorides free sulfur dioxide total sulfur dioxide density
                                                                                        34.0 0.99780 3.51
                     7.4
                               0.700
                                        0.00
                                                           0.076
                                                                          11.0
                                                                                                            0.56
                                                     1.9
                                                                                                                   9.4
                                        0.00
                                                     2.6
                                                           0.098
                                                                          25.0
                                                                                        67.0 0.99680 3.20
                                                                                                            0.68
                     7.8
                               0.880
                                                                                                                   9.8
                                                    2.3
                                                                                        54.0 0.99700 3.26
                     7.8
                                                           0.092
                                                                                                                   9.8
                               0.760
                                        0.04
                                                                          15.0
                                                                                                            0.65
                                                           0.075
                                                                          17.0
                     11.2
                               0.280
                                        0.56
                                                     1.9
                                                                                        60.0 0.99800 3.16
                                                                                                            0.58
                                                                                                                   9.8
                     7.4
                               0.700
                                        0.00
                                                           0.076
                                                                          11.0
                                                                                        34.0 0.99780 3.51
                                                                                                            0.56
                                                     1.9
                                                                                                                   9.4
                                          •••
                                                      ...
          1594
                     6.2
                               0.600
                                        0.08
                                                    2.0
                                                           0.090
                                                                          32.0
                                                                                        44.0 0.99490 3.45
                                                                                                            0.58
                                                                                                                  10.5
                     5.9
                               0.550
                                        0.10
                                                     2.2
                                                           0.062
                                                                          39.0
                                                                                        51.0 0.99512 3.52
                                                                                                            0.76
          1595
                                                                                                                  11.2
                                                                                        40.0 0.99574 3.42
                                                     2.3
                                                           0.076
                               0.510
                                        0.13
                                                                          29.0
                                                                                                            0.75
                                                                                                                  11.0
                                        0.12
                                                                                         44.0 0.99547 3.57
                               0.645
                                                     2.0
                                                                                                            0.71
                                                                          18.0
                     6.0
                               0.310
                                        0.47
                                                           0.067
                                                                                        42.0 0.99549 3.39
                                                                                                            0.66 11.0
                                                    3.6
         1599 rows × 12 columns
 In [5]: 1 ### Your code here
           2 df["residual sugar"].mean()
 Out[5]: 2.53880550343965
         Item 3 b)
In [12]: 1 ### Starter code
             df = pd.read_csv("wine-quality.csv")
           4 ### Your code here
           5 df['quality'].value_counts()
Out[12]: quality
               681
               638
                53
         Name: count, dtype: int64
          Item 3 c)
In [22]: 1 ### Starter code
             df = pd.read_csv("wine-quality.csv")
           4 ### Your code here
             df1 = df.groupby('quality')
             df1['alcohol'].mean()
Out[22]: quality
                9.955000
              10.265094
               9.899706
              10.629519
              11.465913
              12.094444
         Name: alcohol, dtype: float64
          Item 3 d)
In [27]:
          1 ### Starter code
             df = pd.read_csv("wine-quality.csv")
           4 ### Your code here
             df2 = df.groupby('quality')
            7 df1['pH'].mean()
Out[27]: quality
               3.398000
               3.381509
              3.304949
               3.318072
              3.290754
               3.267222
         Name: pH, dtype: float64
         Unsupervised leaning
         Item 4
In [31]:
          1 ### Starter code
           2 from sklearn.preprocessing import StandardScaler
           3 from sklearn.decomposition import PCA
             df = pd.read_csv("wine-quality.csv")
           7 features = ['fixed acidity', 'volatile acidity', 'citric acid', 'residual sugar',
                         'chlorides', 'free sulfur dioxide', 'total sulfur dioxide', 'density',
                         'pH', 'sulphates', 'alcohol']
          10 X = df[features]
          12 ### Your code here
          13 sc = StandardScaler()
          14 X_scaled = sc.fit_transform(X)
          15 pca = PCA(n_components=0.95)
          16 pca.fit(X_scaled)
             pca.n_components_
Out[31]: 9
         Supervised learning
         Item 5 a)
In [35]: 1 ### Starter code
           2 from sklearn.tree import DecisionTreeClassifier
             df = pd.read_csv("wine-quality.csv")
           6 features = ['fixed acidity', 'volatile acidity', 'citric acid', 'residual sugar',
                         'chlorides', 'free sulfur dioxide', 'total sulfur dioxide', 'density',
                         'pH', 'sulphates', 'alcohol']
           9 labels = 'quality'
          11 X = df[features]
          12 y = df[labels]
          14 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.2, random_state = 12)
          16 ### Your code here
          17 dt = DecisionTreeClassifier(max_depth=5, random_state=12)
          18 dt.fit(X_train, y_train)
          19 dt.feature_importances_
Out[35]: array([0.028678 , 0.1071444 , 0.02178991, 0.02959076, 0.01048532,
                 0.01770544, 0.13105678, 0.
                                                   , 0.02876334, 0.20474289,
                 0.42004315])
          Item 5 b)
          1 ### Starter code
In [39]:
           2 from sklearn.neighbors import KNeighborsClassifier
              df = pd.read_csv("wine-quality.csv")
           6 features = ['fixed acidity', 'volatile acidity', 'citric acid', 'residual sugar',
                         'chlorides', 'free sulfur dioxide', 'total sulfur dioxide', 'density',
                         'pH', 'sulphates', 'alcohol']
           9 labels = 'quality'
          11 X = df[features]
          12 y = df[labels]
          14 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.2, random_state = 42)
          16 X_{\text{new}} = [[12.0, 0.6, 0.7, 2.2, 0.074, 10.0, 47.0, 1.0, 3.25, 0.6, 9.0]]
          18 ### Your code here
          19 knn = KNeighborsClassifier(n_neighbors=4)
          20 knn.fit(X_train, y_train)
          21 knn.predict(X_new)
         /Users/janeli/anaconda3/envs/bonding/lib/python3.11/site-packages/sklearn/base.py:439: UserWarning: X does not have
         valid feature names, but KNeighborsClassifier was fitted with feature names
           warnings.warn(
Out[39]: array([3])
         Item 6 a)
In [47]: 1 ### Starter code
           2 from sklearn.linear_model import LinearRegression
             df = pd.read_csv("wine-quality.csv")
           6 X = df[['fixed acidity', 'residual sugar']]
           7 y = df['pH']
           9 # create train-test split
          10 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.2, random_state = 42)
          12 ### Your code here
          13 model = LinearRegression()
          14 model.fit(X_train, y_train)
          15 print(model.intercept_)
          16 print(model.coef_)
          17 print(model.score(X_train, y_train))
         3.813959024718791
          [-0.06042496 \quad 0.00025412]
         0.45724793983222456
          Item 6 b)
In [54]: 1 ### Starter code
           from sklearn.model_selection import train_test_split, cross_val_score
           3 from sklearn.preprocessing import PolynomialFeatures
           4 from sklearn.metrics import mean_squared_error
             df = pd.read_csv("wine-quality.csv")
           8 X = df[['fixed acidity']]
           9 y = df['pH']
          11 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.33, random_state = 50)
          13 ### Your code here
          14 model = PolynomialFeatures(degree=2)
          15 X_train = model.fit_transform(X_train)
          16 X_test = model.fit_transform(X_test)
          18 model = LinearRegression()
          19 model.fit(X_train, y_train)
          20 y_pred = model.predict(X_test)
          21 mean_squared_error(y_test,y_pred)
Out [54]: 0.01206391332110709
          Item 7 a)
         This item uses a different dataset: air_quality.csv . Before starting, inspect the data carefully:
In [55]: 1 df = pd.read_csv('air_quality.csv', parse_dates=True, index_col=0).sort_index().asfreq(freq = 'H')
           2 df.head()
          /var/folders/gj/1nzwkx3x3rg1zwxj2mx5xkrm0000gn/T/ipykernel_95446/2445795387.py:1: FutureWarning: 'H' is deprecated
         and will be removed in a future version, please use 'h' instead.
           df = pd.read_csv('air_quality.csv', parse_dates=True, index_col=0).sort_index().asfreq(freq = 'H')
Out [55]:
                          RADON_pCi/L TEMP_F HUMIDITY PRESSURE_hPa CO2_ppm VOC_ppb High_CO2
                  recorded
          2021-02-01 00:00:00
                                 1.22
                                       63.39
                                                 39.0
                                                              989
                                                                       560
                                                                              75.0
                                                                                        No
          2021-02-01 01:00:00
                                1.19
                                       62.67
                                                 38.5
                                                              989
                                                                       548
                                                                              91.0
                                                                                        No
          2021-02-01 02:00:00
                                       63.55
                                                                             110.0
                                 1.19
                                                 38.0
                                                                       593
                                                                                        No
                                                              989
          2021-02-01 03:00:00
                                 1.22
                                       63.07
                                                 37.5
                                                                             110.0
                                                                      573
                                                              989
                                                                                        No
          2021-02-01 04:00:00
                                 1.24
                                       62.33
                                                 37.5
                                                                       555
                                                                             138.0
                                                              990
                                                                                        No
           1 ### Your code here
In [88]:
           2 import matplotlib.pyplot as plt
              rolling_avg_humidity = df.HUMIDITY.rolling("7d").mean()
              ax = rolling_avg_humidity.plot(figsize=(10, 6), color='blue')
             ax.set_title('Rolling average (over 7 days) of HUMIDITY')
             ax.set_xlabel('Date')
             ax.grid(True)
             ax.set_ylabel('Humidity')
          12 ## You may finde the following lines of code helpful:
          13 fig = ax.get_figure()
          14 fig.savefig("WuzhenLi.png")
                                          Rolling average (over 7 days) of HUMIDITY
              50
                                                                         M
              45
              30
             25
                      Mar
                                                                              Oct
                                                                                      Nov
                                                                                              Dec
                                      May
                                               Jun
                                                              Aug
                                                                                                      2022
                                                              Date
          Item 7 b)
           1 ### Starter code
           2 from statsmodels.tsa.ar_model import AutoReg
           4 df = pd.read_csv('air_quality.csv', parse_dates=True, index_col=0).sort_index().asfreq(freq = 'H')
           6 #train/test splits
             end_train = '2022-1-30 23:00'
           8 begin_test = '2022-1-31 00:00'
           9 end_test = '2022-1-31 23:00' #Last timestamp in the series
          11 df_train = df.loc[:end_train, ['RADON_pCi/L']]
          12 df_test = df.loc[begin_test:, ['RADON_pCi/L']].copy()
          14 ### Your code
          16 model = AutoReg(df_train, lags=100)
          17 model = model.fit()
          18 time = '2022-01-31 10:00:00'
          19 model.predict(start=time, end=time)
         Intel MKL WARNING: Support of Intel(R) Streaming SIMD Extensions 4.2 (Intel(R) SSE4.2) enabled only processors has
         been deprecated. Intel oneAPI Math Kernel Library 2025.0 will require Intel(R) Advanced Vector Extensions (Intel(R)
         AVX) instructions.
```

/var/folders/gj/1nzwkx3x3rg1zwxj2mx5xkrm0000gn/T/ipykernel\_95446/3264789675.py:4: FutureWarning: 'H' is deprecated

df = pd.read\_csv('air\_quality.csv', parse\_dates=True, index\_col=0).sort\_index().asfreq(freq = 'H')

and will be removed in a future version, please use 'h' instead.

1.378586

Do not forget to upload your Jupyter notebook!

Out[78]: 2022-01-31 10:00:00

Freq: h, dtype: float64

Final Exam (Remote)