HW3. by Handan Cetin | USCID: 6074572947 | github: handancetin In [1]: import os import csv import pandas as pd import numpy as np import matplotlib.pyplot as plt import seaborn as sns from sklearn.utils import resample import random random.seed(1234) # for reproducibility import warnings warnings.filterwarnings('ignore') # for plots 1. Time Series Classification Part 1: Feature Creation/Extraction An interesting task in machine learning is classification of time series. In this problem, we will classify the activities of humans based on time series obtained by a Wireless Sensor Network 1.(a) Download the AReM data The dataset contains 7 folders that represent seven types of activities. In each folder, there are multiple files each of which represents an instant of a human performing an activity.1 Each file contains 6 time series collected from activities of the same person, which are called avg\_rss12 var\_rss12 • avg\_rss13 var\_rss13 avg\_rss23 var\_rss23 There are 88 instances in the dataset, each of which contains 6 time series and each time series has 480 consecutive values. In [2]: filepath = '../data/AReM/' 1.(b) Test and Train data division • Test: datasets 1 and 2 in bending1 and bending2, and datasets 1, 2, and 3 in other folders • Train: remaining datasets In [3]: trainFiles = [] testFiles = [] fileCounter = 0# Loop through each file in the folders for root, dir, file in os.walk(filepath): for f in file: fpof = os.path.join(root, f) if 'bending' in fpof: if ('dataset1' in fpof) or ('dataset2' in fpof): testFiles.append(fpof) else: trainFiles.append(fpof) else: if ('dataset1' in fpof) or ('dataset2' in fpof) or ('dataset3' in fpof): testFiles.append(fpof) trainFiles.append(fpof) fileCounter = fileCounter + 1 print('Out of total', fileCounter, 'files:', len(trainFiles), 'datasets added to training list, and', len(testFiles), 'datasets added to testing list.') Out of total 88 files: 39 datasets added to training list, and 49 datasets added to testing list. 1.(c) Feature Extraction Classification of time series usually needs extracting features from them. In this problem, we focus on time-domain features 1.(c)i. Research what types of time-domain features are usually used in time series classification. List them. **Statistical Measures:** Mean Median Standard deviation Variance Skewness • Kurtosis Higher-order moments (e.g., third and fourth moments) • Various percentiles (e.g., 25th, 75th) Minimum and maximum values Range (difference between max and min) Mode (most frequent value) Rolling statistics (Rolling mean, standard deviation, etc.) Autocorrelation at different lags Cross-correlation with other relevant time series Slope (linear regression coefficients, etc) **Time Series Characteristics:**  Length of time series Number of data points Trend, seasonal, or residual components Entropy neasures (Shannon entropy, Spectral entropy, Permutation entropy, Sample entropy, etc.) Time-domain representation of frequency domain features like dominant frequency or peak frequency Central moments (e.g., second, third, and fourth central moments) Histograms of the time series values • Shape-based features like fractal dimensions Time Series Amplitude and Time Derivatives (Amplitude variations, derivatives, etc.) Lagged values or differences at different time intervals Autoregressive coefficients Moving Average coefficients 1.(c)ii. Extract the time-domain features minimum, maximum, mean, median, standard deviation, first quartile, and third quartile for all of the 6 time series in each instance In [4]: trainingSet = [] testingSet = [] for exp in ['train', 'test']: if exp == 'train' : files = trainFiles else: files = testFiles values = [] folderName = [] # Read each file and keep the folder name for f in files: filename\_components = os.path.dirname(f).split("/") folderName.append(filename components[-1]) # bending2/dataset4 is not comma-seperated, so this part is required separator = csv.Sniffer().sniff(open(f).read()).delimiter if separator == ",": fc = pd.read\_csv(f, skiprows=5, header=None, on\_bad\_lines='skip') else: fc = pd.read\_csv(f, skiprows=5, sep="\\s+", header=None, on\_bad\_lines='skip') fc.columns = ['time', 'avg\_rss12', 'var\_rss12', 'avg\_rss13', 'var\_rss13', 'avg\_rss23', 'var\_rss23'] values.append(fc.describe().drop('count').drop(columns='time').T.values.flatten()) In [5]: trainingSet = [] testingSet = [] for exp in ['train', 'test']: if exp == 'train' : files = trainFiles else: files = testFiles values = [] folderName = [] # Read each file and keep the folder name for f in files: filename\_components = os.path.dirname(f).split("/") folderName.append(filename\_components[-1]) # bending2/dataset4 is not comma-seperated, so this part is required separator = csv.Sniffer().sniff(open(f).read()).delimiter if separator == ",": fc = pd.read\_csv(f, skiprows=5, header=None, on\_bad\_lines='skip') else: fc = pd.read\_csv(f, skiprows=5, sep="\\s+", header=None, on\_bad\_lines='skip') fc.columns = ['time','avg\_rss12', 'var\_rss12', 'avg\_rss13', 'var\_rss13', 'avg\_rss23', 'var\_rss23'] values.append(fc.describe().drop('count').drop(columns='time').T.values.flatten()) # Generate statistics df = pd.DataFrame(values) statNames = [] for s in range(1, 7): for stat in ['mean', 'std', 'min', '1st quart', 'median', '3rd quart', 'max']: statNames.append(f'{stat} ({s})') df.columns = statNames # add activity column df['activity'] = pd.Series(folderName) # Save the dataframes if exp == 'train' : trainingSet = df else: testingSet = df trainingSet['isTraining'] = True testingSet['isTraining'] = False In [6]: allData = pd.concat([trainingSet, testingSet]) allData Out[6]: 1st 1st 3rd median 3rd median min min min max max max quart std (2) std (6) std (1) quart mean (2) mean (6) quart activity isTraining mean (1) (1) quart (1) (1) (2) (5) (1) (6) (1) (6) (6) 43.969125 1.618364 36.25 43.310 44.50 44.6700 48.00 0.413125 0.263111 0.0 30.75 0.555312 0.487826 0.0 0.0000 0.49 0.8300 2.96 bending1 True 43.454958 1.386098 37.00 42.500 43.25 45.0000 48.00 0.378083 0.315566 0.0 ... 33.50 0.679646 0.622534 0.0 0.4300 0.50 0.8700 5.26 bending1 True **2** 42.179812 3.670666 33.00 0.0000 39.150 43.50 45.0000 47.75 0.696042 0.630860 0.0 ... 38.67 0.613521 0.524317 0.0 0.50 1.0000 2.18 bending1 True 41.678063 2.243490 33.00 41.330 42.7500 45.75 0.0 ... 37.50 0.383292 0.389164 0.0000 41.75 0.535979 0.405469 0.0 0.43 0.5000 1.79 bending1 True 43.000 43.954500 1.558835 35.00 44.33 45.0000 47.40 0.426250 0.338690 0.0 ... 38.50 0.493292 0.513506 0.0 0.0000 0.43 0.9400 1.79 bending1 True 36.541667 4.018922 27.00 33.250 36.00 39.8125 44.33 2.847958 1.892397 0.0 24.50 3.225458 1.769758 0.0 1.8850 2.87 4.2625 9.18 cycling False **45** 35.879875 4.614878 19.00 33.000 36.00 39.5000 43.75 3.414312 2.237955 0.0 ... 26.50 3.093021 1.626034 0.0 1.8900 2.93 4.0600 8.50 cycling False **46** 35.752354 4.614802 18.50 33.000 36.00 39.3300 44.25 3.328104 2.140576 0.0 ... 24.33 3.069667 1.748326 0.0 1.7975 2.77 4.0600 9.39 cycling False 37.177042 3.581301 24.25 1.601799 0.0 1.5000 34.500 36.25 40.2500 45.00 2.374208 0.0 ... 25.50 2.921729 1.852600 2.50 3.9000 9.34 cycling False 2.737307 36.248768 3.824632 23.33 33.415 39.2500 43.50 2.093999 0.0 ... 27.00 3.532463 1.965267 0.0 2.1700 4.6250 11.15 cycling False 88 rows × 44 columns 1.(c)iii. Estimate the standard deviation of each of the time-domain features you extracted from the data stdTable = pd.DataFrame({ "Standart Deviation of Statistics": allData.describe().loc['std'].index, "All Data" : allData.describe().loc['std'].values, "Train Set Only" : trainingSet.describe().loc['std'].values, "Test Set Only" : testingSet.describe().loc['std'].values}) stdTable Out[7]: **Standart Deviation of Statistics** All Data Train Set Only Test Set Only 0 5.335700 4.850250 5.737933 mean (1) 1.772185 1.854324 1.721256 1 std (1) 2 min (1) 9.569975 9.077137 10.037780 1st quart (1) 6.153874 5.957810 3 6.366662 median (1) 5.440054 4 5.011517 5.809848 5 3rd quart (1) 5.138925 4.653297 5.534215 6 max (1) 4.394362 3.962409 4.731086 7 1.628452 mean (2) 1.574198 1.518276 8 0.884137 0.866277 0.906096 std (2) 0.000000 0.000000 0.000000 min (2) 10 1st quart (2) 0.946386 0.904977 0.983765 1.364829 11 median (2) 1.412293 1.456547 12 3rd quart (2) 2.125399 2.036952 2.210574 max (2) 5.062729 13 5.191506 5.011786 15 std (3) 0.946670 1.098571 0.816820 min (3) 2.956462 2.930310 16 3.019004 4.184895 17 1st quart (3) 4.220658 4.220449 18 median (3) 4.036396 3.898314 4.088755 3rd quart (3) 4.171628 19 3.959001 4.280531 4.951745 20 max (3) 4.875137 4.828204 1.159240 21 mean (4) 1.166178 1.174711 22 std (4) 0.458283 0.468153 0.451051 0.000000 0.000000 23 min (4) 0.000000 1st quart (4) 0.843405 0.830551 24 0.856999 1.166821 median (4) 1.145985 1.122680 25 26 3rd quart (4) 1.552504 1.538679 1.567411 27 max (4) 2.183625 2.308045 2.102071 28 5.941841 5.445393 mean (5) 5.675543 0.941745 29 std (5) 1.024918 1.113642 30 min (5) 6.124001 6.469931 5.863744 5.896663 6.273085 31 1st quart (5) 6.096465 5.505017 32 median (5) 5.813782 6.187095 5.721373 33 3rd quart (5) 5.531720 5.404274 6.430708 34 max (5) 5.741238 5.094435 35 mean (6) 1.154889 1.136198 1.174114 36 std (6) 0.517651 0.506597 0.526738 37 min (6) 0.045838 0.000000 0.061429 0.761659 0.759782 38 1st quart (6) 0.758687 39 median (6) 1.086474 1.070193 1.106415 40 3rd quart (6) 1.523739 1.506349 1.543162 2.655811 2.417709 41 max (6) 2.518921 fig, axes = plt.subplots(1, 3) fig.set\_size\_inches(12,4) sns.boxplot(allData.describe().T, saturation = 1, ax = axes[0], width=0.5, flierprops = dict(marker = '.')) axes[0].set\_title('All Data') sns.boxplot(trainingSet.describe().T, saturation = 1, ax = axes[1], width=0.5, flierprops = dict(marker = '.')) axes[1].set\_title('Training Set') sns.boxplot(testingSet.describe().T, saturation = 1, ax = axes[2], width=0.5, flierprops = dict(marker = '.')) axes[2].set\_title('Testing Set') plt.suptitle('Summary of Summary Statistics of Each Feature') plt.tight\_layout(pad = 1) plt.show() Summary of Summary Statistics of Each Feature Testing Set All Data Training Set 80 50 50 40 40 60 30 30 40 20 20 20 10 10 50% 50% 25% 75% max count mean std 25% 75% 25% 50% 75% count mean std min min count mean std min 1.(c)iii. (cont.) Then, use Python's bootstrapped or any other method to build a 90% bootsrap confidence interval for the standard deviation of each feature. In [9]: bootstrappedSet = [] for sample in range(1000): resampled = resample(allData.iloc[:, :-2]) # Bootstrapping excluding the activity and isTraining columns bootstrappedSet.append(resampled.apply(lambda col: np.std(col))) bootstrappedSet = pd.DataFrame(bootstrappedSet) bootstrappedSet.columns = list(allData.columns[:-2]) confInterval = bootstrappedSet.apply(lambda col : (np.percentile(col, 5), np.percentile(col, 95)), axis=0).apply(np.around, args=(2,)).T confInterval.columns = ['lower bound', 'higher bound'] confInterval['CI range'] = confInterval['higher bound']-confInterval['lower bound'] confInterval Out[9]: lower bound higher bound CI range mean (1) 4.69 5.85 1.16 std (1) 1.56 1.94 0.38 min (1) 8.20 10.67 2.47 1st quart (1) 6.59 1.01 5.58 median (1) 4.81 5.94 1.13 3rd quart (1) 4.30 5.82 1.52 3.37 1.97 max (1) 5.34 mean (2) 1.38 1.69 0.31 std (2) 0.80 0.93 0.13 min (2) 0.00 0.00 0.00 0.21 1st quart (2) 0.82 1.03 median (2) 1.23 1.53 0.30 3rd quart (2) 1.88 2.28 0.40 4.61 5.36 0.75 max (2) mean (3) 3.39 4.45 1.06 std (3) 0.76 1.11 0.35 min (3) 2.75 3.08 0.33 1st quart (3) 1.07 3.60 4.67 median (3) 3.41 4.51 1.10 3rd quart (3) 3.50 4.65 1.15 max (3) 4.14 5.42 1.28 1.07 1.21 mean (4) 0.14 0.42 0.06 std (4) 0.48 min (4) 0.00 0.00 0.00 1st quart (4) 0.77 0.89 0.12 median (4) 1.05 1.20 0.15 3rd quart (4) 1.42 1.62 0.20 max (4) 1.96 2.36 0.40 mean (5) 4.39 6.76 2.37 std (5) 0.81 1.20 0.39 4.29 3.26 min (5) 7.55 1st quart (5) 4.78 7.20 2.42 median (5) 4.48 6.88 2.40 3rd quart (5) 4.30 6.56 2.26 1.79 max (5) 4.73 6.52 mean (6) 1.06 1.21 0.15 std (6) 0.47 0.54 0.07 min (6) 0.00 0.08 80.0 1st quart (6) 0.69 0.80 0.11 median (6) 0.99 1.14 0.15 3rd quart (6) 0.19 1.40 1.59 max (6) 2.22 2.76 0.54 1.(c)iv. Use your judgement to select the three most important time-domain features In [10]: CIRange = confInterval.sort\_values(by=['CI range'])['CI range'] fig, ax = plt.subplots() fig.set\_size\_inches(3,10) ax.barh(range(len(CIRange)), CIRange) ax.set\_yticks(range(len(CIRange))) ax.set\_yticklabels(CIRange.index) ax.invert\_yaxis() # labels read top-to-bottom ax.set\_xlabel('Confidence Interval Range') plt.show() min (4) min (2) std (4) std (6) min (6) 1st quart (6) -1st quart (4) std (2) mean (4) mean (6) median (4) median (6) 3rd quart (6) 3rd quart (4) 1st quart (2) median (2) mean (2) min (3) std (3) std (1) std (5) 3rd quart (2) max (4) max (6) max (2) 1st quart (1) mean (3) 1st quart (3) median (3) median (1) 3rd quart (3) mean (1) max (3) 3rd quart (1) max (5) max (1) 3rd quart (5) mean (5) median (5) 1st quart (5) min (1) min (5) 2 3 Confidence Interval Range I would choose the three most important features by looking at the width of the confidence intervals for each feature as narrower intervals would indicate higher confidence in the estimation. However, this selection would not work in case of overlapping distributions. In [11]: # 'activity' labels as binary classification: 1 if from bending datasets, 0 otherwise. dataSet = allData.copy() dataSet['activity'] = np.where(dataSet['activity'].str.contains('bending'), 'bending', 'not\_bending') dataSubset = dataSet[['max (3)', 'max (4)', 'median (4)', '3rd quart (3)', '3rd quart (4)', '3rd quart (5)', 'activity']] # visualize scatter plots for the features of avg\_rss13, var\_rss13, avg\_rss23. p = sns.pairplot(data = dataSubset, kind = 'scatter', hue='activity', grid\_kws={'despine' : False}, plot\_kws={"s": 5}) p.fig.set\_size\_inches(12,10) sns.move\_legend(p, "lower center", bbox\_to\_anchor=(0.5, -0.1)) 30 max (3) 10 max (4) 5 40 max (5) 02 median (3) Ä median (4) median (5) 20 10 20 3rd quart (3) 10 3rd quart (4) Marie 1 3rd quart (5) 30 20 40 50 0 5 50 20 25 50 0 20 20 10 3rd quart (4) 3rd quart (3) max (4) max (3) max (5) median (3) median (4) median (5) 3rd quart (5) activity bending not\_bending After looking at the various combinations of pairwise scatter plots, I chose the features with relatively clearer seperation: max median 3rd quartile 2. ISLR 3.7.4 Question: I collect a set of data (n = 100 observations) containing a single predictor and a quantitative response. I then fit a linear regression model to the data, as well as a separate cubic regression, i.e.  $Y = \beta_0 + \beta_1 X + \beta_2 X^2 + \beta_3 X^3 + \epsilon$ . (a) Suppose that the true relationship between X and Y is linear, i.e.  $Y = \beta_0 + \beta_1 X + \epsilon$ . Consider the training residual sum of squares (RSS) for the linear regression, and also the training RSS for the cubic regression. Would we expect one to be lower than the other, would we expect them to be the same, or is there not enough information to tell? Justify your answer. • Answer: If we assume that the true relationship is linear, the linear regression model should provide a good fit to the training data. However, we can expect training RSS with cubic regression model be relatively lower than the linear regression model, since the cubic regression is more flexible and may overfit to the data that we are training the model. Another possibility is that, if the relationship is truly linear (without outliers), the cubic regression model would set nonlinear coefficients (for example  $\beta_3$ ) close to zero. (b) Answer (a) using test rather than training RSS. • Answer: We would expect the testing RSS with liner regression model be lower than the cubic regression model. Although the cubic regression may capture some nonlinear patterns due to its flexibity, it might introduce unnecessary complexity that doesn't improve the model's fit. The linear model is definitely a more appropriate choice when the true relationship is linear. (c) Suppose that the true relationship between X and Y is not linear, but we don't know how far it is from linear. Consider the training RSS for the linear regression, and also the training RSS for the cubic regression. Would we expect one to be lower than the other, would we expect them to be the same, or is there not enough information to tell? Justify your answer. • Answer: Similar scenerio as in (a): We would expect training RSS with cubic regression to be lower than the linear regression model as the cubic regression model would capture the flexibility in the training data. (d) Answer (c) using test rather than training RSS. • Answer: There is not enough information to answer this question. The comparison of RSS between linear and cubic regression models depends on the true nature of the data we have. In this scenario, we could try both models to see which one is better fit.