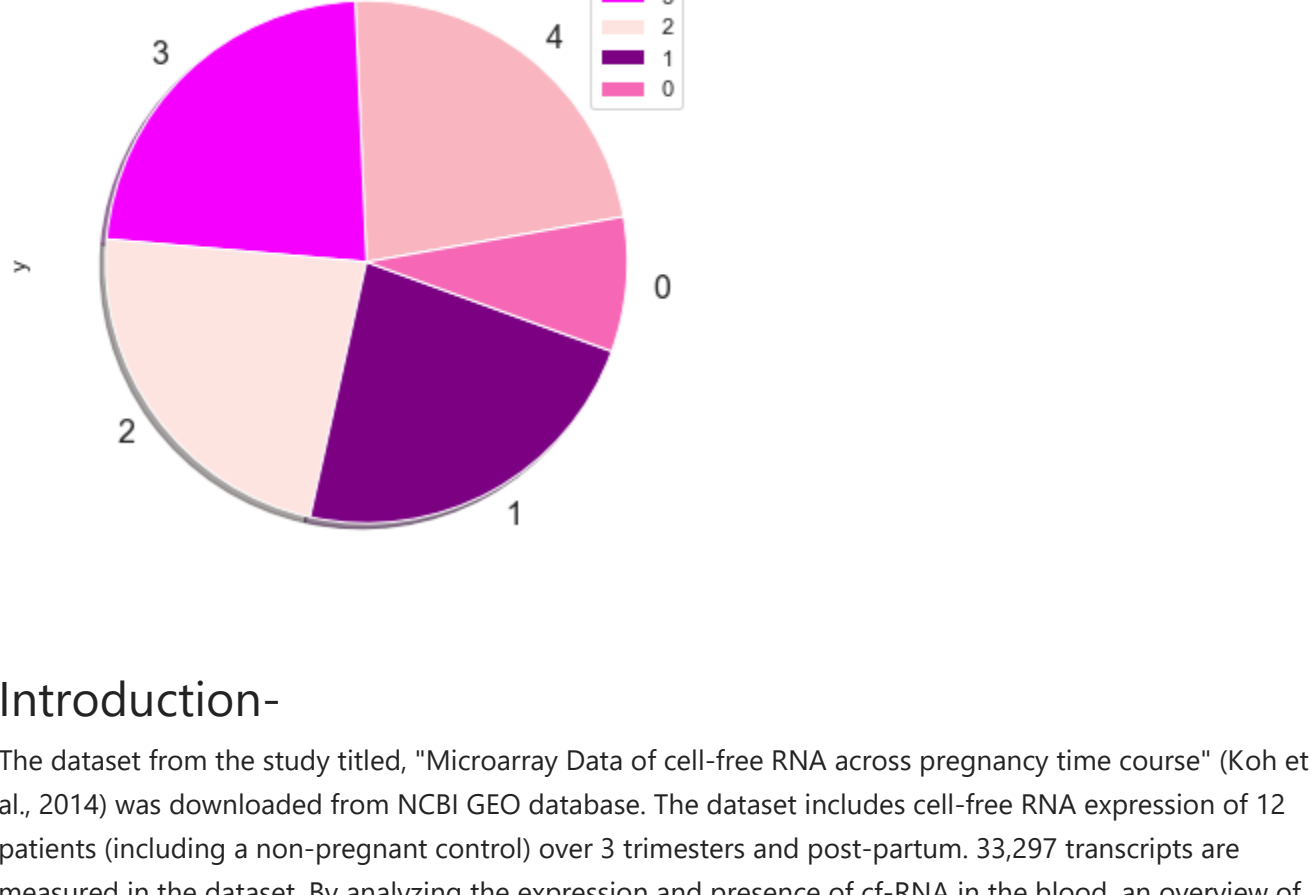


Predicting Pregnancy Stage from Cell free mRNA Levels

Jessie Bologna & Lindsay Reisman
NYU Bioinformatic - Machine Learning
5/15/21



Introduction-

The dataset from the study titled, "Microarray Data of cell-free RNA across pregnancy time course" (Koh et al., 2014) was downloaded from NCBI GEO database. The dataset includes cell-free RNA expression of 12 patients (including a non-pregnant control) over 3 trimesters and post-partum. 33,297 transcripts are measured in the dataset. By analyzing the expression and presence of cf-RNA in the blood, an overview of tissue development and overall health of a fetus can be determined (2014). This study aimed to provide a non-invasive technique for examining fetal tissue development.

Our goal with this dataset is to determine what trimester a patient is in given the expression levels. From this analysis, we are hoping to find which transcript has the highest differential expression during the final stages of pregnancy. Using various Machine Learning models to classify the data to learn and understand how and when these transcripts are expressed, and what up-regulate genes they code for during pregnancy.

Methods -

Overview: To determine which machine learning method worked best for our dataset we begin by downloading the dataset into a pandas dataframe. We next ran various standard data preprocessing steps to get the data ready for classification modeling. We created some charts and graphs to visualize our data. Finally, we created a final dataframe that included all the changes made during our preprocessing steps, we then split that data into a feature matrix (X) and a target matrix (y), and lastly we ran various classification models to see which model was best able to predict which class a patient was in.

Import the Data:

The data consists of 48 samples (12 patients) over the course of 5 stages of pregnancy: first, second, and third trimester, postpartem, and control/non-pregnant microarray data of cell free RNA was collected across 33,297 genes for each patient during each of these time courses

ID_REF	7892501	7892502	7892503	7892504	7892505	7892506	7892507	7892508	7892509	7892510	...	8180409	8180410	8180411
count	48.000000	48.000000	48.000000	48.000000	48.000000	48.000000	48.000000	48.000000	48.000000	48.000000	...	48.000000	48.000000	48.000000
mean	5.004943	2.482006	2.561178	7.310324	2.787235	2.605373	3.165746	2.243372	10.500638	2.562259	...	11.946849	10.312863	8.729550
std	1.853566	0.287033	0.247941	1.055773	0.425361	0.404084	0.727252	0.185912	1.742822	0.266118	...	0.645432	0.727021	1.149331
min	2.089713	1.938133	1.982890	5.428908	2.102378	2.122335	2.367414	1.923975	2.718209	2.022346	...	7.643652	6.219645	5.253879
25%	3.286410	2.287814	2.390578	6.533863	2.540726	2.351941	2.791254	2.106138	10.338530	2.390800	...	11.976620	10.136710	8.073634
50%	5.519676	2.446147	2.514482	7.021481	2.696597	2.464081	2.975391	2.231955	10.952625	2.561447	...	12.035485	10.439085	9.131304
75%	6.390771	2.619133	2.759148	8.047929	2.965962	2.745546	3.281732	2.346940	11.409250	2.682029	...	12.119032	10.698903	9.442569
max	9.105850	3.351991	3.142244	10.706420	4.643989	4.226212	7.068907	2.800776	11.993470	3.234456	...	12.279910	11.043520	10.496080

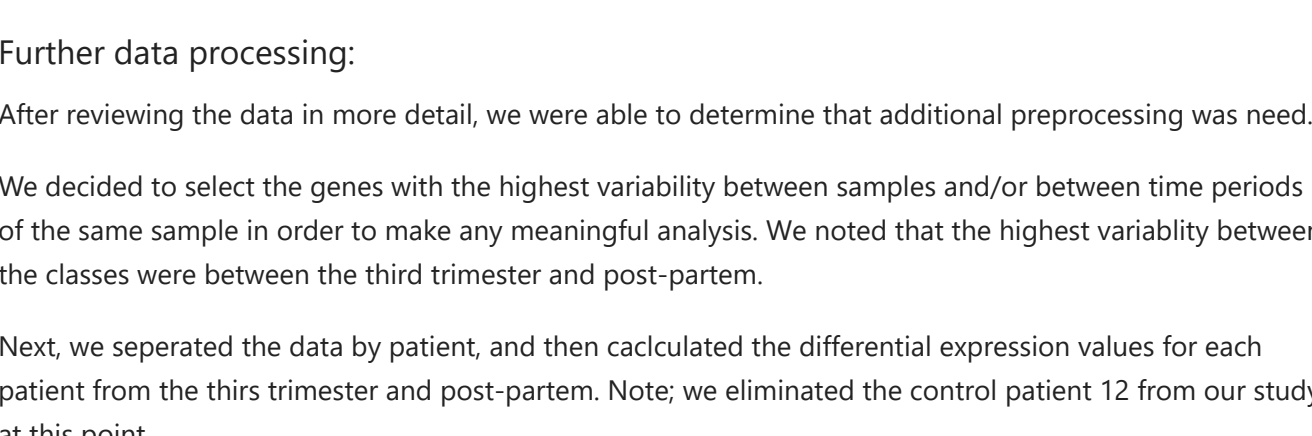
8 rows × 33297 columns

Data Preprocessing :

We started with the standard data preprocessing steps; removing any missing values in the data, and transposing the dataset to be the proper shape needed for running our models

Normalize the data across samples:

In order to be able to compare our samples against each other we needed to normalize the data or account for any outliers. Here we used Log2 base to normalize the data. Log2 fold change was suggested by a classmate during our presentation. Log2 fold is a great tool to help differentiate between original and fold change of a sample using a log2 base for all samples to hover around 0 and 1.



Predicting Pregnancy Stage from Cell free mRNA Levels

Supplementary document: Python code -

Jessie Bologna & Lindsay Reisman

NYU Bioinformatic - Machine Learning

5/15/21

Load the libraries

```
In [5]: import pandas as pd
import numpy as np
from sklearn.metrics import accuracy_score
import matplotlib.pyplot as plt
import seaborn as sns
import warnings
import scipy.stats as stats
from sklearn.preprocessing import MinMaxScaler
from sklearn import svm
from sklearn.linear_model import LogisticRegression
from sklearn.model_selection import cross_val_score
from sklearn.metrics import classification_report
from sklearn.metrics import accuracy_score
from sklearn.svm import SVC
from sklearn.neighbors import KNeighborsClassifier
from sklearn.naive_bayes import GaussianNB
from sklearn.tree import DecisionTreeClassifier
from sklearn.discriminant_analysis import LinearDiscriminantAnalysis
from sklearn.model_selection import GridSearchCV
from sklearn.metrics import confusion_matrix
from sklearn.linear_model import LinearRegression
from sklearn.neural_network import MLPClassifier

In [6]: # Ignore unimportant warnings
def warnings.warn("deprecated", DeprecationWarning)
warnings.filterwarnings("ignore")
with warnings.catch_warnings():
    warnings.simplefilter("ignore")
```

Import the Data

The data consists of 48 samples (12 patients) over the course of 5 stages of pregnancy: first, second, and third trimester, postpartum, and control/non-pregnant

Microarray data of cell free RNA was collected across 33,297 genes for each patient during each of these time courses

```
In [36]: # import the data
df = pd.read_csv('GSE56899_series_matrix.txt', sep = '\t')
# define the column names
columns = ['l1_R2', 'l1_R3', 'l1_R4', 'l1_R5', 'l1_R6', 'l1_R7', 'l1_R8', 'l1_R9', 'l1_R10', 'l1_R11', 'l1_R12', 'l1_R13', 'l1_R14', 'l1_R15', 'l1_R16', 'l1_R17', 'l1_R18', 'l1_R19', 'l1_R20', 'l1_R21', 'l1_R22', 'l1_R23', 'l1_R24', 'l1_R25', 'l1_R26', 'l1_R27', 'l1_R28', 'l1_R29', 'l1_R30', 'l1_R31', 'l1_R32', 'l1_R33', 'l1_R34', 'l1_R35', 'l1_R36', 'l1_R37', 'l1_R38', 'l1_R39', 'l1_R40', 'l1_R41', 'l1_R42', 'l1_R43', 'l1_R44', 'l1_R45', 'l1_R46', 'l1_R47', 'l1_R48', 'l1_R49', 'l1_R50', 'l1_R51', 'l1_R52', 'l1_R53', 'l1_R54', 'l1_R55', 'l1_R56', 'l1_R57', 'l1_R58', 'l1_R59', 'l1_R60', 'l1_R61', 'l1_R62', 'l1_R63', 'l1_R64', 'l1_R65', 'l1_R66', 'l1_R67', 'l1_R68', 'l1_R69', 'l1_R70', 'l1_R71', 'l1_R72', 'l1_R73', 'l1_R74', 'l1_R75', 'l1_R76', 'l1_R77', 'l1_R78', 'l1_R79', 'l1_R80', 'l1_R81', 'l1_R82', 'l1_R83', 'l1_R84', 'l1_R85', 'l1_R86', 'l1_R87', 'l1_R88', 'l1_R89', 'l1_R90', 'l1_R91', 'l1_R92', 'l1_R93', 'l1_R94', 'l1_R95', 'l1_R96', 'l1_R97', 'l1_R98', 'l1_R99', 'l1_R100', 'l1_R101', 'l1_R102', 'l1_R103', 'l1_R104', 'l1_R105', 'l1_R106', 'l1_R107', 'l1_R108', 'l1_R109', 'l1_R110', 'l1_R111', 'l1_R112', 'l1_R113', 'l1_R114', 'l1_R115', 'l1_R116', 'l1_R117', 'l1_R118', 'l1_R119', 'l1_R120', 'l1_R121', 'l1_R122', 'l1_R123', 'l1_R124', 'l1_R125', 'l1_R126', 'l1_R127', 'l1_R128', 'l1_R129', 'l1_R130', 'l1_R131', 'l1_R132', 'l1_R133', 'l1_R134', 'l1_R135', 'l1_R136', 'l1_R137', 'l1_R138', 'l1_R139', 'l1_R140', 'l1_R141', 'l1_R142', 'l1_R143', 'l1_R144', 'l1_R145', 'l1_R146', 'l1_R147', 'l1_R148', 'l1_R149', 'l1_R150', 'l1_R151', 'l1_R152', 'l1_R153', 'l1_R154', 'l1_R155', 'l1_R156', 'l1_R157', 'l1_R158', 'l1_R159', 'l1_R160', 'l1_R161', 'l1_R162', 'l1_R163', 'l1_R164', 'l1_R165', 'l1_R166', 'l1_R167', 'l1_R168', 'l1_R169', 'l1_R170', 'l1_R171', 'l1_R172', 'l1_R173', 'l1_R174', 'l1_R175', 'l1_R176', 'l1_R177', 'l1_R178', 'l1_R179', 'l1_R180', 'l1_R181', 'l1_R182', 'l1_R183', 'l1_R184', 'l1_R185', 'l1_R186', 'l1_R187', 'l1_R188', 'l1_R189', 'l1_R190', 'l1_R191', 'l1_R192', 'l1_R193', 'l1_R194', 'l1_R195', 'l1_R196', 'l1_R197', 'l1_R198', 'l1_R199', 'l1_R200', 'l1_R201', 'l1_R202', 'l1_R203', 'l1_R204', 'l1_R205', 'l1_R206', 'l1_R207', 'l1_R208', 'l1_R209', 'l1_R210', 'l1_R211', 'l1_R212', 'l1_R213', 'l1_R214', 'l1_R215', 'l1_R216', 'l1_R217', 'l1_R218', 'l1_R219', 'l1_R220', 'l1_R221', 'l1_R222', 'l1_R223', 'l1_R224', 'l1_R225', 'l1_R226', 'l1_R227', 'l1_R228', 'l1_R229', 'l1_R230', 'l1_R231', 'l1_R232', 'l1_R233', 'l1_R234', 'l1_R235', 'l1_R236', 'l1_R237', 'l1_R238', 'l1_R239', 'l1_R240', 'l1_R241', 'l1_R242', 'l1_R243', 'l1_R244', 'l1_R245', 'l1_R246', 'l1_R247', 'l1_R248', 'l1_R249', 'l1_R250', 'l1_R251', 'l1_R252', 'l1_R253', 'l1_R254', 'l1_R255', 'l1_R256', 'l1_R257', 'l1_R258', 'l1_R259', 'l1_R260', 'l1_R261', 'l1_R262', 'l1_R263', 'l1_R264', 'l1_R265', 'l1_R266', 'l1_R267', 'l1_R268', 'l1_R269', 'l1_R270', 'l1_R271', 'l1_R272', 'l1_R273', 'l1_R274', 'l1_R275', 'l1_R276', 'l1_R277', 'l1_R278', 'l1_R279', 'l1_R280', 'l1_R281', 'l1_R282', 'l1_R283', 'l1_R284', 'l1_R285', 'l1_R286', 'l1_R287', 'l1_R288', 'l1_R289', 'l1_R290', 'l1_R291', 'l1_R292', 'l1_R293', 'l1_R294', 'l1_R295', 'l1_R296', 'l1_R297', 'l1_R298', 'l1_R299', 'l1_R300', 'l1_R301', 'l1_R302', 'l1_R303', 'l1_R304', 'l1_R305', 'l1_R306', 'l1_R307', 'l1_R308', 'l1_R309', 'l1_R310', 'l1_R311', 'l1_R312', 'l1_R313', 'l1_R314', 'l1_R315', 'l1_R316', 'l1_R317', 'l1_R318', 'l1_R319', 'l1_R320', 'l1_R321', 'l1_R322', 'l1_R323', 'l1_R324', 'l1_R325', 'l1_R326', 'l1_R327', 'l1_R328', 'l1_R329', 'l1_R330', 'l1_R331', 'l1_R332', 'l1_R333', 'l1_R334', 'l1_R335', 'l1_R336', 'l1_R337', 'l1_R338', 'l1_R339', 'l1_R340', 'l1_R341', 'l1_R342', 'l1_R343', 'l1_R344', 'l1_R345', 'l1_R346', 'l1_R347', 'l1_R348', 'l1_R349', 'l1_R350', 'l1_R351', 'l1_R352', 'l1_R353', 'l1_R354', 'l1_R355', 'l1_R356', 'l1_R357', 'l1_R358', 'l1_R359', 'l1_R360', 'l1_R361', 'l1_R362', 'l1_R363', 'l1_R364', 'l1_R365', 'l1_R366', 'l1_R367', 'l1_R368', 'l1_R369', 'l1_R370', 'l1_R371', 'l1_R372', 'l1_R373', 'l1_R374', 'l1_R375', 'l1_R376', 'l1_R377', 'l1_R378', 'l1_R379', 'l1_R380', 'l1_R381', 'l1_R382', 'l1_R383', 'l1_R384', 'l1_R385', 'l1_R386', 'l1_R387', 'l1_R388', 'l1_R389', 'l1_R390', 'l1_R391', 'l1_R392', 'l1_R393', 'l1_R394', 'l1_R395', 'l1_R396', 'l1_R397', 'l1_R398', 'l1_R399', 'l1_R400', 'l1_R401', 'l1_R402', 'l1_R403', 'l1_R404', 'l1_R405', 'l1_R406', 'l1_R407', 'l1_R408', 'l1_R409', 'l1_R410', 'l1_R411', 'l1_R412', 'l1_R413', 'l1_R414', 'l1_R415', 'l1_R416', 'l1_R417', 'l1_R418', 'l1_R419', 'l1_R420', 'l1_R421', 'l1_R422', 'l1_R423', 'l1_R424', 'l1_R425', 'l1_R426', 'l1_R427', 'l1_R428', 'l1_R429', 'l1_R430', 'l1_R431', 'l1_R432', 'l1_R433', 'l1_R434', 'l1_R435', 'l1_R436', 'l1_R437', 'l1_R438', 'l1_R439', 'l1_R440', 'l1_R441', 'l1_R442', 'l1_R443', 'l1_R444', 'l1_R445', 'l1_R446', 'l1_R447', 'l1_R448', 'l1_R449', 'l1_R450', 'l1_R451', 'l1_R452', 'l1_R453', 'l1_R454', 'l1_R455', 'l1_R456', 'l1_R457', 'l1_R458', 'l1_R459', 'l1_R460', 'l1_R461', 'l1_R462', 'l1_R463', 'l1_R464', 'l1_R465', 'l1_R466', 'l1_R467', 'l1_R468', 'l1_R469', 'l1_R470', 'l1_R471', 'l1_R472', 'l1_R473', 'l1_R474', 'l1_R475', 'l1_R476', 'l1_R477', 'l1_R478', 'l1_R479', 'l1_R480', 'l1_R481', 'l1_R482', 'l1_R483', 'l1_R484', 'l1_R485', 'l1_R486', 'l1_R487', 'l1_R488', 'l1_R489', 'l1_R490', 'l1_R491', 'l1_R492', 'l1_R493', 'l1_R494', 'l1_R495', 'l1_R496', 'l1_R497', 'l1_R498', 'l1_R499', 'l1_R500', 'l1_R501', 'l1_R502', 'l1_R503', 'l1_R504', 'l1_R505', 'l1_R506', 'l1_R507', 'l1_R508', 'l1_R509', 'l1_R510', 'l1_R511', 'l1_R512', 'l1_R513', 'l1_R514', 'l1_R515', 'l1_R516', 'l1_R517', 'l1_R518', 'l1_R519', 'l1_R520', 'l1_R521', 'l1_R522', 'l1_R523', 'l1_R524', 'l1_R525', 'l1_R526', 'l1_R527', 'l1_R528', 'l1_R529', 'l1_R530', 'l1_R531', 'l1_R532', 'l1_R533', 'l1_R534', 'l1_R535', 'l1_R536', 'l1_R537', 'l1_R538', 'l1_R539', 'l1_R540', 'l1_R541', 'l1_R542', 'l1_R543', 'l1_R544', 'l1_R545', 'l1_R546', 'l1_R547', 'l1_R548', 'l1_R549', 'l1_R550', 'l1_R551', 'l1_R552', 'l1_R553', 'l1_R554', 'l1_R555', 'l1_R556', 'l1_R557', 'l1_R558', 'l1_R559', 'l1_R560', 'l1_R561', 'l1_R562', 'l1_R563', 'l1_R564', 'l1_R565', 'l1_R566', 'l1_R567', 'l1_R568', 'l1_R569', 'l1_R570', 'l1_R571', 'l1_R572', 'l1_R573', 'l1_R574', 'l1_R575', 'l1_R576', 'l1_R577', 'l1_R578', 'l1_R579', 'l1_R580', 'l1_R581', 'l1_R582', 'l1_R583', 'l1_R584', 'l1_R585', 'l1_R586', 'l1_R587', 'l1_R588', 'l1_R589', 'l1_R590', 'l1_R591', 'l1_R592', 'l1_R593', 'l1_R594', 'l1_R595', 'l1_R596', 'l1_R597', 'l1_R598', 'l1_R599', 'l1_R600', 'l1_R601', 'l1_R602', 'l1_R603', 'l1_R604', 'l1_R605', 'l1_R606', 'l1_R607', 'l1_R608', 'l1_R609', 'l1_R610', 'l1_R611', 'l1_R612', 'l1_R613', 'l1_R614', 'l1_R615', 'l1_R616', 'l1_R617', 'l1_R618', 'l1_R619', 'l1_R620', 'l1_R621', 'l1_R622', 'l1_R623', 'l1_R624', 'l1_R625', 'l1_R626', 'l1_R627', 'l1_R628', 'l1_R629', 'l1_R630', 'l1_R631', 'l1_R632', 'l1_R633', 'l1_R634', 'l1_R635', 'l1_R636', 'l1_R637', 'l1_R638', 'l1_R639', 'l1_R640', 'l1_R641', 'l1_R642', 'l1_R643', 'l1_R644', 'l1_R645', 'l1_R646', 'l1_R647', 'l1_R648', 'l1_R649', 'l1_R650', 'l1_R651', 'l1_R652', 'l1_R653', 'l1_R654', 'l1_R655', 'l1_R656', 'l1_R657', 'l1_R658', 'l1_R659', 'l1_R660', 'l1_R661', 'l1_R662', 'l1_R663', 'l1_R664', 'l1_R665', 'l1_R666', 'l1_R667', 'l1_R668', 'l1_R669', 'l1_R670', 'l1_R671', 'l1_R672', 'l1_R673', 'l1_R674', 'l1_R675', 'l1_R676', 'l1_R677', 'l1_R678', 'l1_R679', 'l1_R680', 'l1_R681', 'l1_R682', 'l1_R683', 'l1_R684', 'l1_R685', 'l1_R686', 'l1_R687', 'l1_R688', 'l1_R689', 'l1_R690', 'l1_R691', 'l1_R692', 'l1_R693', 'l1_R694', 'l1_R695', 'l1_R696', 'l1_R697', 'l1_R698', 'l1_R699', 'l1_R700', 'l1_R701', 'l1_R702', 'l1_R703', 'l1_R704', 'l1_R705', 'l1_R706', 'l1_R707', 'l1_R708', 'l1_R709', 'l1_R710', 'l1_R711', 'l1_R712', 'l1_R713', 'l1_R714', 'l1_R715', 'l1_R716', 'l1_R717', 'l1_R718', 'l1_R719', 'l1_R720', 'l1_R721', 'l1_R722', 'l1_R723', 'l1_R724', 'l1_R725', 'l1_R726', 'l1_R727', 'l1_R728', 'l1_R729', 'l1_R730', 'l1_R731', 'l1_R732', 'l1_R733', 'l1_R734', 'l1_R735', 'l1_R736', 'l1_R737', 'l1_R738', 'l1_R739', 'l1_R740', 'l1_R741', 'l1_R742', 'l1_R743', 'l1_R744', 'l1_R745', 'l1_R746', 'l1_R747', 'l1_R748', 'l1_R749', 'l1_R750', 'l1_R751', 'l1_R752', 'l1_R753', 'l1_R754', 'l1_R755', 'l1_R756', 'l1_R757', 'l1_R758', 'l1_R759', 'l1_R760', 'l1_R761', 'l1_R762', 'l1_R763', 'l1_R764', 'l1_R765', 'l1_R766', 'l1_R767', 'l1_R768', 'l1_R769', 'l1_R770', 'l1_R771', 'l1_R772', 'l1_R773', 'l1_R774', 'l1_R775', 'l1_R776', 'l1_R777', 'l1_R778', 'l1_R779', 'l1_R780', 'l1_R781', 'l1_R782', 'l1_R783', 'l1_R784', 'l1_R785', 'l1_R786', 'l1_R787', 'l1_R788', 'l1_R789', 'l1_R790', 'l1_R791', 'l1_R792', 'l1_R793', 'l1_R794', 'l1_R795', 'l1_R796', 'l1_R797', 'l1_R798', 'l1_R799', 'l1_R800', 'l1_R801', 'l1_R802', 'l1_R803', 'l1_R804', 'l1_R805', 'l1_R806', 'l1_R807', 'l1_R808', 'l1_R809', 'l1_R810', 'l1_R811', 'l1_R812', 'l1_R813', 'l1_R814', 'l1_R815', 'l1_R816', 'l1_R817', 'l1_R818', 'l1_R819', 'l1_R820', 'l1_R821', 'l1_R822', 'l1_R823', 'l1_R824', 'l1_R825', 'l1_R826', 'l1_R827', 'l1_R828', 'l1_R829', 'l1_R830', 'l1_R831', 'l1_R832', 'l1_R833', 'l1_R834', 'l1_R835', 'l1_R836', 'l1_R837', 'l1_R838', 'l1_R839', 'l1_R840', 'l1_R841', 'l1_R842', 'l1_R843', 'l1_R844', 'l1_R845', 'l1_R846', 'l1_R847', 'l1_R848', 'l1_R849', 'l1_R850', 'l1_R851', 'l1_R852', 'l1_R853', 'l1_R854', 'l1_R855', 'l1_R856', 'l1_R857', 'l1_R858', 'l1_R859', 'l1_R860', 'l1_R861', 'l1_R862', 'l1_R863', 'l1_R864', 'l1_R865', 'l1_R866', 'l1_R867', 'l1_R868', 'l1_R869', 'l1_R870', 'l1_R871', 'l1_R872', 'l1_R873', 'l1_R874', 'l1_R875', 'l1_R876', 'l1_R877', 'l1_R878', 'l1_R879', 'l1_R880', 'l1_R881', 'l1_R882', 'l1_R883', 'l1_R884', 'l1_R885', 'l1_R886', 'l1_R887', 'l1_R888', 'l1_R889', 'l1_R890', 'l1_R891', 'l1_R892', 'l1_R893', 'l1_R894', 'l1_R895', 'l1_R896', 'l1_R897', 'l1_R898', 'l1_R899', 'l1_R900', 'l1_R901', 'l1_R902', 'l1_R903', 'l1_R904', 'l1_R905', 'l1_R906', 'l1_R907', 'l1_R908', 'l1_R909', 'l1_R910', 'l1_R911', 'l1_R912', 'l1_R913', 'l1_R914', 'l1_R915', 'l1_R916', 'l1_R917', 'l1_R918', 'l1_R919', 'l1_R920', 'l1_R921', 'l1_R922', 'l1_R923', 'l1_R924', 'l1_R925', 'l1_R926', 'l1_R927', 'l1_R928', 'l1_R929', 'l1_R930', 'l1_R931', 'l1_R932', 'l1_R933', 'l1_R934', 'l1_R935', 'l1_R936', 'l1_R937', 'l1_R938', 'l1_R939', 'l1_R940', 'l1_R941', 'l1_R942', 'l1_R943', 'l1_R944', 'l1_R945', 'l1_R946', 'l1_R947', 'l1_R948', 'l1_R949', 'l1_R950', 'l1_R951', 'l1_R952', 'l1_R953', 'l1_R954', 'l1_R955', 'l1_R956', 'l1_R957', 'l1_R958', 'l1_R959', 'l1_R960', 'l1_R961', 'l1_R962', 'l1_R963', 'l1_R964', 'l1_R965', 'l1_R966', 'l1_R967', 'l1_R968', 'l1_R969', 'l1_R970', 'l1_R971', 'l1_R972', 'l1_R973', 'l1_R974', 'l1_R975', 'l1_R976', 'l1_R977', 'l1_R978', 'l1_R979', 'l1_R980', 'l1_R981', 'l1_R982', 'l1_R983', 'l1_R984', 'l1_R985', 'l1_R986', 'l1_R987', 'l1_R988', 'l1_R989', 'l1_R990', 'l1_R991', 'l1_R992', 'l1_R993', 'l1_R994', 'l1_R995', 'l1_R996', 'l1_R997', 'l1_R998', 'l1_R999', 'l1_R1000', 'l1_R1001', 'l1_R1002', 'l1_R1003', 'l1_R1004', 'l1_R1005', 'l1_R1006', 'l1_R1007', 'l1_R1008', 'l1_R1009', 'l1_R1010', 'l1_R1011', 'l1_R1012', 'l1_R1013', 'l1_R1014', 'l1_R1015', 'l1_R1016', 'l1_R1017', 'l1_R1018', 'l1_R1019', 'l1_R1020', 'l1_R1021', 'l1_R1022', 'l1_R1023', 'l1_R1024', 'l1_R1025', 'l1_R1026', 'l1_R1027', 'l1_R1028', 'l1_R1029', 'l1_R1030', 'l1_R1031', 'l1_R1032', 'l1_R1033', 'l1_R1034', 'l1_R1035', 'l1_R1036', 'l1_R1037', 'l1_R1038', 'l1_R1039', 'l1_R1040', 'l1_R1041', 'l1_R1042', 'l1_R1043', 'l1_R1044', 'l1_R1045', 'l1_R1046', 'l1_R1047', 'l1_R1048', 'l1_R1049', 'l1_R1050', 'l1_R1051', 'l1_R1052', 'l1_R1053', 'l1_R1054', 'l1_R1055', 'l1_R1056', 'l1_R1057', 'l1_R1058', 'l1_R1059', 'l1_R1060', 'l1_R1061', 'l1_R1062', 'l1_R1063', 'l1_R1064', 'l1_R1065', 'l1_R1066', 'l1_R1067', 'l1_R1068', 'l1_R1069', 'l1_R1070', 'l1_R1071', 'l1_R1072', 'l1_R1073', 'l1_R1074', 'l1_R1075', 'l1_R1076', 'l1_R1077', 'l1_R1078', 'l1_R1079', 'l1_R1080', 'l1_R1081', 'l1_R1082', 'l1_R1083', 'l1_R1084', 'l1_R1085', 'l1_R1086', 'l1_R1087', 'l1_R1088', 'l1_R1089', 'l1_R1090', 'l1_R1091', 'l1_R1092', 'l1_R1093', 'l1_R1094', 'l1_R1095', 'l1_R1096', 'l1_R1097', 'l1_R1098', 'l1_R1099', 'l1_R1100', 'l1_R1101', 'l1_R1102', 'l1_R1103', 'l1_R1104', 'l1_R1105', 'l1_R1106', 'l1_R1107', 'l1_R1108', 'l1_R1109', 'l1_R1110', 'l1_R1111', 'l1_R1112', 'l1_R1113', 'l1_R1114', 'l1_R1115', 'l1_R1116', 'l1_R1117', 'l1_R1118', 'l1_R1119', 'l1_R1120', 'l1_R1121', 'l1_R1122', 'l1_R1123', 'l1_R1124', 'l1_R1125', 'l1_R1126', 'l1_R1127', 'l1_R1128', 'l1_R1129', 'l1_R1130', 'l1_R1131', 'l1_R1132', 'l1_R1133', 'l1_R1134', 'l1_R1135', 'l1_R1136', 'l1_R1137', 'l1_R1138', 'l1_R1139', 'l1_R1140', 'l1_R1141', 'l1_R1142', 'l1_R1143', 'l1_R1144', 'l1_R1145', 'l1_R1146', 'l1_R1147', 'l1_R1148', 'l1_R1149', 'l1_R1150', 'l1_R1151', 'l1_R1152', 'l1_R1153', 'l1_R1154', 'l1_R1155', 'l1_R1156', 'l1_R1157', 'l1_R1158', 'l1_R1159', 'l1_R1160', 'l1_R1161', 'l1_R1162', 'l1_R1163', 'l1_R1164', 'l1_R1165', 'l1_R1166', 'l1_R1167', 'l1_R1168', 'l1_R1169', 'l1_R1170', 'l1_R1171', 'l1_R1172', 'l1_R1173', 'l1_R1174', 'l1_R1175', 'l1_R1176', 'l1_R1177', 'l1_R1178', 'l1_R1179', 'l1_R1180', 'l1_R1181', 'l1_R1182', 'l1_R1183', 'l1_R1184', 'l1_R1185', 'l1_R1186', 'l1_R1187', 'l1_R1188', 'l1_R1189', 'l1_R1190', 'l1_R1191', 'l1_R1192', 'l1_R1193', 'l1_R1194', 'l1_R1195', 'l1_R1196', 'l1_R1197', 'l1_R1198', 'l1_R1199', 'l1_R1200', 'l1_R1201', 'l1_R1202', 'l1_R1203', 'l1_R1204', 'l1_R1205', 'l1_R1206', 'l1_R1207', 'l1_R1208', 'l1_R1209', 'l1_R1210', 'l1_R1211', 'l1_R1212', 'l1_R1213', 'l1_R1214', 'l1_R1215', 'l1_R1216', 'l1_R1217', 'l1_R1218', 'l1_R1219', 'l1_R1220', 'l1_R1221', 'l1_R1222', 'l1_R1223', 'l1_R1224', 'l1_R1225', 'l1_R1226', 'l1_R1227', 'l1_R1228', 'l1_R1229', 'l1_R1230', 'l1_R1231', 'l1_R1232', 'l1_R1233', 'l1_R1234', 'l1_R1235', 'l1_R1236', 'l1_R1237', 'l1_R1238', 'l1_R1239', 'l1_R1240', 'l1_R1241', 'l1_R1242', 'l1_R1243', 'l1_R1244', 'l1_R1245', 'l1_R1246', 'l1_R1247', 'l1_R1248', 'l1_R1249', 'l1_R1250', 'l1_R1251', 'l1_R1252', 'l1_R1253', 'l1_R1254', 'l1_R1255', 'l1_R1256', 'l1_R1257', 'l1_R1258', 'l1_R1259', 'l1_R1260', 'l1_R1261', 'l1_R1262', 'l1_R1263', 'l1_R1264', 'l1_R1265', 'l1_R1266', 'l1_R1267', 'l1_R1268', 'l1_R1269', 'l1_R1270', 'l1_R1271', 'l1_R1272', 'l1_R1273', 'l1_R1274', 'l1_R1275', 'l1_R1276', 'l1_R1277', 'l1_R1278', 'l1_R1279', 'l1_R1280', 'l1_R1281', 'l1_R1282', 'l1_R1283', 'l1_R1284', 'l1_R1285', 'l1_R1286', 'l1_R1287', 'l1_R1288', 'l1_R1289', 'l1_R1290', 'l1_R1291', 'l1_R1292', 'l1_R1293', 'l1_R1294', 'l1_R1295', 'l1_R1296', 'l1_R1297', 'l1_R1298', 'l1_R1299', 'l1_R1300', 'l1_R1301', 'l1_R1302', 'l1_R1303', 'l1_R1304', 'l1_R1305', 'l1_R1306', 'l1_R1307', 'l1_R1308', 'l1_R1309', 'l1_R1310', 'l1_R1311', 'l1_R1312', 'l1_R1313', 'l1_R1314', 'l1_R1315', 'l1_R1316', 'l1_R1317', 'l1_R1318', 'l1_R1319', 'l1_R1320', 'l1_R1321', 'l1_R1322', 'l1_R1323', 'l1_R1324', 'l1_R1325', 'l1_R1326', 'l1_R1327', 'l1_R1328', 'l1_R1329', 'l1_R1330', 'l1_R1331', 'l1_R1332', 'l1_R1333', 'l1_R1334', 'l1_R1335', 'l1_R1336', 'l1_R1337', 'l1_R1338', 'l1_R1339', 'l1_R1340', 'l1_R1341', 'l1_R1342', 'l1_R1343', 'l1_R1344', 'l1_R1345', 'l1_R1346', 'l1_R1347', 'l1_R1348', 'l1_R1349', 'l1_R1350', 'l1_R1351', 'l1_R1352', 'l1_R1353', 'l1_R1354', 'l1_R1355', 'l1_R1356', 'l1_R1357', 'l1_R1358', 'l1_R1359', 'l1_R1360', 'l1_R1361', 'l1_R1362', 'l1_R1363', 'l1_R1364', 'l1_R1365', 'l1_R1366', 'l1_R1367', 'l1_R1368', 'l1_R1369', 'l1_R1370', 'l1_R1371', 'l1_R1372', 'l1_R1373', 'l1_R1374', 'l1_R1375', 'l1_R1376', 'l1_R1377', 'l1_R1378', 'l1_R1379', 'l1_R1380', 'l1_R1381', 'l1_R1382', 'l1_R1383', 'l1_R1384', 'l1_R1385', 'l1_R1386', 'l1_R1387', 'l1_R1388', 'l1_R1389', 'l1_R1390', 'l1_R1391', 'l1_R1392', 'l1_R1393', 'l1_R1394', 'l1_R1395', 'l1_R1396', 'l1_R1397', 'l1_R1398', 'l1_R1399', 'l1_R1400', 'l1_R1401', 'l1_R1402', 'l1_R1403', 'l1_R1404', 'l1_R1405', 'l1_R1406', 'l1_R1407', 'l1_R1408', 'l1_R1409', 'l1_R1410', 'l1_R1411', 'l1_R1412', 'l1_R1413', 'l1_R1414', 'l1_R1415', 'l1_R1416', 'l1_R1417', 'l1_R1418', 'l1_R1419', 'l1_R1420', 'l1_R1421', 'l
```



```
8.91575337e-01+0.00000000e+00j 7.5482407e-01+0.00000000e+00j
6.6300059e-01+0.00000000e+00j 5.8498929e-01+0.00000000e+00j
4.3327995e-01+0.00000000e+00j 3.2075315e-01+0.00000000e+00j
2.70728330e-01+0.00000000e+00j 2.1512526e-01+0.00000000e+00j
1.2672830e-01+0.00000000e+00j 9.8822802e-02+0.00000000e+00j
1.36730472e-01+0.00000000e+00j 4.7566831e-02+0.00000000e+00j
7.0326848e-02+0.00000000e+00j 4.05100450e-02+0.00000000e+00j
3.33754136e-02+0.00000000e+00j 1.05100450e-02+0.00000000e+00j
-4.73778712e-16+0.00000000e+00j -4.49437527e-16+0.00000000e+00j
2.89865877e-16+0.00000000e+00j 1.25413641e-16+1.58859276e-16j
1.25413641e-16-1.58879276e-16j 1.26445382e-17+9.0865342e-17j
1.26445382e-17-9.08683342e-17j -1.23028549e-16+0.00000000e+00j
-2.20538142e-16+0.00000000e+00j
```

```
In [67]: tot = sum(eigen_vals)
vc = [(1/tot) for i in sorted(eigen_vals, reverse=True)]
cum_var_exp = np.cumsum(vc)

plt.ylabel('% Variance Explained')
plt.xlabel('% of Features')
plt.title('Score Plot')

plt.plot(cum_var_exp)
```

```
Out[67]: [matplotlib.lines.Line2D at 0x112068327c0>]
```

```
In [68]: eigen_pairs = (np.abs(eigen_vals[1:]),eigen_vecs[1:1]) for i in range(len(eigen_vals)):
eigen_pairs.sort(key=lambda x: x[0],reverse=True)

arr = np.hstack([eigen_pairs[0][1][:n.newaxis],
eigen_pairs[1][1][:n.newaxis]])

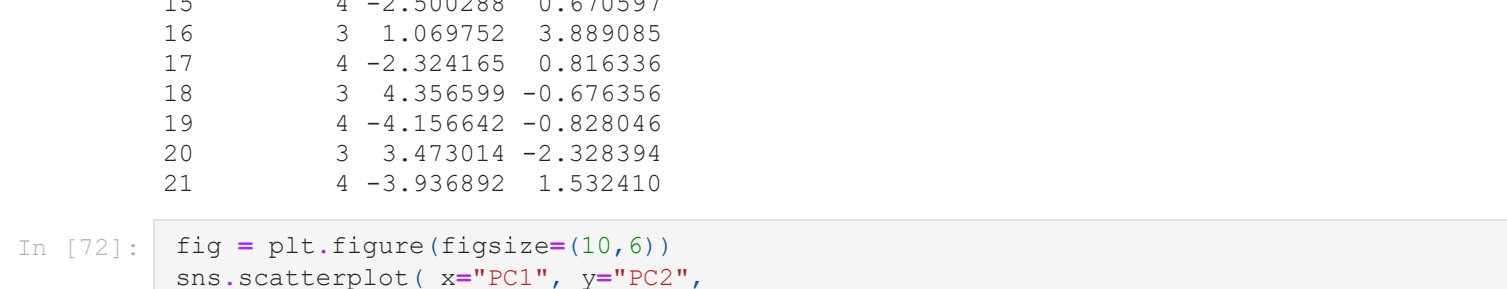
#arr
```

```
In [69]: # PCA across patient groups
X_pca = X.dot(arr)
pca = PCA(n_components=2)
Pc=pca.fit_transform(X)
pca_df = pd.DataFrame(data = Pc, columns = ['PC1', 'PC2'])
```

```
In [70]: d = pd.DataFrame(data=y1, columns = ['patients'])
pca_df1 = pd.concat([d, pca_df, axis=1])
print(pca_df1)
```

```
patients PC1 PC2
0 3 2.430827 -2.543240
1 4 -4.560385 -1.073866
2 3 4.386897 -0.693315
3 4 -3.053299 1.277703
4 3 5.117675 -2.190759
5 4 -3.744242 -0.762400
6 3 3.881749 -1.103985
7 4 -1.632839 0.011237
8 3 4.132520 -0.903661
9 4 -3.816763 -1.272172
10 3 2.641036 -0.398122
11 4 -4.644296 -1.029306
12 3 3.911312 4.909200
13 4 -3.820971 -0.340294
14 3 2.789401 2.973549
15 4 -2.500288 0.670597
16 3 1.069752 3.889085
17 4 -2.324165 0.816336
18 3 4.356599 -0.676356
19 4 -4.156642 -0.828046
20 3 3.473014 -2.328394
21 4 -3.936892 1.532410
```

```
In [72]: fig = plt.figure(figsize=(10,6))
sns.scatterplot(x='PC1', y='PC2',
data=pca_df1,
hue = 'patients',
legend='full') # specify the point size
plt.show()
```



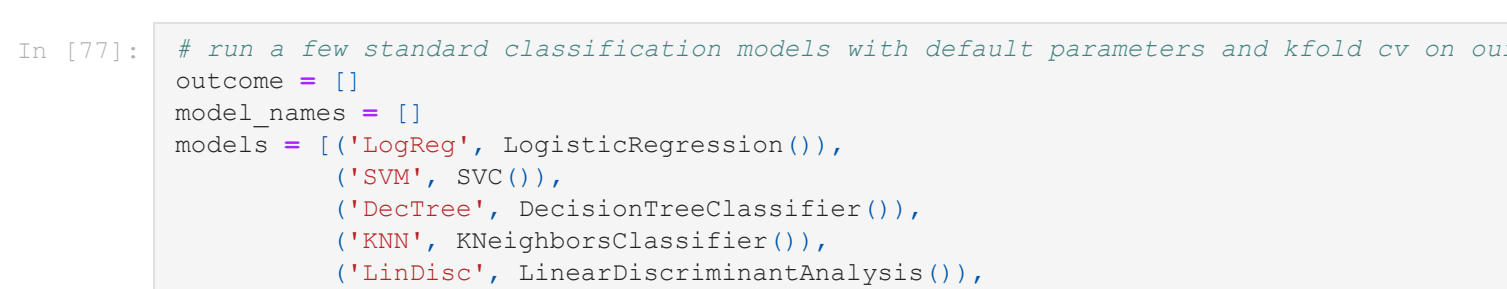
Test some classification models on new dataset -

```
In [77]: # run a few standard classification models with default parameters and kfold cv on our
outcome = []
model_names = []
models = ['LogReg', LogisticRegression()),
('SVM', SVC()),
('DecTree', DecisionTreeClassifier()),
('KNN', KNeighborsClassifier()),
('LinDisc', LinearDiscriminantAnalysis()),
('GaussianNB', GaussianNB()),
('RFC', RandomForestClassifier()),
('AdaBoost', AdaBoostClassifier()),
('MLP', MLPClassifier())

for model_name, model in models:
k_fold_validation = model_selection.KFold(n_splits=5)
results = model_selection.cross_val_score(model, X1, y1, cv=k_fold_validation, scoring='accuracy')
outcome.append(results)
model_names.append(model_name)
output_message = "%s Mean=%f STD=%f" % (model_name, results.mean(), results.std())
print(output_message)
```

```
LogRegl Mean=1.000000 STD=0.000000
SVMl Mean=0.960000 STD=0.080000
DecTree Mean=1.000000 STD=0.000000
KNNl Mean=1.000000 STD=0.000000
LinDiscl Mean=1.000000 STD=0.000000
GaussianNB Mean=1.000000 STD=0.000000
RFCl Mean=1.000000 STD=0.000000
AdaBoostl Mean=1.000000 STD=0.000000
MLPl Mean=1.000000 STD=0.000000
```

```
In [76]: fig = plt.figure(figsize=(12,4))
ax = fig.add_subplot(1,1,1)
models1 = ['LogReg', 'SVM', 'SVM-1', 'DecTree', 'DT_Grid', 'KNN', 'LinDisc', 'GaussianNB',
results = [100,80,100,60, 100, 100,95,95,100,90,100]
ax.bar(models1,results)
plt.title('Machine Learning Model Comparison', fontsize=(16))
plt.show()
```



Results:

Models: Classification and Prediction

Now our data is ready to be used to begin testing various machine learning classification models.

After testing our a few models we were able to determine that most of the models ran with a very high accuracy rate. SVM, decision trees, and random forest all performed the best when the hyper-parameters were carefully selected.

For the SVM models, after adjusting the kernel to be "linear" and the regularization parameter to 0.05 the performance increased from around 88% to 100% accuracy.

The decision tree model performed poorly with the default parameters, but when adjusted using gridsearch to find the best parameters the model improved from around 55% to 100%.

```
In [79]: # Separate into test and train data
from sklearn.model_selection import train_test_split

X_train, X_test, y_train, y_test = train_test_split(
X1, y1, test_size=0.4)
```

```
In [80]: # scaling the features - standardize
from sklearn.preprocessing import StandardScaler

scaler = StandardScaler()
scaler.fit(X_train)
X_train_std = scaler.transform(X_train)
y_test_std = scaler.transform(X_test)
```

```
In [81]: svc = svm.SVC(kernel='rbf', C=1)
svc.fit(X_train_std, y_train)

y_hat = svc.predict(X_test_std)

# Calculate accuracy
print(pd.crosstab(y_test, y_hat, rownames=['Actual'], colnames=['Predicted']))
print("Accuracy = ", accuracy_score(y_test, y_hat))

print(classification_report(y_test, y_hat))
```

```
Predicted 3 4
Actual 3 0
4 0 3
Accuracy = 1.0

precision recall f1-score support

3 1.00 1.00 1.00 6
4 1.00 1.00 1.00 3

accuracy 1.00 9
macro avg 1.00 1.00 1.00 9
weighted avg 1.00 1.00 1.00 9
```

```
In [82]: # Update the SVC model to see if we can get a better performance
svc1 = svm.SVC(kernel='linear', C=0.05, gamma='auto') # changed regularization, kernel
svc1.fit(X_train_std, y_train)

# add in cross validation
scores1 = cross_val_score(svc1, X_test_std, y_test, cv=5)
print(scores1)

print("%0.2f accuracy with a standard deviation of %0.2f" % (scores1.mean(), scores1.std()))

y_hat1 = svc1.predict(X_test_std)

# Calculate accuracy
print(pd.crosstab(y_test, y_hat1, rownames=['Actual'], colnames=['Predicted']))
print("Accuracy = ", accuracy_score(y_test, y_hat1))

print(classification_report(y_test, y_hat1))
```

```
[1. 1. 1. 1. 1.]
1.00 accuracy with a standard deviation of 0.00
Predicted 3 4
Actual 3 0
4 0 3
Accuracy = 1.0

precision recall f1-score support

3 1.00 1.00 1.00 6
4 1.00 1.00 1.00 3

accuracy 1.00 9
macro avg 1.00 1.00 1.00 9
weighted avg 1.00 1.00 1.00 9
```

```
In [83]: # Create decision tree with adaboost classifier
from sklearn.ensemble import AdaBoostClassifier
from sklearn.tree import DecisionTreeClassifier
from sklearn.metrics import confusion_matrix

# Create a decision tree object
tree = DecisionTreeClassifier()

# Fit on training data
tree.fit(X_train_std, y_train)
abc = AdaBoostClassifier(base_estimator=tree)

# Train Adaboost Classifier
model = abc.fit(X_train_std, y_train)

#Predict the response with test dataset
y_pred = model.predict(X_test_std)

print("Accuracy for Decision Tree model: ", accuracy_score(y_test, y_pred)*100) ## Still 55%

# Make a Confusion Matrix
print(pd.crosstab(y_test, y_pred, rownames=['Actual'], colnames=['Predicted']))
```

```
Accuracy for Decision Tree model: 88.88888888888889
Predicted 3 4
Actual 3 5
4 0 3
```

```
In [84]: # Try Gridsearch to tune the hyperparameters to see if we get better accuracy
from sklearn.model_selection import GridSearchCV
tree_param = {'criterion':['gini', 'entropy'],
'max_depth':[3,4,5,6,7,8,9,10,11,12,15,20],
'min_samples_split':[5,10,20,30,40,50,60,70,80,90,100],
'min_samples_leaf':[2,3,4,5,6,7,8,9,10],
'max_features':['list(range(1,X_train_std.shape[1])),
'max_features': ['auto', 'sqrt', 'log2'],
'presort':['True, False]]

grid = GridSearchCV(DecisionTreeClassifier(), tree_param, cv=4)
grid.fit(X_train_std, y_train)
#best_predict = grid.best_score_(y_test)
print("Best Accuracy = ",grid.best_score_)
print(grid.best_params)
print(grid.best_estimator_)
#print(pd.crosstab(y_test, best_predict, rownames=['Actual'], colnames=['Predicted']))
```

```
Best Accuracy = 1.0
['criterion': 'gini', 'max_depth': 3, 'max_features': 'auto', 'min_samples_leaf': 3,
'min_samples_split': 5, 'presort': True]
DecisionTreeClassifier(max_depth=3, max_features='auto', min_samples_leaf=3,
min_samples_split=5, presort=True)
```

```
In [85]: # Now use adaboost on the gridsearch results to see if we can get an even better pred
tree = grid.best_estimator_
tree.fit(X_train_std, y_train)

# Create adaboost classifier object
abc1 = AdaBoostClassifier(n_estimators=100,
learning_rate=5,
base_estimator=tree)

# Train Adaboost Classifier
model3 = abc1.fit(X_train_std, y_train)

# Predict the response for test dataset
y_hat2 = model3.predict(X_test_std)
print("Accuracy:",accuracy_score(y_test, y_hat2))

Accuracy: 0.8888888888888889
```

```
In [86]: from sklearn.ensemble import RandomForestClassifier
from sklearn.datasets import make_classification
classifier = RandomForestClassifier(n_estimators = 20,
criterion = 'entropy')
classifier.fit(X_train_std, y_train)
y_pred = classifier.predict(X_test_std)

# Making the Confusion Matrix
print(pd.crosstab(y_test, y_pred, rownames=['Actual'], colnames=['Predicted']))
print("Accuracy = ", accuracy_score(y_test, y_pred))

Predicted 3 4
Actual 3 1
4 0 3
Accuracy = 0.8888888888888889
```

```
In [87]: model4 = LogisticRegression().fit(X_train_std, y_train)
y_hat4 = model4.predict(X_test_std)

# Making the Confusion Matrix
print(pd.crosstab(y_test, y_hat4, rownames=['Actual'], colnames=['Predicted']))
print("Accuracy = ", accuracy_score(y_test, y_hat4))

Predicted 3 4
Actual 3 0
4 0 3
Accuracy = 1.0
```

References -

Koh W, Pan W, Gawad C, Fan HC et al. Noninvasive in vivo monitoring of tissue-specific global gene expression in humans. Proc Natl Acad Sci U S A 2014 May 20;111(20):7361-6. PMID: 24799715

Scikit-learn: Machine Learning in Python, Pedregosa et al, JMLR 12, pp. 2825-2830, 2011.