	 Converting the remaining columns to floats if necessary. Copying this version of the data (using the copy method) to a variable to preserve it. We will be using it later. data = pd.read_csv("C:\\Users\\rsnen\\Desktop\\IBM\\Wholesale_Customers_Data.csv", sep=',') data.shape (440, 8) data.head() Channel Region Fresh Milk Grocery Frozen Detergents_Paper Delicassen
	1 2 3 7057 9810 9568 1762 3293 1776 2 2 3 6353 8808 7684 2405 3516 7844 3 1 3 13265 1196 4221 6404 507 1788 4 2 3 22615 5410 7198 3915 1777 5185 data = data.drop(['Channel', 'Region'], axis=1) Fresh int64
((((((((((Milk int64 Grocery int64 Frozen int64 Detergents_Paper int64 Delicassen in
e	Scaling data Insuring the data is scaled and (relatively) normally distributed. Examining the correlation and skew. Performing any transformations and scale with a scaling method. Viewing the pairwise correlation plots of the new data. Examining the correlation With data.core we can see the correlation between each one of the different features. So, this will give us for each feature the correlation with all the other features in a square matrix in a square data frame. And just the square of the different features are square matrix in a square data frame.
a fo fr	nat we can get the highest correlation which feature is the highest correlation and because of one feature with itself will always have a correlation of one. We're going to replace that diagonal value which are going to all 1's with all 0's. So we're saying for x in the range of format.shape 0, it's a square matrix we could have called shape zero or shape one. So that's going to be for every single value in our matrix, for every single nur or the range of our matrix, we're going to take the diagonal value. So 0 0,1 1,2 2 and replace that 1 with a 0. And we can see now our correlation matrix has the correlation between fresh and milk and grocery. And the resh and fresh, it's just a 0 across each one of the different diagonals. corr_mat = data.corr() # Strip the diagonal for future examination for x in range(corr_mat.shape[0]): corr_mat.iloc[x,x] = 0.0 corr_mat.
	Fresh 0.00000 0.100510 -0.011854 0.345881 -0.101953 0.244690 Milk 0.100510 0.00000 0.728335 0.123994 0.661816 0.406368 Grocery -0.011854 0.728335 0.000000 -0.040193 0.924641 0.205497 Frozen 0.345881 0.123994 -0.040193 0.000000 -0.131525 0.390947 Detergents_Paper -0.101953 0.661816 0.924641 -0.131525 0.000000 0.069291 Delicassen 0.244690 0.406368 0.205497 0.390947 0.069291 0.000000 # the two categories with their respective most strongly correlated variable. corr_mat.abs().idxmax()
	Fresh Frozen Milk Grocery Grocery Detergents_Paper Frozen Delicassen Detergents_Paper Grocery Delicassen Milk dtype: object Examing the skew values Examing skew for each one of our ten values. And then take the long transformation, if necessary, for those that have higher skew. skew's going to be a value with 0 being no skew,positive value being a right skew,a legative value being a left skew. The higher that value is, the stronger the skew. Those that have a higher skew and we see here we have these values that have tend to have a higher skew. And for those, we're going to to take the log transformation with each, hopefully creating more normally displacements.
g	ata. So for calling each one of these log columns.index, so these are, this is our log columns that we just defined is that pandas series for called the index we get each one of these Delicatessen frozen milk and so color to also match up with each one of our different data columns. So we're going to place those columns in place with the log transformation of those columns. log_columns = data.skew().sort_values(ascending=False) # We sort them from largest to smallest. log_columns = log_columns.loc[log_columns > 0.75] # log columns that are greater than 0.75. log_columns Series([], dtype: float64) # The log transformations for col in log_columns.index: data[col] = np.log1p(data[col])
V	<pre>from sklearn.preprocessing import MinMaxScaler mms = MinMaxScaler() for col in data.columns: data[col] = mms.fit_transform(data[[col]]).squeeze() //sualize the relationship between the variables. sns.set_context('notebook') sns.set_style('white') sns.pairplot(data);</pre>
	1.0 0.8 0.6 0.0 0.0 1.0 0.8
	0.6
	0.2 0.0 1.0 0.8 0.6 0.2 0.0 0.0
	1.0 1.0 1.0 1.0 1.0 1.0 1.0 1.0 1.0 1.0
:	• Use Scikit-learn's pipeline function, recreate the data pre-processing scheme above (transformation and scaling) using a pipeline. If you used a non-Scikit learn function to transform the data (e.g. NumPy's log further class called FunctionTransformer.
u	• Use the pipeline to transform the original data that was stored at the earlier. • Compare the results to the original data to verify that everything worked. **Jote: Scikit-learn has a more flexible Pipeline function and a shortcut version called make_pipeline. Either can be used. Also, if different transformations need to be performed on the data, a FeatureUnio seed. **From sklearn.preprocessing import FunctionTransformer from sklearn.pipeline import Pipeline **The custom NumPy log transformer
T	estimators = [('logip', log_transformer), ('minmaxscale', MinMaxScaler())] pipeline = Pipeline(estimators) # Convert the original data data_pipe = pipeline.fit_transform(data_orig) The results are identical. Note that machine learning models and grid searches can also be added to the pipeline (and in fact, usually are.) #data pipe should equal that data that we just transformed. #So we're going to check that using NumPy.allclose, #which is just going to check that each value within each of our arrays are exactly the same, np.allclose(data_pipe, data)
F	PCA Performing PCA with n_components ranging from 1 to 5. Storing the amount of explained variance for each number of dimensions. Also store the feature importance for each number of dimensions. Hint: PCA doesn't explicitly provide this after a model is fit, but the components_ properties can be used to determine something that approxi importance. How you decided to do so is entirely up to you. Plot the explained variance and feature importances.
11 11 11 11 11 11 11 11 11 11 11 11 11	<pre>from sklearn.decomposition import PCA n = 2 pd.Series({'n':n, 'model':PCAmod,</pre>
(PCAmod.explained_variance_ratiosum() 0.9807028483254034 from sklearn.decomposition import PCA n = 2 pd.Series({'n':n, 'model':PCAmod, 'var': PCAmod.explained_variance_ratiosum()}) weights = PCAmod.explained_variance_ratioreshape(-1,1)/PCAmod.explained_variance_ratiosum() weights = PCAmod.explained_variance_ratioreshape(-1,1)/PCAmod.explained_variance_ratiosum() array([[0.41431035], [0.29503005], [0.1261552],
	[0.09762247], [0.06688193]]) from sklearn.decomposition import PCA n = 2 pd.Series({'n':n, 'model':PCAmod,
p fi T p	[0.02851832, 0.00452192, 0.01154961, 0.06543015, 0.03408488, 0.09749418], [0.08738228, 0.00401993, 0.00147339, 0.02125841, 0.00918582, 0.03660402], [0.09351675, 0.05111478, 0.00046533, 0.0017493, 0.04120455, 0.01212489]]) The larger these absolute values are, the more they contributed to each component and the more important that feature is. What we had here before is we took the absolute value, because we don't care about whether the institute or negative. We just care about how much it affected that principal component. Then we're weighting it according to these weights, weights are going to be how important each one the principal components are stone is going to be multiplied by 0.41, and this second one is going to be multiplied by 0.29, so that we don't put on too much weight. We see here that we use 70 percent of whatever feature this is; this is the fifth then we use 70 percent here in the second feature. The second PCA for a different feature, we want to ensure that these do not get equal weights. This should get a higher weight than this one since this is part of the virincipal component. Weights
	array([[0.41431035],
4	# how much each one of these different features with their weightings were able to comprise these #principal components that we have. array([0.26359241, 0.35535498, 0.15845903, 0.42512385, 0.3825994 ,
	<pre>for n in range(1, 6): # Create and fit the model PCAmod = PCA(n_components=n) PCAmod.fit(data) # Store the model and variance pca_list.append(pd.Series({'n':n, 'model':PCAmod,</pre>
	'features': data.columns, 'values':abs_feature_values.sum()})) pca_df = pd.concat(pca_list, axis=1).T.set_index('n') pca_df model var 1 PCA(n_components=1) 0.406315 2 PCA(n_components=2) 0.695652 3 PCA(n_components=3) 0.819373 4 PCA(n_components=4) 0.915112
	5 PCA(n_components=5) 0.980703 feature_weight_list[1] n features values 0 2 Fresh 0.114290 1 2 Milk 0.194238 2 Grocery 0.090570 3 2 Frozen 0.261864
C	4 2 Detergents_Paper 0.184878 5 2 Delicassen 0.154159 Create a table of feature importances for each data column. features_df = (pd.concat(feature_weight_list)
	1 0.077124 0.328639 0.041701 0.121100 0.15660 0.274776 2 0.154159 0.184878 0.114290 0.261864 0.090570 0.194238 3 0.235332 0.170622 0.115511 0.264779 0.076650 0.137106 4 0.233928 0.145971 0.209356 0.236063 0.061963 0.112720 5 0.211672 0.186937 0.177477 0.196504 0.051589 0.175820 Create a plot of explained variances. Sns. set_context('talk') ax = pca_df['var'].plot(kind='bar')
	ax.set(xlabel='Number of dimensions', ylabel='Percent explained variance vs Dimensions'); Explained Variance vs Dimensions 0.8 0.0 0.0 0.0 0.0 0.0 0.0 0.
A	Number of dimensions and here's a plot of feature importances. ax = features_df.plot(kind='bar', figsize=(13,8)) ax.legend(loc='upper right') ax.set(xlabel='Number of dimensions', ylabel='Relative importance', title='Feature importance vs Dimensions');
	0.30 0.25 Feature importance vs Dimensions Delicassen Detergents_Paper Frozen Grocery Mills
	0.20 0.15 0.10 0.05
Т	Number of dimensions Kernel PCA Fit a KernelPCA model with kernel='rbf'. You can choose how many components and what values to use for the other parameters (rbf refers to a Radial Basis Function kernel, and the gamma parar governs scaling of this kernel and typically ranges between 0 and 1). Several other kernels can be tried, and even passed ss cross validation parameters (see this example). If you want to tinker some more, use GridSearchCV to tune the parameters of the KernelPCA model. The second step is tricky since grid searches are generally used for supervised machine learning methods and rely on scoring metrics, such as accuracy, to determine the best model. However, a custom scoring function we written for GridSearchCV, where larger is better for the outcome of the scoring function.
٧	written for GridSearchCV, where larger is better for the outcome of the scoring function. What would such a metric involve for PCA? What about percent of explained variance? Or perhaps the negative mean squared error on the data once it has been transformed and then inversely transformed? from sklearn.decomposition import KernelPCA from sklearn.model_selection import GridSearchCV from sklearn.metrics import mean_squared_error # Custom scoreruse negative rmse of inverse transform def scorer(pcamodel, X, y=None): try: X_val = X.values except: X_val = X
	<pre># The grid search parameters param_grid = {'gamma':[0.001, 0.01, 0.05, 0.1, 0.5, 1.0],</pre>
L	kernelPCA.best_estimator_ KernelPCA(fit_inverse_transform=True, gamma=1.0, kernel='rbf', n_components=4) et's explore how our model accuracy may change if we include a PCA in our model building pipeline. Let's plan to use sklearn's Pipeline class and create a pipeline that has the following steps: 1. A scaler 2. $PCA(n_e ompo \neq nts = n)$ 3. $LogisticRegression$ • Write a function that takes in a value of n and makes the above pipeline, then predicts the "Activity" column over a 5-fold StratifiedShuffleSplit, and returns the average test accuracy • For various values of n , call the above function and store the average accuracies. • Plot the average accuracy by number of dimensions.
	<pre>Plot the average accuracy by number of dimensions. data = pd.read_csv("C:\\Users\\rsnen\\Desktop\\IBM\\Human_Activity_Recognition_Using_Smartphones_Data.csv", sep=',') data.columns Index(['tBodyAcc-mean()-X', 'tBodyAcc-mean()-Y', 'tBodyAcc-mean()-Z', 'tBodyAcc-std()-X', 'tBodyAcc-std()-Y', 'tBodyAcc-std()-Z', 'tBodyAcc-mad()-X', 'tBodyAcc-mad()-Y', 'tBodyAcc-mad()-Z', 'tBodyAcc-max()-X', 'tBodyAcc-max()-Y', 'tBodyAcc-mad()-Y', 'tBodyAcc-mad()-Z', 'tBodyAcc-max()-X', 'tBodyAcc-max()-X', 'tBodyAcc-max()-X', 'tBodyAcc-max()-X', 'tBodyAcc-max()-X', 'tBodyAcc-max()-X', 'tBodyAcc-max()-Y', 'tBodyA</pre>
le	'angle(Y,gravityMean)', 'angle(Z,gravityMean)', 'Activity'], dtype='object', length=562) data.shape (10299, 562) ets reduce number of columns with PCA from sklearn.pipeline import Pipeline from sklearn.preprocessing import StandardScaler from sklearn.model_selection import StratifiedShuffleSplit from sklearn.linear_model import LogisticRegression from sklearn.metrics import accuracy_score
	<pre>X = data.drop('Activity', axis=1) y = data.Activity sss = StratifiedShuffleSplit(n_splits=5, random_state=42) def get_avg_score(n): pipe = [</pre>
	<pre>0.8836893203883495, 0.9368932038834951, 0.965242718446602, 0.9716504854368931, 0.9786407766990293, 0.9842718446601942, 0.9852427184466019] sns.set_context('talk') ax = plt.axes() ax.plot(ns, score_list) ax.set(xlabel='Number of Dimensions',</pre>
	title='LogisticRegression Accuracy vs Number of dimensions on the Human Activity Dataset') ax.grid(True)
	LogisticRegression Accuracy vs Number of dimensions on the Human Activity Dataset Output Out