# CS 7641 Machine Learning
# Robert Bobkoskie
# \

# Assignment 3: Unsupervised Learning and Dimensionality Reduction

## 1. INTRODUCTION

The goal of this project is to explore the Unsupervised Learning (UL), Clustering: K-Means, Expectation Maximization, and Dimensionality Reduction (DR): Principal Component Analysis (PCA), Independent Component Analysis (ICA), Randomized Projections (RP) and Non-Negative Matrix Factorization (NMF). Unsupervised learning is used to infer some representation of unlabeled data. Clustering groups observations together, showing homogeneity in the data. Both clustering, when applied as a 'feature learning' step, and DR can used for reducing the amount of data. DR reduces the number of dimensions, and is particularly useful in addressing the 'curse of dimensionality', where higher dimensionalities require more training data.

This project will explore clustering in both contexts: inference and feature learning, and DR for feature reduction. It should be noted that unsupervised learning does not have labeled data, thus accuracy, as used in supervised learning is not a relevant metric. Since this project utilized problems that had labels, they were used only to predict and compare accuracy for the unsupervised learning models employed in this project. The clustering, DR and classification algorithms were implemented using scikit-learn[1] (sklearn). Graphing and dataset preprocessing was accomplished using python: pandas, numpy and matplotlib.pyplot.

## 2. Problems

For this assignment, two problems from the first assignment were used. The first problem was the built-in IRIS data set that came with sklearn. The IRIS data set is well known, thus, I will not deeply discuss this data set. Briefly, the data set contains 3 classes of 50 instances each, where each class refers to a type of iris plant. There are four attributes (features): sepal length (cm), sepal width (cm), petal length (cm) and petal width (cm). The objective (problem) is to classify (predict) the species of iris based on the four attributes.

The second problem[2] (SPAM) was also obtained from the **UCI Machine Learning Repository**. The data consists of a collection of attributes from both spam and non-spam (ham) email: 48 features related to word counts and 9 features associated with chars and strings of capital chars. The objective is to classify the e-mail as spam (1) or ham (0).

To minimize issues with feature scale, the data for all problems were normalized and scaled to values between zero and one. This data transformation was applied to all problems and used for all algorithms in this assignment.

## 3. Tune the Hyper-Parameters and Classify

For the classification, a Neural Network was used for all experiments in this assignment. The sklearn hyper-parameters: `activation`, `learning_rate` and `solver` were optimized by applying a ten-fold cross validation using sklearn's GridSearchCV. The number of hidden layers for the Neural Network were not included in the GridSearch, that is, the number of hidden layers was left

at the default setting (`sklearn.neural_network.MLPClassifier: hidden_layer_sizes: tuple, length = default (100,)`). To account for algorithmic processing efficiency, the run-time of both the dimensionality reduction and classification were captured. It should be noted that the computer running the experiments was not exclusively restricted to running experiments, e.g., there were other operations and multi-tasked operations occurring during the experimentation. This could have impacted the time, or at least leas to increases in the deviation of observed run-time values. Also the optimal hyper-parameters for classification change as per the number of dimensions. This also had an impact on the variance of observed classification run-time.

To avoid bias in generalizing to the test data when classification (Supervised Learning, Neural Network) was performed, the problem was partitioned to an 80% test, 20% training for both clustering and dimensionality reduction. The reduced data was then partitioned into 50% test/training for classification.

## 4. DISCUSS THE RESULTS

To provide a baseline for this project, the two problems were classified with a neural network, without clustering and without dimensionality reduction. Classification was performed using all features, and using only two features (chosen at random). This results of this ad-hoc approach of feature reduction will be compared against both clustering and dimensionality reduction. The results in table show a large decrease in test accuracy, while the learning curve for IRIS show a need for more training data using ad-hoc feature reduction.

**Table 1**: Summery of Accuracy for classification Neural Network.

| Table 1 Summary of Accuracy | Nrl Net | | Nrl Net | |
|---|---|---|---|---|
| | I R I S | S P A M | I R I S | S P A M |
| # features | All | all | (0,3) | (0,14) |
| *CLF Run Time (s) | 28 | 205 | 25 | 155 |
| Accuracy Train % | 97 | 92 | 96 | 67 |
| # Misclassified | 2 | 205 | 12 | 748 |
| # Correct | 73 | 2095 | 63 | 1552 |
| Accuracy Test % | 97 | 91 | 84 | 67 |

\* Time for classification includes 10-fold cross validation.

**Figure 1a**: Plots of data classification using only 2 features. The shaded regions denote decision boundaries. The classification of the test data is depicted in the plot by circles.
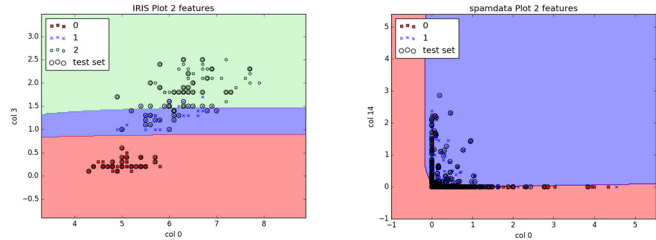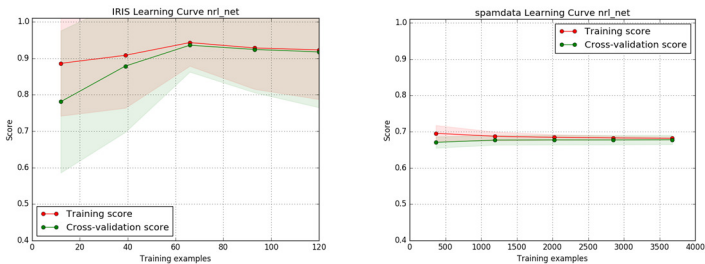


**Figure 1b**: Learning curves for classifying the problems with a Neural Network that uses only two features.

## 4.1. Clustering: k-means, Expectation Maximization (EM)

The first clustering algorithm for discussion is k-means. K-means identifies clusters using an iterative process of assigning and optimizing cluster centers within the data. The optimization step minimizing the sum-of-squares (Euclidian) distances between points within a cluster and the cluster center. When the number of clusters is fixed to k, k-means clustering is an optimization problem: find the k cluster centers and assign the points to the nearest cluster center, such that the squared distances from the cluster are minimized. The optimization problem is known to be NP-hard, and thus the common approach is to search only for approximate solutions. Drawbacks of k-means: k must be specified in advance, bias toward clusters of approximately similar size. Like hill climbing, k-means is susceptible to local minimum, thus, is dependent on where the initial centers for the clusters are placed.

Expectation Maximization (EM), which was implemented using a Gaussian Mixture Model (GMM) for this project, is also an iterative algorithm, with Expectation and Maximization steps. Like k-means, EM exposes clusters in the data. However, where k-means tends to find clusters of comparable spatial profile, EM allows clusters to have different shapes.
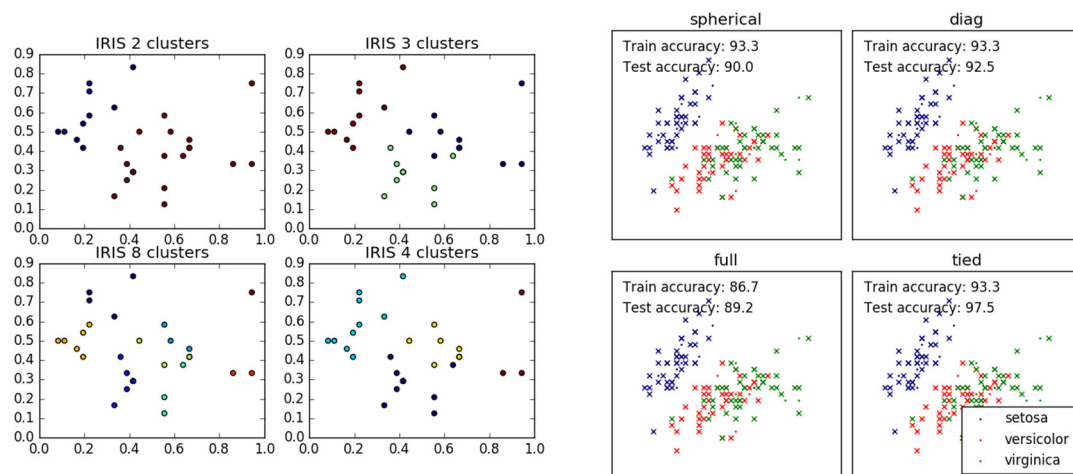
**Table 2**: Summary of results from clustering and expectation maximization. For this experiment, the objective was to observe the performance of k-means across a range: 2 ≤ k ≤ # features. The upper bound for k (# features) is not a good heuristic, it was selected to demonstrate the significance of choosing a 'good' value for k.

| Table 2 Clustering | Num Clstr 2 | Num Clstr 3 | Num Clstr 8 | Num Clstr 4** | Num Clstr 58** | EM | Notes |
|---|---|---|---|---|---|---|---|
| **IRIS:** *Run Time (s)* *Accuracy Train %* *Accuracy Test %* | 0.025 70 66 | 0.026 30 29 | 0.042 47 43 | 0.033 20 23 | | 0.011 93 98 | **k-means**: K=2 performs the best, although the performance is not linear as k increases. **EM cov='tied'**: Best performance. |
| **SPAM:** *Run Time (s)* *Accuracy Train %* *Accuracy Test %* | 0.126 62 59 | 0.122 50 49 | 0.184 9 11 | | 0.555 0 0.33 | 0.019 63 60 | **k-means**: K=2 performs the best, the performance decreases linearly as k increases. **EM cov='spherical'**: Best performance. |

\* Accuracy estimated using labels from problem.
\*\* There values were obtained from the number of features for the problem.

**Figure 2**: k-means plots (left) of IRIS for 2 ≤ k ≤ # features. EM plots for covariance: ['spherical', 'diag','full', 'tied']. Note the clustering appears more pronounced for EM. K-means, k=2 appears to generate the best clusters.

## 4.2. Dimensionality Reduction

There were some challenges for obtaining interesting results when applying dimensionality reduction to the problems. Values for quantifying the dimensionality reduction methods were captured: Eigenvalues, Kurtosis and Reconstruction Error. The results for both problems were similar: PCA Eigenvalues for the first two dimensions were dominant. With increasing dimensionality: ICA Kurtotic distributions plateau at a strongly leptokurtic (tall and narrow curve) value; NMF and RP reconstruction error approach zero.
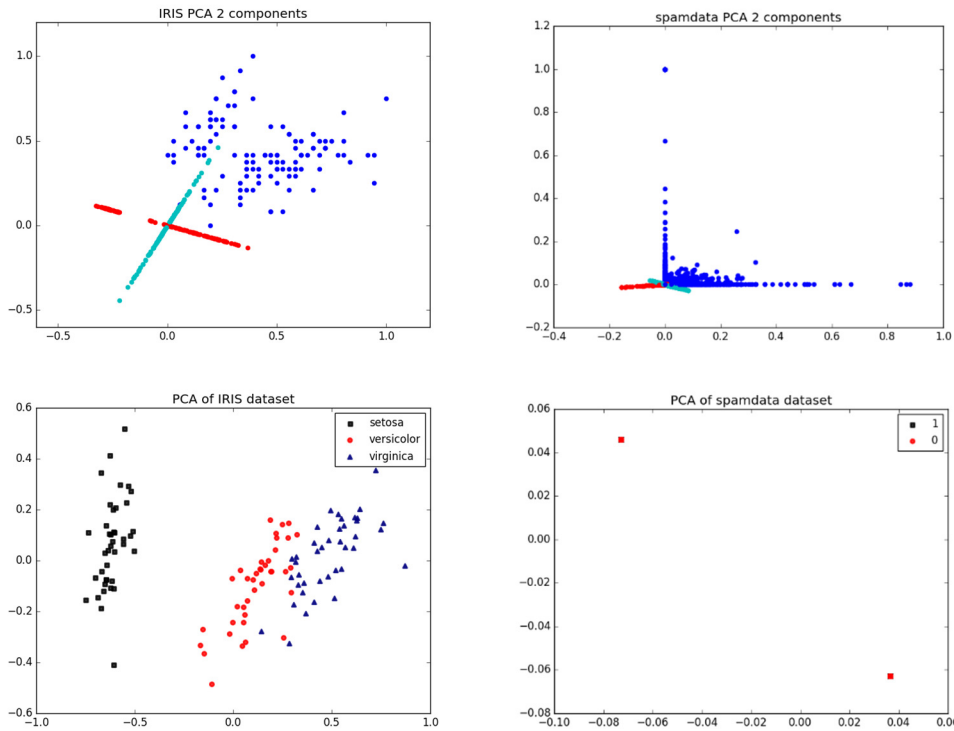


**Figure 3**: The first two principal component vectors are plotted against the original data points (top). Note the principal components are axis parallel to the SPAM data.

The data is transformed using the first two principal components and plotted (bottom). Note the well define clusters for both problems. Interestingly, the data collapsed for the PCA transformation for SPAM.

## 4.3. Clustering with Dimensionality Reduction

This experiment is a repeat of the experiment in section 4.1, but with dimensionality reduction applied prior to clustering. The data was reduced using the first two, then three components.

**Table 3a**: Summary of results from clustering and expectation maximization on dimensionally reduced problem using the first two components from PCA.

| Table 3a Clustering/PCA with (n=2) | Num Clstr 2 | Num Clstr 3 | Num Clstr 8 | Num Clstr 4** | Num Clstr 58** | EM | Notes: The number of components DR keeps=2. |
|---|---|---|---|---|---|---|---|
| **IRIS:** *Run Time (s)* *\*Accuracy Train %* *\*Accuracy Test %* | 0.026 0 0 | 0.026 4.2 3.1 | 0.046 0 0 | 0.03 38 49 | | 0.009 96 96 | k-means+DR: Performs worse than k-means. **EM+DR**: Performs on par with EM. |
| **SPAM:** *Run Time (s)* *\*Accuracy Train %* *\*Accuracy Test %* | 0.033 56 59 | 0.059 49 50 | 0.116 22 22 | | 0.429 2 3 | 0.014 54 57 | k-means+DR: Performs on par with k-means. **EM+DR**: Performs on par with EM. |

\*   Accuracy estimated using labels from problem.
\*\* There values were obtained from the number of features for the problem

**Table 3b**: Summary of results from clustering and expectation maximization on a dimensionally reduced problem using the first three components from PCA.

| Table 3b Clustering/PCA with (n=3) | Num Clstr 2 | Num Clstr 3 | Num Clstr 8 | Num Clstr 4** | Num Clstr 58** | EM | Notes: The number of components DR keeps=3. |
|---|---|---|---|---|---|---|---|
| | | | | | | | |
| **IRIS:** *Run Time (s)* *Accuracy Train %* *Accuracy Test %* | 0.03 75 64 | 0.03 0 7.3 | 0.04 8 7 | 0.03 17 17 | | 0.008 96 99 | **k-means+DR**: Performs better than k-means when #clusters=2. **EM+DR**: Performs better than EM. |
| **SPAM:** *Run Time (s)* *Accuracy Train %* *Accuracy Test %* | 0.06 56 59 | 0.07 80 78 | 0.115 21 20 | | 0.5 2 1 | 0.02 55 56 | **k-means+DR**: Performs better than k-means when #clusters=3. **EM+DR**: Performs worse than EM. |

\* Accuracy estimated using labels from problem.
\*\* There values were obtained from the number of features for the problem

**Figure 3b**: k-means plots (left) of IRIS for 2 ≤ k ≤ # features. EM plots for covariance: ['spherical', 'diag', 'full', 'tied']. PCA was used prior to clustering to reduce the data using the first three components. Note the clustering appears more pronounced for both k-means and EM when DR is applied first.
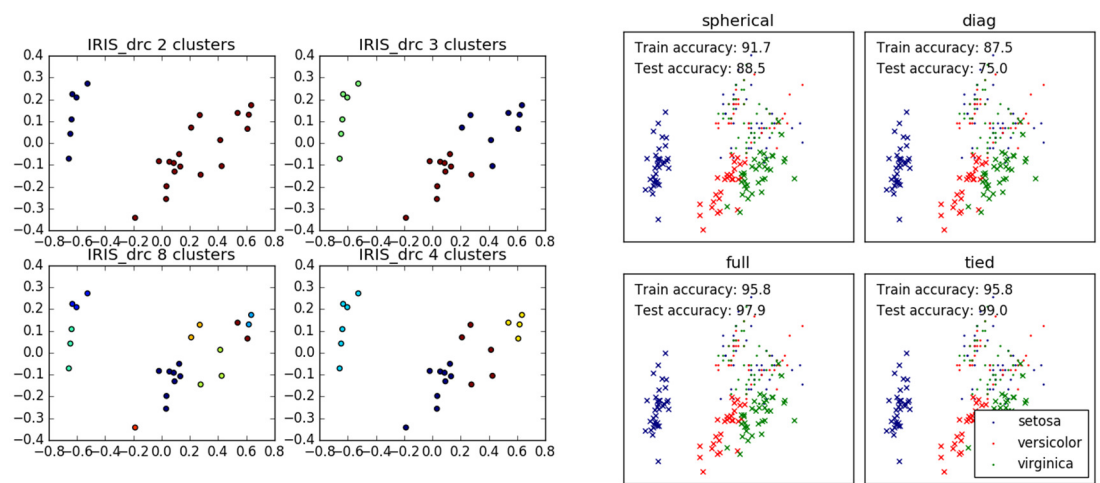


**Table 4a**: Summary of results from clustering and expectation maximization on a dimensionally reduced problem using the first two components from ICA.

| Table 4a Clustering/ICA with (n=2) | Num Clstr 2 | Num Clstr 3 | Num Clstr 8 | Num Clstr 4** | Num Clstr 58** | EM | Notes: The number of components DR keeps=2. |
|---|---|---|---|---|---|---|---|
| **IRIS:** *Run Time (s)* *Accuracy Train %* *Accuracy Test %* | 0.014 0 1 | 0.017 88 76 | 0.042 17 32 | 0.03 21 19 | | 0.009 96 96 | **k-means+DR**: Performs better than k-means with #clusters=3. **EM+DR**: Performs on par with EM. |
| **SPAM:** *Run Time (s)* *Accuracy Train %* *Accuracy Test %* | 0.059 81 79 | 0.06 78 78 | 0.12 34 34 | | 0.61 7 5 | 0.01 44 43 | **k-means+DR**: Performs better than k-means with #clusters=[2, 3]. **EM+DR**: Performs worse than EM. |

\* Accuracy estimated using labels from problem.
\*\* There values were obtained from the number of features for the problem

**Table 4b**: Summary of results from clustering and expectation maximization on a dimensionally reduced problem using the first three components from ICA.

| Table 4b Clustering/ICA with (n=3) | Num Clstr 2 | Num Clstr 3 | Num Clstr 8 | Num Clstr 4** | Num Clstr 58** | EM | Notes: The number of components DR keeps=3. |
|---|---|---|---|---|---|---|---|
| | | | | | | | |
| **IRIS:** *Run Time (s)* *Accuracy Train %* *Accuracy Test %* | 0.03 4 4 | 0.03 88 68 | 0.04 8 7 | 0.03 8 13 | | 0.008 96 99 | **k-means+DR:** Performs on par with k-means. **EM+DR:** Performs better than EM. |
| **SPAM:** *Run Time (s)* *Accuracy Train %* *Accuracy Test %* | 0.03 19 21 | 0.04 78 78 | 0.11 11 10 | | 0.314 1 1 | 0.014 55 58 | **k-means+DR:** Performs better than k-means when #clusters=3. **EM+DR:** Performs worse than EM. |

\* Accuracy estimated using labels from problem.
\*\* There values were obtained from the number of features for the problem

**Figure 4b**: k-means plots (left) of IRIS for 2 ≤ k ≤ # features. EM plots for covariance: ['spherical', 'diag','full', 'tied']. ICA was used prior to clustering to reduce the data using the first three components. Note the clustering appears about the same for both k-means and EM when DR is applied first.
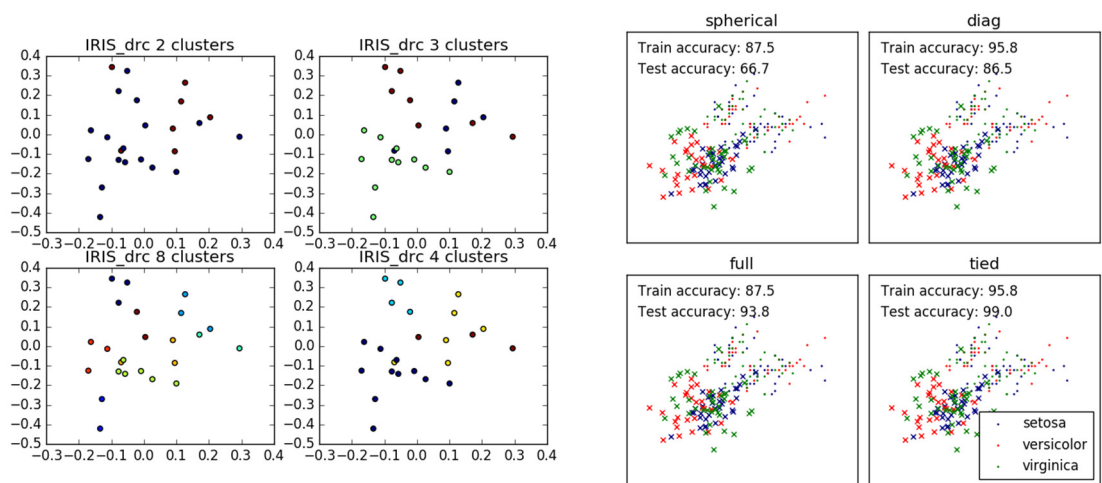


**Table 5a**: Summary of results from clustering and expectation maximization on a dimensionally reduced problem using the first two components from Randomized Projections (RP).

| Table 5a Clustering/RP with (n=2) | Num Clstr 2 | Num Clstr 3 | Num Clstr 8 | Num Clstr 4** | Num Clstr 58** | EM | Notes: The number of components DR keeps=2. |
|---|---|---|---|---|---|---|---|
| **IRIS:** *Run Time (s)* *Accuracy Train %* *Accuracy Test %* | 0.02 42 24 | 0.029 46 42 | 0.041 8 2 | 0.023 17 14 | | 0.009 96 96 | **k-means+DR:** Performs worse than k-means. **EM+DR:** Performs on par with EM. |
| **SPAM:** *Run Time (s)* *Accuracy Train %* *Accuracy Test %* | 0.012 57 58 | 0.072 44 45 | 0.19 28 27 | | 0.41 1 2 | 0.009 54 56 | **k-means+DR:** Performs on par with k-means with #clusters=2. **EM+DR:** Performs worse than EM. |

\* Accuracy estimated using labels from problem.
\*\* There values were obtained from the number of features for the problem

**Table 5b**: Summary of results from clustering and expectation maximization on a dimensionally reduced problem using the first three components from Randomized Projections (RP).

| Table 5b Clustering/RP with (n=3) | Num Clstr 2 | Num Clstr 3 | Num Clstr 8 | Num Clstr 4** | Num Clstr 58** | EM | Notes: The number of components DR keeps=3. |
|---|---|---|---|---|---|---|---|
| | | | | | | | |
| **IRIS:** *Run Time (s)* *\*Accuracy Train %* *\*Accuracy Test %* | *0.02* 75 63 | *0.03* 25 34 | *0.04* 0 0 | *0.03* 17 23 | | *0.006* 96 96 | **k-means+DR:** Performs on par with k-means when #clusters=2. **EM+DR:** Performs on par with EM. |
| **SPAM:** *Run Time (s)* *\*Accuracy Train %* *\*Accuracy Test %* | *0.06* 43 46 | *0.1* 57 55 | *0.2* 26 26 | | *0.5* 4 3 | *0.01* 53 55 | **k-means+DR:** Performs on par with k-means when #clusters=3. **EM+DR:** Performs worse than EM. |

\* Accuracy estimated using labels from problem.
\*\* There values were obtained from the number of features for the problem

**Figure 5b**: k-means plots (left) of IRIS for 2 ≤ k ≤ # features. EM plots for covariance: ['spherical', 'diag','full', 'tied']. RP was used prior to clustering to reduce the data using the first three components. Note the clustering appears about the same for EM, but more pronounced for k-means when DR is applied first.
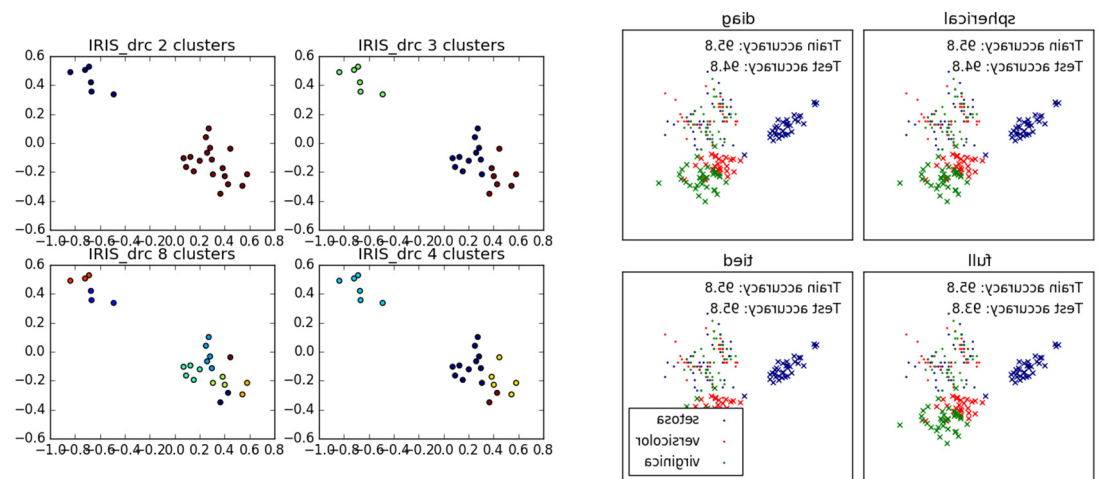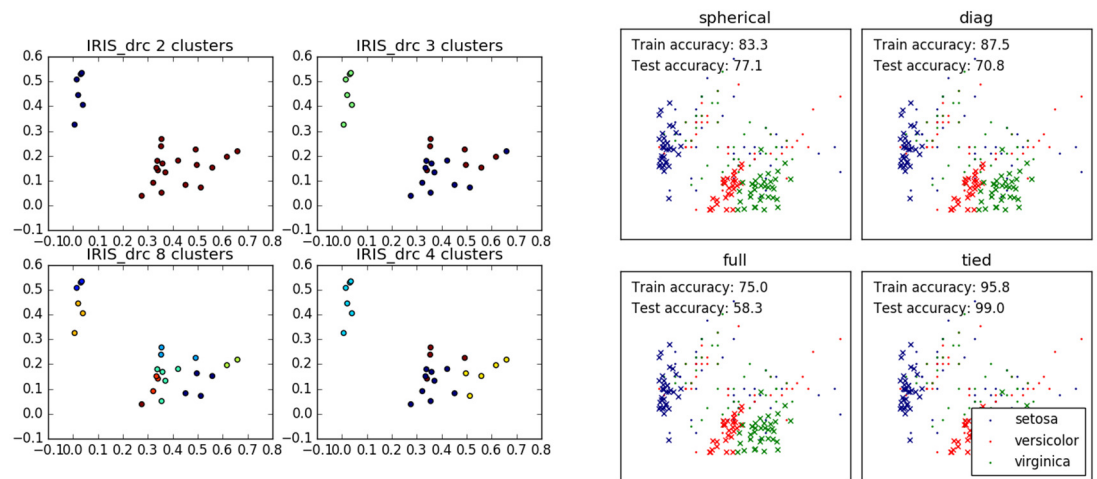


**Table 6a**: Summary of results from clustering and expectation maximization on a dimensionally reduced problem using the first two components from Non-Negative Matrix Factorization (NMF).

| Table 6a Clustering/NMF with (n=2) | Num Clstr 2 | Num Clstr 3 | Num Clstr 8 | Num Clstr 4** | Num Clstr 58** | EM | Notes: The number of components DR keeps=2. |
|---|---|---|---|---|---|---|---|
| **IRIS:** *Run Time (s)* *\*Accuracy Train %* *\*Accuracy Test %* | *0.02* 0 0 | *0.026* 25 29 | *0.04* 0 0 | *0.026* 0 7 | | *0.006* 96 96 | **k-means+DR:** Performs worse than k-means. **EM+DR:** Performs on par with EM. |
| **SPAM:** *Run Time (s)* *\*Accuracy Train %* *\*Accuracy Test %* | *0.04* 54 58 | *0.05* 76 76 | *0.13* 5 6 | | *0.46* 3 3 | *0.02* 53 57 | **k-means+DR:** Performs better than k-means with #clusters=3. **EM+DR:** Performs worse than EM. |

\* Accuracy estimated using labels from problem.
\*\* There values were obtained from the number of features for the problem

**Table 6b**: Summary of results from clustering and expectation maximization on a dimensionally reduced problem using the first three components from Non-Negative Matrix Factorization (NMF).

| Table 6b Clustering/NMF with (n=3) | Num Clstr 2 | Num Clstr 3 | Num Clstr 8 | Num Clstr 4** | Num Clstr 58** | EM | Notes: The number of components DR keeps=3. |
|---|---|---|---|---|---|---|---|
| | | | | | | | |
| **IRIS:** *Run Time (s)* *Accuracy Train %* *Accuracy Test %* | *0.02* 75 64 | *0.03* 17 18 | *0.04* 4 6 | *0.03* 21 22 | | *0.008* 96 99 | **k-means+DR**: Performs on par with k-means when #clusters=2. **EM+DR**: Performs better than EM. |
| **SPAM:** *Run Time (s)* *Accuracy Train %* *Accuracy Test %* | *0.054* 44 41 | *0.05* 54 57 | *0.14* 18 19 | | *0.4* 1 1 | *0.02* 51 53 | **k-means+DR**: Performs on par with k-means when #clusters=3. **EM+DR**: Performs worse than EM. |

\* Accuracy estimated using labels from problem.
\*\* There values were obtained from the number of features for the problem

**Figure 6b**: k-means plots (left) of IRIS for 2 ≤ k ≤ # features. EM plots for covariance: ['spherical', 'diag','full', 'tied']. NMF was used prior to clustering to reduce the data using the first three components. Note the clustering appears about the same for EM, but more pronounced for k-means when DR is applied first.
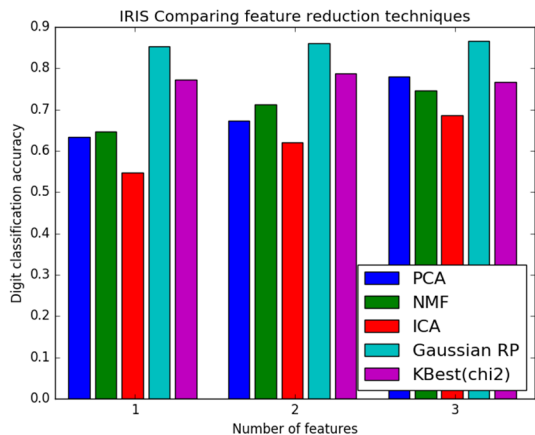


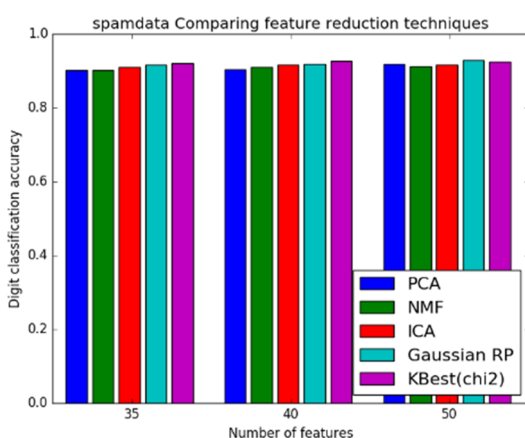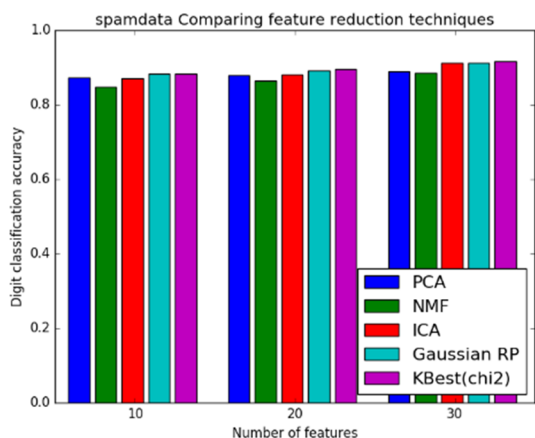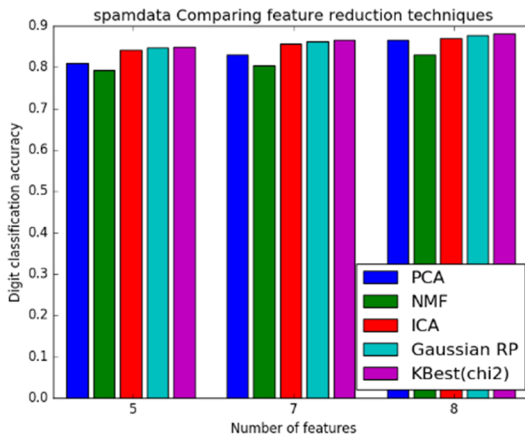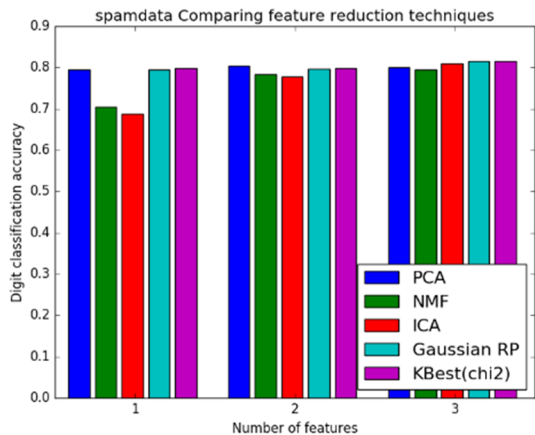## 4.4. Classification on a Dimensionality Reduced Problem

A pipeline in sklearn was utilized to evaluate, and apply the Dimensionality Reduction (DR) algorithms to a problem prior to classification. The pipeline combined DR: PCA, ICA, NMF, Gaussian RP and K-Best (chi2), then performed classification using a neural network with 10-fold cross validation. Recall the number of features (dimensions) in the problems: IRIS = 4, SPAM= 58. The range of dimensionality reduction for IRIS was 1 ≤ DR ≤ 3, and 1 ≤ DR ≤ 5 for SPAM. To baseline DR, a chi-square test was implemented to 'weed out' the features that are the most likely to be independent of class and therefore irrelevant for classification.

**Figure 1**: A comparison of feature (dimensionality) reduction methods. The graphical analysis presents the evaluation of DR algorithms across a range of feature reduction. Observe the accuracy of the classification as number of dimensions for the problem is increased. For IRIS, notice that the accuracy for classification improves for: PCA, NMF and ICA as more dimensions are passed to the neural network. Gaussian RP is the most accurate and consistent across all dimensions, and is superior to the baseline chi-square test. With three dimensions, PCA reaches the chi-square upper-bound as well.

Similar, but not as dramatic results are seen in the SPAM problem. NMF and ICA show strong improvement in classification accuracy when the dimensionality is increased from one to two. All DR methods exhibit improved accuracy as the dimensionality is increased from two to ten. After ten dimensions, a very slow plateau of ~90% accuracy is obtained.



Compare the results of classification using dimensionality reduction to keep two dimensions with our ad-hoc feature reduction using two features for classification (table 1, figure 1). With the ad-hoc approach, the accuracy was: IRIS=84%, SPAM=67%. Although IRIS fared worse for all DR methods except for RP, SPAM performed better across all DR approaches. A greater than 90% classification was achieved with DR on the SPAM problem at ~10 dimensions, matching the results of classification without DR. Applying DR to IRIS was not as successful, as the 97% accuracy without DR was not achieved.

### 4.5. Classification Using a Neural Network on Clustered Data

The section of the project will use a neural network to classify data that has been clustered using: k-means and Expectation Maximization (EM). The approach was to treat the clustered data as features to the neural network. For k-means, since we have labeled data, a simple heuristic for defining the cluster size was to assign the number of classes to the cluster size. The EM feature in sklearn did not have a transformation method, so I used the 'n_samples' attribute set to the size of the test set to generate a transformed data set for classification: 'X = estimator.sample(n_samples=len(X_test))[0]'.

**Table 7**: Summary of results from classification using clustering and expectation for data reduction. The reduced data set was passed to a neural network for classification.

| Table 7 CLF/Clustering | Num Clstr 2** | Num Clstr 3** | EM | Notes |
|---|---|---|---|---|
| | | | | |
| **IRIS:** *CLF Run Time (s)* Accuracy Train % # Misclassified # Correct Accuracy Test % | 57 92 6 54 90 | | 46 40 39 21 35 | **CLF+k-means**: Performs better than CLF+DR. **CLF+EM**: Performs worse than CLF+DR. |
| **SPAM:** *CLF Run Time (s)* Accuracy Train % # Misclassified # Correct Accuracy Test % | | 125 75 457 1383 75 | 114 55 793 1047 57 | **CLF+k-means**: Performs better than CLF+DR. **CLF+EM**: Performs worse than CLF+DR. |

\*   Time for classification includes 10-fold cross validation.

\*\* There values were obtained from the number of classes for the problem

## 5. CONCLUSION

When clustering is performed over a dimensionally reduced problem, choosing the cluster size, dimensionality reducing algorithm and scope of dimensionality reduction relative to the problem is critical. Some clustering algorithms, EM for instance, may be more predictable over a range of dimensionality reduction. Others, like k-means, are less forgiving, giving wildly disparate results for initial cluster sizes and across dimensionality reduction. Classification using dimensionality reduction to reduce the feature space was shown to be more accurate than an ad-hoc approach for the SPAM problem, and performed better for randomized projections for the IRIS problem. Surprisingly, classification using k-means to reduce the feature space performed even better the dimensionality reduction. The methods used to measure dimensionality reduction: Eigenvalues, Kurtosis and Reconstruction Error showed that the optimal value for the number of dimensions to keep was approximate to the number of classes in the problem. Of course, with unsupervised learning, we would not have class information, but this could be estimated from the measurements, or from clustering.

---

[1] http://scikit-learn.org/stable/

[2] https://archive.ics.uci.edu/ml/datasets/Spambase