# Project 3

Purushartha Singh

March 16, 2021

**Abstract**

The project discusses two main methods of feature selection to reduce a problem with a high number of dimensions to the most relevant dimensions.

# Contents

# 1    Introduction

Simply put, feature selection is the process of choosing a subset of the features from the original data according to certain criteria using Filter methods (choosing features using criteria independent of the original problem) or wrapper methods (Solving discriminant problem for each subset of features and selecting the best). It is a part of the feature extraction problem for datasets which have high number of features where classifying using all the features is highly inefficient and would lead to over fitting to the training data due to curse of dimensionality. In contrast, selecting a few of the more helpful features will greatly improve the performance and lead to better results. [1]

# 2    Objective

The aim of this project is feature selection for the classification problem using the high dimensional Taiji dataset provided. The two primary components of the procedure are:

1. search strategy : Identify the subset of most relevant features out of the provided features

2. evaluation function: Measure the discriminative power of the selected feature subset using filter and wrapper methods of feature selection.

3. Effect of M and N: Effect of the number of frames before and after the data which are labelled for that class on the overall classification

# 3    Approach

## 3.1    Data

The data used for the assignment was the new Taiji dataset. The dataset comprised of 1962 features and a total of 21919 data entries for each feature. It must be noted that the last element of the feature set is the frame time which was not relevant to the classification (outside of M and N value) and so, was omitted from the classification process giving a total of 1961 features.

The data points consisted of observations of 10 individuals and this was used to our advantage by using the leave one out methodology for creating 10 sets of training and testing data each containing 19918 training data points and 2001 testing data points. Further description of the features can be found in the project description.

## 3.2    Methodology

The following sections cover a discussion of the two primary methods used for feature selection:

### 3.2.1   Filter method

**Definition:**   Given a data set, select a subset m from n features where m¡n such that the value of criterion function is optimized over m.

Filter methods are "based on performance evaluation metric calculated directly from the data, without direct feedback from predictors that will finally used on the data with reduced number of features."[1]. The two primary methods of filtering discussed in class were variance ratio and augmented variance ratio. The code for the project utilizes augmented variance ratio (AVR) to find the top 2 per cent of the entire feature set which is then passed onto the wrapper method for further selection. The evaluation uses the following formula as the metric such that the features with the highest value of the AVR metric are shortlisted. The AVR for a feature is given by:

$$AVR(F) = \frac{Var(S_F)}{\frac{1}{C}\sum_{k=1,..,C}\frac{Var_k(S_F)}{min_{m\neq k}|\mu_k(S_F)-\mu_m(S_F)|}} \tag{1}$$

such that $S_F$ is the feature being evaluated, $Var(S_F)$ is the variance across the feature, $Var_k(S_F)$ is the variance across the feature for class k, $\mu_k(S_F)$ is the mean across the feature for class k, and $\mu_m(S_F)$ is the mean across the feature for class m such that m is not the same class as k. $C$ is the total number of classes.

Notice that the criteria used for evaluation is independent of the specific problem. The examples are information-theory based measure such as the variance ratio, etc. and the the information-theoretic measures find the correlation, mutual information or information gain of the different feature subsets similar to how Fisher method does.

Once the AVR has been found for each feature, the features are ranked and the highest scoring features are selected.

### Advantages

1. **Fast execution:** Only a subset of the features are used for the rest of the problem, thus, speeding up the computation.

2. **Generality:** Since the evaluation criteria is not data-specific, this method can be generalized over data sets.

### Disadvantages

1. It does not reveal any information about the selected feature aside from ability to split the dataset along the classes for the training dataset. As such, it may not be useful for making informed decisions and the non specific way can lead to loss of information.

### 3.2.2   Wrapper method

**Definition:**   In this method, the features are selected on the basis of quantitative rates. More specifically, the wrapper method implies solving of discriminant problem for each subset of features and selecting the best classification rates. For each subset of features selected by the wrapper methods should separate the class distribution as much as the

original class distribution.

The search strategy is the systematic procedure to choose there candidate feature subsets. These include Complete, greedy and randomized algorithms. The approach taken in the project is the complete approach using Sequential Forward Selection. Since the input is limited to the top AVR selections from the filter method, the NP complete nature of the problem is circumvented by keeping the input size small. The following algorithm is used to recursively find the best fitting accuracy on the training method (using Linear discriminant classification method):

1. Start with an empty set $P$ and a set of features $Q$

2. Select the next best feature $x^+ = argmax[Q(Y_k) + x]$ where Yth feature of Q is the evaluation criterion.

3. Update $P_{k+1} = P_k + x^+$, $Q_{k+1} = Q_k - x^+$ and $k = k + 1$

4. Jump to step 2 until the maximum value of $x^+$ has been achieved and no further addition of features increases the classification accuracy on the training set $\forall x^+ \in Q_{k+1}$.

5. Return $P_{k+1}$

**Advantages**

1. **Accuracy:** The method finds the best set of features from the given set to maximize the training classification accuracy.

2. **Avoids over fitting:** Uses cross-validation for predictive accuracy

**Disadvantages**

1. **Slow execution:** Given the NP complete run time, the problem can be very resource intensive if the input is not managed to very small set of features.

2. **Lack of generality:** The approach is highly dependent on the classifier and feature subset. In fact, during the project, using different classifiers lead to significant increase and decrease in the accuracy of the forwarded features.

## 3.3   Workflow

As requested, the workflow of the project has been provided below in fig 1:

# 4   Results

## 4.1   Feature Selection

This analysis was done with default values of M = 100 and N = 20. The Filter method selected top 1% of the features based on AVR. From the resulting classification, an average accuracy of 0.6958 was achieved over 10 iterations using all but one method to make
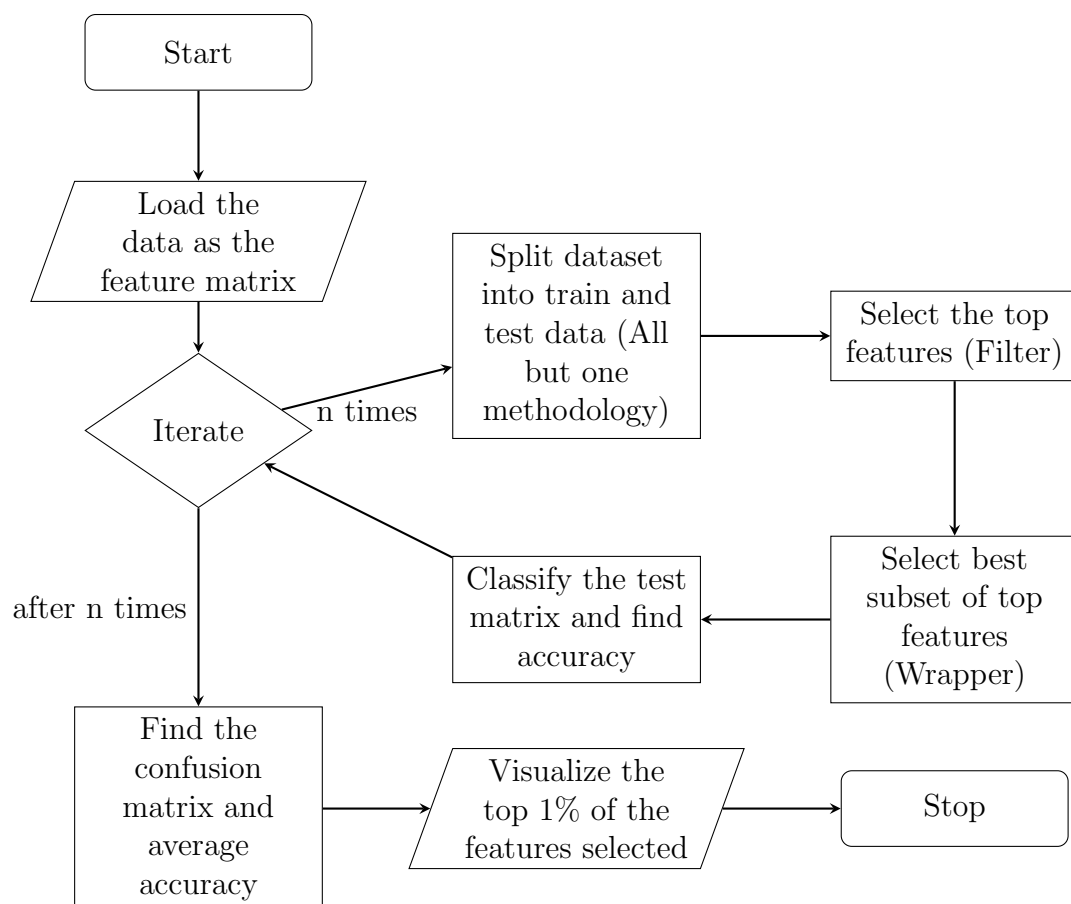
Figure 1: The basic workflow of the code. Details of each step can be found in the methodology section
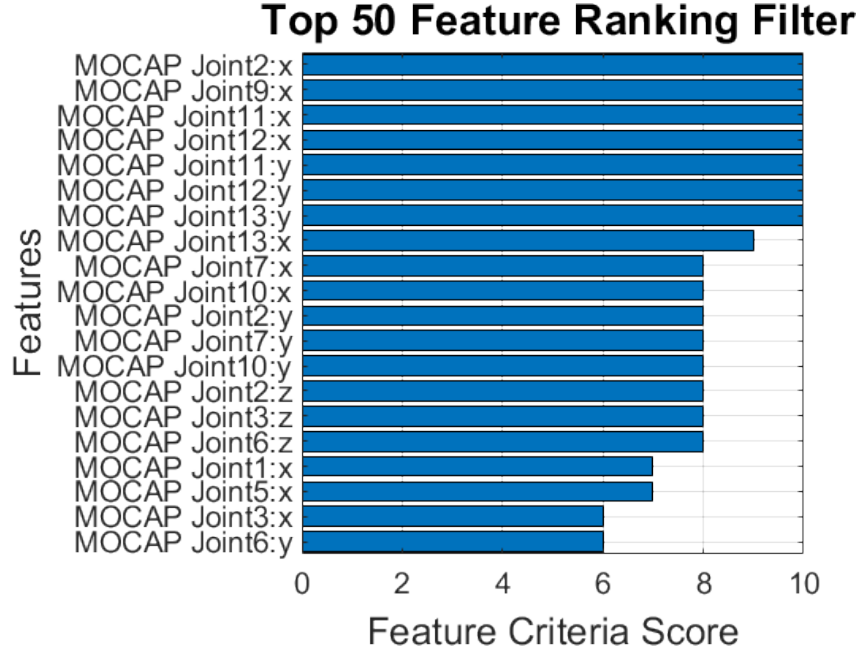
Figure 2: Frequency of appearance of features over 10 iterations in the filter method of feature selection

the training and testing datasets. The 20 most used features in this task for the filter and wrapper methods are shown in figure 2 and 3 respectively. Note that for all the classifications, Linear discriminant classification was used.

## 4.2  Classification Analysis

### 4.2.1  Base run

For the given parameters of M=100, N=20, K=0.01 (Filter % value of 1%), and Linear discriminant classifier, the average accuracy achieved for the run was 0.6983 with a standard deviation of 0.2740. The classification matrix and confusion matrix have been shown in fig 4 and 5 respectively. It must be noted that due to lack of space, the labels in the classification matrix have been omitted. For this run, the classification rate for each class for each iteration and the average classification rate for each class over all the iterations are shown in 6.

### 4.2.2  Final Run

With the values of more optimal parameters known, a final run with M=60, N=400, K=0.05, and Ensemble classifier got average classification accuracy of **0.8342**. with average standard deviation of 0.1538. The classification matrix and confusion matrix have been shown in fig 7 and 8 respectively. For this run, the classification rate for each class for each iteration and the average classification rate for each class over all the iterations are shown in 9. Further optimization may be done to improve the classification accuracy further.

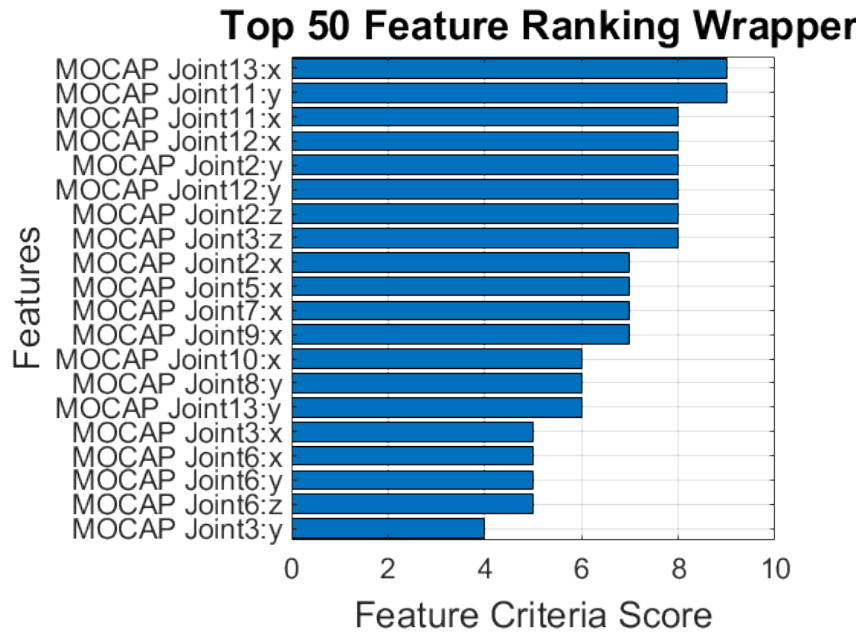## Top 50 Feature Ranking Wrapper

Figure 3: Frequency of appearance of features over 10 iterations in the wrapper method of feature selection

## Average Classification Matrix
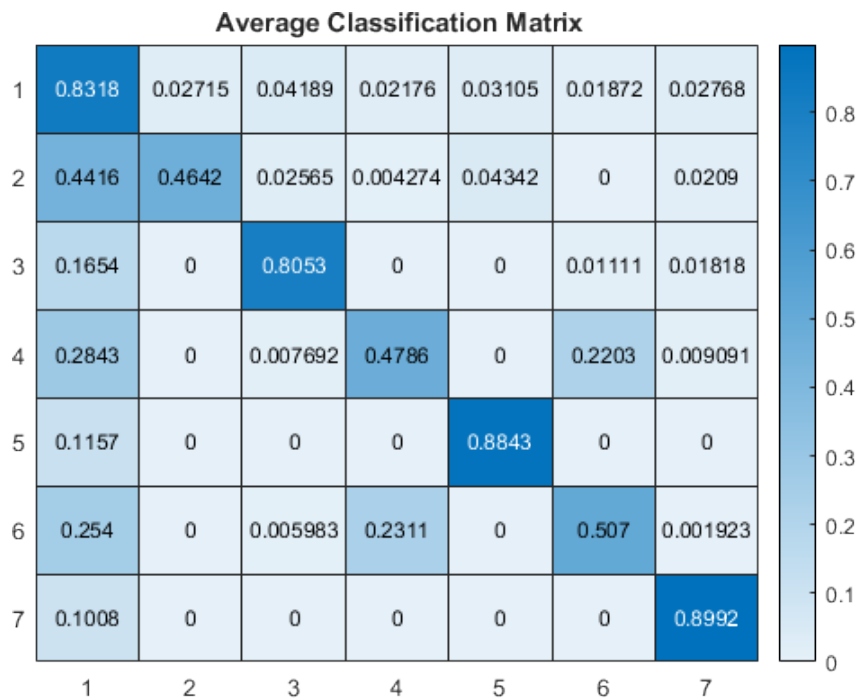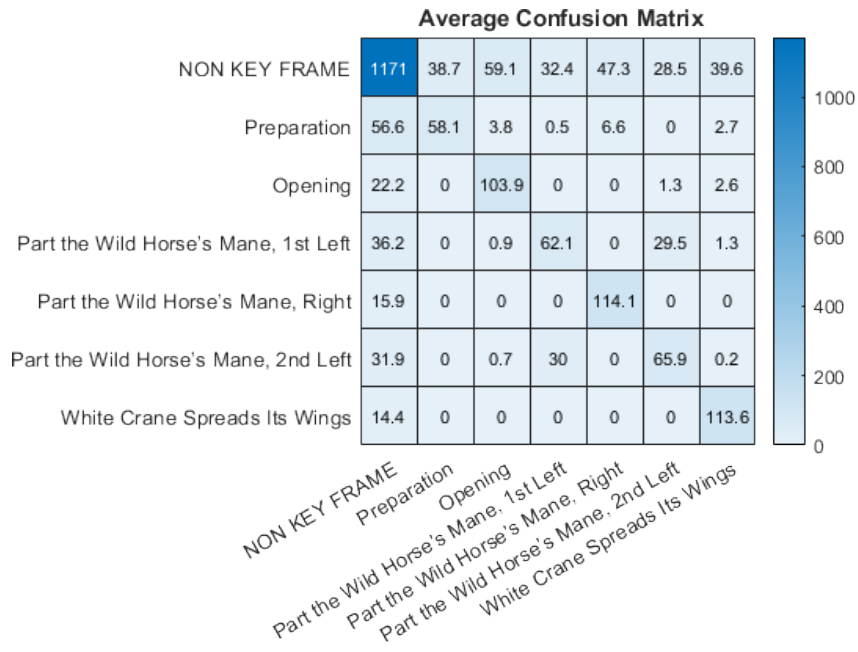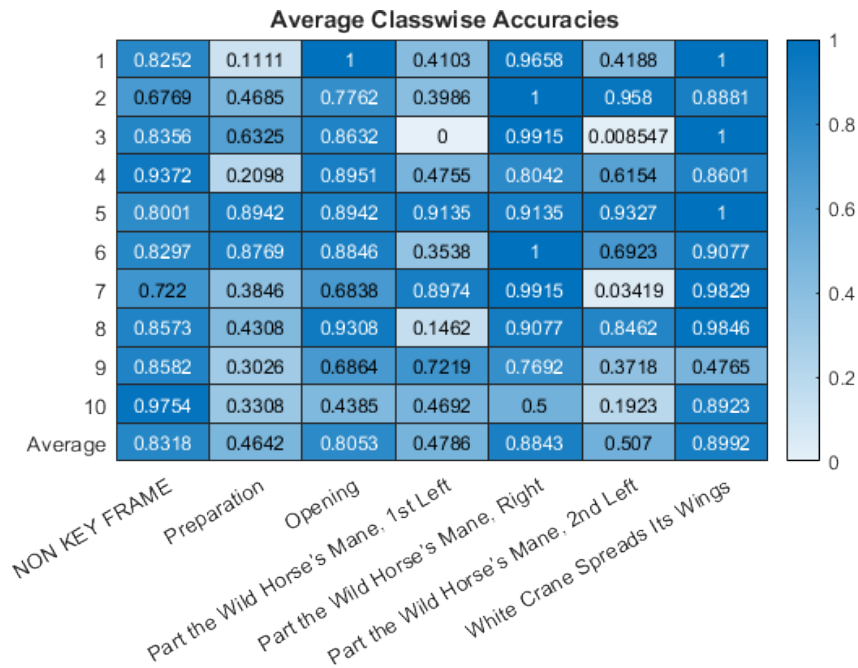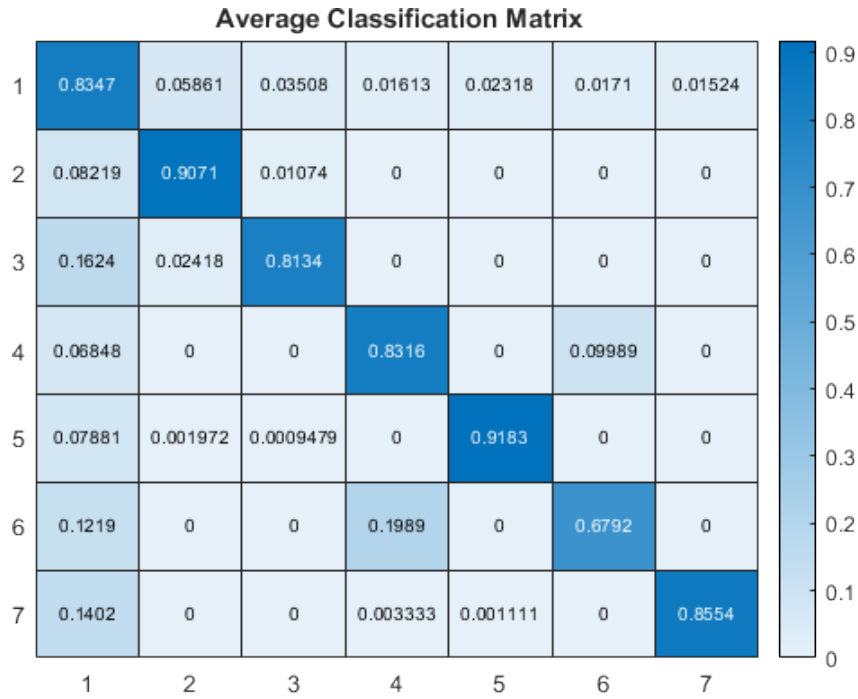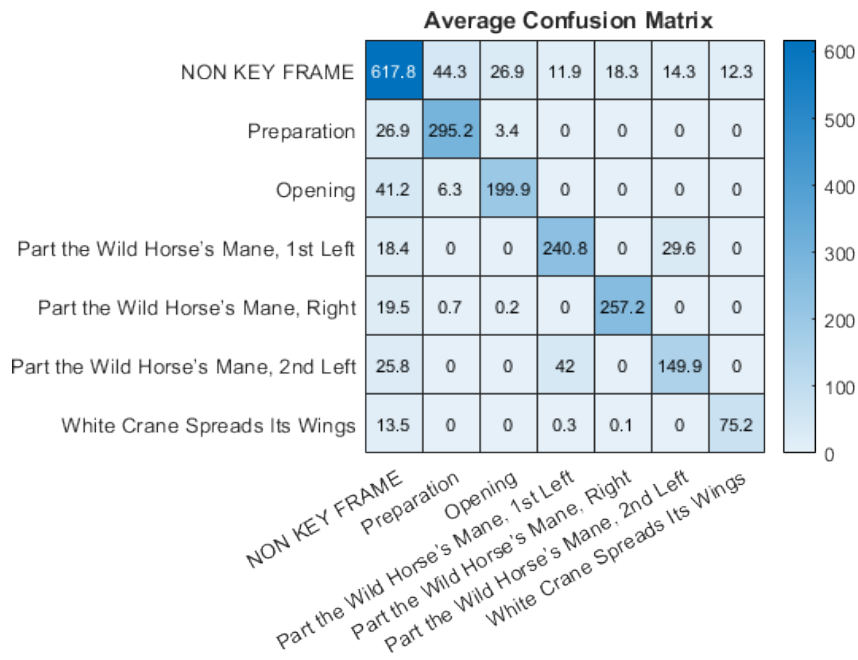
Figure 4: M=100, N=20, K=0.01, Classification Matrix

Figure 5: M=100, N=20, K=0.01, Confusion Matrix



Figure 6: M=100, N=20, K=0.01, Class wise accuracy and average for all iterations

Figure 7: M=60, N=400, K=0.05, Ensemble Classification Matrix


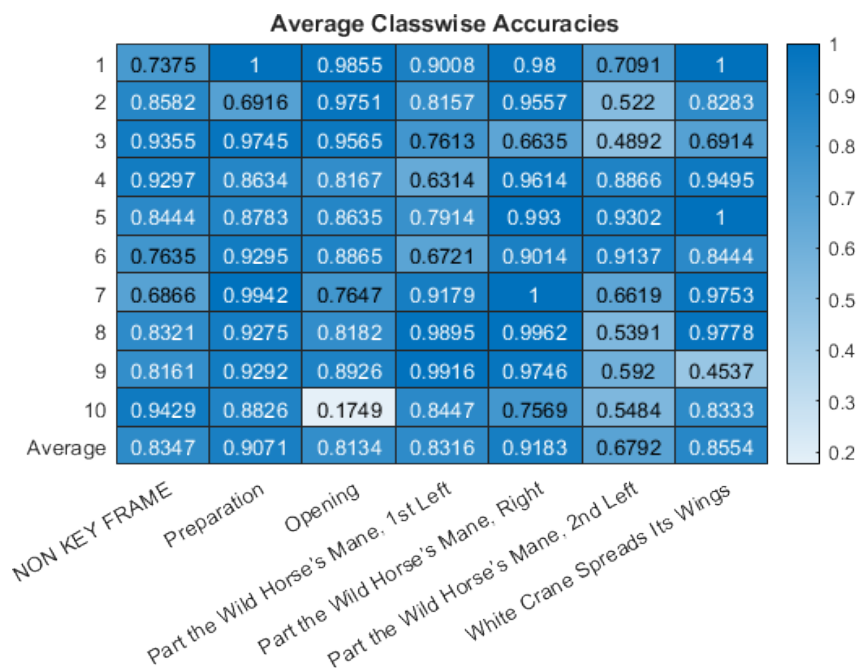
Figure 8: M=60, N=400, K=0.05, Ensemble Confusion Matrix

Figure 9: M=60, N=400, K=0.05, Ensemble Class wise accuracy and average for all iterations

## 4.3   EC: Use of different classifiers and ensemble learning

Use of different classifiers lead to pretty much the same results with the base parameters of M=100, N=20, K=0.01. Here is a list of the accuracy results for KNN, and decision tree classifier:

- K- Nearest Neighbors: 0.7068

- Decision Tree: 0.6412

- Naive Bayes: 0.7026

- Linear Discriminant (Base): 0.6983

Ensemble learning is used with a few parameters to improve the performance compared to these three methods to **0.7295**. This can be attributed to the following parameters:

- First, the bag method is used which is a bootstrap method. This allows for artificial increase of sample data without need for more actual observations.

- The expected value of prior is changed to the empirical values in the training data. This allows for a weight to be assigned to each expected value based on the probability of an item appearing in the provided data, thereby reducing the weight of data of classes which appear more often and boosting the data which comes from classes with fewer datapoints.

- Number of Bins is increased to 50 to improve the speed of the classifier since it allows for multiple values to be calculated in parallel making the process more time efficient.
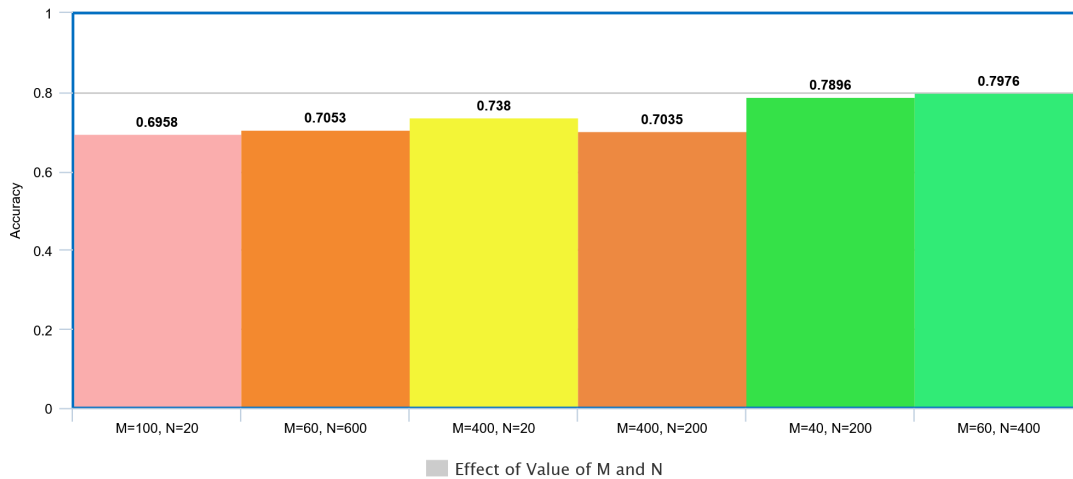
Figure 10: Effect of M and N values on the classification accuracy. K=0.01

## 4.4   Effect of M and N

M stands for the number of frames before the pose and N stands for number of frames after the pose. These give the data appearing at frame values in M-N range the assigned class for the pose during the assignment of classes to the data. The values of M and N seem to have a significant impact on the average accuracy of the classification as can be seen in fig 10. For uniformity, the value of K is kept at 0.01 and Linear discriminant classifier is used. As can be seen in the figure, value of accuracy increases by 10% by tuning the values of M and N.

The reason for this improvement can be due to the fact that increasing the number of labelled frames would decrease the large number of unclassified frames which often end of causing the confusion matrix to skew towards it. Having a more balanced dataset gives better classification accuracy.

# 5   Conclusion

Feature selection is important to reduce the computational efforts and space for working with large data sets. It helps identify features that are most useful and thus, improve the overall accuracy of the algorithm.

Given the novelty of the dataset, it was a very interesting experience having to learn on the go about different parameters and how multiple aspects of reducing dimensionality must be looked at. In addition, this was my first introduction to data specific parameters such as M and N and I was very surprised to look at the drastic effect of the two on the overall accuracy of the data.

This was also a very interesting look at how dimensionality reduction ends up causing a lot of data loss at times but at the same time, keeping the extra dimensions often results in poorer results. Overall, I was a bit sad that the time period for this project

was short and as a result, I was not able to be more thorough with testing and parameter tuning. It would have also been very interesting to have attempted some of the other extra credit activities such as data visualization and running this on the bigger dataset which i was unable to do due to lack of time.

# References

[1] Christopher M. Bishop *Pattern Recognition and Machine Learning.* 2006 2, 3