# Machine Learning Algorithms comparison and analysis

902905501
Yuanheng Wang

## Problem Description:

Two datasets are selected for the purpose of this assignment. One is "Bank Marketing", and the other is "Chess (King-Rook vs. King-Pawn)". Both datasets are acquired from UCI Machine Learning Dataset Repository. Below are the explanations of why to choose these two datasets.

Bank Marketing:

The data is related to a marketing campaign of Portuguese banking institution. The data contains 20 attributes( and the final output result is whether people with certain attributes would choose to subscribe the service provided. The full dataset contains 45211 instances and is a typical classification problem related to business field. I found the problem extremely interesting because the concept behind this dataset can be expanded all sorts of areas in the real world, especially for business undergoing big decisions on whether or not to launch marketing campaigns and relative strategies. For instance, imagine Microsoft is evaluating and researching the market on deciding whether to launch a new bundle of Xbox with a third party bundle game. If information of customers can be retrieved from previous launch events, including their ages, occupations, estimated annual salaries and geographic information, Microsoft would be able to identify the most efficient marketing strategy and tailor the channels of marketing in different areas. The subject can be extended to relevant social studies as well. Government usually have the most thorough information of the residents. If local government wants to push a service or a new proposal, this similar methodology can be applied to give an approximate possibility of passing certain bills.

King-Rook vs. King-Pawn:

This is a typical game problem that I believe will introduce me to the concept of game AI. The situation is as follows: White has King and one Rook left on the chess board, Black has King and one Pawn left on the position of A7, which is one step away from advancement. The data contains 36 attributes, and the $37^{th}$ is the result whether white can win or not. The part I find interesting about this dataset is its attributes. Each attribute is a state description and all those 36 attributes combined describe the entire chessboard. When the system becomes more complicated, for example, Starcraft, there will be more attributes describing the situation of the game. But in principle, I think the concept should be fairly similar. To be more general, not only games, but anything in real world that can be depicted by multiple attributes combined should allow us to apply same thought process.
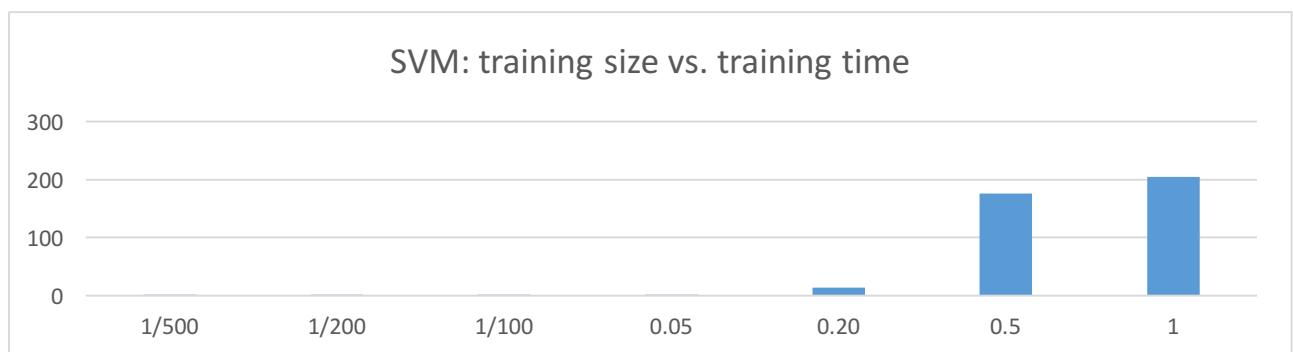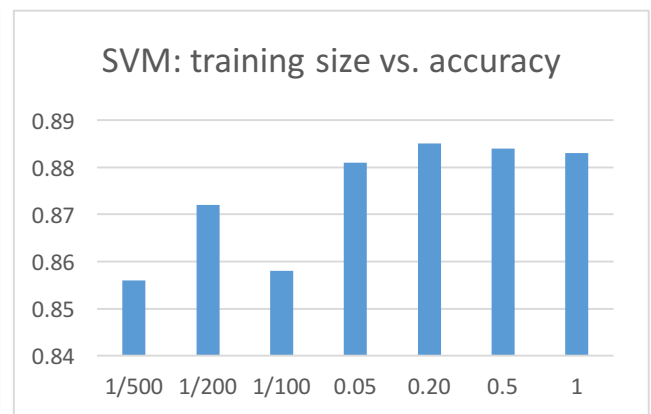
## Algorithm Comparison:

All the results below are attained by tuning the parameters, finding the best parameter settings for one algorithm. By splitting the dataset into training and real-world untouched data, I used cross validating method for the training dataset to get the best model. Then test the model with untouched real world model to get the most unbiased prediction accuracy.

Bank Marketing – SVM:
The SVC classifier in sklearn is implemented.
The tested parameters are
Parameter = {'C':[1,10], 'kernel': ('linear', 'rbf')}. One issue to be mentioned is for this dataset, the 'linear' kernel is extremely slow, and a result can only be obtained by shrinking the training size to 1/400 of the original size and skip the cross validation step. The accuracy is relatively lower than 'rbf' kernel with the same training size, but the training time required is much higher.

Bank Marketing – Decision Tree
The DecisionTreeClassifier in sklearn is implemented in this step. Below is the table of data.
A more direct graph representation.



Decision Tree is doing decently well in this problem. The time it took to train was much faster than SVM, and the accuracy is nearly the same. The model with default setting actually took the longest time to do cross validation and came up with the worst accuracy. The cross validation time and training time were largely deducted when max_depth parameter was specified and with less training and prediction time, setting max_depth equal to 5 improves the accuracy. This indicates that by default decision tree tends to be overfitting with given dataset for this problem.

Bank Marketing – Boosting
The AdaBoostClassifier in sklearn is implemented in this step. Below is the detail data.

## BOOSTING MODEL COMPARISON

Legend: ■default ■n_estimators=20 ■n_estimators=5 ■learning_rate=0.5 ■n_estimators=100,learning_rate=0.5

| | AVERAGE ACCURACY SCORE | CROSS VALIDATION TRAINING TIME | TRAINING TIME | PREDICTION TIME | WORLD ACCURACY | TYPE |
|---|---|---|---|---|---|---|
| default | 0.903 | 7.713 | 1.901 | 0.08 | 0.899 | 0 |
| n_estimators=20 | 0.898 | 3.168 | 0.804 | 0.031 | 0.895 | 0 |
| n_estimators=5 | 0.896 | 0.894 | 0.215 | 0.009 | 0.892 | 0 |
| learning_rate=0.5 | 0.902 | 7.957 | 2.025 | 0.085 | 0.897 | 0 |
| n_estimators=100,learning_rate=0.5 | 0.902 | 15.601 | 3.84 | 0.167 | 0.898 | 0 |

## Boosting: training size vs. training time&accuracy

| training size | training time | Accuracy |
|---|---|---|
| 1 | 1.901 | 0.903 |
| | 1.333 | 0.897 |
| 0.20 | 0.526 | 0.895 |
| | 0.195 | 0.894 |
| 1/100 | 0.105 | 0.871 |
| | 0.103 | 0.867 |
| 1/500 | 0.087 | 0.83 |

■training time ■Accuracy

## Bank Marketing – K-nearest Neighbors

The KNeighborsClassifier in sklearn is implemented in this step. Below is the detail data.

### KNN MODEL COMPARISON

■ Series1   ■ Series2   ■ Series3   ■ Series4

| | AVERAGE ACCURACY SCORE | CROSS VALIDATION TRAINING TIME | TRAINING TIME | PREDICTION TIME | WORLD ACCURACY |
|---|---|---|---|---|---|
| Series1 | 0.883 | 2.18 | 0.178 | 0.643 | 0.876 |
| Series2 | 0.876 | 1.972 | 0.169 | 0.574 | 0.872 |
| Series3 | 0.885 | 2.002 | 0.173 | 0.629 | 0.881 |
| Series4 | 0.887 | 2.276 | 0.18 | 0.728 | 0.883 |

### KNN: training size vs. training time&accuracy

| | training time | Accuracy |
|---|---|---|
| 1 | 0.178 | 0.876 |
| | 0.087 | 0.876 |
| 0.20 | 0.023 | 0.877 |
| | 0.003 | 0.879 |
| 1/100 | 0.001 | 0.876 |
| | 0.001 | 0.877 |
| 1/500 | 0 | 0.872 |

■ training time   ■ Accuracy

## Bank Marketing – Neural networks

The MLPClassifier in sklearn is intended to be implemented in this step. Unfortunately, sklearn 0.17 doesn't support multilayer perceptron. I intended to upgrade sklearn to 0.18 dev version, but Canopy has compatibility issue with the upgrade using git. Due to the time limit and technical capability restraint, this dataset is unable to transfer into arff format to be analyzed in GUI. Future research will be conducted in order to apply neural network algorithm on bank marketing data.

## Chess – SVM



**SVM: MODEL COMPARSION**

Legend: C=1,kernel='rbf' | C=10,kernel='rbf' | C=100,kernel='rbf' | C=1000,kernel='rbf' | C=1000,kernel='linear'

| | AVERAGE ACCURACY SCORE | CROSS VALIDATION TRAINING TIME | TRAINING TIME | PREDICTION TIME | WORLD ACCURACY |
|---|---|---|---|---|---|
| C=1,kernel='rbf' | 0.931 | 0.975 | 0.23 | 0.075 | 0.944 |
| C=10,kernel='rbf' | 0.966 | 0.567 | 0.131 | 0.036 | 0.976 |
| C=100,kernel='rbf' | 0.992 | 0.464 | 0.123 | 0.022 | 0.991 |
| C=1000,kernel='rbf' | 0.993 | 0.453 | 0.112 | 0.017 | 0.994 |
| C=1000,kernel='linear' | 0.971 | 8.913 | 4.215 | 0.011 | 0.976 |

**Training size vs. training time&accuracy**

training time | Accuracy

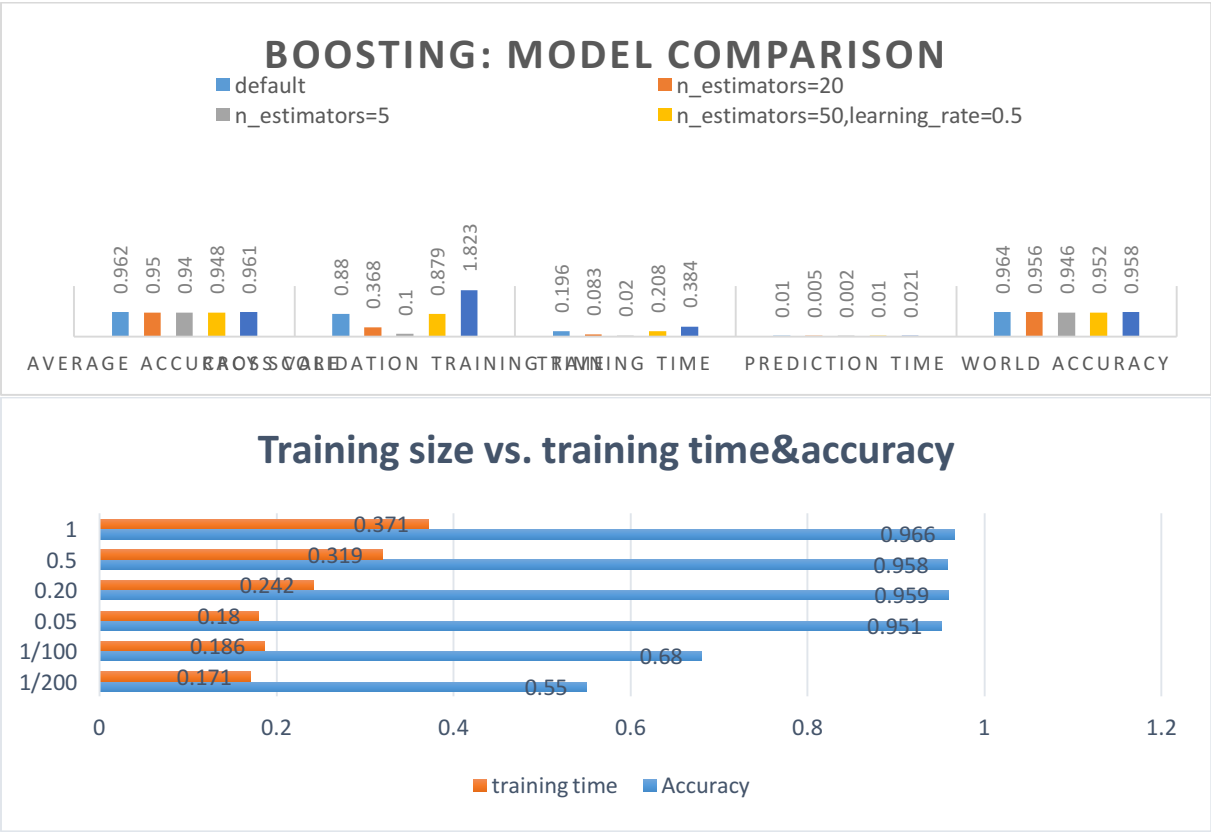| training size | training time | Accuracy |
|---|---|---|
| 0.5 | 0.112 | 0.994 |
| | 0.081 | 0.987 |
| | 0.016 | 0.978 |
| 0.05 | 0.003 | 0.91 |
| | 0.001 | 0.701 |
| 1/200 | 0 | 0.625 |

## Chess - Decision Tree



**DECISION TREE: MODEL COMPARSION**

Legend: default | min_samples_split=10 | min_samples_split=20 | min_samples_split=30,max_depth=5 | min_samples_split=50,max_depth=5

| | AVERAGE ACCURACY SCORE | CROSS VALIDATION TRAINING TIME | TRAINING TIME | PREDICTION TIME | WORLD ACCURACY |
|---|---|---|---|---|---|
| default | 0.99 | 0.038 | 0.007 | 0 | 0.996 |
| min_samples_split=10 | 0.99 | 0.034 | 0.006 | 0 | 0.996 |
| min_samples_split=20 | 0.985 | 0.03 | 0.007 | 0 | 0.992 |
| min_samples_split=30,max_depth=5 | 0.939 | 0.028 | 0.006 | 0 | 0.945 |
| min_samples_split=50,max_depth=5 | 0.939 | 0.029 | 0.005 | 0 | 0.945 |

**Training size vs. training time&accuracy**

Accuracy / training time values by training size:
- 0.5: 0.996, 0.007
- 0.994, 0.004
- 0.98, 0.002
- 0.05: 0.949, 0
- 0.631, 0
- 1/200: 0.555, 0

Legend: training time, Accuracy

Chess – Boosting



**BOOSTING: MODEL COMPARISON**

Legend: default, n_estimators=20, n_estimators=5, n_estimators=50,learning_rate=0.5

| | default | n_estimators=20 | n_estimators=5 | n_estimators=50, learning_rate=0.5 | |
|---|---|---|---|---|---|
| AVERAGE ACCURACY | 0.962 | 0.95 | 0.94 | 0.948 | 0.961 |
| CROSS VALIDATION TRAINING TIME | 0.88 | 0.368 | 0.1 | 0.879 | 1.823 |
| TRAINING TIME | 0.196 | 0.083 | 0.02 | 0.208 | 0.384 |
| PREDICTION TIME | 0.01 | 0.005 | 0.002 | 0.01 | 0.021 |
| WORLD ACCURACY | 0.964 | 0.956 | 0.946 | 0.952 | 0.958 |



**Training size vs. training time&accuracy**

- 1: 0.966, 0.371
- 0.5: 0.958, 0.319
- 0.20: 0.959, 0.242
- 0.05: 0.951, 0.18
- 1/100: 0.68, 0.186
- 1/200: 0.55, 0.171

Legend: training time, Accuracy

Chess – KNN

## KNN: MODEL COMPARISON

Legend: ■ default  ■ n_neighbors=3  ■ n_neighbors=4  ■ n_neighbors=10

| Metric | default | n_neighbors=3 | n_neighbors=4 | n_neighbors=10 |
|---|---|---|---|---|
| AVERAGE ACCURACY SCORE | 0.936 | 0.931 | 0.928 | 0.941 |
| CROSS VALIDATION TRAINING TIME | 0.473 | 0.459 | 0.445 | 0.465 |
| TRAINING TIME | 0.017 | 0.008 | 0.008 | 0.011 |
| PREDICTION TIME | 0.254 | 0.257 | 0.245 | 0.23 |
| WORLD ACCURACY | 0.95 | 0.944 | 0.941 | 0.951 |

## Training size vs. training time&accuracy

| training size | training time | Accuracy |
|---|---|---|
| 0.5 | 0.017 | 0.95 |
| | 0.004 | 0.916 |
| 0.05 | 0.002 | 0.864 |
| | 0.001 | 0.817 |
| 1/200 | 0.001 | 0.648 |
| | 0 | 0.533 |

Legend: ■ training time  ■ Accuracy

Chess – Neural Network

Due to the limit of sklearn version 0.17, MPLClassifier can't be appied on the data. Therefore, weka GUI is used for testing neural network algorithm. Below are detailed data.

## MULTILAYER NETWORK: MODEL COMPARISON

Legend: ■ learning_rate=0.1,hidden_layer='a'  ■ learning_rate=0.3,hidden_layer='a'  ■ learning_rate=0.01,hidden_layer='a'  ■ learning_rate=0.1,hidden_layer='i'  ■ learning_rate=0.1,hidden_layer='o'

| Metric | lr=0.1,'a' | lr=0.3,'a' | lr=0.01,'a' | lr=0.1,'i' | lr=0.1,'o' |
|---|---|---|---|---|---|
| AVERAGE ACCURACY SCORE | 0.98 | 0.977 | 0.97 | 0.977 | 0.977 |
| CROSS VALIDATION TRAINING TIME | 6.5 | 5.98 | 6.22 | 11.61 | 0.92 |

## Training size vs. training time&accuracy



| | training time | Accuracy |
|---|---|---|
| 1 | 0.9 | 0.985 |
| 0.5 | 0.9 | 0.952 |
| 0.20 | 0.89 | 0.928 |
| 0.05 | 0.87 | 0.778 |
| 1/100 | 0.89 | 0.615 |
| 1/200 | 0.88 | 0.513 |

# Analysis of result:

Speed and Accuracy analysis:
From the results above, it's found that Support Vector Machine is the most computation heavy algorithm, especially for Bank marketing data. As mentioned before, in order to target similar prediction accuracy, SVM spends much more time than other algorithms. When kernel was tuned to be "linear", the phenomena was fully exposed. In chess dataset, when "linear" kernel was implemented, it was 400 times slower than using "rbf" kernel. On the other hand, for chess dataset, the calculation time of SVM is much faster than SVM for bank marketing dataset. The reason may be that chess each chess attribute has less possible values. Therefore, SVM only needs to bring up minimum amount of dimensions.

Boosting, decision tree and KNN are much faster. Meantime, they generate about same accurate models as SVM, sometimes even better.

Comparing the performance of decision trees on two datasets, we can find that it indicates overfitting for bank marketing data. When min_samples_split becomes larger (minimum amount of nodes allowed for splitting) and max_depth becomes smaller (the maximum levels of tree), the accuracy tends to go up. That means by default the decision tree generates too many levels and split to too many branches to fit the training dataset. The resulted model doesn't general pattern. On the other hand, the decision tree does very well with chess data set. By default, it renders a 99.6% prediction accuracy. By reducing the tree branch and leaves, its performance gets poorer. One reason is that chess board is very deterministic. When chess board is described by multiple attributes in enough detail, whether White can win or not becomes certain. Decision Tree is extremely good at describing deterministic events.

Boosting and KNN doesn't have too much difference in performance. For both datasets, the more estimators (for boosting) or more neighbors (for KNN), the more accurate prediction will

be. Of course, that means more training time. For boosting, it seems number of estimators and training time have a linear relationship. This relationship for K-nearest neighbors is not obvious.

Learning Curve analysis:
This part is to find the relationship between size of training data and model prediction accuracy. For chess data, 1/20 of the original data (3196 incidents) will bring the forecast accuracy to a reasonably good level. After that, increasing the number of data will slightly improve performance of the model. This holds true for all the tested algorithms.

However, for bank marketing data, an interesting phenomena was observed. No matter how small the size was, the model performance wouldn't be affected by too much for all the algorithms. One possible reason is that the chessboard can't be fully described by too few number of incidents. On the other hand, for bank dataset, several individual data can be used to represent majority of the population. Since whether people will decide to subscribe the service is a stochastic process. Therefore, models can generate decent prediction result based on a few data points but it will be very hard to surpass certain threshold, 95% for example.

Neural network is also a very accurate algorithm with a little longer training time. For five-fold cross validation, it takes on average 5 seconds to form the data. However, if hidden layer 'o' is applied, the training time is magnificently decreased with the accuracy remain unchanged. This can be a good selection when speed is a major concern in the project.

Overall, decision tree is the best algorithm for both datasets. Not only does decision tree require minimum amount of training time and prediction time, but it also maintains as one of the algorithms that has the highest prediction accuracy. When given minimum viable amount of training data, decision tree gets close to optimal forecasting accuracy very fast.

Improvements on models:
Many things could be done to improve the accuracy of the algorithms. In general, if more time is allowed and more computation power is accessible, I would apply GridSearchCV, which is a nice little widget we can use to automatically adjust the parameters in the algorithm to find out the best tuned model for certain algorithm. This will give me more stringent selections of parameter settings instead of trying different settings manually. For instance, for decision tree, instead of guessing the most appropriate min_samples_split or max_depth, if I throw in a list of parameters, the system will generate the best-fitting pruning metrics in a systematic manner.

More robust cross validation should be applied when time permits. Currently for all the algorithms except SVM I applied 5-fold cross validation. A higher order-fold cross validation will minimize the variance and avoid the situation where data is not shuffled and split random enough. This can become an issue, especially when dataset is not big enough.

For Support Vector Machine, due to the limit of the computing power, I was only able to do three-fold for the cross validation. Although it eliminates some of the bias, a more divided cross validation dataset is preferred. One way to solve the computing issue is to run codes in virtual

machines over cloud service, like AWS or Azure. More research will be conducted on how to do that. Another issue is that system always crashes when I try to run 'linear' kernel on the dataset. For boosting, if different forms of weak learners can be specified, the prediction accuracy can be improved and less learners will be needed.

Moreover, as only neural network is tested using weka GUI, the configuration between sklearn and weka might be different and hence result in different result. Given more time, I would test all the other four algorithms on weka and see whether there will be performance change.

In order to understand all the five algorithms more thoroughly, classification problems with more than two categories can be experimented in the future. Both datasets only ask for the classifier to categorize the data into two categories. For bank marketing data, it's whether customers will subscribe; for chess data, it's whether White will win or not. In reality, a machine learning classification problem can be much more complicated with more categories. For instance, depending on the environment, an intelligent robot should be able to respond with actions that's more than just picking two options. For example, for a simple two category problem decision tree can do really well, but maybe neural network will beat all other algorithms when solving a more complex problem.