

Report for Machine Learning Project 1

Yiwei Yan

September 20, 2015

1 Problem Statement and Dataset Selection

The classification methods discussed in this reported are among those simply yet effective algorithms developed mostly in the early phase of machine learning (ML) research, including *K*-Nearest Neighbors (KNN), Decision Trees, Boosting (Adaboost decision trees), kernel Support Vector Machines (SVM), and Neural Networks (NN) [2, 3, 4]. To evaluate these algorithms, I selected the Iris dataset for algorithm illustration, the USPS digit dataset and the Bank Marketing dataset for experimental evaluation.

Table 1 provides some details of the datasets being used. In particular, these datasets and the associated classification problems are picked due to the following reasons:

- **Iris dataset:** Iris flower set is a simple benchmark dataset that I use to illustrate the algorithms. For visualization, I perform a dimensional reduction using principal component analysis (PCA) to reduce the feature dimension to 2. Figure 1 (LEFT) visualizes the Iris data by projecting it to the top 3 principal components.
- **USPS dataset:** USPS digit dataset refers to numeric data obtained from the scanning of handwritten digits from envelopes by USPS. This dataset is especially interesting due to the following reasons. Figure 1 (RIGHT) visualizes some image patches re-generated from the USPS features.
 - (1) Classifying handwritten digits is a very meaningful application in reality.
 - (2) It illustrates the typical difficulties in computer vision (which is often an important application of ML), that the features by simply stacking image pixel intensity only represent very low-level information. Considering that a major research effort in computer vision is devoted to feature engineering, classification directly based on raw pixel intensity is definitely a challenging task.
 - (3) The nature of the dataset itself is interesting to ML, because it is of multi-class (10 classes), high/middle-size feature dimension (256, from 16×16 image patches), but each feature dimension is continuous and well scaled with all values between 0 and 1.
 - (4) Finally, it is a benchmark dataset and thus there are a lot of performance evaluation results as a comparison to the results in this report.
- **Bank Marketing (BM) dataset:** The Bank Marketing dataset [5] is retrieved from the UCI dataset library (<http://archive.ics.uci.edu/ml/datasets/Bank+Marketing>). The classification target is binary, representing “whether a client subscribed a term deposit”, and the features of 20 dimensions represents 20 social and economic attributes of the client, such as age (numeric), type of job (12 categories), marital status (4 categories), number of employees (numeric), etc. This dataset is very interesting because:
 - (1) It is from a real application and has practical meanings for bank industries.
 - (2) This dataset is of a large size, with 41188 samples. This can not only test the algorithm performances in large-scale learning, but also highlights the differences of training/testing runtime.
 - (3) The features are challenging due to its heterogeneous nature. The feature contains both continuous (numeric) and discrete (categorical) attributes. Moreover, even the numeric values are not well scaled. For example, the age attribute has a limited numerical range, but the “number of employees” attribute can differ by thousands. This in fact will highlight the benefits of performing data normalization before learning, as will be discussed later in the report.

Table 1: Dataset descriptions

Dataset name	# samples	Feature dim.	Feature attributes	#classes	Usage
Iris	150	4	continuous	3	illustration
USPS	9298	256	continuous	10	evaluation
Bank Marketing	41188	20	continuous + discrete	2	evaluation

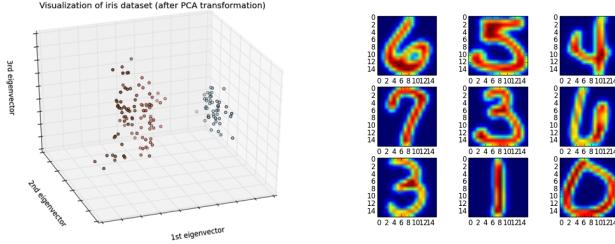


Figure 1: Visualization of Iris dataset (Left) and USPS dataset (Right).

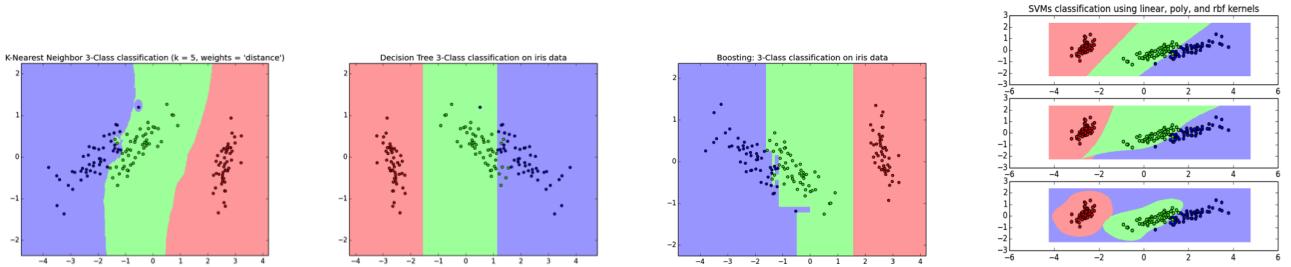


Figure 2: (1) Decision boundaries of the KNN classifier trained on Iris data, visualized by projecting to the top two principal components. (2) Decision Tree. (3) Adaboost Decision Trees. (4) Kernel SVMs.

- (4) The data is “not perfect”. On one hand, some attributes are missing or unknown. Six out of the 20 attributes (such as education, marital status, etc.) can have the value as “unknown”. On the other hand, the data is collected manually from clients’ responses, and thus it may be contaminated due to potential fake responses. Note that such a contamination of data contamination is practical in social/economic science¹.

I wrote the code in Python for all algorithms except Neural Network based on `Scikit-learn` package, and the Neural Network algorithm based on `Pylearn2` package. I integrated all algorithms with Cross-validation into one single API function call “CVtest”, which reports training and testing accuracy as well runtime. The main experiment script calls “CVtest” and plots the results.

2 Methodologies, and Illustration on Iris Dataset

To illustrate the five algorithms, I use the Iris dataset as an example.

2.1 k -Nearest Neighbor (KNN)

KNN is a straightforward method. Given a predefined number of neighbors to retrieve, it find the predefined number of training samples closest in distance to the new point, and predict the label of the testing datum from these. The performance of KNN can be affected by different parameters, including the number of neighbors (the k), the formulation of “distance”, and the neighbor searching/retrieval methods.

Figure 2 (1) depicts an example of KNN using my code, trained on the Iris dataset (note that the training process of KNN is simply storing the training sample in some data structure). In this example, $k = 5$, and Euclidean distance is used. The decision boundary visualized in the subspace spanned by the top two principal components reveals that, the boundary is purely data-driven and determined by the localities of training samples.

2.2 Decision Tree

Decision tree is a non-parametric model for supervised learning, which learns if-then-else decision rules according to certain criteria. The Criteria aim to find the split of feature attributes to obtain optimal impurity defined by different measures. Typical measures of impurity include the entropy and the Gini measure.

Without pruning, the tree grows to reach the optimal split for the **training set**. However, this may lead to overfitting, because the optimal split for the training set may not be the best for the testing set. Thus, restricting the maximum depth D_t of the tree is important to combat overfitting. Figure ?? (2) shows the result of pruning the tree to only three levels. It can be seen that the decision boundary works well for the testing set.

¹I personally encountered similar data contamination when I worked on my previous M.Sc thesis on Applied Economics :-).

Another key observation from Figure 2 (2) is that the decision boundary of decision tree is composed by “hyperplane segments parallel to feature dimensions”, as shown in the figure that the boundaries are all perpendicular to the x -axis. This is because each split only operates on one feature dimension.

2.3 Boosting (Adaboost) Decision Tree

I used the Adaboost for boosting the decision trees. AdaBoost aims to train a sequence of weak classifiers/regressors by repeatedly modifying versions of the data. Adaboost decision tree uses decision trees as the weak classifiers. The predictions from all of them are combined through a weighted majority vote to produce the final prediction.

Figure 2 (3) shows the learned decision boundary using my code for an Adaboost Decision Trees Classifier on Iris dataset, with each decision tree pruned to three levels and with 35 decision trees. The boundary has similar properties to a single decision tree, because the splitting is still with respect to one dimension. However, the boundary is more irregular due to the classifier ensemble.

2.4 (Kernel) Support Vector Machines (SVM)

SVM has been an extremely popular method for a long time due to its various advantages. From my point of view, some major advantages include: the associated optimization programming is usually convex and thus relatively easy to solve by quadratic programming (QP); as a kernel machine, it can capture nonlinearity in the dataset with kernels by mapping data into a higher (maybe infinite) dimensional space; it is sparse, with the boundary depending on the “support vectors”.

Associated with the kernel method are the kernels of different types. Using the code submitted, the decision boundaries of SVMs with different kernels on Iris set are obtained as depicted in Figure 2 (4). It can be observed that the shapes of decision boundaries highly depend on the kernel. A linear kernel gives linear boundaries (straight lines), while polynomial and rbf kernels gives nonlinear boundaries characterized by polynomial functions and radial basis functions respectively.

2.5 Neural Networks (NN)

When using NN for classification, the most important factors are number of hidden layers and the number of nodes/units in each layer. A common observation is that, increasing the number of nodes and layers will reduce the bias of the NN model, which may increase the classification accuracy but in the risk of overfitting.

The activation function is also an important factor. The common choices of activation functions include Sigmoid, Tanh, and Rectifier. In the experiment I will use Sigmoid and focus the analysis on numbers of units and layers.

3 Results and Analysis on USPS Digit Dataset

This section presents the experimental results and analysis on USPS dataset using the methods discussed in section 2 under **Cross-Validation (CV)** framework. In order to differentiate the performance, I randomly selected a subset of data with 1000 samples to perform the experiments (if the whole dataset is used, then almost all methods give 95% to 98% testing accuracy, which makes it difficult to compare).

3.1 k-fold Cross-Validation

I use the k-fold cross-validation framework to evaluate the performances due to the benefits of CV.

To fairly validate the performance of a classifier, it is necessary to partition all available data into two mutually exclusive sets: the training set and the testing set. However, by partitioning the available data, we drastically reduce the number of samples which can be used for learning the model, and the results can depend on a particular random choice for the testing set. The cross-validation procedure aims to solve this problem. Typically in a k-fold CV, the dataset is split into k smaller sets. Then for each of the k folds, a model is trained using k-1 of the folds as training data, and the resulting model is validated on the remaining part of the data. Accordingly, a k-fold CV has k iterations during the whole process.

3.2 Results and Analysis using KNN on USPS

In order to find the best KNN model, I first tested different distance metrics. In details, the general *Minkowski* distance is used. I tested two ways to compute the weights for Minkowski distance: (1) **uniform** weights, with which all points in each neighborhood are weighted equally; (2) **inverse distance** weights, which weight points by the inverse of their distance, such that closer neighbors of a query point will have a greater influence than neighbors which are further away. Moreover, I also tested the norm used for computing the Minkowski distance, including l_1 and l_2 norms.

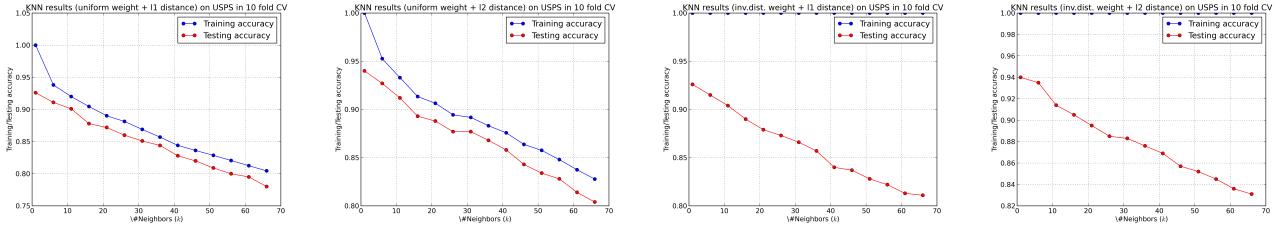


Figure 3: Parameter sweep on k (number of neighbors) in KNN algorithm on USPS dataset. Four combinations of distance metrics are tested, with uniform/ inverse-distance weights combined with $l1$ / $l2$ norm.

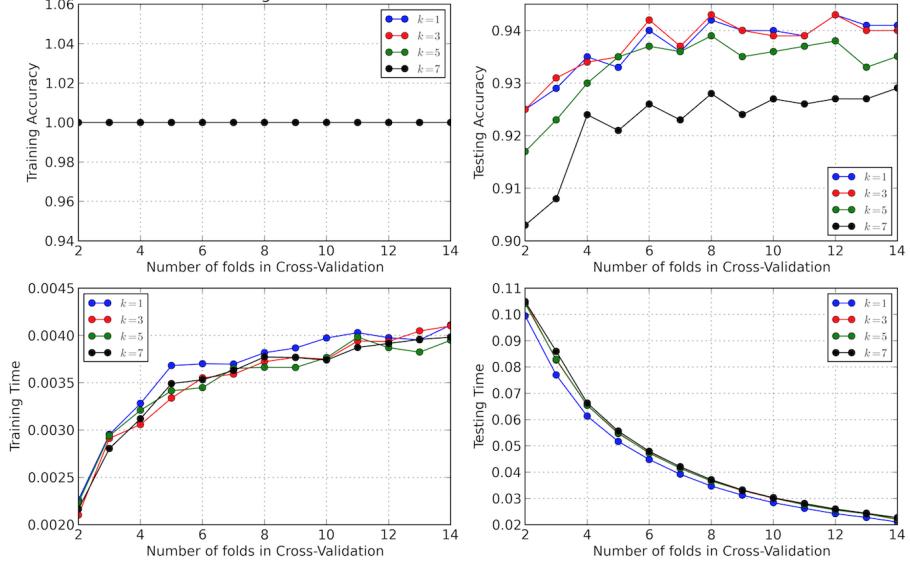


Figure 4: Results of KNN on USPS with different training size versus testing size ratio.

To evaluate the effect of different k values (the number of neighbors to retrieve) and pick the best distance metric at the same time, I perform parameter sweeping experiments on k in the four combinations of distance metrics. The results are showed in Figure 3. From this result we can see that:

- Increasing the value of k does not necessarily increase the accuracy. In fact, on the dataset being tested, increasing k , especially after k is larger than 10, will decease the testing accuracy almost monotonously.
- Inverse-distance weights gives all 100% training accuracy, while uniform weights do not. This is because the training of KNN is simply storing the training data (in a proper data structure). When testing the training accuracy, each query datum has exactly the same point stored in the KNN model, and accordingly the inverse-distance weight assigns a large weight (infinity in ideal cases) which further gives 100% training accuracy.
- Considering both the training and testing accuracy, we can conclude that inverse-weights with $l2$ distances give the best performance in this experiment.

Using inverse-weights with $l2$ distances, I further perform a thorough evaluation of KNN by controlling the sizes of training set and testing set, with $k \in \{1, 3, 5, 7\}$. The control of sizes is done by **controlling the number of folds in the k-fold CV**. Increasing the CV folds will increase the **ratio of training size versus the testing size**. In this way, we can evaluate the performance under different sizes, and **maintain the benefits of CV** for fair comparison.

Furthermore, four key performance metrics are evaluated, including: training accuracy, testing accuracy, training time, testing time. To best minimized the randomized effects of CV partitioning, I evaluated all these metrics by averaging the metrics across all k iterations in one k -fold CV. The results are shown in Figure 4.

From Figure 4, we can see that as the ratio of training size versus testing size increases, the testing accuracy increases, training time increases (in log/log-linear speed), and the testing decreases. This is reasonable, because with larger training dataset, the training should be less biased, and the model captures more the underlying structure embedded in the data. The best testing accuracy is given by $k = 3$.

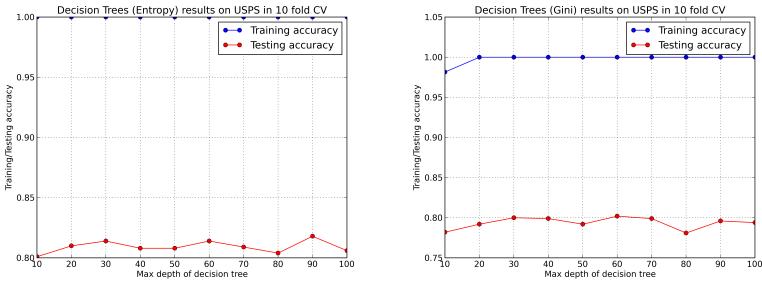


Figure 5: Parameter sweep on D_{DT} (maximum depth allowed) in Decision Trees algorithm on USPS dataset, with entropy and gini criteria respectively.

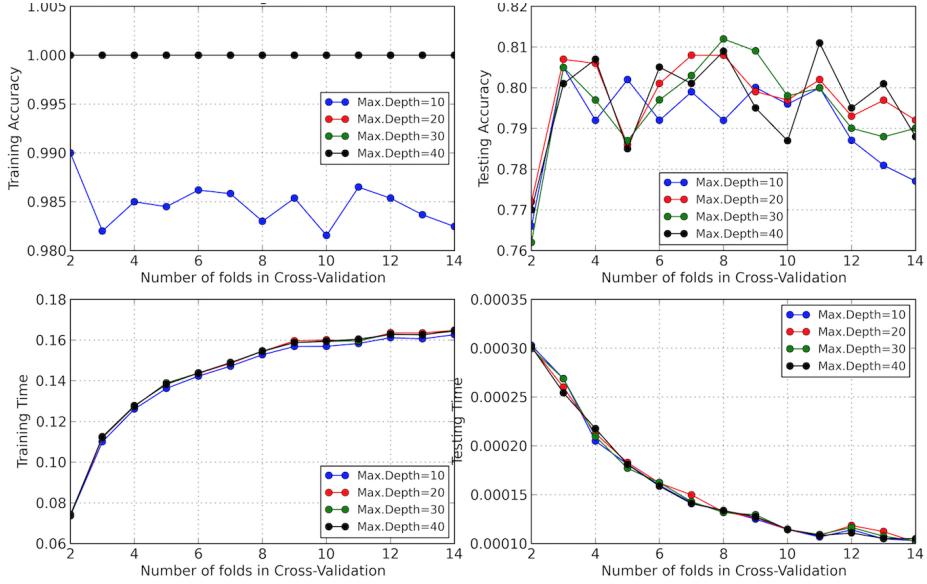


Figure 6: Results of Decision Trees on USPS with different training size versus testing size ratio.

3.3 Results and Analysis using Decision Trees on USPS

As discussed in Section 2.2, an important factor of decision trees is the maximum tree depth D_{DT} . With larger D_{DT} , the decision tree is easier to overfit; however, if D_{DT} is too small, the tree may become overly biased and not flexible enough to model the embedded structure. Therefore, I first perform a parameter sweep on D_{DT} to find a good range of D_{DT} . Furthermore, because there are two popular criteria for attribute splitting (entropy and gini), two experiments are performed with the two criteria respectively.

Figure 5 depicts the results of parameter sweeps. From the figure we can see that the testing accuracy increases monotonously when $D_{DT} \leq 30$, but after that the benefit of increasing D_{DT} is not obvious. Moreover, the entropy criterion and gini criterion give similar results. The mean testing accuracy using entropy is 79.37% and gini is 80.92%. Therefore, I pick gini for the next experiment.

A comprehensive experiment by controlling k-fold CV is then performed. The results are depicted in Figure 6. From Figure 6, we have the following analysis:

- As the size of training set increases and testing set decreases, the testing accuracy first goes up but then drops (slightly). This indicates that overfitting can happen in decision trees when training set is large.
- Excessively restricting the maximum depth can decrease the testing accuracy, because the model itself is more biased. In particular, $D_{DT} = 10$ gives the worst testing accuracy as well as training accuracy in the experiment.
- Depth of the tree does not result in much difference in training/testing time, especially for testing time the deeper tree may result in even faster prediction under the OS scheduling. This is because the prediction of decision trees is extremely efficient, which can be viewed as running through if-else statements.

3.4 Results and Analysis using Boosting Decision Trees on USPS

Based on the decision trees result, I also use the Gini criterion in Adaboost Decision Trees. The number of weak classifiers for ensemble is a trade-off between efficiency and accuracy (more weak classifiers being used is expected

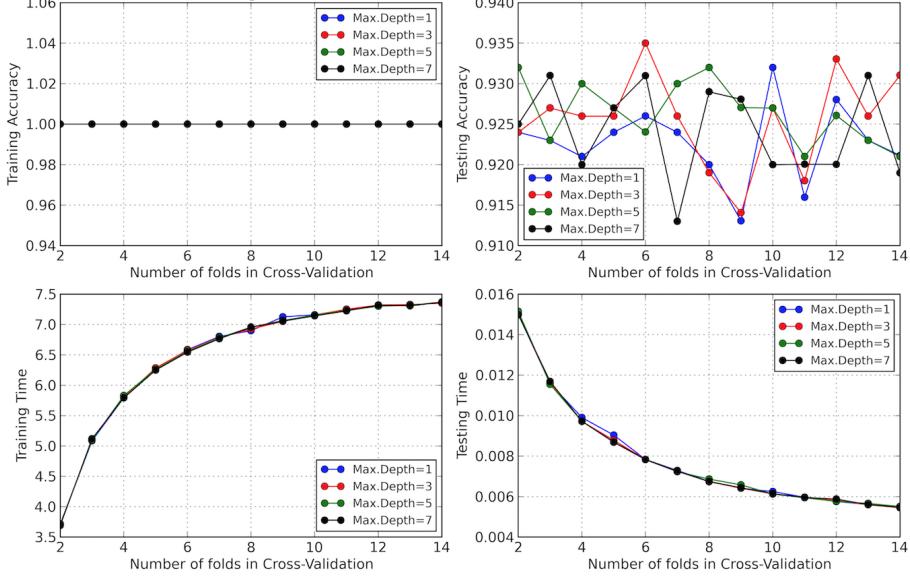


Figure 7: Results of **Adaboost** Decision Trees on USPS with different training size versus testing size ratio.

to give better accuracy but sacrifice efficiency), I maintain this number to be 50. For this multi-class problem, the SAMME [6] method is used for boosting.

The experiment focuses on investigating into the effect of tree-pruning in Adaboost decision trees. Due to the benefits of feature ensemble, the pruning can be more aggressive than a single decision tree. In particular, I tested $D_{DT} \in \{1, 3, 5, 7\}$, and the results are shown in Figure 7.

From Figure 7, we have the following analysis:

- Even with the aggressive pruning, the testing accuracy with Adaboost is greatly improved over a single decision tree. This demonstrates that ensemble methods improve generalizability / robustness over a single estimator.
- With 50 classifiers, the changes of maximum depths in each weak classifier does not result in obvious changes in testing accuracy (the training accuracy is always 100% in this case). This indicates that we can apply very aggressive pruning to prune the tree to 1 level. This also justify why people tend to use “decision stump”, whose depth is 1, in ensemble methods.
- The Adaboost DT is expensive in both training time and testing time, but the different depths result in very similar time. This again demonstrates the depth of each tree does not affect the performance (e.g. training convergence rate, etc.) of the final ensemble classifier that much.

3.5 Results and Analysis using (kernel) SVMs on USPS

One of the major factor of performance in SVMs in the kernel used. Widely used kernels include: (1) linear kernel, $\langle x, x' \rangle$; (2) polynomial kernel, $(\gamma \langle x, x' \rangle + r)^d$; (3) radial basis function (rbf) kernel, $\exp(-\gamma|x - x'|^2)$; (4) sigmoid ($\tanh(\gamma \langle x, x' \rangle + r)$).

I tested all of these kernels in the experiments. Furthermore, the hyper-parameters, such as the bandwidth in rbf kernel, greatly affect the performance. I **tested different values for hyper-parameters** in a factor of 10, i.e. $\gamma \in \{0.001, 0.01, 0.1, 1\}$. Here I only show the best two combinations for the interest of space.

From the results in Figure 8 and 9, we have the following analysis:

- For this dataset, hyper-parameter $\gamma = 0.01$ gives better result than $\gamma = 0.1$. Note that ideally the hyper-parameter should be tuned differently for different kernels.
- When $\gamma = 0.01$, the polynomial kernel gives the best training and testing accuracy, while the rbf kernel also gives similar result. This indicates the USPS data lies on a nonlinear embedded structure.
- When $\gamma = 0.01$, the most expensive to the least in training time is: sigmoid > rbf > polynomial > linear. This is the same for testing time, which is heavily affected by the number of support vectors defining the decision boundaries.
- The training and testing accuracies also increase when training set size increases, but they saturate very fast, indicating that SVMs are actually not very “data-hungry” since they are not models with extremely low bias.

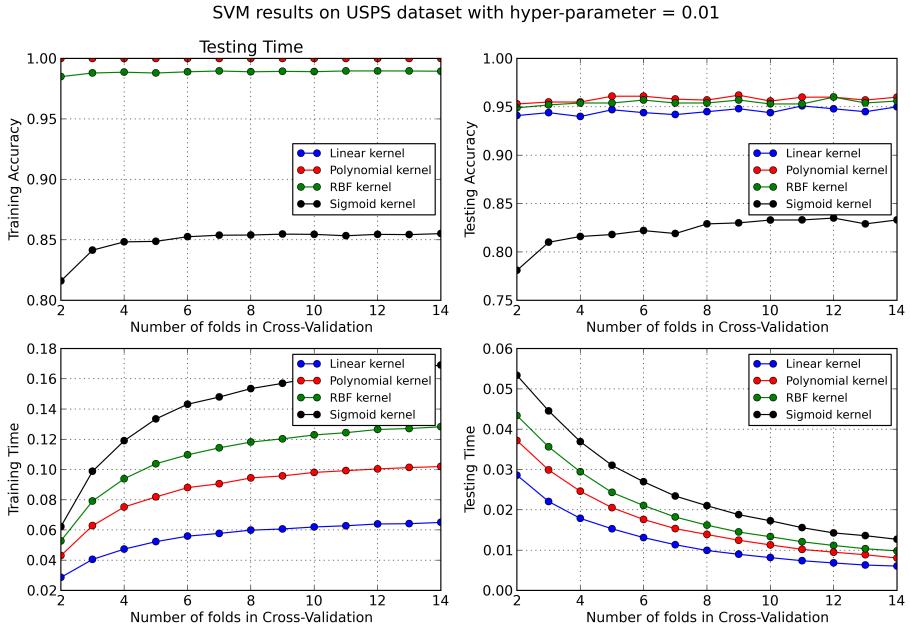


Figure 8: Results of kernel SVMs on USPS with different training size versus testing size ratio, with hyper-parameter $\gamma = 0.01$.

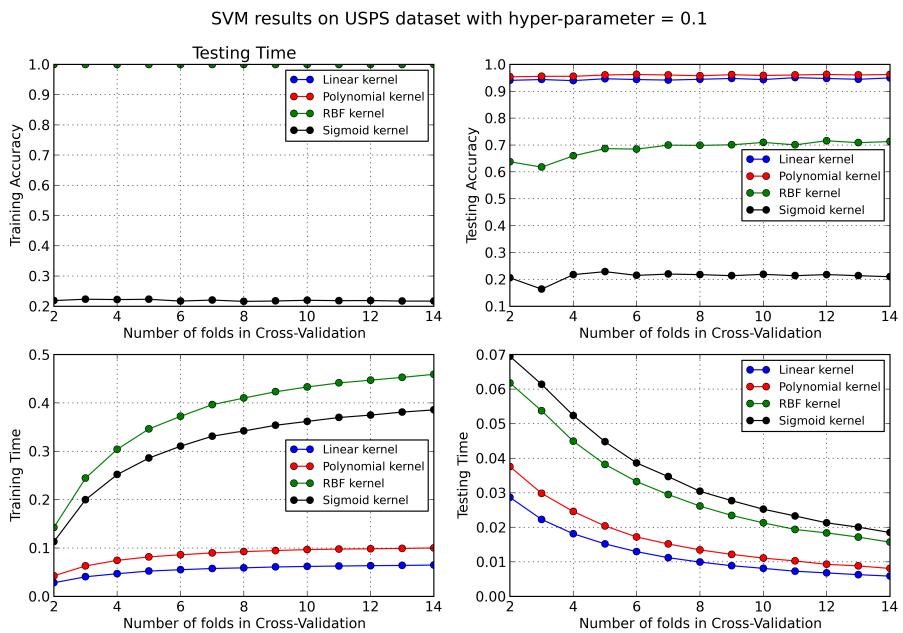


Figure 9: Results of kernel SVMs on USPS, $\gamma = 0.1$.

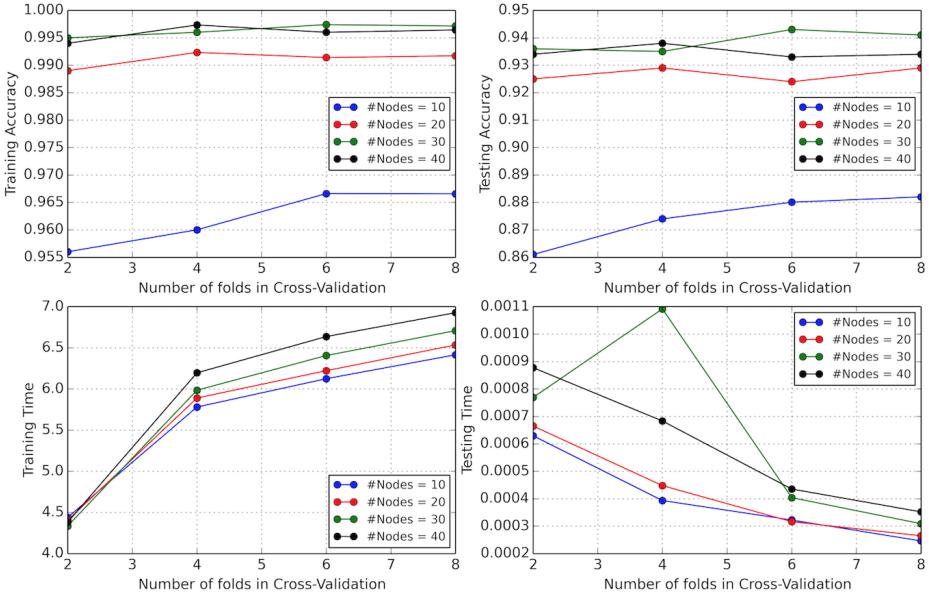


Figure 10: Results of Neural Network with **one hidden layer** on USPS with different training size versus testing size ratio.

3.6 Results and Analysis using Neural Networks on USPS

For the Neural Network models, I tested the structures with one and two hidden layers respectively. Sigmoid functions are used as the activation functions for both the hidden layers and output layer. Furthermore, I evaluated the performance with different number of units/nodes in the hidden layers. The results of Neural Network with one hidden layer (1-HL) are shown in Figure 10, and the results with two hidden layers (2-HL) are shown in Figure 11.

From the results, we have the following analysis:

- For both models with 1-HL and 2-HL, increasing the number of hidden nodes will increase both the training and testing accuracies. However, as this number increase, saturation status may be encountered. For example, in 1-HL model when node number reaches 30, and in 2-HL model when node number reaches 45, the increment of accuracy is not large any more.
- With one more hidden layer, 2-HL has higher testing accuracy than 1-HL when the node number is small. However, in this experiment when node number is large, the benefit of adding more HL is not very evident. In my opinion, this may due to (1) the sample number is not large enough, (2) more optimization iterations should be taken (I constraint the iteration number to 30).
- Training is expensive, especially when node number or layer number is large. However, testing of Neural Network is quite efficient.

I will leave further discussion to section 5.

4 Results and Analysis on Bank Marketing (BM) Dataset

This section presents the experimental results and analysis on BM dataset. Due to similar reasons stated in section 3, a random subset of 5000 samples are used for experiments. **In this experiment, I will not repeat the observations/analysis the same as USPS experiment; rather, I will analyze the *different* observations.**

4.1 Data Preprocessing

The BM dataset has 20 features, but the features are mixed with continuous(numeric) data and discrete(categorical) data. The categorical data are represented by text. In order to apply the learning models, the categorical data is first transformed to numeric representations with natural numbers.

Furthermore, since the range of values in each feature dimension is very different, a **feature normalization** step is performed to re-scale all dimensions. This step is very important for the numerical stability in some algorithms like SVMs, and it speeds up the convergence a lot. In fact, without normalization, the training time of a SVM is more than 10x longer than with normalization (note that the training complexity of SVM can vary between $\mathcal{O}(N^2)$ and $\mathcal{O}(N^3)$, N is the sample number). This timing result is not included in the report for the interest of space, but it can be easily tested using my code submitted.

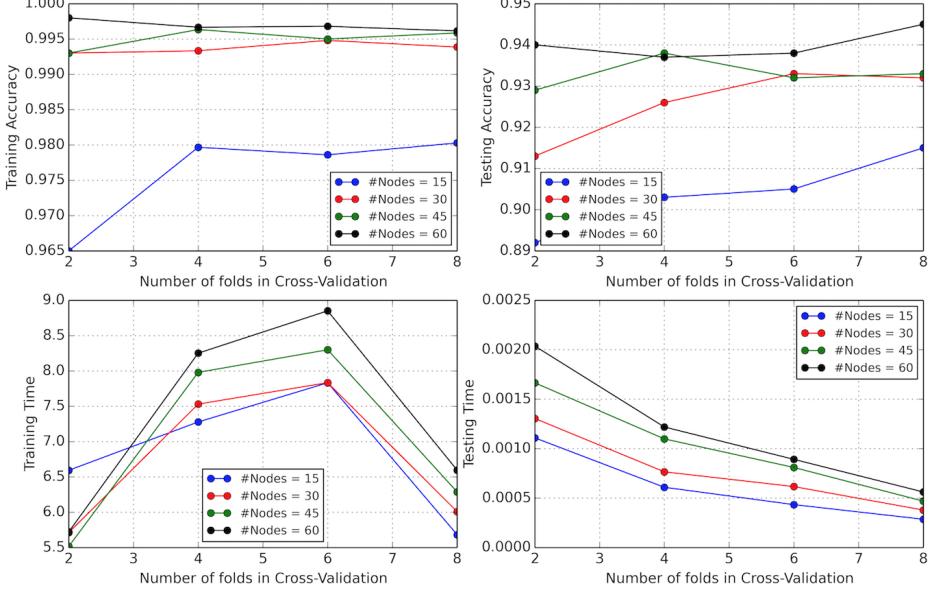


Figure 11: Results of Neural Network with **two hidden layers** on USPS with different training size versus testing size ratio.

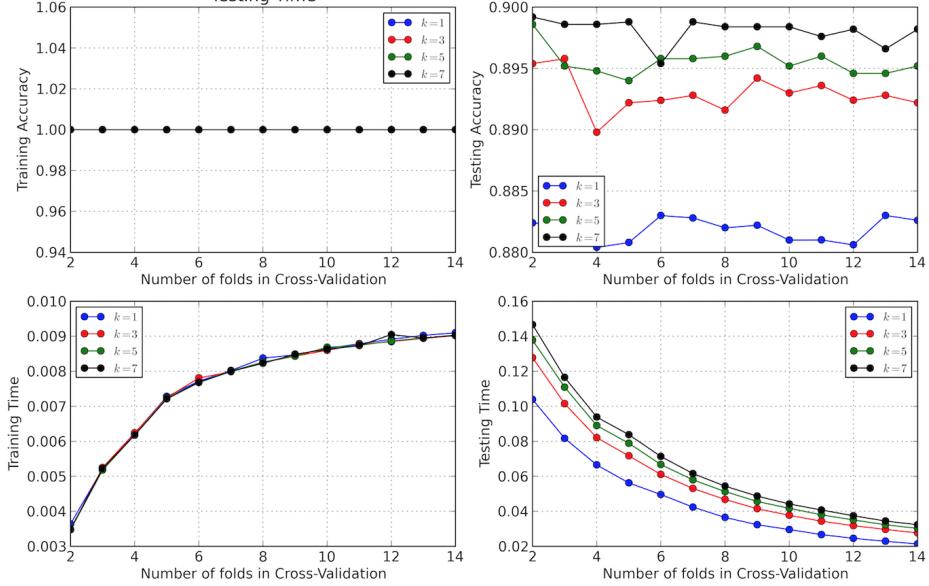


Figure 12: Results of KNN on BM with different training size versus testing size ratio.

4.2 Results and Analysis using KNN on BM

Performing a similar experiment as in section 3.2, I obtain the result in Figure 12.

Observations **different** from USPS experiment are analyzed as below:

- For BM, increasing the k for KNN improves the testing accuracy. The reason may be that the testing data point is quite different from any of individual training data point, and thus instead of using small amount of neighbors to predict the query point, it is better to predict by averaging over a large neighborhood.
- Increasing the training set does not increase the testing accuracy. This can be because redundancy exists in the dataset.

4.3 Results and Analysis using Decision Trees on BM

Similar to section 3.3, the results on BM using a single decision tree are reported in Figure 13.

Observations **different** from USPS experiment are analyzed as below:

- For BM, the most shallow tree with $D_{DT} = 10$ gives best testing accuracy although it does not give 100% training accuracy. Compared to this, $D_{DT} \in \{20, 30, 40\}$ all give similar results lower than $D_{DT} = 10$. This

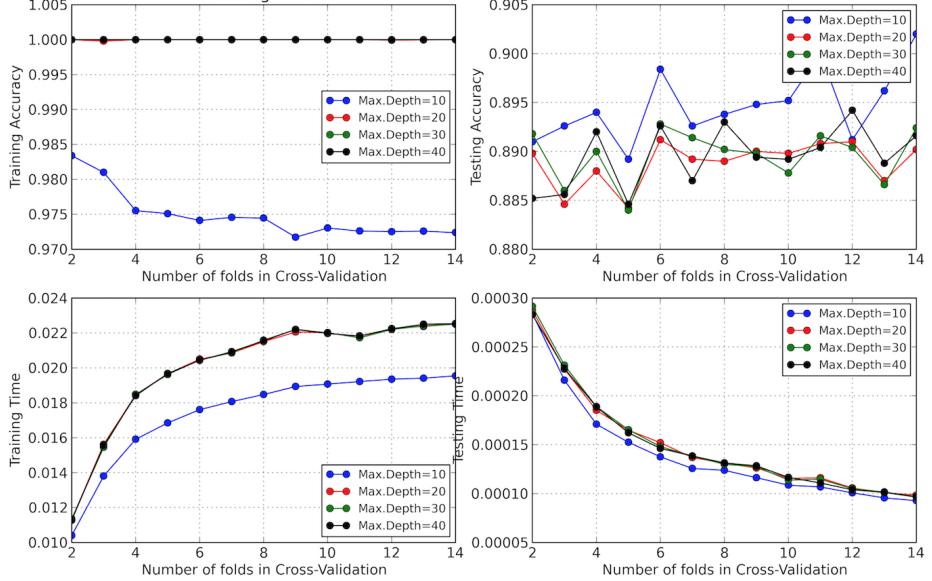


Figure 13: Results of a Decision Tree on BM with different training size versus testing size ratio.

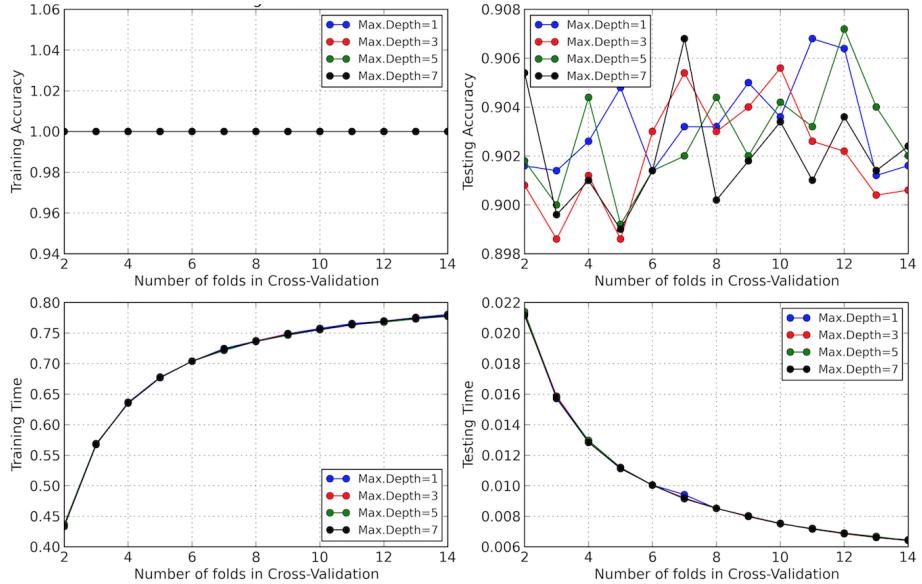


Figure 14: Results of Adaboost Decision Trees on BM with different training size versus testing size ratio.

can be because (1) overfitting, (2) some feature dimensions are redundant or even contain outliers that are detrimental to classification.

- $D_{DT} = 10$ give relative runtime reduction greater than the USPS case, because the training set size is much larger.

Similar to section 3.4, the results on BM using Adaboost decision trees are reported in Figure 14.

The behavior of Adaboost decision trees on BM is pretty similar to the behavior on USPS set. Again, the results demonstrate that “decision stump” is enough to exploit the benefits of the ensemble algorithm.

To apply the SVM on BM, I again performed a parameter sweep on the hyper-parameter γ from 0.0001 to 1 in a multiplicative factor of 10. The best result is given by $\gamma = 0.001$. The result is depicted in Figure 15.

Observations **different** from USPS experiment are analyzed below:

- For BM, the best γ value found is smaller than the one for USPS by a factor of 10. Part of the reason is the scale of values in BM is smaller than in USPS.
- Interestingly, the best testing accuracy on BM is given by linear kernels, while the worst result is from polynomial kernels. Although the absolute differences of testing accuracy is small, it still implies that the BM data lies on an underlying manifold close to linear.
- The fact that linear kernel does not give the best training accuracy reveals that the linear kernel results in a strongest model (in another word, the model is not very flexible due to the linear assumption). It also reveals

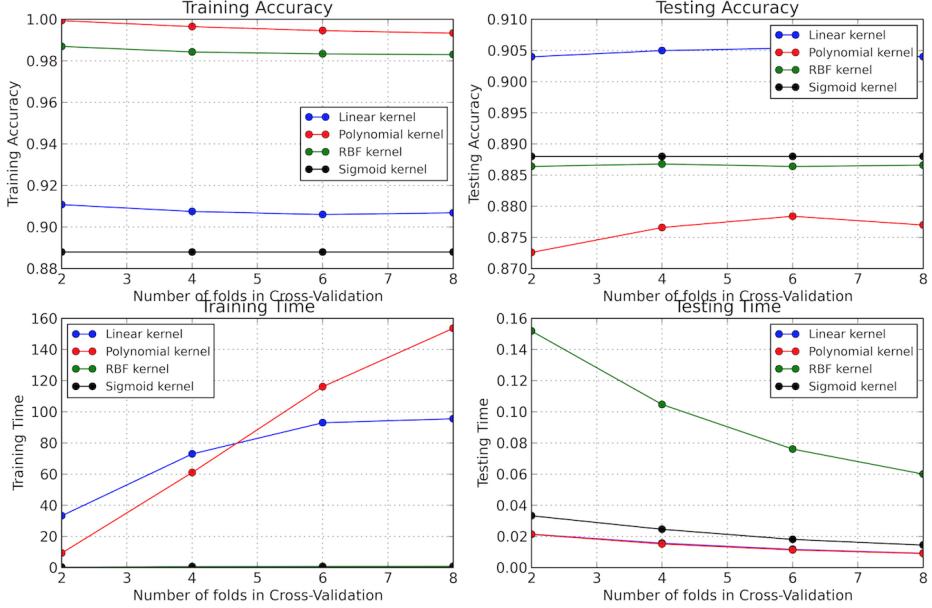


Figure 15: Results of kernel SVMs ($\gamma = 0.001$) on BM with different training size versus testing size ratio.

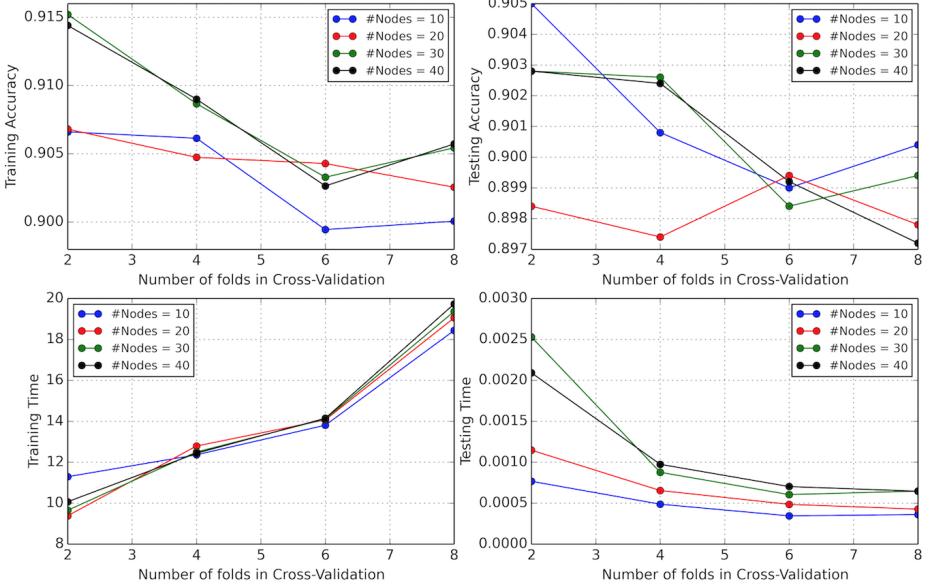


Figure 16: Results of Neural Network (1-HL) on BM with different training size versus testing size ratio.

that if some property of the underlying structure is known, we should exploit this property by using the more bias model, although this may scarify the training accuracy.

4.4 Results and Analysis using Neural Networks on BM

Here both the 1-HL and 2-HL results are reported, in Figure 16 and 17. Different from the USPS result, in 2-HL case number of node ≥ 10 already leads to saturation of testing accuracy. Moreover, unlike the USPS case, the training time in this experiment is very similar under the node numbers tested, but the testing time is in linear order with the node number.

5 Further Discussion and Comparison across Methods

By analyzing the performances across all algorithms and datasets, the following conclusions can be reached:

- Considering the testing accuracy on both datasets, the ranking of algorithms from giving the highest testing accuracy to the lowest testing accuracy is:
 $\text{SVMs} > \text{NN} > \text{AdaboostDT} > \text{KNN} > \text{DT}$,
 where “NN” stands for “Neural Network”, and “DT” stands for “Decision Trees”.

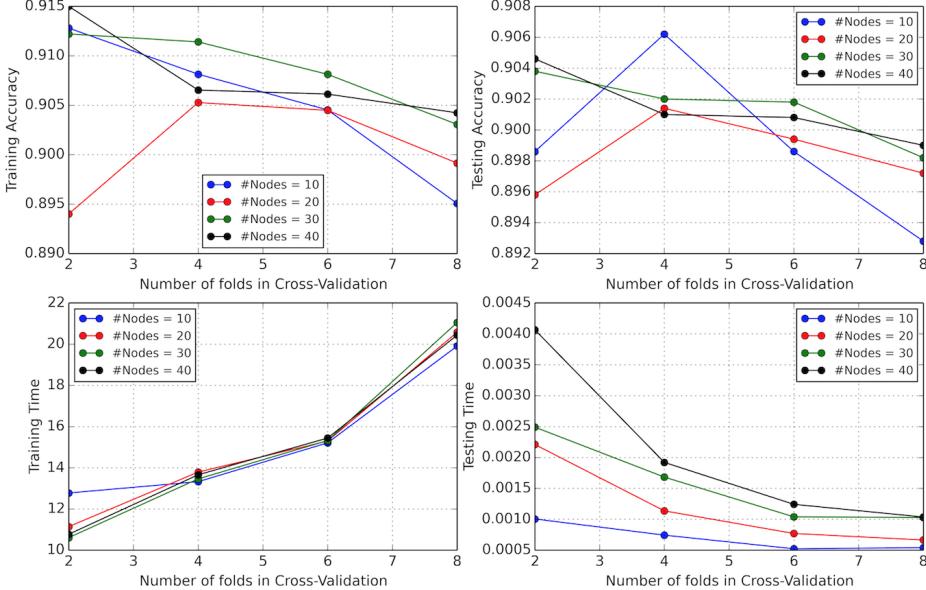


Figure 17: Results of Neural Network (2-HL) on BM.

- Regarding training time, the ranking from the **longest** time to the **shortest** time is: NN > {SVMs or AdaboostDT} > DT > KNN.
- Regarding testing time, the ranking from the **longest** time to the **shortest** time is: {SVMs or AdaboostDT} > NN > KNN > DT.
- Based on the above three conclusions, we can clearly see that different methods have their own benefits. While SVMs in general strike a good balance of overall performance, in the scenarios that accuracy is the most important, AdaboostDT should also be considered; in the scenarios that require fast prediction, DT is the best choice; in the scenarios that needs fast model training, KNN is preferred. NN with more hidden layers (as the structure used in deep learning) is expected to handle (extremely) large-scale and complicated datasets, but it may require more dedicated optimization process. Thus, which one is the “best”, really depends on the application scenario.
- Some other tricks learned from the experiments are:
 1. Usually the more training data, the more accurate the model will be. Also note that though it is a trade-off between the training size and training time, but larger training size may not affect the testing time that much. The testing time is determined more by the model itself.
 2. Overfitting is a practical issue that needs attention. Lower bias model is more flexible to capture different data properties, but at the same time it make result in overfitting. It is necessary to test different models with different degrees of constraints.
 3. Hyper-parameters in some models have a key influence on the final accuracy. In real-world usage, it is important to perform parameter sweep as fine-grain as possible.
 4. Data pre-processing, such as normalization, whitening, re-scaling, can have great impact on both the training time (especially for algorithms relying on numerical optimization) and accuracy.

References

- [1] Michalski, R. S., Carbonell, J. G., and Mitchell, T. M., “Machine learning: An artificial intelligence approach”, *Springer Science & Business Media*, 2013.
- [2] Duda, R. O., Hart, P. E., and Stork, D. G., “Pattern classification”, *John Wiley & Sons*, 2012.
- [3] Bishop, C. M., “Pattern recognition and machine learning”, Springer, 2006.
- [4] Moro, S., Cortez, P., and Rita, P., “A Data-Driven Approach to Predict the Success of Bank Telemarketing”. *Decision Support Systems*, In press, <http://dx.doi.org/10.1016/j.dss.2014.03.001>
- [5] Zhu, J., Zou, H., Rosset, S., and Hastie, T. “Multi-class adaboost”. *Statistics and its Interface*, 2(3), 349-360, 2009.