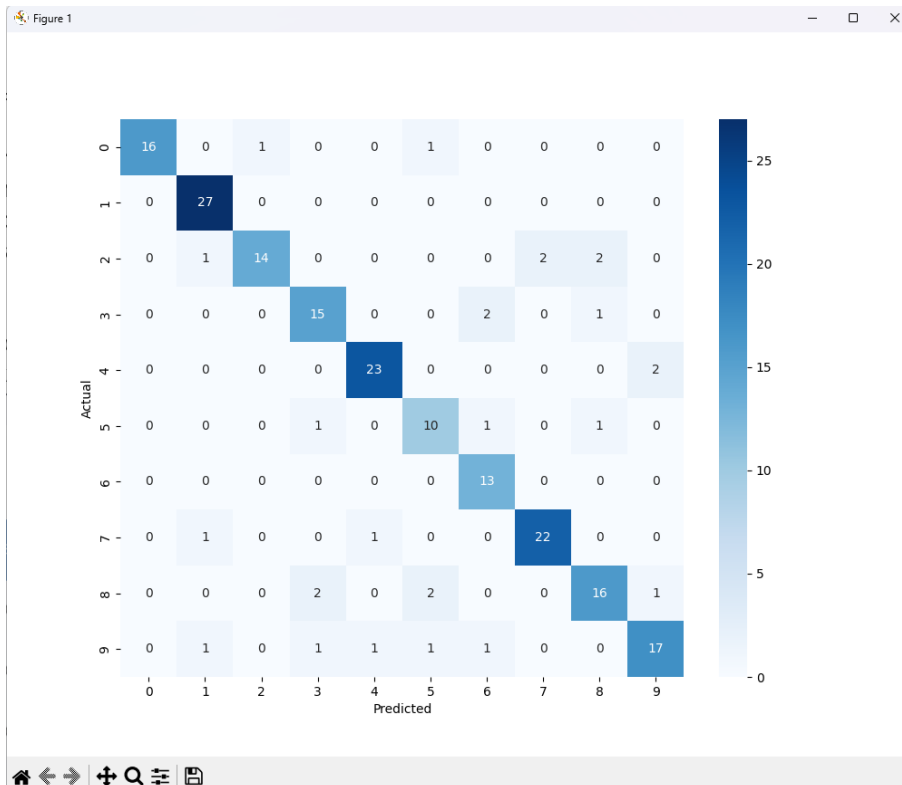


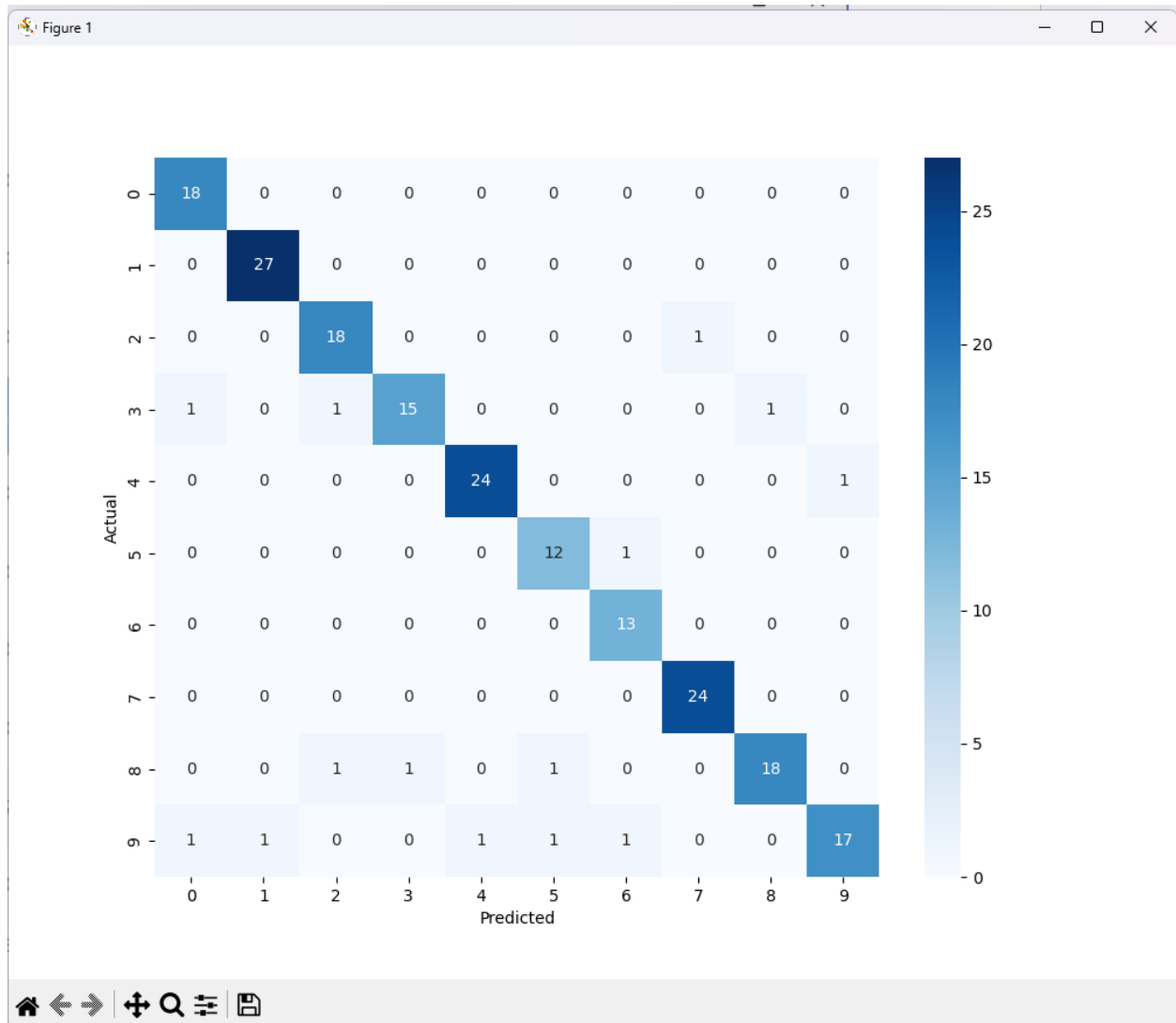
**2. (4.0 points) Implement a k-nearest neighbors classifier for both Euclidean distance and Cosine Similarity using the signature provided in starter.py. Please describe all of your design choices and hyper-parameter selections in a paragraph. Once you are satisfied with performance on the validation set, run your classifier on the test set and summarize results in a 10x10 confusion matrix. Analyze your results in another paragraph.**

When searching for the optimal value of  $k$  in K-Nearest Neighbors (KNN), we typically explore a range of values. In our case, we considered a range from  $k=1$  to  $k=11$ . The inclusion of  $k=11$  is notable because it matches the number of classes in our problem, which allows us to perform a majority vote even when all 10 nearest neighbors are distinct.

After experimentation with the 'euclidean' distance metric on our validation dataset, we discovered that  $k=2$  yielded the best results. Consequently, we adopted  $k=2$  for the KNN model when applied to the testing dataset. This choice resulted in an accuracy of 0.865. The resulting confusion matrix is presented below:



After experimentation with the 'cosim' distance metric on our validation dataset, we discovered that  $k=3$  yielded the best results. Consequently, we adopted  $k=3$  for the KNN model when applied to the testing dataset. This choice resulted in an accuracy of 0.93. The resulting confusion matrix is presented below:



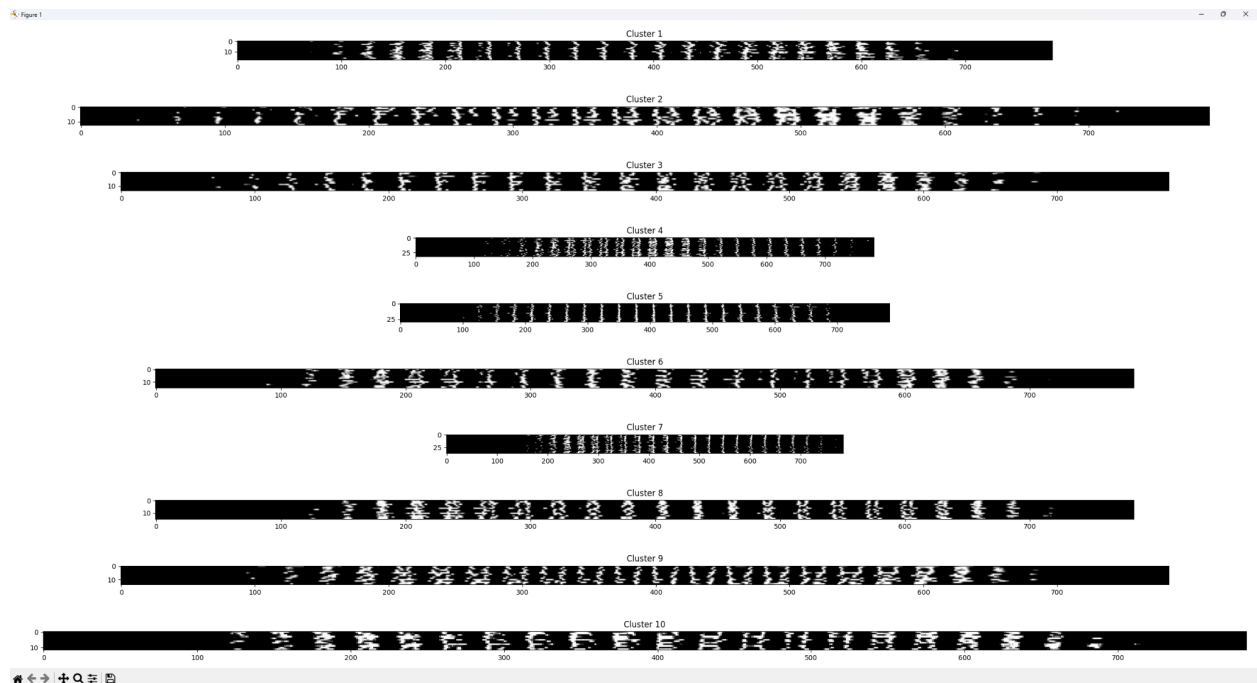
**3. (4.0 points) Implement a k-means classifier in the same manner as described above for the knearest neighbors classifier. The labels should be ignored when training your k-means classifier. Describe your design choices and analyze your results in about one paragraph each.**

We have a total of 10 classes, ranging from 0 to 9. Therefore, we have chosen 'k' as 10 to represent our intention to divide the dataset into 10 clusters.

After experimenting with the 'euclidean' distance metric on our validation dataset, we determined that 20 iterations produced the best results. This resulted in a testing data accuracy of 0.615. It's worth noting that accuracy varies due to the random initialization of centroids at the beginning of the process.

After experimenting with the 'cosim' distance metric on our validation dataset, we determined that 30 iterations produced the best results. This resulted in a testing data accuracy of 0.51. It's worth noting that accuracy varies due to the random initialization of centroids at the beginning of the process.

Here is an example of our test dataset when fit into our clusters.



**4. (1.0 points) Collaborative filters are essentially how recommendation algorithms work on sites like Amazon ("people who bought blank also bought blank") and Netflix ("you watched blank, so you might also like blank"). They work by comparing distances between users. If two users are similar, then items that one user has seen and liked but the other hasn't seen are recommended to the other user. What distance metric should you use to compare user to each other? Given the k-nearest neighbors of a user, how can these k neighbors be used to estimate the rating for a movie that the user has not seen? In about one paragraph describe how you would implement a collaborative filter, or provide pseudo-code.**

In collaborative filtering for recommendation systems, one commonly used distance metric to compare users is cosine similarity. Cosine similarity measures the cosine of the angle between two vectors, representing users' preferences or ratings. It is effective in capturing the similarity in preference patterns between users and is suitable for handling sparse data, which is common in recommendation systems.

Implementation of rating for a movie that a user has not seen based on their k-nearest neighbors:

- Compute the similarity (cosine similarity) between the target user and all other users in the dataset.
- Select the k-nearest neighbors with the highest similarity scores to the target user.
- For the target movie, calculate a weighted average of the ratings given by the k-nearest neighbors, using their similarity scores as weights. This weighted average serves as an estimate of the rating for the target user.

**5. (1.0 bonus points) Prepare a soft k-means classifier using the guideline provided for Question #3 above.**

We developed our soft k-means algorithm as an extension of our k-means algorithm. The key distinction between the two lies in the assignment of data points to clusters. In k-means, each data point is exclusively assigned to a single cluster, whereas in soft k-means, every data point is assigned a degree of membership to all clusters.

We maintained a 'k' value of 10 and performed 20 iterations, consistent with our k-means clustering approach. However, we observed a lower accuracy of 0.135. We attribute this discrepancy to the inherent nature of soft k-means, which tends to distribute data points across multiple clusters. This distribution of data points among various clusters can result in lower accuracy compared to the strict assignments made in traditional k-means clustering. Again, the accuracy varies due to the random initialization of centroids at the beginning of the process.

=====HW1=====

## COMP349 HW1

Chen Si  
Dong Shu  
Ding Zhang

1. (2.0 points) Did you alter the Node data structure? If so, how and why?

- Yes, we did alter the Node data structure.
- Besides the fields labels and children, we also add two fields: attributes and entropy for each node. The field “attribute” of each Node inherits the parent’s attribute, which is an indication that the current node is splitted using that particular attribute. The field “label” is the unique value of that attribute, which serves as edges connecting the nodes if we would draw out the decision tree on paper. If the node is a leaf node, the attribute field would be None, and the value of the label is the predicted class label. We have also added several functions: change\_label, change\_attribute, add\_child and change\_entropy. These methods allow us to change or add values to each of the node values.

2. (2.0 points) How did you handle missing attributes, and why did you choose this strategy?

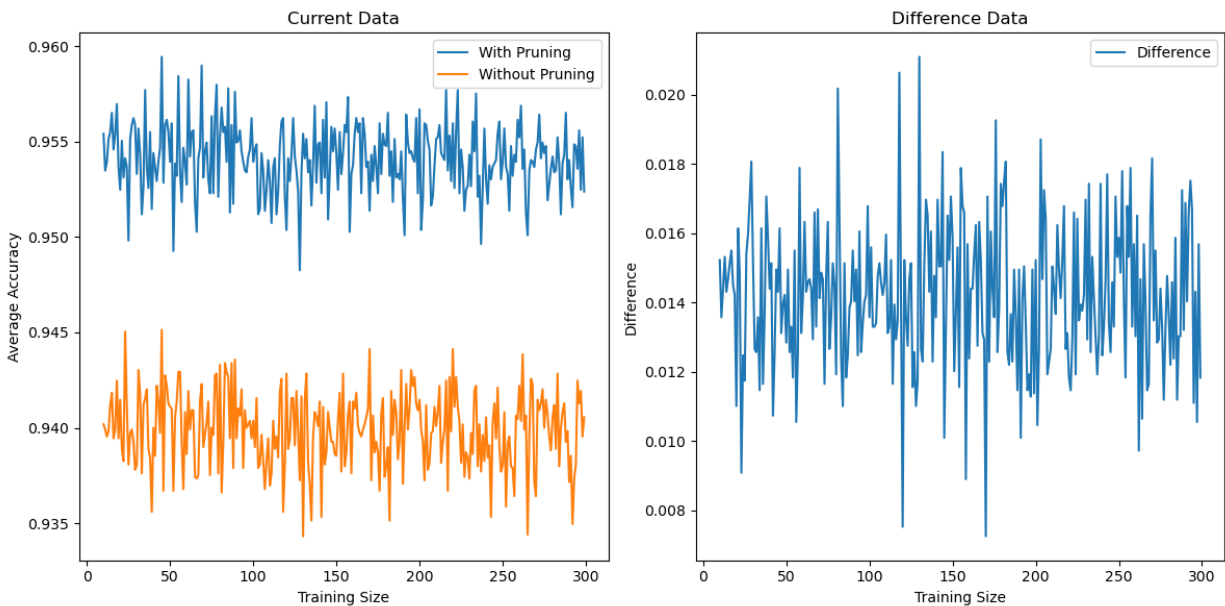
- We replace missing attributes with the most common label in the data rows (the mode of that column). In our opinion, if the data set is large enough, missing values are trivial in a sense that they won’t affect the overall prediction a lot. Thus, we choose to replace missing values with the mode of the column.

3. (2.0 points) How did you perform pruning, and why did you choose this strategy?

1. Pre-Pruning. We perform pruning by calculating the absolute value between the parent’s entropy and the child’s entropy. We prune the node if the entropy difference is greater than some threshold  $\epsilon$  (in our case we set  $\epsilon = 0.1$ ).
2. Post-Pruning. We mainly utilize the reduced error pruning technique. For a given tree and validation set, we first perform the ID3 algorithm on the validation set and record the current precision value. We then find the deepest node of the entire tree and retrieve its parent node. We remove the child nodes from the parent node and recalculate the precision value. If the precision value after pruning is greater than before, we preserve the pruned tree and recursively perform pruning on the next deepest node. If not, we do not perform pruning.
  - a. Two approaches: start performing pruning from the root to the leaf node(work the way down)
  - b. Start from the leaves of the tree and work your way up towards the root.

- i. We perform a pruning function for each node in the tree.
- ii. Temporarily remove the subtree from the validation dataset and calculate the accuracy.
- iii. Compare the accuracy with the accuracy of the original tree (i.e before removing the subtree). If the accuracy remains the same or improves, then we prune the subtree; otherwise, we keep it.
- iv. Repeat this process for all nodes in the tree.

4. (4.0 points) Now you will try your learner on the `house_votes_84.data`, and plot learning curves. Specifically, you should experiment under two settings: with pruning, and without pruning. Use training set sizes ranging between 10 and 300 examples. For each training size you choose, perform 100 random runs, for each run testing on all examples not used for training (see `testPruningOnHouseData` from `unit_tests.py` for one example of this). Plot the average accuracy of the 100 runs as one point on a learning curve (x-axis = number of training examples, y-axis = accuracy on test data). Connect the points to show one line representing accuracy with pruning, the other without. Include your plot in your pdf, and answer two questions:



a. What is the general trend of both lines as training set size increases, and why does this make sense?

- Based on our generated learning curves shown above, we have noticed that the deviation of the pruning curve tends to decrease significantly as the number of training sizes increases. The learning curve for our non-pruned tree, on the other hand, has an increasing (or at least non-decreasing) deviation. Even though it is not significantly clear in this graph, the tree with pruning should have a trend of increasing accuracy as training size increases, and a trend of decreasing accuracy for non-pruned trees.
- This makes sense because the key idea for pruning is trying to reduce overfitting. There is a much larger probability that the tree will overfit on a large dataset, which makes the algorithm harder to generalize to test sets. Pruning will reduce the impact of overfitting and have higher average precision for larger datasets.

**b. How does the advantage of pruning change as the data set size increases? Does this make sense, and why or why not?**

- Similar to the idea above, as the data set size increases, the advantage of pruning will become more significant. A decision tree is very susceptible to overfitting, especially when the size of the dataset increases. This makes it harder for the original tree to generalize to data sets other than the training dataset, as it fits the training set too well. Pruning will remove the observed overfitting and can be better generalized to other testsets. This advantage becomes significantly noticeable as dataset size increases.

*Note: depending on your particular approach, pruning may not improve accuracy consistently or may decrease it (especially for small data set sizes). You can still receive full credit for this as long as your approach is reasonable and correctly Implemented.*

**5. (optional 2.0 points) Use your ID3 code to construct a Random Forest classifier using the candy.data dataset. You can construct any number of random trees using methods of your choosing. Justify your design choices and compare results to a single decision tree constructed using the ID3 algorithm.**

- Using a single decision tree, we achieved an average accuracy of 0.7209090909090909 (with pruning) and 0.6927272727272724 (without pruning) after running the model 100 times on the candy dataset. In contrast, when employing a Random Forest classifier with 100 trees, our test accuracy significantly improved to 0.8636363636363636. Note that this number may vary between different runs, but in general, our Random Forest Predictor had significant improvement on accuracy.
- We opted for a Random Forest size of 100 trees because the testing time for the single decision tree was set to 100 runs. This decision ensures a fair comparison between the two models. It is evident that when comparing a single decision tree to a Random Forest

classifier, the accuracy of the single decision tree, even with pruning, falls behind that of the Random Forest classifier.