

3rd Travis Tran
University of California, Irvine
Irvine, USA
67987698

kNN is an instance-based learning algorithm that essentially speaks for itself; it classifies data points based on the k nearest neighbors in the feature space of the majority class. We experimented with a different number of neighbors (k), using 5-fold cross-validation to determine the optimal k value. To

achieve this result, we used the scikit-learn package to evaluate performance on both training and test sets.

Hyperparameters investigated:

- k: The number of nearest neighbors.
 - Experimented with $k = [1, 5, 10, 50]$

B. Logistic

Logistic regression is a linear model primarily used for binary and multiclass classification. The model applies a logistic function to a linear combination of input features to analyze the probability of a given class. In our implementation, we scaled the features using normalization. Additionally, we experimented with different regularization strength (C) values, using 5-fold cross-validation to determine the optimal C value. To achieve this result, we used the scikit-learn package to evaluate performance on both training and test sets.

Hyperparameters investigated:

- Regularization Strength (C): Smaller values specify stronger regulation
 - Experimented with $C = [0.01, 0.1, 1, 10, 100]$
- Max iterations: Maximum number of iterations for the solver.
 - Experimented with $\text{max iterations} = [500, 1000, 2000, 5000]$

C. Random Forest

The classification made by random forest takes the majority vote of a group of decision trees; each tree trained on a random subset of the given training dataset focuses on random features (aspects of a data point). Inference via decision trees is achieved by choosing junctions depending on whether a threshold is satisfied. To achieve this result, we used the scikit-learn package to evaluate performance on both training and test sets.

Hyperparameter investigated:

- n_estimator: Number of decision trees trained
 - Experimented with $n_estimator = [5, 25, 50, 100]$
- min_sample_split: Minimum amount of data points allowed before splitting further
 - Experimented with $min_sample_split = [2, 5, 10]$
- min_sample_leaf: The smallest amount of data points allowed at a leaf node
 - Experimented with $min_sample_leaf = [2, 5, 10]$

D. Neural Network

For the project, we employed both a Convolutional Neural Networks (CNN) and a Multi layer Perceptrons (MLP). CNNs are deep learning models that consist of convolutional layers that apply filters to extract key features from an image. Additionally, they include pooling layers for cutting dimensionality along with other layers to ultimately classify images. MLPs are a type of feedforward neural network that consists of layers of nodes that are all fully connected to each other. To implement these neural networks, we used scikit-learn for the MLP and PyTorch for the CNN.

Hyperparameters investigated:

- Epochs (for CNN): Number of complete cycles on the entire training dataset
 - Experimented with $\text{epochs} = [1, 5, 10, 20, 30]$
- Batch Size (for MLP): Number of training examples used in one iteration of training
 - Experimented with $\text{batch sizes} = [32, 64, 128, 256]$
- Solver (for MLP): Algorithm used for optimizing weights
 - Experimented with $\text{solver} = ['sgd', 'adam']$
- Max iteration (for MLP): Maximum number of iterations training will run
 - Experimented with $\text{max iteration} = [100, 500, 1000]$

IV. EXPERIMENTAL SETUP

For our project, we decided to use a subset of the data to speed up our models' runtime. As such, we used 10,000 images for the training dataset and 2,000 images for the testing dataset. To evaluate our model's performance, we looked at classification accuracy/error, analyzed the confusion matrix, and measured prediction speed. This would provide a solid benchmark to determine which model achieves high accuracy while maintaining speed.

To experiment with the hyper-parameters of K-Nearest Neighbors (kNN), Logistic Classification, Random Forest, Convolutional Neural Network (CNN), and Multi-Layer Perceptron (MLP), we made sure to implement a seed for reproducible results. For kNN, we tuned the hyper-parameter, k, to identify the best number of neighbors to predict. We did this by using 5-fold cross-validation, where we varied the value of k and returned the k value with the highest accuracy on the test data. We applied a similar approach to our Logistic classifier, using 5-fold cross-validation to get the best regularization strength, C, which helps reduce over-fitting. For Random Forest, we employed a brute-force technique, testing various values for the hyper-parameters, "n_estimators," "min_samples_split," and "min_samples_leaf." Optimizing these would give us the best values for the number of trees, the minimum number of data points to split on, and the minimum number of data points in a leaf node. In the case of CNN, we experimented with different epoch values to determine the best number of complete cycles through the training dataset. Lastly, for MLP, we varied batch sizes to find the most optimal number of training examples used in an iteration. We also experimented with different solvers and maximum iterations to help the model better converge.

For the kNN classifier, we found an optimal value of $k=10$, while for the logistic classifier, we set $C=0.01$ with a maximum iteration of 5000. In the case of Random Forest, we configured the number of estimators to 100, with a minimum sample split of 5 and a minimum sample leaf of 2. For CNN, we found the number of epochs of 20 to be the most optimal. Finally, for the MLP, we opted for a batch size of 128, the solver 'adam' and a maximum iteration of 1000.

V. EXPERIMENTAL RESULTS

These results were obtained by running the fine-tuned models with a test set of 2,000 images.

TABLE I
ACCURACY

Classifier	Accuracy
kNN	28.15%
Logistic	35.10%
Random Forest	42.70%
MLP	40.55%
CNN	54.75%

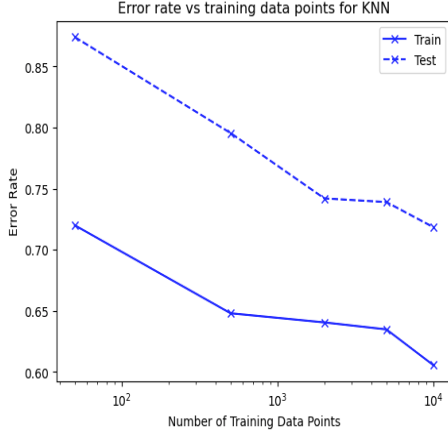


Fig. 3. Learning curve for kNN.

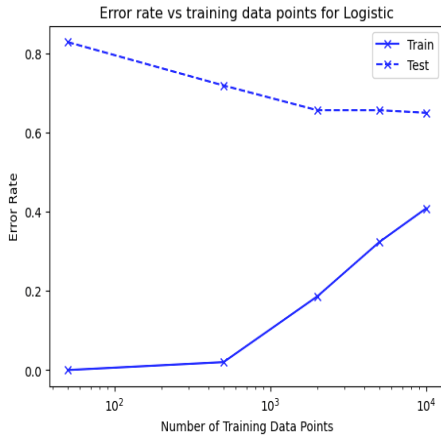


Fig. 4. Learning curve for Logistic.



Fig. 5. Learning curve for RandomForest.



Fig. 6. Learning curve for MLP.

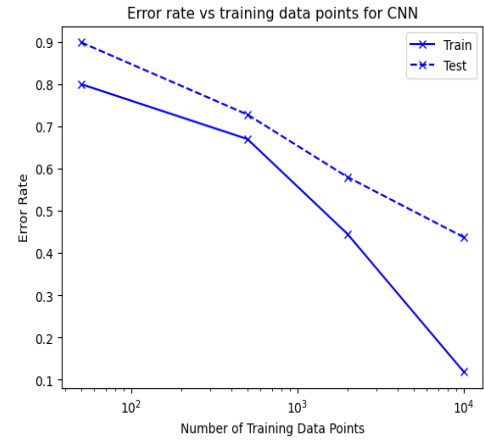


Fig. 7. Learning curve for CNN.

The testing Error rate decreases for all models as the training dataset grows. Among all, CNN showed the steepest decrease in testing error rate. CNN is the only model starting with a high error rate on training data, yet it decreases drastically as

the data size grows. The logistic classifier is the only model whose error rate increases as sample size increases, likely an effect of regularization (which other models do not implement)

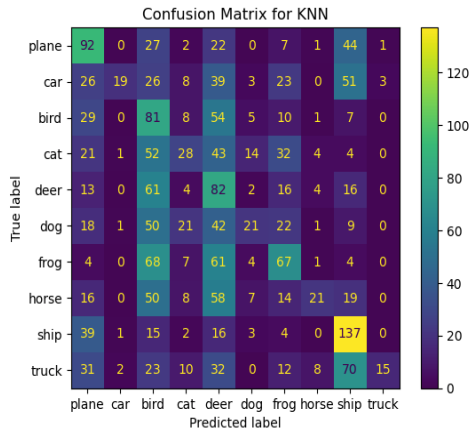


Fig. 8. Confusion matrix for KNN.

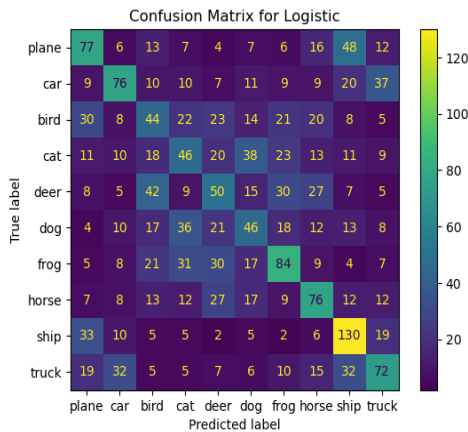


Fig. 9. Confusion matrix for Logistic.

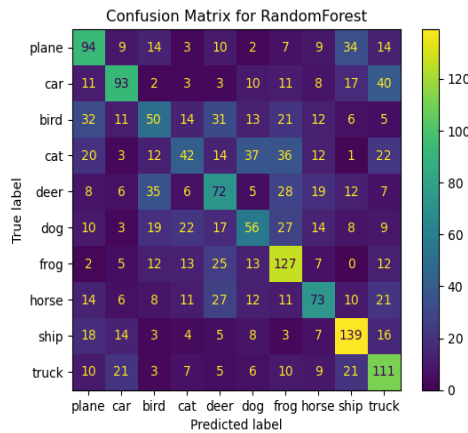


Fig. 10. Confusion matrix for Random Forest.

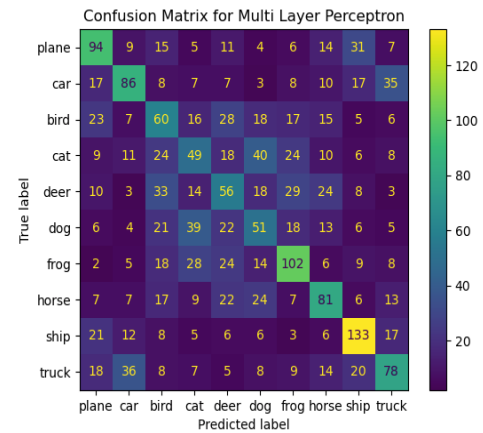


Fig. 11. Example of a figure caption.

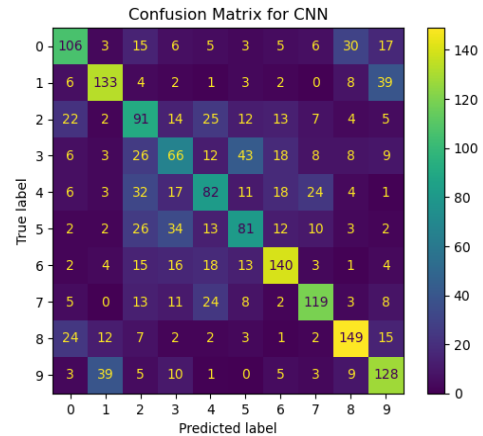


Fig. 12. Example of a figure caption.

KNN performs poorly, especially when misclassifying other animals as birds or deer, which doesn't happen for different classifiers. Logistic, RandomForest, and MLP have similar recognition patterns; they all seem to be confused among animals (the lighter blue square area in the middle covers animal classes: bird, cat, deer...horse). Most models seem to correctly classify non-animal classes at a higher rate.

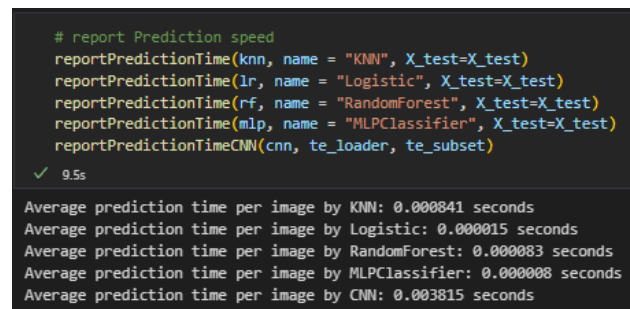


Fig. 13. Predictions speed of classifiers

VI. INSIGHT

- 1) Classifiers have a hard time telling animals apart !
 - Most classifiers have a hard time classifying between animal classes while having a much lower mis-identification rate between non-animal(Vehicle) and telling non-animal apart from animal. Hypothesis being that vehicle's picture tends to have a distinctive background and color while animal pictures are taken with a natural background(green for lawn or forest, gray, blue for sky), and their color is limited (you are unlikely to see a blue cat). However, vehicles, being artificial and their picture taken in many more settings, can have strong contrast colors (paint) and vastly different backgrounds (outdoor vs indoor).
- 2) Strengths and weaknesses of different models
 - CNN took more resources to train and infer, yet it provides the best accuracy. Random Forest and MLP offer a good balance as they provide good accuracy while taking up fewer resources to train and less time to predict. Logistic outperformed kNN in accuracy and prediction time; kNN is likely not an optimal model for this dataset. This could be because the features are too similar, causing more mis-classifications. In a real-world setting, it may be best to use MLP since it provides satisfactory accuracy and is the quickest predictor.
- 3) Discussion
 - Through confusion matrices, learning curves, and measurement of prediction speed, the models' strengths and weaknesses were exposed. This data helped us better understand and learn the intricacies of each model and raise awareness for choosing the adequate model and making trades for real-world scenarios like autonomous vehicles, security and surveillance, and medical imaging analysis, three critical and growing fields in artificial intelligence.

VII. CONTRIBUTIONS

- **Jordan:** Implemented the k-nearest Neighbors and logistic regression classifiers. He was responsible for tuning the hyper-parameters, evaluating the performance of these models, and visualizing the results.
- **Travis:** Designed the neural network structure for the Multi-Layer Perceptron and Convolutional Neural Network. He also developed the architectures for tuning the hyper-parameters and evaluated performance.
- **Shifeng:** Translated the report into Jupyter Notebook to generate graphs and various plots used throughout the sections. He also implemented the random forest classifier and visualized the results in different hyper-parameter settings.