

Finn Thompson

Lakshmanan Sockalingam

Siddhant Rajoriya

Aryan Roy

CMPSC 448

02 December 2023

Mushroom Poisonousness Classification Report

Mushroom classification is a critical task, especially when distinguishing between edible and poisonous varieties. In this project, we employ deep learning techniques, specifically Convolutional Neural Networks (CNN) and Recurrent Neural Networks (RNN), to address this classification challenge. The dataset comprises various attributes of mushrooms, such as cap shape, color, and odor. The goal is to develop models that can accurately predict whether a mushroom is safe to consume or potentially harmful.

The preprocessing phase involves transforming categorical features into numerical representations through label encoding. The dataset is then split into training and testing sets to facilitate model training and evaluation. The subsequent sections delve into the detailed implementation, architectures, training procedures, and results of both CNN and RNN models for mushroom classification. Additionally, the report highlights observed challenges, the

strategies employed to overcome them, and concludes with insights into the overall effectiveness of the deep learning models in this context.

Task

The task at hand is to classify whether a mushroom is poisonous or not based on various features. This is a binary classification problem, where the model needs to predict whether the mushroom belongs to the 'poisonous' or 'edible' class. Classifying whether mushrooms are poisonous or not is often overlooked and can lead to severe health issues and in some cases death when consumed. This application has a strong commercial application especially in mushroom farms where mushrooms are grown for harvesting and cultivation purposes. Therefore this project will aim to be a step towards enhancing the quality control procedure to validating that only edible and non-poisonous mushrooms are sold for consumption.

Dataset

The dataset used for this task is sourced from 'mushrooms.csv.' It includes various attributes such as cap shape, cap color, odor, etc., with the target variable being the class of the mushroom ('poisonous' or 'edible'). These attributes are based on not only the obvious physical attributes such as cap shape and cap color, but also include fine-margin attributes such as the type of stalk surface below the ring and print color of the spores. This gives the dataset a variety of features to include while keeping high generalization.

Preprocessing

To prepare our data in a format suitable for Convolutional and Recurrent Neural Networks, it goes through a few essential preprocessing tasks. First, label encoding is applied to the target variable, 'class', as well as categorical features, converting them into numerical values that are suitable for input into the deep learning models. Next, we use a stratified approach to split the dataset into a training and test set. This maintains a balanced class distribution in each set, and is fundamental for assessing each of the models' performance on new data. Without splitting the data, it would be unclear how the models performed when confronted with new data. Using the test set, we can tune hyperparameters to optimize accuracy for our outputs. Each of these preprocessing stages are critical for transforming the dataset into a format compatible with the neural networks we chose regardless of the specific architecture (CNN or RNN).

Convolutional Neural Network (CNN) Implementation

The CNN model we designed is implemented using the Keras framework. This architecture is identified by a few key components. Firstly, the model consists of an embedding layer. This is responsible for learning a more dense representation of each of the input features. Next, the architecture consists of three core convolutional layers which are each followed by a max-pooling layer. These layers are designed to capture the hierarchical features in the data. The max-pooling layers serve to shrink the size of the information, concentrating on the most important parts. Next a global max-pooling layer is incorporated, which enables our model to capture the most significant feature from the whole sequence of data.

Finally, the architecture concludes with fully connected layers. These layers use a dropout technique in an effort to prevent overfitting and activation functions to make the model more flexible in understanding complex patterns in the data. Though this design is complex, each piece contributes to our CNN's ability to learn and represent the patterns in the Mushroom dataset.

Recurrent Neural Network (RNN) Implementation

Like the CNN model, our RNN model is implemented using the Keras framework. It begins with an embedding layer, similar to the one employed in the CNN model. However, for an RNN model it can capture the relationships and dependencies between sequential elements, something that RNN models excel at. Following that is the simple recurrent layer, which similarly recognizes and understands these dependencies in the data. Next is the dropout layer, used to enhance our model's robustness by preventing overfitting, which is a common challenge in training neural networks.

Finally the architecture concludes with a dense layer. This layer is equipped with a sigmoid activation function. This activation function is a mathematical operation that is applied to each of the outputs in the dense layer. It will squash the output into a value between 0 and 1, which in our case can be translated to whether the mushroom is poisonous or edible. Without this final layer, we cannot perform the transition from our many components to a prediction of yes or no. With that layer implemented it is only a matter of training and testing the RNN on our data to enable the model to start making predictions. The RNN model has many similarities to our CNN model, but where it excels is in dealing with sequential data. Because our data is tabular we do not expect this to dramatically effect the difference in our models' accuracy

CNN Training

The CNN model is trained over 10 epochs with a binary cross-entropy loss function and the Adadelata optimizer. The training accuracy steadily increases, reaching 91.78%, and the validation accuracy peaks at 93.60%. The model effectively learns to differentiate between poisonous and edible mushrooms.

CNN models are trained using standard backpropagation. It is used in feedforward neural networks which are unidirectional. Standard backpropagation works by iteratively propagating the error signal back through the network, beginning from the output layer and progressing back to the input layer. In every layer, the direction which reduces the error signal determines the direction in which the connection weights are updated.

The training procedure chronologically occurs by presenting the CNN model with batches of mushroom data which comprises mushroom description and their classification. The model extracts the features and transforms them into numerical representation.

RNN Training

The RNN model is trained for 10 epochs as well, using binary cross-entropy loss and the Adadelata optimizer. The training accuracy improves over epochs, reaching 90.25%, while the validation accuracy achieves 90.41%. The RNN effectively captures sequential information for classification.

RNN models are trained using backpropagation through time (BPTT). BPTT, similar to standard backpropagation, in addition takes into account the feedback loops (recurrent connections).

BPTT is based on unrolling the RNN over a period of time. The error signal is then propagated back through the unrolled network, beginning from the output layer and progressing back to the input layer. In every layer, the direction which reduces the error signal determines the direction in which the connection weights are updated.

The RNN training process involves feeding the model batches of sequences and their corresponding labels. The weights are updated by the model to minimize the binary cross-entropy loss between the actual labels and the prediction probabilities. It is observed over the course of 10 Epochs that the model progressively becomes better at classifying whether a mushroom is poisonous or not which is reflected by the increase in the model accuracy.

Results

CNN:

Accuracy: 93.60%

Precision: 96.07%

Recall: 90.00%

F1 Score: 92.93%

Macro F1 Score: 93.55%

AUC: 93.38%

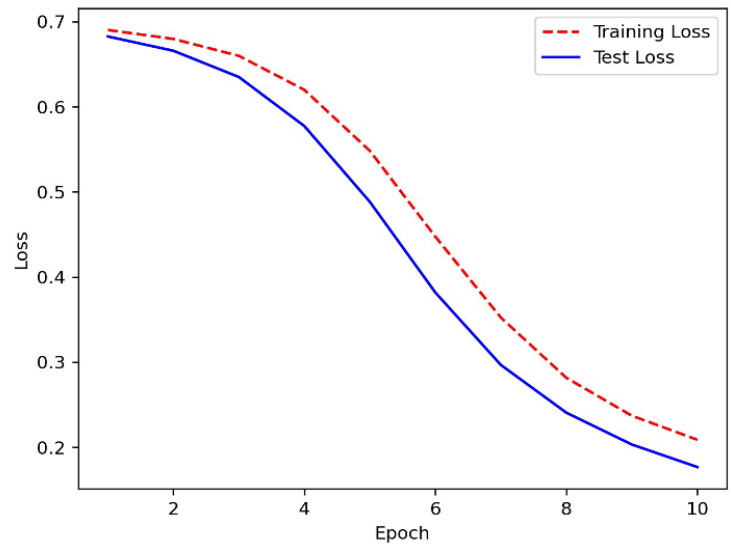


Fig 1. CNN Loss Curves

RNN:

Accuracy: 90.41%

Precision: 96.04%

Recall: 82.89%

F1 Score: 88.98%

Macro F1 Score: 90.24%

AUC: 89.95%

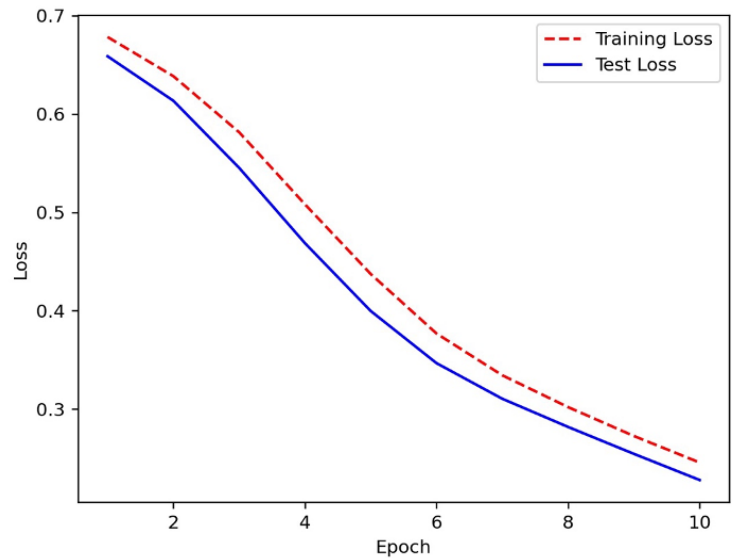


Fig 2. RNN Loss Curves

Both models show excellent performance, with the CNN slightly outperforming the RNN. While precision remains similar for both models, CNN performs better with a difference of ~3% in

accuracy, AUC, and F1 scores. The CNN performs much higher in recall than RNN too. The difference can be representative of CNN's better ability to capture spatial features in the data.

Challenges and Obstacles

Hyperparameter Tuning and Model Development: Tuning several hyperparameters at the same time to get the best accuracy was difficult. In addition, choosing the sizes of the model layers was difficult without going through a lot of training.

- Initially, we tried to solve the issue by conducting a grid search optimization of hyperparameters and model layer numbers. However, we realized that it took too much time. Thus, we decided to do random search optimization to solve hyperparameter tuning issues.

Data Preprocessing: With data in both categorical and numerical values, converting different categorical values columns into numerically representative data and ensuring compatibility with deep learning models was a challenge.

- We overcame the challenge by utilizing label encoding for categorical features and embedding layers for test-based features.

Model Complexity: Choosing appropriate model architectures for the given task was the challenge with several types of basic architectures available in both CNN and RNN models. W

- We experimented with several models and decided on the usage of simpler architectures based on comparatively higher output accuracies.

Maintenance of accuracy: Our training accuracy varied over several training instances. While the variance was very high, it was a challenge to keep it to a minimum. In addition, we had high training and low testing accuracies in the initial tests - showcasing the issue of overfitting, especially in the RNN.

- We solved this problem by adding dropout layers and updating the hyperparameters and the model layer sizes to accommodate the regularization.

Conclusion

The implementation of both CNN and RNN models for mushroom classification was successful. We were able to achieve high accuracies in the early 90% range for both the models. The precision output for both models was especially high at ~96%. The CNN, with its ability to capture spatial patterns, appeared slightly more effective for this task with an AUC difference of around 3%. The choice between these models in real world scenarios could depend on specific requirements, computational resources, and the nature of the dataset. The project allowed us to learn through experience while facing and solving challenges on our own. Overall, the project demonstrated the versatility of deep learning in handling diverse data types for classification tasks.

Works Cited

UCI ML (n.d.), *Mushroom Classification - Safe to eat or deadly poison*,

<https://www.kaggle.com/datasets/uciml/mushroom-classification/data>

Scott, J. (n.d.), *Mushroom Sequence Classification with CNN*,

<https://www.kaggle.com/code/engjay/mushroom-sequence-classification-with-cnn>