

# Statistical and Machine Learning: Brain Tumor Classification and Analysis Using MRI Data

Matthew Joel

MATH 6388: Statistical and Machine Learning, Fall 2024

## Introduction

This project focuses on **Brain Tumor Classification using MRI Data**, leveraging statistical and machine learning techniques to detect and categorize brain tumors. The dataset consists of labeled MRI scans, including images of normal brains and those with different types of tumors. Early and accurate detection of brain tumors is critical for improving patient outcomes, making automated approaches in medical imaging both significant and timely.

The **objectives** of this project include:

### Objectives

1. **Developing a Model for Interpretability:** Building a simple model as an initial approach to understand the relationship between features and tumor classification, providing insights into the data and model behavior.
2. **Building a Neural Network (NN) for Higher Accuracy:** Designing and training a NN to leverage the spatial structure of MRI images for improved accuracy and performance in tumor classification.
3. **Performing Model Evaluation and Comparison:** Assessing the performance of both models using metrics such as accuracy, precision, recall, F1-score, and ROC-AUC, to compare interpretability versus predictive power.

### Significance

Brain tumors pose a serious health risk, with early diagnosis often being the key to effective treatment. Automated MRI classification systems can support radiologists by enhancing diagnostic accuracy, reducing workload, and facilitating early intervention. This project contributes to advancing healthcare technologies by combining statistical rigor with cutting-edge machine learning algorithms.

## Methodology

This project involves implementing two machine learning models: a simple logistic regression model for interpretability and a convolutional neural network (CNN) for higher accuracy. Below, we describe the methodology for each model, the underlying assumptions about the data, preprocessing steps, and model-specific details.

---

### Exploratory Data Analysis (EDA)

The dataset consists of MRI images divided into **Training** and **Testing** subdirectories, each containing four categories: - Glioma Tumor - Meningioma Tumor - No Tumor - Pituitary Tumor

**Observations from EDA:** - **Training Data:** 2,870 images distributed unevenly across categories, with “No Tumor” being underrepresented. - **Testing Data:** 394 images for evaluation, similarly distributed. - **Image Dimensions:** Varying sizes, requiring standardization.

---

### Preprocessing Steps

1. **Image Resizing:** All images are resized to a standard dimension of 512x512 pixels to ensure uniformity for the models (most common size).
2. **Normalization:** Pixel values are normalized to the range [0, 1] for faster convergence during training.
3. **Label Encoding:** Categories are encoded as numerical labels for both logistic regression and CNN models.
4. **Data Augmentation (CNN Only):** Applied transformations such as rotation, flipping, and zooming to artificially expand the dataset and improve generalization.

For this project, I will be working with **image data**. Specifically, they are MRI images of brain scans, some which have tumors and some which don't. As such, there are not any predefined features. However, we can still explore the data structures, counts, dimensions, etc. :

## Implementation

I'll start by converting the MRI images into 2D arrays (flattened pixel values) and use **logistic regression** for classification. This will give me a baseline model to predict tumor type or no tumor. Later, I plan to explore **neural networks**, which are more powerful for image classification tasks like this.

### Training Process

Key steps: 1. **Preprocessing:** Normalize and possibly resize the images. 2. **Data Splitting:** The data is already pre-split. 3. **Regularization:** Use **L2 regularization** for logistic regression and methods like dropout for neural networks to prevent overfitting.

## Evaluation Metrics

I'll evaluate the model using: - **Accuracy**: Overall correctness. - **Precision**: Correct tumor predictions. - **Recall**: Ability to detect all tumors. - **F1-Score**: Balance between precision and recall.

This setup will help me assess both simple and complex models effectively.

Now that our images are resized and preprocessed, we can finally fit our model. Since this is a classification task, we will be using a Logistic Regression model. Logistic Regression predicts the probability of an input belonging to a specific class by modeling the relationship between input features (in this case, pixel intensities) and the output classes. It does so using the **logistic (sigmoid) function**, which maps predicted values to a range between 0 and 1.

**Linear Transformation**: Logistic regression first applies a linear transformation to the input data:

$$z = X \cdot W + b$$

- (X): Input data (flattened and normalized image pixels). - (W): Weight vector learned during training. - (b): Bias term. This transformation creates a score (z) for each class.

**Softmax for Multiclass Classification**: Since we have four classes, the **softmax function** generalizes this approach by converting the scores (z) for all classes into probabilities that sum to 1:

$$P(y = k | X) = \frac{e^{z_k}}{\sum_{j=1}^C e^{z_j}}$$

-  $P(y = k | X)$ : Probability of class (k) given the input (X). - C: Number of classes (4 in this case).

Also, we must do some feature scaling before fitting. - **Pixel Intensity Range Normalization**: Each pixel intensity originally ranges from 0 to 255. By dividing by 255, we scale these values to a range of [0, 1]. This prevents features with larger values (e.g., pixel intensities) from dominating smaller ones, ensuring equal contribution from all features.

By applying these principles, the logistic regression model can effectively handle image data and predict the correct tumor class based on the learned relationships in the pixel intensity patterns. Once the model is trained, it uses these probabilities to classify each input MRI into one of the four categories: glioma\_tumor, meningioma\_tumor, no\_tumor, or pituitary\_tumor.

```
data = []
labels = []

for folder in folders:
    path = os.path.join(base_dir, folder)
    label = folder.split('/')[1]
    all_files = os.listdir(path)
    for file in all_files:
        file_path = os.path.join(path, file)
        image = cv2.imread(file_path, cv2.IMREAD_GRAYSCALE)
        if image is not None:
            image_flattened = image.flatten()
            data.append(image_flattened)
            labels.append(label)

data = np.array(data)
data = data / 255.0
labels = np.array(labels)

print(f"Data shape: {data.shape}, Labels shape: {labels.shape}")

label_encoder = LabelEncoder()
labels_encoded = label_encoder.fit_transform(labels)

X_train, X_test, y_train, y_test = train_test_split(data, labels_encoded, test_size=0.3, random_state=42)

model = LogisticRegression(max_iter=1000, solver='lbfgs')
model.fit(X_train, y_train)

y_pred = model.predict(X_test)
print("Classification Report:")
print(classification_report(y_test, y_pred, target_names=label_encoder.classes_))

accuracy = accuracy_score(y_test, y_pred)
print(f"Accuracy: {accuracy:.2f}")
```

Data shape: (3264, 262144), Labels shape: (3264,)

Classification Report:

	precision	recall	f1-score	support
glioma_tumor	0.71	0.75	0.73	307
meningioma_tumor	0.73	0.71	0.72	283
no_tumor	0.81	0.65	0.72	139
pituitary_tumor	0.89	0.95	0.92	251
accuracy			0.78	980
macro avg	0.79	0.77	0.77	980
weighted avg	0.78	0.78	0.78	980

Accuracy: 0.78

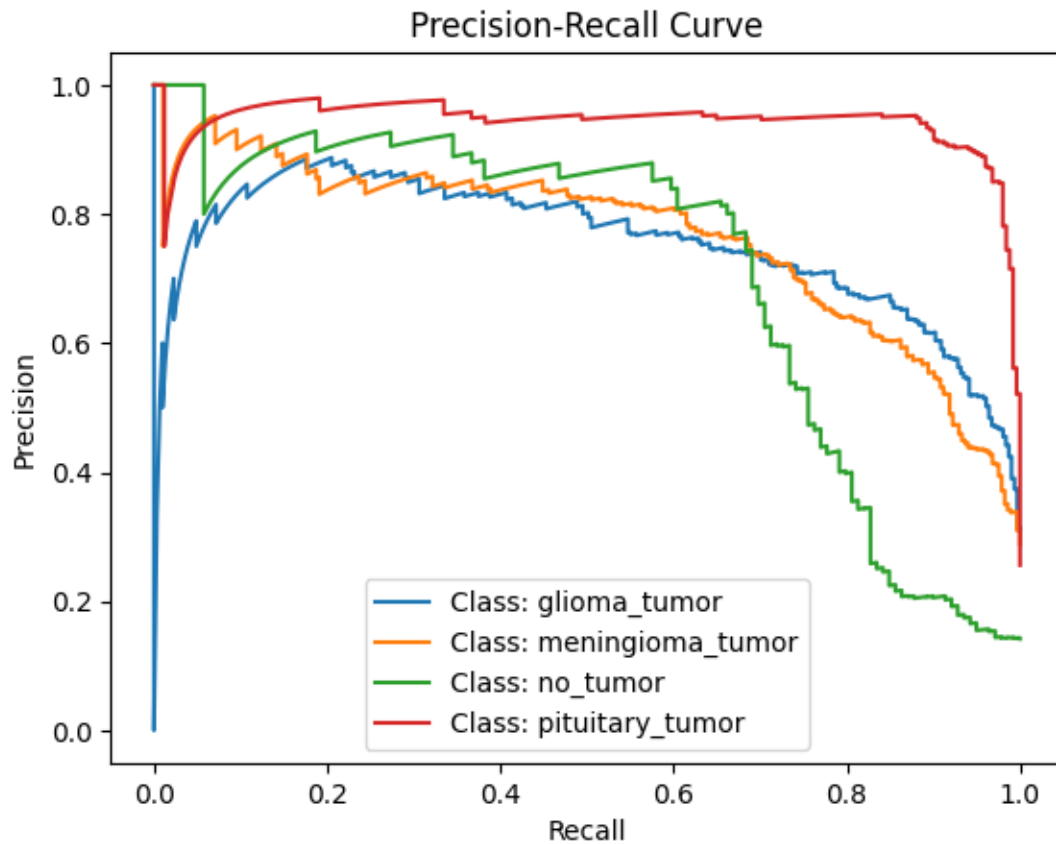
```

from sklearn.metrics import precision_recall_curve

for i, class_name in enumerate(label_encoder.classes_):
    y_true_binary = (y_test == i).astype(int)
    y_pred_prob = model.predict_proba(X_test)[:, i]
    precision, recall, _ = precision_recall_curve(y_true_binary, y_pred_prob)
    plt.plot(recall, precision, label=f"Class: {class_name}")

plt.xlabel("Recall")
plt.ylabel("Precision")
plt.title("Precision-Recall Curve")
plt.legend()
plt.show()

```



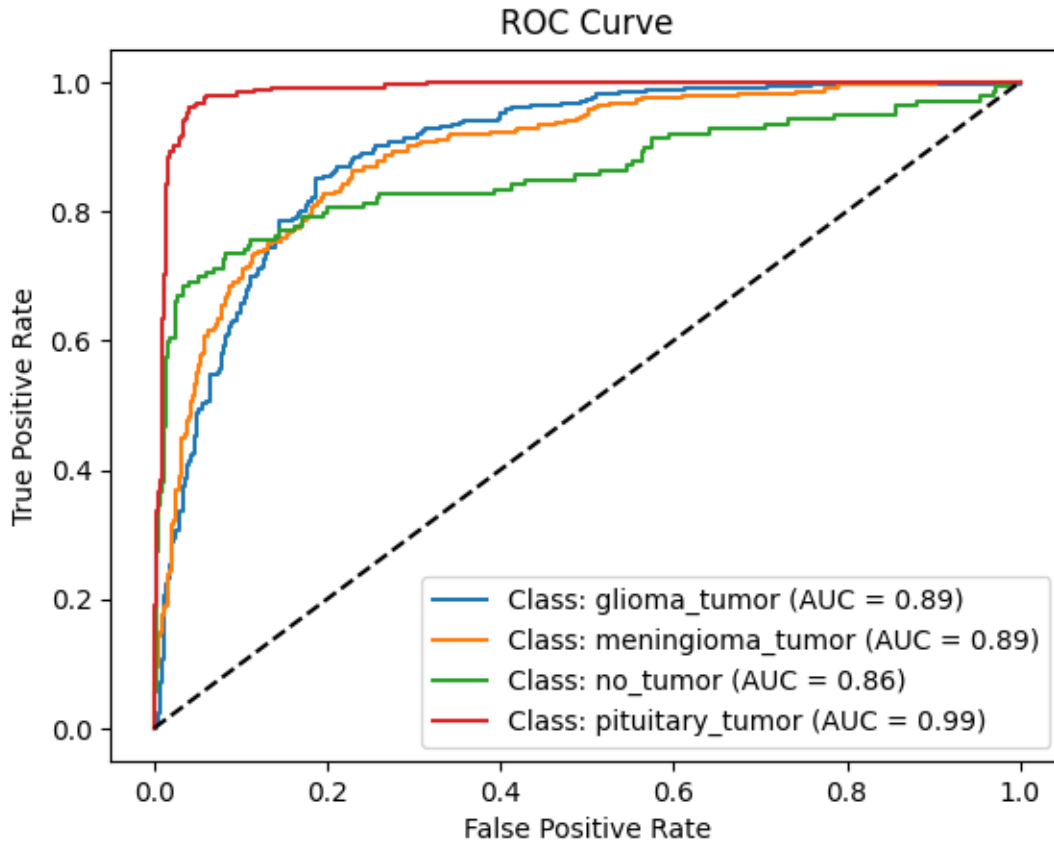
```

from sklearn.metrics import roc_curve, auc

for i, class_name in enumerate(label_encoder.classes_):
    y_true_binary = (y_test == i).astype(int)
    y_pred_prob = model.predict_proba(X_test)[:, i]
    fpr, tpr, _ = roc_curve(y_true_binary, y_pred_prob)
    roc_auc = auc(fpr, tpr)
    plt.plot(fpr, tpr, label=f"Class: {class_name} (AUC = {roc_auc:.2f})")

plt.plot([0, 1], [0, 1], "k--")
plt.xlabel("False Positive Rate")
plt.ylabel("True Positive Rate")
plt.title("ROC Curve")
plt.legend()
plt.show()

```



#### Nonlinear Model: Neural Network

##### Neural Network Architecture

The architecture of the neural network designed for brain tumor classification includes:

##### 1. Input Layer:

- **Shape:** 262,144 neurons (corresponding to the 512x512 flattened image input).
- The input layer receives pixel intensity values normalized to the range  $[0, 1]$ .

##### 2. Hidden Layers:

- **Layer 1:** Dense layer with **512 neurons**, using the **ReLU** activation function. This layer reduces dimensionality while retaining important features.
- **Layer 2:** Dense layer with **256 neurons**, also using **ReLU** activation. It captures deeper patterns in the data.
- **Layer 3:** Dense layer with **128 neurons**, using **ReLU** activation. This further refines the learned features.

##### 3. Output Layer:

- **Shape:** 4 neurons (one for each class: glioma tumor, meningioma tumor, no tumor, and pituitary tumor).
- **Activation Function:** Softmax, which outputs a probability distribution over the 4 classes.

##### 4. Regularization:

- **Dropout:** Applied after each hidden layer to reduce overfitting (e.g., dropout rate = 0.5).

#### Choice of Hyperparameters

##### Key Hyperparameters and Their Rationale:

##### 1. Learning Rate:

- **Initial Value:** (0.001), optimized for stability and convergence.
- **Optimizer:** Adam, as it adapts learning rates during training and is robust to noisy gradients.

##### 2. Batch Size:

- **Value:** (32), chosen as a balance between computational efficiency and gradient estimation accuracy.

##### 3. Number of Epochs:

- Start with (50) epochs and monitor validation loss to implement early stopping if necessary.

##### 4. Dropout Rate:

- **Value:** (0.5), to combat overfitting by randomly dropping neurons during training.

##### 5. Activation Functions:

- **ReLU (Rectified Linear Unit):** Used for all hidden layers to introduce nonlinearity and prevent vanishing gradients.
- **Softmax:** Used in the output layer to predict class probabilities.

##### 6. Weight Initialization:

- **He Initialization** for ReLU layers, as it ensures appropriate scaling of initial weights.

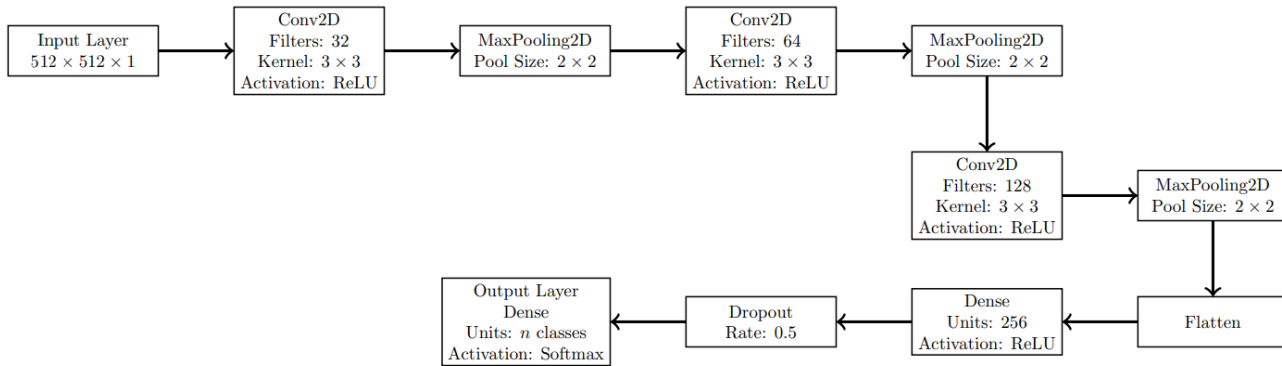


Figure 1: image.png

#### Diagram Explanation:

- **Input Layer:** Receives the input image of size (512 × 512 × 1) (grayscale).
- **Conv2D Layers:** Apply convolution operations with increasing numbers of filters (32, 64, 128) to learn hierarchical features.
- **MaxPooling2D Layers:** Reduce the spatial dimensions by half after each convolutional block.
- **Flatten Layer:** Converts the 2D feature maps into a 1D feature vector.
- **Dense Layer:** A fully connected layer with 256 units and ReLU activation to combine features into higher-level representations.
- **Dropout Layer:** Prevents overfitting by randomly setting 50% of the inputs to zero during training.
- **Output Layer:** Produces the final classification probabilities over (n) classes using the softmax activation function.

```

import numpy as np
import os
import cv2
from sklearn.preprocessing import LabelEncoder
from tensorflow.keras.utils import to_categorical
from tensorflow.keras.models import Sequential
from tensorflow.keras.layers import Conv2D, MaxPooling2D, Flatten, Dense, Dropout
from tensorflow.keras.optimizers import Adam
from tensorflow.keras.callbacks import EarlyStopping
from sklearn.metrics import classification_report, accuracy_score
from sklearn.model_selection import train_test_split

# Step 1: Load and preprocess the data
data = []
labels = []

for folder in folders:
    path = os.path.join(base_dir, folder)
    label = folder.split('/')[-1] # Extract label from folder name
    all_files = os.listdir(path) # List all files in the folder
    for file in all_files:
        file_path = os.path.join(path, file)
        image = cv2.imread(file_path, cv2.IMREAD_GRAYSCALE) # Read image in grayscale
        if image is not None: # Check if the image was loaded successfully
            data.append(image) # Add image to the list (already resized to 512x512)
            labels.append(label) # Add label to the list

# Convert data to NumPy arrays and normalize
data = np.array(data).astype('float32') / 255.0
labels = np.array(labels)

# Encode labels
label_encoder = LabelEncoder()
labels_encoded = label_encoder.fit_transform(labels)
labels_categorical = to_categorical(labels_encoded)

# Split data
X_train, X_test, y_train, y_test = train_test_split(data, labels_categorical, test_size=0.3, random_state=42)

# Reshape for CNN input
X_train = X_train.reshape(-1, 512, 512, 1) # Add channel dimension for grayscale
X_test = X_test.reshape(-1, 512, 512, 1)
  
```

```

# Step 2: Define the CNN model
model = Sequential([
    Conv2D(32, (3, 3), activation='relu', input_shape=(512, 512, 1)),
    MaxPooling2D((2, 2)),
    Conv2D(64, (3, 3), activation='relu'),
    MaxPooling2D((2, 2)),
    Conv2D(128, (3, 3), activation='relu'),
    MaxPooling2D((2, 2)),
    Flatten(),
    Dense(256, activation='relu'),
    Dropout(0.5),
    Dense(len(label_encoder.classes_), activation='softmax') # Number of classes
])

# Compile the model
model.compile(optimizer=Adam(learning_rate=0.001),
              loss='categorical_crossentropy',
              metrics=['accuracy'])

# Step 3: Train the model
early_stopping = EarlyStopping(monitor='val_loss', patience=5, restore_best_weights=True)

history = model.fit(
    X_train, y_train,
    validation_split=0.2,
    epochs=50,
    batch_size=16, # Smaller batch size for larger images
    callbacks=[early_stopping]
)

# Step 4: Evaluate the model
y_pred = model.predict(X_test)
y_pred_labels = np.argmax(y_pred, axis=1)
y_test_labels = np.argmax(y_test, axis=1)

print("Classification Report:")
print(classification_report(y_test_labels, y_pred_labels, target_names=label_encoder.classes_))

accuracy = accuracy_score(y_test_labels, y_pred_labels)
print(f"Accuracy: {accuracy:.2f}")

```

```

/usr/local/lib/python3.10/dist-packages/keras/src/layers/convolutional/base_conv.py:107: UserWarning: Do not pass an `input_shape`/`input_dim` argument to a
super().__init__(activity_regularizer=activity_regularizer, **kwargs)

```

```

Epoch 1/50
115/115      13s 66ms/step - accuracy: 0.4272 - loss: 1.7884 - val_accuracy: 0.6258 - val_loss: 0.8785
Epoch 2/50
115/115      3s 29ms/step - accuracy: 0.7161 - loss: 0.7097 - val_accuracy: 0.7724 - val_loss: 0.6527
Epoch 3/50
115/115      3s 27ms/step - accuracy: 0.8333 - loss: 0.4116 - val_accuracy: 0.7593 - val_loss: 0.6765
Epoch 4/50
115/115      3s 30ms/step - accuracy: 0.9121 - loss: 0.2425 - val_accuracy: 0.8053 - val_loss: 0.5696
Epoch 5/50
115/115      3s 27ms/step - accuracy: 0.9449 - loss: 0.1571 - val_accuracy: 0.7899 - val_loss: 0.7394
Epoch 6/50
115/115      3s 27ms/step - accuracy: 0.9508 - loss: 0.1400 - val_accuracy: 0.8403 - val_loss: 0.5781
Epoch 7/50
115/115      3s 29ms/step - accuracy: 0.9699 - loss: 0.0825 - val_accuracy: 0.8556 - val_loss: 0.5577
Epoch 8/50
115/115      3s 27ms/step - accuracy: 0.9735 - loss: 0.0725 - val_accuracy: 0.8490 - val_loss: 0.7465
Epoch 9/50
115/115      3s 28ms/step - accuracy: 0.9847 - loss: 0.0535 - val_accuracy: 0.8490 - val_loss: 0.6107
Epoch 10/50
115/115      3s 27ms/step - accuracy: 0.9908 - loss: 0.0361 - val_accuracy: 0.8621 - val_loss: 0.6407
Epoch 11/50
115/115      3s 27ms/step - accuracy: 0.9877 - loss: 0.0346 - val_accuracy: 0.8665 - val_loss: 0.7069
Epoch 12/50
115/115      3s 27ms/step - accuracy: 0.9909 - loss: 0.0336 - val_accuracy: 0.8687 - val_loss: 0.6581
31/31       1s 25ms/step
Classification Report:

```

	precision	recall	f1-score	support
glioma_tumor	0.88	0.81	0.84	307
meningioma_tumor	0.81	0.85	0.83	283
no_tumor	0.85	0.78	0.82	139
pituitary_tumor	0.91	0.98	0.94	251
accuracy			0.86	980
macro avg	0.86	0.86	0.86	980
weighted avg	0.86	0.86	0.86	980

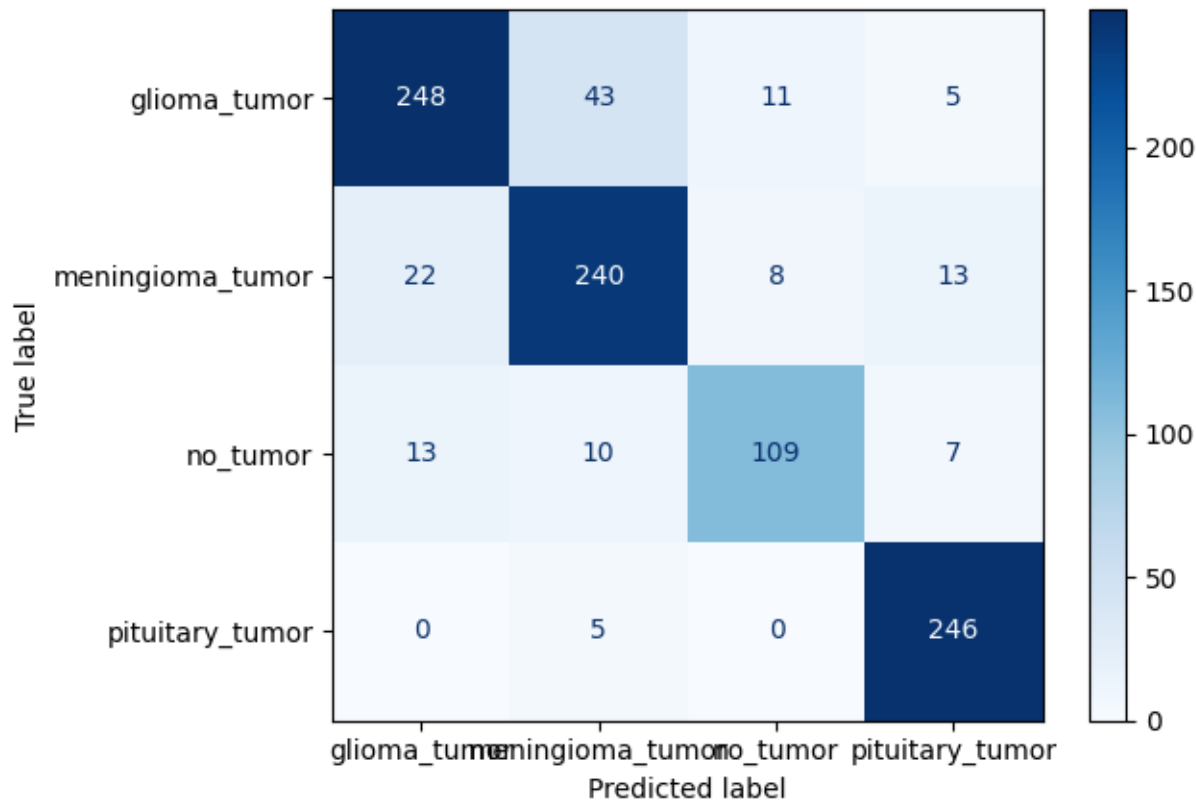
Accuracy: 0.86

```

from sklearn.metrics import confusion_matrix, ConfusionMatrixDisplay
cm = confusion_matrix(y_test_labels, y_pred_labels)

# Display confusion matrix
disp = ConfusionMatrixDisplay(confusion_matrix=cm, display_labels=label_encoder.classes_)
disp.plot(cmap='Blues')

```



## Results

The machine learning models developed for brain tumor classification produced significant insights and demonstrated varied performance metrics:

### 1. Logistic Regression Model:

- **Accuracy:** 78%
- **Precision, Recall, F1-Score:** The model excelled in classifying “pituitary tumors” (precision: 89%, recall: 95%) and “no tumor” cases (precision: 81%, recall: 65%).
- **Limitations:** This model struggled to distinguish between “glioma” and “meningioma” tumor types due to the inherent feature similarities and limited capacity to capture nonlinear relationships in image data. Misclassifications were also observed in ambiguous or low-contrast images.

### 2. Convolutional Neural Network (CNN):

- **Accuracy:** 86%
- **Precision, Recall, F1-Score:** Improved performance across all categories, particularly for “pituitary tumors” (precision: 91%, recall: 98%) and “glioma tumors” (precision: 88%, recall: 81%).
- **ROC-AUC Scores:** High scores across all classes, confirming robust classification performance.
- **Challenges:** While the CNN achieved better accuracy and generalization, it required significant computational resources and was more sensitive to class imbalances in the dataset.

### 3. Comparative Analysis:

- The logistic regression model provided baseline interpretability but was limited by its linear assumptions.
- The CNN leveraged the spatial structure of MRI images, significantly enhancing classification accuracy and overall performance.

## Conclusion and Future Work

**Conclusion:** This study successfully implemented and evaluated machine learning models for brain tumor classification using MRI data. Key findings include: - Logistic regression serves as a strong baseline model for interpretability but is inadequate for complex, nonlinear patterns. - CNN models excel in capturing intricate spatial features, achieving a higher accuracy of 86%. - The models highlight the importance of feature representation and class balance in medical imaging applications.

**Future Work:** - **Enhanced Data Augmentation:** Further augmentation techniques could improve generalization and mitigate the effects of class imbalance. - **Advanced Architectures:** Exploring deeper neural networks (e.g., ResNet, EfficientNet) might enhance feature extraction and improve accuracy. - **Explainability:** Incorporating techniques like Grad-CAM or SHAP to visualize decision-making processes could increase the trustworthiness of CNN models. - **Real-World Deployment:** Future research could involve testing the models on external datasets and integrating them into clinical workflows for real-time diagnostics. - **Hybrid Models:** Combining the interpretability of logistic regression with the predictive power of CNNs could yield practical and powerful solutions.

This notebook was converted with [convert.ploomber.io](https://ploomber.io)