

# IMAGE SCRAPING AND CLASSIFICATION PROJECT



Submitted by:

PRIYANKA

## TABLE OF CONTENTS

1.	Acknowledgement
2.	Introduction
3.	Analytical problem framing
4.	Data Visualization
5.	Model Training and Testing
6.	Conclusion

# ACKNOWLEDGEMENT

I would like to express my gratitude to Data Trained for providing me with several learning opportunities. I would like to thank Flip Robo for giving me this platform to intern and learn in the vast field of Data Science. Our mentor Miss. Khushboo Garg has always been a support throughout the entire project, I would like to thank her for patiently resolving every query I landed up into. I have been learning through various online modes during this project. I would like to mention Mr. Krishh Naik, his YouTube videos have been a major learning platform to me.

# INTRODUCTION

## What is Image Classification?

When we human look at a picture we can easily identify what it depicts without any problem. Humans learn this skill of identification at early stages of life. This image identification problem though is not an easy task for a machine or a computer. They don't see the world the way humans do. Image classification thus becomes a challenge for computers or machines to learn, which is where the role of deep learning comes in.

Image classification is the ability of a computer to analyze an image and identify the class the image falls under. The class can be anything from automobiles to animals etc.

For instance, we input an image of a bus into the computer then the ability of computer to tell that the loaded image is a bus or the probability that the image is of a bus is what image classification is.

Image classification is a supervised learning problem.

## Background of the domain problem

Earlier image classification used raw pixel data which means that computer would break down the image into individual pixels. Problem with this approach was that two similar product images could look way distinct because of the different background, colours angles, poses and many other factors. This made it a task for computers to identify the images.

## Review of literature

As earlier method of classification of image was not reliable this is where deep learning entered. Deep learning is a subset of machine learning that allows machine to learn from data. Deep learning is known to use neural networks for this work.

In Neural Networks, the input gets filtered through hidden layers of nodes. Each node processes the input and passes their individual result to the next layer of node. This process repeats until it reaches an output layer and machine is able to answer.

There are different types of neural networks based on how hidden layers work. Image classification with deep learning is mostly done using CNN which is convolutional neural network. In CNN the nodes in the hidden layers don't always share their output with every node in the next layer known as convolutional layer.

## Motivation of the problem undertaken

Deep learning allows machine to identify and extract features from images. This means they can learn the features to look for in images by analyzing a lot of pictures. User does not need to enter the filters manually. Deep learning has made this problem of image classification quite easy.

## **ANALYTICAL PROBLEM FACING**

### Data sources and their formats

Data used in the project is images of saree, men jeans and men trouser which is scraped from e-commerce amazon. Scraping of data is done using selenium. After scraping the images in three different folders. The images were separated into training and testing folders each consisting three sub folders belonging to one of the three- Saree, Jeans and Trouser. These subfolders would serve as three categories in the project.

We will build a model that would take images as input and will classify them in one of the three categories- Saree, Jeans or Trouser.

### Visualization

To check whether the three folders of saree, jeans and trouser have the right images we visualize some images randomly using 'random'.

The code to load the images is given below:

```
#Randomly displaying the pictures from our file of saree from train folder
plt.figure(figsize=(20,20))
train_folder="C:/Users/Priyanka/Downloads/amazon_images/train/saree"
for i in range(5):
    file = random.choice(os.listdir(train_folder))
    image_path = os.path.join(train_folder, file)
    img=mpimg.imread(image_path)
    ax=plt.subplot(1,5,i+1)
    ax.title.set_text(file)
    plt.imshow(img)
```

The output of the code printed five random images from the folder of saree.



Similarly, we randomly loaded images of trouser and jeans.

The loaded images are given below:

### Jeans Men:

*#Randomly displaying the pictures from our file of jeans from train folder*

```
plt.figure(figsize=(20,20))
train_folder="C:/Users/Priyanka/Downloads/amazon_images/train/jeans"
for i in range(5):
    file = random.choice(os.listdir(train_folder))
    image_path= os.path.join(train_folder, file)
    img=mpimg.imread(image_path)
    ax=plt.subplot(1,5,i+1)
    ax.title.set_text(file)
    plt.imshow(img)
```



### Trouser Men:

*#Randomly displaying the pictures from our file of trouser from train folder*

```
plt.figure(figsize=(20,20))
train_folder=r"C:/Users/Priyanka/Downloads/amazon_images/train/trouser"
for i in range(5):
    file = random.choice(os.listdir(train_folder))
    image_path= os.path.join(train_folder, file)
    img=mpimg.imread(image_path)
    ax=plt.subplot(1,5,i+1)
    ax.title.set_text(file)
    plt.imshow(img)
```



## Loading the Images in the python notebook

Train and test dataset was separately loaded using the open cv python library. There are many ways in which we can load the images, one can use PIL, open cv, IPython, Keras to load images in python. Open cv python library is the mostly used among all to load and resize the images. In the project we have also used open cv to load the train and test data.

We have defined a function to load the datasets. Let us have a look at the function:

*#defining a function to load dataset which will return image data in array form and class name.*

```
def create_dataset(img_folder):

    img_data_array=[] #empty list to store image data
    class_name=[] #empty list to store class name

    for dir1 in os.listdir(img_folder):
        for file in os.listdir(os.path.join(img_folder, dir1)):

            image_path= os.path.join(img_folder, dir1, file)
            image= cv2.imread( image_path, cv2.COLOR_BGR2RGB)
            image=cv2.resize(image, (IMG_HEIGHT, IMG_WIDTH),interpolation = cv2.INTER_AREA)
            image=np.array(image)
            image = image.astype('float32')
            image /= 255
            img_data_array.append(image)
            class_name.append(dir1)

    return img_data_array, class_name
# extracting the image array and class name
img_data, class_name =create_dataset(r"C:/Users/Priyanka/Downloads/amazon_images/train/")
```

A function is defined which will return us the class name of each image and each image in the form of an array.

1. We have used os.listdir method to list all the files in the train folder. So as output it will give us three files of Saree, Jeans and Trousers folders.

2. We have used `os.path.join` method to join the paths of the folders listed by `os.listdir` with the train folder to get into the folders of Saree, Jeans and trouser in order to download the images.
3. Further, we have downloaded the images using open cv python method.
4. Image augmentation is performed on downloaded image in which the images are re-sized, converted into an array and normalized in the range of 0-1.
5. Image array and class to which the image belongs are returned.

The `class_name` returned by the function is a list containing the word saree, jeans and trouser. We need to convert this into integer type classification 0,1, and 2. That way we will be able to use the classes as label output.

6. Following steps were performed to convert the keywords saree, jeans and trouser to integer.

```
#creating a dictionary with Label and respective numerical association
```

```
target_dict={k: v for v, k in enumerate(np.unique(class_name))}  
target_dict  
{'jeans': 0, 'saree': 1, 'trouser': 2}
```

```
#setting the numerical value of each label
```

```
target_val= [target_dict[class_name[i]] for i in range(len(class_name))]
```

Here, we have built a dictionary in which we have the name of classes along with a integer representation. Further this dictionary is used to put label corresponding each image as 0,1 or 2.

7. Creating training data

```
# creating x_train and y_train
```

```
x_train=np.array(img_data, np.float32)  
y_train=np.array(list(map(int,target_val)), np.float32)
```

We have used the above downloaded image array and `target_val` list and stored the data under the names `x_train` and `y_train` respectively.

Similarly we have loaded the test data and split the data into `x_test` and `y_test`.

## Hardware and Software Requirements and Tools Used

Hardware specifications-

Processor- Intel(R) Core(TM) i5-10210U CPU @ 1.60GHz 2.11 GHz

RAM-8.00 GB

Edition-Windows10

This project was completed using Python 3.0 which is an interpreted high-level and general-purpose programming language. Jupyter notebook which is an open-source web application provided GUI environment for the python notebook.

Important libraries that were used in the completion of this project are:

1. Numpy- Also known as Numerical python is a python library used to perform algebraic and statical analysis.

2. Matplotlib.pyplot- Matplotlib is a plotting library in python used for creating visualizations in python. Pyplot is the sub-module of matplotlib library which is the collection of functions that makes matplotlib works like pyplot. Pyplot in simple language helps in making the figure decorative and informative. We can put labels with the help of pyplot, we can insert some lines and points in the figure using pyplot etc.

3. Tensor flow- It is a free and open-source library for machine learning. It is used in deep learning. It allows programmers to create large scale neural networks with dense layers.

4. Keras - Keras is a free open-source python library that works on top of tensor flow. It allows user to develop and evaluate deep learning models.

5. OS – The main purpose of the OS module is to interact with your operating system. We can use OS to create folders, remove folders, move folders, and sometimes change the working directory. We have used it to access the names of files within a file path by doing `listdir()` and also we have used `os.path.join` to join different path easily.

6. Random – Random is a python module which is used to generate pseudo-random numbers. It is used to perform some action randomly. We have used random to randomly pick image paths to load random images of our data.

7. OpenCV-Python – We have used this module to load images and resizing the image.

## Model Training and Testing

We have a built a convolutional neural network.

```
cnn=tf.keras.Sequential([tf.keras.layers.InputLayer(input_shape=(IMG_HEIGHT,IMG_WIDTH, 3)),
                        tf.keras.layers.Conv2D(filters=32, kernel_size=3, strides=(2, 2), activation='relu'),
                        tf.keras.layers.MaxPool2D(2,2),
                        tf.keras.layers.Conv2D(filters=64, kernel_size=3, strides=(2, 2), activation='relu'),
                        tf.keras.layers.MaxPool2D(2,2),
                        tf.keras.layers.Flatten(),
                        tf.keras.layers.Dense(64,activation='relu'),
                        tf.keras.layers.Dense(3, activation='softmax')])
cnn.compile(optimizer='adam', loss='sparse_categorical_crossentropy', metrics=['accuracy'])
```

The model type that we will be using is Sequential. Sequential is the easiest way to build a model in Keras. It allows you to build a model layer by layer.

We have 2 Conv2D layers. Convolution layer will deal with our input images, which are seen as 2-dimensional matrices.

32 in the first layer and 64 in the second layer are the number of nodes in each layer. This number can be adjusted to be higher or lower, depending on the size of the dataset. In our case, 32 and 64 works well, so we will stick with this for now.

Kernel size is the size of the filter matrix for our convolution. So a kernel size of 3 means we will have a 3x3 filter matrix.

Activation is the activation function for the layer. The activation function we will be using for convolutional layers is ReLU, or Rectified Linear Activation. This activation function has been proven to work well in neural networks.

In between the Conv2D layers and the dense layer, there is a 'Flatten' layer. Flatten serves as a connection between the convolution and dense layers.



'Dense' is the layer type we will use in for our output layer. Dense is a standard layer type that is used in many cases for neural networks.

We will have 3 nodes in our output layer, one for each possible outcome (0-2).

The activation is 'softmax'. Softmax makes the output sum up to 1 so the output can be interpreted as probabilities. The model will then make its prediction based on which option has the highest probability.

## Compiling the model

Next, we need to compile our model. Compiling the model takes three parameters: optimizer, loss and metrics.

The optimizer controls the learning rate. We have used 'adam' as our optimizer. Adam is generally a good optimizer to use for many cases. The Adam optimizer adjusts the learning rate throughout training.

The learning rate determines how fast the optimal weights for the model are calculated. A smaller learning rate may lead to more accurate weights (up to a certain point), but the time it takes to compute the weights will be longer.

We will use 'sparse\_categorical\_crossentropy' for our loss function. This is the most common choice for integer type classification. One advantage of using sparse categorical cross entropy is it saves time in memory as well as computation because it simply uses a single integer for a class, rather than a whole vector. A lower score indicates that the model is performing better.

To make things even easier to interpret, we will use the 'accuracy' metric to see the accuracy score when we train the model.

## Training the model

Now we will train our model. To train, we will use the 'fit()' function on our model with the following parameters: training data (x\_train), target data (y\_train) and the number of epochs.

The number of epochs is the number of times the model will cycle through the data. The more epochs we run, the more the model will improve, up to a certain point. After that point, the model will stop improving during each epoch. For our model, we will set the number of epochs to 25.

```
#fitting the training data
cnn.fit(x_train, y_train, epochs=25)
```

```
Epoch 17/25
28/28 [=====] - 3s 98ms/step - loss: 0.0266 - accuracy: 0.9966
Epoch 18/25
28/28 [=====] - 3s 96ms/step - loss: 0.0197 - accuracy: 0.9966
Epoch 19/25
28/28 [=====] - 3s 96ms/step - loss: 0.0182 - accuracy: 0.9977
Epoch 20/25
28/28 [=====] - 3s 94ms/step - loss: 0.0097 - accuracy: 1.0000
Epoch 21/25
28/28 [=====] - 3s 94ms/step - loss: 0.0080 - accuracy: 1.0000
Epoch 22/25
28/28 [=====] - 3s 94ms/step - loss: 0.0075 - accuracy: 1.0000
Epoch 23/25
28/28 [=====] - 3s 97ms/step - loss: 0.0060 - accuracy: 1.0000
Epoch 24/25
28/28 [=====] - 3s 95ms/step - loss: 0.0060 - accuracy: 1.0000
Epoch 25/25
28/28 [=====] - 3s 95ms/step - loss: 0.0038 - accuracy: 1.0000

Out[43]: <tensorflow.python.keras.callbacks.History at 0x1f5219befd0>
```

After 19 epoch we started getting an accuracy of 1 that means our model was able to understand all the images of the three classes.

## Using our model to make predictions

To see the actual predictions that our model has made for the test data, we are using the predict function. The predict function will give an array with 3 numbers. These numbers are the probabilities that the input image represents each class (0–2). The array index with the highest number represents the model prediction. The sum of each array equals 1 (since each number is a probability).

To show this, we will show the predictions for the first 5 images in the test set.

```
#predicting the output label for test data
y_pred=cnn.predict(x_test)

#displaying first five predictions
y_pred[:5]

array([[6.8443668e-01, 2.4977615e-07, 3.1556305e-01],
       [5.2379489e-01, 5.2998044e-12, 4.7620511e-01],
       [9.9977106e-01, 2.0983339e-12, 2.2890595e-04],
       [9.9980539e-01, 7.0145140e-10, 1.9461283e-04],
       [9.9696285e-01, 2.0464207e-10, 3.0370930e-03]], dtype=float32)
```

We can see the index of the highest numbers are 0, 0, 0, 0, 0.

Let us compare these with actual results.

```
#checking the actual values of the label
y_test[:5]

array([0., 0., 0., 0., 0.], dtype=float32)
```

Actual results too showed 0, 0, 0, 0, 0 for the first five images. Our model was correct in classifying the images.

### Performance metrics used

We have a classification type data so performance metrics like accuracy score, confusion matrix, classification report, f1 score were used to study the performance of the models.

### Interpretation of the result

Error Metrics	Result
Accuracy Score	0.9366
F1 Score	0.9367

Our model was able to classify the images into the categories fairly well as we got the accuracy score and f1 score of 0.94 each.

## **CONCLUSION**

We were successful in building an end-to-end deep learning model to classify the three categories of images.