**FLIP ROBO**

# *IMAGE SCRAPPING AND CLASSIFICATION PROJECT*

**Submitted by:**
**NIPAM GOGOI**

*Flip Robo Technologies*

## *ACKNOWLEDGMENT*

I would like to thank Flip Robo Technologies for providing me with the opportunity to work on this project from which I have learned a lot. I am also grateful to Miss Khushboo Garg for her constant guidance and support.

*Flip Robo Technologies*

# *INTRODUCTION*

Images are one of the major sources of data in the field of data science and AI. This field is making appropriate use of information that can be gathered through images by examining its features and details

The idea behind this project is to build a deep learning-based Image Classification model on images that will be scraped from e-commerce portal. This is done to make the model more and more robust.

## BUSINESS PROBLEM FRAMING

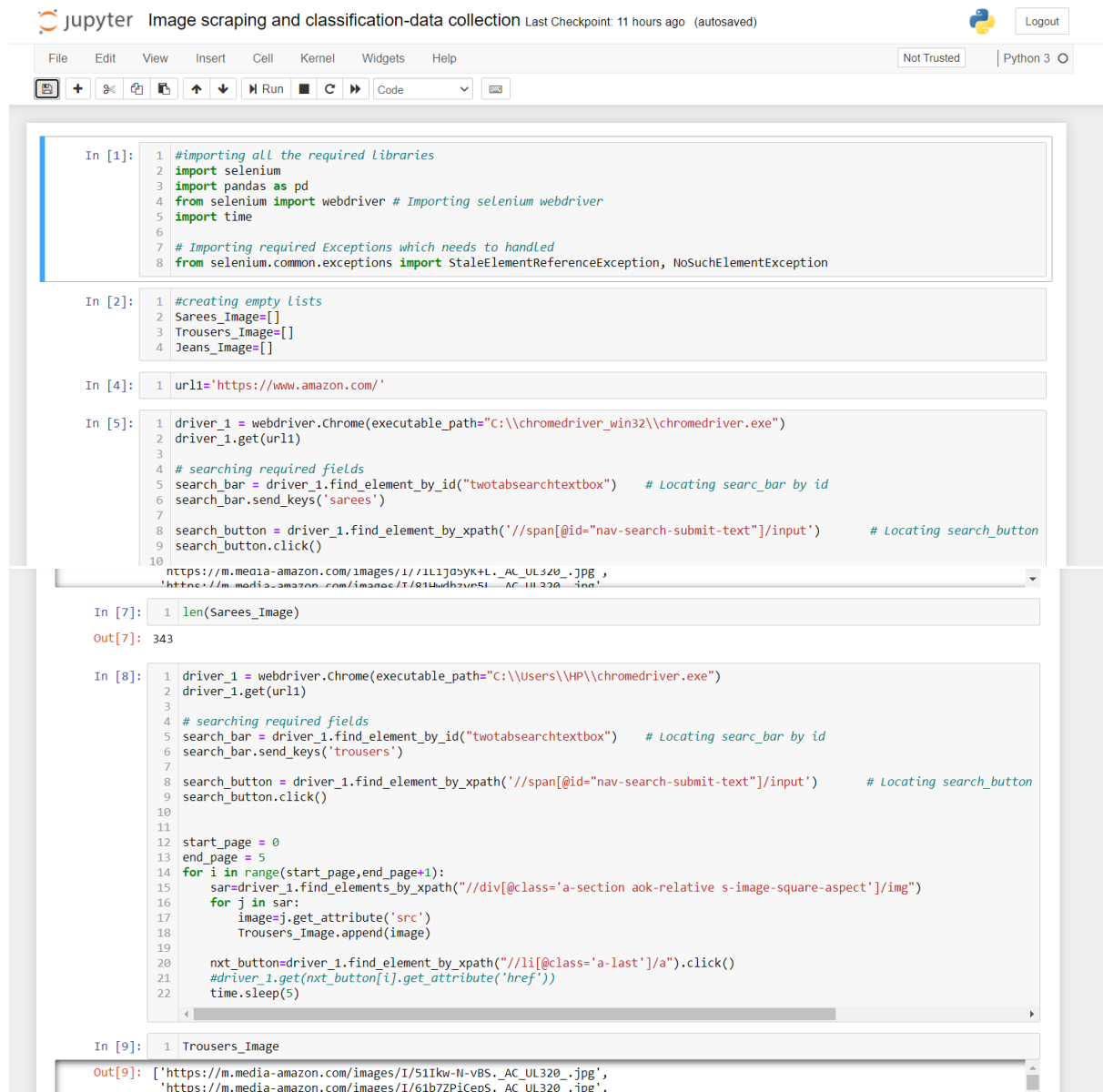This task is divided into two phases: Data Collection and Mode Building.

We need to scrape images from e-commerce portal, Amazon.com. The clothing categories used for scraping will be:

-   Sarees (women)
-   Trousers (men)
-   Jeans (men)

You need to scrape images of these 3 categories and build your data from it. That data will be provided as an input to your deep learning problem. You need to scrape minimum 200 images of each categories. There is no maximum limit to the data collection. You are free to apply image augmentation techniques to increase the size of your data but make sure the quality of data is not compromised.

# *ANALYTICAL PROBLEM FRAMING*

## DATA COLLECTION - WEBSCRAPPING

File    Edit    View    Insert    Cell    Kernel    Widgets    Help                                    Not Trusted    Python 3 ○

```python
In [1]:  1  #importing all the required libraries
         2  import selenium
         3  import pandas as pd
         4  from selenium import webdriver # Importing selenium webdriver
         5  import time
         6
         7  # Importing required Exceptions which needs to handled
         8  from selenium.common.exceptions import StaleElementReferenceException, NoSuchElementException
```

```python
In [2]:  1  #creating empty lists
         2  Sarees_Image=[]
         3  Trousers_Image=[]
         4  Jeans_Image=[]
```

```python
In [4]:  1  url1='https://www.amazon.com/'
```

```python
In [5]:  1  driver_1 = webdriver.Chrome(executable_path="C:\\chromedriver_win32\\chromedriver.exe")
         2  driver_1.get(url1)
         3
         4  # searching required fields
         5  search_bar = driver_1.find_element_by_id("twotabsearchtextbox")    # Locating searc_bar by id
         6  search_bar.send_keys('sarees')
         7
         8  search_button = driver_1.find_element_by_xpath('//span[@id="nav-search-submit-text"]/input')        # Locating search_button
         9  search_button.click()
        10
```
https://m.media-amazon.com/images/I/71L1ja5yK+L._AC_UL320_.jpg ,
'https://m.media-amazon.com/images/I/81Hwdhzyr5L._AC_UL320_.jpg'

```python
In [7]:  1  len(Sarees_Image)
Out[7]:  343
```

```python
In [8]:  1  driver_1 = webdriver.Chrome(executable_path="C:\\Users\\HP\\chromedriver.exe")
         2  driver_1.get(url1)
         3
         4  # searching required fields
         5  search_bar = driver_1.find_element_by_id("twotabsearchtextbox")    # Locating searc_bar by id
         6  search_bar.send_keys('trousers')
         7
         8  search_button = driver_1.find_element_by_xpath('//span[@id="nav-search-submit-text"]/input')        # Locating search_button
         9  search_button.click()
        10
        11
        12  start_page = 0
        13  end_page = 5
        14  for i in range(start_page,end_page+1):
        15      sar=driver_1.find_elements_by_xpath("//div[@class='a-section aok-relative s-image-square-aspect']/img")
        16      for j in sar:
        17          image=j.get_attribute('src')
        18          Trousers_Image.append(image)
        19
        20      nxt_button=driver_1.find_element_by_xpath("//li[@class='a-last']/a").click()
        21      #driver_1.get(nxt_button[i].get_attribute('href'))
        22      time.sleep(5)
```

```python
In [9]:  1  Trousers_Image
Out[9]:  ['https://m.media-amazon.com/images/I/51Ikw-N-vBS._AC_UL320_.jpg',
         'https://m.media-amazon.com/images/I/61b7ZPiCepS._AC_UL320_.jpg',
```

```
In [18]:   1  for index, link in enumerate(Trousers):
           2      print('Downloading {0} of 343 trouser images'.format(index+1))
           3      response=requests.get(link)
           4      with open('Trousers_Images/img{0}.jpeg'.format(index+1),"wb") as file:
           5          file.write(response.content)
```

```
Downloading 1 of 343 trouser images
Downloading 2 of 343 trouser images
Downloading 3 of 343 trouser images
Downloading 4 of 343 trouser images
Downloading 5 of 343 trouser images
Downloading 6 of 343 trouser images
Downloading 7 of 343 trouser images
Downloading 8 of 343 trouser images
Downloading 9 of 343 trouser images
Downloading 10 of 343 trouser images
Downloading 11 of 343 trouser images
Downloading 12 of 343 trouser images
Downloading 13 of 343 trouser images
Downloading 14 of 343 trouser images
Downloading 15 of 343 trouser images
Downloading 16 of 343 trouser images
Downloading 17 of 343 trouser images
Downloading 18 of 343 trouser images
Downloading 19 of 343 trouser images
Downloading 20 of 343 trouser images
```

```
In [25]:   1  for index, link in enumerate(Jeans):
           2      print('Downloading {0} of 343 jeans images'.format(index+1))
           3      response=requests.get(link)
           4      with open('Jeans_Images/img{0}.jpeg'.format(index+1),"wb") as file:
           5          file.write(response.content)
```

```
Downloading 1 of 343 jeans images
Downloading 2 of 343 jeans images
```

# MODEL BUILDING

```python
In [1]:  1  import os
         2  import numpy as np
```

```python
In [2]:  1  # Directory with Jeans pictures
         2  jeans_dir = os.path.join('Clothes/Train/Jeans_Images')
         3
         4  # Directory with Saree pictures
         5  saree_dir = os.path.join('Clothes/Train/Sarees_Images')
         6
         7  # Directory with Trouser pictures
         8  trousers_dir = os.path.join('Clothes/Train/Trousers_Images')
```

```python
In [3]:  1  train_jeans_names = os.listdir(jeans_dir)
         2  print(train_jeans_names[:5])
         3
         4  train_saree_names = os.listdir(saree_dir)
         5  print(train_saree_names[:5])
         6
         7  train_trousers_names = os.listdir(trousers_dir)
         8  print(train_trousers_names[:5])
```

```
['img1.jpeg', 'img10.jpeg', 'img100.jpeg', 'img101.jpeg', 'img102.jpeg']
['img1.jpeg', 'img10.jpeg', 'img100.jpeg', 'img101.jpeg', 'img102.jpeg']
['img1.jpeg', 'img10.jpeg', 'img100.jpeg', 'img101.jpeg', 'img102.jpeg']
```

```python
In [4]:  1  print('total jeans images:', len(os.listdir(jeans_dir)))
         2  print('total saree images:', len(os.listdir(saree_dir)))
         3  print('total trousers images:', len(os.listdir(trousers_dir)))
```

```
total jeans images: 300
total saree images: 300
total trousers images: 300
```

```python
In [5]:  1  import matplotlib.pyplot as plt
         2  import matplotlib.image as mpimg
```

```python
In [6]:  1  # Parameters for our graph; we'll output images in a 4x4 configuration
         2  nrows = 4
         3  ncols = 4
         4
         5  # Index for iterating over images
         6  pic_index = 0
         7
         8  fig = plt.gcf()
         9  fig.set_size_inches(ncols * 4, nrows * 4)
        10
        11  pic_index += 8
        12  next_jeans_pix = [os.path.join(jeans_dir, fname)
        13                    for fname in train_jeans_names[pic_index-8:pic_index]]
        14  next_saree_pix = [os.path.join(saree_dir, fname)
        15                    for fname in train_saree_names[pic_index-8:pic_index]]
        16  next_trousers_pix = [os.path.join(trousers_dir, fname)
        17                    for fname in train_trousers_names[pic_index-8:pic_index]]
        18
        19  print ("Showing some jeans pictures...")
        20  print()
        21  for i, img_path in enumerate(next_jeans_pix):
        22      # Set up subplot; subplot indices start at 1
        23      sp = plt.subplot(nrows, ncols, i + 1)
        24      sp.axis('Off') # Don't show axes (or gridlines)
        25
        26      img = mpimg.imread(img_path)
        27      plt.imshow(img)
        28
        29  plt.show()
        30
        31  print ("Showing some Sarees pictures...")
```

```
43  plt.show()
44
45  print ("Showing some trousers pictures...")
46  print()
47  fig = plt.gcf()
48  fig.set_size_inches(ncols * 4, nrows * 4)
49  for i, img_path in enumerate(next_trousers_pix):
50      # Set up subplot; subplot indices start at 1
51      sp = plt.subplot(nrows, ncols, i + 1)
52      sp.axis('Off') # Don't show axes (or gridlines)
53
54      img = mpimg.imread(img_path)
55      plt.imshow(img)
56
57  plt.show()
```

Showing some jeans pictures...



Showing some Sarees pictures...

```
1  print("Count of Training Images")
2  print("No.of Images of Sarees in train dataset :",len(os.listdir('Clothes/Train/Sarees_Images')))
3  print("No.of Images of Jeans in train dataset :",len(os.listdir('Clothes/Train/Jeans_Images')))
4  print("No.of Images of Trousers in train dataset :",len(os.listdir('Clothes/Train/Trousers_Images')))
5  print()
6  print("Count of Test Images")
7  print("No.of Images of Sarees in test dataset :",len(os.listdir('Clothes/Test/Sarees_Images')))
8  print("No.of Images of Jeans in test dataset :",len(os.listdir('Clothes/Test/Jeans_Images')))
9  print("No.of Images of Trousers in test dataset :",len(os.listdir('Clothes/Test/Trousers_Images')))
```

```
Count of Training Images
No.of Images of Sarees in train dataset : 300
No.of Images of Jeans in train dataset : 300
No.of Images of Trousers in train dataset : 300

Count of Test Images
No.of Images of Sarees in test dataset : 43
No.of Images of Jeans in test dataset : 43
No.of Images of Trousers in test dataset : 43
```

```
1  train_data='Clothes/Train'
2  test_data='Clothes/Test'
```

```
1  input_shape=(200,320,3)
2  batch_size=12
```

```
1  from tensorflow.keras.preprocessing.image import ImageDataGenerator
```

```
1  # Data Augmentation on Training Images
2
3  Train_datagen=ImageDataGenerator(rescale=1./255,
4                                      zoom_range=0.2,
5                                      rotation_range=30,
6                                      horizontal_flip=True)
7  Training_set=Train_datagen.flow_from_directory(train_data,
8                                      target_size=(200,320),
9                                      batch_size=batch_size,
10                                     class_mode='categorical')
11
```

```
Found 900 images belonging to 3 classes.
```

```
1  # Test Data Generator
2  Test_datagen=ImageDataGenerator(rescale=1./255)
3  Test_set=Test_datagen.flow_from_directory(test_data,
4                                      target_size=(200,320),
5                                      batch_size=batch_size,
6                                      class_mode='categorical')
```

```
Found 129 images belonging to 3 classes.
```

```
1  import tensorflow as tf
```

```
1  model = tf.keras.models.Sequential([
2      # Note the input shape is the desired size of the image 200x 320 with 3 bytes color
3      # The first convolution
4      tf.keras.layers.Conv2D(16, (3,3), activation='relu', input_shape=(200, 320, 3)),
5      tf.keras.layers.MaxPooling2D(2, 2),
6      # The second convolution
7      tf.keras.layers.Conv2D(32, (3,3), activation='relu'),
8      tf.keras.layers.MaxPooling2D(2,2),
9      # The third convolution
10     tf.keras.layers.Conv2D(64, (3,3), activation='relu'),
11     tf.keras.layers.MaxPooling2D(2,2),
12     # The fourth convolution
13     tf.keras.layers.Conv2D(64, (3,3), activation='relu'),
14     tf.keras.layers.MaxPooling2D(2,2),
15     # The fifth convolution
16     tf.keras.layers.Conv2D(64, (3,3), activation='relu'),
17     tf.keras.layers.MaxPooling2D(2,2),
18     # Flatten the results to feed into a dense layer
19     tf.keras.layers.Flatten(),
20     # 128 neuron in the fully-connected layer
21     tf.keras.layers.Dense(128, activation='relu'),
22     # 3 output neurons for 3 classes with the softmax activation
23     tf.keras.layers.Dense(3, activation='softmax')
24 ])
```

```
[15]:   1  model.summary()
```

Model: "sequential"

| Layer (type)                      | Output Shape          | Param # |
|-----------------------------------|-----------------------|---------|
| conv2d (Conv2D)                   | (None, 198, 318, 16)  | 448     |
| max_pooling2d (MaxPooling2D)      | (None, 99, 159, 16)   | 0       |
| conv2d_1 (Conv2D)                 | (None, 97, 157, 32)   | 4640    |
| max_pooling2d_1 (MaxPooling2      | (None, 48, 78, 32)    | 0       |
| conv2d_2 (Conv2D)                 | (None, 46, 76, 64)    | 18496   |
| max_pooling2d_2 (MaxPooling2      | (None, 23, 38, 64)    | 0       |
| conv2d_3 (Conv2D)                 | (None, 21, 36, 64)    | 36928   |
| max_pooling2d_3 (MaxPooling2      | (None, 10, 18, 64)    | 0       |
| conv2d_4 (Conv2D)                 | (None, 8, 16, 64)     | 36928   |
| max_pooling2d_4 (MaxPooling2      | (None, 4, 8, 64)      | 0       |
| flatten (Flatten)                 | (None, 2048)          | 0       |
| dense (Dense)                     | (None, 128)           | 262272  |
| dense_1 (Dense)                   | (None, 3)             | 387     |

Total params: 360,099
Trainable params: 360,099
Non-trainable params: 0

```
In [16]:   1  from tensorflow.keras.optimizers import RMSprop
           2
           3  model.compile(loss='categorical_crossentropy',
           4            optimizer=RMSprop(lr=0.001),
           5            metrics=['acc'])
```

```
In [17]:   1  total_sample = Training_set.n
```
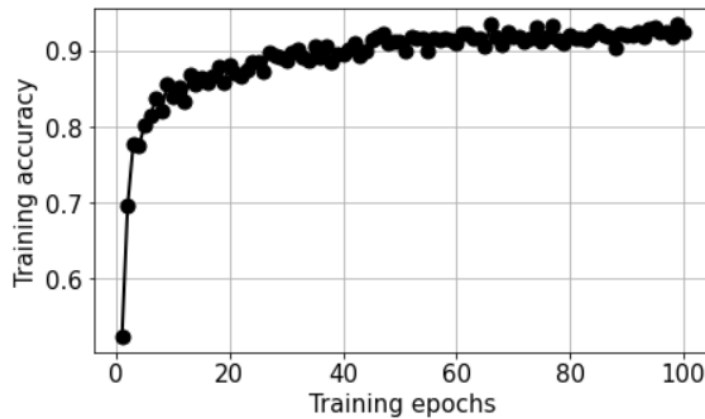
```
In [18]:   1  n_epochs = 100
```

```
In [19]:   1  history = model.fit(
           2        Training_set,
           3        steps_per_epoch=int(total_sample/batch_size),
           4        epochs=n_epochs,
           5        verbose=1)
```

```
75/75 [==============================] - 74s 984ms/step - loss: 0.5941 - acc: 0.7916
Epoch 4/100
75/75 [==============================] - 69s 915ms/step - loss: 0.5365 - acc: 0.7674
Epoch 5/100
75/75 [==============================] - 69s 919ms/step - loss: 0.5040 - acc: 0.8008
Epoch 6/100
75/75 [==============================] - 69s 917ms/step - loss: 0.4706 - acc: 0.8065
Epoch 7/100
75/75 [==============================] - 71s 940ms/step - loss: 0.3977 - acc: 0.8448
Epoch 8/100
75/75 [==============================] - 69s 913ms/step - loss: 0.4285 - acc: 0.8135
Epoch 9/100
75/75 [==============================] - 69s 916ms/step - loss: 0.4094 - acc: 0.8503
Epoch 10/100
75/75 [==============================] - 69s 922ms/step - loss: 0.3892 - acc: 0.8363
Epoch 11/100
75/75 [==============================] - 69s 917ms/step - loss: 0.3691 - acc: 0.8626
Epoch 12/100
```

```
1  plt.figure(figsize=(7,4))
2  plt.plot([i+1 for i in range(n_epochs)],history.history['acc'],'-o',c='k',lw=2,markersize=9)
3  plt.grid(True)
4  plt.title("Training accuracy with epochs\n",fontsize=18)
5  plt.xlabel("Training epochs",fontsize=15)
6  plt.ylabel("Training accuracy",fontsize=15)
7  plt.xticks(fontsize=15)
8  plt.yticks(fontsize=15)
9  plt.show()
```
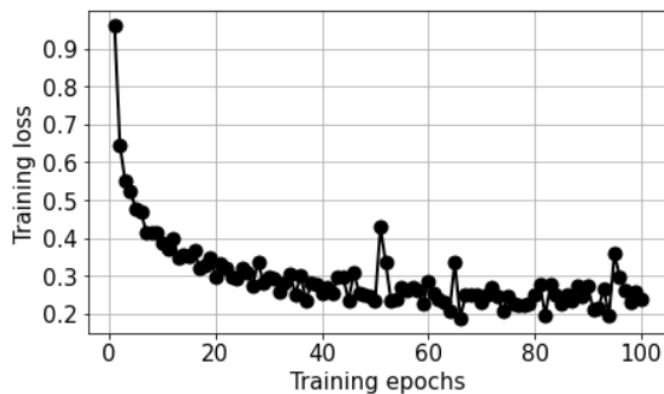

Training accuracy with epochs

```
1  plt.figure(figsize=(7,4))
2  plt.plot([i+1 for i in range(n_epochs)],history.history['loss'],'-o',c='k',lw=2,markersize=9)
3  plt.grid(True)
4  plt.title("Training loss with epochs\n",fontsize=18)
5  plt.xlabel("Training epochs",fontsize=15)
6  plt.ylabel("Training loss",fontsize=15)
7  plt.xticks(fontsize=15)
8  plt.yticks(fontsize=15)
9  plt.show()
```


Training loss with epochs

```
from PIL import Image
```

```
test_jeans = os.listdir('Clothes/Test/Jeans_Images/')[:5]
print(test_jeans)
test_sarees = os.listdir('Clothes/Test/Sarees_Images/')[:5]
print(test_sarees)
test_trousers = os.listdir('Clothes/Test/Trousers_Images/')[:5]
print(test_trousers)
```

```
['img301.jpeg', 'img302.jpeg', 'img303.jpeg', 'img304.jpeg', 'img305.jpeg']
['img301.jpeg', 'img302.jpeg', 'img303.jpeg', 'img304.jpeg', 'img305.jpeg']
['img301.jpeg', 'img302.jpeg', 'img303.jpeg', 'img304.jpeg', 'img305.jpeg']
```

```
Categories = ["Jeans", "Sarees", "Trousers"]
```

```
set(Test_set.classes)  # Unique Values of the Classes to be predicted
```
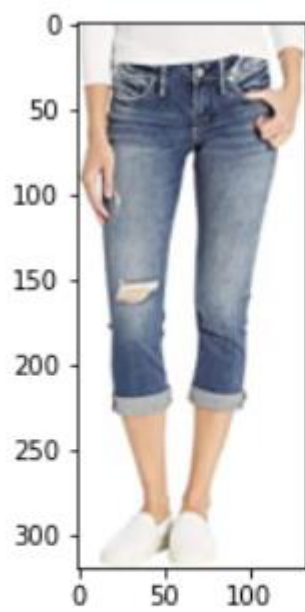
```
{0, 1, 2}
```

```
# Defining Function to plot the test image and predict class using the model trained
def ploting_predict(path,img_name):
    img_test=Image.open(path + img_name)
    plt.imshow(img_test)
    plt.show()  # Plotting the input Test Image
    img_test = img_test.resize((200,320))
    img_test=np.expand_dims(img_test,axis=0) # Expand dimensions for proper prediction
    print("Predicted as: ",Categories[int(np.argmax(model.predict(img_test)))])  # Predicting Classes using the Model
    print()
```

```
print()
```

In [30]:
```
path = 'Clothes/Test/Jeans_Images/'
for i in test_jeans:
    ploting_predict(path,i)
```

Predicted as:   Jeans

```
1  path = 'Clothes/Test/Sarees_Images/'
2  for i in test_sarees:
3      ploting_predict(path,i)
```



Predicted as:   Sarees

```
1  path = 'Clothes/Test/Trousers_Images/'
2  for i in test_trousers:
3      ploting_predict(path,i)
```

## *CONCLUSION*

- The Image Data was collected using webscrapping from Amazon for Jeans, Sarees and Trousers.

- We have used Deep Learning model – Convolutional Neural Network for the project

- The model is working well and was able to classify the three clothing items properly with 92% overall accuracy.

- We can improve the classification of Jeans and Trousers by increasing the training dataset. For Sarees it was accurately predicted in test data.

- Since in all three categories there were some extra/unnecessary items other than the main items hence, it could have been removed and we could have got better result. Moreover, training data could have been increased.