

DATA SCIENCE, MACHINE LEARNING

# How to Create a New Custom Dataset from Images



Uday Sai · Follow

Published in Towards AI · 7 min read · Aug 24, 2020

If you're a person like me, trying to build your custom image dataset out of raw images, then this article is just for you!

We all have learned how to build machine learning models on the classic MNIST/Fashion MNIST datasets. But, what if you want to train a model to recognize your friends' faces? A dataset for that purpose is not readily available on the internet.

After working on public datasets for months, I wanted to create a custom dataset of my face images and use them for face identification.

Real expertise is demonstrated by using machine learning to solve your own problems. Building *your own image dataset* is a non-trivial task by itself. Surprisingly, it is covered far less comprehensively in most online courses.

I searched for ways to do it and finally figured it out.

In this article, you will learn how to prepare your own dataset of raw images.

which you can then use for your own image classification/computer vision projects.

## Steps

1. Gather images for your dataset
2. Rename the pictures according to their classes
3. Merge them into one folder
4. Resize the pictures
5. Convert all images into the same file format
6. Convert images into a CSV file
7. A few tweaks to the CSV file
8. Load the CSV (BONUS)

### Gather images for your dataset

As an example, let's say that I want to build a model that can differentiate between Keanu Reeves and me XD.

If you need to create a dataset of your own face or bulk download images from google, [this article](#) from *pyimagesearch* walks you through it.

After getting the images, sort the images into different folders according to their classes. For the sake of simplicity, I'm going to use just five images per class (You can use as many as you want. The more, the better).

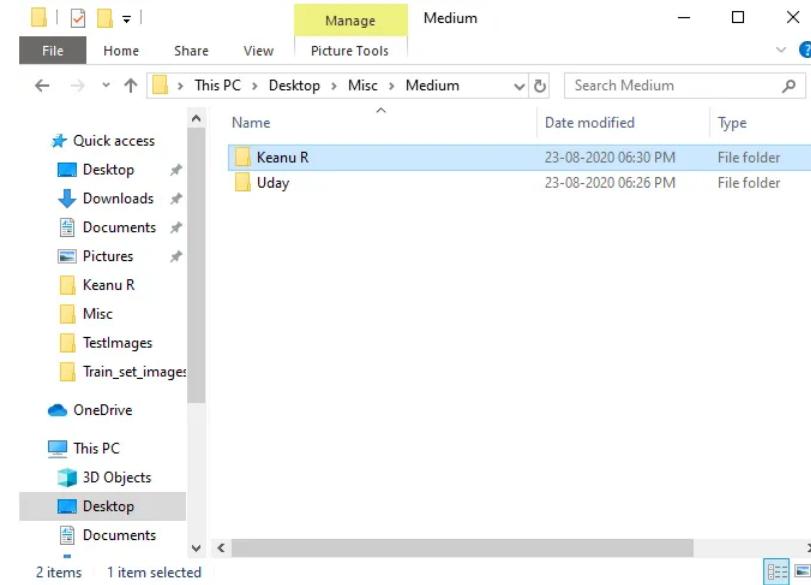
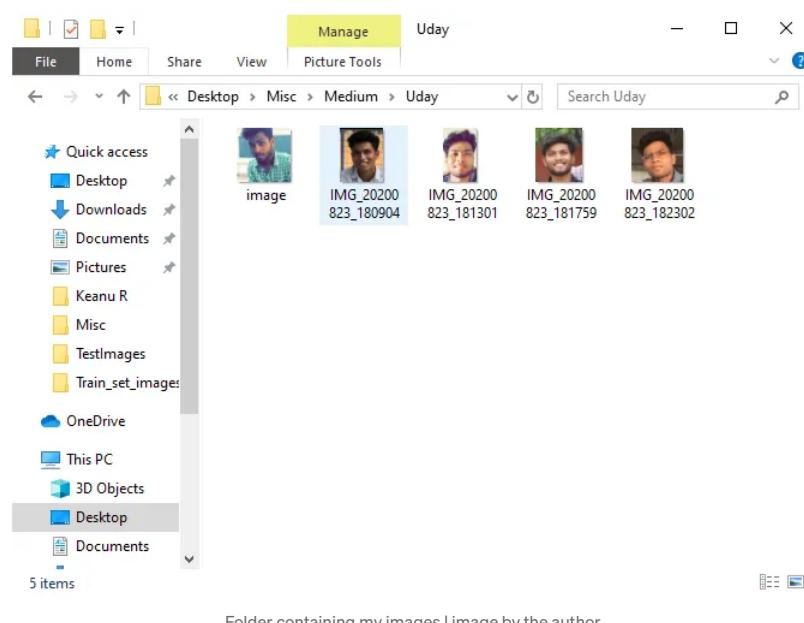
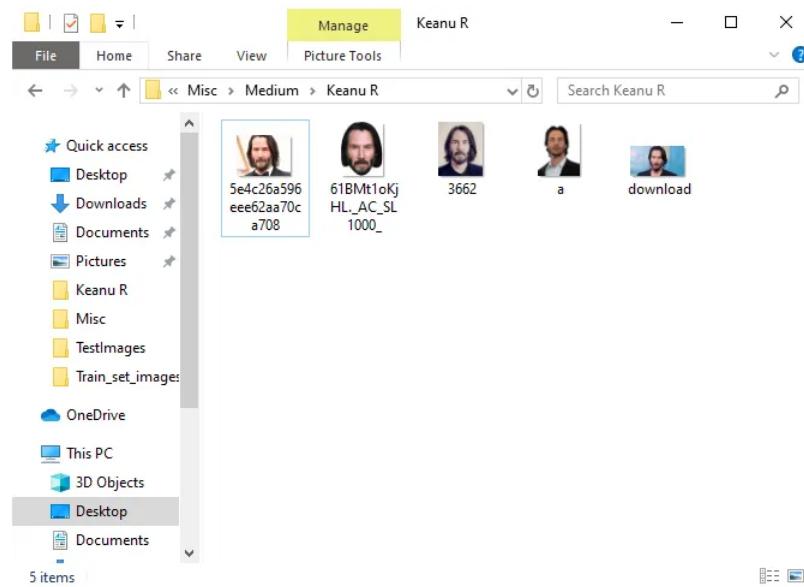


Image by the [author](#)

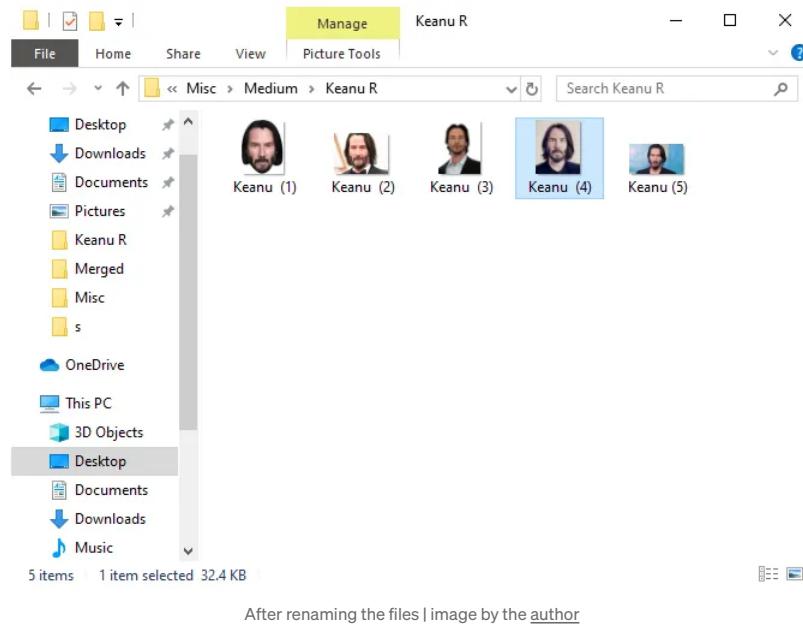


A machine learning model is only as good as the data we put into it.

Clean the data. Remove duplicates. Crop the images around your point of interest (in this example, faces of Keanu and me) to make the most of your data.

#### **Rename the pictures according to their classes**

1. Open the folder and select all images.
2. Right-click on them.
3. Rename all of them by their class.



4. Repeat it for all the remaining classes. Name the classes with at least one different alphabet (this is needed in the latter part of the process).

#### **Merge all the images into a single folder.**

#### **Resize the pictures**

The tool we use for this is [\*Image Resizer for Windows\*](#). It's free, small, and completely malware-free.



Lets you resize images by right-clicking

Image by the [author](#)

Once it's downloaded, click **Install**.

Once the program is installed on your computer, you're good to go. Now, go to the folder containing the photos that you want to resize.

Select your photos. Then right-click on them and choose to **Resize pictures** from the options.

A window will then pop up. Here, you can modify the basic settings for the pictures that will be processed.

You can select the size for the pictures. In this case, I resize the images to 48 x 48 pixels.

## Resize your pictures

Select a size.

- Small (fits within 854 × 480 pixels)
- Medium (fits within 1366 × 768 pixels)
- Large (fits within 1920 × 1080 pixels)
- Phone (fits within 320 × 569 pixels)
- Custom Fit  ×  Pixels

- Make pictures smaller but not larger
- Resize the original pictures (don't create copies)
- Ignore the orientation of pictures

[Advanced options...](#)

**Resize**

Cancel

Don't forget to change from "Fit" to "Stretch." | [image by the author](#)

Note: Sometimes, smaller pictures get ignored by the resizer. After Resizing, select all the images and verify if all the images are of the same size.

### Convert all images into the same file format

Here is a neat trick to do this easily and efficiently. You could either choose .png or .jpg format.

**Step 1** — Type cmd on the taskbar search field and jointly press **Ctrl + Shift + Enter** keys. If you come across UAC prompt, click **Yes**.

**Step 2** — In the Command Prompt, first input the path of the new folder where you stored the files (images of Spotlight). To do so, type in –

```
cd path of the folder
```

**Note** — Please replace the **pathofthefolder** with the actual path.

```
cd C:\Users\Uday\Desktop\Misc\Medium\Merged
```

**Step 3** — To Change the images to JPG format, type in the given batch command, and press **Enter**.

```
Ren *.* *.jpg
```

**Step 4** — To convert the images to PNG format, use the following batch command –

```
Ren *.* *.png
```

### Convert the images into a CSV

Run the following code to convert all the images into a CSV and label them accordingly.

```
from PIL import Image
import numpy as np
import sys
import os
import csv

# default format can be changed as needed
def createFileList(myDir, format='.jpg'):
    fileList = []
    print(myDir)
    labels = []
    names = []
    keywords = {"K": "1", "U": "0",} # keys and values to be changed
    as needed

    for root, dirs, files in os.walk(myDir, topdown=True):
        for name in files:
            if name.endswith(format):
                fullName = os.path.join(root, name)
                fileList.append(fullName)
                for keyword in keywords:
                    if keyword in name:
                        labels.append(keywords[keyword])
                    else:
                        continue
                names.append(name)
    return fileList, labels, names

# load the original image
myFileList, labels, names = createFileList('/content/')
i = 0
for file in myFileList:
    print(file)
    img_file = Image.open(file)
    # img_file.show()

    # get original image parameters...
    width, height = img_file.size
    format = img_file.format
    mode = img_file.mode

    # Make image Greyscale
    img_grey = img_file.convert('L')
    #img_grey.save('result.png')
    #img_grey.show()

    # Save Greyscale values
    value = np.asarray(img_grey.getdata(),
                      dtype=np.int).reshape((width, height))
    value = value.flatten()

    value = np.append(value, labels[i])
    i += 1

    print(value)
    with open("name_you_want.csv", 'a') as f:
        writer = csv.writer(f)
        writer.writerow(value)
```

1. I've used K and U alphabets as keys to recognize the classes from the file names (Keanu has K in it and Uday has U in it). Change it as per your needs.
2. To keep images in color instead of greyscale images replace 'L' with 'RGB.' Also, add depth value before saving the image. Depth = 3 representing the number of color channels (Red, Green, Blue).

```
img_grey = img_file.convert('L')# replace L with RGB
value = np.asarray(img_grey.getdata(), dtype=np.int).reshape((width,
height, 3))
```

3. *name\_you\_want* will be the name of the CSV file created. Feel free to

change it.

You have your dataset ready. Well, almost ready.

A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U
129	197	399	202	203	204	206	207	205	208	209	208	209	210	212	210	213	211	212	210	
2	78	56	77	67	67	67	63	69	55	66	62	70	60	69	74	62	72	69	61	77
3	11	9	10	12	12	8	10	10	5	10	10	3	8	6	7	12	1	7	9	
4	201	204	208	211	211	209	207	205	204	211	209	197	187	187	192	196	194	176	170	198
5	255	254	255	255	255	255	254	254	254	247	254	252	252	253	254	231	187	165	150	148
6	237	238	235	235	227	189	136	105	83	82	81	80	74	68	69	74	72	72	73	74
7	232	232	232	232	231	229	230	232	231	234	231	232	231	238	230	235	232	233	238	226
8	251	251	251	251	251	251	251	251	251	251	251	251	251	251	251	251	251	251	251	251
9	286	225	245	247	248	248	250	249	250	250	247	253	254	223	124	118	103	95	93	73
10	113	102	105	113	105	100	113	124	134	131	129	91	99	103	104	101	94	137	161	143
11																				
12																				
13																				
14																				
15																				
16																				
17																				
18																				
19																				
20																				
21																				

And the last column in the sheet are the labels | image by the [author](#)

### A few tweaks to the CSV file

1. Scroll to the end, click on an empty cell and use the following Excel formula to concatenate the pixel values.

```
=TEXTJOIN(" ", TRUE, (A1:B1))  
#replace B1 with last but one column name
```

CJJ	CJK	CJL	CJM	CIN	CJO	CJP	CJQ	CJR	CJS	CJT	CJU	CJV	CJW	CJX
1	18	19	15	15	22	26	25	1	=TEXTJOIN(" ",TRUE,(A1:CJ1))					
2	156	162	156	155	126	73	60	0						
3	168	164	215	204	198	191	193	0						
4	31	123	128	131	152	140	139	1						
5	254	254	254	255	253	253	253	1						
6	94	87	93	87	74	103	114	0						
7	79	64	70	55	54	100	91	0						
8	15	15	13	11	9	8	8	1						
9	70	73	79	83	82	83	91	1						
10	54	41	46	55	94	99	62	0						
11														

Image by the [author](#)

2. Drag the formula to the remaining rows.

A screenshot of Microsoft Excel showing a table with data. The formula `=TEXTJOIN(",\",TRUE,A1:CJ1)` is entered in cell CJ1. The data consists of 10 rows of numerical values from columns CJJ to CJX.

Image by the author

3. Copy that column values to the notepad. Re-copy them and paste them back. This way, you will retain the pixel values and not the formula.

4. Now select all cells except the labels and concatenated values and delete them.

A screenshot of Microsoft Excel showing the same table after deleting the labels and concatenated values. Only the numerical data remains in the cells.

Image by the author

5. Cut the remaining columns and paste them at the beginning of the sheet.

A screenshot of Microsoft Excel showing the table with the columns rearranged. The first 10 columns (CJJ to CJX) have been moved to the end of the sheet, and new columns A through O are now at the beginning.

Image by the author

6. Name the columns accordingly.

Image by the author

Aaaaanndddd, we're DONE!

Congratulations! You've created a brand new custom image dataset from scratch.

## Bonus

## Load the CSV

Load the CSV and run this following code snippet and you all good to good.

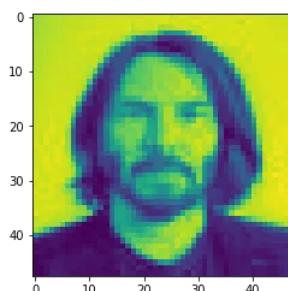
```
import pandas as pd
import cv2
import matplotlib.pyplot as plt
import numpy as np
from sklearn.model_selection import train_test_split

dataset_path = '/content/Keanu&Uday.csv'
image_size=(48,48) #add 3 if RGB image

def load():
    data = pd.read_csv(dataset_path)
    pixels = data['Pixels'].tolist()
    width, height= 48, 48 ,# add depth 3 if RGB image
    faces = []
    for pixel_sequence in pixels:
        face = [int(pixel) for pixel in pixel_sequence.split(' ')]
        face = np.asarray(face).reshape(width, height,) #add depth if
RGB image
        a = face
        face = np.resize(face.astype('uint8'),image_size)
        faces.append(face.astype('float32'))

    faces = np.asarray(faces)
    A = faces
    faces = np.expand_dims(faces, -1)
    return faces, A

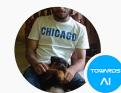
faces,A = load()
plt.imshow(A[0].astype("uint8"))
```



Thanks for reading! I hope you found this article useful. Here's the link to the Colab notebook.

Resources: [GitHub repository](#) and [Google Colab](#)

Machine Learning   Deep Learning   Neural Networks   Data Science  
Data Visualization



## Written by Uday Sai

20 Followers · Writer for Towards AI

Pursuing excellence and adding value to people's lives | CSE Student at CMR Technical Campus, Hyderabad

Follow

### More from Uday Sai and Towards AI



in Towards AI

#### How To Do RAG Without Vector Databases

Knowledge Graphs Are All You Need 😊

◆ · 10 min read · Mar 2, 2024

451 1

+ ...



Ruben Aster in Towards AI

#### LangChain SQL Agent for Massive Documents Interaction

Leverage LangChain SQL Agent and GPT for Document Analysis and Interaction

18 min read · Mar 7, 2024

801 5

+ ...



 IVAN ILIN in Towards AI

## Advanced RAG Techniques: an Illustrated Overview

A comprehensive study of the advanced retrieval augmented generation techniques...

19 min read · Dec 17, 2023

 See all from Uday Sai

 See all from Towards AI

 Boris Meinardus in Towards AI

## Machine Learning Was Hard Until I Learned These 5 Secrets!

The secrets no one tells you but make learning ML a lot easier and enjoyable.

◆ 10 min read · Mar 28, 2024

 1K

 8



...

## Recommended from Medium



 Maahi Patel

### The Complete Guide to Image Preprocessing Techniques in...

Have you ever struggled with poor quality images in your machine learning or compute...

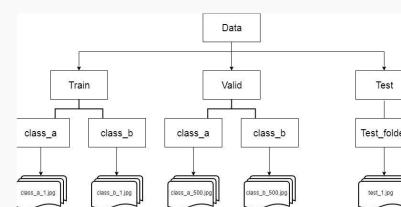
11 min read · Oct 23, 2023

 144

 3



...



 Md Shahbaz Alam

### How to Pass Image Datasets to CNN Models using Image Data...

“Data will talk to you if you’re willing to listen.”  
—Jim Bergeson

4 min read · Nov 2, 2023

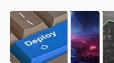
 93





...

## Lists



### Predictive Modeling w/ Python

20 stories · 1071 saves



### Practical Guides to Machine Learning

10 stories · 1286 saves



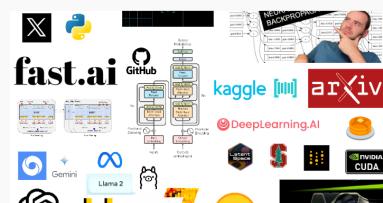
### Natural Language Processing

1355 stories · 838 saves



### data science and AI

40 stories · 124 saves



 Benedict Neo in bitgrit Data Science Publication

### Roadmap to Learn AI in 2024

A free curriculum for hackers and programmers to learn AI

11 min read · Mar 10, 2024

 9.6K

 103



...



 Anushree Das

### Image Dataset Analysis using Python Libraries—Pandas and...

A tutorial on how to use pandas and matplotlib for data analysis.

4 min read · Oct 10, 2023

 1





...



 Sohaib Zafar

## Choosing the Right Pre-Trained Model: A Guide to VGGNet, ResN...

Introduction

13 min read · Nov 12, 2023

 29  1



 Pallavi Vangari

## Exploratory Data Analysis on CIFAR-10 Dataset using Python

Introduction

2 min read · Oct 31, 2023

[See more recommendations](#)