

✓ importing libraries

```
import pandas as pd
```

✓ reading only required data

```
df = pd.read_csv("https://raw.githubusercontent.com/ieee8023/covid-chestxray-dataset/refs/heads/master/metadata.csv")
```

✓ selecting required columns

df

	view	filename	grid icon
0	PA	auntninnie-a-2020_01_28_23_51_6665_2020_01_28...	edit icon
1	PA	auntninnie-b-2020_01_28_23_51_6665_2020_01_28...	
2	PA	auntninnie-c-2020_01_28_23_51_6665_2020_01_28...	
3	PA	auntninnie-d-2020_01_28_23_51_6665_2020_01_28...	
4	PA	nejmc2001573_f1a.jpeg	
...	...	...	
945	AP	072ecaf8c60a81980abb57150a8016_jumbo-9.jpeg	
946	AP	ff33c406392b968d483174c97eb857_jumbo-9.jpeg	
947	PA	000001-266.jpg	
948	AP	000001-272.jpg	
949	L	000002-268.jpg	

950 rows × 2 columns

Next steps: [Generate code with df](#) [New interactive sheet](#)

✓ metadata info

df.info()

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 950 entries, 0 to 949
Data columns (total 2 columns):
 #   Column      Non-Null Count  Dtype  
--- 
 0   view        950 non-null    object 
 1   filename    950 non-null    object 
dtypes: object(2)
memory usage: 15.0+ KB
```

✓ metadata description

df.describe()

	view	filename	grid icon
count	950	950	
unique	7	950	
top	PA	000002-268.jpg	
freq	344	1	

✓ selecting only AP and PA value views

```
final_df = df[df["view"].isin(["AP", "PA"])]
```

final\_df

	view	filename	grid icon
0	PA	auntninnie-a-2020_01_28_23_51_6665_2020_01_28_...	edit icon
1	PA	auntninnie-b-2020_01_28_23_51_6665_2020_01_28_...	
2	PA	auntninnie-c-2020_01_28_23_51_6665_2020_01_28_...	
3	PA	auntninnie-d-2020_01_28_23_51_6665_2020_01_28_...	
4	PA	nejmc2001573_f1a.jpeg	
...	...	...	
943	AP	02b973e10caa192fd4e6825ad4aeaf_jumbo-10.jpeg	
945	AP	072ecacf8c60a81980abb57150a8016_jumbo-9.jpeg	
946	AP	ff33c406392b968d483174c97eb857_jumbo-9.jpeg	
947	PA	000001-266.jpg	
948	AP	000001-272.jpg	

547 rows × 2 columns

Next steps: [Generate code with final\\_df](#) [New interactive sheet](#)

```
final_df.info()
final_df.describe()
```

```
<class 'pandas.core.frame.DataFrame'>
Index: 547 entries, 0 to 948
Data columns (total 2 columns):
 #   Column      Non-Null Count  Dtype  
--- 
 0   view        547 non-null    object 
 1   filename    547 non-null    object 
dtypes: object(2)
memory usage: 12.8+ KB
```

	view	filename	grid icon
count	547	547	
unique	2	547	
top	PA	000001-272.jpg	
freq	344	1	

```
ap = df[df["view"] == "AP"]
pa = df[df["view"] == "PA"]
ap_filename = ap["filename"]
pa_filename = pa["filename"]
```

```
print("AP dataframe")
ap.info()
print("\n")
print("PA dataframe")
pa.info()
```

```
AP dataframe
<class 'pandas.core.frame.DataFrame'>
Index: 203 entries, 9 to 948
Data columns (total 2 columns):
 #   Column      Non-Null Count  Dtype  
--- 
 0   view        203 non-null    object 
 1   filename    203 non-null    object 
dtypes: object(2)
memory usage: 4.8+ KB
```

```
PA dataframe
<class 'pandas.core.frame.DataFrame'>
Index: 344 entries, 0 to 947
Data columns (total 2 columns):
 #   Column      Non-Null Count  Dtype  
--- 
 0   view        344 non-null    object 
 1   filename    344 non-null    object 
dtypes: object(2)
```

```
dtypes: object(2)
memory_usage: 8.1+ KB
```

## ✓ downloading images

```
import urllib
import os
from concurrent.futures import ThreadPoolExecutor
from functools import partial

url_base = "https://raw.githubusercontent.com/ieee8023/covid-chestxray-dataset/refs/heads/master/images/"

def downloadImg(imgName:str,imgType:str):
    fullUrl = f"{url_base}{imgName}"
    fileDest = f"./images/{imgType}/{imgName}"
    if os.path.exists(fileDest):
        return
    try:
        urllib.request.urlretrieve(fullUrl, fileDest)
    except Exception as e:
        return f"Error downloading {imgName}:{e}"
```

## mounting google drive

## ✓ making directories to save images

```
saveDirs = ["./images","./images/AP","./images/PA"]
for sd in saveDirs:
    os.makedirs(sd,exist_ok=True)
```

## ✓ parallelly downloading images

```
with ThreadPoolExecutor(max_workers=10) as executor:
    func_with_args = partial(downloadImg,imgType="AP")
    executor.map(func_with_args,ap_filename.tolist())

with ThreadPoolExecutor(max_workers=10) as executor:
    func_with_args = partial(downloadImg,imgType="PA")
    executor.map(func_with_args,pa_filename.tolist())
```

## ✓ Summary

```
print("Total images in original dataset :",df.filename.count())
print("No of images after preprocessing :",final_df.filename.count())
print("\nDistribution of views\n",final_df.view.value_counts())

Total images in original dataset : 950
No of images after preprocessing : 547

Distribution of views
view
PA    344
AP    203
Name: count, dtype: int64
```