

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import requests
import json
import seaborn as sns
import re

from sklearn.model_selection import train_test_split
from sklearn.linear_model import LogisticRegression
from sklearn.neighbors import KNeighborsClassifier
from sklearn.ensemble import RandomForestClassifier
from sklearn.model_selection import RandomizedSearchCV, GridSearchCV

import sklearn
%matplotlib inline

import pandas
import numpy

from sklearn.metrics import accuracy_score


pd.set_option('display.max_rows', 10)
pd.set_option('display.max_columns', 10)
def build_dataframe(file):
    """This function will create Pandas DataFrame"""
    df = pd.read_csv(file)
    return df


def draw_histogram(dataframe):
    """create histogram for dataframe passed"""
    dataframe.isna().sum().plot(kind='bar')
    plt.show()
```

```

def type_categorization(df1, df2):
    """create histogram to categorize types of two dataframes"""
    df1['type'].hist(bins=80, alpha=0.8)
    df2['type'].hist(bins=80, alpha=0.8)
    plt.legend(['Disney_Plus', 'Hulu'])
    plt.show()

def missing_values_check(dataframe):
    """This will check missing values within dataframe"""
    movie_missing = dataframe.isnull().sum()
    return movie_missing

def count_entries(file_name, chunk_size, colname):
    """ This will return a dictionary with counts of occurrences as value for key"""
    counts_dict = {}
    for chunk in pd.read_csv(file_name, chunksize=chunk_size):
        for entry in chunk[colname]:
            if entry in counts_dict.keys():
                counts_dict[entry] += 1
            else:
                counts_dict[entry] = 1
    return counts_dict

# loading data from csv file
df_disney_plus= build_dataframe('C:/Users/Swara/FidelityA/UCD Project/disney_plus_titles.csv')
df_hulu = build_dataframe('C:/Users/Swara/FidelityA/UCD Project/hulu_titles.csv')

# Code to demonstrate basic operations on pandas dataframe like find missing & duplicate

```

```

disney_missing = missing_values_check(df_disney_plus)
hulu_missing = missing_values_check(df_hulu)
print('Display Missing Values - Disney Plus: ', disney_missing)
print('Display Missing Values - Hulu : ', hulu_missing)
disney_dups = df_disney_plus.duplicated(subset=['title'], keep=False)
print(df_disney_plus[disney_dups])
hulu_dups = df_hulu.duplicated(subset=['title'], keep=False)
print(df_hulu[hulu_dups])
# Merging dataframes and removing duplicates

df_merge = pd.concat([df_disney_plus, df_hulu])
print(df_merge.info())
duplicates_combined = df_merge.duplicated(subset=['title'], keep=False)
print(df_merge[duplicates_combined].info())
#Drop duplicates
distinct_movies = df_merge.drop_duplicates(subset=['title','director','cast'])
print(distinct_movies.info())

# Code to demonstrate creating function, dictionary, using iterators

result_counts = count_entries('C:/Users/Swara/FidelityA/UCD Project/disney_plus_titles.csv', 10,
'release_year')
print(result_counts)

# histogram for two dataframe

draw_histogram(df_disney_plus)
draw_histogram(df_hulu)

# Compare type of each dat set(TV,Movie Show)

type_categorization(df_disney_plus, df_hulu)

```

```
# Code to demonstrate loading data from API

response_API = requests.get('https://api.covid19india.org/state_district_wise.json')

#print(response_API.status_code)

data = response_API.text

parse_json = json.loads(data)

active_case = parse_json['Andaman and Nicobar Islands']['districtData']['South Andaman']['active']

print("Active cases in South Andaman:", active_case)
```

```
df_tweets = build_dataframe('C:/Users/Swara/FidelityA/UCD Project/Tweets.csv')

hash_tags = []

regex = r"#\b\w\w+\b"

for i in df_tweets['text']:

    word = re.findall(regex, i)

    if len(word) > 0:

        hash_tags.append(word)

print(hash_tags)
```

```
# Code to demonstrate regular expression in python, Extracting all hash tags from Tweets

Tweet= pandas.read_csv("C:/Users/Swara/FidelityA/UCD Project/Tweets.csv")

Tweet.head()
```

```
(len(Tweet)-Tweet.count())/len(Tweet)

del Tweet['tweet_coord']

del Tweet['airline_sentiment_gold']

del Tweet['negativereason_gold']

Mood_count=Tweet['airline_sentiment'].value_counts()

Index = [1,2,3]

plt.bar(Index,Mood_count)

plt.xticks(Index,['negative','neutral','positive'],rotation=45)

plt.ylabel('Mood Count')

plt.xlabel('Mood')
```

```

plt.title('Count of Moods')

#Machine Learning Code snippet

#Load Data

df = pd.read_csv("C:/Users/Swara/FidelityA/UCD Project/heart.csv")

df.shape

df.head()

train_df.prognosis.value_counts().head()

#how many class of one feature or target

df["target"].value_counts()

#bar chart

df["target"].value_counts().plot(kind='bar', color=["salmon", "lightblue"])

df.info()

#check missing values of all features

df.isna().sum()

df.describe()

#Heart Disease Frequency according to Sex

df.sex.value_counts()

#Compare target and sex column

pd.crosstab(df.target, df.sex)

#Create plot of crosstab

pd.crosstab(df.target, df.sex).plot(kind="bar", figsize=(10,6), color=["salmon", "lightblue"])

plt.title("Heart Disease Frequency for Sex")

plt.xlabel("0 = No Disease, 1=Disease")

plt.ylabel("Amount")

plt.legend(["Female", "Male"]);

plt.xticks(rotation=0);

#Create new figure

plt.figure(figsize=(10,6))

#Scatter with positive examples

```

[illegible]

```
cv=5,  
n_iter=20,  
verbose=True)
```

```
# Fit random hyperparameter search model for LogisticRegression  
rs_log_reg.fit(X_train, y_train)  
print('Best Params :- ', end="")
```