


Upload the Dataset

```
from google.colab import files
uploaded = files.upload()
```

 Choose Files

netflix\_titles.csv


- netflix\_titles.csv(text/csv) - 3399671 bytes, last modified: 5/8/2025 - 100% done

Saving netflix\_titles.csv to netflix\_titles.csv

Load the Dataset

```
import pandas as pd

df = pd.read_csv('/content/netflix_titles.csv')
df.head()
```



	show_id	type	title	director	cast	country	date_added	release_year	rating	duration	listed_ir
0	s1	Movie	Dick Johnson Is Dead	Kirsten Johnson	NaN	United States	September 25, 2021	2020	PG-13	90 min	Documentaries
1	s2	TV Show	Blood & Water	NaN	Ama Qamata, Khosi Ngema, Gail Mabalane, Thaban...	South Africa	September 24, 2021	2021	TV-MA	2 Seasons	International TV Shows, TV Dramas, TV Mysteries
2	s3	TV Show	Ganglands	Julien Leclercq	Sami Bouajila, Tracy Gotoas, Samuel Jouy, Nabi...	NaN	September 24, 2021	2021	TV-MA	1 Season	Crime TV Shows, International TV Shows, TV Act..
3	s4	TV Show	Jailbirds New Orleans	NaN	NaN	NaN	September 24, 2021	2021	TV-MA	1 Season	Docuseries, Reality TV
4	s5	TV Show	Kota Factory	NaN	Mayur More, Jitendra Kumar, Ranjan Raj, Alam K...	India	September 24, 2021	2021	TV-MA	2 Seasons	International TV Shows, Romantic TV Shows, TV ..

Next steps:

[Generate code with df](#)

 [View recommended plots](#)

[New interactive sheet](#)

Data Exploration

```
df.info()
df.describe(include='all')
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 8807 entries, 0 to 8806
Data columns (total 12 columns):
#   Column          Non-Null Count  Dtype
---  -
0   show_id         8807 non-null   object
1   type            8807 non-null   object
2   title           8807 non-null   object
3   director        6173 non-null   object
4   cast            7982 non-null   object
5   country         7976 non-null   object
6   date_added      8797 non-null   object
7   release_year    8807 non-null   int64
8   rating          8803 non-null   object
9   duration        8804 non-null   object
10  listed_in       8807 non-null   object
11  description      8807 non-null   object
dtypes: int64(1), object(11)
memory usage: 825.8+ KB
```

	show_id	type	title	director	cast	country	date_added	release_year	rating	duration	listed
count	8807	8807	8807	6173	7982	7976	8797	8807.000000	8803	8804	8
unique	8807	2	8807	4528	7692	748	1767	NaN	17	220	
top	s8807	Movie	Zubaan	Rajiv Chilaka	David Attenborough	United States	January 1, 2020	NaN	TV-MA	1 Season	Dram Internatic Mo
freq	1	6131	1	19	19	2818	109	NaN	3207	1793	
mean	NaN	NaN	NaN	NaN	NaN	NaN	NaN	2014.180198	NaN	NaN	1
std	NaN	NaN	NaN	NaN	NaN	NaN	NaN	8.819312	NaN	NaN	1
min	NaN	NaN	NaN	NaN	NaN	NaN	NaN	1925.000000	NaN	NaN	1
25%	NaN	NaN	NaN	NaN	NaN	NaN	NaN	2013.000000	NaN	NaN	1
50%	NaN	NaN	NaN	NaN	NaN	NaN	NaN	2017.000000	NaN	NaN	1
75%	NaN	NaN	NaN	NaN	NaN	NaN	NaN	2019.000000	NaN	NaN	1
max	NaN	NaN	NaN	NaN	NaN	NaN	NaN	2021.000000	NaN	NaN	1

Check for Missing Values and Duplicates

```
print(df.isnull().sum())
print("\nDuplicate Rows:", df.duplicated().sum())
```

```
show_id      0
type         0
title        0
director     2634
cast         825
country      831
date_added   10
release_year  0
rating       4
duration     3
listed_in    0
description  0
dtype: int64

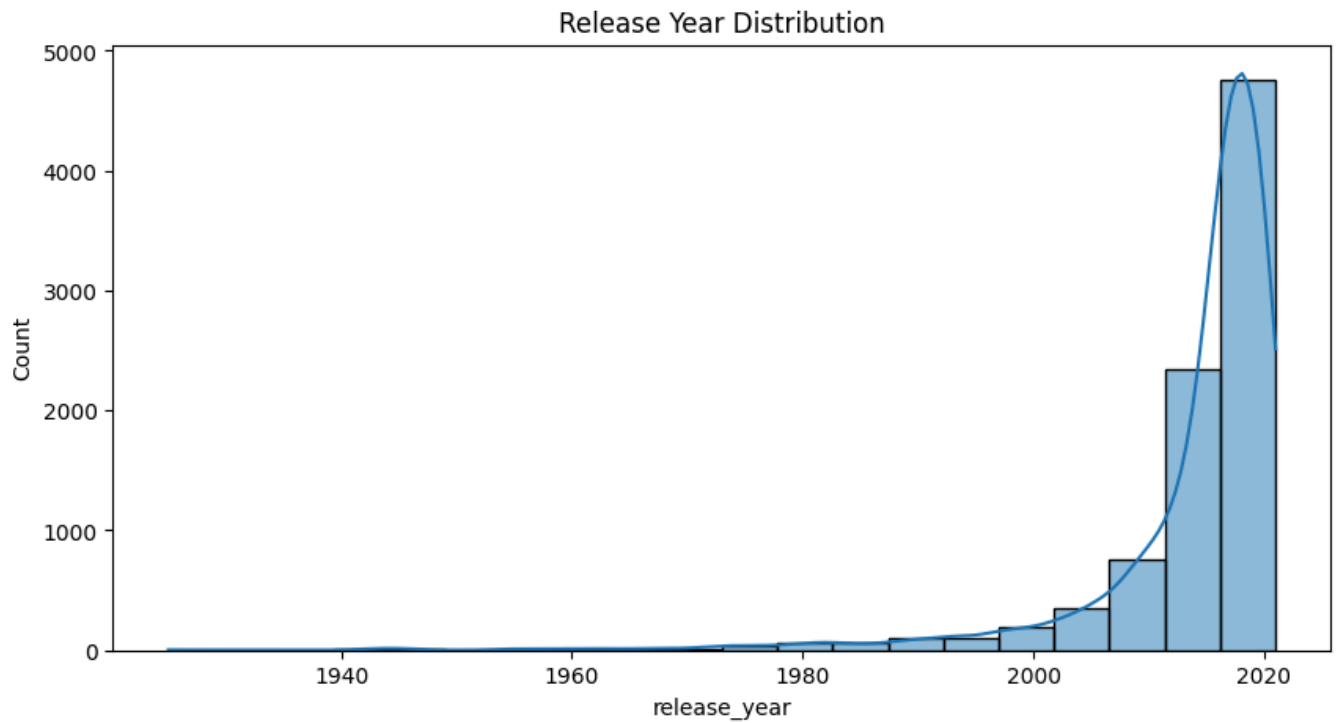
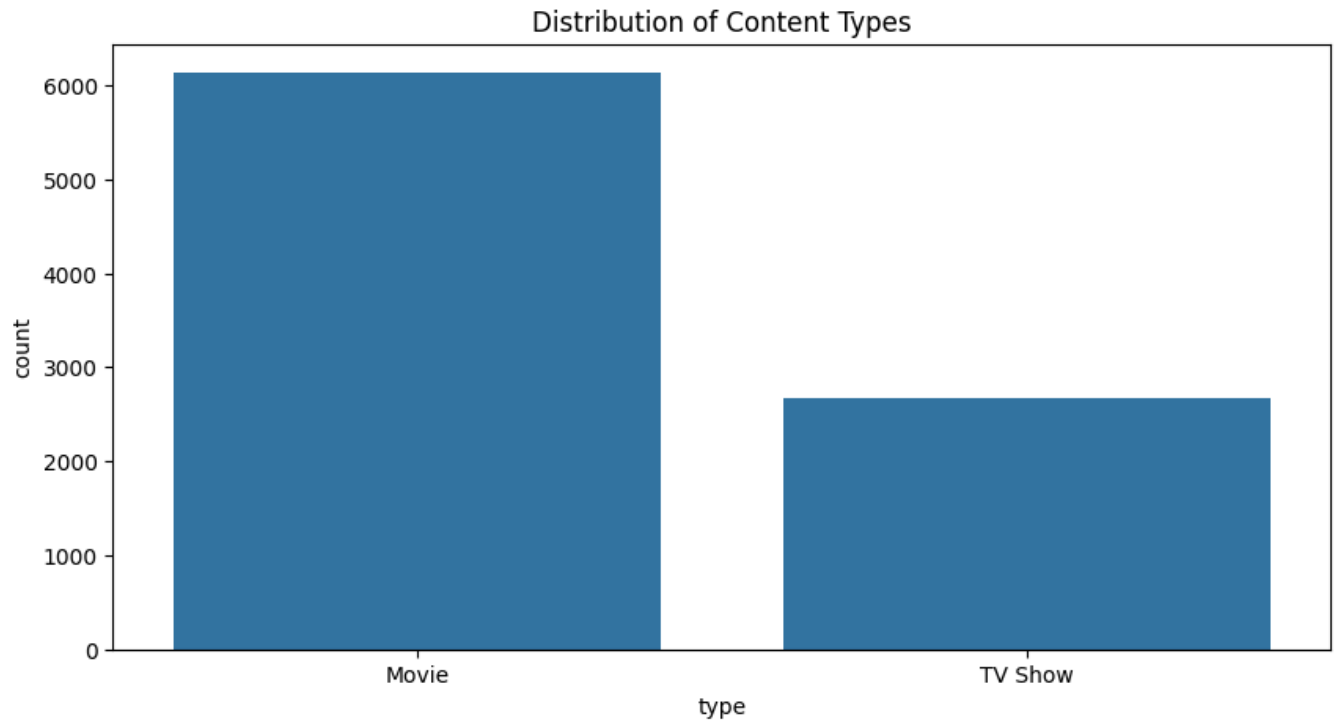
Duplicate Rows: 0
```

## Visualize a Few Features

```
import seaborn as sns
import matplotlib.pyplot as plt

plt.figure(figsize=(10,5))
sns.countplot(data=df, x='type')
plt.title("Distribution of Content Types")
plt.show()

plt.figure(figsize=(10,5))
sns.histplot(df['release_year'], bins=20, kde=True)
plt.title("Release Year Distribution")
plt.show()
```



#### Identify Target and Features

```
# This dataset doesn't have a traditional 'target' column.  
# For example, we could try to predict the 'type' (Movie/TV Show) based on other features.  
target = 'type'  
features = df.drop(columns=['type', 'show_id', 'title', 'description'])
```

#### Convert Categorical Columns to Numerical

```
df['date_added'] = pd.to_datetime(df['date_added'], format='mixed', errors='coerce')
df['year_added'] = df['date_added'].dt.year
df['month_added'] = df['date_added'].dt.month
df = df.drop(columns=['date_added'])
```

## One-Hot Encoding

```
df_encoded = pd.get_dummies(df, columns=['rating', 'country', 'listed_in'], drop_first=True)
```

## Feature Scaling

```
print(df_encoded.columns.tolist())
```

```
df_encoded.rename(columns={
    'Year Added': 'year_added',
    'Month Added': 'month_added'
}, inplace=True)
```

```
missing_cols = [col for col in numerical_cols if col not in df_encoded.columns]
if missing_cols:
    print(f"Missing columns: {missing_cols}")
else:
    df_encoded[numerical_cols] = scaler.fit_transform(df_encoded[numerical_cols])
```

```
➡ ['show_id', 'type', 'title', 'director', 'cast', 'date_added', 'release_year', 'duration', 'description', 'rating_
Missing columns: ['year_added', 'month_added']
```

## Train-Test Split

```
from sklearn.model_selection import train_test_split

X = df_encoded.drop(columns=['type'])
y = df['type']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
```

## Model Building

```
import pandas as pd
from sklearn.model_selection import train_test_split
from sklearn.ensemble import RandomForestClassifier

# Load your dataset
df = pd.read_csv('/content/netflix_titles.csv') # Adjust path as needed

# Save target separately
target = df['type']

# Drop unnecessary columns
df = df.drop(columns=['show_id', 'title', 'description', 'cast', 'director', 'type'], errors='ignore')

# Convert date_added
df['date_added'] = pd.to_datetime(df['date_added'], format='mixed', errors='coerce')
df['year_added'] = df['date_added'].dt.year
df['month_added'] = df['date_added'].dt.month
df = df.drop(columns=['date_added'])
```




```
# Drop rows with missing values in critical columns
df = df.dropna(subset=['rating', 'country', 'release_year'])

# Encode categorical features
df_encoded = pd.get_dummies(df, drop_first=True)

# Align target with the filtered DataFrame
target = target.loc[df_encoded.index].astype('category').cat.codes

# Train-test split
X_train, X_test, y_train, y_test = train_test_split(df_encoded, target, test_size=0.2, random_state=42)

# Train model
model = RandomForestClassifier()
model.fit(X_train, y_train)
```


 ▼ RandomForestClassifier  

RandomForestClassifier()

## Evaluation

```
from sklearn.metrics import classification_report, accuracy_score

y_pred = model.predict(X_test)
print(classification_report(y_test, y_pred))
print("Accuracy:", accuracy_score(y_test, y_pred))
```



	precision	recall	f1-score	support
0	0.98	0.99	0.98	1150
1	0.97	0.94	0.96	445
accuracy			0.98	1595
macro avg	0.98	0.97	0.97	1595
weighted avg	0.98	0.98	0.98	1595

Accuracy: 0.9768025078369906

## Make Predictions from New Input

```
sample = X_test.iloc[0:1]
prediction = model.predict(sample)
print("Predicted Type:", prediction[0])
```

 Predicted Type: 0

## Convert to DataFrame and Encode

```
# Step 1: Fit the scaler on your training data
scaler = StandardScaler()
scaler.fit(df_encoded[numerical_cols]) # df_encoded is your training dataset

# Step 2: Now you can transform new data safely
new_df[numerical_cols] = scaler.transform(new_df[numerical_cols])

import joblib
```

```

joblib.dump(model, 'model.pkl')
joblib.dump(scaler, 'scaler.pkl')

model = joblib.load('model.pkl')
scaler = joblib.load('scaler.pkl')

```

Predict the Final Grade (assuming 'type')

```

feature_columns = df_encoded.columns # Save this right after encoding


# Create your new input as a DataFrame
new_data = {
    'country': ['India'],
    'rating': ['TV-MA'],
    'release_year': [2021],
    'duration': ['1 Season'],
    'year_added': [2022],
    'month_added': [7],
    'listed_in': ['Dramas, International TV Shows']
}
new_df = pd.DataFrame(new_data)

# One-hot encode new_df using same logic
new_df_encoded = pd.get_dummies(new_df, drop_first=True)

# Align new input with training columns
new_df_encoded = new_df_encoded.reindex(columns=feature_columns, fill_value=0)

final_prediction = model.predict(new_df_encoded)
print("Final Prediction:", "TV Show" if final_prediction[0] == 1 else "Movie")

```


 Final Prediction: TV Show

Deployment - Building an Interactive App python Copy Edit

```

!pip install gradio
import gradio as gr

```

 Collecting gradio

```

  Downloading gradio-5.29.0-py3-none-any.whl.metadata (16 kB)
Collecting aiofiles<25.0,>=22.0 (from gradio)
  Downloading aiofiles-24.1.0-py3-none-any.whl.metadata (10 kB)
Requirement already satisfied: anyio<5.0,>=3.0 in /usr/local/lib/python3.11/dist-packages (from gradio) (4.9.0)
Collecting fastapi<1.0,>=0.115.2 (from gradio)
  Downloading fastapi-0.115.12-py3-none-any.whl.metadata (27 kB)
Collecting ffmpy (from gradio)
  Downloading ffmpy-0.5.0-py3-none-any.whl.metadata (3.0 kB)
Collecting gradio-client==1.10.0 (from gradio)
  Downloading gradio_client-1.10.0-py3-none-any.whl.metadata (7.1 kB)
Collecting groovy~=0.1 (from gradio)
  Downloading groovy-0.1.2-py3-none-any.whl.metadata (6.1 kB)
Requirement already satisfied: httpx>=0.24.1 in /usr/local/lib/python3.11/dist-packages (from gradio) (0.28.1)
Requirement already satisfied: huggingface-hub>=0.28.1 in /usr/local/lib/python3.11/dist-packages (from gradio) (3.1.6)
Requirement already satisfied: jinja2<4.0 in /usr/local/lib/python3.11/dist-packages (from gradio) (3.1.6)
Requirement already satisfied: markupsafe<4.0,>=2.0 in /usr/local/lib/python3.11/dist-packages (from gradio) (3.1.0)
Requirement already satisfied: numpy<3.0,>=1.0 in /usr/local/lib/python3.11/dist-packages (from gradio) (2.0.2)
Requirement already satisfied: orjson~=3.0 in /usr/local/lib/python3.11/dist-packages (from gradio) (3.10.18)

```

Requirement already satisfied: packaging in /usr/local/lib/python3.11/dist-packages (from gradio) (24.2)  
 Requirement already satisfied: pandas<3.0,>=1.0 in /usr/local/lib/python3.11/dist-packages (from gradio) (2.2.2)  
 Requirement already satisfied: pillow<12.0,>=8.0 in /usr/local/lib/python3.11/dist-packages (from gradio) (11.2.  
 Requirement already satisfied: pydantic<2.12,>=2.0 in /usr/local/lib/python3.11/dist-packages (from gradio) (2.1  
 Collecting pydub (from gradio)  
 Downloading pydub-0.25.1-py2.py3-none-any.whl.metadata (1.4 kB)  
 Collecting python-multipart<=0.0.18 (from gradio)  
 Downloading python\_multipart-0.0.20-py3-none-any.whl.metadata (1.8 kB)  
 Requirement already satisfied: pyyaml<7.0,>=5.0 in /usr/local/lib/python3.11/dist-packages (from gradio) (6.0.2)  
 Collecting ruff<=0.9.3 (from gradio)  
 Downloading ruff-0.11.8-py3-none-manylinux\_2\_17\_x86\_64.manylinux2014\_x86\_64.whl.metadata (25 kB)  
 Collecting safehttpx<0.2.0,>=0.1.6 (from gradio)  
 Downloading safehttpx-0.1.6-py3-none-any.whl.metadata (4.2 kB)  
 Collecting semantic-version<=2.0 (from gradio)  
 Downloading semantic\_version-2.10.0-py2.py3-none-any.whl.metadata (9.7 kB)  
 Collecting starlette<1.0,>=0.40.0 (from gradio)  
 Downloading starlette-0.46.2-py3-none-any.whl.metadata (6.2 kB)  
 Collecting tomlkit<0.14.0,>=0.12.0 (from gradio)  
 Downloading tomlkit-0.13.2-py3-none-any.whl.metadata (2.7 kB)  
 Requirement already satisfied: typer<1.0,>=0.12 in /usr/local/lib/python3.11/dist-packages (from gradio) (0.15.3  
 Requirement already satisfied: typing-extensions<=4.0 in /usr/local/lib/python3.11/dist-packages (from gradio) (  
 Collecting uvicorn<=0.14.0 (from gradio)  
 Downloading uvicorn-0.34.2-py3-none-any.whl.metadata (6.5 kB)  
 Requirement already satisfied: fsspec in /usr/local/lib/python3.11/dist-packages (from gradio-client==1.10.0->gr  
 Requirement already satisfied: websockets<16.0,>=10.0 in /usr/local/lib/python3.11/dist-packages (from gradio-cl  
 Requirement already satisfied: idna>=2.8 in /usr/local/lib/python3.11/dist-packages (from anyio<5.0,>=3.0->gradi  
 Requirement already satisfied: sniffio>=1.1 in /usr/local/lib/python3.11/dist-packages (from anyio<5.0,>=3.0->gr  
 Requirement already satisfied: certifi in /usr/local/lib/python3.11/dist-packages (from httpx>=0.24.1->gradio) (  
 Requirement already satisfied: httpcore==1.\* in /usr/local/lib/python3.11/dist-packages (from httpx>=0.24.1->gra  
 Requirement already satisfied: h11>=0.16 in /usr/local/lib/python3.11/dist-packages (from httpcore==1.\*->httpx>=  
 Requirement already satisfied: filelock in /usr/local/lib/python3.11/dist-packages (from huggingface-hub>=0.28.1  
 Requirement already satisfied: requests in /usr/local/lib/python3.11/dist-packages (from huggingface-hub>=0.28.1  
 Requirement already satisfied: tqdm>=4.42.1 in /usr/local/lib/python3.11/dist-packages (from huggingface-hub>=0.  
 Requirement already satisfied: python-dateutil>=2.8.2 in /usr/local/lib/python3.11/dist-packages (from pandas<3.  
 Requirement already satisfied: pytz>=2020.1 in /usr/local/lib/python3.11/dist-packages (from pandas<3.0,>=1.0->g  
 Requirement already satisfied: tzdata>=2022.7 in /usr/local/lib/python3.11/dist-packages (from pandas<3.0,>=1.0-  
 Requirement already satisfied: annotated-types>=0.6.0 in /usr/local/lib/python3.11/dist-packages (from pydantic<  
 Requirement already satisfied: pydantic-core==2.33.2 in /usr/local/lib/python3.11/dist-packages (from pydantic<2

## Create a Prediction Function

```
def predict_type(release_year, year_added, month_added):
    input_data = pd.DataFrame([[release_year, year_added, month_added]],
                               columns=['release_year', 'year_added', 'month_added'])
    input_data[numerical_cols] = scaler.transform(input_data[numerical_cols])
    return model.predict(input_data)[0]
```

## Create the Gradio Interface

```
import gradio as gr
import pandas as pd

# Save the feature columns after training
feature_columns = df_encoded.columns

# Define the prediction function
def predict_netflix_type(country, rating, release_year, duration, year_added, month_added, listed_in):
    # Create a new DataFrame for input
    new_input = pd.DataFrame({
        'country': [country],
        'rating': [rating],
        'release_year': [release_year],
        'duration': [duration],
        'year_added': [year_added],
        'month added': [month added].
```



```

        'listed_in': [listed_in]
    })

# Encode the input similar to training
new_encoded = pd.get_dummies(new_input, drop_first=True)

# Reindex to match training features
new_encoded = new_encoded.reindex(columns=feature_columns, fill_value=0)

# Predict
prediction = model.predict(new_encoded)[0]
return "TV Show" if prediction == 1 else "Movie"

# Create the Gradio interface
interface = gr.Interface(
    fn=predict_netflix_type,
    inputs=[
        gr.Textbox(label="Country"),
        gr.Textbox(label="Rating (e.g. TV-MA)"),
        gr.Number(label="Release Year"),
        gr.Textbox(label="Duration (e.g. 1 Season or 90 min)"),
        gr.Number(label="Year Added"),
        gr.Number(label="Month Added"),
        gr.Textbox(label="Genres (Listed In)"),
    ],
    outputs="text",
    title="Netflix Content Type Predictor",
    description="Predict whether the content is a Movie or TV Show based on input features."
)

interface.launch()

```

➡ It looks like you are running Gradio on a hosted Jupyter notebook. For the Gradio app to work, sharing must be e

Colab notebook detected. To show errors in colab notebook, set debug=True in launch()

\* Running on public URL: <https://b4e268096ec7ded9d5.gradio.live>

This share link expires in 1 week. For free permanent hosting and GPU upgrades, run `gradio deploy` from the termi