

PROJECT PROPOSAL

PREDICTING MOVIE SUCCESS USING LOGISTIC REGRESSION



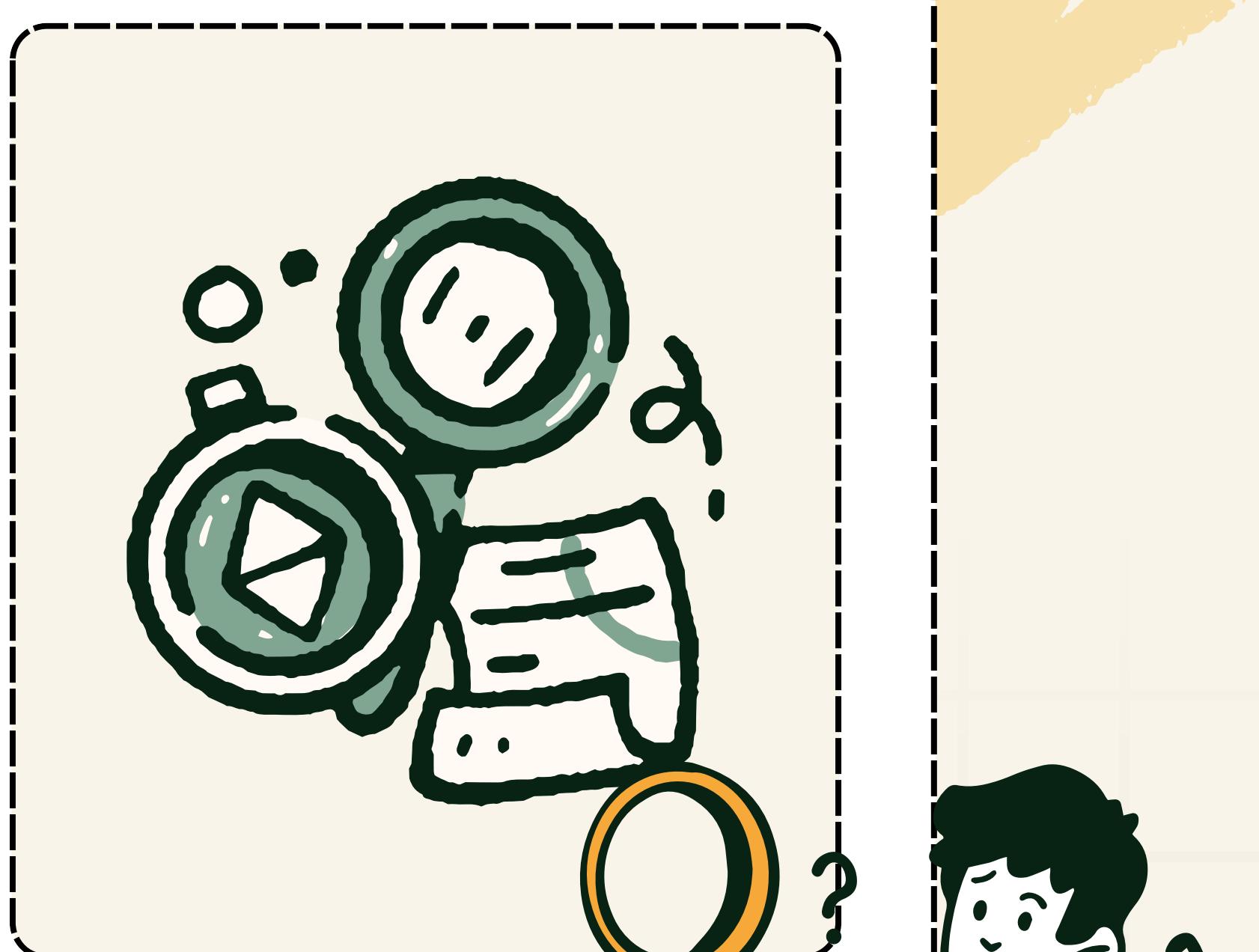
PRESENTED BY

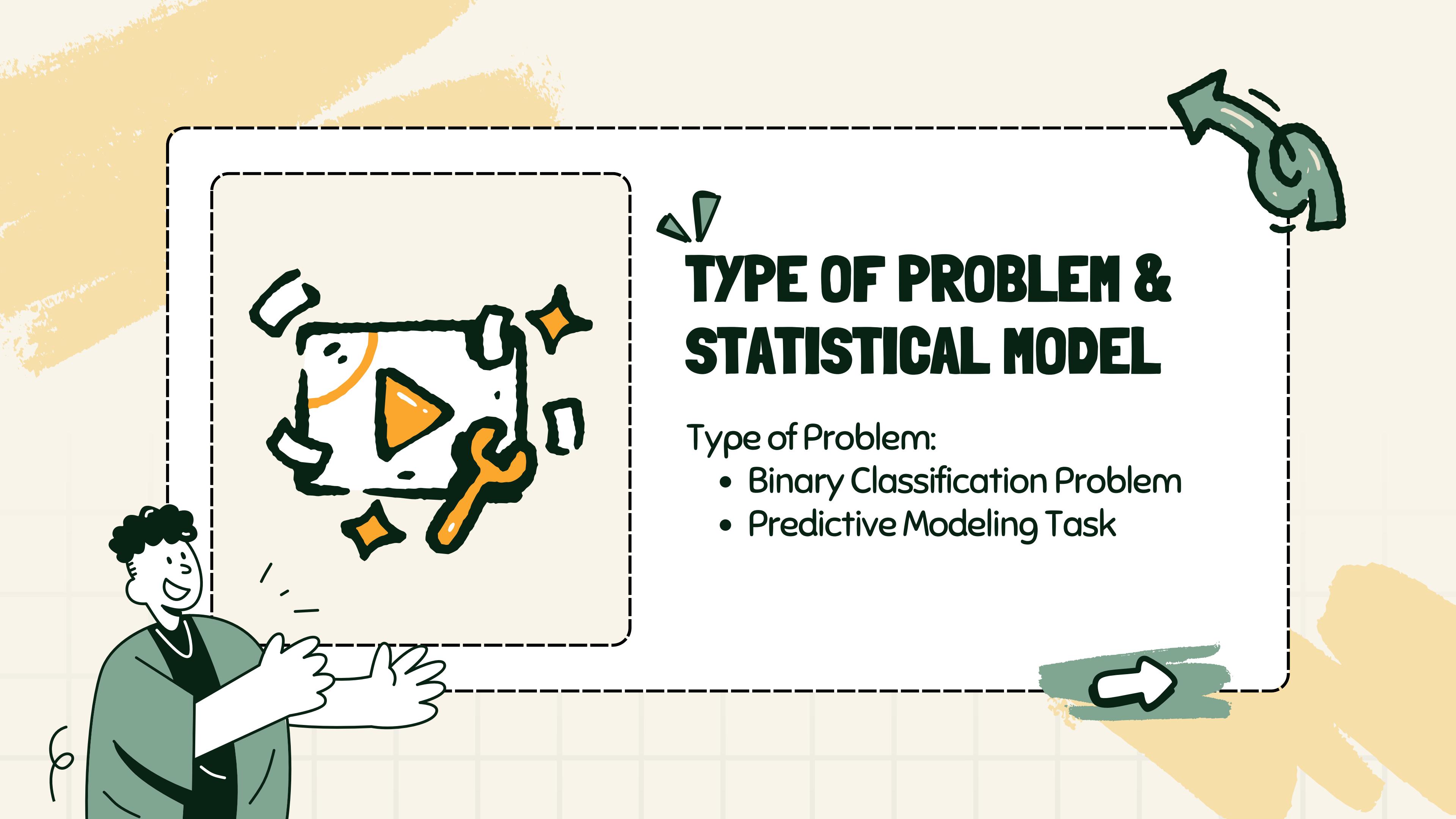
PRASHANT KOIRALA
AASKA KOIRALA
AISHMITA YONZAN



PROBLEM STATEMENT

Can we predict whether a movie will be successful or not based on its metadata (e.g., budget, runtime, genre, and other pre-release attributes)?





TYPE OF PROBLEM & STATISTICAL MODEL

Type of Problem:

- Binary Classification Problem
- Predictive Modeling Task

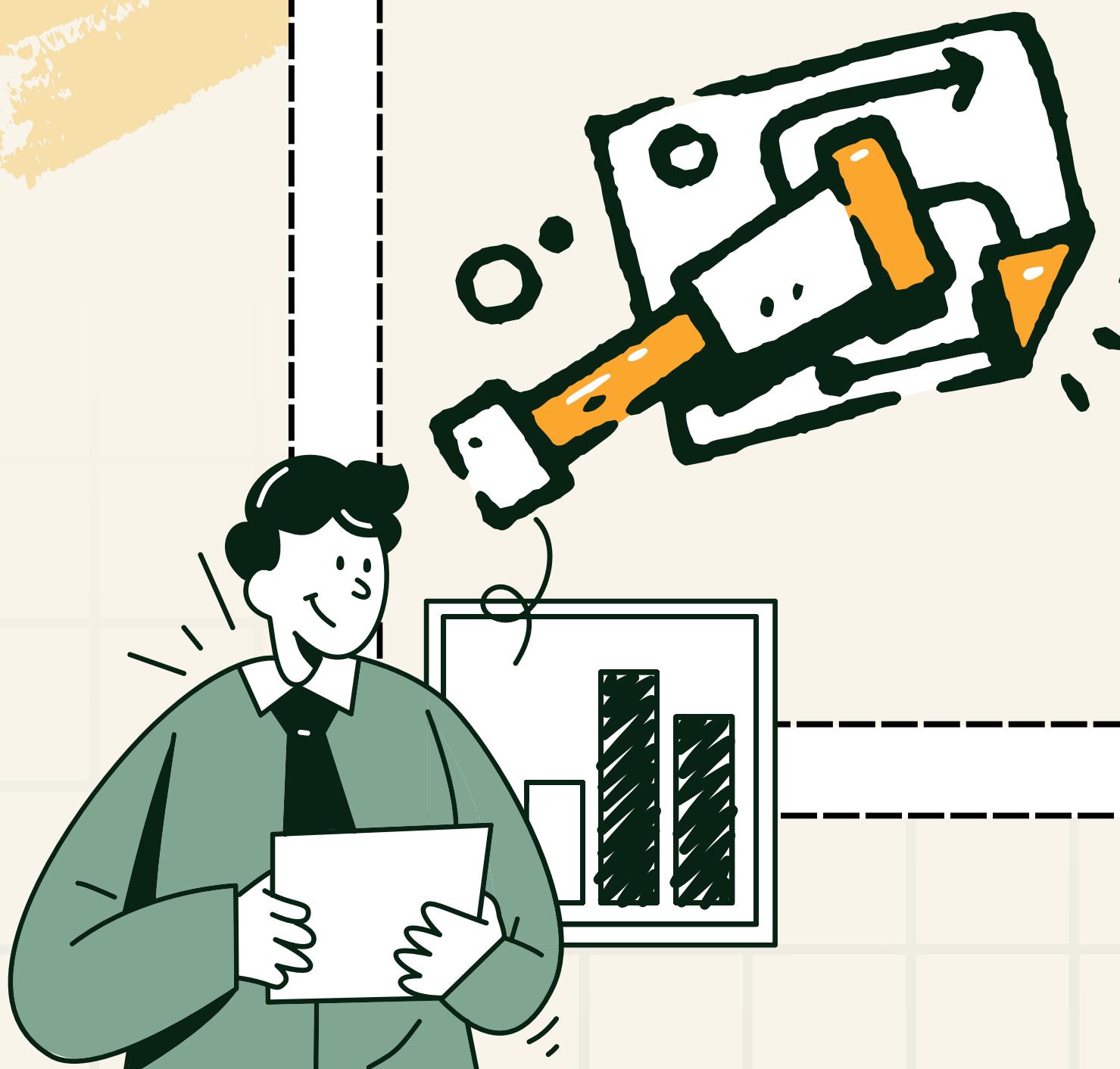
STATISTICAL TECHNIQUE CHOSEN

Logistic Regression



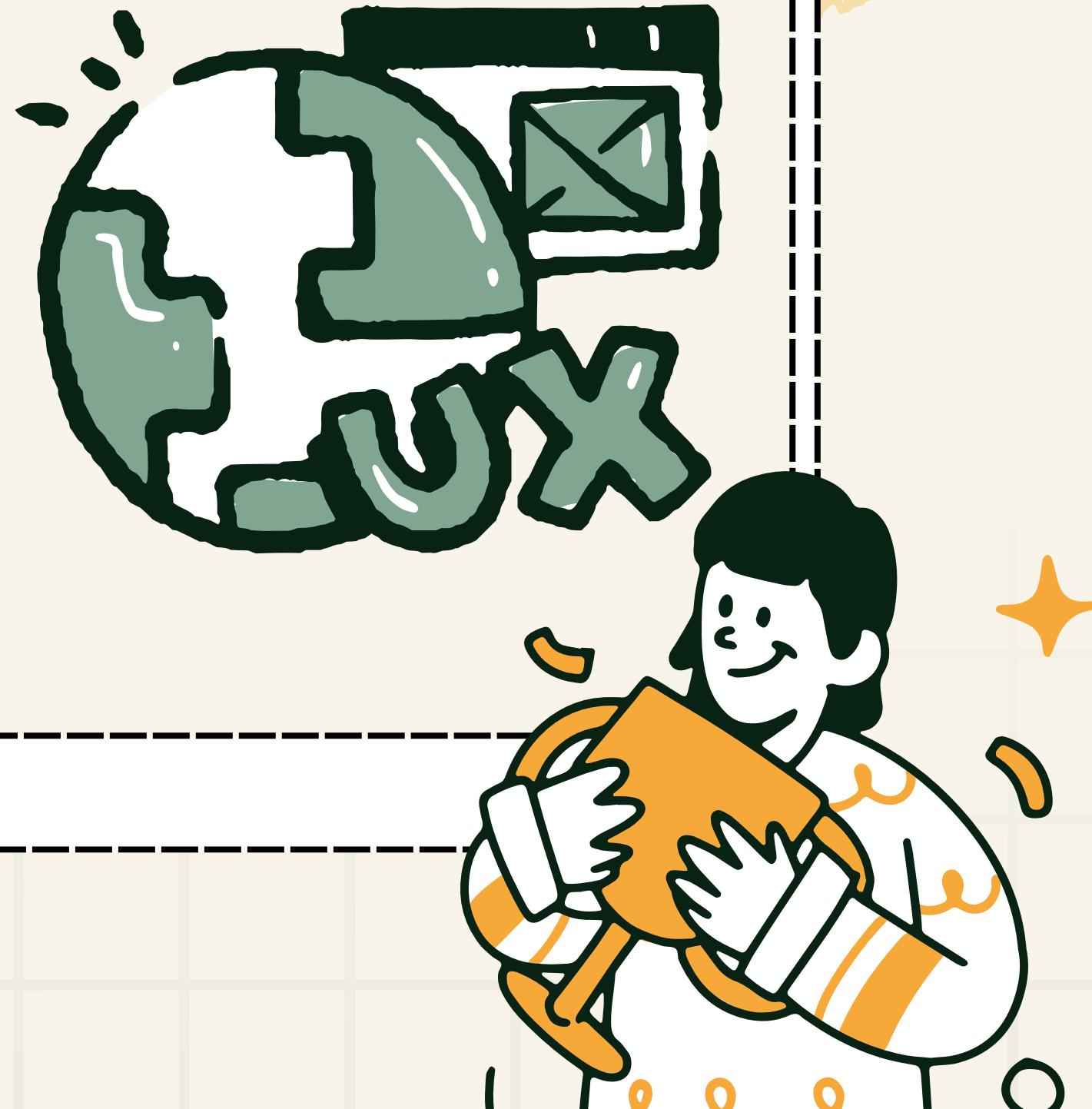
LOGISTIC REGRESSION

- A statistical model used for binary classification.
- Estimates the probability that a given input belongs to a particular category.
- Uses the sigmoid function to output a value between 0 and 1.
- Predicts outcomes like “Success (1)” or “Not Success (0)”.



KEY FEATURES

- It is ideal for binary outcomes (success vs. not success).
- It provides clear coefficients to understand how each feature (e.g., budget) impacts the probability of success.
- Suitable for a student project, as it's straightforward to implement and interpret.



LITERATURE REVIEW & DATASET SELECTION

DA, O., OM, S., AK, F., O, A., A, O., A, O., A, W., & M, Y. (2021). Movie success prediction using data mining. *British Journal of Computer Networking and Information Technology*, 4(2), 22–30. <https://doi.org/10.52589/bjcnit-cqocirec>

Zheng, Y. (2024). Predicting movie box office based on machine learning, deep learning, and statistical methods. *Applied and Computational Engineering*, 94(1), 20–32. <https://doi.org/10.54254/2755-2721/94/2024melb0069>

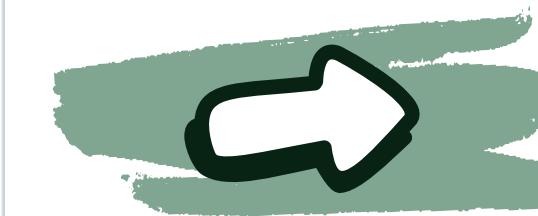
Subramaniyaswamy, V., Vaibhav, M. V., Prasad, R. V., & Logesh, R. (2017). Predicting movie box office success using multiple regression and SVM. *Journal of Science & Statistics*, 182–186. <https://doi.org/10.1109/iss1.2017.8389394>

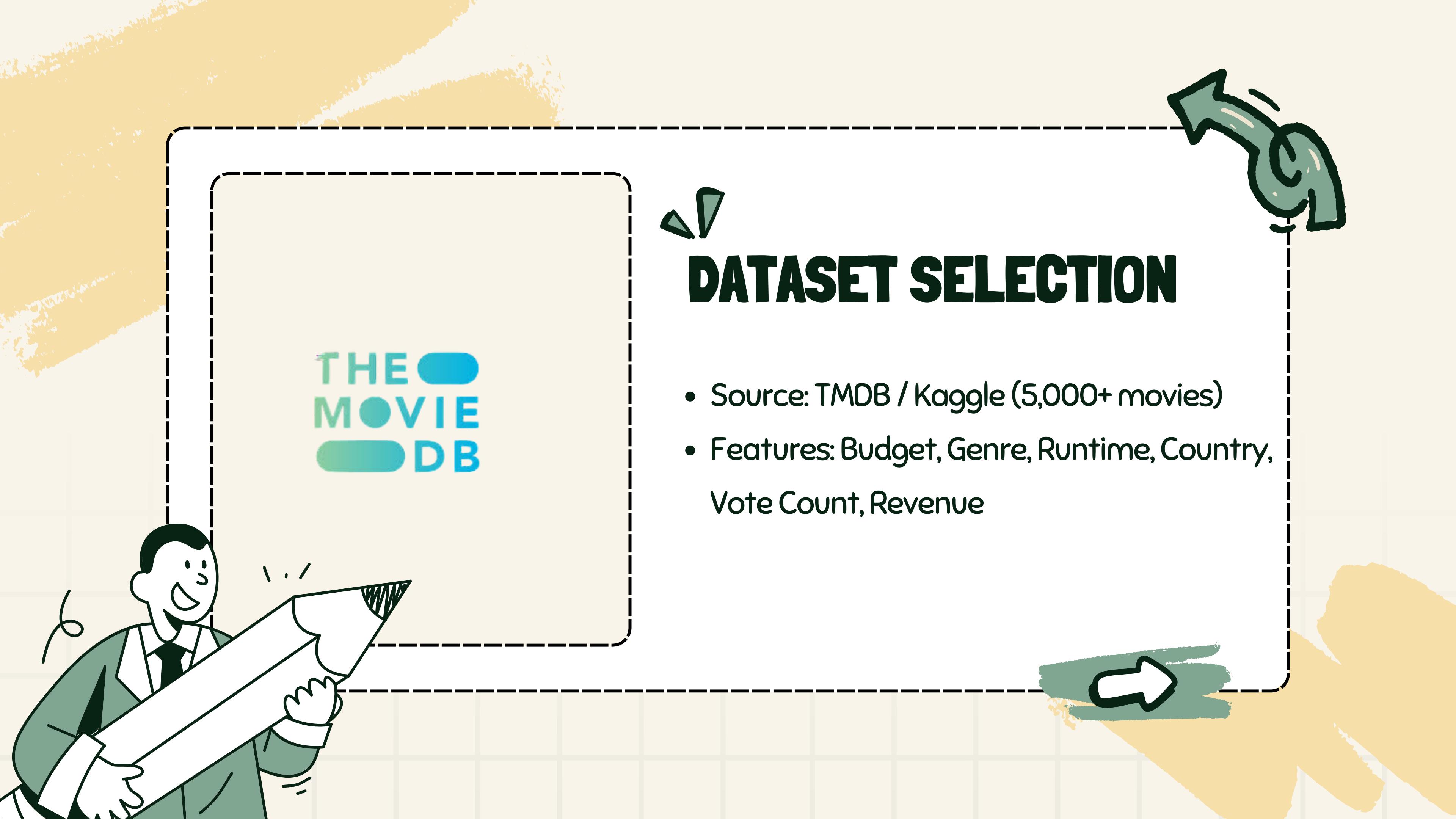
Velingkar, G., Varadarajan, R., Lanka, S., & M, A. K. (2022). Movie Box-Office success Prediction using Machine Learning. 2022 Second International Conference on Power, Control and Computing Technologies (ICPC2T), 9, 1–6. <https://doi.org/10.1109/icpc2t53885.2022.9776798>

Mr, N. V., Mr, P. M., & Pb, S. B. (2014). Predicting movie success based on IMDB data. *International Journal of Business Intelligents*, 3(2). <https://doi.org/10.20894/ijbi.105.003.002.004>

BRIEF OVERVIEW OF LITERATURE REVIEW

<p>Proceedings Article • 10.1109/ISS1.2017.8389394</p> <p><input type="checkbox"/> 3. Predicting movie box office success using multiple regression and SVM</p> <p>V. Subramaniyaswamy, M. Viginesh Vaibhav, R. Vishnu Prasad +1 more</p> <p>1 Dec 2017</p> <p>66 26 Request PDF Podcast Chat</p> <p> 66</p>	<p>The paper you referenced is not the same as the one discussed here. This research focuses on predicting movie box office success using multiple regression and SVM, analyzing factors like trailer views, Wikipedia page views, and critic ratings.</p>
<p>Proceedings Article • 10.1109/ICPC2T53885.2022.9776798</p> <p><input type="checkbox"/> 4. Movie Box-Office Success Prediction Using Machine Learning</p> <p>Gaurang Velingkar, Rakshita Varadarajan, Sidharth Lanka +1 more</p> <p>1 Mar 2022</p> <p>66 4 Request PDF Podcast Chat</p> <p> 66</p>	<p>The paper focuses on predicting movie box office success using machine learning, specifically employing a Random Forest Regression model to analyze factors like genre, budget, and star power, aiding investors in making informed decisions before a movie's release.</p>
<p> . Proceedings Article • 10.1109/icpc2t53885.2022.9776798</p> <p><input type="checkbox"/> 5. Movie Box-Office Success Prediction Using Machine Learning</p> <p>1 Mar 2022</p> <p>66 8 PDF Podcast Chat </p>	<p>The paper focuses on predicting movie box office success using machine learning, specifically employing a Random Forest Regression model to analyze factors like genre, budget, and star power, aiding investors in making informed decisions before a movie's release.</p>





DATASET SELECTION

- Source: TMDB / Kaggle (5,000+ movies)
- Features: Budget, Genre, Runtime, Country, Vote Count, Revenue

EXPLORATORY DATA ANALYSIS (EDA)

tmdb_5000_movies

Search (Cmd + Ctrl + U)

Home Insert Draw Page Layout Formulas Data Review View Acrobat

Paste Calibri (Body) 12 A General Conditional Formatting Insert Delete Sort & Filter Add-ins

B I U A Format as Table Cell Styles Find & Select Create PDF and share link

U1 x ✓ ffx

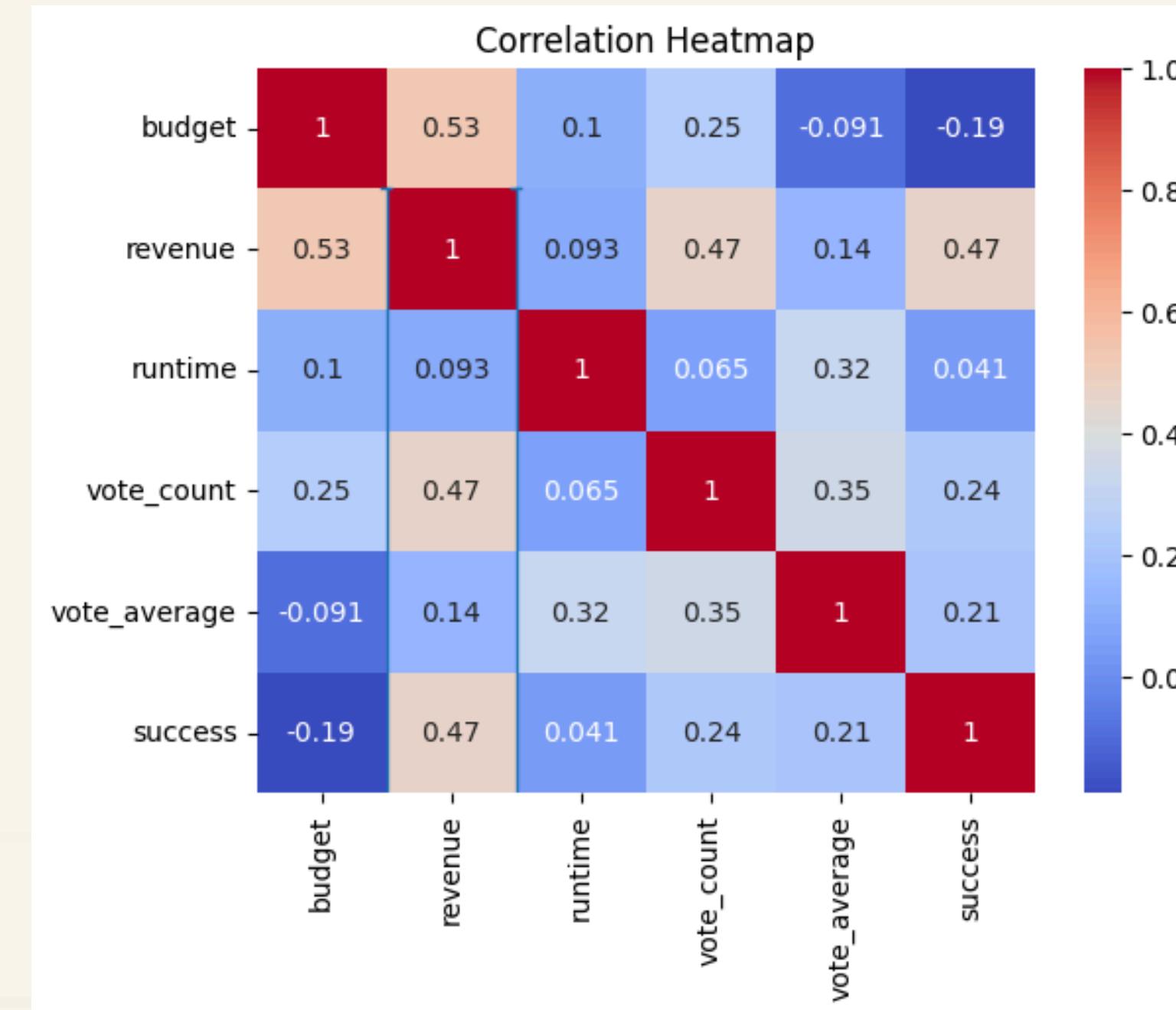
	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	V
1	budget	genres	homepage	id	keywords	original_lang	original_title	overview	popularity	production_c	production_c	release_date	revenue	runtime	spoken_lang	status	tagline	title	vote_average	vote_count		
2	237000000	[{"id": 28, "n http://www.	19995	[{"id": 1463, "en	Avatar	In the 22nd c	150.437577	[{"name": "I	[{"iso_3166_	12/10/09	2787965087		162	[{"iso_639_1 Released	Enter the W	Avatar		7.2	11800			
3	300000000	[{"id": 12, "n http://disney	285	[{"id": 270, "en	Pirates of the	Captain Barb	139.082615	[{"name": "V	[{"iso_3166_	5/19/07	961000000		169	[{"iso_639_1 Released	At the end of	Pirates of the		6.9	4500			
4	245000000	[{"id": 28, "n http://www.	206647	[{"id": 470, "en	Spectre	A cryptic me	107.376788	[{"name": "C	[{"iso_3166_	10/26/15	880674609		148	[{"iso_639_1 Released	A Plan No Or	Spectre		6.3	4466			
5	250000000	[{"id": 28, "n http://www.	49026	[{"id": 849, "en	The Dark Kni	Following the	112.31295	[{"name": "L	[{"iso_3166_	7/16/12	1084939099		165	[{"iso_639_1 Released	The Legend E	The Dark Kni		7.6	9106			
6	260000000	[{"id": 28, "n http://movie	49529	[{"id": 818, "en	John Carter	John Carter i	43.926995	[{"name": "V	[{"iso_3166_	3/7/12	284139100		132	[{"iso_639_1 Released	Lost in our w	John Carter		6.1	2124			
7	258000000	[{"id": 14, "n http://www.	559	[{"id": 851, "en	Spider-Man	The seeming	115.699814	[{"name": "C	[{"iso_3166_	5/1/07	890871626		139	[{"iso_639_1 Released	The battle w	Spider-Man		5.9	3576			
8	260000000	[{"id": 16, "n http://disney	38757	[{"id": 1562, "en	Tangled	When the kir	48.681969	[{"name": "V	[{"iso_3166_	11/24/10	591794936		100	[{"iso_639_1 Released	They're takin	Tangled		7.4	3330			
9	280000000	[{"id": 28, "n http://marve	99861	[{"id": 8828, "en	Avengers: Ag	When Tony S	134.279229	[{"name": "V	[{"iso_3166_	4/22/15	1405403694		141	[{"iso_639_1 Released	A New Age	Avengers: Ag		7.3	6767			
10	250000000	[{"id": 12, "n http://harryp	767	[{"id": 616, "en	Harry Potter	As Harry beg	98.885637	[{"name": "V	[{"iso_3166_	7/7/09	933959197		153	[{"iso_639_1 Released	Dark Secrets	Harry Potter		7.4	5293			
11	250000000	[{"id": 28, "n http://www.	209112	[{"id": 849, "en	Batman v Su	Fearing the a	155.790452	[{"name": "C	[{"iso_3166_	3/23/16	873260194		151	[{"iso_639_1 Released	Justice or rev	Batman v Su		5.7	7004			
12	270000000	[{"id": 12, "n http://www.	1452	[{"id": 83, "en	Superman R	Superman re	57.925623	[{"name": "C	[{"iso_3166_	6/28/06	391081192		154	[{"iso_639_1 Released	Superman Re			5.4	1400			
13	200000000	[{"id": 12, "n http://www.	10764	[{"id": 627, "en	Quantum of	Quantum of	107.928811	[{"name": "E	[{"iso_3166_	10/30/08	586090727		106	[{"iso_639_1 Released	For love, for	Quantum of		6.1	2965			
14	200000000	[{"id": 12, "n http://disney	58	[{"id": 616, "en	Pirates of the	Captain Jack	145.847379	[{"name": "V	[{"iso_3166_	6/20/06	1065659812		151	[{"iso_639_1 Released	Jack is back!	Pirates of the		7	5246			
15	255000000	[{"id": 28, "n http://disney	57201	[{"id": 1556, "en	The Lone Rai	The Texas Ra	49.046956	[{"name": "V	[{"iso_3166_	7/3/13	89289910		149	[{"iso_639_1 Released	Never Take C	The Lone Rai		5.9	2311			
16	225000000	[{"id": 28, "n http://www.	49521	[{"id": 83, "en	Man of Steel	A young boy	99.398009	[{"name": "L	[{"iso_3166_	6/12/13	662845518		143	[{"iso_639_1 Released	You will belie	Man of Steel		6.5	6359			
17	225000000	[{"id": 12, "name": "Adver	2454	[{"id": 818, "en	The Chronicl	One year aft	53.978602	[{"name": "V	[{"iso_3166_	5/15/08	419651413		150	[{"iso_639_1 Released	Hope has a n	The Chronicle		6.3	1630			
18	220000000	[{"id": 878, "n http://marve	24428	[{"id": 242, "en	The Avenger	When an un	144.448633	[{"name": "P	[{"iso_3166_	4/25/12	1519557910		143	[{"iso_639_1 Released	Some assem	The Avenger		7.4	11776			
19	380000000	[{"id": 12, "n http://disney	1865	[{"id": 658, "en	Pirates of th	Captain Jack	135.413856	[{"name": "V	[{"iso_3166_	5/14/11	1045713802		136	[{"iso_639_1 Released	Live Forever	Pirates of the		6.4	4948			
20	225000000	[{"id": 28, "n http://www.	41154	[{"id": 4379, "en	Men in Black	Agents J (Wi	52.035179	[{"name": "A	[{"iso_3166_	5/23/12	624026776		106	[{"iso_639_1 Released	They are bac	Men in Black		6.2	4160			
21	250000000	[{"id": 28, "n http://www.	122917	[{"id": 417, "en	The Hobbit: T	Immediately	120.965743	[{"name": "V	[{"iso_3166_	12/10/14	956019788		144	[{"iso_639_1 Released	Witness the	The Hobbit: T		7.1	4760			
22	215000000	[{"id": 28, "n http://www.	1930	[{"id": 1872, "en	The Amazing	Peter Parker	89.866276	[{"name": "C	[{"iso_3166_	6/27/12	752215857		136	[{"iso_639_1 Released	The untold st	The Amazing		6.5	6586			
23	200000000	[{"id": 28, "n http://www.	20662	[{"id": 4147, "en	Robin Hood	When soldie	37.668301	[{"name": "I	[{"iso_3166_	5/12/10	310669540		140	[{"iso_639_1 Released	Rise and rise	Robin Hood		6.2	1398			
24	250000000	[{"id": 12, "n http://www.	57158	[{"id": 603, "en	The Hobbit: T	The Dwarves	94.370564	[{"name": "V	[{"iso_3166_	12/11/13	958400000		161	[{"iso_639_1 Released	Beyond dark	The Hobbit: T		7.6	4524			
25	180000000	[{"id": 12, "n http://www.	2268	[{"id": 392, "en	The Golden C	After overha	42.990906	[{"name": "N	[{"iso_3166_	12/4/07	372234864		113	[{"iso_639_1 Released	There are w	The Golden C		5.8	1303			
26	207000000	[{"id": 12, "name": "Adver	254	[{"id": 774, "en	King Kong	In 1933 New	61.22601	[{"name": "V	[{"iso_3166_	12/14/05	550000000		187	[{"iso_639_1 Released	The eighth w	King Kong		6.6	2337			
27	200000000	[{"id": 18, "n http://www.	597	[{"id": 2580, "en	Titanic	84 years late	100.025899	[{"name": "P	[{"iso_3166_	11/18/97	1845034188		194	[{"iso_639_1 Released	Nothing on E	Titanic		7.5	7562			
28	250000000	[{"id": 12, "n http://marve	271110	[{"id": 393, "en	Captain Ame	Following th	198.372395	[{"name": "S	[{"iso_3166_	4/27/16	1153304495		147	[{"iso_639_1 Released	Divided We	Captain Ame		7.1	7241			
29	209000000	[{"id": 53, "name": "Thriller	44833	[{"id": 1721, "en	Battleship	When manki	64.928382	[{"name": "L	[{"iso_3166_	4/11/12	303025485		131	[{"iso_639_1 Released	The Battle	f Battleship		5.5	2114			
30	150000000	[{"id": 28, "n http://www.	135397	[{"id": 1299, "en	Jurassic Wor	Twenty-two	418.708552	[{"name": "L	[{"iso_3166_	6/9/15	1513528810		124	[{"iso_639_1 Released	The park is o	Jurassic Wor		6.5	8662			
31	200000000	[{"id": 28, "n http://www.	37724	[{"id": 470, "en	Skyfall	When Bond's	93.004993	[{"name": "C	[{"iso_3166_	10/25/12	1108561013		143	[{"iso_639_1 Released	Think on you	Skyfall		6.9	7604			
32	200000000	[{"id": 28, "n http://www.	558	[{"id": 851, "en	Spider-Man	Peter Parker	35.149586	[{"name": "C	[{"iso_3166_	6/25/04	783766341		127	[{"iso_639_1 Released	There's a hei	Spider-Man		6.7	4321			
33	200000000	[{"id": 28, "n http://marve	68721	[{"id": 949, "en	Iron Man 3	When Tony S	77.68208	[{"name": "N	[{"iso_3166_	4/18/13	1215439994		130	[{"iso_639_1 Released	Unleash the	Iron Man 3		6.8	8806			
34	200000000	[{"id": 10751, "n http://disney	12155	[{"id": 818, "en	Alice in Won																	

EXPLORATORY DATA ANALYSIS (EDA)

Features Used for Analysis

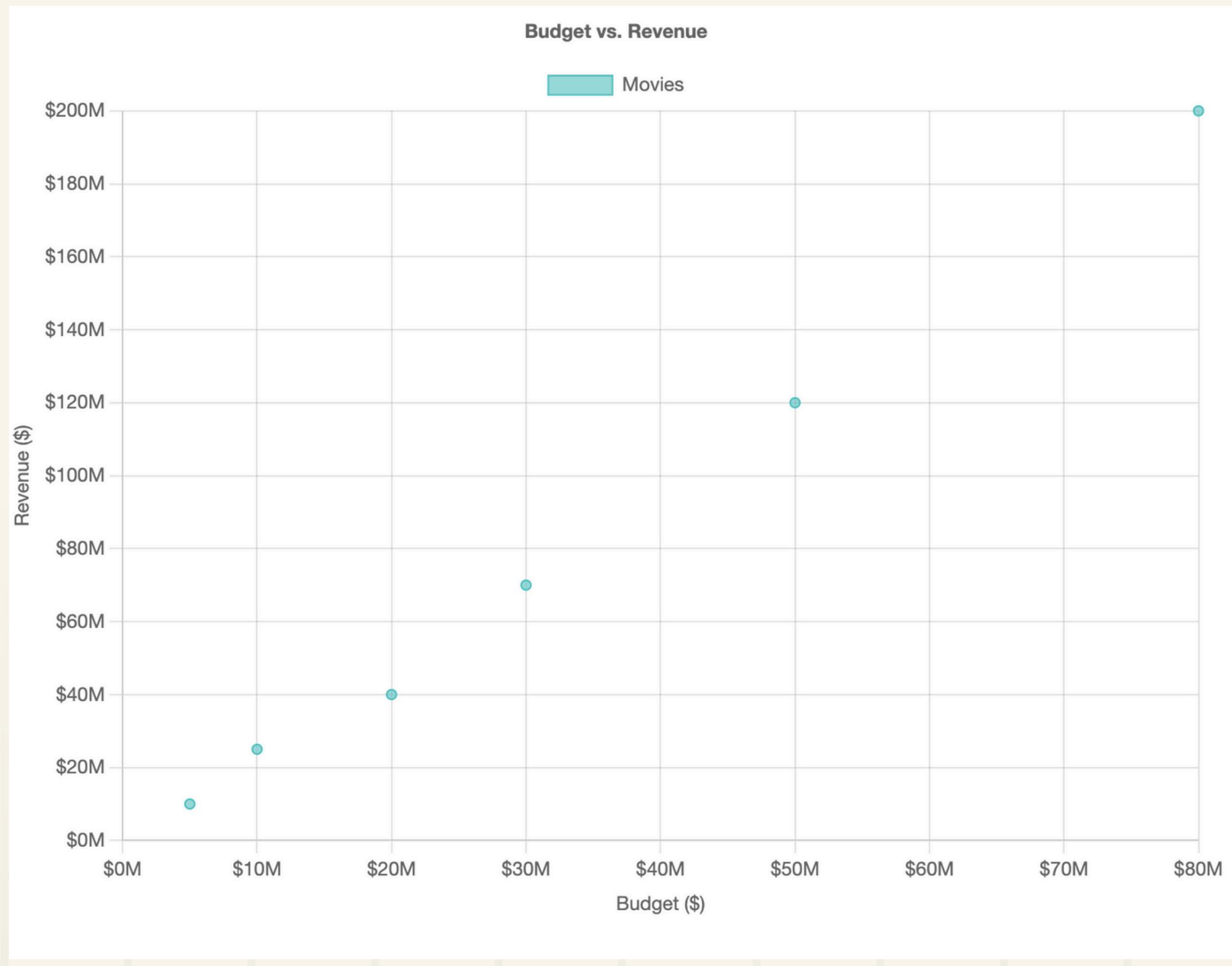
- **Numerical:** budget, revenue, runtime, year, vote_count (and possibly vote_average if not dropped).
- **Categorical:** certification_US, genre, country.
- **Target:** success (True/False, where True means $\text{revenue} \geq 2 \times \text{budget}$).

WHY THOSE FEATURES?



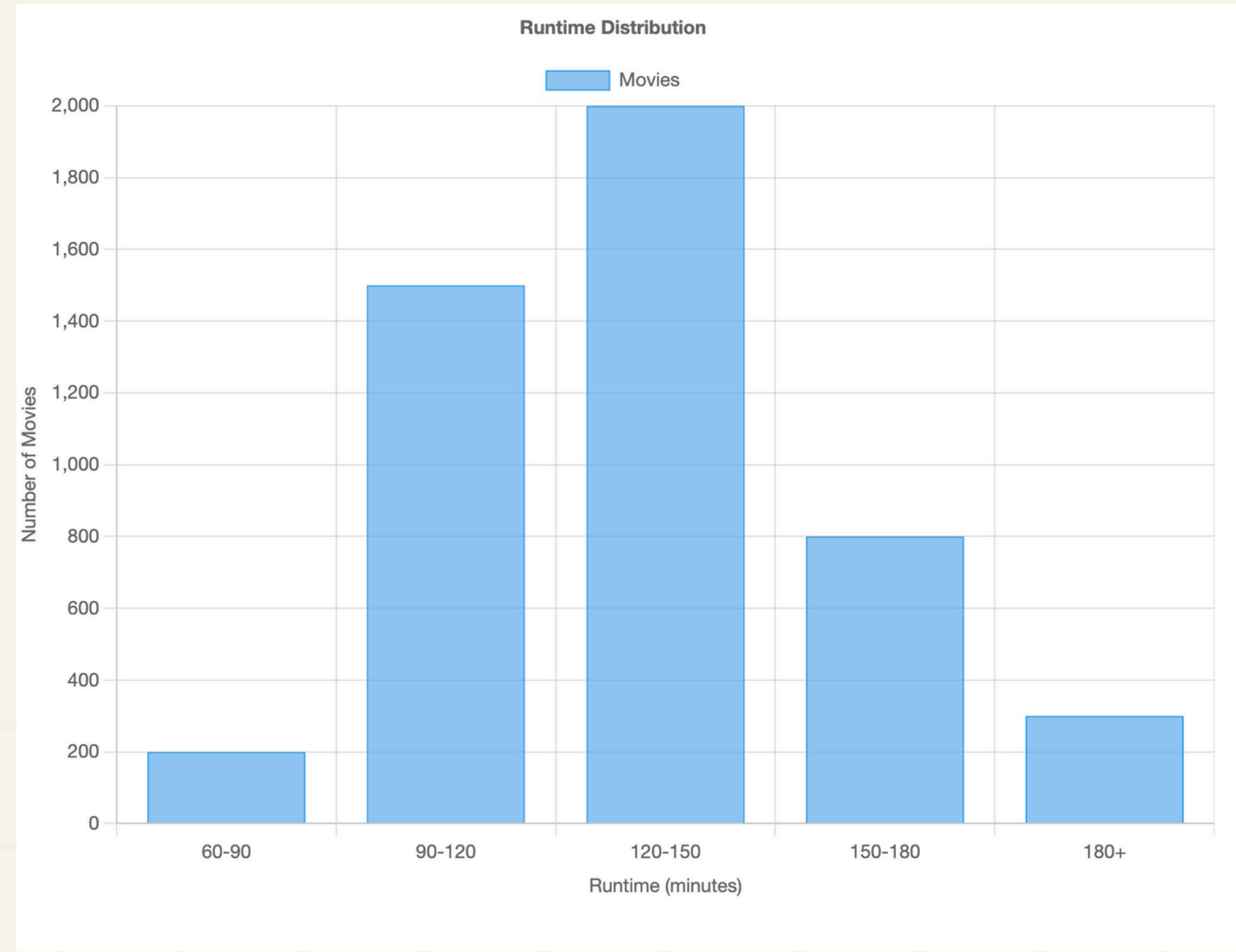
1. SUCCESS IS MOST TIED TO REVENUE (0.47) AND VOTE_COUNT (0.24), SO THESE FEATURES MIGHT PREDICT SUCCESS WELL.
2. BUDGET AND REVENUE ARE RELATED (0.53), MEANING BIGGER BUDGETS OFTEN LEAD TO HIGHER REVENUE.
3. RUNTIME HAS LITTLE CONNECTION TO OTHER FEATURES (ALL NEAR 0), SO IT MIGHT NOT BE A STRONG PREDICTOR OF SUCCESS.

WHY THOSE FEATURES?



PS: SHOWS IF HIGHER BUDGETS LEAD TO HIGHER REVENUE.

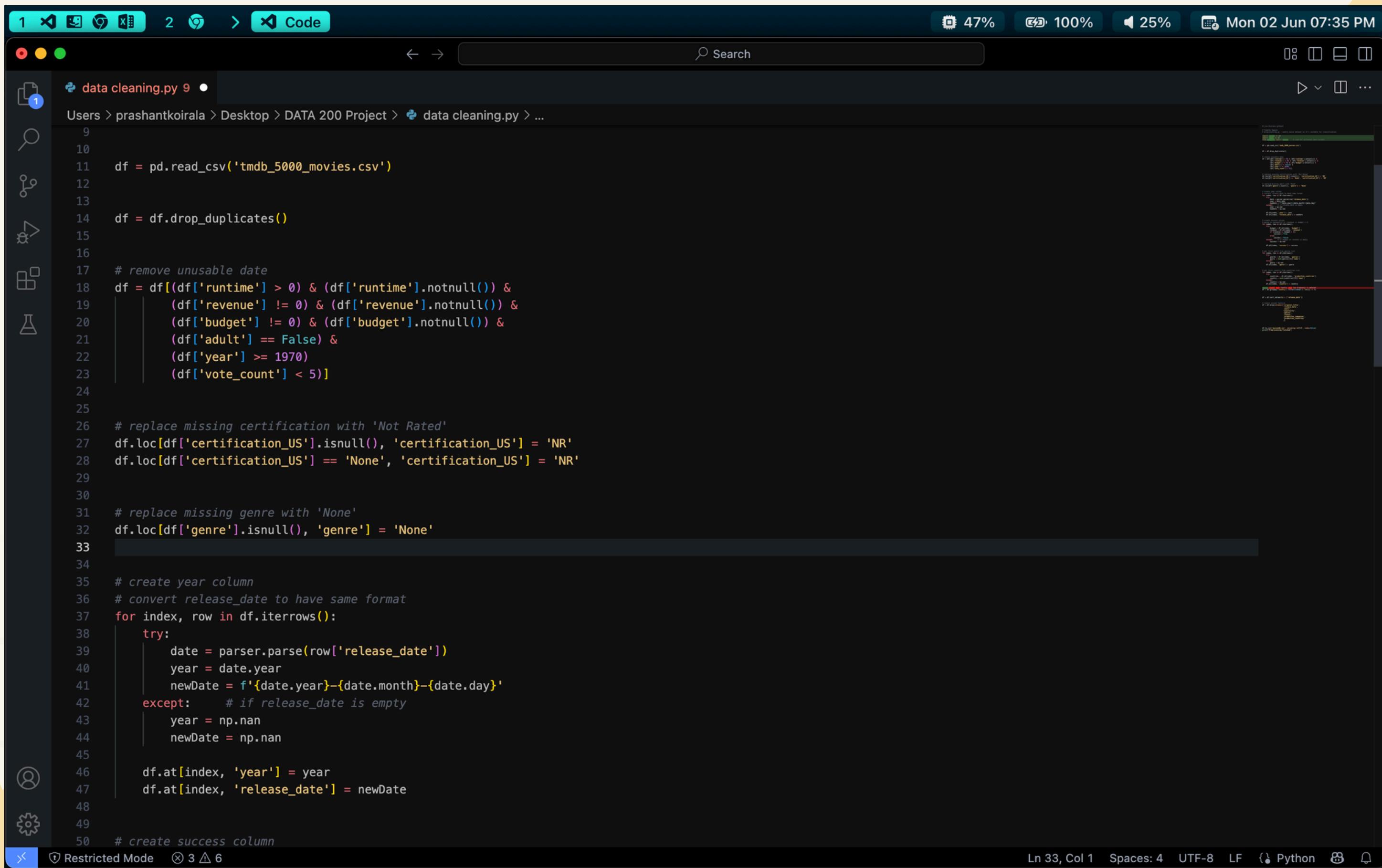
WHY THOSE FEATURES?



PS: SHOWS HOW RUNTIMES ARE DISTRIBUTED.

DATA CLEANING

DATA CLEANING PROGRAM



The screenshot shows a dark-themed code editor window titled "data cleaning.py 9". The file path is "Users > prashantkoirala > Desktop > DATA 200 Project > data cleaning.py > ...". The code is written in Python and performs the following tasks:

- Imports pandas: `import pandas as pd`
- Reads a CSV file: `df = pd.read_csv('tmdb_5000_movies.csv')`
- Drops duplicates: `df = df.drop_duplicates()`
- Filters rows based on runtime, revenue, budget, adult status, year, and vote count.
- Replaces missing certification values with "Not Rated":
 - For US certification: `df.loc[df['certification_US'].isnull(), 'certification_US'] = 'NR'`
 - For other certification: `df.loc[df['certification'] == 'None', 'certification'] = 'NR'`
- Replaces missing genre values with "None": `df.loc[df['genre'].isnull(), 'genre'] = 'None'`
- Creates a "year" column by extracting the year from the release date. If the release date is empty, it sets the year to np.nan.
 - Converts release_date to datetime: `date = parser.parse(row['release_date'])`
 - Extracts the year: `year = date.year`
 - Creates a new date string: `newDate = f'{date.year}-{date.month}-{date.day}'`
 - Handles empty release dates: `except: # if release_date is empty`
 - Assigns np.nan to year and newDate if release_date is empty.
- Updates the DataFrame with the extracted year and new date.
- Creates a "success" column.

The code editor interface includes a sidebar with icons for file operations, search, and help. The status bar at the bottom shows "Ln 33, Col 1" and "Python".

```
1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47 48 49 50
```

```
1 import pandas as pd
2
3 df = pd.read_csv('tmdb_5000_movies.csv')
4
5 df = df.drop_duplicates()
6
7 # remove unusable data
8 df = df[(df['runtime'] > 0) & (df['runtime'].notnull()) &
9         (df['revenue'] != 0) & (df['revenue'].notnull()) &
10        (df['budget'] != 0) & (df['budget'].notnull()) &
11        (df['adult'] == False) &
12        (df['year'] >= 1970) &
13        (df['vote_count'] < 5)]
14
15
16 # replace missing certification with 'Not Rated'
17 df.loc[df['certification_US'].isnull(), 'certification_US'] = 'NR'
18 df.loc[df['certification_US'] == 'None', 'certification_US'] = 'NR'
19
20
21 # replace missing genre with 'None'
22 df.loc[df['genre'].isnull(), 'genre'] = 'None'
23
24
25 # create year column
26 # convert release_date to have same format
27 for index, row in df.iterrows():
28     try:
29         date = parser.parse(row['release_date'])
30         year = date.year
31         newDate = f'{date.year}-{date.month}-{date.day}'
32     except: # if release_date is empty
33         year = np.nan
34         newDate = np.nan
35
36         df.at[index, 'year'] = year
37         df.at[index, 'release_date'] = newDate
38
39 # create success column
```

DATA CLEANING

1. Removed duplicates and movies with missing/zero budget, revenue, or runtime.
2. Replaced missing certification_US with “NR”; extracted first genre and country.
3. Filtered movies before 1970 and countries with < 5 occurrences.
4. Created success (True if revenue \geq 2x budget) and year columns.

AFTER DATA CLEANING

CONCLUSION

We successfully cleaned the TMDB dataset, creating a success variable ($\text{revenue} \geq 2x \text{ budget}$) and identifying key features like budget, revenue, and vote count. The correlation analysis showed revenue (0.47) and vote count (0.24) as strong predictors of success, with budget (0.19) having a weaker link. This supports our logistic regression approach, and further analysis will refine these findings.

REFERENCES

- DA, O., OM, S., AK, F., O, A., A, O., A, O., A, W., & M, Y. (2021). Movie success prediction using data mining. *British Journal of Computer Networking and Information Technology*, 4(2), 22–30. <https://doi.org/10.52589/bjcnit-cqocirec>
- Zheng, Y. (2024). Predicting movie box office based on machine learning, deep learning, and statistical methods. *Applied and Computational Engineering*, 94(1), 20–32. <https://doi.org/10.54254/2755-2721/94/2024melb0069>
- Subramaniyaswamy, V., Vaibhav, M. V., Prasad, R. V., & Logesh, R. (2017). Predicting movie box office success using multiple regression and SVM. *Journal of Science & Statistics*, 182–186. <https://doi.org/10.1109/iss1.2017.8389394>
- Velingkar, G., Varadarajan, R., Lanka, S., & M, A. K. (2022). Movie Box-Office success Prediction using Machine Learning. 2022 Second International Conference on Power, Control and Computing Technologies (ICPC2T), 9, 1–6.
<https://doi.org/10.1109/icpc2t53885.2022.9776798>
- Mr, N. V., Mr, P. M., & Pb, S. B. (2014). Predicting movie success based on IMDB data. *International Journal of Business Intelligents*, 3(2). <https://doi.org/10.20894/ijbi.105.003.002.004>



THANK YOU

