

- <https://www>

- 'Name' - Book title.
 - 'Author' - Person who wrote the book.
 - 'User Rating' - Book rating out of 5.
 - 'Reviews' - Number of people/readers reviewing the book via. rating.
 - 'Price' - Cost of each book in US dollars.
 - 'Year' - Book launch year.
 - 'Genre' - Categories of books
 - To analyse the data using python libraries pandas and numpy and for data visualisation using plotly library

- The books we
- Visualize the
- Top 10 book
- Does a high
- Is the mean
- Mean price

- ```
: import pandas as pd
import numpy as np
import plotly.express as px
import plotly.graph_objects as go
import warnings
warnings.filterwarnings('ignore')
```

|   |             |
|---|-------------|
| 0 | 10-D        |
| 1 |             |
| 2 | 12 Rules fo |
| 3 |             |

- ## Distribution of genre in complete dataset.

```
: print(data.filter(['Name', 'Genre']).shape)
print(data.filter(['Name', 'Genre']).drop_duplicates().shape)

(550, 2)
(351, 2)

: category = data.filter(['Name', 'Genre']).drop_duplicates()
category.isnull().sum()
category.drop_duplicates(inplace=True)
category = category.groupby(['Genre']).agg(count_genre=('Genre', 'count'))
category.reset_index(level=0, inplace=True)
fig = px.pie(category, values='count_genre', names='Genre', title='Distribution of genres')
fig.show()
```

### Distribution of genre

A pie chart titled "Distribution of genre" showing the proportion of two categories. The chart is split into two main segments: a blue segment representing 54.4% and an orange segment representing 45.6%. The segments are labeled with their respective percentages.

| Genre  | Percentage |
|--------|------------|
| Blue   | 54.4%      |
| Orange | 45.6%      |

- Pie chart shows the distribution between fiction and non-fiction books.
- 160 books genre is fiction and 191 genre is non-fiction.
- Non-fiction books is 8.8 % higher than fiction.

- ```
: year_genre = data.filter(['Name', 'Year', 'Genre'])
: # year_genre.drop_duplicates(inplace=True)
: year_genre = year_genre.groupby(['Year', 'Genre'])
: year_genre.reset_index(level=[0,1], inplace=True)
: fig = px.bar(x=year_genre.Year, y=year_genre['Total'])
: fig.show()
```

Genre comparison per year

The chart displays the total number of books published per year, categorized by genre. The non-fiction genre (blue bars) shows a general upward trend from approximately 24 books in 2009 to about 29 books in 2013. The fiction genre (orange bars) also shows a general upward trend, starting at approximately 26 books in 2009 and reaching about 29 books in 2013. Both genres show a slight dip in 2010.

Year	Non-fiction	Fiction
2009	24	26
2010	20	30
2011	21	29
2012	21	29
2013	24	26

- Sales data is not available in a dataset.

Assumption

- If we assume all readers have provided a review then we can assume the number of reviews as number of books sold.

Hence, we can compute profit as the product of reviews and price for a given year.

```
In [7]: bestselling = data.filter(['Name', 'Author', 'Reviews', 'Price', 'Genre'])
bestselling.drop_duplicates(subset=['Name'], inplace=True)
bestselling['selling_price'] = bestselling['Reviews'] * bestselling['Price']

bestselling_fiction = bestselling[bestselling.Genre == 'Fiction']
bestselling_fiction['rank_fiction_price'] = bestselling_fiction.selling_price.rank(method='first', ascending=False).astype(np.int32)
bestselling_fiction = bestselling_fiction[bestselling_fiction.rank_fiction_price < 11].copy()

bestselling_non_fiction = bestselling[bestselling.Genre == 'Non_Fiction']
```

```
copy()

fig = go.Figure(data=[
    go.Bar(name='Fiction', x=bestselling_fiction.Author, y=bestselling_fiction.selling_price),
    go.Bar(name='Non-Fiction', x=bestselling_non_fiction.Author, y=bestselling_non_fiction.selling_price)
])
fig.update_layout(barmode='group', title='Top 10 profitable author per genre')
fig.update_xaxes(title_text="Authors")
```

Top 10 profitable author per genre

1.5M

Fiction

Non-Fiction

Author	Genre	Profit (M)
J.K. Rowling	Fiction	1.5
Agatha Christie	Fiction	1.4
Dan Brown	Fiction	1.3
Stephen King	Fiction	1.2
Harlan Coben	Fiction	1.1
E.L. James	Fiction	1.0
J.R.R. Tolkien	Fiction	0.9
J.D. Salinger	Fiction	0.8
William Shakespeare	Fiction	0.7
Charles Dickens	Fiction	0.6

Item	Price
George	~0.55M
E L Jam	~0.6M
Gillian F	~0.55M
J.K. Row	~0.55M
John Gr	~0.6M
Paula H	~0.7M
Delia Ov	~0.6M
Donna T	~0.6M
Jordan H	~0.25M
Michelle	~0.65M
David GG	~0.25M
Tara We	~0.4M
Lin-Manu	~0.35M
Mark Må	~0.4M
Americas	~0.4M
Cheryl S	~0.5M
Laura H	~0.5M
Cheryl S	~0.3M

Most profitable au



Top 1

- | Book Title | Author |
|---|--|
| Fifty Shades of Grey: Book One of the Fifty Shades Trilogy (Fifty Shades of Grey) | J. K. Rowling |
| The Fault in Our Stars | Percy Jackson & the Olympians: The Lightning Thief (Percy Jackson & the Olympians) |
| The Nightingale | Stephenie Meyer |
| All the Light We Cannot See | Mark Haddon |
| The Goldfinch | Colson Whitehead |

- X-axis represent the total number of reviews.
- Y-axis represent the books name.
- The number of reviews ranges between 37 and 87,841.
- By far the most reviews have been given to 'Where the
Train' by Paula Hawking with a user rating of 4.1.

The
most
most
most
most

- ```
most_rvs_year['rank_Reviews'] = most_rvs_year.groupby(['Year']).Reviews.rank(method='first', ascending=False).astype(np.int32)
most_rvs_year = most_rvs_year[most_rvs_year.rank_Reviews == 1]
most_rvs_year
fig = px.bar(x=most_rvs_year.Year, y=most_rvs_year.Reviews, color=most_rvs_year.Name, title='Books with the maximum number of reviews per year',
 labels={'x': 'Years', 'y': 'Total Reviews'})
fig.show()
```

A bar chart titled "Color" showing the count of users for different colors. The y-axis is labeled "Users" and ranges from 0 to 80k. The x-axis categories are color names: Red, Green, Blue, Orange, Purple, Yellow, and Cyan. The bars are colored according to the categories they represent. The chart shows that the "Green" category has the highest user count, followed by "Orange" and "Red".

| Color Category | Users (approx.) |
|----------------|-----------------|
| Red            | 79,000          |
| Green          | 85,000          |
| Blue           | 5,000           |
| Orange         | 79,000          |
| Purple         | 5,000           |
| Yellow         | 5,000           |
| Cyan           | 5,000           |

| Year | Reviews (k) |
|------|-------------|
| 2010 | 20          |
| 2011 | 30          |
| 2012 | 28          |
| 2013 | 25          |
| 2014 | 25          |
| 2015 | 28          |
| 2016 | 30          |
| 2017 | 28          |
| 2018 | 30          |
| 2019 | 32          |

**Visualize the distribution of genre with respect to reviews.**

### Review comparision per genre

A pie chart comparing the percentage of reviews for Fiction and Non Fiction genres. The chart is divided into two segments: a blue segment representing Fiction at 57.3% and an orange segment representing Non Fiction at 42.7%. The Non Fiction segment is labeled with its percentage value, 42.7%, inside it.

| Genre       | Percentage |
|-------------|------------|
| Fiction     | 57.3%      |
| Non Fiction | 42.7%      |

A pie chart illustrating the distribution of reviews by genre. The chart is divided into two main segments: a large blue segment representing Non Fiction reviews (2,810,195) and a smaller red segment representing Fiction reviews (3,764,110). The blue segment occupies approximately three-quarters of the circle, while the red segment occupies the remaining quarter.

| Genre       | Number of Reviews |
|-------------|-------------------|
| Non Fiction | 2,810,195         |
| Fiction     | 3,764,110         |

In [13]:

```
max_rating = data.filter(['Name', 'Author', 'User Rating', 'Reviews']).drop_duplicates()
max_rating['rank_rating'] = max_rating['User Rating'].rank(method='first', ascending=False).astype(np.int32)
max_rating = max_rating[max_rating.rank_rating < 11]
max_rating['rank_review'] = max_rating['Reviews'].rank(method='first', ascending=False).astype(np.int32)
max_rating = max_rating[max_rating.rank_review < 11]
print(max_rating.rank_rating.min(), max_rating.rank_rating.max())
print(max_rating.rank_review.min(), max_rating.rank_review.max())
max_rating['total_rank'] = (max_rating['rank_rating'] + max_rating['rank_review'])/2
max_rating['rank_total'] = max_rating['total_rank'].rank(method='first', ascending=True).astype(np.int32)
)
max_rating = max_rating[max_rating.rank_total < 11]
fig = px.bar(y=max_rating.Name, x=max_rating['User Rating'], labels={'x':'Rating', 'y':'Book Title'}, title='Top 10 books with maximum rating')
fig.show()
```

### Top 10 books with maximum rating

The figure is a horizontal bar chart titled "Top 10 books with maximum rating". The y-axis is labeled "Book Title" and lists the following titles from top to bottom: "The Very Hungry Caterpillar", "The Nightingale: A Novel", "The 5 Love Languages: The Secret to Love that Lasts", "Oh, the Places You'll Go!", "Jesus Calling: Enjoying Peace in His Presence (with Scripture References)", "Harry Potter and the Chamber of Secrets: The Illustrated Edition (Harry Potter, Book 2)", "Dog Man: Fetch-22: From the Creator of Captain Underpants (Dog Man #8)", and "Can't Hurt Me: Master Your Mind and Defy the Odds". Each title is accompanied by a solid blue horizontal bar.

| Book Title                                                                              |
|-----------------------------------------------------------------------------------------|
| The Very Hungry Caterpillar                                                             |
| The Nightingale: A Novel                                                                |
| The 5 Love Languages: The Secret to Love that Lasts                                     |
| Oh, the Places You'll Go!                                                               |
| Jesus Calling: Enjoying Peace in His Presence (with Scripture References)               |
| Harry Potter and the Chamber of Secrets: The Illustrated Edition (Harry Potter, Book 2) |
| Dog Man: Fetch-22: From the Creator of Captain Underpants (Dog Man #8)                  |
| Can't Hurt Me: Master Your Mind and Defy the Odds                                       |

| Book Title                           | Rating |
|--------------------------------------|--------|
| Becoming                             | 4.9    |
| The Handmaid's Tale                  | 4.9    |
| The Brief Wondrous Life of Oscar Wao | 4.8    |
| The Road                             | 4.8    |
| The Art of the Novel                 | 4.8    |
| The Goldfinch                        | 4.8    |

```
print(rating_price.isnull().sum().sum())
rating_price = pd.DataFrame(rating_price.groupby(['User Rating']).agg({'Price': 'mean', 'Name': 'nunique'}))
rating_price.reset_index(level=0,inplace=True)
fig = px.scatter(rating_price, x="User Rating", y="Price", title='Higher rating affect on price', size='Name')
fig.update_xaxes(title_text="User Rating")
fig.update_yaxes(title_text="Price in $")
fig.show()
```

0

A bubble chart illustrating the relationship between User Rating (X-axis) and Price in dollars (Y-axis). The X-axis ranges from approximately 3.2 to 5.0, and the Y-axis ranges from 11 to 16. Data points are represented by blue circles, where the size of the circle corresponds to a third variable, likely popularity or volume. The chart shows a general trend where higher user ratings are associated with higher prices.

| User Rating | Price (\$) | Bubble Size (approx.) |
|-------------|------------|-----------------------|
| 3.2         | 12.0       | Small                 |
| 3.8         | 14.0       | Medium                |
| 4.0         | 11.8       | Very Small            |
| 4.1         | 11.6       | Very Small            |
| 4.2         | 11.6       | Very Small            |
| 4.3         | 11.8       | Medium                |
| 4.3         | 14.2       | Large                 |
| 4.5         | 16.0       | Very Large            |
| 4.6         | 11.8       | Large                 |
| 4.6         | 12.0       | Large                 |
| 4.7         | 14.0       | Large                 |
| 4.8         | 11.8       | Large                 |
| 4.9         | 12.8       | Medium                |

- X-axis represent user rating.
- Y-axis represent price in dollars.
- There is no clear relationship between user rating and price.
- We have seen that the number of books are more in higher rating.

## Is the mean price is changing over the years?

```
In [15]: price_year = data.filter(['Price', 'Year'])
price_year = pd.DataFrame(price_year.groupby(['Year']).Price.mean())
price_year.reset_index(level=0, inplace=True)
fig = px.line(price_year, x="Year", y="Price", title='Mean price change over the years')
fig.update_xaxes(title_text="Year")
fig.update_yaxes(title_text="Price in $")
fig.show()
```

A line graph illustrating price fluctuations over time. The vertical axis represents the price in dollars, ranging from 10 to 15. The horizontal axis represents time. The price starts at approximately 15.2, drops to a low of about 13.5, rises to a peak of about 15.3, dips to 14.7, stays flat at 14.7, drops sharply to 10.5, rises to a peak of about 13.2, falls to 11.5, stays flat at 11.5, and finally drops to a low of about 10.2.

| Time | Price (\$) |
|------|------------|
| 1    | 15.2       |
| 2    | 13.5       |
| 3    | 15.3       |
| 4    | 14.7       |
| 5    | 14.7       |
| 6    | 10.5       |
| 7    | 13.2       |
| 8    | 11.5       |
| 9    | 11.5       |
| 10   | 10.2       |

## Observation

- X-axis represent years.
- Y-axis represent price in dollars.
- Sudden fall in price between 2014-2015 but we don't know the reason because unavailability of data.

## Mean price per genre.

```
In [16]: genre_price = data.filter(['Price','Genre'])
genre_price = pd.DataFrame(genre_price.groupby(['Genre']).Price.mean())
genre_price.reset_index(level=0,inplace=True)
fig = px.bar(genre_price, x='Genre', y='Price',title='Mean price per genre')
fig.update_xaxes(title_text="Genre")
```

Mean price per genre

| Genre     | Mean Price (\$) |
|-----------|-----------------|
| Action    | ~10.8           |
| Adventure | ~14.5           |
| RPG       | ~14.5           |
| Strategy  | ~14.5           |

The figure is a horizontal box plot titled 'Genre' on the x-axis, which has two categories: 'Fiction' and 'Non Fiction'. The y-axis is labeled 'Price' and has a single tick mark at 0. The box plot shows the distribution of price for each genre. The median for Non Fiction is approximately 14.84, while for Fiction it is approximately 10.85. Both distributions are skewed to the right, with many outliers at higher price points.

| Genre       | Median | Approx. Range (Min - Max) |
|-------------|--------|---------------------------|
| Fiction     | 10.85  | ~5.00 - 25.00             |
| Non Fiction | 14.84  | ~10.00 - 25.00            |