

IMPORTING CSV FILE

In [4]:

```

1 # import pandas as pd
2 import pandas as pd
3 df = pd.read_csv (r"C:\Users\Anusha V\Desktop\YouTube.csv")
4 print(df)

```

| | Rank | Username | Categories | Suscribers | Country | \ |
|-----|------|---------------------|------------|------------|----------------|---|
| 0 | 1 | tseries | True | 249500000 | India | |
| 1 | 2 | MrBeast | False | 183500000 | Estados Unidos | |
| 2 | 3 | CoComelon | True | 165500000 | Unknown | |
| 3 | 4 | SETIndia | False | 162600000 | India | |
| 4 | 5 | KidsDianaShow | True | 113500000 | Unknown | |
| .. | ... | ... | ... | ... | ... | |
| 194 | 195 | dianaandromahin8102 | True | 26200000 | Unknown | |
| 195 | 196 | nickiminaj | False | 26100000 | Estados Unidos | |
| 196 | 197 | mariliamendoncareal | True | 26100000 | Brasil | |
| 197 | 198 | TheLallantop | False | 26000000 | India | |
| 198 | 199 | AGT | True | 25900000 | Estados Unidos | |

| | Visits | Likes | Comments | \ |
|-----|-----------|---------|----------|---|
| 0 | 86200 | 2700 | 78 | |
| 1 | 117400000 | 5300000 | 18500 | |
| 2 | 7000000 | 24700 | 0 | |
| 3 | 15600 | 166 | 9 | |
| 4 | 3900000 | 12400 | 0 | |
| .. | ... | ... | ... | |
| 194 | 109000 | 322 | 0 | |
| 195 | 1600000 | 98300 | 7600 | |
| 196 | 0 | 0 | 0 | |
| 197 | 30800 | 822 | 52 | |
| 198 | 186800 | 3100 | 197 | |

Links

| | |
|-----|--|
| 0 | http://youtube.com/channel/UCq-Fj5jknLsUf-MWSy... (http://youtube.co m/channel/UCq-Fj5jknLsUf-MWSy...) |
| 1 | http://youtube.com/channel/UCX60Q3DkcsbYNE6H8u... (http://youtube.co m/channel/UCX60Q3DkcsbYNE6H8u...) |
| 2 | http://youtube.com/channel/UCbCmjCuTUZos6Inko4... (http://youtube.co m/channel/UCbCmjCuTUZos6Inko4...) |
| 3 | http://youtube.com/channel/UCpEhnqL0y41EpW2TvW... (http://youtube.co m/channel/UCpEhnqL0y41EpW2TvW...) |
| 4 | http://youtube.com/channel/UCK8GzjMOrta8yxDcKf... (http://youtube.co m/channel/UCK8GzjMOrta8yxDcKf...) |
| .. | ... |
| 194 | http://youtube.com/channel/UC5ma-WCc8jNCV0WyRn... (http://youtube.co m/channel/UC5ma-WCc8jNCV0WyRn...) |
| 195 | http://youtube.com/channel/UC3j0d7GUMhpgJRBhiL... (http://youtube.co m/channel/UC3j0d7GUMhpgJRBhiL...) |
| 196 | http://youtube.com/channel/UCwfEOOn001DWcyTgzVV... (http://youtube.co m/channel/UCwfEOOn001DWcyTgzVV...) |
| 197 | http://youtube.com/channel/UCx8Z14PpntdaxCt2ha... (http://youtube.co m/channel/UCx8Z14PpntdaxCt2ha...) |
| 198 | http://youtube.com/channel/UCT2X19JJJaJGUN7mrYu... (http://youtube.co m/channel/UCT2X19JJJaJGUN7mrYu...) |

Loading [MathJax]/extensions/MathML/content/mathml.js
119 rows x 9 columns

HEAD FUNCTION

In [288]: 1 df.head(20)

Out[288]:

| | Rank | Username | Categories | Suscribers | Country | Visits | Likes | Comments |
|----|------|----------------------|------------|------------|----------------|-----------|---------|----------|
| 0 | 1 | tseries | True | 249500000 | India | 86200 | 2700 | |
| 1 | 2 | MrBeast | False | 183500000 | Estados Unidos | 117400000 | 5300000 | 1 |
| 2 | 3 | CoComelon | True | 165500000 | Unknown | 7000000 | 24700 | |
| 3 | 4 | SETIndia | False | 162600000 | India | 15600 | 166 | |
| 4 | 5 | KidsDianaShow | True | 113500000 | Unknown | 3900000 | 12400 | |
| 5 | 6 | PewDiePie | False | 111500000 | Estados Unidos | 2400000 | 197300 | |
| 6 | 7 | LikeNastyofficial | True | 107500000 | Unknown | 2600000 | 28000 | |
| 7 | 8 | VladandNiki | False | 101400000 | Unknown | 4100000 | 22100 | |
| 8 | 9 | zeemusiccompany | True | 99700000 | India | 74300 | 2600 | |
| 9 | 10 | WWE | False | 97200000 | Estados Unidos | 184500 | 6300 | |
| 10 | 11 | BLACKPINK | True | 91300000 | Estados Unidos | 863200 | 146900 | |
| 11 | 12 | GoldminesTelefilms | False | 89700000 | India | 34600 | 421 | |
| 12 | 13 | SonySAB | True | 85400000 | India | 35500 | 615 | |
| 13 | 14 | 5MinuteCraftsYouTube | False | 80300000 | India | 164600 | 703 | |
| 14 | 15 | BTS | True | 76500000 | India | 969700 | 180300 | |
| 15 | 16 | HYBELABELS | False | 76500001 | India | 969700 | 180300 | |
| 16 | 17 | HYBELABELS | True | 76500002 | India | 969700 | 180300 | |
| 17 | 18 | Pinkfong | False | 69600000 | Unknown | 506400 | 397 | |
| 18 | 19 | ChuChuTV | True | 67500000 | Unknown | 652100 | 17800 | |
| 19 | 20 | KondZilla | False | 66600000 | Brasil | 30400 | 1100 | |

TAIL FUNCTION

In [289]: 1 df.head(20)

| | Rank | Username | Categories | Suscribers | Country | Visits | Likes | Comments |
|----|------|----------------------|------------|------------|----------------|-----------|---------|----------|
| 0 | 1 | tseries | True | 249500000 | India | 86200 | 2700 | |
| 1 | 2 | MrBeast | False | 183500000 | Estados Unidos | 117400000 | 5300000 | 1 |
| 2 | 3 | CoComelon | True | 165500000 | Unknown | 7000000 | 24700 | |
| 3 | 4 | SETIndia | False | 162600000 | India | 15600 | 166 | |
| 4 | 5 | KidsDianaShow | True | 113500000 | Unknown | 3900000 | 12400 | |
| 5 | 6 | PewDiePie | False | 111500000 | Estados Unidos | 2400000 | 197300 | |
| 6 | 7 | LikeNastyofficial | True | 107500000 | Unknown | 2600000 | 28000 | |
| 7 | 8 | VladandNiki | False | 101400000 | Unknown | 4100000 | 22100 | |
| 8 | 9 | zeemusiccompany | True | 99700000 | India | 74300 | 2600 | |
| 9 | 10 | WWE | False | 97200000 | Estados Unidos | 184500 | 6300 | |
| 10 | 11 | BLACKPINK | True | 91300000 | Estados Unidos | 863200 | 146900 | |
| 11 | 12 | GoldminesTelefilms | False | 89700000 | India | 34600 | 421 | |
| 12 | 13 | SonySAB | True | 85400000 | India | 35500 | 615 | |
| 13 | 14 | 5MinuteCraftsYouTube | False | 80300000 | India | 164600 | 703 | |
| 14 | 15 | BTS | True | 76500000 | India | 969700 | 180300 | |
| 15 | 16 | HYBELABELS | False | 76500001 | India | 969700 | 180300 | |
| 16 | 17 | HYBELABELS | True | 76500002 | India | 969700 | 180300 | |
| 17 | 18 | Pinkfong | False | 69600000 | Unknown | 506400 | 397 | |
| 18 | 19 | ChuChuTV | True | 67500000 | Unknown | 652100 | 17800 | |
| 19 | 20 | KondZilla | False | 66600000 | Brasil | 30400 | 1100 | |

DESCRIBE FUNCTION

In [290]: 1 df.describe()

| | Rank | Suscribers | Visits | Likes | Comments |
|-------|------------|--------------|--------------|--------------|--------------|
| count | 199.000000 | 1.990000e+02 | 1.990000e+02 | 1.990000e+02 | 199.000000 |
| mean | 100.000000 | 4.405226e+07 | 2.261227e+06 | 1.047747e+05 | 2096.356784 |
| std | 57.590508 | 2.753718e+07 | 1.084192e+07 | 5.321669e+05 | 7913.381544 |
| min | 1.000000 | 2.590000e+07 | 0.000000e+00 | 0.000000e+00 | 0.000000 |
| 25% | 50.500000 | 3.030000e+07 | 3.225000e+04 | 5.845000e+02 | 0.000000 |
| 50% | 100.000000 | 3.550000e+07 | 2.257000e+05 | 3.900000e+03 | 28.000000 |
| 75% | 149.500000 | 4.500000e+07 | 1.050000e+06 | 3.465000e+04 | 375.000000 |
| max | 199.000000 | 2.495000e+08 | 1.174000e+08 | 5.300000e+06 | 82800.000000 |

Loading [MathJax]/extensions/MathML/content-mathml.js

PIVOT FUNCTION

```
In [227]: 1 pivot_table = df.pivot_table(index = 'Likes', aggfunc='mean')
2 pivot_table
```

Out[227]:

| | Categories | Rank |
|------------|------------|---------|
| Likes | | |
| 0 | 4567.0 | 124.125 |
| 1000 | 4567.0 | 161.000 |
| 10100 | 4567.0 | 70.000 |
| 10200 | 4567.0 | 36.000 |
| 106500 | 4567.0 | 157.000 |
| ... | ... | ... |
| 983 | 4567.0 | 23.000 |
| 98300 | 4567.0 | 196.000 |
| 993 | 4567.0 | 192.000 |
| 994 | 4567.0 | 64.000 |
| HYBELABELS | 4567.0 | 16.500 |

167 rows × 2 columns

MELT FUNCTION

```
In [228]: 1 melt_df = pd.melt(df,id_vars=['Likes'], var_name='visits')
2 melt_df
```

Out[228]:

| | Likes | visits | value |
|------|---------|-------------|-------|
| 0 | 2700 | Rank | 1 |
| 1 | 5300000 | Rank | 2 |
| 2 | 24700 | Rank | 3 |
| 3 | 166 | Rank | 4 |
| 4 | 12400 | Rank | 5 |
| ... | ... | ... | ... |
| 2383 | 322 | Unnamed: 12 | NaN |
| 2384 | 98300 | Unnamed: 12 | NaN |
| 2385 | 0 | Unnamed: 12 | NaN |
| 2386 | 822 | Unnamed: 12 | NaN |
| 2387 | 3100 | Unnamed: 12 | NaN |

2388 rows × 3 columns

Loading [MathJax]/extensions/MathML/content-mathml.js

REDUCE FUNCTION

In [175]:

```
1 from functools import reduce
2 def multiply(x, y):
3     return x * y
4 column_to_reduce = 'Rank'
5 result = reduce(multiply, df['Rank'])
6 print(result)
```

```
39432893368239525177618160696609253114756798884358663164737126662217972498
17016714601521420059923119520886060694598194151288213951213185525309633124
76414965556731428635381658618698494471961222810725832120127016645932065613
71414742663876212120378695162016062870278978433011301595208516203117585042
939808946111139481185194868736000000000000000000000000000000000000000000000000000
000
```

MAPPING FUNCTION

Loading [MathJax]/extensions/MathML/content-mathml.js

```
In [5]: 1 def map_grades(grade):
2     grade_maping = {
3         'A':90,
4         'B':80,
5         'C':70,
6         'D':60,
7         'E':50,
8         'F':40
9     }
10 LeagueIndex = 'Grade'
11 df['Likes'] = df['Rank'].map(map_grades)
12 print(df)
```

| | Rank | Username | Categories | Suscribers | Country | \ |
|-----|------|---------------------|------------|------------|----------------|---|
| 0 | 1 | tseries | True | 249500000 | India | |
| 1 | 2 | MrBeast | False | 183500000 | Estados Unidos | |
| 2 | 3 | CoComelon | True | 165500000 | Unknown | |
| 3 | 4 | SETIndia | False | 162600000 | India | |
| 4 | 5 | KidsDianaShow | True | 113500000 | Unknown | |
| .. | .. | ... | ... | ... | ... | |
| 194 | 195 | dianaandromahin8102 | True | 26200000 | Unknown | |
| 195 | 196 | nickiminaj | False | 26100000 | Estados Unidos | |
| 196 | 197 | mariliamendoncareal | True | 26100000 | Brasil | |
| 197 | 198 | TheLallantop | False | 26000000 | India | |
| 198 | 199 | AGT | True | 25900000 | Estados Unidos | |

| | Visits | Likes | Comments | \ |
|-----|-----------|-------|----------|---|
| 0 | 86200 | None | 78 | |
| 1 | 117400000 | None | 18500 | |
| 2 | 7000000 | None | 0 | |
| 3 | 15600 | None | 9 | |
| 4 | 3900000 | None | 0 | |
| .. | .. | .. | .. | |
| 194 | 109000 | None | 0 | |
| 195 | 1600000 | None | 7600 | |
| 196 | 0 | None | 0 | |
| 197 | 30800 | None | 52 | |
| 198 | 186800 | None | 197 | |

Links

| | | |
|-----|---|---|
| 0 | http://youtube.com/channel/UCq-Fj5jknLsUF-MWSy... | (http://youtube.co |
| 1 | http://youtube.com/channel/UCX60Q3DkcsbYNE6H8u... | (http://youtube.co |
| 2 | http://youtube.com/channel/UCbCmjCuTUZos6Inko4... | (http://youtube.co |
| 3 | http://youtube.com/channel/UCpEhnqL0y41EpW2TvW... | (http://youtube.co |
| 4 | http://youtube.com/channel/UCk8GzjMOrta8yxDcKf... | (http://youtube.co |
| .. | .. | .. |
| 194 | http://youtube.com/channel/UC5ma-WCc8jNCV0WyRn... | (http://youtube.co |
| 195 | http://youtube.com/channel/UC3j0d7GUMhpgJRBhiL... | (http://youtube.co |
| 196 | http://youtube.com/channel/UCwfEOOn001DWcyTgzVV... | (http://youtube.co |
| 197 | http://youtube.com/channel/UCx8Z14PpntdaxCt2ha... | (http://youtube.co |
| 198 | http://youtube.com/channel/UCT2X19JJaJGUN7mrYu... | (http://youtube.co |

[199 rows x 9 columns]

Filter FUNCTION

In [8]:

```

1 def score_condition(score):
2     return score > 5
3 Likes = 'score_condition'
4 filtered_df = df[df['Likes'].apply(score_condition)]
5 print(filtered_df)

```

```

-
TypeError Traceback (most recent call last)
t)
~\AppData\Local\Temp\ipykernel_10508\466627290.py in <module>
      2     return score > 5
      3 Likes = 'score_condition'
----> 4 filtered_df = df[df['Likes'].apply(score_condition)]
      5 print(filtered_df)

~\anaconda3\lib\site-packages\pandas\core\series.py in apply(self, func, convert_dtype, args, **kwargs)
    4431         dtype: float64
    4432         """
-> 4433         return SeriesApply(self, func, convert_dtype, args, kwargs).apply()
    4434
    4435     def _reduce(

~\anaconda3\lib\site-packages\pandas\core\apply.py in apply(self)
    1086         return self.apply_str()
    1087
-> 1088         return self.apply_standard()
    1089
    1090     def agg(self):

~\anaconda3\lib\site-packages\pandas\core\apply.py in apply_standard(self)
    1141             # List[Union[Callable[..., Any], str]]]]"; expected
ed
    1142             # "Callable[[Any], Any]"
-> 1143             mapped = lib.map_infer(
    1144                 values,
    1145                 f, # type: ignore[arg-type]

~\anaconda3\lib\site-packages\pandas\_libs\lib.pyx in pandas._libs.lib.map_infer()

~\AppData\Local\Temp\ipykernel_10508\466627290.py in score_condition(score)
      1 def score_condition(score):
----> 2     return score > 5
      3 Likes = 'score_condition'
      4 filtered_df = df[df['Likes'].apply(score_condition)]
      5 print(filtered_df)

TypeError: '>' not supported between instances of 'NoneType' and 'int'

```

LAMBDA FUNCTION

```
In [178]: 1 modify_score = lambda x: x * 2
2 LeagueIndex = 'score'
3 df['Likes'] = df['Rank'].apply(modify_score)
4 print(df)
```

Loading [MathJax]/extensions/MathML/content-mathml.js

| | Rank | Username | Categories | Suscribers | \ |
|-----|------|---------------------|---------------------------|------------|-----|
| 0 | 1 | tseries | MÁsica y baile | 249500000 | |
| 1 | 2 | MrBeast | Videojuegos, Humor | 183500000 | |
| 2 | 3 | CoComelon | EducaciÃ³n | 165500000 | |
| 3 | 4 | SETIndia | EducaciÃ³n | 162600000 | |
| 4 | 5 | KidsDianaShow | AnimaciÃ³n, Juguetes | 113500000 | |
| .. | ... | ... | ... | ... | ... |
| 194 | 195 | dianaandromahin8102 | AnimaciÃ³n, Juguetes | 26200000 | |
| 195 | 196 | nickiminaj | MÁsica y baile | 26100000 | |
| 196 | 197 | mariliamendoncareal | MÁsica y baile | 26100000 | |
| 197 | 198 | TheLallantop | Noticias y PolÃtica | 26000000 | |
| 198 | 199 | AGT | MÁsica y baile, PelÃculas | 25900000 | |

| | Country | Visits | Likes | Comments | \ |
|-----|----------------|-----------|-------|----------|-----|
| 0 | India | 86200 | 2 | 78 | |
| 1 | Estados Unidos | 117400000 | 4 | 18500 | |
| 2 | Unknown | 7000000 | 6 | 0 | |
| 3 | India | 15600 | 8 | 9 | |
| 4 | Unknown | 3900000 | 10 | 0 | |
| .. | ... | ... | ... | ... | ... |
| 194 | Unknown | 109000 | 390 | 0 | |
| 195 | Estados Unidos | 1600000 | 392 | 7600 | |
| 196 | Brasil | 0 | 394 | 0 | |
| 197 | India | 30800 | 396 | 52 | |
| 198 | Estados Unidos | 186800 | 398 | 197 | |

Links Unnamed: 9 Unnamed:

| | | | | | |
|-----|---|---|-----|-----|--|
| 10 | \ | | | | |
| 0 | http://youtube.com/channel/UCq-Fj5jknLsUf-MWSy... | (http://youtube.co | NaN | NaN | |
| 1 | http://youtube.com/channel/UCX60Q3DkcsbYNE6H8u... | (http://youtube.co | NaN | NaN | |
| 2 | http://youtube.com/channel/UCbCmjCuTUZos6Inko4... | (http://youtube.co | NaN | NaN | |
| 3 | http://youtube.com/channel/UCpEhnqL0y41EpW2TvW... | (http://youtube.co | NaN | NaN | |
| 4 | http://youtube.com/channel/UCk8GzjM0rta8yxDcKf... | (http://youtube.co | NaN | NaN | |
| .. | .. | .. | .. | .. | |
| 194 | http://youtube.com/channel/UC5ma-WCc8jNCV0WyRn... | (http://youtube.co | NaN | NaN | |
| 195 | http://youtube.com/channel/UC3j0d7GUMhpgJRBhil... | (http://youtube.co | NaN | NaN | |
| 196 | http://youtube.com/channel/UCwfEOn001DWcyTgzVV... | (http://youtube.co | NaN | NaN | |
| 197 | http://youtube.com/channel/UCx8Z14PpntdaxCt2ha... | (http://youtube.co | NaN | NaN | |
| 198 | http://youtube.com/channel/UCT2X19JJJaJGUN7mrYu... | (http://youtube.co | NaN | NaN | |

Unnamed: 11 Unnamed: 12

| | | |
|-----|-----|-----|
| 0 | NaN | NaN |
| 1 | NaN | NaN |
| 2 | NaN | NaN |
| 3 | NaN | NaN |
| 4 | NaN | NaN |
| .. | .. | .. |
| 194 | NaN | NaN |
| 195 | NaN | NaN |
| 196 | NaN | NaN |

Loading [MathJax]/extensions/MathML/content-mathml.js

```
197      NaN      NaN  
198      NaN      NaN
```

[199 rows x 13 columns]

SORTED FUNCTION

```
In [179]: 1 sorted_df = df.sort_values(by='Rank', ascending=False)
2 print(sorted_df)
```

| Rank | Username | Categories | Suscribers | \ |
|------|-------------------------|---------------------------|------------|-----|
| 198 | 199 AGT | Música y baile, Películas | 25900000 | |
| 197 | 198 TheLallantop | Noticias y Política | 26000000 | |
| 196 | 197 mariliamendoncareal | Música y baile | 26100000 | |
| 195 | 196 nickiminaj | Música y baile | 26100000 | |
| 194 | 195 dianaandromahin8102 | Animación, Juguetes | 26200000 | |
| .. | ... | ... | ... | ... |
| 4 | 5 KidsDianaShow | Animación, Juguetes | 113500000 | |
| 3 | 4 SETIndia | Educación | 162600000 | |
| 2 | 3 CoComelon | Educación | 165500000 | |
| 1 | 2 MrBeast | Videojuegos, Humor | 183500000 | |
| 0 | 1 tseries | Música y baile | 249500000 | |

| | Country | Visits | Likes | Comments | \ |
|-----|----------------|-----------|-------|----------|---|
| 198 | Estados Unidos | 186800 | 398 | 197 | |
| 197 | India | 30800 | 396 | 52 | |
| 196 | Brasil | 0 | 394 | 0 | |
| 195 | Estados Unidos | 1600000 | 392 | 7600 | |
| 194 | Unknown | 109000 | 390 | 0 | |
| .. | ... | ... | ... | ... | |
| 4 | Unknown | 3900000 | 10 | 0 | |
| 3 | India | 15600 | 8 | 9 | |
| 2 | Unknown | 7000000 | 6 | 0 | |
| 1 | Estados Unidos | 117400000 | 4 | 18500 | |
| 0 | India | 86200 | 2 | 78 | |

Links Unnamed: 9 Unnamed:

| | | | | |
|------|--|-----|-----|--|
| 10 \ | | | | |
| 198 | http://youtube.com/channel/UCT2X19JJJaJGUN7mrYu... (http://youtube.co m/channel/UCT2X19JJJaJGUN7mrYu...) | NaN | NaN | |
| 197 | http://youtube.com/channel/UCx8Z14PpntdaxCt2ha... (http://youtube.co m/channel/UCx8Z14PpntdaxCt2ha...) | NaN | NaN | |
| 196 | http://youtube.com/channel/UCwfEOn001DWcyTgzVV... (http://youtube.co m/channel/UCwfEOn001DWcyTgzVV...) | NaN | NaN | |
| 195 | http://youtube.com/channel/UC3j0d7GUMhpgJRBhiL... (http://youtube.co m/channel/UC3j0d7GUMhpgJRBhiL...) | NaN | NaN | |
| 194 | http://youtube.com/channel/UC5ma-WCc8jNCV0WyRn... (http://youtube.co m/channel/UC5ma-WCc8jNCV0WyRn...) | NaN | NaN | |
| .. | ... | ... | ... | |
| .. | ... | ... | ... | |
| 4 | http://youtube.com/channel/Uck8GzjMOrta8yxDcKf... (http://youtube.co m/channel/Uck8GzjMOrta8yxDcKf...) | NaN | NaN | |
| 3 | http://youtube.com/channel/UCpEhnqL0y41EpW2TvW... (http://youtube.co m/channel/UCpEhnqL0y41EpW2TvW...) | NaN | NaN | |
| 2 | http://youtube.com/channel/UCbCmjCuTUZos6Inko4... (http://youtube.co m/channel/UCbCmjCuTUZos6Inko4...) | NaN | NaN | |
| 1 | http://youtube.com/channel/UCX60Q3DkcsbYNE6H8u... (http://youtube.co m/channel/UCX60Q3DkcsbYNE6H8u...) | NaN | NaN | |
| 0 | http://youtube.com/channel/UCq-Fj5jknLsUf-MWSy... (http://youtube.co m/channel/UCq-Fj5jknLsUf-MWSy...) | NaN | NaN | |

Unnamed: 11 Unnamed: 12

| | | |
|-----|-----|-----|
| 198 | NaN | NaN |
| 197 | NaN | NaN |
| 196 | NaN | NaN |
| 195 | NaN | NaN |
| 194 | NaN | NaN |
| .. | ... | ... |
| 4 | NaN | NaN |
| 3 | NaN | NaN |
| 2 | NaN | NaN |

Loading [MathJax]/extensions/MathML/content-mathml.js

2 NaN NaN

```
1           NaN      NaN  
0           NaN      NaN
```

[199 rows x 13 columns]

GROUPBY FUNCTION

```
In [180]: 1 grouped_data = df.groupby('Likes')  
2 result = grouped_data['Likes'].mean()  
3 print(result)
```

```
Likes  
2        2.0  
4        4.0  
6        6.0  
8        8.0  
10       10.0  
...  
390      390.0  
392      392.0  
394      394.0  
396      396.0  
398      398.0  
Name: Likes, Length: 199, dtype: float64
```

AGGREGATION FUNCTIONS

```
In [181]: 1 sum_result = df['Likes'].sum()  
2 sum_result
```

Out[181]: 39800

```
In [182]: 1 sum_result = df['Likes'].min()  
2 sum_result
```

Out[182]: 2

```
In [183]: 1 sum_result = df['Likes'].median()  
2 sum_result
```

Out[183]: 200.0

```
In [184]: 1 sum_result = df['Likes'].median()  
2 sum_result
```

Out[184]: 200.0

```
In [185]: 1 sum_result = df['Likes'].max()  
2 sum_result
```

Out[185]: 398

```
In [186]: 1 sum_result = df['Likes'].count()
2 sum_result
```

Out[186]: 199

```
In [187]: 1 sum_result = df['Likes'].std()
2 sum_result
```

Out[187]: 115.18101695447331

```
In [188]: 1 sum_result = df['Likes'].var()
2 sum_result
```

Out[188]: 13266.666666666666

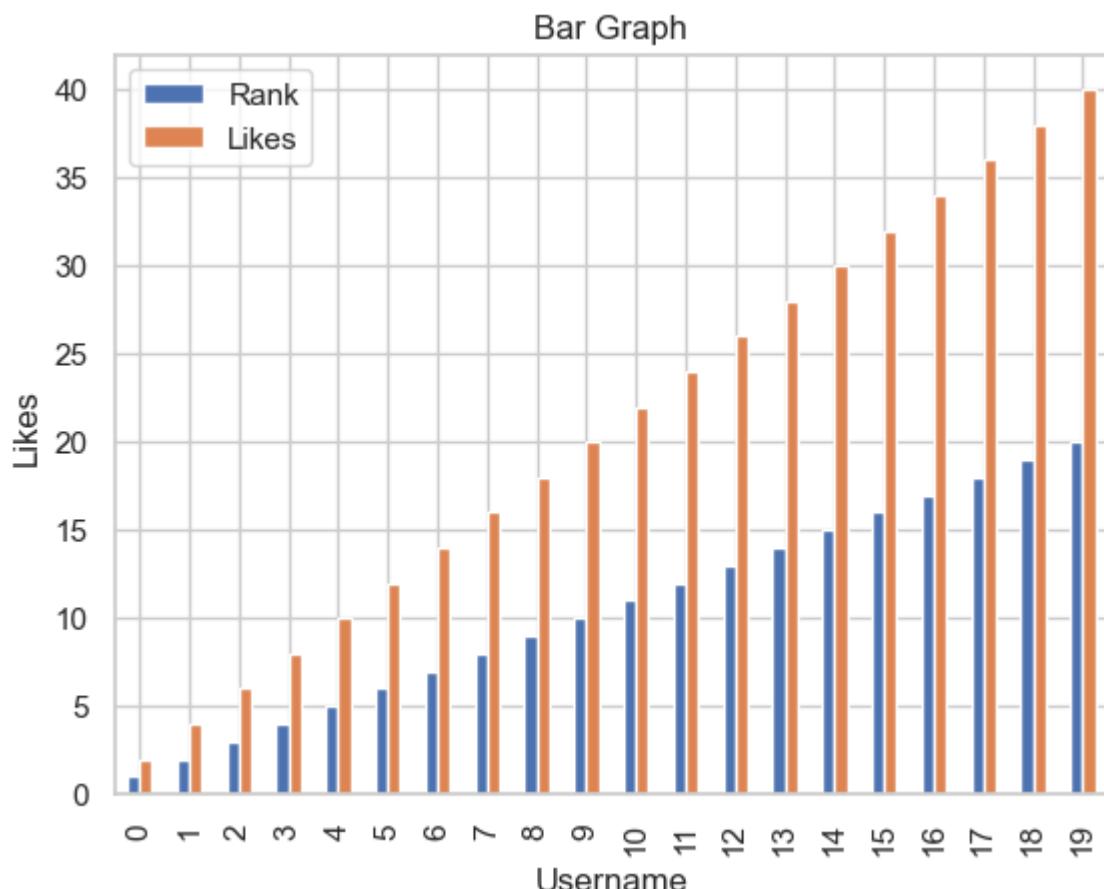
```
In [189]: 1 sum_result = df['Likes'].mode()
2 sum_result
3
```

Out[189]: 0 2
1 4
2 6
3 8
4 10
...
194 390
195 392
196 394
197 396
198 398
Name: Likes, Length: 199, dtype: int64

BAR CHART

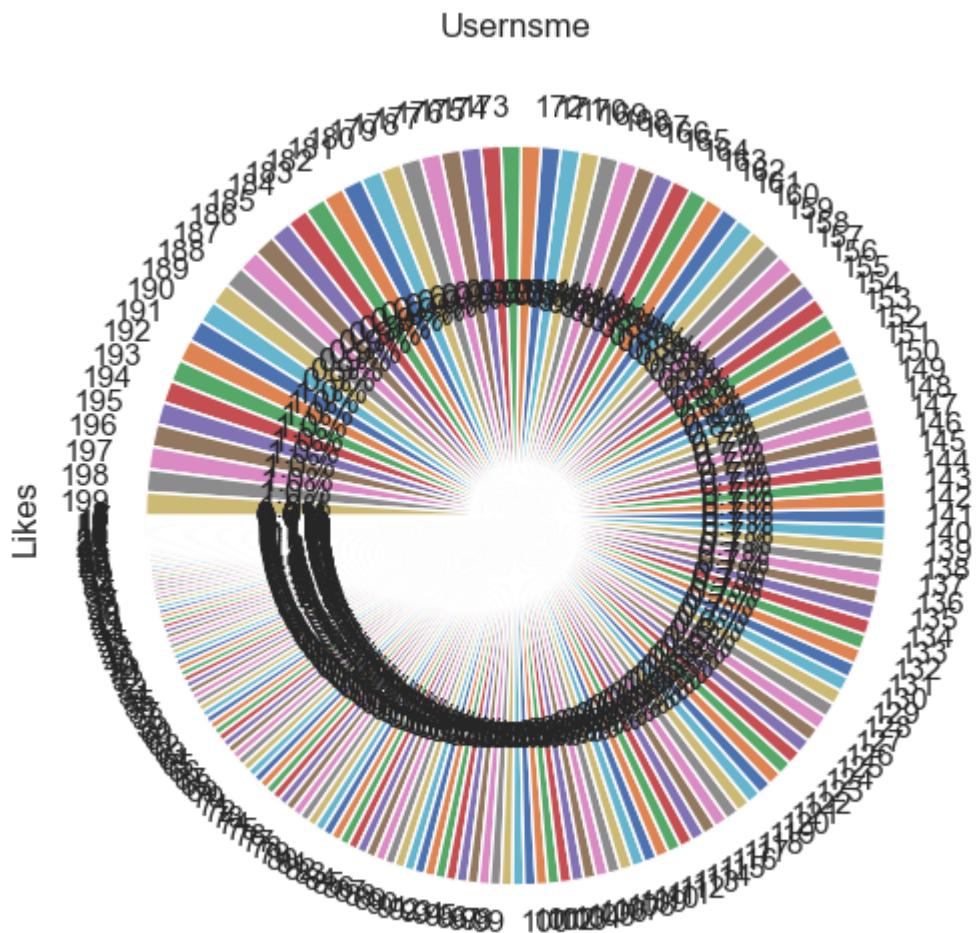
```
In [190]:  
1 import matplotlib.pyplot as plt  
2 bar_data = df.head(20)  
3 plt.figure(figsize=(8, 6))  
4 bar_data.plot(kind='bar')  
5 plt.xlabel('Username')  
6 plt.ylabel('Likes')  
7 plt.title('Bar Graph')  
8 plt.show()
```

<Figure size 800x600 with 0 Axes>



PIE CHART

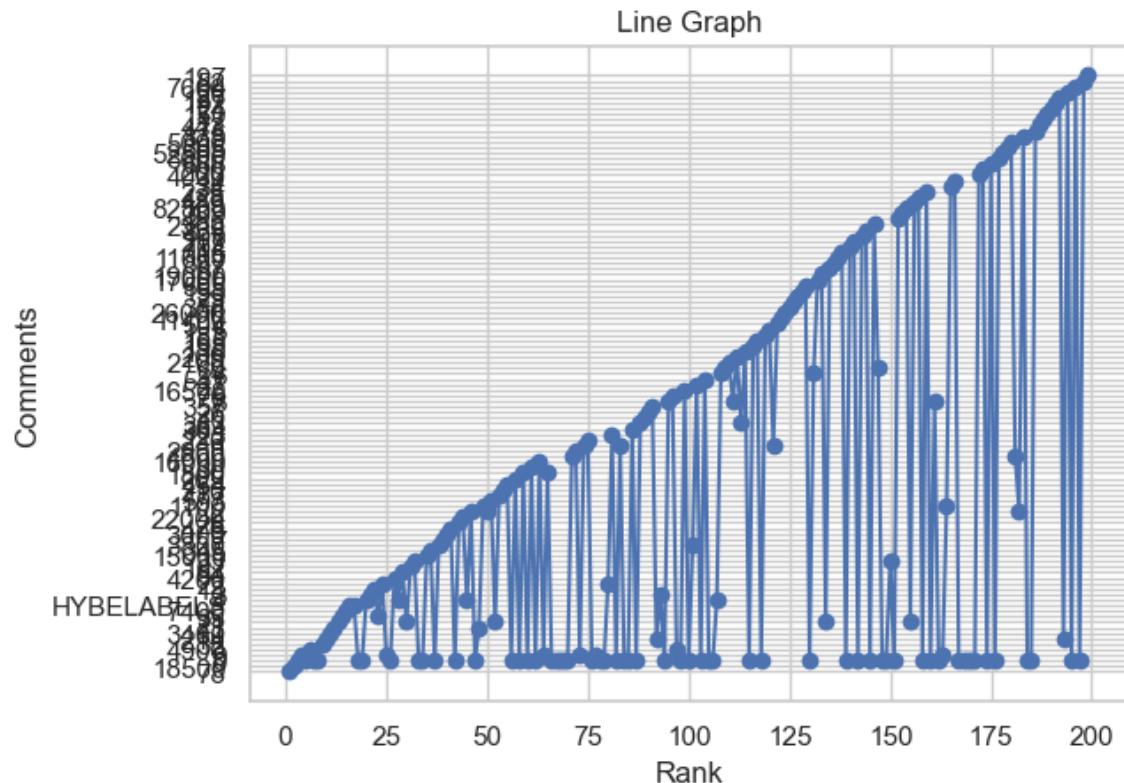
```
In [191]: 1 pie_data = df.groupby('Rank')['Likes'].sum()
2 plt.figure(figsize=(8, 6))
3 pie_data.plot(kind='pie', autopct='%1.1f%%', startangle=180)
4 plt.title('Usernsme')
5 plt.show()
```



LINE GRAPH

In [192]:

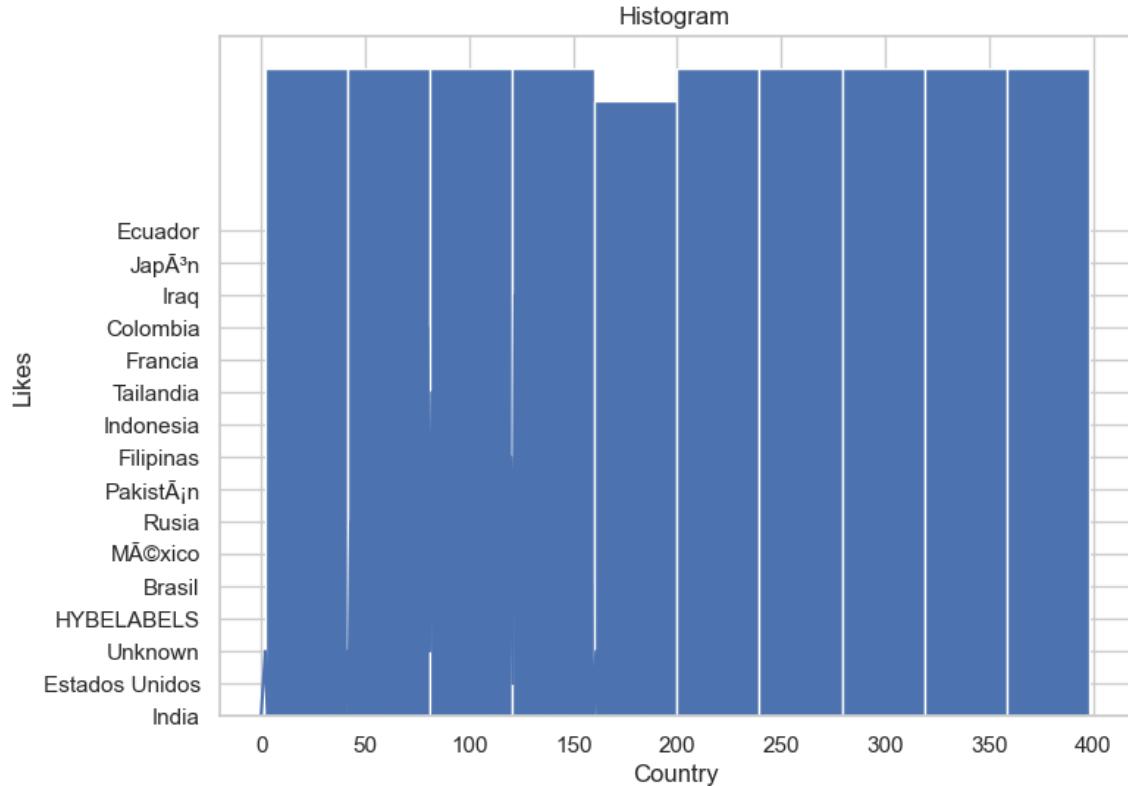
```
1 x_line = df['Rank']
2 y_line = df['Comments']
3 plt.plot(x_line, y_line, marker='o', linestyle='--')
4 plt.xlabel('Rank')
5 plt.ylabel('Comments')
6 plt.title('Line Graph')
7 plt.grid(True)
8 plt.show()
```



HISTOGRAM

In [193]:

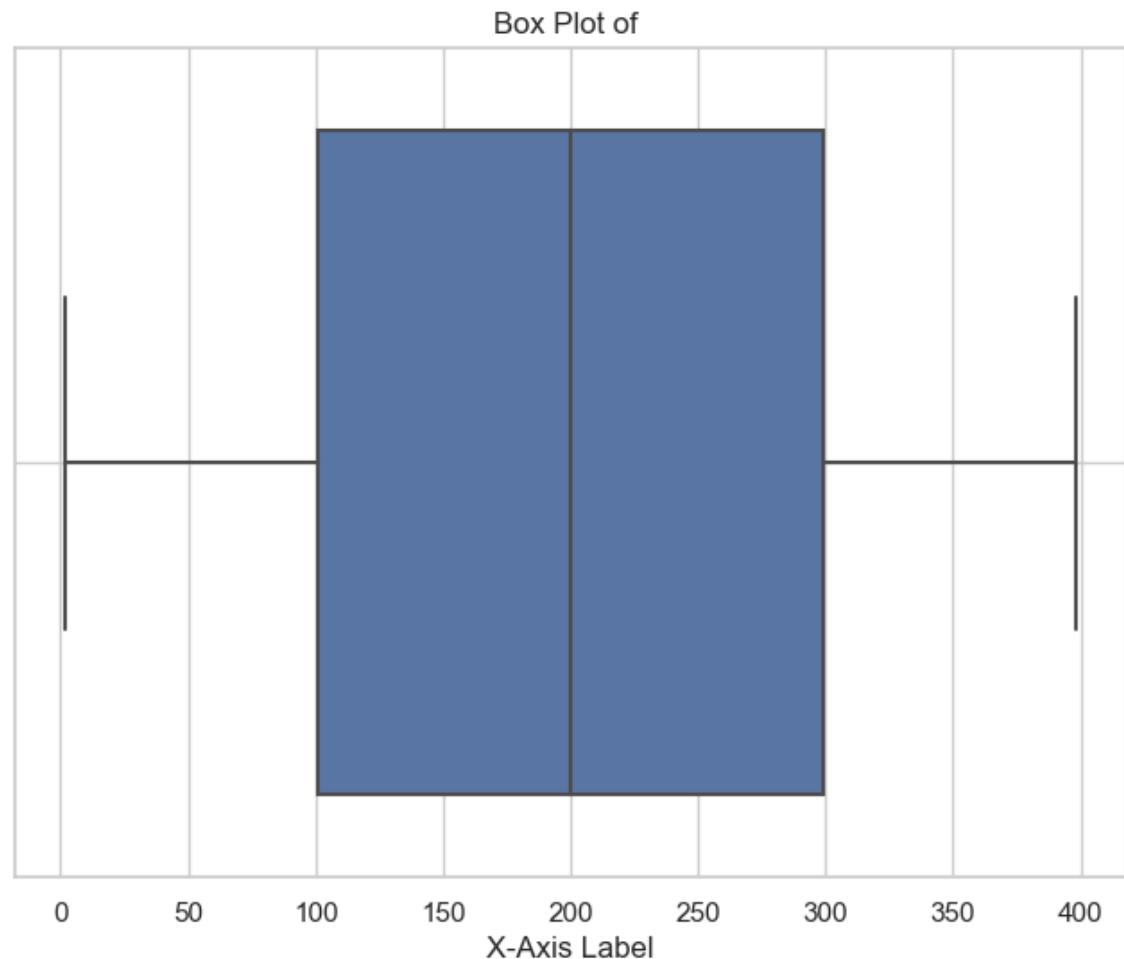
```
1 data_hist =df['Likes']
2 plt.figure(figsize=(8, 6))
3 plt.hist(data_hist, bins=10)
4 plt.xlabel('Country')
5 plt.ylabel('Likes')
6 plt.plot(df.index, df['Country'], label='Likes', color='b')
7
8 # Display the plot
9 plt.grid(True) # Optional: Add grid line
10 plt.title('Histogram')
11 plt.show()
12
```



BOX PLOT

In [194]:

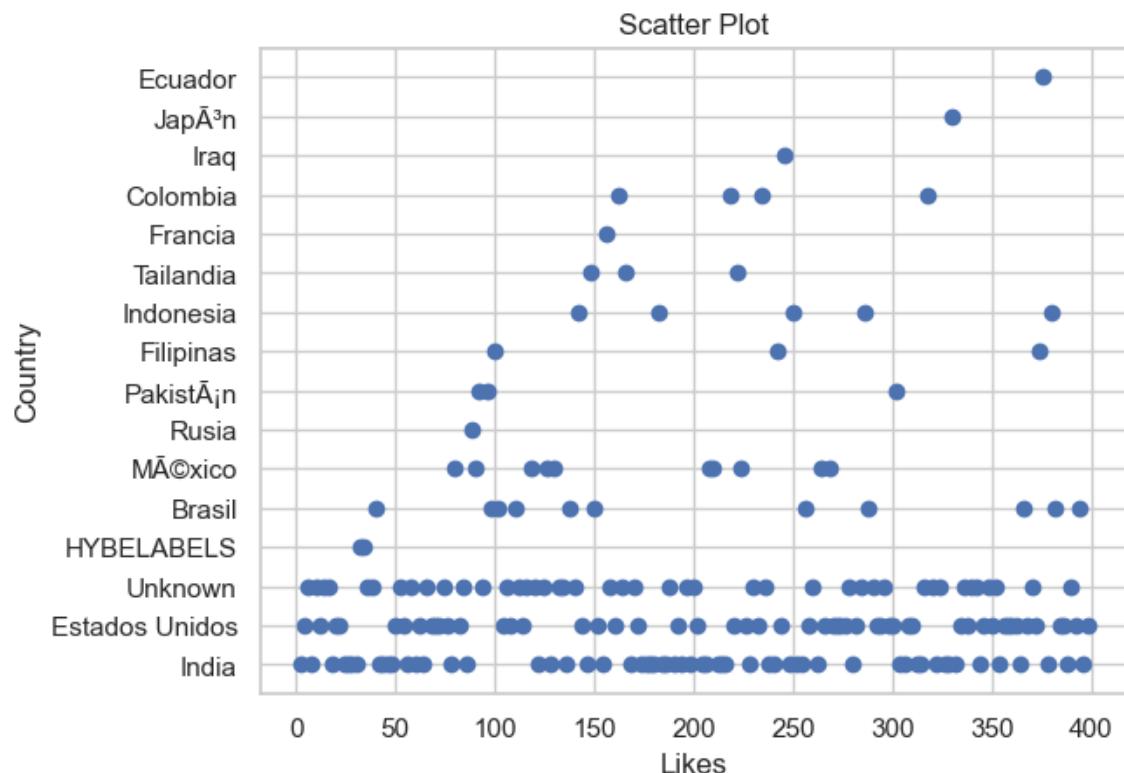
```
1 import seaborn as sns
2 import matplotlib.pyplot as plt
3
4 # Create a box plot
5 plt.figure(figsize=(8, 6))
6 sns.boxplot(x='Likes', data=df)
7 plt.xlabel('X-Axis Label')
8 plt.title('Box Plot of ')
9 plt.grid(True)
10 plt.show()
```



SCATTER PLOT

In [195]:

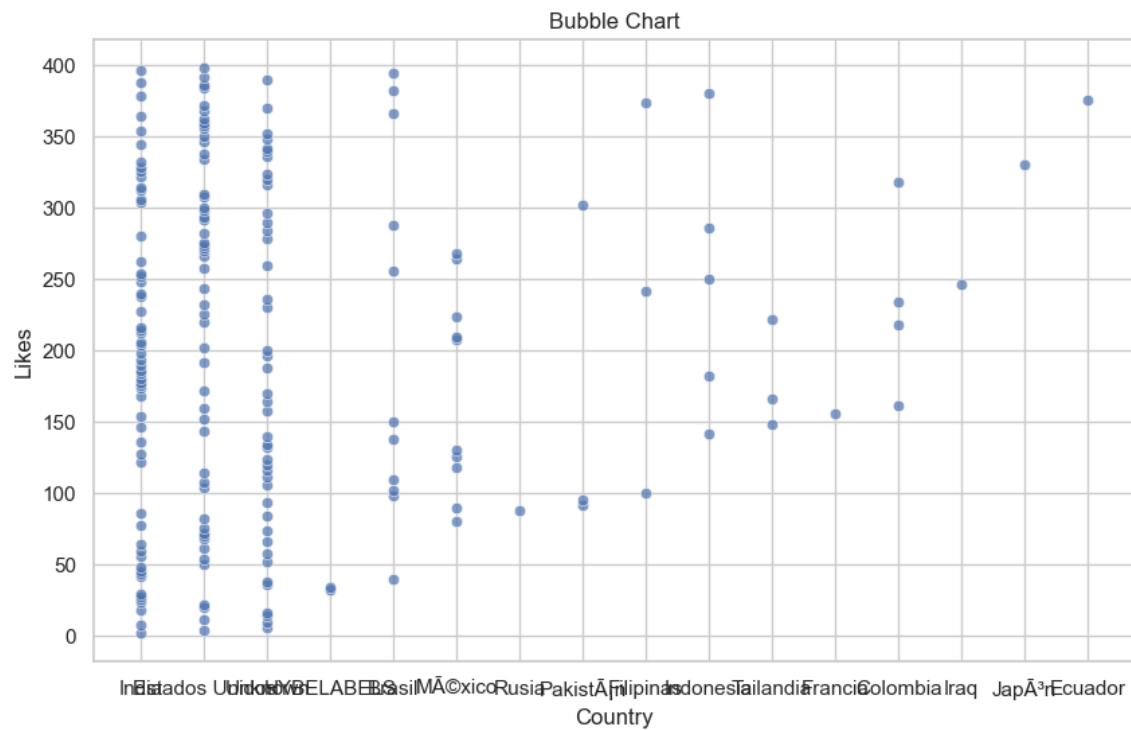
```
1 import pandas as pd
2 import matplotlib.pyplot as plt
3 x_column = 'Likes'
4 y_column = 'Country'
5 # Create a scatter plot
6 plt.scatter(df['Likes'], df['Country'])
7
8 # Add Labels and title
9 plt.xlabel('Likes')
10 plt.ylabel('Country')
11 plt.title('Scatter Plot')
12
13 # Show the plot
14 plt.show()
15
```



BUBBLE CHART

In [196]:

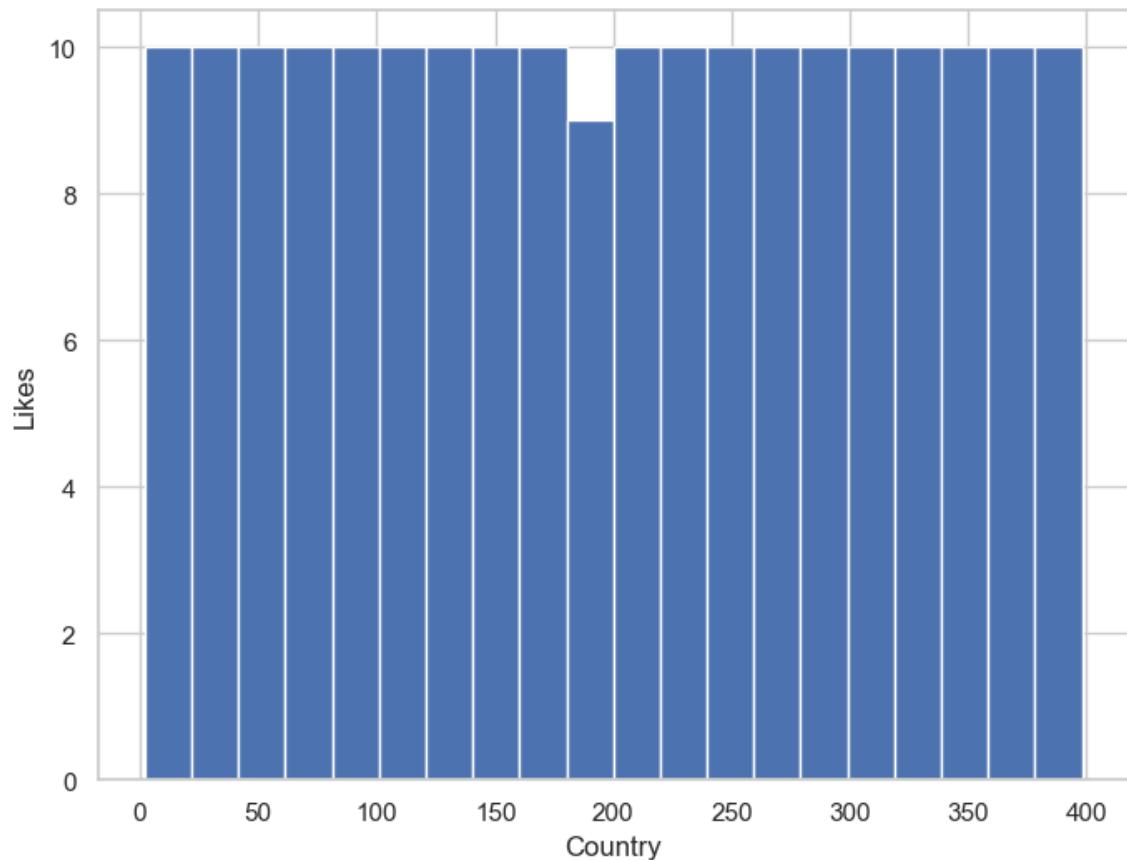
```
1 import pandas as pd
2 import matplotlib.pyplot as plt
3 import seaborn as sns
4
5 # Define the columns for x, y, and bubble size
6 x_column = 'Country' # Replace with the actual column name
7 y_column = 'Likes' # Replace with the actual column name
8 size_column = 'Size_Column' # Replace with the actual column name
9
10 # Create a bubble chart
11 plt.figure(figsize=(10, 6))
12 sns.scatterplot(x=x_column, y=y_column, data=df, sizes=(20, 200), alpha=0.5)
13
14 # Customize the plot properties (title, labels, etc.)
15 plt.title('Bubble Chart')
16 plt.xlabel('Country')
17 plt.ylabel('Likes')
18
19 # Show the bubble chart
20 plt.show()
```



HISTOGRAM

In [197]:

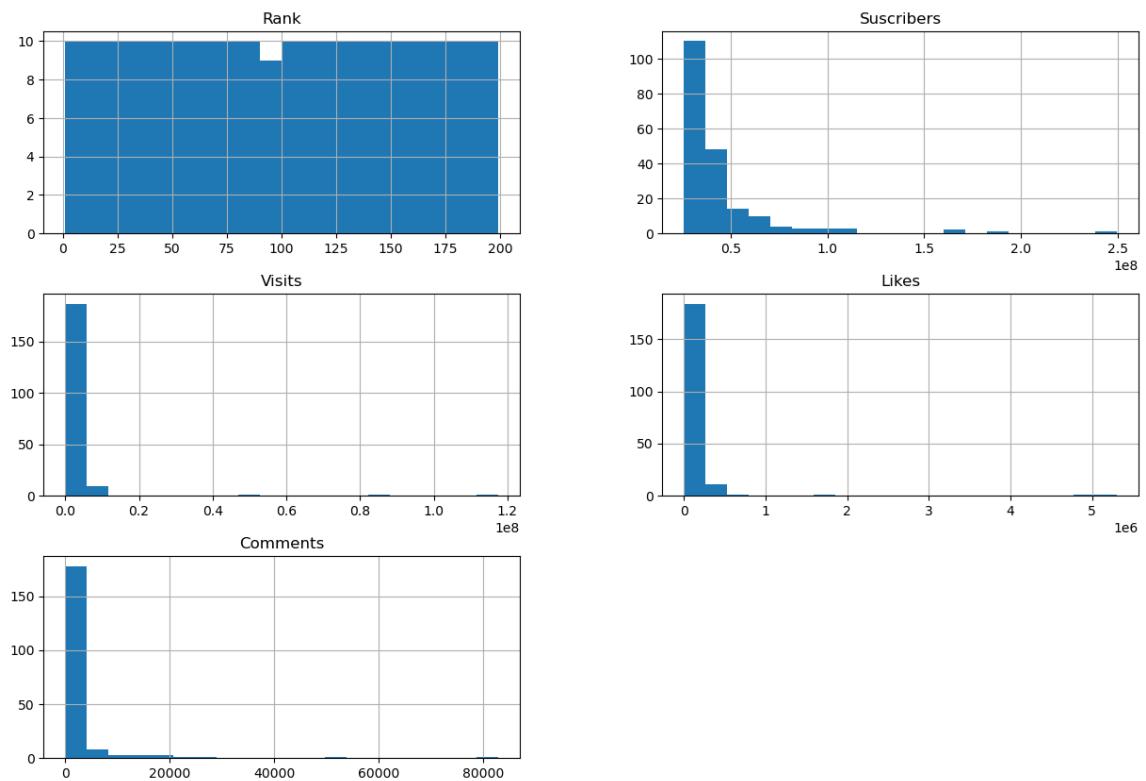
```
1 passenger_age = df['Likes'].dropna()
2 plt.figure(figsize=(8, 6))
3 plt.hist(passenger_age, bins=20)
4 plt.xlabel("Country")
5 plt.ylabel("Likes")
6 plt.grid(True)
7 plt.show()
```



Data exploration

```
In [2]:  
1 # Import necessary libraries  
2 import pandas as pd  
3 import matplotlib.pyplot as plt  
4  
5 # Check the first few rows of the dataset  
6 print(df.head())  
7  
8 # Get an overview of the dataset  
9 print(df.info())  
10  
11 # Generate summary statistics  
12 print(df.describe())  
13  
14 # Count missing values in each column  
15 print(df.isnull().sum())  
16  
17 # Visualize data: Create histograms for numeric columns  
18 df.hist(bins=20, figsize=(15, 10))  
19 plt.show()  
20  
21 # Visualize data: Create scatter plots or other visualizations as needed  
22 # Example: data.plot(x='column_name', y='another_column', kind='scatter')  
# plt.show()  
24
```

```
\_ 0    1      tseries      True  249500000          India   86200
  1    2      MrBeast     False  183500000  Estados Unidos 117400000
  2    3      CoComelon    True  165500000          Unknown 7000000
  3    4      SETIndia    False  162600000          India   15600
  4    5  KidsDianaShow  True  113500000          Unknown 3900000
    Likes  Comments          Links
  0  2700      78  http://youtube.com/channel/UCq-Fj5jknLsUF-MWSy... (h
  http://youtube.com/channel/UCq-Fj5jknLsUF-MWSy...)
  1 5300000     18500  http://youtube.com/channel/UCX60Q3DkcsbYNE6H8u... (h
  http://youtube.com/channel/UCX60Q3DkcsbYNE6H8u...)
  2  24700        0  http://youtube.com/channel/UCbCmjCuTUZos6Inko4... (h
  http://youtube.com/channel/UCbCmjCuTUZos6Inko4...)
  3    166        9  http://youtube.com/channel/UCpEhnqL0y41EpW2TvW... (h
  http://youtube.com/channel/UCpEhnqL0y41EpW2TvW...)
  4  12400        0  http://youtube.com/channel/Uck8GzjMOrta8yxDcKf... (h
  http://youtube.com/channel/Uck8GzjMOrta8yxDcKf...)
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 199 entries, 0 to 198
Data columns (total 9 columns):
 #  Column      Non-Null Count  Dtype  
--- 
 0  Rank         199 non-null    int64  
 1  Username     199 non-null    object  
 2  Categories   199 non-null    bool    
 3  Suscribers  199 non-null    int64  
 4  Country      199 non-null    object  
 5  Visits       199 non-null    int64  
 6  Likes        199 non-null    int64  
 7  Comments     199 non-null    int64  
 8  Links        199 non-null    object  
dtypes: bool(1), int64(5), object(3)
memory usage: 12.8+ KB
None
      Rank  Suscribers  Visits  Likes  Comments
count  199.00000  1.990000e+02  1.990000e+02  1.990000e+02  199.00000
mean   100.00000  4.405226e+07  2.261227e+06  1.047747e+05  2096.356784
std    57.590508  2.753718e+07  1.084192e+07  5.321669e+05  7913.381544
min    1.000000  2.590000e+07  0.000000e+00  0.000000e+00  0.000000
25%    50.500000  3.030000e+07  3.225000e+04  5.845000e+02  0.000000
50%    100.00000  3.550000e+07  2.257000e+05  3.900000e+03  28.000000
75%    149.50000  4.500000e+07  1.050000e+06  3.465000e+04  375.000000
max    199.00000  2.495000e+08  1.174000e+08  5.300000e+06  82800.00000
Rank      0
Username  0
Categories 0
Suscribers 0
Country    0
Visits     0
Likes      0
Comments   0
Links      0
dtype: int64
```



Correlation & Covariance

Loading [MathJax]/extensions/MathML/content-mathml.js

```
In [7]: 1 import pandas as pd
2 import numpy as np
3
4 # Calculate Correlation
5 correlation_matrix = df.corr()
6 print("Correlation Matrix:")
7 print(correlation_matrix)
8
9 # Calculate Covariance
10 covariance_matrix = df.cov()
11 print("\nCovariance Matrix:")
12 print(covariance_matrix)
```

Correlation Matrix:

| | Rank | Categories | Suscribers | Visits | Likes | \ |
|------------|---------------|---------------|------------|-----------|-----------|---|
| Rank | 1.000000e+00 | 4.338515e-16 | -0.677636 | -0.044936 | -0.032294 | |
| Categories | 4.338515e-16 | 1.000000e+00 | 0.017224 | -0.048342 | -0.020632 | |
| Suscribers | -6.776364e-01 | 1.722411e-02 | 1.000000 | 0.255249 | 0.217098 | |
| Visits | -4.493650e-02 | -4.834173e-02 | 0.255249 | 1.000000 | 0.977950 | |
| Likes | -3.229386e-02 | -2.063200e-02 | 0.217098 | 0.977950 | 1.000000 | |
| Comments | 6.149461e-02 | -8.054422e-02 | 0.003749 | 0.460614 | 0.405251 | |

| | Comments |
|------------|-----------|
| Rank | 0.061495 |
| Categories | -0.080544 |
| Suscribers | 0.003749 |
| Visits | 0.460614 |
| Likes | 0.405251 |
| Comments | 1.000000 |

Covariance Matrix:

| | Rank | Categories | Suscribers | Visits | \ |
|------------|---------------|---------------|---------------|---------------|---|
| Rank | 3.316667e+03 | 5.382900e-16 | -1.074650e+09 | -2.805797e+07 | |
| Categories | 5.382900e-16 | 2.512563e-01 | 2.377468e+05 | -2.627161e+05 | |
| Suscribers | -1.074650e+09 | 2.377468e+05 | 7.582961e+14 | 7.620610e+13 | |
| Visits | -2.805797e+07 | -2.627161e+05 | 7.620610e+13 | 1.175472e+14 | |
| Likes | -9.897345e+05 | -5.503610e+03 | 3.181439e+12 | 5.642487e+12 | |
| Comments | 2.802529e+04 | -3.194883e+02 | 8.168526e+08 | 3.951894e+10 | |

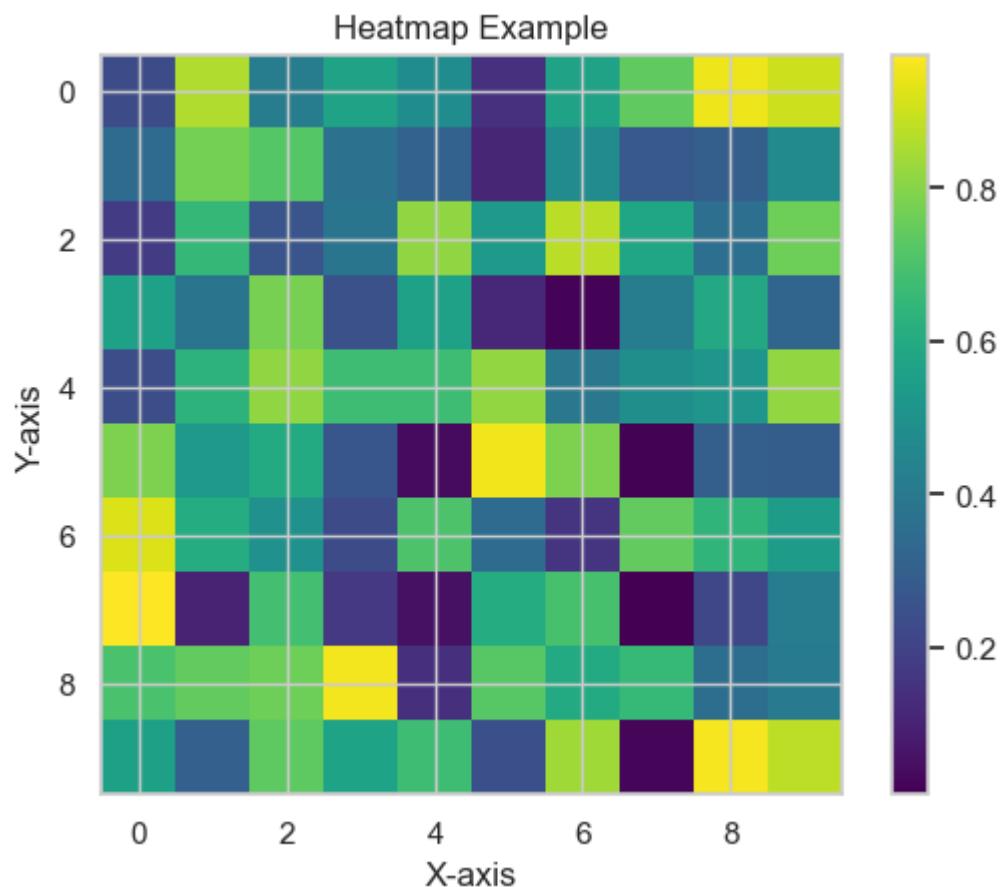
| | Likes | Comments |
|------------|---------------|---------------|
| Rank | -9.897345e+05 | 2.802529e+04 |
| Categories | -5.503610e+03 | -3.194883e+02 |
| Suscribers | 3.181439e+12 | 8.168526e+08 |
| Visits | 5.642487e+12 | 3.951894e+10 |
| Likes | 2.832016e+11 | 1.706609e+09 |
| Comments | 1.706609e+09 | 6.262161e+07 |

In []: 1

HEATMAP

In [229]:

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3
4 # Sample data for the heatmap
5 data = np.random.random((10, 10)) # Replace this with your own data
6
7 # Create a heatmap
8 plt.imshow(data, cmap='viridis') # You can choose a different colormap
9 plt.colorbar() # Add a colorbar for reference
10
11 plt.title('Heatmap Example')
12 plt.xlabel('X-axis')
13 plt.ylabel('Y-axis')
14
15 # Show the heatmap
16 plt.show()
```

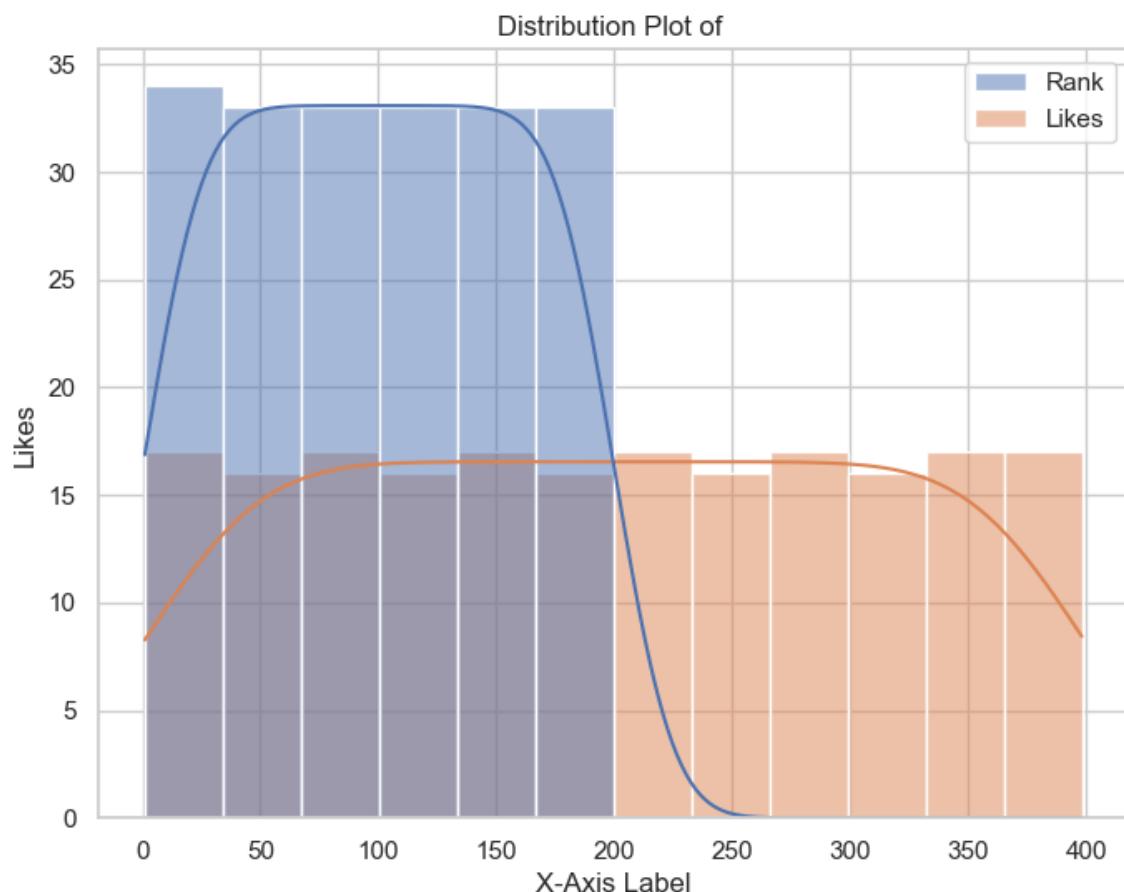


DISTREBUTION PLOT

Loading [MathJax]/extensions/MathML/content-mathml.js

In [199]:

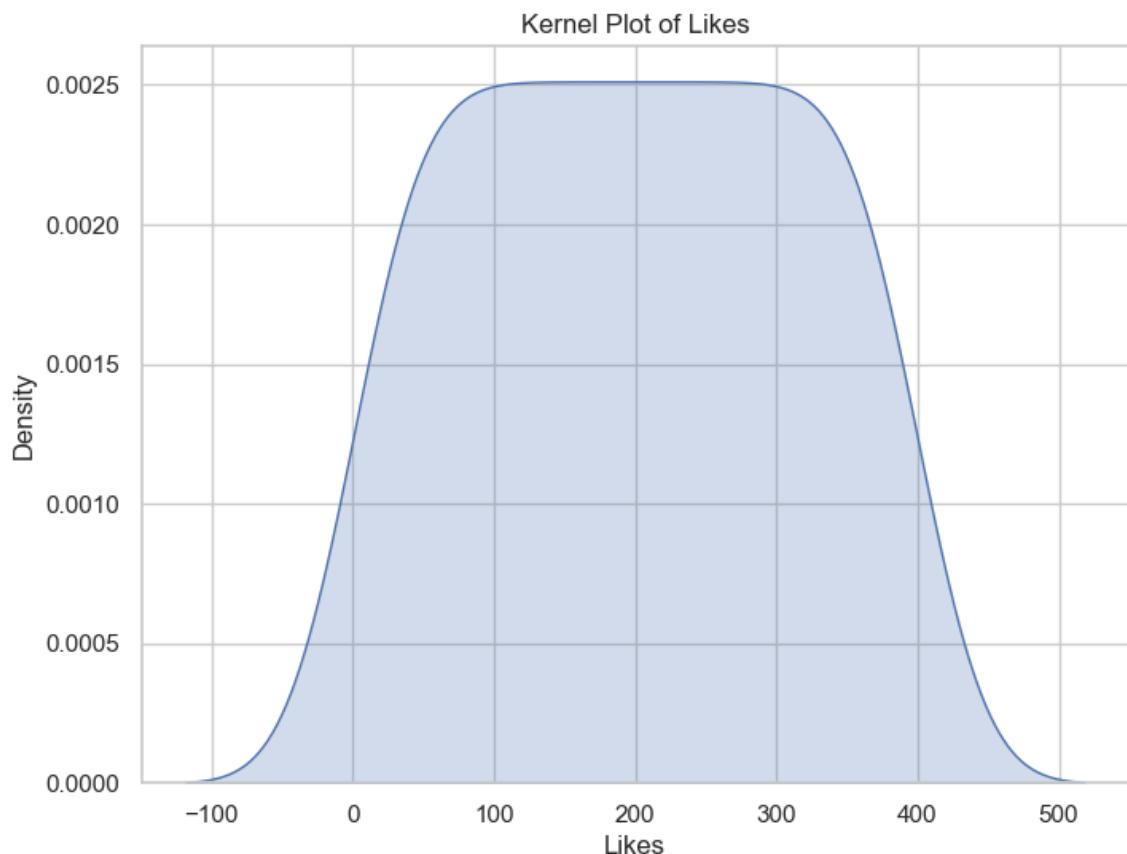
```
1 import pandas as pd
2 import seaborn as sns
3 import matplotlib.pyplot as plt
4
5 # Extract the variable for which you want to create a distribution plot
6 variable_to_plot = 'Rank' # Replace 'ColumnName' with the actual column
7
8 # Create a distribution plot using Seaborn
9 plt.figure(figsize=(8, 6)) # Optional: Set the figure size
10 sns.histplot(data=df, kde=True) # 'kde=True' adds a kernel density estimate
11
12 # Customize the plot (labels, title, etc.)
13 plt.xlabel('X-Axis Label')
14 plt.ylabel('Likes')
15 plt.title('Distribution Plot of ')
16
17 # Show the distribution plot
18 plt.show()
```



KERNEL PLOT

In [200]:

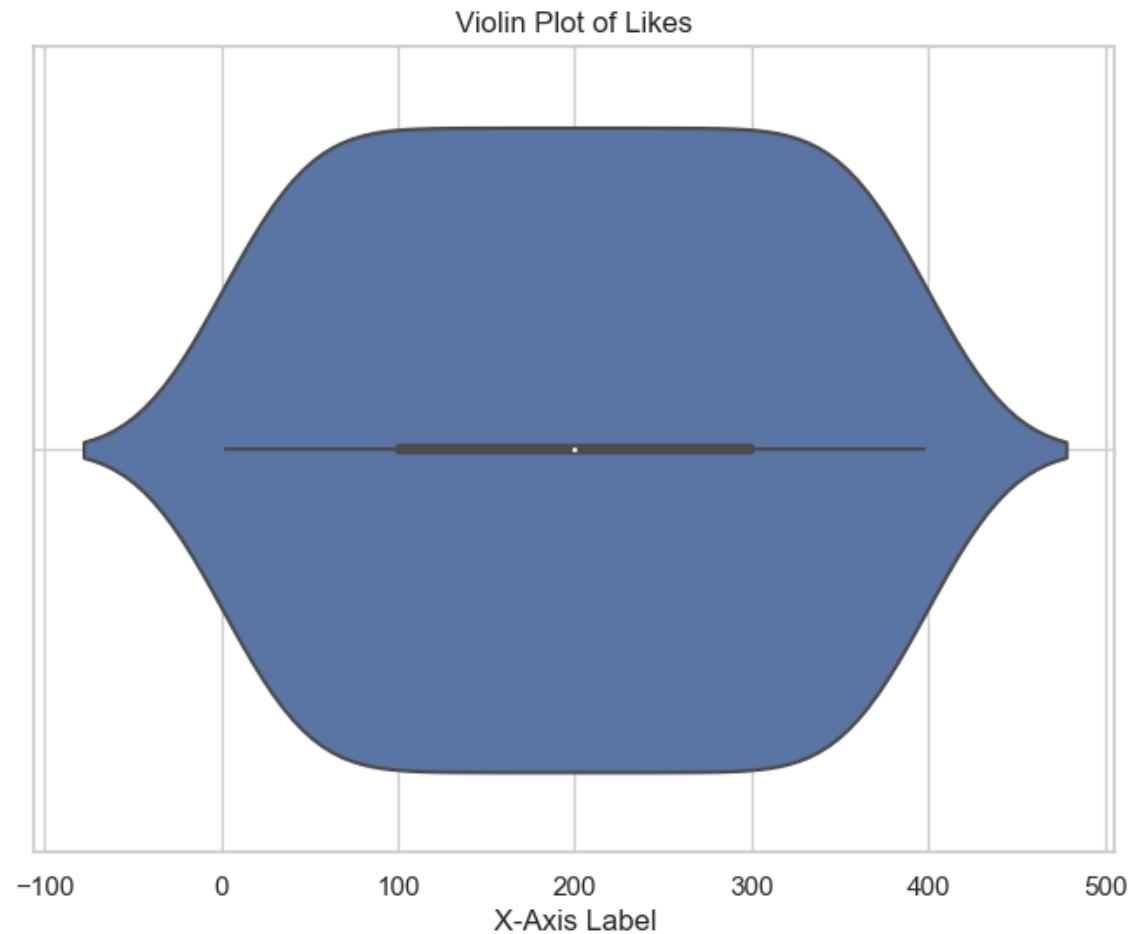
```
1 import seaborn as sns
2 import matplotlib.pyplot as plt
3
4 # Extract the variable for which you want to create a KDE plot
5 variable_to_plot = 'Likes' # Replace 'ColumnName' with the actual column name
6
7 # Create a KDE plot
8 plt.figure(figsize=(8, 6))
9 sns.kdeplot(df['Likes'], shade=True)
10 plt.xlabel('Likes')
11 plt.ylabel('Density')
12 plt.title('Kernel Plot of ' + variable_to_plot)
13 plt.grid(True)
14 plt.show()
```



VIOLIN PLOT

In [201]:

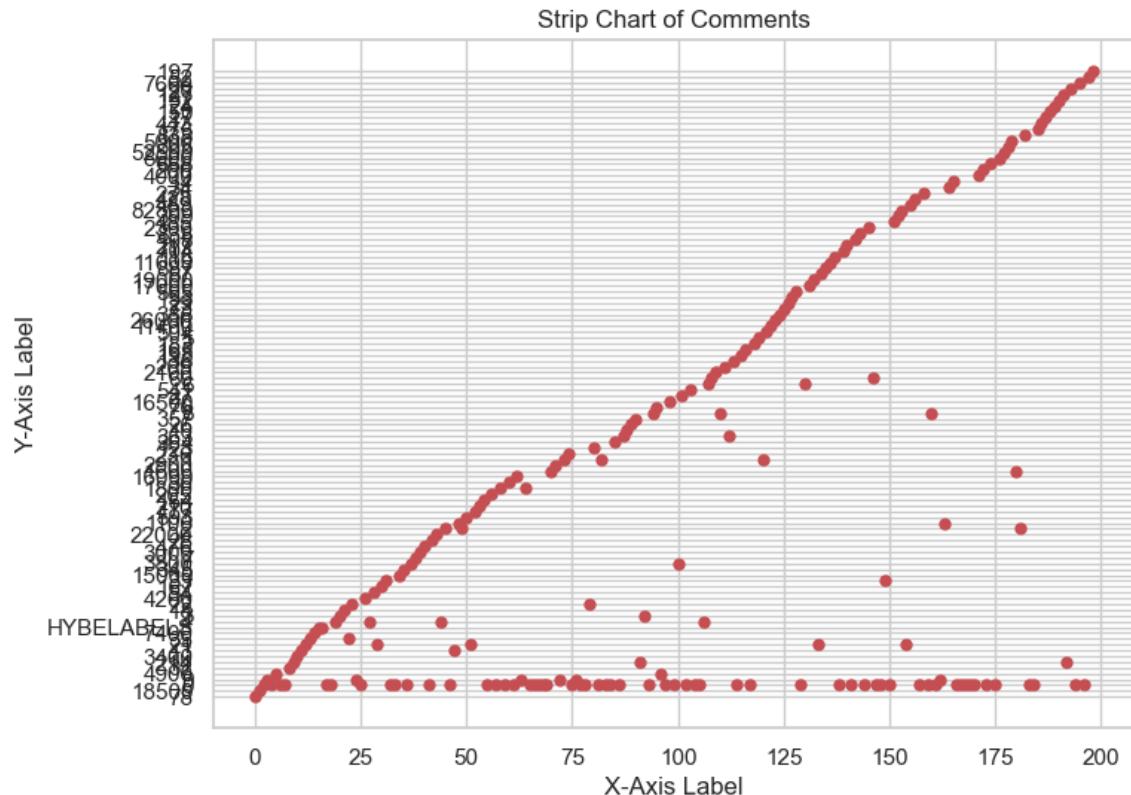
```
1 import seaborn as sns
2 import matplotlib.pyplot as plt
3
4 # Create a violin plot
5 plt.figure(figsize=(8, 6))
6 sns.violinplot(x='Likes', data=df)
7 plt.xlabel('X-Axis Label')
8 plt.title('Violin Plot of ' + variable_to_plot)
9 plt.grid(True)
10 plt.show()
```



STRIP CHART

In [202]:

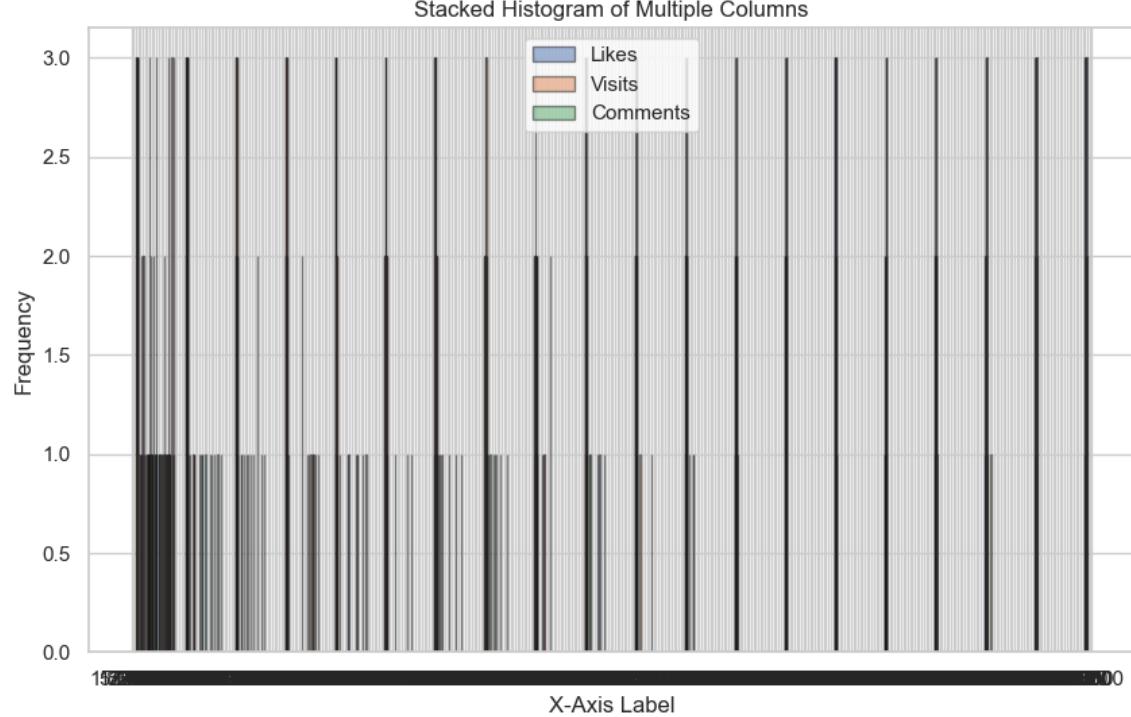
```
1 import pandas as pd
2 import matplotlib.pyplot as plt
3 # Extract the variable for the strip chart
4 variable_to_plot = 'Comments' # Replace 'ColumnName' with the actual column name
5
6 # Create a strip chart
7 plt.figure(figsize=(8, 6))
8 plt.plot(df[variable_to_plot], 'ro', markersize=5) # 'ro' means red circles
9 plt.xlabel('X-Axis Label')
10 plt.ylabel('Y-Axis Label')
11 plt.title('Strip Chart of ' + variable_to_plot)
12 plt.grid(True)
13 plt.show()
```



In [230]:

```
1 import pandas as pd
2 import matplotlib.pyplot as plt
3
4 # Define the columns you want to create a stacked histogram for
5 columns_to_plot = ['Likes', 'Visits', 'Comments'] # Replace with your
6
7 # Create a stacked histogram
8 plt.figure(figsize=(10, 6))
9 plt.hist(df[columns_to_plot].values.T, bins=20, edgecolor='k', alpha=0.
10 plt.xlabel('X-Axis Label')
11 plt.ylabel('Frequency')
12 plt.title('Stacked Histogram of Multiple Columns')
13 plt.legend()
14 plt.grid(True)
15 plt.show()
```

C:\Users\Anusha V\anaconda3\lib\site-packages\IPython\core\pylabtools.py:1
51: UserWarning: Creating legend with loc="best" can be slow with large amounts of data.
fig.canvas.print_figure(bytes_io, **kw)

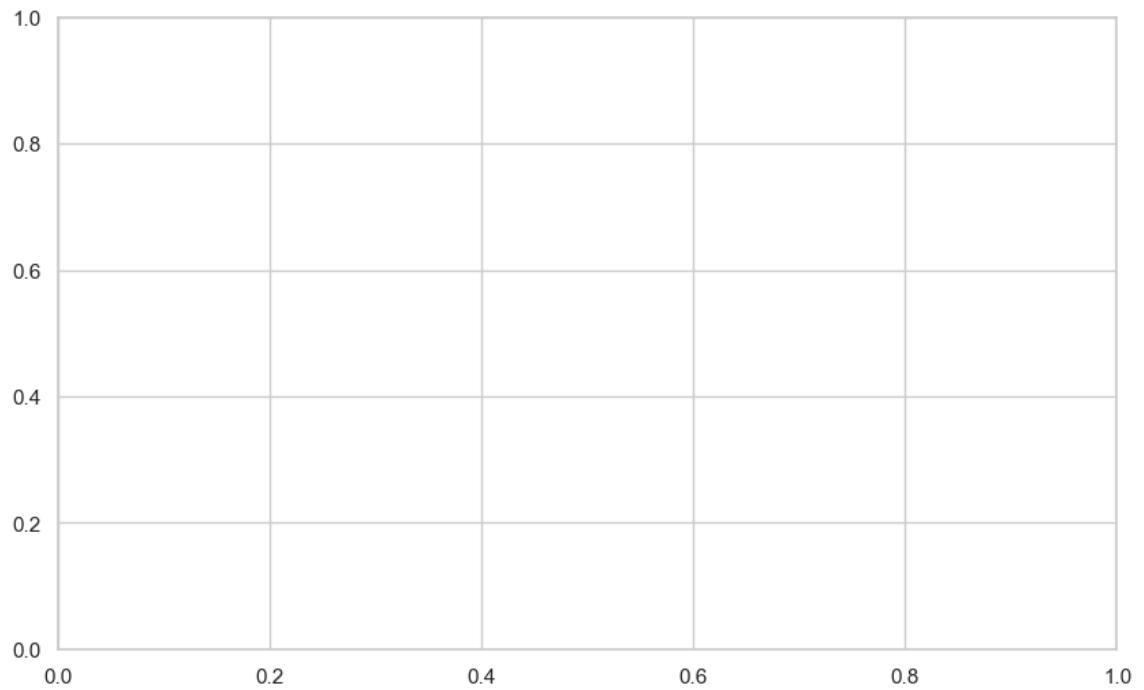


Overlapping Density Plots

In [231]:

```
1 import pandas as pd
2 import seaborn as sns
3 import matplotlib.pyplot as plt
4 # Define the columns you want to create overlapping density plots for
5 columns_to_plot = ['Likes', 'Rank'] # Replace with your column names
6
7 # Create overlapping density plots
8 plt.figure(figsize=(10, 6))
9 for column in columns_to_plot:
10     sns.kdeplot(df[column], label=column, shade=True)
11
12 plt.xlabel('X-Axis Label')
13 plt.ylabel('Density')
14 plt.title('Overlapping Density Plots of Multiple Columns')
15 plt.legend()
16 plt.grid(True)
17 plt.show()
```

```
-  
TypeError  
t)  
~\AppData\Local\Temp\ipykernel_12064\2074536078.py in <module>  
    8 plt.figure(figsize=(10, 6))  
    9 for column in columns_to_plot:  
---> 10     sns.kdeplot(df[column], label=column, shade=True)  
    11  
    12 plt.xlabel('X-Axis Label')  
  
~\anaconda3\lib\site-packages\seaborn\_decorators.py in inner_f(*args, **k  
wargs)  
    44             )  
    45         kwargs.update({k: arg for k, arg in zip(sig.parameters, ar  
gs)})  
---> 46     return f(**kwargs)  
    47     return inner_f  
    48  
  
~\anaconda3\lib\site-packages\seaborn\distributions.py in kdeplot(x, y, sh  
ade, vertical, kernel, bw, gridsize, cut, clip, legend, cumulative, shade_  
lowest, cbar, cbar_ax, cbar_kws, ax, weights, hue, palette, hue_order, hue  
_norm, multiple, common_norm, common_grid, levels, thresh, bw_method, bw_a  
djust, log_scale, color, fill, data, data2, warn_singular, **kwargs)  
    1760     )  
    1761  
-> 1762     p._attach(ax, allowed_types=["numeric", "datetime"], log_scale  
=log_scale)  
    1763  
    1764     if p.univariate:  
  
~\anaconda3\lib\site-packages\seaborn\_core.py in _attach(self, obj, allow  
ed_types, log_scale)  
    1123             f"{allowed_types} is required"  
    1124         )  
-> 1125         raise TypeError(err)  
    1126  
    1127         # Register with the matplotlib unit conversion machine  
ry  
  
TypeError: The x variable is categorical, but one of ['numeric', 'datetim  
e'] is required
```

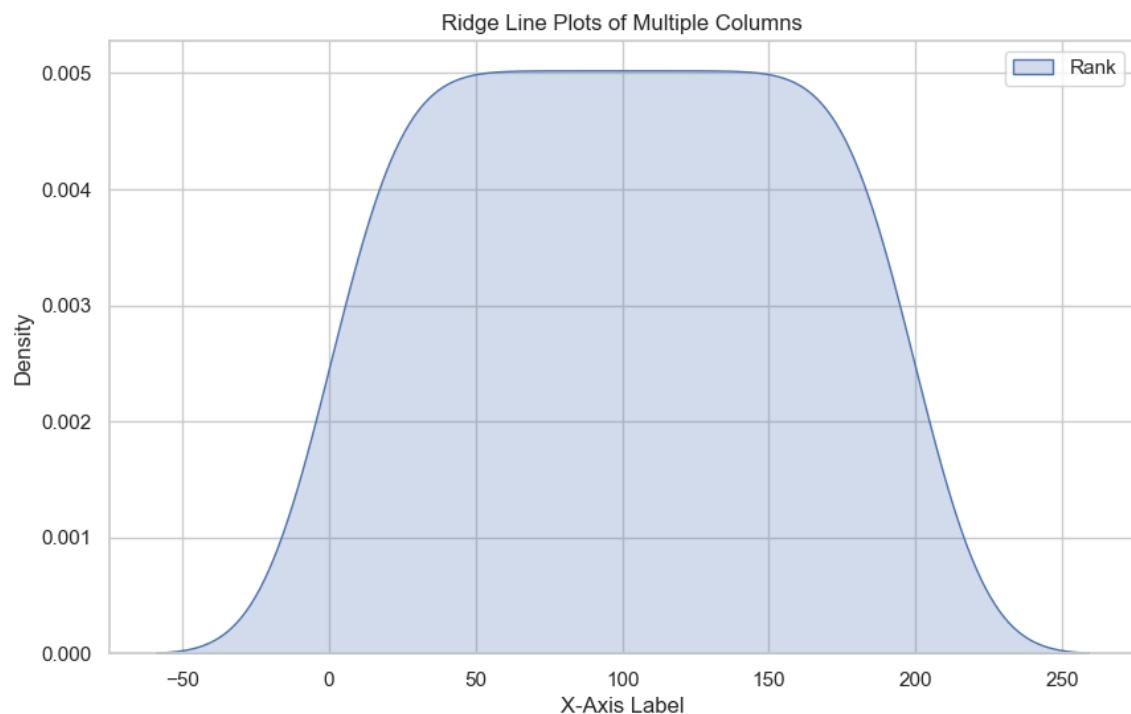


Ridge Line Plots

Loading [MathJax]/extensions/MathML/content-mathml.js

In [235]:

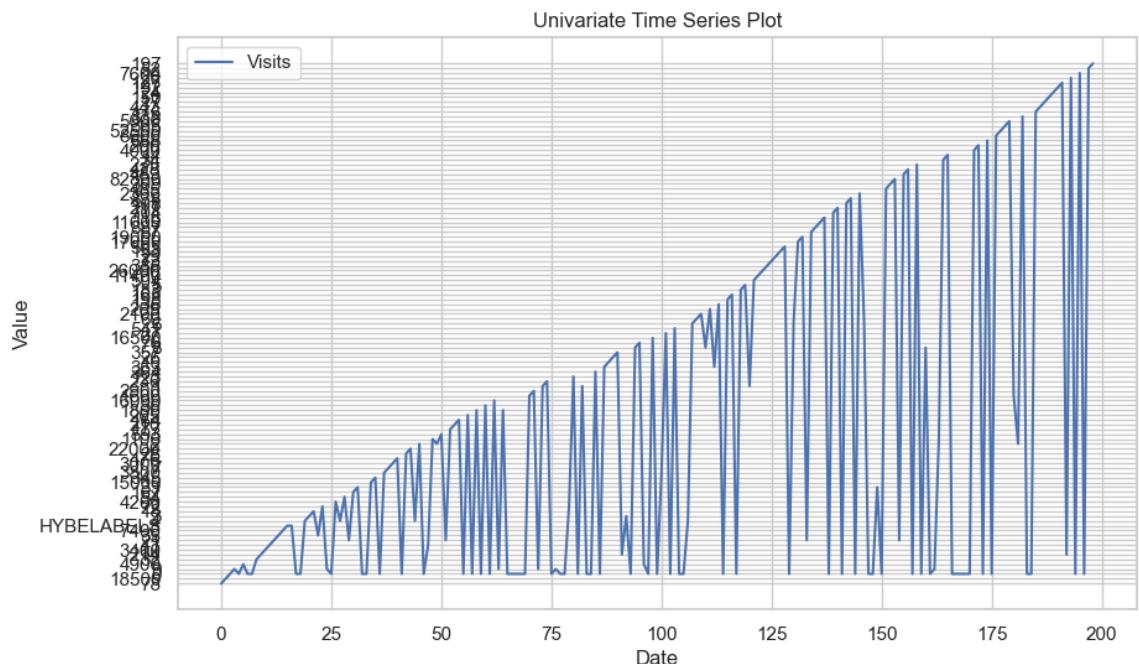
```
1 import pandas as pd
2 import seaborn as sns
3 import matplotlib.pyplot as plt
4 # Define the columns you want to create ridge line plots for
5 columns_to_plot = ['Rank'] # Replace with your column names
6
7 # Create ridge line plots
8 plt.figure(figsize=(10, 6))
9 sns.set(style='whitegrid')
10 for column in columns_to_plot:
11     sns.kdeplot(df[column], label=column, shade=True)
12
13 plt.xlabel('X-Axis Label')
14 plt.ylabel('Density')
15 plt.title('Ridge Line Plots of Multiple Columns')
16 plt.legend()
17 plt.grid(True)
18 plt.show()
```



Univariate Time Series Plot

In [232]:

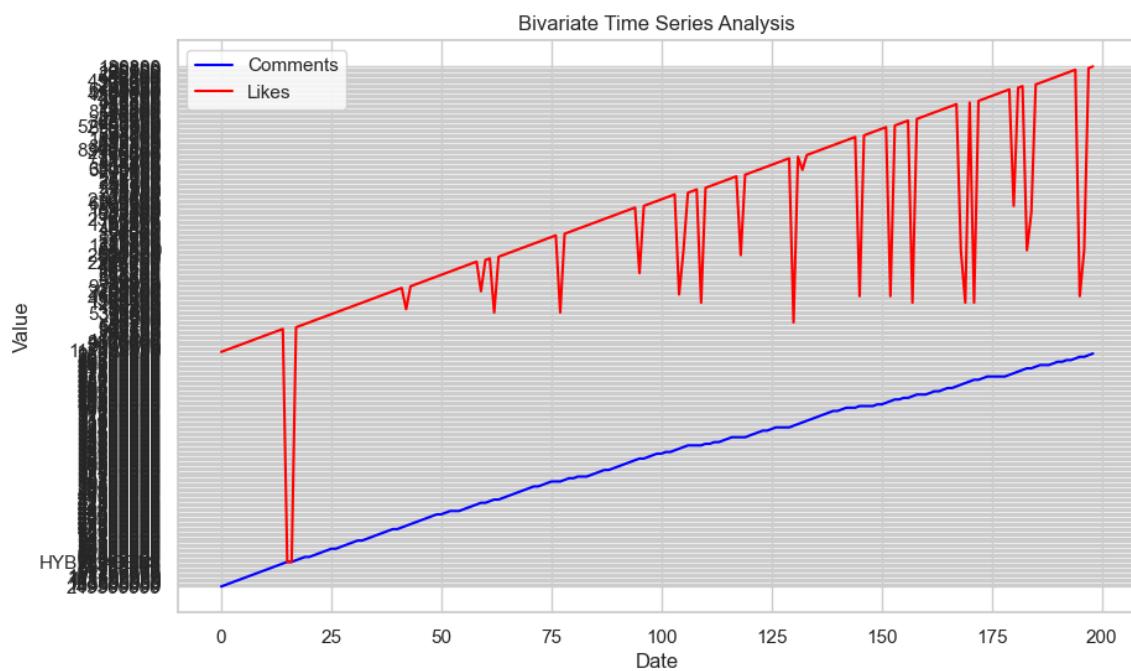
```
1 plt.figure(figsize=(10, 6)) # Optional: Set the figure size
2
3 # Plot the time series data
4 plt.plot(df.index, df['Comments'], label='Visits', color='b')
5
6 # Optional: Add Labels and a title
7 plt.xlabel('Date')
8 plt.ylabel('Value')
9 plt.title('Univariate Time Series Plot')
10
11 # Optional: Customize other plot properties
12
13 # Display the plot
14 plt.grid(True) # Optional: Add grid lines
15 plt.legend() # Optional: Show the Legend if you have multiple series
16 plt.show()
```



Bivariate Time Series Analysis

In [233]:

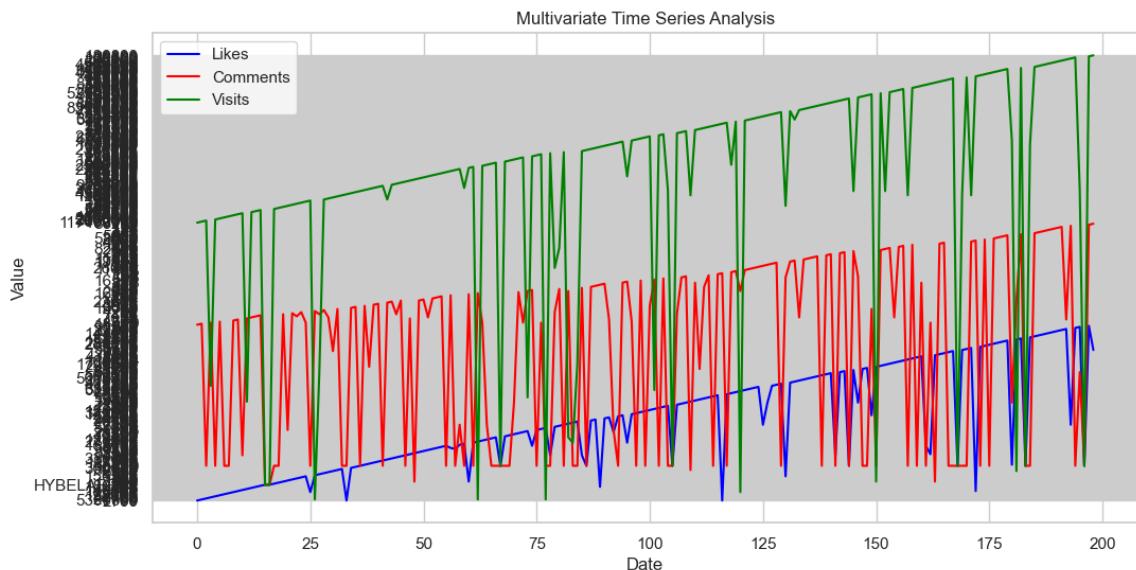
```
1 import pandas as pd
2 import matplotlib.pyplot as plt
3
4
5 # Plot the time series data
6 plt.figure(figsize=(10, 6))
7 plt.plot(df.index, df['Suscribers'], label='Comments', color='blue')
8 plt.plot(df.index, df['Visits'], label='Likes', color='red')
9
10 # Customize the plot
11 plt.xlabel('Date')
12 plt.ylabel('Value')
13 plt.title('Bivariate Time Series Analysis')
14 plt.legend()
15 plt.grid(True)
16
17 # Show the plot
18 plt.show()
```



Multivariate Time Series Analysis

In [243]:

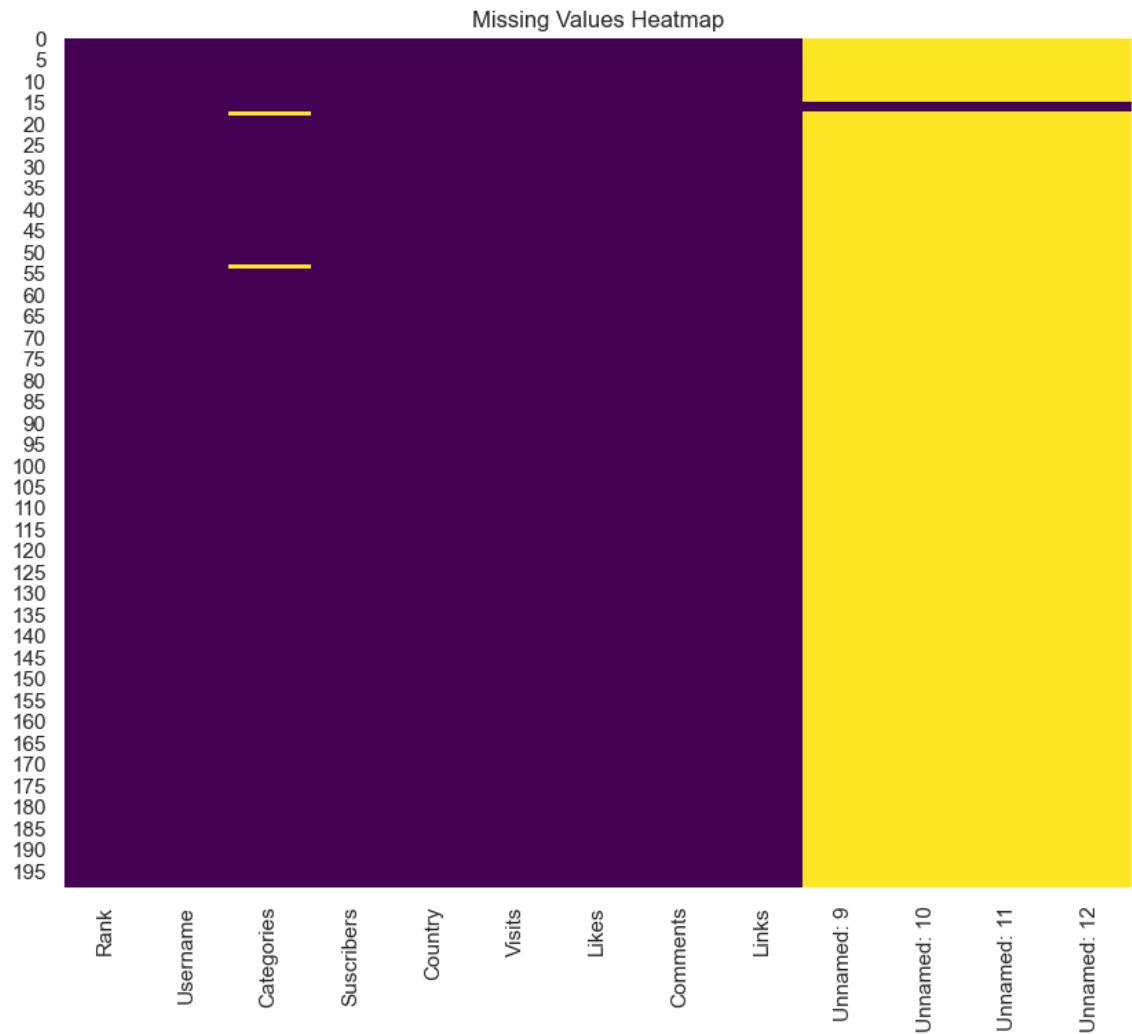
```
1 import pandas as pd
2 import matplotlib.pyplot as plt
3
4
5 # Plot the multivariate time series data
6 plt.figure(figsize=(12, 6))
7 plt.plot(df.index, df['Likes'], label='Likes', color='blue')
8 plt.plot(df.index, df['Comments'], label='Comments', color='red')
9 plt.plot(df.index, df['Visits'], label='Visits', color='green')
10
11 # Customize the plot
12 plt.xlabel('Date')
13 plt.ylabel('Value')
14 plt.title('Multivariate Time Series Analysis')
15 plt.legend()
16 plt.grid(True)
17
18 # Show the plot
19 plt.show()
```



Missing Values Heatmap

In [237]:

```
1 import pandas as pd
2 import seaborn as sns
3 import matplotlib.pyplot as plt
4 # Load your CSV file
5 file_path = r"C:\Users\Anusha V\Desktop\YouTube.csv"
6 df = pd.read_csv(file_path)
7 # Plot the heatmap
8 plt.figure(figsize=(10, 8))
9 sns.heatmap(df.isnull(), cmap='viridis', cbar=False)
10 plt.title('Missing Values Heatmap')
11 plt.show()
```



isnull method

In [239]: 1 df.isnull()

Out[239]:

| | Rank | Username | Categories | Suscribers | Country | Visits | Likes | Comments | Links | Unr |
|-----|-------|----------|------------|------------|---------|--------|-------|----------|-------|-------|
| 0 | False | False | False | False | False | False | False | False | False | False |
| 1 | False | False | False | False | False | False | False | False | False | False |
| 2 | False | False | False | False | False | False | False | False | False | False |
| 3 | False | False | False | False | False | False | False | False | False | False |
| 4 | False | False | False | False | False | False | False | False | False | False |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 194 | False | False | False | False | False | False | False | False | False | False |
| 195 | False | False | False | False | False | False | False | False | False | False |
| 196 | False | False | False | False | False | False | False | False | False | False |
| 197 | False | False | False | False | False | False | False | False | False | False |
| 198 | False | False | False | False | False | False | False | False | False | False |

199 rows × 13 columns



Number of missing values

In [240]: 1 df.isnull().sum().sum()

Out[240]: 790

In [241]: 1 df.isnull().sum()

Out[241]:

| | |
|--------------|-----|
| Rank | 0 |
| Username | 0 |
| Categories | 2 |
| Suscribers | 0 |
| Country | 0 |
| Visits | 0 |
| Likes | 0 |
| Comments | 0 |
| Links | 0 |
| Unnamed: 9 | 197 |
| Unnamed: 10 | 197 |
| Unnamed: 11 | 197 |
| Unnamed: 12 | 197 |
| dtype: int64 | |

```
In [242]: 1 import pandas as pd
2
3
4 # Identify missing values
5 missing_data = df.isnull()
6
7 # Count missing values per column
8 missing_count = missing_data.sum()
9
10 # Count missing values per row
11 missing_count_per_row = missing_data.sum(axis=1)
12
13 # Print or analyze the missing data
14 print("Missing Data Count per Column:")
15 print(missing_count)
16
17 print("Missing Data Count per Row:")
18 print(missing_count_per_row)
```

Missing Data Count per Column:

```
Rank          0
Username      0
Categories    2
Suscribers   0
Country       0
Visits        0
Likes         0
Comments      0
Links         0
Unnamed: 9    197
Unnamed: 10   197
Unnamed: 11   197
Unnamed: 12   197
dtype: int64
```

Missing Data Count per Row:

```
0      4
1      4
2      4
3      4
4      4
..
194    4
195    4
196    4
197    4
198    4
Length: 199, dtype: int64
```

```
In [244]: 1 import pandas as pd
2
3
4 # Check for missing values
5 missing_values = df.isnull().sum()
6
7 # Print the columns with missing values
8 print("Columns with missing values:")
9 print(missing_values[missing_values > 0])
```

Columns with missing values:

| Categories | 2 |
|-------------|-------|
| Unnamed: 9 | 197 |
| Unnamed: 10 | 197 |
| Unnamed: 11 | 197 |
| Unnamed: 12 | 197 |
| dtype: | int64 |

fillna method

```
In [218]: 1 d_f = df.fillna(100)
2 d_f
```

Out[218]:

| | Rank | Username | Categories | Suscribers | Country | Visits | Likes | Con |
|-----|------|---------------------|----------------------------------|------------|----------------|-----------|---------|-----|
| 0 | 1 | tseries | Møsica y baile | 249500000 | India | 86200 | 2700 | |
| 1 | 2 | MrBeast | Videojuegos, Humor | 183500000 | Estados Unidos | 117400000 | 5300000 | |
| 2 | 3 | CoComelon | Educaci&on | 165500000 | Unknown | 7000000 | 24700 | |
| 3 | 4 | SETIndia | Educaci&on | 162600000 | India | 15600 | 166 | |
| 4 | 5 | KidsDianaShow | Animaci&on, Juguetes | 113500000 | Unknown | 3900000 | 12400 | |
| ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 194 | 195 | dianaandromahin8102 | Animaci&on, Juguetes | 26200000 | Unknown | 109000 | 322 | |
| 195 | 196 | nickiminaj | Møsica y baile | 26100000 | Estados Unidos | 1600000 | 98300 | |
| 196 | 197 | mariliamendoncareal | Møsica y baile | 26100000 | Brasil | 0 | 0 | |
| 197 | 198 | TheLallantop | Noticias y Pol&tica | 26000000 | India | 30800 | 822 | |
| 198 | 199 | AGT | Møsica y baile, Pel&culas | 25900000 | Estados Unidos | 186800 | 3100 | |

199 rows × 13 columns

```
In [245]: 1 df_filled = df.fillna(1212)
2 df_filled
```

Out[245]:

| | Rank | Username | Categories | Suscribers | Country | Visits | Likes | Comments |
|-----|------|---------------------|------------|------------|----------------|-----------|---------|----------|
| 0 | 1 | tseries | 4567.0 | 249500000 | India | 86200 | 2700 | |
| 1 | 2 | MrBeast | 4567.0 | 183500000 | Estados Unidos | 117400000 | 5300000 | |
| 2 | 3 | CoComelon | 4567.0 | 165500000 | Unknown | 7000000 | 24700 | |
| 3 | 4 | SETIndia | 4567.0 | 162600000 | India | 15600 | 166 | |
| 4 | 5 | KidsDianaShow | 4567.0 | 113500000 | Unknown | 3900000 | 12400 | |
| ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 194 | 195 | dianaandromahin8102 | 4567.0 | 26200000 | Unknown | 109000 | 322 | |
| 195 | 196 | nickiminaj | 4567.0 | 26100000 | Estados Unidos | 1600000 | 98300 | |
| 196 | 197 | mariliamendoncareal | 4567.0 | 26100000 | Brasil | 0 | 0 | |
| 197 | 198 | TheLallantop | 4567.0 | 26000000 | India | 30800 | 822 | |
| 198 | 199 | AGT | 4567.0 | 25900000 | Estados Unidos | 186800 | 3100 | |

199 rows × 13 columns



dropna method

```
In [247]: 1 d_f = df.dropna(0)
2 d_f
```

C:\Users\Anusha V\AppData\Local\Temp\ipykernel_12064\1675708418.py:1: FutureWarning: In a future version of pandas all arguments of DataFrame.dropna will be keyword-only.

```
d_f = df.dropna(0)
```

Out[247]:

| | Rank | Username | Categories | Suscribers | Country | Visits | Likes |
|----|------|------------|------------|------------|------------|------------|------------|
| 15 | 16 | HYBELABELS | 4567.0 | HYBELABELS | HYBELABELS | HYBELABELS | HYBELABELS |
| 16 | 17 | HYBELABELS | 4567.0 | HYBELABELS | HYBELABELS | HYBELABELS | HYBELABELS |



Replace

In [248]:

```
1 import numpy as np
2 df12 = df.replace(to_replace=np.nan, value='not found')
3 df12
```

Out[248]:

| | Rank | Username | Categories | Suscribers | Country | Visits | Likes | Com |
|-----|------|---------------------|------------|------------|----------------|-----------|---------|-----|
| 0 | 1 | tseries | 4567.0 | 249500000 | India | 86200 | 2700 | |
| 1 | 2 | MrBeast | 4567.0 | 183500000 | Estados Unidos | 117400000 | 5300000 | |
| 2 | 3 | CoComelon | 4567.0 | 165500000 | Unknown | 7000000 | 24700 | |
| 3 | 4 | SETIndia | 4567.0 | 162600000 | India | 15600 | 166 | |
| 4 | 5 | KidsDianaShow | 4567.0 | 113500000 | Unknown | 3900000 | 12400 | |
| ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 194 | 195 | dianaandromahin8102 | 4567.0 | 26200000 | Unknown | 109000 | 322 | |
| 195 | 196 | nickiminaj | 4567.0 | 26100000 | Estados Unidos | 1600000 | 98300 | |
| 196 | 197 | mariliamendoncareal | 4567.0 | 26100000 | Brasil | 0 | 0 | |
| 197 | 198 | TheLallantop | 4567.0 | 26000000 | India | 30800 | 822 | |
| 198 | 199 | AGT | 4567.0 | 25900000 | Estados Unidos | 186800 | 3100 | |

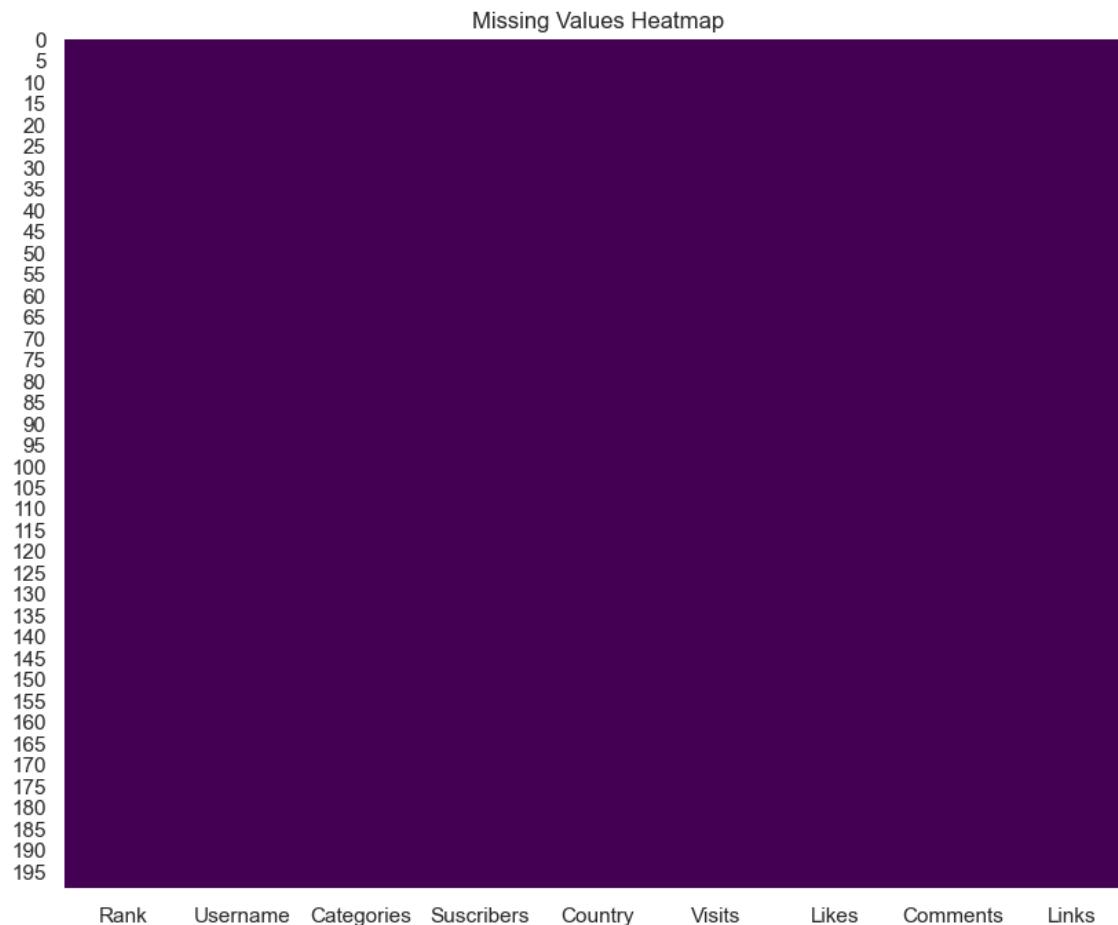
199 rows × 13 columns



HEATMAP

In [251]:

```
1 import pandas as pd
2 import seaborn as sns
3 import matplotlib.pyplot as plt
4 file_path = r"C:\Users\Anusha V\Desktop\YouTube.csv"
5 df1 = pd.read_csv(file_path)
6 # Fill missing values
7 df1['Likes'].fillna(df1['Likes'].mean(), inplace=True)
8 df1['Comments'].fillna(df1['Comments'].mean(), inplace=True)
9 df1['Visits'].fillna(df1['Visits'].mean(), inplace=True)
10
11
12 # Heatmap
13 plt.figure(figsize=(10, 8))
14 sns.heatmap(df1.isnull(), cmap='viridis', cbar=False)
15 plt.title('Missing Values Heatmap')
16 plt.show()
```



In [291]:

```
1 import pandas as pd
2 import seaborn as sns
3 import matplotlib.pyplot as plt
4 # Load your CSV file
5 file_path = r"C:\Users\Anusha V\Desktop\YouTube.csv"
6 df = pd.read_csv(file_path)
7 # Plot the heatmap
8 plt.figure(figsize=(10, 8))
9 sns.heatmap(df.isnull(), cmap='viridis', cbar=False)
10 plt.title('Missing Values Heatmap')
11 plt.show()
```

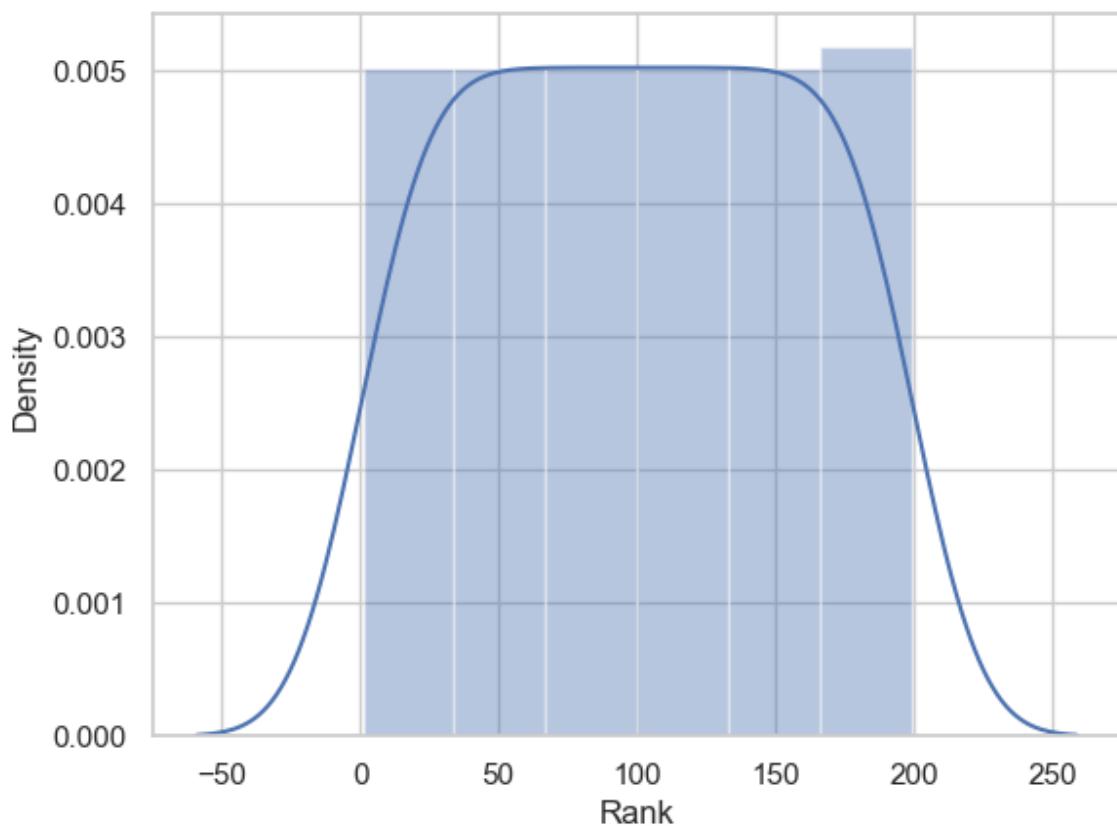


Distplot

```
In [276]: 1 import seaborn as sns  
2 sns.distplot(df['Rank'])
```

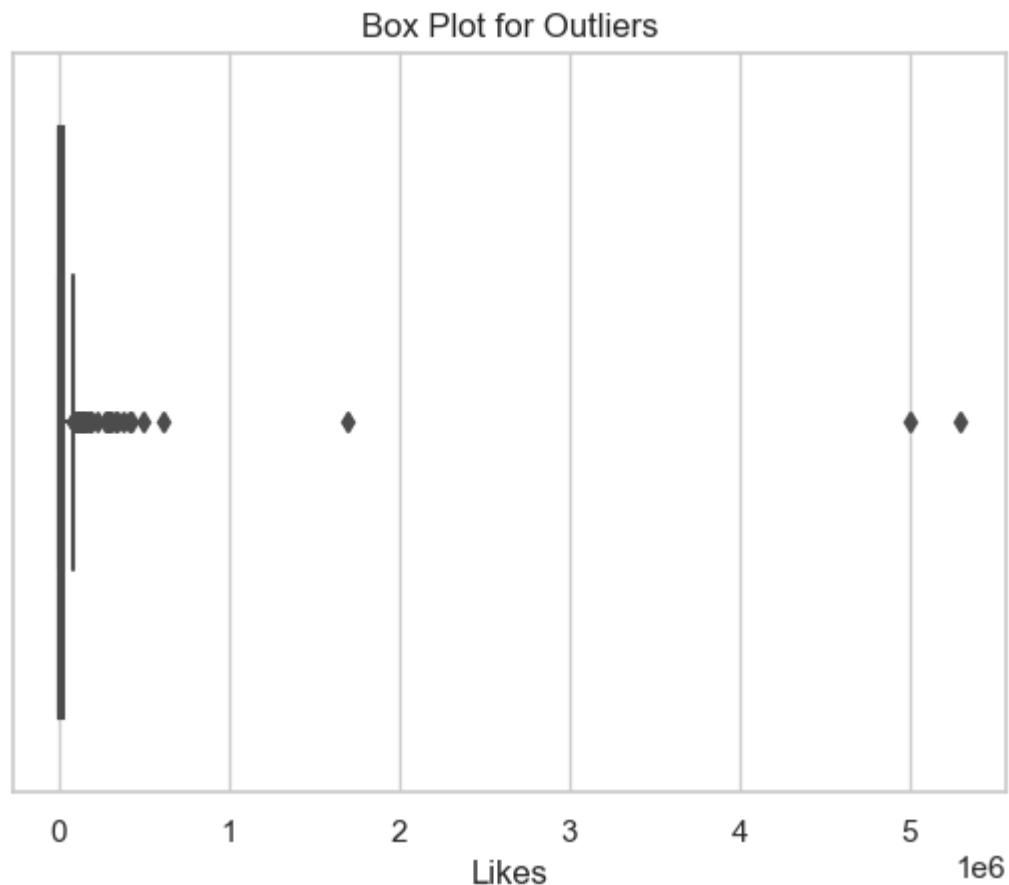
C:\Users\Anusha V\anaconda3\lib\site-packages\seaborn\distributions.py:261
9: FutureWarning: `distplot` is a deprecated function and will be removed
in a future version. Please adapt your code to use either `displot` (a fig
ure-level function with similar flexibility) or `histplot` (an axes-level
function for histograms).
warnings.warn(msg, FutureWarning)

Out[276]: <AxesSubplot:xlabel='Rank', ylabel='Density'>



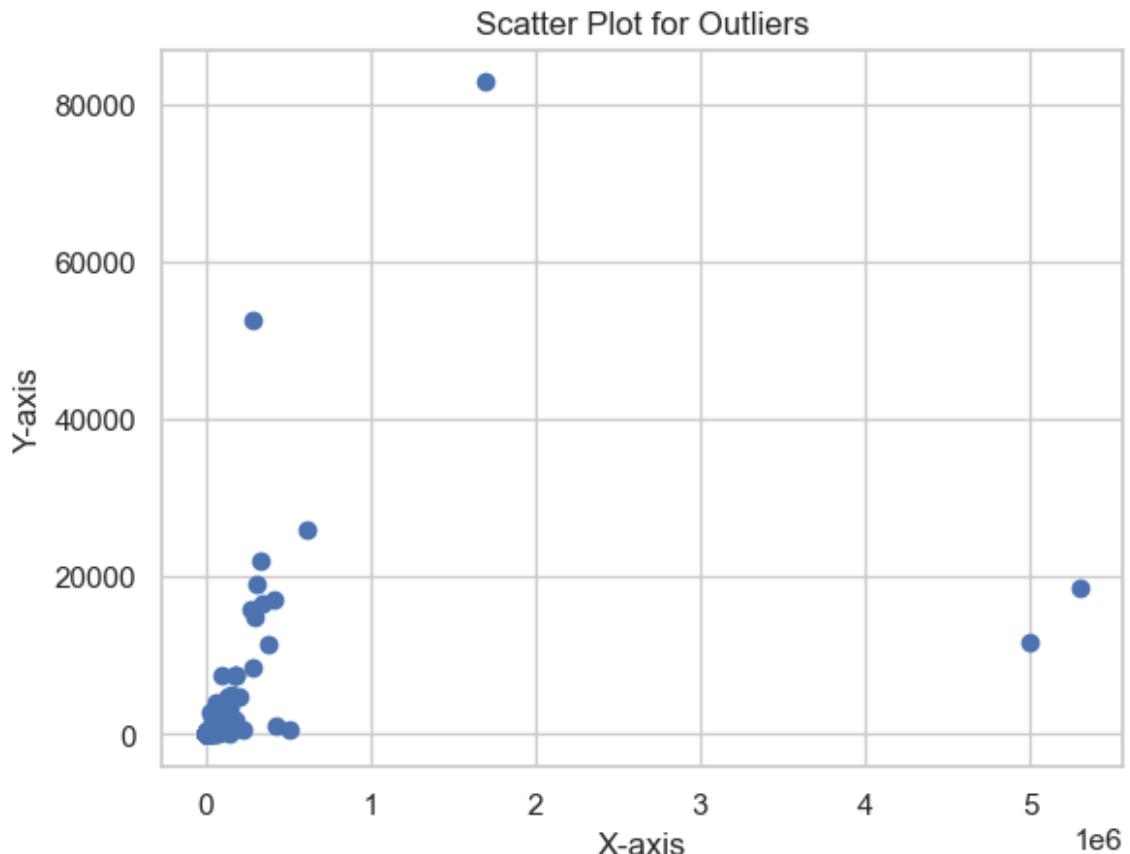
Box Plot for Outliers

```
In [256]:  
1 import matplotlib.pyplot as plt  
2 import seaborn as sns  
3  
4 # Assuming 'df' is your DataFrame  
5 sns.boxplot(x=df['Likes'])  
6 plt.title("Box Plot for Outliers")  
7 plt.show()
```



Scatter Plot for Outliers

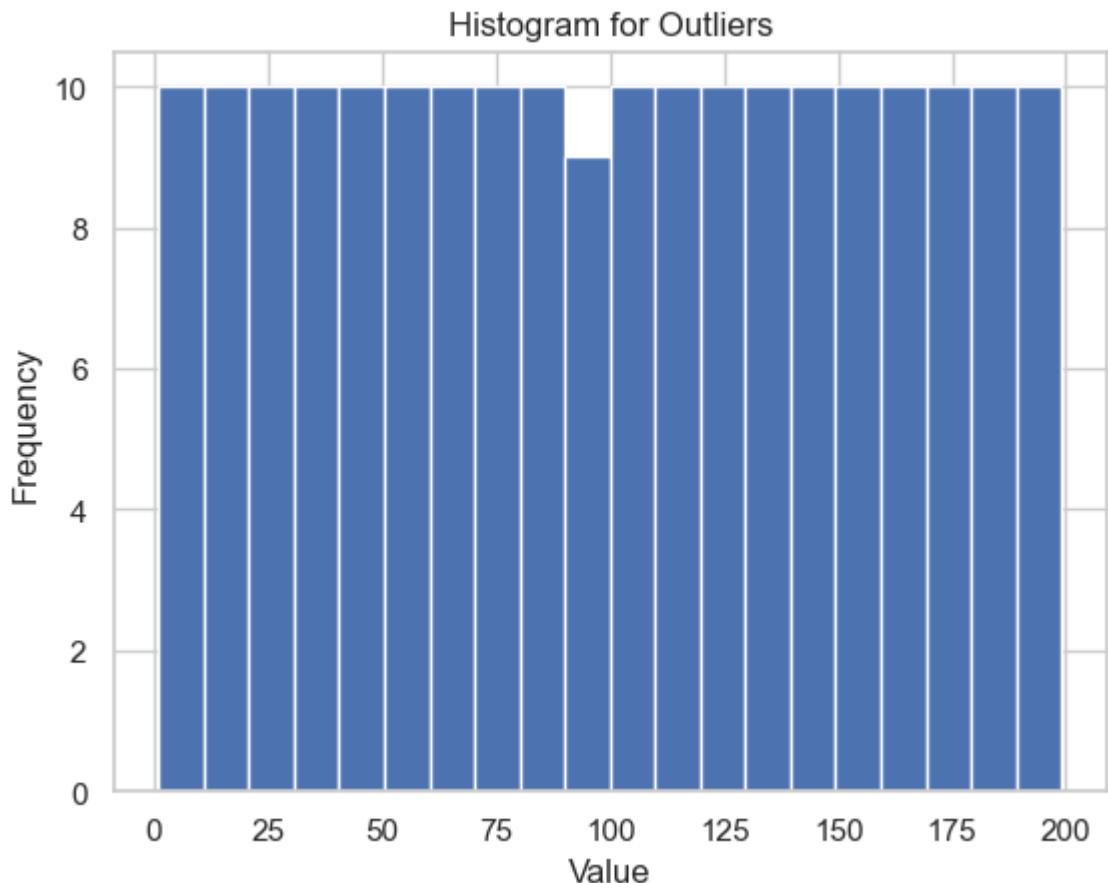
```
In [257]:  
1 import matplotlib.pyplot as plt  
2  
3 # Assuming 'df' is your DataFrame and 'x' and 'y' are columns  
4 plt.scatter(df['Likes'], df['Comments'])  
5 plt.title("Scatter Plot for Outliers")  
6 plt.xlabel("X-axis")  
7 plt.ylabel("Y-axis")  
8 plt.show()  
9
```



Histogram for Outliers

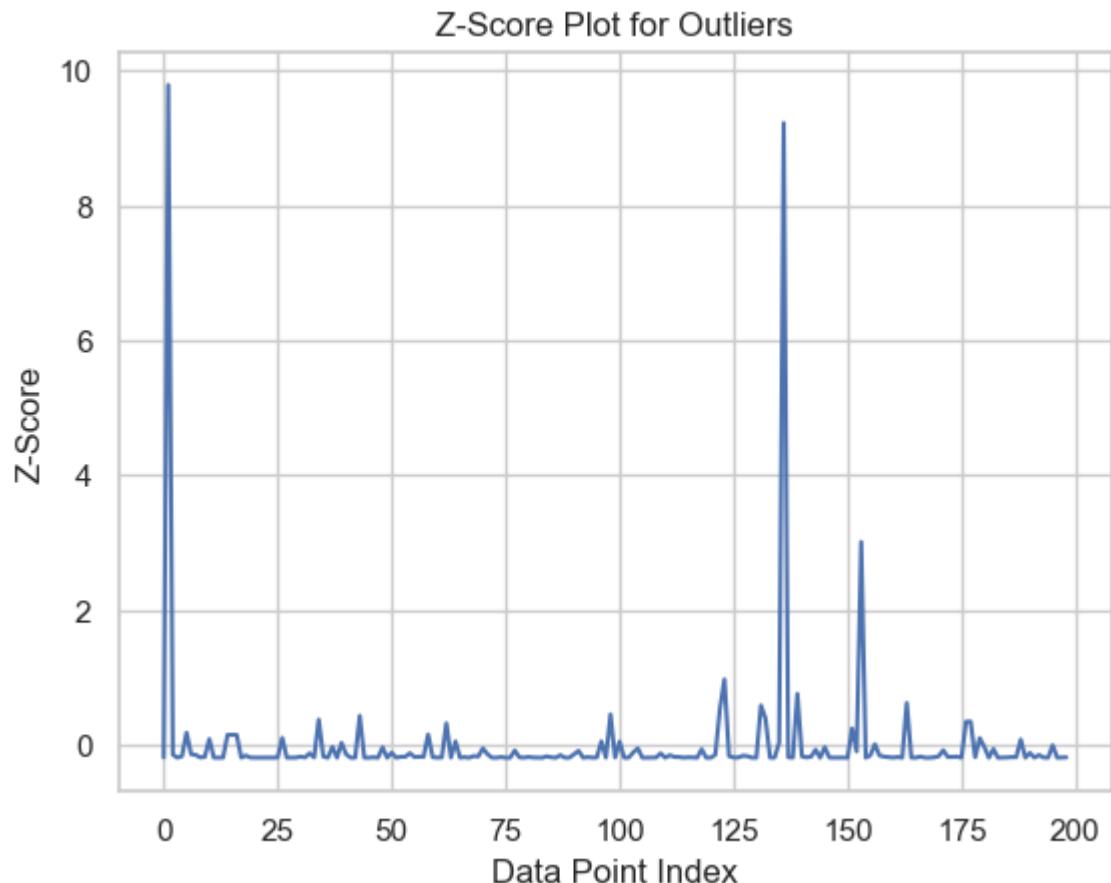
In [260]:

```
1 import matplotlib.pyplot as plt
2
3 # Assuming 'df' is your DataFrame and 'column_name' is the variable of
4 plt.hist(df['Rank'], bins=20) # Adjust the number of bins as needed
5 plt.title("Histogram for Outliers")
6 plt.xlabel("Value")
7 plt.ylabel("Frequency")
8 plt.show()
9
```



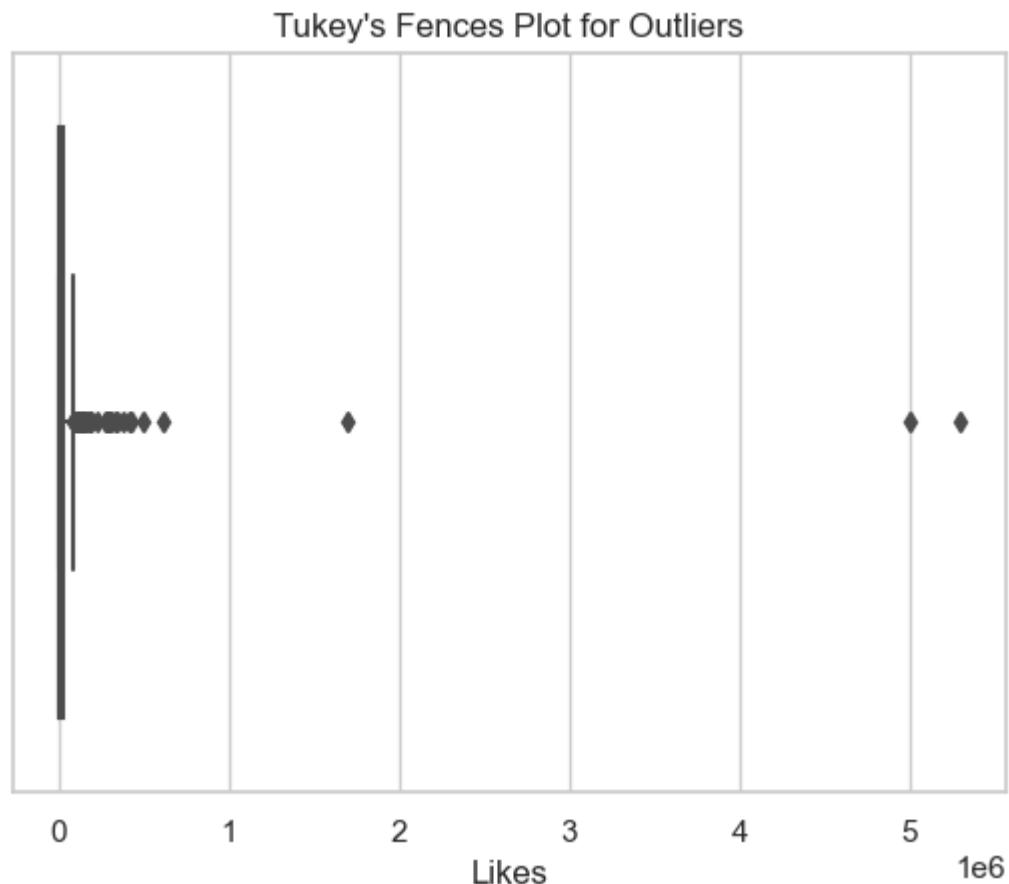
Z-Score Plot for Outliers

```
In [261]: 1 from scipy import stats  
2  
3 # Assuming 'df' is your DataFrame and 'column_name' is the variable of  
4 z_scores = stats.zscore(df['Likes'])  
5 plt.plot(z_scores)  
6 plt.title("Z-Score Plot for Outliers")  
7 plt.xlabel("Data Point Index")  
8 plt.ylabel("Z-Score")  
9 plt.show()  
10
```



Tukey's Fences Plot for Outliers

```
In [263]:  
1 import matplotlib.pyplot as plt  
2 import seaborn as sns  
3  
4 # Assuming 'df' is your DataFrame and 'column_name' is the variable of  
5 sns.boxplot(x=df['Likes'], whis=1.5) # Adjust the whisker length as ne  
6 plt.title("Tukey's Fences Plot for Outliers")  
7 plt.show()  
8
```



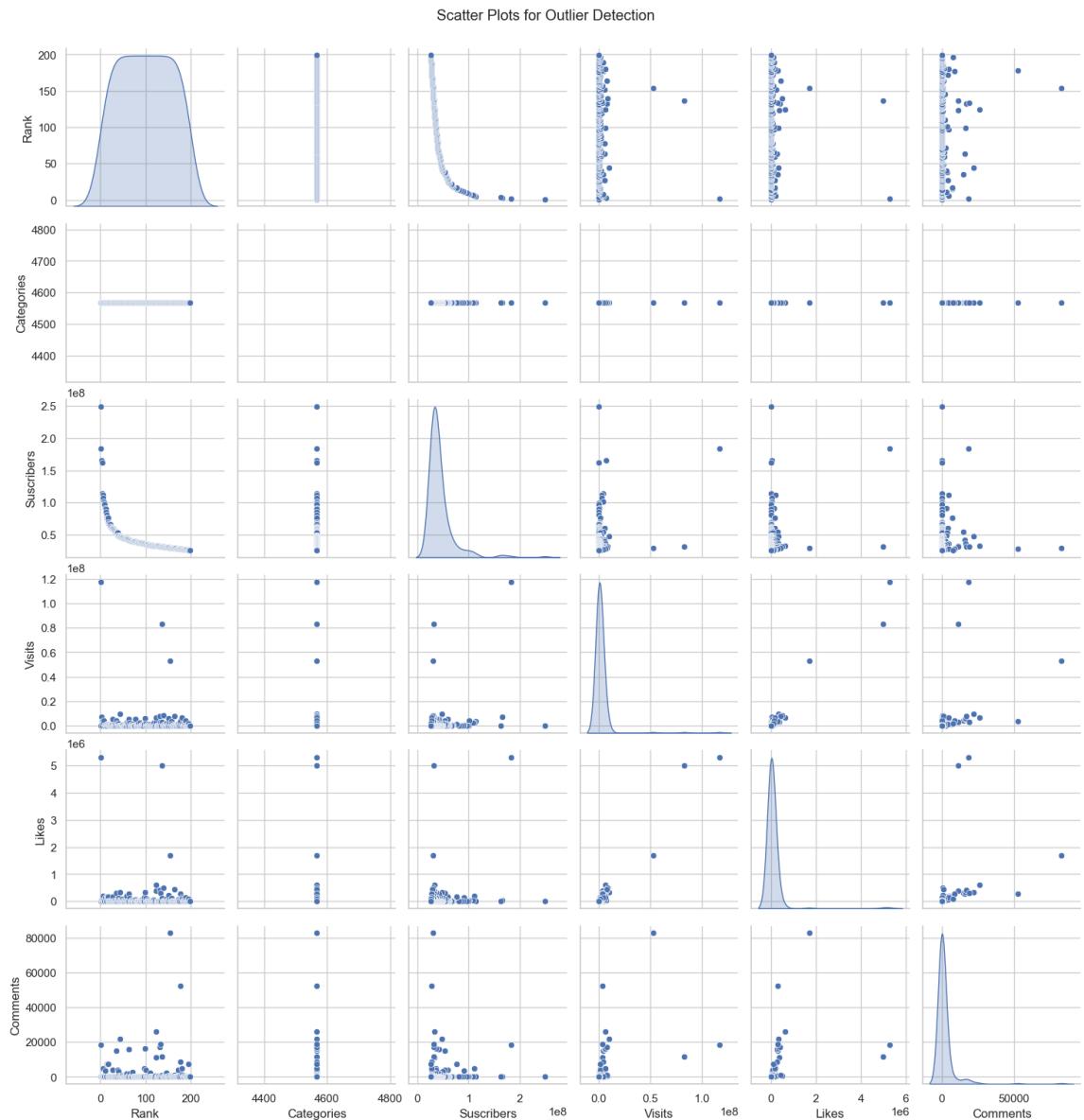
Scatter Plots for Outlier

In [266]:

```

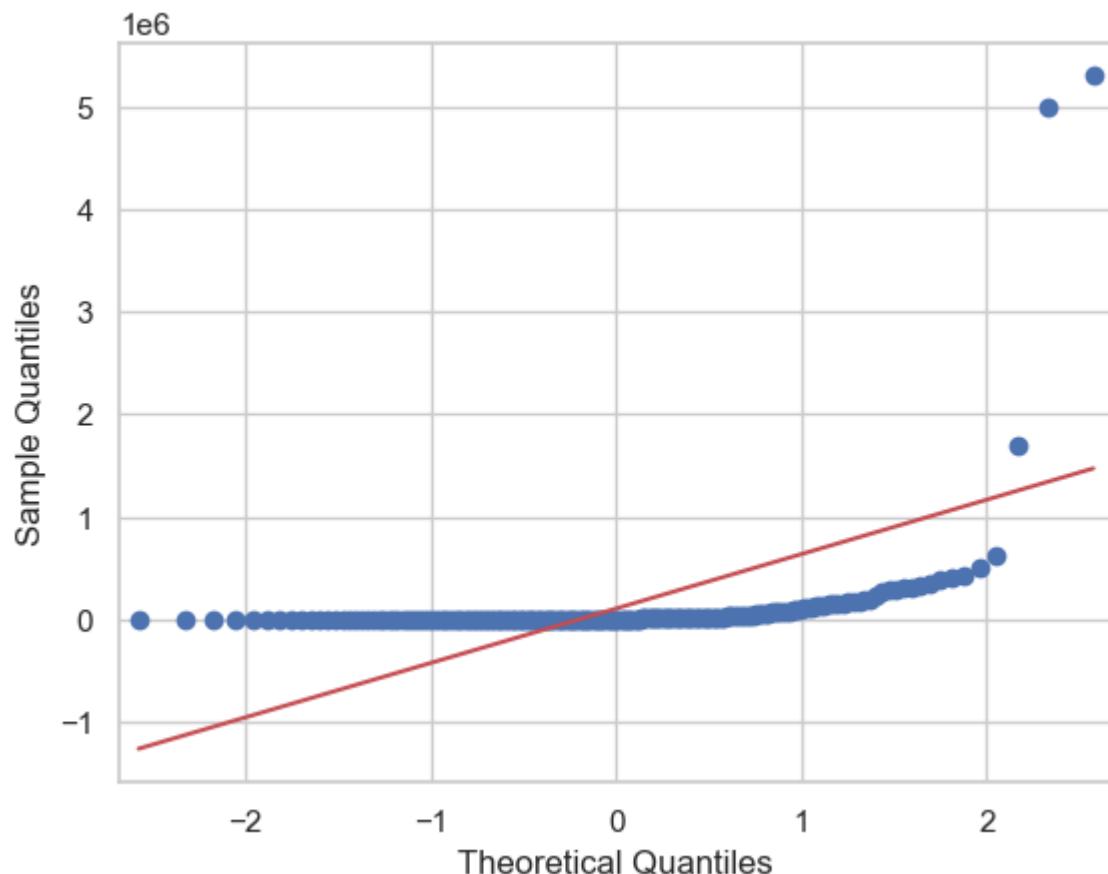
1 # Create scatter plots for pairs of numeric columns
2 numeric_columns = df.select_dtypes(include=['float64', 'int64'])
3 sns.pairplot(data=df, vars=numeric_columns.columns, diag_kind='kde')
4 plt.suptitle('Scatter Plots for Outlier Detection', y=1.02)
5 plt.show()

```

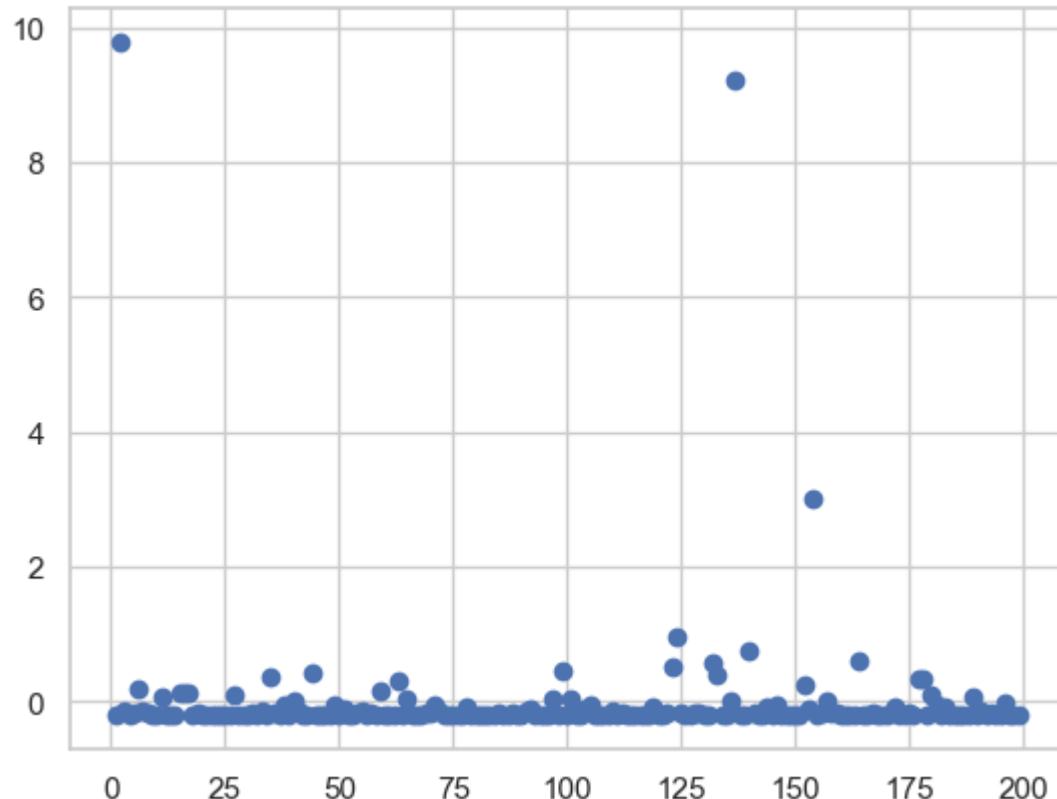


Outlier

```
In [267]: 1 from statsmodels.graphics.gofplots import qqplot  
2  
3 # Create a QQ plot  
4 qqplot(df['Likes'], line='s')  
5 plt.show()
```



```
In [270]: 1 from scipy import stats  
2  
3 z_scores = stats.zscore(df['Likes'])  
4 plt.scatter(df['Rank'], z_scores)  
5 plt.show()
```



Upper limit & Lower limit

```
In [279]: 1 upper_limit = df['Likes'].mean() + 3*df['Likes'].std()  
2 lower_limit = df['Rank'].mean() + 3*df['Rank'].std()  
3 print('upper limit:', upper_limit)  
4 print('lower limit:', lower_limit)
```

upper limit: 1701275.3320238253
lower limit: 272.77152543170996

In [281]: 1 df.loc[(df['Likes'] < upper_limit) | (df['Likes'] > lower_limit)]

Out[281]:

| | Rank | Username | Categories | Suscribers | Country | Visits | Likes | Comments |
|-----|------|---------------------|------------|------------|----------------|-----------|---------|----------|
| 0 | 1 | tseries | 4567 | 249500000 | India | 86200 | 2700 | |
| 1 | 2 | MrBeast | 4567 | 183500000 | Estados Unidos | 117400000 | 5300000 | |
| 2 | 3 | CoComelon | 4567 | 165500000 | Unknown | 7000000 | 24700 | |
| 3 | 4 | SETIndia | 4567 | 162600000 | India | 15600 | 166 | |
| 4 | 5 | KidsDianaShow | 4567 | 113500000 | Unknown | 3900000 | 12400 | |
| ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 194 | 195 | dianaandromahin8102 | 4567 | 26200000 | Unknown | 109000 | 322 | |
| 195 | 196 | nickiminaj | 4567 | 26100000 | Estados Unidos | 1600000 | 98300 | |
| 196 | 197 | mariliamendoncareal | 4567 | 26100000 | Brasil | 0 | 0 | |
| 197 | 198 | TheLallantop | 4567 | 26000000 | India | 30800 | 822 | |
| 198 | 199 | AGT | 4567 | 25900000 | Estados Unidos | 186800 | 3100 | |

199 rows × 9 columns



In [293]: 1 from pandas.api.types import CategoricalDtype
2 df['Categories'] = df['Categories'].astype(CategoricalDtype(categories

In [294]: 1 df['Categories'].unique()

Out[294]: array([True, False])

In [296]: 1 import pandas as pd
2
3 # Specify the number of bins or the bin edges
4 # For simplicity, let's assume you want to create three bins for a column
5 num_bins = 3
6
7 # Perform data discretization
8 d_df=df['Likes'] = pd.cut(df['Comments'], bins=num_bins, labels=False)
9 d_df

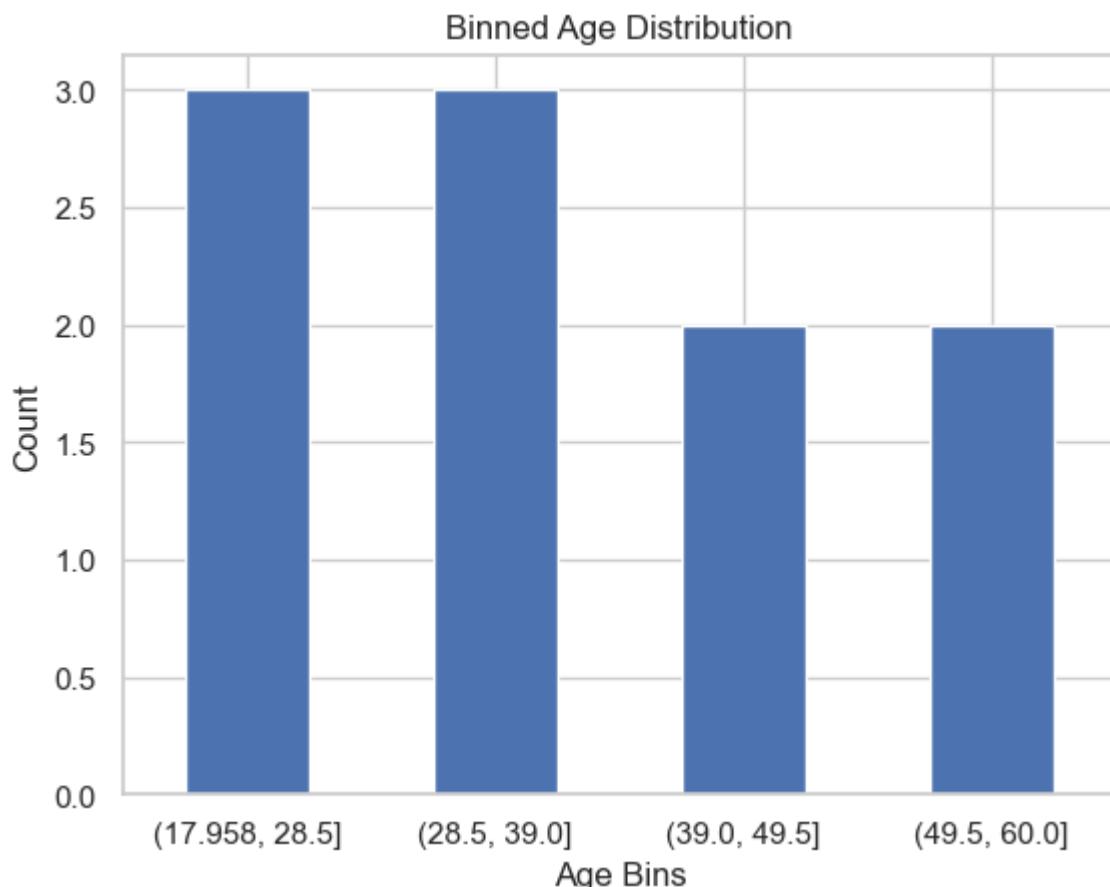
Out[296]: 0 0
1 0
2 0
3 0
4 0
..
194 0
195 0
196 0
197 0
198 0
Name: Comments, Length: 199, dtype: int64

Loading [MathJax]/extensions/MathML/content-mathml.js

equal-width binning

In [297]:

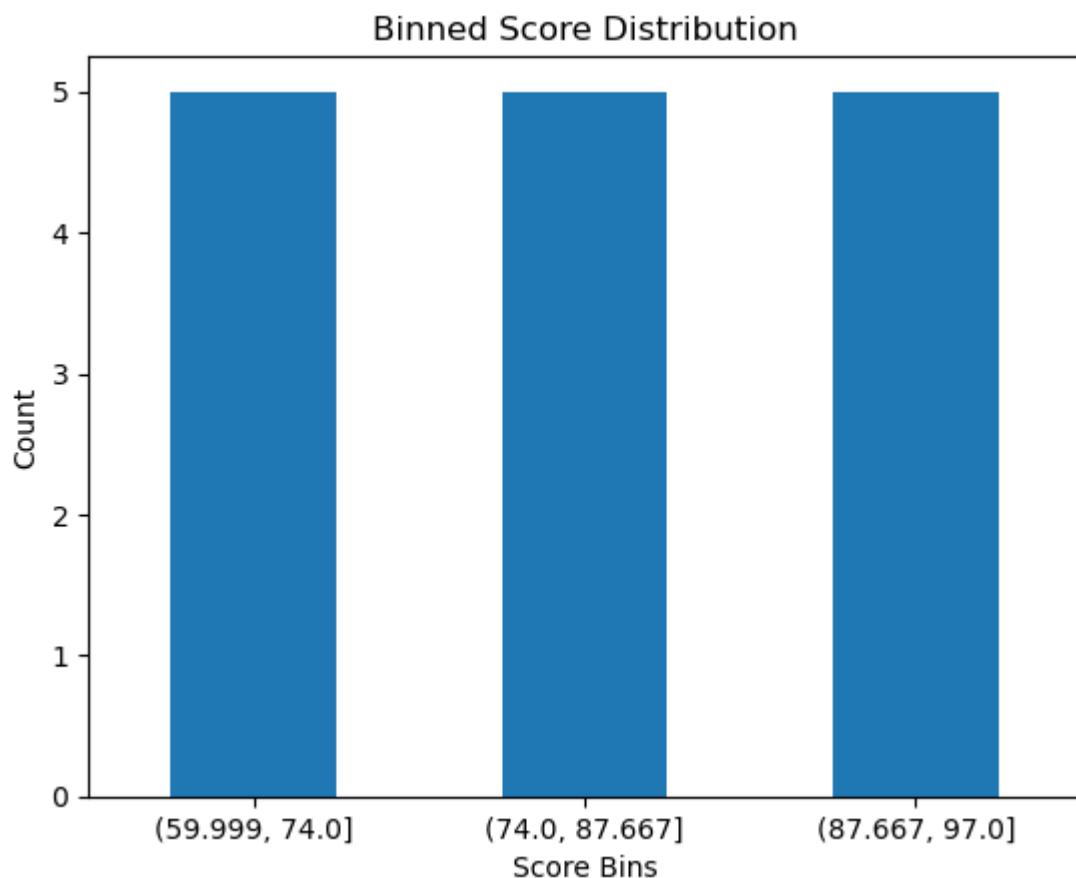
```
1 import pandas as pd
2 import matplotlib.pyplot as plt
3
4 # Create a DataFrame with continuous data
5 data = {'Age': [25, 32, 45, 52, 60, 18, 35, 42, 28, 38]}
6 df = pd.DataFrame(data)
7
8 # Perform equal-width binning
9 num_bins = 4
10 df['Age_Binned'] = pd.cut(df['Age'], bins=num_bins)
11
12 # Count the number of data points in each bin
13 bin_counts = df['Age_Binned'].value_counts().sort_index()
14
15 # Create a bar graph
16 bin_counts.plot(kind='bar')
17 plt.xlabel('Age Bins')
18 plt.ylabel('Count')
19 plt.title('Binned Age Distribution')
20 plt.xticks(rotation=0) # Ensure x-axis labels are not rotated
21 plt.show()
```



equal-frequency binning

Loading [MathJax]/extensions/MathML/content-mathml.js

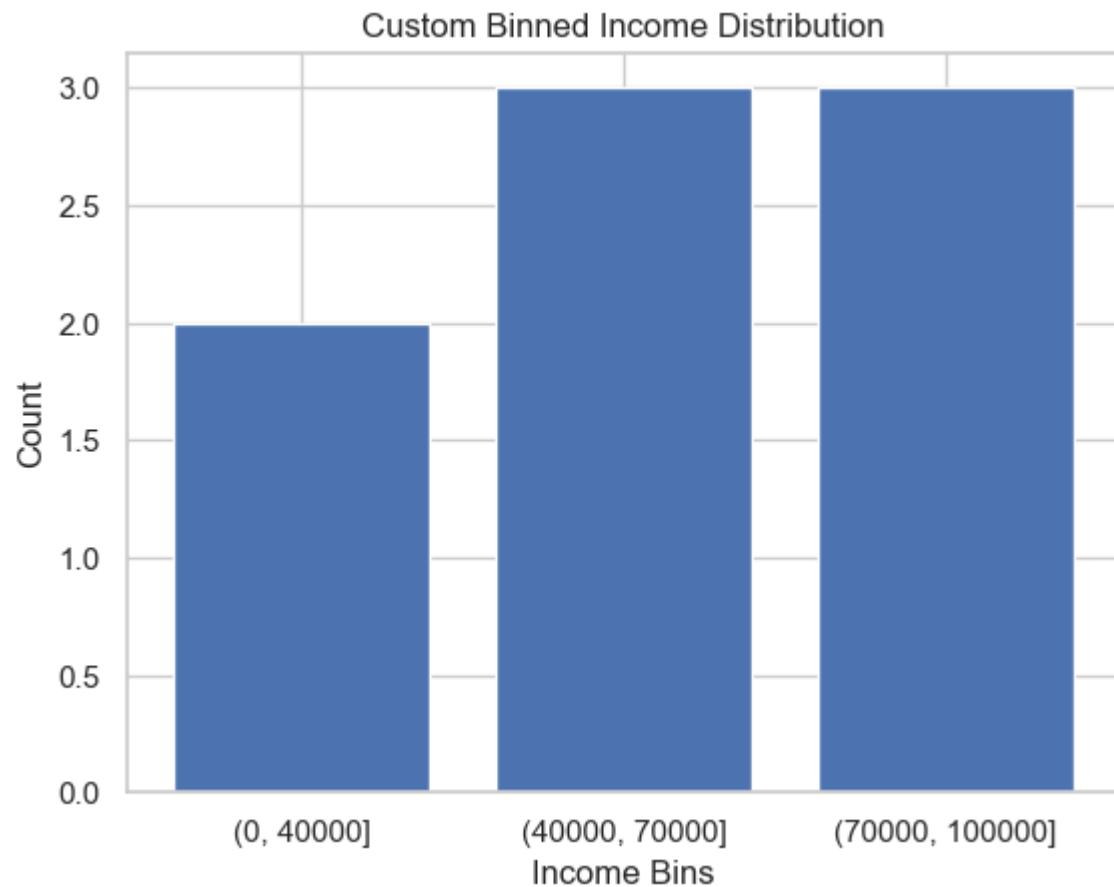
```
In [4]:  
1 # import pandas as pd  
2 import matplotlib.pyplot as plt  
3  
4 # Create a DataFrame with continuous data  
5 data = {'Score': [82, 91, 75, 68, 95, 87, 63, 78, 89, 72, 84, 97, 60, 9  
6 df = pd.DataFrame(data)  
7  
8 # Define the number of bins  
9 num_bins = 3  
10  
11 # Perform equal-frequency binning  
12 df['Score_Binned'] = pd.qcut(df['Score'], q=num_bins)  
13  
14 # Count the number of data points in each bin  
15 bin_counts = df['Score_Binned'].value_counts().sort_index()  
16  
17 # Create a bar graph  
18 bin_counts.plot(kind='bar')  
19 plt.xlabel('Score Bins')  
20 plt.ylabel('Count')  
21 plt.title('Binned Score Distribution')  
22 plt.xticks(rotation=0) # Ensure x-axis labels are not rotated  
23 plt.show()
```



custom binning

In [299]:

```
1 import pandas as pd
2 import matplotlib.pyplot as plt
3
4 # Create a DataFrame with continuous data
5 data = {'Income': [25000, 35000, 45000, 55000, 65000, 75000, 85000, 95000]}
6 df = pd.DataFrame(data)
7
8 # Define custom bin edges
9 custom_bins = [0, 40000, 70000, 100000]
10
11 # Perform custom binning
12 df['Income_Binned'] = pd.cut(df['Income'], bins=custom_bins)
13
14 # Count the number of data points in each bin
15 bin_counts = df['Income_Binned'].value_counts().sort_index()
16
17 # Create a histogram
18 plt.bar(bin_counts.index.astype(str), bin_counts.values)
19 plt.xlabel('Income Bins')
20 plt.ylabel('Count')
21 plt.title('Custom Binned Income Distribution')
22 plt.show()
```



cluster-based binning

In [353]: 1 pip install scikit-learn

```
Requirement already satisfied: scikit-learn in c:\users\anusha v\anaconda3
\lib\site-packages (1.0.2)
Requirement already satisfied: joblib>=0.11 in c:\users\anusha v\anaconda3
\lib\site-packages (from scikit-learn) (1.1.0)
Requirement already satisfied: threadpoolctl>=2.0.0 in c:\users\anusha v\ana
naconda3\lib\site-packages (from scikit-learn) (2.2.0)
Requirement already satisfied: numpy>=1.14.6 in c:\users\anusha v\anaconda
3\lib\site-packages (from scikit-learn) (1.21.5)
Requirement already satisfied: scipy>=1.1.0 in c:\users\anusha v\anaconda3
\lib\site-packages (from scikit-learn) (1.9.1)
Note: you may need to restart the kernel to use updated packages.
```

In [354]: 1 pip install KMeans

Note: you may need to restart the kernel to use updated packages.

```
WARNING: Retrying (Retry(total=4, connect=None, read=None, redirect=None,
status=None)) after connection broken by 'NewConnectionError('<pip._vend
r.urllib3.connection.HTTPSConnection object at 0x0000019A50BAB310>: Failed
to establish a new connection: [Errno 11001] getaddrinfo failed')': /simp
le/kmeans/
WARNING: Retrying (Retry(total=3, connect=None, read=None, redirect=None,
status=None)) after connection broken by 'NewConnectionError('<pip._vend
r.urllib3.connection.HTTPSConnection object at 0x0000019A50BAB640>: Failed
to establish a new connection: [Errno 11001] getaddrinfo failed')': /simp
le/kmeans/
WARNING: Retrying (Retry(total=2, connect=None, read=None, redirect=None,
status=None)) after connection broken by 'NewConnectionError('<pip._vend
r.urllib3.connection.HTTPSConnection object at 0x0000019A50BAB940>: Failed
to establish a new connection: [Errno 11001] getaddrinfo failed')': /simp
le/kmeans/
WARNING: Retrying (Retry(total=1, connect=None, read=None, redirect=None,
status=None)) after connection broken by 'NewConnectionError('<pip._vend
r.urllib3.connection.HTTPSConnection object at 0x0000019A50BABA0>: Failed
to establish a new connection: [Errno 11001] getaddrinfo failed')': /simp
le/kmeans/
WARNING: Retrying (Retry(total=0, connect=None, read=None, redirect=None,
status=None)) after connection broken by 'NewConnectionError('<pip._vend
r.urllib3.connection.HTTPSConnection object at 0x0000019A50BABCA0>: Failed
to establish a new connection: [Errno 11001] getaddrinfo failed')': /simp
le/kmeans/
ERROR: Could not find a version that satisfies the requirement KMeans (fro
m versions: none)
ERROR: No matching distribution found for KMeans
```

```
In [355]: 1 import pandas as pd
2 import matplotlib.pyplot as plt
3
4 # Create a DataFrame with continuous data
5 data = {'Value': [12, 15, 18, 25, 35, 42, 55, 65]}
6 df = pd.DataFrame(data)
7
8 # Define the number of desired bins
9 num_bins = 3
10
11 # Create a K-means clustering model
12 kmeans = KMeans(n_clusters=num_bins, random_state=42)
13
14 # Fit the model to the data
15 df['Cluster'] = kmeans.fit_predict(df[['Value']])
16
17 # Calculate cluster centers
18 cluster_centers = kmeans.cluster_centers_
19
20 # Sort cluster centers to determine bin edges
21 cluster_centers.sort(axis=0)
22
23 # Create bins based on cluster centers
24 bin_edges = cluster_centers.flatten().tolist()
25 bin_edges.append(float('inf'))
26
27 # Perform cluster-based binning
28 df['Value_Binned'] = pd.cut(df['Value'], bins=bin_edges)
29
30 # Count the number of data points in each bin
31 bin_counts = df['Value_Binned'].value_counts().sort_index()
32
33 # Create a bar graph
34 bin_counts.plot(kind='bar')
35 plt.xlabel('Value Bins')
36 plt.ylabel('Count')
37 plt.title('Cluster-Based Binned Value Distribution')
38 plt.xticks(rotation=0) # Ensure x-axis labels are not rotated
39 plt.show()
```

```
-  
NameError: name 'KMeans' is not defined
```

Traceback (most recent call last)
t)
~\AppData\Local\Temp\ipykernel_12064\3132253234.py in <module>
 10
 11 # Create a K-means clustering model
---> 12 kmeans = KMeans(n_clusters=num_bins, random_state=42)
 13
 14 # Fit the model to the data

LabelEncoder

```
In [305]: 1 from sklearn.preprocessing import LabelEncoder  
2  
3 # Sample data  
4 data = ['low', 'medium', 'high', 'low', 'high', 'medium']  
5  
6 # Create a label encoder  
7 label_encoder = LabelEncoder()  
8  
9 # Fit the Label encoder on the data and transform the data  
10 encoded_data = label_encoder.fit_transform(data)  
11  
12 # The encoded data can now be used in your analysis or machine Learning  
13 print(encoded_data)  
14
```

```
[1 2 0 1 0 2]
```

Handling outliers

```
In [364]: 1 import pandas as pd  
2 from scipy import stats  
3 column_name = 'Comments'  
4 df['Z-score'] = stats.zscore(df["Comments"])  
5 output_csv_file = 'z_scores_output.csv'  
6 df.to_csv(output_csv_file, index=False)  
7 print(f'Z-scores calculated and saved to {output_csv_file}')
```

```
Z-scores calculated and saved to z_scores_output.csv
```

In [363]:

```
1 import pandas as pd
2 Q1 = df['Likes'].quantile(0.2)
3 Q3 = df['Comments'].quantile(0.20)
4 IQR = Q3 - Q1
5 lower_bound = Q1 - 1.5 * IQR
6 upper_bound = Q3 + 1.5 * IQR
7 outliers = df[(df['Likes'] < lower_bound) | (df['Comments'] > upper_bound)]
8 outliers
9
```

Out[363]:

| | Rank | Username | Categories | Suscribers | Country | Visits | Likes | Comments |
|-----|------|---------------------|------------|------------|----------------|-----------|---------|----------|
| 0 | 1 | tseries | True | 249500000 | India | 86200 | 2700 | |
| 1 | 2 | MrBeast | False | 183500000 | Estados Unidos | 117400000 | 5300000 | |
| 2 | 3 | CoComelon | True | 165500000 | Unknown | 7000000 | 24700 | |
| 3 | 4 | SETIndia | False | 162600000 | India | 15600 | 166 | |
| 4 | 5 | KidsDianaShow | True | 113500000 | Unknown | 3900000 | 12400 | |
| ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 194 | 195 | dianaandromahin8102 | True | 26200000 | Unknown | 109000 | 322 | |
| 195 | 196 | nickiminaj | False | 26100000 | Estados Unidos | 1600000 | 98300 | |
| 196 | 197 | mariliamendoncareal | True | 26100000 | Brasil | 0 | 0 | |
| 197 | 198 | TheLallantop | False | 26000000 | India | 30800 | 822 | |
| 198 | 199 | AGT | True | 25900000 | Estados Unidos | 186800 | 3100 | |

199 rows × 10 columns



In [321]:

```
1 import pandas as pd
2 df_no_duplicates = df.drop_duplicates()
3 df_no_duplicates.to_csv('cleaned_data.csv', index=False)
4 df_no_duplicates
```

Out[321]:

| | Rank | Username | Categories | Suscribers | Country | Visits | Likes | Com |
|-----|------|---------------------|------------|------------|----------------|-----------|---------|-----|
| 0 | 1 | tseries | True | 249500000 | India | 86200 | 2700 | |
| 1 | 2 | MrBeast | False | 183500000 | Estados Unidos | 117400000 | 5300000 | |
| 2 | 3 | CoComelon | True | 165500000 | Unknown | 7000000 | 24700 | |
| 3 | 4 | SETIndia | False | 162600000 | India | 15600 | 166 | |
| 4 | 5 | KidsDianaShow | True | 113500000 | Unknown | 3900000 | 12400 | |
| ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 194 | 195 | dianaandromahin8102 | True | 26200000 | Unknown | 109000 | 322 | |
| 195 | 196 | nickiminaj | False | 26100000 | Estados Unidos | 1600000 | 98300 | |
| 196 | 197 | mariliamendoncareal | True | 26100000 | Brasil | 0 | 0 | |
| 197 | 198 | TheLallantop | False | 26000000 | India | 30800 | 822 | |
| 198 | 199 | AGT | True | 25900000 | Estados Unidos | 186800 | 3100 | |

199 rows × 10 columns



Univariate analysis

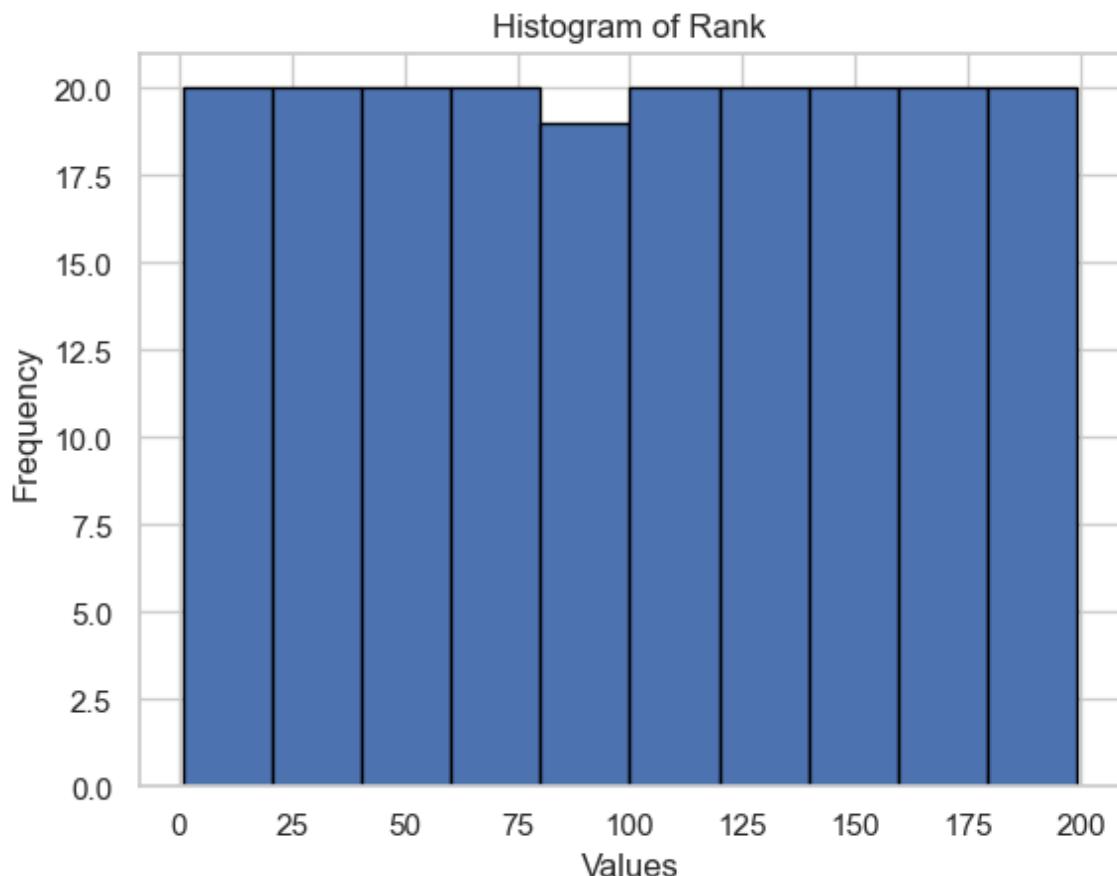
In [326]:

```
1 import pandas as pd
2 import matplotlib.pyplot as plt
3
4 selected_column = 'Rank' # Replace 'column_name' with the name of the column
5 mean = df["Likes"].mean()
6 median = df["Comments"].median()
7 mode = df["Visits"].mode().iloc[0]
8 print(f"Mean: {mean}")
9 print(f"Median: {median}")
10 print(f"Mode: {mode}")
11 plt.hist(df["Rank"], bins=10, edgecolor='black')
12 plt.xlabel("Values")
13 plt.ylabel('Frequency')
14 plt.title(f'Histogram of {selected_column}')
15 plt.show()
16
```

Mean: 104774.6783919598

Median: 28.0

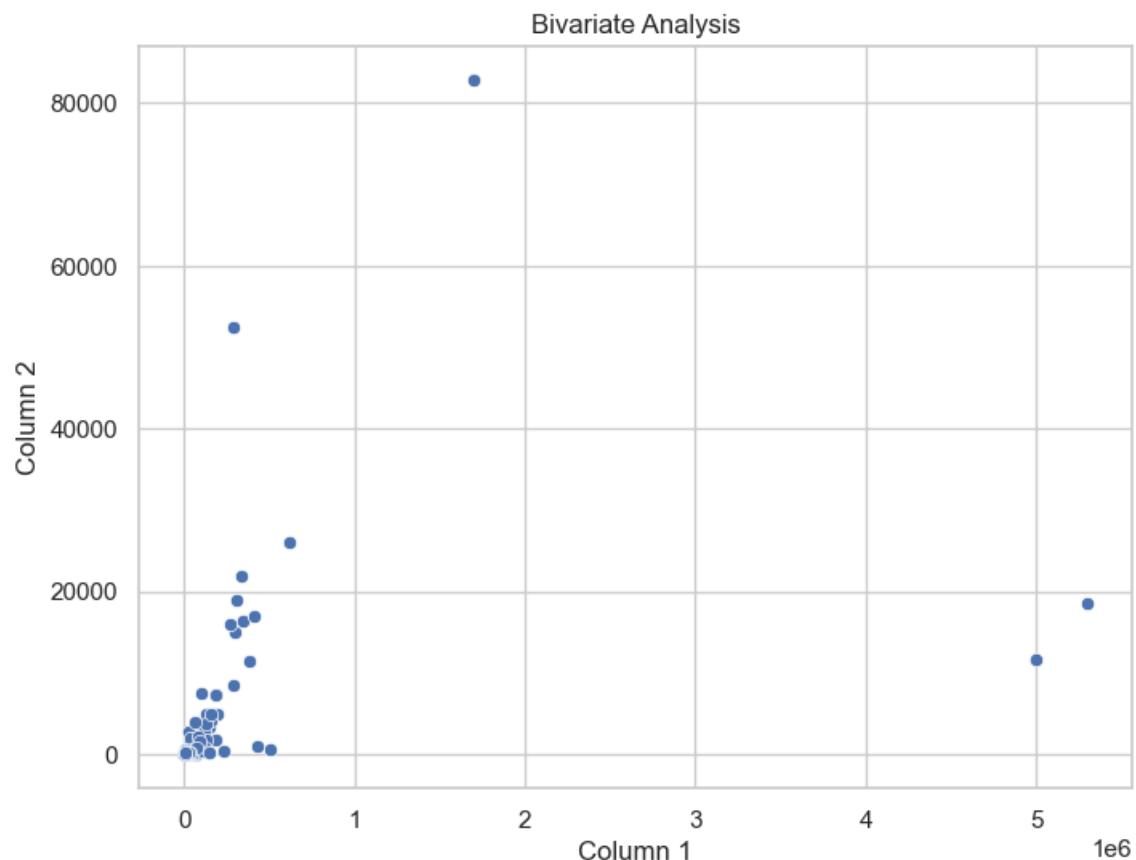
Mode: 0



Bivariate Analysis

In [327]:

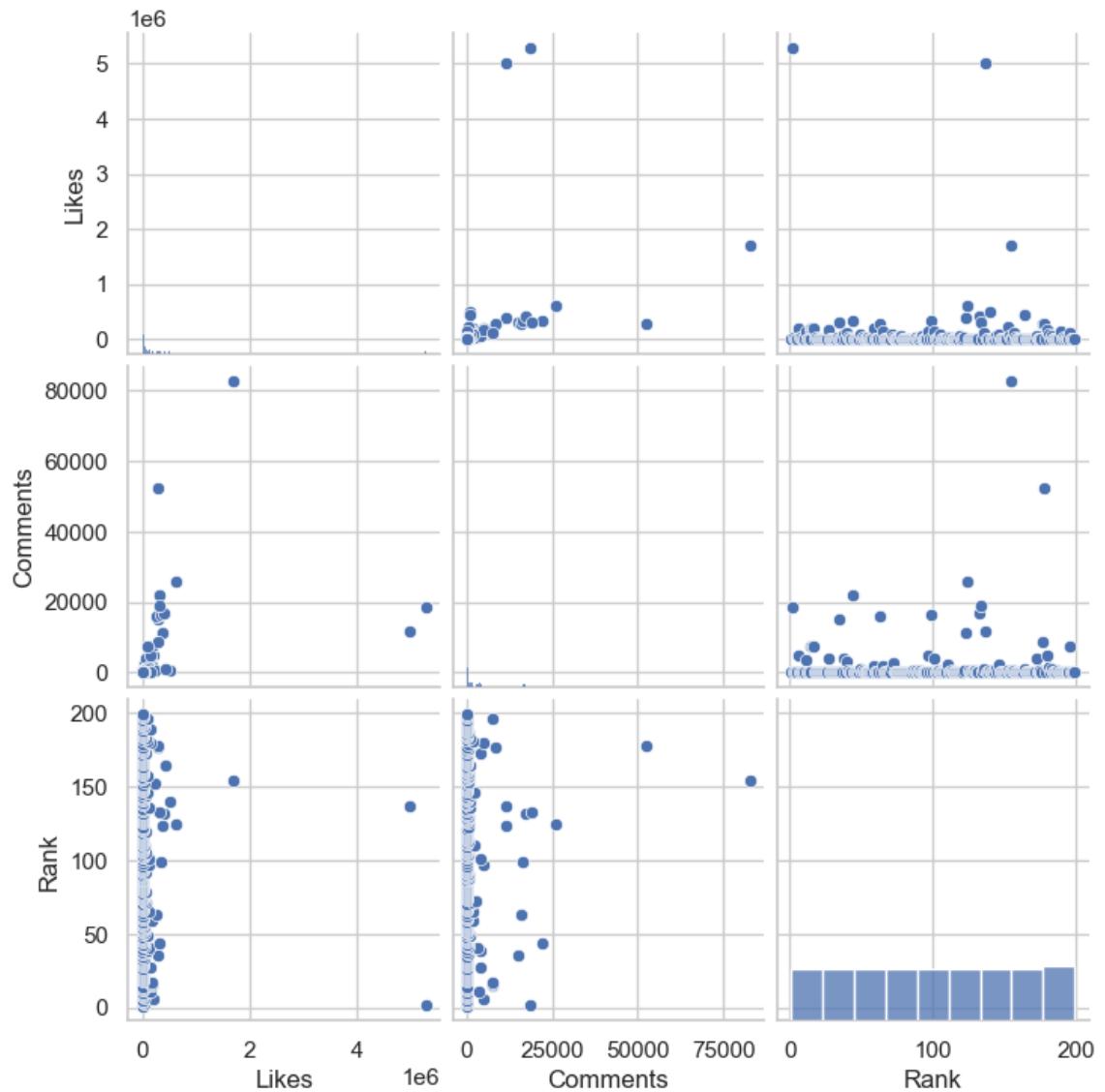
```
1 import pandas as pd
2 import matplotlib.pyplot as plt
3 import seaborn as sns
4 # Read the CSV file into a DataFrame
5
6 # Select two columns for analysis (replace 'column1' and 'column2' with
7 column1 = df['Likes']
8 column2 = df['Comments']
9 # Create scatter plot
10 plt.figure(figsize=(8, 6))
11 sns.scatterplot(x=column1, y=column2)
12 plt.title('Bivariate Analysis')
13 plt.xlabel('Column 1')
14 plt.ylabel('Column 2')
15 plt.show()
```



Multivariate analysis

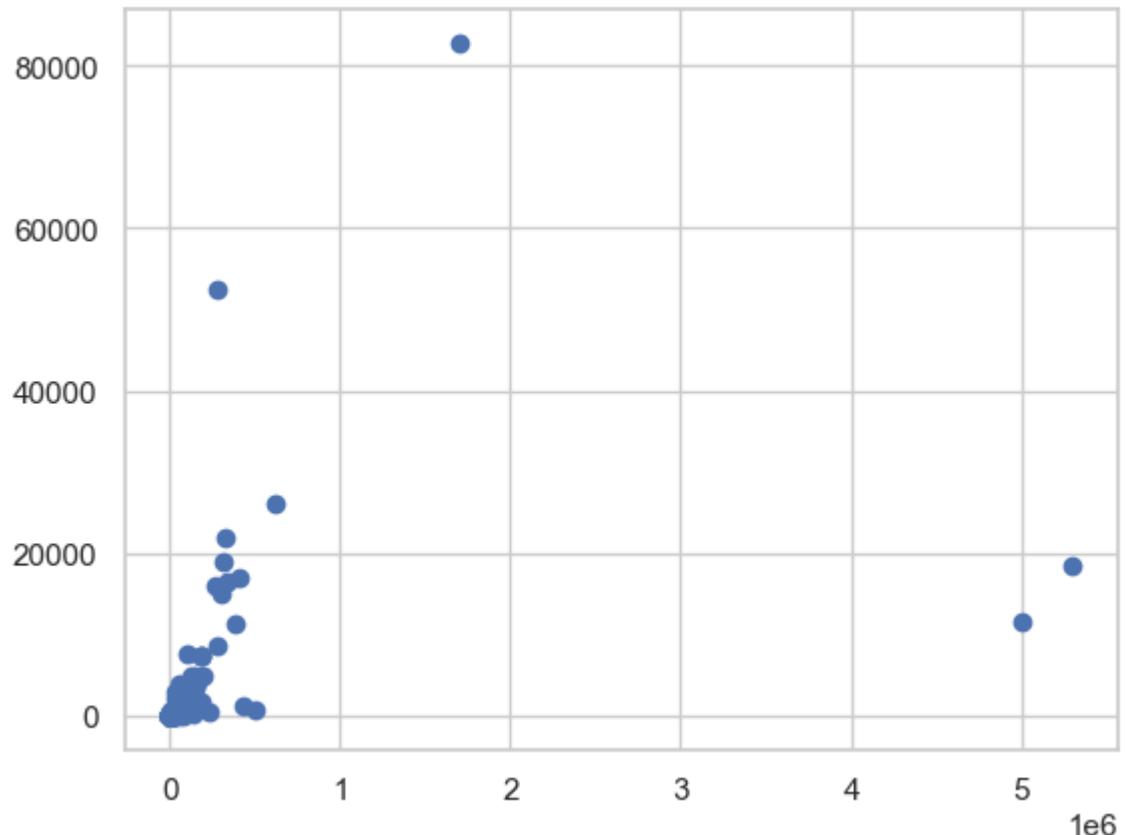
In [328]:

```
1 import pandas as pd
2 import seaborn as sns
3 numeric_data = df[['Likes', 'Comments', 'Rank']]
4 sns.pairplot(numeric_data)
5 plt.show()
```



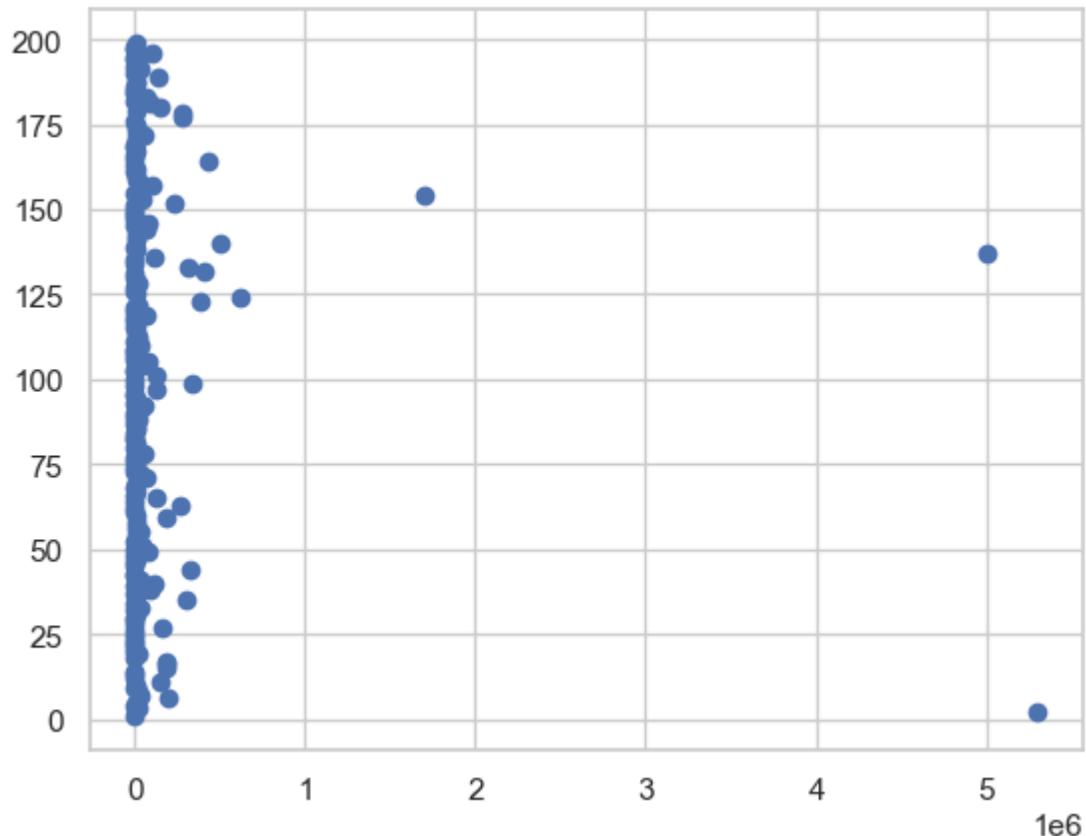
```
In [332]: 1 import matplotlib.pyplot as plt  
2 %matplotlib inline  
3 plt.scatter(df['Likes'], df['Comments'])
```

Out[332]: <matplotlib.collections.PathCollection at 0x242f0c7cc10>



```
In [335]: 1 plt.scatter(df["Likes"],df["Rank"])
```

```
Out[335]: <matplotlib.collections.PathCollection at 0x242ee336bb0>
```



X-Test Y-Test

```
In [336]: 1 x = df[['Likes','Comments']]
 2 y = df[['Rank']]
 3 x
```

```
Out[336]: Likes  Comments
```

| | Likes | Comments |
|-----|---------|----------|
| 0 | 2700 | 78 |
| 1 | 5300000 | 18500 |
| 2 | 24700 | 0 |
| 3 | 166 | 9 |
| 4 | 12400 | 0 |
| ... | ... | ... |
| 194 | 322 | 0 |
| 195 | 98300 | 7600 |
| 196 | 0 | 0 |
| 197 | 822 | 52 |
| 198 | 3100 | 197 |

199 rows × 2 columns

Loading [MathJax]/extensions/MathML/content-mathml.js

In [337]: 1 x

Out[337]:

| | Likes | Comments |
|-----|---------|----------|
| 0 | 2700 | 78 |
| 1 | 5300000 | 18500 |
| 2 | 24700 | 0 |
| 3 | 166 | 9 |
| 4 | 12400 | 0 |
| ... | ... | ... |
| 194 | 322 | 0 |
| 195 | 98300 | 7600 |
| 196 | 0 | 0 |
| 197 | 822 | 52 |
| 198 | 3100 | 197 |

199 rows × 2 columns

In [338]: 1 y

Out[338]:

| | Rank |
|-----|------|
| 0 | 1 |
| 1 | 2 |
| 2 | 3 |
| 3 | 4 |
| 4 | 5 |
| ... | ... |
| 194 | 195 |
| 195 | 196 |
| 196 | 197 |
| 197 | 198 |
| 198 | 199 |

199 rows × 1 columns

In [339]:

```
1 from sklearn.model_selection import train_test_split
2 x_train, x_test, y_train, y_test = train_test_split(x,y,test_size=0.2)
```

In [340]:

```
1 len(x_train)
```

Out[340]:

159

In [341]:

```
1 len(x_test)
```

Out[341]:

40

Loading [MathJax]/extensions/MathML/content-mathml.js

In [342]: 1 x_train

Out[342]:

| | Likes | Comments |
|------------|---------|----------|
| 57 | 7500 | 0 |
| 193 | 1800 | 20 |
| 51 | 430 | 21 |
| 25 | 2600 | 0 |
| 187 | 4300 | 17 |
| ... | ... | ... |
| 190 | 37700 | 192 |
| 136 | 5000000 | 11600 |
| 163 | 431100 | 1100 |
| 2 | 24700 | 0 |
| 125 | 1800 | 73 |

159 rows × 2 columns

In [343]: 1 y_train

Out[343]:

| | Rank |
|------------|------|
| 57 | 58 |
| 193 | 194 |
| 51 | 52 |
| 25 | 26 |
| 187 | 188 |
| ... | ... |
| 190 | 191 |
| 136 | 137 |
| 163 | 164 |
| 2 | 3 |
| 125 | 126 |

159 rows × 1 columns

In [344]: 1 x_test

Out[344]:

| | Likes | Comments |
|------------|--------|----------|
| 79 | 433 | 22 |
| 92 | 255 | 3 |
| 181 | 62 | 2 |
| 173 | 6600 | 0 |
| 31 | 2100 | 33 |
| 3 | 166 | 9 |
| 63 | 994 | 9 |
| 39 | 117100 | 3000 |
| 49 | 180 | 2 |
| 83 | 59 | 0 |
| 102 | 641 | 0 |
| 50 | 39800 | 593 |
| 56 | 8200 | 205 |
| 101 | 601 | 27 |
| 81 | 1800 | 0 |
| 189 | 471 | 24 |
| 158 | 8500 | 275 |
| 122 | 382000 | 11400 |
| 192 | 21300 | 214 |
| 14 | 180300 | 7400 |
| 84 | 10800 | 0 |
| 120 | 83 | 1 |
| 26 | 156500 | 4200 |
| 68 | 13100 | 0 |
| 124 | 13200 | 386 |
| 38 | 163 | 7 |
| 121 | 23300 | 504 |
| 156 | 106500 | 428 |
| 180 | 85800 | 1600 |
| 64 | 128900 | 1800 |
| 46 | 3400 | 0 |
| 12 | 615 | 21 |
| 76 | 211 | 9 |
| 172 | 6300 | 200 |
| 5 | 197300 | 4900 |
| 107 | 267 | 6 |
| 191 | 993 | 127 |
| 65 | 184 | 0 |
| 7 | 22100 | 0 |

Loading [MathJax]/extensions/MathML/content-mathml.js

| Likes | Comments |
|-------|----------|
| 117 | 88 |

Loading [MathJax]/extensions/MathML/content-mathml.js

In [345]: 1 y_test

Out[345]:

| | Rank |
|------------|------|
| 79 | 80 |
| 92 | 93 |
| 181 | 182 |
| 173 | 174 |
| 31 | 32 |
| 3 | 4 |
| 63 | 64 |
| 39 | 40 |
| 49 | 50 |
| 83 | 84 |
| 102 | 103 |
| 50 | 51 |
| 56 | 57 |
| 101 | 102 |
| 81 | 82 |
| 189 | 190 |
| 158 | 159 |
| 122 | 123 |
| 192 | 193 |
| 14 | 15 |
| 84 | 85 |
| 120 | 121 |
| 26 | 27 |
| 68 | 69 |
| 124 | 125 |
| 38 | 39 |
| 121 | 122 |
| 156 | 157 |
| 180 | 181 |
| 64 | 65 |
| 46 | 47 |
| 12 | 13 |
| 76 | 77 |
| 172 | 173 |
| 5 | 6 |
| 107 | 108 |
| 191 | 192 |
| 65 | 66 |

Loading [MathJax]/extensions/MathML/content-mathml.js

7 8

Rank

117 118

In [346]:

```
1 from sklearn.linear_model import LinearRegression
2 clf = LinearRegression()
3 clf.fit(x_train,y_train)
```

Out[346]: LinearRegression()

In [347]:

```
1 clf.predict(x_test)
```

Out[347]: array([[101.30787323],
[101.29595555],
[101.29674181],
[101.24501668],
[101.30273441],
[101.30083679],
[101.29446324],
[102.49251971],
[101.29583351],
[101.29536619],
[101.29088623],
[101.40417926],
[101.3760694],
[101.31007684],
[101.2819648],
[101.30897944],
[101.42271533],
[106.32806798],
[101.28152608],
[105.0832193],
[101.21268709],
[101.29588081],
[103.02846903],
[101.19498278],
[101.46416591],
[101.29946116],
[101.46894538],
[100.77536008],
[101.75434849],
[101.56245719],
[101.26964876],
[101.30577292],
[101.3004904],
[101.3871979],
[103.20396193],
[101.29796126],
[101.37699539],
[101.294404],
[101.12570507],
[101.29514296]])

In [348]:

```
1 clf.score(x_test, y_test)
```

Out[348]: -0.03599704198236631

In [349]: 1 df.info()

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 199 entries, 0 to 198
Data columns (total 10 columns):
 #   Column      Non-Null Count  Dtype  
--- 
 0   Rank        199 non-null    int64  
 1   Username    199 non-null    object  
 2   Categories  199 non-null    bool    
 3   Suscribers  199 non-null    int64  
 4   Country     199 non-null    object  
 5   Visits      199 non-null    int64  
 6   Likes       199 non-null    int64  
 7   Comments    199 non-null    int64  
 8   Links       199 non-null    object  
 9   Z-score     199 non-null    float64 
dtypes: bool(1), float64(1), int64(5), object(3)
memory usage: 14.3+ KB
```

```
In [350]: 1 x_train, x_test, y_train, y_test = train_test_split(x,y,test_size=0.3,r  
2 x_test
```

Out[350]:

| | Likes | Comments |
|------------|---------|----------|
| 59 | 4600 | 0 |
| 5 | 197300 | 4900 |
| 20 | 413 | 3 |
| 124 | 13200 | 386 |
| 52 | 7500 | 177 |
| 19 | 1100 | 4 |
| 161 | 5700 | 0 |
| 55 | 5700 | 0 |
| 69 | 10100 | 0 |
| 2 | 24700 | 0 |
| 98 | 341800 | 16500 |
| 10 | 146900 | 3400 |
| 75 | 3 | 0 |
| 134 | 1400 | 67 |
| 193 | 1800 | 20 |
| 63 | 994 | 9 |
| 110 | 442 | 8 |
| 78 | 3400 | 0 |
| 178 | 4100 | 305 |
| 114 | 623 | 0 |
| 149 | 1200 | 33 |
| 129 | 508 | 0 |
| 61 | 2400 | 0 |
| 87 | 20100 | 363 |
| 102 | 641 | 0 |
| 120 | 83 | 1 |
| 168 | 0 | 0 |
| 1 | 5300000 | 18500 |
| 47 | 242 | 13 |
| 171 | 58500 | 4000 |
| 185 | 2200 | 176 |
| 39 | 117100 | 3000 |
| 76 | 211 | 9 |
| 91 | 54100 | 214 |
| 35 | 10200 | 345 |
| 121 | 23300 | 504 |
| 169 | 3100 | 0 |
| 162 | 242 | 9 |
| 46 | 3400 | 0 |

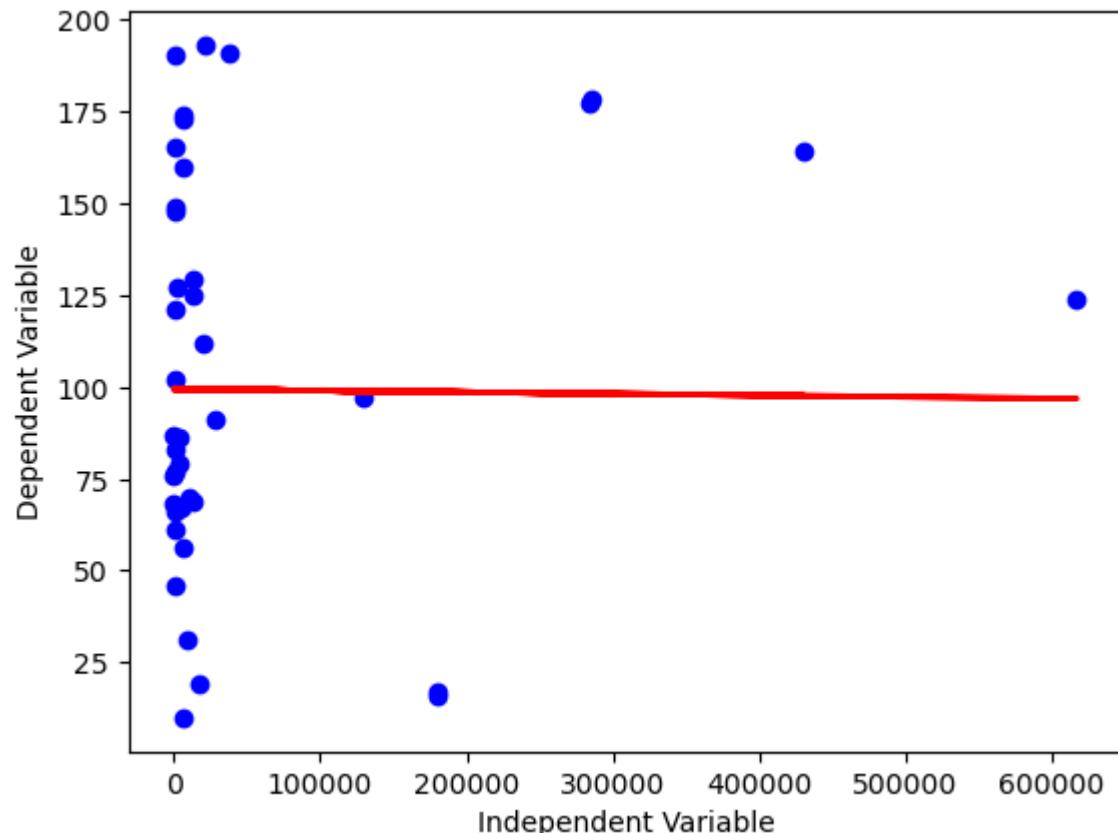
Loading [MathJax]/extensions/MathML/content-mathml.js

| | Likes | Comments |
|------------|--------|----------|
| 173 | 6600 | 0 |
| 189 | 471 | 24 |
| 7 | 22100 | 0 |
| 26 | 156500 | 4200 |
| 148 | 798 | 0 |
| 58 | 183400 | 1800 |
| 72 | 255 | 9 |
| 103 | 42400 | 541 |
| 198 | 3100 | 197 |
| 56 | 8200 | 205 |
| 184 | 651 | 0 |
| 24 | 231 | 9 |
| 43 | 330400 | 22000 |
| 101 | 601 | 27 |
| 142 | 7600 | 108 |
| 21 | 392 | 18 |
| 60 | 1100 | 30 |
| 174 | 8700 | 668 |
| 70 | 73500 | 1600 |
| 90 | 27500 | 357 |
| 49 | 180 | 2 |

Simple Linear Regression

```
In [3]:  
1 import pandas as pd  
2 from sklearn.model_selection import train_test_split  
3 from sklearn.linear_model import LinearRegression  
4 from sklearn.metrics import mean_squared_error  
5 import matplotlib.pyplot as plt  
6  
7 X = df[['Likes']] # Independent variable  
8 y = df['Rank'] # Dependent variable  
9  
10 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2)  
11  
12 model = LinearRegression()  
13 model.fit(X_train, y_train)  
14  
15 y_pred = model.predict(X_test)  
16 mse = mean_squared_error(y_test, y_pred)  
17 print(f"Mean Squared Error: {mse}")  
18  
19 plt.scatter(X_test, y_test, color='blue')  
20 plt.plot(X_test, y_pred, color='red', linewidth=2)  
21 plt.xlabel('Independent Variable')  
22 plt.ylabel('Dependent Variable')  
23 plt.show()  
24  
25
```

Mean Squared Error: 2893.7832522878493



Multiple Linear Regression

In [36]:

```

1 # Import the necessary libraries
2 import numpy as np
3 import pandas as pd
4 from sklearn.model_selection import train_test_split
5 from sklearn.linear_model import LinearRegression
6 from sklearn.metrics import mean_squared_error, r2_score
7
8 # Create or Load a dataset (example data)
9 data = {
10     'Feature1': [1.2, 2.2, 3.1, 4.5, 5.0],
11     'Feature2': [2.4, 3.9, 4.8, 5.9, 7.0],
12     'Feature3': [1.0, 2.0, 3.0, 4.0, 5.0],
13     'Target': [2.3, 3.6, 5.0, 5.5, 7.8]
14 }
15
16 df = pd.DataFrame(data)
17
18 # Split the data into training and testing sets
19 X = df[['Feature1', 'Feature2', 'Feature3']]
20 y = df['Target']
21 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2)
22
23 # Create a Linear Regression model
24 model = LinearRegression()
25
26 # Fit the model to the training data
27 model.fit(X_train, y_train)
28
29 # Make predictions on the test data
30 y_pred = model.predict(X_test)
31
32 # Evaluate the model
33 mse = mean_squared_error(y_test, y_pred)
34 r2 = r2_score(y_test, y_pred)
35
36 print(f"Mean Squared Error: {mse:.2f}")
37 print(f"R-squared: {r2:.2f}")
38
39 # Coefficients and intercept
40 coefficients = model.coef_
41 intercept = model.intercept_
42
43 print("Coefficients:", coefficients)
44 print("Intercept:", intercept)
45

```

Mean Squared Error: 0.04

R-squared: nan

Coefficients: [-2. -0.5 3.85]

Intercept: 2.04999999999978

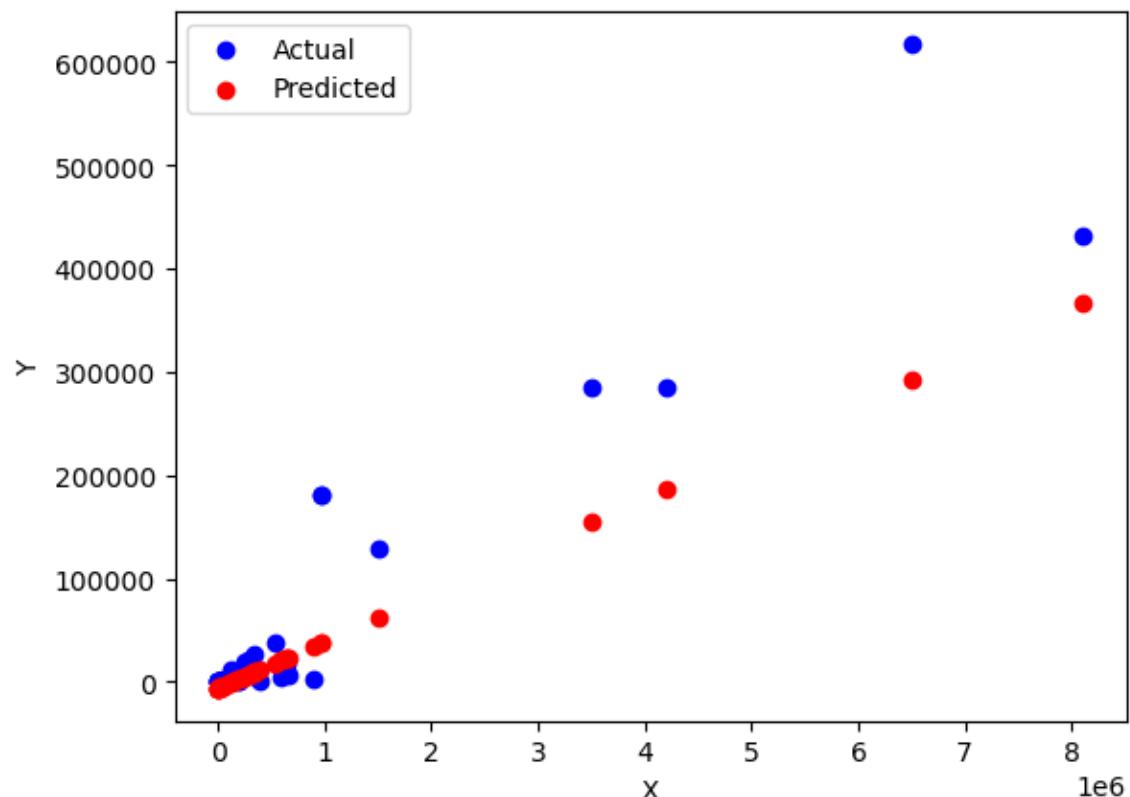
C:\Users\Anusha V\anaconda3\lib\site-packages\sklearn\metrics_regression.py:796: UndefinedMetricWarning: R^2 score is not well-defined with less than two samples.
 warnings.warn(msg, UndefinedMetricWarning)

polynomial Linear Regression

Loading [MathJax]/extensions/MathML/content-mathml.js

In [45]:

```
1 import pandas as pd
2 import numpy as np
3 from sklearn.model_selection import train_test_split
4 from sklearn.preprocessing import PolynomialFeatures
5 from sklearn.linear_model import LinearRegression
6 import matplotlib.pyplot as plt
7
8
9 X = df[['Visits']] # Features (input variable)
10 y = df['Likes'] # Target variable (output variable)
11
12 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2)
13
14 degree = 2 # Set the degree of the polynomial
15 poly = PolynomialFeatures(degree=degree)
16 X_train_poly = poly.fit_transform(X_train)
17 X_test_poly = poly.transform(X_test)
18
19 model = LinearRegression()
20 model.fit(X_train_poly, y_train)
21
22 y_pred = model.predict(X_test_poly)
23
24 # Plot the actual vs. predicted values
25 plt.scatter(X_test, y_test, color='blue', label='Actual')
26 plt.scatter(X_test, y_pred, color='red', label='Predicted')
27 plt.xlabel('X')
28 plt.ylabel('Y')
29 plt.legend()
30 plt.show()
31
```



Cross Validation

In [49]:

```
1 import numpy as np
2 from sklearn.model_selection import KFold
3 from sklearn.metrics import accuracy_score
4 from sklearn.datasets import load_iris
5 from sklearn.svm import SVC
6
7 # Load your dataset (replace this with your data)
8 data = load_iris()
9 X = data.data
10 y = data.target
11
12 # Set the number of folds (k)
13 k = 5
14 kf = KFold(n_splits=k, shuffle=True, random_state=42)
15
16 # Initialize a list to store the accuracy scores for each fold
17 accuracy_scores = []
18
19 # Iterate over each fold
20 for train_index, test_index in kf.split(X):
21     X_train, X_test = X[train_index], X[test_index]
22     y_train, y_test = y[train_index], y[test_index]
23
24     # Create and train your machine learning model (replace with your model)
25     model = SVC()
26     model.fit(X_train, y_train)
27
28     # Make predictions on the test set
29     y_pred = model.predict(X_test)
30
31     # Calculate the accuracy for this fold
32     accuracy = accuracy_score(y_test, y_pred)
33     accuracy_scores.append(accuracy)
34
35 # Calculate the mean and standard deviation of the accuracy scores
36 mean_accuracy = np.mean(accuracy_scores)
37 std_accuracy = np.std(accuracy_scores)
38
39 # Print the results
40 print(f"Accuracy scores for {k}-fold cross-validation: {accuracy_scores}")
41 print(f"Mean Accuracy: {mean_accuracy:.2f}")
42 print(f"Standard Deviation: {std_accuracy:.2f}")
43
```

Accuracy scores for 5-fold cross-validation: [1.0, 1.0, 0.9333333333333333, 0.9333333333333333, 0.9666666666666667]

Mean Accuracy: 0.97

Standard Deviation: 0.03

Evaluation regression model

In [37]:

```

1 # Import the necessary libraries
2 import numpy as np
3 import pandas as pd
4 from sklearn.model_selection import train_test_split
5 from sklearn.linear_model import LinearRegression
6 from sklearn.metrics import mean_squared_error, r2_score, explained_variance
7
8 # Create or Load a dataset (example data)
9 data = {
10     'Feature1': [1.2, 2.2, 3.1, 4.5, 5.0],
11     'Feature2': [2.4, 3.9, 4.8, 5.9, 7.0],
12     'Feature3': [1.0, 2.0, 3.0, 4.0, 5.0],
13     'Target': [2.3, 3.6, 5.0, 5.5, 7.8]
14 }
15
16 df = pd.DataFrame(data)
17
18 # Split the data into training and testing sets
19 X = df[['Feature1', 'Feature2', 'Feature3']]
20 y = df['Target']
21 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2)
22
23 # Create a Linear Regression model
24 model = LinearRegression()
25
26 # Fit the model to the training data
27 model.fit(X_train, y_train)
28
29 # Make predictions on the test data
30 y_pred = model.predict(X_test)
31
32 # Evaluate the model
33 mse = mean_squared_error(y_test, y_pred) # Mean Squared Error
34 r2 = r2_score(y_test, y_pred) # R-squared (coefficient of determination)
35 explained_variance = explained_variance_score(y_test, y_pred) # Explained Variance Score
36 mae = mean_absolute_error(y_test, y_pred) # Mean Absolute Error
37
38 print(f"Mean Squared Error: {mse:.2f}")
39 print(f"R-squared: {r2:.2f}")
40 print(f"Explained Variance Score: {explained_variance:.2f}")
41 print(f"Mean Absolute Error: {mae:.2f}")
42

```

Mean Squared Error: 0.04

R-squared: nan

Explained Variance Score: 1.00

Mean Absolute Error: 0.20

C:\Users\Anusha V\anaconda3\lib\site-packages\sklearn\metrics_regression.py:796: UndefinedMetricWarning: R^2 score is not well-defined with less than two samples.

```
warnings.warn(msg, UndefinedMetricWarning)
```

Random Forest Regressor

```
In [4]: 1 from sklearn.ensemble import RandomForestRegressor
2 from sklearn.datasets import make_regression
3 X, y = make_regression(n_features=4, n_informative=2, random_state=0, s
4 rfr = RandomForestRegressor(max_depth=3)
5 rfr.fit(X, y)
6 print(rfr.predict([[0, 1, 0, 1]]))
```

[32.85592054]

Lasso Regression

```
In [15]: 1 from sklearn import linear_model
2 import numpy as np
3 rng = np.random.RandomState(1)
4 X = np.sort(5 * rng.rand(80, 1), axis=0)
5 y = np.sin(X).ravel()
6 y[::5] += 3 * (0.5 - rng.rand(16))
7 # Fit regression model
8 lassoReg = linear_model.Lasso(alpha=0.1)
9 lassoReg.fit(X,y)
10 # Predict
11 X_test = np.arange(0.0, 5.0, 1)[:, np.newaxis]
12 lassoReg.predict(X_test)
13
```

Out[15]: array([0.78305084, 0.49957596, 0.21610108, -0.0673738 , -0.35084868])

In [9]: 1 X

```
Out[9]: array([[5.71874087e-04],
   [9.14413867e-02],
   [9.68347894e-02],
   [1.36937966e-01],
   [1.95273916e-01],
   [2.49767295e-01],
   [2.66812726e-01],
   [4.25221057e-01],
   [4.61692974e-01],
   [4.91734169e-01],
   [5.11672144e-01],
   [5.16130033e-01],
   [6.50142861e-01],
   [6.87373521e-01],
   [6.96381736e-01],
   [7.01934693e-01],
   [7.33642875e-01],
   [7.33779454e-01],
   [8.26770986e-01],
   [8.46150000e-01],
```

In [11]: 1 y

```
Out[11]: array([
  5.71874056e-04,  9.13140084e-02,  9.66835240e-02,  1.36510390e-01,
  1.94035253e-01,  2.47178482e-01,  2.63658285e-01,  4.12522163e-01,
  4.45464463e-01,  4.72155294e-01,  4.89635918e-01,  4.93517994e-01,
  6.05300129e-01,  6.34509327e-01,  6.41446076e-01,  6.45696215e-01,
  6.69579781e-01,  6.69681218e-01,  7.35748344e-01,  7.50720534e-01,
  8.02397078e-01,  8.36304301e-01,  8.53289303e-01,  8.71444956e-01,
  9.70605329e-01,  9.85824622e-01,  9.91311217e-01,  9.94728365e-01,
  9.98252126e-01,  9.99993246e-01,  9.99977003e-01,  9.87699652e-01,
  9.85915646e-01,  9.15904395e-01,  9.14069859e-01,  8.77830014e-01,
  8.77534605e-01,  8.70630546e-01,  8.69934082e-01,  8.65235592e-01,
  8.60400242e-01,  7.84646075e-01,  6.31686553e-01,  5.37227991e-01,
  4.58019714e-01,  4.45838347e-01,  4.32721158e-01,  3.41153283e-01,
  2.67699531e-01,  2.07303185e-01,  1.93830272e-01,  -1.76451855e-01,
  -2.09188375e-01,  -2.49907779e-01,  -2.80682239e-01,  -2.86826010e-01,
  -3.12470754e-01,  -3.14585948e-01,  -3.24429114e-01,  -3.49651545e-01,
  -4.43974821e-01,  -4.69211146e-01,  -5.64011437e-01,  -5.72153265e-01,
  -5.74888492e-01,  -7.20694771e-01,  -7.59230658e-01,  -7.80436587e-01,
  -8.58088588e-01,  -9.45898629e-01,  -9.48667092e-01,  -9.48706725e-01,
  -9.71490934e-01,  -9.80974083e-01,  -9.85684166e-01,  -9.97200340e-01,
  -9.99943982e-01,  -9.97032450e-01,  -9.91701462e-01,  -9.73227696e-01
])
```

Support Vector Regression

In [2]:

```
1 from sklearn.svm import SVR
2 import numpy as np
3 rng = np.random.RandomState(1)
4 X = np.sort(5 * rng.rand(80, 1), axis=0)
5 y = np.sin(X).ravel()
6 y[::5] += 3 * (0.5 - rng.rand(16))
7 # Fit regression model
8 svr = SVR().fit(X, y)
9 # Predict
10 X_test = np.arange(0.0, 5.0, 1)[:, np.newaxis]
11 svr.predict(X_test)
12
```

Out[2]: array([-0.07840308, 0.78077042, 0.81326895, 0.08638149, -0.6928019])

In [44]:

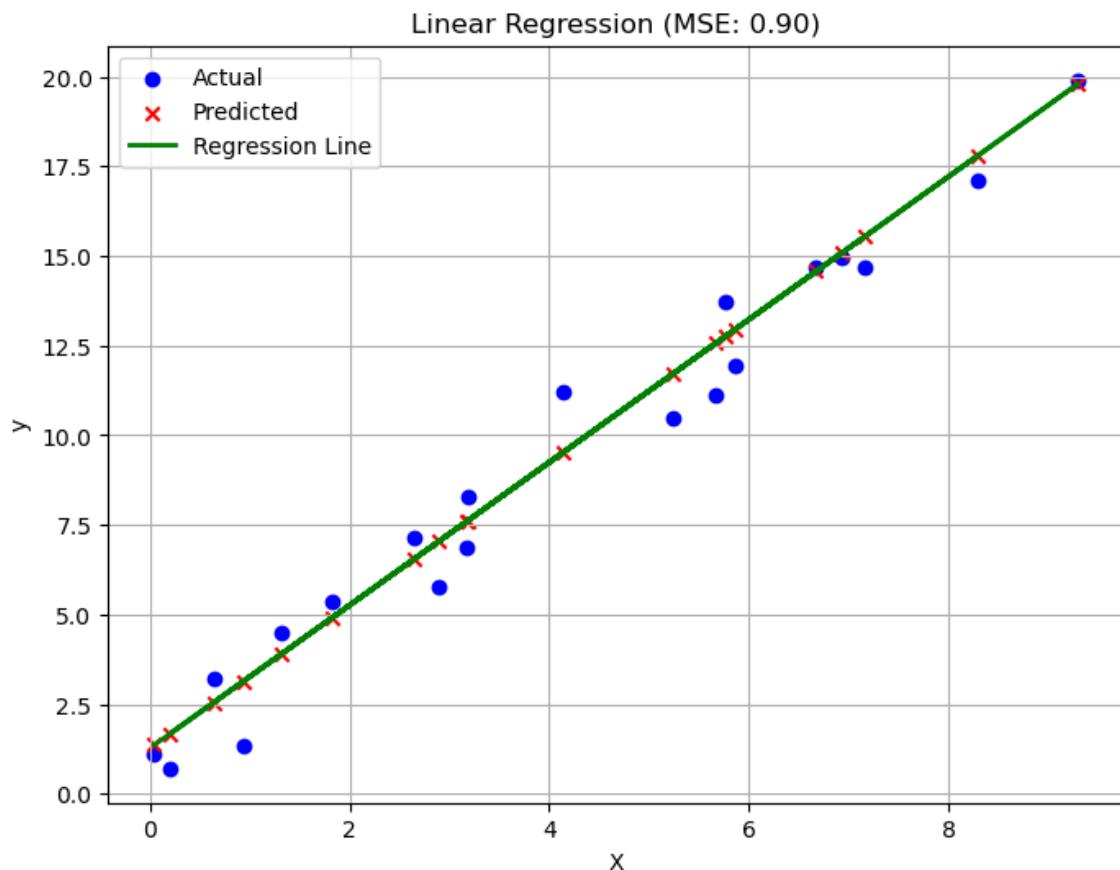
```
1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 from sklearn.model_selection import train_test_split
5 from sklearn.svm import SVR
6 from sklearn.metrics import mean_squared_error, r2_score
7 # Load your dataset from a CSV file
8 df = pd.read_csv (r"C:\Users\Anusha V\Desktop\YouTube.csv")
9 # Split the data into features (X) and the target variable (y)
10 X = data[['Likes']] # Replace 'feature_column' with your feature column
11 y = data['Rank'] # Adjust 'target_column' to your target variable
12 # Split the dataset into a training set and a testing set
13 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3)
14 # Create and train the Support Vector Regression model
15 model = SVR(kernel='linear') # You can choose the kernel type (e.g., 'L
16 model.fit(X_train, y_train)
17 # Make predictions on the testing set
18 y_pred = model.predict(X_test)
19 # Evaluate the model's performance using regression metrics
20 mse = mean_squared_error(y_test, y_pred)
21 r2 = r2_score(y_test, y_pred)
22 # Plot the true values and predicted values
23 plt.scatter(X_test, y_test, color='b', label='True')
24 plt.scatter(X_test, y_pred, color='r', label='Predicted')
25 plt.title('SVR - True vs. Predicted')
26 plt.xlabel('Feature')
27 plt.ylabel('Target')
28 plt.legend(loc='upper left')
29 plt.show()
30 print(f'Mean Squared Error (MSE): {mse:.4f}')
31 print(f'R-squared (R2): {r2:.4f}')
32
```

```
File "C:\Users\Anusha V\AppData\Local\Temp\ipykernel_8816\2609282417.py", line 15
    model = SVR(kernel='linear') # You can choose the kernel type (e.g.,
^line
SyntaxError: positional argument follows keyword argument
```

linear regression model

In [28]:

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 from sklearn.linear_model import LinearRegression
4 from sklearn.metrics import mean_squared_error
5
6 # Generate synthetic data
7 np.random.seed(0)
8 X = np.random.rand(100, 1) * 10
9 y = 2 * X + 1 + np.random.randn(100, 1)
10
11 # Split the data into training and testing sets
12 X_train, X_test, y_train, y_test = X[:80], X[80:], y[:80], y[80:]
13
14 # Fit a linear regression model
15 model = LinearRegression()
16 model.fit(X_train, y_train)
17
18 # Make predictions on the test set
19 predictions = model.predict(X_test)
20
21 # Calculate the overall MSE
22 mse = mean_squared_error(y_test, predictions)
23
24 # Create a scatter plot for the linear regression
25 plt.figure(figsize=(8, 6))
26 plt.scatter(X_test, y_test, color='blue', label='Actual')
27 plt.scatter(X_test, predictions, color='red', marker='x', label='Predictions')
28 plt.plot(X_test, predictions, color='green', linewidth=2, label='Regression Line')
29 plt.xlabel('X')
30 plt.ylabel('y')
31 plt.title(f'Linear Regression (MSE: {mse:.2f})')
32 plt.legend()
33 plt.grid()
34 plt.show()
35
```



Mean Absolute Error(MAE)

In [13]:

```

1 # from sklearn.metrics import mean_absolute_error
2
3 # Actual values (ground truth)
4 actual = [10, 20, 30, 40, 50]
5
6 # Predicted values from your model
7 predicted = [12, 18, 28, 37, 45]
8
9 # Calculate the Mean Absolute Error
10 mae = mean_absolute_error(actual, predicted)
11
12 print(f"Mean Absolute Error: {mae}")

```

Mean Absolute Error: 2.8

In [12]:

```

1 import pandas as pd
2 from sklearn.metrics import mean_absolute_error
3 df = pd.read_csv ('C:\Users\Anusha V\Desktop\YouTube.csv')
4 actual = df['Likes']
5 predicted = df['Rank']
6 mae = mean_absolute_error(actual, predicted)
7 print(f"Mean Absolute Error (MAE): {mae}")

```

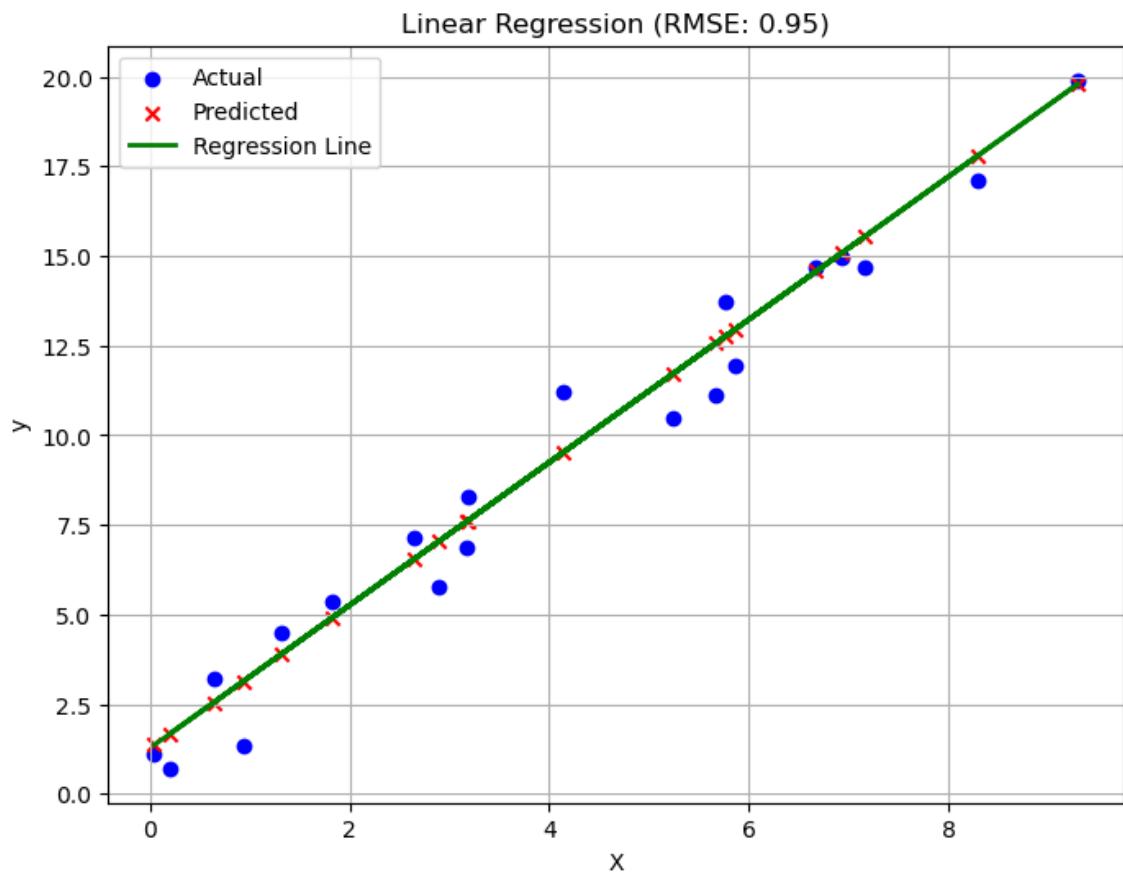
Mean Absolute Error (MAE): 104689.26130653266

root mean squared error

Loading [MathJax]/extensions/MathML/content-mathml.js

In [33]:

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 from sklearn.linear_model import LinearRegression
4 from sklearn.metrics import mean_squared_error
5
6 # Generate synthetic data
7 np.random.seed(0)
8 X = np.random.rand(100, 1) * 10
9 y = 2 * X + 1 + np.random.randn(100, 1)
10
11 # Split the data into training and testing sets
12 X_train, X_test, y_train, y_test = X[:80], X[80:], y[:80], y[80:]
13
14 # Fit a linear regression model
15 model = LinearRegression()
16 model.fit(X_train, y_train)
17
18 # Make predictions on the test set
19 predictions = model.predict(X_test)
20
21 # Calculate squared errors for each prediction
22 squared_errors = (predictions - y_test) ** 2
23
24 # Calculate the overall RMSE
25 rmse = np.sqrt(mean_squared_error(y_test, predictions))
26
27 # Create a scatter plot for the linear regression with RMSE
28 plt.figure(figsize=(8, 6))
29 plt.scatter(X_test, y_test, color='blue', label='Actual')
30 plt.scatter(X_test, predictions, color='red', marker='x', label='Predictions')
31 plt.plot(X_test, predictions, color='green', linewidth=2, label='Regression Line')
32 plt.xlabel('X')
33 plt.ylabel('y')
34 plt.title(f'Linear Regression (RMSE: {rmse:.2f})')
35 plt.legend()
36 plt.grid()
37 plt.show()
```



In [9]:

```

1 import numpy as np
2 from sklearn.metrics import mean_squared_error
3 from math import sqrt
4
5 # Actual values
6 actual = [10, 12, 15, 18, 20]
7
8 # Predicted values
9 predicted = [9, 11, 14, 17, 19]
10
11 # Calculate Mean Squared Error (MSE)
12 mse = mean_squared_error(actual, predicted)
13
14 # Calculate RMSE
15 rmse = sqrt(mse)
16
17 print("RMSE:", rmse)

```

RMSE: 1.0

R-squared (R²) score

```
In [49]: 1 from sklearn.metrics import r2_score
2
3 # Assuming you have observed actual and predicted values
4 actual_values = [10, 20, 30, 40, 50]
5 predicted_values = [12, 18, 28, 38, 52]
6
7 # Calculate the R-squared (R2) score
8 r2 = r2_score(actual_values, predicted_values)
9
10 print(f"R-squared (R2) Score: {r2}")
```

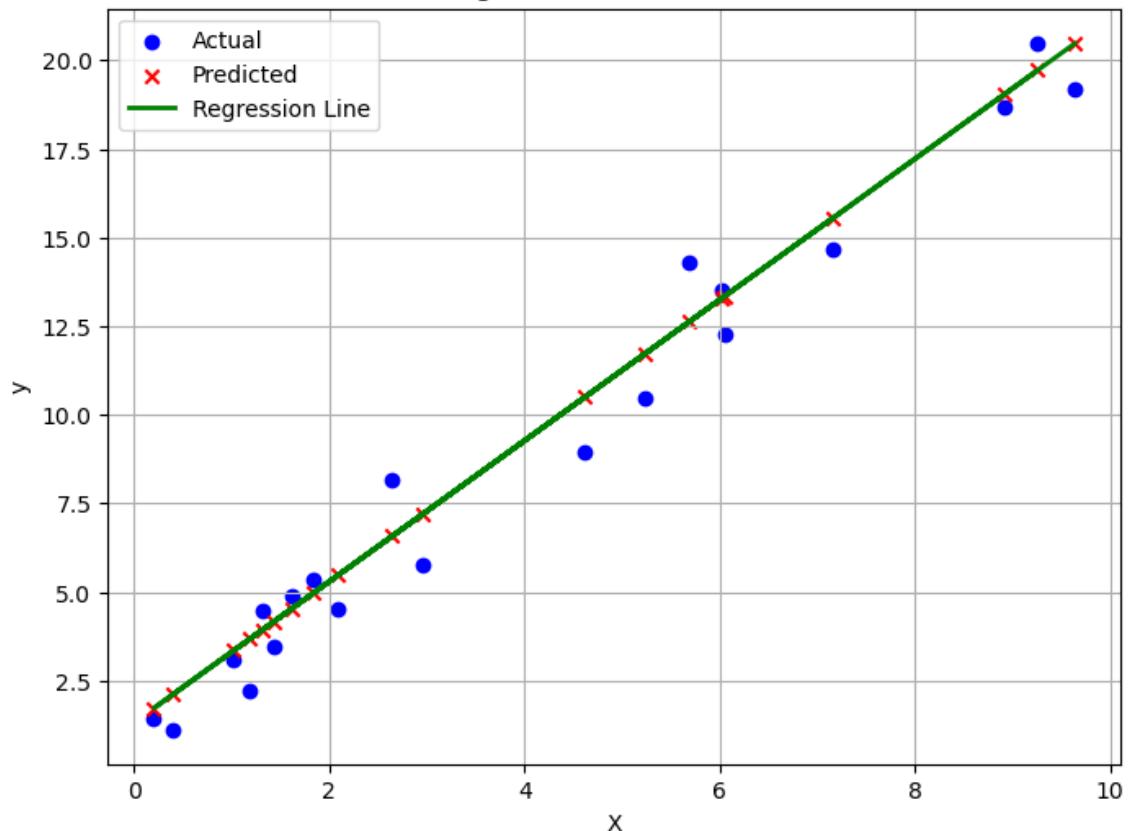
R-squared (R2) Score: 0.98

holdout method

In [35]:

```
1 import matplotlib.pyplot as plt
2 from sklearn.model_selection import train_test_split
3 from sklearn.linear_model import LinearRegression
4
5 # Generate synthetic data
6 np.random.seed(0)
7 X = np.random.rand(100, 1) * 10
8 y = 2 * X + 1 + np.random.randn(100, 1)
9
10 # Split the data into a training and testing set
11 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2)
12
13 # Fit a Linear regression model on the training data
14 model = LinearRegression()
15 model.fit(X_train, y_train)
16
17 # Make predictions on the test set
18 predictions = model.predict(X_test)
19
20 # Create a scatter plot for the linear regression
21 plt.figure(figsize=(8, 6))
22 plt.scatter(X_test, y_test, color='blue', label='Actual')
23 plt.scatter(X_test, predictions, color='red', marker='x', label='Predicted')
24 plt.plot(X_test, predictions, color='green', linewidth=2, label='Regression Line')
25 plt.xlabel('X')
26 plt.ylabel('y')
27 plt.title('Linear Regression with Holdout Method')
28 plt.legend()
29 plt.grid()
30 plt.show()
```

Linear Regression with Holdout Method



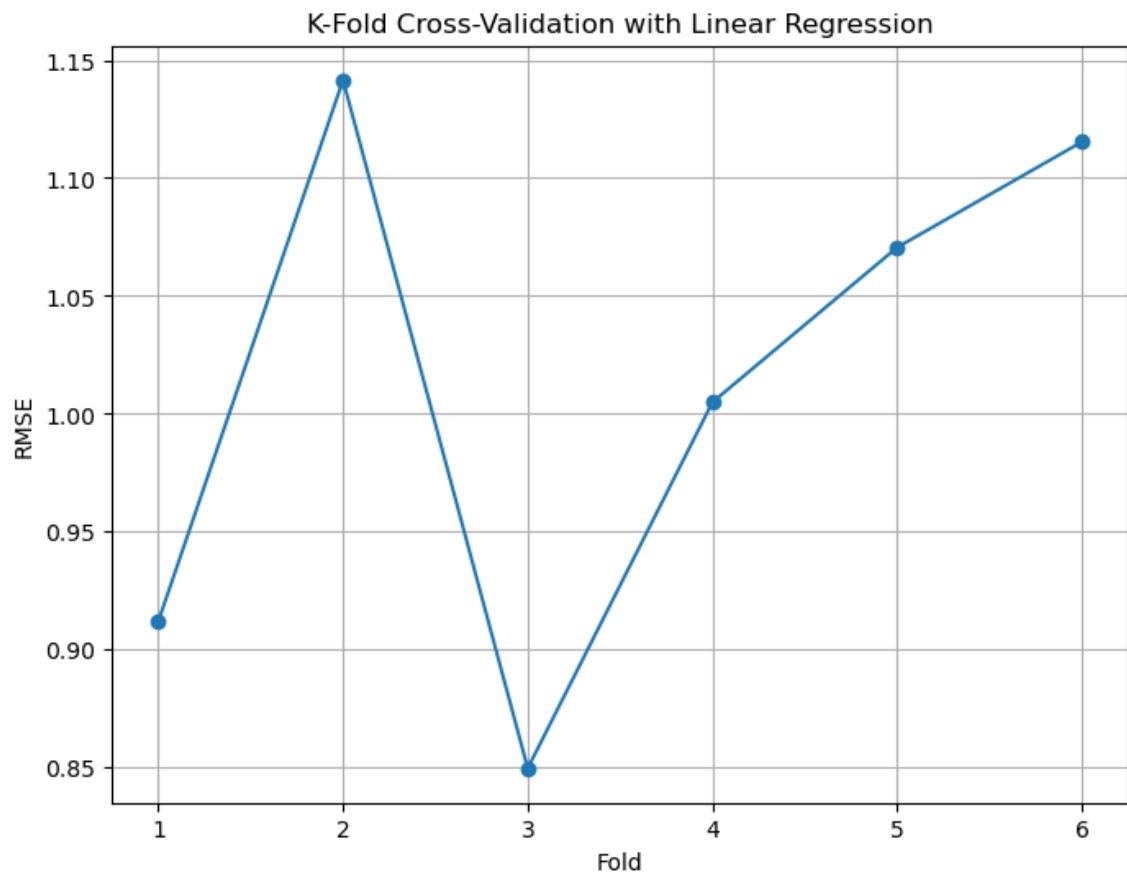
```
In [12]: 1 from sklearn.model_selection import train_test_split
2 from sklearn.datasets import load_iris # You can replace this with your
3
4 # Load your dataset (replace with your data loading code)
5 data = load_iris()
6
7 # Split the dataset into features (X) and target labels (y)
8 X = data.data
9 y = data.target
10
11 # Split the data into a training set and a test set
12 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2)
13 X
14
15 # Now, you can train your machine learning model on X_train and y_train
16 # and test its performance on X_test and y_test.
```

```
Out[12]: array([[5.1, 3.5, 1.4, 0.2],
   [4.9, 3. , 1.4, 0.2],
   [4.7, 3.2, 1.3, 0.2],
   [4.6, 3.1, 1.5, 0.2],
   [5. , 3.6, 1.4, 0.2],
   [5.4, 3.9, 1.7, 0.4],
   [4.6, 3.4, 1.4, 0.3],
   [5. , 3.4, 1.5, 0.2],
   [4.4, 2.9, 1.4, 0.2],
   [4.9, 3.1, 1.5, 0.1],
   [5.4, 3.7, 1.5, 0.2],
   [4.8, 3.4, 1.6, 0.2],
   [4.8, 3. , 1.4, 0.1],
   [4.3, 3. , 1.1, 0.1],
   [5.8, 4. , 1.2, 0.2],
   [5.7, 4.4, 1.5, 0.4],
   [5.4, 3.9, 1.3, 0.4],
   [5.1, 3.5, 1.4, 0.3],
   [5.7, 3.8, 1.7, 0.3],
```

K-Fold Cross-Validation

In [37]:

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 from sklearn.model_selection import KFold
4 from sklearn.linear_model import LinearRegression
5
6 # Generate synthetic data
7 np.random.seed(0)
8 X = np.random.rand(100, 1)
9 y = 2 * X + 1 + np.random.randn(100, 1)
10
11 # Number of folds for cross-validation
12 n_folds = 6
13
14 # Initialize a KFold cross-validation splitter
15 kf = KFold(n_splits=n_folds, shuffle=True, random_state=0)
16
17 # Initialize an array to store RMSE for each fold
18 rmse_values = []
19
20 # Perform K-fold cross-validation
21 for train_index, test_index in kf.split(X):
22     X_train, X_test = X[train_index], X[test_index]
23     y_train, y_test = y[train_index], y[test_index]
24
25     # Fit a Linear regression model on the training data
26     model = LinearRegression()
27     model.fit(X_train, y_train)
28
29     # Make predictions on the test set
30     predictions = model.predict(X_test)
31
32     # Calculate RMSE for this fold
33     rmse = np.sqrt(np.mean((predictions - y_test) ** 2))
34     rmse_values.append(rmse)
35
36 # Plot the RMSE values for each fold
37 plt.figure(figsize=(8, 6))
38 plt.plot(range(1, n_folds + 1), rmse_values, marker='o', linestyle='--')
39 plt.xlabel('Fold')
40 plt.ylabel('RMSE')
41 plt.title('K-Fold Cross-Validation with Linear Regression')
42 plt.grid()
43 plt.show()
```

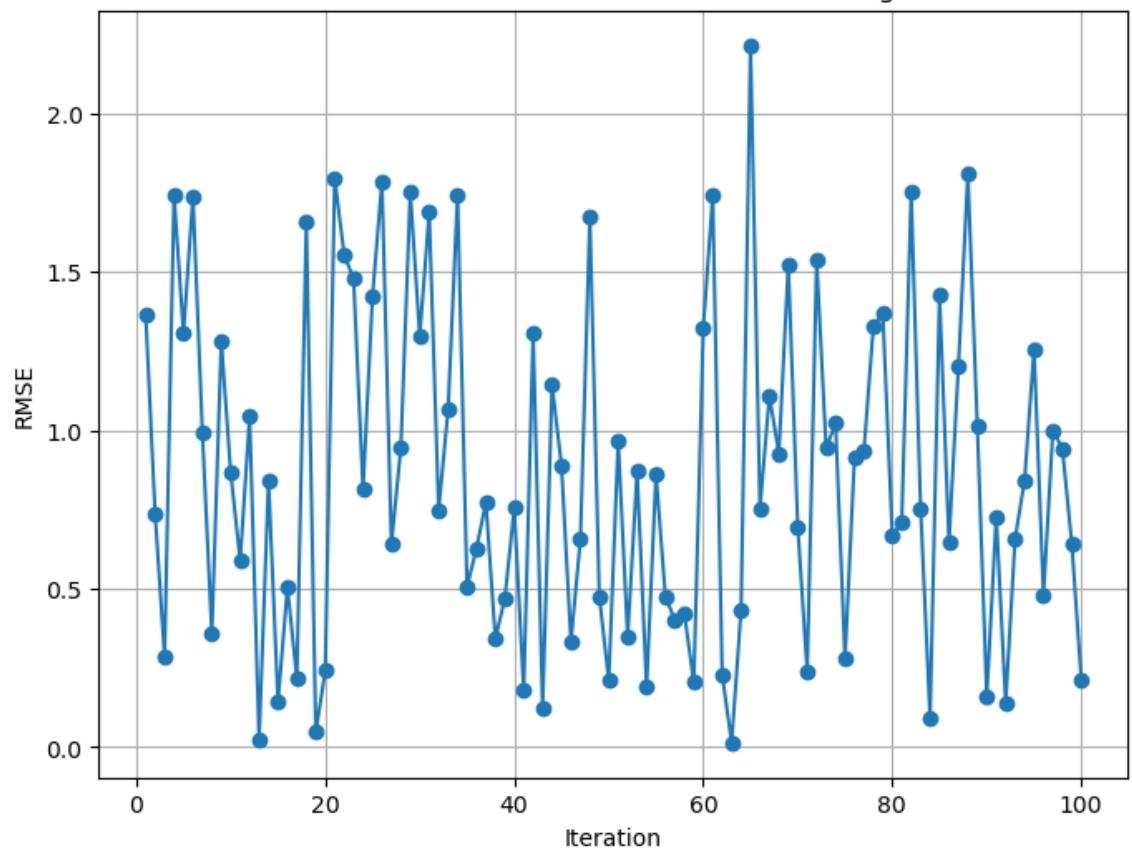


Leave-one-out cross validation

In [31]:

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 from sklearn.model_selection import LeaveOneOut
4 from sklearn.linear_model import LinearRegression
5
6 # Generate synthetic data
7 np.random.seed(0)
8 X = np.random.rand(100, 1) * 10
9 y = 2 * X + 1 + np.random.randn(100, 1)
10
11 # Initialize Leave-One-Out Cross-Validation
12 loo = LeaveOneOut()
13
14 # Initialize an array to store RMSE for each iteration
15 rmse_values = []
16
17 # Perform Leave-One-Out Cross-Validation
18 for train_index, test_index in loo.split(X):
19     X_train, X_test = X[train_index], X[test_index]
20     y_train, y_test = y[train_index], y[test_index]
21
22     # Fit a Linear regression model on the training data
23     model = LinearRegression()
24     model.fit(X_train, y_train)
25
26     # Make predictions on the test data point
27     prediction = model.predict(X_test)
28
29     # Calculate RMSE for this iteration
30     rmse = np.sqrt((prediction - y_test) ** 2)
31     rmse_values.append(rmse)
32
33 # Flatten the RMSE values array
34 rmse_values = np.array(rmse_values).flatten()
35
36 # Plot the RMSE values for each iteration
37 plt.figure(figsize=(8, 6))
38 plt.plot(range(1, len(rmse_values) + 1), rmse_values, marker='o', lines
39 plt.xlabel('Iteration')
40 plt.ylabel('RMSE')
41 plt.title('Leave-One-Out Cross-Validation with Linear Regression')
42 plt.grid()
43 plt.show()
44 # Calculate the average RMSE
45 average_rmse = np.mean(rmse_values)
46 print(f'Average RMSE: {average_rmse:.2f}')
```

Leave-One-Out Cross-Validation with Linear Regression



Average RMSE: 0.87

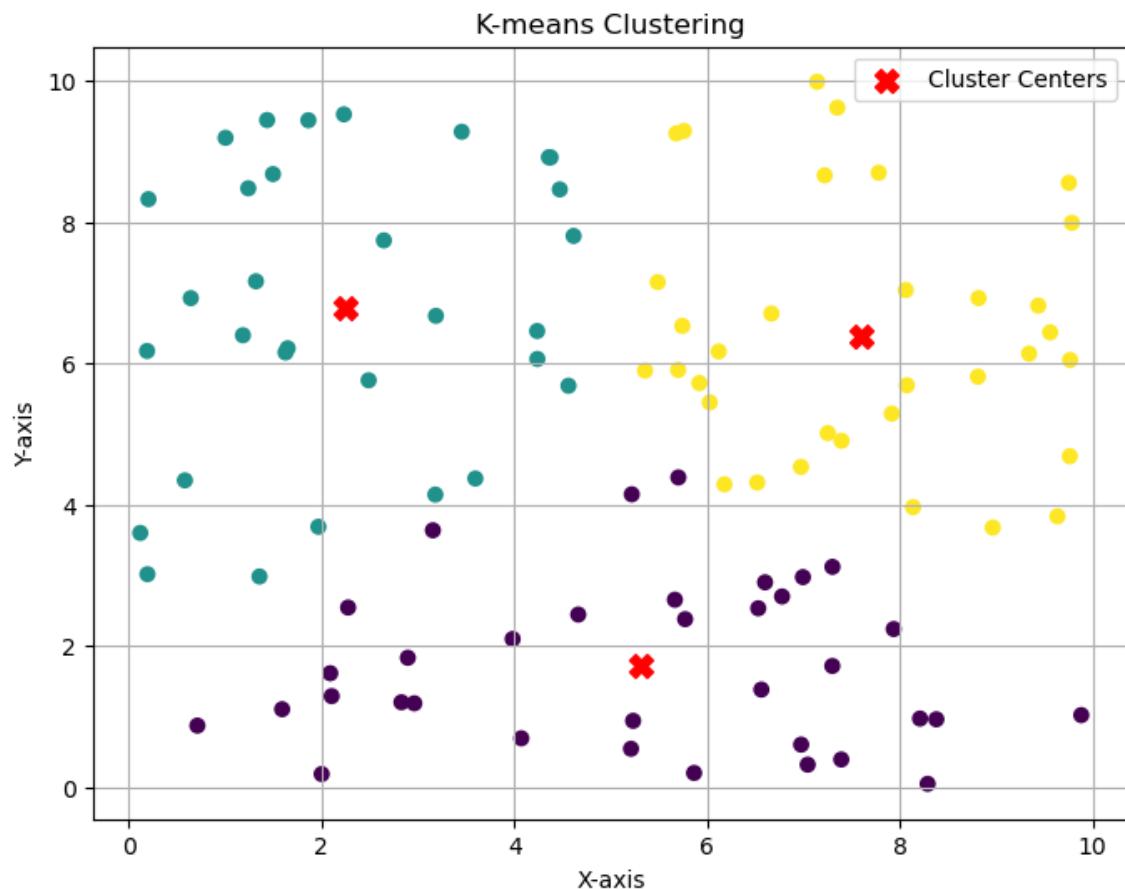
K-means

In [30]:

```

1 import numpy as np
2 import matplotlib.pyplot as plt
3 from sklearn.cluster import KMeans
4
5 # Generate synthetic data
6 np.random.seed(0)
7 X = np.random.rand(100, 2) * 10
8
9 # Create a K-means model with 3 clusters
10 kmeans = KMeans(n_clusters=3, random_state=0)
11
12 # Fit the K-means model to the data
13 kmeans.fit(X)
14
15 # Get cluster labels and cluster centers
16 cluster_labels = kmeans.labels_
17 cluster_centers = kmeans.cluster_centers_
18
19 # Plot the data points and cluster centers
20 plt.figure(figsize=(8, 6))
21 plt.scatter(X[:, 0], X[:, 1], c=cluster_labels, cmap='viridis')
22 plt.scatter(cluster_centers[:, 0], cluster_centers[:, 1], c='red', marker='x')
23 plt.xlabel('X-axis')
24 plt.ylabel('Y-axis')
25 plt.title('K-means Clustering')
26 plt.legend()
27 plt.grid()
28 plt.show()

```

**R-SQUAR**

Loading [MathJax]/extensions/MathML/content-mathml.js

In [39]:

```

1 import numpy as np
2 from sklearn.linear_model import LinearRegression
3 from sklearn.metrics import r2_score
4
5 # Generate synthetic data
6 np.random.seed(0)
7 X = np.random.rand(100, 1) * 10
8 y = 2 * X + 1 + np.random.randn(100, 1)
9
10 # Split the data into training and testing sets
11 X_train, X_test, y_train, y_test = X[:80], X[80:], y[:80], y[80:]
12
13 # Fit a Linear regression model
14 model = LinearRegression()
15 model.fit(X_train, y_train)
16
17 # Make predictions on the test set
18 predictions = model.predict(X_test)
19
20 # Calculate R-squared
21 r_squared = r2_score(y_test, predictions)
22
23 print(f'R-squared: {r_squared:.2f}')

```

R-squared: 0.97

In [4]:

```

1 from sklearn.metrics import r2_score
2
3 # Assuming you have observed actual and predicted values
4 actual_values = [10, 20, 30, 40, 50]
5 predicted_values = [12, 18, 28, 38, 52]
6
7 # Calculate the R-squared (R2) score
8 r2 = r2_score(actual_values, predicted_values)
9
10 print(f"R-squared (R2) Score: {r2}")
11

```

R-squared (R2) Score: 0.98

Mean absolute percentage error

In [38]:

```

1 import numpy as np
2 absolute_percentage_errors = np.abs((y_test - y_pred) / y_test)
3 mape = np.mean(absolute_percentage_errors) * 100 # Multiply by 100 to e
4 mape

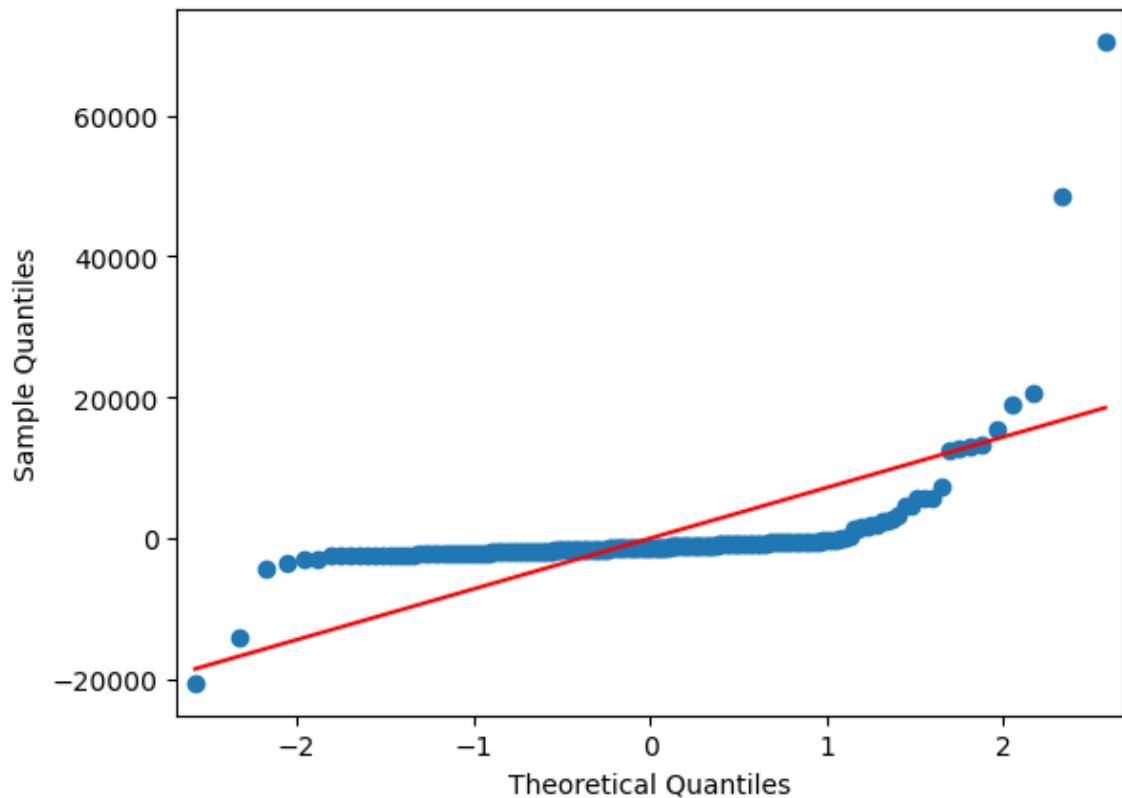
```

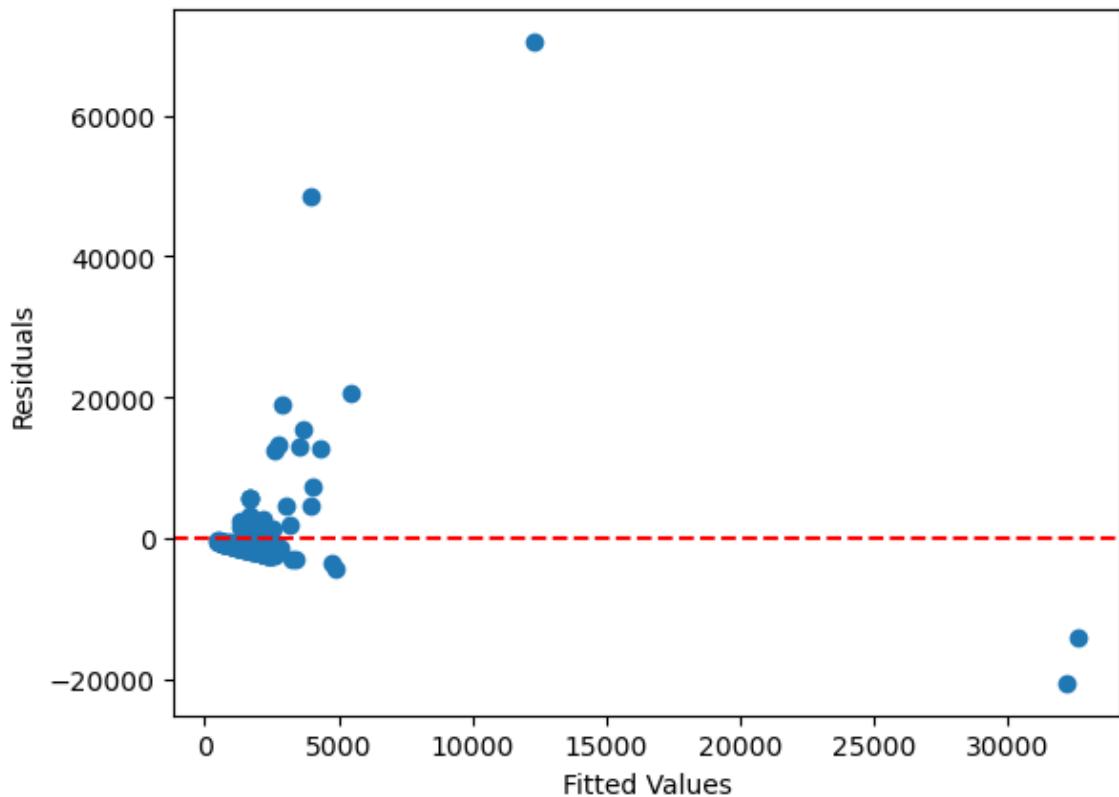
Out[38]: 5.55555555555338

Residual analysis

In [42]:

```
1 import pandas as pd
2 import statsmodels.api as sm
3 import matplotlib.pyplot as plt
4 df = pd.read_csv (r"C:\Users\Anusha V\Desktop\YouTube.csv")
5 X = df[['Likes', 'Rank']] # Features
6 y = df['Comments'] # Target variable
7 X = sm.add_constant(X)
8 model = sm.OLS(y, X).fit()
9 residuals = model.resid
10 sm.qqplot(residuals, line='s')
11 plt.show()
12 plt.scatter(model.fittedvalues, residuals)
13 plt.xlabel('Fitted Values')
14 plt.ylabel('Residuals')
15 plt.axhline(y=0, color='red', linestyle='--')
16 plt.show()
17
```





decision tree

In [2]:

```
1 # Import the necessary Libraries
2 from sklearn import tree
3
4 # Create a dataset for training the decision tree
5 # In this example, we have two features (X) and a target variable (y)
6 X = [[0, 0], [1, 1]]
7 y = [0, 1]
8
9 # Create a decision tree classifier
10 clf = tree.DecisionTreeClassifier()
11
12 # Train the classifier on the dataset
13 clf = clf.fit(X, y)
14
15 # Make predictions using the trained decision tree
16 predictions = clf.predict([[2, 2], [0.5, 0.5]])
17
18 # Print the predictions
19 print(predictions)
20
```

[1 0]

Decision Tree Classifier

```
In [3]:  
1 # Import the necessary libraries  
2 from sklearn import datasets  
3 from sklearn.model_selection import train_test_split  
4 from sklearn.tree import DecisionTreeClassifier  
5 from sklearn.metrics import accuracy_score  
6  
7 # Load a sample dataset (Iris dataset as an example)  
8 iris = datasets.load_iris()  
9 X = iris.data  
10 y = iris.target  
11  
12 # Split the data into training and testing sets  
13 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2)  
14  
15 # Create a Decision Tree Classifier  
16 clf = DecisionTreeClassifier()  
17  
18 # Fit the classifier to the training data  
19 clf.fit(X_train, y_train)  
20  
21 # Make predictions on the test data  
22 y_pred = clf.predict(X_test)  
23  
24 # Calculate the accuracy of the classifier  
25 accuracy = accuracy_score(y_test, y_pred)  
26 print(f"Accuracy: {accuracy * 100:.2f}%")  
27
```

Accuracy: 100.00%

Binary Classifier

In [15]:

```
1 # Import the necessary libraries
2 from sklearn.model_selection import train_test_split
3 from sklearn.linear_model import LogisticRegression
4 from sklearn.metrics import accuracy_score
5
6 # Generate a sample binary classification dataset (you can replace this
7 # Here, we'll create a simple dataset with two features (X) and binary
8 X = [[1.2, 2.3], [2.4, 3.5], [3.6, 4.7], [4.8, 5.9], [5.0, 6.1]]
9 y = [0, 1, 0, 1, 1]
10
11 # Split the data into training and testing sets
12 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2)
13
14 # Create a Binary Classifier (Logistic Regression is used in this example)
15 clf = LogisticRegression()
16
17 # Fit the classifier to the training data
18 clf.fit(X_train, y_train)
19
20 # Make predictions on the test data
21 y_pred = clf.predict(X_test)
22
23 # Calculate the accuracy of the classifier
24 accuracy = accuracy_score(y_test, y_pred)
25 print(f"Accuracy: {accuracy * 100:.2f}%")
26
```

Accuracy: 0.00%

multi class classifier

Support Vector Machine (SVM)

```
In [17]: 1 # Import the necessary libraries
2 from sklearn import datasets
3 from sklearn.model_selection import train_test_split
4 from sklearn.svm import SVC
5 from sklearn.metrics import accuracy_score
6
7 # Load a sample dataset (Iris dataset as an example)
8 iris = datasets.load_iris()
9 X = iris.data
10 y = iris.target
11
12 # Split the data into training and testing sets
13 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2)
14
15 # Create a Support Vector Machine (SVM) classifier
16 clf = SVC(kernel='linear', C=1)
17
18 # Fit the classifier to the training data
19 clf.fit(X_train, y_train)
20
21 # Make predictions on the test data
22 y_pred = clf.predict
23 y_pred
```

Out[17]: <bound method BaseSVC.predict of SVC(C=1, kernel='linear')>

```
In [18]: 1 # Import the necessary libraries
2 from sklearn import datasets
3 from sklearn.model_selection import train_test_split
4 from sklearn.svm import SVC
5 from sklearn.metrics import accuracy_score
6
7 # Load a sample dataset (Iris dataset as an example)
8 iris = datasets.load_iris()
9 X = iris.data
10 y = iris.target
11
12 # Split the data into training and testing sets
13 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2)
14
15 # Create a Support Vector Machine (SVM) classifier
16 clf = SVC(kernel='linear', C=1.0)
17
18 # Fit the classifier to the training data
19 clf.fit(X_train, y_train)
20
21 # Make predictions on the test data
22 y_pred = clf.predict(X_test)
23
24 # Calculate the accuracy of the classifier
25 accuracy = accuracy_score(y_test, y_pred)
26 print(f"Accuracy: {accuracy * 100:.2f}%")
```

Accuracy: 100.00%

multi-class classification

Loading [MathJax]/extensions/MathML/content-mathml.js

In [19]:

```

1 # Import necessary libraries
2 from sklearn.datasets import load_iris
3 from sklearn.model_selection import train_test_split
4 from sklearn.ensemble import RandomForestClassifier
5 from sklearn.metrics import accuracy_score, classification_report
6
7 # Load a sample dataset (Iris dataset as an example)
8 iris = load_iris()
9 X = iris.data
10 y = iris.target
11
12 # Split the data into training and testing sets
13 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2)
14
15 # Create a Random Forest Classifier
16 clf = RandomForestClassifier(n_estimators=100, random_state=42)
17
18 # Fit the classifier to the training data
19 clf.fit(X_train, y_train)
20
21 # Make predictions on the test data
22 y_pred = clf.predict(X_test)
23
24 # Calculate and print the accuracy of the classifier
25 accuracy = accuracy_score(y_test, y_pred)
26 print(f"Accuracy: {accuracy * 100:.2f}%")
27
28 # Print classification report (precision, recall, F1-score for each class)
29 report = classification_report(y_test, y_pred, target_names=iris.target_names)
30 print(report)
31

```

Accuracy: 100.00%

| | precision | recall | f1-score | support |
|--------------|-----------|--------|----------|---------|
| setosa | 1.00 | 1.00 | 1.00 | 10 |
| versicolor | 1.00 | 1.00 | 1.00 | 9 |
| virginica | 1.00 | 1.00 | 1.00 | 11 |
| accuracy | | | 1.00 | 30 |
| macro avg | 1.00 | 1.00 | 1.00 | 30 |
| weighted avg | 1.00 | 1.00 | 1.00 | 30 |

multi-class classification

```
In [22]: 1 # Import the necessary libraries
2 from sklearn import datasets
3 from sklearn.model_selection import train_test_split
4 from sklearn.svm import SVC
5 from sklearn.metrics import accuracy_score
6
7 # Load a sample dataset (Iris dataset as an example)
8 iris = datasets.load_iris()
9 X = iris.data
10 y = iris.target
11
12 # Split the data into training and testing sets
13 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2)
14
15 # Create a Support Vector Machine (SVM) classifier
16 clf = SVC(kernel='linear', C=1.0)
17
18 # Fit the classifier to the training data
19 clf.fit(X_train, y_train)
20
21 # Make predictions on the test data
22 y_pred = clf.predict(X_test)
23
24 # Calculate the accuracy of the classifier
25 accuracy = accuracy_score(y_test, y_pred)
26 print(f"Accuracy: {accuracy * 100:.2f}%")
27
```

Accuracy: 100.00%

Confusion Matrix:

In [27]:

```
1 import numpy as np
2
3 # Given confusion matrix
4 confusion_matrix = np.array([
5     [13, 45],
6     [0, 15]
7 ])
8
9 # (i) Accuracy
10 total_samples = np.sum(confusion_matrix)
11 correct_predictions = np.trace(confusion_matrix)
12 accuracy = correct_predictions / total_samples
13 print(f"(i) Accuracy: {accuracy:.2f}")
14
15 # (ii) Precision for Class 1
16 precision_class_1 = confusion_matrix[1, 1] / (confusion_matrix[0, 1] +
17 print(f"(ii) Precision for Class 1: {precision_class_1:.2f}")
18
19 # (iii) Recall for Class 1
20 recall_class_1 = confusion_matrix[1, 1] / (confusion_matrix[1, 0] + con
21 print(f"(iii) Recall for Class 1: {recall_class_1:.2f}")
22
23 # (iv) Specificity for Class 0
24 specificity_class_0 = confusion_matrix[0, 0] / (confusion_matrix[0, 0]
25 print(f"(iv) Specificity for Class 0: {specificity_class_0:.2f}")
26
27 # (v) F1-Score for Class 1
28 f1_score_class_1 = 2 * (precision_class_1 * recall_class_1) / (precision_
29 print(f"(v) F1-Score for Class 1: {f1_score_class_1:.2f}")
30
```

(i) Accuracy: 0.38
(ii) Precision for Class 1: 0.25
(iii) Recall for Class 1: 1.00
(iv) Specificity for Class 0: 0.22
(v) F1-Score for Class 1: 0.40

In [22]:

```

1 import numpy as np
2
3 # Given confusion matrix
4 confusion_matrix = np.array([[0, 0, 0],
5                             [13, 45, 8],
6                             [0, 15, 32]])
7
8 # (i) Accuracy
9 total_samples = np.sum(confusion_matrix)
10 correct_predictions = np.trace(confusion_matrix)
11 accuracy = correct_predictions / total_samples
12 print(f"(i) Accuracy: {accuracy:.2f}")
13
14 # (ii) Precision for Class 1
15 precision_class_1 = confusion_matrix[1, 1] / (confusion_matrix[1, 0] +
16 print(f"(ii) Precision for Class 1: {precision_class_1:.2f}")
17
18 # (iii) Recall for Class 1
19 recall_class_1 = confusion_matrix[1, 1] / (confusion_matrix[0, 1] + con
20 print(f"(iii) Recall for Class 1: {recall_class_1:.2f}")
21
22 # (iv) Specificity for Class 0
23 specificity_class_0 = (confusion_matrix[0, 0] + confusion_matrix[2, 0])
24 print(f"(iv) Specificity for Class 0: {specificity_class_0:.2f}")
25
26
27 # (v) F1-Score for Class 1
28 f1_score_class_1 = 2 * (precision_class_1 * recall_class_1) / (precision_
29 print(f"(v) F1-Score for Class 1: {f1_score_class_1:.2f}")
30
31
32

```



(i) Accuracy: 0.68
(ii) Precision for Class 1: 0.68
(iii) Recall for Class 1: 0.75
(iv) Specificity for Class 0: 0.00
(v) F1-Score for Class 1: 0.71

In [23]:

```

1 # Import the necessary Libraries
2 from sklearn.metrics import confusion_matrix
3 import numpy as np
4
5 # Actual and predicted Labels (replace these with your actual data)
6 actual = np.array([1, 0, 1, 1, 0, 1, 0, 1, 0, 0])
7 predicted = np.array([1, 0, 1, 0, 1, 1, 0, 1, 0, 1])
8
9 # Create the confusion matrix
10 cm = confusion_matrix(actual, predicted)
11
12 # Print the confusion matrix
13 print("Confusion Matrix:")
14 print(cm)
15

```

Confusion Matrix:

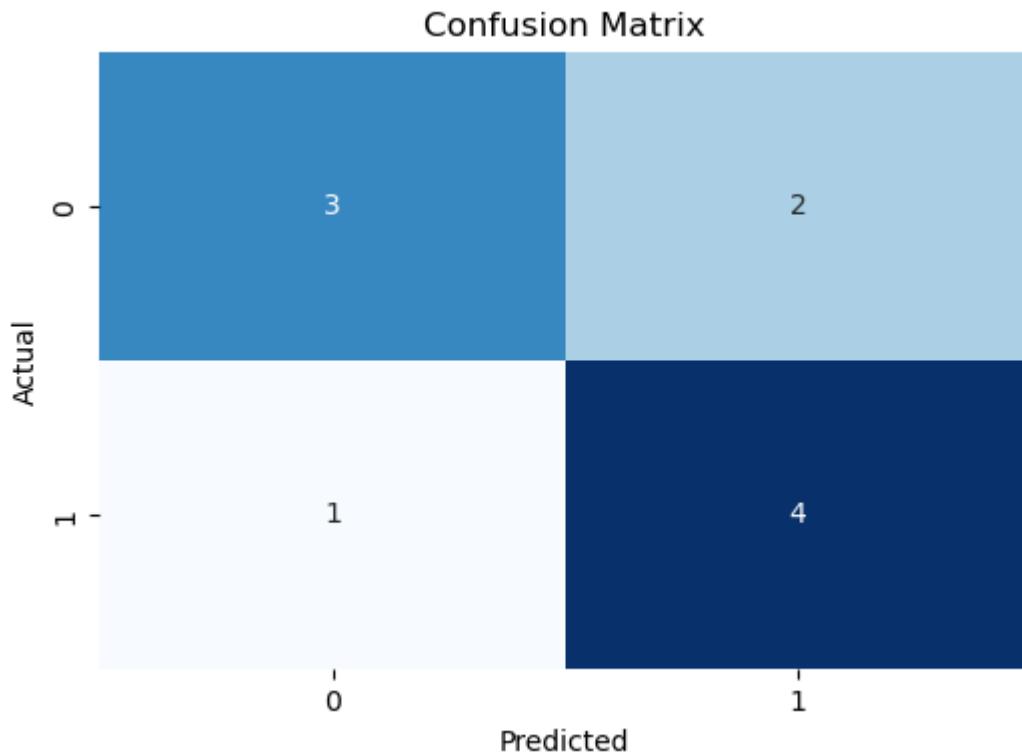
[[3 2]

[1 4 1]

Loading [MathJax]/extensions/MathML/content-mathml.js

In [20]:

```
1 from sklearn.metrics import confusion_matrix
2 import matplotlib.pyplot as plt
3 import seaborn as sns
4
5 # Example ground truth (actual) Labels and predicted labels
6 true_labels = [1, 0, 1, 1, 0, 1, 0, 0, 1, 0]
7 predicted_labels = [1, 1, 0, 1, 0, 1, 0, 0, 1, 1]
8
9 # Create a confusion matrix
10 cm = confusion_matrix(true_labels, predicted_labels)
11
12 # Display the confusion matrix as a heatmap
13 plt.figure(figsize=(6, 4))
14 sns.heatmap(cm, annot=True, fmt="d", cmap="Blues", cbar=False)
15 plt.xlabel('Predicted')
16 plt.ylabel('Actual')
17 plt.title('Confusion Matrix')
18 plt.show()
19
```



Creating a Confusion Matrix:

```
In [24]: 1 from sklearn.metrics import confusion_matrix
2
3 # True Labels and predicted Labels
4 true_labels = [0, 1, 1, 0, 2, 1, 2, 0]
5 predicted_labels = [0, 1, 0, 1, 2, 2, 2, 0]
6
7 # Generate the confusion matrix
8 cm = confusion_matrix(true_labels, predicted_labels)
9 print(cm)
10
```

```
[[2 1 0]
 [1 1 1]
 [0 0 2]]
```

Evaluating Performance Metrics:

```
In [27]: 1 from sklearn.metrics import classification_report
2
3 # Assuming 'true_labels' and 'predicted_labels' are defined
4 report = classification_report(true_labels, predicted_labels)
5 print(report)
6
```

| | precision | recall | f1-score | support |
|--------------|-----------|--------|----------|---------|
| 0 | 0.67 | 0.67 | 0.67 | 3 |
| 1 | 0.50 | 0.33 | 0.40 | 3 |
| 2 | 0.67 | 1.00 | 0.80 | 2 |
| accuracy | | | 0.62 | 8 |
| macro avg | 0.61 | 0.67 | 0.62 | 8 |
| weighted avg | 0.60 | 0.62 | 0.60 | 8 |

Custom Confusion Matrix Analysis:

```
In [31]: 1 tp = cm[1, 1] # True Positives
2 fn = cm[1, 0] # False Negatives
3 fp = cm[0, 1] # False Positives
4 tn = cm[0, 0] # True Negatives
5
6 sensitivity = tp / (tp + fn) # Recall
7 specificity = tn / (tn + fp)
```

```
In [32]: 1 sensitivity
```

Out[32]: 0.5

```
In [33]: 1 specificity
```

Out[33]: 0.6666666666666666

Multi-Class Confusion Matrix:

In [34]:

```

1 from sklearn.metrics import confusion_matrix
2
3 # Assuming 'true_labels' and 'predicted_labels' for multi-class
4 cm = confusion_matrix(true_labels, predicted_labels)
5 print(cm)
6

```

```
[[2 1 0]
 [1 1 1]
 [0 0 2]]
```

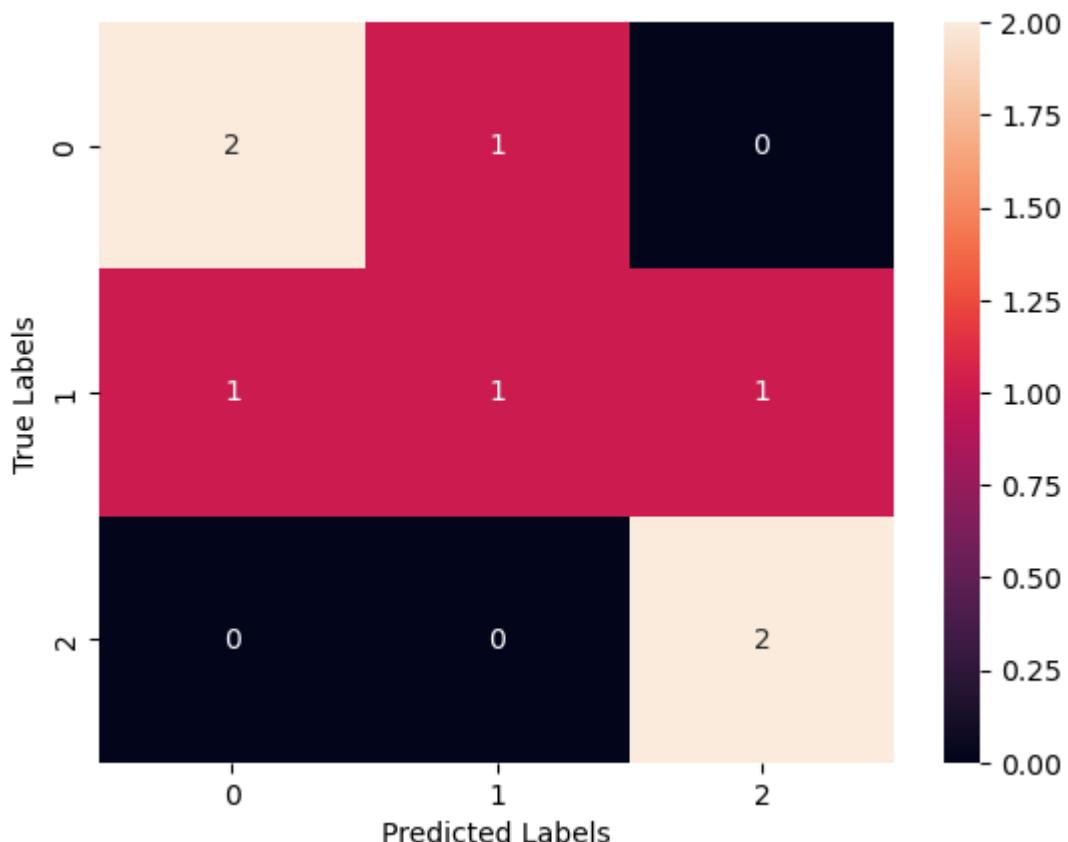
Visualizing a Confusion Matrix:

In [35]:

```

1 import seaborn as sns
2 import matplotlib.pyplot as plt
3
4 # Assuming 'cm' is the confusion matrix from previous code
5 sns.heatmap(cm, annot=True, fmt="d")
6 plt.xlabel("Predicted Labels")
7 plt.ylabel("True Labels")
8 plt.show()
9

```



F1 score

Loading [MathJax]/extensions/MathML/content-mathml.js

```
In [50]: 1 from sklearn.metrics import f1_score
2
3 # Example true labels and predicted labels
4 true_labels = [1, 0, 1, 1, 0, 1, 0, 1]
5 predicted_labels = [1, 0, 0, 1, 0, 1, 1, 1]
6
7 # Calculate the F1 score
8 f1 = f1_score(true_labels, predicted_labels)
9 print(f"F1 Score: {f1:.2f}")
10
```

F1 Score: 0.80

DecisionTree Classifier hyperparameters

```
In [3]: 1 # Import necessary Libraries
2 from sklearn.tree import DecisionTreeClassifier
3 from sklearn.model_selection import GridSearchCV
4 from sklearn.datasets import load_iris
5
6 # Load a sample dataset (you can replace this with your own dataset)
7 data = load_iris()
8 X, y = data.data, data.target
9
10 # Create a DecisionTreeClassifier
11 clf = DecisionTreeClassifier()
12
13 # Define a dictionary of hyperparameters and their possible values
14 param_grid = {
15     'criterion': ['gini', 'entropy'],
16     'max_depth': [None, 10, 20, 30],
17     'min_samples_split': [2, 5, 10],
18     'min_samples_leaf': [1, 2, 4]
19 }
20
21 # Create a GridSearchCV object to find the best hyperparameters
22 grid_search = GridSearchCV(clf, param_grid, cv=5, n_jobs=-1)
23
24 # Fit the grid search to the data
25 grid_search.fit(X, y)
26
27 # Print the best hyperparameters and the corresponding accuracy score
28 print("Best hyperparameters: ", grid_search.best_params_)
29 print("Best accuracy score: ", grid_search.best_score_)
30
```

Best hyperparameters: {'criterion': 'gini', 'max_depth': None, 'min_samples_leaf': 1, 'min_samples_split': 5}
 Best accuracy score: 0.9666666666666668

Decision Tree Visualization

In [4]:

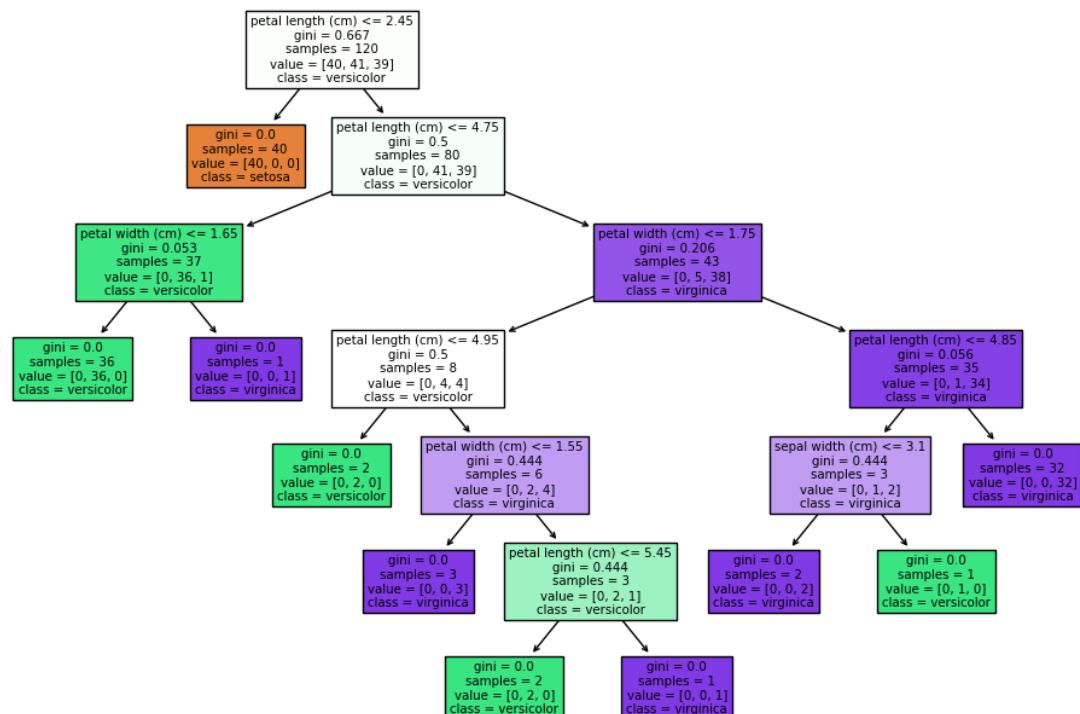
```

1 import numpy as np
2 import matplotlib.pyplot as plt
3 from sklearn.datasets import load_iris
4 from sklearn.tree import DecisionTreeClassifier
5 from sklearn.model_selection import train_test_split
6 from sklearn.metrics import accuracy_score
7
8 data = load_iris()
9 X = data.data
10 y = data.target
11
12 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2)
13
14 # Create a Decision Tree classifier
15 clf = DecisionTreeClassifier(criterion='gini', max_depth=None, random_state=42)
16
17 # Train the classifier
18 clf.fit(X_train, y_train)
19
20 y_pred = clf.predict(X_test)
21
22 # Calculate the accuracy of the model
23 accuracy = accuracy_score(y_test, y_pred)
24 print(f"Accuracy: {accuracy:.2f}")
25
26 from sklearn.tree import plot_tree
27
28 # Plot the decision tree
29 plt.figure(figsize=(12, 8))
30 plot_tree(clf, filled=True, feature_names=data.feature_names, class_names=['Setosa', 'Versicolor', 'Virginica'])
31 plt.title("Decision Tree Visualization")
32 plt.show()
33

```

Accuracy: 1.00

Decision Tree Visualization



Loading [MathJax]/extensions/MathML/content-mathml.js

binary classification

In [42]:

```

1 from sklearn.model_selection import train_test_split
2 from sklearn.linear_model import LogisticRegression
3
4 X = [[1.2, 2.3], [2.4, 3.5], [3.6, 4.7], [4.8, 5.9], [5.0, 6.1]]
5 y = [0, 1, 0, 1, 1]
6 # Load and preprocess your data
7 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2
8
9 # Create and train a binary classification model (e.g., logistic regression)
10 model = LogisticRegression()
11 model.fit(X_train, y_train)
12
13
14 # Make predictions
15 predictions = model.predict(X_test)
16 predictions

```

Out[42]: array([1])

In [39]:

```

1 # Import necessary Libraries
2 import numpy as np
3 from sklearn.model_selection import train_test_split
4 from sklearn.linear_model import LogisticRegression
5 from sklearn.metrics import accuracy_score
6
7 # Create synthetic data for binary classification
8 X = np.array([[1, 2], [2, 3], [3, 4], [4, 5], [2, 2], [3, 3], [4, 4], [5, 5], [1, 1], [0, 0], [1, 0], [0, 1]])
9 y = np.array([0, 0, 0, 0, 1, 1, 1, 1, 0, 1, 0, 1])
10
11 # Split the data into training and testing sets
12 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2)
13
14 # Create a logistic regression model
15 model = LogisticRegression()
16
17 # Train the model on the training data
18 model.fit(X_train, y_train)
19
20 # Make predictions on the test data
21 y_pred = model.predict(X_test)
22
23
24
25 # Calculate the accuracy of the model
26 accuracy = accuracy_score(y_test, y_pred)
27 print(f"Accuracy: {accuracy * 100:.2f}%")
28
29 predictions = model.predict(X_test)
30 predictions
31

```

Accuracy: 100.00%

Out[39]: array([0, 1])

Loading [MathJax]/extensions/MathML/content-mathml.js

multi-class classification

In [8]:

```

1 import numpy as np
2 from sklearn import datasets
3 from sklearn.model_selection import train_test_split
4 from sklearn.svm import SVC
5 from sklearn.metrics import accuracy_score, classification_report
6
7 # Load the Iris dataset
8 iris = datasets.load_iris()
9 X = iris.data
10 y = iris.target
11
12 # Split the data into training and testing sets
13 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3)
14
15 # Initialize the Support Vector Machine classifier
16 svm_classifier = SVC(kernel='linear', C=1)
17
18 # Fit the classifier to the training data
19 svm_classifier.fit(X_train, y_train)
20
21 # Make predictions on the test data
22 y_pred = svm_classifier.predict(X_test)
23
24 # Evaluate the classifier's performance
25 accuracy = accuracy_score(y_test, y_pred)
26 report = classification_report(y_test, y_pred, target_names=iris.target)
27 -/
28 # Print the results
29 print("Accuracy:", accuracy)
30 print("Classification Report:\n", report)
31

```

Accuracy: 1.0

Classification Report:

| | precision | recall | f1-score | support |
|--------------|-----------|--------|----------|---------|
| setosa | 1.00 | 1.00 | 1.00 | 19 |
| versicolor | 1.00 | 1.00 | 1.00 | 13 |
| virginica | 1.00 | 1.00 | 1.00 | 13 |
| accuracy | | | 1.00 | 45 |
| macro avg | 1.00 | 1.00 | 1.00 | 45 |
| weighted avg | 1.00 | 1.00 | 1.00 | 45 |

multi-label classification

In [48]:

```

1 from sklearn.model_selection import train_test_split
2 from sklearn.preprocessing import MultiLabelBinarizer
3 from sklearn.multiclass import OneVsRestClassifier
4 from sklearn.svm import SVC
5 from sklearn.metrics import accuracy_score, hamming_loss
6
7 data = [
8     ("This is a cat", {"animal": 1, "pet": 1}),
9     ("This is a dog", {"animal": 1, "pet": 1}),
10    ("This is a fish", {"animal": 1, "pet": 0}),
11    ("This is a book", {"animal": 0, "pet": 0}),
12    ("This is a rabbit", {"animal": 1, "pet": 1}),
13]
14 # Split data into X (input) and y (labels)
15 X = [item[0] for item in data]
16 y = [item[1] for item in data]
17
18 # Convert labels to a binary matrix
19 mlb = MultiLabelBinarizer()
20 y_bin = mlb.fit_transform(y)
21
22 # Split data into training and testing sets
23 X_train, X_test, y_train, y_test = train_test_split(X, y_bin, test_size=0.2)
24
25 # Create a classifier (e.g., Support Vector Machine)
26 classifier = OneVsRestClassifier(SVC(kernel="linear"))
27
28 # Fit the model to the training data
29 classifier.fit(X_train, y_train)
30
31
32 # Predict labels for the test data
33 y_pred = classifier.predict(X_test)
34
35 # Inverse transform to get labels
36 y_pred_labels = mlb.inverse_transform(y_pred)
37
38 # Evaluate the model
39 accuracy = accuracy_score(y_test, y_pred)
40 hamming = hamming_loss(y_test, y_pred)
41
42 print("Accuracy: ", accuracy)
43 print("Hamming Loss: ", hamming)
44

```

Accuracy: 1.0

Hamming Loss: 0.0

C:\Users\Anusha V\anaconda3\lib\site-packages\sklearn\multiclass.py:79: UserWarning: Label 0 is present in all training examples.
 warnings.warn(
 C:\Users\Anusha V\anaconda3\lib\site-packages\sklearn\multiclass.py:79: UserWarning: Label 1 is present in all training examples.
 warnings.warn(

```
In [8]: 1 import sys
2 from sklearn.feature_extraction.text import TfidfVectorizer
3 from sklearn.multiclass import OneVsRestClassifier
4 from sklearn.svm import SVC
5 from sklearn.preprocessing import MultiLabelBinarizer
6
7 # Example data for training (adjust as needed)
8 X_train = ["text1", "text2", "text3"]
9 y_train = [["label1", "label2"], ["label2", "label3"], ["label1"]]
10
11 # Vectorize the text
12 vectorizer = TfidfVectorizer()
13 X_train = vectorizer.fit_transform(X_train)
14
15 # Binarize the labels
16 mlb = MultiLabelBinarizer()
17 y_train = mlb.fit_transform(y_train)
18
19 # Train a multi-label classifier (e.g., SVM)
20 classifier = OneVsRestClassifier(SVC())
21 classifier.fit(X_train, y_train)
22
23
24 # Output the result
25 print("Predicted Labels: ", predicted_labels)
26
```

Predicted Labels: classifier.fit

imbalanced classification

```
In [17]: 1 pip install -U imbalanced-learn
2
```

```
Collecting imbalanced-learn
  Downloading imbalanced_learn-0.11.0-py3-none-any.whl (235 kB)
----- 235.6/235.6 kB 720.1 kB/s eta 0:00:00
Requirement already satisfied: threadpoolctl>=2.0.0 in c:\users\anusha v\anaconda3\lib\site-packages (from imbalanced-learn) (2.2.0)
Requirement already satisfied: numpy>=1.17.3 in c:\users\anusha v\anaconda3\lib\site-packages (from imbalanced-learn) (1.21.5)
Requirement already satisfied: scikit-learn>=1.0.2 in c:\users\anusha v\anaconda3\lib\site-packages (from imbalanced-learn) (1.0.2)
Requirement already satisfied: scipy>=1.5.0 in c:\users\anusha v\anaconda3\lib\site-packages (from imbalanced-learn) (1.9.1)
Collecting joblib>=1.1.1
  Downloading joblib-1.3.2-py3-none-any.whl (302 kB)
----- 302.2/302.2 kB 1.3 MB/s eta 0:00:00
Installing collected packages: joblib, imbalanced-learn
Attempting uninstall: joblib
  Found existing installation: joblib 1.1.0
  Uninstalling joblib-1.1.0:
    Successfully uninstalled joblib-1.1.0
Successfully installed imbalanced-learn-0.11.0 joblib-1.3.2
Note: you may need to restart the kernel to use updated packages.
```

```
In [19]: 1 from imblearn.over_sampling import RandomOverSampler  
2 from sklearn.ensemble import RandomForestClassifier  
3  
4 # Apply random oversampling to balance the data  
5 oversampler = RandomOverSampler()  
6 X_resampled, y_resampled = oversampler.fit_resample(X, y)  
7  
8 # Create a Random Forest classifier  
9 model = RandomForestClassifier()  
10  
11 # Train the model on the resampled data  
12 model.fit(X_resampled, y_resampled)  
13  
14 # Make predictions on new data  
15 y_pred = model.predict(X_test)  
16 y_pred
```

```
Out[19]: array([1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0,  
    1, 1, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,  
    1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,  
    1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,  
    1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,  
    1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,  
    1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,  
    1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,  
    1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,  
    1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,  
    1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,  
    1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1])
```

Housing Price dataset

In [31]:

```

1 from sklearn.datasets import load_boston
2 import pandas as pd
3
4 # Load the Boston Housing Price dataset
5 boston = load_boston()
6
7 # Create a pandas DataFrame to explore the data
8 boston_df = pd.DataFrame(data=boston.data, columns=boston.feature_names)
9 boston_df['target'] = boston.target
10
11 # Display the first few rows of the dataset
12 print(boston_df.head())
13
14 # Split the data into features (X) and target variable (y)
15 X = boston.data
16 y = boston.target
17
18 # You can use X and y for machine Learning tasks such as regression:
19
20 # Example of training a Linear regression model
21 from sklearn.linear_model import LinearRegression
22 from sklearn.model_selection import train_test_split
23 from sklearn.metrics import mean_squared_error, r2_score
24
25 # Split the data into a training and test set
26 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2)
27
28 # Create and fit a Linear regression model
29 model = LinearRegression()
30 model.fit(X_train, y_train)
31
32 # Make predictions on the test set
33 y_pred = model.predict(X_test)
34
35 # Evaluate the model
36 mse = mean_squared_error(y_test, y_pred)
37 r2 = r2_score(y_test, y_pred)
38 print(f"Mean Squared Error: {mse:.2f}")
39 print(f"R-squared (R2): {r2:.2f}")
40

```

| | CRIM | ZN | INDUS | CHAS | NOX | RM | AGE | DIS | RAD | TAX | \ |
|---|---------|------|-------|------|-------|-------|------|--------|-----|-------|---|
| 0 | 0.00632 | 18.0 | 2.31 | 0.0 | 0.538 | 6.575 | 65.2 | 4.0900 | 1.0 | 296.0 | |
| 1 | 0.02731 | 0.0 | 7.07 | 0.0 | 0.469 | 6.421 | 78.9 | 4.9671 | 2.0 | 242.0 | |
| 2 | 0.02729 | 0.0 | 7.07 | 0.0 | 0.469 | 7.185 | 61.1 | 4.9671 | 2.0 | 242.0 | |
| 3 | 0.03237 | 0.0 | 2.18 | 0.0 | 0.458 | 6.998 | 45.8 | 6.0622 | 3.0 | 222.0 | |
| 4 | 0.06905 | 0.0 | 2.18 | 0.0 | 0.458 | 7.147 | 54.2 | 6.0622 | 3.0 | 222.0 | |

| | PTRATIO | B | LSTAT | target |
|---|---------|--------|-------|--------|
| 0 | 15.3 | 396.90 | 4.98 | 24.0 |
| 1 | 17.8 | 396.90 | 9.14 | 21.6 |
| 2 | 17.8 | 392.83 | 4.03 | 34.7 |
| 3 | 18.7 | 394.63 | 2.94 | 33.4 |
| 4 | 18.7 | 396.90 | 5.33 | 36.2 |

Mean Squared Error: 24.29

R-squared (R2): 0.67

```
C:\Users\Anusha V\anaconda3\lib\site-packages\sklearn\utils\deprecation.py:87: FutureWarning: Function load_boston is deprecated; `load_boston` is deprecated in 1.0 and will be removed in 1.2.
```

The Boston housing prices dataset has an ethical problem. You can refer to the documentation of this function for further details.

The scikit-learn maintainers therefore strongly discourage the use of this dataset unless the purpose of the code is to study and educate about ethical issues in data science and machine learning.

In this special case, you can fetch the dataset from the original source::

```
import pandas as pd
import numpy as np

data_url = "http://lib.stat.cmu.edu/datasets/boston"
raw_df = pd.read_csv(data_url, sep="\s+", skiprows=22, header=None)
data = np.hstack([raw_df.values[::2, :], raw_df.values[1::2, :2]])
target = raw_df.values[1::2, 2]
```

Alternative datasets include the California housing dataset (i.e. `:func:`~sklearn.datasets.fetch_california_housing``) and the Ames housing dataset. You can load the datasets as follows::

```
from sklearn.datasets import fetch_california_housing
housing = fetch_california_housing()

for the California housing dataset and::

from sklearn.datasets import fetch_openml
housing = fetch_openml(name="house_prices", as_frame=True)

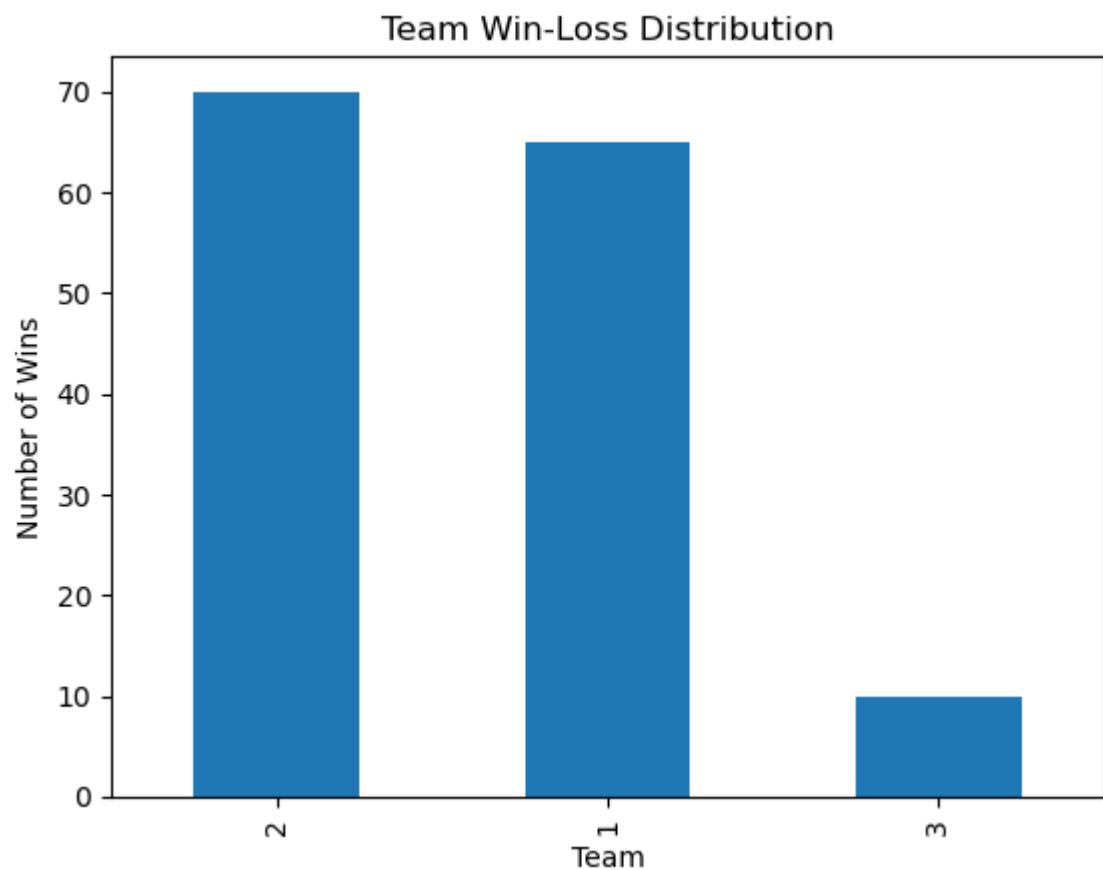
for the Ames housing dataset.
```

```
warnings.warn(msg, category=FutureWarning)
```

cricket match data

In [36]:

```
1 import pandas as pd
2 import matplotlib.pyplot as plt
3
4 # Load historical cricket match data
5 cricket_data = pd.read_csv(r"C:\Users\Anusha V\Downloads\DATA (1).csv")
6
7 # Perform EDA
8 # Example: Plot a bar chart to visualize the win-loss distribution of t
9 team_wins = cricket_data['wins'].value_counts()
10 team_wins.plot(kind='bar', title='Team Win-Loss Distribution')
11 plt.xlabel('Team')
12 plt.ylabel('Number of Wins')
13 plt.show()
14
15
```



In []:

1

gini

```
In [6]: 1 def gini_index(probabilities):
2     """
3         Calculate the Gini index for a set of class probabilities.
4
5         :param probabilities: A list of class probabilities (e.g., [0.2, 0.
6         :return: The Gini index for the given probabilities.
7         """
8         gini = 1.0 # Initialize the Gini index
9
10        for p in probabilities:
11            gini -= p ** 2
12
13        return gini
14
15    # Example usage:
16    class_probabilities = [0.2, 0.3, 0.5]
17    gini = gini_index(class_probabilities)
18    print(f"Gini Index: {gini}")
19
```

Gini Index: 0.62

Entropy

```
In [1]: 1 import math
2
3 def calculate_entropy(data):
4     """
5         Calculate the entropy of a set of data.
6
7     Parameters:
8         data (list): A list of labels or classes.
9
10    Returns:
11        float: The entropy of the data.
12    """
13    total_samples = len(data)
14    class_counts = {}
15
16    # Count the occurrences of each class in the data
17    for label in data:
18        if label in class_counts:
19            class_counts[label] += 1
20        else:
21            class_counts[label] = 1
22
23    entropy = 0.0
24
25    # Calculate entropy using the formula: -p(x) * Log2(p(x))
26    for label in class_counts:
27        probability = class_counts[label] / total_samples
28        entropy -= probability * math.log2(probability)
29
30    return entropy
31
32 # Example usage
33 data = ['A', 'B', 'A', 'C', 'B', 'A', 'A', 'B']
34 entropy = calculate_entropy(data)
35 print(f'Entropy of the data: {entropy}')
36
```

Entropy of the data: 1.4056390622295665

StandardScaler(Standardization)

```
In [25]: 1 x = df.iloc[:, [7, 6]]
2 y = df.iloc[:, 5]
```

```
In [26]: 1 x.head()
```

| | Comments | Likes |
|---|----------|---------|
| 0 | 78 | 2700 |
| 1 | 18500 | 5300000 |
| 2 | 0 | 24700 |
| 3 | 9 | 166 |
| 4 | 0 | 12400 |

Loading [MathJax]/extensions/MathML/content-mathml.js

```
In [27]: 1 from sklearn.model_selection import train_test_split
2 from sklearn.preprocessing import StandardScaler
```

```
In [28]: 1 x_train, x_test, y_train, y_test = train_test_split(x, y, random_state
```

```
In [29]: 1 x_train.head()
```

Out[29]:

| | Comments | Likes |
|-----|----------|--------|
| 84 | 0 | 10800 |
| 16 | 7400 | 180300 |
| 19 | 4 | 1100 |
| 2 | 0 | 24700 |
| 142 | 108 | 7600 |

```
In [30]: 1 scaler = StandardScaler(). fit(x_train)
```

```
In [31]: 1 scaler
```

Out[31]: StandardScaler()

```
In [32]: 1 scaler.mean_
```

Out[32]: array([2221.26174497, 89513.84563758])

```
In [33]: 1 scaler.scale_
```

Out[33]: array([-8663.65510725, 436331.85622113])

```
In [37]: 1 scaler.transform(x_train)
```

Out[37]: array([[-2.56388524e-01, -1.80399035e-01],
 [5.97754434e-01, 2.08066757e-01],
 [-2.55926825e-01, -2.02629820e-01],
 [-2.56388524e-01, -1.48542548e-01],
 [-2.43922654e-01, -1.87732902e-01],
 [-2.49463041e-01, -2.01483904e-01],
 [-2.53964605e-01, -2.02941510e-01],
 [-2.53964605e-01, -2.03759694e-01],
 [-2.55349701e-01, -2.04621424e-01],
 [-2.55349701e-01, -2.04667260e-01],
 [-2.56388524e-01, -2.05150837e-01],
 [3.09192624e-01, 2.47027928e-01],
 [-2.56388524e-01, -1.78336384e-01],
 [-2.33303579e-01, -1.90712286e-01],
 [-2.56388524e-01, -1.99192070e-01],
 [-2.55349701e-01, -2.04596214e-01],
 [1.82225427e-01, -8.51151614e-03],
 [-2.15181897e-01, -1.42125414e-01],
 [-2.50155589e-01, -2.03335705e-01],

```
In [39]: 1 x_train_scales = scaler.transform(x_train)
2 x_train_scales
```

```
Out[39]: array([[-2.56388524e-01, -1.80399035e-01],
   [ 5.97754434e-01,  2.08066757e-01],
   [-2.55926825e-01, -2.02629820e-01],
   [-2.56388524e-01, -1.48542548e-01],
   [-2.43922654e-01, -1.87732902e-01],
   [-2.49463041e-01, -2.01483904e-01],
   [-2.53964605e-01, -2.02941510e-01],
   [-2.53964605e-01, -2.03759694e-01],
   [-2.55349701e-01, -2.04621424e-01],
   [-2.55349701e-01, -2.04667260e-01],
   [-2.56388524e-01, -2.05150837e-01],
   [ 3.09192624e-01,  2.47027928e-01],
   [-2.56388524e-01, -1.78336384e-01],
   [-2.33303579e-01, -1.90712286e-01],
   [-2.56388524e-01, -1.99192070e-01],
   [-2.55349701e-01, -2.04596214e-01],
   [ 1.82225427e-01, -8.51151614e-03],
   [-2.15181897e-01, -1.42125414e-01],
   [-2.50155589e-01, -2.03335705e-01],
   [-2.50155589e-01, -2.03335705e-01]]
```

```
In [41]: 1 print(x_train_scales.mean(axis=0))
[ 0.0000000e+00 -2.9804645e-18]
```

```
In [43]: 1 print(x_train_scales.std(axis=0))
[1. 1.]
```

```
In [48]: 1 scaler = StandardScaler(). fit(x_test)
2 scaler.mean_
```

```
Out[48]: array([-1724.14, 150251.96])
```

```
In [11]: 1 import pandas as pd
          2 data = pd.read_csv (r"C:\Users\Anusha V\Downloads\archive (8)\APY.csv")
          3 data.head(10)
```

Out[11]:

| | State | District | Crop | Crop_Year | Season | Area | Production | Yield |
|---|----------------------------|----------|----------|-----------|------------|--------|------------|-------|
| 0 | Andaman and Nicobar Island | NICOBARS | Arecanut | 2007 | Kharif | 2439.6 | 3415.0 | 1.40 |
| 1 | Andaman and Nicobar Island | NICOBARS | Arecanut | 2007 | Rabi | 1626.4 | 2277.0 | 1.40 |
| 2 | Andaman and Nicobar Island | NICOBARS | Arecanut | 2008 | Autumn | 4147.0 | 3060.0 | 0.74 |
| 3 | Andaman and Nicobar Island | NICOBARS | Arecanut | 2008 | Summer | 4147.0 | 2660.0 | 0.64 |
| 4 | Andaman and Nicobar Island | NICOBARS | Arecanut | 2009 | Autumn | 4153.0 | 3120.0 | 0.75 |
| 5 | Andaman and Nicobar Island | NICOBARS | Arecanut | 2009 | Summer | 4153.0 | 2080.0 | 0.50 |
| 6 | Andaman and Nicobar Island | NICOBARS | Arecanut | 2000 | Kharif | 1254.0 | 2000.0 | 1.59 |
| 7 | Andaman and Nicobar Island | NICOBARS | Arecanut | 2001 | Kharif | 1254.0 | 2061.0 | 1.64 |
| 8 | Andaman and Nicobar Island | NICOBARS | Arecanut | 2002 | Whole Year | 1258.0 | 2083.0 | 1.66 |
| 9 | Andaman and Nicobar Island | NICOBARS | Arecanut | 2003 | Whole Year | 1261.0 | 1525.0 | 1.21 |

```
In [14]: 1 X = data[['Crop_Year', 'Season', 'Yield']]
          2 y = data['Crop_Year']
```

k mean

```
In [15]: 1 import numpy as np
          2 from sklearn.neighbors import KNeighborsClassifier
          3 from sklearn.model_selection import train_test_split
          4 from sklearn.metrics import accuracy_score
          5 from sklearn.cluster import KMeans
          6
```

```
In [17]: 1 from sklearn.datasets import load_iris
          2 iris = load_iris()
          3 X = iris.data
          4 y = iris.target
```

```
In [18]: 1 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2)
```

```
In [19]: 1 k = 3 # Set the value of k
          2 knn_classifier = KNeighborsClassifier(n_neighbors=k)
          3 knn_classifier.fit(X_train, y_train)
```

Out[19]: KNeighborsClassifier(n_neighbors=3)

Loading [MathJax]/extensions/MathML/content-mathml.js

In [20]: 1 y_pred = knn_classifier.predict(X_test)

```
C:\Users\Anusha V\anaconda3\lib\site-packages\sklearn\neighbors\_classification.py:228: FutureWarning: Unlike other reduction functions (e.g. `skew`, `kurtosis`), the default behavior of `mode` typically preserves the axis it acts along. In SciPy 1.11.0, this behavior will change: the default value of `keepdims` will become False, the `axis` over which the statistic is taken will be eliminated, and the value None will no longer be accepted. Set `keepdims` to True or False to avoid this warning.
    mode, _ = stats.mode(_y[neigh_ind, k], axis=1)
```

In [21]: 1 accuracy = accuracy_score(y_test, y_pred)
2 print(f"Accuracy: {accuracy * 100:.2f}%")
3

Accuracy: 100.00%

k MEAN

In [22]: 1 data=pd.read_csv(r"C:\Users\Anusha V\Downloads\archive (8)\APY.csv")
2 selected_columns = ['Crop_Year', 'Season', 'Yield']
3 X=data[selected_columns]
4

In [23]: 1 X.isnull().sum()

Out[23]: Crop_Year 0
Season 0
Yield 0
dtype: int64

In [9]: 1 import numpy as np
2 import pandas as pd
3 from sklearn.datasets import make_classification

In [10]:

```
1 import pandas as pd
2 from sklearn.datasets import make_classification
3
4 # Generate synthetic data
5 X, y = make_classification(n_features=5, n_redundant=0, n_informative=5
6
7 # Create a DataFrame
8 df = pd.DataFrame(X, columns=["Feature1", "Feature2", "Feature3", "Feat
9 df['target'] = y
10
11 # Print the shape and the first few rows of the DataFrame
12 print(df.shape)
13 print(df.head())
14
```

(100, 6)

| | Feature1 | Feature2 | Feature3 | Feature4 | Feature5 | target |
|---|-----------|-----------|-----------|-----------|-----------|--------|
| 0 | -0.528476 | 2.049389 | 0.704279 | -3.275012 | 3.555149 | 1 |
| 1 | 0.515171 | -0.698543 | 0.194818 | -3.008275 | -0.890609 | 1 |
| 2 | 0.051224 | 1.454693 | -1.665245 | -1.872938 | 0.638409 | 1 |
| 3 | 1.920039 | 0.189161 | 1.949257 | -0.355041 | -2.558094 | 0 |
| 4 | 0.398833 | -1.520716 | 0.397772 | -0.063983 | 0.245394 | 0 |

In [15]:

```
1 import random
2
3 # function for row sampling
4 def sample_rows(df, percent):
5     return df.sample(int(percent * df.shape[0]), replace=True)
6
7 # function for feature sampling
8 def sample_features(df, percent):
9     cols = random.sample(df.columns.tolist()[:-1], int(percent * df.shape[1]))
10    return df[cols]
11
12 # Example usage
13 # df is your DataFrame
14 # Sample 50% of rows
15 sampled_rows = sample_rows(df, 0.5)
16
17 # Sample 50% of features (columns)
18 sampled_features = sample_features(df, 0.5)
19 sampled_features
20
21
```

Out[15]:

| | Feature3 | Feature1 | Feature5 |
|-----|-----------|-----------|-----------|
| 0 | 0.704279 | -0.528476 | 3.555149 |
| 1 | 0.194818 | 0.515171 | -0.890609 |
| 2 | -1.665245 | 0.051224 | 0.638409 |
| 3 | 1.949257 | 1.920039 | -2.558094 |
| 4 | 0.397772 | 0.398833 | 0.245394 |
| ... | ... | ... | ... |
| 95 | 0.580032 | -0.974712 | 0.495764 |
| 96 | 2.427803 | -0.627004 | 3.220083 |
| 97 | -0.059303 | -0.836839 | 0.086734 |
| 98 | 1.429864 | -0.874207 | -1.519296 |
| 99 | 0.378503 | 0.544607 | -2.403189 |

100 rows × 3 columns

In [22]:

```
1 def combined_sampling (df,row_percent, col_percent):
2     new_df = sample_rows(df, row_percent)
3     return sample_features (new_df,col_percent)
4 import random
5
6 # Example usage
7 # df is your DataFrame
8 # Sample 50% of rows and 50% of features (columns)
9 sampled_data = combined_sampling(df, 0.5, 0.5)
10 sampled_data
```

Out[22]:

| | Feature5 | Feature2 | Feature1 |
|---|-----------|-----------|-----------|
| 25 | 2.383718 | -0.194155 | 1.953884 |
| 31 | 0.162926 | -0.245893 | 2.985782 |
| 11 | 1.199305 | 3.189073 | 1.504901 |
| 90 | 0.660209 | -2.729023 | -0.753172 |
| 9 | -2.700445 | 2.972274 | 0.472425 |
| 81 | -4.573964 | 2.196432 | 2.500614 |
| 89 | -4.038451 | 2.387250 | 1.216692 |
| 37 | 1.463461 | -1.402099 | 0.998925 |
| 67 | 2.544064 | -0.610241 | 0.347334 |
| 54 | -2.098093 | 1.744157 | 0.603254 |
| 77 | 2.639847 | -0.982625 | 1.127300 |
| 80 | -0.606813 | -0.964081 | 1.735420 |
| 88 | 0.099062 | -0.056112 | 1.762495 |
| 97 | 0.086734 | -1.007120 | -0.836839 |
| 51 | 0.581332 | -0.986208 | 0.246330 |
| 97 | 0.086734 | -1.007120 | -0.836839 |
| 95 | 0.495764 | -2.351240 | -0.974712 |
| 78 | -2.006040 | 1.333253 | 0.764933 |
| 95 | 0.495764 | -2.351240 | -0.974712 |
| 45 | 0.047452 | -1.057637 | 1.744233 |
| 19 | 0.008011 | -1.770500 | 2.099001 |
| 7 | 1.537535 | -1.629851 | 0.325215 |
| 9 | -2.700445 | 2.972274 | 0.472425 |
| 78 | -2.006040 | 1.333253 | 0.764933 |
| 0 | 3.555149 | 2.049389 | -0.528476 |
| 72 | -1.738130 | -1.777368 | 1.360138 |
| 26 | 0.646512 | 0.972441 | 2.132425 |
| 57 | -1.258582 | -1.667245 | 1.509598 |
| 96 | 3.220083 | -2.182029 | -0.627004 |
| 1 | -0.890609 | -0.698543 | 0.515171 |
| 4 | 0.245394 | -1.520716 | 0.398833 |
| 25 | 2.383718 | -0.194155 | 1.953884 |
| 1 | -0.890609 | -0.698543 | 0.515171 |
| 92 | 2.917935 | 2.023404 | -1.276507 |
| 31 | 0.162926 | -0.245893 | 2.985782 |
| 71 | -0.543031 | -0.412917 | -2.981880 |
| 14 | -0.654628 | 0.111394 | -0.737606 |
| 51 | 0.581332 | -0.986208 | 0.246330 |
| Loading [MathJax]/extensions/MathML/content-mathml.js | | | |
| 80 | -0.606813 | -0.964081 | 1.735420 |

| | Feature5 | Feature2 | Feature1 |
|----|-----------|-----------|-----------|
| 9 | -2.700445 | 2.972274 | 0.472425 |
| 27 | 1.832087 | -1.387378 | 2.833763 |
| 52 | 3.108148 | 1.295792 | 1.440234 |
| 2 | 0.638409 | 1.454693 | 0.051224 |
| 24 | 2.737663 | -2.721685 | 2.150645 |
| 51 | 0.581332 | -0.986208 | 0.246330 |
| 26 | 0.646512 | 0.972441 | 2.132425 |
| 65 | -2.930517 | -0.468903 | 3.043369 |
| 72 | -1.738130 | -1.777368 | 1.360138 |
| 49 | -0.716775 | -1.915392 | -2.627376 |
| 78 | -2.006040 | 1.333253 | 0.764933 |

In [23]:

```

1 df1 = sample_rows(df, 0.1)
2 df1.shape
3

```

Out[23]: (10, 6)

In [24]:

```

1 df1 = sample_rows(df, 0.1)
2 df1.shape
3

```

Out[24]: (10, 6)

In [25]:

```

1 df3 = sample_rows(df, 0.1)
2 df3.shape
3

```

Out[25]: (10, 6)

In [26]:

```

1 from sklearn.tree import DecisionTreeClassifier
2 clf1 = DecisionTreeClassifier()
3 clf2 = DecisionTreeClassifier()
4 clf3 = DecisionTreeClassifier()
5

```

In [27]:

```

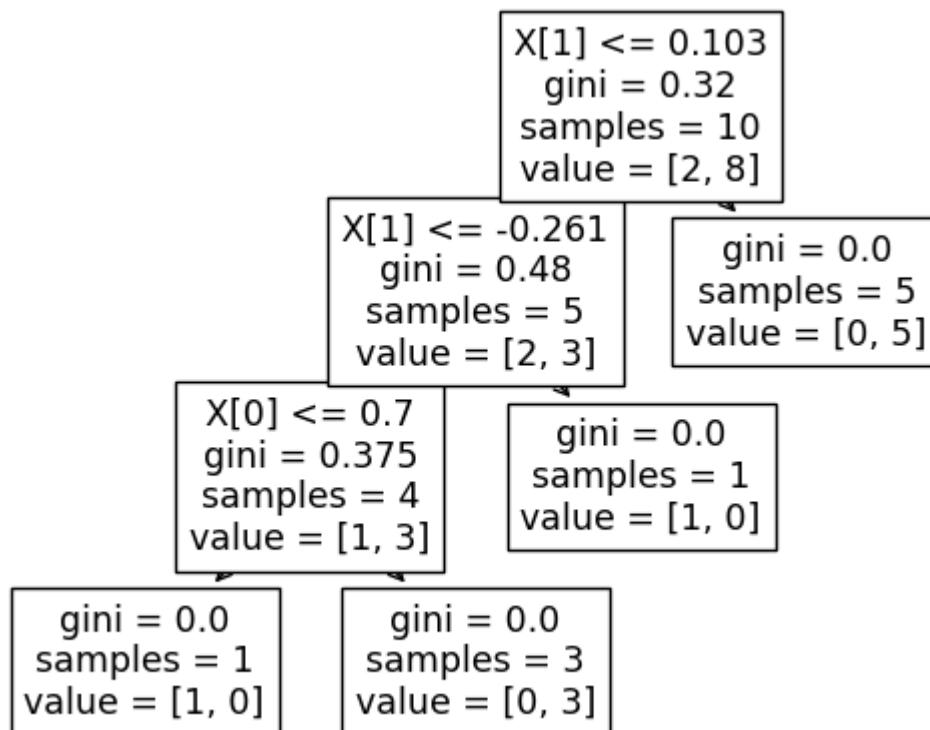
1 clf1.fit(df1.iloc[:, 0:2], df1.iloc[:, -1])
2 clf2.fit(df2.iloc[:, 0:2], df2.iloc[:, -1])
3 clf3.fit(df3.iloc[:, 0:2], df3.iloc[:, -1])

```

Out[27]: DecisionTreeClassifier()

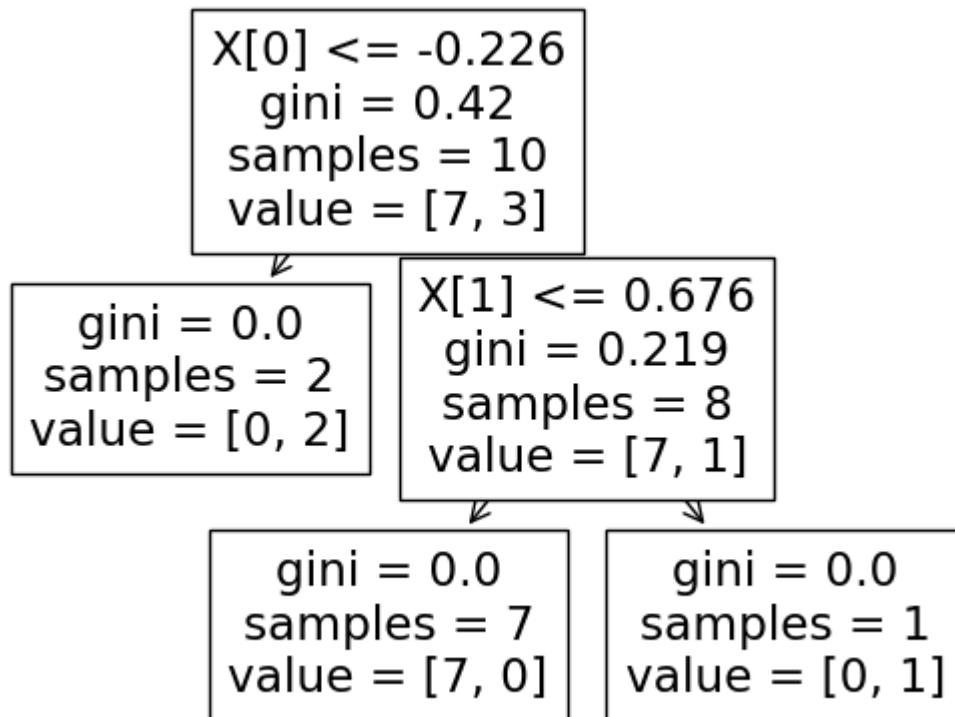
```
In [28]: 1 from sklearn.tree import plot_tree
2 plot_tree(clf1)
```

```
Out[28]: [Text(0.6666666666666666, 0.875, 'X[1] <= 0.103\ngini = 0.32\nsamples = 10\nvalue = [2, 8']),
Text(0.5, 0.625, 'X[1] <= -0.261\ngini = 0.48\nsamples = 5\nvalue = [2, 3']),
Text(0.3333333333333333, 0.375, 'X[0] <= 0.7\ngini = 0.375\nsamples = 4\nvalue = [1, 3']),
Text(0.1666666666666666, 0.125, 'gini = 0.0\nsamples = 1\nvalue = [1, 0']),
Text(0.5, 0.125, 'gini = 0.0\nsamples = 3\nvalue = [0, 3']),
Text(0.6666666666666666, 0.375, 'gini = 0.0\nsamples = 1\nvalue = [1, 0']),
Text(0.8333333333333334, 0.625, 'gini = 0.0\nsamples = 5\nvalue = [0, 5])]
```



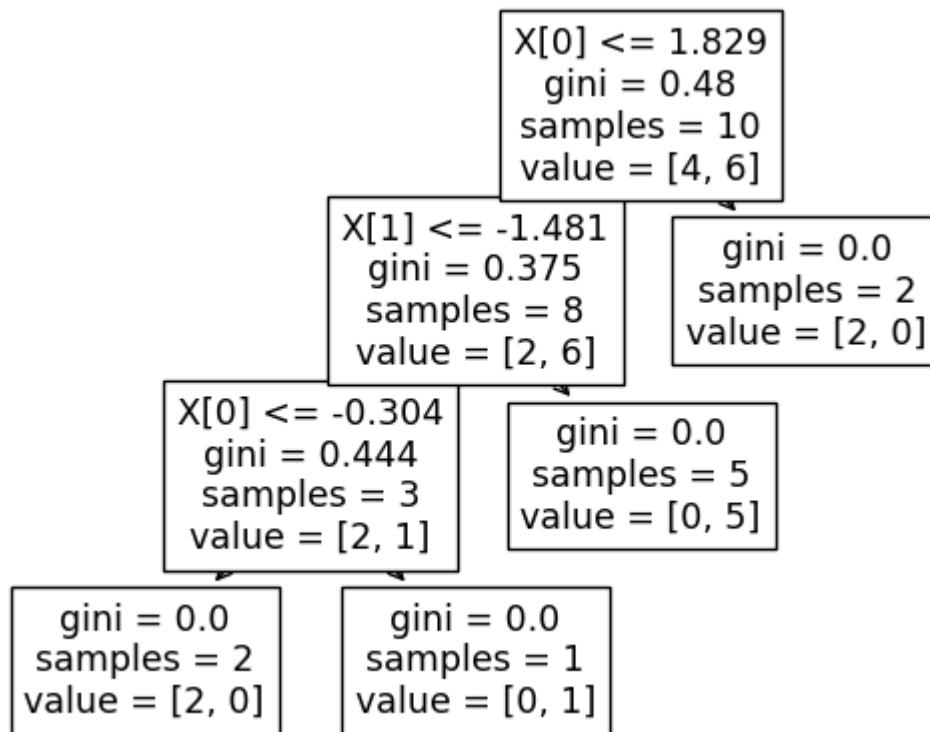
```
In [29]: 1 plot_tree(c1f2)
```

```
Out[29]: [Text(0.4, 0.8333333333333334, 'X[0] <= -0.226\n gini = 0.42\n samples = 10\n nvalue = [7, 3]'),  
 Text(0.2, 0.5, ' gini = 0.0\n samples = 2\n nvalue = [0, 2]'),  
 Text(0.6, 0.5, 'X[1] <= 0.676\n gini = 0.219\n samples = 8\n nvalue = [7, 1]'),  
 Text(0.4, 0.1666666666666666, ' gini = 0.0\n samples = 7\n nvalue = [7, 0]'),  
 Text(0.8, 0.1666666666666666, ' gini = 0.0\n samples = 1\n nvalue = [0, 1]')]
```



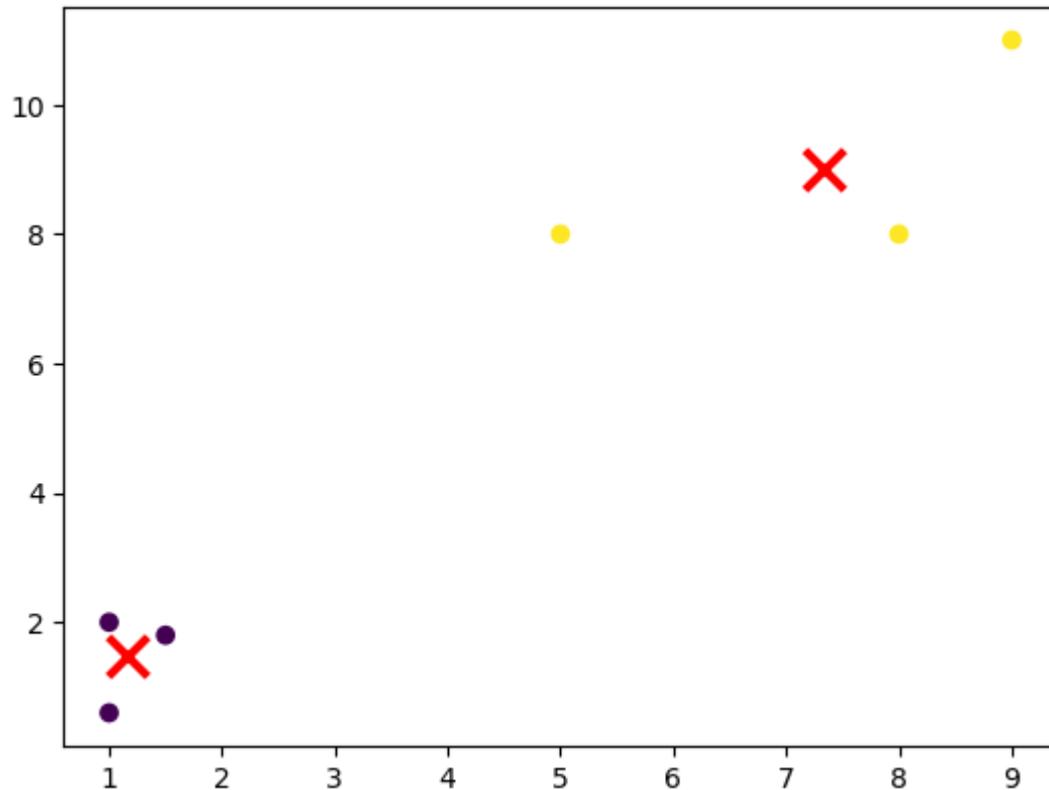
```
In [30]: 1 plot_tree(c1f3)
```

```
Out[30]: [Text(0.6666666666666666, 0.875, 'X[0] <= 1.829\ngini = 0.48\nsamples = 10\nvalue = [4, 6]'),  
Text(0.5, 0.625, 'X[1] <= -1.481\ngini = 0.375\nsamples = 8\nvalue = [2, 6]'),  
Text(0.3333333333333333, 0.375, 'X[0] <= -0.304\ngini = 0.444\nsamples = 3\nvalue = [2, 1]'),  
Text(0.1666666666666666, 0.125, 'gini = 0.0\nsamples = 2\nvalue = [2, 0]'),  
Text(0.5, 0.125, 'gini = 0.0\nsamples = 1\nvalue = [0, 1]'),  
Text(0.6666666666666666, 0.375, 'gini = 0.0\nsamples = 5\nvalue = [0, 5]'),  
Text(0.8333333333333334, 0.625, 'gini = 0.0\nsamples = 2\nvalue = [2, 0])]
```



In [31]:

```
1 from sklearn.cluster import KMeans
2 import matplotlib.pyplot as plt
3
4 # Generate some example data (replace this with your own data)
5 import numpy as np
6 X = np.array([[1, 2], [5, 8], [1.5, 1.8], [8, 8], [1, 0.6], [9, 11]])
7
8 # Create a K-Means model with the desired number of clusters (e.g., 2)
9 kmeans = KMeans(n_clusters=2)
10
11 # Fit the model to your data
12 kmeans.fit(X)
13
14 # Get the cluster labels for each data point
15 labels = kmeans.labels_
16
17 # Get the coordinates of the cluster centers
18 centers = kmeans.cluster_centers_
19
20 # Plot the data points and cluster centers
21 plt.scatter(X[:, 0], X[:, 1], c=labels)
22 plt.scatter(centers[:, 0], centers[:, 1], marker='x', s=200, linewidths=2)
23 plt.show()
24
```



hierarchical clustering

```
In [33]: 1 from sklearn.cluster import AgglomerativeClustering  
2 clustering = AgglomerativeClustering(n_clusters=3)  
3 labels = clustering.fit_predict(X)  
4 labels
```

```
Out[33]: array([1, 0, 1, 0, 1, 2], dtype=int64)
```

density based clustering

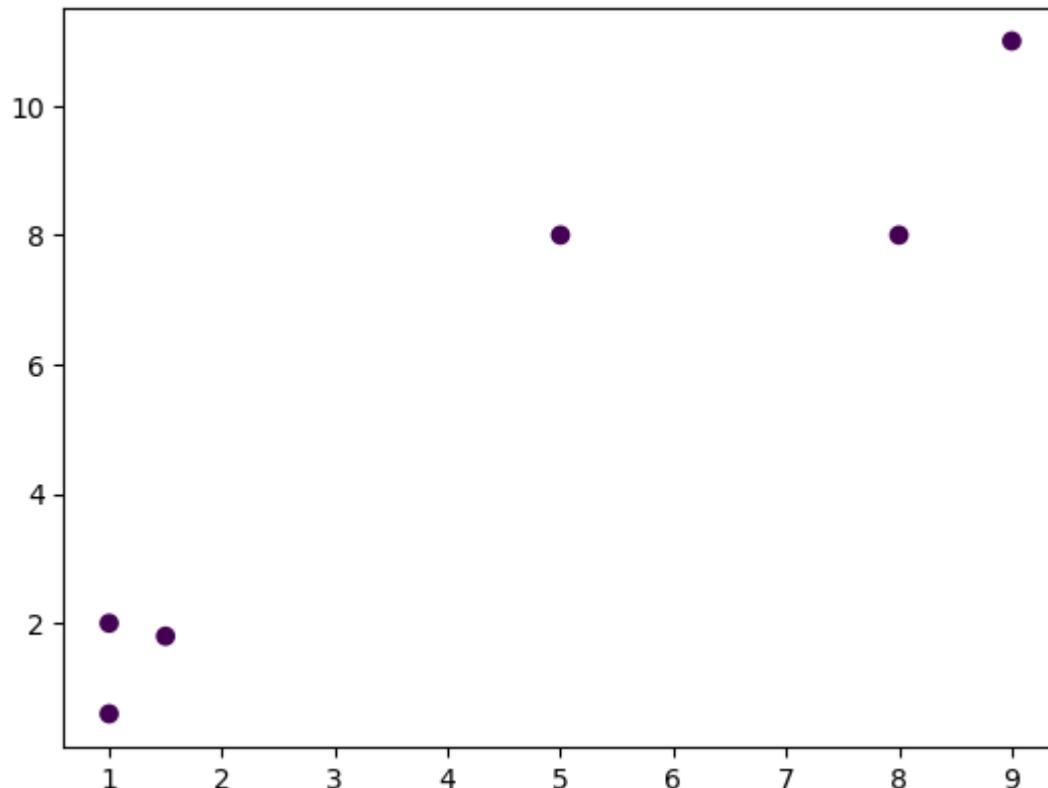
```
In [34]: 1 from sklearn.cluster import DBSCAN  
2 dbscan = DBSCAN(eps=0.3, min_samples=5)  
3 labels = dbscan.fit_predict(X)  
4 labels
```

```
Out[34]: array([-1, -1, -1, -1, -1, -1], dtype=int64)
```

partitional clustering

In [36]:

```
1 from sklearn.cluster import DBSCAN
2 import matplotlib.pyplot as plt
3
4 # Generate some example data (replace this with your own data)
5 import numpy as np
6 X = np.array([[1, 2], [5, 8], [1.5, 1.8], [8, 8], [1, 0.6], [9, 11]])
7
8 # Create a DBSCAN model
9 dbSCAN = DBSCAN(eps=0.3, min_samples=2)
10
11 # Fit the model to your data
12 dbSCAN.fit(X)
13
14 # Get the cluster labels for each data point (-1 represents noise)
15 labels = dbSCAN.labels_
16
17 # Plot the data points and clusters (including noise)
18 plt.scatter(X[:, 0], X[:, 1], c=labels)
19 plt.show()
20
```



distribution model based of cluster

In [37]:

```

1 import numpy as np
2 import matplotlib.pyplot as plt
3 from sklearn.mixture import GaussianMixture
4
5 # Generate some example data (replace this with your own data)
6 data = np.random.randn(300, 2)
7 data[100:200, 0] += 5
8 data[200:300, 1] += 5
9
10 # Create a K-Means model with the desired number of clusters (e.g., 3)
11 kmeans = KMeans(n_clusters=3)
12 kmeans
13

```

Out[37]: KMeans(n_clusters=3)

fuzzy clustering

In [39]:

```

1 pip install -U scikit-fuzzy

Collecting scikit-fuzzy
  Downloading scikit-fuzzy-0.4.2.tar.gz (993 kB)
  ----- 994.0/994.0 kB 31.2 kB/s eta 0: 00:00
    Preparing metadata (setup.py): started
    Preparing metadata (setup.py): finished with status 'done'
Requirement already satisfied: numpy>=1.6.0 in c:\users\anusha v\anaconda3\lib\site-packages (from scikit-fuzzy) (1.21.5)
Requirement already satisfied: scipy>=0.9.0 in c:\users\anusha v\anaconda3\lib\site-packages (from scikit-fuzzy) (1.9.1)
Requirement already satisfied: networkx>=1.9.0 in c:\users\anusha v\anaconda3\lib\site-packages (from scikit-fuzzy) (2.8.4)
Building wheels for collected packages: scikit-fuzzy
  Building wheel for scikit-fuzzy (setup.py): started
  Building wheel for scikit-fuzzy (setup.py): finished with status 'done'
  Created wheel for scikit-fuzzy: filename=scikit_fuzzy-0.4.2-py3-none-any.whl size=894075 sha256=4343413b7cb37e9505f40ff3362f46f146bfd963813d0a4de12125836171b977
  Stored in directory: c:\users\anusha v\appdata\local\pip\cache\wheels\32\2c\aa\90a7d7dd8448ec029f298a61f3490275e99b17aa348be675c
Successfully built scikit-fuzzy
Installing collected packages: scikit-fuzzy
Successfully installed scikit-fuzzy-0.4.2
Note: you may need to restart the kernel to use updated packages.

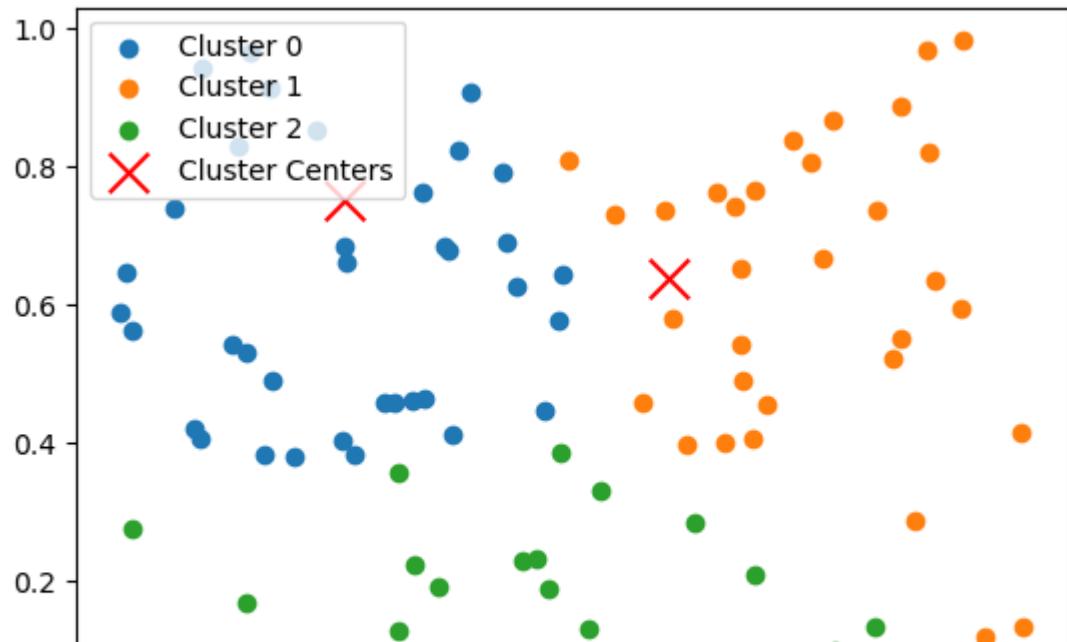
WARNING: Retrying (Retry(total=4, connect=None, read=None, redirect=None, status=None)) after connection broken by 'ReadTimeoutError("HTTPSConnectionPool(host='pypi.org', port=443): Read timed out. (read timeout=15)": /simple/scikit-fuzzy)'
WARNING: Retrying (Retry(total=3, connect=None, read=None, redirect=None, status=None)) after connection broken by 'ReadTimeoutError("HTTPSConnectionPool(host='pypi.org', port=443): Read timed out. (read timeout=15)": /simple/scikit-fuzzy)'
WARNING: Retrying (Retry(total=2, connect=None, read=None, redirect=None, status=None)) after connection broken by 'ReadTimeoutError("HTTPSConnectionPool(host='pypi.org', port=443): Read timed out. (read timeout=15)": /simple/scikit-fuzzy)'

```

Loading [MathJax]/extensions/MathML/content-mathml.js

In [40]:

```
1 import numpy as np
2 import skfuzzy as fuzz
3 from skfuzzy.cluster import cmeans
4 import matplotlib.pyplot as plt
5
6 # Generate some example data (replace this with your own data)
7 data = np.random.rand(2, 100)
8
9 # Specify the number of clusters
10 n_clusters = 3
11
12 # Specify the fuzziness exponent (typically 2.0 for FCM)
13 m = 2.0
14
15 # FCM clustering
16 cntr, u, u0, d, jm, p, fpc = cmeans(data, n_clusters, m, error=0.005, m=2.0)
17
18 # Get the cluster memberships for each data point
19 cluster_membership = np.argmax(u, axis=0)
20
21 # Visualize the data points and cluster centers
22 for i in range(n_clusters):
23     plt.scatter(data[0], cluster_membership == i], data[1], cluster_membership == i]
24
25 plt.scatter(cntr[0], cntr[1], s=200, c='red', marker='x', label='Cluster Centers')
26 plt.legend()
27 plt.show()
28
```



knn classification

```
In [2]: 1 # Import necessary Libraries
2 from sklearn.neighbors import KNeighborsClassifier
3 from sklearn.datasets import load_iris
4 from sklearn.model_selection import train_test_split
5
6 # Load a sample dataset (Iris dataset)
7 iris = load_iris()
8 X, y = iris.data, iris.target
9
10 # Split the dataset into training and testing sets
11 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2)
12
13 # Create a k-NN classifier
14 knn = KNeighborsClassifier(n_neighbors=3)
15
16 # Train the classifier on the training data
17 knn.fit(X_train, y_train)
18
19 # Predict the class labels for the test data
20 y_pred = knn.predict(X_test)
21 y_pred
```

C:\Users\Anusha V\anaconda3\lib\site-packages\sklearn\neighbors_classification.py:228: FutureWarning: Unlike other reduction functions (e.g. `skew`, `kurtosis`), the default behavior of `mode` typically preserves the axis it acts along. In SciPy 1.11.0, this behavior will change: the default value of `keepdims` will become False, the `axis` over which the statistic is taken will be eliminated, and the value None will no longer be accepted. Set `keepdims` to True or False to avoid this warning.

```
    mode, _ = stats.mode(_y[neigh_ind, k], axis=1)
```

Out[2]: array([1, 0, 2, 1, 1, 0, 1, 2, 1, 1, 2, 0, 0, 0, 0, 1, 2, 1, 1, 2, 0, 2,
0, 2, 2, 2, 2, 2, 0, 0])

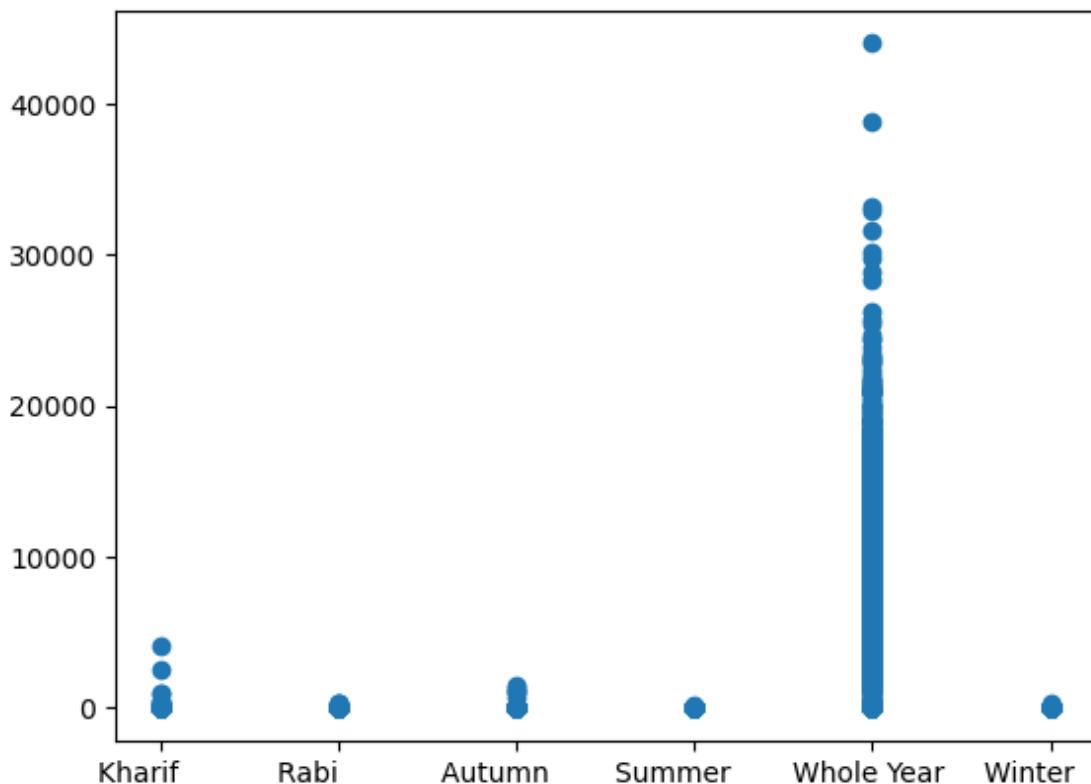
```
In [18]: 1 df=pd.read_csv(r"C:\Users\Anusha V\Downloads\archive (8)\APY.csv")
2 df.head(10)
```

Out[18]:

| | State | District | Crop | Crop_Year | Season | Area | Production | Yield |
|---|----------------------------|----------|----------|-----------|------------|--------|------------|-------|
| 0 | Andaman and Nicobar Island | NICOBARS | Arecanut | 2007 | Kharif | 2439.6 | 3415.0 | 1.40 |
| 1 | Andaman and Nicobar Island | NICOBARS | Arecanut | 2007 | Rabi | 1626.4 | 2277.0 | 1.40 |
| 2 | Andaman and Nicobar Island | NICOBARS | Arecanut | 2008 | Autumn | 4147.0 | 3060.0 | 0.74 |
| 3 | Andaman and Nicobar Island | NICOBARS | Arecanut | 2008 | Summer | 4147.0 | 2660.0 | 0.64 |
| 4 | Andaman and Nicobar Island | NICOBARS | Arecanut | 2009 | Autumn | 4153.0 | 3120.0 | 0.75 |
| 5 | Andaman and Nicobar Island | NICOBARS | Arecanut | 2009 | Summer | 4153.0 | 2080.0 | 0.50 |
| 6 | Andaman and Nicobar Island | NICOBARS | Arecanut | 2000 | Kharif | 1254.0 | 2000.0 | 1.59 |
| 7 | Andaman and Nicobar Island | NICOBARS | Arecanut | 2001 | Kharif | 1254.0 | 2061.0 | 1.64 |
| 8 | Andaman and Nicobar Island | NICOBARS | Arecanut | 2002 | Whole Year | 1258.0 | 2083.0 | 1.66 |
| 9 | Andaman and Nicobar Island | NICOBARS | Arecanut | 2003 | Whole Year | 1261.0 | 1525.0 | 1.21 |

```
In [19]: 1 import matplotlib.pyplot as plt
2 plt.scatter(df["Season"],df["Yield"])
```

Out[19]: <matplotlib.collections.PathCollection at 0x2126a43d8e0>



```
In [23]: 1 from sklearn.cluster import KMeans
2 import numpy as np
3
4 # Sample data
5 data = np.array([[1, 2], [1, 4], [1, 0], [4, 2], [4, 4], [4, 0]])
6
7 # Create a KMeans object with the desired number of clusters
8 km = KMeans(n_clusters=3)
9
10 # Fit the data and predict the cluster labels
11 cluster_labels = km.fit_predict(data)
12
13 # The cluster_labels now contain the assigned cluster labels for each d
14 print(cluster_labels)
15
```

[0 1 0 2 1 2]

In []: 1

Breast Cancer Wisconsin (diagnostic)

```
In [3]: 1 import numpy as np
2 import pandas as pd
3 from sklearn.datasets import load_breast_cancer
4 from sklearn.model_selection import train_test_split
5 from sklearn.ensemble import RandomForestClassifier
6 from sklearn.metrics import accuracy_score, classification_report
```

```
In [4]: 1 data = load_breast_cancer()
2 X = data.data
3 y = data.target
```

```
In [6]: 1 rf_classifier = RandomForestClassifier(n_estimators=100, random_state=42)
2 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2)
```

```
In [7]: 1 rf_classifier. fit(X_train, y_train)
2
```

Out[7]: RandomForestClassifier(random_state=42)

```
In [9]: 1 y_pred = rf_classifier.predict(X_test)
2 accuracy = accuracy_score(y_test, y_pred)
3 accuracy
```

Out[9]: 0.9649122807017544

In [10]:

```

1 report = classification_report(y_test, y_pred)
2 print(report)
3

```

| | precision | recall | f1-score | support |
|--------------|-----------|--------|----------|---------|
| 0 | 0.98 | 0.93 | 0.95 | 43 |
| 1 | 0.96 | 0.99 | 0.97 | 71 |
| accuracy | | | 0.96 | 114 |
| macro avg | 0.97 | 0.96 | 0.96 | 114 |
| weighted avg | 0.97 | 0.96 | 0.96 | 114 |

Principal Component Analysis (PCA)

In [17]:

```

1 from sklearn.decomposition import PCA
2 import numpy as np
3
4 # Sample data
5 data = np.array([[1, 2, 3], [4, 5, 6], [7, 8, 9]])
6
7 # Create a PCA object with the desired number of components
8 pca = PCA(n_components=2)
9
10 # Fit the data and transform it to the reduced dimension
11 transformed_data = pca.fit_transform(data)
12
13 # The transformed_data now contains the reduced dimension representation
14 print(transformed_data)
15

```

```

[[ 5.19615242  0.        ]
 [-0.          0.        ]
 [-5.19615242  0.        ]]

```

In [30]:

```

1 import pandas as pd
2 import numpy as np
3 import random as rd
4 from sklearn.decomposition import PCA
5 from sklearn import preprocessing
6 import matplotlib.pyplot as plt

```

In [6]:

```

1 import sys
2 import pandas as pd
3 from sklearn.decomposition import PCA
4
5
6
7

```

```
In [ ]: 1 import sys
2 import pandas as pd
3 from sklearn.decomposition import PCA
4
5 # Extract the user's input from the command line arguments (e.g., the d
6 data_filename = sys.argv[1]
7
8 # Load the data from the CSV file
9 data = pd.read_csv(data_filename)
10
11 # Perform PCA on the data
12 pca = PCA(n_components=2)
13 pca_result = pca.fit_transform(data)
14
15 # Process the PCA result as needed and return it
16 print("PCA Result:")
17 print(pca_result)
18
```

```
In [7]: 1 data_filename = sys.argv[1]
2
```

```
In [14]: 1 data = pd.read_csv(r"C:\Users\Anusha V\Desktop\YouTube.csv")
2
3
```

```
In [16]: 1 from sklearn.decomposition import PCA
2
3 # Assuming 'data' is your data in the form of a NumPy array or a Pandas
4 pca = PCA(n_components=2)
5 pca_result = pca
6
7
```

```
In [17]: 1 print("PCA Result:")
2 print(pca_result)
3
```

PCA Result:
PCA(n_components=2)

```
In [19]: 1 import pandas as pd
2 data = pd.read_csv("C:/Users/Anusha V/Downloads/archive (12)/diabetes.c
3 data.head()
```

| | Pregnancies | Glucose | BloodPressure | SkinThickness | Insulin | BMI | DiabetesPedigreeFunct |
|---|-------------|---------|---------------|---------------|---------|------|-----------------------|
| 0 | 6 | 148 | 72 | 35 | 0 | 33.6 | 0.1 |
| 1 | 1 | 85 | 66 | 29 | 0 | 26.6 | 0.1 |
| 2 | 8 | 183 | 64 | 0 | 0 | 23.3 | 0.1 |
| 3 | 1 | 89 | 66 | 23 | 94 | 28.1 | 0.0 |
| 4 | 0 | 137 | 40 | 35 | 168 | 43.1 | 2.4 |

In [21]: 1 data.isnull().sum()

Out[21]: Pregnancies 0
Glucose 0
BloodPressure 0
SkinThickness 0
Insulin 0
BMI 0
DiabetesPedigreeFunction 0
Age 0
Outcome 0
dtype: int64

In [22]: 1 data.describe()

Out[22]:

| | Pregnancies | Glucose | BloodPressure | SkinThickness | Insulin | BMI | Diab |
|--------------|-------------|------------|---------------|---------------|------------|------------|------------|
| count | 768.000000 | 768.000000 | 768.000000 | 768.000000 | 768.000000 | 768.000000 | 768.000000 |
| mean | 3.845052 | 120.894531 | 69.105469 | 20.536458 | 79.799479 | 31.992578 | |
| std | 3.369578 | 31.972618 | 19.355807 | 15.952218 | 115.244002 | 7.884160 | |
| min | 0.000000 | 0.000000 | 0.000000 | 0.000000 | 0.000000 | 0.000000 | |
| 25% | 1.000000 | 99.000000 | 62.000000 | 0.000000 | 0.000000 | 27.300000 | |
| 50% | 3.000000 | 117.000000 | 72.000000 | 23.000000 | 30.500000 | 32.000000 | |
| 75% | 6.000000 | 140.250000 | 80.000000 | 32.000000 | 127.250000 | 36.600000 | |
| max | 17.000000 | 199.000000 | 122.000000 | 99.000000 | 846.000000 | 67.100000 | |

In [23]: 1 data.Outcome.value_counts()

Out[23]: 0 500
1 268
Name: Outcome, dtype: int64

In [24]: 1 268/500

Out[24]: 0.536

In [25]: 1 x = data.drop('Outcome', axis='columns')

In [26]: 1 y = data.Outcome

In [27]: 1 from sklearn.preprocessing import StandardScaler

In [28]: 1 scaler = StandardScaler()

```
In [33]: 1 X_scaled = scaler.fit_transform(x)
          2 X_scaled[:3]
```

```
Out[33]: array([[ 0.63994726,  0.84832379,  0.14964075,  0.90726993, -0.69289057,
                  0.20401277,  0.46849198,  1.4259954 ],
                 [-0.84488505, -1.12339636, -0.16054575,  0.53090156, -0.69289057,
                  -0.68442195, -0.36506078, -0.19067191],
                 [ 1.23388019,  1.94372388, -0.26394125, -1.28821221, -0.69289057,
                  -1.10325546,  0.60439732, -0.10558415]])
```

```
In [34]: 1 from sklearn.model_selection import train_test_split
```

```
In [35]: 1 X_train,X_test,y_train,y_test = train_test_split(X_scaled,y, stratify=y)
```

```
In [38]: 1 X_train.shape
```

```
Out[38]: (576, 8)
```

```
In [39]: 1 X_test.shape
```

```
Out[39]: (192, 8)
```

```
In [41]: 1 y_train.value_counts
```

```
Out[41]: <bound method IndexOpsMixin.value_counts of 745>
          0    67
          1    570
          2    517
          3    717
          ..
          4    735
          5    475
          6    46
          7    767
          8    540
          Name: Outcome, Length: 576, dtype: int64>
```

```
In [42]: 1 201/375
```

```
Out[42]: 0.536
```

```
In [43]: 1 451/854
```

```
Out[43]: 0.5281030444964872
```

```
In [44]: 1 from sklearn.tree import DecisionTreeClassifier
```

```
In [45]: 1 from sklearn.model_selection import cross_val_score
```

```
In [49]: 1 scores = cross_val_score(DecisionTreeClassifier(), x, y, cv=5)
          2 scores
```

```
Out[49]: array([0.66233766, 0.68181818, 0.68831169, 0.79084967, 0.74509804])
```

Loading [MathJax]/extensions/MathML/content-mathml.js

In [50]: 1 scores.mean()

Out[50]: 0.7136830489771666

In [57]:

```

1 from sklearn.ensemble import BaggingClassifier
2 BaggingClassifier(
3     base_estimator=DecisionTreeClassifier(),
4     n_estimators=0,
5     max_samples=0.8,
6     oob_score=True,
7     random_state=0
8 )
9
10

```

Out[57]: BaggingClassifier(base_estimator=DecisionTreeClassifier(), max_samples=0.8, n_estimators=0, oob_score=True, random_state=0)

In [2]:

```

1 import pandas as pd
2 df = pd.read_csv(r"C:\Users\Anusha V\Downloads\archive (13)\healthcare-
3 df

```

Out[2]:

| | id | gender | age | hypertension | heart_disease | ever_married | work_type | Residence_ |
|-------------|-----------|---------------|------------|---------------------|----------------------|---------------------|------------------|-------------------|
| 0 | 9046 | Male | 67.0 | 0 | 1 | Yes | Private | L |
| 1 | 51676 | Female | 61.0 | 0 | 0 | Yes | Self-employed | I |
| 2 | 31112 | Male | 80.0 | 0 | 1 | Yes | Private | I |
| 3 | 60182 | Female | 49.0 | 0 | 0 | Yes | Private | L |
| 4 | 1665 | Female | 79.0 | 1 | 0 | Yes | Self-employed | I |
| ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 5105 | 18234 | Female | 80.0 | 1 | 0 | Yes | Private | L |
| 5106 | 44873 | Female | 81.0 | 0 | 0 | Yes | Self-employed | L |
| 5107 | 19723 | Female | 35.0 | 0 | 0 | Yes | Self-employed | I |
| 5108 | 37544 | Male | 51.0 | 0 | 0 | Yes | Private | I |
| 5109 | 44679 | Female | 44.0 | 0 | 0 | Yes | Govt_job | L |

5110 rows × 12 columns

In [4]: 1 df.isnull().sum().sum

Out[4]: <bound method NDFrame._add_numeric_operations.<locals>.sum of id
0
gender 0
age 0
hypertension 0
heart_disease 0
ever_married 0
work_type 0
Residence_type 0
avg_glucose_level 0
bmi 201
smoking_status 0
stroke 0
dtype: int64>

In [6]: 1 means = df.mean()
2 df.fillna(means, inplace = True)
3 df.to_csv(r"C:\Users\Anusha V\Downloads\archive (13)\healthcare-dataset
4 means

C:\Users\Anusha V\AppData\Local\Temp\ipykernel_22100\4153111646.py:1: FutureWarning: Dropping of nuisance columns in DataFrame reductions (with 'numeric_only=None') is deprecated; in a future version this will raise TypeError. Select only valid columns before calling the reduction.

means = df.mean()

Out[6]: id 36517.829354
age 43.226614
hypertension 0.097456
heart_disease 0.054012
avg_glucose_level 106.147677
bmi 28.893237
stroke 0.048728
dtype: float64

In [7]: 1 df.head()

Out[7]:

| | id | gender | age | hypertension | heart_disease | ever_married | work_type | Residence_type |
|---|-------|--------|------|--------------|---------------|--------------|---------------|----------------|
| 0 | 9046 | Male | 67.0 | 0 | 1 | Yes | Private | Urban |
| 1 | 51676 | Female | 61.0 | 0 | 0 | Yes | Self-employed | Rural |
| 2 | 31112 | Male | 80.0 | 0 | 1 | Yes | Private | Rural |
| 3 | 60182 | Female | 49.0 | 0 | 0 | Yes | Private | Urban |
| 4 | 1665 | Female | 79.0 | 1 | 0 | Yes | Self-employed | Rural |

```
In [11]: 1 RM = df.drop_duplicates()
2 RM.to_csv(r"C:\Users\Anusha V\Downloads\archive (13)\healthcare-dataset"
3 RM
```

Out[11]:

| | id | gender | age | hypertension | heart_disease | ever_married | work_type | Residence_ |
|-------------|-----------|---------------|------------|---------------------|----------------------|---------------------|------------------|-------------------|
| 0 | 9046 | Male | 67.0 | 0 | 1 | Yes | Private | L |
| 1 | 51676 | Female | 61.0 | 0 | 0 | Yes | Self-employed | I |
| 2 | 31112 | Male | 80.0 | 0 | 1 | Yes | Private | I |
| 3 | 60182 | Female | 49.0 | 0 | 0 | Yes | Private | L |
| 4 | 1665 | Female | 79.0 | 1 | 0 | Yes | Self-employed | I |
| ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 5105 | 18234 | Female | 80.0 | 1 | 0 | Yes | Private | L |
| 5106 | 44873 | Female | 81.0 | 0 | 0 | Yes | Self-employed | L |
| 5107 | 19723 | Female | 35.0 | 0 | 0 | Yes | Self-employed | I |
| 5108 | 37544 | Male | 51.0 | 0 | 0 | Yes | Private | I |
| 5109 | 44679 | Female | 44.0 | 0 | 0 | Yes | Govt_job | L |

5110 rows × 12 columns



```
In [15]: 1 encoder = df['gender'] = df['gender'].replace({'Male': 0 , 'Female': 1})
2 encoder
```

Out[15]:

| | |
|------|---|
| 0 | 0 |
| 1 | 1 |
| 2 | 0 |
| 3 | 1 |
| 4 | 1 |
| .. | |
| 5105 | 1 |
| 5106 | 1 |
| 5107 | 1 |
| 5108 | 0 |
| 5109 | 1 |

Name: gender, Length: 5110, dtype: object

In [17]: 1 df.describe()

| | id | age | hypertension | heart_disease | avg_glucose_level | bn |
|--------------|--------------|-------------|---------------------|----------------------|--------------------------|-------------|
| count | 5110.000000 | 5110.000000 | 5110.000000 | 5110.000000 | 5110.000000 | 5110.000000 |
| mean | 36517.829354 | 43.226614 | 0.097456 | 0.054012 | 106.147677 | 28.89323 |
| std | 21161.721625 | 22.612647 | 0.296607 | 0.226063 | 45.283560 | 7.69801 |
| min | 67.000000 | 0.080000 | 0.000000 | 0.000000 | 55.120000 | 10.30000 |
| 25% | 17741.250000 | 25.000000 | 0.000000 | 0.000000 | 77.245000 | 23.80000 |
| 50% | 36932.000000 | 45.000000 | 0.000000 | 0.000000 | 91.885000 | 28.40000 |
| 75% | 54682.000000 | 61.000000 | 0.000000 | 0.000000 | 114.090000 | 32.80000 |
| max | 72940.000000 | 82.000000 | 1.000000 | 1.000000 | 271.740000 | 97.60000 |



In [19]: 1 groupby = df.groupby("age").count()
2 groupby

| | id | gender | hypertension | heart_disease | ever_married | work_type | Residence_type | a |
|--------------|-----------|---------------|---------------------|----------------------|---------------------|------------------|-----------------------|----------|
| age | | | | | | | | |
| 0.08 | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 2 |
| 0.16 | 3 | 3 | 3 | 3 | 3 | 3 | 3 | 3 |
| 0.24 | 5 | 5 | 5 | 5 | 5 | 5 | 5 | 5 |
| 0.32 | 5 | 5 | 5 | 5 | 5 | 5 | 5 | 5 |
| 0.40 | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 2 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 78.00 | 102 | 102 | 102 | 102 | 102 | 102 | 102 | 102 |
| 79.00 | 85 | 85 | 85 | 85 | 85 | 85 | 85 | 85 |
| 80.00 | 70 | 70 | 70 | 70 | 70 | 70 | 70 | 70 |
| 81.00 | 60 | 60 | 60 | 60 | 60 | 60 | 60 | 60 |
| 82.00 | 56 | 56 | 56 | 56 | 56 | 56 | 56 | 56 |

104 rows × 11 columns



```
In [23]: 1 pivot = df.pivot(columns = 'age')
2 pivot
```

Out[23]:

| | age | 0.08 | 0.16 | 0.24 | 0.32 | 0.40 | 0.48 | 0.56 | 0.64 | 0.72 | 0.80 | ... | 73.00 | 74.00 | 75.00 | 76.00 |
|------|-----|------|------|------|------|------|------|------|------|------|------|-----|-------|-------|-------|-------|
| 0 | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | ... | NaN | NaN | NaN | NaN |
| 1 | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | ... | NaN | NaN | NaN | NaN |
| 2 | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | ... | NaN | NaN | NaN | NaN |
| 3 | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | ... | NaN | NaN | NaN | NaN |
| 4 | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | ... | NaN | NaN | NaN | NaN |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 5105 | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | ... | NaN | NaN | NaN | NaN |
| 5106 | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | ... | NaN | NaN | NaN | NaN |
| 5107 | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | ... | NaN | NaN | NaN | NaN |
| 5108 | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | ... | NaN | NaN | NaN | NaN |
| 5109 | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | ... | NaN | NaN | NaN | NaN |

5110 rows × 1144 columns



```
In [22]: 1 df.dtypes
```

```
Out[22]: id                  int64
gender               object
age                  float64
hypertension        int64
heart_disease      int64
ever_married        object
work_type            object
Residence_type      object
avg_glucose_level  float64
bmi                 float64
smoking_status      object
stroke              int64
dtype: object
```

```
In [24]: 1 df = df.drop(['work_type'], axis=1)
2 df
```

Out[24]:

| | id | gender | age | hypertension | heart_disease | ever_married | Residence_type | avg_c |
|------|-------|--------|------|--------------|---------------|--------------|----------------|-------|
| 0 | 9046 | 0 | 67.0 | 0 | 1 | Yes | Urban | |
| 1 | 51676 | 1 | 61.0 | 0 | 0 | Yes | Rural | |
| 2 | 31112 | 0 | 80.0 | 0 | 1 | Yes | Rural | |
| 3 | 60182 | 1 | 49.0 | 0 | 0 | Yes | Urban | |
| 4 | 1665 | 1 | 79.0 | 1 | 0 | Yes | Rural | |
| ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 5105 | 18234 | 1 | 80.0 | 1 | 0 | Yes | Urban | |
| 5106 | 44873 | 1 | 81.0 | 0 | 0 | Yes | Urban | |
| 5107 | 19723 | 1 | 35.0 | 0 | 0 | Yes | Rural | |
| 5108 | 37544 | 0 | 51.0 | 0 | 0 | Yes | Rural | |
| 5109 | 44679 | 1 | 44.0 | 0 | 0 | Yes | Urban | |

5110 rows × 11 columns

In []:

1