

1 Exercici 1

Realitza la pràctica del notebook a GitHub "03 EXAMINING DATA" amb seaborn i el dataset

In [1]:

```
1 import numpy as np
2 import pandas as pd
3 import matplotlib.pyplot as plt
4 import seaborn as sns
```

In [2]:

```
1 df = pd.read_csv("https://raw.githubusercontent.com/IT-Academy-BCN/Data-Science/main/datasets/03_examining_data/03_examining_data.csv")
```

In [4]:

```
1 df.shape
```

(244, 7)

In [5]:

```
1 df.columns
```

Index(['total_bill', 'tip', 'sex', 'smoker', 'day', 'time', 'size'], dtype='object')

In [24]:

```
1 df.head()
```

| | total_bill | tip | sex | smoker | day | time | size |
|---|------------|------|--------|--------|-----|--------|------|
| 0 | 16.99 | 1.01 | Female | No | Sun | Dinner | 2 |
| 1 | 10.34 | 1.66 | Male | No | Sun | Dinner | 3 |
| 2 | 21.01 | 3.50 | Male | No | Sun | Dinner | 3 |
| 3 | 23.68 | 3.31 | Male | No | Sun | Dinner | 2 |
| 4 | 24.59 | 3.61 | Female | No | Sun | Dinner | 4 |

In [27]:

```
1 df.time.unique()
```

```
array(['Dinner', 'Lunch'], dtype=object)
```

In [28]:

```
1 df.time.nunique()
```

```
2
```

In [37]:

```
1 df.smoker.unique()
```

```
array(['No', 'Yes'], dtype=object)
```

In [36]:

```
1 df.sex.unique()
```

```
array(['Female', 'Male'], dtype=object)
```

In [67]:

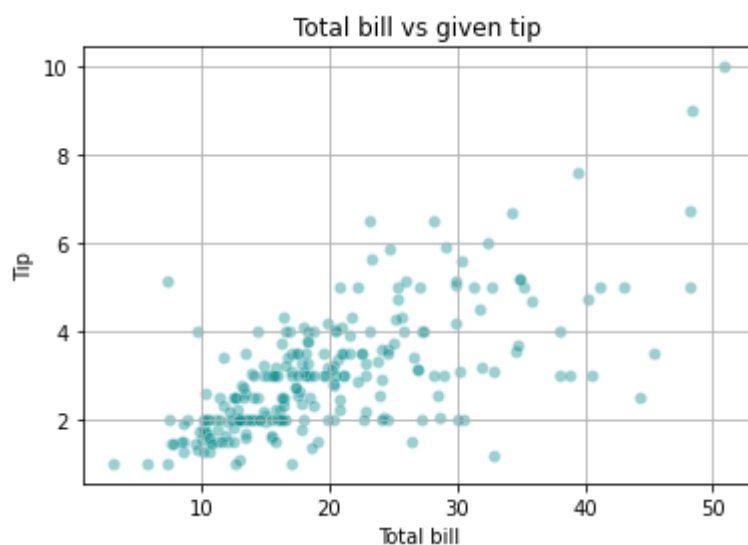
```
1 df.describe().round(2)
```

| | total_bill | tip | party |
|-------|------------|--------|--------|
| count | 244.00 | 244.00 | 244.00 |
| mean | 19.79 | 3.00 | 2.57 |
| std | 8.90 | 1.38 | 0.95 |
| min | 3.07 | 1.00 | 1.00 |
| 25% | 13.35 | 2.00 | 2.00 |
| 50% | 17.80 | 2.90 | 2.00 |
| 75% | 24.13 | 3.56 | 3.00 |
| max | 50.81 | 10.00 | 6.00 |

1.1 Scatterplots

In [23]:

```
1 #plot
2 sns.scatterplot(x = df.total_bill, y = df.tip, color = "DarkCyan", alpha = 0.4)
3
4 plt.title("Total bill vs given tip")
5 plt.ylabel("Tip")
6 plt.xlabel("Total bill")
7 plt.grid()
8 plt.show()
```



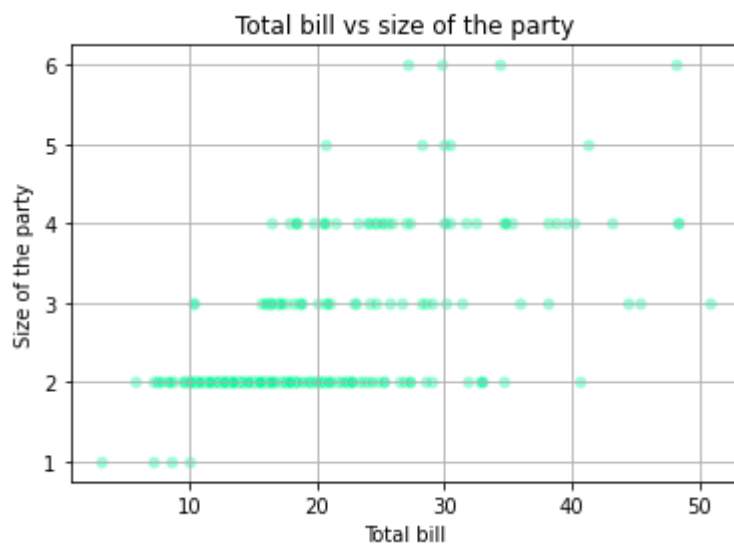
There is linear correlation between the total bill and the tip.

In [3]:

```
1 df.rename(columns = {"size" : "party"}, inplace = True)
```

In [41]:

```
1 sns.scatterplot(x = df.total_bill, y = df.party, color = "MediumSpringGreen", alph
2
3 plt.title("Total bill vs size of the party")
4 plt.ylabel("Size of the party")
5 plt.xlabel("Total bill")
6 plt.grid()
7 plt.show()
```



In [42]:

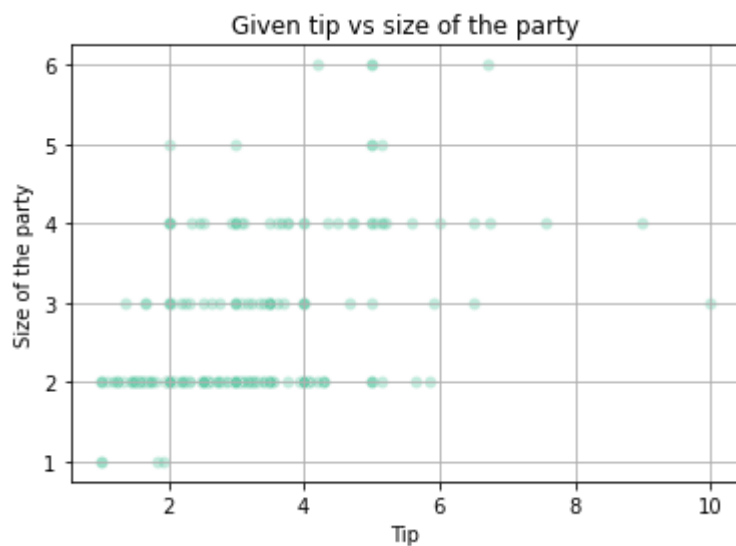
```
1 df.total_bill.corr(df.party)
```

0.5983151309049022

There is some correlaion between the total bill and the size of the party.

In [43]:

```
1 sns.scatterplot(x = df.tip, y = df.party, color = "MediumAquaMarine", alpha = 0.4)
2
3 plt.title("Given tip vs size of the party")
4 plt.ylabel("Size of the party")
5 plt.xlabel("Tip")
6 plt.grid()
7 plt.show()
```



In [44]:

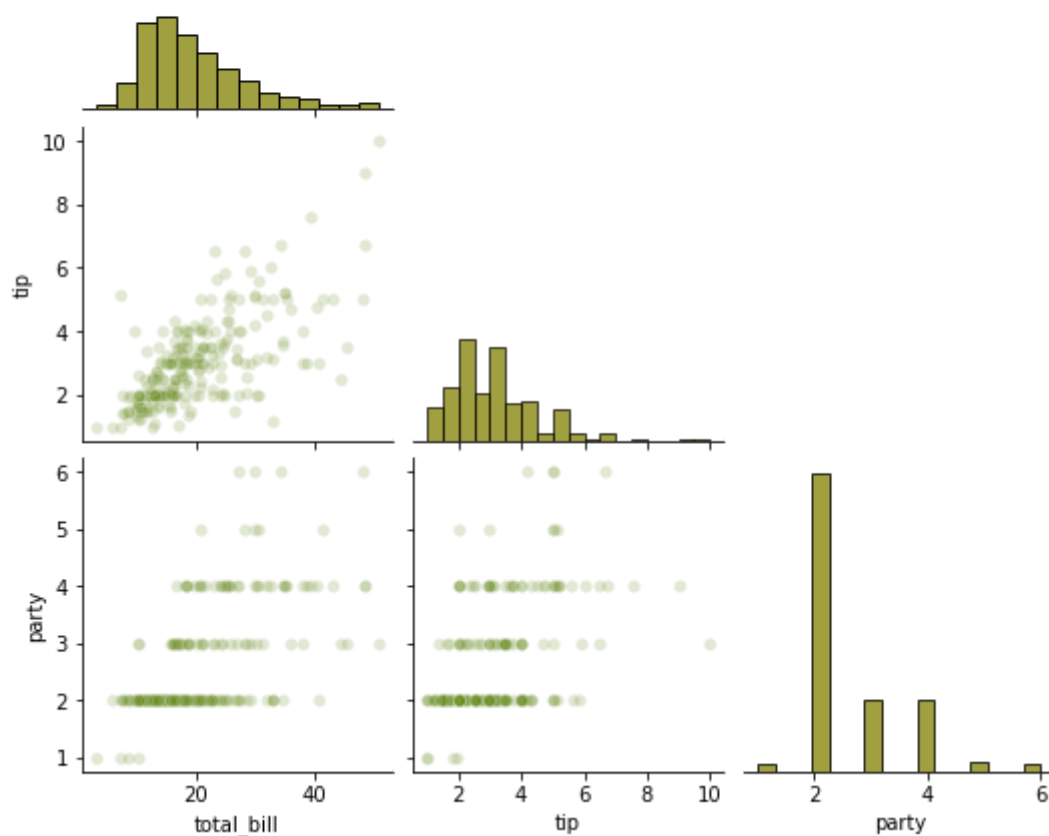
```
1 df.tip.corr(df.party)
```

0.4892987752303577

There is less correlation between the size of the party and the tip given. It is not a good pred

In [102]:

```
1 #matrix plot
2 sns.pairplot(df, diag_kind = "hist", plot_kws = {"alpha":0.2, "color":"OliveDrab"})
3     diag_kws = { "color":"Olive"}
```



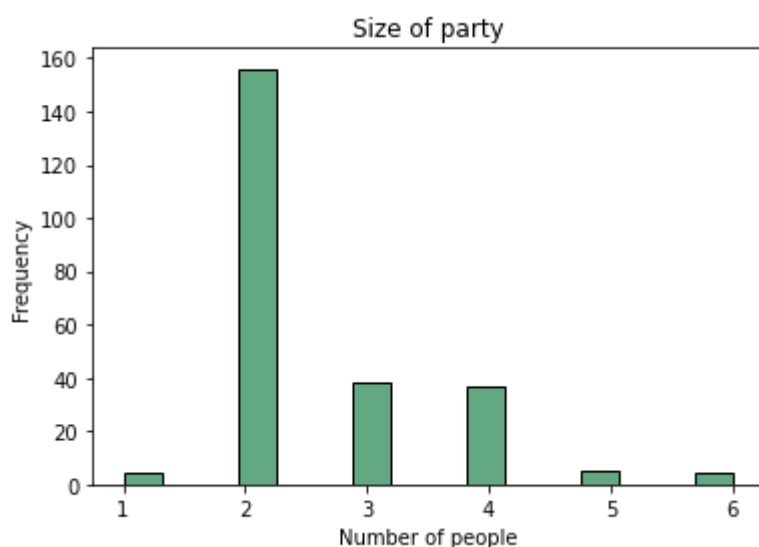
Matrix plot that gathers the above plots.

1.2 Histograms

In [104]:

```
1 sns.histplot(data = df, x = "party", color = "SeaGreen")
2 plt.title("Size of party")
3 plt.ylabel("Frequency")
4 plt.xlabel("Number of people")
```

Text(0.5, 0, 'Number of people')



The majority of people going to the restaurant go in groups of 2. That means it is unimodal a

In [78]:

```
1 df.total_bill.mean() #tip average
```

19.785942622950824

In [77]:

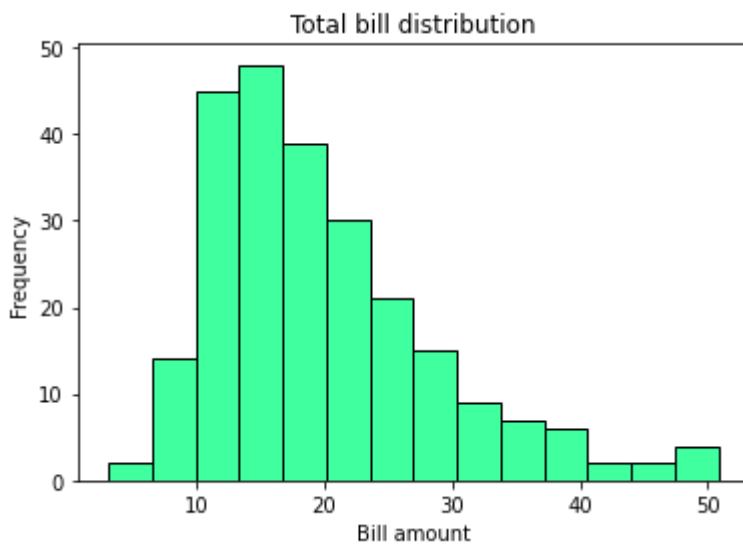
```
1 df.total_bill.std() #tip standard deviation
```

8.902411954856856

In [82]:

```
1 sns.histplot(data = df, x = "total_bill", color = "SpringGreen")
2 plt.title("Total bill distribution")
3 plt.ylabel("Frequency")
4 plt.xlabel("Bill amount")
```

Text(0.5, 0, 'Bill amount')



This distribution is unimodal, and is very skewed to the left. Most of the counts are close to the mean, which is within less than a standard deviation away from it.

In [79]:

```
1 df.head()
```

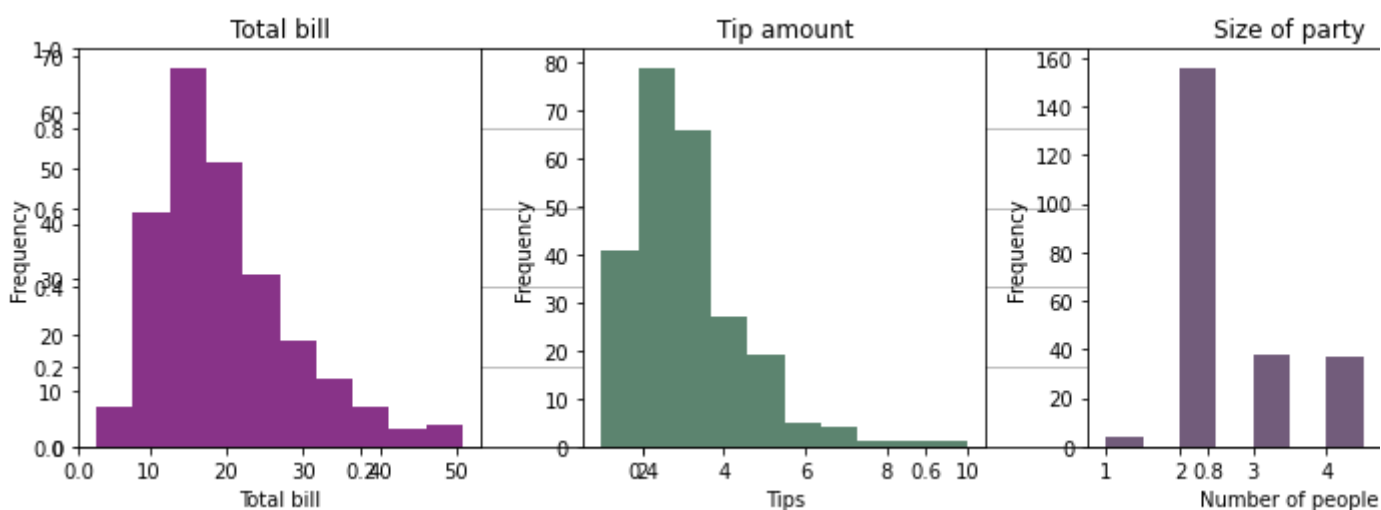
| | total_bill | tip | sex | smoker | day | time | party |
|---|------------|------|--------|--------|-----|--------|-------|
| 0 | 16.99 | 1.01 | Female | No | Sun | Dinner | 2 |
| 1 | 10.34 | 1.66 | Male | No | Sun | Dinner | 3 |
| 2 | 21.01 | 3.50 | Male | No | Sun | Dinner | 3 |
| 3 | 23.68 | 3.31 | Male | No | Sun | Dinner | 2 |
| 4 | 24.59 | 3.61 | Female | No | Sun | Dinner | 4 |

In [106]:

```

1  fig = plt.figure(figsize = (11, 4))
2
3  ax1 = fig.add_subplot(1, 3, 1)
4
5  ax1.hist(df["total_bill"], color = "#883388")
6  plt.title("Total bill")
7  plt.ylabel("Frequency")
8  plt.xlabel("Total bill")
9
10 ax2 = fig.add_subplot(1, 3, 2)
11
12 ax2.hist(df["tip"], color = "#5C846F")
13 plt.title("Tip amount")
14 plt.ylabel("Frequency")
15 plt.xlabel("Tips")
16
17 ax3 = fig.add_subplot(1, 3, 3)
18
19 ax3.hist(df["party"], color = "#725C7B")
20 plt.title("Size of party")
21 plt.ylabel("Frequency")
22 plt.xlabel("Number of people")
23
24 plt.tight_layout()

```



All of the above plots show that they are all very skewed to the left and unimodal. The mode

cases and the great majority of values are comprised in the range of the mean \pm one standa

In [96]:

```
1 df.shape
```

```
(244, 7)
```

In [99]:

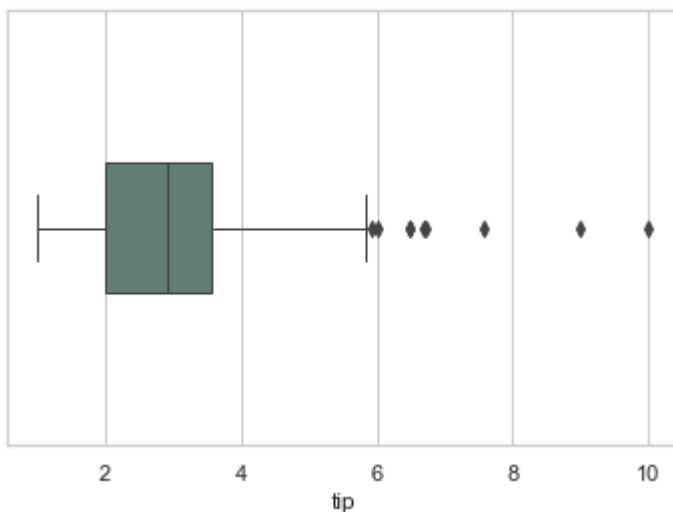
```
1 df["tip"].describe().round(2)
```

```
count    244.00
mean       3.00
std        1.38
min         1.00
25%        2.00
50%        2.90
75%        3.56
max        10.00
Name: tip, dtype: float64
```

In [127]:

```
1 sns.set(style = "whitegrid")
2 sns.boxplot(x = df["tip"], color = "#5C846F", fliersize = 5, linewidth = 1, width
```

<AxesSubplot:xlabel='tip'>

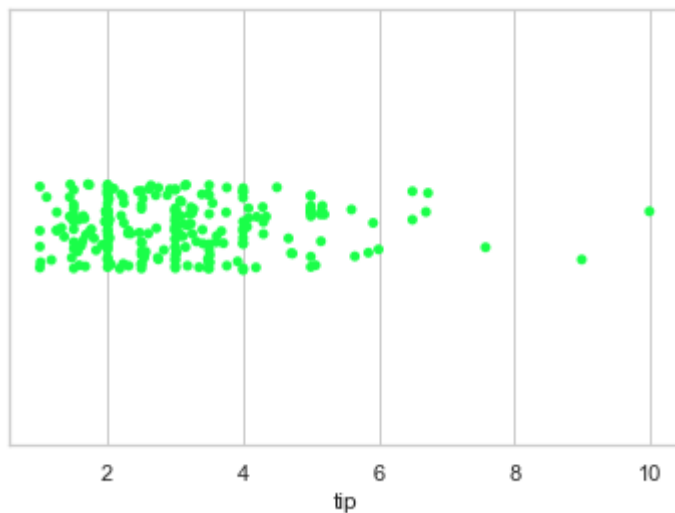


50 % of all values go between 2 and 3.56. Although the median is a bit closer to the third quartile, on the interquartile range tend to be below it, not above it, as evidenced by the shorter whisker. There are outliers that are very high and far from the other values.

```
In [131]:
```

```
1 sns.stripplot(x = df["tip"], color = "#1AFE49")
```

```
<AxesSubplot:xlabel='tip'>
```

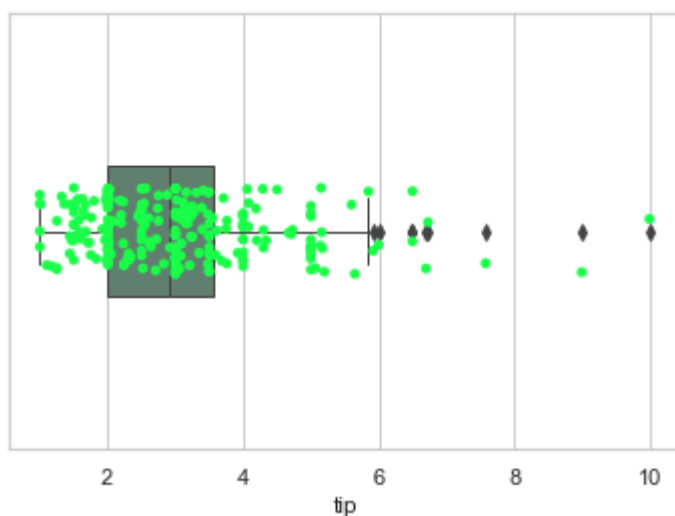


As the boxplot, this scatter shows that the majority of values are in a small range and there are a few outliers.

```
In [128]:
```

```
1 sns.boxplot(x = df["tip"], color = "#5C846F", fliersize = 5, linewidth = 1, width = 0.5)  
2 sns.stripplot(x = df["tip"], color = "#1AFE49")
```

```
<AxesSubplot:xlabel='tip'>
```

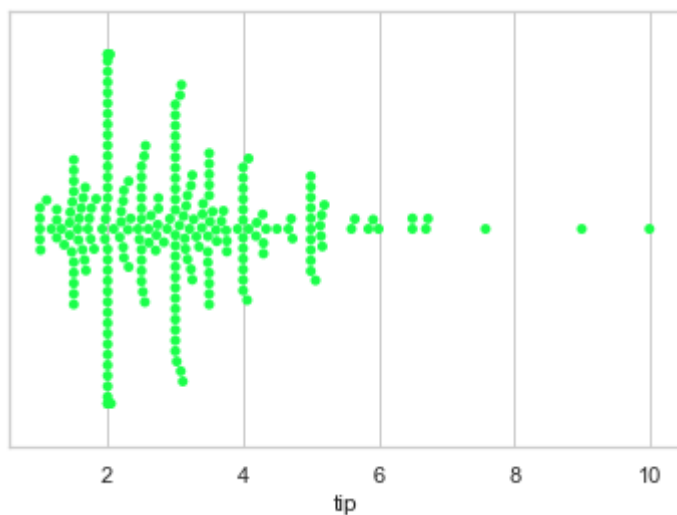


Both plots together confirm that they overlap and show that most of the values are in the interval [2, 6].

In [133]:

```
1 sns.swarmplot(x = df["tip"], color = "#1AFE49")
```

<AxesSubplot:xlabel='tip'>

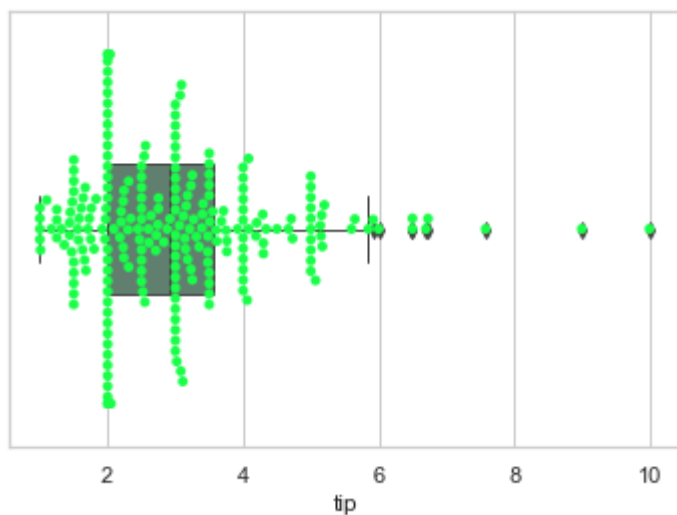


This plot allows us to see more clearly where the majority of the values fall, and shows a rather the median.

In [134]:

```
1 sns.boxplot(x = df["tip"], color = "#5C846F", fliersize = 5, linewidth = 1, width  
2 sns.swarmplot(x = df["tip"], color = "#1AFE49")
```

<AxesSubplot:xlabel='tip'>



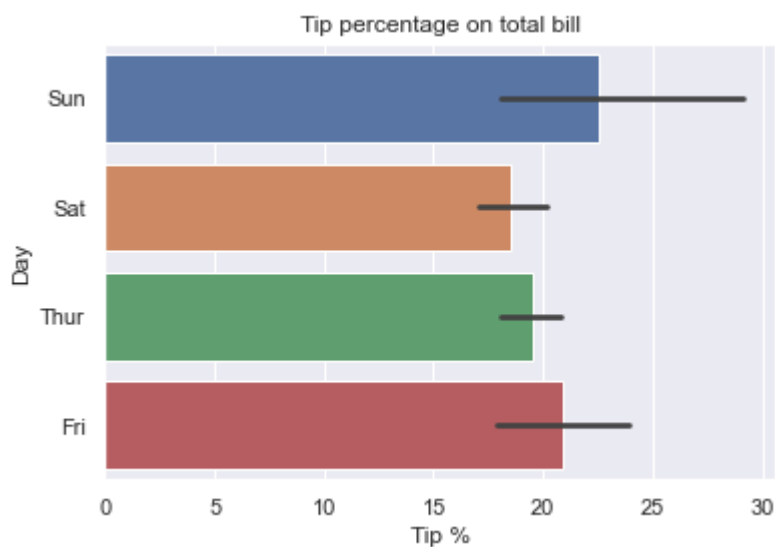
In [7]:

```
1 df['tip_pct'] = round((df['tip'] / (df['total_bill'] - df['tip'])) * 100, 2)
2 df.head()
```

| | total_bill | tip | sex | smoker | day | time | party | tip_pct |
|---|------------|------|--------|--------|-----|--------|-------|---------|
| 0 | 16.99 | 1.01 | Female | No | Sun | Dinner | 2 | 6.32 |
| 1 | 10.34 | 1.66 | Male | No | Sun | Dinner | 3 | 19.12 |
| 2 | 21.01 | 3.50 | Male | No | Sun | Dinner | 3 | 19.99 |
| 3 | 23.68 | 3.31 | Male | No | Sun | Dinner | 2 | 16.25 |
| 4 | 24.59 | 3.61 | Female | No | Sun | Dinner | 4 | 17.21 |

In [11]:

```
1 sns.barplot(x = 'tip_pct', y = 'day', data = df, orient = "h")
2 plt.title("Tip percentage on total bill")
3 plt.ylabel("Day")
4 plt.xlabel("Tip %")
5 plt.show()
6 sns.set(style = "darkgrid")
```



In [13]:

```
1 df.describe().round(2)
```

| | total_bill | tip | party | tip_pct |
|-------|------------|--------|--------|---------|
| count | 244.00 | 244.00 | 244.00 | 244.00 |
| mean | 19.79 | 3.00 | 2.57 | 20.21 |
| std | 8.90 | 1.38 | 0.95 | 16.34 |
| min | 3.07 | 1.00 | 1.00 | 3.70 |
| 25% | 13.35 | 2.00 | 2.00 | 14.83 |
| 50% | 17.80 | 2.90 | 2.00 | 18.31 |
| 75% | 24.13 | 3.56 | 3.00 | 23.68 |
| max | 50.81 | 10.00 | 6.00 | 245.24 |

In [16]:

```
1 round(df.describe(include = 'all'), 2)
```

| | total_bill | tip | sex | smoker | day | time | party | tip_pct |
|--------|------------|--------|------|--------|-----|--------|--------|---------|
| count | 244.00 | 244.00 | 244 | 244 | 244 | 244 | 244.00 | 244.00 |
| unique | NaN | NaN | 2 | 2 | 4 | 2 | NaN | NaN |
| top | NaN | NaN | Male | No | Sat | Dinner | NaN | NaN |
| freq | NaN | NaN | 157 | 151 | 87 | 176 | NaN | NaN |
| mean | 19.79 | 3.00 | NaN | NaN | NaN | NaN | 2.57 | 20.21 |
| std | 8.90 | 1.38 | NaN | NaN | NaN | NaN | 0.95 | 16.34 |
| min | 3.07 | 1.00 | NaN | NaN | NaN | NaN | 1.00 | 3.70 |
| 25% | 13.35 | 2.00 | NaN | NaN | NaN | NaN | 2.00 | 14.83 |
| 50% | 17.80 | 2.90 | NaN | NaN | NaN | NaN | 2.00 | 18.31 |
| 75% | 24.13 | 3.56 | NaN | NaN | NaN | NaN | 3.00 | 23.68 |
| max | 50.81 | 10.00 | NaN | NaN | NaN | NaN | 6.00 | 245.24 |

In [17]:

```
1 df.isnull().sum() / len(df) # % of null values
```

```
total_bill    0.0  
tip           0.0  
sex           0.0  
smoker        0.0  
day           0.0  
time          0.0  
party         0.0  
tip_pct       0.0  
dtype: float64
```

1.3 One variable

In [22]:

```
1 sns.boxplot(y = "tip_pct", data = df[df.tip < 10], color = "#FFBA49", fliersize =  
2             orient = 'v', linewidth = 1 , width = 0.3)  
3 plt.title("Tip percentage on total bill")  
4 plt.ylabel("Tip %")  
5 plt.show()
```

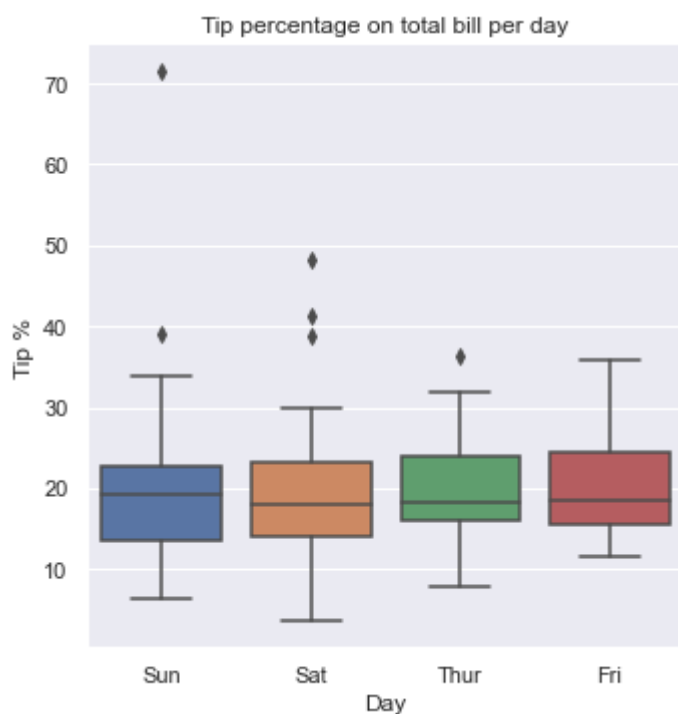


There's a value that is extremely far off from the rest, including the other outliers. Since it says bill, it could mean it was wrongly introduced into the dataset.

1.4 Two variables

In [25]:

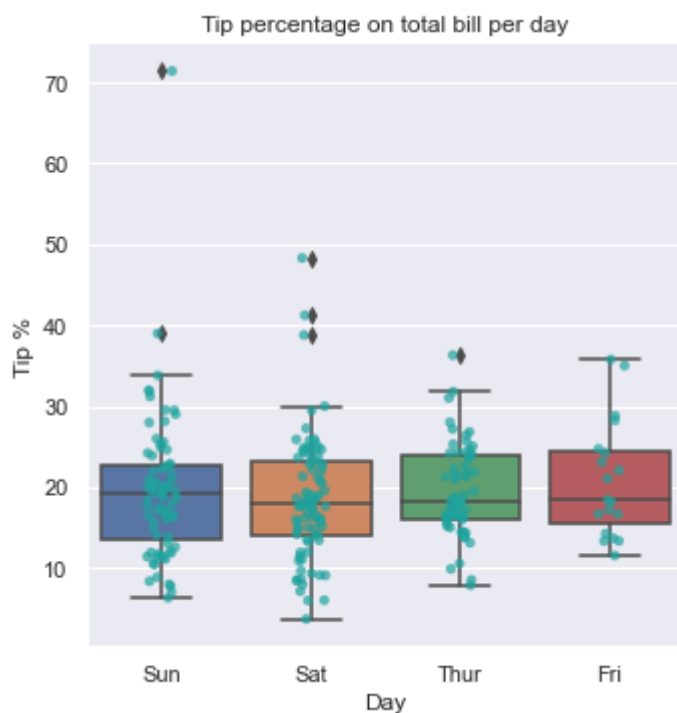
```
1 sns.catplot(x = 'day', y = 'tip_pct', kind = 'box', data = df[df.tip_pct < 245])
2 plt.title("Tip percentage on total bill per day")
3 plt.ylabel("Tip %")
4 plt.xlabel("Day")
5 plt.show()
```



Having removed that extreme outlier, the values are closer to the median, although there are day, the median does not vary much, but the bottom whisker, that means the lower values, d Saturdays and the highest is on Fridays.

In [28]:

```
1 sns.catplot(x = 'day', y = 'tip_pct', kind = 'box', data = df[df.tip_pct < 245])
2 sns.stripplot(x = 'day', y = 'tip_pct', data = df[df.tip_pct < 245], orient = 'v',
3               color = "#20A39E", alpha = 0.7)
4 plt.title("Tip percentage on total bill per day")
5 plt.ylabel("Tip %")
6 plt.xlabel("Day")
7 plt.show()
```

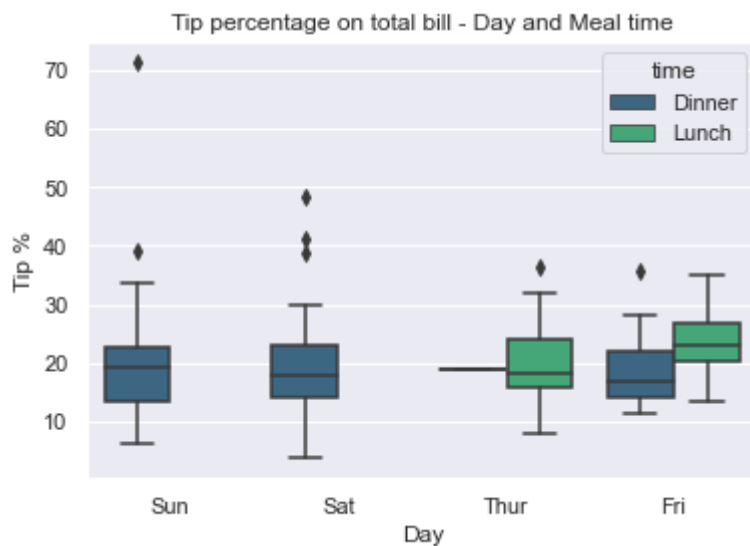


We can see the scatter shows that most of the values fall between the whiskers of the box plot

1.5 Three variables

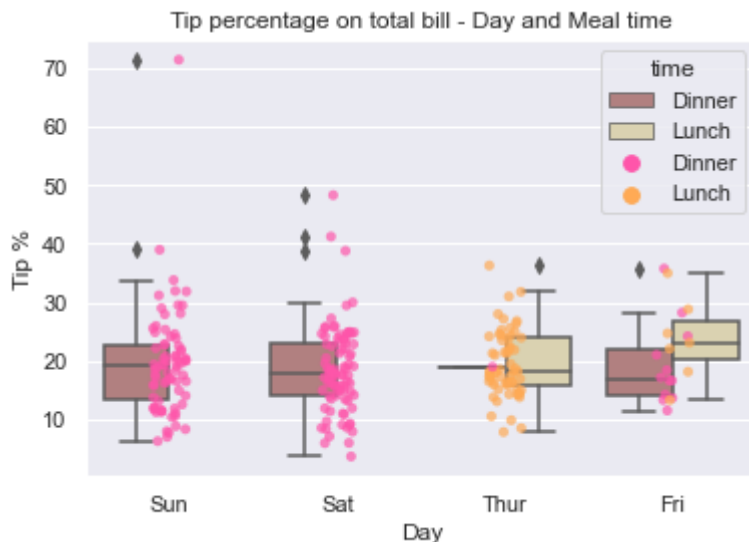
In [33]:

```
1 sns.boxplot(x = 'day', y = 'tip_pct', hue = 'time', data = df[df.tip_pct < 245], p
2 plt.title("Tip percentage on total bill - Day and Meal time")
3 plt.ylabel("Tip %")
4 plt.xlabel("Day")
5 plt.show()
```



In [41]:

```
1 sns.boxplot(x = 'day', y = 'tip_pct', hue = 'time', data = df[df.tip_pct < 245], p
2 sns.stripplot(x = 'day', y = 'tip_pct', hue = 'time', data = df[df.tip_pct < 245],
3               palette = "spring", alpha = 0.7)
4 plt.title("Tip percentage on total bill - Day and Meal time")
5 plt.ylabel("Tip %")
6 plt.xlabel("Day")
7 plt.show()
```

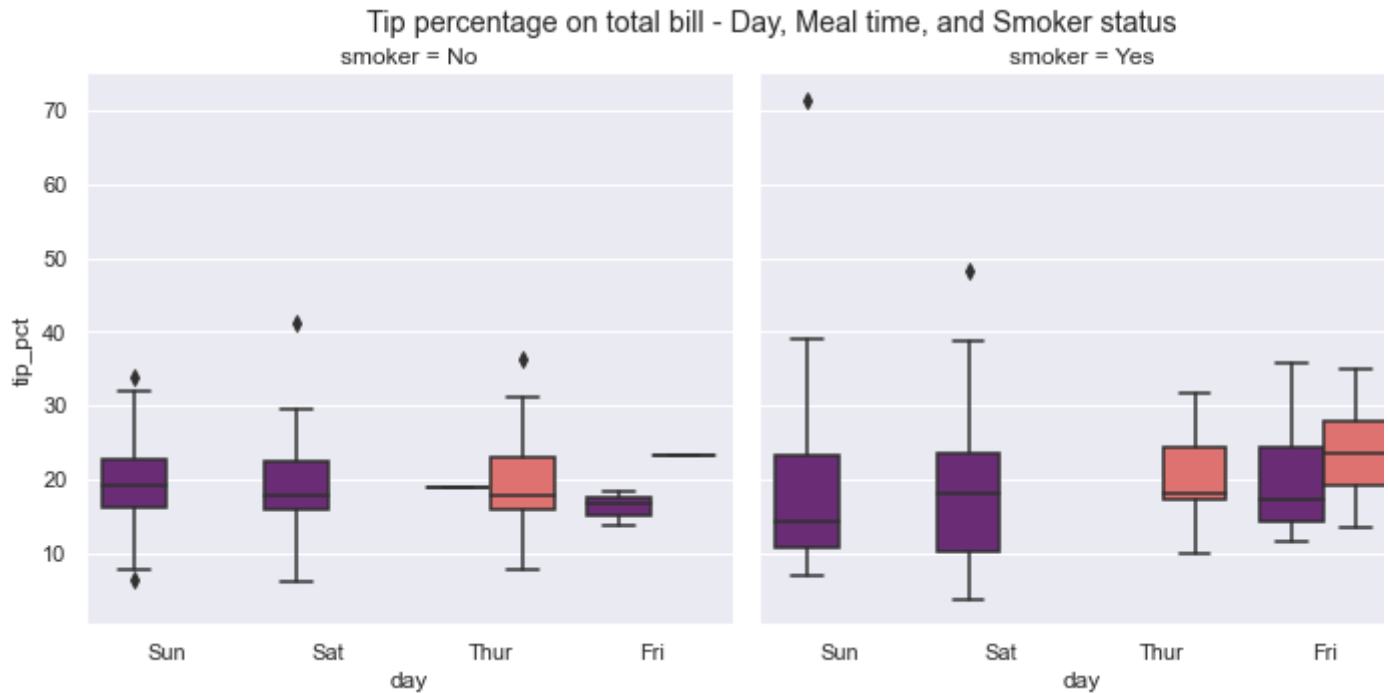


It seems there were no lunches on Sunday and Saturday, and hardly any dinners on Thursday. On average, people are more generous on a time, compared to another, but the outliers tend to be higher on dinner time.

1.6 Four variables

In [53]:

```
1 plot = sns.catplot(x = 'day', y = 'tip_pct', hue = 'time', col = 'smoker', kind =  
2 data = df[df.tip_pct < 245], palette = "magma")  
3 plot.fig.suptitle("Tip percentage on total bill - Day, Meal time, and Smoker statu  
4 plot.fig.subplots_adjust(top = 0.89)  
5 plt.show()
```



There is more variability on tip generosity on smokers than in the the non smoker group.

2 Exercici 2

Repeteix l'exercici 1 amb el dataset que disposem en el repositori de GitHub PRE-PROCES

In [13]:

```
1 movies = pd.read_table(  
2     "https://raw.githubusercontent.com/IT-Academy-BCN/Data-Science/main/Pre-proces  
3     sep = ":", engine = "python", encoding = "ISO-8859-1", header = None,  
4     names = ["id", "title", "genres"], index_col = 0)
```

In [4]:

```
1 movies.head()
```

| | title | genres |
|----|------------------------------------|------------------------------|
| id | | |
| 1 | Toy Story (1995) | Animation Children's Comedy |
| 2 | Jumanji (1995) | Adventure Children's Fantasy |
| 3 | Grumpier Old Men (1995) | Comedy Romance |
| 4 | Waiting to Exhale (1995) | Comedy Drama |
| 5 | Father of the Bride Part II (1995) | Comedy |

In [14]:

```
1 #split name and year into two columns  
2 movies[["title", "year"]] = movies["title"].str.rsplit("(", expand = True, n = 1)
```

In [15]:

```
1 #remove end parenthesis  
2 movies["year"] = movies["year"].str.rstrip(")")
```

In [16]:

```
1 #create dummy list for genre
2 dmovies = movies.genres.str.get_dummies(sep = "|")
3 dmovies.head()
```

| | Action | Adventure | Animation | Children's | Comedy | Crime | Documentary | Drama | Fantasy | Film Noi |
|----|--------|-----------|-----------|------------|--------|-------|-------------|-------|---------|-------------|
| id | | | | | | | | | | |
| 1 | 0 | 0 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 |
| 2 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 0 |
| 3 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 |
| 4 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 |
| 5 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 |

In [17]:

```
1 #add dummy list to main dataset
2 movies = movies.join(dmovies)
```

In [15]:

```
1 movies.head()
```

| | title | genres | year | Action | Adventure | Animation | Children's | Comedy | Cri |
|----|-----------------------------------|------------------------------|------|--------|-----------|-----------|------------|--------|-----|
| id | | | | | | | | | |
| 1 | Toy Story | Animation Children's Comedy | 1995 | 0 | 0 | 1 | 1 | 1 | 0 |
| 2 | Jumanji | Adventure Children's Fantasy | 1995 | 0 | 1 | 0 | 1 | 0 | 0 |
| 3 | Grumpier Old Men | Comedy Romance | 1995 | 0 | 0 | 0 | 0 | 1 | 0 |
| 4 | Waiting to Exhale | Comedy Drama | 1995 | 0 | 0 | 0 | 0 | 1 | 0 |
| 5 | Father of the Bride Part II | Comedy | 1995 | 0 | 0 | 0 | 0 | 1 | 0 |

5 rows x 21 columns

In [9]:

```
1 movies.shape
```

(3883, 21)

2.1 One variable

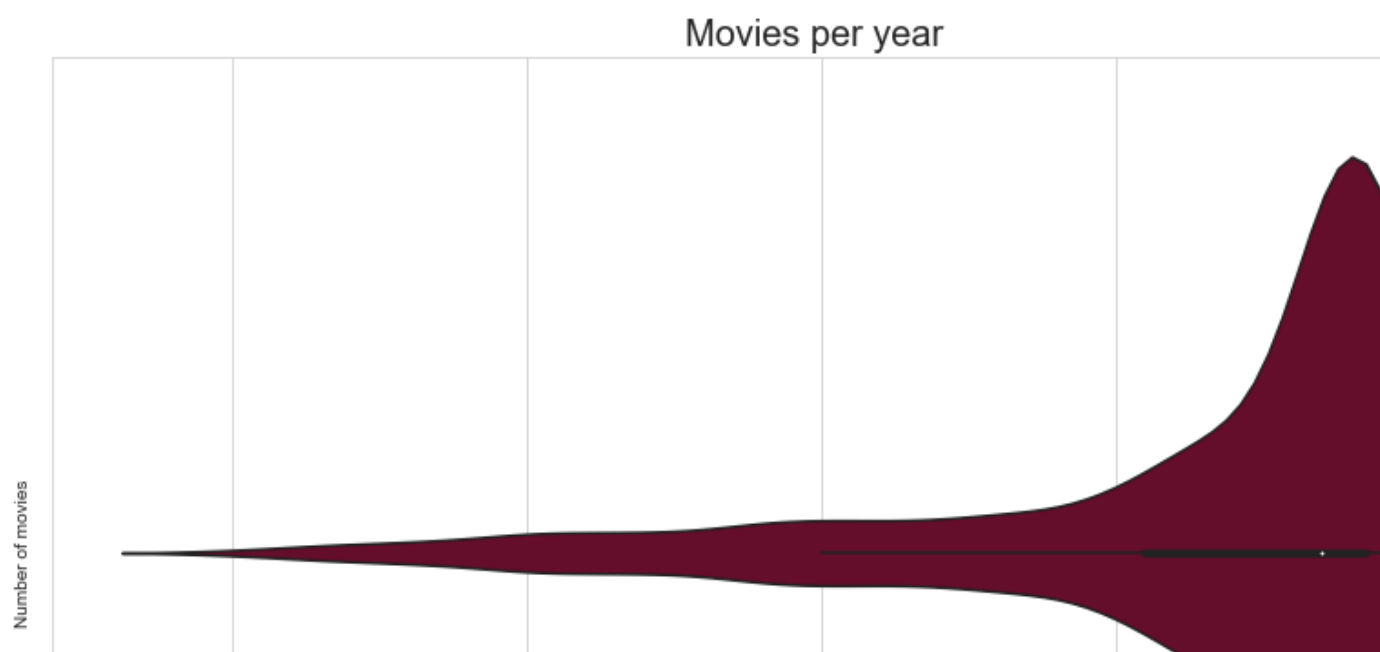
In [18]:

```
1 #make year variables into int
2 movies["year"] = movies["year"].astype(str).astype(int)
```

In [49]:

```
1 sns.set_style("whitegrid")
2 plt.figure(figsize = (15, 10))
3 sns.violinplot(x = "year", data = movies, color = "#720026")
4 plt.title("Movies per year", fontsize = 20)
5 plt.ylabel("Number of movies")
6 plt.xlabel("Year")
```

Text(0.5, 0, 'Year')



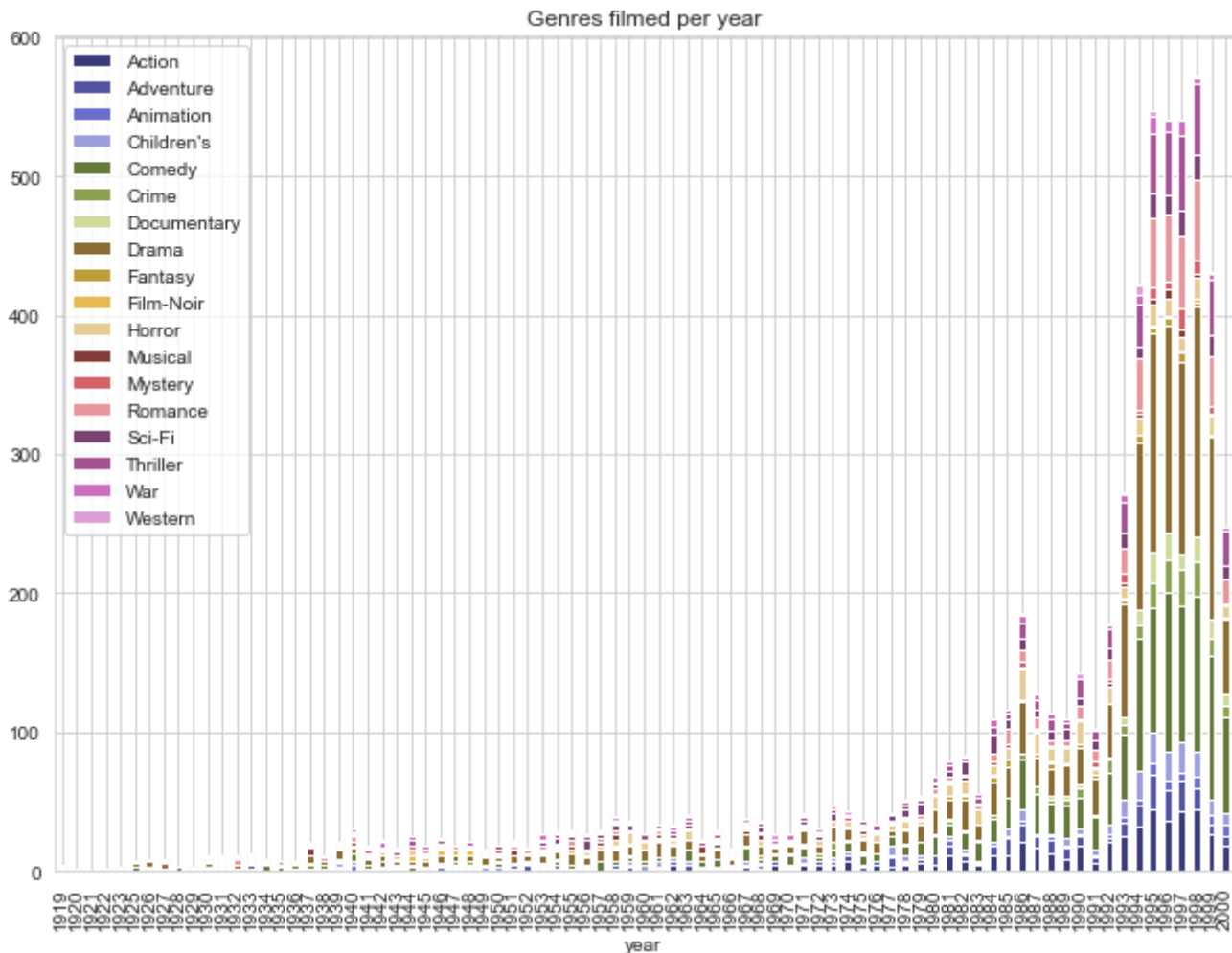
The majority of the movies in the dataset were made in the 1990s.

2.2 Two variables

In [32]:

```
1 movies.groupby("year").sum().plot(kind = "bar", stacked = True, title = "Genres fi
2                               figsize = (11, 8), colormap = "tab20b")
```

```
<AxesSubplot:title={'center':'Genres filmed per year'}, xlabel='year'>
```



It seems that throughout the years, Drama is the most filmed genre, followed by Comedy, with a rise in action movies in the 1990s.

3 Exercici 3

En aquest exercici no us donarem gaires indicacions perquè volem que ens mostreu la vostra

gràfiques i interpretacions del dataset "movies.dat" del exercici anterior.

```
In [21]:  
1 df = movies.groupby("year").sum().reset_index()
```

```
In [23]:  
1 df.tail()
```

| Children's | Comedy | Crime | Documentary | Drama | Fantasy | Film-Noir | Horror | Musical | Mystery | Romance | Sci-Fi |
|------------|--------|-------|-------------|-------|---------|-----------|--------|---------|---------|---------|--------|
| 115 | | 23 | 19 | 150 | 5 | 2 | 12 | 7 | 6 | 48 | 14 |
| 98 | | 26 | 11 | 139 | 6 | 2 | 10 | 5 | 15 | 52 | 18 |
| 112 | | 25 | 18 | 166 | 2 | 3 | 15 | 3 | 10 | 58 | 17 |
| 103 | | 12 | 15 | 130 | 2 | 0 | 14 | 1 | 5 | 37 | 15 |
| 69 | | 8 | 8 | 55 | 1 | 0 | 8 | 1 | 1 | 17 | 10 |

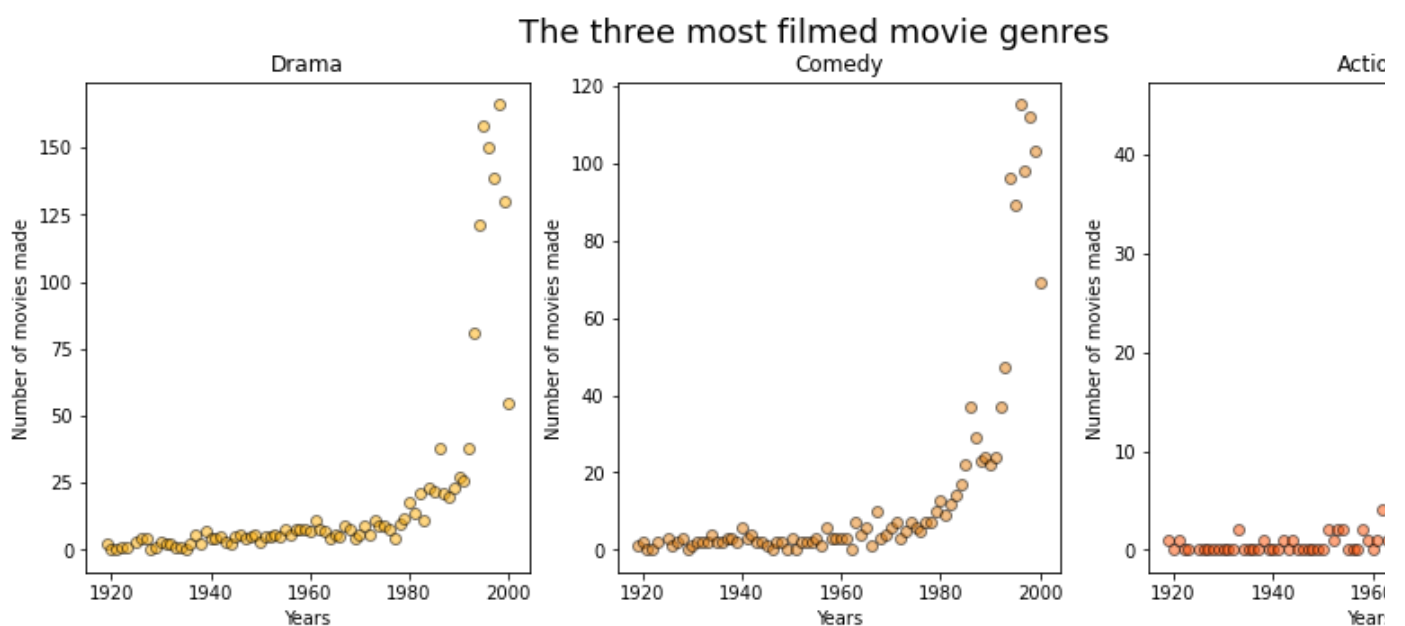
In [27]:

```

1 fig = plt.figure(figsize = (15, 5))
2 plt.suptitle("The three most filmed movie genres", size = 18)
3
4 ax1 = fig.add_subplot(1, 3, 1)
5 ax1.scatter(x = df.year, y = df.Drama, c = "#FFAE03", alpha = 0.5, edgecolors = "b")
6 plt.title("Drama")
7 plt.ylabel("Number of movies made")
8 plt.xlabel("Years")
9
10 ax2 = fig.add_subplot(1, 3, 2)
11 ax2.scatter(x = df.year, y = df.Comedy, c = "#E67F0D", alpha = 0.5, edgecolors = "b")
12 plt.title("Comedy")
13 plt.ylabel("Number of movies made")
14 plt.xlabel("Years")
15
16 ax3 = fig.add_subplot(1, 3, 3)
17 ax3.scatter(x = df.year, y = df.Action, c = "#FE4E00", alpha = 0.5, edgecolors = "b")
18 plt.title("Action")
19 plt.ylabel("Number of movies made")
20 plt.xlabel("Years")

```

Text(0.5, 0, 'Years')



In all of these genres there's an exponential rise of the number of movies, although this could be more movies from the 80s and 90s in the dataset.

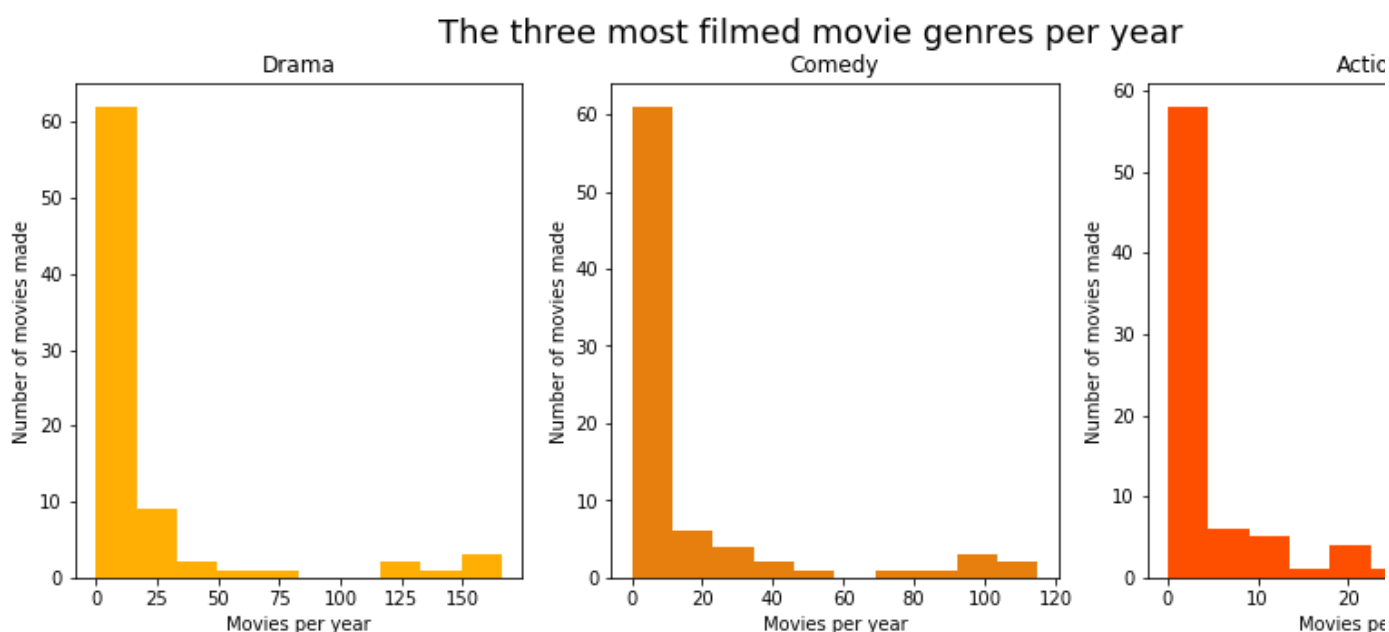
In [30]:

```

1 fig = plt.figure(figsize = (15, 5))
2 plt.suptitle("The three most filmed movie genres per year", size = 18)
3
4 ax1 = fig.add_subplot(1, 3, 1)
5 ax1.hist(x = df.Drama, color = "#FFAE03")
6 plt.title("Drama")
7 plt.ylabel("Number of movies made")
8 plt.xlabel("Movies per year")
9
10 ax2 = fig.add_subplot(1, 3, 2)
11 ax2.hist(x = df.Comedy, color = "#E67F0D")
12 plt.title("Comedy")
13 plt.ylabel("Number of movies made")
14 plt.xlabel("Movies per year")
15
16 ax3 = fig.add_subplot(1, 3, 3)
17 ax3.hist(x = df.Action, color = "#FE4E00")
18 plt.title("Action")
19 plt.ylabel("Number of movies made")
20 plt.xlabel("Movies per year")

```

```
Text(0.5, 0, 'Movies per year')
```



The high number of movies made for each genre per year is not the norm. In the Drama and less than 20 movies with those genres made. As for Action, that number falls to 10.

