

## Libraries

```
In [97]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
import plotly.express as px
import plotly.graph_objects as go
from plotly.subplots import make_subplots
```

## Load the data

```
In [22]: df=pd.read_csv('C:/Users/SLEPFL/Desktop/USA/data_analyst/udemy_dataset/udemy_courses.csv')
```

	course_id	course_title	url	is_paid	price	num_subscribers	num_reviews	num_lectures	level	content_duration	published_timestamp	subject
0	1070968	Ultimate Investment Banking Course	https://www.udemy.com/ultimate-investment-bank...	True	200	2147	23	51	All Levels	1.5	2017-01-18T20:58:56Z	Business Finance
1	1113822	Complete GST Course & Certification - Gov. You...	https://www.udemy.com/goods-and-services-tax/...	True	75	2792	923	274	All Levels	39.0	2017-03-09T16:34:20Z	Business Finance
2	1006314	Financial Modeling for Business Analysts and C...	https://www.udemy.com/financial-modeling-for-b...	True	45	2174	74	51	Intermediate Level	2.5	2016-12-19T19:26:30Z	Business Finance
3	1212580	Beginner to Pro - Financial Analysis in Excel ...	https://www.udemy.com/complete-excel-finance-c...	True	95	2451	11	36	All Levels	3.0	2017-05-30T20:07:24Z	Business Finance
4	1010568	How to Maximize Your Profits Trading Options	https://www.udemy.com/how-to-maximize-your-pro...	True	200	1276	45	26	Intermediate Level	2.0	2016-12-13T14:57:18Z	Business Finance
5	192970	Trading Penny Stocks: A Guide for All Levels L...	https://www.udemy.com/trading-penny-stocks-a-0...	True	150	9221	138	25	All Levels	3.0	2014-05-02T15:13:30Z	Business Finance
6	739964	Investing And Trading For Beginners: Mastering...	https://www.udemy.com/investing-and-trading-fo...	True	65	1540	178	26	Beginner Level	1.0	2016-02-21T18:23:12Z	Business Finance
7	403100	Trading Stock Chart Patterns For Immediate, Ex...	https://www.udemy.com/trading-chart-patterns-f...	True	95	2917	148	23	All Levels	2.5	2016-01-30T22:13:03Z	Business Finance
8	476268	Options Trading 3: Advanced Stock Profit and ...	https://www.udemy.com/day-trading-stock-optio...	True	195	5172	34	38	Expert Level	2.5	2015-05-28T00:14:03Z	Business Finance
9	1167710	The Only Investment Strategy You Need For Your...	https://www.udemy.com/the-only-investment-str...	True	200	827	14	15	All Levels	1.0	2017-04-18T18:13:32Z	Business Finance

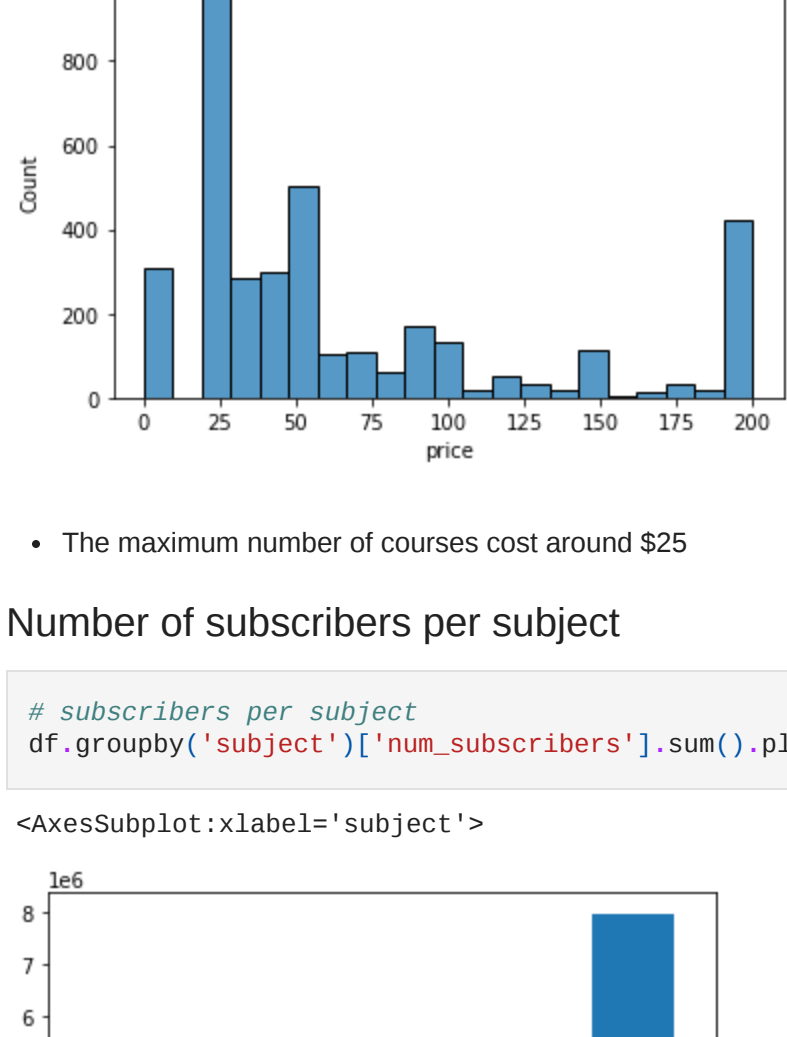
```
In [42]: df['published_date']>df['published_timestamp'].str.split('').str.get(8)
df['published_date'].pd.to_datetime(df['published_date'], format='%Y-%m-%d')
df['year']=df['published_date'].dt.year
```

### What is the distribution of subjects

```
In [116]: # Distribution of subjects
df['subject'].unique()
```

```
Out[116]: array(['Business Finance', 'Graphic Design', 'Musical Instruments',
       'Web Development'], dtype=object)
```

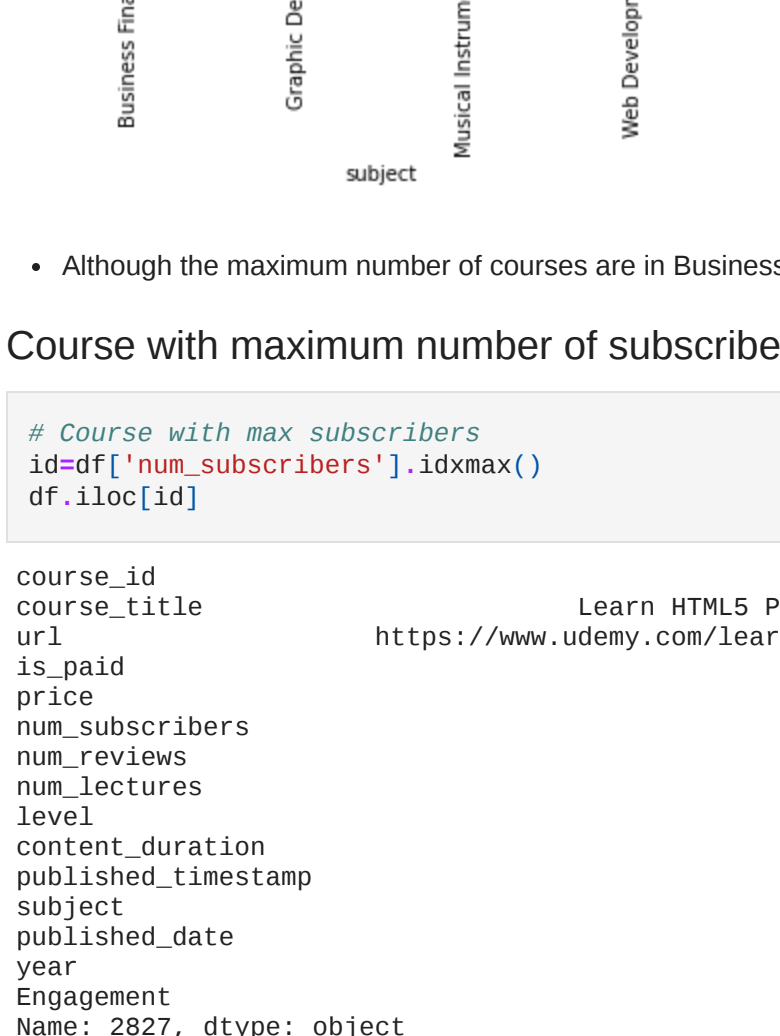
```
In [115]: # Distribution of subjects(course-wise)
print(df['subject'].value_counts())
plt.figure(figsize=(5,3))
df['subject'].value_counts().plot(kind='pie', autopct='%1.1f%%', shadow=True)
plt.show()
```



- It can be seen here, that 4 subjects in total are present and most number of courses are in Web Development

### Price for different courses

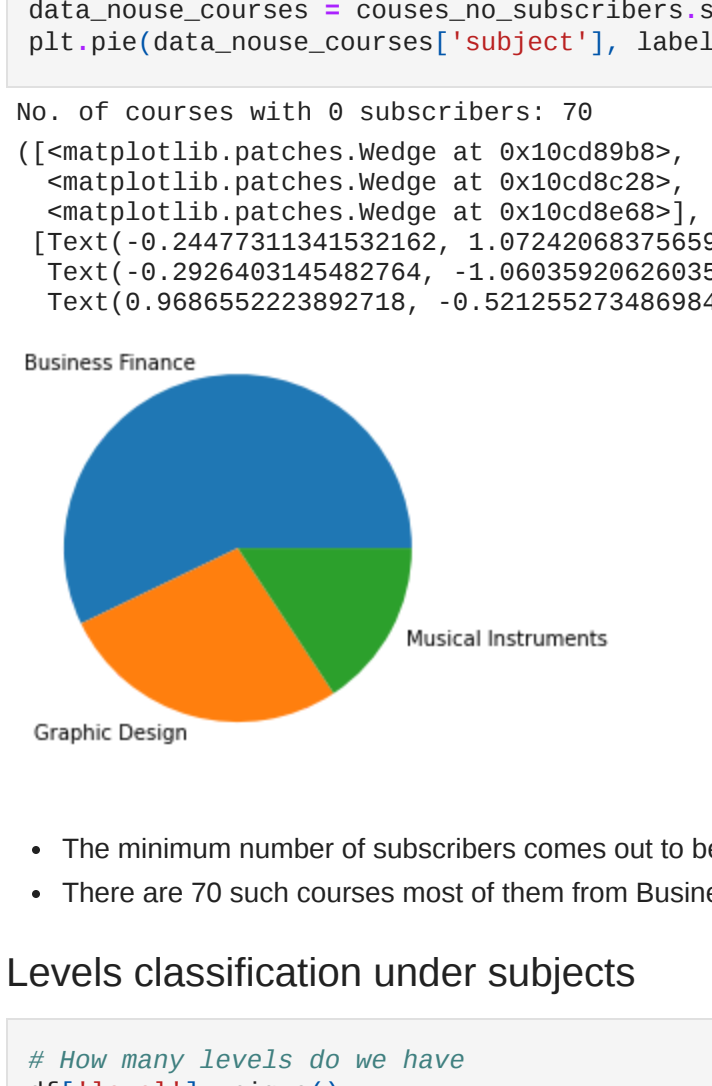
```
In [119]: # Price range of courses
sns.histplot(data=df, x='price')
```



- The maximum number of courses cost around \$25

### Number of subscribers per subject

```
In [91]: # subscribers per subject
df.groupby('subject')['num_subscribers'].sum().plot(kind='bar')
```



- Although the maximum number of courses are in Business Finance, the number of subscribers are maximum in Web Development.

### Course with maximum number of subscribers

```
In [101]: # Course with max subscribers
idx=df['num_subscribers'].idxmax()
df.iloc[idx]
```

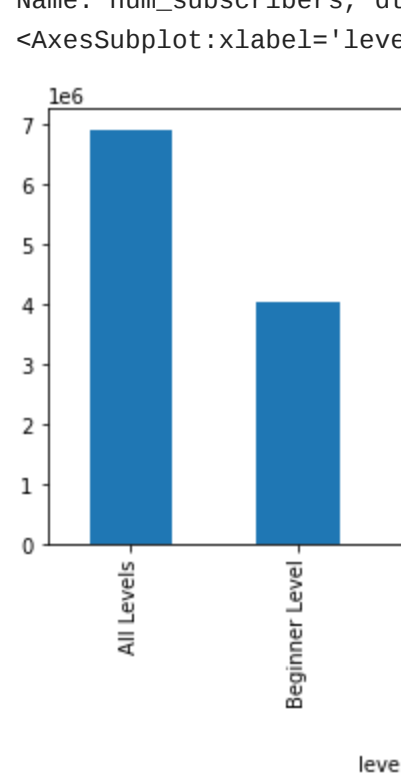
```
Out[101]: course_id      41295
course_title      Learn HTML5 Programming From Scratch
url               https://www.udemy.com/learn-html5-programming-...
is_paid           False
price              8
num_subscribers   268923
num_reviews       8629
num_lectures      45
level             All Levels
content_duration  18.5
published_timestamp 2013-02-14T07:43:41Z
subject           Web Development
published_date     2813-02-14 08:09:08
year              2013
engagement        2913
Name: 2897, dtype: object
```

- course with maximum number of subscribers comes out to be Web Development course titled 'Learn HTML5 Programming From Scratch' with 268923 subscribers

### Course with minimum number of subscribers

```
In [102]: # Courses with 0 subscribers
courses_no_subscribers = df.loc[df.num_subscribers == 0]
print('No. of courses with 0 subscribers: ',str(courses_no_subscribers.num_subscribers.value_counts()[0]))
data_nouse_courses = courses_no_subscribers.subject.value_counts().reset_index()
plt.plot(data_nouse_courses['subject'], labels = data_nouse_courses['index'])
```

```
Out[102]: No. of courses with 0 subscribers: 78
[<matplotlib.patches.Wedge at 8x18cd898ab>,
 <matplotlib.patches.Wedge at 8x18cd8c2b>,
 <matplotlib.patches.Wedge at 8x18cd868d>]
[Text(0.2447311345152162, 1.0724266837565984, 'Business Finance'),
 Text(-0.2505460134527164, -1.06893926093025, 'Graphic Design'),
 Text(0.9686552223892718, -0.5212652734889842, 'Musical Instruments')]
```



- The minimum number of subscribers comes out to be 0.
- There are 78 such courses most of them from Business Finance. It can be noted that each Web development course has atleast one subscriber

### Levels classification under subjects

```
In [45]: # How many levels do we have
df['level'].unique()
```

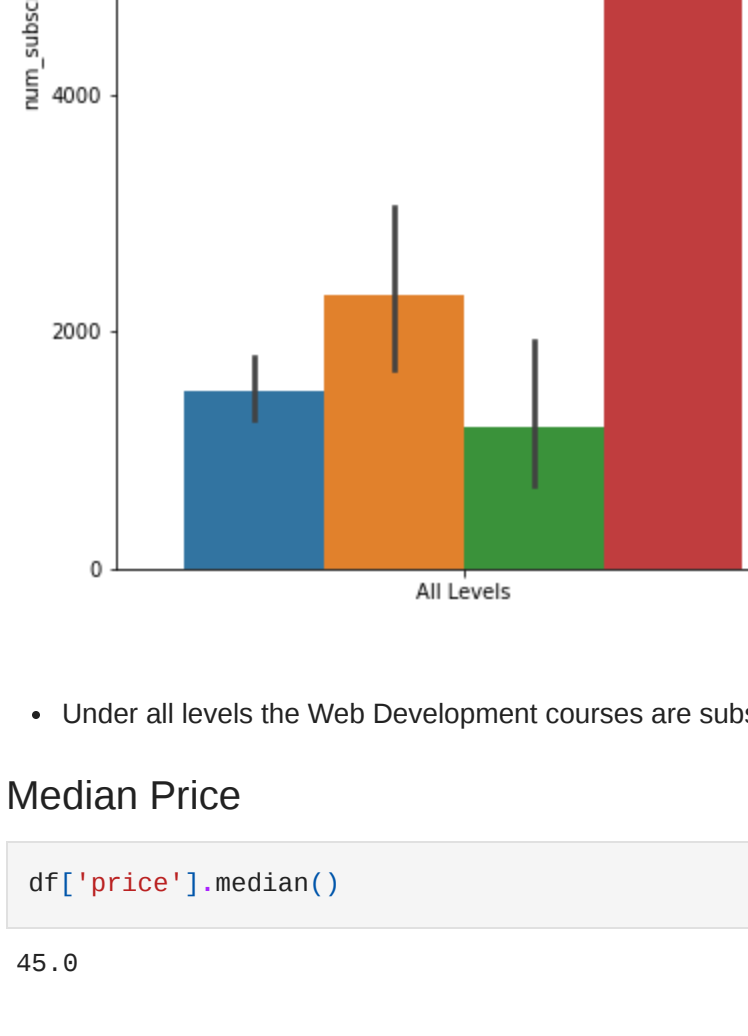
```
Out[45]: array(['All Levels', 'Intermediate Level', 'Beginner Level',
       'Expert Level'], dtype=object)
```

- We have 4 different levels:

- All Levels
- Intermediate Levels
- Beginner Levels
- Expert Level

### Number of courses per level

```
In [48]: # distribution of courses per level
df['level'].value_counts().plot(kind='bar')
```

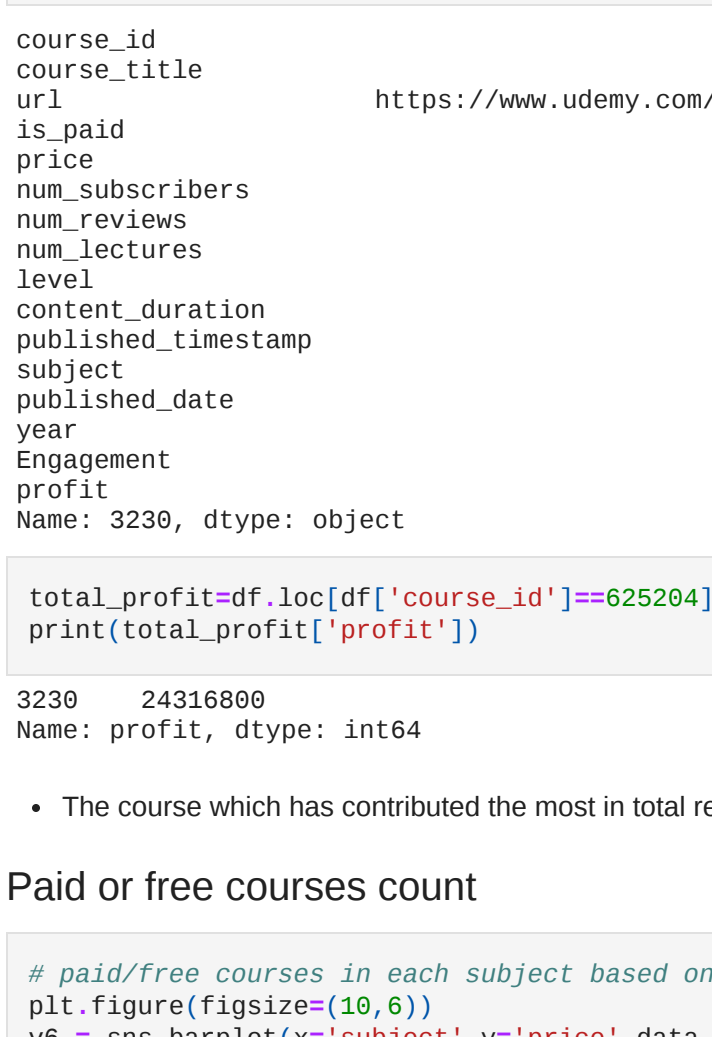


- Maximum number of courses lie under All Levels and there are minimum courses for Expert Level

### Subscribers for each level

```
In [58]: # How many subscribers per level
print(df.groupby('level')['num_subscribers'].sum())
df.groupby('level')['num_subscribers'].sum().plot(kind = 'bar')
```

```
Out[58]: level
All Levels      6915676
Beginner Level  4655843
Expert Level    58286
Intermediate Level 42085
Name: num_subscribers, dtype: int64
```



- Maximum number of courses are subscribed under All Level and the Expert Level courses are least subscribed

### Subscribers for each subject under different levels

```
In [62]: # subscribers comparison per subject for each level
df.groupby('subject')['level'].value_counts()
plt.figure(figsize=(25,18))
sns.barplot(x='level', y='num_subscribers', hue='subject', data=df)
```



- Under all levels the Web Development courses are subscribed the most

### Median Price

```
In [121]: df['price'].median()
```

```
Out[121]: 45.8
```

- The median price comes out to be \$45

### Total revenue collected

```
In [109]: # total earning
df['profit']=df['price']*df['num_subscribers']
sum=df['profit'].sum()
```

```
Out[109]: 884921315
```

- The total profit comes out to be \$884921315 till date for all the udemy courses purchased.

### Highest profitable course

```
In [115]: #highest profitable course
idx=df['profit'].idxmax()
df.iloc[idx]
```

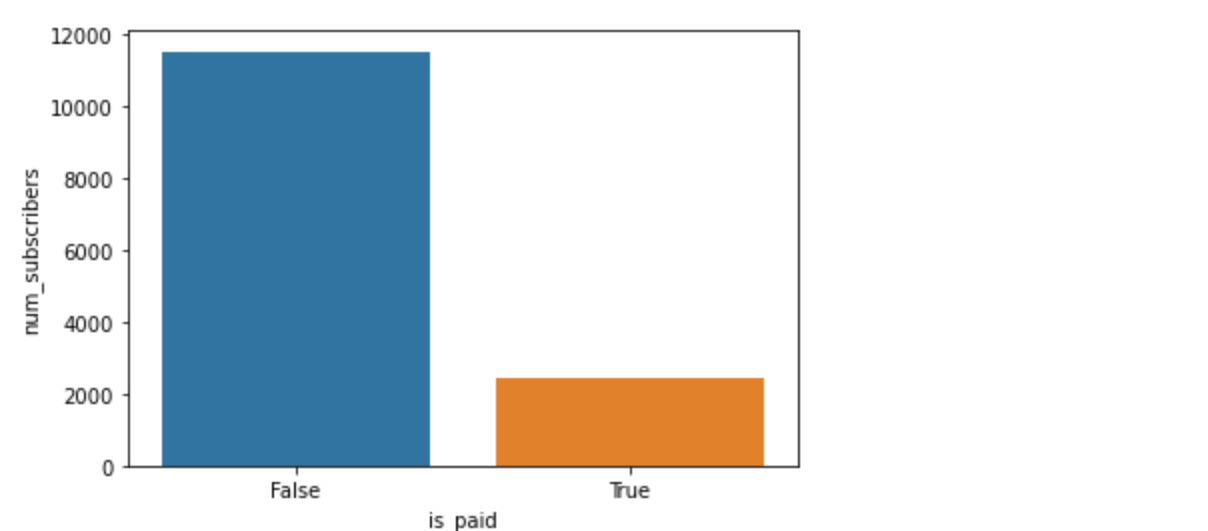
```
Out[115]: course_id      625204
course_title      The Web Developer Bootcamp
url               https://www.udemy.com/the-web-developer-bootcamp/...
is_paid           True
price              200
num_subscribers   121094
num_reviews       27445
num_lectures      342
level             All Levels
content_duration  42.8
published_timestamp 2015-11-02T21:13:27Z
subject           Web Development
published_date     2915-11-02 08:58:08
year              2015
engagement        149029
profit             24316800
Name: 3238, dtype: object
```

- The course which has contributed the most in total revenue is 'The Web Developer Bootcamp' and has brought a revenue of \$24316800

### Paid or free courses count

```
In [92]: # paid/free courses in each subject based on level
plt.figure(figsize=(18,6))
ve = sns.barplot(x='subject', y='price', data = df, hue = 'level', ci=False)
```

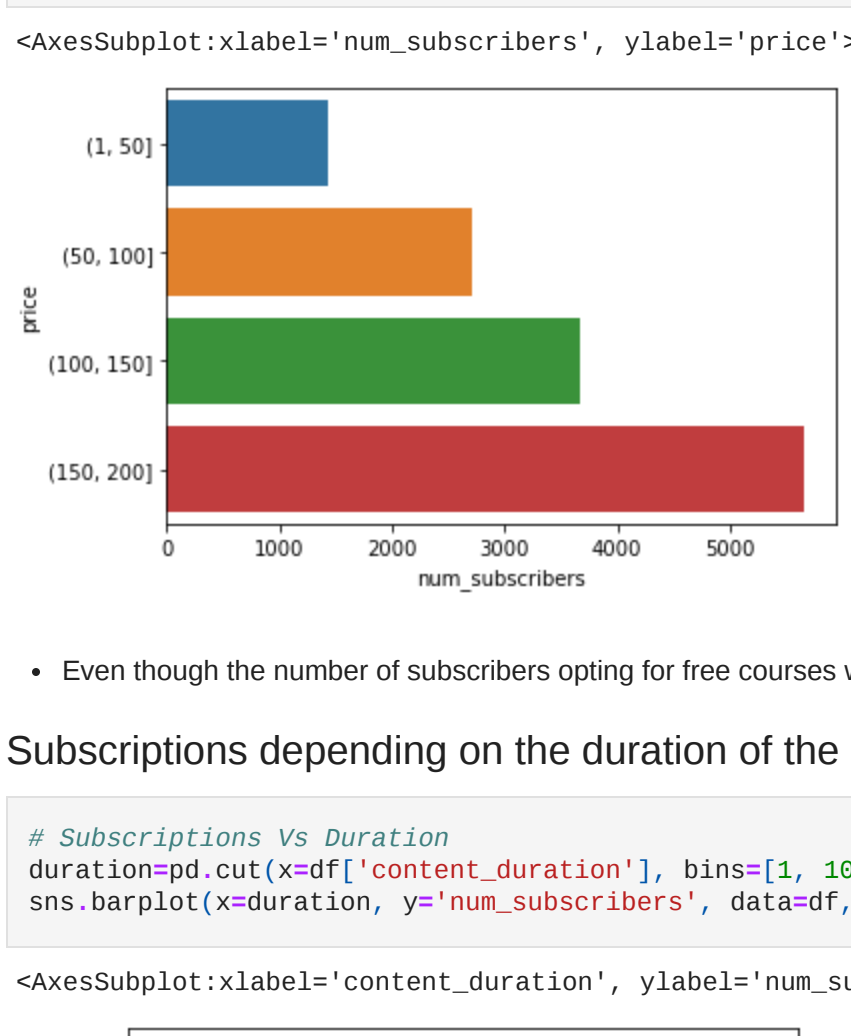
```
# count of paid/unpaid courses
price_list=df['is_paid'].unique()
price_count=df['is_paid'].value_counts().reset_index()
fig=px.bar(price_count, x='index', y='is_paid', text='is_paid', color='is_paid',
           title='level count of courses paid and unpaid for ',
           labels={'index':'paid/unpaid courses', 'is_paid':'count of paid/unpaid courses'})
fig.update_layout(showlegend=False, width=600)
fig.show()
```



- Number of paid courses are 3398 and free courses are 310
- The Expert Level course in all the subjects leaving Musical Instruments subject are the most expensive ones.

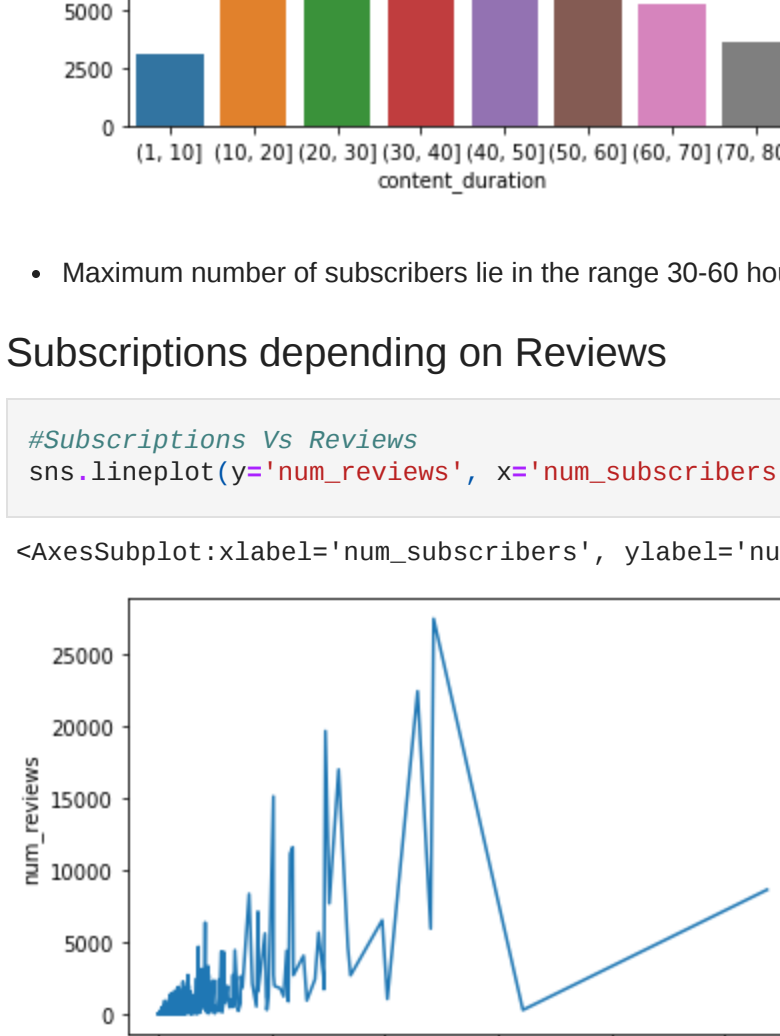
### Subscriptions dependency on price

```
In [96]: # price vs user
sns.scatterplot(data=df, x='price', y='num_subscribers')
```



### Subscribers dependencies on course is paid/free

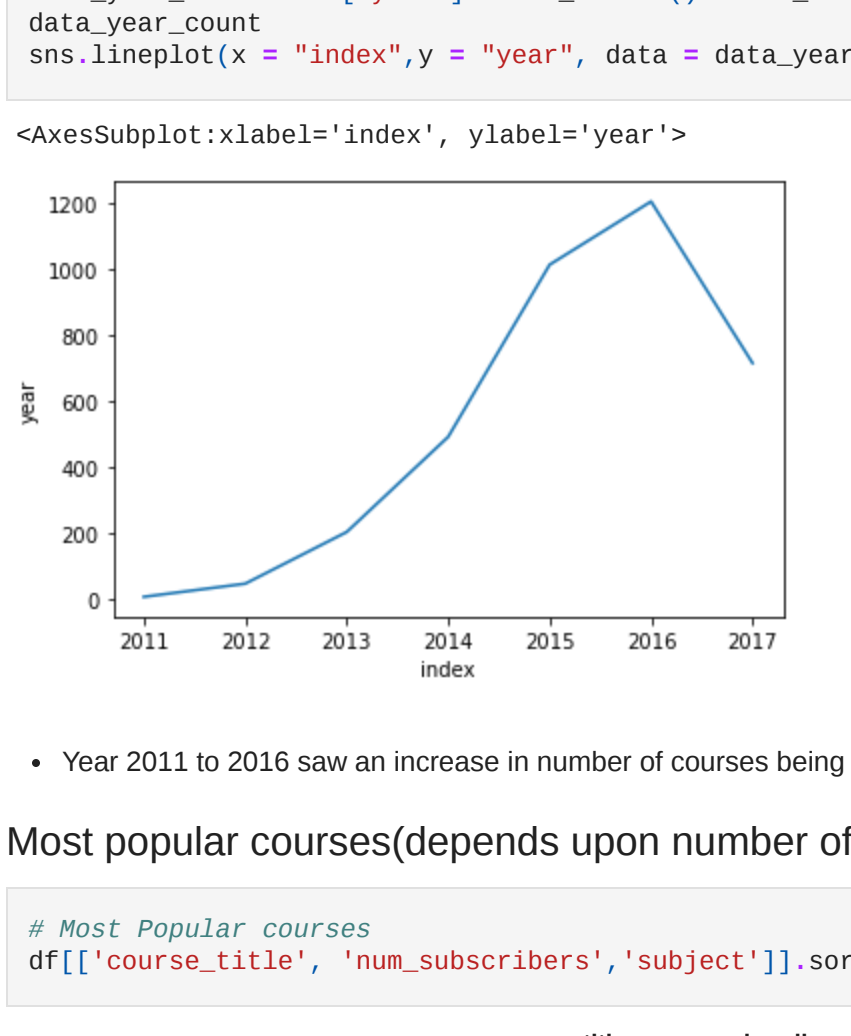
```
In [91]: # subscription Vs paid
sns.barplot(x='is_paid', y='num_subscribers', data=df, ci=False)
```



- Subscription is more for free courses rather than the paid ones

### Subscriptions in different price ranges

```
In [36]: # Subscriptions Vs price
price_range_cut=df['price'].bins[1, 50, 100, 150, 200]
sns.barplot(y=price, x='num_subscribers', data=df, ci=False)
```



- Even though the number of subscribers opting for free courses was much higher than the paid ones, those who opted for the paid ones chose the course in range \$150-200 the most expensive range.

### Subscriptions depending on the duration of the course

```
In [27]: # Subscriptions Vs Duration
duration_range_cut=df['content_duration'].bins[1, 10, 20, 30, 40, 50, 60, 70, 80]
sns.barplot(x=duration, y='num_subscribers', data=df, ci=False)
```



- Maximum number of subscribers lie in the range 30-60 hours range

### Subscriptions depending on Reviews

```
In [38]: # Subscriptions Vs Reviews
sns.lineplot(y='num_reviews', x='num_subscribers', data=df, ci=False)
```



- No particular pattern can be seen here

### Courses published per year

```
In [69]: # Publications per year
data_year_count = df['year'].value_counts().reset_index()
data_year_count
sns.lineplot(x = "index", y = "year", data = data_year_count)
```



- Year 2011 to 2016 saw an increase in number of courses being published every year. But there was a sudden fall in 2017.

### Most popular courses(depends upon number of subscribers)

```
In [71]: # Most Popular courses
df[['course_title', 'df.num_subscribers', 'subject']].sort_values('num_subscribers', ascending = False).head(5)
```

	course_title	num_subscribers	subject
2827	Learn HTML5 Programming From Scratch	268923	Web Development
3032	Coding for Entrepreneurs Basic	161326	Web Development
3230	The Web Developer Bootcamp	121584	Web Development
2783	Build Your First Website in 1 Week with HTML5 ...	120291	Web Development
3232	The Complete Web Developer Course 2.0	114612	Web Development

- All of the courses belong to web development
- Most engaging courses(depends on number of reviews and subscribers)

```
In [89]: # Most engaging Courses
df[['engagement', 'df.num_subscribers']] = df['num_reviews']
df[['course_title', 'engagement', 'subject']].sort_values('engagement', ascending = False).head(5)
```

	course_title	Engagement	subject
2827	Learn HTML5 Programming From Scratch	277552	Web Development
3032	Coding for Entrepreneurs Basic	161326	Web Development
3230	The Web Developer Bootcamp	149029	Web Development
3232	The Complete Web Developer Course 2.0	118624	Web Development
2783	Build Your First Website in 1 Week with HTML5 ...	116215	Web Development

- The most engaging course belongs to Web Development.

### Best cost beneficial course

```
In [130]: # What courses offer the best cost benefit
# we can determine this by finding the course with the least amount paid, but has the highest number of subscribers.
cost_benefit = df[(df['price']<df['price'].mean()) &
(df['num_subscribers']>df['num_subscribers'].mean())].sort_values('num_subscribers', ascending=False)[['course_title']].head(1).unique()
```

The course which offers the best cost benefit is : 'Learn HTML5 Programming From Scratch'

- 'Learn HTML5 Programming From Scratch' comes out to be the most beneficial course.

```
In [130]: jupyter nbconvert --execute --to pdf udemy_analysis.ipynb
SyntaxError: invalid syntax
```