

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

Load Data

```
In [2]: 1 df=pd.read_csv("udemy_courses.csv")
```

Basic data exploration

```
In [3]: 1 df.head()
```

Out[3]:

	course_id	course_title	url	is_paid	price	num_subscribers	num_re
0	1070968	Ultimate Investment Banking Course	https://www.udemy.com/ultimate-investment-bank...	True	200	2147	
1	1113822	Complete GST Course & Certification - Grow You...	https://www.udemy.com/goods-and-services-tax/	True	75	2792	
2	1006314	Financial Modeling for Business Analysts and C...	https://www.udemy.com/financial-modeling-for-b...	True	45	2174	
3	1210588	Beginner to Pro - Financial Analysis in Excel ...	https://www.udemy.com/complete-excel-finance-c...	True	95	2451	
4	1011058	How To Maximize Your Profits Trading Options	https://www.udemy.com/how-to-maximize-your-pro...	True	200	1276	

```
In [4]: 1 df.level.value_counts()
```

```
Out[4]: All Levels          1929
Beginner Level             1270
Intermediate Level         421
Expert Level                58
Name: level, dtype: int64
```

```
In [5]: 1 ##value counts of all level of course
        2 df.level.value_counts()
```

```
Out[5]: All Levels          1929
        Beginner Level      1270
        Intermediate Level   421
        Expert Level        58
        Name: level, dtype: int64
```

Basic Data Cleaning

```
In [6]: 1 ##Use LabelEncoder for encoding of is_paid column
        2 from sklearn.preprocessing import LabelEncoder
        3 labelencoder = LabelEncoder()
        4 df['is_paid'] = labelencoder.fit_transform(df['is_paid'])
        5 df.head()
```

```
Out[6]:
```

	course_id	course_title	url	is_paid	price	num_subscribers	num_re
0	1070968	Ultimate Investment Banking Course	https://www.udemy.com/ultimate-investment-bank...	1	200	2147	
1	1113822	Complete GST Course & Certification - Grow You...	https://www.udemy.com/goods-and-services-tax/	1	75	2792	
2	1006314	Financial Modeling for Business Analysts and C...	https://www.udemy.com/financial-modeling-for-b...	1	45	2174	
3	1210588	Beginner to Pro - Financial Analysis in Excel ...	https://www.udemy.com/complete-excel-finance-c...	1	95	2451	
4	1011058	How To Maximize Your Profits Trading Options	https://www.udemy.com/how-to-maximize-your-pro...	1	200	1276	

```
In [7]: 1  ##Custom LabelEncoder for encoding level column
2  cleanup_nums = {"level":      {"All Levels": 0, "Beginner Level": 1, "Intermed
3                                     "Expert Level":3}}
4  df = df.replace(cleanup_nums)
5  df.head()
```

Out[7]:

	course_id	course_title	url	is_paid	price	num_subscribers	num_re
0	1070968	Ultimate Investment Banking Course	https://www.udemy.com/ultimate-investment-bank...	1	200	2147	
1	1113822	Complete GST Course & Certification - Grow You...	https://www.udemy.com/goods-and-services-tax/	1	75	2792	
2	1006314	Financial Modeling for Business Analysts and C...	https://www.udemy.com/financial-modeling-for-b...	1	45	2174	
3	1210588	Beginner to Pro - Financial Analysis in Excel ...	https://www.udemy.com/complete-excel-finance-c...	1	95	2451	
4	1011058	How To Maximize Your Profits Trading Options	https://www.udemy.com/how-to-maximize-your-pro...	1	200	1276	

```
In [8]: 1  df.isnull().sum()
2  ##no null value in the data set
```

```
Out[8]: course_id      0
course_title    0
url             0
is_paid         0
price           0
num_subscribers 0
num_reviews     0
num_lectures    0
level           0
content_duration 0
published_timestamp 0
subject         0
dtype: int64
```

```
In [9]: 1  ##showing datatypes of all column  
        2  df.dtypes
```

```
Out[9]: course_id          int64  
        course_title      object  
        url               object  
        is_paid           int64  
        price             int64  
        num_subscribers    int64  
        num_reviews        int64  
        num_lectures       int64  
        level             int64  
        content_duration   float64  
        published_timestamp object  
        subject            object  
        dtype: object
```

```
In [10]: 1  df.is_paid.value_counts()
```

```
Out[10]: 1    3368  
        0    310  
        Name: is_paid, dtype: int64
```

```
In [11]: 1  df.subject.value_counts()
```

```
Out[11]: Web Development      1200  
        Business Finance      1195  
        Musical Instruments     680  
        Graphic Design         603  
        Name: subject, dtype: int64
```

```
In [12]: 1  df.level.value_counts()
```

```
Out[12]: 0    1929  
        1    1270  
        2     421  
        3      58  
        Name: level, dtype: int64
```

All Tasks And Answers

```
In [13]: 1  ## creating a new column name per_course_total_revenue  
        2  df["per_course_total_revenue"]=df["price"]*df["num_subscribers"]
```

In [14]:

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

Out[14]:

	course_id	course_title	url	is_paid	price	num_subscribers	num_re
0	1070968	Ultimate Investment Banking Course	https://www.udemy.com/ultimate-investment-bank...	1	200	2147	
1	1113822	Complete GST Course & Certification - Grow You...	https://www.udemy.com/goods-and-services-tax/	1	75	2792	
2	1006314	Financial Modeling for Business Analysts and C...	https://www.udemy.com/financial-modeling-for-b...	1	45	2174	
3	1210588	Beginner to Pro - Financial Analysis in Excel ...	https://www.udemy.com/complete-excel-finance-c...	1	95	2451	
4	1011058	How To Maximize Your Profits Trading Options	https://www.udemy.com/how-to-maximize-your-pro...	1	200	1276	

In [15]:

```
1 ##for subject
2 ##Which subject got highest revenue?
3 df.loc[df['per_course_total_revenue'].idxmax()]
4 ##we see that webdevelopment got highest revenue at about rupees 24316800/-
```

Out[15]:

```
course_id          625204
course_title       The Web Developer Bootcamp
url                https://www.udemy.com/the-web-developer-bootcamp/
                  (https://www.udemy.com/the-web-developer-bootcamp/)
is_paid            1
price              200
num_subscribers    121584
num_reviews        27445
num_lectures       342
level              0
content_duration   43
published_timestamp 2015-11-02T21:13:27Z
subject            Web Development
per_course_total_revenue 24316800
Name: 3230, dtype: object
```

```
In [16]: 1  ##See all column's unique values
2  def unique_value(df):
3      for column in df:
4          print(f'{column}:{df[column].unique()}')
5
6  ##unique_value(df)
```

```
In [17]: 1  ## All subjects name and their total courses counts
2  df.subject.value_counts()
```

```
Out[17]: Web Development      1200
Business Finance      1195
Musical Instruments      680
Graphic Design      603
Name: subject, dtype: int64
```

Revenue differences among all subjects with the help of pie chart

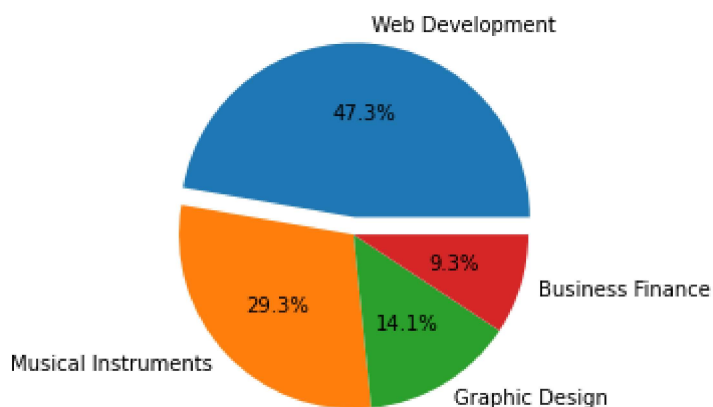
```
In [18]: 1  web=df[df['subject']=="Web Development"]
2  web_max=web.loc[web['per_course_total_revenue'].idxmax()]
3
4  Business_fin=df[df['subject']=="Business Finance"]
5  Business_fin_max=Business_fin.loc[Business_fin['per_course_total_revenue'].i
6
7  Musical_Instruments=df[df['subject']=="Musical Instruments"]
8  Musical_Instruments_max=Musical_Instruments.loc[Musical_Instruments['per_cou
9
10 Graphic_Design=df[df['subject']=="Graphic Design"]
11 Graphic_Design_max=Graphic_Design.loc[Graphic_Design['per_course_total_reven
12
13 ##concatination of all this dataframes
14 revenue_dif=pd.concat([web_max,Business_fin_max,Musical_Instruments_max,Grap
15 revenue_dif=revenue_dif.T
16
```

```
In [19]: 1 ##Showing the values from high to Low
2 revenue_dif=revenue_dif.sort_values(by=['per_course_total_revenue'],ascending=False)
3 revenue_dif.head()
```

Out[19]:

	course_id	course_title	url	is_paid	price	num_subscribers	n
3230	625204	The Web Developer Bootcamp	https://www.udemy.com/the-web-developer-bootcamp/	1	200	121584	
1979	238934	Pianoforall - Incredible New Way To Learn Pian...	https://www.udemy.com/pianoforall-incredible-n...	1	200	75499	
1213	820194	Photoshop for Entrepreneurs - Design 11 Practi...	https://www.udemy.com/photoshop-for-entreprene...	1	200	36288	
40	648826	The Complete Financial Analyst Course 2017	https://www.udemy.com/the-complete-financial-a...	1	195	24481	

```
In [20]: 1 ## plotting pie chart acoording to max
2 sujet=['Web Development','Musical Instruments','Graphic Design','Business F
3 exp=[0.1,0,0,0]
4 plt.pie(revenue_dif.per_course_total_revenue,labels=sujet,autopct='%2.1f%%'
5 plt.show()
```



```
In [21]: 1 w,b,m,g=df.subject.value_counts()
```

In [22]: 1 df.subject.value_counts()

Out[22]: Web Development 1200
 Business Finance 1195
 Musical Instruments 680
 Graphic Design 603
 Name: subject, dtype: int64

In [23]: 1 df['per_course_total_revenue'].dtypes

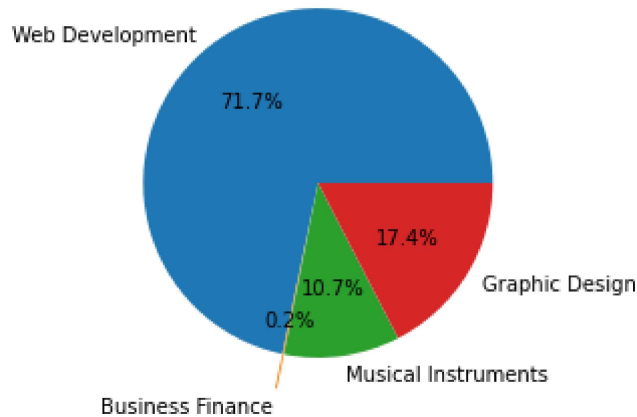
Out[23]: dtype('int64')

In [24]: 1 *## For Average*
 2 web=df[df['subject']=='Web Development']['per_course_total_revenue'].sum()
 3 web=web/w
 4 bus=df[df['subject']=='Business Finance']['per_course_total_revenue'].sum()
 5 bus=b
 6 mus=df[df['subject']=='Musical Instruments']['per_course_total_revenue'].sum()
 7 mus=mus/m
 8 grap=df[df['subject']=='Graphic Design']['per_course_total_revenue'].sum()
 9 grap=grap/g
 10 data=[web,bus,mus,grap]
 11 subject=['Web Development','Business Finance','Musical Instruments','Graphic Design']
 12 index=[0,1,2,3]
 13 q1=pd.DataFrame(data,index)
 14 headers = ["Avg_revenue_per_subject"]
 15 q1.columns = headers
 16 q1['subject']=subject
 17
 18 q1.head()

Out[24]:

	Avg_revenue_per_subject	subject
0	525703.145833	Web Development
1	1195.000000	Business Finance
2	78469.198529	Musical Instruments
3	127666.948590	Graphic Design


```
In [25]: 1  ## plotting pie chart according to Average
2  subject1=['Web Development','Business Finance','Musical Instruments','Graphic
3  exp=[0,0.2,0,0]
4  plt.pie(q1.Avg_revenue_per_subject ,labels=subject1,autopct='%2.1f%%',explod
5  plt.show()
```



Which course is popular according to reviews and create a graph of that distribution

```
In [26]: 1  ## Highest rating course
2  df.loc[df['num_reviews'].idxmax()]
```

```
Out[26]: course_id                625204
course_title                The Web Developer Bootcamp
url            https://www.udemy.com/the-web-developer-bootcamp/
            (https://www.udemy.com/the-web-developer-bootcamp/)
is_paid                      1
price                        200
num_subscribers             121584
num_reviews                 27445
num_lectures                342
level                       0
content_duration             43
published_timestamp          2015-11-02T21:13:27Z
subject                     Web Development
per_course_total_revenue      24316800
Name: 3230, dtype: object
```

```

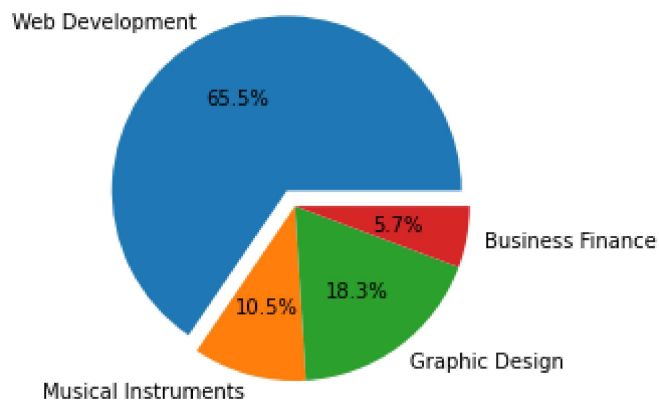
In [27]: 1  ## for counting max
          2  web=df[df['subject']=="Web Development"]
          3  web_review_max=web.loc[web['num_reviews'].idxmax()]
          4
          5  Business_fin=df[df['subject']=="Business Finance"]
          6  Business_review_fin_max=Business_fin.loc[Business_fin['num_reviews'].idxmax()]
          7
          8  Musical_Instruments=df[df['subject']=="Musical Instruments"]
          9  Musical_review_Instruments_max=Musical_Instruments.loc[Musical_Instruments['
10
11  Graphic_Design=df[df['subject']=="Graphic Design"]
12  Graphic_review_Design_max=Graphic_Design.loc[Graphic_Design['num_reviews'].i
13
14  ##concatination of all this dataframes
15  review_df=pd.concat([web_review_max,Business_review_fin_max,Musical_review_I
16  review_df=review_df.T

```

```

In [28]: 1  ## plotting pie chart according to max
          2  subject=['Web Development','Musical Instruments','Graphic Design','Business F
          3  exp=[0.1,0,0,0]
          4  plt.pie(review_df.num_reviews,labels=subject,autopct='%2.1f%%',explode=exp)
          5  plt.show()

```



```

In [29]: 1  ## For Average
2  web=df[df['subject']=='Web Development']['num_reviews'].sum()
3  web=web/w
4  bus=df[df['subject']=='Business Finance']['num_reviews'].sum()
5  bus=b
6  mus=df[df['subject']=='Musical Instruments']['num_reviews'].sum()
7  mus=mus/m
8  grap=df[df['subject']=='Graphic Design']['num_reviews'].sum()
9  grap=grap/g
10 data=[web,bus,mus,grap]
11 subject=['Web Development','Business Finance','Musical Instruments','Graphic
12 index=[0,1,2,3]
13 q2=pd.DataFrame(data,index)
14 headers = ["Avg_review_per_subject"]
15 q2.columns = headers
16 q2['subject']=subject
17
18 q2.head()

```

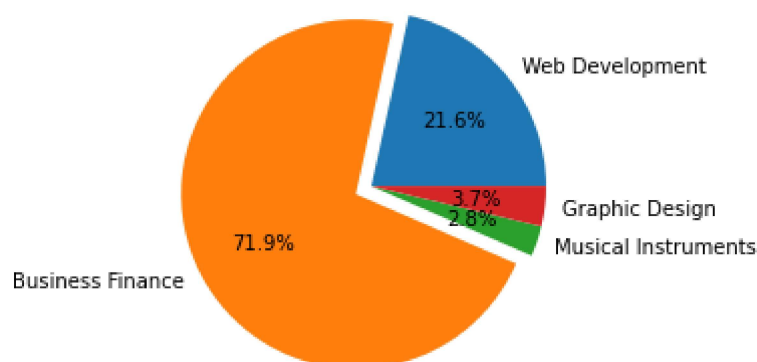
Out[29]:

	Avg_review_per_subject	subject
0	358.354167	Web Development
1	1195.000000	Business Finance
2	46.652941	Musical Instruments
3	61.475954	Graphic Design

```

In [30]: 1  ## plotting pie chart according to Average
2  subject1=['Web Development','Business Finance','Musical Instruments','Graphic
3  exp=[0,0.1,0,0]
4  plt.pie(q2.Avg_review_per_subject,labels=subject1,autopct='%2.1f%',explode=e
5  plt.show()

```



Which course is popular according to subscribers and create a graph of that distribution

```

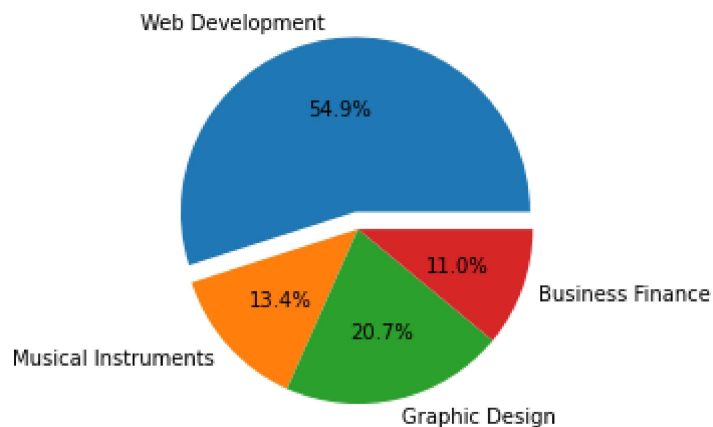
In [31]: 1 web=df[df['subject']=="Web Development"]
2 web_sub_max=web.loc[web['num_subscribers'].idxmax()]
3
4 Business_fin=df[df['subject']=="Business Finance"]
5 Business_sub_fin_max=Business_fin.loc[Business_fin['num_subscribers'].idxmax()]
6
7 Musical_Instruments=df[df['subject']=="Musical Instruments"]
8 Musical_sub_Instruments_max=Musical_Instruments.loc[Musical_Instruments['num_subscribers'].idxmax()]
9
10 Graphic_Design=df[df['subject']=="Graphic Design"]
11 Graphic_sub_Design_max=Graphic_Design.loc[Graphic_Design['num_subscribers'].idxmax()]
12
13 ##concatination of all this dataframes
14 review_sub=pd.concat([web_sub_max,Business_sub_fin_max,Musical_sub_Instruments_max,Graphic_sub_Design_max])
15 review_sub=review_sub.T

```

```

In [32]: 1 ## plotting pie chart acoording to max
2 subject=['Web Development','Musical Instruments','Graphic Design','Business Finance']
3 exp=[0.1,0,0,0]
4 plt.pie(review_sub.num_subscribers,labels=subject,autopct='%2.1f%%',explode=explode)
5 plt.show()

```



```

In [33]: 1  ## For Average
          2  web=df[df['subject']=='Web Development']['num_subscribers'].sum()
          3  web=web/w
          4  bus=df[df['subject']=='Business Finance']['num_subscribers'].sum()
          5  bus=b
          6  mus=df[df['subject']=='Musical Instruments']['num_subscribers'].sum()
          7  mus=mus/m
          8  grap=df[df['subject']=='Graphic Design']['num_subscribers'].sum()
          9  grap=grap/g
         10  data=[web,bus,mus,grap]
         11  subject=['Web Development','Business Finance','Musical Instruments','Graphic
         12  index=[0,1,2,3]
         13  q3=pd.DataFrame(data,index)
         14  headers = ["Avg_subscribers_per_subject"]
         15  q3.columns = headers
         16  q3['subject']=subject
         17
         18  q3.head()

```

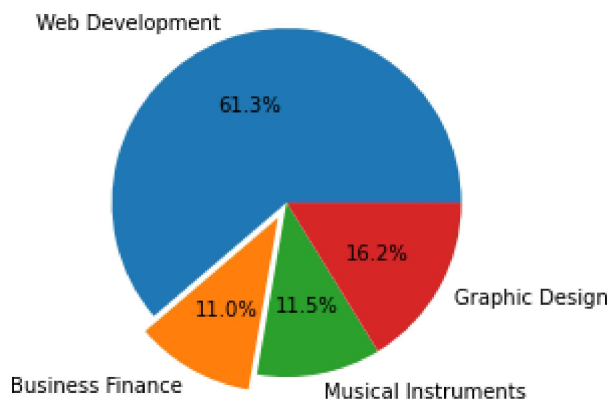
Out[33]:

	Avg_subscribers_per_subject	subject
0	6650.476667	Web Development
1	1195.000000	Business Finance
2	1245.130882	Musical Instruments
3	1763.097844	Graphic Design

```

In [34]: 1  ## plotting pie chart acoording to Average
          2  subject1=['Web Development','Business Finance','Musical Instruments','Graphic
          3  exp=[0,0.1,0,0]
          4  plt.pie(q3.Avg_subscribers_per_subject,labels=subject1,autopct='%2.1f%%',expl
          5  plt.show()

```



show the relation between course duration and subscribers

In [35]:

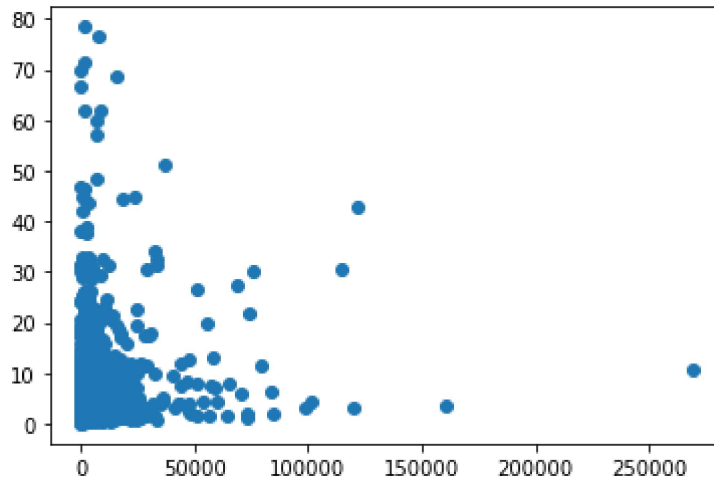
```
1 df.corr()
```

Out[35]:

	course_id	is_paid	price	num_subscribers	num_reviews	num_lect
course_id	1.000000	-0.013679	0.142319	-0.167856	-0.058550	-0.000000
is_paid	-0.013679	1.000000	0.328513	-0.266159	-0.087471	0.100000
price	0.142319	0.328513	1.000000	0.050769	0.113696	0.300000
num_subscribers	-0.167856	-0.266159	0.050769	1.000000	0.649946	0.100000
num_reviews	-0.058550	-0.087471	0.113696	0.649946	1.000000	0.200000
num_lectures	-0.024646	0.112574	0.330160	0.157746	0.243029	1.000000
level	0.078451	-0.033156	-0.073219	-0.062092	-0.055649	-0.100000
content_duration	-0.057223	0.094417	0.293450	0.161839	0.228889	0.800000
per_course_total_revenue	-0.053973	0.072902	0.346617	0.557175	0.769948	0.300000

In [36]:

```
1 ##Correlation coefficient values less than +0.8 or greater than -0.8 are not
2 ##scatter plot to show the correlation
3 plt.scatter(df['num_subscribers'], df['content_duration'])
4 plt.show()
```



In [37]:

```
1 from scipy.stats import linregress
2 linregress(df['num_subscribers'], df['content_duration'])
3 ## here p value is 5.2486796944911873e-23 which is less than 0.05 so relation
4 ## see the r value
```

Out[37]: LinregressResult(slope=0.00010308643347981602, intercept=3.764934092352409, rvalue=0.16183867741001334, pvalue=5.2486796944911873e-23, stderr=1.0367359669694272e-05)

```
In [38]: 1  ## Showing in different way
2  ##  $\rho$  = population correlation coefficient (unknown)
3  ##  $r$  = sample correlation coefficient (known; calculated from sample data)
4
5  np.corrcoef(df['num_subscribers'], df['content_duration'])
6
7  ## Here r value is 0.16183868 which shows a positive corelation not more
```

```
Out[38]: array([[1.          , 0.16183868],
                [0.16183868, 1.          ]])
```

```
In [39]: 1  ##some notes
2
3  ##If the p-value is less than the significance Level ( $\alpha = 0.05$ )
4  ## Decision: Reject the null hypothesis.
5  ##Conclusion: "There is sufficient evidence to conclude that there is a sign
6
7  ## If the p-value is NOT less than the significance Level ( $\alpha = 0.05$ )
8  ## Then Decision: DO NOT REJECT the null hypothesis.
9  ##Conclusion: "There is insufficient evidence to conclude that there is a si
10
```

which course is popular according to level and create a graph of that distribution

```
In [40]: 1  pd.crosstab(df.level, df.subject)## it shows no of courses according to Leve
```

```
Out[40]:
```

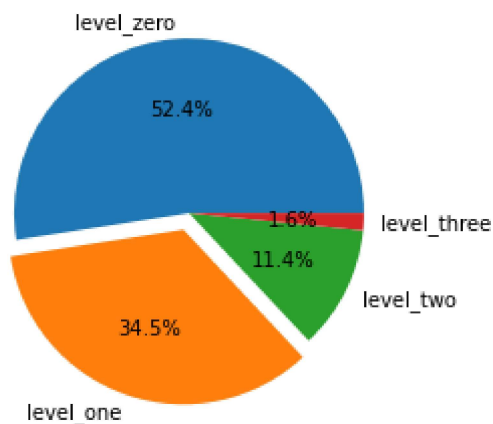
	subject	Business Finance	Graphic Design	Musical Instruments	Web Development
level					
0		696	298	276	659
1		340	243	296	391
2		128	57	101	135
3		31	5	7	15

```
In [41]: 1  df.level.value_counts()
```

```
Out[41]: 0    1929
1    1270
2     421
3      58
Name: level, dtype: int64
```

```
In [42]: 1  df3=df.level.value_counts()
```

```
In [43]: 1 subject1=['level_zero','level_one','level_two','level_three']
2 exp=[0,0.1,0,0]
3 plt.pie(df3,labels=subject1,autopct='%2.1f%%',explode=exp)
4 plt.show()
```



Best unpaid course according to subscribers

```
In [44]: 1 df_unpaid=df[df['is_paid']==0]
2 df.loc[df_unpaid['num_subscribers'].idxmax()]
```

```
Out[44]: course_id                                41295
course_title                                Learn HTML5 Programming From Scratch
url                                https://www.udemy.com/learn-html5-programming-...
      (https://www.udemy.com/learn-html5-programming-...)
is_paid                                          0
price                                          0
num_subscribers                                268923
num_reviews                                    8629
num_lectures                                    45
level                                          0
content_duration                                10.5
published_timestamp                2013-02-14T07:03:41Z
subject                                Web Development
per_course_total_revenue                                0
Name: 2827, dtype: object
```



```
In [50]: 1  ## Highest paid course according to subscribers
          2  paid=df[df['is_paid']==1]
          3  paid.loc[paid[paid['subject']=='Web Development']['per_course_total_revenue'
```

```
Out[50]: course_id                                625204
          course_title                            The Web Developer Bootcamp
          url                                     https://www.udemy.com/the-web-developer-bootcamp/
          (https://www.udemy.com/the-web-developer-bootcamp/)
          is_paid                                  1
          price                                    200
          num_subscribers                         121584
          num_reviews                             27445
          num_lectures                             342
          level                                    0
          content_duration                         43
          published_timestamp                     2015-11-02T21:13:27Z
          subject                                  Web Development
          per_course_total_revenue                 24316800
          Name: 3230, dtype: object
```

```
In [ ]: 1
```