

Cdemy

This project is based on the dataset for the Udemy Data Analytics Project published on Medium. A member of the data team at Udemy had worked on the project with excel and

The aim is to explore the number of courses for each subject area, the number of subscribers, how well the courses were rated and how much revenue is being generated etc. For

The scope of the project covers only four course categories: Business Finance, Graphics Design, Musical Instruments and Web Development.

each subject category, we will also identify the key words for the most best rated and most sought after courses..

visualised with PowerBI but I will be using python and PowerBI.

We will find answers to the following questions and more:

· What levels of courses are free?

• Which subject area has the highest number of subscribers?

• Which subject area has the highest and lowest number of published courses?

2011 Name: year, dtype: int64 Let us clean the text columns a bit. In [119... Let us create a regex function to clean the text column off: - punctuation - special characters def clean(text): text = str(text).title() text = re.sub('\[.*?\]', '', text) text = re.sub('https?://\S+|www\.\S+', '', text) #remove url text = re.sub('[%s]' % re.escape(string.punctuation), '', text) #remove punctuations text = re.sub('\n', '', text)
text = re.sub('[0-9]', '', text) text = re.sub('<.*?>+', '', text) return text #apply the function created df['course_title'] = df['course_title'].apply(clean) In [120... #clean off punctautions from subject column df['subject']=df['subject'].str.split(': ').str[-1].str.lstrip() df['subject'].value_counts() Web Development Out[120]: **Business Finance** 1191 Musical Instruments 680 Graphic Design 602 Name: subject, dtype: int64 Let us create the Revenue column by multiplying the num_subscribers by price of the course. df['revenue'] = df['num_subscribers'] * df['price'] In [121... In [123... #separate the component dataframes df_biz = df.query('subject == "Business Finance"') df_gfx = df.query('subject == "Graphic Design"') df_mus = df.query('subject == "Musical Instruments"') df_dev = df.query('subject == "Web Development"') print('There are {} {} courses'.format(df_biz.shape[0], df_biz.subject[0])) print('There are {} {} courses'.format(df_gfx.shape[0], df_gfx.subject[0])) print('There are {} {} courses'.format(df_mus.shape[0], df_mus.subject[0])) print('There are {} {} courses'.format(df_dev.shape[0], df_dev.subject[0])) There are 1191 Business Finance courses There are 602 Graphic Design courses There are 680 Musical Instruments courses There are 1199 Web Development courses **Data Cleaning Steps:** Checked datatypes and missing values. • Dropped null values from the dataframe. Created a column to categorize the courses into free and paid. • Extracted and stored the date only data from the datetime column and created the year column. · Created a function with regex to clean up the text columns. Dropped off columns that are not necessary for the analysis. • Created the revenue column. **Exploratory Data Analysis** From the pie chart below, Web Development leads with 1203 courses while Graphics Designs at 602 courses is the least. result = df['price_group'].value_counts().reset_index() result.columns = ['price_group', 'count'] fig = px.bar(result, x='count', y='price_group', orientation='h', labels={'count': 'Count', 'price_group': 'Price Group'}, title='Count of Price Groups', color='price_group') fig.show() Count of Price Groups Price Group Paid free Paid Price Group free 0 500 1000 1500 2000 2500 3000 3500 Count

2017

2014

2013

2012

713

490

201

45

df['subject'].value_counts() In [133... Web Development 1199 Out[133]: Business Finance 1191 Musical Instruments 680 Graphic Design 602 Name: subject, dtype: int64 result = df['subject'].value_counts().reset_index() In [164... result.columns = ['subject', 'count'] fig = px.pie(result, names='subject', values='count', title='Distribution of Subjects', labels={'count': 'Count'}) fig.show() 0 Distribution of Subjects Web Development **Business Finance Musical Instruments** Graphic Design 32.4% 32.7% 16.4% df.groupby('subject')['revenue'].sum().sort_values(ascending=False) In [141.. subject Out[141]: Web Development 627597400.0 123735315.0 Business Finance 76983170.0 Graphic Design 53359055.0 Musical Instruments Name: revenue, dtype: float64 result = df.groupby('subject')['revenue'].sum().reset_index() In [144... fig = px.bar(result, x='subject', y='revenue', labels={'revenue': 'Total Revenue', 'subject': 'Subject'}, title='Total Revenue by Subject') fig.show() Total Revenue by Subject 600M 500M Total Revenue 400M 300M 200M 100M 0 **Musical Instruments Business Finance** Graphic Design Web Development Subject In [166... fig, ((ax0, ax1), (ax2, ax3)) = plt.subplots(2,2,figsize=(15,10))labels = [0,1,2,3,4,5]ax0.bar(df_biz.year.value_counts().index,df_biz.year.value_counts().values) ax0.set_title('Business Finance') ax0.set_ylabel('Number of Courses') ax1.bar(df_gfx.year.value_counts().index,df_gfx.year.value_counts().values,color='r') ax1.set_title('Graphic Design') ax2.bar(df_mus.year.value_counts().index,df_mus.year.value_counts().values,color='y') ax2.set_title('Musical Instruments') ax2.set_ylabel('Number of Courses') ax3.bar(df_dev.year.value_counts().index,df_dev.year.value_counts().values, color='c') ax3.set_title('Web Development')

fig.tight_layout() plt.show() Graphic Design **Business Finance** 350 175 300 150 125 Number of Courses 100 150 75 100 50 50 25 2013 2013 2014 2014 2015 2017 2012 2015 2016 2017 2012 2016 Musical Instruments Web Development 400 200 Number of Courses 300 200 100 2017 2012 2013 2014 2015 2016 2013 2015 2016 2017 2011 2012 2014 result = df.groupby('year')['subject'].value_counts().reset_index(name='count') heatmap_data = result.pivot(index='year', columns='subject', values='count') plt.figure(figsize=(12, 8)) sns.heatmap(heatmap_data, annot=True, cmap='YlGnBu', fmt='g', linewidths=.6) plt.title('Subject Counts by Year') plt.show() Subject Counts by Year 2011 - 400 2012 6 10 10 19 - 350 - 300 84 23 39 55 - 250 year 2014 192 65 120 113 - 200 339 168 171 336 - 150 2016 347 181 228 448 - 100 - 50 2017 155 112 **Business Finance** Graphic Design Web Development Musical Instruments subject import plotly.graph_objects as go from plotly.subplots import make_subplots # Create subplot grid fig = make_subplots(rows=2, cols=2, subplot_titles=['Business Finance', 'Graphic Design', 'Musical Instruments', 'Web Development'], shared_yaxes=True, horizontal_spacing=0.1, vertical_spacing=0.15) # Plot for Business Finance fig.add_trace(go.Bar(x=df_biz.groupby('year')['revenue'].sum().index, y=df_biz.groupby('year')['revenue'].sum(), name='Business Finance'), row=1, co # Plot for Graphic Design

fig.add_trace(go.Bar(x=df_gfx.groupby('year')['revenue'].sum().index, y=df_gfx.groupby('year')['revenue'].sum(), name='Graphic Design', marker_color

fig.add_trace(go.Bar(x=df_mus.groupby('year')['revenue'].sum().index, y=df_mus.groupby('year')['revenue'].sum(), name='Musical Instruments', marker_

fig.add_trace(go.Bar(x=df_dev.groupby('year')['revenue'].sum().index, y=df_dev.groupby('year')['revenue'].sum(), name='Web Development', marker_colo

Graphic Design

Web Development

2014

Rating

mean min max

ngular Formerly Angular The Complete Guide

Javascript Understanding The Weird Parts

• The dataset is limited to only four subject areas while leaving out a ton of interesting areas including Data Science, Cyber Security, Cloud Computing, Digital Marketing.

• For courses with high rating and subscriptions, they are courses on Accounting, Forex, Stock in Business Finance; Photoshop, Adobe Illustrator in Graphics Design; Piano and

• For revenue generation, 2014 was the big year for Musical Instruments, 2015 for Business Finance and Web Development while Graphics Design made the

• The dataset is not updated to the current year, 2022. It would have been nice to see how digital learning grew in the Covid-19 era.

• Web Development and Graphics Designs have the highest and lowest number of published courses respectively.

• Subscribers are more interested in Web Development courses and least in Musical Instruments courses.

• Web Development courses are more expensive, attractive more subscribers; hence generate the highest revenue.

• It was not explained what the scale of the Rating was; does a rating of 0.0 mean there was no rating at all or the course received the least score.

• There is no data on the content creators; it would have been great to know which tutors make the most subscribed and best rated contents.

• There are no demographic data on the subscribers and the duration of their learning of the chosen subject.

• The rating of the courses didnt provide details of the comments made by the subscribers.

• It was not expressely stated what the unit of the course duration was - whether hours or minutes.

• Graphics Design courses followed by Business Finance received the best ratings.

Guitar in Musical Instruments and HTML, CSS, building a Website in Web Development.

• The year 2016 has highest number of courses published across the different subjects.

• The % of the courses that are free are more in expert level.

The courses received high ratings across board.

most revenue in 2016.

2016

star_Rating

5 3.878254

5 3.978405

5 2.039706

5 3.635530

mean min

0.0

0.0 15099800.0

Learn And Understand Nodejs

revenue

mean

Learn And Understand Angularis Modern React With Redux

max

4773795.0 103891.952141

7257600.0 127879.019934

0.0 24316800.0 523434.028357

2012

df.groupby('subject')[['price', 'num_subscribers', 'Rating', 'star_Rating', 'revenue']].agg(['min', 'max', 'mean'])

mean min max

65576 1569.026868 0.00 1.00 0.690353

53851 1766.026578 0.01 0.99 0.730382

19 268923 6619.922435 0.00 1.00 0.642127

fig = px.treemap(top_courses, path=['course_title'], values='revenue', title='Top 10 Courses by Revenue')

num_subscribers

101154 1245.130882 0.00

top_courses = df[['course_title', 'revenue']].sort_values('revenue', ascending=False).head(10)

The Complete Web Developer Course

max

fig.update_layout(height=600, width=800, title_text='Revenue Across Subjects and Years', showlegend=False)

Plot for Musical Instruments

Revenue Across Subjects and Years

Business Finance

Musical Instruments

2014

max

Graphic Design 0.0 200.0 57.890365

Musical Instruments 0.0 200.0 49.558824

Web Development 0.0 200.0 77.035029

Top 10 Courses by Revenue

The Web Developer Bootcamp

0.0 200.0 68.694374

2016

price

mean min

Plot for Web Development

Update layout

Show the plot
fig.show()

40M

30M

20M

10M

0

250M

200M

150M

100M

50M

Out[183]:

In [255...

0

2012

subject

Top 10 courses by revenue

Business Finance

Create a treemap

Show the plot
fig.show()

Comments

Conclusions

2012