# Udemy Course Analysis
## -Group 21



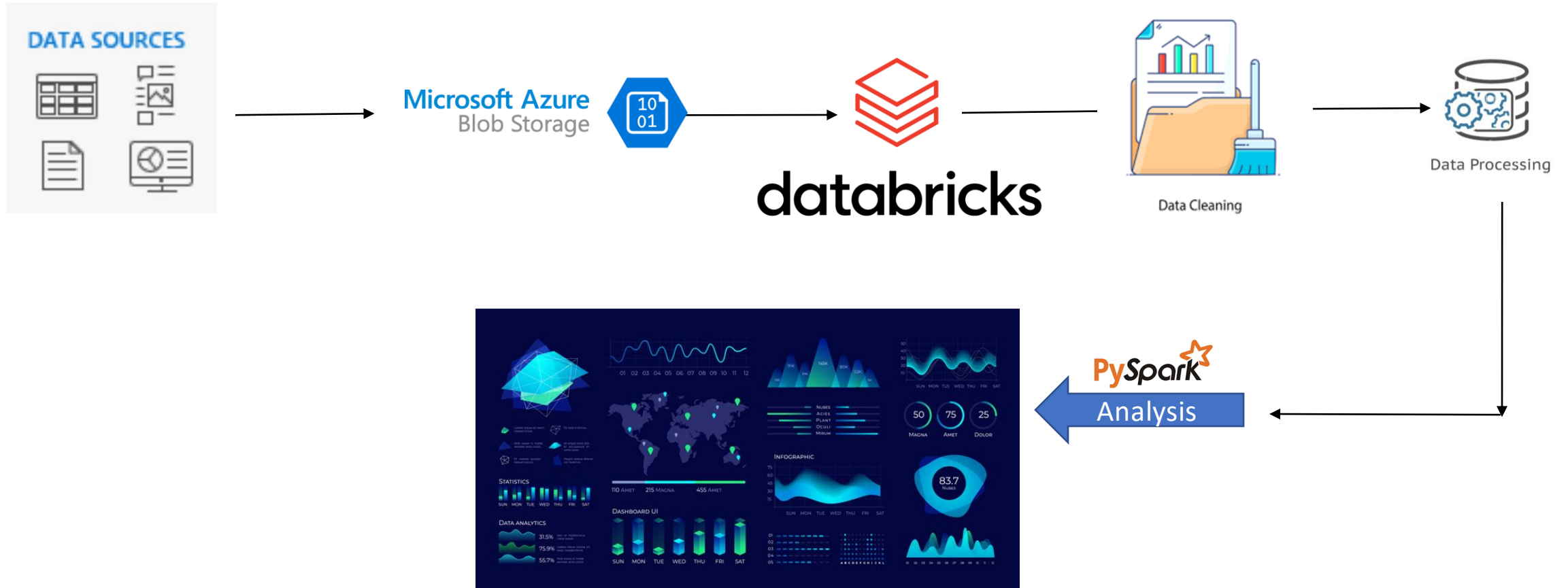**Bindu Parvati, Jonnala Gadda –16338568**

# Introduction

- Udemy, Inc. operates an edtech platform that delivers online learning and teaching services.

- The platform features a broad range of courses spanning various fields such as business, academics, the arts, health and fitness, language, music, and technology.

- The primary objective of the project is to conduct statistical analysis on the course data available on Udemy.

UMKC

## About dataset
Name: Udemy courses
Data size: ~1.6GB
Date – October 2022
No of files: 2
No of columns: 23

It has detailed information on all the Udemy courses. The data is segregated into 2 CSV files:
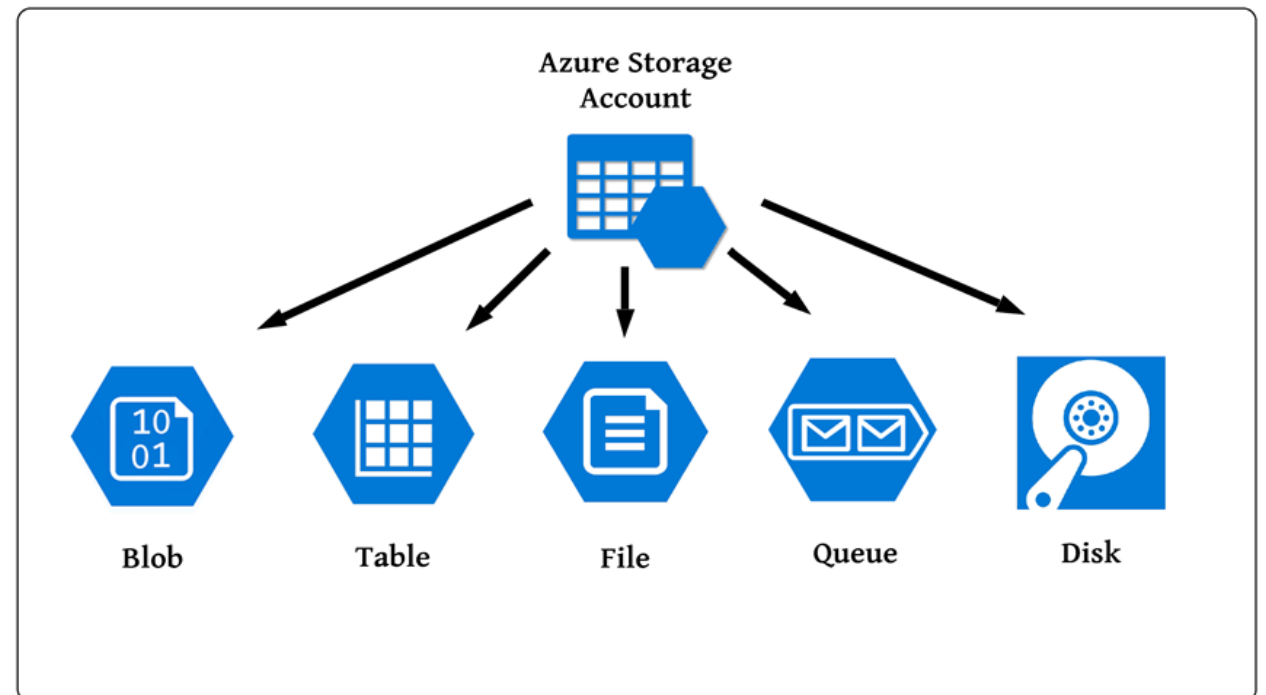
- Comments: This file contains over 9 million comments and ratings. It has 6 columns (comment_id, course_id, rate, date , display_name, comment)

- Course_info: This file holds the information of 209,734 courses offered by Udemy and 73,514 instructors teaching in 79 languages in 13 different categories. It has 17 columns(course_id, course_title, is_paid, price, headline, num_subscribers, num_lectures, content_length_min, published_time, last_update_date, category, subcategory, topic, language, course_url, instructor_name, instructor_url

- Data Source: Kaggle

# Architecture

# Storage Account

- All Azure Storage data objects, such as blobs, file shares, queues, tables, and disks, are included in an Azure storage account.

- A storage account offers a distinct namespace and is reachable through HTTP or HTTPS from anywhere in the globe.

- robust, very accessible, safe, and incredibly expandable.

- Blob storage is a type of cloud storage that is designed to hold large volumes of unstructured data.

# Storage Account Creation in Azure
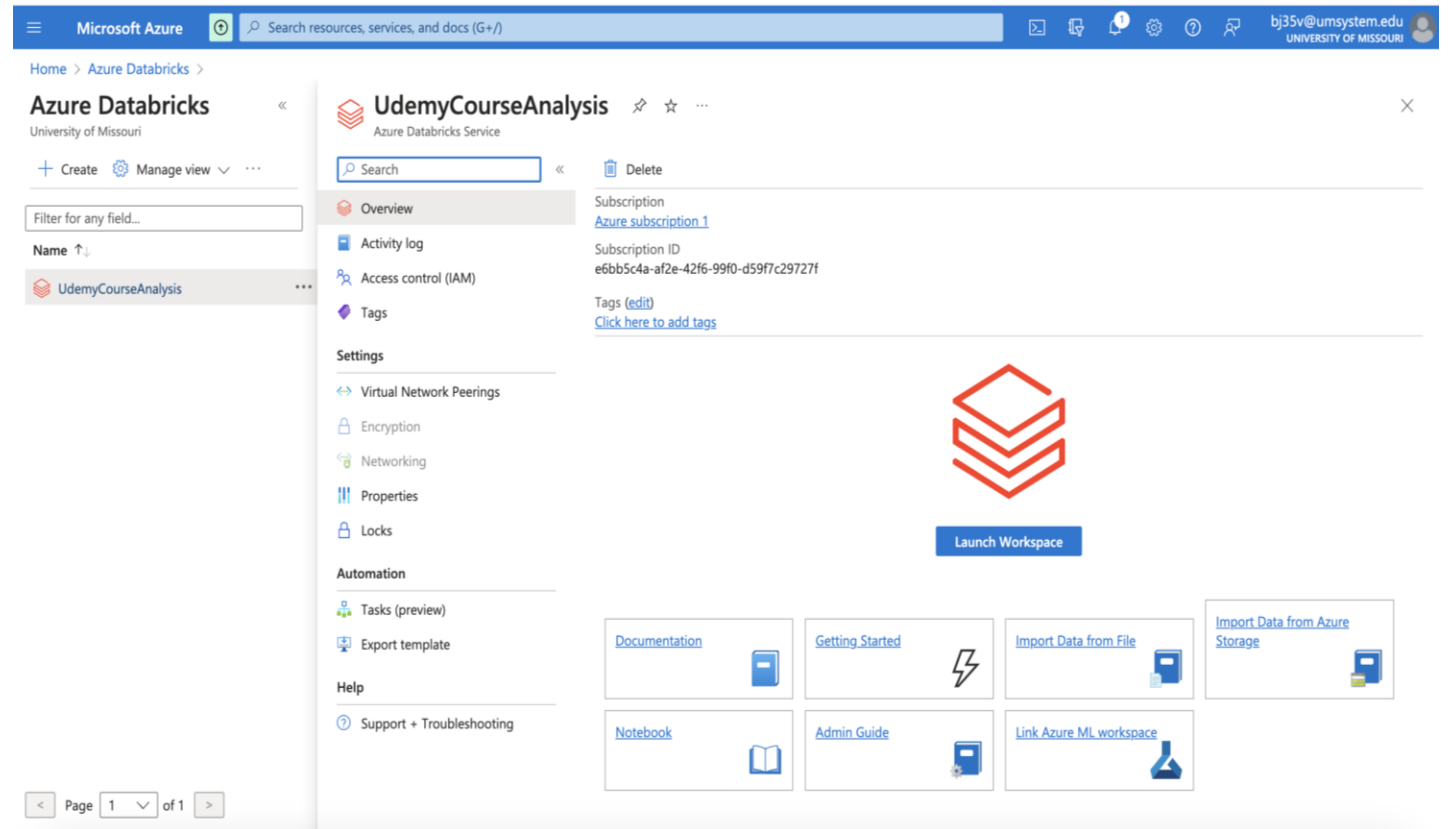
# Azure Databricks

For creating, implementing, sharing, and maintaining organization-grade data solutions at scale, Azure Databricks provides a standardized collection of tools.

The Microsoft Azure Databricks Lakehouse Platform integrates with your cloud account's cloud storage and security while managing and deploying cloud infrastructure on your behalf.

# Environment setup and Tasks performed

- Mounting Blob storage with Databricks Notebook

- Data cleaning by removal of nulls

- Identified and re-created the messy data

- Re-indexing and re-arranging the columns

- Defining Schema and Data filtering

- Basic analysis of the data



```
Cmd 1

1  dbutils.fs.mount(
2    source = 'wasbs://rawdata@udemyprojectdata.blob.core.windows.net',
3    mount_point = '/mnt/team21data',
4    extra_configs =
   {'fs.azure.account.key.udemyprojectdata.blob.core.windows.net':'ws0gfWpjj1nRRDAi0q9XKcn78bTUSNzMrLktEvEFFXIrCPNDHijYR32lFi/Pf90nORqF5
   J18EbLq+AStNvuplg=='}
5  )
```

Out[48]: True

Command took 20.74 seconds -- by bj35v@umsystem.edu at 4/3/2023, 12:45:54 PM on Jonnala Parvati's Cluster

```
Cmd 2

1  dbutils.fs.ls("/mnt/datateam21")
```

Out[26]: [FileInfo(path='dbfs:/mnt/datateam21/Comments.csv', name='Comments.csv', size=1609425094, modificationTime=1680543721000),
 FileInfo(path='dbfs:/mnt/datateam21/Course_info.csv', name='Course_info.csv', size=76151841, modificationTime=1680542792000)]

Command took 0.36 seconds -- by bj35v@umsystem.edu at 4/4/2023, 11:26:10 AM on Jonnala Parvati's Cluster

UMKC

# Defining Schema

```
1   from pyspark.sql.functions import *
2   from pyspark.sql.types import *
3
4   newDF=[StructField('id',IntegerType(),True),
5         StructField('course_id',IntegerType(),True),
6         StructField('rate',StringType(),True),
7         StructField('date', TimestampType(),True),
8         StructField('display_name',StringType(),True),
9         StructField('comment',StringType(),True)
10        ]
11        ]
12  finalStruct=StructType(fields=newDF)
13
```

Command took 0.10 seconds -- by bj35v@umsystem.edu at 4/5/2023, 10:24:24 AM on Jonnala Parvati's Cluster

```
1   comments = spark.read.schema(finalStruct).option("mode", "DROPMALFORMED").
2   options(header= 'TRUE', delimiter = ',').csv('dbfs:/mnt/team21data/Comments.csv')
```

▼ 🗏 comments: pyspark.sql.dataframe.DataFrame
        id: integer
        course_id: integer
        rate: string
        date: timestamp
        display_name: string
        comment: string

Command took 0.17 seconds -- by bj35v@umsystem.edu at 4/5/2023, 10:24:26 AM on Jonnala Parvati's Cluster

```
1   import pyspark
2   from pyspark.sql.functions import *
3   from pyspark.sql.types import *
4
5   newDF1=[StructField('id',FloatType(),True),
6         StructField('title',StringType(),True),
7         StructField('is_paid',StringType(),True),
8         StructField('price',FloatType(),True),
9         StructField('headline',StringType(),True),
10        StructField('num_subscribers',FloatType(),True),
11        StructField('avg_rating',FloatType(),True),
12        StructField('num_reviews',FloatType(),True),
13        StructField('num_comments',FloatType(),True),
14        StructField('num_lectures',FloatType(),True),
15        StructField('content_length_min',FloatType(),True),
16        StructField('published_time', TimestampType(),True),
17        StructField('last_update_date',TimestampType(),True),
18        StructField('category',StringType(),True),
19        StructField('subcategory',StringType(),True),
20        StructField('topic',StringType(),True),
21        StructField('language',StringType(),True),
22        StructField('course_url',StringType(),True),
23        StructField('instructor_name',StringType(),True),
24        StructField('instructor_url',StringType(),True),
25        ]
26  finalStruct_course = StructType(fields=newDF1)
```

# Data



```python
1  course_info.show(10)
```
Python

▸ (1) Spark Jobs

```
+------+--------------------+-------+------+--------------------+---------------+----------+-----------+------------+------------+----------------+--------------------+---------+
|    id|               title|is_paid| price|            headline|num_subscribers|avg_rating|num_reviews|num_comments|num_lectures|content_length_min|          published|
_time|    last_update_date|          category|          subcategory|           topic|language|          course_url|     instructor_name|          instructor_url|
+------+--------------------+-------+------+--------------------+---------------+----------+-----------+------------+------------+----------------+--------------------+---------+
|4715.0|Online Vegan Vege...|   True| 24.99|Learn to cook del...|         2231.0|      3.75|      134.0|        42.0|        37.0|          1268.0|2010-08-05 22:
06:13|2020-11-06 00:00:00|          Lifestyle|     Food & Beverage|Vegan Cooking| English|/course/vegan-veg...|         Angela Poch|     /user/angelapoch/|
|1769.0|The Lean Startup ...|  False|   0.0|"Debunking Myths ...|        26474.0|       4.5|      709.0|       112.0|         9.0|            88.0|2010-01-12 18:
09:46|               null|          Business|     Entrepreneurship| Lean Startup| English|/course/the-lean-...|          Eric Ries|     /user/ericries/|
|5664.0|How To Become a V...|   True| 19.99|Get the tools you...|         1713.0|       4.4|       41.0|        13.0|        14.0|            82.0|2010-10-13 18:
07:17|2019-10-09 00:00:00|          Lifestyle|     Other Lifestyle|Vegan Cooking| English|/course/see-my-pe...|         Angela Poch|     /user/angelapoch/|
|7723.0|How to Train a Puppy|   True|199.99|Train your puppy ...|         4988.0|       4.8|      395.0|        88.0|        36.0|          1511.0|2011-06-20 20:
08:38|2016-01-13 00:00:00|          Lifestyle| Pet Care & Training| Pet Training| English|/course/complete-...|          Ian Dunbar|     /user/ian-dunbar/|
|8157.0|Web Design from t...|   True|159.99|Learn web design ...|         1266.0|      4.75|       38.0|        12.0|        38.0|           569.0|2011-06-23 18:
31:20|               null|             Design|          Web Design|   Web Design| English|/course/web-desig...|     E Learning Lab|     /user/edwin-ang-2/|
|8139.0|14-Day Yoga Detox...|   True| 29.99|Lose weight, get ...|        20505.0| 4.5301204|      796.0|       135.0|        31.0|          1163.0|2011-07-15 04:
13:24|2018-05-22 00:00:00|   Health & Fitness|               Yoga|         Yoga| English|/course/yoga-for-...|       Sadie Nardini|     /user/sadienardini/|
|2762.0|Simple Strategy f...|   True| 39.99|Use my favorite T...|         3309.0|      3.85|      958.0|       241.0|         8.0|            80.0|2010-04-14 16:
32:46|2019-03-07 00:00:00|Finance & Accounting| Investing & Trading|Swing Trading| English|/course/swing-tra...|         Tom Watson|     /user/tomwatson/|
|8082.0|Ruby Programming ...|   True| 74.99|Learn Ruby Progra...|        28824.0|       4.0|      741.0|       189.0|        56.0|           363.0|2011-07-08 21
```

Command took 0.51 seconds -- by bj35v@umsystem.edu at 4/5/2023, 11:41:43 AM on Jonnala Parvati's Cluster

```python
1
2  rows = course_info.count()
3  print("number of rows in course_info file = ", rows)
4
5  cols = len(course_info.columns)
6  print("number of columns in course_info file = ", cols)
```

▸ (2) Spark Jobs

```
number of rows in course_info file =  209734
number of columns in course_info file =  20
```

```python
1  comments.show(5)
```

▸ (1) Spark Jobs

```
+---------+---------+----+-------------------+--------------+--------------------+
|       id|course_id|rate|               date|  display_name|             comment|
+---------+---------+----+-------------------+--------------+--------------------+
| 88962892|  3173036| 1.0|2021-06-30 01:54:25|         Rahul|I think a beginne...|
|125535470|  4913148| 5.0|2022-10-07 18:17:41|         Marlo|Aviva is such a n...|
| 68767147|  3178386| 3.5|2020-10-19 13:35:37| Yamila Andrea|Muy buena la intr...|
|125029758|  3175814| 5.0|2022-10-01 04:13:49|     Jacqueline|This course is th...|
| 76584052|  3174896| 4.5|2021-01-30 16:45:11|       Anthony|I found this cour...|
+---------+---------+----+-------------------+--------------+--------------------+
only showing top 5 rows
```

```python
1  rows = comments.count()
2  print("number of rows in comments file = ", rows)
3
4  cols = len(comments.columns)
5  print("number of columns in comments file = ", cols)
```
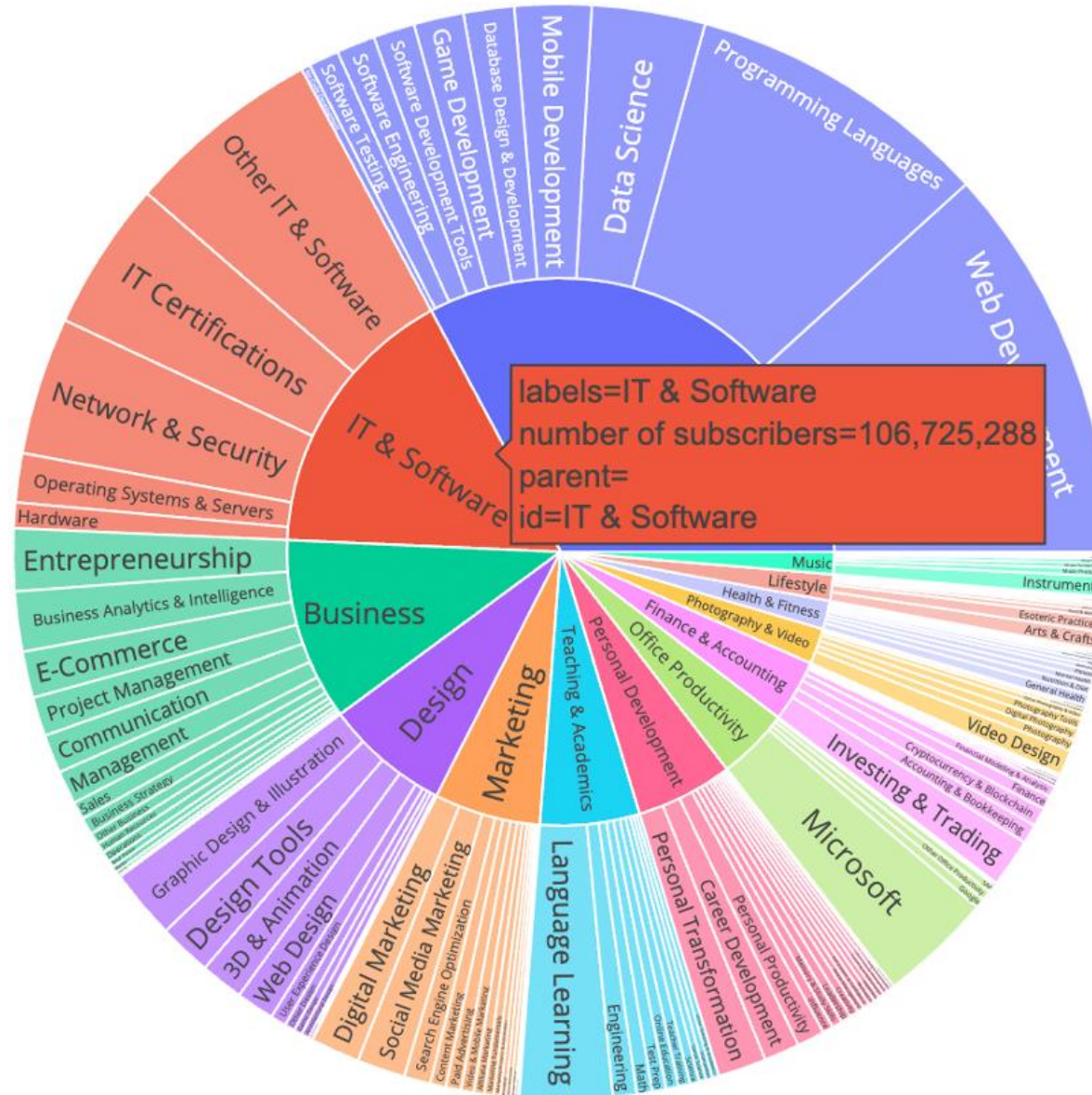
▸ (2) Spark Jobs

```
number of rows in comments file =  10813069
number of columns in comments file =  6
```
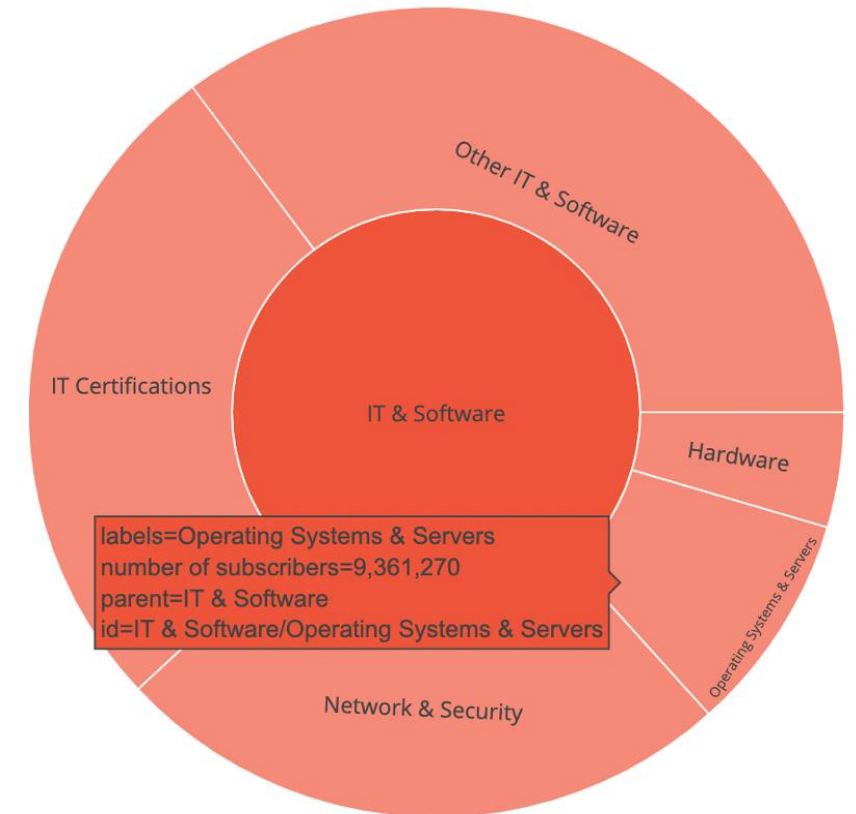
UMKC

# Exploratory Analysis

- Number of categories: 13

- Number of subcategories: 258

- Number of topics:3925

- Number of paid courses:187959

- Number of free courses:21733

- Number of Instructors:72893

```
+----------------------------+          +--------------------+
|count(DISTINCT category)|              |count(DISTINCT topic)|
+----------------------------+          +--------------------+
|                         13|           |                3925|
+----------------------------+          +--------------------+

+----------------------------+          +----------------------------+
|count(DISTINCT subcategory)|           |count(DISTINCT instructor_name)|
+----------------------------+          +----------------------------+
|                        258|           |                       72893|
+----------------------------+          +----------------------------+
```

| 3 | unpaid_courses |
|---|----------------|

| 2 | paid_courses |
|---|--------------|

▸ (2) Spark Jobs          ▸ (2) Spark Jobs

Out[191]: 21733          Out[190]: 187959

# Number of subscribers in different subcategories



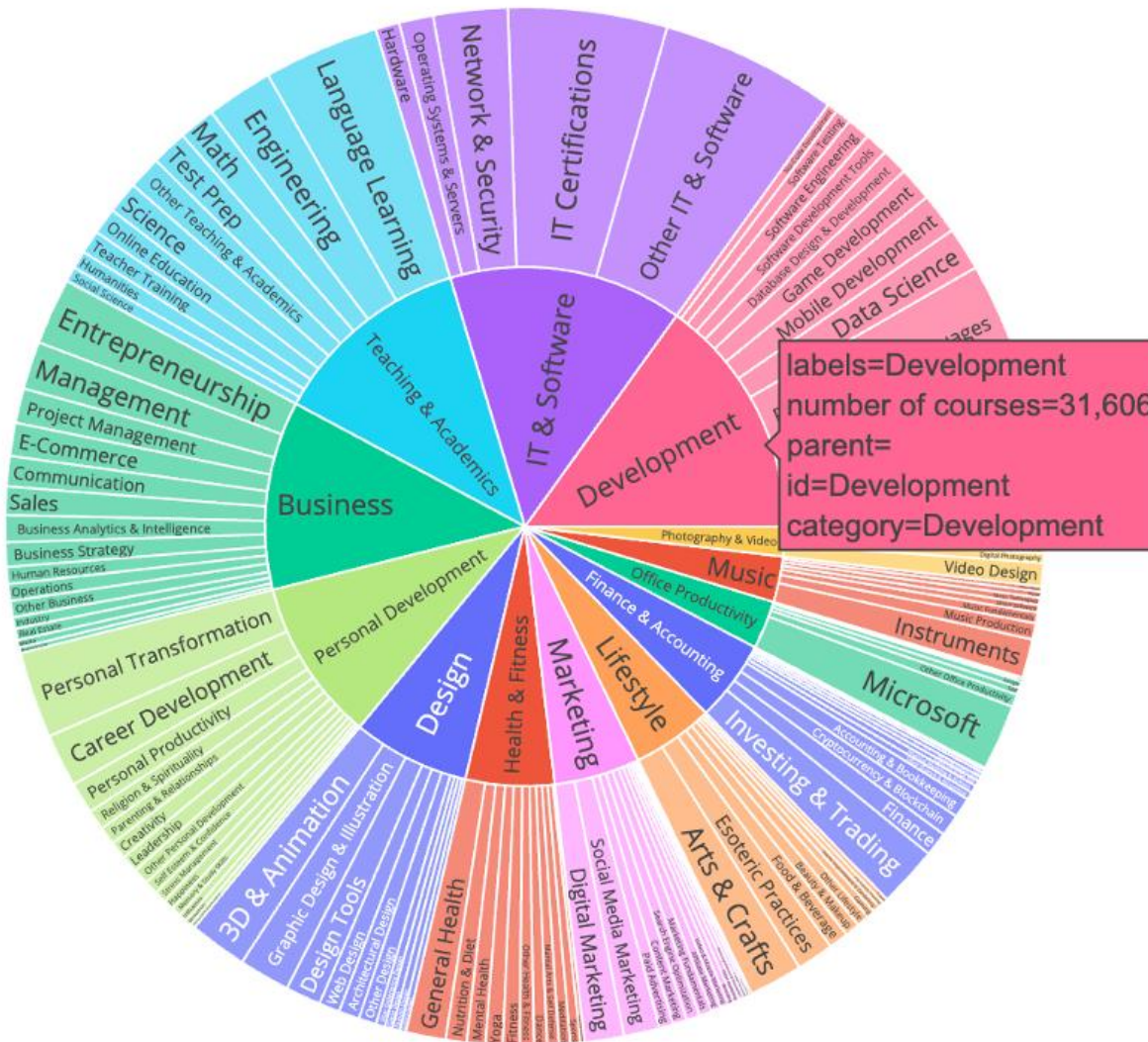labels=IT & Software
number of subscribers=106,725,288
parent=
id=IT & Software

labels=Operating Systems & Servers
number of subscribers=9,361,270
parent=IT & Software
id=IT & Software/Operating Systems & Servers

Number of courses offered in different subcategories

labels=Development
number of courses=31,606
parent=
id=Development
category=Development

Number of courses offered in different subcategories

labels=Web Development
number of courses=10,382
parent=Development
id=Development/Web Development
category=Development

# Growth Analysis

| published_year | id | num_subscribers | num_reviews | num_comments | num_lectures |
|---|---|---|---|---|---|
| 2010 | 4 | 33727.0 | 1842.0 | 408.0 | 68.0 |
| 2011 | 57 | 1328159.0 | 60041.0 | 13170.0 | 4555.0 |
| 2012 | 463 | 7208781.0 | 356369.0 | 70807.0 | 25517.0 |
| 2013 | 1771 | 1.4778123E7 | 1120787.0 | 232587.0 | 76614.0 |
| 2014 | 3392 | 2.7486349E7 | 2164269.0 | 455951.0 | 139467.0 |
| 2015 | 7071 | 5.6369593E7 | 6031783.0 | 1111070.0 | 296750.0 |
| 2016 | 7967 | 6.7377486E7 | 7735998.0 | 1445077.0 | 366585.0 |
| 2017 | 12241 | 8.1335593E7 | 9476263.0 | 1710549.0 | 560803.0 |
| 2018 | 20615 | 9.227419E7 | 8876357.0 | 1567888.0 | 940099.0 |
| 2019 | 23501 | 8.8881482E7 | 6496135.0 | 1125479.0 | 973949.0 |
| 2020 | 44864 | 1.20561779E8 | 5553367.0 | 984244.0 | 1571114.0 |
| 2021 | 51391 | 7.2717205E7 | 2573612.0 | 502163.0 | 1615536.0 |
| 2022 | 36095 | 1.770909E7 | 543948.0 | 150609.0 | 1081437.0 |

UMKC

Number of Categories and Subcategories Over the Years



Number of comments posted per year

| year | num_categories | num_subcategories |
|------|----------------|-------------------|
| 2010 | 3 | 4 |
| 2011 | 9 | 25 |
| 2012 | 13 | 95 |
| 2013 | 13 | 121 |
| 2014 | 13 | 129 |
| 2015 | 13 | 129 |
| 2016 | 13 | 130 |
| 2017 | 13 | 130 |
| 2018 | 13 | 130 |
| 2019 | 13 | 130 |
| 2020 | 13 | 130 |
| 2021 | 13 | 130 |
| 2022 | 13 | 130 |

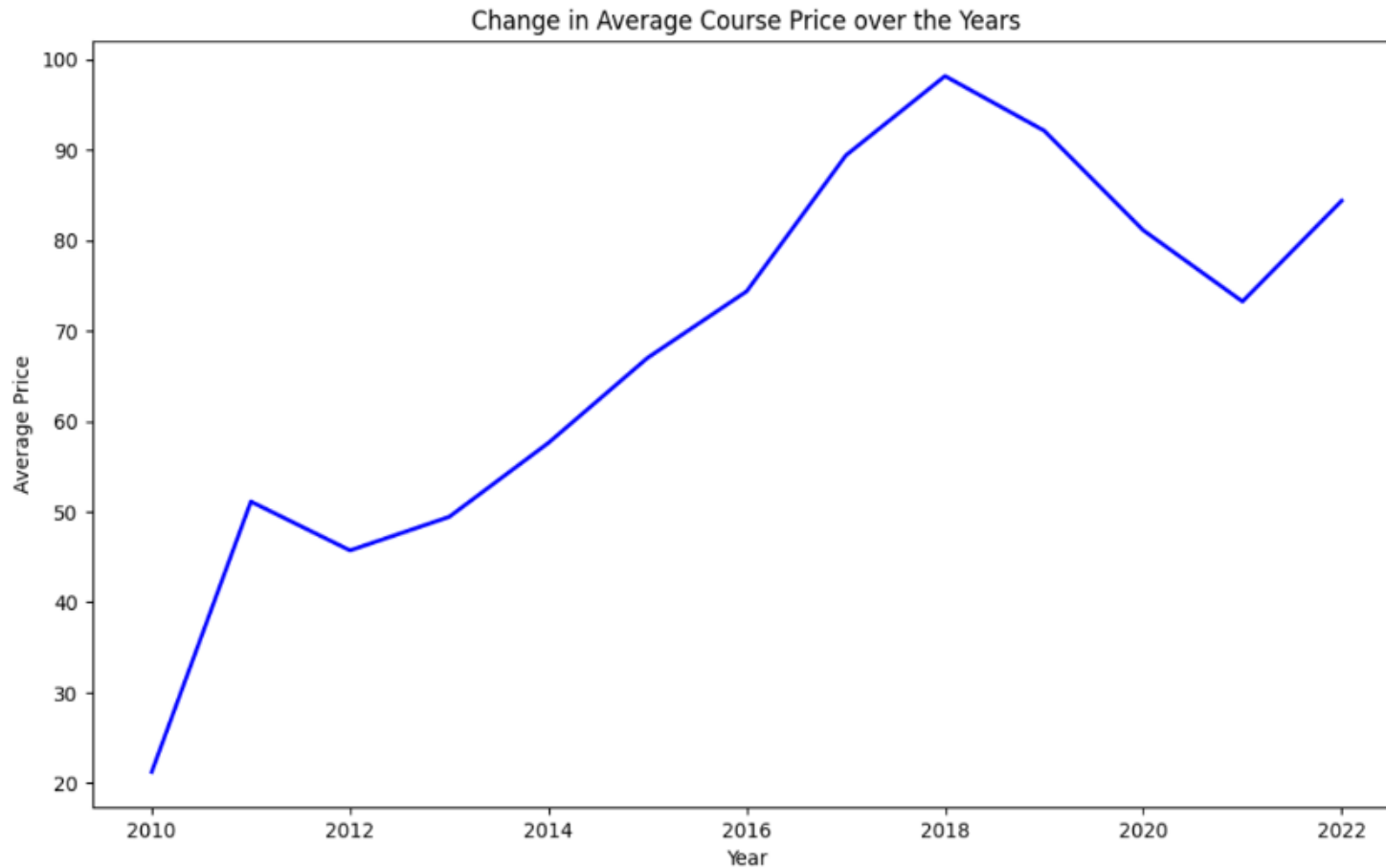Udemy Course Growth Over Time

- The dataset indicates that the first course was released in the first quarter of 2010.
- In 2010, one course was released each quarter.
- Up until the first quarter of 2020, there were roughly 7000 released courses every quarter.
- More than 13,200 courses were published in the second quarter of 2020, roughly doubling the number of courses published in a single quarter.
- The COVID-19 pandemic, which began in March 2020, can be a reason for this increase.
- Many people have turned to online education as a result of the pandemic, staying at home. Following that, the quantity of published courses stayed within the same range and varied between 11,000 and 14,000 courses each quarter.
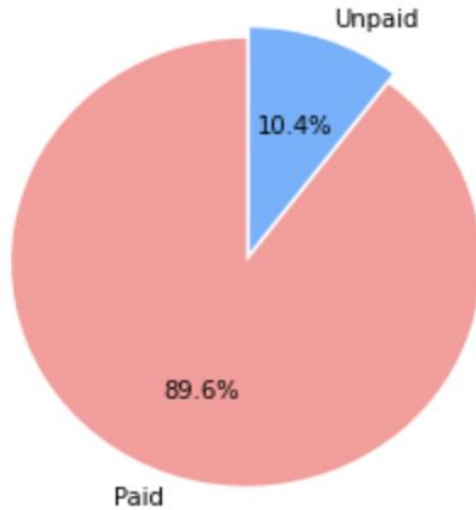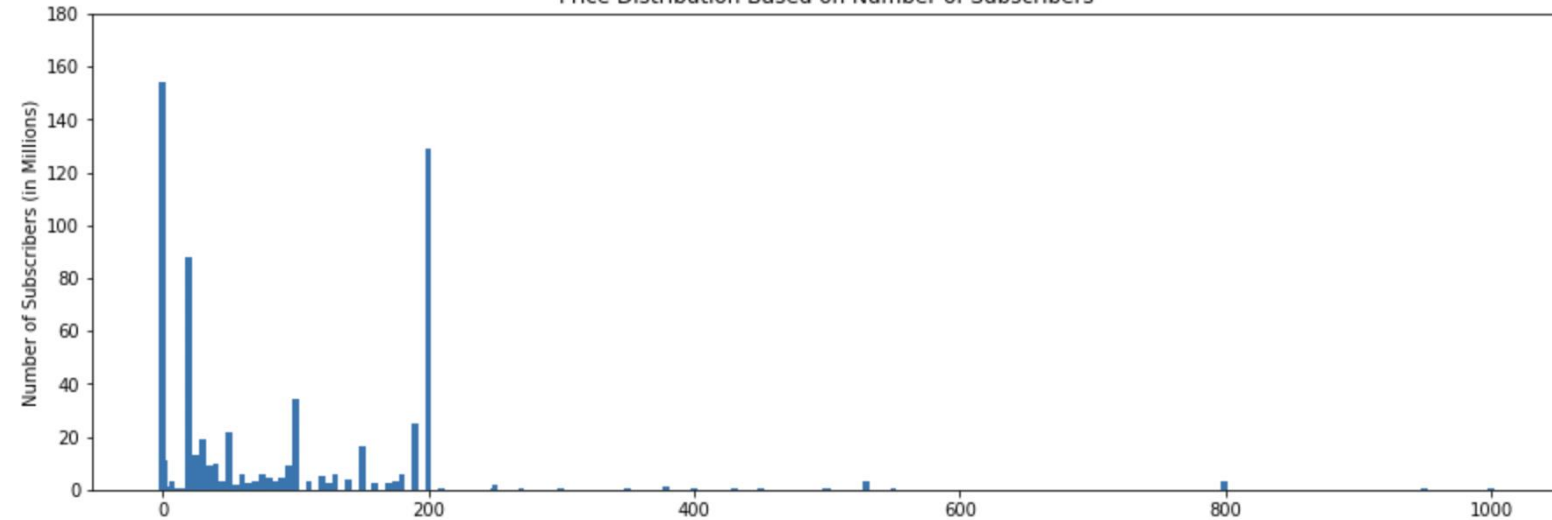
Year-wise count of courses updated



Year-over-year Growth Rates of Udemy Subscribers and Revenue

## Change in Average Course Price over the Years



| year | avg_price | price_change |
|------|-----------|--------------|
| 2010 | 21.24250030517578 | null |
| 2011 | 51.11403602466249 | 140.62 |
| 2012 | 45.717106322747846 | -10.56 |
| 2013 | 49.426245583308614 | 8.11 |
| 2014 | 57.59842641691168 | 16.53 |
| 2015 | 66.99407659489368 | 16.31 |
| 2016 | 74.35113059477587 | 10.98 |
| 2017 | 89.37098373886212 | 20.2 |
| 2018 | 98.13552318094766 | 9.81 |
| 2019 | 92.09890771377194 | -6.15 |
| 2020 | 81.08784826585313 | -11.96 |
| 2021 | 73.22534042648067 | -9.7 |
| 2022 | 84.3582485774976 | 15.2 |

# Overall Udemy Analysis
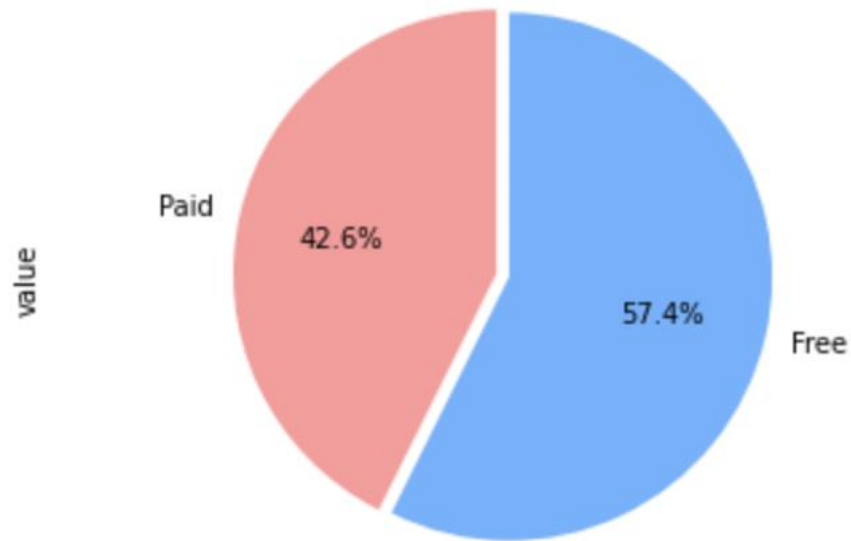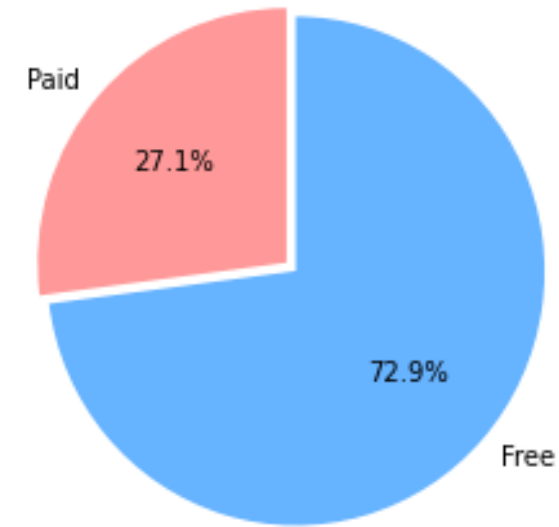
Ratio of Paid and Free Courses based on Number of Comments

Ratio of Subscribers per Course for Paid and Free Courses

- 89.6% are paid courses and 10.4% are unpaid courses.
- Despite the fact that Udemy offers a wide selection of courses ranging in price from $0.1 to $1000 US, as shown, the majority of them had original prices of less than $200.
- The number of subscribers enrolling the free courses has increased. People are willing to invest up to $200 for good content.
- Student engagement is almost similar for both paid and free courses.

## Top 10 Instructors by Earnings



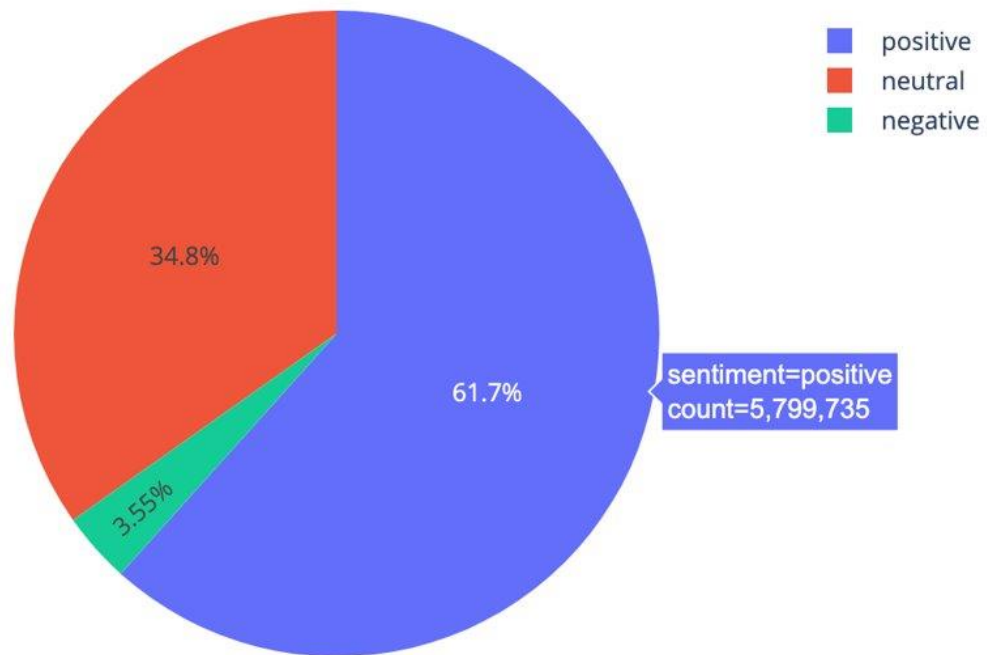| instructor_name | earned | avg_length | avg_price |
|---|---|---|---|
| Srinidhi Ranganathan | 1.73513164065625E9 | 74.85 | 273.00174986521404 |
| Learn Tech Plus | 1.1983609094453125E9 | 177.5759162303665 | 150.98476833323534 |
| TJ Walker | 1.1716159671679688E9 | 443.6555023923445 | 174.9900046626917 |
| Jose Portilla | 8.189048139404297E8 | 1025.5510204081634 | 172.74510624943946 |
| YouAccel Training | 8.01222185265625E8 | 329.64 | 133.8900029373169 |
| Creative Online School | 6.38075221E8 | 201.75 | 187.49000511169433 |
| Robert (Bob) Steele | 6.297034091914062E8 | 1374.439393939394 | 165.482504223333226 |
| Kirill Eremenko | 5.435664651762695E8 | 680.4444444444445 | 171.0736149681939 |
| Joseph Delgadillo | 5.434852525E8 | 748.0 | 165.37615849421576 |
| 365 Careers | 5.358052063120117E8 | 458.51807228915663 | 179.92988453141177 |
| Ing. Tomas Moravek | Facebook Ads Expert | 4.532673521875E8 | 1342.9772727272727 | 186.35432330044833 |
| Chris Haroun | 4.3560850329248047E8 | 463.2903225806452 | 124.90968033575243 |
| Dr. Angela Yu | 3.9845071703125E8 | 2941.5 | 184.36500453948975 |
| Academind by Maximilian Schwarzmüller | 3.782185150383301E8 | 1465.969696969697 | 157.41424629905006 |
| Stephane Maarek | AWS Certified Cloud Practitioner,Solutions Architect,Developer | 3.7408264850097656E8 | 573.0 | 114.85111342536078 |
| Joe Parys | 3.594630980625E8 | 353.05263157894734 | 191.6347420090123 |
| Sandor Kiss | 3.443059305E8 | 240.1 | 195.8233388264974 |

- A variety of course lengths were offered by Udemy.
- On this platform, short courses are more typical.
- There are courses with 0 to 1,000 minutes of content.
- It's noteworthy that 9,373 courses contain 0 lectures or videos because their content duration is nil.
- Most of these courses consist of practice exams.
- Additionally, graph demonstrates that while courses with 40 to 60 minutes of content are the most common, the number of courses with less than 40 minutes of content is minuscule.

```
+---------+---------+----+-------------------+-------------+--------------------+----------------+---------+
|       id|course_id|rate|               date| display_name|             comment|comment_language|sentiment|
+---------+---------+----+-------------------+-------------+--------------------+----------------+---------+
| 88962892|  3173036| 1.0|2021-06-30 01:54:25|        Rahul|I think a beginne...|              en| positive|
|125535470|  4913148| 5.0|2022-10-07 18:17:41|        Marlo|Aviva is such a n...|              en| positive|
| 68767147|  3178386| 3.5|2020-10-19 13:35:37|Yamila Andrea|Muy buena la intr...|              es| negative|
|125029758|  3175814| 5.0|2022-10-01 04:13:49|    Jacqueline|This course is th...|             en| positive|
| 76584052|  3174896| 4.5|2021-01-30 16:45:11|      Anthony|I found this cour...|              en| positive|
|124129784|  4693438| 1.0|2022-09-20 18:30:29|         Jiang|nothing informati...|             en|  neutral|
|121769970|  4693272| 3.5|2022-08-22 18:48:56|      Kenneth|Multiple spelling...|              en| positive|
| 57260120|  3168632| 5.0|2020-06-03 05:56:44|         Tony|Very unique way o...|              en| positive|
| 77427106|  3188362| 4.0|2021-02-10 05:19:29|       HIROKO|グルテンフリーのポイントがよくわか...|             ja|  neutral|
|103846020|  4164550| 4.5|2022-01-02 17:04:53|         Jess|Good Course!  Inf...|              en| positive|
+---------+---------+----+-------------------+-------------+--------------------+----------------+---------+
```
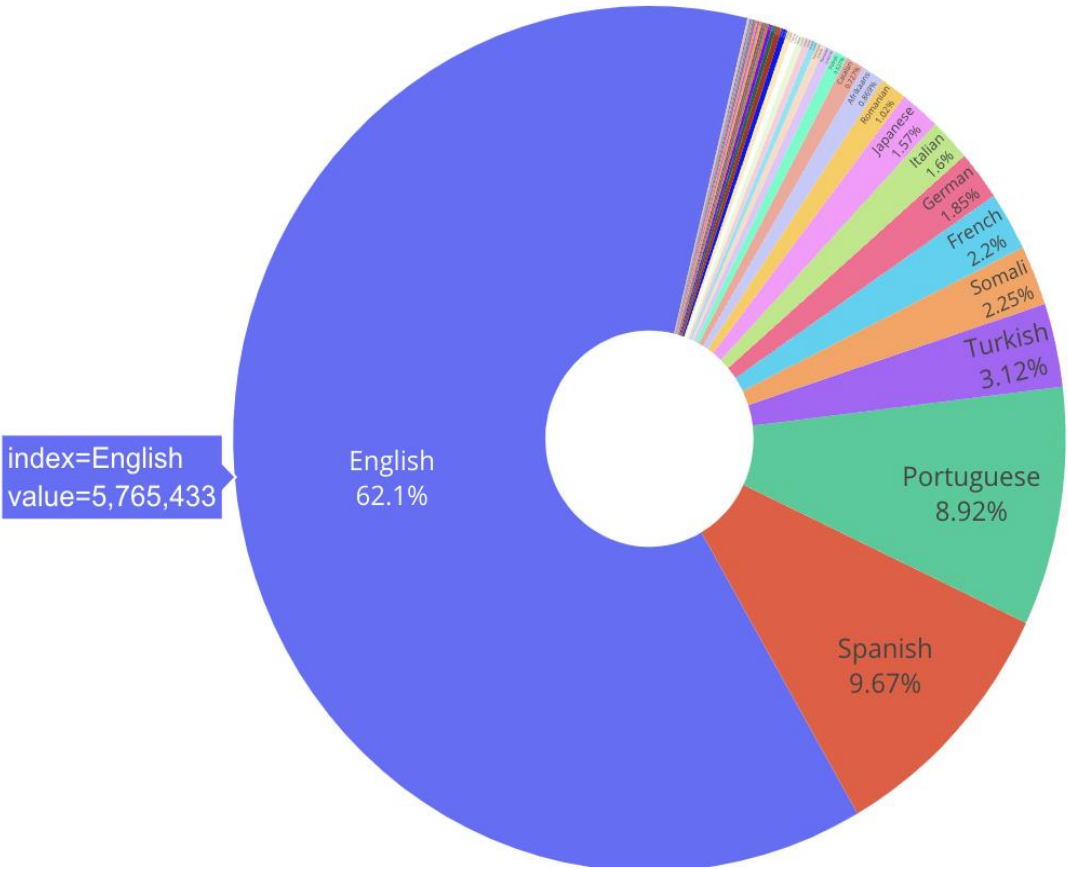
| instructor_name | title | positive_comment_count | negative_comment_count | neutral_comment_count |
|---|---|---|---|---|
| Jose Portilla | 2022 Complete Pyt... | 32355 | 1523 | 5147 |
| Kyle Pew | Microsoft Excel -... | 30148 | 1111 | 4815 |
| Colt Steele | The Web Developer... | 26805 | 959 | 3231 |
| Dr. Angela Yu | The Complete 2022... | 23671 | 873 | 3176 |
| Rob Percival | The Complete Digi... | 22761 | 896 | 3876 |
| Jonas Schmedtmann | The Complete Java... | 19538 | 753 | 2345 |
| Jose Portilla | The Complete SQL ... | 19350 | 751 | 2022 |
| Stephane Maarek |... | Ultimate AWS Cert... | 18625 | 405 | 2116 |
| Kirill Eremenko | Machine Learning ... | 18170 | 1228 | 3163 |
| Jaysen Batchelor | The Ultimate Draw... | 17507 | 848 | 3984 |



Top 10 Instructors with higest Number of Comments



Share of Languages in Comments

| instructor_name | title | is_paid | avg_rating | num_subscribers | price | course_id | positive_comment_count | negative_comment_count | neutral_comment_count |
|---|---|---|---|---|---|---|---|---|---|
| Sounds True | Sounds True Prese... | True | 5.0 | 9053.0 | 199.99 | 787672.0 | 720 | 18 | 106 |
| N.T. Wright | Paul: A Biography | True | 5.0 | 2181.0 | 54.99 | 1343784.0 | 250 | 4 | 35 |
| Samira Mian | Draw Islamic Geom... | True | 5.0 | 2987.0 | 19.99 | 819676.0 | 180 | 3 | 26 |
| Estelle Black | How To Master The... | False | 5.0 | 17883.0 | 0.0 | 1219774.0 | 176 | 9 | 27 |
| António Araújo | How to get Paid t... | True | 5.0 | 17390.0 | 99.99 | 914070.0 | 173 | 13 | 21 |
| Greg Reverdiau | Part 1 FAA Privat... | True | 5.0 | 2223.0 | 19.99 | 1798036.0 | 125 | 4 | 9 |
| N.T. Wright | The Resurrection ... | True | 5.0 | 1560.0 | 79.99 | 1048016.0 | 121 | 3 | 18 |
| Eric Arceneaux | Phase 2 - Becomin... | True | 5.0 | 2434.0 | 99.99 | 437148.0 | 118 | 5 | 20 |
| Michael C. Bush | THE 8 FACTORS: Ga... | True | 5.0 | 3204.0 | 49.99 | 80938.0 | 117 | 7 | 8 |
| Sorin Dumitrascu | PMI-ACP Certifica... | True | 5.0 | 13082.0 | 174.99 | 529326.0 | 114 | 3 | 5 |

only showing top 10 rows

# Overall Insights

- Udemy's membership base and number of published courses more than doubled during the Covid-19 pandemic. With 33% of all Udemy enrollments (468 million), the "Development" category has by far the most subscribers.
- Despite the affordable course costs, competent instructors made a large profit by offering top-notch courses.
- More subscribers signed up for courses that weren't paid for and costs around $20.
- If the information is good, it doesn't matter how long the course is; subscribers are willing to pay about $200.
- Most of the comments are positive, only a few offer neutral suggestions, and barely any have anything bad to say.
- English was the language that received the most comments, followed by Spanish and Portuguese.
- Jose Portilla, who is among the top 10 earners, received the most positive comments.

UMKC

# Thank You