



# Building a Course Recommend System



## Team Members:

- Safinaz Ahmed
- Mina Waheed
- Kerolos Adel

# 💡 Project Overview



## Main Idea: Course Recommender System

Build a Course Recommender System that helps users discover online courses based on their interests by analyzing course content and user behavior.



## Goal: Relevant & High-Quality Courses

Help users find **relevant** and **high-quality** courses easily. We use **content-based filtering** based on features like course title, subject, level and price.

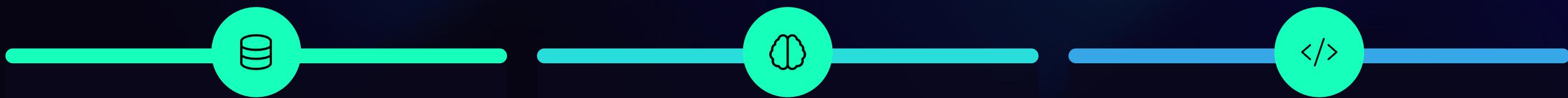


## Why This Matters: Personalized Discovery

Online learners face **too many choices**. Our system makes discovery **faster**, **smarter**, and **personalized**, streamlining the learning journey.

# Agenda

Today, we'll explore the journey of building a personalized course recommender system.



## Preprocessing & Visualizations

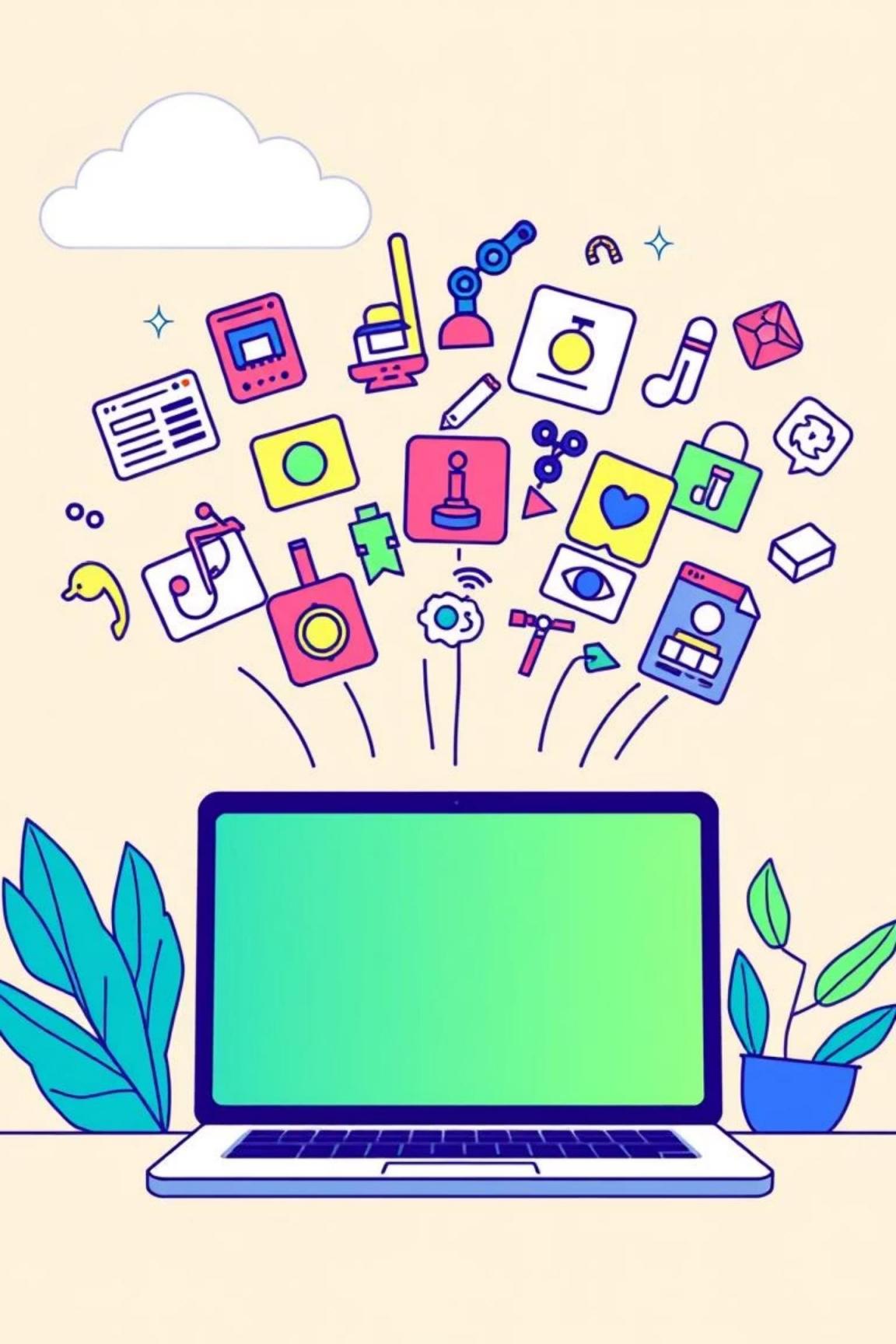
Preparing and understanding the Udemy course dataset.

## Modeling

Developing the content-based recommendation engine.

## Deployment & Interface

Bringing the system to users with an intuitive Streamlit app.



# Dataset Overview: Udemy Courses

Our foundation is a rich dataset of over 3,600 Udemy courses, providing a comprehensive view of online learning content.

- **Source:** Publicly available Udemy courses dataset.
- **Key Features:** Course title, description, subject, level, price, number of subscribers, reviews, and duration.
- **Diversity:** Spanning various subjects from web development to music, offering a broad spectrum for recommendations.

# Dataset Overview

A	B	C	D	E	F	G	H	I	J	K	
course_id	course_title	url	is_paid	price	num_subscribers	num_reviews	num_lectures	level	content_duration	published_timestamp	subject
1070968	Ultimate Investment Strategy Course	https://www.udemy.com/ultimate-investment-strategy-course/	TRUE	200	2147	23	51	All Levels	1.5	2017-01-18T20:58:58Z	Business Finance
1113822	Complete GST Course for Accountants	https://www.udemy.com/complete-gst-course-for-accountants/	TRUE	75	2792	923	274	All Levels	39	2017-03-09T16:34:20Z	Business Finance
1006314	Financial Modeling for Investment Banking	https://www.udemy.com/financial-modeling-for-investment-banking/	TRUE	45	2174	74	51	Intermediate Level	2.5	2016-12-19T19:26:30Z	Business Finance
1210588	Beginner to Pro - Python Programming	https://www.udemy.com/beginner-to-pro-python-programming-course/	TRUE	95	2451	11	36	All Levels	3	2017-05-30T20:07:24Z	Business Finance
1011058	How To Maximize Your Trading Profits	https://www.udemy.com/how-to-maximize-your-trading-profits/	TRUE	200	1276	45	26	Intermediate Level	2	2016-12-13T14:57:18Z	Business Finance
192870	Trading Penny Stocks for Profit	https://www.udemy.com/trading-penny-stocks-for-profit/	TRUE	150	9221	138	25	All Levels	3	2014-05-02T15:13:30Z	Business Finance
739964	Investing And Trading Like A Pro	https://www.udemy.com/investing-and-trading-like-a-pro/	TRUE	65	1540	178	26	Beginner Level	1	2016-02-21T18:23:12Z	Business Finance
403100	Trading Stock Charts For Profit	https://www.udemy.com/trading-stock-charts-for-profit/	TRUE	95	2917	148	23	All Levels	2.5	2015-01-30T22:13:03Z	Business Finance
476268	Options Trading For Dummies	https://www.udemy.com/options-trading-for-dummies/	TRUE	195	5172	34	38	Expert Level	2.5	2015-05-28T00:14:03Z	Business Finance
1167710	The Only Investment Strategy Course You'll Ever Need	https://www.udemy.com/the-only-investment-strategy-course-youll-ever-need/	TRUE	200	827	14	15	All Levels	1	2017-04-18T18:13:32Z	Business Finance
592338	Forex Trading Secrets Revealed	https://www.udemy.com/forex-trading-secrets-revealed/	TRUE	200	4284	93	76	All Levels	5	2015-09-11T16:47:02Z	Business Finance
975046	Trading Options For Dummies	https://www.udemy.com/trading-options-for-dummies/	TRUE	200	1380	42	17	All Levels	1	2016-10-18T22:52:31Z	Business Finance
742602	Financial Management For Dummies	https://www.udemy.com/financial-management-for-dummies/	TRUE	30	3607	21	19	All Levels	1.5	2016-02-03T18:04:01Z	Business Finance
794151	Forex Trading Course For Beginners	https://www.udemy.com/forex-trading-course-for-beginners/	TRUE	195	4061	52	16	All Levels	2	2016-03-16T15:40:19Z	Business Finance
1196544	Python Algo Trading Course	https://www.udemy.com/python-algo-trading-course/	TRUE	200	294	19	42	All Levels	7	2017-04-28T16:41:44Z	Business Finance
504036	Short Selling: Learn How To Short Sell Stocks	https://www.udemy.com/short-selling-learn-how-to-short-sell-stocks/	TRUE	75	2276	106	19	Intermediate Level	1.5	2015-06-22T21:18:35Z	Business Finance
719698	Basic Technical Analysis For Dummies	https://www.udemy.com/basic-technical-analysis-for-dummies/	TRUE	20	4919	79	16	Beginner Level	1.5	2016-01-08T17:21:26Z	Business Finance
564966	The Complete Python Developer Course	https://www.udemy.com/the-complete-python-developer-course/	TRUE	200	2666	115	52	All Levels	4	2015-08-10T21:07:35Z	Business Finance
606928	7 Deadly Mistakes When Trading Forex	https://www.udemy.com/7-deadly-mistakes-when-trading-forex/	TRUE	50	5354	24	23	All Levels	1.5	2015-09-21T18:10:34Z	Business Finance
58977	Financial Statement Analysis For Dummies	https://www.udemy.com/financial-statement-analysis-for-dummies/	TRUE	95	8095	249	12	Beginner Level	0.5833333333333333	2013-06-09T00:21:26Z	Business Finance
1242604	Winning Forex Trading Strategies	https://www.udemy.com/forex-trading-strategies/	TRUE	200	809	3	25	All Levels	2	2017-06-06T02:54:04Z	Business Finance
798740	Forex Traders - Complete Guide To Trading	https://www.udemy.com/forex-traders-complete-guide-to-trading/	TRUE	200	2295	84	39	All Levels	4	2016-05-02T19:26:48Z	Business Finance
506568	Create A Business Plan For Your Startup	https://www.udemy.com/create-a-business-plan-for-your-startup/	TRUE	75	10149	83	16	All Levels	2	2015-05-26T17:25:46Z	Business Finance

Activate Windows

# Exploratory Data Analysis (EDA)

understand key characteristics of each course

Column Name	Description
course_id	Unique identifier for each course
course_title	Title of the course
url	Link to the course on Udemy
is_paid	Indicates if the course is free or paid (TRUE = Paid)
price	Course price in USD
num_subscribers	Total number of students enrolled in the course
num_reviews	Number of student reviews submitted
num_lectures	Total number of video lectures in the course
level	Course difficulty level(e.g., Beginner, Intermediate, Expert, All Levels)
content_duration	Duration of the course content in hours
published_timestamp	Date and time when the course was published
subject	Category or domain of the course (e.g., Business Finance, Web Development)

# Data Cleaning

## Removed Duplicates

Ensured each course is unique by checking `course_id` and `course_title`

## Handled Missing Values

Checked for nulls in key columns (`price`, `num_subscribers`, `level`, etc.)

## Check Outliers

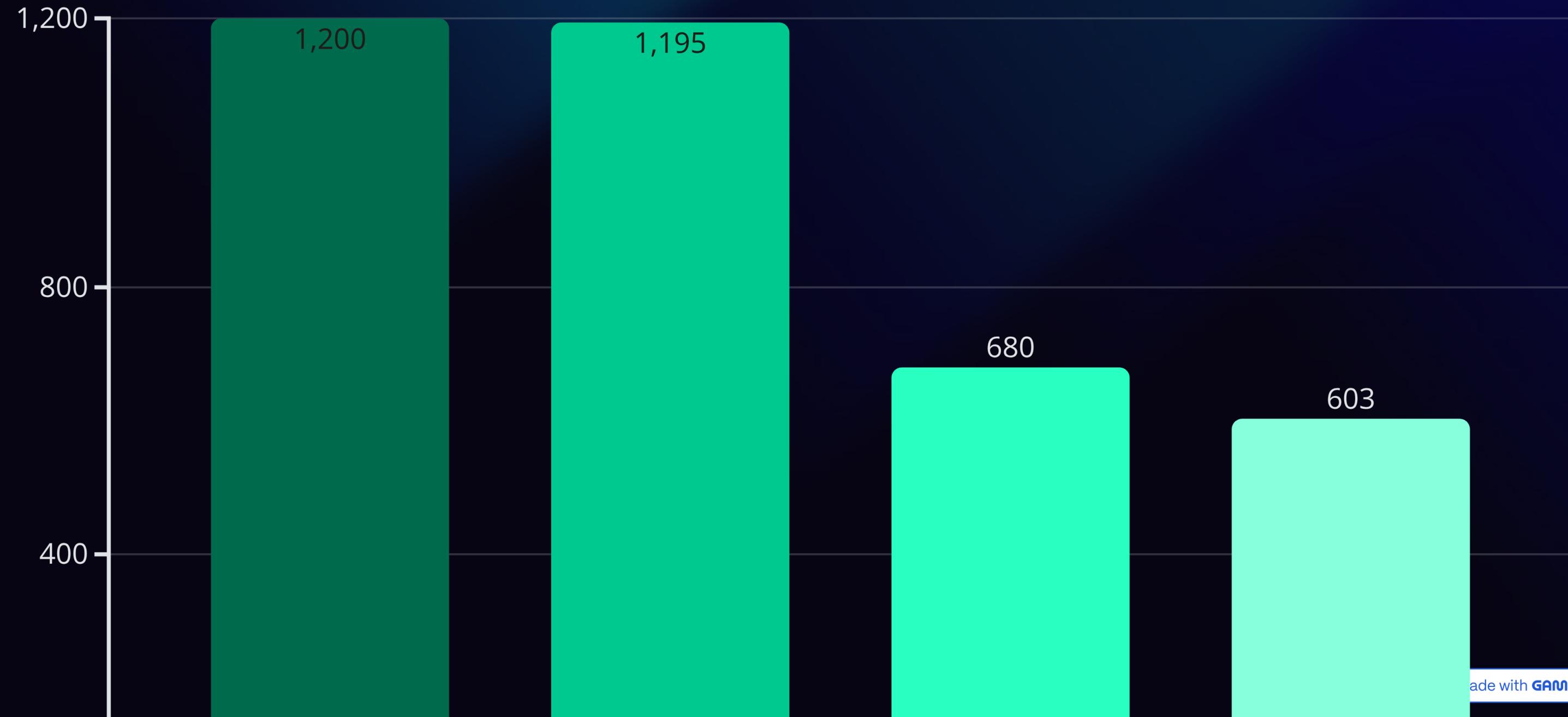
Used IQR to detect/remove unrealistic values

The dataset became **clean**, **consistent**, and **ready for preprocessing and modeling**

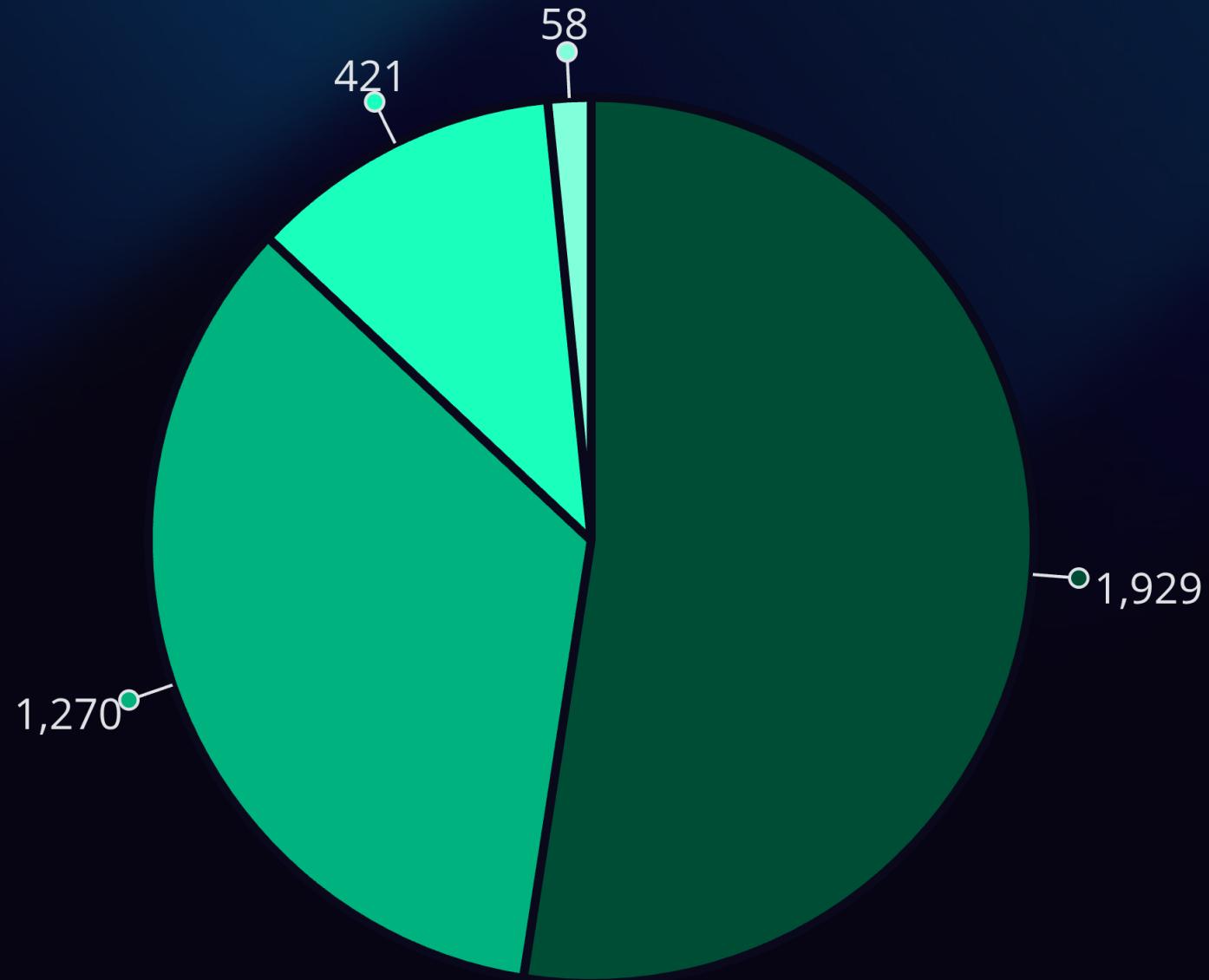
# Key Data Visualizations

Visualizing the dataset helps us understand distributions and potential correlations, guiding our feature selection.

Number of Courses per Subject



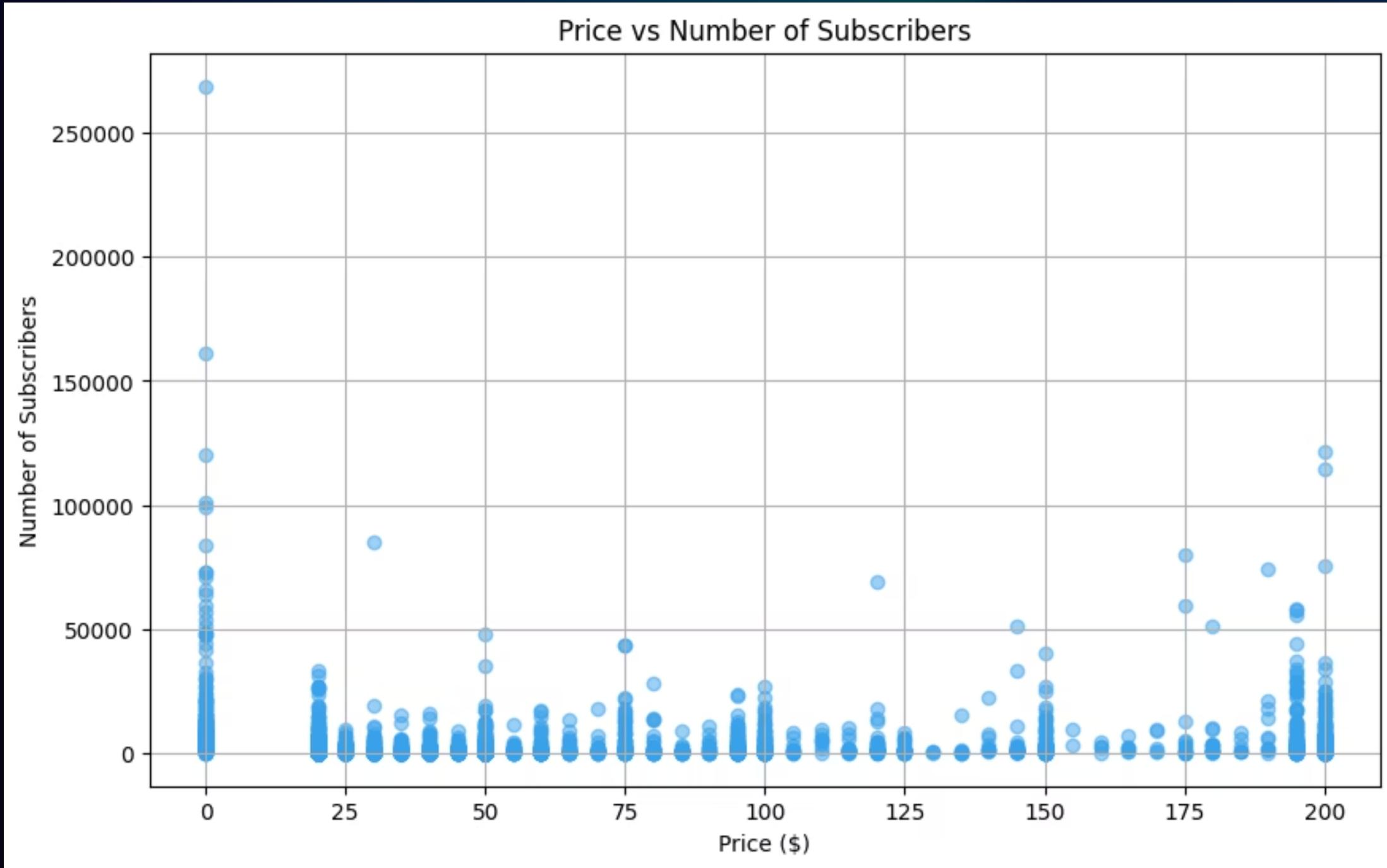
## Distribution of Courses by Level



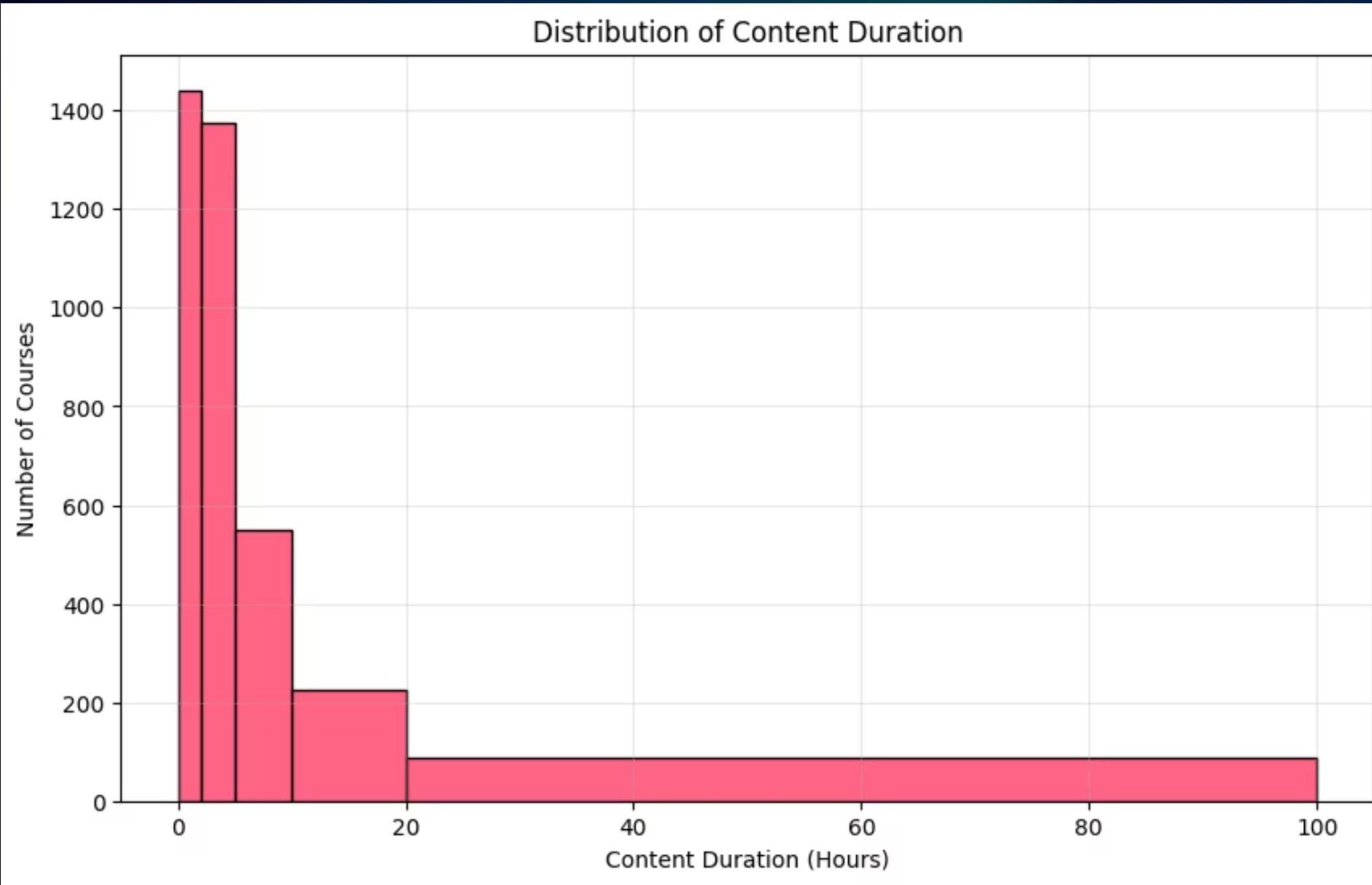
## Distribution of Subscribers by Subject



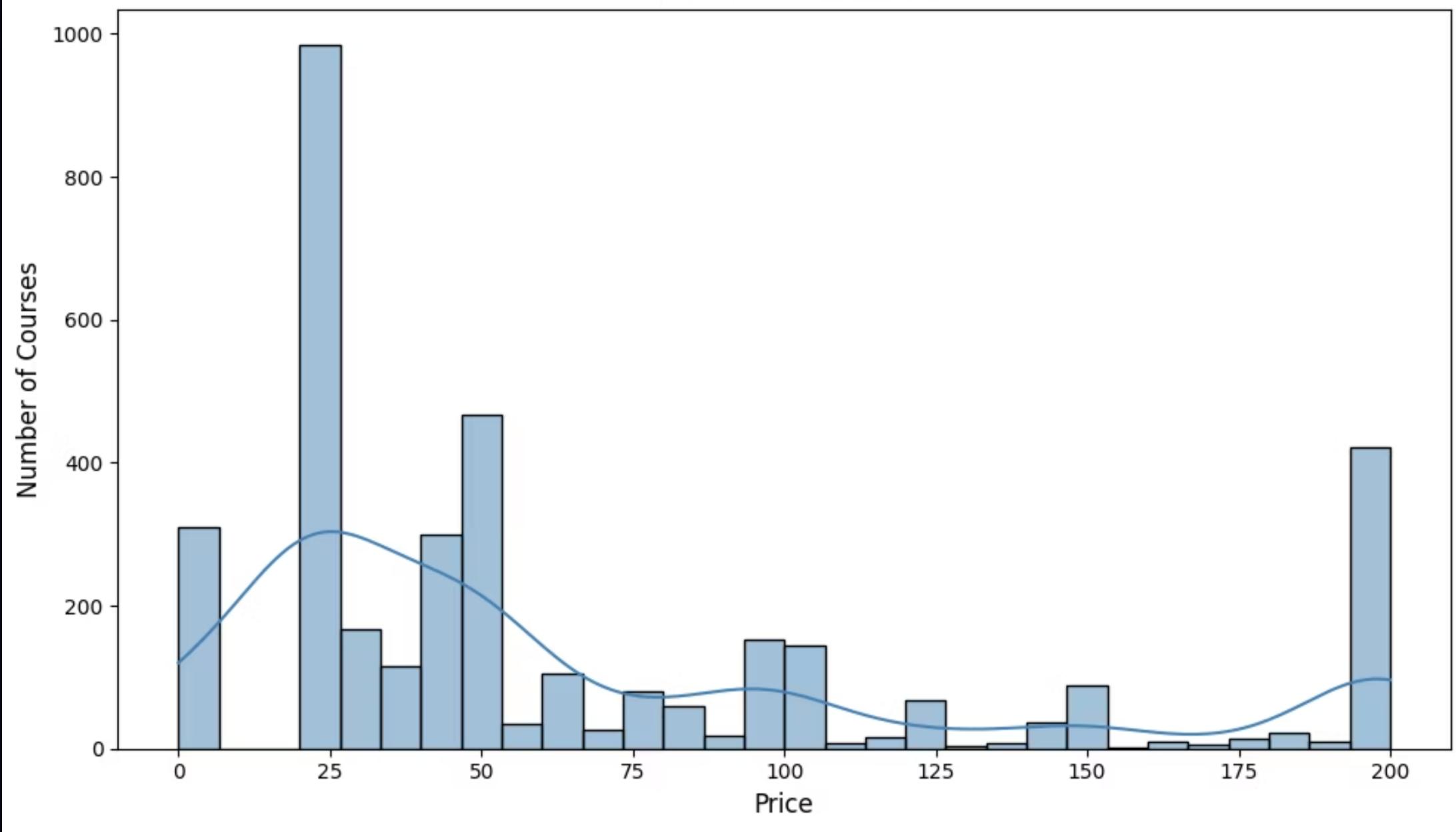
### Price vs Number of Subscribers

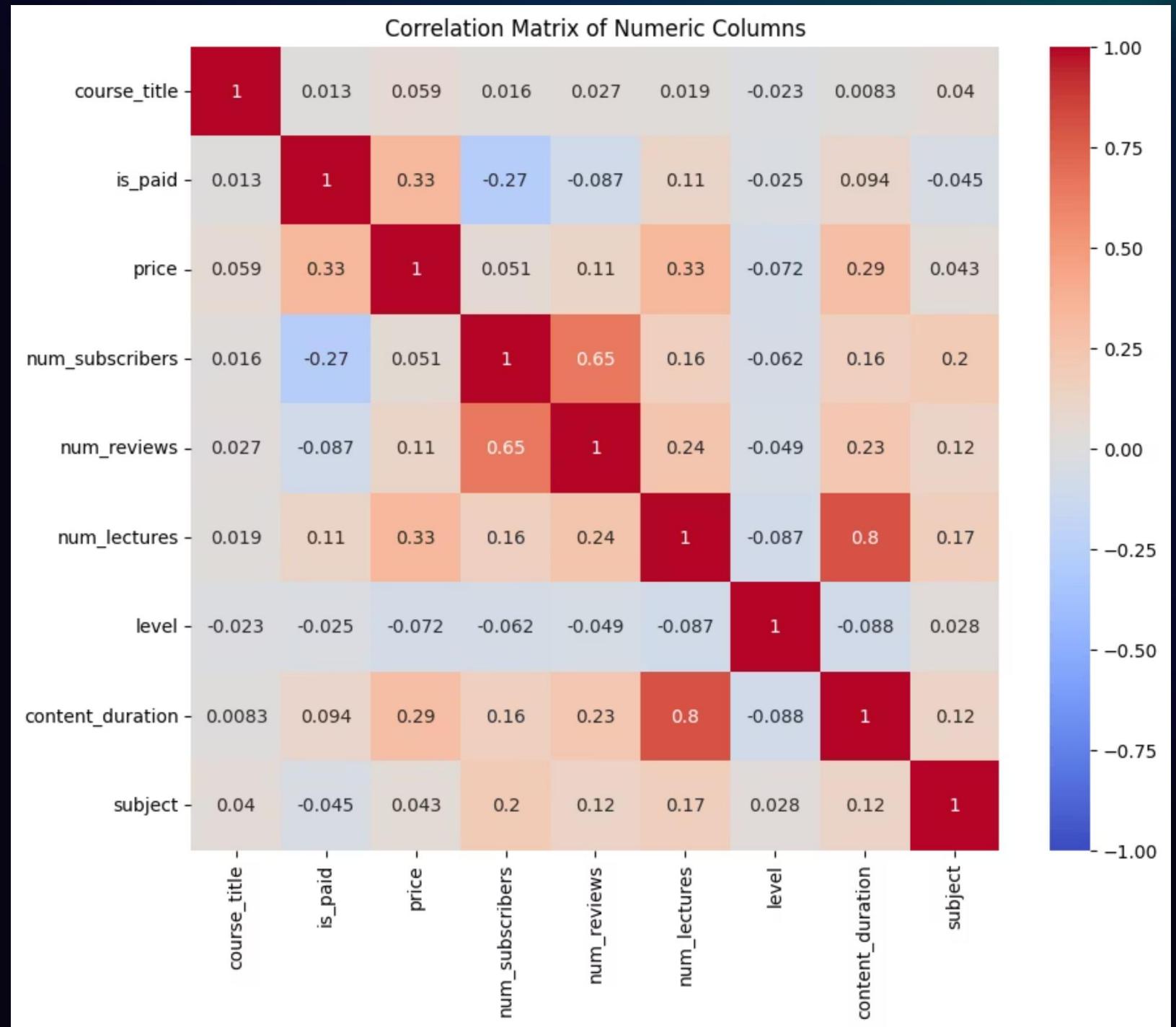


### Distribution of Content Duration



## Distribution of Course Prices





# Features Used for Clustering

## Textual Feature:

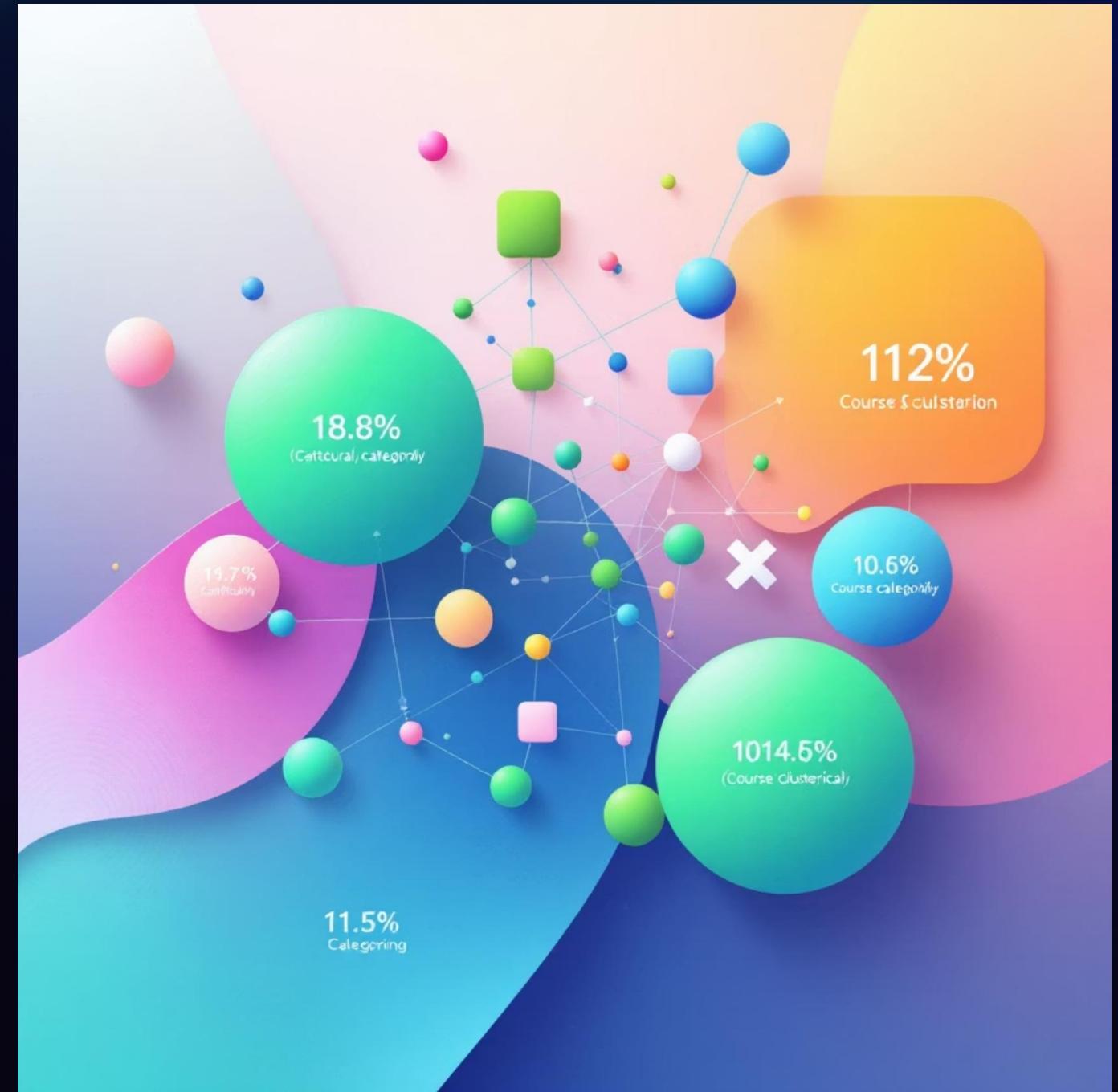
- course\_title → Processed using TF-IDF Vectorization
- Helps capture topic similarities between courses

## Categorical Features(Encoded):

- subject
- level

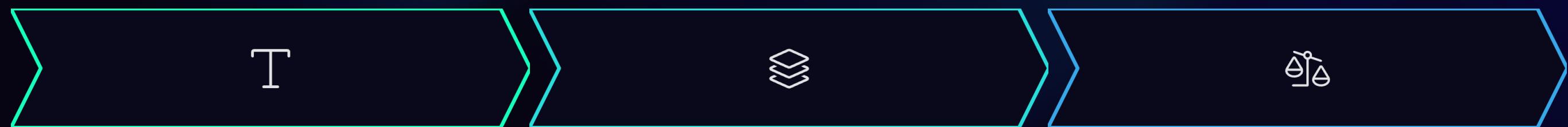
## Numeric Features(Scaled):

- price
- num\_subscribers
- num\_reviews



# Data Preprocessing Pipeline

Transforming raw course data into a structured format for effective modeling is crucial.



## Text Handling (TF-IDF)

Course titles and descriptions converted into numerical vectors for semantic understanding.

## Categorical Encoding

Subjects and levels are one-hot encoded to represent distinct categories numerically.

## Numerical Scaling

Subscriber counts and prices normalized to prevent dominance of larger values.

# Modeling: Content-Based Filtering

Our recommender system uses content-based filtering, matching courses based on their inherent characteristics.

1

## Core Algorithm

Cosine similarity to measure similarity between course feature vectors.

2

## Feature Weighting

Experimented with weighting features like course description (TF-IDF) higher for semantic relevance.

3

## Feature Experimentation

Tested combinations of text features (title, description) and numeric features (price, subscribers) for optimal performance.

First, I tried **K-Means** with **PCA** to reduce the dimensionality of the features—like price, number of subscribers, reviews, lectures, and content duration, plus encoded subject and level. I got some pretty good **Silhouette Scores**, which showed the clusters were tight and well-separated. Then, I went a step further and used **embeddings** for the course titles, generated with a model like all-MiniLM-L6-v2, to capture the semantic meaning of the titles. I combined these embeddings with the other features and ran K-Means again.

For the feature weights, I tested two approaches. One, I used **fixed weights**, giving more importance to stuff like num\_subscribers and the title embeddings. Two, I ran **Grid Search** to find the best weights for the features, testing different combinations. The Silhouette Scores came out really strong—honestly, they looked great, sometimes above 0.6, which is solid for clustering.

But, here's the thing: when I tested the clusters in a practical setting—like recommending courses based on them—the results weren't that accurate. Some predictions were off, and the recommendations didn't always make sense, like suggesting courses that weren't really similar.

So, I decided to try **DBSCAN** instead, thinking it might handle outliers better or find more natural clusters. I used the same setup: PCA in one run, embeddings in another, fixed weights, and Grid Search for weights, plus tuning eps and min\_samples. Again, the Silhouette Scores were good, especially with the best parameters from Grid Search. But, same problem—when I applied it practically, the recommendations were still off. Some courses ended up in the noise cluster, and others didn't group as expected.

Here's a sample of the results: for **K-Means** with **PCA**, I got a Silhouette Score around 0.5449 with 2 clusters. For **DBSCAN**, the run had a Silhouette Score of about 0.567 with 3 clusters. But in both cases, when I checked the actual course recommendations, they didn't align well with what users would expect.

# Appling Kmean With PCA

Explained Variance Ratio: [0.45419456 0.27698876 0.16038283 0.06905668 0.03937718]

Grid Search Results:

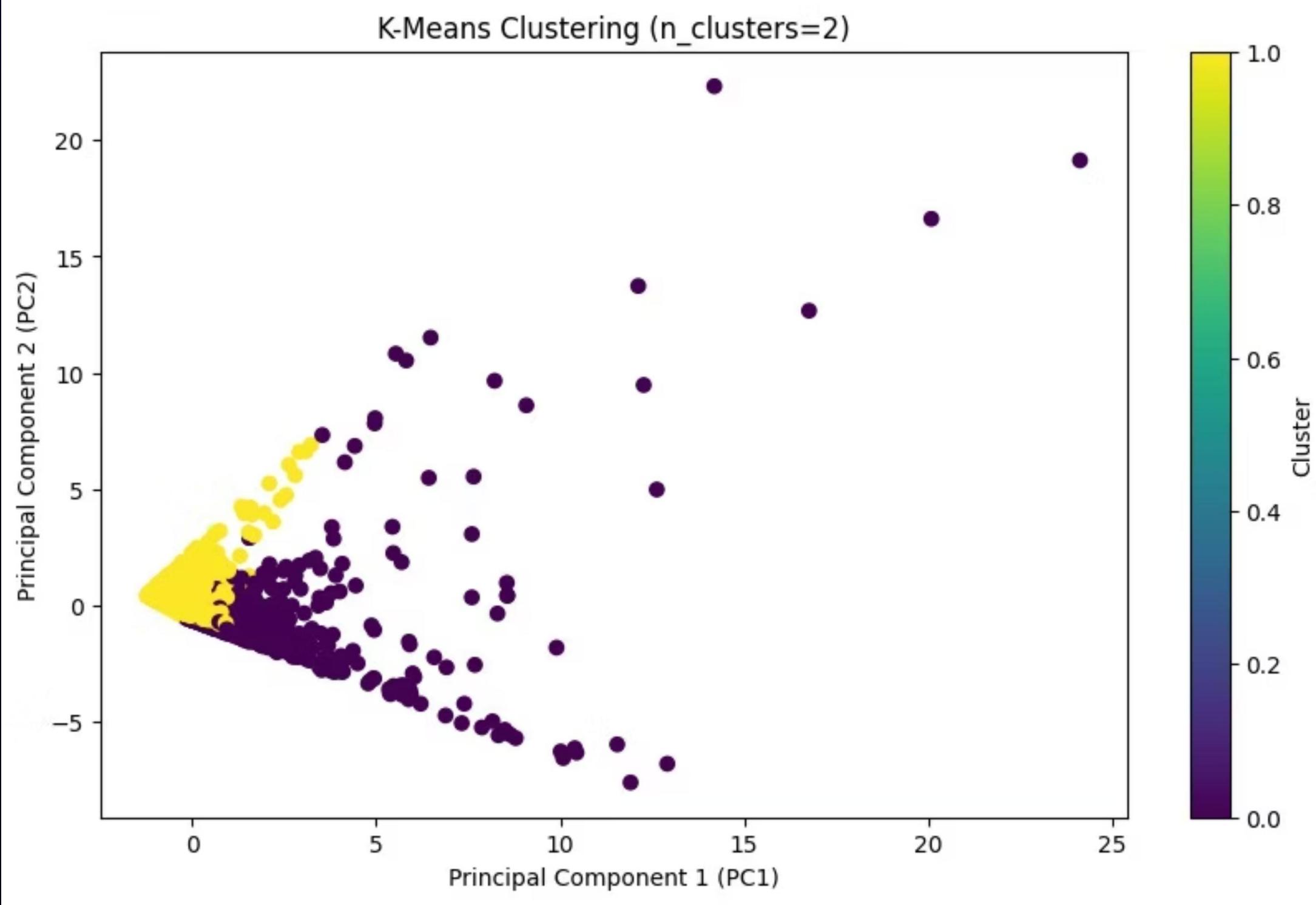
	n_clusters	silhouette_score	inertia
0	2	0.544908	14192.473969
1	3	0.506211	10787.611718
2	4	0.508255	7876.805244
3	5	0.522939	6751.834561
4	6	0.339572	6252.126414
5	7	0.366010	5004.939774
6	8	0.375513	4681.739173
7	9	0.380105	4492.191780
8	10	0.413572	4171.267535

Best Parameters:

n\_clusters: 2

Best silhouette Score: 0.5449084797102958

K-Means Clustering (n\_clusters=2)



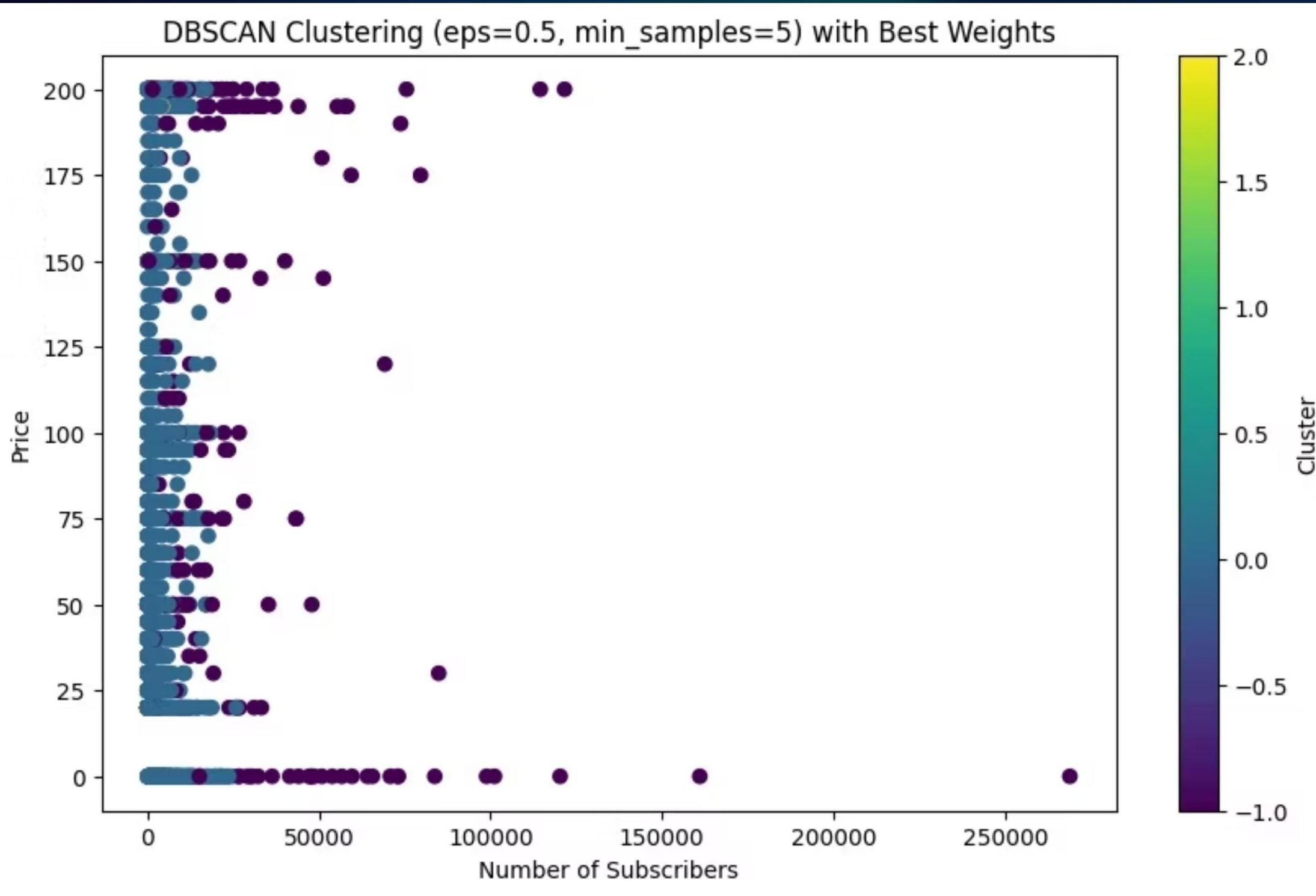
# Appling DBSCAN With grid search for Weights and parameters

```
Grid Search Results (Top 5 by silhouette Score):
    eps  min_samples                               weights \
4   0.5          5  {'price': 0.5, 'num_subscribers': 1.0, 'num_re...
8   0.7          7  {'price': 0.5, 'num_subscribers': 1.0, 'num_re...
7   0.7          5  {'price': 0.5, 'num_subscribers': 1.0, 'num_re...
17  0.7          7  {'price': 0.5, 'num_subscribers': 1.0, 'num_re...
16  0.7          5  {'price': 0.5, 'num_subscribers': 1.0, 'num_re...

    n_clusters  n_noise  silhouette_score
4            3      274        0.567055
8            2      198        0.562350
7            2      182        0.560877
17           2      229        0.560262
16           2      203        0.556389

Best Parameters:
eps: 0.5
min_samples: 5
Weights: {'price': 0.5, 'num_subscribers': 1.0, 'num_reviews': 1.0, 'num_lectures': 1.0, 'content_duration': 1.0}
Best Silhouette Score: 0.5670546299913514
Number of Clusters: 3
Number of Noise Points: 274
```

DBSCAN Clustering ( $\text{eps}=0.5$ ,  $\text{min\_samples}=5$ ) with Best Weights



But here's the problem: when I tested these models practically—like trying to recommend courses based on the clusters—the predictions were off. For example, with K-Means, I'd get recommendations that didn't really match the course's topic or user expectations. Same thing with DBSCAN; some courses ended up in the noise cluster, and others were grouped in ways that didn't make sense practically.

So, I switched to **Cosine Similarity** for recommendations. I used the same features—price, subscribers, reviews, and the title embeddings—and calculated the similarity between courses. This approach worked much better! The recommendations were more accurate, like suggesting Python courses for a Python course, or matching based on price and popularity. For instance, when I tested Cosine Similarity on a course, it recommended courses with similar titles and features, and they felt way more relevant compared to K-Means or DBSCAN.

Here's a sample of the results:

# Feature Weighting

To control the influence of each feature (textual, numeric, categorical) on the recommendation results – making the similarity more balanced and meaningful.

## \* Feature Groups & Weights

Feature Group	Features Included	Applied Weight
Textual	<code>course_title</code> (TF-IDF vectorized)	 High
Categorical	<code>subject, level</code> (one-hot encoded)	 Medium
Numeric	<code>price, num_subscribers,</code> <code>num_reviews</code> (scaled)	 Low

### Choosed Course

123

course_id	1100054
course_title	FOREX TRADING - Learn in a quick + profitable ...
url	<a href="https://www.udemy.com/forex-the-only-simple-tr...">https://www.udemy.com/forex-the-only-simple-tr...</a>
is_paid	True
price	100
num_subscribers	271
num_reviews	48
num_lectures	47
level	All Levels
content_duration	3.5
published_timestamp	2017-02-20T23:45:20Z
subject	Business Finance

# Choosed Course

Cosine Similarity (Original):

	course_id	course_title	url	is_paid	price	num_subscribers	num_reviews	num_lectures	level	content_duration	published_timestamp	subject	similarity
7	964828	Forex Trading - Learn An Effective Forex Trad...	https://www.udemy.com/forex-trading-strategy-t...	True	95	223	34	16	All Levels	1.5	2016-12-14T22:45:02Z	Business Finance	0.955458
9	324656	Forex Trading	https://www.udemy.com/how-to-trade-forex/	True	95	136	14	25	All Levels	4.5	2014-11-18T10:52:45Z	Business Finance	0.953922
3	575476	Forex trading made simple	https://www.udemy.com/forex-mentor-online-school/	True	100	45	4	31	All Levels	1.5	2015-09-25T19:08:09Z	Business Finance	0.947554

## Cosine Similarity

KMeans Recommendation (PCA):												
	course_id	course_title	url	is_paid	price	num_subscribers	num_reviews	num_lectures	level	content_duration	published_timestamp	subject
34	1041926	How to Trade Forex like a Hedge Fund: Long FX ...	https://www.udemy.com/how-to-trade-forex-like-a-hedge-fund-long-fx-trading/	True	115	504	5	32	All Levels	1.5	2017-03-23T16:23:33Z	Business Finance
44	876646	Direction-Independ Trading - Elite Forex Trade...	https://www.udemy.com/direction-independent-trading-elite-forex-trader/	True	120	973	15	28	All Levels	2.5	2016-06-21T22:46:41Z	Business Finance
15	1189592	Learn How To Successfully Trade Forex: In 5 Si...	https://www.udemy.com/learn-how-to-successfully-trade-forex-in-5-steps/	True	95	64	7	44	All Levels	8.5	2017-05-03T22:09:03Z	Business Finance

## Kmeans with PCA

	course_id	course_title	url	is_paid	price	num_subscribers	num_reviews	num_lectures	level	content_duration	published_timestamp	subject
10	944534	Forex:Trade Management & Psychology	https://www.udemy.com/forexmoney-management-ps...	True	150	13	0	19	Intermediate Level	2.0	2016-09-21T19:07:06Z	Business Finance
6	1247394	Cryptocurrency Trading: Complete Guide To Trad...	https://www.udemy.com/cryptocurrency-trading/	True	95	367	42	35	All Levels	5.0	2017-06-21T23:18:47Z	Business Finance
72	581598	Commodity Futures Day Trading Strategies	https://www.udemy.com/futures-day-trade-course/	True	95	91	5	28	All Levels	5.0	2015-08-28T23:47:29Z	Business Finance

## DBSCAN With PCA

# Deployment with Streamlit

The recommender system is deployed as an interactive web application using Streamlit, providing a seamless user experience.

The screenshot shows a Streamlit application interface titled "Course Recommender System". On the left, there is a sidebar titled "Filter Courses" with dropdown menus for "Choose Subject" (set to "Business Finance"), "Choose Level" (set to "All"), and a search bar for "Search by Title". The main area features a title "Course Recommender System" with a target icon. Below it is a "Choose Page" input set to "1" with a page navigation bar. Three course cards are displayed in a grid:

- Ultimate Investment Banking Course**  
Business Finance  
Duration: 1.5 hours  
Subscribers: 2,147  
Reviews: 23  
Price: 200
- Complete GST Course & Certification - Grow Your CA**  
Business Finance  
Duration: 39.0 hours  
Subscribers: 2,792  
Reviews: 923  
Price: 75
- Financial Modeling for Business Analysts and Consu**  
Business Finance  
Duration: 2.5 hours  
Subscribers: 2,174  
Reviews: 74  
Price: 45

At the bottom right, there is a message: "Go to Settings to activate Windows." A "Select" button is located at the bottom of each course card. In the top right corner of the main area, there are "Deploy" and three-dot buttons.



Ready to find your next course?



Let's Try the Interface

# Key Takeaways & Next Steps

This project demonstrates a robust approach to personalized course recommendations.

## Personalized Learning

Content-based filtering provides highly relevant course suggestions.

## Scalable Deployment

Streamlit offers a rapid and interactive deployment solution.

## Future Enhancements

Incorporating collaborative filtering and real-time user feedback for even smarter recommendations.