

ISBESTSELLER PREDICTION

PREPEARD BY:

RANYA ALKHTANI

EMAIL & GITHUB ACCOUNTS:

RANYAALKHTANI@GMAIL.COM

CLASSIFICATION PROBLEM
PROJECT REPORT



» ABSTRACT

The goal of this project was to make classification on sentiment analyzed data in order to predict whether the course will be bestseller or not. I worked with Udemy courses Dataset which I scribed. after exploring the data I come up with a target and an idea of model. I wanted to make prove of concept using Udemy courses Dataset. the classifier that I chosen as the Best model, have been going through a lot of experiences with different features and algorithms to improve its performance.

» DESIGN

This model is to provide the courses required in the market place, especially as they rely on the specific criteria to meet and succeed in this task. Predict whether the course is a best seller or not may be one of the most important criteria that motivates an instructors to offer that course. By using AI-based classification models, we can help the teacher to cover such a need.

» DATA

The data obtained by performing web scarbing on (<https://www.udemy.com>) website which include many courses with 65 languages and covering almost all subjects that anyone needs. The dataset consists of approximately 20,000 entries linked to 11 features. The features we concern with and get from this website for the courses are ["Title", "Instructor", "Description", "Rating", "Participants", "Original_Praice", "Discount_Price", "Hours", "num_lectures", "Level", "Catogory"].

» ALGORITHMS

- **Data Collection:** web scribing.
- **EDA: plotting :** null, correlation.
- **Data Cleaning:** Handling null, stop-word removal, punctuation removal, text lowercase, numbers removal.
- **Pre-Processing:**
 - **Feature preparation:** categorical feature encoding.
 - **Feature Engineering:** Tf_idf feature, word2vec.
- **Modeling:** Logistic Regression: KNN, Naive Bayes, Gradient Boosting, Random Forest.



› TOOLS

Technologies: python, google colab

Libraries: SKlearn, StatesModel,neighbors, svm, BeautifulSoup, Selenium, NumPy, Pandas, Matplotlib, Seaborn, langdetect, LinearRegression, time, os, Request, IPython.core.display.

› COMMUNICATION

in this part of the report I will show and document all of communication work has been shown both of visualisations and slide show:

AGENDA

- Goal
- Data Collection
- EDA
- Data Cleaning
- Pre-Processing
- Modeling
- Models Evaluation

GOAL

- Predicting the isBestseller.

DATA COLLECTION

- Web Scrapping "Udemy"
- 12 Features -> 8 Features
- Almost 20K Record -> 13K 'en'

EDA

• Impalacement checking

DATA CLEANING

• Null Handling

PRE-PROCESSING

• Feature preparation

• Feature Engineering

ED A

• Null Values

ED A

• Title word cloud

ED A

• Description word cloud

MODELING

FITTED DATA

MODELS

OverSampled

Tomeklink

SMOTE

ENN

TF_NF

Word2Vec

Logistic Regression

KNN

Naive Bayes

Random Forest

MODELS EVALUATION

Data	Models	F1	Precision	Recall	Accuracy
OverSampled	Logistic Regression	0.340	0.221	0.737	0.577
	Naive Bayes	0.350	0.246	0.744	0.586
	Random Forest	0.347	0.349	0.349	0.606
Tomeklink	Logistic Regression	0.150	0.009	0.387	0.178
	Naive Bayes	0.349	0.334	0.366	0.739
	Random Forest	0.322	0.485	0.241	0.660
SMOTE	Logistic Regression	0.361	0.248	0.601	0.604
	Naive Bayes	0.353	0.293	0.449	0.721
	Random Forest	0.360	0.347	0.364	0.608
ENN	Logistic Regression	0.360	0.228	0.612	0.528
	Naive Bayes	0.361	0.239	0.737	0.613
	Random Forest	0.361	0.225	0.620	0.632
TF_NF	Logistic Regression	0.293	0.212	0.476	0.601
	Naive Bayes	0.300	0.000	0.000	0.602
	Word2Vec	0.002	0.260	0.001	0.607

MODELS EVALUATION

ROC Curve - AUR value

Baseline Model No AUR = 0.5

LR sampled with ENN AUR = 0.788

KNN sampled with SMOTE AUR = 0.7

Naive Bayes sampled with ENN AUR = 0.701

RF sampled with over sampling AUR = 0.942

THANK YOU FOR LISTENING