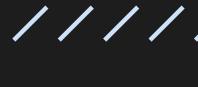
## Influencer Marketing

/////

ML-BASED INFLUENCER RECOMMENDER SYSTEM

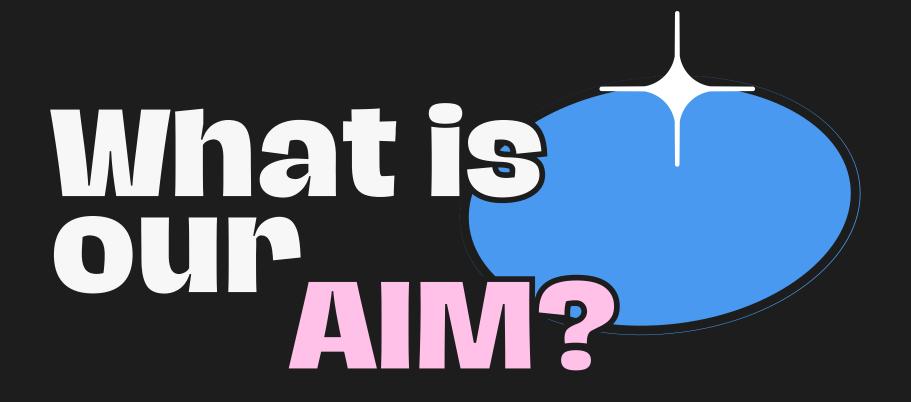
GROUP 14 - ADITYA, RUHI, DAKSH



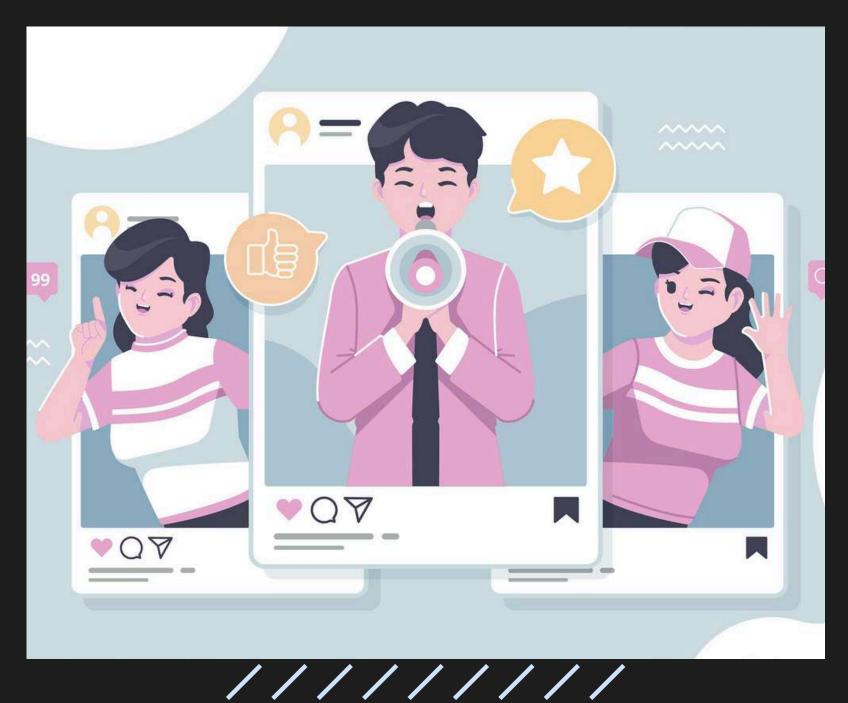




# "SECTION 1"



We aim to build a machine learning-powered recommender system to help brands find the most relevant influencers based on bio similarity, engagement, and past sponsorship data.





Over 70% of brands in India are currently investing in influencer marketing.

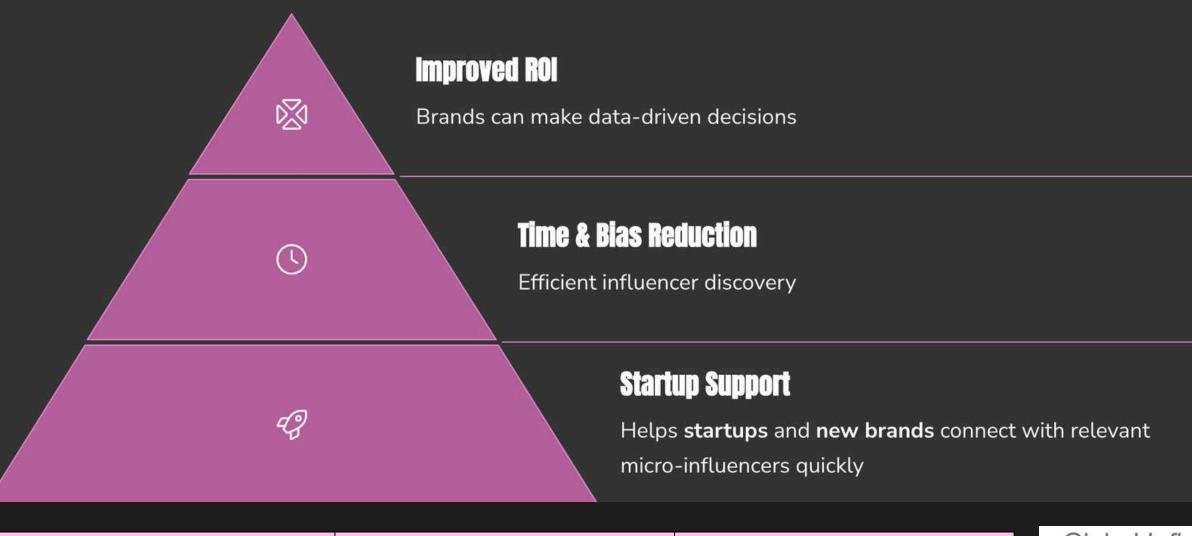
Influencer marketing is powerful but inefficient.

Brands struggle to find the right influencer from thousands of options.

Simple metrics (like follower count) aren't enough to guarantee campaign success.

Manual selection is time-consuming and often ineffective.

This project boosts cost-effective marketing by identifying high-impact influencers and detecting fake engagement. It empowers small businesses and refines brand messaging using NLP-driven social media analysis.



- 71% of Users more likely to purchase a given product if recommended by the right influencer
- This solution is scalable and can be deployed across various industries using influencer marketing.

Metric	Insight	Source
ROI Boost	Up to \$5.78 return for every \$1 spent with the right influencer.	Influencer Marketing Hub (2023)
Cost Savings	Relevant micro/nano influencers can cut campaign costs by 30–40%.	Later x Fohr (2023)
Engagement	Well-matched influencers deliver 2x higher engagement rates.	SocialPubli (2022)
Conversion Rate	Proper alignment drives 43% more conversions than untargeted outreach.	Kantar x CreatorIQ (2022)



# "SECTION 2"

## Literature Review

### How prior work tackles influencer selection?

Prior research has applied statistical regression, clustering, deep NLP, multi-task learning, and ML to identify and rank influencers. We build on these by combining lightweight transformers, CatBoost being our gradient-boosted trees, and audience sentiment into a production-ready pipeline.



## Modeling Influencer Marketing Campaigns in Social Networks

#### **Developed Solution:**

- Campaign Simulation Model
- Create a fake social network
- Set parameters (category, budget, engagement)
- Allow hiring to happen
- Simulate a campaign
- Measure the outcome

#### **Shortcomings:**

- No semantic alignment between brand and influencer
- Lack of audience sentiment analysis.
- No predictive performance metrics on real-world data, instead of training on past data, they built a "what-if" model that works on assumptions (category, engagement,etc)

Our solution predicts which influencers will most likely drive sponsorship success - grounded in semantic, sentiment, and engagement signals.

### Modeling Influencer Marketing Campaigns in Social Networks

Ronak Doshi\* , Ajay Ramesh\* , Student Member, IEEE, and Shrisha Rao , Senior Member, IEEE

Abstract—Social media are extensively used in today's world, and facilitate quick and easy sharing of information, which makes them a good way to advertise products. Influencers of a social media network, owing to their massive popularity, provide a huge potential customer base. However, it is not straightforward to decide which influencers should be selected for an advertizing campaign that can generate high returns with low investment. In this work, we present an agent-based model (ABM) that can simulate the dynamics of influencer advertizing campaigns in a variety of scenarios and can help to discover the best influencer marketing strategy. Our system is a probabilistic graph-based model that provides the additional advantage to incorporate real-world factors such as customers'

Doshi, R., Ramesh, A., & Rao, S. (2022). Modeling influencer marketing campaigns in social networks (arXiv:2106.01750v3). arXiv. https://doi.org/10.48550/arXiv.2106.01750

### Enhancing Influencer Marketing Strategies

#### **Developed Solution:**

- Takes influencer and brand features along with sponsorship history
- Data is cleaned and encoded
- Predicts sponsorship success based on sponsorships that already took place.

#### **Shortcomings:**

- Only caters to large influencers
- Lack of sentiment analysis
- Does not account for new influencers or brands

These gaps motivate our pipeline's use of explicit semantic and sentiment features, a transparent CatBoost model for clear feature importances, and a design that scales efficiently to larger, more diverse influencer populations.

Rivera, O. (2022). Enhancing influencer marketing strategies through machine learning: Predictive analysis of influencer-generated interactions (Master's thesis, KTH Royal Institute of Technology). Retrieved from https://kth.diva-portal.org/smash/get/diva2%3A1783645/FULLTEXTO1.pdf

## **Enhancing Influencer Marketing Strategies through Machine Learning**

Predictive Analysis of Influencer-Generated Interactions

#### **OLIMPIA RIVERA**

Date: July 1, 2023

Supervisor: Haibo Li

Examiner: Anders Hedman

School of Electrical Engineering and Computer Science

Swedish title: Förbättra Marknadsföringsstrategier Genom Maskininlärning

The field of influencer marketing has experienced rapid growth in recent years. However, uncovering the true effectiveness of this marketing approach remains a significant challenge. This thesis addresses the challenge of predicting the effectiveness of influencer marketing campaigns by employing advanced machine learning techniques, specifically the Auto Machine Learning framework Autogluon. With the aim of democratizing machine learning and empowering businesses in the influencer marketing domain, this work leverages Autogluon to predict the interactions generated by influencers when posting affiliate links. By evaluating various settings of AutoGluon and assessing the performance using metrics such as R-squared score, we observed promising results with good predictive accuracy. The findings from our study contribute to critical discussions in the field. This research offers a streamlined

## Using Machine Learning to Connect Brands with Influencers

#### **Developed Solution:**

- Looks at all previous "successful" brand-influencer pairings
- Analyzes patterns in those past collaborations
- Clusters influencers from similar collaborations together
- A brand fills out a three-step form (budget, category, reach)
- Model picks an influencer that falls in the matching cluster and recommends them

#### **Shortcomings:**

- No test set or quantitative metrics to validate clustering quality or recommendation precision.
- Unsupervised clustering can lead to groupings without clear labels so we dont really know that new recommendations truly match brand needs or not.
- It doesn't address each brand's unique category or needs
- Lack of audience sentiment analysis.

Master of Science and Engineering in Interaction Technology and Design

Using Machine Learning to Connect Brands with Influencers

by Jonathan HEDLUND

Master Thesis

February 9, 2022

Fall 2021

Master Thesis, 30 ECTS

Master of Science in Interaction Technology and Design, 300 ECTS

Examiner – Thomas Mejtoft Supervisor at UmU – Kalle Prorok

We frame it as a supervised CatBoost classifier with clear evaluation metrics, expose feature importances for explainabiliy.

Hedlund, J. (2021). Using machine learning to connect brands with influencers (Master's thesis, Umeå University). Retrieved from https://www.diva-portal.org/smash/record.jsf?pid=diva2%3A1636469

### Ranking Micro-Influencers

#### **Developed Solution:**

- Given a brand's profile, which micro-influencers should I partner with?
- Learns what makes a good match by looking at both their post captions and images.
- Gives a ranked list of most likely to fit the brand.

#### **Shortcomings:**

- It only counts likes/comments, but never looks at whether those reactions were positive or negative.
- Does not cosider an influener or brand's bio
- Only trained for small influencers.

#### **Performance Metrics:**

- Recall@10 ≥ 0.19
- Recall@50 = 0.55

Our model adds sentiment scores and deep text understanding, and uses CatBoost classifier to give clear influencer matches.

#### Ranking Micro-Influencers: a Novel Multi-Task Learning and Interpretable Framework

Adam Elwood

lastminute.com

Chiasso, Switzerland
adam.elwood@lastminute.com

Alberto Gasparin

lastminute.com

Chiasso, Switzerland
alberto.gasparin@lastminute.com

Alessandro Rozza

lastminute.com

Chiasso, Switzerland
alessandro.rozza@lastminute.com

Abstract—With the rise in use of social media to promote branded products, the demand for effective influencer marketing has increased. Brands are looking for improved ways to identify valuable influencers among a vast catalogue; this is even more challenging with "micro-influencers", which are more affordable than mainstream ones but difficult to discover. In this paper, we propose a novel multi-task learning framework to improve the state of the art in micro-influencer ranking based on multimedia content. Moreover, since the visual congruence between a brand and influencer has been shown to be good measure of compatibility, we provide an effective visual method for interpreting our models' decisions, which can also be used to inform brands' media strategies. We compare with the current state-of-the-art on a recently constructed public dataset and we show significant improvement both in terms of accuracy and model complexity. The techniques for ranking and interpretation presented in this work can be generalised to arbitrary multimedia ranking tasks

Elwood, A., Gasparin, A., & Rozza, A. (2021). Ranking micro-influencers: A novel multi-task learning and interpretable framework (arXiv:2107.13943v1). arXiv. https://doi.org/10.48550/arXiv.2107.13943



#### **Data Collection**

Publicly available data of 25,282 brands and 38,113 influencers and 16,01,074 posts.

#### **Pre-processing**

Handle missing data, Standardize features, Merge datasets, Compute metrics

#### **Feature Engineering**

Engagement metrics, Influencer Sentiment, Brand & influencer categories, Bio similarity, Comment Sentiments.

#### **Model Prediction**

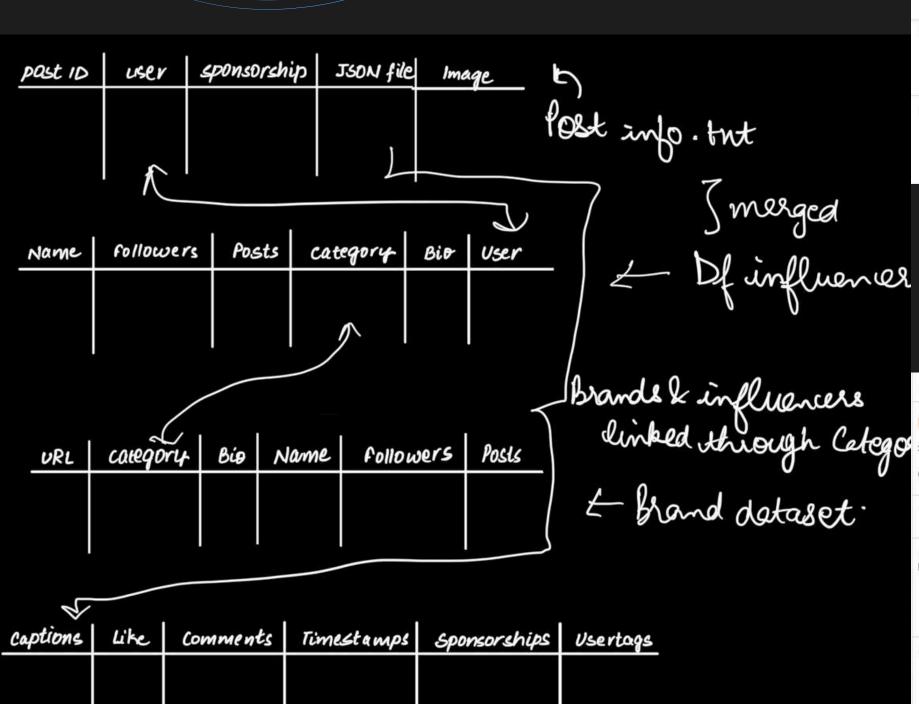
The CatBoostClassifier model (trained on historical sponsorship data) predicts the probability of sponsorship for each influencer.

#### **Recommendation System**

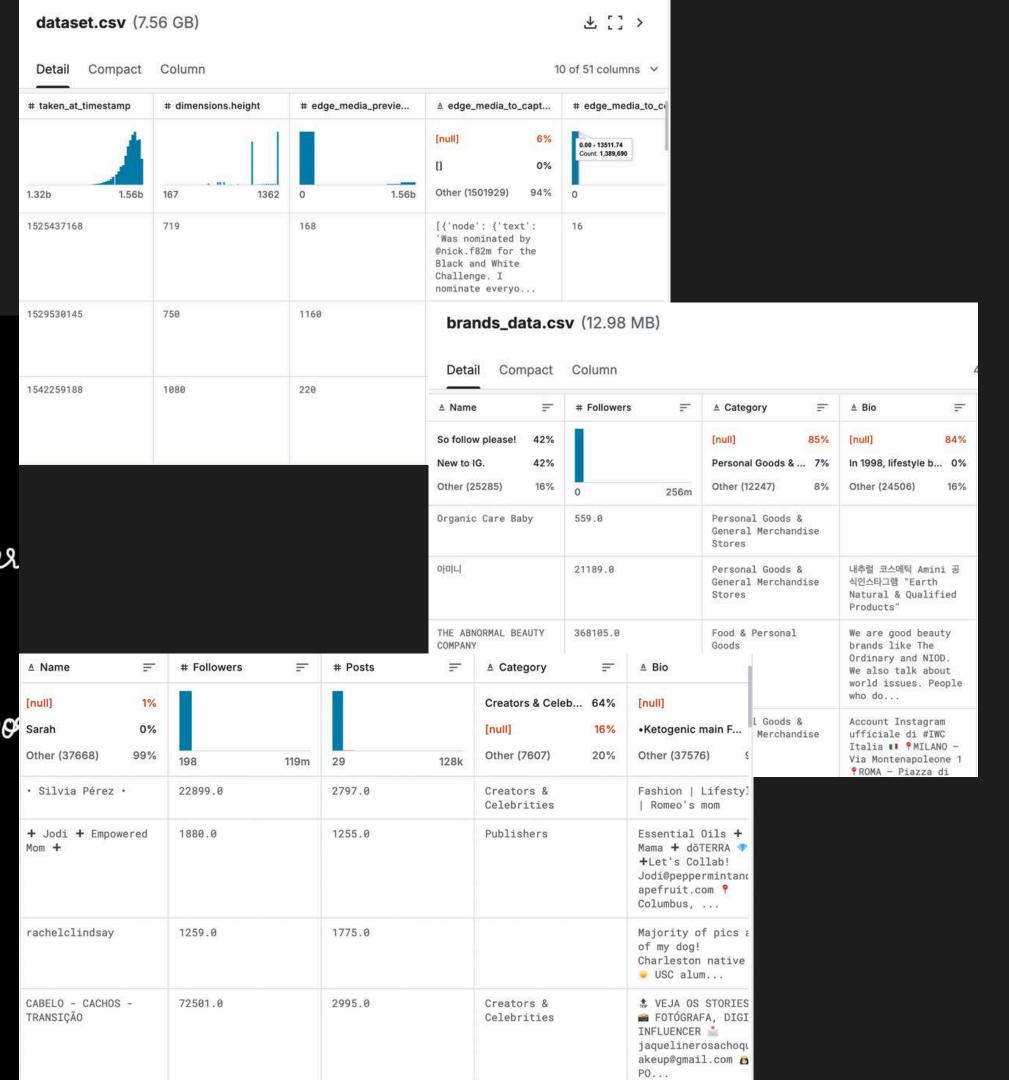
Predict sponsorship probability for each influencer, Rank influencers by predicted probability and recommend top relevant influencers (hopefully).

# "SECTION 3"





15son files

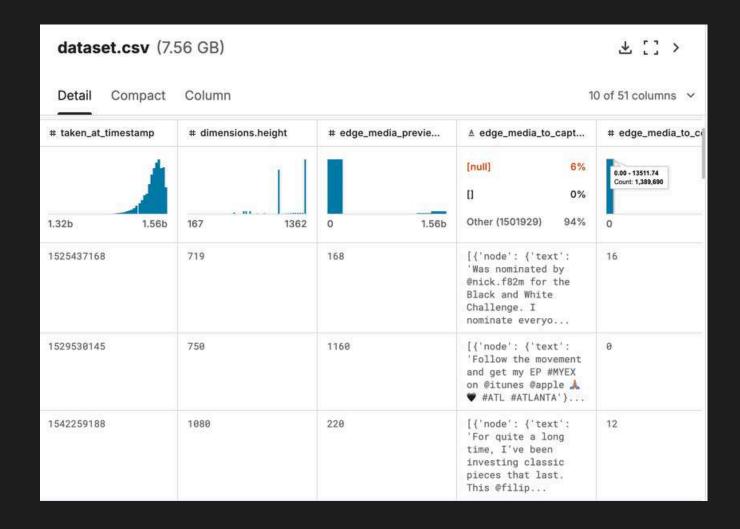




- Brand Username
- Brand Bio
- Brand Category
- Influencer Username
- Influencer Name
- Influencer Category
- Influencer Bio

- Influencer Followers
- Engagement Rate
- Sponsorship
- Bio Similarity
- Brand Name
- Followers
- Sentiment Score

	Brand Username	Brand Bio	Brand Category	Influencer Username	Influencer Name	Influencer Category	Influencer Bio	Influencer Followers	Influencer Engagement Rate	sponsorship_final	bio_similarity	Brand Name	Followers	sentiment_score
0	pukkaherbs	Celebrating how delicious, organic herbal teas	Grocery & Convenience Stores	wellness_ed	Amy Hopkinson	Publishers	I can't dance but I can burpee 🜉 Digital Ed @W	46736.0	1.076717	1	0.268495	Pukka Herbs	61353.0	0.9783
1	pukkaherbs	Celebrating how delicious, organic herbal teas	Grocery & Convenience Stores	niomismart	Niomi Smart	Unknown		1685570.0	0.421798	1	0.334331	Pukka Herbs	61353.0	0.9993
2	pukkaherbs	Celebrating how delicious, organic herbal teas	Grocery & Convenience Stores	bbaldz	brooke baldwin	Creators & Celebrities	LA. ▼: brookebaldwin99@gmail.com	26254.0	31.429285	0	0.067289	Pukka Herbs	61353.0	0.0000
3	pukkaherbs	Celebrating how delicious, organic herbal teas	Grocery & Convenience Stores	lauraschra	Lau	Home Goods Stores	-	31572.0	31.429285	0	0.054145	Pukka Herbs	61353.0	0.0000
4	pukkaherbs	Celebrating how delicious, organic herbal teas	Grocery & Convenience Stores	beautyandthebambino	ALIAEL	Unknown	Living + Loving + Life aliajacksonc@gmail.com	9219.0	31.429285	0	0.205278	Pukka Herbs	61353.0	0.0000
				***			***			1.1.			***	
38230	europeanwax	Revealing	Transportation & Accomodation Services	_davidecorsini_	Davide Corsini	Creators & Celebrities	#lifestyle #sportaddicted #mediainfluencer '**	62339.0	3.185205	0	0.217633	European Wax Center	35206.0	0.9775
38231	europeanwax	Pevesling	Transportation & Accomodation Services	rheamendezona	Maria Regina Ramas Mendezona	Creators & Celebrities	@thesocialtumbleweed Wife to @marcomendezona	6014.0	29.794501	0	0.175660	European Wax Center	35206.0	0.9972



Detail Compact Column							
∆ Name =	# Followers	△ Category =	∆ Bio =				
So follow please! 42%  New to IG. 42%  Other (25285) 16%	0 256m	[null] 85% Personal Goods & 7% Other (12247) 8%	[null] 84% In 1998, lifestyle b 0% Other (24506) 16%				
Organic Care Baby	559.0	Personal Goods & General Merchandise Stores					
아미니	21189.0	Personal Goods & General Merchandise Stores	내추럴 코스메틱 Amini 공 식인스타그램 "Earth Natural & Qualified Products"				
THE ABNORMAL BEAUTY COMPANY	368105.0	Food & Personal Goods	We are good beauty brands like The Ordinary and NIOD. We also talk about world issues. People who do				
IWC Schaffhausen	12315.0	Personal Goods &	Account Instagram				

- Contains JSON files of the 16,01,074 posts.
- JSON files have various information such as captions, likes, comments, timestamps, sponsorship, usertags, etc.
- Narrowed down to relevant attributes (Name, Followers, Category, Bio)
- Data preprocessing done on Kaggle.

- Contains Instagram profiles of the 25,282 brands scraped by authors.
- A dataset of 10 brand attributes (eg. Name, Followers, Following, No. of posts, Category, URLs, Bio, Email, Phone no., Profile picture.)
- Narrowed down to 4 relevant attributes (Name, Followers, Category, Bio)

A Name	# Followers =	# Posts =	A Category =	A Bio
[null] 1% Sarah 0% Other (37668) 99%	198 119m	29 128k	Creators & Celeb 64% [null] 16% Other (7607) 20%	[null] •Ketogenic main F Other (37576)
• Silvia Pérez •	22899.0	2797.0	Creators & Celebrities	Fashion   Lifesty]   Romeo's mom
+ Jodi + Empowered Mom +	1880.0	1255.0	Publishers	Essential Oils + Mama + dōTERRA * +Let's Collab! Jodi@peppermintanc apefruit.com * Columbus,
rachelclindsay	1259.0	1775.0		Majority of pics a of my dog! Charleston native
CABELO - CACHOS - TRANSIÇÃO	72501.0	2995.0	Creators & Celebrities	<pre></pre>

- Contains Instagram profiles of the 38,113 influencers scraped by author.
- A dataset of 11 influencer attributes (Name, followers, No. of posts, URLs, category, bio, email, phone no., profile picture, username).
- Narrowed down to 5 relevant attributes (Name, Followers, No. of posts, Category, Bio)

	Brand Username	Brand Bio	Brand Category	Influencer Username	Influencer Name	Influencer Category	Influencer Bio	Influencer Followers	Influencer Engagement Rate	sponsorship_final	bio_similarity	Brand Name	Followers	sentiment_score
o	pukkaherbs	Celebrating how delicious, organic herbal teas	Grocery & Convenience Stores	wellness_ed	Amy Hopkinson	Publishers	I can't dance but I can burpee ➡️ Digital Ed @W	46736.0	1.076717	1	0.268495	Pukka Herbs	61353.0	0.9783
1	pukkaherbs	Celebrating how delicious, organic herbal teas	Grocery & Convenience Stores	niomismart	Niomi Smart	Unknown	□ Lifestyle Vlogger       □ Author of Eat Smart       □ C	1685570.0	0.421798	1	0.334331	Pukka Herbs	61353.0	0.9993
2	pukkaherbs	Celebrating how delicious, organic herbal teas	Grocery & Convenience Stores	bbaldz	brooke baldwin	Creators & Celebrities	LA. ▼: brookebaldwin99@gmail.com	26254.0	31.429285	0	0.067289	Pukka Herbs	61353.0	0.0000
3	pukkaherbs	Celebrating how delicious, organic herbal teas	Grocery & Convenience Stores	lauraschra	Lau	Home Goods Stores	- e - our house in pictures '1904' - Drenthe	31572.0	31.429285	0	0.054145	Pukka Herbs	61353.0	0.0000
4	pukkaherbs	Celebrating how delicious, organic herbal teas	Grocery & Convenience Stores	beautyandthebambino	ALIAEL LISON	Unknown	Living + Loving + Life  aliajacksonc@gmail.com	9219.0	31.429285	0	0.205278	Pukka Herbs	61353.0	0.0000
***	***	***	(1999)	***	(99)	***		•••	***	//**	***	,,,	***	***
38230	europeanwax	Revealing You. Revealing Beautiful Skin.	Transportation & Accomodation Services	_davidecorsini_	Davide Corsini	Creators & Celebrities	#lifestyle #sportaddicted #mediainfluencer 🤾	62339.0	3.185205	0	0.217633	European Wax Center	35206.0	0.9775
38231	europeanwax	Revealing You. Revealing Beautiful Skin.	Transportation & Accomodation Services	rheamendezona	Maria Regina Ramas Mendezona	Creators & Celebrities	@thesocialtumbleweed Wife to @marcomendezona	6014.0	29.794501	0	0.175660	European Wax Center	35206.0	0.9972

- Obtained after merging and preprocessing the initial datasets.
- For every brand, there are all sponsored influencers + 10 unsponsored influencers

## Ethical considerations

### The dataset used in our project was originally collected by the authors through web scraping from Instagram

- Scraping enabled collection of attributes like followers, posts, bios, and engagement metrics that are not easily accessible via APIs.
- Violates Instagram's (Meta's) Terms of Service, which strictly prohibit unauthorised automated access.
- Unstable as Instagram can change its structure or block IPs at any time

- To avoid these issues, the author used the official Instagram Graph API instead of scraping, which requires approval and access tokens, for legitimate and scalable data collection.
- Access only allowed to Instagram Business or Creator accounts



```
edge_media_to_tagged_user.edges \
     sponsor_user_username
          drogariasaopaulo [{'node': {'user': {'full name': 'Vichy Brasil...
122
           uniquemonbijoux [{'node': {'user': {'full_name': 'Maysa Leão '...
140
           homesensecanada [{'node': {'user': {'full_name': 'Style at Hom...
580
                   zalando [{'node': {'user': {'full_name': 'ALLSAINTS', ...
600
                     stage [{'node': {'user': {'id': '1298729246', 'usern...
710
           outletdelcorso [{'node': {'user': {'full_name': 'OUTLET DEL C...
806
3250
                    iorane
                  shopiixc [{'node': {'user': {'full name': 'Erika Munro ...
3658
             veuveclicquot [{'node': {'user': {'full_name': 'Veuve Clicqu...
4428
        constance calcados [{'node': {'user': {'full name': 'Constance Ca...
4698
              sponsor_user_full_name
122
                        Vichy Brasil
140
      Maysa Leão 🦙 Marketing Digital
                       Style at Home
580
600
                           ALLSAINTS
710
                                None
806
                    OUTLET DEL CORSO
3250
                                None
          Erika Munro Kennerly, Esq.
3658
                  Veuve Clicquot USA
4428
                  Constance Calcados
4698
```

- Cleaned the tagged user text to remove unwanted characters.
- Isolated the relevant part of the tagged text.
- Extracted sponsor full names from the tagged users column.
- Matched the extracted names with the sponsor usernames already provided.

```
Category \
0 Personal Goods & General Merchandise Stores
  Personal Goods & General Merchandise Stores
                        Food & Personal Goods
  Personal Goods & General Merchandise Stores
        Transportation & Accomodation Services
                                               Bio
  내추럴 코스메틱 Amini 공식인스타그램 "Earth Natural & Oualif...
2 We are good beauty brands like The Ordinary an...
3 Account Instagram ufficiale di #IWC Italia 💵 ...
4 Welcome to the Instagram account for Four Seas...
Final Dataset Sample:
                            JSON_File
                                            Name Followers Posts \
        User
0 alisasia 1309041812857818435.json Alisa Sia
                                                  41099.0 691.0
  alisasia 1311846669234786866.json Alisa Sia
                                                  41099.0 691.0
2 alisasia 1315560311952470229.json Alisa Sia
                                                  41099.0 691.0
3 alisasia 1318531733175899446.json Alisa Sia
                                                  41099.0 691.0
4 alisasia 1343280729400051114.json Alisa Sia
                                                  41099.0 691.0
                Category
O Creators & Celebrities Los Angeles Snapchat & Twitter: AlisaSia alisa...
1 Creators & Celebrities Los Angeles Snapchat & Twitter: AlisaSia alisa...
2 Creators & Celebrities Los Angeles Snapchat & Twitter: AlisaSia alisa...
3 Creators & Celebrities Los Angeles Snapchat & Twitter: AlisaSia alisa...
4 Creators & Celebrities Los Angeles Snapchat & Twitter: AlisaSia alisa...
```

Cleaning missing data

					+ D-+	/64 - 1 - 12 - 13	×
	owner.username	Engagem			it kate	(Standardized)	1
0	NaN		Na			NaN	
1	NaN	20	Na			NaN	
2	essiesophie		1.46078			-0.045592	
3	essiesophie		6.74495			0.041928	
4	essiesophie		7.09253			-0.052507	
5	essiesophie		9.11225			-0.049310	
6	essiesophie		0.79849			-0.062471	
7	essiesophie		4.20854			0.196221	
8	essiesophie		1.03334			-0.062100	
9	ellaxarnold		3.57940	1		-0.058069	
10	ellaxarnold		2.56331	3		-0.059677	
11	madison_silotti		5.48916	3		-0.055046	
12	madison_silotti	1	0.73526	4		-0.046741	
13	madison_silotti	6	7.95624	9		0.043845	
14	madison_silotti	6	6.21430	0		0.041088	
						40000 000 500 1044	
2000	Average Engageme		Averag	e Engagement	Rate	(Standardized)	
0		NaN				NaN	
1		NaN				NaN	
2		.207274				0.157577	
3		.207274				0.157577	
4		.207274				0.157577	
5	37	.207274				0.157577	
6	37	.207274				0.157577	
7	37	.207274				0.157577	
8	37	.207274				0.157577	
9	54	.962384				0.344262	
10	54	.962384				0.344262	
11	84	.151615				0.651171	
12		.151615				0.651171	
13		.151615				0.651171	
14		.151615				0.651171	

- Computed engagement rate by (Likes + 2\*comments)/followers
- Any missing engagement values were filled by calculating the mean

```
# Display 10-15 samples from both datasets
sample_influencers = influencers_df.head(7) # First 7 influencer bios
sample_brands = brands_df.head(8) # First 8 brand bios
# Print processed samples
print("\nInfluencer Bios (Processed):\n", sample_influencers[['Bio', 'Processed_Bio']])
print("\nBrand Bios (Processed):\n", sample_brands[['Bio', 'Processed_Bio']])
[nltk_data] Error loading punkt: <urlopen error [Errno −3] Temporary
               failure in name resolution>
[nltk data]
<ipython-input-18-4f6f14a264ba>:10: DtypeWarning: Columns (6) have mixed types. Specify dtype option on import or set low_me
 brands_df = pd.read_csv('/kaggle/input/influencer-dataset/brands_data.csv', usecols=['Bio'])
Influencer Bios (Processed):
                                                 Bio \
                  Fashion | Lifestyle | Romeo's mom
  Essential Oils + Boy Mama + doTERRA * +Let's C...
  Majority of pics are of my dog!
  🎄 VEJA OS STORIES 🏝 FOTÓGRAFA, DIGITAL INFLUE...
  *Auckland, New Zealand WHonest reviews / Routi...
  I can't dance but I can burpee 🟴 Digital Ed @W...
  Editing a novel previously oldfashionedsus Bo...
                                      Processed_Bio
                       fashion lifestyle romeos mom
  essential oils boy mama dterra lets collab jod...
  majority of pics are of my dog charleston nati...
  veja os stories fotgrafa digital influencer ja...
  auckland new zealand honest reviews routines f...
  i cant dance but i can burpee digital ed @ wom...
  editing a novel previously oldfashionedsus boo...
```

Cleaning comments are other text data

```
import pandas as pd
from nltk.sentiment import SentimentIntensityAnalyzer
import nltk
from langdetect import detect, DetectorFactory
# Ensure consistent language detection
DetectorFactory.seed = 0
# Download VADER lexicon
nltk.download('vader_lexicon')
# Load data
df = pd.read_csv("/kaggle/input/comments/merged_df-2.csv")
df['extracted_comments'] = df['extracted_comments'].fillna("")
# Initialize VADER
sia = SentimentIntensityAnalyzer()
# Function: Detect language and apply sentiment
def get_vader_sentiment(text):
        if detect(text) == 'en':
            return sia.polarity_scores(text)['compound']
        else:
            return 0 # Non-English
        return 0 # Empty or undetectable text
# Apply sentiment analysis
df['sentiment_score'] = df['extracted_comments'].apply(get_vader_sentiment)
# Average sentiment per influencer
influencer_scores = df.groupby('User').agg({
     'sentiment_score': 'mean',
     'extracted_comments': lambda x: list(x)[:3] # first 3 comments for display
}).reset_index()
# Sort by sentiment (ascending)
influencer_scores = influencer_scores.sort_values(by='sentiment_score', ascending=True)
```

- VADER LEXICON for sentiment analysis.
- Widely used for sentiment analysis on social media text

```
import pandas as pd
import numpy as np
from catboost import CatBoostClassifier
from sklearn.preprocessing import StandardScaler
from sklearn.metrics.pairwise import cosine_similarity
from sentence_transformers import SentenceTransformer
```

- Standardisation done using StandardScaler
- SentenceTransformer and Cosine similarity were used for bio similarty

# "SECTION 4"

#### Semantic Similarity using Sentence Transformers and Cosine Similarity

#### Why?

- Bio Similarity Score: to quantify alignment
- Captures Meaning, Not Just Keywords: Unlike TextBlob (which is rule-based and keyword-driven)

#### How?

- Sentence Embeddings: Convert brand and influencer bios into dense vectors.
- Cosine Similarity: Measure semantic closeness between each brand-influencer bio pair.

#### CatBoostClassifier

#### Why?

• Trains a CatBoost classifier to learn patterns that distinguish sponsored vs. non-sponsored influencer-brand pairs.

#### How?

CatBoost uses gradient boosting — that builds a sequence of decision trees, where each new tree tries to fix the errors of the previous trees.

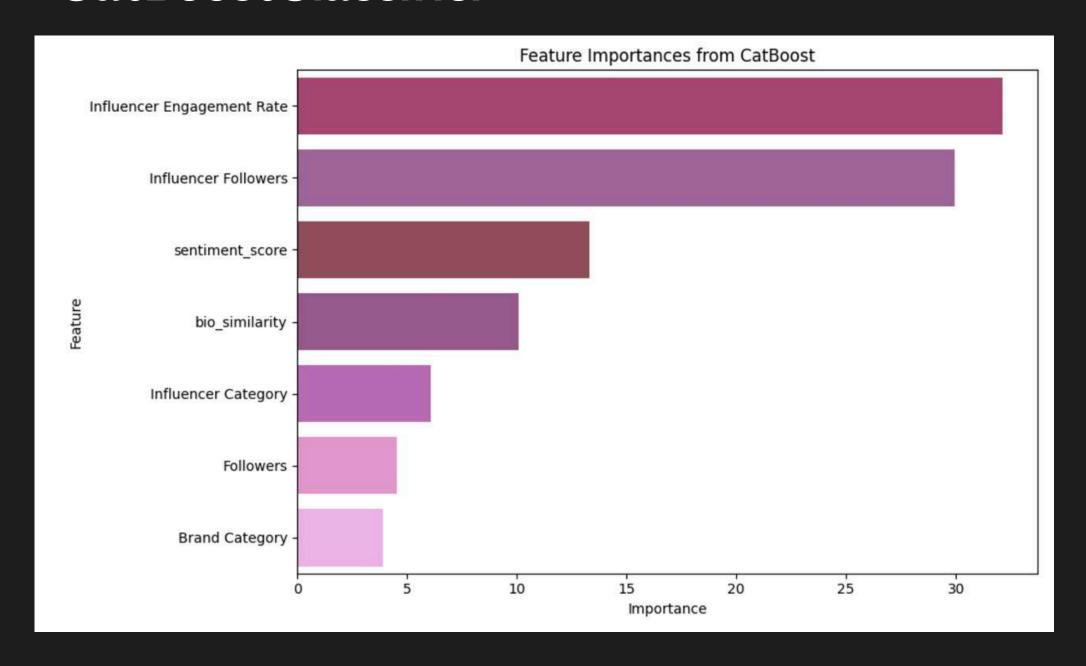
- Start with a basic prediction (like the average label).
- Compute residuals (the mistakes).
- Train new trees to correct those mistakes.
- Repeat for many rounds each tree gradually improves the model.

#### CatBoostClassifier

Feature	Benefit
Native categorical handling	No need for encoding; better performance
High accuracy on tabular data	Usually outperforms random forests
Low overfitting	Thanks to ordered boosting
Easy to use	Minimal preprocessing

• Traditional Boosting uses the full dataset to train at each stage — leading to overfitting whereas CatBoost's Ordered Boosting trains each sample using only the samples before it, preventing data leakage and improving generalization.

#### CatBoostClassifier



## Challenges Faced & How We Tackled Them

#### **Dataset Challenges**

- Brand usernames missing & Sponsorships not directly mentioned → Extracted from custom preprocessing.
- Massive dataset size (~7 crore records) → Couldn't load in Kaggle due to RAM limits.

#### Solution:

 Used smart pairing: for each brand, included all sponsored influencers + 10 random unsponsored ones → reduced data size without loss of logic.

#### **Algorithmic & Modeling Decisions**

- Initially used TF-IDF for bio similarity → Sentence Transformers for deeper semantic understanding after testing.
- Initially considered XGBoost → CatBoost after deeper domain research

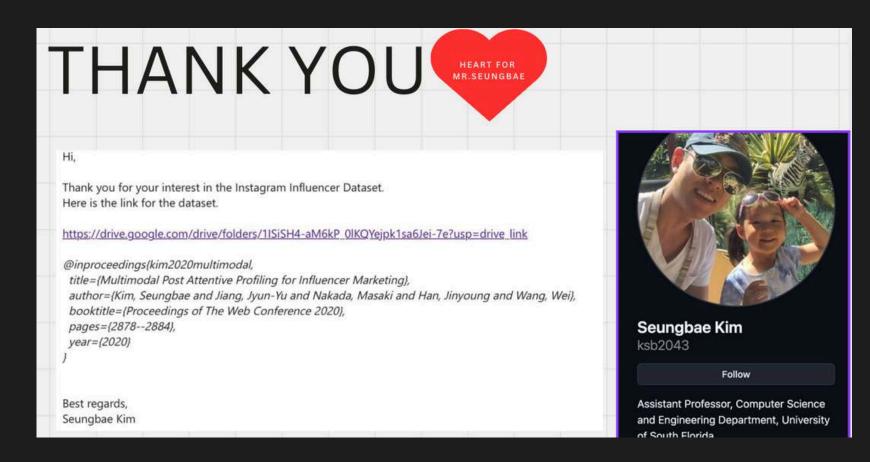
# Challenges Faced & How We Tackled Them

#### Data Availability & Quality

- Dataset acquired after extensive research and cold emailing.
- Limited sponsorship history → Augmented slightly but avoided overdoing to maintain consistency.
- Some influencer usernames and metrics changed over time (e.g., follower count) → Accepted as part of dataset limitation.

#### Solution:

- Treated historical sponsorships as ground truth.
- Acknowledge that with an up-to-date dataset, these issues can be fully resolved.



### The Fun Part

```
# === 2. Input new brand data ===
brand_bio = "we have the best furniture and interior design in the industry"
brand_category = "furniture"
brand_followers = 5000
```

Batches	: 100%		1/1 [00:00<00:00, 15
	Influencer Username	Influencer Followers	sponsorship_prob
10308	hommeboys	-0.079969	0.999412
8256	anthonygeorgehome	-0.099726	0.999356
35281	astoldbymichelle	-0.085095	0.999113
10854	centered_by_design	-0.088690	0.998724
4910	coco.and.jack	-0.075875	0.998537



861 posts 777K followers 905 following

(astoldbymichelle

Lifestyle & Home Decor Blogger

:• Weekly Decorating, DIY, & Hosting Ideas

TX, Orlando, FL :+ Mom 3 5 5

@ www.shopltk.com/explore/astoldbymichelle + 3

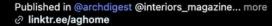


anthonygeorgehome

anthonygeorgehome

(anthonygeorgehome

interior designer | 🛋 ceramicist | 🏺 animal lover and 🗂 | 🗫 🦘 🦮























**⊞ POSTS** TREELS 2 TAGGED



REELS

2 TAGGED







**⊞ POSTS** 

REELS

2 TAGGED



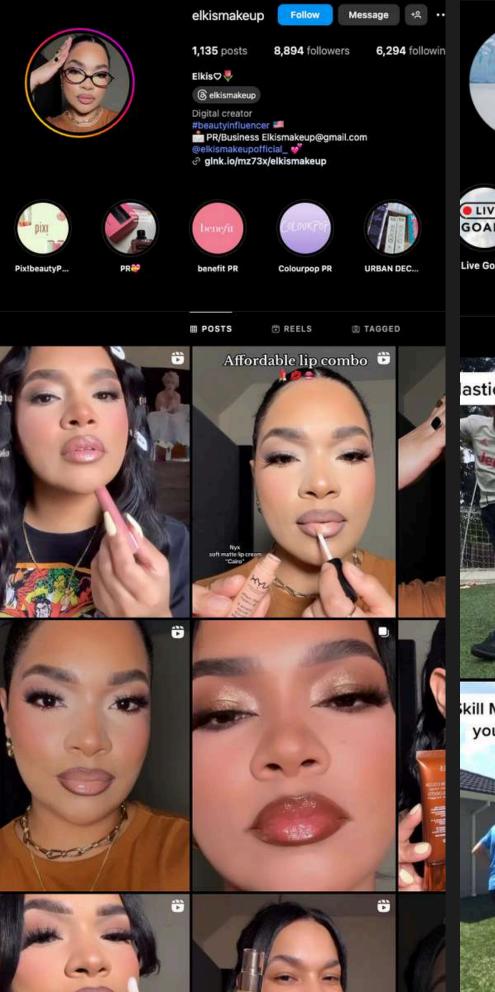


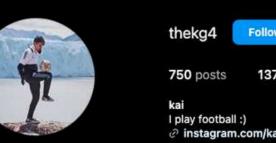
## The Fun Part

```
# === 2. Input new brand data ===
brand_bio = "beauty brand focuses on makeup and specializes in lipsticks"
brand_category = "makeup"
brand_followers = 5000
```

. 63		
.nfluencer Username	Influencer Follow	ers sponsorship_prob \
wong.jpg	-0.077	716 0.999188
elkismakeup	-0.100	610 0.998936
_make_up_looks	-0.099	955 0.998865
tulipheels95	-0.096	729 0.998858
mandiglitter	-0.095	591 0.998313
	wong.jpg elkismakeup _make_up_looks tulipheels95	elkismakeup -0.100 _make_up_looks -0.099 tulipheels95 -0.096

```
# === 2. Input new brand data ===
brand_bio = "passioante brand about sports football specefically"
brand_category = "football"
brand_followers = 5000
22914
                    thekg4
                                        0.033219
                                                          0.996388
14447
                                       -0.054909
                                                          0.991702
               _han_banan_
                                                          0.987175
28511
            angeliquefiske
                                       -0.100603
       thefootballrepublic
5639
                                       -0.012542
                                                          0.986366
```





I play football :)
⊘ instagram.com/kaigalia



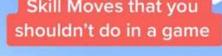




**⊞ POSTS** D TAGGED









# "SECTION 5"

## Performance Metrics of our model

#### Mean Reciprocal Rank (MRR): 0.473

• On average, the first true sponsored influencer appears around the 2nd position in our ranked list.

#### Recall@3: 0.8902

• In nearly 9 out of 10 cases, at least one correct influencer is in the top three recommendations.

```
return sum(mrr_list) / len(mrr_list)
mrr_score = calculate_mrr(final_df)
print("MRR Score:", mrr_score)
|
MRR Score: 0.4729697599196747
```

```
k = 3
recall_scores = final_df.groupby('Brand Username').apply(lambda x: recall_at_k(x, k))
overall_recall_at_k = recall_scores.mean()
print(f"Overall Recall@{k}: {overall_recall_at_k:.4f}")

Overall Recall@3: 0.8902
```

#### **Classification Metrics**

- Precision ≈ 0.83 Recall ≈ 0.87 (sponsored)
- Balanced Accuracy  $\approx$  0.9081, F1  $\approx$  0.93

#### **Evaluation Protocol**

 Computed on a stratified 80/20 hold-out split to preserve sponsorship ratios and ensure robust performance estimates.

Why it matters: High MRR and Recall@3 minimize the list brands need to review & Balanced Precision/Recall shows we avoid both irrelevant picks and missed opportunities.

#### Deploying Our Influencer Recommendation Engine: Transforming Influencer Marketing

#### Vision:

• Tap into the booming \$40B influencer marketing industry with a first-of-its-kind ML-powered recommendation system that helps brands find the most effective, cost-efficient influencers.

#### Why Now?

- No commercialized influencer recommendation model exists yet.
- Backed by our passion, rigorous development, and support from Prof.
   Brainerd — we are ready to scale this innovation to market.

Business Model	Description
B2B SaaS Platform	Subscription-based web interface for startups and agencies to find influencers.
	Operate as a dedicated agency recommending influencers and earning:
Full-Service Marketing Agency	Commission from brand-influencer deals
	Revenue share from successful collaborations

## Scaling Considerations for our model

**Global Expansion** 

Add multilingual support and integrate additional platforms (e.g., TikTok, YouTube) to serve new markets.

**Advanced Capabilities** 

Incorporate visual-content matching and fraud/bot detection as load and scope grow.

**Go-to-Market Scaling** 

Launch a freemium tier to drive early adoption and forge agency partnerships for rapid traction.

**Data Freshness** 

Recompute embeddings and sentiment scores daily or weekly; monitor input distributions and retrain when performance degrades.

## "The End"

```
# === 2. Input new brand data ===
 brand_bio = "Machine Learning & Pattern Recognition"
 brand_category = "AI"
 brand_followers = 5000
                                                1/1 [00:00<00:00, 3
Batches: 100%
      Influencer Username Influencer Followers
                                                sponsorship_prob
35902
        Mr. Siddharth
                                    1.00
                                                        1.00
27841
12539
        Mr. Pushpinder Singh
                                    1.00
                                                        1.00
30536
```



### References:

- https://doi.org/10.48550/arXiv.2106.01750
- https://kth.diva-portal.org/smash/get/diva2%3A1783645/FULLTEXT01.pdf
- https://www.diva-portal.org/smash/record.jsf?pid=diva2%3A1636469
- <a href="https://www.statista.com/topics/2496/influence-marketing/#topicOverview">https://www.statista.com/topics/2496/influence-marketing/#topicOverview</a>
- https://doi.org/10.48550/arXiv.2107.13943
- <a href="https://www.tcs.com/what-we-do/industries/consumer-goods-distribution/white-paper/reimagining-influencer-marketing-with-machine-learning-nlp">https://www.tcs.com/what-we-do/industries/consumer-goods-distribution/white-paper/reimagining-influencer-marketing-with-machine-learning-nlp</a>
- https://prohibitionpr.co.uk/digital-marketing/influencer-marketing/the-impact-of-working-with-the-wrong-influencers-in-marketing-campaigns/
- <a href="https://dl.acm.org/doi/abs/10.1145/3366423.3380052">https://dl.acm.org/doi/abs/10.1145/3366423.3380052</a>