

Applied Machine Learning

PROJECT PRESENTATION

TEAM : Longshot

MAY 2025

Content

A Team

B Problem Statement

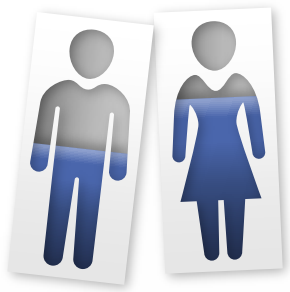
C Data and Design

D Experiments

E Results

F Diagnostics

G Conclusions and Recommendations



Team and Roles



Gauranga

gauranga.mds2023@cmi.ac.in



Nikita

nikitakumari@study.iitm.ac.in



Dhruv

22f3001413@ds.study.iitm.ac.in



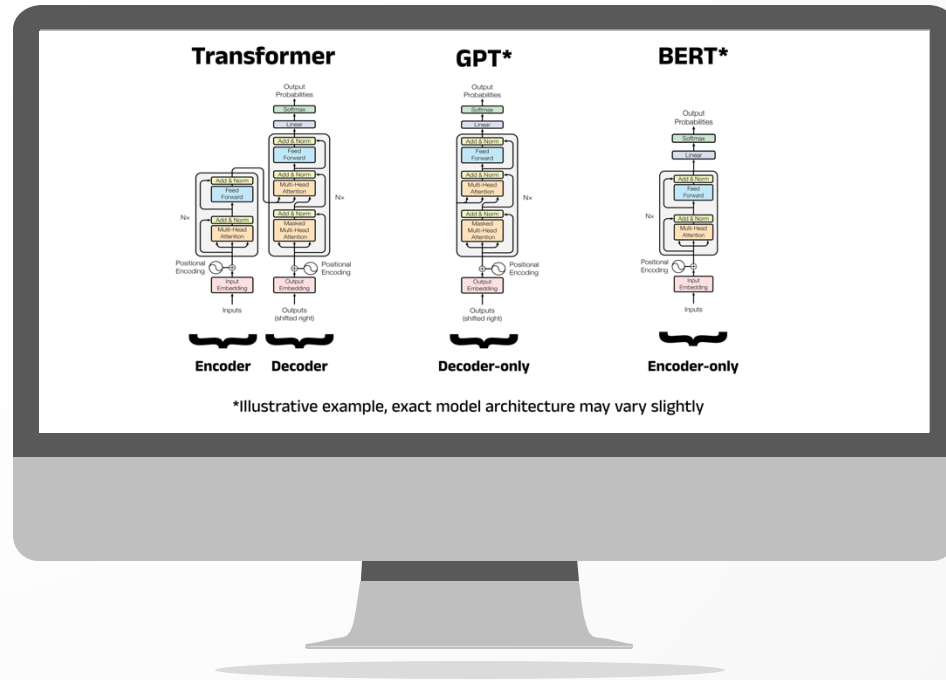
Vasu

21f3002975@ds.study.iitm.ac.in



Member	Responsibilities
Nikita	Background research, Data pre-processing and Design
Dhruv	Code base and ANN Modelling
Gauranga	ML Modelling and Experiments
Vasu	Interpretation and Reporting

Backdrop



Situation

There is a strong belief that Deep Learning based models can solve any problem way better than traditional ML models.

Question

Are the modern AI modeling architectures the silver bullet?

Are there any considerations to keep in mind?

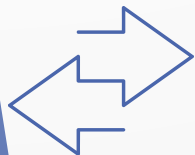
Answer

Systematic comparison of Classical Models and Modern Deep Learning based Architectures.

Problem of choice

- I. Good Reads is a very popular source of book reviews for users. Book reviews impact reader perceptions and can directly impact sales.
- II. Currently likes and comments on reviews are being used to determine the popularity of a review which in turn determines the placement of the review. Higher the popularity, higher is placement.
- III. However, this is our best guess and a black box. To verify this, we have reconstructed the popularity metric and tested whether a model can be trained to determine which reviews will be popular vs. not so popular.
- IV. While doing so, we are evaluating whether simpler, faster classical methods can match modern deep learning models at this classification task.

Study design



	Choices	Metrics	Experiments
Data	<ul style="list-style-type: none">15 Mn reviews data set (2 Mn books, 465 k users)Construct book review popularity measure based on likes and commentsMeta data of reviews (length of the review, % verbs etc.)English reviews	<ul style="list-style-type: none">Class imbalance : ~15% reviews are popular	<ul style="list-style-type: none">Down samplingTfidf, CBoW, Skip-Gram and BERT embeddings.Normalization of meta data
Models	<ul style="list-style-type: none">Model = f (Review + Review's Meta data) -> Classification	<ul style="list-style-type: none">Precision RecallFN is costlier => Recall	<ul style="list-style-type: none">LogisticXGBoostANNTransformer
Action	<ul style="list-style-type: none">If prob(popular) > default threshold push the review up the order on the website	<ul style="list-style-type: none">Update threshold to improve recall	<ul style="list-style-type: none">Pick highest recall subject to acceptable precision score

Data and feature engineering



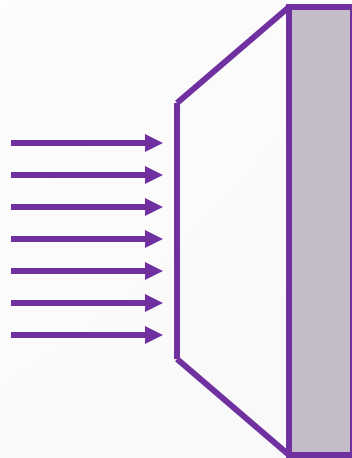
15 Mn
multilingual reviews

2 Mn books

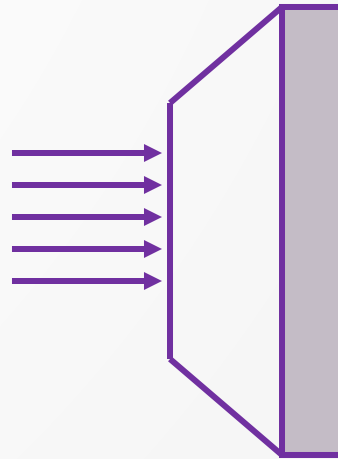
465k users

- English reviews
- Review reactions = Likes + Comments
- At least 100 reactions

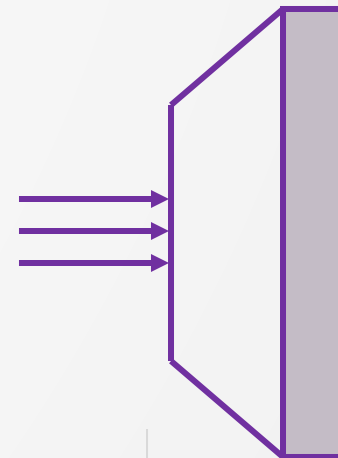
Cleaning



Ft. Engg.



Tokenization



- No of reviews by reviewer
- # of sentences, length # of words, # and % verbs, %, # and % of nouns # and % of adjectives nouns,
- VADER sentiment score
- User rating, Diff from average

- NLTK
- Stop words removal
- Tokenization

POPULAR REVIEW IS
DEFINED AS:

Share (Comment + )

> 2 %

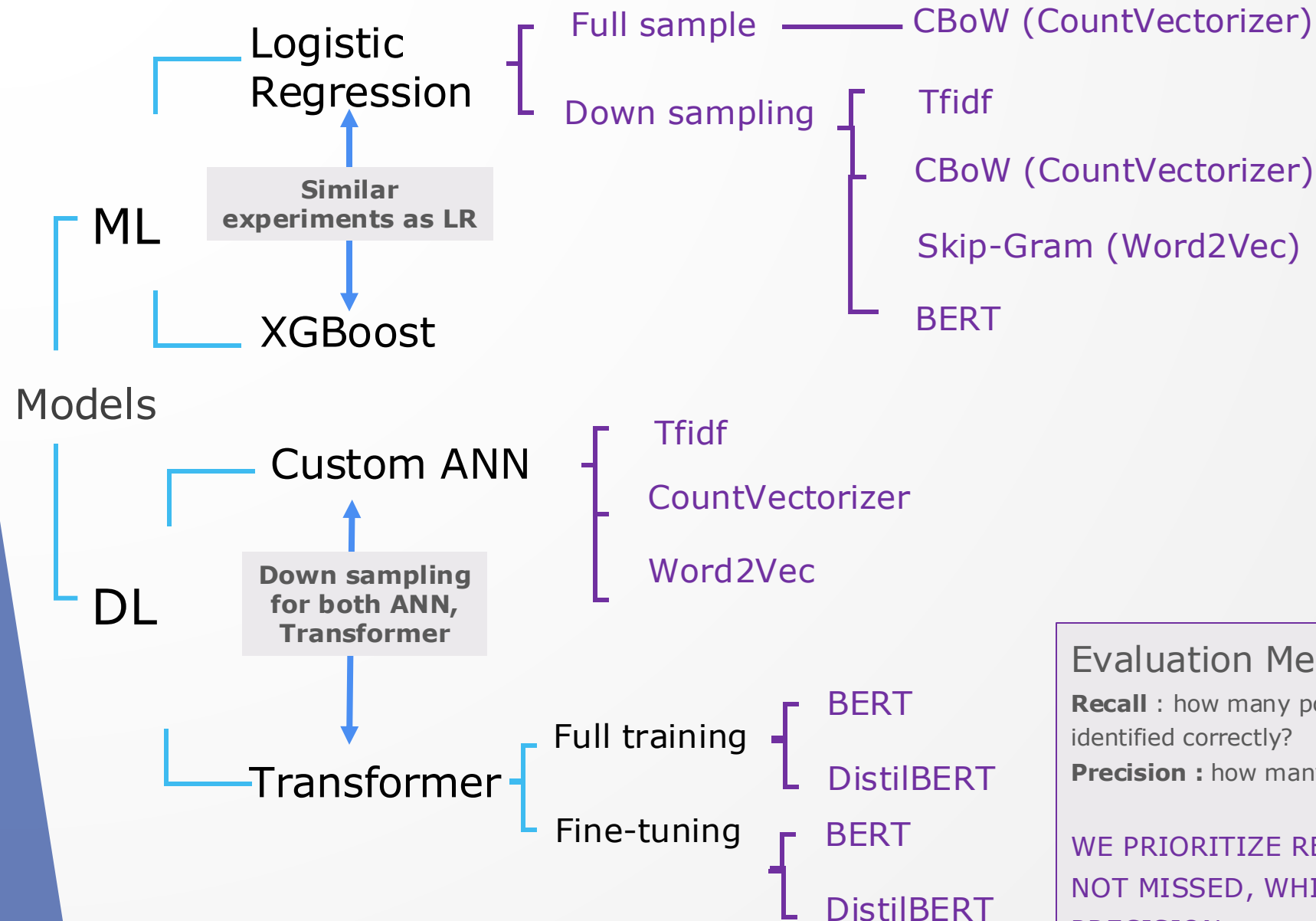
83%

17%

Not Popular Popular

1.7 Mn Reviews

Key Experiments#



Hyper-Parameters

LR

Regularization strength
(C = 10, 1, 0.01, 0.001)

XGBoost

Estimators : 100,1000
Depth: 4,6
Lr : 0.1, 0.3

BERT/ DistilBERT

No of tokens : 64, 128
Epochs: 1,3

Evaluation Metrics

Recall : how many popular reviews have been identified correctly?

Precision : how many identified as popular are correct?

WE PRIORITIZE RECALL TO ENSURE POPULAR REVIEWS ARE NOT MISSED, WHILE MAINTAINING AN ACCEPTABLE LEVEL OF PRECISION



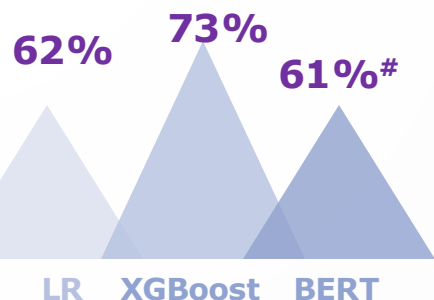
Not all experiments are mentioned here

Results

Model	Down Sampling	Embeddings	Accuracy	Precision	Recall/ Sensitivity	f1	Specificity	AUC (ROC)		REMARKS
LR	Yes	Tfidf	0.74	0.32	0.66	0.43	0.75	0.70		
LR	Yes #	CBoW	0.77	0.34	0.62	0.44	0.79	0.71		Acceptable
LR	No	CBoW	0.86	0.60	0.18	0.28	0.98	0.58		Very poor recall
LR	Yes	Word2Vec	0.74	0.30	0.64	0.41	0.74	0.69		
LR	Yes	BERT	0.71	0.29	0.66	0.41	0.72	0.69		
XGBoost	Yes	Tfidf	0.72	0.31	0.73	0.44	0.72	0.72		
XGBoost	Yes	CBoW	0.72	0.31	0.73	0.44	0.72	0.72		XGBoost does better !
XGBoost	No	CBoW	0.68	0.68	0.22	0.34	0.98	0.60		
XGBoost	Yes	Word2Vec	0.71	0.30	0.72	0.42	0.71	0.72		
ANN	Yes	Tfidf	0.72	0.31	0.70	0.43	0.73	0.71		
ANN	Yes	CBoW	0.71	0.30	0.74	0.43	0.70	0.72		No upside vs. ML models
BERT	Yes	BERT	0.74	0.30	0.58	0.39	0.77	0.73		Not the best as expected
BERT- Ft*	Yes	BERT	0.64	0.22	0.57	0.32	0.65	0.65		

* Ft – Fine tune, # 15% of the full sample are positive reviews, train sample down sized to reflect 50-50% split

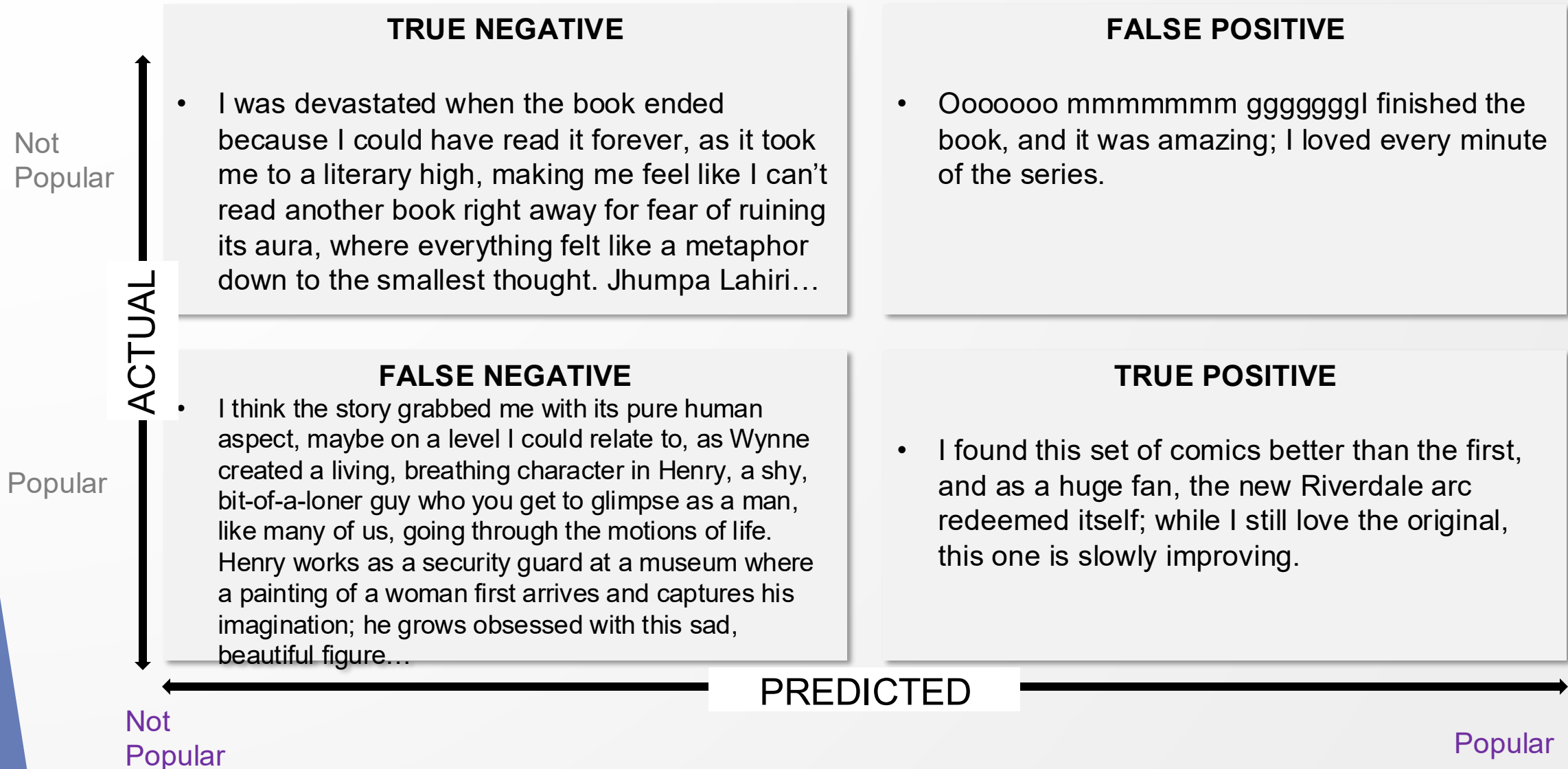
Why is RECALL not improving with DL models?



- I. Are we measuring the right thing?
- II. Is embeddings an issue?
- III. Is dataset imbalance making learning hard?

Popularity	→ Embeddings	→ Imbalanced
<ul style="list-style-type: none">• Popularity (dependent/outcome) metric was created via feature engineering and not independently or directly measured from users• Correlation between popularity and features is low (0.3) to begin with	<p>Diagnostics on TP and FN reviews revealed the following:</p> <ol style="list-style-type: none">1. When emotional subtlety is missing, the model is failing to learn low key appreciation and classifying it as unpopular2. Deeply thoughtful and philosophical reviews with no overt positivity are being deemed unpopular3. Positivity and sentiment are being confused. Could be a popular review but highly critical and negative in sentiment	<ul style="list-style-type: none">• Yes, imbalance is an issue evidenced by poor model performance on full sample• Down sampling is showing better modeling outcomes (70% vs. ~ 20% recall scores)

A closer look at TP and FN reviews



Conclusions and recommendations



- Define popularity clearly, what is it supposed to capture?
- Re-evaluate how to measure popularity metric
- Ideally an independent measure of popularity is required and not feature engineered
 - E.g. : review and recommendation
- And if required, capture from users (maintain a stratified sample)



- Domain Adaptive Pre-Training : Train BERT further using unsupervised masked language modelling (MLM) on reviews corpus
- Test for Sentiment \neq Popularity
- Stop-words and tokenization should not remove evocative differentiators



- Deal with data set imbalance; down sampling has shown promise in mitigating this risk
- Deep Learning models do not provide an automatic advantage, this is a cognitive bias
- Do not overlook traditional ML models (particularly Boosting algorithms), more interpretable and computationally friendly
- Follow Occam's Razor : adding meta data and additional features may not automatically improve model performance

For suggestion and questions



21f3002975@ds.study.iitm.ac.in,
gauranga.mds2023@cmi.ac.in

TEAM : Longshot

MAY 2025