

## Problem Statement

Design a Machine Learning for QnA Ranking

→ FB Community  
→ Quora  
→ Reddit

## Objectives :

1. Given a Q and answer  $\{A_1, A_2, \dots, A_n\}$   
predict a score, such that when the answers  
are sorted by score, the best takes the first place.

In Practice → Maximize the relevance metric. (NDCG)  
click/engagement rate.

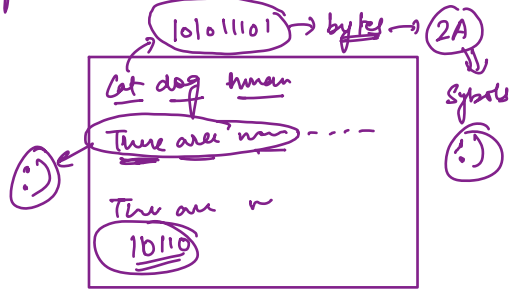
2. Scale → Mn of QnA  
Latency < 100 ms

Data Collection & feature engineering Layer.

→ Historical data on QnA & platform logs



- + upvotes, likes, timestamp, edit history, accepted
- + user's reputation (warning / bans / removed content)
- + clicks, impressions
- + sub-replies & main replies.
- + grading (0-4)
- + dwell time



## Feature Engineering

### + Text processing (Tokenisation)

- + lowercase
- + stripped

### + Contextual embedding → Transformer Sentence-BERT

### + Metadata

## Modelling Strategy

# Baseline model

+ simple supervised model

+ votes  
+ impression  
+ dwell

} logistic Reg  
XGBoost.

<https://medium.com/acing-ai/quora-ml-platform-for-ranking-answers-baaf00cc97e8>

+ SBERT

[https://sbert.net/examples/sentence\\_transformer/applications/retrieve\\_rerank/README.html](https://sbert.net/examples/sentence_transformer/applications/retrieve_rerank/README.html)

LLMs → Prompt ( — )

Score

+ explain → { }

[https://arxiv.org/html/](https://arxiv.org/html/2411.00142v1#:~:text=Most%20recently%2C%20utilizing%20Large%20Language,approaches%20utilize%20an%20LLM%20either)

2411.00142v1#:~:text=Most%20recently%2C%20utilizing%20Large%20Language,approaches%20utilize%20an%20LLM%20either

# Evaluation

+ Baseline  
F score, AUC ...  
NDCG, MRR

Given a query, and a list of answers, the objective of the ranking model is to rank the answers with optimal rank related metrics, such as NDCG

Target / Amazon / Flipkart

Promotions & Campaigns

20 years  $\rightarrow$   $P_1, P_2, \dots, P_{1000}$  Items  $[I_1, I_2, I_3, \dots]$

Promotion  $\rightarrow$  Predict Units

Strategy design

$\rightarrow$   
 $\rightarrow$   
 $\rightarrow$

Task  $\rightarrow$  5 mins

Break  $\rightarrow$  5 mins

Put answers on chat

Promotion → June & July 2025



June & July 2024 | 2023 | 2022 ...

Step 1

Most similar promotion from the past year.

Hadoop

↳ HQL (↪ SQL)

Historical data.

Promotion	Items	Unit Sales	Timestamp. (Fiscal weekend)
P <sub>A</sub> ↖	I <sub>1</sub> ↗	100	2020-08-01
P <sub>A</sub> ✓	I <sub>2</sub> ✓	200	2020-08-07
P <sub>A</sub>	I <sub>3</sub>		2020-08-14
P <sub>B</sub>	I <sub>x</sub>		
P <sub>n</sub>	I <sub>y</sub>		

P<sub>C</sub> ← I<sub>1</sub>

P<sub>XA</sub> ⇒ I<sub>XA</sub> [ I<sub>1</sub>, I<sub>2</sub>, I<sub>3</sub> ]

P<sub>XA</sub> → P<sub>A</sub> ⇒ Metric

nesting

Similarity

① Item Similarity

→  $\frac{A \cap B}{B}$  ⇒  $\frac{\text{Common Items}}{\text{Items in the current promo}}$

$P_A [I_1, I_2, I_3, I_4]$   
 $P_{xA} [I_1, I_2, I_5]$

$\Rightarrow \boxed{\frac{2}{3}} \Rightarrow \underline{\underline{\text{Jaarud}}}$

② Brand Jaarud

③ Promo length  $\rightarrow$  20 days (Past promotion)

10 days

④ Location Jaarud

$\Rightarrow \frac{200}{300} \Rightarrow \underline{\underline{0.66}}$

$\Rightarrow \textcircled{2}$

Current Promo	Matched Promo	Item Jaarud	Brand	Promo length	Location	Weighted Score
$P_A$	$\rightarrow P_A$	0.66	0.78	2	0.66	1
$P_{xA}$	$\rightarrow P_B$	-	-	-	-	-
$P_A$	$\rightarrow P_C$	-	-	-	-	-
$P_{xA}$	$\rightarrow P_D$	-	-	-	-	-
...						

$P_{xA} \rightarrow P_A$

$(a \times 0.66 + b \times 0.78 + c \times 2 + d \times 0.66)$

Max(Weighted)

[End of story  $\rightarrow$  Most similar promotion]

Predict the sales of promotional units for the current promotion

1) Target variable  $\rightarrow$  Promotional units from past promotions.

2) Macro metric → Count of Promotional Campaigns One item is used in.



3) Trend based feature → Item wise sales in last quarter, last 6 months.

Step 1

Already given

Current Promo	Matched promo	Items	$F_x$	Last Qtr Sales	Last 6 months Sales	Tricomp	Seasonal	Item
$\frac{P_{xA}}{P_{xA}}$	$\frac{P_A}{P_A}$	$I_1$ $I_2$ ...	30 22				Multiplier	Jacard

Brand	locapm	Promotion	Units
Jacard	Jacard	length	

↓

Treatment

Null / outliers / Negatives

↓

Train Test Split

Time based

Train  $\rightarrow$  2023-01-01  $\rightarrow$  2024-01-01  
Test  $\rightarrow$  2024-01-02  $\rightarrow$  —



Model

XGBoost  $\rightarrow$

Random Forest

$\downarrow$

Evaluation | Hyperparameter tuning



Post-Model Analysis  $\rightarrow$  SHAPley Value