# VIDEO RECOMMENDATION

ONLINE METRICS:
1. CTR
2. Reaction e.g. like, upvote/downvote, rating
3. Num. of completed videos
4. Watch time as a % of total session time

LABEL CRITERIA:   +ve:
1. clicked
2. >80% watched
3. +ve reaction e.g. like, upvote, high rating

−ve:
1. Not clicked (but have impression)
2. <10% watched
3. −ve reaction e.g. report, dislike, downvote, low rating

FEATURES:
→ all item interaction
→ historic interaction
→ user's context when interaction happened

| User Features | User-Item Interaction features | Item Features | Context Features | Sparse Features |
|---|---|---|---|---|
| Demographic info: | − avg. session time | − Lang | − Time of day | − User Id |
| − age | − Seq. features | − Titles & Tags | − Day of week | − Item Id |
| − gender | − search history | − Duration | − is_holiday | |
| − location | − Liked videos | − genre | − device | |
| − Language | − Watched videos | − time elapsed since release | | |
| Account info: | | − time on platform | | |
| − account age | | − engagement (e.g. likes, upvote, watched) (rate & counter feature) | | |
| | | − maturity rating | | |
| | | − Modality of data features (e.g. video feature-engineering) | | |

− Rate & Counter features
− Context Features:  Sine & Cosine fn for feature-engg.
   (to show 23 is closer to 1
− Sparse Features:  Deep hash Embedding

# MODELLING :

- →Candidate Generation : — Multiple Ways
- →Filtering — Geography restriction
  - — Age restriction
  - — Already watched
- → Ranking — Re-ranking — Multi-task learning
- →Re-ranking : — Diversity
  - — Freshness

## CANDIDATE GENERATION :

— 2 Tower: User features + User-item historical interaction
(User Tower) + Context + Sparse features
related to User
eg. User id

(Item Tower) : Item features

— Sampling for 2 tower & Bias removal — Popularity
— Fairness
— Privacy
— Position

— Matrix Factorization Vs 2-tower

— ANN: Quantization & Indexing

— Evaluation Metric: Recall based

## RANKING

— Pointwise, Pairwise & Listwise : Pros & Cons
— Single Task vs Multi-task learning: Pros & Cons
— Multi-task learning: Shared bottom Vs MMoE Machines/Depth
— Single Task or Expert NW in MMoE: Deep and Wide NW or Factorization — Popularity
— Sampling : Balancing Bias removal — Fairness
Train-Test split — time, User basis — Privacy
basis — Position
— Loss fn : Focal loss

— Features: Cand. generation score + Avg. Embedding
+ Any specific user-item
interaction feature

— Evaluation Metric: Precision based — MAP
— NDCG
Diversity

SPECIAL ISSUES IN RECSYS: — Cold start of users & items
— Serving bias, position bias, popularity bias
— Calibration

# FEED RANKING

→ No. of times feed need to be refreshed each day for each user

→ Any specific engagement to be optimized or weighted engagement (e.g. like, comment, share, etc.)

MODEL ARCHITECTURE

→ Multi-stage ∴ 1) Candidate Selection (not 2-tower, heuristic based)
2) Ranking (multi-task classification)

→ ML Objective : Maximize wt. Score based on both explicit & implicit reaction

→ ONLINE METRICS :
1) Impression duration (e.g. time spent viewing the post)
2) Reaction rate
3) Hiding/reporting rate
4) Sessions with reaction / Total # of sessions
5) Retention: re-visit after x duration
6) Avg. session duration
7) No. of sessions in x days

→ LABEL :  Posts with explicit reactions = +ve
(like, comment, share, etc.)

Posts with implicit reactions = +ve
(impression duration > 2 sec for passive users)

|  | Like | Comm | Share | Subsc |
|---|---|---|---|---|
| Post 1 → | 1 | 1 | 0 | 0 |
| Post 2 → | 1 | 1 | 1 | 0 |
| Post 3 | 0 | 0 | 0 | 1 |

Posts with hide/report = -ve

→ CANDIDATE SELECTION : 
+ (FILTERING)
1) Posts from 1st degree connections (unseen or impression duration < t∞)
2) Post from pages/people you follow
3) Post from groups joined
4) Trending & popular posts in your geography, age group
5) Posts with high engagement in your network
6) Posts related to your interest

→ FEATURE-ENGINEERING (FOR RANKING MODEL): Can use cross-features

Item: Textual Content, Image Or Videos, Reactions, Hashtags, Post age
Post popularity, Post length

User: Demographics: age, gender, lang, location, account age
Context: time/day of week, device, is_holiday
User-item historical engagement: Rate & Counter features
"all" e.g. reaction to text/content
video
image
avg. no. of hashtags used in last 3 months

User-"this" item engagement areas: top topics interested in
topic similarity, posted by friend/family
mentioned in post

User-author engagement: Like/comment/share rate
posts Length of friendship

Author: Degree of influence, historical trend of engagement on author's post

Sparse: Post id, Author Id, User id

→ RANKING MODEL: Multi-task (Shared layers + Head for each engagement
(+ RE-RANKING) Loss fn: sum of loss for each task)

→ OFFLINE METRICS: NDCG@K, MAP@K

# FRAUD/ABUSE/RISK ML SYSTEM DESIGN

**Issues:**

1. Delayed Labels / concept drift / continual learning
   - proxy label
   - Re-formulate the ML problem with time horizon
   - Domain Experts/Manual annotation
   - Unsupervised Anomaly Detection: Auto-Encoder (reconstruction loss above threshold)
     (also if labels are unavailable)
   - Supervised Anomaly Detection: one-class SVM with one-class Classification (if labels of one class known) Isolation forest
     training only on majority class labels

2. Imbalanced data:
   - Under-sample majority class / over-sample minority class
     ↓ SMOTE
   - Metrics: Precision/Recall/F1 instead of Accuracy
     PR AUC instead of ROC AUC
   - Cost sensitive fn learning    e.g. scale-pos wt in Xgboost

3. Explainability    (Human in the loop — rules → Association rule mining
   bigness
   supervised ML cannot stop new, unseen or evolved fraud patterns

4. Anonymize data for sensitive info

5. Point-in-time values: Using external APIs to fetch data, e.g. email related APIs
   Need to make sure that at the time of training, the values of these APIs should be from the point in time when event occurred (e.g. transaction happened)

# TYPES OF FEATURES → Identity ↳ Behavior

- Affinity features: based on counters: how many times two features with certain values were seen together (e.g. how many times email + IP was seen together in last 15 mins)

- Velocity features: no. of transactions from this IP in ~~last~~ 15 mins

- Reputation features: reputation of email domain no. of tickets purchased in last 6 months

- External API based features: e.g. credit score of user, lat-long of IP

- Profile based features: e.g. zip code

# MODELING:

i) Supervised : Xgboost / LightGBM

ii) Temporal / Time-series: Convert it into supervised ML by adding temporal info as features e.g. day, time, week etc + using a time-window and crafting features e.g. refresh

iii) Graph ML: GraphSage ← can be trained in batches ← can predict on unseen node

# HARMFUL CONTENT/COMMENT

→ Define 'harmful': different categories

→ Business obj : increase platform safety

→ Some categories of harm need real-time system and some categories could be batch

→ ML Objective: accurately predict harmful post/comment

→ Real-time or combination of realtime + batch

→ If real-time and # of post/comment huge ‼ (Multi-stage)
filter by 1. new accounts post/comment per unit time
2. Old accounts post/comment: historical reporting

→ Multi-task classification (new post posted on platform → classify into diff. catg. of harm)

→ ONLINE METRICS : 
1) Missed/FN: # of posts reported by users
2) FP: # of posts reported which were found to be not harmful (valid appeal)
3) Impression count of posts that were reported
4) Some ratio of above

→ LABEL : 
i) Annotated
ii) User reports (threshold by certain # of reports)
- Sample to make sure all categories of harm are appropriately represented
- wt examples by recency (constantly need to catch up)

→ **FEATURES:** Poster, Post, Comment, Commenter

(identity features: email, ip, phone no. etc.

Velocity features (rate & counter features): #times transaction happened in last x hrs etc.

**Poster**
1. Violation history
 - # of violations
 - # of user reports (by time range)
 # Profane words (by time range)
2. Account age
 # of followers, followees

**Post**
- Embedding for diff. modalities
→ Profane word used
→ sentiment
→ Context features for post
→ Lang.
→ Reactions

**Comment**
→ Embedding
→ Profane words
→ sentiment of each comm & aggregated
→ Relevance of comment to post
→ Comment creation velocity
→ Lang.
→ Reactions

**Commenter**
1. Violation history
2. Account age
 # of followers, followees

Poster - Commenter historical interaction
Poster - Comment → reaction of poster to each comment
Post - Comment → relevance of comment to post

→ **MODELING:** Multi-Stage: Pros & Cons w.r.t. alternatives

1. Single binary classifier: i) cannot determine Catg.
(one model)   ii) diff. to improve in a category if not doing well in that category

2. Multiple binary classifier i) Train & Serving multiple
(multiple models)   models is expensive/inefficient

3. Multi-label: Same feature transformation for all
(one model)   categories → not ideal

✓4. Multi-task: 1) Training & Serving: not expensive
(one model)   2) Shared layers transform features beneficial to all categories (prevents redundancy)
   3) Training data for each category contributes to learning of other tasks (especially useful if labeled data for one category is very limited)

→ **OFFLINE METRICS:** ✓Recall + False +ve Rate
  ✓AUCPR

⇒ **SERVING:** Usually in 2 forms: ⌈Low confidence → Manual Review
Model Score  ✓prediction
Use  ⌊High confidence → Remove
  ✓prediction

# PERSONALIZED SEARCH

**ONLINE METRICS:**
1. CTR (Click/impression)
2. Average Position of Click
3. Further engagement e.g. — dwell time > t sec
   — content watch duration
   — purchase/conversion
4. Time to success (click + dwell time/content played/purchased)
5. # of similar searches in a session (Counter metric)

**LABELS:** +ve: 1. Clicked
2. Clicked OR + further engagement above a threshold

−ve: 1. Not clicked (but impression)
2. Clicked + further engagement below a threshold

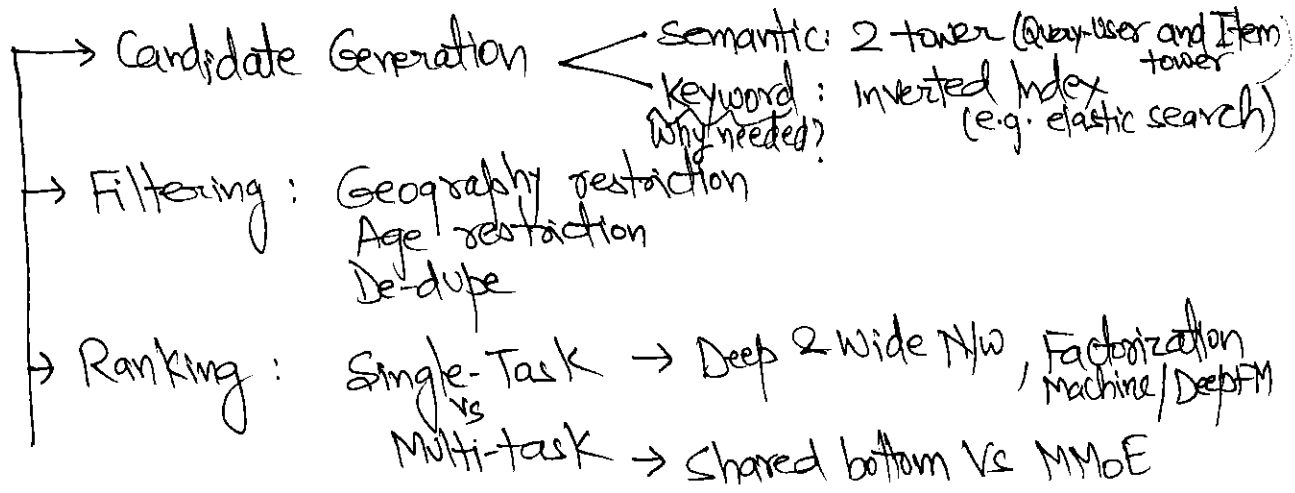— Human labellers/annotated data

**FEATURES:**

Query Processing: Typo Correction/Fuzzy Matching, Stemming/Lemmatization, Synonyms, Removing irrelevant/stop words,

> context when the query happened   Not used in 2 tower

| Query Features | Item Features | Searcher Features | Context Features | Query-Doc | Searcher-Doc |
|---|---|---|---|---|---|
| — Intent | — Keywords/Terms (emb.) BERT | — Demographic | — time of day | — Similarity measure | — Dist. b/w searcher & item |
| — Keywords/Term (BERT embedd.) | — Location | -age | — day of week | — cosine similarity b/w emb.| |
| — Location (e.g. lat, long from where Query originated) | — Reputation features: | -gender | — device | | |
| | — PageRank | — location | — is holiday | — TF-IDF match | |
| | — Author's reputation | — lang | | — Historical engagement for the (query-doc) pair | |
| (rate + counter features) | — Engagement features | | | | |
| | — clicks received | | | | |
| | — dwell time/purchased/watched | | | | |

**IMP:** Prev. historical queries of user does not matter as it is as user's query diff. things. However, some characteristics e.g. intent, radius/dist; type of content with high engagement is useful

# MODELING

→ Candidate Generation ← Semantic: 2 tower (Query-User and Item tower)
  Keyword: Inverted Index (e.g. elastic search)
  Why needed?

→ Filtering: Geography restriction
  Age restriction
  De-dupe

→ Ranking: Single-Task → Deep & Wide N/w, Factorization Machine/DeepFM
  vs
  Multi-task → Shared bottom vs MMoE

## CANDIDATE GENERATION:

- 2 tower → Sampling
  → Bias Removal (position, popularity, fairness, privacy)
  → How features are in which tower
    which
  ↓ ANN: Quantization & Inverted Index

— Eval Metric: Recall based

## RANKING:

— Pointwise, Pairwise & Listwise: Pros & Cons

— Xgboost vs DNN: Pros & Cons

— Single Task vs Multi-task (Usually single-task)

— Training Data: Balancing & Bias Removal & Train-test split

— Cross-features inclusion (if single task)

— Loss fn: Focal loss

— Evaluation Metric: — Top-K
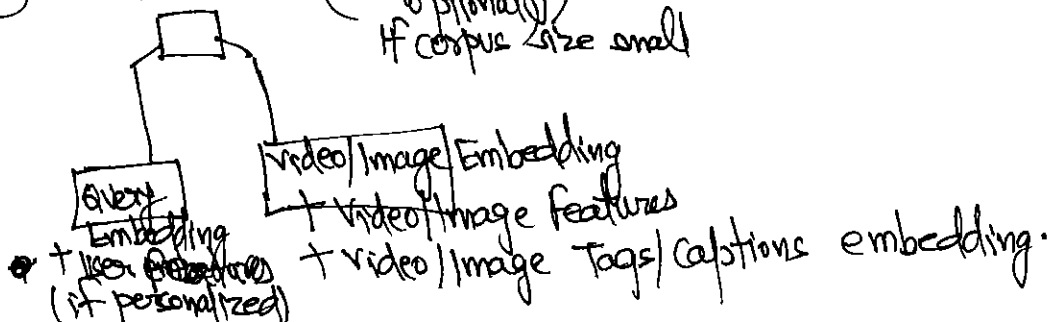  — Precision based ← MAP
    NDCG

# IMAGE/VIDEO SEARCH

→ Online Metrics : CTR, Click + dwell time > t sec, Click + dwell time > t sec
     video watched + download (of link
     # of similar searches in a session      (if not visit link + address
     Sessions with atleast 1 click       copyright → then only
     Time to success (Click + dwell time > t sec +       download)
                                any further actions)
                              like video watch duration etc.
          Further engagement: Total watch duration
                              Video completion rate

→ Labels :   Annotated
          +ve : Click + dwell time > t sec (+ any further
                                            engagement / user actions)

          -ve : Impressions only
          Sampling , focal loss (multiple modalities or imbalanced)

→ Features:   Image: — Captions derived from image (git-base
                                                    model)
                    — Tag supplied by uploader/article linked
                    — Image classification: list of all objects
                                    + overall scene description
                                    + characteristics of image:
                    — Image embeddings (ResNet) (e.g. close-up, black-
                                                          and white)
          Video: (Same as Image)
                    + Video embeddings (Video MVIT)
                    + duration of video
                    + uploader
                    + Video popularity

→ Model :   2 Approaches: (Can use combination)
               Query : word embedding       (blue balloon
                                            blue, balloon → all n-gram)

               ① 2 tower   (+ Ranking)
                              optional)
                           if corpus size small



               Query          Video/Image Embedding
               Embedding   + Video/Image features
          • + user features   + Video/Image Tags/Captions embedding.
             (if personalized)

(2) Query : word2vector (Deep box approach)

(Lightweight approach) Image/Video: all tags derived from image classification/video segmentation
+
textual features of video/image → Word2Vec

Cosine similarity b/w Query embedding vector + Image/Video
(word2vec)                                            emb. vector
                                                     (Word2Vec)

(Optional) formal training not required as both are use same
word2vec embedding techniques

→ Loss Fn: Focal loss
→ OFFLINE METRICS: MAP, Precision@K, NDCG, Top-K

→ Serving/Inference: Video/Image Indexing pipeline (to compute embeddings & ANN)
+
Text Indexing pipeline (to compute embeddings) of textual tags, titles,

# TIME SERIES ANOMALY DETECTION

→ Data is non iid (values in series are co-related)

→ Difficult to get labels i.e. Unavailability of labelled data
(few $\frac{or}{or}$)

→ STL Decomposition of Time-Series Data:
- Trend
- Seasonality
- Residual (Use this component to detect anomalies)
(or intrusion detection)

→ Features for cybersecurity Attack (or intrusion detection)
(Can convert time-series into supervised ML problem using feature-engineering e.g. tsfresh package)
- Headers of protocol
- Count features / Velocity features
- Metrics , Events , Logs , Traces
  e.g. CPU, e.g. login e.g. process (end-to-end request used attempt app logs path thou n/w)

→ N/w Features:
- I/o packets
- I/o drop
- I/o errors

→ App/Process Features
- CPU util
- I/o memory
- Threads
- file descriptors

System Features
- System disk I/o
- System CPU util.
- System CPU state (context switches, system calls, interrupts)
- System memory

→ Train-test split by time window

→ Balance false +ve with false -ve

→ Eval Metrics: MSE, MAE, RMSE, MAPE (Mean Abs % Error)

→ Model:
- Rule-based
- Statistical
- Traditional ML : Isolation forest / one Class SVM / K-Means
- DL : Auto-encoders

# PEOPLE YOU MAY KNOW

→ Goal: grow the network

→ Assumption: abt avg. no. of connections per user = 1000

→ Assumption: social graph of most users not very dynamic i.e. connections don't change significantly over a short period

→ ONLINE METRICS:
- Primary: # of conn. req. accepted in the last X days
- Secondary: # of conn. req. sent in the last X days
  (drawback: does not reflect real growth of n/w since acceptance is also needed)
- Counter: ignore or hide suggestions in the last X days

→ LABELS:  +ve: sent + accepted      -ve: hide or ignore suggestions or anything that is not +ve label

→ FEATURE - ENG:   Actors: Sender & Receiver

## Affinity Features

| User Features | Demographic Affinity | Education & Work Affinity | Social Affinity |
|---|---|---|---|
| - # of connections | - age | with - School, College, Major | - # of common connections |
| - # of followers, following | - gender | overlapping Job / Company / Industry | - Profile visits |
| - # of pending requests | - lang | time-period | - Common connections |
| - Account's age | - city, country | | wt. by time i.e. how long they have existed |
| - # of received reactions (influence) | | | |

→ CANDIDATE GENERATION: i) Friends of friends
(FoF) (2-hop neighbors)

→ Limit cand.: threshold on # of common connections (e.g. 10, 20 etc.)
→ Metric: Recall oriented

ii) Users from common groups

Avg 1000 conn per user
→ 1000 × 1000 of FoF on average
(search space reduction from Billion to Million)

→ Cold-start of new user: i) Based on profile info
ii) Influential users
popular

→ RANKING: i) Single Task : (Binary classification)
ii) Factorization Machines : Adv.
iii) GNN: GraphSage, Link prediction, Node embeddings
adv.

→ OFFLINE METRICS: MAP & NDCG

→ ONLINE vs BATCH PREDICTION:

<u>ONLINE</u>

→ Pro : Computed only for users who login or visit homepage (fraction of total users)

→ Con : Since recomm. are calc. on the fly, if takes long time then user experience is poor

DECISION: BATCH ← no delay better user experience
calc. it only for active users
social graph does not change quickly ie,
recommendations remain relevant for extended period

<u>BATCH</u>

→ Pro : Recommendations calc. beforehand, no delay

→ Con: Need to calc. it for all users

→ RE-RANKER: Diversity + Bump-up new user

→ POPULAR BIAS : Limit users who have high no. of connections

→ ~~ACCEPTANCE~~

→ DELAYED LABELS: Acceptance can take time, data analysis on historical acceptance time period should inform when a recommended connection is -ve label

# AUTOCOMPLETE OR TYPEAHEAD SYSTEM

→ Latency: Model inference on almost every keystroke  P90 < 50 ms

→ Personalization: Users have their own style, need to capture uniqueness of their personal style

→ Fairness & Privacy: Sensitive & confidential info should not be part of suggestions

→ GOAL:   i) Min. keystrokes needed to reach first relevant item
          ii) Rank user's INTENDED Query at top

→ ONLINE METRICS:   i) Avg. keystroke to get first relevant/click item
                    ii) Avg. position of click on the suggestion list
          Counter    iii) Avg. No. of similar searches before first click
          Metrics:

→ LABELS: — Query log of all users (decayed by time recency)  Prefix  suffix
          — Convert into prefix-suffix pairs  e.g. gdp | of Europe
            (multiple prefix-suffix pairs for a single query log)

          — Remove sensitive/confidential info from the query log

→ FEATURES: (Mostly used in Ranking phase)

| Prefix Features | Candidate Features | User Features ➡ | Previous Search History of User Features |
|---|---|---|---|
| —Prefix length | — suggestion length (char & words) | —Typing speed | (Cand.Prefix/ & search History) |
| —Does it end with a space char? | —Cand. freq. in query log for diff. time periods e.g. 1 week, 1 month | —Location of user  } Demographic features | candidate |
| (users most likely to peer at suggestions after a full word) | — Location Notion? | — Age | session history (short-term)   all prev. search history across session (long term) |
| | — Time Notion? | — Gender | —n-gram similarity |
| —Time of typing | —Cand. freq in querylog submitted by users in same age/gender/region group | — Language | —embedding similarity |
| | | —Avg. length of suggestion user clicked in the past | $q_1 → q_2 → q_3 → \ldots$ |
| | | —Similarity b/w suggestion words & words in prev. queries in same session | (session history) |
| | | e.g. cosine similarity, edit distance | —Term combination: added/removed  —Query similarity: cosine similarity  —CTR |

→ MODELING:

CANDIDATE GENERATION
i) Trie & inverted index (term → location of content)
(historical queries along with popularity)
issue! A significant proportion of queries issued daily have never been seen previously

ii) From query log, build prefix-suffix pairs → seq2seq Model using LSTM

iii) Using LLMs

RANKING: Factorization Machines
(After adding features) to candidates

→ OFFLINE METRICS : Cand. Generation: Recall based
Ranking: MRR
Minimum keystroke length (MKS)

→ RE-RANKING / FILTERING PHASE: Remove suggestions having
i) Age restriction   ii) Geography restriction
iii) Profane words
Bump up suggestions   i) Location-sensitive
ii) Time-sensitive

→ SPELL CORRECTION OR NON-PREFIX MATCHING:
i) Build a model with common misspelt words → correct words mapped
ii) While searching in trie, if no suggestions or less popular suggestions jump to a diff. node paying a cost dictated from conversion table. Conversion Table is created with misspelt words → correct words along with penalty

iii) Find other words with minimum edit distance