

ISSUES IN RECOMMENDATION SYSTEMS

1. COLD START PROBLEM: (user-item interaction not known)

a) COLD START OF ITEMS:

- We usually know the metadata of items (e.g. color category, price, etc.) hence find similar items and show them together (content based filtering)
- Learning period: boost the presence of items in recommendations for a certain period of time and decay it over time (randomly or after content based filtering)

b) COLD START OF USERS:

- Don't know anything about user
 - Popular content ~~2/3~~ ~~1/2~~ ~~1/3~~ ~~1/4~~ ~~1/5~~ ~~1/6~~ ~~1/7~~ ~~1/8~~ ~~1/9~~ ~~1/10~~ ~~1/11~~ ~~1/12~~ ~~1/13~~ ~~1/14~~ ~~1/15~~ ~~1/16~~ ~~1/17~~ ~~1/18~~ ~~1/19~~ ~~1/20~~ ~~1/21~~ ~~1/22~~ ~~1/23~~ ~~1/24~~ ~~1/25~~ ~~1/26~~ ~~1/27~~ ~~1/28~~ ~~1/29~~ ~~1/30~~ ~~1/31~~ ~~1/32~~ ~~1/33~~ ~~1/34~~ ~~1/35~~ ~~1/36~~ ~~1/37~~ ~~1/38~~ ~~1/39~~ ~~1/40~~ ~~1/41~~ ~~1/42~~ ~~1/43~~ ~~1/44~~ ~~1/45~~ ~~1/46~~ ~~1/47~~ ~~1/48~~ ~~1/49~~ ~~1/50~~ ~~1/51~~ ~~1/52~~ ~~1/53~~ ~~1/54~~ ~~1/55~~ ~~1/56~~ ~~1/57~~ ~~1/58~~ ~~1/59~~ ~~1/60~~ ~~1/61~~ ~~1/62~~ ~~1/63~~ ~~1/64~~ ~~1/65~~ ~~1/66~~ ~~1/67~~ ~~1/68~~ ~~1/69~~ ~~1/70~~ ~~1/71~~ ~~1/72~~ ~~1/73~~ ~~1/74~~ ~~1/75~~ ~~1/76~~ ~~1/77~~ ~~1/78~~ ~~1/79~~ ~~1/80~~ ~~1/81~~ ~~1/82~~ ~~1/83~~ ~~1/84~~ ~~1/85~~ ~~1/86~~ ~~1/87~~ ~~1/88~~ ~~1/89~~ ~~1/90~~ ~~1/91~~ ~~1/92~~ ~~1/93~~ ~~1/94~~ ~~1/95~~ ~~1/96~~ ~~1/97~~ ~~1/98~~ ~~1/99~~ ~~1/100~~ ~~1/101~~ ~~1/102~~ ~~1/103~~ ~~1/104~~ ~~1/105~~ ~~1/106~~ ~~1/107~~ ~~1/108~~ ~~1/109~~ ~~1/110~~ ~~1/111~~ ~~1/112~~ ~~1/113~~ ~~1/114~~ ~~1/115~~ ~~1/116~~ ~~1/117~~ ~~1/118~~ ~~1/119~~ ~~1/120~~ ~~1/121~~ ~~1/122~~ ~~1/123~~ ~~1/124~~ ~~1/125~~ ~~1/126~~ ~~1/127~~ ~~1/128~~ ~~1/129~~ ~~1/130~~ ~~1/131~~ ~~1/132~~ ~~1/133~~ ~~1/134~~ ~~1/135~~ ~~1/136~~ ~~1/137~~ ~~1/138~~ ~~1/139~~ ~~1/140~~ ~~1/141~~ ~~1/142~~ ~~1/143~~ ~~1/144~~ ~~1/145~~ ~~1/146~~ ~~1/147~~ ~~1/148~~ ~~1/149~~ ~~1/150~~ ~~1/151~~ ~~1/152~~ ~~1/153~~ ~~1/154~~ ~~1/155~~ ~~1/156~~ ~~1/157~~ ~~1/158~~ ~~1/159~~ ~~1/160~~ ~~1/161~~ ~~1/162~~ ~~1/163~~ ~~1/164~~ ~~1/165~~ ~~1/166~~ ~~1/167~~ ~~1/168~~ ~~1/169~~ ~~1/170~~ ~~1/171~~ ~~1/172~~ ~~1/173~~ ~~1/174~~ ~~1/175~~ ~~1/176~~ ~~1/177~~ ~~1/178~~ ~~1/179~~ ~~1/180~~ ~~1/181~~ ~~1/182~~ ~~1/183~~ ~~1/184~~ ~~1/185~~ ~~1/186~~ ~~1/187~~ ~~1/188~~ ~~1/189~~ ~~1/190~~ ~~1/191~~ ~~1/192~~ ~~1/193~~ ~~1/194~~ ~~1/195~~ ~~1/196~~ ~~1/197~~ ~~1/198~~ ~~1/199~~ ~~1/200~~ ~~1/201~~ ~~1/202~~ ~~1/203~~ ~~1/204~~ ~~1/205~~ ~~1/206~~ ~~1/207~~ ~~1/208~~ ~~1/209~~ ~~1/210~~ ~~1/211~~ ~~1/212~~ ~~1/213~~ ~~1/214~~ ~~1/215~~ ~~1/216~~ ~~1/217~~ ~~1/218~~ ~~1/219~~ ~~1/220~~ ~~1/221~~ ~~1/222~~ ~~1/223~~ ~~1/224~~ ~~1/225~~ ~~1/226~~ ~~1/227~~ ~~1/228~~ ~~1/229~~ ~~1/230~~ ~~1/231~~ ~~1/232~~ ~~1/233~~ ~~1/234~~ ~~1/235~~ ~~1/236~~ ~~1/237~~ ~~1/238~~ ~~1/239~~ ~~1/240~~ ~~1/241~~ ~~1/242~~ ~~1/243~~ ~~1/244~~ ~~1/245~~ ~~1/246~~ ~~1/247~~ ~~1/248~~ ~~1/249~~ ~~1/250~~ ~~1/251~~ ~~1/252~~ ~~1/253~~ ~~1/254~~ ~~1/255~~ ~~1/256~~ ~~1/257~~ ~~1/258~~ ~~1/259~~ ~~1/260~~ ~~1/261~~ ~~1/262~~ ~~1/263~~ ~~1/264~~ ~~1/265~~ ~~1/266~~ ~~1/267~~ ~~1/268~~ ~~1/269~~ ~~1/270~~ ~~1/271~~ ~~1/272~~ ~~1/273~~ ~~1/274~~ ~~1/275~~ ~~1/276~~ ~~1/277~~ ~~1/278~~ ~~1/279~~ ~~1/280~~ ~~1/281~~ ~~1/282~~ ~~1/283~~ ~~1/284~~ ~~1/285~~ ~~1/286~~ ~~1/287~~ ~~1/288~~ ~~1/289~~ ~~1/290~~ ~~1/291~~ ~~1/292~~ ~~1/293~~ ~~1/294~~ ~~1/295~~ ~~1/296~~ ~~1/297~~ ~~1/298~~ ~~1/299~~ ~~1/300~~ ~~1/301~~ ~~1/302~~ ~~1/303~~ ~~1/304~~ ~~1/305~~ ~~1/306~~ ~~1/307~~ ~~1/308~~ ~~1/309~~ ~~1/310~~ ~~1/311~~ ~~1/312~~ ~~1/313~~ ~~1/314~~ ~~1/315~~ ~~1/316~~ ~~1/317~~ ~~1/318~~ ~~1/319~~ ~~1/320~~ ~~1/321~~ ~~1/322~~ ~~1/323~~ ~~1/324~~ ~~1/325~~ ~~1/326~~ ~~1/327~~ ~~1/328~~ ~~1/329~~ ~~1/330~~ ~~1/331~~ ~~1/332~~ ~~1/333~~ ~~1/334~~ ~~1/335~~ ~~1/336~~ ~~1/337~~ ~~1/338~~ ~~1/339~~ ~~1/340~~ ~~1/341~~ ~~1/342~~ ~~1/343~~ ~~1/344~~ ~~1/345~~ ~~1/346~~ ~~1/347~~ ~~1/348~~ ~~1/349~~ ~~1/350~~ ~~1/351~~ ~~1/352~~ ~~1/353~~ ~~1/354~~ ~~1/355~~ ~~1/356~~ ~~1/357~~ ~~1/358~~ ~~1/359~~ ~~1/360~~ ~~1/361~~ ~~1/362~~ ~~1/363~~ ~~1/364~~ ~~1/365~~ ~~1/366~~ ~~1/367~~ ~~1/368~~ ~~1/369~~ ~~1/370~~ ~~1/371~~ ~~1/372~~ ~~1/373~~ ~~1/374~~ ~~1/375~~ ~~1/376~~ ~~1/377~~ ~~1/378~~ ~~1/379~~ ~~1/380~~ ~~1/381~~ ~~1/382~~ ~~1/383~~ ~~1/384~~ ~~1/385~~ ~~1/386~~ ~~1/387~~ ~~1/388~~ ~~1/389~~ ~~1/390~~ ~~1/391~~ ~~1/392~~ ~~1/393~~ ~~1/394~~ ~~1/395~~ ~~1/396~~ ~~1/397~~ ~~1/398~~ ~~1/399~~ ~~1/400~~ ~~1/401~~ ~~1/402~~ ~~1/403~~ ~~1/404~~ ~~1/405~~ ~~1/406~~ ~~1/407~~ ~~1/408~~ ~~1/409~~ ~~1/410~~ ~~1/411~~ ~~1/412~~ ~~1/413~~ ~~1/414~~ ~~1/415~~ ~~1/416~~ ~~1/417~~ ~~1/418~~ ~~1/419~~ ~~1/420~~ ~~1/421~~ ~~1/422~~ ~~1/423~~ ~~1/424~~ ~~1/425~~ ~~1/426~~ ~~1/427~~ ~~1/428~~ ~~1/429~~ ~~1/430~~ ~~1/431~~ ~~1/432~~ ~~1/433~~ ~~1/434~~ ~~1/435~~ ~~1/436~~ ~~1/437~~ ~~1/438~~ ~~1/439~~ ~~1/440~~ ~~1/441~~ ~~1/442~~ ~~1/443~~ ~~1/444~~ ~~1/445~~ ~~1/446~~ ~~1/447~~ ~~1/448~~ ~~1/449~~ ~~1/450~~ ~~1/451~~ ~~1/452~~ ~~1/453~~ ~~1/454~~ ~~1/455~~ ~~1/456~~ ~~1/457~~ ~~1/458~~ ~~1/459~~ ~~1/460~~ ~~1/461~~ ~~1/462~~ ~~1/463~~ ~~1/464~~ ~~1/465~~ ~~1/466~~ ~~1/467~~ ~~1/468~~ ~~1/469~~ ~~1/470~~ ~~1/471~~ ~~1/472~~ ~~1/473~~ ~~1/474~~ ~~1/475~~ ~~1/476~~ ~~1/477~~ ~~1/478~~ ~~1/479~~ ~~1/480~~ ~~1/481~~ ~~1/482~~ ~~1/483~~ ~~1/484~~ ~~1/485~~ ~~1/486~~ ~~1/487~~ ~~1/488~~ ~~1/489~~ ~~1/490~~ ~~1/491~~ ~~1/492~~ ~~1/493~~ ~~1/494~~ ~~1/495~~ ~~1/496~~ ~~1/497~~ ~~1/498~~ ~~1/499~~ ~~1/500~~ ~~1/501~~ ~~1/502~~ ~~1/503~~ ~~1/504~~ ~~1/505~~ ~~1/506~~ ~~1/507~~ ~~1/508~~ ~~1/509~~ ~~1/510~~ ~~1/511~~ ~~1/512~~ ~~1/513~~ ~~1/514~~ ~~1/515~~ ~~1/516~~ ~~1/517~~ ~~1/518~~ ~~1/519~~ ~~1/520~~ ~~1/521~~ ~~1/522~~ ~~1/523~~ ~~1/524~~ ~~1/525~~ ~~1/526~~ ~~1/527~~ ~~1/528~~ ~~1/529~~ ~~1/530~~ ~~1/531~~ ~~1/532~~ ~~1/533~~ ~~1/534~~ ~~1/535~~ ~~1/536~~ ~~1/537~~ ~~1/538~~ ~~1/539~~ ~~1/540~~ ~~1/541~~ ~~1/542~~ ~~1/543~~ ~~1/544~~ ~~1/545~~ ~~1/546~~ ~~1/547~~ ~~1/548~~ ~~1/549~~ ~~1/550~~ ~~1/551~~ ~~1/552~~ ~~1/553~~ ~~1/554~~ ~~1/555~~ ~~1/556~~ ~~1/557~~ ~~1/558~~ ~~1/559~~ ~~1/560~~ ~~1/561~~ ~~1/562~~ ~~1/563~~ ~~1/564~~ ~~1/565~~ ~~1/566~~ ~~1/567~~ ~~1/568~~ ~~1/569~~ ~~1/570~~ ~~1/571~~ ~~1/572~~ ~~1/573~~ ~~1/574~~ ~~1/575~~ ~~1/576~~ ~~1/577~~ ~~1/578~~ ~~1/579~~ ~~1/580~~ ~~1/581~~ ~~1/582~~ ~~1/583~~ ~~1/584~~ ~~1/585~~ ~~1/586~~ ~~1/587~~ ~~1/588~~ ~~1/589~~ ~~1/590~~ ~~1/591~~ ~~1/592~~ ~~1/593~~ ~~1/594~~ ~~1/595~~ ~~1/596~~ ~~1/597~~ ~~1/598~~ ~~1/599~~ ~~1/600~~ ~~1/601~~ ~~1/602~~ ~~1/603~~ ~~1/604~~ ~~1/605~~ ~~1/606~~ ~~1/607~~ ~~1/608~~ ~~1/609~~ ~~1/610~~ ~~1/611~~ ~~1/612~~ ~~1/613~~ ~~1/614~~ ~~1/615~~ ~~1/616~~ ~~1/617~~ ~~1/618~~ ~~1/619~~ ~~1/620~~ ~~1/621~~ ~~1/622~~ ~~1/623~~ ~~1/624~~ ~~1/625~~ ~~1/626~~ ~~1/627~~ ~~1/628~~ ~~1/629~~ ~~1/630~~ ~~1/631~~ ~~1/632~~ ~~1/633~~ ~~1/634~~ ~~1/635~~ ~~1/636~~ ~~1/637~~ ~~1/638~~ ~~1/639~~ ~~1/640~~ ~~1/641~~ ~~1/642~~ ~~1/643~~ ~~1/644~~ ~~1/645~~ ~~1/646~~ ~~1/647~~ ~~1/648~~ ~~1/649~~ ~~1/650~~ ~~1/651~~ ~~1/652~~ ~~1/653~~ ~~1/654~~ ~~1/655~~ ~~1/656~~ ~~1/657~~ ~~1/658~~ ~~1/659~~ ~~1/660~~ ~~1/661~~ ~~1/662~~ ~~1/663~~ ~~1/664~~ ~~1/665~~ ~~1/666~~ ~~1/667~~ ~~1/668~~ ~~1/669~~ ~~1/670~~ ~~1/671~~ ~~1/672~~ ~~1/673~~ ~~1/674~~ ~~1/675~~ ~~1/676~~ ~~1/677~~ ~~1/678~~ ~~1/679~~ ~~1/680~~ ~~1/681~~ ~~1/682~~ ~~1/683~~ ~~1/684~~ ~~1/685~~ ~~1/686~~ ~~1/687~~ ~~1/688~~ ~~1/689~~ ~~1/690~~ ~~1/691~~ ~~1/692~~ ~~1/693~~ ~~1/694~~ ~~1/695~~ ~~1/696~~ ~~1/697~~ ~~1/698~~ ~~1/699~~ ~~1/700~~ ~~1/701~~ ~~1/702~~ ~~1/703~~ ~~1/704~~ ~~1/705~~ ~~1/706~~ ~~1/707~~ ~~1/708~~ ~~1/709~~ ~~1/710~~ ~~1/711~~ ~~1/712~~ ~~1/713~~ ~~1/714~~ ~~1/715~~ ~~1/716~~ ~~1/717~~ ~~1/718~~ ~~1/719~~ ~~1/720~~ ~~1/721~~ ~~1/722~~ ~~1/723~~ ~~1/724~~ ~~1/725~~ ~~1/726~~ ~~1/727~~ ~~1/728~~ ~~1/729~~ ~~1/730~~ ~~1/731~~ ~~1/732~~ ~~1/733~~ ~~1/734~~ ~~1/735~~ ~~1/736~~ ~~1/737~~ ~~1/738~~ ~~1/739~~ ~~1/740~~ ~~1/741~~ ~~1/742~~ ~~1/743~~ ~~1/744~~ ~~1/745~~ ~~1/746~~ ~~1/747~~ ~~1/748~~ ~~1/749~~ ~~1/750~~ ~~1/751~~ ~~1/752~~ ~~1/753~~ ~~1/754~~ ~~1/755~~ ~~1/756~~ ~~1/757~~ ~~1/758~~ ~~1/759~~ ~~1/760~~ ~~1/761~~ ~~1/762~~ ~~1/763~~ ~~1/764~~ ~~1/765~~ ~~1/766~~ ~~1/767~~ ~~1/768~~ ~~1/769~~ ~~1/770~~ ~~1/771~~ ~~1/772~~ ~~1/773~~ ~~1/774~~ ~~1/775~~ ~~1/776~~ ~~1/777~~ ~~1/778~~ ~~1/779~~ ~~1/780~~ ~~1/781~~ ~~1/782~~ ~~1/783~~ ~~1/784~~ ~~1/785~~ ~~1/786~~ ~~1/787~~ ~~1/788~~ ~~1/789~~ ~~1/790~~ ~~1/791~~ ~~1/792~~ ~~1/793~~ ~~1/794~~ ~~1/795~~ ~~1/796~~ ~~1/797~~ ~~1/798~~ ~~1/799~~ ~~1/800~~ ~~1/801~~ ~~1/802~~ ~~1/803~~ ~~1/804~~ ~~1/805~~ ~~1/806~~ ~~1/807~~ ~~1/808~~ ~~1/809~~ ~~1/810~~ ~~1/811~~ ~~1/812~~ ~~1/813~~ ~~1/814~~ ~~1/815~~ ~~1/816~~ ~~1/817~~ ~~1/818~~ ~~1/819~~ ~~1/820~~ ~~1/821~~ ~~1/822~~ ~~1/823~~ ~~1/824~~ ~~1/825~~ ~~1/826~~ ~~1/827~~ ~~1/828~~ ~~1/829~~ ~~1/830~~ ~~1/831~~ ~~1/832~~ ~~1/833~~ ~~1/834~~ ~~1/835~~ ~~1/836~~ ~~1/837~~ ~~1/838~~ ~~1/839~~ ~~1/840~~ ~~1/841~~ ~~1/842~~ ~~1/843~~ ~~1/844~~ ~~1/845~~ ~~1/846~~ ~~1/847~~ ~~1/848~~ ~~1/849~~ ~~1/850~~ ~~1/851~~ ~~1/852~~ ~~1/853~~ ~~1/854~~ ~~1/855~~ ~~1/856~~ ~~1/857~~ ~~1/858~~ ~~1/859~~ ~~1/860~~ ~~1/861~~ ~~1/862~~ ~~1/863~~ ~~1/864~~ ~~1/865~~ ~~1/866~~ ~~1/867~~ ~~1/868~~ ~~1/869~~ ~~1/870~~ ~~1/871~~ ~~1/872~~ ~~1/873~~ ~~1/874~~ ~~1/875~~ ~~1/876~~ ~~1/877~~ ~~1/878~~ ~~1/879~~ ~~1/880~~ ~~1/881~~ ~~1/882~~ ~~1/883~~ ~~1/884~~ ~~1/885~~ ~~1/886~~ ~~1/887~~ ~~1/888~~ ~~1/889~~ ~~1/890~~ ~~1/891~~ ~~1/892~~ ~~1/893~~ ~~1/894~~ ~~1/895~~ ~~1/896~~ ~~1/897~~ ~~1/898~~ ~~1/899~~ ~~1/900~~ ~~1/901~~ ~~1/902~~ ~~1/903~~ ~~1/904~~ ~~1/905~~ ~~1/906~~ ~~1/907~~ ~~1/908~~ ~~1/909~~ ~~1/910~~ ~~1/911~~ ~~1/912~~ ~~1/913~~ ~~1/914~~ ~~1/915~~ ~~1/916~~ ~~1/917~~ ~~1/918~~ ~~1/919~~ ~~1/920~~ ~~1/921~~ ~~1/922~~ ~~1/923~~ ~~1/924~~ ~~1/925~~ ~~1/926~~ ~~1/927~~ ~~1/928~~ ~~1/929~~ ~~1/930~~ ~~1/931~~ ~~1/932~~ ~~1/933~~ ~~1/934~~ ~~1/935~~ ~~1/936~~ ~~1/937~~ ~~1/938~~ ~~1/939~~ ~~1/940~~ ~~1/941~~ ~~1/942~~ ~~1/943~~ ~~1/944~~ ~~1/945~~ ~~1/946~~ ~~1/947~~ ~~1/948~~ ~~1/949~~ ~~1/950~~ ~~1/951~~ ~~1/952~~ ~~1/953~~ ~~1/954~~ ~~1/955~~ ~~1/956~~ ~~1/957~~ ~~1/958~~ ~~1/959~~ ~~1/960~~ ~~1/961~~ ~~1/962~~ ~~1/963~~ ~~1/964~~ ~~1/965~~ ~~1/966~~ ~~1/967~~ ~~1/968~~ ~~1/969~~ ~~1/970~~ ~~1/971~~ ~~1/972~~ ~~1/973~~ ~~1/974~~ ~~1/975~~ ~~1/976~~ ~~1/977~~ ~~1/978~~ ~~1/979~~ ~~1/980~~ ~~1/981~~ ~~1/982~~ ~~1/983~~ ~~1/984~~ ~~1/985~~ ~~1/986~~ ~~1/987~~ ~~1/988~~ ~~1/989~~ ~~1/990~~ ~~1/991~~ ~~1/992~~ ~~1/993~~ ~~1/994~~ ~~1/995~~ ~~1/996~~ ~~1/997~~ ~~1/998~~ ~~1/999~~ ~~1/1000~~ ~~1/1001~~ ~~1/1002~~ ~~1/1003~~ ~~1/1004~~ ~~1/1005~~ ~~1/1006~~ ~~1/1007~~ ~~1/1008~~ ~~1/1009~~ ~~1/1010~~ ~~1/1011~~ ~~1/1012~~ ~~1/1013~~ ~~1/1014~~ ~~1/1015~~ ~~1/1016~~ ~~1/1017~~ ~~1/1018~~ ~~1/1019~~ ~~1/1020~~ ~~1/1021~~ ~~1/1022~~ ~~1/1023~~ ~~1/1024~~ ~~1/1025~~ ~~1/1026~~ ~~1/1027~~ ~~1/1028~~ ~~1/1029~~ ~~1/1030~~ ~~1/1031~~ ~~1/1032~~ ~~1/1033~~ ~~1/1034~~ ~~1/1035~~ ~~1/1036~~ ~~1/1037~~ ~~1/1038~~ ~~1/1039~~ ~~1/1040~~ ~~1/1041~~ ~~1/1042~~ ~~1/1043~~ ~~1/1044~~ ~~1/1045~~ ~~1/1046~~ ~~1/1047~~ ~~1/1048~~ ~~1/1049~~ ~~1/1050~~ ~~1/1051~~ ~~1/1052~~ ~~1/1053~~ ~~1/1054~~ ~~1/1055~~ ~~1/1056~~ ~~1/1057~~ ~~1/1058~~ ~~1/1059~~ ~~1/1060~~ ~~1/1061~~ ~~1/1062~~ ~~1/1063~~ ~~1/1064~~ ~~1/1065~~ ~~1/1066~~ ~~1/1067~~ ~~1/1068~~ ~~1/1069~~ ~~1/1070~~ ~~1/1071~~ ~~1/1072~~ ~~1/1073~~ ~~1/1074~~ ~~1/1075~~ ~~1/1076~~ ~~1/1077~~ ~~1/1078~~ ~~1/1079~~ ~~1/1080~~ ~~1/1081~~ ~~1/1082~~ ~~1/1083~~ ~~1/1084~~ ~~1/1085~~ ~~1/1086~~ ~~1/1087~~ ~~1/1088~~ ~~1/1089~~ ~~1/1090~~ ~~1/1091~~ ~~1/1092~~ ~~1/1093~~ ~~1/1094~~ ~~1/1095~~ ~~1/1096~~ ~~1/1097~~ ~~1/1098~~



2. EXPLORATION - EXPLOITATION / DIVERSITY | SERVING - BIAS / FEEDBACK LOOP ISSUE

- Labels generated from recommendations (ranking model) Algorithmic Bias of RecSys becomes training data for next iteration of RecSys Training
 - Feedback loop issue
 - Serving bias
 - Diversity \downarrow

- To mitigate above issues, we need to have strategy of exploration vs exploitation

1) Multiple Candidate generators

2) ϵ -greedy approach to select when to explore and when to exploit

3) Heuristics: (certain % of times in the recommended list use O/P of RecSys and remaining $(100 - x)\%$ of times use random items)

4) Re-ranking: Penalize items in the final ranked list which are above a certain threshold in Similarity to ranked above items

Measuring Diversity: Avg pairwise similarity of an item to all other items ranked above it \downarrow \rightarrow high diversity

Theoretically, ranking models should be trained solely on unbiased datasets. However, the exploration dataset may be too small.

3. BIAS IN RECSYS

1. POSITION BIAS: Occurs when items at the top of the list are more likely to be selected by users

MITIGATION:

1) Add positional features:

Training: Add position of items as feature

Serving: set all items to position = 1 as feature for ranking (to negate impact of position)

2) Model position bias: Model position and relevancy of item separately and then normalizes using position

$$P_{i,p} = R_i \times P_p$$

Observed prob. of click on item i on position p True relevance of item i Prob of click on position p

$$\therefore R_i = \frac{P_{i,p}}{P_p}$$

3) Use Interleaving in A/B Tests

2. POPULARITY BIAS: Final Ranked list shown to user has large number of popular items

MITIGATION:

- 1) Penalize items in training set by their popularity i.e. # of times bought, clicked, etc (sort all items in training set by # of times engaged with and select limited number of pop. items)
- 2) Add diversity in training set / candidate generators
- 3) GF classification algo used in ranking phase penalize the O/P with freq. of item engagement i.e. bought, clicked, etc
- 4) Re-ranking: only keep a certain % of items with popularity above a certain threshold

3. SERVING BIAS - (Mentioned in 2)

4. CLICKBAIT BIAS: Label based only on clicks X
MITIGATION: watch length or combination of factors
(e.g. likes, comment, etc.)

5. DURATION BIAS: If watched or not watched \rightarrow label criteria
then shorter videos more likely to be watched
than longer videos

MITIGATION: watch length in quantile buckets
& predict watch quantile

4. NEGATIVE SAMPLING: (More # of -ve than +ve labels available then how to sample -ve labels)

— Popularity based sampling: Sample -ve items based on their popularity where less popular items are preferred over highly popular items.

— Hard -ve sampling: Selects -ve items that are similar to +ve items

— Cluster -ve items based on criteria such as genre, popularity, item creator, etc and pick certain % from each cluster

FRESHNESS: Item age as feature

Re-ranking stage \rightarrow bump up ranking based on freshness

CANDIDATE GENERATION

SOURCES :

1) Non-personalized sources

- Popular content
- Trending content
- New content
- Trending/Popular content in a geography, age grp or topics followed

2) Personalized

~~networks~~ sources:

- In-network content: items generated by users/businesses that the user is currently following/connected to

- Historical content/
content based filtering

Items that are similar to items that the user has liked, previously bought/engaged with or shown interest in. E.g. based on characteristics such as genre, topic, keywords, media, etc.

- Collaborative filtering: Items that are recommended based on the similarity of the user's preferences to those of other users in the system

ALGORITHMS FOR GENERATING CANDIDATES: (can be pre-computed as well)

1. Matrix Factorization
2. 2 Tower Neural Network (either only user features or item features or both)
3. Graph based
4. Neighborhood based (with or without embedding of user or items)
interaction matrix only } & using Pearson correlation coeff

LIMITING PRUNING CANDIDATES:

1. Heuristic: Ratio of candidates from each cand. generation method based on business objectives (e.g. product diversity, market penetration, sales etc.)
2. Thresholding: Limit # of candidates from each cand. gen. method based on thresholding of scoring function from that algo
3. Universal Score: Calibrate all algorithms score and put a universal threshold

SAMPLING-NEGATIVE EXAMPLES IN RETRIEVAL PHASE FOR 2 TOWER

1. Batch Random Negative:

For a query, an item is clicked and rest of the items that are not clicked (-ve)

If sample from these not clicked items randomly: random -ve

However, there is computational cost attached to

- Fetching features for each of -ve items
- Each -ve item has to go through item tower

For a given batch of queries, if the ~~same~~ random -ve set is used for all +ve examples = Batch Random Negative

Training speed is much faster

2. Batch Negative:

query 1 → clicked item 1 : +ve example
query 2 → clicked item 2 : -ve example (for query 1)

— We can use clicked item 2, etc as -ve example for query 1 - clicked item 1 : +ve example

— Each item embedding in the batch serves once as +ve example and $|B|-1$ times -ve example for other queries in the batch ($|B|$ = size of batch)

Adv: Same advantage as batch random -ve

Disadv: Training data has only clicked items and users ~~not~~ interacting with small set of popular items leading to bias (not so popular items will be at disadvantage)

Items that are not favored by the existing system are less likely to get user feedback. Accordingly, sampling only batch -ves will end up with a model lacking resolution for long-tail items, which seldom appear in training data.

3. Hard -ve Mining: (Batch Hard Negatives)

- To make model better at differentiating b/w easy and hard examples, an iterative approach can be utilized

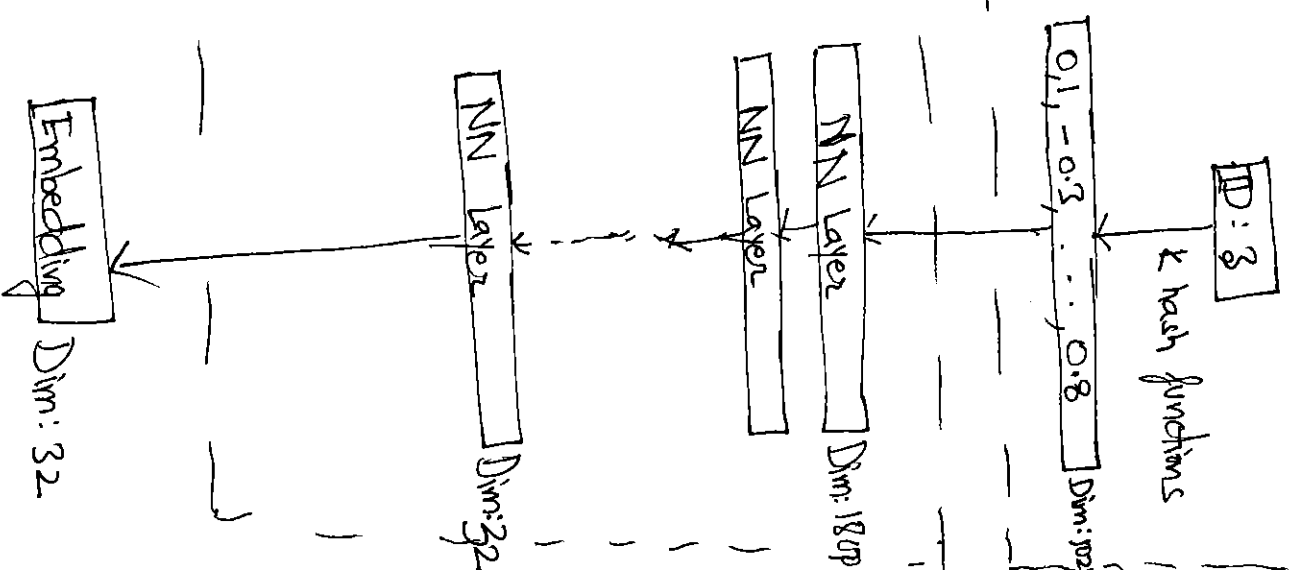
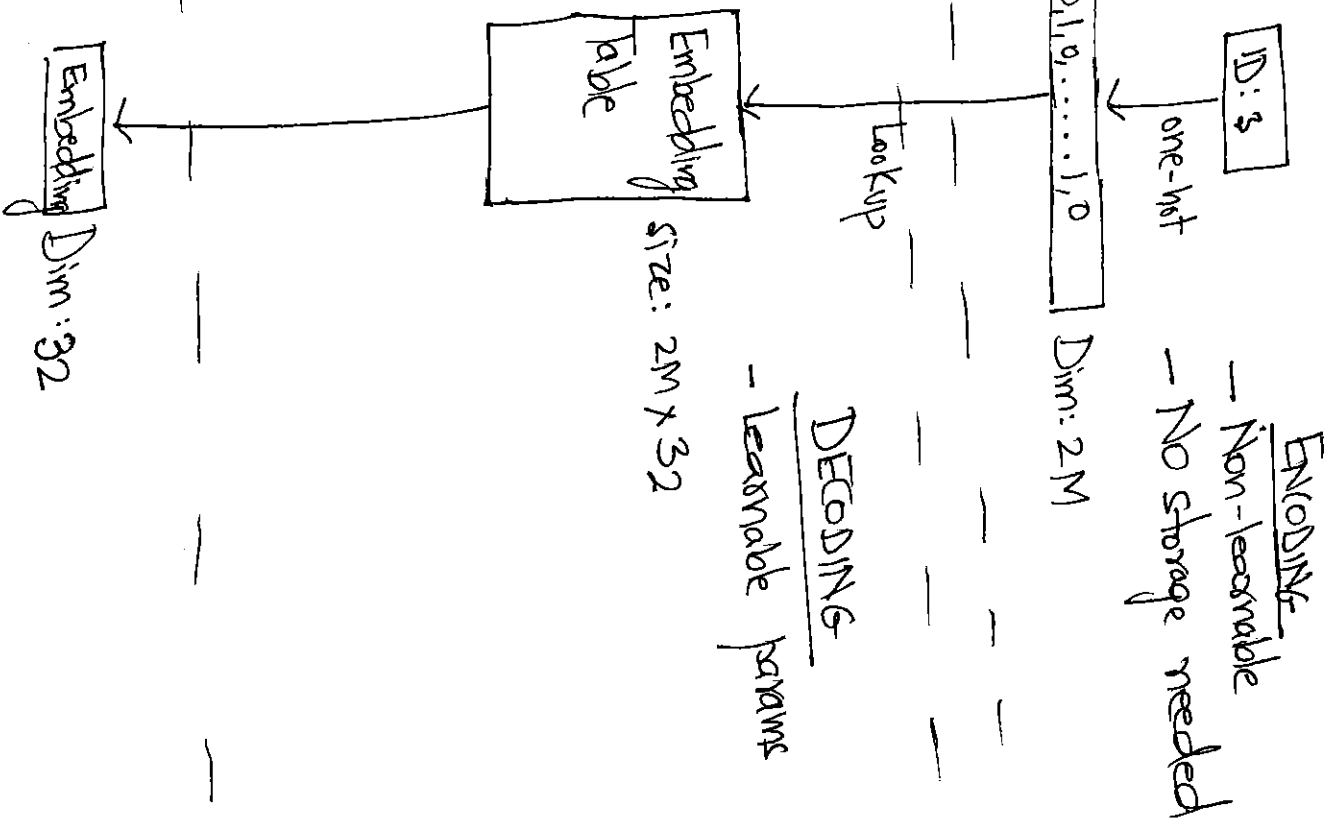
Pass 1: Generate -ve candidate/item set by using in-batch random -ve and/or in-batch -ve and training the two-tower (along with +ve candidate/item set) encoders

Pass 2: Use ANN to get -ve examples/items which are closer to input query (hard -ves) and re-train

NOTE: The dataset need not be balanced at this stage. Can use focal loss or balanced cross-entropy to account for it.

DEEP HASH EMBEDDING

- Embedding tables or look-up tables that map sparse id features to low-dimensional dense vectors are core components of modern recsys.
- We first hash the id into a one-hot hash vector, use an embedding layer (fully connected layer) ^{during training} and convert into low-dimensional dense vector.
- At serving time, this embedding layer can be used as a lookup table i.e. embedding table for sparse ids.
- By design, embedding tables are memory-hungry.
e.g. A feature with 5000 sparse ids connected to 128 dimension vector is of size (5000, 128)
- For fast inference, these embedding tables reside in primary memory (RAM, ^{L1/L2 cache} or SSDs)
CPU memory
- The key idea behind Deep Hash Embedding (DHE) is to replace embedding tables with an encoder-decoder architecture that has a much smaller memory footprint.



→ Encoding step: Input id → apply K -hash fns → K -dimensional dense vector normalized in $[-1, 1]$

Decoding step: K -dimensional dense vector normalized in range $[-1, 1]$ → small NN (snaps) → semantically meaningful embeddings

→ Advantages: 1) No hash collisions (due to large no. of hash fns e.g. 1024)
equal to no. of dimensions

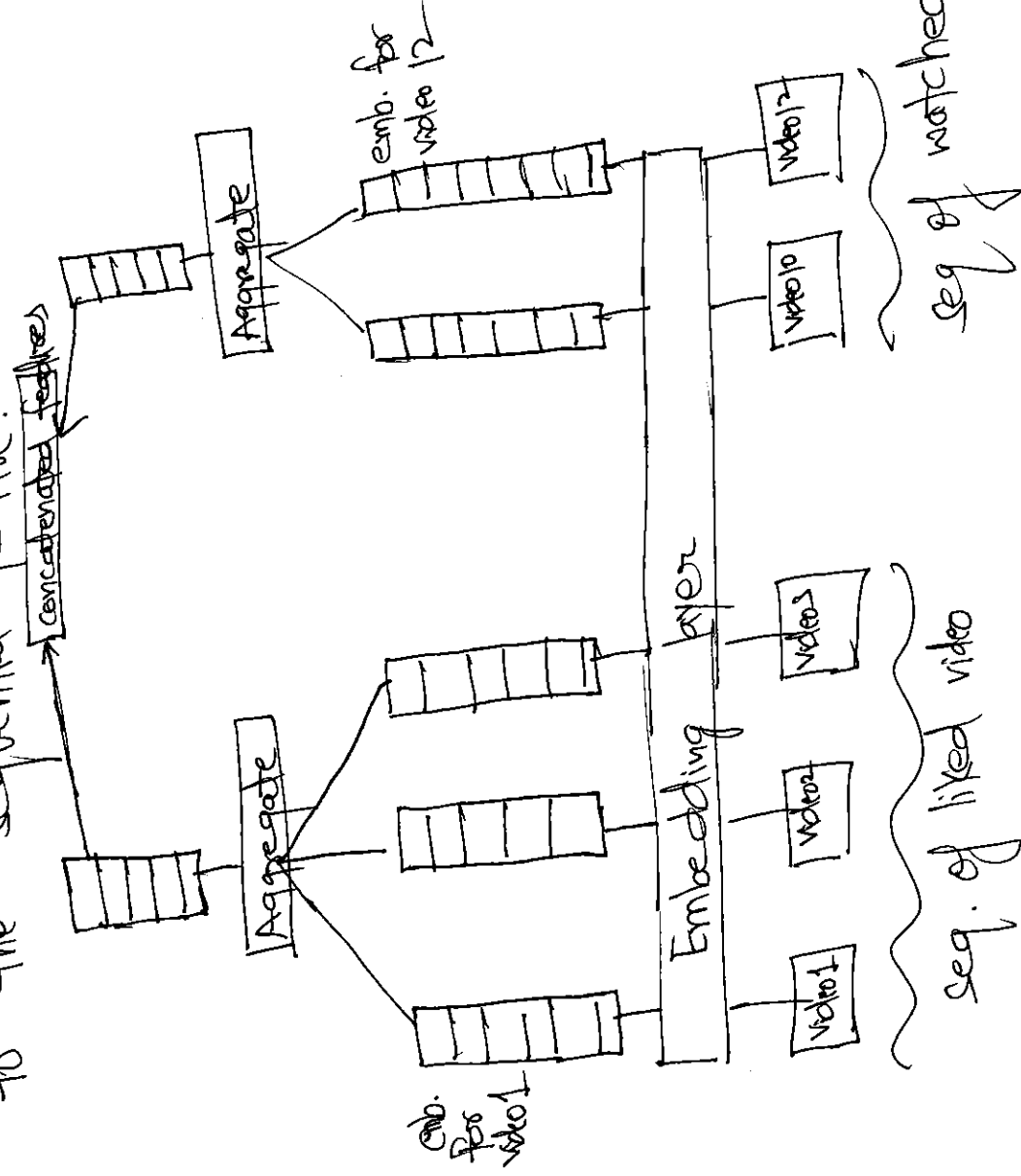
2) Low memory footprint since no sparse required for embedding table

WAYS TO ENCODE SEQUENCE FEATURES

1. Seq. features e.g. list of liked videos, list of watched videos.

Each item id is a sparse feature encoded via an embedding layer i.e. o/p of embedding layer is the dense, low-dimensional representation of item Id.

Now for a sequence of item ids, they are fed as separate feature into embedding layer and then aggregated (avg.) to get embedding vector corresponding to the sequential feature.



2. If 2 tower is used in Card. generation
then we have embeddings corresponding to each item
in the item seq. \rightarrow aggregate them
(avg.)

3. Use an attention block and the O/P of attention
block $(K, N) \rightarrow$ concatenate them to get one
feature vector. (E.g. as used in MMoE network)

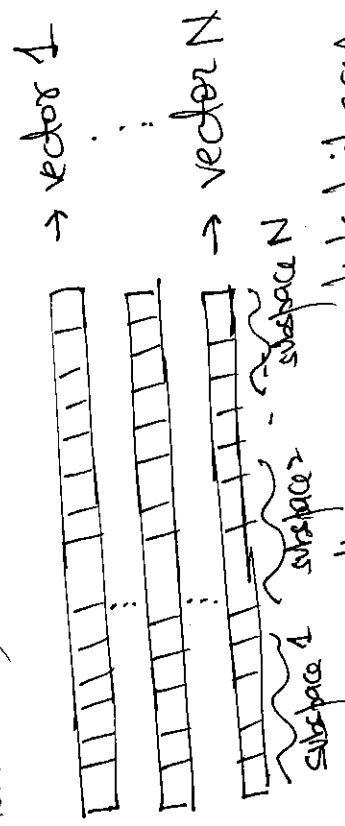
QUANTIZATION AND INDEXING FOR ANN

→ If the no. of items or users is huge for a use case then Storage and search for vector embeddings is problematic and needs scaling

→ STORAGE SOLVE: 1) Use less no. of bits to ~~store~~ embeddings (e.g. float 16 etc.)
2) Product Quantization (PQ)

→ SEARCH SOLVE : INDEXES

Combination of { 1) Inverted Index
yields highly scalable 2) HNSW (Hierarchical Navigable Small World)
Search (e.g. > billion item search)



PRODUCT QUANTIZATION:

→ Divide each vector into smaller segments/subspaces

→ Run K-means clustering on each subspace
(e.g. $K=256$ centroids)


→ For each subspace/segment of vector, find the nearest Centroid (e.g. centroid 240) and substitute it with that Centroid ID

→ Example: if 256 Centroids then 8 bits can represent the indices of these 256 centroids ($2^8=256$)


→ Essentially, we have encoded our original vectors with Centroid IDs (i.e. query of centroid ids = PA code for the vector)

INVERTED INDEX: → Run K-mean clustering on PA codes of vectors & the centroids so obtained are indexed/


stored in memory

→  → List of PA codes for vectors

stored in database

 → List of PA codes for vectors

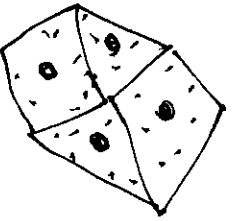
⋮

 → List of PA codes for vectors

~~~~~

Partition (Centroids)  
Inverted Index

→ This structure is also referred to as "Voronoi cell"



• Partition (Centroid)  
⋮  
~~PA~~ code of vectors (usually in database)

2. imp. params:

i) no. of cells to create: (K in K-means)

if it is high, more cells are created but fewer vectors to search within each cell → prioritizes search-speed

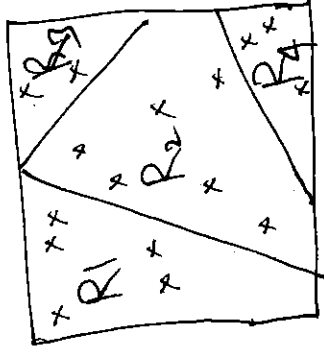
ii) no. of cells to search:

if it is high, more time needed for search but better quality, especially for edge query vectors → prioritizes search quality

# APPROXIMATE NEAREST NEIGHBOR (ANN)

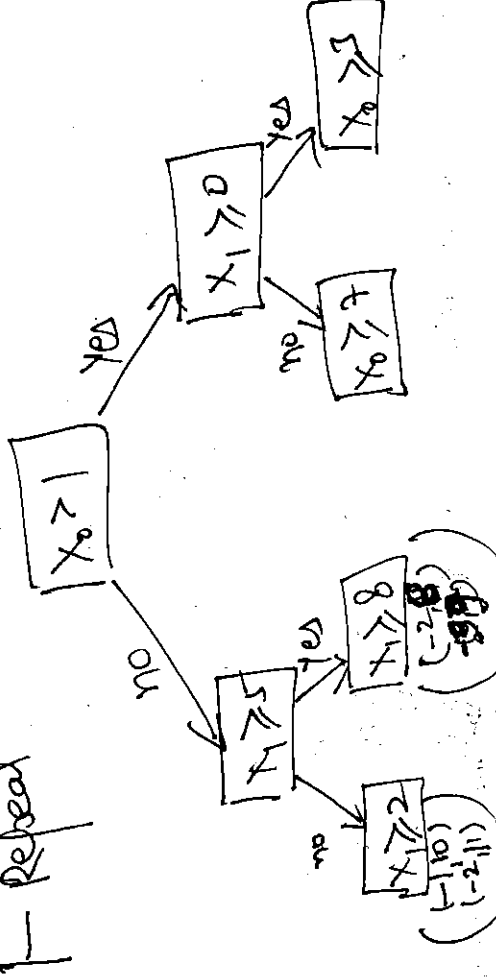
## 1. Tree based ANN (ANN - SPOTIFY)

- Split the space into multiple partitions using tree structure & leverage characteristics of tree to perform faster search
- Non-leaf nodes split the space into two partitions given the criterion
- Leaf nodes represent a particular region in space (partition)
- To find nearest neighbors, the algo only searches the partition the query point belongs to



Example: K-d tree, Annoy

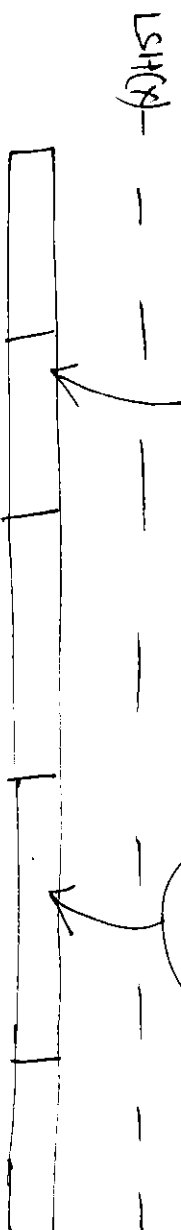
- Pick a random dimension from your K-dim vector
- Find the median of that dimension across all of the vectors in your current collection
- Split the vector on the median value
- Repeat



Partitioned vectors

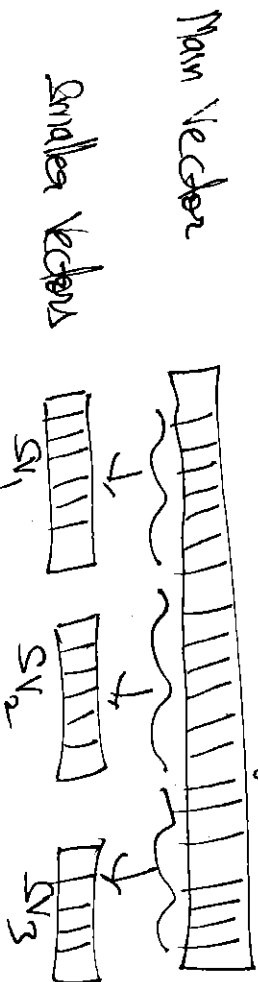
2. Locality Sensitive Hashing  $\rightarrow$  FFAISS (Facebook)

- Used hash fn to bucket pts
- Pts in the same bucket are closer to each other than pts in another bucket
- LSH only searches pts belonging to the same bucket as query pt
- Instead of minimizing hash collisions, "maximize hash collisions" so that pts/rects in same bucket are closer



3. Quantization based  $\rightarrow$  SVM (Google) / FACS (Facebook)

- Represent pts closer to each other by a "representative vector"
- Only performing KNN on all pts and representing pts belonging to the cluster by its centroid  $\rightarrow$  accuracy
- Instead divide the vector into smaller vectors and find centroids for each of the smaller vectors and represent the vector by a combination of the centroids of these smaller vectors



Find K centroids/clusters of

clusters of all  $SN_1$   $SN_2$   $SN_3$   $SN_4$  (3 centroids)

$SK_{11}$   $SK_{12}$   $SK_{13}$   $SK_{14}$

$SK_{21}$   $SK_{22}$   $SK_{23}$   $SK_{24}$  (4 centroids)

$SK_{31}$   $SK_{32}$

$2 \times 4 \times 2 = 16$  (9 centroids)

$\equiv$  comb. of 24 cent. to be examined

## CHALLENGE IN TRAINING MULTI-MODAL SYSTEMS:

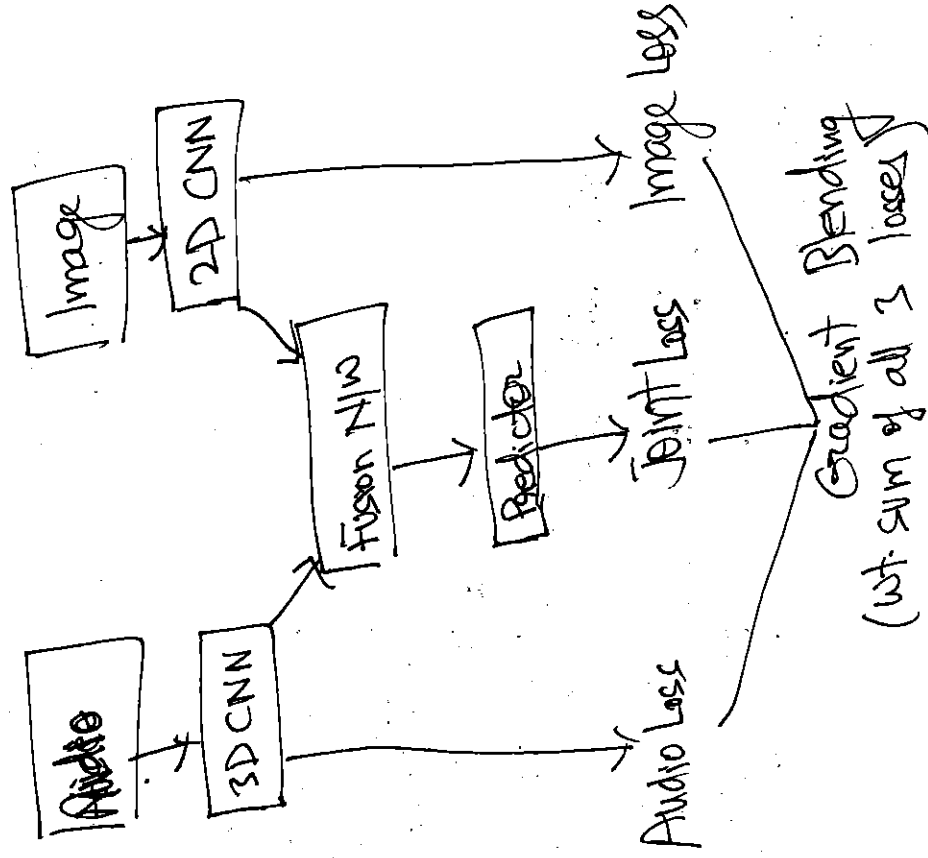
— Over-fitting (when the learning speed varies across different modalities, one modality (e.g. image) can dominate the learning process)

— Techniques to overcome overfitting in multi-model settings

— Gradient Blending

— Focal Loss

1) Gradient Blending: Adjusting loss function to take into account each model's individual wts as well as combined wts.



## 2) Focal loss: (Better than Cross Entropy Loss)

- Cases where cross-entropy loss performs badly:

1) Class imbalance: - Majority class examples will dominate the loss function and gradient descent, causing it to be updated in the direction of the model becoming more confident in predicting majority class while putting less emphasis on minority class.

- Balanced Cross-Entropy solves above problem by adding a wt. factor to each class (could be inverse class freq. or a hyper-parameter determined by cross-validation): Same is done by focal loss

2) Fails to distinguish between hard and easy examples

- Hard examples are those in which the model repeatedly makes huge error whereas easy examples are those which are easily classified

Focal loss =

$$-CE = -p_t$$

$$-FL = -(1-p_t)^\gamma \log p_t$$

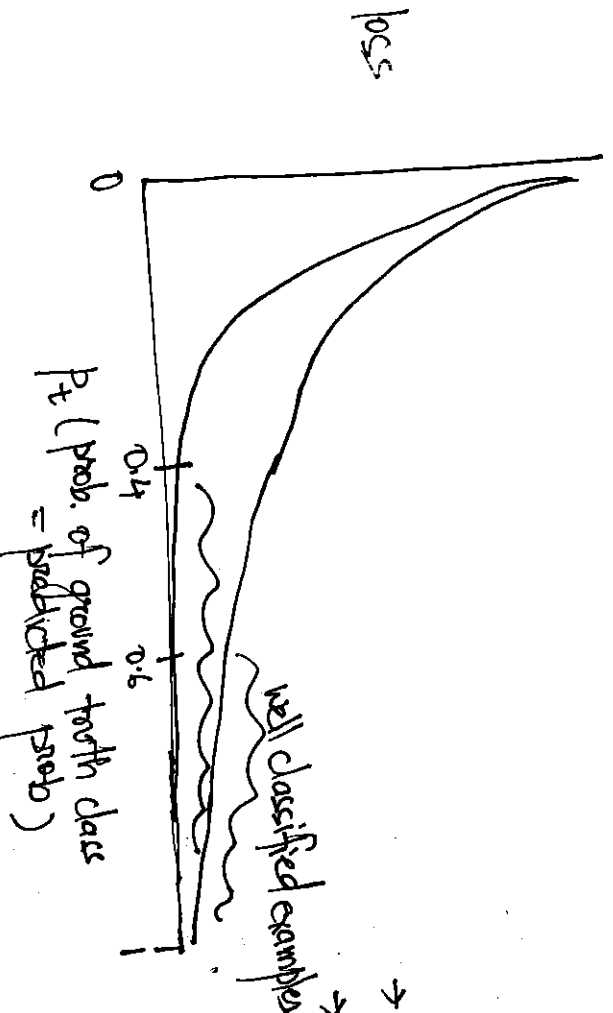
where  $\gamma$  = modulating factor

→ when  $\gamma = 0$ ,  $FL = CE$

→ loss is low when  $p_t \in [0.1, 1]$

area of well classified easy examples

even examples that are easily classified ( $p \gg 0.5$ ), incur a loss with non-trivial magnitude





→ Now, as we increase the value of  $\gamma$ , we slowly extend that range of predicted prob. where loss is kw to  $\sim [0.4, 1]$

"

This means we are "extending" or "relaxing" our criteria of well classified examples

→ This means with standard CE loss, the model will push scores even for well-classified examples further & further away till the predicted prob reaches 1

Whereas

the model trained with FL will not care much abt well classified examples (will not push their prob to 1), instead will focus on improving/reducing loss for hard examples (whose error is higher)

wt. factor  $\nearrow$  modulating factor

$$\text{Hence } FL(p) = \begin{cases} -\alpha (1-p)^{\gamma} \log p & ; \gamma = 1 \\ -(1-\alpha) p^{\gamma} \log(1-p) & ; \text{otherwise} \end{cases}$$



# MULTIPLE OBJECTIVE CLASSIFICATION SYSTEM

→ Example: 1) Detect harmful content

- Nudity
- Violence
- Hate

2) Increase Engagement

- Likes
- Comment
- share
- click
- Dwell time

→ OPTIONS:

1) Single binary classifier:

Harmful or Engagement Prob

Model

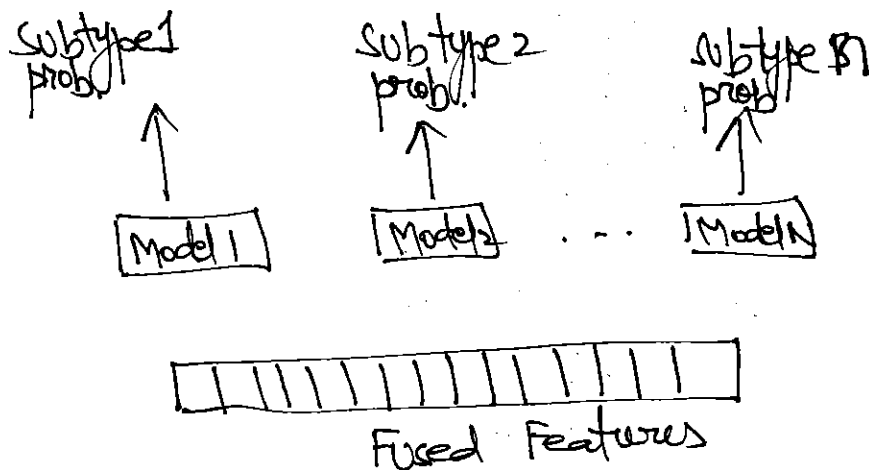


Fused features

Issues: Difficult to determine which subtype of harm or engagement  
— Cannot inform users or business which subtype resulted in the outcome

— Cannot improve performance in a subtype, if the system is bad in some subtypes

2.) One binary classifier per subtype

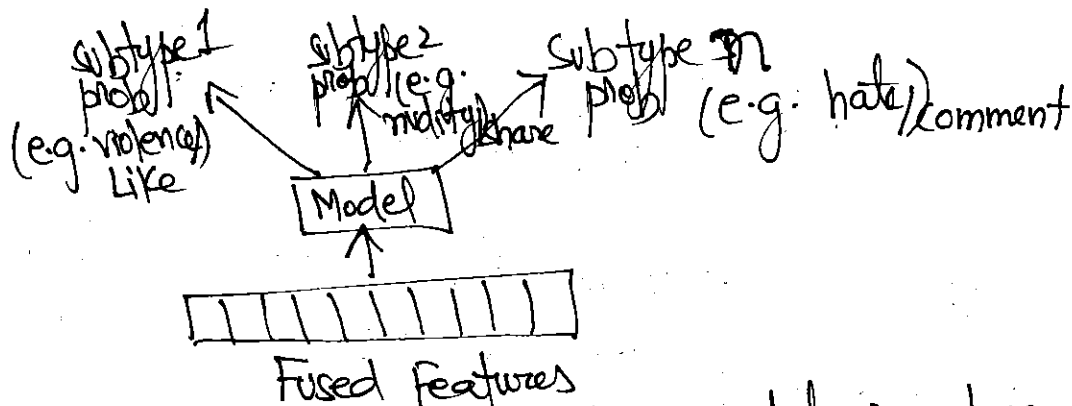


Adv. 1) Can explain the outcome in terms of subtype  
e.g. why a post was taken down

2) Can improve the model for a particular subtype,  
if the subtype model performance is poor

Disadv. 1) Training & maintaining  $n$  models is time consuming & expensive

3) Multi label classifier



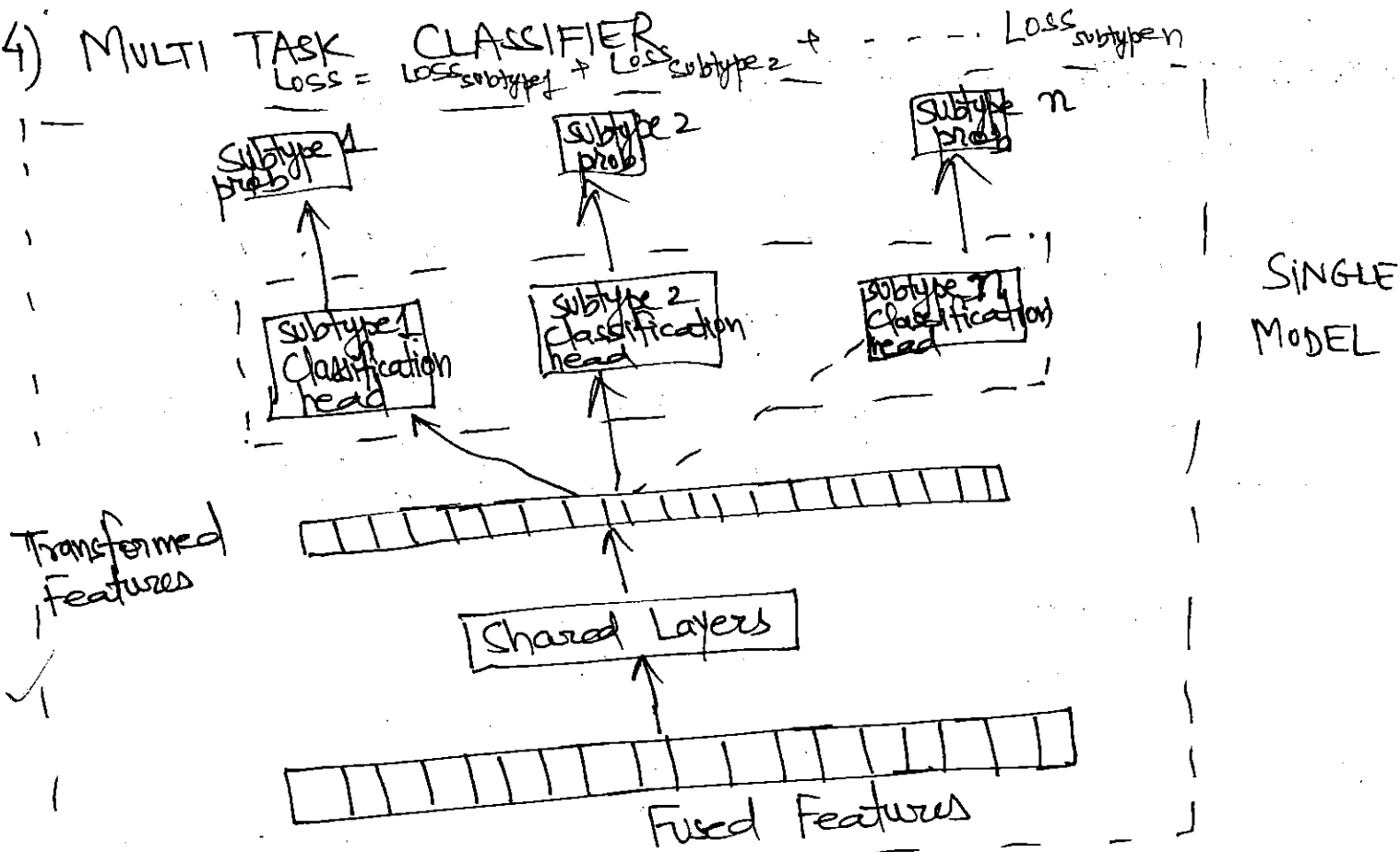
Adv. 1) Training & Maintaining the model is less costly

Disadv. 1) Each observation in training data ~~needs~~ have to be labeled for all subtypes

2) If diff. subtypes require diff. feature transformations, not possible

#### 4) MULTI TASK CLASSIFIER

$$\text{Loss} = \text{Loss}_{\text{subtype 1}} + \text{Loss}_{\text{subtype 2}} + \dots + \text{Loss}_{\text{subtype n}}$$



Shared Layers: Common across all subtypes

Task Specific Layers: Set of independent ML Layers (Classification head)

Each classification head transforms features in a way that is optimal for predicting a specific subtype (eg. harm) probability

Adv: 1) Not expensive to train or maintain since it is single model

2) Shared layers transform features in a way that is beneficial for each task (prevents redundant feature transformations) (↓ decreases training time compared to separate layers)

3) Each obs. in training data "need" not be labeled for all subtypes (e.g. weather and animal type prediction, but needed in posts → comments, like, ret) and training data of each subtype contributes to the learning of other subtypes (especially useful if one subtype has less training data)

## NOTE ON MULTI-TASK CLASSIFICATION:

Usually employed for scenarios where multi-label for different tasks (usually related tasks) on the same input.  
 E.g. predict type of animal and weather on a given pic  
 cat/dog/rat snake      cloudy/sunny/raining

Concatenate both dataset (182) dataset 1 dataset 2  
 - For feed-based systems: Diff. tasks like comment share on the same

## HOW TO COMPILE TRAINING DATASET FOR MULTI-TASK LEARNING

Usually employed in scenarios where multiple labels for multiple tasks on the same input. E.g. predict animal and weather on an image  
 cat rat bat lion cloudy sunny rain

- Training Dataset:

① Animal Dataset (or better) + Weather Dataset (concatenated)

(note: animal dataset images should not show weather and vice-versa) but at prediction time, prediction needs to happen for both tasks

Animal Dataset example 1  
 Animal Dataset example 2  
 Weather Dataset example 1  
 Weather Dataset example 2

② Feed based systems: Like, Comment, Share, etc. on a post

|         | Like | Comment | Share |
|---------|------|---------|-------|
| Post #1 | 1    | 2       | 0     |
| Post #2 | 0    | 0       | 1     |

(similar to multi-label dataset)

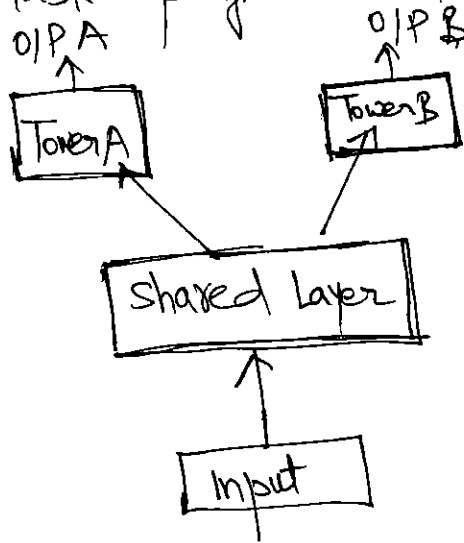
(Can have multiple reactions on a post)

# MULTI-TASK LEARNING (AND MULTI-TOWER)

## SHARED BOTTOM MODEL:

→ Same/shared base

→ Task specific head



→ Issue: If the tasks are not co-related then the performance is sub-opt as the shared layer with common representation introduces noise

## MIXTURE OF EXPERTS (MoE):

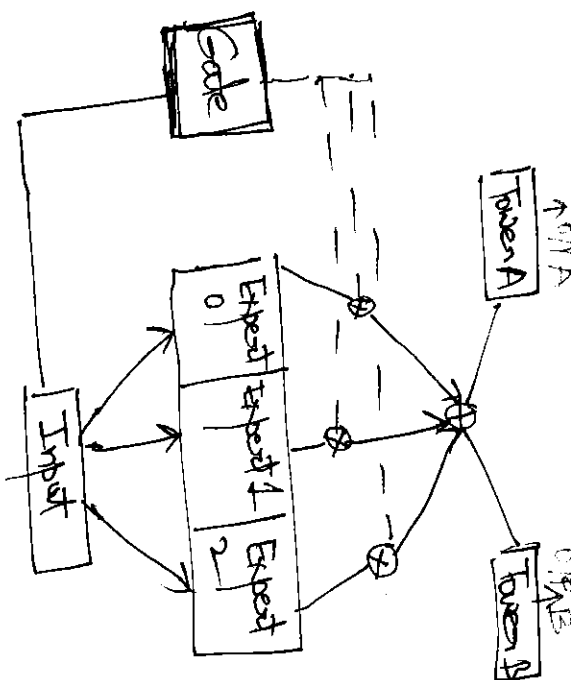
→ solves the issue ~~of~~ shared bottom model i.e. performs better even if the tasks are not co-related

→ ~~Multiple expert n/w~~ + a gating mechanism to weight ~~each expert's~~ <sup>o/p</sup>

→ Each expert in the n/w is able to learn different patterns in data and focus on diff. things

→ The gating n/w acts as a weighting scheme and takes the wt. avg. of expert n/w conditioned on input data i.e.

the model is able to activate parts of n/w differently on a per sample basis.



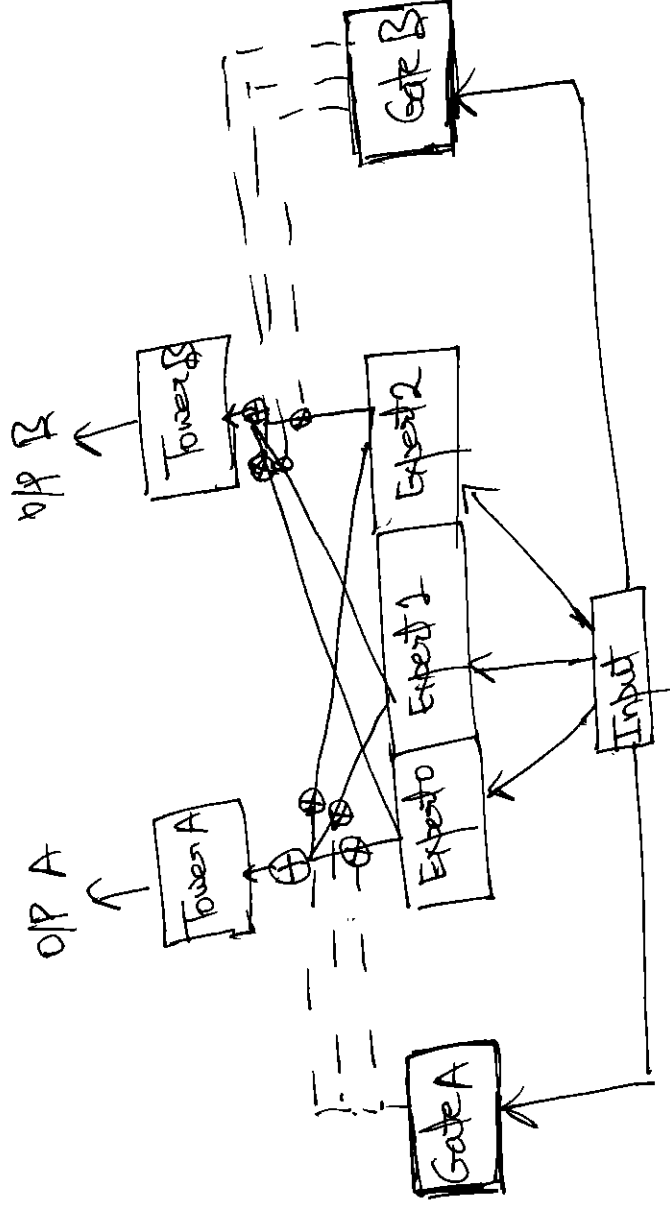
MoE

## MULTI-GATE MIXTURE OF EXPERTS (MMoE):

- Performs even if the tasks are not correlated and better than MoE
- MMoE is diff. than MoE : MMoE has a gating m/w adaptation of MoE to multi-task learning. for each task

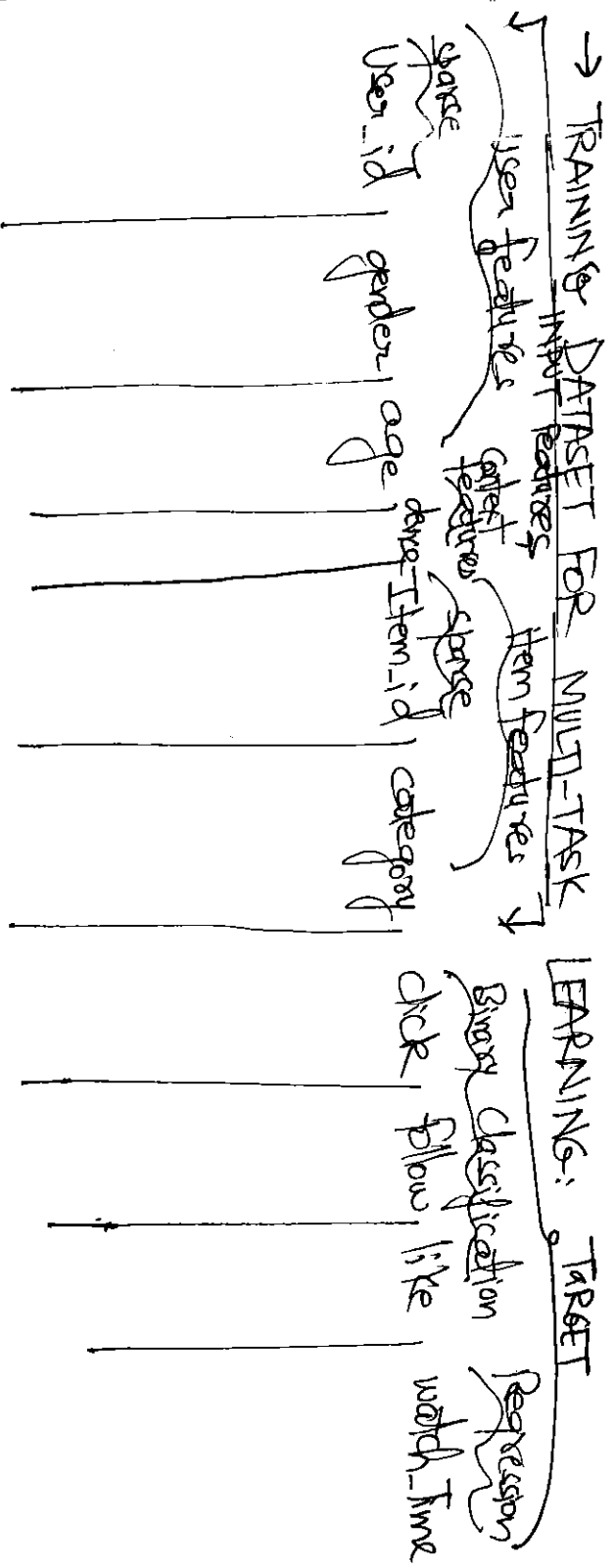
- This allows the model to learn per task and per sample weighting of each expert m/w (instead of just per sample weighting like MoE)
- MMoE essentially learns to model relationships b/w different tasks i.e. tasks which have little in common with each other will result in the gating m/w of each task learning to use diff. expert m/w.
- Experts are shared across all tasks and the gating m/w is task-specific. Each task also has a task-specific "tower" to decouple the optimization for tasks





MMoE

Note: Expert N/w is any DNN architecture



→ Coos-features b/w user-item not needed

→ user features (dense + sparse) ✓

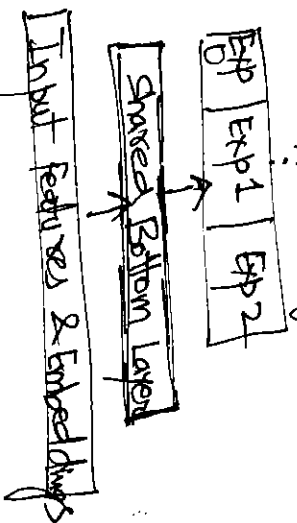
Item features (dense + sparse) ✓

Context features ✓

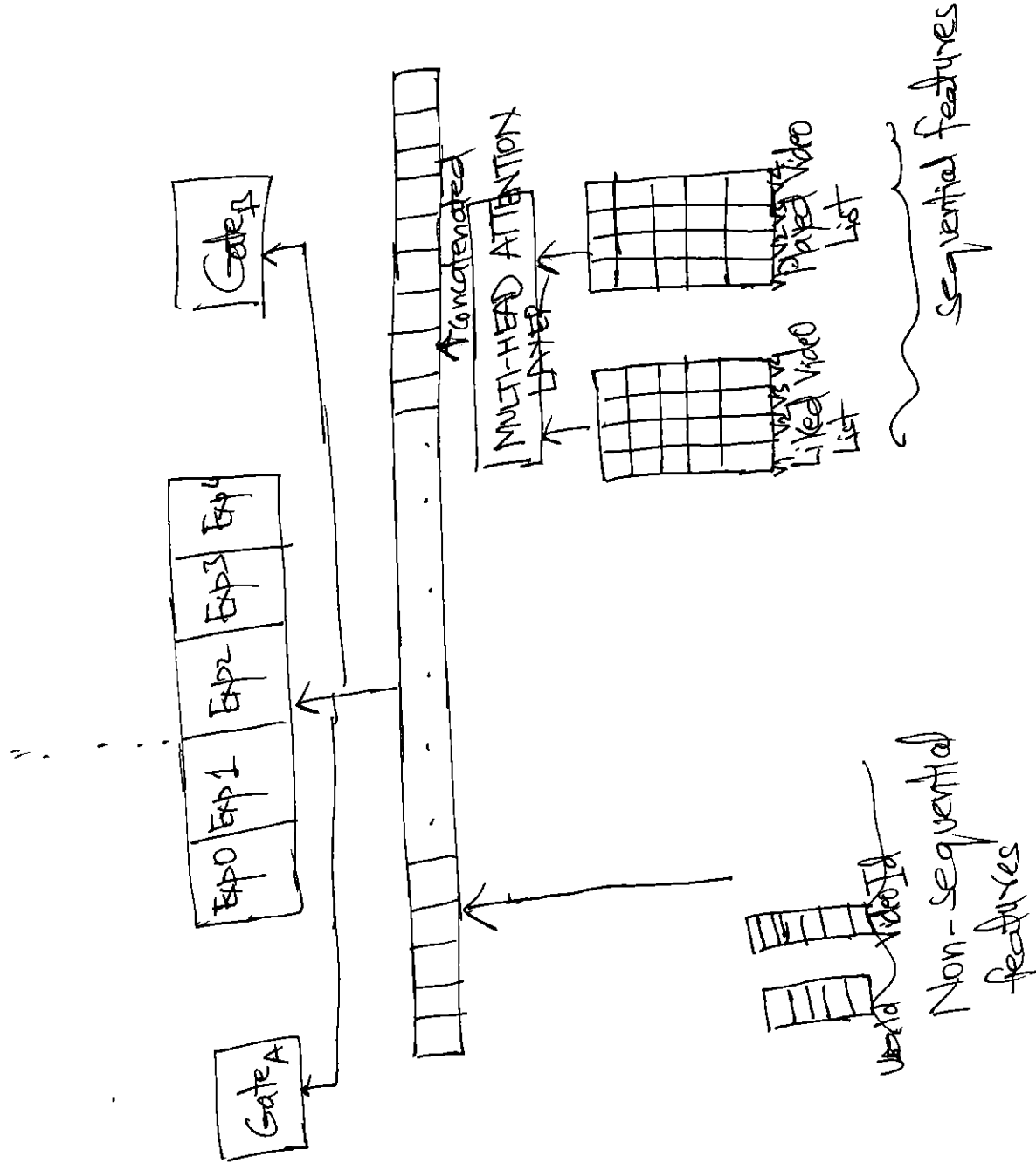
→ Target can be a combination of classification & regression

→ Loss FN =  $\sum$  (cross-entropy loss for each task (or MSE loss))

→ In practice, the lowest layer in MMoE is kept as 'shared' bcz the high dimensionality of I/P could lead to significant model training & serving costs.



# How to INGEST SEQ FEATURES INTO MMbE ARCHITECTURE



→ For each sequential feature, a multi-head attention module is used to fuse the embeddings into a concatenation of vector ops

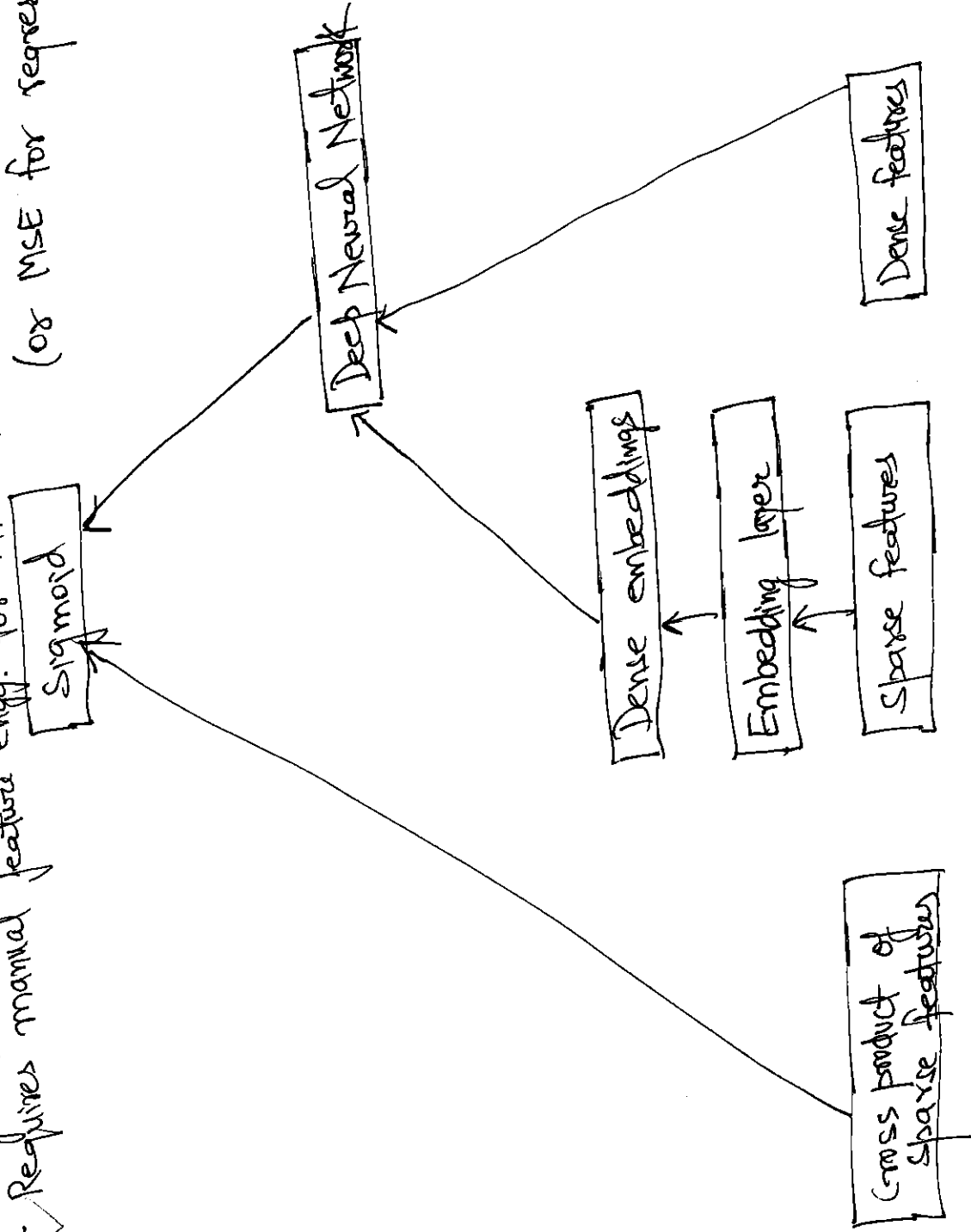
→ These sequential and non-sequential vectors are then concatenated and fed into MMbE layer



# WIDE AND DEEP MODEL (Google Play)

- Wide part can "memorize" seen feature interactions using cross-product feature transformations (need to manually hand craft these cross-features)  
unlike DeepFM or DON

- Deep part can "generalize" to previously unseen feature interactions
  - Jointly trains wide linear model and Deep NN
  - Useful if dataset has several sparse features
  - Requires manual feature engg. for interaction terms (or MSE for regression)





## FACTORIZATION MACHINES

- Extension of linear model that is designed to capture interactions between features within high-dimensional sparse datasets

- Mathematically,

$$\hat{y} = w_0 + \sum_i w_i x_i + \sum_{i > j} \langle v_i, v_j \rangle x_i x_j$$

global bias      linear term      factorization term

→ factorization term that models pairwise interaction between  $i^{\text{th}}$  and  $j^{\text{th}}$  variable (dot product of  $i^{\text{th}}$  and  $j^{\text{th}}$  feature embeddings)

- Requires no manual feature-engg for interaction terms (unlike wide & Deep Model)

- Can be used for classification as well as regression

- Pros: 1) The interaction between features can be estimated even under high sparsity

2) The number of params as well as time for prediction and learning is linear

3) Diminishes cold start of user or item i.e. if you know either some user features or some item features or some auxiliary features, FM can kick start recommendations (unlike Matrix Factorization where you need user interaction for items)

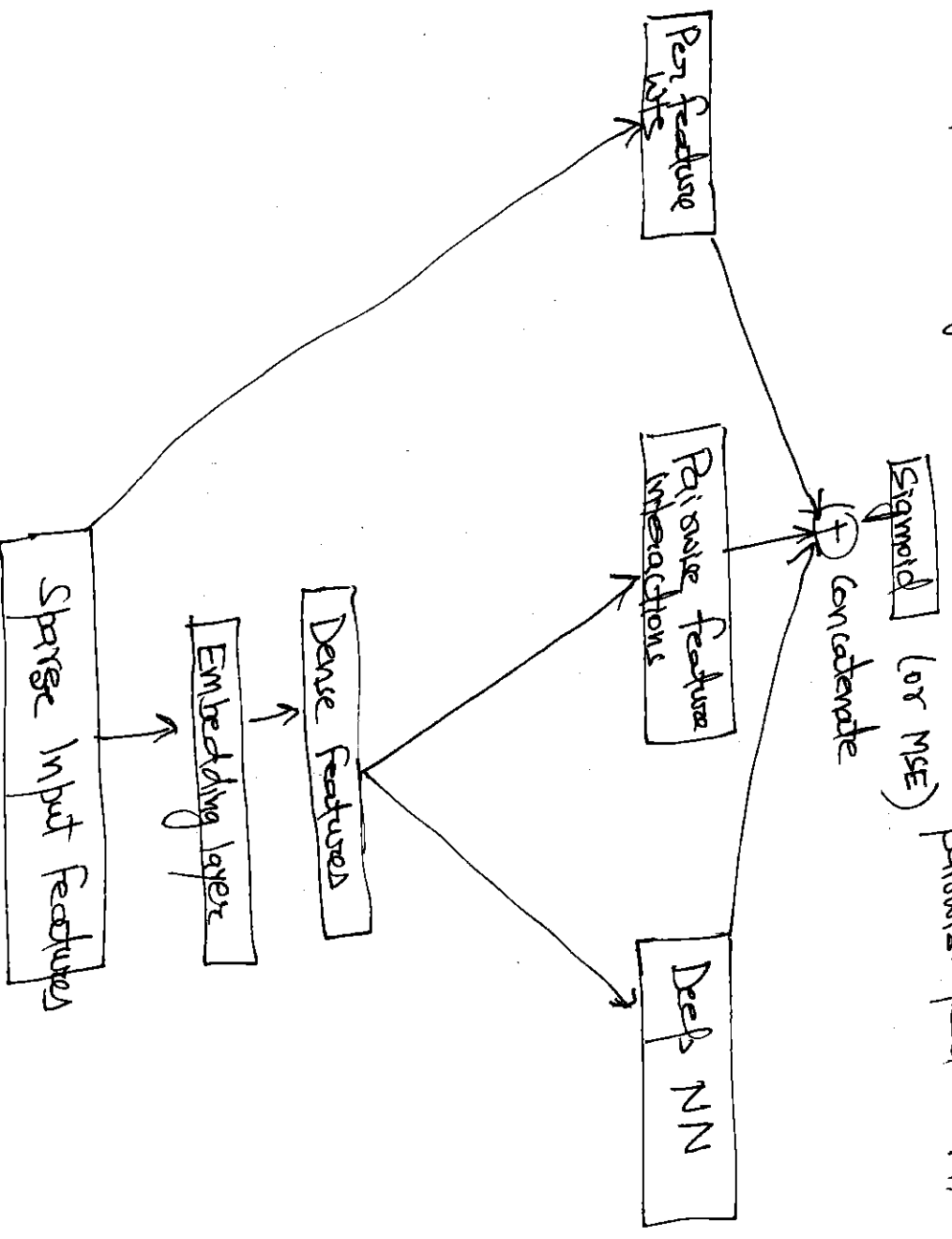
CONS: Only captures low level pairwise feature interactions (most implementations only support pairwise interactions e.g. Sagemaker's FM)

- Dataset looks something like this

| User features |   |   | Item features |   |   | Context features |   |   | Pairwise features |   |   | Label (e.g. rating or binary label) |   |
|---------------|---|---|---------------|---|---|------------------|---|---|-------------------|---|---|-------------------------------------|---|
| 1             | 0 | 0 | 0             | 1 | 1 | 1                | 0 | 0 | 1                 | 1 | 0 | 0                                   | 5 |
| 0             | 0 | 1 | 0             | 1 | 0 | 0                | 0 | 0 | 0                 | 0 | 1 | 1                                   | 4 |

## DEEP FACTORIZATION MACHINES (DeepFM)

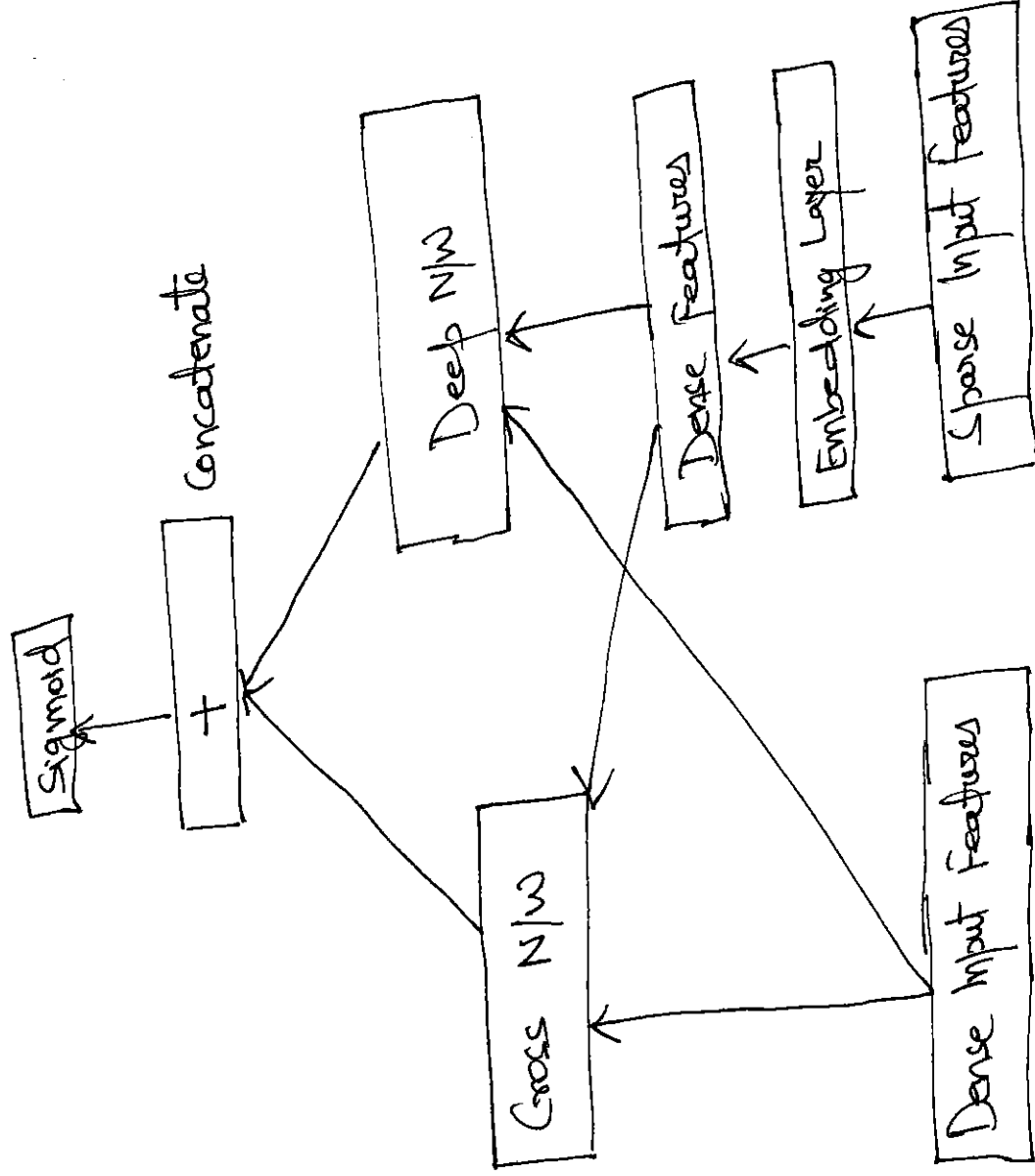
To capture higher order features along with FMs low level pairwise feature interactions





# DEEP AND CROSS NETWORK (DCN)

- Improvement over Wide & Deep Network
  - no manual feature-engg required to find feature interactions
- Diff from DeepFM:
  - The polynomial degree of interaction term increases at each layer and is determined by layer depth (Generalization of FM to higher degree)
- Cross N/W: Learns feature crosses
- Deep N/W: Learns complex and generalizable features

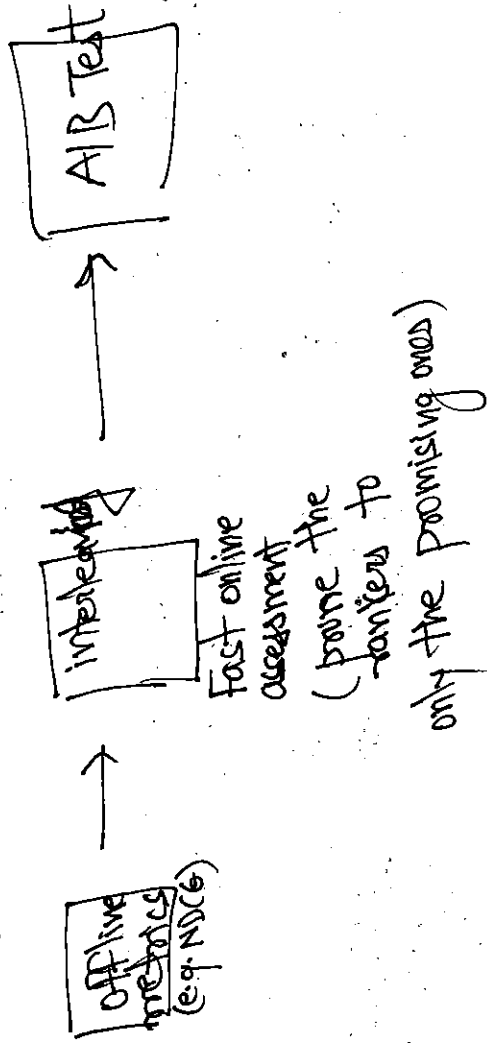




## ONLINE EVALUATION OF RANKERS - INTERLEAVING

- Interleaving speeds up the process of determining which ranker is better since it requires smaller sample size (same user is presented interleaved result from multiple rankers)

- Usually, the following process is used



WHY

- i) Interleaving is highly sensitive to ranking algo quality i.e. it reliably identifies the best algo with considerably smaller sample size compared to traditional A/B testing
- ii) Interleaving is predictive of success in second stage i.e. A/B test

# - APPROACH (Team Drafting and Competitive Pairs)

How  
Ranker 1  
(ranked list)

|   |
|---|
| A |
| B |
| C |
| D |

(Team A)

Ranker 2  
(ranked list)

|   |
|---|
| X |
| B |
| A |
| Y |

(Team B)

Flip a coin  
to determine  
which ranker  
goes first,  
then alternate b/w  
both rankers

A } competitive  
X } Pair  
B  
C  
Y } competitive  
D } Pair

i) From both rankers, only select the highest ranked result that has not been selected yet

ii) If the items in both list are different  
→ create a competitive pair

(metrics will be  
calc. only on competitive  
pairs)

iii) If the items in both list are same, then select it but do not assign to either team (metrics on this item will not be taken into account to determine which ranker is better)

Metrics such as viewing hrs, clicks etc. are calculated for both teams from competitive pairs and statistical tests (such as t-test) are used to determine verdict

## - NOTEWORTHY POINTS: BENEFITS

1) Flipping a coin to determine which ranker goes first and then alternating the turns  $\rightarrow$  reduction of position bias since both rankers have equal prob. of going first

2) Only calc. metrics on competitive pairs and leaving non-competitive items in list improves sensitivity in determining better ranker. E.g. if 2 rankers produce the exact ~~ranklist~~ list, the entire search query can be ignored since the preference is exactly zero. In traditional interleaving would still associate clicks to teams and add noise to the result

3) Same user sees results from both rankers hence variation in user characteristics is minimized compared to traditional A/B testing where one grp of users see variant A and another grp of users (with possible diff. characteristics) see variant B.

4) Needs engineering infrastructure to run this setup, if not already present



## A/B TESTING

### EXPERIMENTAL DESIGN:

- ① → Metric: Pick a metric to test along with guardrail metrics
- ② → Decide on significance level ( $\alpha$ ), power threshold ( $1 - \beta$ ) and Minimum Detectable Effect (MDE)

Significance level ( $\alpha$ ) [or  $(1 - \alpha)$  confidence level]: Percent of the time a difference will be detected, assuming one does not exist (Type I or false +ve error). Usually set to 0.05

Power: Percent of time minimum effect size will be detected assuming it exists.

$\beta$  → Type II error or false -ve error  
 $1 - \beta$  → Power (prob. of not making a type II error)  
Usually set to 0.8

Minimum Detectable Effect (MDE)

- Represents relative minimum improvement over the baseline that you are willing to detect in an experiment to a certain degree of statistical significance
- It is the smallest effect that will be detected  $(1 - \beta)\%$  of the time
- The smaller your MDE is, the larger the sample size required to reach statistical significance
- MDE is decided in consultation with stakeholders (e.g. 5%).
- The smaller your baseline is, the larger the sample size required to detect the same MDE

Note: effect = metric delta b/w control & treatment

### ③ Stopping Cond. for A/B Test: Sample size & Experiment Duration

#### SAMPLE SIZE:

- Depends on a) baseline metric, b) statistical significance level ( $\alpha$ ), c) Statistical power ( $1-\beta$ ) and d) MDE
- Sample size calculators can determine sample size needed per variation

#### EXPERIMENT LENGTH:

- Depends on daily traffic ( $\frac{\text{total traffic for the feature under test}}{\text{daily traffic for the feature under test}}$ )
- Should account for differences in usage pattern e.g. day of week, holidays etc. (usually set at 2 weeks)

### ④ Assign Expts:

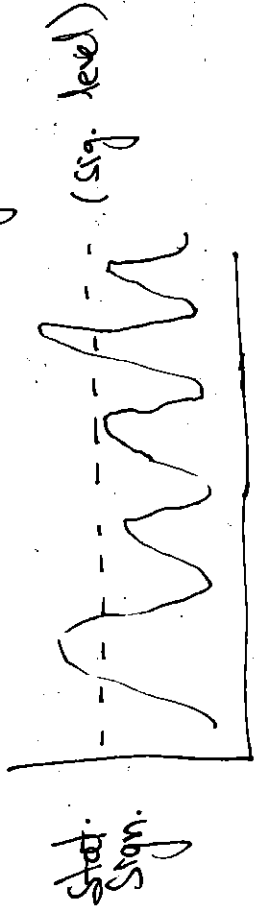
Randomize control & treatment sufficiently to account for confounding variable (doon the line (e.g. day of week, season, word of mouth publicity, press (+ve or -ve)))



## A/B TESTING MISTAKES

→ Stopping A/B Tests Early:

Reaching statistical significance alone is not enough to declare a winner. bc stat. sig. varies for the duration of test



Duration of test

Declare a winner only if:

- 1) Sample size requirement is met
- 2) Duration of test covered variations in day of week, holidays (etc) i.e. usually one business cycle (usually 2 weeks)

→ Running multiple tests at the same time on overlapping traffic

If you run multiple tests at the same time, some tests will succeed with statistical significance much larger than  $\alpha$

Family-wise Error Rate (FWER)

$$\text{FWER or } \alpha = 0.05$$

Prob. of not committing Type I error:  $1 - 0.05$

Prob. of not committing Type I error for 10 tests =  $(1 - 0.05)^{10}$

Prob. of committing Type I error at least once for 10 tests

$$= 1 - (1 - 0.05)^{10}$$

$$= 0.40$$

$$= 40\%$$

## Sol<sup>n</sup> :- 1) Bonferroni correction (BC)

Adjusts  $\alpha$  level required based on total number of tests running  
i.e.  $\alpha/n$  where  $n$  = no. of tests running

E.g. we conduct multiple comparison of 10 metrics  
By default  $\alpha \pm 0.05$  for each test

Using Bonferroni correction new  $\alpha = \frac{\text{old } \alpha}{n}$

$$= \frac{0.05}{10} = 0.005$$

Now, if we want to claim any statistically significant results, the respective p-value has to be smaller than 0.005

As seen, we have lowered the rejection bar from 0.05 to 0.005, making it much harder to reject Null hypothesis.

This method is criticized for being too conservative & lacking statistical power

## 2) Holm Bonferroni (HB)

Due to conservative nature of BC method, rarely used in practice

At controls FWER (family-wise error rate) by

- i) ranking the p-values from each test
- ii) adjusting rejection criteria for each null hypothesis

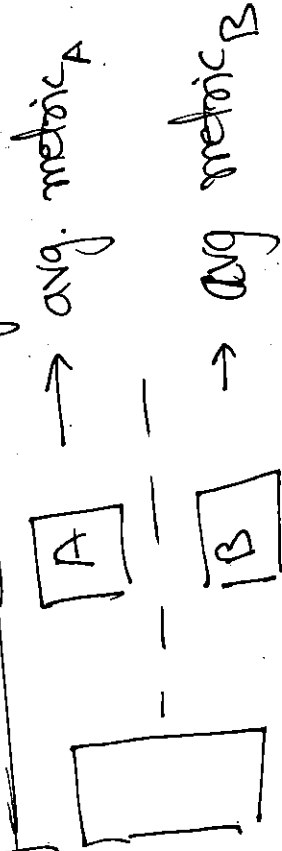
Diff b/w BC & HB: HB's rejection threshold depends on all p-values (unlike BC method)

HB has more statistical power than BC method

### 3) FDR (False Discovery Rate)

- Proportion of false ~~res~~ among all significant results
- Eg. if you have made 100 significant results and set FDR to 0.05 then 5 of them are false discoveries (or false rejection)

→ Dealing with non-normality



AB tests are dressed up Z-tests and t-tests and these tests assume the var. of interest (metric we are testing for) is normally distributed.

This assumption is usually true due to large sample size + Central Limit Theorem.

But if the normality assumption of metric does not hold true then we have to use other techniques:

- Running alternate test: Wilcoxon rank-sum test does not assume normality

- Gather more data to invoke CLT

## → Dealing with Novelty Effect / Primary Effects,

A/B tests may have an exaggerated initial effect due to novelty effect where a new feature attracts social media attention and PR hype causing inquisitive users to rush to checkout new changes and thereby leading metric to jump spuriously

In the opposite way, there could be primary effect where users are fixed on what is familiar and have an aversion to changes making test metrics to be stagnant / low initially.

To detect novelty or primary effect, slice & dice the A/B test results on new users only

## → Dealing with network effects

Generally assumed that each user in control & treatment groups are independent i.e. no interference b/w control & treatment. However, this assumption does not hold true all the time.

Example: in social n/ws - the action of users are likely impacted by those around them (who may be in other groups) (control/treatment)

Facebook like: treatment group users may find it entertaining & post about it & send links to friends (who may start using the feature) → can lead to increased engagement in control group

To solve for it → create clusters of connected people and assign them as control or treatment group

→ Nuances of taking action on A/B Test Results:

- sample size & duration of test ✓
- Statistical significance reached ✓
- No significant drop in guardrail metrics ✓

Also, time & effort to build & scale feature maintain < business impact of feature

→ A/B Test Universal Holdouts

Useful in quantifying the business impact of entire org

E.g.

Team A → 5% ↑ in conversion

B → 8% ↑ in conversion

...

N → N% ↑ in conversion

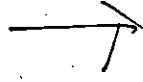
Universal Holdout: No tests are run for the grp of users

Gap of users part of any test



conversion

Universal Holdout



conversion

$\Delta$  = o/p of entire org  
in conversion metrics

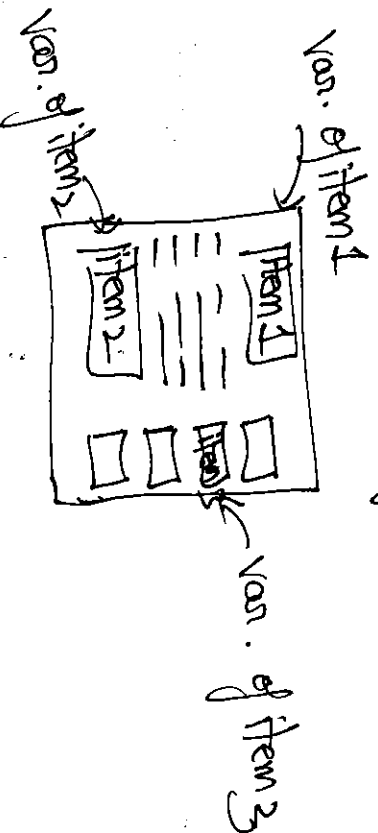
## BANDIT BASED A/B TESTING: (MAB)

- Bandit algorithm starts sending ~~meta~~ traffic to "winning" variant to minimize "regret"
- Useful in scenarios where you can't wait to finish regular A/B test to determine decrease in metric or the metric ↓ has a huge cost
- Explore vs Exploit  
↓  
reg. traffic allocation to "winning" variant with pt. in time
- Significance bel may not be reached

## MULTIVARIATE TESTS

these are (not A/B tests)

- Test multiple versions of ~~diff elements~~ <sup>same page</sup> to isolate which attributes cause the largest impact



- Ton of traffic needed to perform multivariate test
- They test original against variation but each variation contains diff. design elements (e.g. which combination of item 1, 2 and 3 variation ↑ metric) [figures out interaction effect]

## MONITORING

\* - COST  
INFRASTRUCTURE: - Amount of CPU/GPU

- Memory
- Disk I/O
- N/w I/O
- Throughput (Mbps) (no. of concurrent users)
- Latency (client  $\leftrightarrow$  server)

## PREDICTION/SCORING PIPELINE:

- Kafka  $\rightarrow$  Flink consumer lag
- Real-time processing pipeline (Flink): Error handling Timeout
- Populating into feature store: Distribution drift b/w training & prediction data (data drift) training-serving skew
- Featureization code similar to training pipeline
- Outages in pipelines
- Upstream schema changes in data stores

MODEL: - Model performance metric change: precision, recall (if label is available in an acceptable time frame)

- Model prediction metrics (i.e. flag 99% of users) (if label not available in an acceptable time frame)

TRAINING PIPELINE: Base table schema changes / Unavailability (e.g. using DDL to monitor)

MODEL TRAINING MONITOR: - Vanishing gradient, overfitting, poor initialization, dead-rell, exploding tensor, loss not decreasing (E.g. Scikitlearn Debugger)

- Data leakage

## RE-TRAINING FREQUENCY

1. Periodic (e.g. weekly)
2. When model metrics degrade

## REDUCING LATENCY (WHILE SCORING)

1. PRE-COMPUTING FEATURES, EMBEDDINGS
2. CACHING FEATURES, EMBEDDINGS AND PREDICTIONS  
if used by several downstream models
3. MODEL COMPRESSION :- Quantization  
(reduce ints and activations from float32 to int 8)
  - Knowledge distillation
  - Pruning weights and connections that contribute least to model performance

4. USE SIMPLER FEATURE ENGINEERING
5. USE MORE HARDWARE / INFRASTRUCTURE RESOURCES FOR INFERENCE PIPELINE

## SCALING SCORING

- Docker Containers + Kubernetes, Load balancers + Auto-scaling (e.g.  $\uparrow$  Noisy request threshold)



## REDUCING TRAINING TIME:

1. Distributed Training
2. Large Batch sizes (reduce the no. of times parameters are updated)
3. Utilize GPU acceleration: GPUs are optimized for parallel processing, have large memory bandwidth and are designed specifically for matrix calc.
4. Downsample training data
5. Use optimized implementation of algorithms e.g. tensorflow, pytorch
6. Optimize learning rate schedule: fine tune learning rate to balance speed of convergence and stability of optimization losses



## MODEL CALIBRATION: (Especially useful in multi stage recommendation system)

- Necessary when

- 1) Comparing experimental vs production model (or any two versions of model)
- 2) Model O/P is consumed by a downstream model (or subsequent models in multi-stage settings)

- Example: Consider 2 models; (e.g. for CTR prediction)

- 1<sup>st</sup> model : score varies b/w 0 to 0.2

- 2<sup>nd</sup> model : score varies b/w 0 to 1

Both models have same test set performance (AUC)

bcz their relative ordering of predictions is the same

Their scores are not directly comparable, a score of 0.2 means a very high CTR for the first model and relatively low CTR for the second model

If production and experimental models are not calibrated to the same distribution then downstream systems (e.g. ranking or re-ranking functions) that rely on model score distributions will need to be re-tuned (unsustainable in the long run)

If calibrated then model experiments and launches can be decoupled from the changes to the rest of recommendation system

CALIBRATION: Connecting model scores to realworld probability in the range  $[0, 1]$

To all the samples that a classifier gave 0.8 (pred. prob) appear 80% actually belong to true class  
→ well-calibrated binary classifier

Calibrating a classifier consists of fitting a regressor that maps O/P of classifier → calibrated prob in  $[0, 1]$

Two methods

— Platt Scaling: fitting a sigmoid fn  
— Isotonic regression: fitting a discrete, stepwise, monotonically inc. fn

(We know O/P of uncalibrated classifier & true label hence can fit a fn)

# PROS AND CONS OF ML ALGOS

## LOGISTIC REGRESSION:

Pros:

- Fast training
- Fast inference
- Interpretable
- Can be updated thru continuous learning (online learning)

Cons:

- Non-linear problems can't be solved
- Inability to capture feature interactions

## DECISION TREES:

Pros:

- Fast training
- Fast inference
- Little to no data prep : Normalizing x  
Scaling x  
Missing values x
- Interpretable

Cons:

- Overfitting

# RANDOM FOREST (BAGGING + "random" selection of features)

Pros:

- Does not significantly increase training time bcz DT can be trained in parallel
- Does not add much latency at inference time bcz DT can process input in parallel
- Reduces ~~over~~overfitting (high variance)
- Interpretable (feature importance  $\checkmark$ , coeff.  $\times$ )
- Non-linear decision boundary
- Online learning possible (Mondrian Forests)
- Easy data grab.

Cons:

- Not helpful in high bias scenarios

## BOOSTING (GBDT, XGBoost)

Pros:

- Though it builds trees sequentially, training time can be reduced with efficient implementations like XGBoost (builds multiple branches of tree in parallel) + distributed training possible
- Implementations available for fast inference (catboost, lightGBM)
- Reduces both bias & variance
- Easy data prep: Normalization  $\times$   
Scaling  $\times$   
Missing values  $\times$
- Works well with structured data

CONS:

- Unsuited for continuous learning from streaming data
- Does not work well on unstructured data
- Lots of hyperparams to tune

## Neural NETWORKS

Pros

- Continual learning possible
- Works well with unstructured data + embeddings
- Can learn very complex tasks and non-linear decision boundaries

CONS:

- Large training data is required
- Computationally expensive to train
- Some effort in data prep. Eg. if input features are in very different ranges, the model  
↓  
normalization,  
log-scaling, one-hot encoding, may converge very slowly.  
embedding
- ~~Black~~ - box nature in interpretability





## DELAYED LABELS

Issue → model becomes stale  
increases ~~if~~ false -ves

### SOLUTION:

- i) Proxy label: If there is some info/feature that co-relates strongly with label and is available before
- ii) Re-formulate the ML problem with a time horizon instead of whether fraud will be committed by a particular user?  
reformulate it as whether fraud will be committed in the next 30 days by this user?

### iii) Supervised Anomaly Detection:

Using historical data of bad actors,  
used anomaly detection algo like  
one-class SVM (trained only on one class)  
and essentially learns the decision boundary/distribution of one class

### iv) Unsupervised: Auto-Encoders

Re-construction loss above a threshold

### v) Domain Experts: craft/annotate labels

