

# Two Tower Model with Self-Supervision Learning

Friday, March 15, 2024

3:50 PM

## Two Tower Model Issues :

- ① A few popular items takes most clicks
- ② Large number of items get a few clicks.
- ③ The representation vector of popular items is learned well  
However, the vector of unpopular items is poorly learned.  
Take large fraction

How to solve? data augmentation

## Two Tower Model Training :

Batch Examples :

users	items
#1	#1
#2	#2
⋮	⋮
#N	#N

Positive Samples :

$(1, 1), (2, 2), \dots, (N, N)$

Negative Samples :

$(1, 2), (1, 3) \dots (1, N)$

$(2, 1), (2, 3) \dots (2, N)$

$(N, 1), (N, 2) \dots (N, N-1)$

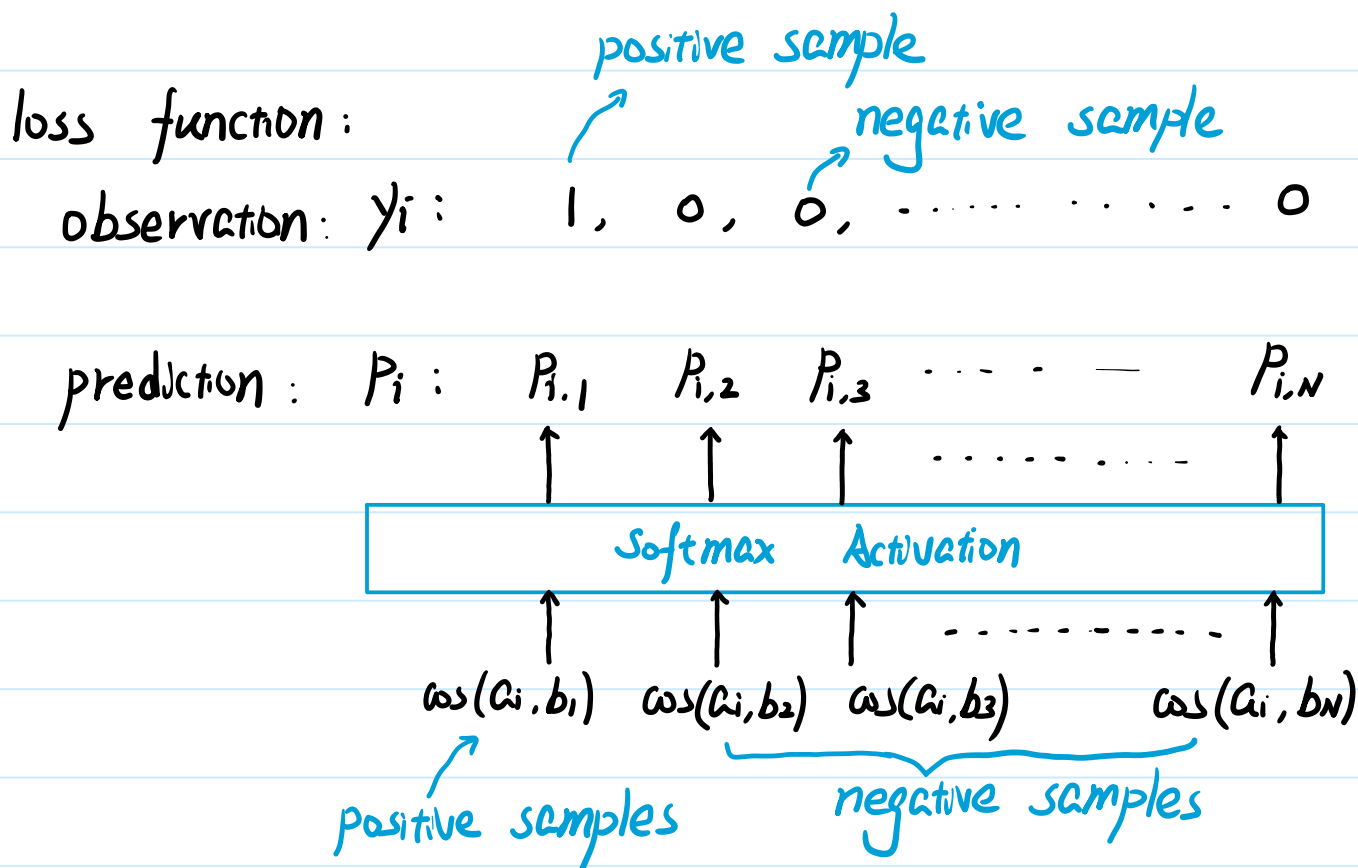
# of positive samples :  $N$       # of negative :  $N(N-1)$

objective: make  $\cos(a_i, b_+)$  large ; make  $\cos(a_i, b_-)$  small

# Two Tower Model with Self-Supervision Learning

Friday, March 15, 2024

3:50 PM



Cross Entropy Loss ( $y_i, P_i$ )

$$L = -\log P_{i,i} = -\log \frac{\exp(\cos(a_i, b_j))}{\sum_{j=1}^n \exp(\cos(a_i, b_j))}$$

Probability of item  $j$  being selected:

$$P_j \propto (\# \text{ of click})$$

Estimated like from user  $i$  to item  $j$ :  $\cos(a_i, b_j)$

Corrected like:  $\cos(a_i, b_j) - \log P_j$

# Two Tower Model with Self-Supervision Learning

Friday, March 15, 2024

3:50 PM

## Training Process:

① Randomly select  $N$  users and their interacted items  $(a_1, b_1), (a_2, b_2), \dots, (a_N, b_N) \rightarrow$  one batch

② Loss function:

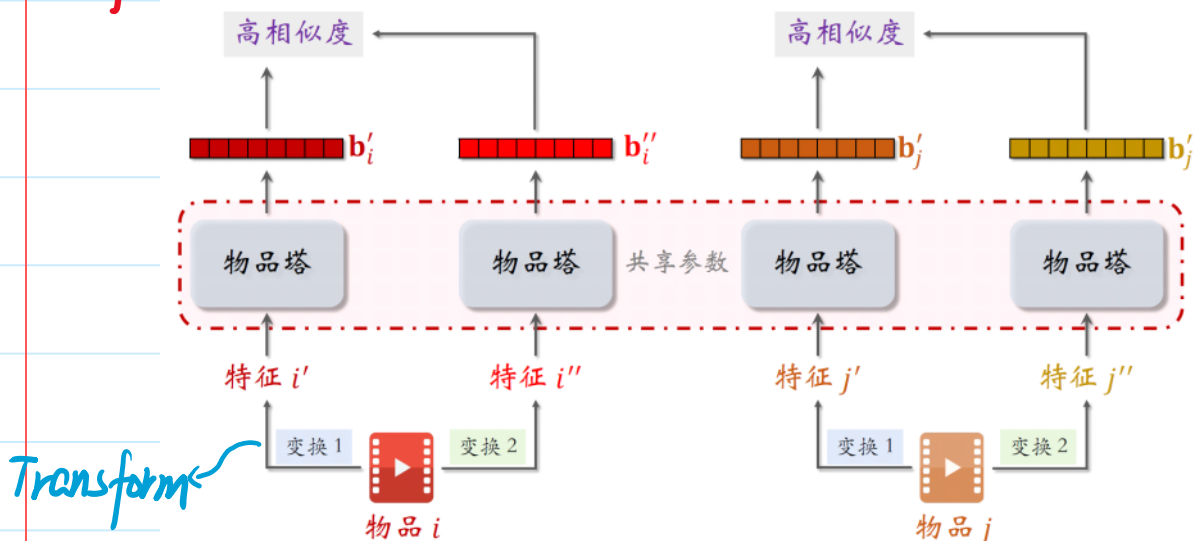
$$L_{\text{main}}[i] = -\log \frac{\exp[\omega(a_i, b_i) - \log P_i]}{\sum_{j=1}^n \exp[\omega(a_i, b_j) - \log P_j]}$$

for user  $i$

③ gradient descent:

$$\frac{1}{N} \sum_{i=1}^N L_{\text{main}}[i]$$

## Self-supervision:



(Picture from Shusen Wang on Youtube/Bilibili)

① make  $\cos(b'_i, b''_i)$  and  $\cos(b'_j, b''_j)$  large

② make  $\cos(b'_i, b'_j)$  and  $\cos(b'_i, b''_j)$  small

same item: high similarity; different item: low similarity



# Two Tower Model with Self-Supervision Learning

Friday, March 15, 2024

3:50 PM

## Transform #1: Random mask

① select some categorical features and mask them

e.g. item #1: genre = (photo, beauty, travel)  
→ mask: genre = ("missing")

## Transform #2: Drop out.

① randomly drop part of features

e.g. item # genre = (photo, beauty, travel)  
→ drop: genre = (photo)

## Transform #3: Complementary:

e.g. items: (ID, genre, hashtag, city)

two groups { (ID, hashtag)  
(genre, city)

each item has two representation:

(ID, "missing", hashtag, "missing") → item tower

("missing", genre, "missing", city) → item tower

vector representation similar

# Two Tower Model with Self-Supervision Learning

Friday, March 15, 2024

3:50 PM

## Transform #4: mask correlated features.

- ① gender:  $U = \{\text{male, female, neutral}\}$
- ② genre:  $V = \{\text{cosmetics, digit, football, tech}\}$
- ③  $u = \text{male}$  and  $v = \text{digit}$  probability is large
- ④  $u = \text{male}$  and  $v = \text{cosmetics}$  probability is small.

mutual information:

$$MI(U, V) = \sum_{u \in U} \sum_{v \in V} p(u, v) \cdot \log \frac{p(u, v)}{p(u) \cdot p(v)}$$

$k$  features  $\rightarrow k \times k$  matrix for MI

- ① random select one feature as "seed feature"
- ② get most relevant  $k/2$  features
- ③ mask "seed" and  $k/2$  features

Pros: works best

Cons: implementation is difficult; hard to maintain

Overall: NOT effie

# Two Tower Model with Self-Supervision Learning

Friday, March 15, 2024

3:50 PM

## Model Training: for self supervision

- ① uniformly select  $m$  items (NOT based on # of clicks)
- ② perform two types of transformations  
item tower generates two groups of vectors  
 $(b_1', b_2', \dots, b_m')$   $(b_1'', b_2'', \dots, b_m'')$
- ③ loss function for item  $i$ :

$$L_{\text{self}}[i] = -\log \frac{\exp(\omega(b_i', b_i''))}{\sum_{j=1}^m \exp(\omega(b_i', b_j''))}$$

$i$  and  $j$  similarity

- ④ objective function:

$$\frac{1}{m} \sum_{i=1}^m L_{\text{self}}[i]$$

## Self supervision + Two Tower Model:

- ① random sample  $n$  (user, item) from clicked items
- ② uniform sample  $m$  items from all items
- ③ gradient descent:

$$\underbrace{\frac{1}{n} \sum_{i=1}^n L_{\text{main}}[i]}_{\text{two tower}} + \underbrace{\alpha}_{\text{hyper parameter}} \underbrace{\frac{1}{m} \sum_{j=1}^m L_{\text{self}}[j]}_{\text{self supervision}}$$