

1. Kd树原理介绍:

1.1. 构建: Kd是二叉搜索树的一种.

1.2. 构建方式: 每个节点按照某一个维度.

进行切分. 按照

维度的选择的方式一般

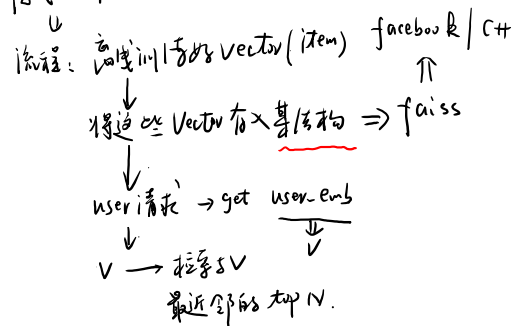
是看方差的大小.

选方差大的维度. 进行

对数据的切分. 因为区分度大.

1. 向量检索
2. YouTube 召回、YouTube Rank
3. 模型部署常用方法。○

1. 隐式召回/embedding 召回。



索引结构: faiss - kd-tree - Annoy - HNSW.

vector 的维度:

1. 较低时.

2. 远超过 30 维度.

1. id: 800 万个, dim: 200.

faiss: 30ms 以内.

向量检索的流程:

1. 离线训练 (idV, 双塔).

2. 双塔会已经有了. 可推荐的 item \rightarrow get-id \rightarrow get-特征

\rightarrow 过 item 相似的双塔 \rightarrow

get item 的 embedding.

\rightarrow 放入 faiss 做 train.

3. user 请求 \rightarrow get-user-embedding

(过 user 侧塔)

\rightarrow faiss 检索.

\rightarrow 业务侧的过筛.

} serving

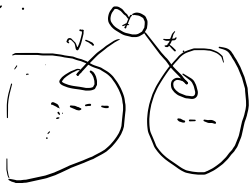
1. KNN, kd-tree.

kd-tree: 二叉树. 二叉搜索树.

[2, 3] → kd-tree

float: 2, 3 → 查找

$\log(N)$



- [1, 2, 3, 4, 5]
- [2, 3, 3, 4, 5]
- [2, 3, 3, 4, 5]
- [1, 2, 3, 4, 5]

1
2
3 查找?
4
5

eg. [(2, 3), (5, 4), (9, 6), (4, 7), (8, 1), (7, 2)]

