



Chapter 4: Finding Similar Items

-Locality Sensitive Hashing

崔金华

电子邮箱: jhcui@hust.edu.cn

个人主页: <https://csjhcui.github.io/>

Contents

提纲

4.1

集合相似度的应用

Applications of Set Similarity

4.2

文档的Shingling

Shingling of Documents

4.3

最小哈希

Minhashing

4.4

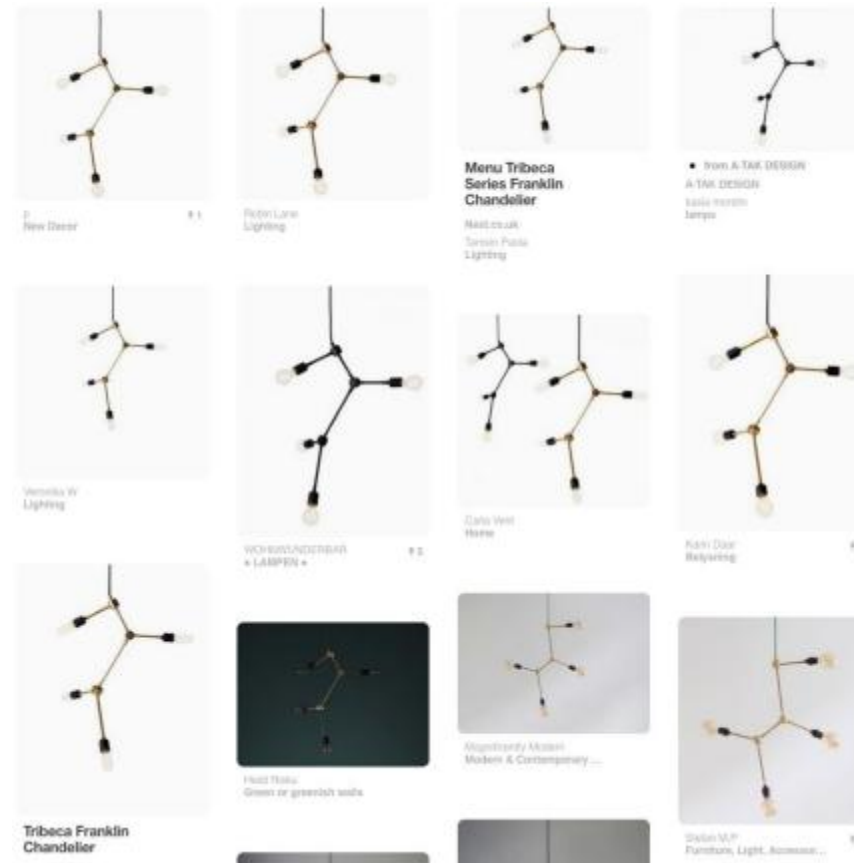
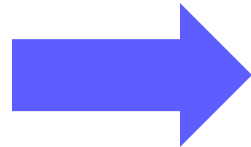
局部敏感哈希算法

Locality-Sensitive Hashing

4.1 Pinterest Visual Search

□ Given a query image patch, find similar images

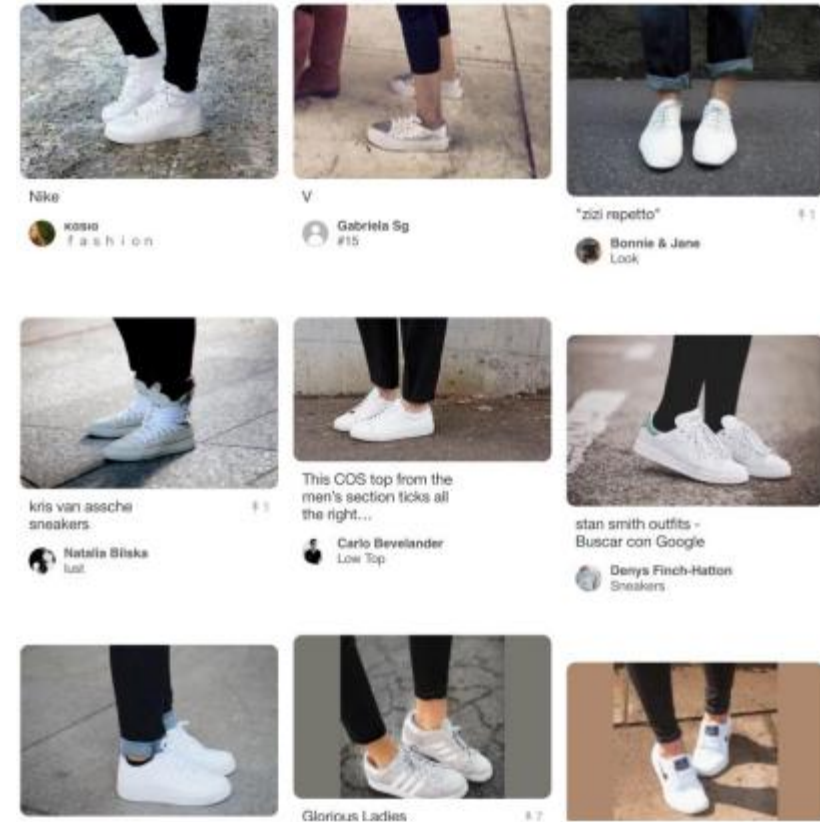
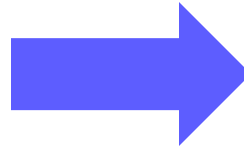
Visually similar results



4.1 Pinterest Visual Search

□ Given a query image patch, find similar images

Visually similar results



4.1 A Common Metaphor

□ Many problems can be expressed as finding “similar” sets:

➤ Find near-neighbors in high-dimensional space

□ Examples:

➤ Customers who purchased similar products

- Products with similar customer sets

➤ Images with similar features

- Users who visited similar websites

➤ Pages with similar words

- For duplicate detection, classification by topic

4.1 Task: Finding Similar Documents

□ **Goal:** Given a large number (N in the millions or billions) of documents, find “near duplicate” pairs

□ **Applications:**

- Mirror websites, or approximate mirrors
 - Don't want to show both in search results
- Similar news articles at many news sites
 - Cluster articles by “same story”

焦点新闻 >

法制网

中共中央政治局召开会议 中共中央总书记习近平主持会议

4 小时前

国际金融报

中共中央政治局召开会议

12 小时前

SOH_NEWS_CN

大风暴正来临！习近平强调「忠诚」李希首创巡视「三板斧」

13 小时前

瞭望东方周刊

联播+ | 如何开展这项主题教育？中央明确方向

10 小时前

完整报道

人民网

在中央党校建校90周年庆祝大会暨2023年春季学期开学典礼上的讲话--新闻报道-中国共产党新闻网

2 小时前



4.1 Task: Finding Similar Documents

□ Problems:

- Many small pieces of one document can appear out of order in another
- Too many documents to compare all pairs
- Documents are so large or so many that they cannot fit in main memory