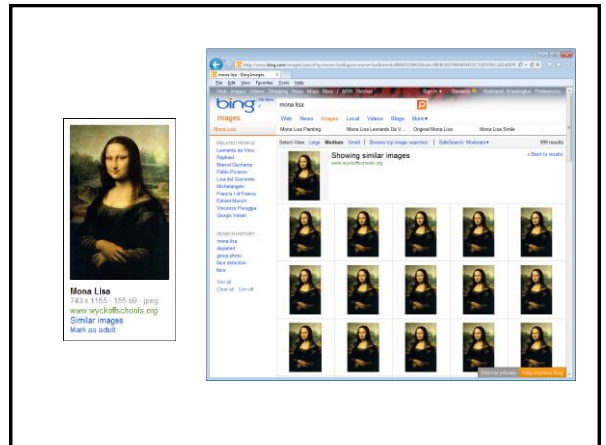
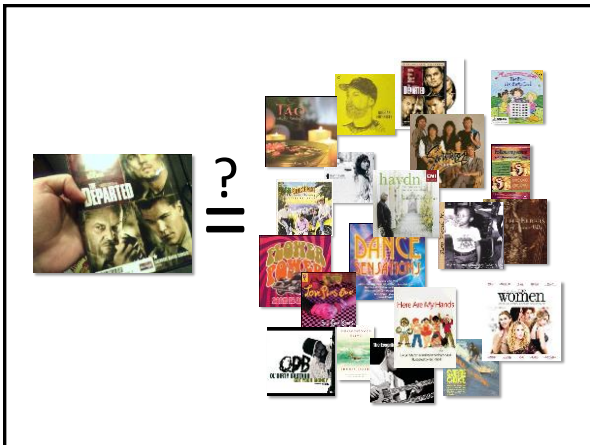


# Large-scale matching

CSE P 576

Larry Zitnick ([larryz@microsoft.com](mailto:larryz@microsoft.com))



## Large scale matching

How do we match millions or billions of images in under a second?

Is it even possible to store the information necessary?

1 image (640x480 jpg) = 100 kb  
 1 million images = 100 gigabytes  
 1 billion images = 100 terabytes  
 100 billion images = 10,000 terabytes (Flickr has 5 billion)

## Interest points

Currently, interest point techniques are the main method for scaling to large databases.



Find interest points



Extract patches

## Searching interest points

How do we find similar descriptors across images?

Nearest neighbor search:

Linear search:

1 million images x 1,000 descriptors = 1 billion descriptors (Too slow!)

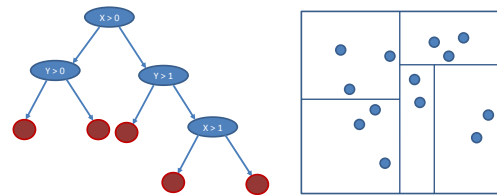
Instead use approximate nearest neighbor:

- KD-tree
- Locality sensitive hashing

## KD-tree

Short for "k-dimensional tree."

Creates a binary tree that splits the data along one dimension:



## KD-tree

Algorithm for creating tree:

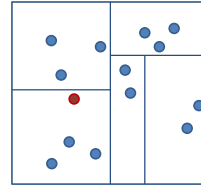
1. Find dimension with highest variance (sometimes cycle through dimensions).
2. Split at median.
3. Recursively repeat steps 1 and 2, until less than "n" data points exist in leaves.

Searching for approximate nearest neighbor:

1. Traverse down the tree until you reach the leaf node.
2. Linearly search for nearest neighbor among all data points in leaf node.

## Problem:

The nearest neighbor may not be in the "found" leaf:

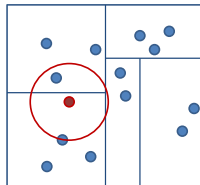


## Backtracking (or priority search)

The nearest neighbor may not be in the "found" leaf:

Backtrack to see if any decision boundaries are closer than your current "nearest neighbor."

In high dimensional space, the number of "backtracks" are typically limited to a fixed number. The closest decision boundaries are stored and sorted.



## High-dimensional space

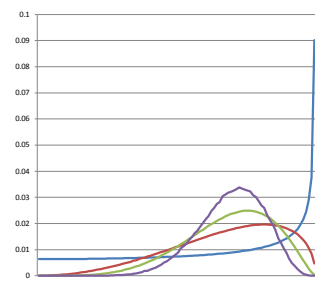
How far away are two random points on an n-dimensional sphere?

2D

4D

6D

10D



Don't follow your "2D intuitions"

## KD-tree

Other variations:

- Use Principal Component Analysis to align the principal axes of the data with the coordinate axes
- Use multiple randomized KD-trees (by rotating data points)

**Optimised KD-trees for fast image descriptor matching**  
Chanop Silpa-Anan Richard Hartley, CVPR 2008

## Storing the descriptors

Storing the descriptors is expensive:

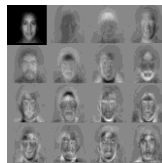
1000 descriptors x 128 dimensions x 1 byte =  
128,000 bytes per image

1 million images = 120 gigabytes

## PCA

Reduce the dimensionality of the descriptor using PCA.

Pick the "n" orthogonal dimensions with highest variance.



Eigenvectors of face images

Storing the descriptors is still expensive:

1000 descriptors x 32 dimensions x 1 byte =  
32,000 bytes per image

1 million images = 30 gigabytes

## Locality sensitive hashing

- Assume points are embedded in Euclidean space



- How to binarize so Hamming distance approximates Euclidean distance?

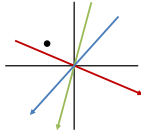
`Ham_Dist(10001010, 11101110) = 3`

1000 descriptors x 64 dimensions x 1 bit =  
64,000 bytes per image

1 million images = 7.5 gigabytes **We're in RAM!**

## Finding the binary code

- For each bit:
  - Compute a random unit vector "r".
  - For input vector "v", set the bit equal to:
 
$$h(v) = \text{sign}(v \cdot r)$$



The following holds:

$$Pr[h(u) = h(v)] = 1 - \frac{\theta(u, v)}{\pi}$$

$$1 - \frac{\theta(u, v)}{\pi} \text{ is closely related to } \cos(\theta(u, v)).$$

1 0 0

## P-stable distributions

- Projects onto an integer instead of bit:
  - Compute a random Gaussian vector "a".
  - For input vector "v", set the index equal to:

$$h_{a,b}(v) = \left\lfloor \frac{a \cdot v + b}{r} \right\rfloor$$

- "b" and "r" are chosen by hand.

Datar, M.; Immorlica, N., Indyk, P., Mirrokni, V.S. (2004). "Locality-Sensitive Hashing Scheme Based on p-Stable Distributions". *Proceedings of the Symposium on Computational Geometry*.

## Other methods

Spectral hashing (uses thresholded eigenvectors) to find binary codes:

$$h(v) = \text{sign}(\cos(kw \cdot v))$$

"w" is a principal component, "k" is chosen based on data.

Spectral Hashing, Yair Weiss, Antonio Torralba, Rob Fergus, NIPS 2008

Locality-sensitive binary codes:

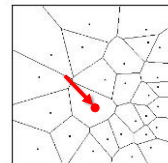
$$h(v) = \text{sign}(\cos(kw \cdot v))$$

"w" is a randomly sampled vector, "k" is chosen based on data.

Locality-Sensitive Binary Codes from Shift-Invariant Kernels, Maxim Raginsky, Svetlana Lazebnik, NIPS 2009

## Visual words

What if we just quantize the descriptors to create "visual words."



Each descriptor = one integer

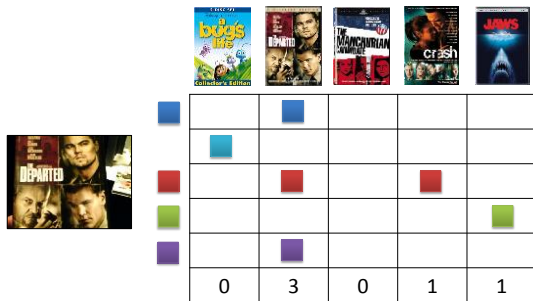
1000 descriptors x 32 bits =

4,000 bytes per image

**We're in RAM (on my laptop)!**

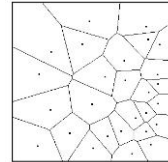
1 million images = 3.7 gigabytes

## Inverse lookup table



## Creating vocabulary

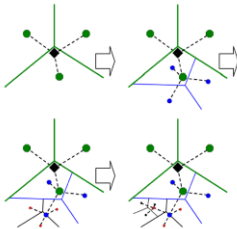
Naïve method is to use k-means clustering on the descriptors.



But this is slow to assign new descriptors to visual words. Need to match the descriptor to every cluster mean = expensive when the vocabulary has 1,000,000 words.

## Vocabulary tree

Recursively split descriptor space using k-means, with  $k \in [3,10]$



Only need 60 comparisons for  $k = 10$  with 1,000,000 visual words.

Scalable Recognition with a Vocabulary Tree, D. Nister, H. Stewenius, CVPR 2006.

## Stop words

If a visual word commonly occurs in many images remove it. You don't want too many images returned for a single word in the inverse look-up table.

It is common practice to have a few thousand stop words for large vocabulary sizes.

It's why search engines don't use the words "a" and "the" in their search queries...

## Weighting visual words

Some visual words are more informative than others.

Use TF-IDF weighting. If a visual word occurs frequently in an image but is rare in other images give it a higher weight.

$$\text{TF (term frequency)} \quad \text{tf}_{i,j} = \frac{n_{i,j}}{\sum_k n_{k,j}}$$

$$\text{IDF (inverse document frequency)} \quad \text{idf}_i = \log \frac{|D|}{|\{j : t_i \in d_j\}|}$$

$$\text{TF-IDF} \quad (\text{tf-idf})_{i,j} = \text{tf}_{i,j} \times \text{idf}_i$$

Commonly used in many types of document retrieval.

## Reducing # of visual words

What if storing a single integer per visual word is too much?

1000 descriptors x 32 bits =  
4,000 bytes per image

1 million images = 3.7 gigabytes

How can we reduce the number of visual words?

100 visual words x 32 bits =  
400 bytes per image

1 million images = 380 megabytes **We're in RAM (on smartphone)!**

## Randomly removing visual words

$$P(v \in I_1, v \in I_2) = 0.5$$

Randomly remove 2/3s of visual words:

$$P(v \in I_1, v \in I_2) = 0.5 * 0.33 * 0.33 = 0.0555$$

Image 1	Image 2	Image 1	Image 2
2	2	2	85
21	32	67	231
23	85	321	678
67	86	363	743
105	35		
231	105		
321	231		
363	321		
375	375		
578	678		
586	743		
745	745		

**Not a good idea!**

## Randomly remove specific visual words

For example: Remove all even visual words.

$$P(v \in I_1, v \in I_2) = 0.5$$

Some images may not have any visual words remaining:

Image 1	Image 2	Image 1	Image 2
2	2	63	85
22	32	325	35
24	85	745	105
63	86		231
104	35		321
232	105		743
325	231		745
366	321		
378	378		
578	678		
586	743		
745	745		

## Min-hash

Maintain Jaccard similarity while keeping a constant number of visual words per image.

Jaccard similarity:

$$\text{sim}_J(C_i, C_j) = \frac{|C_i \cap C_j|}{|C_i \cup C_j|}$$

$C_1$	$C_2$	
0	1	$\text{sim}_J(C_1, C_2) = 2/5 = 0.4$
1	0	
1	1	
0	0	
1	1	
0	1	

\*Minhash slides based on material from Rajeev Motwani and Jeff Ullman

## Key Observation

- For columns  $C_i, C_j$ , four types of rows

	$C_i$	$C_j$
<b>A</b>	1	1
<b>B</b>	1	0
<b>C</b>	0	1
<b>D</b>	0	0

- Overload notation:** A = # of rows of type A

- Claim**  $\text{sim}_J(C_i, C_j) = \frac{A}{A+B+C}$

## Min Hashing

- Randomly **permute** rows
- Hash  $h(C_i) =$  index of first row with 1 in column  $C_i$
- Surprising Property**
- Why?**  $P[h(C_i) = h(C_j)] = \text{sim}_J(C_i, C_j)$ 
  - Both are  $A/(A+B+C)$
  - Look down columns  $C_i, C_j$  until first **non-Type-D** row
  - $h(C_i) = h(C_j) \iff$  type A row

## Min-Hash Signatures

- Pick** – P random row permutations
- MinHash Signature**
  - $\text{sig}(C) =$  list of P indexes of first rows with 1 in column C
- Similarity of signatures**
  - Let  $\text{sim}_H(\text{sig}(C_i), \text{sig}(C_j)) =$  fraction of permutations where MinHash values agree
  - Observe  $E[\text{sim}_H(\text{sig}(C_i), \text{sig}(C_j))] = \text{sim}_J(C_i, C_j)$



## Sketches

What if hashes aren't unique enough? I.e., we return too many possible matches per hash in an inverse look-up table?

Concatenate the hashes into "sketches" of size "k".

$$h_1 = 23, h_2 = 243, h_3 = 598 \rightarrow s_1 = 598,243,023$$

The probability of two sketches colliding is:

$$\text{sim}_j(C_i, C_j)^k$$

Typically you have to balance the precision/recall tradeoffs when picking the sketch size and number of sketches.

## Overview

1 million images	100 GB
1 million images (descriptors)	120 GB
1 million images (descriptors PCA)	30 GB
1 million images (binary descriptors)	7.5 GB
1 million images (visual words)	3.7 GB
1 million images (hashed visual words)	380 MB
10 billion images (hashed visual words)	3.6 TB