Visual Recognition and Machine Learning Summer School Paris 2011

Large scale visual search – part 1

Josef Sivic

http://www.di.ens.fr/~josef INRIA, WILLOW, ENS/INRIA/CNRS UMR 8548 Laboratoire d'Informatique, Ecole Normale Supérieure, Paris

With slides from: O. Chum, K. Grauman, I. Laptev, S. Lazebnik, B. Leibe, D. Lowe, J. Philbin, J. Ponce, D. Nister, C. Schmid, N. Snavely, A. Zisserman

Outline

- 1. Local invariant features (45 mins, C. Schmid)
- Matching and recognition with local features (45 mins, J. Sivic)

3. Efficient visual search (45 mins, J. Sivic)

4. Very large scale visual indexing – recent work (45 mins, C. Schmid)

Practical session – Panorama stitching (60 mins) Download: http://www.di.ens.fr/willow/events/cvml2011/mosaic.zip

Example II: Two images again





1000+ descriptors per image



Match regions between frames using SIFT descriptors and spatial consistency



Multiple regions overcome problem of partial occlusion

Approach - review

1. Establish tentative (or putative) correspondence based on local appearance of individual features (now)

2. Verify matches based on semi-local / global geometric relations (Part 2).

What about multiple images?

• So far, we have seen successful matching of a query image to a single target image using local features.

• How to generalize this strategy to multiple target images with reasonable complexity?

• 10, 10², 10³, ..., 10⁷, ... 10¹⁰ images?

History of "large scale" visual search with local regions

Schmid and Mohr '97 Sivic and Zisserman'03 Nister and Stewenius'06 Philbin et al.'07 Chum et al.'07 + Jegou et al.'07 Chum et al.'08 Jegou et al. '09 Jegou et al. '10

- 1k images
- 5k images
- 50k images (1M)
- 100k images
- 1M images
- 5M images
- 10M images
- ~100M images

All on a single machine in \sim 1 second!

Two strategies

- 1. Efficient approximate nearest neighbour search on local feature descriptors.
- 2. Quantize descriptors into a "visual vocabulary" and use efficient techniques from text retrieval.

(Bag-of-words representation)

Strategy I: Efficient approximate NN search



- 1. Compute local features in each image independently (Part 1)
- 2. "Label" each feature by a descriptor vector based on its intensity (Part 1)
- 3. Finding corresponding features is transformed to finding nearest neighbour vectors
- 4. Rank matched images by number of (tentatively) corresponding regions
- 5. Verify top ranked images based on spatial consistency (Part 2)

Finding nearest neighbour vectors

Establish correspondences between object model image and images in the database by **nearest neighbour matching** on SIFT vectors



Solve following problem for all feature vectors, $\mathbf{x}_j \in \mathcal{R}^{128}$, in the query image: $\forall j \ NN(j) = \arg\min_i ||\mathbf{x}_i - \mathbf{x}_j||$ where, $\mathbf{x}_i \in \mathcal{R}^{128}$, are features from all the database images.

Quick look at the complexity of the NN-search

N ... images

- M ... regions per image (~1000)
- D ... dimension of the descriptor (~128)

Exhaustive linear search: O(M NMD)

Example:

- Matching two images (N=1), each having 1000 SIFT descriptors Nearest neighbors search: 0.4 s (2 GHz CPU, implemenation in C)
- Memory footprint: 1000 * 128 = 128kB / image

# of images	CPU time	Memory req.	
N = 1,000	. ~7min	(~1	00MB)
N = 10,000	. ~1h7min	(~	1GB)
N = 10 ⁷	~115 days	(~	1TB)
All images on Facebook: $N = 10^{10} \dots \sim 300$ years (~ 1PB)			

Nearest-neighbor matching

Solve following problem for all feature vectors, \mathbf{x}_{i} , in the query image:

$$\forall j \ NN(j) = \arg\min_i ||\mathbf{x}_i - \mathbf{x}_j||$$

where x_i are features in database images.

Nearest-neighbour matching is the major computational bottleneck

- Linear search performs *dn* operations for *n* features in the database and *d* dimensions
- No exact methods are faster than linear search for d>10
- Approximate methods can be much faster, but at the cost of missing some correct matches. Failure rate gets worse for large datasets.

Indexing local features: approximate nearest neighbor search



Best-Bin First (BBF), a variant of k-d trees that uses priority queue to examine most promising branches first [Beis & Lowe, CVPR 1997]



Locality-Sensitive Hashing (LSH), a randomized hashing technique using hash functions that map similar points to the same bin, with high probability [Indyk & Motwani, 1998]

K-d tree

• K-d tree is a binary tree data structure for organizing a set of points in a K-dimensional space.

• Each internal node is associated with an axis aligned hyper-plane splitting its associated points into two sub-trees.

• Dimensions with high variance are chosen first.

• Position of the splitting hyper-plane is chosen as the mean/median of the projected points.



Images: Anna Atramentov

K-d tree construction

Simple 2D example



K-d tree query



Slide credit: Anna Atramentov

K-d tree: Backtracking

Backtracking is necessary as the true nearest neighbor may not lie in the query cell.

But in some cases, almost all cells need to be inspected.





A bad distribution which forces almost all nodes to be inspected.

Figure: A. Moore

Solution: Approximate nearest neighbor K-d tree

Key ideas:

- Search k-d tree bins in order of distance from query
- Requires use of a priority queue
- Limit the number of neighbouring k-d tree bins to explore: only approximate NN is found



Reduce the boundary effects by randomization

Randomized K-d trees

- How to choose the dimension to split and the splitting point?
 - Pick dimension with the highest variance
 - Split at the mean/median

 Multiple randomized trees increase the chances of finding nearby points



Randomized K-d trees: discussion

- Find approximate nearest neighbor in O(logN) time, where N is the number of data points.
- Increased memory requirements: needs to store multiple (~8) trees
- Good performance in practice for recognition problems (NN-search for SIFT descriptors and image patches).
- Code available online:

http://people.cs.ubc.ca/~mariusm/index.php/FLANN/FLANN

Indexing local features: approximate nearest neighbor search



Best-Bin First (BBF), a variant of k-d trees that uses priority queue to examine most promising branches first [Beis & Lowe, CVPR 1997]



Locality-Sensitive Hashing (LSH), a randomized hashing technique using hash functions that map similar points to the same bin, with high probability [Indyk & Motwani, 1998] Locality Sensitive Hashing (LSH)

Idea: construct hash functions g: $\mathbb{R}^d \rightarrow \mathbb{Z}^k$ such that

for any points p,q:

If $||p-q|| \le r$, then Pr[g(p)=g(q)] is "high" or "not-so-small" If ||p-q|| > cr, then Pr[g(p)=g(q)] is "small"

Example of g: linear projections

 $g(p) = \langle h_1(p), h_2(p), ..., h_k(p) \rangle$, where $h_{X,b}(p) = \lfloor (p^*X+b)/w \rfloor$

[.] is the "floor" operator.
X_i are sampled from a Gaussian.
w is the width of each quantization bin.
b is sampled from uniform distr. [0,w]. [Datar-Im

[Datar-Immorlica-Indyk-Mirrokni'04]

Locality Sensitive Hashing (LSH)

- Choose a random projection
- Project points
- Points close in the original space remain close under the projection
- Unfortunately, converse not true



 Answer: use multiple quantized projections which define a high-dimensional "grid" Locality Sensitive Hashing (LSH)

- Cell contents can be efficiently indexed using a hash table
- Repeat to avoid quantization errors near the cell boundaries



- Point that shares at least one cell = potential candidate
- Compute distance to all candidates

LSH: discussion

In theory, query time is O(kL), where k is the number of projections and L is the number of hash tables

I.e. independent of the number of points, N.

In practice, LSH has high memory requirements as large number of projections/hash tables are needed.

Code and more materials available online: http://www.mit.edu/~andoni/LSH/

See also: <u>http://cobweb.ecn.purdue.edu/~malcolm/yahoo/</u> <u>Slaney2008(LSHTutorialDraft).pdf</u>

Comparison of approximate NN-search methods

http://www.cs.ubc.ca/~lowe/papers/09muja.pdf

FAST APPROXIMATE NEAREST NEIGHBORS WITH AUTOMATIC ALGORITHM CONFIGURATION

Marius Muja, David G. Lowe

Computer Science Department, University of British Columbia, Vancouver, B.C., Canada mariusm@cs.ubc.ca, lowe@cs.ubc.ca

Keywords: nearest-neighbors search, randomized kd-trees, hierarchical k-means tree, clustering.

Abstract: For many computer vision problems, the most time consuming component consists of nearest neighbor matching in high-dimensional spaces. There are no known exact algorithms for solving these high-dimensional problems that are faster than linear search. Approximate algorithms are known to provide large speedups with only minor loss in accuracy, but many such algorithms have been published with only minimal guidance on selecting an algorithm and its parameters for any given problem. In this paper, we describe a system that answers the question, "What is the fastest approximate nearest-neighbor algorithm for my data?" Our system will take any given dataset and desired degree of precision and use these to automatically determine the best algorithm and parameter values. We also describe a new algorithm that applies priority search on hierarchical k-means trees, which we have found to provide the best known performance on many datasets. After testing a range of alternatives, we have found that multiple randomized k-d trees provide the best performance for other datasets. We are releasing public domain code that implements these approaches. This library provides about one order of magnitude improvement in query time over the best previously available software and provides Comparison of approximate NN-search methods

Dataset: 100K SIFT descriptors



Code for all methods available online, see Muja&Lowe'09

Figure: Muja&Lowe'09

Approximate nearest neighbour search (references)

- J. L. Bentley. Multidimensional binary search trees used for associative searching. Comm. ACM, 18(9), 1975.
- Freidman, J. H., Bentley, J. L., and Finkel, R. A. An algorithm for finding best matches in logarithmic expected time. *ACM Trans. Math. Softw., 3:209–226, 1977.*
- Arya, S., Mount, D. M., Netanyahu, N. S., Silverman, R., and Wu, A. Y. An optimal algorithm for approximate nearest neighbor searching in fixed dimensions. *Journal of the ACM*, 45:891–923, 1998.
- C. Silpa-Anan and R. Hartley. Optimised KD-trees for fast image descriptor matching. In CVPR, 2008.
- M. Muja and D. G. Lowe. Fast approximate nearest neighbors with automatic algorithm configuration. In VISAPP, 2009.
- P. Indyk and R. Motwani, "Approximate nearest neighbors: towards removing the curse of dimensionality," in *Proc. of 30th ACM Symposium on Theory of Computing, 1998*
- G. Shakhnarovich, P. Viola, and T. Darrell, "Fast pose estimation with parametersensitive hashing," in *Proc. of the IEEE International Conference on Computer Vision,* 2003.
- R. Salakhutdinov and G. Hinton, "Semantic Hashing," ACM SIGIR, 2007.
- Y. Weiss, A. Torralba, and R. Fergus, "Spectral hashing," in NIPS, 2008.

ANN - search (references continued)

- O. Chum, J. Philbin, and A. Zisserman. Near duplicate image detection: min-hash and tfidf weighting. BMVC., 2008.
- M. Raginsky and S. Lazebnik, "Locality-Sensitive Binary Codes from Shift-Invariant Kernels," in *Proc. of Advances in neural information processing systems, 2009.*
- B. Kulis and K. Grauman, "Kernelized locality-sensitive hashing for scalable image search," Proc. of the IEEE International Conference on Computer Vision, 2009.
- J. Wang, S. Kumar, and S.-F. Chang, "Semi-supervised hashing for scalable image retrieval," in IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR), 2010.
- J. Wang, S. Kumar, and S.-F. Chang, "Sequential projection learning for hashing with compact codes," in Proceedings of the 27th International Conference on Machine Learning, 2010.

So far ...

- Linear exhaustive search can be prohibitively expensive for large image collections
- Answer (so far): approximate NN search methods
 - Randomized KD-trees
 - Locality sensitive hashing
- However, memory footprint can be still high.
 Example: N = 10⁷ images, 10¹⁰ SIFT features with 128B per feature > 1TB of memory

Look how text-based search engines (Google) index documents – **inverted files**.

Indexing text with inverted files



Need to map feature descriptors to "visual words".









Visual words

Example: each group of patches belongs to the same visual word


Samples of visual words (clusters on SIFT descriptors):





More specific example

Samples of visual words (clusters on SIFT descriptors):





More specific example

Visual words

- First explored for texture and material representations
- *Texton* = cluster center of filter responses over collection of images
- Describe textures and materials based on distribution of prototypical texture elements.



Leung & Malik 1999; Varma & Zisserman, 2002; Lazebnik, Schmid & Ponce, 2003;

Inverted file index for images comprised of visual words



- Score each image by the number of common visual words (tentative correspondences)
- Worst case complexity is linear in the number of images N
- In practice, it is linear in the length of the lists (<< N)

Another interpretation: Bags of visual words

Summarize entire image based on its distribution (histogram) of visual word occurrences.

Analogous to bag of words representation commonly used for documents.







Another interpretation: the bag-of-words model

For a vocabulary of size K, each image is represented by a K-vector

$$\mathbf{v}_d = (t_1, \dots, t_i, \dots, t_K)^\top$$

where t_i is the number of occurrences of visual word i.

Images are ranked by the normalized scalar product between the query vector v_a and all vectors in the database v_d :

$$f_d = \frac{\mathbf{v}_q^{\top} \mathbf{v}_d}{\|\mathbf{v}_q\|_2 \|\mathbf{v}_d\|_2}$$

Scalar product can be computed efficiently using inverted file.

What if vectors are binary? What is the meaning of $\mathbf{v}_q^{\top} \mathbf{v}_d$?

Strategy I: Efficient approximate NN search



- 1. Compute local features in each image independently (offline)
- 2. "Label" each feature by a descriptor vector based on its intensity (offline)
- 3. Finding corresponding features is transformed to finding nearest neighbour vectors
- 4. Rank matched images by number of (tentatively) corresponding regions
- 5. Verify top ranked images based on spatial consistency (Part 2)

Strategy II: Match histograms of visual words



- 1. Compute affine covariant regions in each frame independently (offline)
- 2. "Label" each region by a vector of descriptors based on its intensity (offline)
- 3. Build histograms of visual words by descriptor quantization (offline)
- 4. Rank retrieved frames by matching vis. word histograms using inverted files.
- 5. Verify retrieved frame based on spatial consistency (Part 2)

Visual words: discussion I.

Efficiency – cost of quantization

 Need to still assign each local descriptor to one of the cluster centers. Could be prohibitive for large vocabularies (K=1M)

- Approximate NN-search still needed
- True also for building the vocabulary

Visual words: discussion II.

Generalization

 Is vocabulary/quantization learned on one dataset good for searching another dataset?

• Experimentally observe a loss in performance.

But, see recent work by Jegou et al.:

Hamming Embedding and Weak Geometry Consistency for Large Scale Image Search, ECCV'2008

http://lear.inrialpes.fr/pubs/2008/JDS08a/

Visual words: discussion III.

• Need to determine the size of the vocabulary, K.

• Other algorithms for building vocabularies, e.g. agglomerative clustering / mean-shift, but typically more expensive.

Supervised quantization?

Also give examples of images / descriptors which should and should not match.

E.g.:

Philbin et al. ECCV'10, http://www.robots.ox.ac.uk/~vgg/publications/html/philbin10b-bibtex.html

Visual search using local regions (references)

- C. Schmid, R. Mohr, Local Greyvalue Invariants for Image Retrieval, PAMI, 1997
- J. Sivic, A. Zisserman, Text retrieval approach to object matching in videos, ICCV, 2003
- D. Nister, H. Stewenius, Scalable Recognition with a Vocabulary Tree, CVPR, 2006.
- J. Philbin, O. Chum, M. Isard, J. Sivic, A. Zisserman, Object retrieval with large vocabularies and fast spatial matching, CVPR, 2007
- O. Chum, J. Philbin, M. Isard, J. Sivic, A. Zisserman, Total Recall: Automatic Query Expansion with a Generative Feature Model for Object Retrieval, ICCV, 2007
- H. Jegou, M. Douze, C. Schmid, Hamming embedding and weak geometric consistency for large scale image search, ECCV'2008
- O. Chum, M. Perdoch, J. Matas: Geometric min-Hashing: Finding a (Thick) Needle in a Haystack, CVPR 2009
- H. Jégou, M. Douze and C. Schmid, On the burstiness of visual elements, CVPR, 2009
- H. Jégou, M. Douze, C. Schmid and P. Pérez, Aggregating local descriptors into a compact image representation, CVPR'2010

Demo

Oxford Buildings Search

http://www.robots.ox.ac.uk/~vgg/research/oxbuildings/ index.html

Oxford buildings dataset

- Automatically crawled from **flickr**
- Consists of:

Dataset	Resolution	# images	# features	Descriptor size
i	1024×768	$5,\!062$	$16,\!334,\!970$	1.9 GB
ii	1024×768	99,782	$277,\!770,\!833$	$33.1~\mathrm{GB}$
iii	500×333	$1,\!040,\!801$	$1,\!186,\!469,\!709$	141.4 GB
Total		$1,\!145,\!645$	$1,\!480,\!575,\!512$	$176.4~\mathrm{GB}$



Oxford buildings dataset

Landmarks plus queries used for evaluation



- Ground truth obtained for 11 landmarks
- Evaluate performance by mean Average Precision

Measuring retrieval performance: Precision - Recall

- Precision: % of returned images that
 are relevant
- Recall: % of relevant images that are returned





Average Precision



- A good AP score requires both high recall and high precision
- Application-independent

Performance measured by mean Average Precision (mAP) over 55 queries on 100K or 1.1M image datasets







- high precision at low recall (like google)
- variation in performance over query
- none retrieve all instances

Why aren't all objects retrieved?



Obtaining visual words is like a sensor measuring the image

"noise" in the measurement process means that some visual words are missing or incorrect, e.g. due to

- Missed detections
- Changes beyond built in invariance
- Quantization effects

Query expansion
 Better quantization

Consequence: Visual word in query is missing in target image

Automatic query expansion

Visual word representations of two images of the same object may differ (due to e.g. detection/quantization noise) resulting in missed returns

Initial returns may be used to add new relevant visual words to the query

Strong spatial model prevents 'drift' by discarding false positives

[Chum, Philbin, Sivic, Isard, Zisserman, ICCV'07; Chum, Mikulik, Perdoch, Matas, CVPR'11]

Visual query expansion - overview

1. Original query





3. Spatial verification

4. New enhanced query





















Demo

Query Expansion

Query image

Originally retrieved







































Quantization errors

Typically, quantization has a significant impact on the final performance of the system [Sivic03,Nister06,Philbin07]

Quantization errors split features that should be grouped together and confuse features that should be separated



Overcoming quantization errors

• Soft-assign each descriptor to multiple cluster centers [Philbin et al. CVPR 2008, Van Gemert et al. ECCV 2008]



Learning a vocabulary to overcome quantization errors [Mikulik et al. ECCV 2010, Philbin et al. ECCV 2010] Beyond bag-of-visual-words.

Locality-constrained linear coding. [Wang et al. CVPR 2010]

- Represent data point as a linear combination of nearby cluster centers.
- Store the coefficients of linear combination.

Used for category-level classification.



Other recent work

Learning a vocabulary to overcome quantization errors [Mikulik et al. ECCV 2010, Philbin et al. ECCV 2010]

Large scale image clustering [Chum et al. CVPR 2009, Philbin et al. IJCV 2010, Li et al., ECCV 2008]

Matching in structured datasets (3D landmarks or street-view images) [Knopp et al. ECCV 2010, Zamir&Shah ECCV 2010, Li et al. ECCV 2010, Baatz et al. ECCV 2010]

Putting geometry into index

[Perdoch et al. CVPR 2009, Jegou et al. ECCV 2008]

Improving efficiency by min-hash and geometric min-hash [Chum et al. 07-09] and compressing bag-of-visual-word vectors [Jegou et al. 09].

What objects/scenes local regions do not work on?



What objects/scenes local regions do not work on?



E.g. texture-less objects, objects defined by shape, deformable objects, wiry objects.

Example applications of large scale visual search and matching Application: Internet-based inpainting Photo-editing using images of the same place [Whyte, Sivic and Zisserman, 2009], but see also [Hays and Efros, 2007].





Application: place recognition (retrieval in a structured (on a map) database)



[Knopp, Sivic, Pajdla, ECCV 2010]

Correctly recognized examples













More correctly recognized examples


Matching and 3D reconstruction in large unstructured datasets.



Building Rome in a Day, Sameer Agarwal, Noah Snavely, Ian Simon, Steven M. Seitz and Richard Szeliski, International Conference on Computer Vision, 2009 http://grail.cs.washington.edu/rome/

See also [Havlena, Torrii, Knopp and Pajdla, CVPR 2009].

Figure: N. Snavely

Example of the final 3D point cloud and cameras

57,845 downloaded images, 11,868 registered images. This video: 4,619 images.



Mobile visual search

Bing visual scan





Google Goggles

<complex-block>



PLINKART

Plink Art is an app for your mobile phone that lets you identify almost any work of art just by taking a photo of it.

Snap pictures of objects (media covers including books, CDs, DVDs, games, and newspapers and magazines) receive information, price comparisons, and reviews. Consists of an iPhone application and a web-based tool, which remembers all your requests.

Gives instant information on what you see

Visual Search

Κ



and others... Snaptell.com, Moodstocks.com

Example







k Visual Search

Learn more

kooaba Paperboy delivers digital extras for print



What next?

Searchable visual memory?



Figure: iphoneography.com

What next?

Visual search for texture-less, wiry, or deformable objects..



.. also faces/people and their actions [see I. Laptev's lecture on Friday]



The IKEA problem

Recognize all IKEA furniture in all YouTube videos



Category-level visual search [See lectures by C. Schmid and A. Zisserman tomorrow]



same category











See also e.g. [Torresani et al. ECCV 2010]

Recognize and match non-photographic depictions









Automatic alignment of paintings and photographs depicting the same 3D scene [Russell, Sivic, Ponce and Dessales, 2011]

Inputs





Photographs





3D model



Painting



Viewpoint of painting



Outline

- 1. Local invariant features (45 mins, C. Schmid)
- 2. Matching and recognition with local features (45 mins, J. Sivic)
- 3. Efficient visual search (45 mins, J. Sivic)
- 4. Very large scale visual indexing recent work (45 mins, C. Schmid)

Practical session (60 mins)