

Similarity and Invariance; Searching for Similar Images

36-350: Data Mining

2 September 2009

READING: Section 14.5 in the textbook.

So far, we have seen how to search and categorize texts by representing them as feature vectors — bag-of-word vectors, appropriately scaled and weighted. Except for page-rank, all of the procedures we have used have depended *only* on the feature vectors, not on the underlying text. (Page-rank doesn't even depend on the text, just on the links.) We can use the same techniques on any other kind of data, if we can just define numerical features for them — and as long as we are content to *just* use numerical features. The nature of the underlying objects, and even the method of representation, play no role in things like similarity searching or multi-dimensional scaling. They don't care where vectors come from, or what they actually represent, just that they're vectors. Using feature vectors rather than raw data is a kind of **abstraction**, in the sense in which we use that term in math or computer programming, and one of the advantages of abstraction is, exactly, facilitating the re-use of methods and procedures. Whether recycling our techniques from text to some other domain will work well depends largely on how well we chose our features, i.e., the mode of representation we use to hide data and produce abstract vectors. As an example, let's look at images.

1 Image Representation

At the lowest level, computers represent images as matrices of pixels; the matrix entries give the color of each pixel.¹ So let's talk about colors.

1.1 Colors

Normal human eyes contain three different sorts of color-sensitive cells, each tuned to a different part of the visible spectrum. These give us the three primary

¹Actually, most graphics formats, like `jpeg`, use other representations, which are more compact for typical images. All formats get translated into pixel matrices when they're displayed, however, so we'll ignore this complication.

colors, red, green and blue, and others colors are combinations of them.² The **RGB** encoding of a color is a three-dimensional vector indicating how much of red, green and blue needs to be blended to make the color, from 0 to 1 on each axis. Black is (0, 0, 0) and white is (1, 1, 1); all the other colors fall somewhere in the cube. There are other representations of color space, but RGB is by far the most common, and we'll stick with it.³

Many physically distinct colors, from different points in the RGB cube, are psychologically indistinguishable⁴, or differ only trivially. Rather than use the full color space available to the computer, it is often desirable to **quantize** color down to a smaller number of values, effectively merging similar colors (so that we just have one “red”, and not millions of shades of red). Geometrically, this means **partitioning** the RGB cube, cutting it into cells, and treating all colors in a cell as equivalent, calling them with a single name. (This corresponds to stemming for words, or, if the cells are very big and coarse, going beyond stemming to lumping together words with similar meanings, like “steel” and “iron”.)

1.2 Back to Images

Just as with documents, we would now like to do similarity searching for images. The user should be able to give the system an initial image and tell it “find me more which look like this”.

The most straightforward approach would use the computer’s representation of the image directly. If it’s a color image of $M \times N$ pixels, we have a list of $3MN$ numbers; call this our feature vector, and calculate Euclidean distances as before. This can actually work OK for some purposes, like undoing the effects of some kinds of camera distortions, but it’s not very good at finding meaningfully similar images. For instance, in Figure 1, `flower2` is a close-up of the middle flower from `flower1`, but their Euclidean distance will be very large. (Why?)

Since the bag-of-words representation served us well with documents, we might try the same thing here; it’s called the bag of colors. We first chose a quantization of the color space, and then, for each quantized color, count the number of pixels of that color in the image (possibly zero). That is, our procedure is as follows:

1. The user gives us an image Q . Convert it into a vector of color counts.
2. For each image in the collection, measure the distance to Q .
3. Return the k images closest to Q .

²Red-green color-blind people effectively have only two primary colors, because of mutations changing the tuning of their color-sensitive cells. Some other species of animals have more than three primary colors.

³Color-encoding schemes for *printing* are different, because the physics of blending inks and dyes is different from that of blending light.

⁴Indistinguishable by whom? On average, women tend to be better at color discrimination than men, and there is some evidence that this is due to hormonal differences, but the overlap is large, and most people can get better with practice.

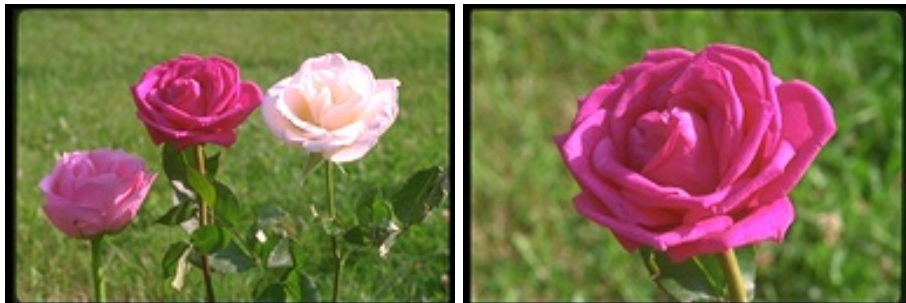


Figure 1: The images `flower1` (left) and `flower2` (right) from our data set. `flower2` is a close-up photo of the middle flower from `flower1`.

Today’s data are photographs of flowers, photographs of tigers, and photographs of the ocean. Table 1 is part of the color-count representation of the data.

Just as with documents, we found it useful to emphasize rare words, it can be helpful to emphasize rare colors, by using the **inverse picture frequency**, *IPF*:

$$IPF(c) = \log \left(\frac{N}{n_c} \right)$$

Colors which occur frequently, like a black border, will have a small *IPF*(*c*). (Notice that both `flower1` and `flower2` have such a border.)

Taking a set of 30 pictures, equally divided into scenes of the ocean, of flowers and of tigers, and ask how often a picture’s nearest neighbor is of the same type, we get the results in Table 2 depending on how we normalize the bags-of-colors (cf. Figure 2).

As a sheer classification task, this example is a bit trivial. Our ancestors spent the last forty-odd million years getting very good at sorting flowers from big cats at a glance.⁵ On the other hand, we do not have the same evolutionary priming for other kinds of image classification, like looking at a microscope slide and telling whether the cells on it are normal or cancerous...

2 Invariance

When trying to decide whether or not a certain representation is going to be useful, it’s often helpful to think about the **invariances** of the representation. That is, what sorts of changes could we make to the object, without changing its representation? (The representation is **invariant under** those changes.) Invariances tell us what the representation ignores, what concrete aspects of the object it abstracts away from. Abstraction as such is neither good nor bad; particular abstractions are more or less useful for particular tasks.

⁵More precisely, those of our relatives who could not do this are tragically under-represented among our ancestors.

	violetred4	goldenrod	lightskyblue4	gray58.2
flower1	59	0	0	0
flower2	128	0	0	0
flower3	166	75	0	0
tiger1	0	4	457	85
tiger2	0	0	0	2
tiger3	0	0	0	115
ocean1	0	0	4326	433
ocean2	0	0	2761	142
ocean3	0	0	1179	8596

Table 1: Part of the bag-of-colors vector for the data.

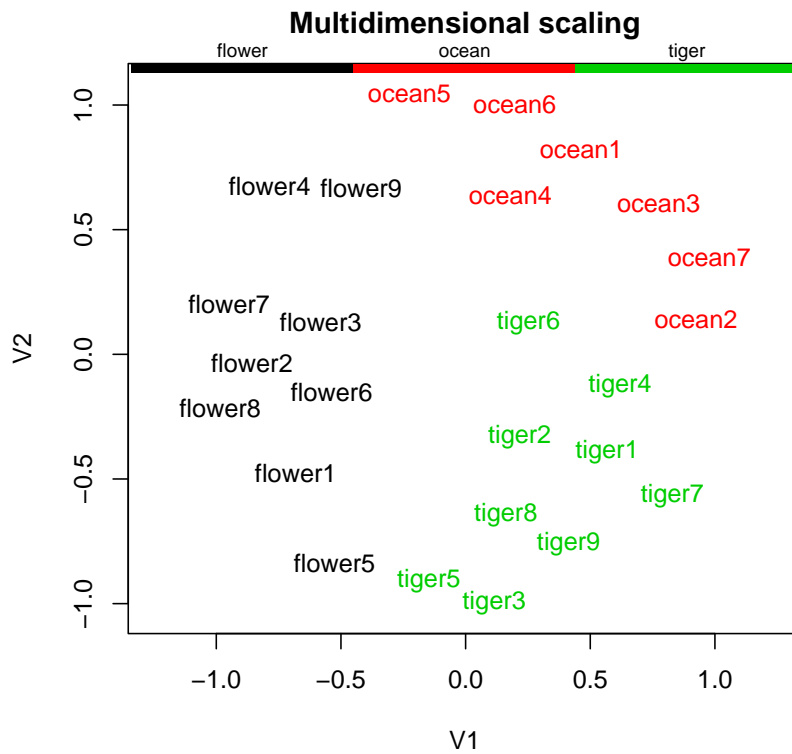


Figure 2: Multi-dimensional scaling of the flower/tiger/ocean images, using bag-of-color vectors, Euclidean distance with length normalization, and inverse-picture-frequency weighting.

Normalization	Equal weight	IPF weight
None	8	4
Picture size	8	0
Euclidean length	7	0

Table 2: How often an image’s nearest neighbor belongs to a different class, measuring distance between bag-of-color vectors with different normalizations and weightings of colors. There are 30 images total.

The bag-of-words representation is invariant to changes in punctuation, word order, and general grammar. If it is normalized, it is also invariant to changes in document length. If inverse-document-frequency weighting is used, it is invariant under changes in the exact counts of the most common words (“the”, “of”, “and”, etc.) which appear in any text written in English. Stemming, and grouping together synonyms, would add extra invariances. These invariances make the representations good at finding documents with similar topics, even though they might be written by different people. (By the same token, these invariances make it hard to, say, find all the stories by a given author.)

What invariances do we get for images in the bag-of-colors representation?

- Position of objects, pose, small changes in camera angle
- Small amounts of zooming
- Small amounts of unusual colors
- Differences in texture: we can scramble the pixels in any region, even the whole image, without effect. (Plaids and stripes seem the same.)

Some other invariances we might like, but don’t get with color counts:

- Lighting and shadows; time of day
- Occlusion (one object covers another), placement in 3D
- Blurring, changes in focus

If there is an invariance we don’t want, we can “break” that invariance by adding additional features. There are various schemes for representing textures, so that plaids, stripes, spots, “snow”, solid colors, etc. can be distinguished, and many systems for content-based image searching combine color-counting with textures.

It’s harder to *add* an invariance than to break one; generally you need to go back to the beginning and re-design your representation. The sorts of invariances we I just said we’d like to have, but don’t, with this representation are all very hard to achieve; roughly 40% of the human cortex is devoted to visual processing for a very good reason.

You should think about what invariances the following representations have, and how they compare to the invariances of the bag-of-colors representation. The first one is much smaller than the bag of colors, and the second one is larger.


1. Represent an image by its average pixel color.
2. Represent an image by dividing it into four quarters, and doing color-counts in each quarter. (Is there any reason why four equal parts should be special?)

3 Practice

Content-based image retrieval systems are used commercially. All the ones I know of include the bag-of-colors representation, but generally also a number of other features. For example, you should go to <http://labs.systemone.at/retrievr/> and think about what features they are using, and how well they work. (Cf. Figure 3.)

Google's image search is very popular, but it's a total cheat, and doesn't analyze the images at all. Instead, it searches on the phrases you give it, takes the pages which get a high rank, and assumes that images which appear on those pages close to your query phrase are probably pictures of what you are looking for. This works well for a lot of things, but can give curious results when there aren't, actually, many pictures of what you are looking for.

Search by:
[Sketch](#) • [Image](#)




Search for similar images by

- [uploading](#) an image file
- or [entering the URL](#) of an image.


Now! Please help us [rate sketches](#) for The Art of retrieval! Only a million to go ...

Still new! You can search by [uploading images](#) now as well. Also: Your sketches have URLs! Send 'em around.


This is an experimental service. Please treat it nicely and send copious amounts of feedback! For some background, [read here](#).




From [karavekov](#)




From [pupanna](#)




From [Feuillu](#)




From [Rebsana](#)




From [kelsana](#)




From [Katemina](#)




From [Fodd](#)




From [mkaggen](#)




From [Lydia Pinkham](#)




From [Holly Yvonne](#)




From [Mi Xavier](#)




From [alecani](#)




From [pasodoble](#)



From [rebekka](#)



From [LoveYourDog.com](#)



From [heavenuphere](#)

Figure 3: Example of similarity-search for images in action, using <http://labs.systemone.at/>. The query image, on the left, is an old snapshot of my cat. The returned images were obtained from flickr.com. What features are being used here? Are they working?