Learning Globally-Consistent Local Distance Functions for Shape Based Image Retrieval and Classification

Andrea Frome, Yoram Singer, Fei Sha, Jitendra Malik ICCV 2007

Presented by Roie Kliper

Learning Globally-Consistent Local Distance Functions for Shape-Based Retrieval and Classification

> Andrea Frome EECS, UC Berkeley antrea, fromegonail.com

Google, Inc. singergoogle.com

Fei Sha EECS, UC Berkeley feishages.berkeley.edu

EECS, UC Berkeley nalik@cs.berkeley.edu

Yoram Singer

Abstract

We address the problem of visual category recogni by learning an image-to-image distance function that attempts to satisfy the following property: the distance betoren images from the same category should be less than the distance between images from different categories. We use patch-based feature vectors common in object recognition work as a basis for our image-to-image distance functions. Our large-margin formulation for learning the distance functions is similar to formulations used in the machine learning literature on distance metric learning, how ever we differ in that we learn local distance functionsa different parameterized function for every image of our training set-whereas typically a single global distance function is learned. This was a neural approach first introduced in Frome, Singer, & Malik, NIPS 2006. In that work we learned the local distance functions independently, and the outputs of these functions could not be compared at test time without the use of additional heuristics or training. Here we introduce a different approach that has the advantoge that it learns distance functions that are globally consistent in that they can be directly compared for purposes of retrieval and classification. The output of the learning algorithm are weights assigned to the image features, which is installively oppositing in the computer vision setting: some features are more salient than others, and which are more salient depends on the category, or image, being considered. We train and test using the Galtech 101 object recognition benchmark. Using Afteen training images per cate gory, we achieved a mean recognition rate of 63.2% and using twenty images per category, a rate of 66.6%.



Figure 1. There images from the Calaecht01 data set, to drog category, one from the **72605** category. We we distance functions between pairs of images such that to from *j* to 1(D_{jn}) is smaller than from *k* to 1(D_{kn}). It this one form the basis of our learning algorithm.

1. Introduction

Coasider the triplet of images, drawn tech101 dataset [4], shown in Figure 1. We wa sify a query image 1, and we have stored exempt t and k. Let D_{in} be the distance from image D_{kk} be the distance from image k to i where DThen a nearest neighbor classifier which assign gory of the query image (based on which of L imalier, would trivially do the right thing. No this to work, the distance function need not be ric, in general $D_{jk} \neq D_{kj}.$ To approach this pr parameterize the image to image dotance funct weighted linear combination of distances betw based shape feature descriptors, such as SIFT [] metric blur [2]. These features characterize inaby fixed length vectors, which can be compared u L2 metrics. One possible approach to computing to-image distance is to attempt to solve the correproblem by taking into account both distances be

Uuune

- What are Distance Functions and what are the good for?
- Global vs. Local Distance Functions.
- Algorithm Outline & Details...
- Results
- Some conclusions

= (Income the intermediation on Talana (none Energy (Theory 2007))

Distance i unchuns

- A distance function is a function defined over pairs of data points. $D: X \times X \rightarrow \Re$
- A distance metric is a mapping to the nonnegative real number $D: X \times X \rightarrow \Re^+$, that also obeys the following properties:
 - Isolation.
 - Symmetry.
 - Triangular Inequality.



- There are several intuitive ways of transforming a similarity function into a distance function and vice-versa.
- One important and widely used type of similarity functions are kernel functions.
 - A kernel function K can be used to define a Gram matrix (also known as a kernel matrix).
 - The Gram matrix is in essence a similarity matrix of the set of vectors S.
 - It also has several other appealing properties: it is symmetric and positive semidefinite. (In fact it can be shown that any positive semi-definite symmetric matrix corresponds to some kernel function k.)

important?

- A distance functions defines a topology over the space.
- Relieves interesting manifolds (Manifold learning)
- Many Distance-based algorithms:
 - NN classification
 - Radial basis function network
 - K-means clustering,
 - Various graph based methods (linkage clustering)
 - kernel methods.
- Many others....

oalegonzation scenario



Ochicial luca.

 Distance between images from the same category should be less than distance between images from different categories.



Olobal v S Local

• Global: (Ambitious) In general try to capture the global dimensions that are responsible to the data categorization: $D: X \times X \rightarrow \Re$

• e.g: Mahalanobis

• Local: (modest) Learn a distance that is specific for a given focal point (image). $D_a: X \to \Re$

• Make it globally consistent....

Algonum Outime

- Input: patch-based feature vectors such as SIFT or geometric blur. (Using geometric blur and simple color features)
- Parametrize the image-to-image distance functions using a weighted linear combination of distances between patch based shape feature descriptors.
- Generate Triplet based constraints.
- Large margin optimization formulation
- Solve (of the shelf solver)
- Output:
 - weights assigned to the image features
 - Local distance function which are globally consistent (thus can be compared)

descriptors extraction and use: SIFT, Geometric Blur

- Robust under affine transformations.
- These features characterize image patches by fixed length vectors, which can be compared using L1 or L2 Norms.
- Different features can be combined into a sing bag of feature.



• Berg & Malik, CVPR 01

COIDI I Caluico

- Color features are histograms of eightpixel radius patches also centered at edge pixels in the image.
- Any "pixels" in a patch off the edge of the image are counted in a "undefined" bin, and we convert the HSV coordinates of the remaining points to a Cartesian space where the z direction is value and (x, y) is the Cartesian projection of the hue/saturation dimensions.
- Divide the (x, y) space into an 11 °— 11 grid, and make three divisions in the z direction.



routare te mage alotarioee



f

Combining real are types

































A DISTAILLE I UNCTION

 The image to image distance function is define as follows:

$$D_{j \to i} = \sum_{m=1}^{M} w_{j \to i,m} \cdot d_{j \to i,m} = \left\langle \begin{matrix} \mathbf{r} & \mathbf{r} \\ W_{j}, d_{j \to i} \end{matrix} \right\rangle$$

• A Triplet constraint will look like this:

 $\left\langle \vec{W}_{i}, d_{i \to i} \right\rangle > \left\langle W_{k}, d_{k \to i} \right\rangle$

Al alta Set of Thplets





Chousing inpicts

- Choosing an exhaustive set of triplets in impractical.
 - Data is too large to fit in the memory.
 - Iteration run time increases linearly with the number of triplets.
 - Only few triplets are informative.
- Solution: pruning:
 - For each of the M elementary distances in focal image F, we find the top K closest images
 - If K contains both in-class and out-of-class images, create all triplets from (F,inclass,out-class) combinations
 - If K are all in-class images, get closest out-class image then make K triplets (reverse all K are out-class)

ine opunization

- In an idealistic setting we can hope for the following to hold: $\forall_{i,j,k\in T} w_k \cdot d_{k\rightarrow j} > w_i \cdot d_{i\rightarrow j}$
- Rescaling w we can get $\forall_{i,j,k\in T} w_k \cdot d_{k\to j} > w_i \cdot d_{i\to j} + \psi_{i,j,k\in T} w_k \cdot d_{k\to j} > w_i \cdot d_{i\to j} + \psi_{i,j,k\in T} w_k \cdot d_{k\to j} > \psi_i \cdot d_{i\to j} + \psi_{i,j,k\in T} w_k \cdot d_{k\to j} > \psi_i \cdot d_{i\to j} + \psi_{i,j,k\in T} w_k \cdot d_{k\to j} > \psi_i \cdot d_{i\to j} + \psi_{i,j,k\in T} w_k \cdot d_{k\to j} > \psi_i \cdot d_{i\to j} + \psi_{i,j,k\in T} w_k \cdot d_{k\to j} > \psi_i \cdot d_{i\to j} + \psi_{i,j,k\in T} w_k \cdot d_{k\to j} > \psi_i \cdot d_{i\to j} + \psi_{i,j,k\in T} w_k \cdot d_{k\to j} > \psi_i \cdot d_{i\to j} + \psi_{i,j,k\in T} w_k \cdot d_{k\to j} > \psi_i \cdot d_{i\to j} + \psi_{i,j,k\in T} w_k \cdot d_{k\to j} > \psi_i \cdot d_{i\to j} + \psi_{i,j,k\in T} w_k \cdot d_{k\to j} > \psi_i \cdot d_{i\to j} + \psi_{i,j,k\in T} w_k \cdot d_{k\to j} > \psi_i \cdot d_{i\to j} + \psi_i +$
- In reality a relaxation of the above is needed: Using large margin terminology.(hinge loss) we get a linear penalty for deviation:

 $\sum \left[1 - W \cdot X_{i, i, k}\right]_{+}$

The optimization oont...

 Adding a Regularization term: We get the following convex problem:

 $\frac{1}{2} \|W\|^2 + C \sum_{i,j,k\in T} [1 - W \cdot X_{i,j,k}]_+$

Standard soft margin

positive W forces p the distance function

Linear objective function with multip The Dual problem: constraints:

 $\underset{W_{\mathcal{X},\xi_{\mathcal{U}}}}{\underset{W_{\mathcal{X},\xi_{\mathcal{U}}}}{\max}} \frac{1}{2} \left\| \frac{1}{2} \sum_{(ijk)\in T} C\alpha_{ijk} \sum_{ijk} \xi_{ijk} \mu \right\|^{2} + \sum_{(ijk)\in T} \alpha_{ijk} \xi_{ijk} \mu + \sum_{(ijk)\in T} \alpha_{ijk} \mu + \sum_{(ijk)\in T} \alpha_{$



$$\begin{split} & \overleftarrow{\psi}_{i,j,k} : \overleftarrow{\psi}_{k} \leq \widehat{\mathcal{A}}_{ijk} \leq C \\ & \overleftarrow{\psi}_{i,j,k} : \overleftarrow{\mu}_{m}^{W} \geq \widehat{\mathcal{A}}_{ijk} \geq 1 - \xi_{ijk} \\ & \overleftarrow{\psi}_{m} : W_{m} \geq 0 \end{split}$$

Vancon IVI



















itesuits.



• (Taken from Frome's Thesis, 2007)

ILESUILS.



11620112



FASAS FASAS FASAS REGRETARIA RECERTARIA R

i and i futte message

- Seems like for different visual object categories, different notions of similarity are needed: color informative for some classes, local shape for others, geometry for others.
- Though Global learning is an attractive idea, Local learning might be more adequate and meaningful.
- The visual features of an image are not equally important in determining its similarity to some prototype.
- Which features are more salient depend on the category, or even the image, being considered. (hence locality)
- The described method allows for a combination of different features in a single bag of features