Supervised Set-to-Set Hashing in Visual Recognition

I-Hong Jhuo

CODAIT, **IBM**

Outline

- Introduction
- Set-to-Set Hashing
- Results
- Summary

Outline

- Introduction
- Set-to-Set Hashing
- Results
- Summary

Point to Point (P2P) Hashing



Introduction

- P2P (ANN) search characteristics:
 - Computational time (sublinear search complexity)
 - Storage space

What are the unsolved problems ?

 How to measure distance between sets
 Is there any useful information

Outline

- Introduction
- Set-to-Set Hashing
- Results
- Summary

Set to Set (S2S) Hashing

Image set, ID2

Query image set

Video Database



S01E04: 03'42"



Image set, ID1



S01E06: 05'22"



S01E02: 04'20"





Image set, ID3



S01E02: 00'21"

Supervised Information



Image Set Hashing



Reference:

- Deep Image Set Hashing, [arXiv:1606.05381v2, 2016]
- End-to-End image Set hashing, but No structural Information

Introduction

• What are we proposing?

Set-to-Set Hashing => binary codes for a set

 Multiple information => statistical & structural information to improve the retrieval performance

o State-of-the-art results

Multiple Information for Image Set Hashing

Statistical information:
 Gaussian-logarithm kernel

Structural information:

 Graph kernel to capture topological (clique decomposition) similarity



Structural Information and Similarity Measure



- Structural information: $K_g(\mathbf{x}_i, \mathbf{x}_j) = \frac{\sum_{p=1}^{n_i} \sum_{q=1}^{n_j} A_{ip} A_{jq} g(x_{ip}, x_{jq})}{\sum_{r=1}^{n_i} A_{ip} \sum_{q=1}^{n_j} A_{iq}}$
- Statistical information: $K_s(\mathbf{x}_i, \mathbf{x}_j) = \phi(C_i)^\top \phi(C_j) \\ = exp(-||log(C_i) log(C_j)||_F^2/2\gamma_s^2)$

Structural Information and Similarity Measure



- Structural information: $K_g(\mathbf{x}_i, \mathbf{x}_j) = \frac{\sum_{p=1}^{n_i} \sum_{q=1}^{n_j} A_{ip} A_{jq} g(x_{ip}, x_{jq})}{\sum_{n=1}^{n_i} A_{ip} \sum_{q=1}^{n_j} A_{jq}}$
- Statistical information: $K_s(\mathbf{x}_i, \mathbf{x}_j) = \phi(C_i)^\top \phi(C_j) \\ = exp(-||log(C_i) log(C_j)||_F^2/2\gamma_s^2)$

Structural Information and Similarity Measure



- Structural information: $K_g(\mathbf{x}_i, \mathbf{x}_j) = \frac{\sum_{p=1}^{n_i} \sum_{q=1}^{n_j} A_{ip} A_{jq} g(x_{ip}, x_{jq})}{\sum_{r=1}^{n_i} A_{ip} \sum_{q=1}^{n_j} A_{iq}}$
- Statistical information: $K_s(\mathbf{x}_i, \mathbf{x}_j) = \phi(C_i)^\top \phi(C_j) \\ = exp(-||log(C_i) log(C_j)||_F^2/2\gamma_s^2)$

The Mapping Function

Weak Learners:

$$\begin{split} f(\mathbf{x}) &= \mathrm{sign}(K_m(\mathbf{x}_a,\mathbf{x}) - K_m(\mathbf{x}_b,\mathbf{x}) + \varepsilon) \\ & \text{positive} & \text{The } \mathbf{m}_{\mathrm{th}} \, \mathrm{kernel} & \text{negative} \end{split}$$

A strong split (bit):

$$F(\mathbf{x}) = \operatorname{sign}(\sum_{t=1}^{T} \lambda^{t} f^{t}(\mathbf{x}))$$

The best weak learner at iteration t

Image Set Hashing via Multiple Information

Problem Formulation



Experiment Settings

- Dataset:
 - CIFAR-10: 60, 000 images are randomly selected to 195 image sets for training process; 100 image sets for test query, the left images are around 50,000 for test database (around 1700 image sets)
 - Big Bang Theory video (image set) benchmark (BBT): 3,341 videos with 12 characters, 100 image sets for training, 100 image sets for test query, 3041 image sets for test database

Retrieval Performance on

CIFAR-10

		Our Method					
bits	LSH [13]	SH [49]	SSH [41]	KLSH [16]	KSH [21]	HER [18]	ISH
8bits	0.1063	0.1279	0.1165	0.1067	0.2363	0.2284	0.2822
12bits	0.1073	0.1330	0.1416	0.1210	0.2486	0.2672	0.3069
24bits	0.1086	0.1317	0.1512	0.1420	0.2680	0.2772	0.3159
32bits	0.1194	0.1322	0.1574	0.1501	0.2818	0.2937	0.3292
48 bits	0.1105	0.1352	0.1629	0.1622	0.3003	0.3154	0.3334

Performance on CIFAR-10. KLSH, KSH, HER, ISH use image sets as input



Performance on BBT TV-series

		TV drama: the Big Bang Theory					
Method	8 bits	16 bits	32 bits	64 bits	128 bits		
LSH [13]	0.2109	0.2086	0.2092	0.1963	0.1994		
ITQ [10]	0.2935	0.3025	0.2989	0.3029	0.3060		
SH [49]	0.2377	0.2652	0.2665	0.2623	0.2673		
DBC [30]	0.4489	0.4495	0.4235	0.4005	0.3867		
SSH [42]	0.2716	0.2855	0.2662	0.2584	0.3003		
MM-NN [25]	0.3752	0.3955	0.4664	0.5124	0.4922		
KLSH [16]	0.2450	0.2498	0.2381	0.2256	0.2325		
KSH (point)	0.4090	0.4366	0.4454	0.4567	0.4604		
KSH [21] (set)	0.4590	0.4619	0.4534	0.4685	0.4631		
HER [18]	0.4606	0.5049	0.5227	0.5490	0.5539		
ISH ⁰	0.4833	0.5279	0.5359	0.5501	0.5712		
ISH	0.5018	0.5592	0.5864	0.6007	0.6280		

Performance on BBT. KLSH, KSH, HER, ISH⁰, ISH use image sets as input





 S2S provides the promising applications/results comparing to the traditional hashing methods (P2P)

 S2S hashing improves the retrieval performance by utilizing multiple (structural and statistical) information

- State-of-the-art results in retrieval datasets
 - No complex deep learning frameworks

MAX for Developers



IBM <u>D</u>ata <u>A</u>sset e<u>X</u>change (**DAX**)

- Curated free and open datasets under open data licenses
- Standardized dataset formats and metadata
- Ready for use in enterprise AI applications
- Complement to the Model Asset eXchange (MAX)

Data Asset eXchange ibm.biz/data-asset-exchange Model Asset eXchange ibm.biz/model-exchange



(Natural Language Processing) (Language Modeling) (Contracts

CDLA-Sharing | CSV | JSON

This dataset consists of raw and processed execution

logs generated from two versions of Nutch, an open

source web crawler application.

Nutch

CDLA-Sharing | CoNLL-U Finance Proposition Bank

Text from approximately 1000 english sentences obtained from IBM's public annual financial reports, annotated with a laver of "universal" semantic role A dataset of online discussion threads crawled from Ubuntu Forums, with associated subjectivity labels.

CC BY-SA 4.0 | XML

Forum Subjectivity



CODAIT, IBM