Multimedia Analysis:

Marriage of Signal Processing and Machine Learning

Tsuhan Chen 陈祖翰 Professor, Carnegie Mellon University tsuhan@cmu.edu



A 10-Year Journey...

- IEEE Multimedia Signal Processing (MMSP) Technical Committee, 1996~
- MMSP Workshops
 - Princeton 1997, Los Angeles 1998, Copenhagen 1999, Cannes 2001, St. Thomas 2002, Siena 2004, Shanghai 2005, Victoria 2006
- International Conf on Multimedia (ICME)
 - New York 2000, Tokyo 2001, Lausanne 2002, Baltimore 2003, Taipei 2004, Amsterdam 2005, Toronto 2006, Beijing 2007
- *IEEE Transactions on Multimedia*, March 1999~
 - Special issues: networked multimedia 2001, multimedia database 2002, multimodal interface 2003, streaming media 2004, MPEG-21 2005
- SPS Distinguished Lecturer 2007



Carnegie Mellon

Multimedia Analysis: Bits vs. Content



[Baker/Kanade]

It's Bayesian in machine learning, a.k.a., *prior*... Number of all possible 16×12 images = $2^{16\times12\times8}$ >> $30\times60\times60\times24\times365\times$ human history×world population >> number of all possible face images

Thoughts

 "The most compelling shapes are those near to our hearts: people's faces, a gracefully moving body, a natural scene with rustling leaves and flowing water. Evolution has tuned us to these sights..."

[Lengyel, 1998]

 Multimedia analysis... more than processing bits... it's all about the *content* ... it's signal processing + machine learning [Chen, 2006]



Content-Based Information Retrieval

Many Interesting Applications...



Carnegie Mellon

LOGOS [Leung&Chen ICME'02]





Hand-Drawn Sketches



3D Objects [Zhang&Chen ACM MM'01]



ModelRetrieval - C:\ChaZhang\demo\newtest1_0600.mdf

Database View Adjust Experiment Options Help



Carnegie Mellon

3D Protein Structures [Chen&Chen ICIP'02]



Content-Based Information Retrieval (CBIR)



Advanced Multimedia Processing Lab

High-Level Semantics



Possible Solutions



Advanced Multimedia Processing Lab

Semantic Information

- Hidden annotation
 - Object #*n* has Attribute #*k* \rightarrow explicit
- Relevance feedback
 - Objects #m and #n are (not) similar \rightarrow implicit
- Q: How to represent and propagate semantic information?
- Q: How to use explicit/implicit semantic information to improve retrieval?



Semantic Information as Probabilities

	Attribute 1	Attribute 2	Attribute 3	 Attribute <i>K</i>
Object 1	ρ_{11}	$p_{_{12}}$	P_{13}	 P_{1K}
Object 2	<i>P</i> ₂₁	<i>p</i> ₂₂	<i>P</i> ₂₃	 <i>р</i> _{2К}
•••	:	• • •	•	•
Object <i>N</i>	$\rho_{_{N1}}$	p_{N2}	p_{N3}	 P _{NK}

 p_{nk} : Attribute Probabilities



Annotate one object...

	Attribute 1	Attribute 2	Attribute 3	 Attribute <i>K</i>
Object 1				
Object 2	1	0	1	0
•••	:	• •	• •	:
Object <i>N</i>				

When an object is annotated, p_{nk} is set to 0/1

Q: How to propagate? A: Based on low level features

Advanced Multimedia Processing Lab

Carnegie Mellon

Semantic Propagation





Carnegie Mellon

Semantic Propagation (cont.)





Semantic Propagation (cont.)

	Attribute 1	Attribute 2	Attribute 3	 Attribute <i>K</i>
Object 1	$ ho_{11}$	$p_{_{12}}$	p_{13}	 p_{1K}
Object 2	1	0	1	0
•••	:	•	•	:
Object <i>N</i>	$\rho_{_{N1}}$	p_{N2}	p_{N3}	 P _{NK}

Q: Which to annotate next?





- Choose the most uncertain object to annotate
 - Uncertainty determined by the entropy of attribute probabilities
- "Selective sampling"
 - May want to consider density in feature space too

Recap...

- Maintain attribute probabilities of each object
- Set an attribute probability to 1/0 when annotated
- Propagate probabilities to non-annotated objects
- Choose the most uncertain object in the database to annotate next
 - Use probabilities to estimate uncertainty
- Use probabilities to measure semantic distance...







Relevance Feedback

- Relevance feedback
 - Ask for user's feedback during the retrieval
 - "Object #i is (not) similar to the query"
 - "Objects #m and #n are (not) similar"
 - Implicit semantic information
- Use feedback to improve retrieval
 - Way 1: Move the query point
 - Way 2: Weigh the features
 - Way 3: "Warp" the feature space



An Example





An Example





Move the Query Point





Feature Weighting





Feature Space Warping





Feature Space Warping



$$v_{pi} = \left[\gamma \sum_{j=1}^{M} u_i \exp\left(-c \left|v_{ij}\right|\right) \right] v_{iq}$$

This is also semantic propagation!!!

Advanced Multimedia Processing Lab

Experiment Result [Bang&Chen ICIP'02]



Advanced Multimedia Processing Lab

Semantic Propagation is the Key

- Without semantic propagation, hidden annotation and relevance feedback are not very useful
- With enough relevance feedback, can we can accomplish information retrieval without lowlevel features at all?



Pushing Content to Extreme

Content-Free Information Retrieval...



Content-Free Information Retrieval (CFIR)

- With enough relevance feedback, retrieval is based more and more on feedback, less and less on features
- In the extreme case, retrieval based on feedback only
 - Retrieval based on user history
- e.g., Amazon.com



Carnegie Mellon

Example -- How CFIR Works

User History				
	1	0	0	1
	0	1	1	0
	0	1	1	?



Carnegie Mellon

Example -- How CFIR Works



CBIR vs. CFIR

- Will user *U* like image *X*?
- Two different approaches:
 - Look at what U likes
 - \rightarrow Characterize images \rightarrow Content-based IR
 - Look at which users like X
 - \rightarrow Characterize users \rightarrow Content-free IR





Experiment Results [Liu&Chen ICASSP'05]



Content without User Feedback

Extracting content from nothing...



Unsupervised Image Categorization



[Caltech face + background dataset]

Unsupervised Image Categorization



[UIUC car dataset]

"Bag of Words" Representation



DoG interest point detector + SIFT descriptor [Lowe]



Codebook



Graphical Model



Need to handle background...





Solution: Add a hidden layer \rightarrow PLSA

Probabilistic Latent Semantic Analysis (PLSA)



- Hofmann 01, Monay and Gatica-Perez
 04, Sivic et al. 05, Quelhas et al. 05
- Can model complex scenes

Document **Chpia** Appsticance

P(d, z, w) = P(d)P(z|d)P(w|z)

Inferance: Maximum likelihood
 estimationrizinion EMPa(gorith)

- Segmentation P(z|d,w)

d : image z : topic

w : word

Problem with PLSA



"Bag of Words" is the problem...

• As long as the parts are present, the exact position does not matter too much



Picasso, 1943

Dali, 1936







Not so for general objects!

Enforcing Clustering

[Liu&Chen ICIP'06]

• A number of S = 10 fixed spatial distributions





 $s_1 \, t_0 \, s_9$

Enforcing Clustering: Semantic-Shift



[Liu&Chen CVPR'06]

- d : image
- z : topic
- w : word appearance
- x : word position

Document Character Approximation Semantics p(d, z, w, x) = P(d) P(z|d) P(w|z) p(x|z, d)

Representing Location Semantics



- Assume single foreground object
- Location semantics p(x|z,d)
 - Foreground:

$$p(x|z_{\mathsf{FG}}, d_i) \equiv \mathcal{N}(x|\mu_i, \Sigma_i)$$

– Background:

Complement of foreground distribution



$$\begin{array}{c} \text{Learning in Semantic-Shift} \\ \hline P(z_k | d_i, w_j, x_p^{d_i}) \propto p(x_p^{d_i} | z_k, d_i) P(z_k | d_i) P(w_j | z_k) \\ posterio \\ \hline \mathcal{N}(x | \mu_i, \Sigma_i) \\ \hline$$

Learning in Semantic-Shift

p(x|z,d)Location semantics



This is why "semantic-shift"

P(z|d, w, x)

posterior



P(w|z)

Topic appearance



Learning all 3 terms simultaneously... Completely unsupervised...

Results

[Liu&Chen CVPR'06]





0.9 → 0.98

Future Work

- Intra doject modeling
- Video: spatial temporal modeling
- Training the codewords in the loop
- Multiple ojects

Conclusions

- Machine learning can bridge the gap between low-level features (bits) and high-level semantics (content)
- Hidden annotation and relevance feedback can help; semantic propagation is the key
 - Active learning
- "Content-free" information retrieval is possible
 - Bayesian framework
- Content extraction without user feedback is possible
 - Unsupervised learning; graphical models



Afterthoughts...

- Feng-Shui (风水)
 - Ancient Chinese room arrangement technique
- Way 1 (low-level):
 - Write down all the rules
 - Too many and do not generalize
- Way 2 (high-level):
 - Imagine how a dragon would move through the room to arrange it in a livable manner
 - Intuitive and creative
 - Done by some Feng-Shui masters







Advanced Multimedia Processing Lab

Please visit us at: http://amp.ece.cmu.edu

