



Semantic Image Segmentation and Web-Supervised Visual Learning

Florian Schroff Andrew Zisserman University of Oxford, UK

Antonio Criminisi

Microsoft Research Ltd, Cambridge, UK

• • Outline

• Part I: Semantic Image Segmentation

- Goal: automatic segmentation into object regions
- Texton-based Random Forest classifier

• Part II: Web-Supervised Visual Learning

- Goal: harvest class specific images automatically
 - Use text & metadata from web-pages
 - Learn visual model
- Part III: Learn segmentation model from harvested images

Goal: Classification & Segmentation



Image

Classification/Segmentation



Challenges in Object Recognition







Intra-class variations:
 appearance
 differences/similarities among
 objects of the same class

o Inter-class variations:

appearance differences/similarities between objects of different classes

o Lighting and viewpoint

Importance of Context



Oliva and Torralba (2007)

- Context often delivers important cues
- Human recognition heavily relies on context
- In ambiguous cases context is crucial for recognition

System Overview

- Treat object recognition as supervised classification problem:
 - Train classifier on labeled training data

f∩r

- Apply to new unseen test images
- Feature extraction/description
 - Crucial to have a discriminative feature representation



feature extraction classifier (SVM, NN, Random Forest)

training

images

Part I: Image Segmentation

Supervised classification problem: Classify each pixel in the image



Image Segmentation

- Introduction to textons and single-class
 histogram models (SCHM)
- Comparison of nearest neighbour (NN) and Random Forest
- Show strength of Random Forests to combine multiple features

Background: Feature Extraction



Background: Texton Vocabulary



Map Features to Textons



Training Images

Feature Vectors per pixel

Map to textons (pre-clustered)

Resulting texton-maps

Texton-Based Class Models

- Learn texton histograms given class regions
- Represent each class as a set of texton histograms
- Commonly used for texture classification (region ⇔ whole image)

(Leung&Malik ICCV99, Varma&Zisserman CVPR03,

Cula&Dana SPIE01, Winn et al. ICCV05)



Exemplar based class models (Nearest Neighbour or SVM classifier)



Model each class by a single model! (Schroff *et al.* ICVGIP 06) (rediscovered by Boiman, Shechtman, Irani CVPR 08)

(SHCM improve generalization and speed)

Pixelwise Classification (NN)



Sheep model

Kullback-Leibler Divergence: Testing

- KL does not penalize zero bins in the test histogram which are non-zero in the model histogram
- Thus, KL is better suited for singlehistogram class models, which have many non-zero bins due to different class appearances
- This better suitability was shown by our experiments



$$D_{KL}(\mathbf{h}||\mathbf{q}) = \sum \mathbf{h}_i \log(\frac{\mathbf{h}_i}{q_i})$$



Combine Single Histogram Class Model and Random Forest

Random Forest (Training)

• During training each node "selects" the feature from a precompiled feature pool that optimizes the information gain





• Combination of independent decision trees

• Emperical class posteriors in leaf nodes are averaged

- Kleinberg, Stochastic Discrimination 90
- Amit & Geman, Neural Computation 97; Breiman 01
- Lepetit & Fua, PAMI06; Winn et al, CVPR06; Moosman et al., NIPS06





window size\offset	0	10	20	30
25:25	76.2%	78.8%	81.0%	82.2%
10:35	77.7%	79.3%	80.8%	82.0%

• Learning of offset and rectangle shapes/sizes, as well as the channels improves performance

More Feature Types

RGB



Textons





Pixel to be classified

Weighted sum of textons

Difference of HOG responses

- Compute differences over various responses (RGB, textons, HOG)
- Use difference of rectangle responses together with a threshold as node-test $t_p < \lambda$?

Feature Response: Example



- Example of centered rectangle response:
 - Red-channel
 - Green-channel
 - Blue-channel
- Example of rectangle difference (red- and green-channel)

Features: HOG Detailed



 Each pixel is discribed by a "stacked" hog descriptor with different parameters

 Difference computed over responses of one gradient bin with respect to a certain normalization and cellsize

c=cellsize

Importance of different feature types







HOG

HOG & RGB





HOG & RGB

Importance of different feature types







HOG

RGB





HOG & RGB

Importance of different feature types

RGB









HOG & RGB



HOG



 Use global energy minimization instead of maximum a posteriori (MAP) estimate

Image Segmentation using Energy Minimization

Conditional Random Field (CRF) • energy minimization using, e.g. Graph-Cut or TRW-S



 $\mathbf{g}_{ij} \in \mathbb{R}^d$ Colour difference vector

c_i = binary variable representing label ('fg' or 'bg') of pixel i



Labelling problem

Graph Cut

CRF and Colour-Model



- CRF as commonly used (e.g. Shotton et al. ECCV06: TextonBoost)
- TRW-S is used to maximize this CRF
- Perform two iterations: one with one w/o colour model

MSRC-Databases

9-classes: building, grass, tree, cow, sky, airplane, face, car, bicycle

120 training-120 testimages

Similar: 21-classes



Images

Groundtruth

Images

Groundtruth

Segmentation Results (MSRC-DB) with Colour-Model



w/o CRF Class posteriors only





Segmentation Results (MSRC-DB) with Colour-Model

Classification



Image



Classification Quality









Segmentation Results (MSRC-DB 21 classes)

Classification

Image overlay

Classification Quality



CRF













21-class MSCR dataset



VOC2007-Database

20 classes:

Aeroplane Bicycle Bird **Boat Bottle** Bus Car Cat Chair Cow Diningtable Dog Horse **Motorbike** Person Pottedplant Sheep Sofa Train **Tymonitor**









Images









Groundtruth













Images





Groundtruth

VOC 2007





	textons	RGB	HOG	HOG,RGB	HOG,RGB,	HOG,RGB,	CRF	Others
					F17	SHCMs		
9-class (pixelwise)	72.8%	72.2%	75.9%	84.5%	86.0%	84.9%	87.3%	84.9% [1]
21-class (pixelwise)	40.3%	54.0%	56.3%	69.9%	71.5%	71.7%	73.7%	72.7% [2]
VOC2007 (class avg.)	13.6%	10.1%	17.4%	17.7%	19.7%	21.4%	24.4 %	20.0% [3]

[1] Verbeek et al. NIPS2008; [2] Shotton et al. ECCV2006;[3] Shotton et al. CVPR 2008 (raw results w/o image level prior)

Combination of features improves performance

 CRF improves performance and most importantly visual quality

• • • Summary

- Discriminative learning of rectangle shapes and offsets improves performance
- Different feature types can easily be combined in the random forest framework
- Combining different feature types improves performance

Part II: Web-Supervised Visual Learning

• Goal: retrieve class specific images from the web

- No user interaction (fully automatic)
- Images are ranked using a multi-modal approach:
 - Text & metadata from the web-pages
 - Visual features
- Previous work on learning relationships between words and images:
 - Barnard et al. JMLR 03 (Matching Words and Pictures)
 - Berg et al. CVPR 04, CVPR 06



learn text ranker once





Text&Metadata Ranker

Why don't we start with Google image search?
Limited return (only 1000 images)

- Goal: object class independent ranker
- Rank images using Bayes model on binary feature vector:

a=(context10, context50, filename, filedir, imagealt, imagetitle, websitetitle)

$$P(y|\mathbf{a}) = P(a_1, \dots, a_4|y) \prod_{j=1}^{7} P(a_i|y) \cdot P(y)/Z$$



Visual Ranking

How to learn visual model from these noisy images?

• Where do we get the training data from?



• Train on top text ranked images \rightarrow positive data

- Randomly sample images \rightarrow negative data
- Support Vector Machine (SVM)
 - robust to noise

Filter drawings & abstract Images drawing&symbolic





With 7.3 million unique visitors per month as audited by the ABC, SFGate is the leading news and information Web site for the tian Francisco flay Area. Reflecting the diverse spirit of the region 50 Galo dolivers the most up to the minute stories, in-depth special reports, unboatable local sports coverage, the best regional listings and cutting edge anterbainment coverage

























o Gradient- & colour-histograms o RBF-SVM





- 400 visual-words from four interest point detectors
- HOG descriptor to represent shape
- o RBF-SVM on "stacked" feature vector

Example: Penguin

- 1. Enter "penguin"
- 2. Retrieve images from web pages returned by Google web search on **penguin**
 - 522 in-class, 1771 non-class



- 3. Remove drawings & abstract images
 - 391 in-class, 784 non-class



Example: Penguin continued

4. rank images using naïve Bayes metadata ranker



5. Train SVM on visual features using ranked images as noisy training data

6. Final re-ranking using trained SVM







Text+visual ranked images

- Text ranker:
 - rank images for new requested object-class
- Visual ranker:
 - Train visual classifier and re-rank images



airplane: HOG+BOW t+v

Examples continued



shark: HOG+BOW t+v

Examples continued





giraffe: HOG+BOW t+v



zebra: HOG+BOW t+v



car: re-ranking Google images using our visual re-ranking



leopard: re-ranked Berg (2006) images

• • • Summary

- Use object-class independent text ranker to retrieve training data
- Train visual classifier on top text ranked images
- Show applicability on different datasets
 - Google image search
 - Berg et al. (Animals on the Web)

Part III: Segmentation from Harvested Images

- Random Forest pixelwise classification
- Use weak supervision
 - No segmented training data
 - Per image classlabels are used



- Segment images in 21-class MSRC dataset
 - Weak supervision: 52.1% (w/o CRF)
 - Strong supervision: 71.5% (w/o CRF)

(following images with CRF)





















image











weak supervision





strong supervision

Learn Segmentation Model

- Train Random Forest on top ranked
 100 car images and
 200 randomly sampled background images
- Segment images in 21-class MSRC dataset (using CRF with colour-model)



car: HOG+BOW t+v



















automatic segmentation

image

 $\operatorname{groundtruth}$

• • • Summary

- Show that Random Forest can be trained on weakly labelled training data
- Combine strong Random Forest segmentation with unsupervised visual learning
- This allows learning of segmentation models w/o requiring manually labeled training data

Discussion & Future Work

- Image level class priors (Shotton *et al.* CVPR08) can improve performance dramatically
- Incorporate a more global shape into the decision trees
- Hierarchy of trees
 - Top trees classifying interesting image subareas
 - Subsequent trees perform fine grained segmentation