On the use of Regions for Semantic Image Segmentation

Rui Hu¹, <u>Diane Larlus</u>², Gabriela Csurka² ¹ University of Surrey, UK ² Xerox Research Center Europe

> Tuesday December 18th 2012 ICVGIP 2012



Table of contents



- 2 Proposed Benchmark
- 3 Recognition model
- 4 Image level prior
- **5** Conditional Random Field
- 6 Conclusions



• Unsupervised Image segmentation







• Unsupervised Image segmentation





• Semantic Image segmentation







- State-of-the-art semantic segmentation methods usually leverage
 - Local appearance of objects (class likelihood maps)
 - Local consistency (constraining neighboring labels)
 - Global consistency (image level priors)
- That are combined in
 - unified CRF framework [Verbeek & Triggs 2007, Kohli et al 2009, Ladicky et al 2009]
 - sequential framework

[Yang et al 2007, Csurka & Perronnin 2008]



Recent methods use unsupervised image partition in Regions, or Super-Pixels to enhance semantic segmentation:

- Local appearance is predicted based on region descriptors [Gu et al 2009, Lim et al 2009, Vijayanarasimhan & Grauman 2011, Lucchi et al 2011]
- Local consistency is enforced within regions:
 - in a post processing step [Csurka and Perronnin 2008]
 - or using higher order potentials in the CRF [Kohli et al 2009, Ladicky et al 2009, Gonfaus et al 2010]



What is the best way to use regions ?

We propose a benchmark studying the role and benefit of regions at different stages of the segmentation process.



Table of contents

Semantic Segmentation

2 Proposed Benchmark

3 Recognition model

Image level prior

5 Conditional Random Field

6 Conclusions



We propose a benchmark based on 3 components

• A standard dataset:

MSRC-21 dataset

• A standard super-pixel method:

Berkeley segmentation approach

• A standard pipeline:

Fisher-Vector based patch classification Condition Random Field



MRSC-21 dataset

- Standard benchmark, 591 images:
 - 276 images for training
 - 59 images for validation
 - 275 images for testing
- 21 classes:

building, grass, tree, cow, sheep, sky, aeroplane, water, face, car, bicycle, flower, sign, bird, book, chair, road, cat, dog, body, boat

- Evaluate pixel-level classification
 - Average class-based accuracy

[Shotton et al, IJCV 2009]

































Patch-based Fisher Vector Representation



[Csurka and Perronnin, IJCV 2011]

Dense patch extraction at single scale or at 5 different scales, described using:



Conditional Random Field model

- Dense CRF model
 - [Krähenbühl and Koltun, NIPS 2011]
 - Model with unary and pairwise potentials
 - Unary term: based on the patch-based FV classification
 - Pairwise term: all pairwise pixel connections are considered (not only 4 or 8 neighborhood systems)



Table of contents

Semantic Segmentation

- 2 Proposed Benchmark
- 3 Recognition model
- Image level prior
- **5** Conditional Random Field
- 6 Conclusions



Appearance model

- Patch-based system: PB-SIS
 - Classify each patch individually
 - Accumulate patch probabilities at the pixel level



- Region-based system: RB-SIS
 - · Aggregation of patches for each region of the hierarchy
 - Classify each region individually
 - Accumulate region information at the pixel level



Recognition model

Appearance only

- Patch-based semantic image segmentation: PB-SIS
- Region-based semantic image segmentation: RB-SIS

| | One scale (1S) | | Multi scale (MS) | |
|----------------|----------------|--------|------------------|--------|
| | PB-SIS | RB-SIS | PB-SIS | RB-SIS |
| COL | 55.72 | 62.84 | 62.31 | 65.94 |
| SIFT | 46.10 | 61.98 | 54.29 | 65.44 |
| APP (COL+SIFT) | 63.63 | 70.24 | 69.98 | 72.90 |

Regions are great assets that improve local appearance based prediction.



Exploiting the shape and the hierarchy of regions

For RB-SIS using regions, we can:

use gPb as shape descriptor

[Gu et al CVPR 2009, Lim et al ICCV 2009]

• exploit partially the hierarchy through Bags-of-Triplets





Exploiting the shape and the hierarchy of regions

RB-SIS: shape and bags-of-triplets

| | shape only | +APP(1S) | +APP(MS) |
|-----------|------------|----------|----------|
| BoR | 34.77 | 70.35 | 71.85 |
| BoR + BoT | 42.70 | 71.18 | 72.99 |

- Shape alone performs poorly
- Hierarchy helps a lot for shape alone, but less when appearance is present



Table of contents

Semantic Segmentation

- 2 Proposed Benchmark
- 3 Recognition model
- 4 Image level prior
- **5** Conditional Random Field
- 6 Conclusions



Appearance based predictions are combined with

- Global image classification (global Fisher Vector + SVM)
- Location prior (object location likelihood prior from training)

| | REC | + GL |
|--------|-------|-------|
| PB-SIS | 69.98 | 75.20 |
| RB-SIS | 72.99 | 75.88 |

Recognition (REC) is enhanced with global and location (GL) priors



Table of contents

Semantic Segmentation

- Proposed Benchmark
- **3** Recognition model
- 4 Image level prior
- **5** Conditional Random Field
- 6 Conclusions



Conditional Random Field (CRF)

We use a dense CRF formulation

- unary potential: best recognition model enhanced with global and location priors
- pairwise potential: all pixel pairs are connected with pairwise
 - middle range regularization
 - longer range color-dependent regularization





Conditional Random Field

We extend the dense CRF to use region information

- unary potential: best recognition model
- pairwise potential
 - middle range regularization
 - longer range color-dependent regularization
 - additional potential using leaf regions





Conditional Random Field

• Dense CRF results without (dCRF) and with (dCRFSP) region-based regularization

| | REC | + GL | dCRF | dCRFSP |
|--------|-------|-------|-------|--------|
| PB-SIS | 69.98 | 75.20 | 76.69 | 77.25 |
| RB-SIS | 72.99 | 75.88 | 75.80 | 76.02 |

- CRF regularization brings little improvement to RB-SIS
- PB-SIS benefits more from CRF, and outperforms RB-SIS



Qualitative results

test image - groundtruth - PB-prior - RB-prior - PB-dCRFSP - RB-dCRFSP





Table of contents

Semantic Segmentation

- 2 Proposed Benchmark
- 3 Recognition model
- 4 Image level prior
- **5** Conditional Random Field





Conclusions

Proposed framework allows to evaluate the contribution of each component

Take Home Message:

- Simple recognition model using regions and global prior is already very competitive, no need for regularization
- When a CRF is considered, the patch-based model is enough, and regions could be used only at a later stage



Thanks for your attention ! Questions ?



Backup-slides



Main limitation of an image partitioned into regions:

- No possible recovery if a region groups multiple classes. Possible solutions:
 - Multiple segmentation to obtain overlapping sets of regions [Pantofaru et al 2008, Gould et al 2009]
 - Exploiting a hierarchy of regions

[Ladicky et al 2009, Gu et al 2009, Lim et al 2009, Munoz et al 2010]

Graph of regions

[Chen et al 2011]



Patch-based Fisher Vector Representation

No regularization: simple patch voting



Conditional Random Field (CRF)

We use a dense CRF formulation:

• CRF based regularization: dCRF

$$E(\mathbf{x}) = \sum_{i} \psi_{u}(\mathbf{x}_{i}) + \sum_{i < j} \delta_{\mathbf{x}_{i},\mathbf{x}_{j}} \psi_{p}(\mathbf{x}_{i},\mathbf{x}_{j}),$$

Pairwise potential

$$\begin{split} \psi_{p}(x_{i}, x_{j}) &= \omega_{1} \exp\left(-\frac{|p_{i} - p_{j}|^{2}}{2\theta_{\alpha}^{2}} - \frac{|I_{i} - I_{j}|^{2}}{2\theta_{\beta}^{2}}\right) \\ &+ \omega_{2} \exp\left(-\frac{|p_{i} - p_{j}|^{2}}{2\theta_{\gamma}^{2}}\right) \end{split}$$

with p_i and l_i being the position and RGB value of pixel x_i respectively.

Conditional Random Field

• CRF based regularization: dCRF





Conditional Random Field (CRF)

We extend the dense CRF to use region information:

• CRF based regularization: dCRFSP

$$E(\mathbf{x}) = \sum_{i} \psi_{u}(\mathbf{x}_{i}) + \sum_{i < j} \delta_{\mathbf{x}_{i},\mathbf{x}_{j}} \hat{\psi}_{\rho}(\mathbf{x}_{i},\mathbf{x}_{j}),$$

Pairwise potential

$$\hat{\psi}_{p}(x_{i}, x_{j}) = \psi_{p}(x_{i}, x_{j}) + \omega_{3} \exp\left(-\frac{|p_{i} - p_{j}|^{2}}{2\theta_{\alpha}^{2}} - \frac{|R_{i} - R_{j}|^{2}}{2\theta_{\delta}^{2}}\right)$$

with position p_i , RGB value of pixel I_i and the leaf region that contains x_i , R_i .



Conditional Random Field

• CRF based regularization: dCRFSP



