Fisher and VLAD with FLAIR Koen E. A. van de Sande¹², Cees G. M. Snoek¹², Arnold W. M. Smeulders¹²³ ¹University of Amsterdam, ²Euvision Technologies, ³Centrum Wiskunde & Informatica

Object Recognition

Object recognition seeks to answer 2 questions:

- What is it?
- Where is it?



Box Encodings

State-of-the-art encodings:

- VLAD [JegouTPAMI12]
- Fisher [PerronninECCV10]
- Current state-of-the-art for recognition:
- VOC'10 [CinbisICCV13]
- VOC'12 [SandeVOC12]
- ImageNet-DET'13 (this work)
- Drawback: Computationally expensive

Key Problems



Large boxes are slower to encode than small

Existing approaches:

- Integral images compute the sum of an area in O(1)
- Multidimensional integral images [PorikliCVPR05]
- extend this to histograms
- Approximate at significant accuracy costs

For VLAD encoding, integral histograms require 14GB of memory and minutes to compute per image: unusable in practice.

For Fisher, there are additional complexities with multiple assignment, power and L_2 normalizations

FLAIR solves these problems. It was our secret ingredient to winning the 2013 IM GENET 200 Object Detection Challenge

References

- Cinbis et al. Segmentation driven object detection with fisher vectors (ICCV2013) Jegou et al. Aggregating local image descriptors into compact codes (TPAMI2012)
- Perronnin et al. Improving the fisher kernel for large-scale image classification
- (ECCV2010) Porikli. Integral histogram: A fast way to extract histograms in cartesian spaces (CVPR2005)
- Uijlings et al. Selective Search for Object Recognition (IJCV2013)
- Van de Sande et al. Hybrid Coding for Selective Search (PASCAL VOC2012)
- Wang et al. Regionlets for Generic Object Detection (ICCV2013)

The *k*-th feature vector element:

VLAD

Evaluate VLAD independent of box area: *O(KWHD+KDBS)*

Decomposition over k still possible, even with multiple assignment. Consider:

Goal: Construct feature vector for *B* boxes with *S* spatial pyramid cells N descriptors of length D in the image of size (W,H)

$$\phi(\vec{v}_n, \vec{c}_k) = \begin{cases} 1 & \text{if } k \text{ is the closest codeword for} \\ 0 & \text{otherwise} \end{cases}$$



One descriptor now affects *D* elements of feature vector:

$$\vec{F}_k = \sum_{n=1}^N (\vec{v}_n - \vec{c}_k) \phi(\vec{v}_n, \vec{c}_k)$$

Uses multidimensional integral images



Fisher is the normalized gradient of the loglikelihood under a Mixture-of-Gaussians.

$$\mathbf{k}: \quad \frac{\Delta \ln p}{\Delta \vec{\mu}_k} = \frac{1}{N} \sum_{n=1}^N \frac{p(k|\vec{v}_n)}{\sqrt{\pi_k}} \left(\frac{\vec{v}_n - \vec{\mu}_k}{\vec{\sigma}_k}\right),$$
$$\frac{\Delta \ln p}{\sqrt{\pi_k}} = \frac{1}{N} \sum_{n=1}^N \frac{p(k|\vec{v}_n)}{\sqrt{\pi_k}} \left(\frac{(\vec{v}_n - \vec{\mu}_k)^2}{\vec{\sigma}_k}\right) = \frac{1}{N}$$

$$\frac{\ln p}{\Delta \vec{\sigma}_k} = \frac{1}{N} \sum_{n=1}^{N} \frac{p(k|\vec{v}_n)}{\sqrt{\pi_k}} \left(\frac{(\vec{v}_n - \vec{\mu}_k)^2}{\vec{\sigma}_k^2} - 1 \right)$$

$$S_0(k) = \sum_{n=1}^N p(k|\vec{v}_n), \qquad \vec{S}_1(k) = \sum_{n=1}^N p(k|\vec{v}_n)\vec{v}_n, \qquad \vec{S}_2(k) = \sum_{n=1}^N p(k|\vec{v}_n)\vec{v}_n^2$$

Rewrite using these to:

$$\begin{array}{l} \vdots \quad \frac{\Delta \ln p}{\Delta \vec{\mu}_k} = \frac{1}{N\sqrt{\pi_k}\vec{\sigma}_k} \left(\vec{S}_1(k) - \vec{\mu}_k \cdot S_0(k)\right), \\ \\ \frac{\Delta \ln p}{\Delta \vec{\sigma}_k} = \frac{1}{N\sqrt{\pi_k}} \left(\frac{\vec{S}_2(k) - 2\vec{\mu}_k \vec{S}_1(k) + \vec{\mu}_k^2 \cdot S_0(k)}{\vec{\sigma}_k^2} - S_0(k)\right) \end{array}$$

The scalar integral image $S_0(k)$ and the multidimensional integral images $S_1(k)$ and $S_2(k)$ are supplemented by a scalar integral image holding the number of descriptors N in an area. With these four integral images the gradients for a single codeword evaluate in O(D) and independent of box area, similar to VLAD





UNIVERSITEIT VAN AMSTERDAM



plane	bike	bird	boat	bottle	bus	car	cat	chair
52.4	54.3	13.0	15.6	35.1	54.2	49.1	31.8	15.5
56.2	42.4	15.3	12.6	21.8	49.3	36.8	46.1	12.9
61.3	46.4	21.1	21.0	18.1	49.3	45.0	46.9	12.8
53.3	55.3	19.2	21.0	30.0	54.4	46.7	41.2	20.0
65.0	48.9	25.9	24.6	24.5	56.1	54.5	51.2	17.0
61.3	52.3	27.8	25.7	21.3	54.0	45.6	54.0	15.5
	52.4 56.2 61.3 53.3 65.0 61.3	52.454.356.242.461.346.453.355.365.048.961.352.3	52.454.313.056.242.415.361.346.421.153.355.319.265.048.925.961.352.327.8	52.454.313.015.656.242.415.312.661.346.421.121.053.355.319.221.065.048.925.924.661.352.327.825.7	52.454.313.015.6 35.1 56.242.415.312.621.861.346.421.121.018.153.3 55.3 19.221.030.0 65.0 48.925.924.624.561.352.3 27.825.7 21.3	52.454.313.015.6 35.1 54.256.242.415.312.621.849.361.346.421.121.018.149.353.3 55.3 19.221.030.054.4 65.0 48.925.924.624.5 56.1 61.352.3 27.825.7 21.354.0	52.454.313.015.6 35.1 54.249.156.242.415.312.621.849.336.861.346.421.121.018.149.345.053.3 55.3 19.221.030.054.446.7 65.0 48.925.924.624.5 56.154.5 61.352.3 27.825.7 21.354.045.6	52.454.313.015.6 35.1 54.249.131.856.242.415.312.621.849.336.846.161.346.421.121.018.149.345.046.953.3 55.3 19.221.030.054.446.741.2 65.0 48.925.924.624.5 56.154.5 51.261.352.3 27.825.7 21.354.045.6 54.0



http://koen.me/research/flair