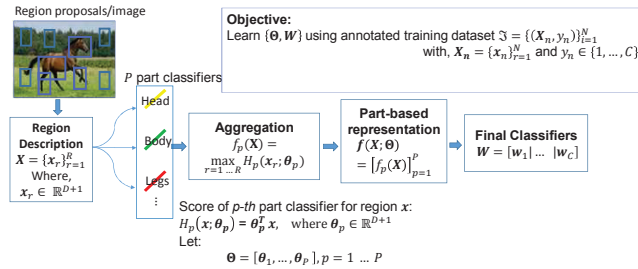


# SPLeaP: Soft Pooling of Learned Parts for Image Classification

## Introduction

### Part-Based Models (PBMs) in the context of Image Classification.



## Motivation

### Motivation for using PBMs:

- By breaking object / image into parts at least some of these are visible and recognizable.
- Parts can be recombined in many ways to construct an object / scene.
- Parts are distinctive of a particular class.
- Part-based representation are compact.

### Challenges:

- How to train the part classifiers ?
- Difficult problem – No annotations of relevant regions within an image.
- Chicken and egg problem:  
Discriminative / relevant regions ↔ Part Classifiers
- Initialization critical.

## Approach

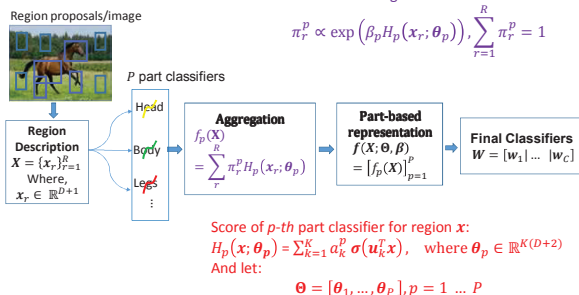
### Novel joint learning framework for learning parameters in PBMs:

- We describe each part classifier as a linear combination of weak non-linear classifiers.
- We introduce a parameter, referred as the "per-part softness pooling coefficient" inside the optimization process.
- Learning is done by a block-wise Stochastic Gradient Descent (SGD).

### Contributions:

#### 1. A boosted non-linear part classifiers.

#### 2. Parametric soft-max Pooling.



**Objective:**  
Learn  $\{\theta, W, \beta\}$  using  $\mathcal{S} = \{(X_n, y_n)\}_{n=1}^N$   
with,  $X_n = \{x_r\}_{r=1}^R, y_n \in \{1, \dots, C\}$  and pooling coefficient,  $\beta = [\beta_p]_{p=1}^P$

### Multi-Class Classification Loss:

Using Logistic regression, class label for an input image  $X$  is predicted as:

$$Pr(y = c | X; \theta, \beta, W) = \frac{\exp(w_c^T f(X; \theta, \beta))}{\sum_{c=1}^C \exp(w_c^T f(X; \theta, \beta))}$$

$$\min_{\theta, \beta, W} - \sum_n \sum_c [y_n = c] \ln Pr(c | X_n; \theta, \beta, W) + \mu \|\theta\|^2 + \delta \|W\|^2$$

## Results

### Experimental Settings:

#### Datasets:

- Pascal-VOC-2007 (20 classes) - animals (e.g., dog, cat), vehicles (e.g., aeroplane, car) and other manufactured objects (e.g., tv monitor, chair).
- MIT Indoor 67 (67 classes) - e.g., nursery, movie theater, casino or meeting room.
- Willow dataset (7 classes) - e.g. play instrument, walk.

#### Performance measure:

- mean Average Precision (mAP) for Pascal-VOC and Willow.
- Accuracy (acc) for MIT dataset.

#### Region proposal scheme: Selective search.

#### Region feature extraction:

- 128-D from 13-layer architecture (VGG-128)
- 4096-D from 16-layer architecture (VD-16)
- Krizhevsky's architecture pre-trained using ImageNet + Places database (HybridCNN).

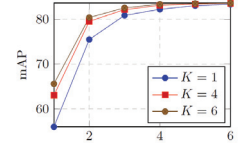
### Comparison with state-of-the-art methods:

Pascal-VOC-2007				MIT indoor67 (left) and Willow (right)			
Methods	mAP	Methods	mAP	Methods	acc	Methods	mAP
VGG128-G	75.5	VD-16-G	81.73	Gong et al.	68.80	Sicre et al.	81.90
Oquab et al.	77.31	Cimpoi et al.	85.1	Li et al.	69.69	VD-16-G	85.12
Li et al.	77.90	SPLeaP-VD-16	<b>88.01</b>	HybridC NN-G	72.54	VD-16 (dense evaluation)	85.67
Cimpoi et al.	79.50	16-layer architecture		Parizi et al.	73.30	SPLeaP - VD-16	<b>88.47</b>
CNN-S fine tuned [5]	82.42			SPLeaP-HybridCN	<b>73.45</b>		
SPLeaP-VGG-128	<b>84.68</b>						
		13-layer architecture					

### 1. Parametric Soft-Max Pooling:

Average pooling	Max Pooling	Cross Valid. $\beta_p = \beta$	Learned $\beta_p$
80.77	83.23	84.31	<b>84.68</b>

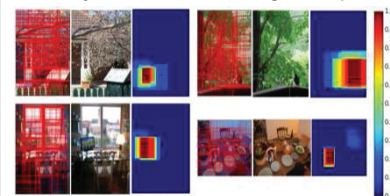
### 2. A boosted non-linear part classifiers:



Plot of test mAP vs number of training epochs

### Specialization of Part Classifiers to discriminative regions:

#### Selectivity of Part Classifiers using heatmaps



Heatmaps for images Pascal-VOC-2007 of classes (clockwise from top-left) "potted plant", "bird", "bottle" and "TV monitor".



Discriminative parts for the four classes (clockwise from top-left) "horse", "motorbike", "dining table", and "potted plant".

## Conclusion

- We introduce SPLeaP, a novel part-based model for image classification.
- Based on non-linear part classifiers combined with part-dependent soft pooling - both being trained jointly with the image classifiers surpasses standard pooling approaches and other PBMs on several challenging classification tasks.
- Our method does not need any particular initialization of the parts.

## References

- Parizi et al.: Automatic Discovery and Optimization of Parts for Image Classification. In: ICLR (2015)
- Li et al.: Mid-level deep pattern mining. In: CVPR (2015).
- Song et al.: Multi-scale orderless pooling of deep convolutional activation features ECCV (2014).
- Cimpoi et al.: Deep filter banks for texture recognition and segmentation. In: CVPR (2015)
- Chaffield et al.: Return of the Devil in the Details: Delving Deep into Convolutional Nets. In: BMVC 2014.
- Oquab et al.: Learning and transferring mid-level image representations using convolutional neural networks.
- Sicre et al.: Discovering and aligning discriminative mid-level features for image classification