

## - Abstract

- ⇒ A framework for semantic segmentation.
- ⇒ Completely unsupervised: Inputs are only images. No labels.
- ⇒ Classes are discovered by clustering of per-pixel features learned by deep networks
- ⇒ Clustering pixel features alone does not work
  - Fail to learn high level concepts
  - Overfits to low-level visual cues
- ⇒ Introduce inductive biases to make it all work out.
  - ① Invariance to photometric transformations
  - ② Equivariance to geometric transformations
- ⇒ Results
  - ① Better results on COCO and Cityscapes datasets.
    - ⇒ Accuracy +17.5
    - ⇒ mIoU +4.5
  - ② Better initialization for supervised learning

## - Glossary

- ⇒ Thing = foreground object
- ⇒ Stuff = background object.

## - Problem Statement

- ⇒ Input = a set of unlabeled images in a domain  $D$ .
- ⇒ Outputs
  - A set of visual classes  $C$
  - A semantic segmentation function  $f$

- A set of visual classes  $C$
  - A semantic segmentation function  $f_\theta$
- $\Rightarrow$  Given an image from  $D$ , assign to each pixel a class in  $C$ .

### - Baseline clustering approach

$\Rightarrow$  Based on DeepCluster <https://arxiv.org/abs/1807.05520>

① Start with an NN that outputs some pixel-level features

② Cluster with current features

③ Use cluster labels as pseudo-labels to train a feature embedder

④ Repeat ② and ③ until satisfied.

$\Rightarrow$  DeepCluster, however, only works at image level, not pixel level

$\Rightarrow$  Let  $\{x_1, x_2, \dots, x_n\}$  be the set of training images.

$\Rightarrow$  For each image  $x_i$ ,  $f_\theta(x_i)$  is a feature tensor containing embedding for each pixel.

$\Rightarrow$  Let  $f_\theta(x_i)[p]$  denote the embedding of pixel  $p$  of image  $x_i$ .

$\Rightarrow$  Let  $g_w(\cdot)$  = a classifier operating on the  $f_\theta(x_i)[p]$ 's.

$\Rightarrow$  The baseline approach alternates between the following two steps.

① Use the current embedding to cluster the pixels by k-means

$$\min_{y, \mu} \sum_{i, p} \|f_\theta(x_i)[p] - \mu_{y, ip}\|^2$$

Here,  $y_{ip}$  = class label of  $x_i[p]$ .

$\mu_k$  = the mean vector associated with the  $k^{\text{th}}$  class.

k-means is computed with minibatch k-means

<https://openreview.net/forum?id=B1EPCL-d-B>

- ② Use the cluster labels to train a pixel classifier using standard cross entropy loss

$$\min_{\theta, w} \sum_{i, p} \mathcal{L}_{CE}(g_w(f_{\theta}(x_i)[p]), y_{ip})$$

$$\mathcal{L}_{CE}((s_1, s_2, \dots, s_k), s_{y_{ip}}) = -\log \frac{e^{s_{y_{ip}}}}{\sum_k e^{s_k}}$$

Here,  $(s_1, s_2, \dots, s_k) \leftarrow g_w(f_{\theta}(x_i)[p])$  is the class scores outputted by  $g_w$ .

### - A modification to the baseline

- $\Rightarrow$  To prevent noisy gradients and noisy clusters, the paper removes the classifier  $g_w$ .
- $\Rightarrow$  Instead, it labels each pixels with the centroid of the nearest cluster.
- $\Rightarrow$  The score of a pixel embedding is the cosine distance from the embedding to the centroid.

$$\mathcal{L}_{clust}(f_{\theta}(x_i)[p], y_{ip}, \mu) = -\log \left( \frac{e^{-d(f_{\theta}(x_i)[p], \mu_{y_{ip}})}}{\sum_l e^{-d(f_{\theta}(x_i)[p], \mu_l)}} \right)$$

$$d(v_1, v_2) = \text{cosine-distance}(v_1, v_2)$$

$$= 1 - \frac{v_1 \cdot v_2}{\|v_1\| \|v_2\|}$$

### - Inductive biases

- $\Rightarrow$  So far, there is no guarantee that the features would be semantic.
- $\Rightarrow$  To get such features, the paper introduces two inductive biases.

① Invariance to photometric transformations: The labels should not

- ① Invariance to photometric transformations: The labels should not change if pixel colors are jittered slightly.
- ② Equivariance to geometric transformations: If the image is warped geometrically, then the labels should undergo the same transformation too.

Mathematically:  $G(P(x)) = G(Y)$

where  $x$  = input color image,

$Y$  = image of labels

$P$  = photometric transformation

$G$  = geometric transformation

- Invariance to photometric transformation, operationally

$\Rightarrow$  For each image  $x_i$ , sample two random photometric transformations  $P_i^{(1)}, P_i^{(2)}$ .

$\Rightarrow$  This yields two feature embeddings:

$$z_{ip}^{(1)} = f_{\theta}(P_i^{(1)}(x_i)) [p]$$

$$z_{ip}^{(2)} = f_{\theta}(P_i^{(2)}(x_i)) [p]$$

$\Rightarrow$  We perform two k-means clustering separately, resulting in two labels and two centroids

$$y^{(1)}, \mu^{(1)} = \operatorname{argmin}_{y, \mu} \|z_{ip}^{(1)} - \mu_{y_{ip}}\|$$

$$y^{(2)}, \mu^{(2)} = \operatorname{argmin}_{y, \mu} \|z_{ip}^{(2)} - \mu_{y_{ip}}\|$$

$\Rightarrow$  To make sure that the features satisfy the clustering property, we minimize the clustering loss:

1. For each image  $x_i$ , we sample a random photometric transform  $P_i^{(1)}$ .

$$\mathcal{L}_{\text{within}} = \mathcal{L}_{\text{clust}}(f_\theta(P_i^{(1)}(x_i))(p), y_{ip}^{(1)}, \mu^{(1)}) \\ + \mathcal{L}_{\text{clust}}(f_\theta(P_i^{(2)}(x_i))(p), y_{ip}^{(2)}, \mu^{(2)})$$

$\Rightarrow$  Because the labels should be invariant to photometric transformation, the following should also be minimized,

$$\mathcal{L}_{\text{cross}} = \mathcal{L}_{\text{clust}}(f_\theta(P_i^{(1)}(x_i))(p), y_{ip}^{(2)}, \mu^{(2)}) \\ + \mathcal{L}_{\text{clust}}(f_\theta(P_i^{(2)}(x_i))(p), y_{ip}^{(1)}, \mu^{(1)})$$

$\Rightarrow$  The two losses, when minimize together, should encourage the network to learn an embedding function that

- ① produce features that will be labelled identically, and
- ② produce the same cluster centroids regardless of the color transformations applied,

- Equivariance to geometric transformations, operationally.

$\Rightarrow$  In addition to the  $P_i^{(1)}$  and  $P_i^{(2)}$ , we sample a random geometric transformation  $T_i$  for each image  $i$ .

$\Rightarrow$  Now, we compute

$$z_{ip}^{(1)} = f_\theta(G_i(P_i^{(1)}(x_i)))(p)$$

$$z_{ip}^{(2)} = G_i(f_\theta(P_i^{(2)}(x_i)))(p)$$

$\Rightarrow$  We then cluster and minimize  $\mathcal{L}_{\text{within}} + \mathcal{L}_{\text{cross}}$  with the formula in the last section,

- Pseudocode

for  $x_i \sim D$  do

$P_i^{(1)}, P_i^{(2)} \sim \text{RandomPhotometricTransforms}$

$P_i^{(1)}, P_i^{(2)} \sim \text{RandomPhotometricTransforms}$

$G_i \sim \text{RandomGeometricTransform}$

$Z_{i,:}^{(1)} \leftarrow G_i (f_\theta(P_i^{(1)}(x_i)))[:, :]$

$Z_{i,:}^{(2)} \leftarrow f_\theta(G_i(P_i^{(2)}(x_i)))[:, :]$

end for

$\mu^{(1)}, \gamma^{(1)} \leftarrow \text{KMeans}(\{Z_{ip}^{(1)}\})$

$\mu^{(2)}, \gamma^{(2)} \leftarrow \text{KMeans}(\{Z_{ip}^{(2)}\})$

for  $x_i \sim D$  do

$\mathcal{L}_{\text{within}} \leftarrow \sum_p (\mathcal{L}_{\text{clust}}(Z_{ip}^{(1)}, \gamma_{ip}^{(1)}, \mu^{(1)}) + \mathcal{L}_{\text{clust}}(Z_{ip}^{(2)}, \gamma_{ip}^{(2)}, \mu^{(2)}))$

$\mathcal{L}_{\text{cross}} \leftarrow \sum_p (\mathcal{L}_{\text{clust}}(Z_{ip}^{(1)}, \gamma_{ip}^{(2)}, \mu^{(2)}) + \mathcal{L}_{\text{clust}}(Z_{ip}^{(2)}, \gamma_{ip}^{(1)}, \mu^{(1)}))$

$\mathcal{L}_{\text{total}} \leftarrow \mathcal{L}_{\text{within}} + \mathcal{L}_{\text{cross}}$

$f_\theta \leftarrow \text{backward}(\mathcal{L}_{\text{total}})$

end

### - Training details

$\Rightarrow$  The paper use Feature Pyramid Network <https://arxiv.org/abs/1612.03144>  
with ResNet-18 backbone.

$\Rightarrow$  ResNet-18 is pretrained on ImageNet.

$\Rightarrow$  Fusion dimension of FPN = 128

$\Rightarrow$  L2 normalization on the feature map.

$\Rightarrow$  Images are resized and center cropped to  $320 \times 320$ .

$\Rightarrow$  Overclustering

- Previous works show that jointly optimizing for another higher cluster number leads to better stability and accuracy.

- The paper takes two cluster numbers together  
 $\Rightarrow K_1$  = one set by the user,  
 $\Rightarrow K_2$  is set to a fixed value of 100,
- The paper optimizes

$$\mathcal{L} = \frac{\log K_2}{\log K_1 + \log K_2} \mathcal{L}_{K_1} + \frac{\log K_1}{\log K_1 + \log K_2} \mathcal{L}_{K_2}$$

- Why these weights?

$\Rightarrow$  Magnitude of cross entropy depends on  $\log(K)$ .

$\Rightarrow$  The paper says it also applies a balance term based on cluster size. However, it does not give enough details.

- Results can be found in the paper

$\Rightarrow$  The new algo compares favorably to

- Deep Clusters <https://arxiv.org/abs/1807.05520>
- Invariant information clustering <https://arxiv.org/abs/1807.06653>