

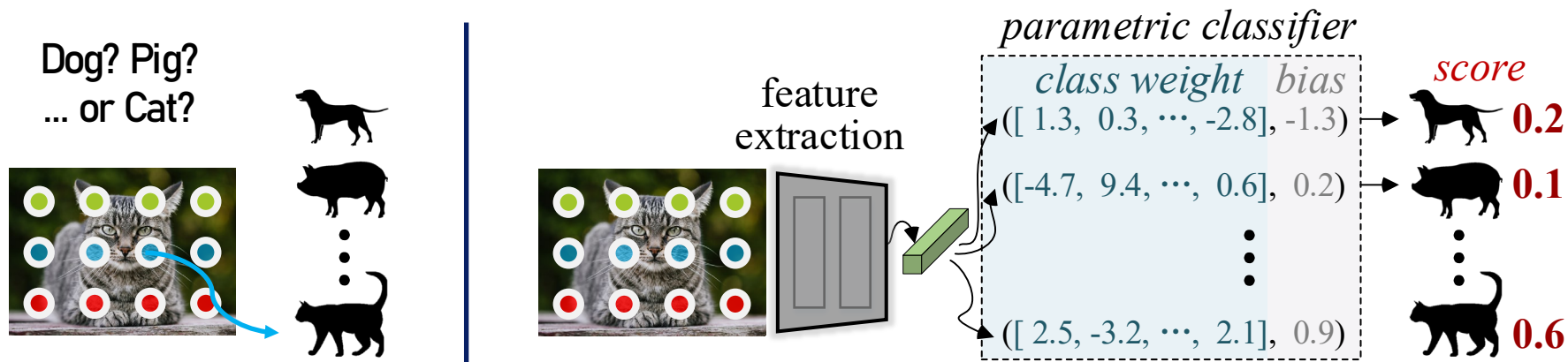
GMMSeg: Gaussian Mixture based Generative Semantic Segmentation Models

Chen Liang, Wenguan Wang, Jiaxu Miao, Yi Yang

Zhejiang University & University of Technology Sydney

NeurIPS 2022

Semantic Segmentation



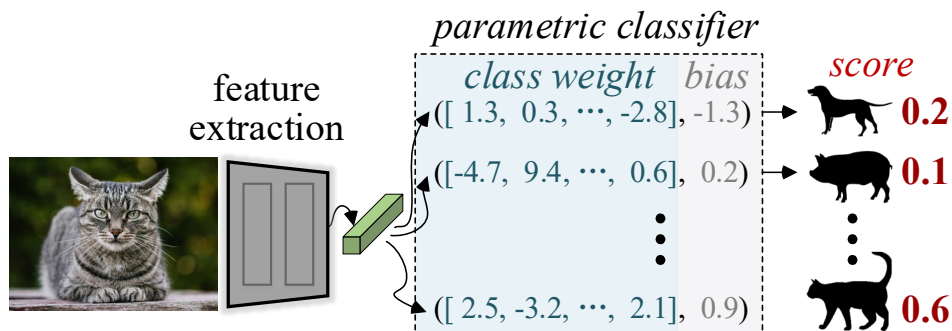
Semantic Segmentation:

Dense Feature extractor + **Parametric softmax classifier**

Deficiency of Parametric Softmax Classifier

Parametric softmax classifier:

$$p(c|\mathbf{x}; \boldsymbol{\omega}, \boldsymbol{\theta}) = \frac{\exp(y_c)}{\sum_{c'} \exp(y_{c'})} = \frac{\exp(\mathbf{w}_c^\top \mathbf{x} + b_c)}{\sum_{c'} \exp(\mathbf{w}_{c'}^\top \mathbf{x} + b_{c'})}$$



- Only learning **decision boundaries**; Ignoring underlying **data distribution**.

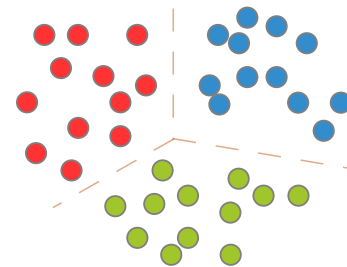


Straightforward: Only learn **decision boundaries**



Fail to capture the intrinsic class characteristics;

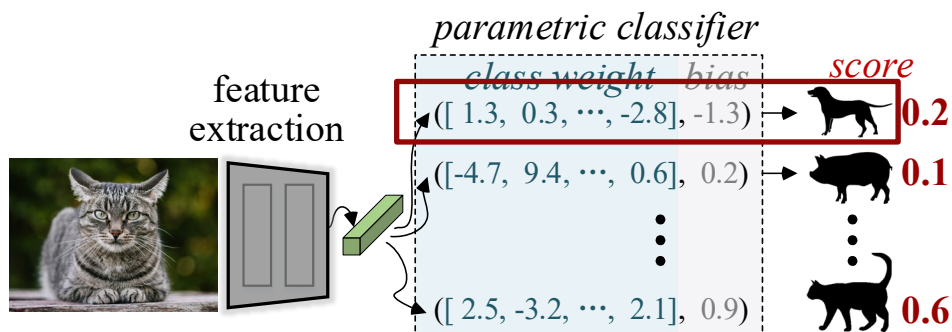
Ignore underlying **data structure**



Deficiency of Parametric Softmax Classifier

Parametric softmax classifier:

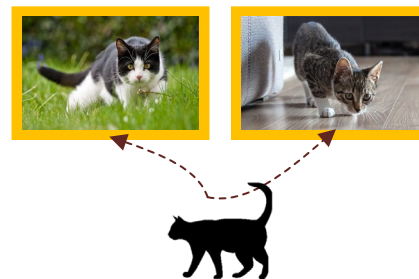
$$p(c|\mathbf{x}; \boldsymbol{\omega}, \boldsymbol{\theta}) = \frac{\exp(y_c)}{\sum_{c'} \exp(y_{c'})} = \frac{\exp(\mathbf{w}_c^\top \mathbf{x} + b_c)}{\sum_{c'} \exp(\mathbf{w}_{c'}^\top \mathbf{x} + b_{c'})}$$



- Implicit **unimodality** assumption; Bearing **no within-class variation**.

The unimodality assumption is rarely the case in real-world scenarios.

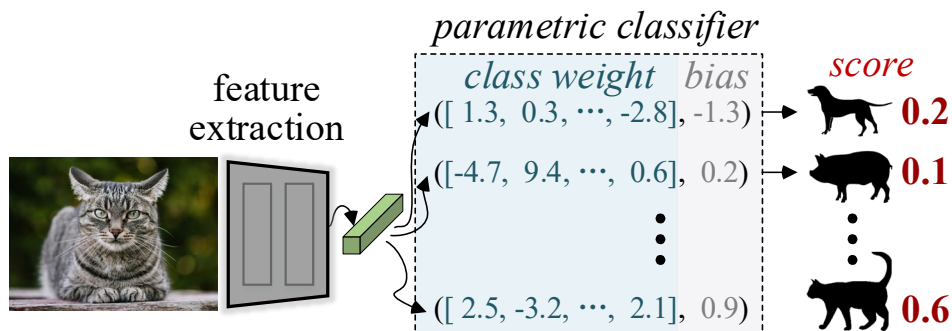
-  The model **less tolerant of intra-class variances**.



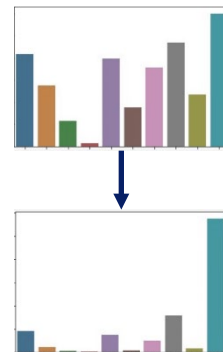
Deficiency of Parametric Softmax Classifier

Parametric softmax classifier:

$$p(c|\mathbf{x}; \boldsymbol{\omega}, \boldsymbol{\theta}) = \frac{\exp(y_c)}{\sum_{c'} \exp(y_{c'})} = \frac{\exp(\mathbf{w}_c^\top \mathbf{x} + b_c)}{\sum_{c'} \exp(\mathbf{w}_{c'}^\top \mathbf{x} + b_{c'})}$$



- Inferior **robustness to out-of-distribution inputs**; Poorly **calibrated**.
- 🗨 Prediction score is useless besides its comparative value against others;
Struggling to recognize out-of-distribution data.
- 🗨 Model accuracy deteriorates rapidly away from the decision boundaries;
yields poorly calibrated predictions.



Rethinking *de facto* Paradigm

- Only learning **decision boundaries**; Ignoring underlying **data structure**.
- Implicit **unimodality** assumption; Bearing **no within-class variation**.
- Inferior **robustness to out-of-distribution inputs**; Poorly **calibrated**.

Is there any way to address the **limitations** of *de facto* segmentation regime?



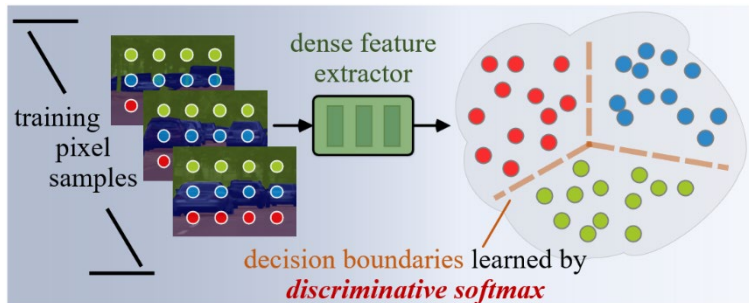
Our answer:

Generative Gaussian Mixture Classifier (GMMSeg)

Discriminative vs. Generative

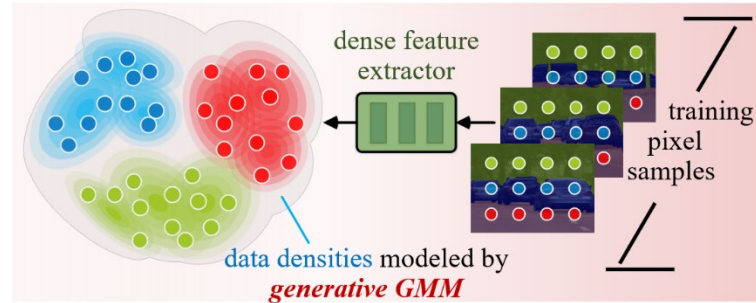
Discriminative Classifier

- Model $p(c|\mathbf{x})$ directly.
- Optimization: $\Pi_{(x,c) \in \mathcal{D}} p(c|\mathbf{x})$
- Model **Decision Boundaries** only
- Example: **Parametric softmax**



Generative Classifier

- Model joint distribution $p(\mathbf{x}, c)$.
Then deduce $p(c|\mathbf{x}) = \frac{p(c)p(\mathbf{x}|c)}{\sum_{c'} p(c')p(\mathbf{x}|c')}$.
Uniform Prior
- Optimization: $\Pi_{(x,c) \in \mathcal{D}} p(\mathbf{x}|c)$
- Model entire **Data Distribution**
- Example: **GMMSeg**



GMMSeg: Distribution Modeling

Use **Gaussian Mixtures** to model arbitrary feature distribution:

- Simple, elegant and powerful.

$$p(\mathbf{x}|c; \phi_c) = \sum_{m=1}^M p(m|c; \boldsymbol{\pi}_c) p(\mathbf{x}|c, m; \boldsymbol{\mu}_c, \boldsymbol{\Sigma}_c) = \sum_{m=1}^M \pi_{cm} \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_{cm}, \boldsymbol{\Sigma}_{cm}).$$

Parameter optimization: Maximizing the log-likelihood over all feature-label pairs $\{(\mathbf{x}_n, c_n)\}_{n=1}^N$

$$\phi_c^* = \arg \max_{\phi_c} \sum_{\mathbf{x}_n: c_n=c} \log p(\mathbf{x}_n|c; \phi_c) = \arg \max_{\phi_c} \sum_{\mathbf{x}_n: c_n=c} \log \sum_{m=1}^M p(\mathbf{x}_n, m|c; \phi_c)$$

$m|c \sim \text{Multinomial}(\boldsymbol{\pi}_c)$: prior of mixture components, $\sum_m \pi_{cm} = 1$.

$\boldsymbol{\mu}_{cm} \in \mathbb{R}^D$, $\boldsymbol{\Sigma}_{cm} \in \mathbb{R}^{D \times D}$: mean and covariance matrix for component m in class c .

$\phi_c = \{\boldsymbol{\pi}_c, \boldsymbol{\mu}_c, \boldsymbol{\Sigma}_c\}$: set of all parameters.

GMMSeg: Sinkhorn EM Optimization

Optimization through vanilla EM (F-function form¹):

$$\mathbf{E}\text{-Step: } q_c^{(t)} = \arg \max_{q_c} F(q_c, \phi_c^{(t-1)}), \quad \mathbf{M}\text{-Step: } \phi_c^{(t)} = \arg \max_{\phi_c} F(q_c^{(t)}, \phi_c).$$

F-function is defined as: $F(q_c, \phi_c) = \mathbb{E}_{q_c}[\log p(\mathbf{x}, m|c; \phi_c)] + H(q_c)$

EM starts with some **initial guess** at parameters $\phi_c^{(0)}$, then iterates over

- **E-Step:** given $\phi_c^{(t-1)}$, compute the posterior $q_c^{(t)}$ over the M components.
 - **M-Step:** given **soft cluster assignment** $q_c^{(t)}$, the parameters are updated as $\phi_c^{(t)}$ such that the F function is maximized.
-

$q_c[m] = p(m|\mathbf{x}, c; \phi_c)$: the probability that data \mathbf{x} is **assigned** to component m .

$H(q_c) = -\mathbb{E}_{q_c}[\log q_c[m]]$: the entropy of q_c .

[1] A view of the em algorithm that justifies incremental, sparse, and other variants. In *Learning in graphical models*, 1998.

GMMSeg: Sinkhorn EM Optimization

Problem: Standard EM suffers from **slow convergence**; Delivers **unsatisfactory results**, potentially due to the parameter sensitivity of EM.

GMMSeg: Sinkhorn EM Optimization

Sinkhorn EM¹ with a **uniform prior** on mixture weight, i.e. $\forall c, m : \pi_{cm} = \frac{1}{M}$:

$$\text{E-Step: } q_c^{(t)} = \arg \max_{q_c \in \mathcal{Q}_c} F(q_c, \phi_c^{(t-1)})$$

$$\text{restricted by a constraint: } \mathcal{Q}_c = \{q_c : \frac{1}{N_c} \sum_{\mathbf{x}_n : c_n = c} p(m | \mathbf{x}_n, c) = \frac{1}{M}\}$$

Analogous to **entropy-regularized OT**: **←** Solved by **Sinkhorn-Knopp** algorithm.

$$\min_{\mathbf{Q}_c \in \mathcal{Q}'_c} \sum_{n,m} \mathbf{Q}_c(n, m) \mathbf{O}_c(n, m) + \epsilon H(\mathbf{Q}_c), \quad \mathcal{Q}'_c = \{\mathbf{Q}_c \in \mathbb{R}_+^{N_c \times M} : \mathbf{Q}_c \mathbf{1}^M = \mathbf{1}^{N_c}, (\mathbf{Q}_c)^\top \mathbf{1}^{N_c} = \frac{N_c}{M} \mathbf{1}^M\}$$

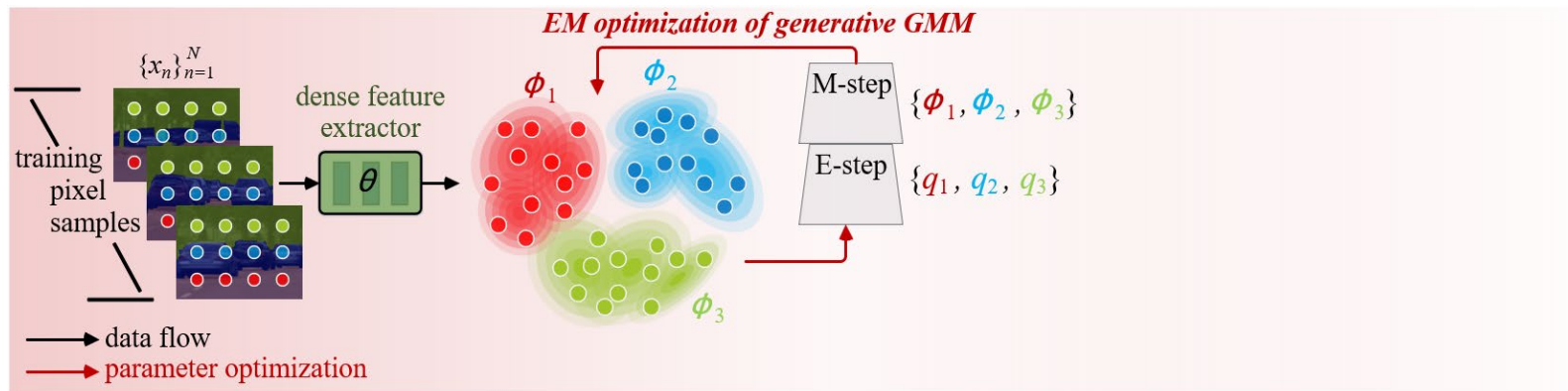
- Sinkhorn EM is proved to have the same global optimum with the standard EM yet is less prone to getting stuck in local optima¹.

$\mathbf{Q}_c(n, m) = q_{cn}[m]$: posterior distribution over the M components (target solution).

$\mathbf{O}_c(n, m) = -\log p(\mathbf{x}_n | c, m)$: negative log-likelihood (cost matrix).

[1] Sinkhorn em: an expectation-maximization algorithm based on entropic optimal transport. *arXiv*, 2021.

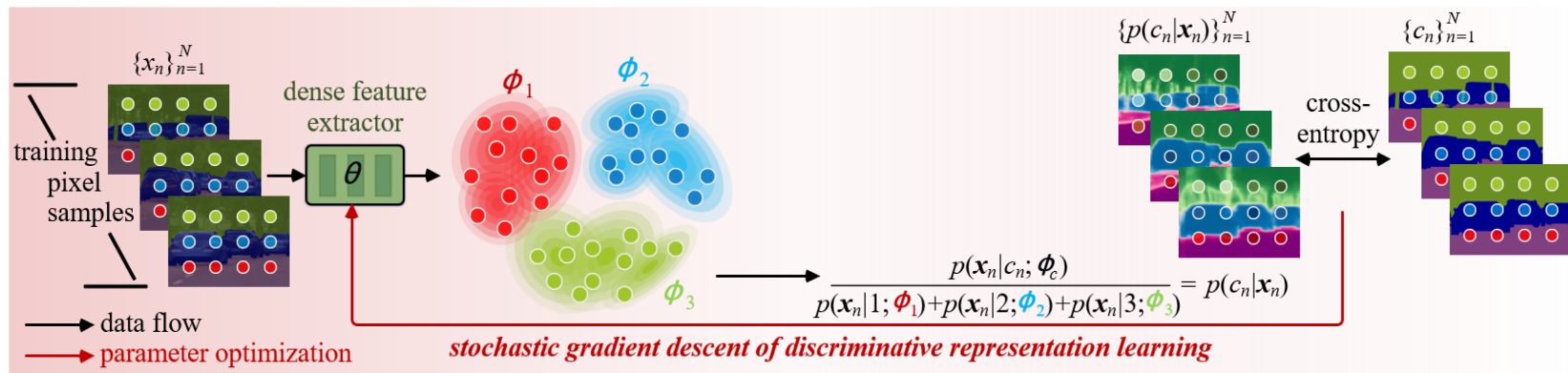
GMMSeg: Hybrid Optimization



Generative Optimization (Sinkhorn EM) of GMM Classifier: $\{\phi_c^*\}_{c=1}^C =$

$$\{\arg \max_{\phi_c} \sum_{\mathbf{x}_n: c_n=c} \log p(\mathbf{x}_n | c; \phi_c)\}_{c=1}^C = \{\arg \max_{\phi_c} \sum_{\mathbf{x}_n: c_n=c} \log \sum_{m=1}^M \pi_{cm} \mathcal{N}(\mathbf{x}_n; \boldsymbol{\mu}_{cm}, \boldsymbol{\Sigma}_{cm})\}_{c=1}^C,$$

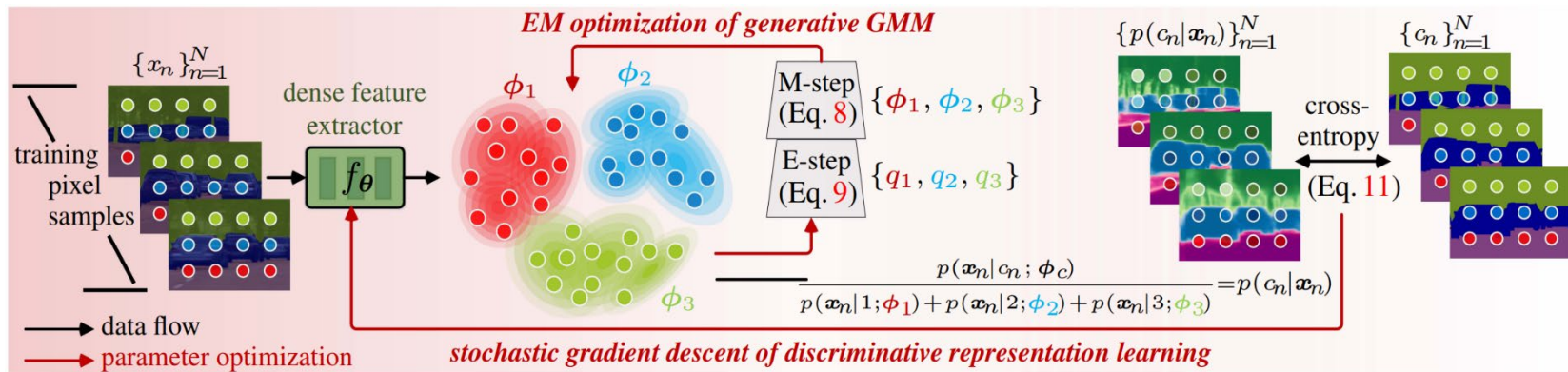
GMMSeg: Hybrid Optimization



Discriminative Learning (Cross-Entropy Loss) of Dense Representation: $\theta^* =$

$$\arg \min_{\theta} - \sum_{(x,c) \in \mathcal{D}} \log p(c|\mathbf{x}; \{\phi_c^*\}_{c=1}^C, \theta) = \arg \min_{\theta} - \sum_{(x,c) \in \mathcal{D}} \log \left(\frac{\sum_{m=1}^M \pi_{cm} \mathcal{N}(f_{\theta}(x); \boldsymbol{\mu}_{cm}, \boldsymbol{\Sigma}_{cm})}{\sum_{c'=1}^C \sum_{m=1}^M \pi_{c'm} \mathcal{N}(f_{\theta}(x); \boldsymbol{\mu}_{c'm}, \boldsymbol{\Sigma}_{c'm})} \right)$$

GMMSeg: Hybrid Optimization



Generative Optimization (Sinkhorn EM) of GMM Classifier: $\{\phi_c^*\}_{c=1}^C =$


$$\{\arg \max_{\phi_c} \sum_{x_n: c_n=c} \log p(x_n|c; \phi_c)\}_{c=1}^C = \{\arg \max_{\phi_c} \sum_{x_n: c_n=c} \log \sum_{m=1}^M \pi_{cm} \mathcal{N}(x_n; \mu_{cm}, \Sigma_{cm})\}_{c=1}^C,$$

Discriminative Learning (Cross-Entropy Loss) of Dense Representation: $\theta^* =$

$$\arg \min_{\theta} - \sum_{(x,c) \in \mathcal{D}} \log p(c|x; \{\phi_c^*\}_{c=1}^C, \theta) = \arg \min_{\theta} - \sum_{(x,c) \in \mathcal{D}} \log \left(\frac{\sum_{m=1}^M \pi_{cm} \mathcal{N}(f_\theta(x); \mu_{cm}, \Sigma_{cm})}{\sum_{c'=1}^C \sum_{m=1}^M \pi_{c'm} \mathcal{N}(f_\theta(x); \mu_{c'm}, \Sigma_{c'm})} \right)$$

Summary: Generative Gaussian Mixture Based Classifier

- **Versatility.**

 GMMSeg is a **principled** framework, fully compatible with modern seg. architectures.

- Best of Both Worlds: **Strong discriminative performance.**

 GMMSeg achieves the merits of both generative and discriminative learning.

Fit data distribution on evolving feature space with **online EM generative optimization.**

Discriminatively end-to-end trained data space under the guidance of the GMM classifier.

- Best of Both Worlds: **Distribution-preserving in nature.**

 GMMSeg can **naturally reject abnormal inputs**, without any post-processing;

 **Meaningful likelihood** of the example fitting each class GMM distribution; **Well-calibrated.**

Experiments: Semantic Segmentation

Quantitative results on ADE20K val, Cityscapes val, and COCO-Stuff test with mean IoU as the evaluation metric.

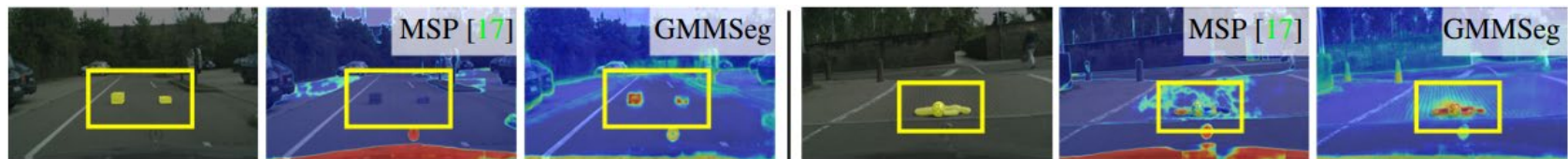
Method	Backbone	ADE _{20K}	Citys.	COCO.
FCN [CVPR15] [1]	ResNet ₁₀₁	39.9	75.5	32.6
PSPNet [CVPR17] [3]	ResNet ₁₀₁	44.4	79.8	37.8
SETR [CVPR21] [9]	†ViT _{Large}	48.2	79.2	-
Segmenter [ICCV21] [8]	†ViT _{Large}	‡51.8	79.1	-
MaskFormer [NeurIPS21] [81]	†Swin _{Base}	‡52.7	-	-
DeepLab _{v3+} [ECCV18] [47]	ResNet ₁₀₁	45.5	80.6	33.8
GMMSeg		46.7 ↑1.2	81.1 ↑0.5	35.5 ↑1.7
OCRNNet [ECCV20] [48]	HRNet _{v2w48}	43.3	80.4	37.6
GMMSeg		44.8 ↑1.5	81.2 ↑0.8	39.2 ↑1.6
UPerNet [ECCV18] [49]	Swin _{Base}	48.0	81.1	43.4
GMMSeg		49.0 ↑1.0	81.8 ↑0.7	44.3 ↑0.9
SegFormer [NeurIPS21] [7]	MiT _{B5}	50.0	82.0	44.0
GMMSeg		50.8 ↑0.8	82.6 ↑0.6	44.7 ↑0.7

†: pretrained on ImageNet_{22K}; ‡: using larger crop-size, *i.e.*, 640×640

Experiments: Anomaly Segmentation

Fishyscapes Lost&Found test and Static test.

Method	Re-training	Extra Network	OoD Data	FS Lost&Found		FS Static	
				AP \uparrow	FPR ₉₅ \downarrow	AP \uparrow	FPR ₉₅ \downarrow
Density - Single-layer NLL [12]	✗	✓	✗	3.01	32.9	40.86	21.29
Density - Minimum NLL [12]	✗	✓	✗	4.25	47.15	62.14	17.43
Density - Logistic Regression [12]	✗	✓	✓	4.65	24.36	57.16	13.39
Image Resynthesis [15]	✗	✓	✗	5.70	48.05	29.6	27.13
Bayesian Deeplab [16]	✓	✗	✗	9.81	38.46	48.70	15.05
OoD Training - Void Class [17]	✓	✗	✓	10.29	22.11	45.00	19.40
Discriminative Outlier Detection Head [18]	✓	✓	✓	31.31	19.02	96.76	0.29
Dirichlet Deeplab [19]	✓	✗	✓	34.28	47.43	31.30	84.60
SynBoost [20]	✗	✓	✓	43.22	15.79	72.59	18.75
MSP [21]	✗	✗	✗	1.77	44.85	12.88	39.83
Entropy [22]	✗	✗	✗	2.93	44.83	15.41	39.75
kNN Embedding - density [12]	✗	✗	✗	3.55	30.02	44.03	20.25
SML [14]	✗	✗	✗	31.05	21.52	53.11	19.64
GMMSeg-DeepLabv3+	✗	✗	✗	55.63	6.61	76.02	15.96



Summary: Rethinking *de facto* Paradigm in Segmentation

- Only learning **decision boundaries**; Ignoring underlying **data structure**.
- Implicit **unimodality** assumption; Bearing **no within-class variation**.
- Inferior **robustness to out-of-distribution inputs**; Poorly **calibrated**.

Our answer:

Generative Classifier: **Gaussian Mixture Model based Classifier (GMMSeg)**.

Paper



Code



Thanks!