



# Deep DPM: Deep Clustering With an Unknown Number of Clusters

2023.01.18

Presenter : Taewoong Kang  
twk@deepnoid.com

## Abstract

- **Effective Deep-clustering method**
  - (When K is unknown, model become computationally expensive)
- Does **not require** knowing the value of **K**
- Using **split/merging framework, a dynamic architecture, novel loss**



Figure 3. Examples of ImageNet images clustered together by DeepDPM. Each panel stands for a different cluster.

## Architecture

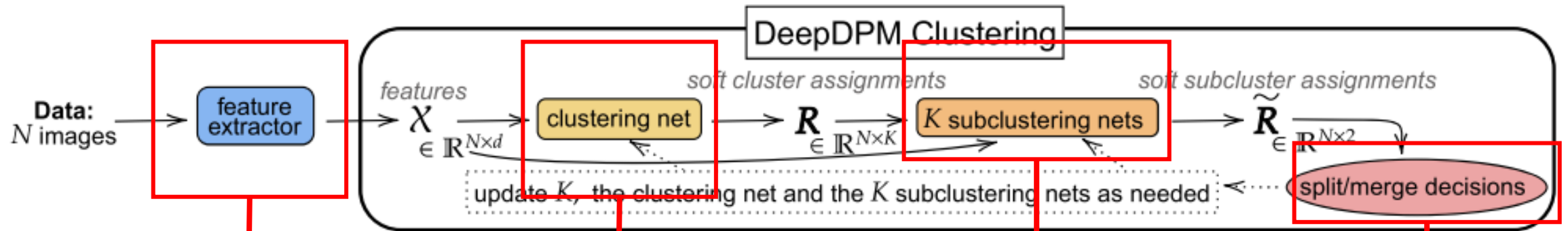


Figure 2. DeepDPM’s pipeline: given features  $\mathcal{X}$ , the clustering net outputs cluster assignments,  $\mathbf{R}$ , while the subclustering nets generate subcluster assignments,  $\tilde{\mathbf{R}}$ . Upon the acceptance of split/merge proposals, all those nets are updated during the learning.

Embedding  
(use MOCO)

Cluster A, B,..., K

Cluster A1, A2  
...  
Cluster K1, K2

Change K

## Notations

$\mathcal{X} = (\mathbf{x}_i)_{i=1}^N$  denote  $N$  data points in  $\mathbb{R}^d$

$z_i$  is the point-to-cluster assignment Cluster label  $(\mathbf{x}_i)_{i:z_i=k}$ .

$$\bar{K} \triangleq |\{k : k \in \mathbf{z}\}|$$

$$\mathbf{z} = (z_i)_{i=1}^N.$$

## DPGMM (the Dirichlet Process Gaussian Mixture Model)

- BNP (Bayesian nonparametric) extension of GMM

$$p(\mathbf{x} | (\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k, \pi_k)_{k=1}^{\infty}) = \sum_{k=1}^{\infty} \pi_k \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

## Notations

- Gaussian pdf (probability density function)  $\mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$
- mean  $\boldsymbol{\mu}_k \in \mathbb{R}^d$  and a  $d$ -by- $d$  covariance matrix  $\boldsymbol{\Sigma}_k$
- $\mathbf{x} \in \mathbb{R}^d$ ,  $\pi_k > 0 \forall k$ , and  $\sum_{k=1}^{\infty} \pi_k = 1$
- $\boldsymbol{\theta}_k = (\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$
- $\boldsymbol{\theta} = (\boldsymbol{\theta}_k)_{k=1}^{\infty}$  are iid, draws from their prior, typically a NIW distribution
- $\boldsymbol{\pi} = (\pi_k)_{k=1}^{\infty}$  are drawn using the GEM (the Griffiths-Engen-McCloskey stick-breaking process)
- $\alpha > 0$  the concentration parameter

## Metropolis-Hastings framework

- Hasting Ratio

$$H_s = \frac{\alpha \Gamma(N_{k,1}) f_{\mathbf{x}}(\mathcal{X}_{k,1}; \lambda) \Gamma(N_{k,2}) f_{\mathbf{x}}(\mathcal{X}_{k,2}; \lambda)}{\Gamma(N_k) f_{\mathbf{x}}(\mathcal{X}_k; \lambda)}$$

$\Gamma$  is the Gamma function

$$N_k = |\mathcal{X}_k|$$

$f_{\mathbf{x}}(\cdot; \lambda)$  Is the marginal likelihood where  $\lambda$  represents the NIW hyperparameters

## NIW

- A **conjugate prior** to the multivariate normal distribution with **an unknown mean** and **unknown covariance matrix**  $\rightarrow$  posterior probability will be in the same distribution
- Algebraically convenient to inference

- Pdf of Inverse-Wishart(IW)

$$\mathcal{W}^{-1}(\Sigma_k; \nu, \Psi) = \frac{|\nu \Psi|^{\frac{\nu}{2}}}{2^{\frac{\nu d}{2}} \Gamma_d(\frac{\nu}{2})} |\Sigma_k|^{-\frac{\nu+d+1}{2}} e^{-\frac{1}{2} \text{tr}(\nu \Psi \Sigma_k^{-1})}$$

$\nu > d - 1$ ,  $\Psi \in \mathbb{R}^{d \times d}$  is SPD, and  $\Gamma_d$  is the ( $d$ -dimensional) multivariate gamma function.

The positive real number  $\nu$  and the SPD matrix  $\Psi$  are called the hyperparameters

- Symmetric and Positive Definite (SPD)

## NIW

- Pdf of Normal-Inverse-Wishart (NIW)

$$p(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k; \kappa, \boldsymbol{m}, \nu, \boldsymbol{\Psi}) = \text{NIW}(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k; \kappa, \boldsymbol{m}, \nu, \boldsymbol{\Psi}) \triangleq \underbrace{\mathcal{N}(\boldsymbol{\mu}_k; \boldsymbol{m}, \frac{1}{\kappa} \boldsymbol{\Sigma}_k)}_{p(\boldsymbol{\mu}_k | \boldsymbol{\Sigma}_k; \kappa, \boldsymbol{m})} \underbrace{\mathcal{W}^{-1}(\boldsymbol{\Sigma}_k; \nu, \boldsymbol{\Psi})}_{p(\boldsymbol{\Sigma}_k; \nu, \boldsymbol{\Psi})}$$

$\boldsymbol{m} \in \mathbb{R}^d$  and  $\kappa > 0$  (while  $\nu$  and  $\boldsymbol{\Psi}$  are as before) and  $\mathcal{N}(\boldsymbol{\mu}_k; \boldsymbol{m}, \frac{1}{\kappa} \boldsymbol{\Sigma}_k)$  is a  $d$ -dimensional Gaussian pdf

$$(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \sim \text{NIW}(\boldsymbol{m}, \kappa, \boldsymbol{\Psi}, \nu) \quad \lambda \triangleq (\boldsymbol{m}, \kappa, \boldsymbol{\Psi}, \nu)$$

- Posterior hyperparameters

$$p(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k | \mathcal{X}_k) = \text{NIW}(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k; \kappa^*, \boldsymbol{m}_k^*, \nu^*, \boldsymbol{\Psi}_k^*)$$

$$\kappa_k^* = \kappa + N_k$$

$$\boldsymbol{m}_k^* = \frac{1}{\kappa_k^*} \left[ \kappa \boldsymbol{m} + \sum_{i: z_i = k} \boldsymbol{x}_i \right]$$

$$\nu_k^* = \nu + N_k$$

$$\boldsymbol{\Psi}_k^* = \frac{1}{\nu_k^*} \left[ \nu \boldsymbol{\Psi} + \kappa \boldsymbol{m} \boldsymbol{m}^T + \left( \sum_{i: z_i = k} \boldsymbol{x}_i \boldsymbol{x}_i^T \right) - \kappa_k^* \boldsymbol{m}_k^* (\boldsymbol{m}_k^*)^T \right]$$



## ▀ Marginal Likelihood Function

- When marginalizing over the parameters of the Gaussian, one obtains the marginal data likelihood

$$\begin{aligned} f_{\mathbf{x}}((\mathbf{x}_i)_{i=1}^N; \lambda) &= f_{\mathbf{x}}((\mathbf{x}_i)_{i=1}^N; \mathbf{m}, \kappa, \boldsymbol{\Psi}, \nu) = \int p((\mathbf{x}_i)_{i=1}^N | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) p(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k; \lambda) d(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \\ &= \frac{1}{\pi^{\frac{Nd}{2}}} \frac{\Gamma_d(\nu^*/2)}{\Gamma_d(\nu/2)} \frac{|\nu \boldsymbol{\Psi}|^{\nu/2}}{|\nu^* \boldsymbol{\Psi}_k^*|^{\nu^*/2}} \left( \frac{\kappa}{\kappa^*} \right)^{d/2} \end{aligned}$$

where  $\Gamma_d$  is the  $d$ -dimensional Gamma function.

## DeepDPM

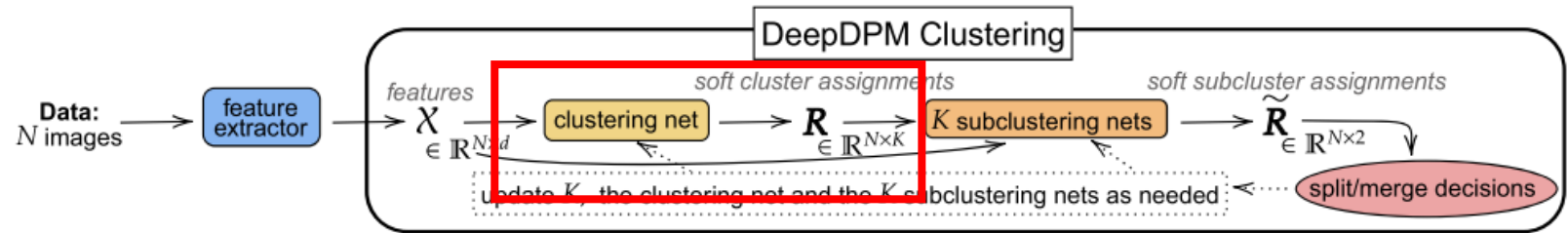


Figure 2. DeepDPM's pipeline: given features  $\mathcal{X}$ , the clustering net outputs cluster assignments,  $\mathbf{R}$ , while the subclustering nets generate subcluster assignments,  $\tilde{\mathbf{R}}$ . Upon the acceptance of split/merge proposals, all those nets are updated during the learning.

### Clustering net

- Given the current value  $K$ , the data is passed to the clustering net.
- Generates  $K$  soft cluster assignments

$$f_{\text{cl}}(\mathcal{X}) = \mathbf{R} = (\mathbf{r}_i)_{i=1}^N \quad \mathbf{r}_i = (r_{i,k})_{k=1}^K$$

- We compute the hard assignments

$$\mathbf{z} = (z_i)_{i=1}^N \text{ by } z_i = \arg \max_k r_{i,k}.$$

$$p(z_k = 1 | \mathbf{x}_n) = \frac{\pi_k \mathcal{N}(\mathbf{x}_n | \mu_k, \Sigma_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x}_n | \mu_j, \Sigma_j)} = \gamma(z_{nk})$$

## DeepDPM

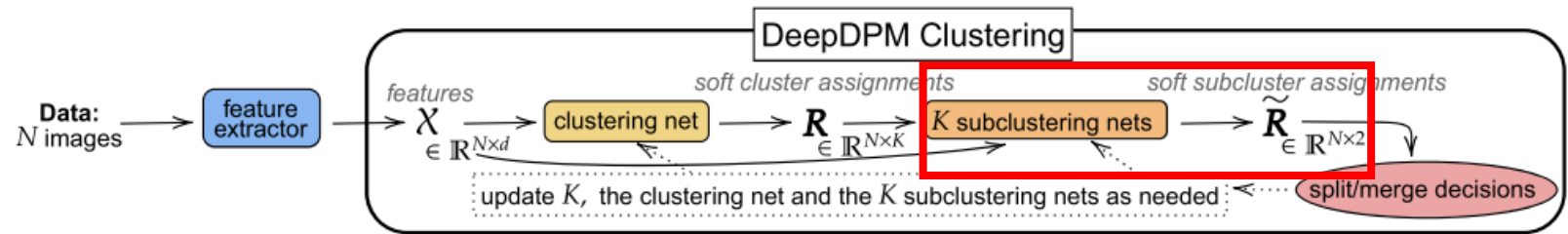


Figure 2. DeepDPM's pipeline: given features  $\mathcal{X}$ , the clustering net outputs cluster assignments,  $\mathbf{R}$ , while the subclustering nets generate subcluster assignments,  $\tilde{\mathbf{R}}$ . Upon the acceptance of split/merge proposals, all those nets are updated during the learning.

- Subclustering net
  - Each subclustering net is fed with the data (hard-) assigned to its respective cluster
  - Generates soft subcluster assignments

$$f_{\text{sub}}^k(\mathcal{X}_k) = \tilde{\mathbf{R}}_k = (\tilde{\mathbf{r}}_i)_{i:z_i=k} \quad \tilde{\mathbf{r}}_i = (\tilde{r}_{i,j})_{j=1}^2$$

## ■ Splits and Merges

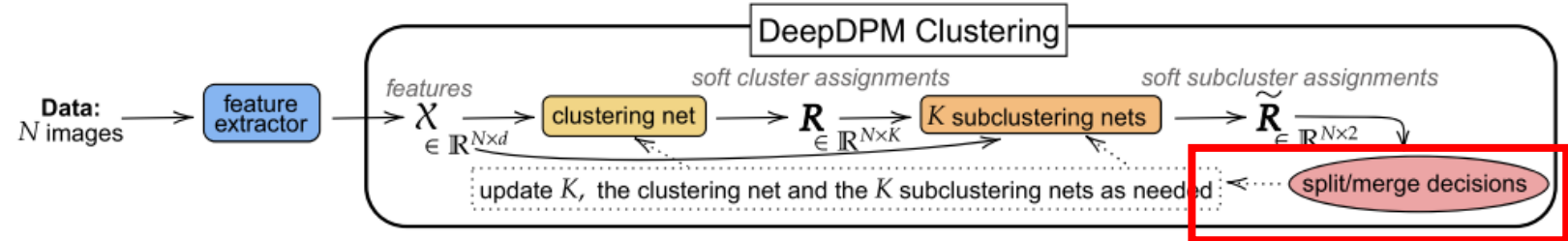


Figure 2. DeepDPM's pipeline: given features  $X$ , the clustering net outputs cluster assignments,  $R$ , while the subclustering nets generate subcluster assignments,  $\tilde{R}$ . Upon the acceptance of split/merge proposals, all those nets are updated during the learning.

### ■ Changing K via Splits and Merges

#### ■ Splits

- A split proposal is accepted stochastically with probability  $\min(1, H_s)$

#### ■ Merges

- Must ensure we never mistakenly merging three clusters together
- Consider the merges of each cluster with only its 3 nearest neighbors
- Merge proposal is accepted/rejected using a Hasting ratio  $H_m = 1/H_s$
- The parameter and the weight of the new-cluster are initialized using the weighted MAP estimation

$$\begin{aligned} \mu_{k_1} &\leftarrow \tilde{\mu}_{k,1}, & \Sigma_{k_1} &\leftarrow \tilde{\Sigma}_{k,1}, & \pi_{k_1} &\leftarrow \pi_k \times \tilde{\pi}_{k,1} \\ \mu_{k_2} &\leftarrow \tilde{\mu}_{k,2}, & \Sigma_{k_2} &\leftarrow \tilde{\Sigma}_{k,2}, & \pi_{k_2} &\leftarrow \pi_k \times \tilde{\pi}_{k,2} \end{aligned}$$

## Loss function

- Motivated by EM (Expectation Maximization) in the Bayesian GMM
- E step
  - Responsibility (책임값)
  - Each E step is followed by a standard M step in the Bayesian GMM (except soft assignment used in MAP)

$$r_{i,k}^E = \frac{\pi_k \mathcal{N}(\mathbf{x}_i; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{k'=1}^K \pi_{k'} \mathcal{N}(\mathbf{x}_i; \boldsymbol{\mu}_{k'}, \boldsymbol{\Sigma}_{k'})} \quad k \in \{1, \dots, K\} \quad \begin{array}{l} (\pi_k, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)_{k=1}^K \text{ from the previous epoch} \\ \sum_{k=1}^K r_{i,k}^E = 1 \end{array}$$

$$\mathcal{L}_{\text{cl}} = \sum_{i=1}^N \text{KL}(\mathbf{r}_i \| \mathbf{r}_i^E) \quad \mathcal{L}_{\text{sub}} = \sum_{k=1}^K \sum_{i=1}^{N_k} \sum_{j=1}^2 \tilde{r}_{i,j} \|\mathbf{x}_i - \tilde{\boldsymbol{\mu}}_{k,j}\|_{\ell_2}^2$$

## ▀ Loss function

- Bayesian M step
  - Use weighted versions of the MAP estimates of  $(\pi_k, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)_{k=1}^K$
  - Instead of  $r_{i,k}^{\text{E}}$  we use  $r_{i,k}$  as weights

## Weighted MAP Estimates

- In the unweighted case, we use

$$\begin{aligned} \Sigma_k &= \frac{\nu^* \Psi_k^*}{\nu^* - d + 1} \\ \mu_k &= m_k^* . \end{aligned} \quad \text{Eq1}$$

- In weighted MAP, we still use Eq1, but instead of the posterior hyperparameters we use their weighted versions

$$\begin{aligned} \kappa_k^* &= \kappa + N_k \\ m_k^* &= \frac{1}{\kappa_k^*} \left[ \kappa m + \sum_{i: z_i = k} \mathbf{x}_i \right] \\ \nu_k^* &= \nu + N_k \\ \Psi_k^* &= \frac{1}{\nu^*} \left[ \nu \Psi + \kappa m m^T + \left( \sum_{i: z_i = k} \mathbf{x}_i \mathbf{x}_i^T \right) - \kappa_k^* m_k^* (m_k^*)^T \right] . \end{aligned}$$



$$\begin{aligned} \kappa_k^* &= \kappa + \sum_{i=1}^N r_{i,k} \\ m_k^* &= \frac{1}{\kappa_k^*} \left[ \kappa m + \sum_{i=1}^N r_{i,k} \mathbf{x}_i \right] \\ \nu_k^* &= \nu + \sum_{i=1}^N r_{i,k} \\ \Psi_k^* &= \frac{1}{\nu^*} \left[ \nu \Psi + \kappa m m^T + \left( \sum_{i=1}^N r_{i,k} \mathbf{x}_i \mathbf{x}_i^T \right) - \kappa_k^* m_k^* (m_k^*)^T \right] \end{aligned}$$

## Amortized EM Inference

- Our method still yields results that are usually **better than the standard EM**
- By the virtue of the smoothness of the function learned by the deep net, we **improve the prediction** for the points in not only the **current batch** but also **other batches**
- The smoothness serves as an inductive bias**, such that points should have similar labels
- When using the **GMM negative log likelihood** instead of our loss, empirically that **led to unstable optimization and/or poor results**

	ACC		
	$K_{\text{init}}=3$	$K_{\text{init}}=10$	$K_{\text{init}}=30$
No splits/merges	.29±.01	.59±.03	.46±.01
No splits	.29±.01	.59±.02	.45±.03
No merges	.46±.00	.58±.01	.47±.01
2-means instead of $f_{\text{sub}}$	.61±.00	.59±.02	.56±.02
No priors in the $M$ step	.58±.01	.57±.02	.58±.01
Isotropic loss instead of $\mathcal{L}_{\text{cl}}$	.58±.00	.58±.00	.58±.02
DeepDPM (full method)	<b>.62±.03</b>	<b>.61±.00</b>	<b>.62±.01</b>

- On Fashion-MNIST

Table 6. DeepDPM's performance under different ablations.



## ▀ A Weak Prior : Letting the Data Speak for Itself

- The inferred  $K$  depends on  $\mathcal{X}$ ,  $\alpha$ , and the NIW hyperparameters
- Intentionally choose the prior to be very weak
- When doing posterior calculation, if  $\alpha \ll N$ , where  $N$  is the number of data points, then the importance of  $\alpha$  diminishes
- For example,  $\nu$  is very high and  $\Psi$  is small  $\rightarrow$  favor small clusters, thus  $K$  is likely to be high
- $\nu$  is very high and  $\Psi$  is large  $\rightarrow$  favor large clusters, so  $K$  will tend to be small
- Our  $\alpha$ ,  $\nu$  and  $\kappa$  are all much smaller than  $N$  in all the datasets

## Classical methods

- Nonparametric ones are less affected by the imbalance.

	NMI	ARI	ACC	NMI	ARI	ACC	NMI	ARI	ACC
	MNIST [18]			USPS [35]			Fashion-MNIST [69]		
$K$ -means <sup>p</sup>	.90±.02	.84±.05	.85±.06	.86±.01	.79±.05	.80±.06	.67±.01	.50±.03	.60±.04
GMM <sup>p</sup>	<b>.94±.00</b>	<b>.95±.00</b>	<b>.98±.00</b>	.86±.02	.79±.05	.81±.06	.66±.01	.49±.02	.58±.03
DBSCAN	.92±0	.86±0	.89±0	.72±0	.46±0	.57±0	.63±0	-.32±0	.39±0
DPM Sampler	.92±.01	.91±.04	.93±.05	.87±.01	.82±.02	.83±.03	.67±.01	.49±.02	.59±.03
moVB	.93±.00	.94±.00	.97±.00	.87±.02	<b>.86±.04</b>	<b>.90±.04</b>	.66±.02	.47±.03	.55±.03
DeepDPM (Ours)	<b>.94±.00</b>	<b>.95±.00</b>	<b>.98±.00</b>	<b>.88±.00</b>	<b>.86±.01</b>	.89±.2	<b>.68±.01</b>	<b>.51±.02</b>	<b>.62±.03</b>
	MNIST <sup>imb</sup>			USPS <sup>imb</sup>			Fashion-MNIST <sup>imb</sup>		
$K$ -means <sup>p</sup>	.89±.03	.84±.06	.83±.06	.82±.02	.71±.05	.71±.05	.62±.01	.46±.02	.56±.03
GMM <sup>p</sup>	.94±.02	.95±.03	.96±.04	.83±.01	.74±.05	.76±.05	.62±.01	.46±.02	.57±.03
DBSCAN	.93±0	.92±0	.94±0	.84±0	.79±0	.80±0	.62±0	.35±0	.46±0
DPM Sampler	.93±.01	.94±.02	.96±.02	.89±.02	.89±.06	.91±.04	<b>.66±.01</b>	<b>.50±.01</b>	<b>.61±.01</b>
moVB	.94±.00	.95±.00	.96±.00	.88±.01	.89±.02	.91±.02	.63±.01	.44±.02	.53±.02
DeepDPM (Ours)	<b>.95±.01</b>	<b>.97±.01</b>	<b>.98±.01</b>	<b>.90±.00</b>	<b>.92±.00</b>	<b>.94±.00</b>	.65±.00	<b>.50±.00</b>	<b>.61±.00</b>

Table 1. Comparing the mean results ( $\pm$ std. dev.) of DeepDPM with classical clustering methods. The results are the mean of 10 independent runs. Methods marked with <sup>p</sup> are parametric (require  $K$ ). Datasets marked with <sup>imb</sup> are imbalanced ones.

## Nonparametric method

- DPM's inferred  $K$  is the closest to GT  $K$

Method	Inferred $K$		
	MNIST	USPS	Fashion-MNIST
DBSCAN	$9.0 \pm 0.00$	$6.0 \pm 0.00$	$4.0 \pm 0.00$
DPM Sampler	$11.3 \pm 0.82$	$8.5 \pm 0.85$	$12.4 \pm 0.97$
moVB	$14 \pm 1.00$	$11.2 \pm 1.08$	$16.9 \pm 2.30$
DeepDPM (Ours)	<b><math>10 \pm 0.00</math></b>	<b><math>9.2 \pm 0.42</math></b>	<b><math>10.2 \pm 0.79</math></b>

Table 2. Comparing the mean inferred value ( $\pm$ std. dev.) for  $K$  of 10 runs among **nonparametric methods** GT  $K = 10$ .

## Deep Nonparametric Methods

	MNIST [18]			STL-10 [15]			Reuters10k [43]		
Method	NMI	ARI	ACC	NMI	ARI	ACC	NMI	ARI	ACC
AdapVAE <sup>†</sup> [74] <i>avg</i>	.86±1.02	.84±2.35	N/A	.75±0.53	<b>.71±0.81</b>	N/A	.45±1.79	.43±5.73	N/A
DCC <sup>†</sup> [52] <i>best</i>	.912	N/A	.96	N/A	N/A	N/A	.59	N/A	.60
DCC <sup>‡</sup> [52] <i>avg</i>	.90±.02	.89±.07	.91±.07	.22±.00	.01±.00	.04±.00	.25±.00	.00±.00	.00±.00
DeepDPM (ours) <i>avg</i>	.90±.01	.91±.02	.93±.03	.78±.004	.70±.01	.84±.01	.61±.00	.64±.01	.83±.00
DeepDPM (ours) <i>best</i>	<b>.92</b>	<b>.93</b>	<b>.96</b>	<b>.79</b>	<b>.71</b>	<b>.85</b>	<b>.61</b>	<b>.64</b>	<b>.83</b>

Table 3. Comparing deep nonparametric methods. <sup>†</sup>: reported in the papers. <sup>‡</sup>: obtained using their code. *avg*: mean ( $\pm$ std. dev.) of 5 runs.

## ▀ The value of Deep Nonparametric Methods

Method	NMI	ARI	ACC
ImageNet-50: Balanced			
DBSCAN	.52±.00	.09±.00	.24±.00
moVB	.70±.01	.38±.01	.55±.02
DPM Sampler	.72±.00	.43±.01	.57±.01
DeepDPM (ours)	.75±.00	.49±.01	.64±.00
DeepDPM (ours)*	<b>.77±.00</b>	<b>.54±.01</b>	<b>.66±.01</b>
ImageNet-50: Imbalanced			
DBSCAN	.33±.00	.04±.00	.24±.00
moVB	.68±.01	.44±.03	.52±.03
DPM Sampler	.70±.00	.40±.01	.51±.00
DeepDPM (ours)	.74±.01	.48±.02	.58±.01
DeepDPM (ours)*	<b>.75±.00</b>	<b>.51±.01</b>	<b>.60±.01</b>

Table 4. Comparison of **nonparametric** methods on ImageNet-50 and its imbalanced version. \* marks results with AE alternation.

Method	Final/best $K$ : balanced	Final/best $K$ : imbalanced
$K$ -means <sup>p</sup>	40	20
DCN++ <sup>p</sup>	60	40
SCAN <sup>p</sup>	70	40
DBSCAN	16	13
moVB	46.2±1.3	<b>46.4±1.1</b>
DPM Sampler	72.0±2.6	70.3±4.6
DeepDPM (ours)	<b>52.0±1.0</b>	43.67±1.2
DeepDPM (ours)*	55.3±1.5	46.3±2.5

Table 5. Comparing the mean ( $\pm$ std. dev.) value for  $K$  found on ImageNet-50 of 3 runs. For the **parametric** methods (marked with <sup>p</sup>) we use the  $K$  value with the best silhouette score. \* marks results obtained with AE alternation.

## ▀ Ablation Study and Robustness to the Initial $K$

	ACC		
	$K_{\text{init}}=3$	$K_{\text{init}}=10$	$K_{\text{init}}=30$
No splits/merges	.29±.01	.59±.03	.46±.01
No splits	.29±.01	.59±.02	.45±.03
No merges	.46±.00	.58±.01	.47±.01
2-means instead of $f_{\text{sub}}$	.61±.00	.59±.02	.56±.02
No priors in the $M$ step	.58±.01	.57±.02	.58±.01
Isotropic loss instead of $\mathcal{L}_{\text{cl}}$	.58±.00	.58±.00	.58±.02
DeepDPM (full method)	<b>.62±.03</b>	<b>.61±.00</b>	<b>.62±.01</b>

Table 6. DeepDPM's performance under different ablations.

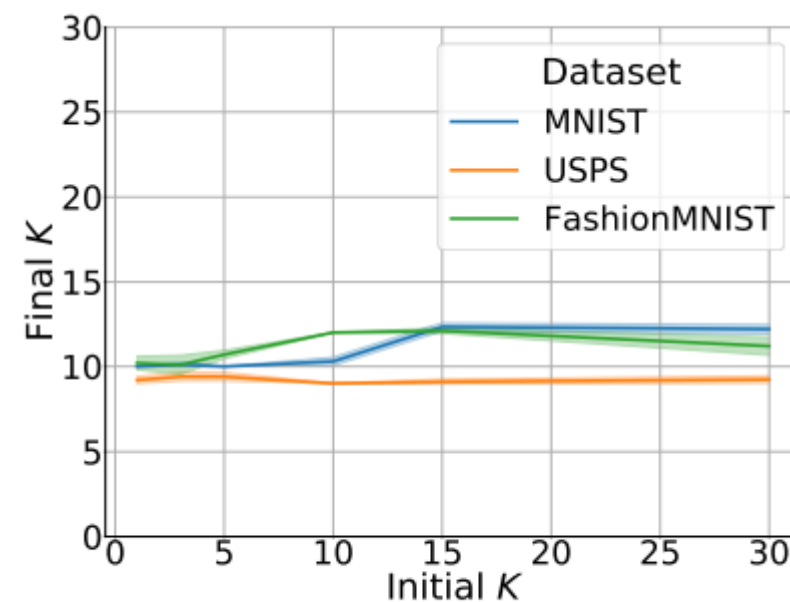


Figure 4. Robustness to the initial  $K$ . GT  $K = 10$  in all datasets.

- On Fashion-MNIST

## Limitations

- If Deep-DPM's input features are poor it would struggle to recover.
- $K$  is known and dataset is balanced, parametric methods may be a slightly better choice.

## Summary

- Outperforms deep and non-deep nonparametric methods and achieves SOTA results
- Demonstrated the added value the nonparametric approach brings to deep clustering

**Clustering the Entire ImageNet Dataset.** On ImageNet, we obtained the following results: ACC: 0.25, NMI: 0.65, ARI: 0.14. Our method was initialized with  $K = 200$  and converged into 707 clusters (GT=1000). These are first results on ImageNet reported for deep nonparametric clustering. [Figure 3](#) shows examples of images clustered together.



# Q & A



감사합니다

## ▀ Feature extraction

- Changing K via Splits and Merges
  - Splits
    - A split proposal is accepted stochastically with probability
  - Merges
    - Must ensure we never mistakenly