Dataset Distillation as Optimal Quantisation

M4DL Conference on Inverse Problems and DL

Hong Ye Tan 7 July 2025

Joint work with Emma Slade @ GSK.ai

Intro

Imagine this scenario:

You have a big dataset.

Imagine this scenario:

- · You have a big dataset.
- · You need to train a lot of models.

Imagine this scenario:

- · You have a big dataset.
- · You need to train a lot of models.
- You can take a small performance hit as long as it trains quickly.

1

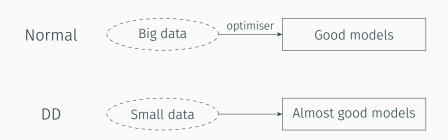
Imagine this scenario:

- · You have a big dataset.
- · You need to train a lot of models.
- · You can take a small performance hit as long as it trains quickly.



Imagine this scenario:

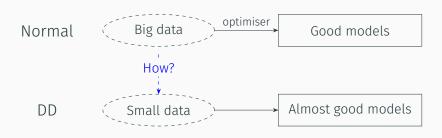
- · You have a big dataset.
- · You need to train a lot of models.
- · You can take a small performance hit as long as it trains quickly.



1

Imagine this scenario:

- · You have a big dataset.
- · You need to train a lot of models.
- · You can take a small performance hit as long as it trains quickly.



1

Problem: find a small surrogate dataset, such that training on it gives performance similar to training on the whole dataset

Problem: find a small surrogate dataset, such that training on it gives performance similar to training on the whole dataset

Problem: find a small surrogate dataset, such that training on it gives good performance

Problem: find a small surrogate dataset, such that training on it gives good performance

Intuition: you want the most important data.

Main acceleration: since your surrogate dataset is small, you train faster.

Problem: find a small surrogate dataset, such that training on it gives good performance

Intuition: you want the most important data.

Main acceleration: since your surrogate dataset is small, you train faster.

Two main archetypes:

- 1. Bi-level approaches
 - · Based on bi-level optimisation
- 2. Disentangled approaches
 - Based on approximation with empirical distributions

Problem: find a small surrogate dataset, such that training on it gives good performance

Intuition: you want the most important data.

Main acceleration: since your surrogate dataset is small, you train faster.

Two main archetypes:

- 1. Bi-level approaches
 - · Based on bi-level optimisation
- 2. Disentangled approaches
 - Based on approximation with empirical distributions

Dataset distillation is:

Dataset distillation is:

 \cdot Optimise over the distilled dataset \mathcal{S} :

 $\mathop{\text{arg min}}_{\mathcal{S}}$

Dataset distillation is:

- Optimise over the distilled dataset S:
- such that when you train on S,

 $\underset{\mathcal{S}}{\operatorname{arg\,min}} \operatorname{TrainingLoss}_{\mathcal{S}}(\theta))$

Dataset distillation is:

```
• Optimise over the distilled dataset S: (Outer loop)
```

- · you get good performance. (Test error)

```
\operatorname*{arg\,min}_{\mathcal{S}} \mathrm{TestError}(\operatorname*{arg\,min}_{\theta} \mathrm{TrainingLoss}_{\mathcal{S}}(\theta))
```

Dataset distillation is:

- Optimise over the distilled dataset S: (Outer loop)
- such that when you train on S, (Inner loop)
- · you get good performance. (Test error)

$$\operatorname*{arg\,min}_{\mathcal{S}}\operatorname{TestError}(\operatorname*{arg\,min}_{\theta}\operatorname{TrainingLoss}_{\mathcal{S}}(\theta))$$

Limitations

- 1. Computational scaling
- 2. Poor mathematical interpretation
- 3. Questionable architecture generalisation
- 4. No mathematical guarantees

$$\underset{\mathcal{S}}{\operatorname{arg\,min}}\operatorname{TestError}(\underset{\theta}{\operatorname{arg\,min}}\operatorname{TrainingLoss}_{\mathcal{S}}(\theta)) \tag{DD}$$

Table 1: Test accuracy with 10 images per class.

| Dataset | Original | |
|-----------|----------|--|
| CIFAR-10 | 84.8 | |
| CIFAR-100 | 56.2 | |

$$\underset{S}{\operatorname{arg \, min \, TestError}}(\underset{\theta}{\operatorname{arg \, min \, TrainingLoss}}_{\mathcal{S}}(\theta)) \tag{DD}$$

Gradient matching: replace the (intractable) training loss minimiser with a normal training regime

$$\underset{\mathcal{S}}{\operatorname{arg\,min\,TestError}}(\operatorname{Train_N_Epochs}_{\mathcal{S}}(\theta)) \tag{GM}$$

Table 1: Test accuracy with 10 images per class.

| Dataset | Original | GM | |
|-----------|----------|------|--|
| CIFAR-10 | 84.8 | 44.9 | |
| CIFAR-100 | 56.2 | 25.3 | |

Distribution matching: match the distributions of the synthetic \mathcal{S} and the training \mathcal{T} after passing through randomly initialised neural networks

$$\arg\min_{\mathcal{S}} \mathbb{E}_{\theta \sim p_{\theta}} \| \frac{1}{|\mathcal{T}|} \sum_{\mathbf{x} \in \mathcal{T}} \psi_{\theta}(\mathbf{x}) - \frac{1}{|\mathcal{S}|} \sum_{\mathbf{x} \in \mathcal{S}} \psi_{\theta}(\mathbf{x}) \|^{2}$$
 (DM)

Table 1: Test accuracy with 10 images per class.

| Dataset | Original | GM | DM |
|-----------|----------|------|------|
| CIFAR-10 | 84.8 | 44.9 | 48.9 |
| CIFAR-100 | 56.2 | 25.3 | 29.7 |

Distribution matching: match the distributions of the synthetic \mathcal{S} and the training \mathcal{T} after passing through randomly initialised neural networks

 Inspired by maximum mean discrepancy of the "RKHS of neural networks"

$$\arg\min_{\mathcal{S}} \mathbb{E}_{\theta \sim p_{\theta}} \| \frac{1}{|\mathcal{T}|} \sum_{\mathbf{x} \in \mathcal{T}} \psi_{\theta}(\mathbf{x}) - \frac{1}{|\mathcal{S}|} \sum_{\mathbf{x} \in \mathcal{S}} \psi_{\theta}(\mathbf{x}) \|^{2}$$
 (DM)

Table 1: Test accuracy with 10 images per class.

| Dataset | Original | GM | DM | |
|-----------|----------|------|------|--|
| CIFAR-10 | 84.8 | 44.9 | 48.9 | |
| CIFAR-100 | 56.2 | 25.3 | 29.7 | |

Matching training trajectories: minimise the parameter distance when training with ${\cal S}$ vs full training ${\cal T}$

$$\operatorname*{arg\,min}_{\mathcal{S}} \mathbb{E}_{\theta_0 \sim p_\theta} \sum_{t=1}^{\mathcal{T}-M} \frac{\|\theta_{t+N}^{\mathcal{S}} - \theta_{t+M}^{\mathcal{T}}\|^2}{\|\theta_{t+M}^{\mathcal{T}} - \theta_{t}^{\mathcal{T}}\|^2}$$

Table 1: Test accuracy with 10 images per class.

| Dataset | Original | GM | DM | MTT |
|-----------|----------|------|------|------|
| CIFAR-10 | 84.8 | 44.9 | 48.9 | 65.3 |
| CIFAR-100 | 56.2 | 25.3 | 29.7 | 40.1 |

Matching training trajectories: minimise the parameter distance when training with ${\cal S}$ vs full training ${\cal T}$

· Requires pretraining "expert models" from many initialisations

$$\operatorname*{arg\,min}_{\mathcal{S}} \mathbb{E}_{\theta_0 \sim p_{\theta}} \sum_{t=1}^{T-M} \frac{\|\theta_{t+N}^{\mathcal{S}} - \theta_{t+M}^{\mathcal{T}}\|^2}{\|\theta_{t+M}^{\mathcal{T}} - \theta_{t}^{\mathcal{T}}\|^2}$$

Table 1: Test accuracy with 10 images per class.

| Dataset | Original | GM | DM | MTT |
|-----------|----------|------|------|------|
| CIFAR-10 | 84.8 | 44.9 | 48.9 | 65.3 |
| CIFAR-100 | 56.2 | 25.3 | 29.7 | 40.1 |

Synthetic images from bi-level methods

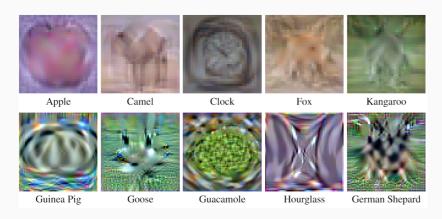


Figure 1: What are these? Directly optimised in image space.

Feature distribution after distillation

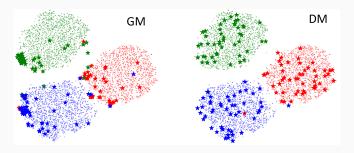


Figure 2: Distribution of synthetic images in CIFAR-10 [Zhao and Bilen '23].

Clustering effect?

• Features from distribution matching are "more uniform" than gradient matching.

Seems like dealing with the data is better than forcing a particular architecture.

If distribution matching is "like clustering" $\stackrel{?}{\rightarrow}$ use clustering for DD?

If distribution matching is "like clustering" $\stackrel{?}{\rightarrow}$ use clustering for DD?

Yes!

- 1. Computational scaling
 - · Removes the inner or outer problem.
 - Uses only one expert model.

If distribution matching is "like clustering" $\stackrel{?}{\rightarrow}$ use clustering for DD?

Yes!

- 1. Computational scaling
 - · Removes the inner or outer problem.
 - · Uses only **one** expert model.
- 2. Poor mathematical interpretation
 - · As a distribution approximation problem.

If distribution matching is "like clustering" $\stackrel{?}{\rightarrow}$ use clustering for DD?

Yes!

- 1. Computational scaling
 - · Removes the inner or outer problem.
 - Uses only **one** expert model.
- 2. Poor mathematical interpretation
 - · As a distribution approximation problem.
- 3. Questionable architecture generalisation
 - · Will be architecture independent.

If distribution matching is "like clustering" $\stackrel{?}{\rightarrow}$ use clustering for DD?

Yes!

- 1. Computational scaling
 - · Removes the inner or outer problem.
 - Uses only **one** expert model.
- 2. Poor mathematical interpretation
 - · As a distribution approximation problem.
- 3. Questionable architecture generalisation
 - · Will be architecture independent.
- 4. No mathematical guarantees
 - · Using distributional convergence.

If distribution matching is "like clustering" $\stackrel{?}{\rightarrow}$ use clustering for DD?

Yes!

- 1. Computational scaling
 - · Removes the inner or outer problem.
 - · Uses only **one** expert model.
- 2. Poor mathematical interpretation
 - · As a distribution approximation problem.
- 3. Questionable architecture generalisation
 - · Will be architecture independent.
- 4. No mathematical guarantees
 - · Using distributional convergence.
- 5. Strange looking distilled images
 - Force the distilled images to look good using generative models.

"Disentangled dataset distillation"

How to use clustering for DD?

 Direct clustering of the big images does not work. (curse of dimensionality)

"Disentangled dataset distillation"

How to use clustering for DD?

 Direct clustering of the big images does not work. (curse of dimensionality)

How about using an encoder-decoder structure? Then you can cluster in the latent space, and you have generative capabilities!

"Disentangled dataset distillation"

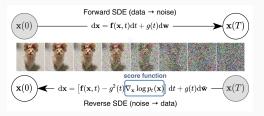
How to use clustering for DD?

 Direct clustering of the big images does not work. (curse of dimensionality)

How about using an encoder-decoder structure? Then you can cluster in the latent space, and you have generative capabilities!

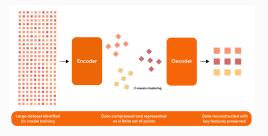
What you want: latent diffusion models

· Decoders given by conditional diffusion models



Algorithm: you have a encoder-decoder pair $(\mathcal{E}, \mathcal{D})$.

1. Take your entire training dataset \mathcal{T} , and encode it to get a dataset of latent variables $\mathcal{Z} = \mathcal{E}(\mathcal{T})$



- 1. Take your entire training dataset \mathcal{T} , and encode it to get a dataset of latent variables $\mathcal{Z} = \mathcal{E}(\mathcal{T})$
- 2. Cluster your latent variables in each class, say with k-means. You get some "centroid" latents in each class $\hat{\mathcal{Z}}$.



Algorithm: you have a encoder-decoder pair $(\mathcal{E}, \mathcal{D})$.

- 1. Take your entire training dataset \mathcal{T} , and encode it to get a dataset of latent variables $\mathcal{Z} = \mathcal{E}(\mathcal{T})$
- 2. Cluster your latent variables in each class, say with k-means. You get some "centroid" latents in each class $\hat{\mathcal{Z}}$.
- 3. Decode these clustered latents \Rightarrow Synthetic dataset $\mathcal{S} = \mathcal{D}(\hat{\mathcal{Z}})$.

 $S = \mathsf{decode} \circ \mathsf{cluster} \circ \mathsf{encode}(\mathcal{T})$



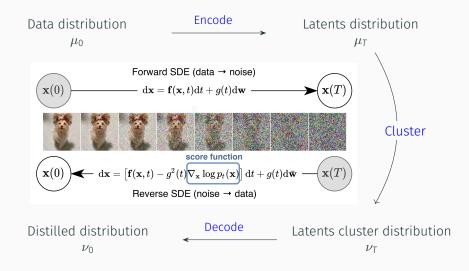
Algorithm: you have a encoder-decoder pair $(\mathcal{E}, \mathcal{D})$.

- 1. Take your entire training dataset \mathcal{T} , and encode it to get a dataset of latent variables $\mathcal{Z} = \mathcal{E}(\mathcal{T})$
- 2. Cluster your latent variables in each class, say with k-means. You get some "centroid" latents in each class $\hat{\mathcal{Z}}$.
- 3. Decode these clustered latents \Rightarrow Synthetic dataset $\mathcal{S} = \mathcal{D}(\hat{\mathcal{Z}})$.

 $S = \mathsf{decode} \circ \mathsf{cluster} \circ \mathsf{encode}(\mathcal{T})$

Oops, someone has done this before [Su et al., CVPR'24].





DDOQ

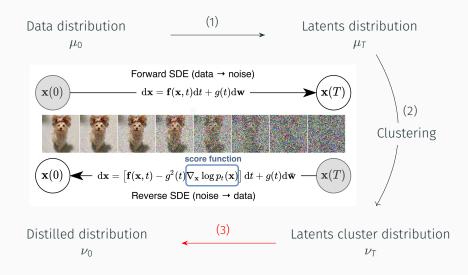
A closer look: decoding

3. Decode these clustered latents $S = \mathcal{D}(\hat{Z})$.

What is the decoding the clusters actually doing?

A look at diffusion models.

Generation from latents



Consistency

Proposition

For the VESDE/VPSDE and any initial data distribution $\mu \in \mathcal{P}_2(\mathbb{R}^d)$ with compact support bounded by R>0, the backwards diffusion process is well posed. Suppose that there are two distributions μ_T, ν_T at time T that undergo the reverse diffusion process (with fixed initial reference measure μ) up to time $t=\delta\in(0,T)$ to produce distributions μ_δ, ν_δ . There exists a (universal explicit) constant $C=C(\delta,T,R,d)\in(0,+\infty)$ such that if $f:\mathbb{R}^d\to\mathbb{R}^n$ is an L-Lipschitz function, then the difference in expectation satisfies

$$\|\mathbb{E}_{\mu_{\delta}}[f] - \mathbb{E}_{\nu_{\delta}}[f]\| \leq CLW_2(\mu_T, \nu_T).$$

Proof: basically stochastic Gronwall's inequality.

Consistency

Proposition (informal)

Suppose you have a bounded data distribution μ_0 , say of images. Then you pass it through your forward stochastic process (noising), up to μ_T . If you have another distribution ν_T that approximates μ_T , then after passing both through the reverse stochastic process (generation), your generated distributions ν_δ will approximate μ_δ .

Consistency

Proposition (informal)

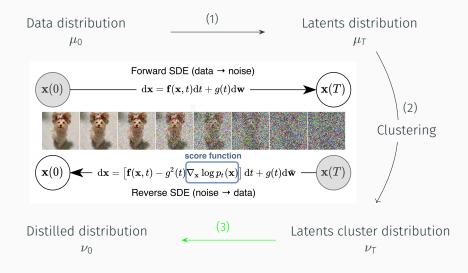
Suppose you have a bounded data distribution μ_0 , say of images. Then you pass it through your forward stochastic process (noising), up to μ_T . If you have another distribution ν_T that approximates μ_T , then after passing both through the reverse stochastic process (generation), your generated distributions ν_δ will (weakly) approximate μ_δ .

If you are close in the latent space, then you are close after decoding/generation.

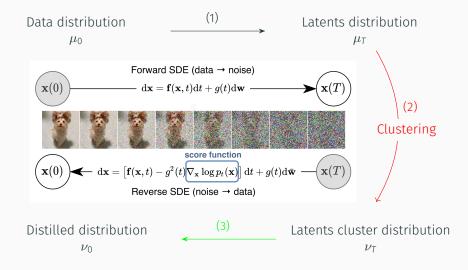
Consequences:

- Gradient steps when training with the distilled dataset are automatically similar to those on the full training dataset (supposedly).
- · No need to enforce through training!

Generation from latents



Generation from latents



An even closer look: clustering

2. Cluster your latent variables in each class, say with k-means. You get some "centroid" latents in each class $\hat{\mathcal{Z}}$.

What is clustering actually doing?

Clustering is optimal quantisation!

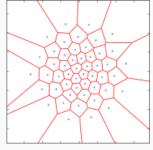
Definition

For a probability measure μ , the **quadratic distortion** of a set of points $\{x_1, ..., x_K\}$ is the μ -average distance to the set:

$$\mathcal{G}: (x_1, ..., x_K) \mapsto \int_{\mathbb{R}^d} \min_i \|x - x_i\|^2 \, \mu(\mathrm{d}x) = \mathbb{E}_{X \sim \mu}[\min_i \|X - x_i\|^2]. \tag{1}$$

The **optimal quantisation** of μ (at level K) is a set of points $\{x_1,...,x_K\}$ that minimizes the quadratic distortion.

Fact: there is a 1-1 correspondence with optimal quantisers $\{x_1,...,x_K\}$ and Wasserstein-minimisers with finite support.



Example quantisation of a Gaussian

Relating optimal quantisers to Wasserstein minimisers

Proposition

For a set of points $\{x_1,...,x_K\}$, the empirical measure $\nu = \sum_i w_i \delta(x_i)$ minimising $\mathcal{W}_2(\mu,\nu)$ satisfying supp $\nu \subset \{x_1,...,x_K\}$ is given by $\nu = \sum_i \mu(C_i)\delta(x_i)$.

- Sets of points
 ↔ approximating measures (automatically)
- Rates (Zador's theorem): $W_2(\nu_K, \mu) \sim \mathcal{O}(K^{-1/d})$ as number of points $K \to \infty$.

 $^{{}^{1}}C_{i} \subset \mathbb{R}^{d}$ are the Voronoi cells, set of points closest to x_{i} .

Clustering is optimal quantisation!

Guess what is an algorithm for solving optimal quantisation?

 $^{^2}$ Actually, a slight modification thereof: the competitive learning vector quantisation (CLVQ) algorithm.

Clustering is optimal quantisation!

Guess what is an algorithm for solving optimal quantisation? k-means clustering!² Produces $\{x_1, ..., x_K\}$.

Previous approach: take the approximating measure to be

$$\nu = \frac{1}{K} \sum_{i=1}^{K} \delta(\mathbf{x}_i) \quad (\approx \frac{1}{|\mathcal{T}|} \sum_{u \in \mathcal{T}} \delta(u) \approx \mu)$$

Key: finding $\{x_i\}$ is not enough! Actually, we need to optimise for $w_i > 0, x_i \in \mathbb{R}^d$ in

$$\nu = \arg\min_{w,x} \mathcal{W}_2 \left(\sum_{i=1}^K w_i \delta(x_i), \mu \right).$$

²Actually, a slight modification thereof: the competitive learning vector quantisation (CLVQ) algorithm.

A tiny modification

So we want small W_2 distance to give good approximation of the latents $\nu_T \approx \mu_T$:

$$\nu = \arg\min_{w,x} \mathcal{W}_2 \left(\sum_{i=1}^K w_i \delta(x_i), \mu \right).$$

How do we get the $w_i \approx \mu(C_i)$?

A tiny modification

So we want small W_2 distance to give good approximation of the latents $\nu_T \approx \mu_T$:

$$\nu = \underset{w,x}{\operatorname{arg\,min}} \mathcal{W}_2\left(\sum_{i=1}^K w_i \delta(x_i), \mu\right).$$

How do we get the $w_i \approx \mu(C_i)$?

Answer: k-means does this automatically! Produces online approximations to $\mu(C_i)$.

We get better approximation (in measures) with no additional computation.

Why do we need the coefficients?

If w_i are uniform $w_i \equiv 1/K$, then this is the Wasserstein barycentre problem (as opposed to optimal quantisation)³

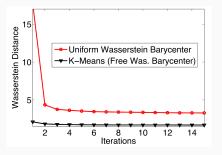
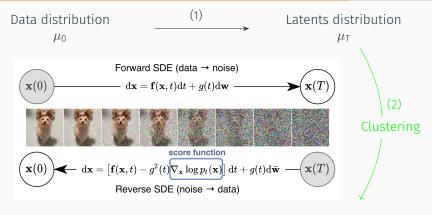


Figure 3: The Wasserstein distance W_2 is higher for the Wasserstein barycentre. Courtesy of ³.

³[Cuturi & Doucet ICML '14].

Generation from latents



Distilled distribution
$$\begin{array}{c} \text{ (3)} \\ \nu_0 \end{array} \qquad \text{Latents cluster distribution}$$

We have the optimal quantisers $\nu_T \approx \mu_T$ from (2), which implies $\nu_0 \approx \mu_0$ from (3).

Algorithm: you have a encoder-decoder pair $(\mathcal{E}, \mathcal{D})$.

1. Take your entire training dataset \mathcal{T} , and encode it to get a dataset of latent variables $\mathcal{Z} = \mathcal{E}(\mathcal{T})$

- 1. Take your entire training dataset \mathcal{T} , and encode it to get a dataset of latent variables $\mathcal{Z} = \mathcal{E}(\mathcal{T})$
- 2. Cluster your latent variables in each class, say with k-means. You get some "centroid" latents in each class $\hat{\mathcal{Z}}$, and some corresponding weights $W = \{w_z\}_{z \in \hat{\mathcal{Z}}}$.

- 1. Take your entire training dataset \mathcal{T} , and encode it to get a dataset of latent variables $\mathcal{Z} = \mathcal{E}(\mathcal{T})$
- 2. Cluster your latent variables in each class, say with k-means. You get some "centroid" latents in each class $\hat{\mathcal{Z}}$, and some corresponding weights $W = \{w_z\}_{z \in \hat{\mathcal{Z}}}$.
- 3. Decode these clustered latents $S = \mathcal{D}(\hat{Z})$.

- 1. Take your entire training dataset \mathcal{T} , and encode it to get a dataset of latent variables $\mathcal{Z} = \mathcal{E}(\mathcal{T})$
- 2. Cluster your latent variables in each class, say with k-means. You get some "centroid" latents in each class $\hat{\mathcal{Z}}$, and some corresponding weights $W = \{w_z\}_{z \in \hat{\mathcal{Z}}}$.
- 3. Decode these clustered latents $S = \mathcal{D}(\hat{Z})$.
- 4. When training with the distilled dataset (S, W), multiply the gradients of each training point with the corresponding weight.

$$\mathcal{T} \underset{1}{\mapsto} \mathcal{Z} \underset{2}{\mapsto} (\hat{\mathcal{Z}}, W) \underset{3}{\mapsto} (\mathcal{D}(\hat{\mathcal{Z}}), W) = (\mathcal{S}, W)$$

Experiments

ImageNet-1K: 1.2M images, 1000 classes, 150GB size.

Bi-level approaches: not possible.

- Previous maximum on 100K 128 \times 128 images split into 200 classes, \sim 600MB size.
- · Already requires 80GB of GPU memory for 10 images per class.

Experiments

Comparing the addition of the *k*-means weights:

| IPC | <i>k</i> -means weights? | ResNet-18 | ResNet-50 | ResNet-101 |
|-----|--------------------------|-----------|-----------|------------|
| 10 | Х | 27.9 | 33.5 | 34.2 |
| | ✓ | 33.1 | 34.4 | 36.7 |
| 50 | X | 55.2 | 62.4 | 63.4 |
| | ✓ | 56.2 | 62.5 | 63.6 |
| 100 | X | 59.3 | 65.4 | 66.5 |
| | ✓ | 60.1 | 65.9 | 66.7 |
| 200 | X | 62.6 | 67.8 | 68.1 |
| | ✓ | 63.4 | 68.0 | 68.6 |
| | | | | |

Table 2: Dataset distillation test accuracy without and with the *k*-means weights. Adding weights makes it uniformly better.

Comparing against other methods

Table 3: Comparison of top-1 classification performance on the ImageNet-1K dataset for baselines versus the proposed DDOQ method at various IPCs. We observe that DDOQ outperforms the previous SOTA method D⁴M, due to the addition of weights to the synthetic data. The maximum performance for all methods should be 69.8 as the soft labels are computed using a pre-trained ResNet-18 model.

| IPC | Method | ResNet-18 | ResNet-50 | ResNet-101 | IPC | Method | ResNet-18 | ResNet-50 | ResNet-101 |
|-----|--------------------|-----------------------|----------------------|----------------------|-----|--------------------|----------------------|----------------------------|-----------------------|
| 10 | TESLA | 7.7 | | | 50 | SRe ² L | 46.8 _{±0.2} | 55.6 _{±0.3} | 60.8 _{±0.5} |
| | SRe ² L | $21.3_{\pm 0.6}$ | $28.4_{\pm 0.1}$ | $30.9_{\pm 0.1}$ | | CDA | 53.5 | 61.3 | 61.6 |
| | RDED | $42.0_{\pm 0.1}$ | - | $48.3_{\pm 1.0}$ | | RDED | $56.5_{\pm 0.1}$ | - | $61.2_{\pm 0.4}$ |
| | D^4M | 27.9 | 33.5 | 34.2 | | D^4M | 55.2 | 62.4 | 63.4 |
| | DDOQ | $33.1_{\pm 0.60}$ | $34.4_{\pm 0.99}$ | $36.7_{\pm 0.80}$ | | DDOQ | $56.2_{\pm 0.07}$ | $\textbf{62.5}_{\pm 0.24}$ | $63.6_{\pm 0.13}$ |
| 100 | SRe ² L | 52.8 _{±0.3} | 61.0 _{±0.4} | 62.8 _{±0.2} | 200 | SRe ² L | 57.0 _{±0.4} | 64.6 _{±0.3} | 65.9 _{±0.3} |
| | CDA | 58.0 | 65.1 | 65.9 | | CDA | 63.3 | 67.6 | 68.4 |
| | D^4M | 59.3 | 65.4 | 66.5 | | D^4M | 62.6 | 67.8 | 68.1 |
| | DDOQ | 60.1 _{±0.15} | $65.9_{\pm 0.15}$ | $66.7_{\pm 0.06}$ | | DDOQ | $63.4_{\pm 0.08}$ | $68.0_{\pm 0.05}$ | 68.6 _{±0.08} |

What do the weights look like?

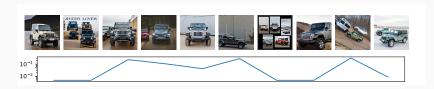


Figure 4: Example distilled images of the "jeep" class in ImageNet-1K along with their *k*-means weights. There is little to no features that can be used to differentiate the low and high weighted images, mainly due to the high fidelity of the diffusion model. However, the weights are indicative of the distribution of the training data in the latent space of the diffusion model.

Summary

- 1. Dataset distillation: turning big data into small data
- 2. By thinking a bit, we get
 - i) ... rid of any bi-level formulations;
 - ii) Theoretical perspective as to why clustering-style methods work
 - iii) State of the art performance for free

Summary

- 1. Dataset distillation: turning big data into small data
- 2. By thinking a bit, we get
 - i) ... rid of any bi-level formulations;
 - ii) Theoretical perspective as to why clustering-style methods work
 - iii) State of the art performance for free

What we really use/assume:

- Generative (diffusion) models as implicitly modelling the underlying data distribution (model expressiveness/trainability)
- Faithfulness of data distribution ↔ latent space mapping (manifold hypothesis)

Appendix

Generalisation

 Table 4: Generalization performance.

| Teacher Network | | Student Network | | | | | |
|-----------------|--------------------------|---------------------|---------------------|---------------------|------------------|--|--|
| | | ResNet-18 | MobileNet-V2 | EfficientNet-B0 | Swin-T | | |
| ResNet-18 | D ⁴ M DDOQ | 55.2 56.2 | 47.9 52.1 | 55.4 58.0 | 58.1 57.4 | | |
| MobileNet-V2 | D ⁴ M DDOQ | 47.6 47.7 | 42.9 45.6 | 49.8 52.5 | 58.9 56.3 | | |
| Swin-T | D ⁴ M DDOQ | 27.5 28.5 | 21.9 24.1 | 26.4 29.3 | 38.1 36.0 | | |

CLVQ

Algorithm 1: CLVQ

Data: initial cluster centers $x_1^{(0)}, ..., x_K^{(0)}$, step-sizes $(\gamma_i)_{i \geq 0}$, $i \leftarrow 0$ Initialize weights $\mathbf{w} = (w_1, ..., w_K) = (1/K, ..., 1/K)$; while not converged do

Sample
$$X_i \sim \mu$$
;
Select "winner" $k_{\min} \in \arg\min_{1 \leq k \leq K} \|X_i - X_k^{(i)}\|$;
Update $X_k^{(i+1)} \leftarrow (1 - \gamma_i) X_k^{(i)} + \gamma_i X_i$ if $k = k_{\min}$, otherwise $X_k^{(i+1)} \leftarrow X_k^{(i)}$;
Update weights $w_k \leftarrow (1 - \gamma_i) w_k + \gamma_i \mathbf{1}_{k=k_{\min}}$; $i \leftarrow i+1$;

end

Result: quantization $\nu_K = \sum_{i=1}^K w_i \delta(x_i^*)$