

# Provable and Efficient Dataset Distillation for Kernel Ridge Regression

Yilan Chen<sup>1</sup>, Wei Huang<sup>2</sup>, Tsui-Wei Weng<sup>1</sup>

<sup>1</sup>UCSD, <sup>2</sup>RIEKN AIP

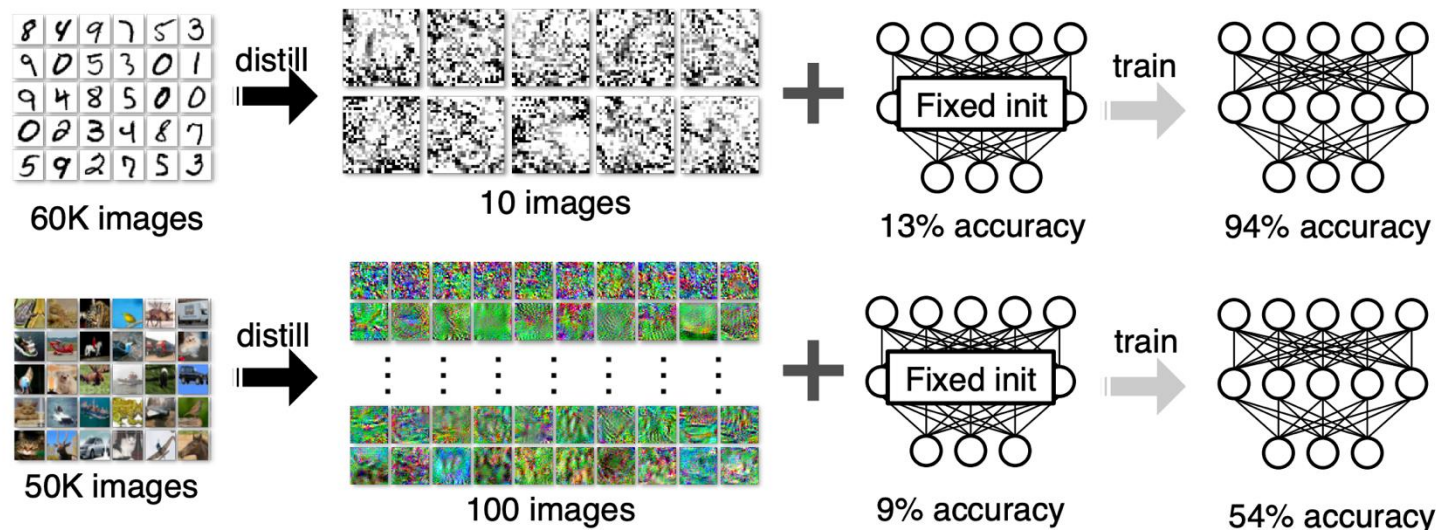
NeurIPS, December 2024

# Dataset Distillation

Motivation: Deep learning models are trained on increasingly larger datasets. It is crucial to reduce computational costs and improve data quality

- LLaMa 3 was pre-trained on over 15 trillion tokens
- Cost of training GPT-4 exceeded \$100 million

Dataset distillation: distill a large dataset into a small synthesized dataset such that models trained on it can achieve similar performance to those trained on the original dataset.



(a) Dataset distillation on MNIST and CIFAR10

# Dataset Distillation

Many empirical works:

- Meta-model Matching
- Gradient Matching
- Trajectory Matching
- Distribution Matching

# Dataset Distillation

Few theoretical analysis:

- [Izzo and Zou, 2023] linear ridge regression (LRR) needs  $d$  data points to recover original model's performance. Kernel ridge regression (KRR) needs  $n$  data
- [Maalouf et al., 2023] use Random Fourier Features (RFF) to approximate shift-invariant kernels and construct  $p$  distilled data for such RFF model.  $p$ : dimension of the RFF

This work:

- for KRR, one data point per class is already necessary and sufficient to recover the original model's performance in many settings
- analytical solutions for distilled dataset

# Problem formulation

Original dataset:  $(\mathbf{X}, \mathbf{Y}) \in \mathbb{R}^{d \times n} \times \mathbb{R}^{k \times n}$

Distilled dataset:  $(\mathbf{X}_S, \mathbf{Y}_S) \in \mathbb{R}^{d \times m} \times \mathbb{R}^{k \times m}$

Kernel ridge regression:  $f(\mathbf{x}) = \mathbf{W} \phi(\mathbf{x})$ , where  $\phi: \mathbb{R}^d \rightarrow \mathbb{R}^p$  and  $\mathbf{W} \in \mathbb{R}^{k \times p}$

$$\min_{\mathbf{W}} \|\mathbf{W} \phi(\mathbf{X}) - \mathbf{Y}\|_F + \lambda \|\mathbf{W}\|_F$$

Analytical solution:  $\mathbf{W} = \mathbf{Y} \phi_\lambda(\mathbf{X})^+$

$$\phi_\lambda(\mathbf{X})^+ = \begin{cases} (\phi(\mathbf{X})^T \phi(\mathbf{X}) + \lambda \mathbf{I}_n)^{-1} \phi(\mathbf{X}), & \lambda > \mathbf{0} \\ \phi(\mathbf{X})^+, & \lambda = \mathbf{0} \end{cases}$$

Similarly for distilled dataset model:  $f_S(\mathbf{x}) = \mathbf{W}_S \phi(\mathbf{x})$

Goal of dataset distillation: find  $(\mathbf{X}_S, \mathbf{Y}_S)$  such that  $\mathbf{W}_S = \mathbf{W}$

# Provable and Efficient Dataset Distillation for Kernel Ridge Regression

## Comparison with existing work

Table 1: Comparison with existing theoretical analysis of dataset distillation. The number of distilled data needed to recover original model’s performance and models analyzed. “-” means not applicable. For linear ridge regression (LRR) and kernel ridge regression (KRR) with subjective feature mapping, our results only need one distilled data per class ( $k \leq d$  in our setting), which is far less than the existing work [9, 21] that require  $n$  or  $p$  points. As an example,  $k = 10, d = 3072, n = 50000$  for CIFAR-10. The  $k, d, n$  of standard datasets are listed in Table 2.  $p$  is the dimension of feature mapping  $\phi : \mathbb{R}^d \mapsto \mathbb{R}^p$ .

	<b>LRR</b>	<b>Kernel ridge regression (KRR)</b>	
		surjective $\phi$	non-surjective $\phi$
Izzo and Zou [9]	$d$	-	$n$ (Gaussian Kernel)
Maalouf et al. [21]	-	-	$p$ (Shift-invariant Kernels)
<b>Our work</b>	$k, (k \leq d)$	$k, (k \leq p)$ (Invertible NNs, FCNN, CNN, Random Fourier Features)	$p$ in general (Deep nonlinear NNs). $k + 1$ for deep linear NNs

# Main result 1: Linear Ridge Regression

## Linear ridge regression

### Theorem

- When  $m < k$  there is no  $\mathbf{X}_S$  can guarantee  $\mathbf{W}_S = \mathbf{W}$  unless the columns of  $\mathbf{W}$  are in the range space of  $\mathbf{Y}_S$ .
- When  $m \geq k$  and  $\mathbf{Y}_S$  is rank  $k$ , let  $r = \min(m, d)$  and take  $\mathbf{D} = \mathbf{Y}_S^+ \mathbf{W} + (\mathbf{I}_m - \mathbf{Y}_S^+ \mathbf{Y}_S) \mathbf{Z}$ , where  $\mathbf{Z} \in \mathbb{R}^{d \times m}$  is any matrix. Suppose the reduced SVD of  $\mathbf{D}$  is  $\mathbf{D} = \mathbf{V} \text{diag}(\sigma'_1, \dots, \sigma'_r) \mathbf{U}^T$  with  $\sigma'_1 \geq \dots \geq \sigma'_r \geq 0$ .

1.  $\lambda_S > 0$ :  $\mathbf{W}_S = \mathbf{W}$  if and only if for any  $\mathbf{D}$  defined above,  $\lambda_S < \frac{1}{4\sigma'_1}$  and  $\mathbf{X}_S =$

$$\mathbf{U} \text{diag}(\sigma_1, \dots, \sigma_r) \mathbf{V}^T \text{ where } \sigma_i = \begin{cases} 0, & \text{if } \sigma'_i = 0, \\ \frac{1 \pm \sqrt{1 - 4\lambda_S \sigma'_i}}{2\sigma'_i}, & \text{otherwise.} \end{cases}$$

2.  $\lambda_S = 0$ :  $\mathbf{W}_S = \mathbf{W}$  if and only if  $\mathbf{X}_S = \mathbf{D}^+$  for any  $\mathbf{D}$  defined above.

# Main result 1: Linear Ridge Regression

- One distilled data per class is necessary and sufficient for  $\mathbf{W}_S = \mathbf{W}$ .
- Intuitively, original dataset  $(\mathbf{X}, \mathbf{Y})$  is compressed into  $\mathbf{X}_S$  through original model's parameter  $\mathbf{W}$ .
  - When  $m = k$ , only one solution. When  $\lambda_S = 0$ ,  $\mathbf{X}_S = (\mathbf{Y}_S^+ \mathbf{W})^+$ .
  - When  $m > k$ , infinitely many distilled datasets since  $\mathbf{Z}$  is a free variable to choose
  - When  $m = n$ ,  $(\mathbf{X}, \mathbf{Y})$  is a distilled dataset for itself.
  - When  $m > n$ , can generate more data than original dataset.

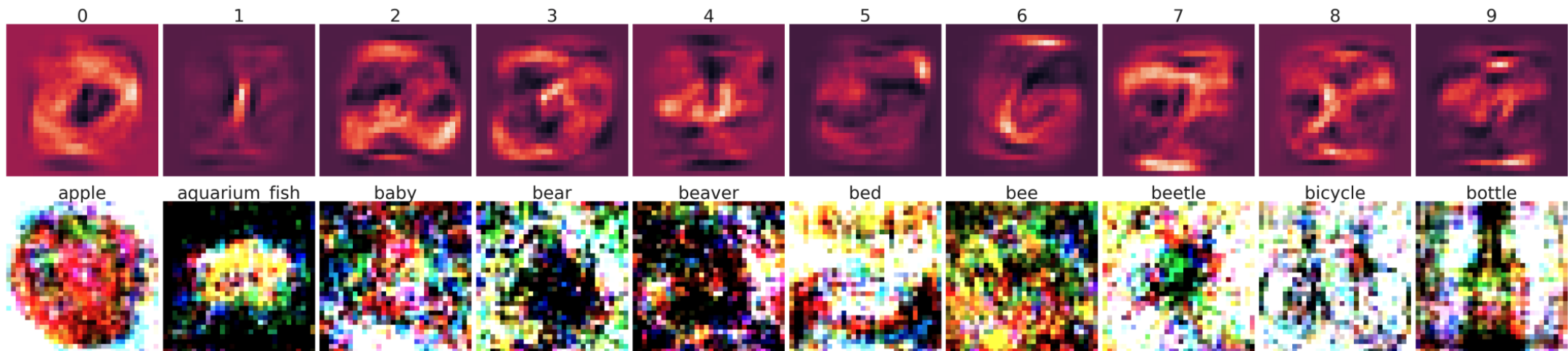
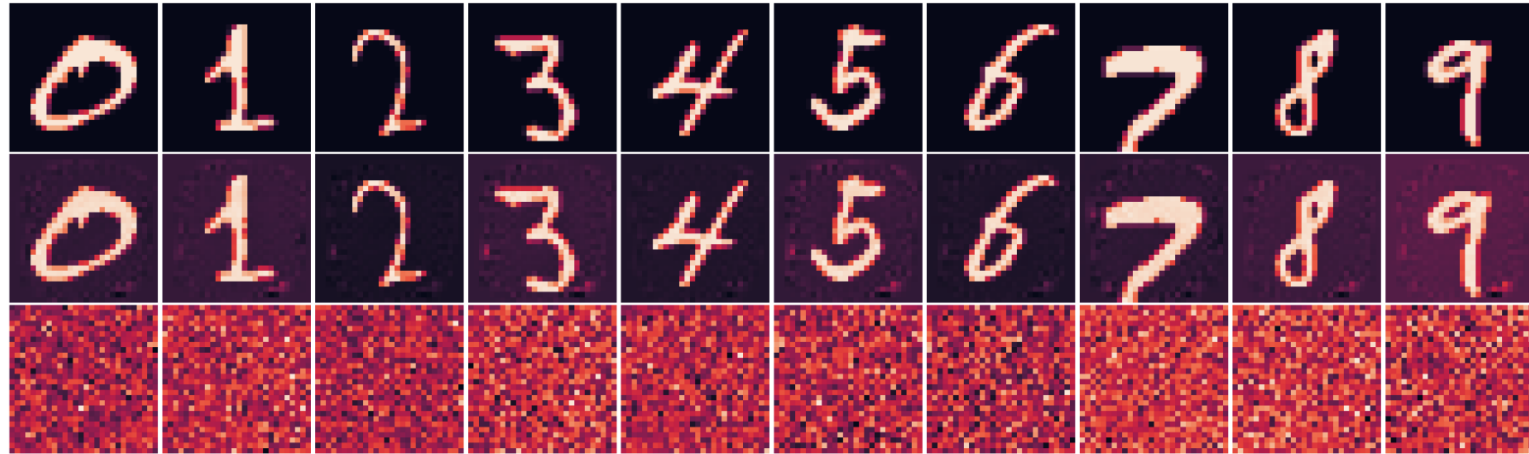


Figure 1: Distilled data of MNIST (first row) and CIFAR-100 (second row) for LRR when  $m = k$ .



# Main result 1: Linear Ridge Regression

1. Find realistic distilled data that is close to original data by solving closed-form solution
2. Generate distilled data from random noise



(a) MNIST with IPC=50.



(b) CIFAR-100 with IPC=5.

Figure 2: Initialized data (first row), distilled data generated from real images using techniques in Sec 4.2 (second row), and distilled data generated from random noise using techniques in Sec 4.1 (third row) for a LRR with  $m = 500$  on MNIST and CIFAR-100. IPC: images per class.

# Main result 2: Kernel Ridge Regression

- The results can be extended to KRR by replacing  $\mathbf{X}_S$  with  $\phi(\mathbf{X}_S)$
- When  $\phi$  is surjective or bijective, can always find a  $\mathbf{X}_S$  for a desired  $\phi(\mathbf{X}_S)$

Examples of surjective  $\phi: p \leq d$

1. Invertible NNs
2. Fully-connected NN (FCNN)
3. Convolutional Neural Network (CNN)
4. Random Fourier Features (RFF)

# Main result 2: Kernel Ridge Regression

## Non-surjective Feature Mapping

- For non-surjective  $\phi$  such as deep nonlinear NNs, one data per class is generally not sufficient as long as  $(\mathbf{Y}_S^+ \mathbf{W})^+$  is not in the range space of  $\phi$
- For deep linear NNs, we show  $m = k + 1$  can be sufficient under certain conditions

For surjective  $\phi$ , our algorithm outperforms previous work such as KIP while being significantly more efficient

Table 4: Comparison between our algorithm and KIP.

Dataset	IPC	KIP [25]		Ours		Speedup over KIP $\uparrow$
		Accuracy $\uparrow$	Cost $\downarrow$ (GPU Sec.)	Accuracy $\uparrow$	Cost $\downarrow$ (GPU Sec.)	
MNIST	1	93.44 $\pm$ 0.17	159	<b>93.72<math>\pm</math>0.14</b>	16	<b>9.9<math>\times</math></b>
	10	<b>93.75<math>\pm</math>0.10</b>	554	93.69 $\pm$ 0.17	16	<b>34.6<math>\times</math></b>
	50	<b>93.72<math>\pm</math>0.11</b>	3114	93.62 $\pm$ 0.24	16	<b>194.6<math>\times</math></b>
CIFAR-10	1	45.83 $\pm$ 0.29	225	<b>47.85<math>\pm</math>0.10</b>	21	<b>10.7<math>\times</math></b>
	10	47.50 $\pm$ 0.29	594	<b>47.76<math>\pm</math>0.12</b>	20	<b>29.7<math>\times</math></b>
	50	47.48 $\pm$ 0.20	3510	<b>47.77<math>\pm</math>0.06</b>	20	<b>175.5<math>\times</math></b>
CIFAR-100	1	20.08 $\pm$ 0.20	616	<b>21.58<math>\pm</math>0.15</b>	20	<b>30.8<math>\times</math></b>
	10	21.56 $\pm$ 0.16	9323	<b>21.59<math>\pm</math>0.15</b>	20	<b>466.1<math>\times</math></b>
	50	-	$\sim$ 396000	<b>21.58<math>\pm</math>0.13</b>	25	<b><math>\sim</math>15840.0<math>\times</math></b>

# Future works

- Determining minimum number of distilled data points required for non-surjective deep neural networks
- Extending our analysis to learnable feature mappings