

Course on Learning Low-Dim Representations for High-Dim Data
Design Deep Networks for Pursuing Low-D Representations

Yi Ma

EECS UCB and IDS/CS HKU

with Sam Buchanan, Qing Qu, Atlas Wang
John Wright, Yuqian Zhang, Zhihui Zhu

June 7, 2023

*“Everything should be made as simple as possible,
but not any simpler.”*

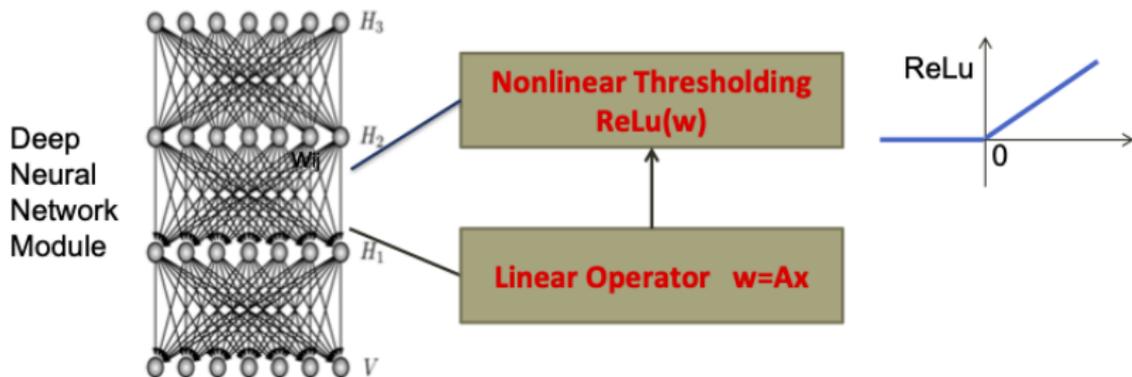
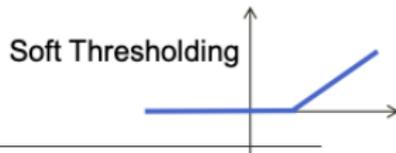
– *Albert Einstein*

Outline

- ① Objectives for Learning from Data
 - Precursors and Motivations
 - Linear and Discriminative Representation (LDR)
- ② Closed-Form Information-Theoretical Measure for LDR
 - Principle of Maximizing Coding Rate Reduction (MCR²)
 - Experimental Verification
- ③ White-Box Deep Networks from Optimizing Rate Reduction
 - Deep Networks as Projected Gradient Ascent
 - Convolution Networks from Shift Invariance
 - Experimental Results
- ④ Closed-Loop Transcription to an LDR (CTRL)
- ⑤ Conclusions and Open Problems

ISTA: Sparse Recovery via ℓ^1 (Wright and Ma, 2022)**CONTEXT – Basic algorithm (ISTA)****Algorithm 8.1** Iterative Soft-Thresholding Algorithm (ISTA) for BPDN

- 1: **Problem:** $\min_{\mathbf{x}} \frac{1}{2} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2 + \lambda \|\mathbf{x}\|_1$, given $\mathbf{y} \in \mathbb{R}^d$, $\mathbf{A} \in \mathbb{R}^{d \times n}$.
- 2: **Input:** $\mathbf{x}_0 \in \mathbb{R}^n$ and $L \geq \lambda_{\max}(\mathbf{A}^T \mathbf{A})$.
- 3: **while** \mathbf{x}_k not converged ($k = 1, 2, \dots$) **do**
- 4: $\mathbf{w}_k \leftarrow \mathbf{x}_k - \frac{1}{L} \mathbf{A}^T (\mathbf{A}\mathbf{x}_k - \mathbf{y})$.
- 5: $\mathbf{x}_{k+1} \leftarrow \text{soft}(\mathbf{w}_k, \lambda/L)$.
- 6: **end while**
- 7: **Output:** $\mathbf{x}_* \leftarrow \mathbf{x}_k$.



Learned ISTA (Gregor and LeCun, ICML 2010)

CONTEXT – Learned ISTA (LISTA)

If only interested in one instance: $y = Ax$ AND with many training data: $\{(y_i, x_i)\}$.
We can **optimize the optimization path** of ISTA using supervised learning:

Algorithm 3 LISTA::fprop

```

LISTA :: fprop( $X, Z, W_e, S, \theta$ )
;; Arguments are passed by reference.
;; variables  $Z(t), C(t)$  and  $B$  are saved for bprop.
 $B = W_e X; Z(0) = h_\theta(B)$ 
for  $t = 1$  to  $T$  do
   $C(t) = B + SZ(t-1)$ 
   $Z(t) = h_\theta(C(t))$ 
end for
 $Z = Z(T)$ 

```

Algorithm 4 LISTA::bprop

```

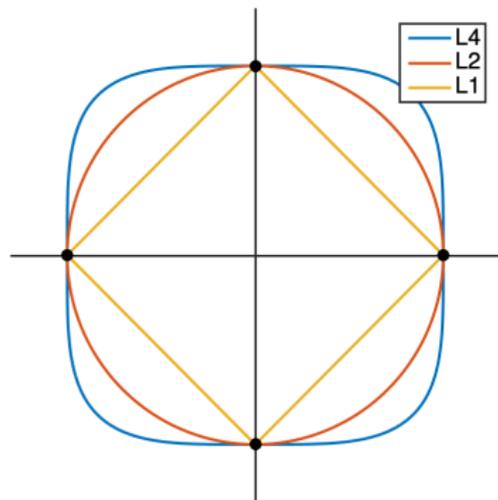
LISTA :: bprop( $Z^*, X, Z, W_e, S, \theta, \delta X, \delta W_e, \delta S, \delta \theta$ )
;; Arguments are passed by reference.
;; Variables  $Z(t), C(t)$ , and  $B$  were saved in fprop.
Initialize:  $\delta B = 0; \delta S = 0; \delta \theta = 0$ 
 $\delta Z(T) = (Z(T) - Z^*)$ 
for  $t = T$  down to 1 do
   $\delta C(t) = h'_\theta(C(t)) \cdot \delta Z(t)$ 
   $\delta \theta = \delta \theta - \text{sign}(C(t)) \cdot \delta C(t)$ 
   $\delta B = \delta B + \delta C(t)$ 
   $\delta S = \delta S + \delta C(t) Z(t-1)^T$ 
   $\delta Z(t-1) = S^T \delta C(t)$ 
end for
 $\delta B = \delta B + h'_\theta(B) \cdot \delta Z(0)$ 
 $\delta \theta = \delta \theta - \text{sign}(B) \cdot h'_\theta(B) \delta Z(0)$ 
 $\delta W_e = \delta B X^T; \delta X = W_e^T \delta B$ 

```

Sparsity: ℓ^0 or ℓ^1 Minimization versus ℓ^4 Maximization

Given $\mathbf{Y} = \mathbf{A}\mathbf{X}$, find \mathbf{A} such that \mathbf{X} is the sparsest:

$$\min_{\mathbf{A}} \|\mathbf{X}\|_0, \quad \mathbf{X} \doteq \mathbf{A}\mathbf{Y}, \quad \text{s.t.} \quad \mathbf{A} \in O(n; \mathbb{R}). \quad (1)$$



Maximizing ℓ^4 norm promotes sparsity or spikiness):

$$\arg \max_{\mathbf{q} \in \mathbb{S}^n} \|\mathbf{q}\|_4 \Leftrightarrow \arg \min_{\mathbf{q} \in \mathbb{S}^n} \|\mathbf{q}\|_0.$$

Figure: ℓ^1 -, ℓ^2 -, and ℓ^4 -spheres in \mathbb{R}^2

Sparse Dictionary Learning via ℓ^4 : the MSP Algorithm

Solve a sparsifying transform \mathbf{A} for $\mathbf{X} = \mathbf{A}\mathbf{Y}$ by solving the program:

$$\max_{\mathbf{A}} \phi(\mathbf{A}) \doteq \|\mathbf{A}\mathbf{Y}\|_4^4, \quad \text{s.t. } \mathbf{A} \in O(n; \mathbb{R}), \quad (2)$$

via projected gradient descent: $\mathbf{A}_{t+1} = \mathcal{P}_{O(n)}[\mathbf{A}_t + \alpha \nabla_{\mathbf{A}} \phi(\mathbf{A}_t)]$.

The Matching, Stretching, and Projection (MSP) Algorithm:

- 1: **Initialize** $\mathbf{A}_0 \in O(n, \mathbb{R})$ initialize \mathbf{A}_0 for iteration
- 2: **for** $t = 0, 1, \dots$
- 3: $\nabla_{\mathbf{A}} \phi(\mathbf{A}_t) = 4(\mathbf{A}_t \mathbf{Y})^{\circ 3} \mathbf{Y}^*$ Matching and Stretching
- 4: $\mathbf{U}\Sigma\mathbf{V}^* = \text{SVD}[\nabla_{\mathbf{A}} \phi(\mathbf{A}_t)]$
- 5: $\mathbf{A}_{t+1} = \mathbf{U}\mathbf{V}^*$ Project \mathbf{A} onto orthogonal group
- 6: **end for**
- 7: **output:** \mathbf{A}_* .

Complete dictionary learning via L4-Norm maximization over the orthogonal group, Zhai et. al. JMLR 2020.

“Deep Networks” as Power Iteration for Low-Dim

“Fixed point” (not gradient descent) interpretation:

$$\mathbf{A}_{t+1} = \mathcal{P}_{O(n)}[\nabla_{\mathbf{A}}\phi(\mathbf{A}_t)] = \mathcal{P}_{O(n)}[(\mathbf{A}_t\mathbf{Y})^{\circ 3}\mathbf{Y}^*].$$

Define “layer” operators and states iteratively:

$$\delta\mathbf{A}_{t+1} \doteq \mathbf{A}_{t+1}\mathbf{A}_t^* \quad \text{and} \quad \mathbf{Z}_{t+1} \doteq \mathbf{A}_{t+1}\mathbf{Y} = \delta\mathbf{A}_{t+1}\mathbf{Z}_t.$$

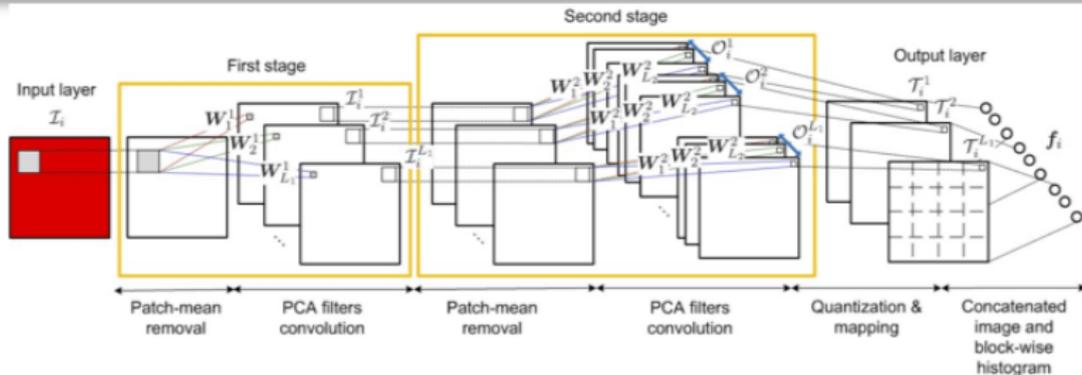
Forward-constructed “deep network” interpretation:

$$\delta\mathbf{A}_{t+1} = \mathcal{P}_{O(n)}[(\mathbf{Z}_t)^{\circ 3}\mathbf{Z}_t^*], \quad \mathbf{X} \leftarrow \delta\mathbf{A}_{t+1}\delta\mathbf{A}_t \cdots \delta\mathbf{A}_1\mathbf{Y}.$$

Table: Fixed Points: PCA (Power iteration), ICA (FastICA), and DL (MSP).

	Objectives	Constraint Sets	Algorithms
Power Iter.	$\varphi(\mathbf{a}) \doteq \frac{1}{2} \ \mathbf{a}^*\mathbf{Y}\ _2^2$	$\mathbf{w} \in \mathbb{S}^{n-1}$	$\mathbf{a}_{t+1} = \mathcal{P}_{\mathbb{S}^{n-1}}[\nabla_{\mathbf{a}}\varphi(\mathbf{a}_t)]$
FastICA	$\psi(\mathbf{a}) \doteq \frac{1}{4} \text{kurt}[\mathbf{a}^*\mathbf{y}]$	$\mathbf{a} \in \mathbb{S}^{n-1}$	$\mathbf{a}_{t+1} = \mathcal{P}_{\mathbb{S}^{n-1}}[\nabla_{\mathbf{a}}\psi(\mathbf{a}_t)]$
MSP	$\phi(\mathbf{A}) \doteq \frac{1}{4} \ \mathbf{A}^*\mathbf{Y}\ _4^4$	$\mathbf{A} \in St(k, n; \mathbb{R})$	$\mathbf{A}_{t+1} = \mathcal{P}_{St(k, n; \mathbb{R})}[\nabla_{\mathbf{A}}\phi(\mathbf{A}_t)]$

The Simplest Network Models - PCANet



Two layers forward-constructed without back-propagation:

- **simplest data-adaptive mapping (PCA),**
- **simplest nonlinear activation (binary),**
- **simplest pooling (histogram).**

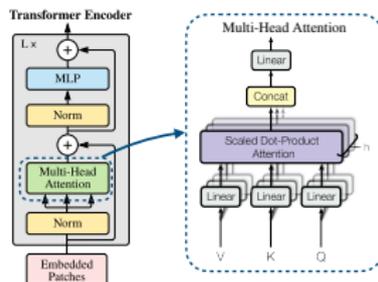
Recognition rates (%) on FERET dataset.

Probe sets	<i>Fb</i>	<i>Fc</i>	<i>Dup-I</i>	<i>Dup-II</i>	<i>Avg.</i>
LBP [18]	93.00	51.00	61.00	50.00	63.75
DMMA [25]	98.10	98.50	81.60	83.20	89.60
P-LBP [21]	98.00	98.00	90.00	85.00	92.75
POEM [26]	99.60	99.50	88.80	85.00	93.20
G-LQP [27]	99.90	100	93.20	91.00	96.03
LGBP-LGXP [28]	99.00	99.00	94.00	93.00	96.25
sPOEM+POD [29]	99.70	100	94.90	94.00	97.15
GOM [30]	99.90	100	95.70	93.10	97.18
PCANet-1 (Irn. CD)	99.33	99.48	88.92	84.19	92.98
PCANet-2 (Irn. CD)	99.67	99.48	95.84	94.02	97.25
PCANet-1	99.50	98.97	89.89	86.75	93.78
PCANet-2	99.58	100	95.43	94.02	97.26

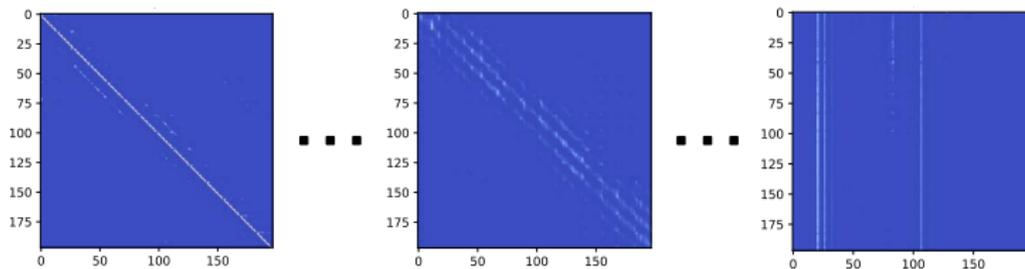
PCANet, Chan and Ma et. al. in *IEEE Trans. On Image Processing*, 2015

Pursuit of Sparsity via Transformer Layers

[Vaswani et al., 2017, Dosovitskiy et al., 2020] (see Sam Buchanan's lecture)



“Attention” matrix



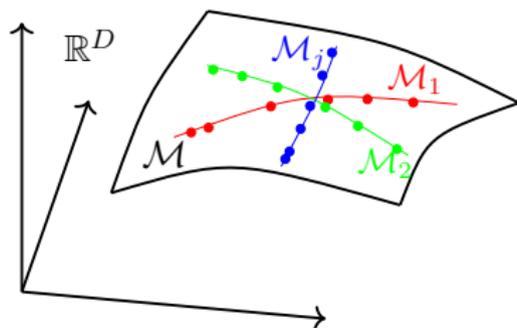
Earlier
Layers

Later
Layers

$$A = \text{Softmax} \left(\frac{qk^T}{\sqrt{D_h}} \right)$$

High-Dim Data with Mixed **Nonlinear** Low-Dim Structures

Figure: High-dimensional Real-World Data: data samples $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_m]$ in \mathbb{R}^D lying on a mixture of low-dimensional submanifolds $\mathbf{X} \subset \cup_{j=1}^k \mathcal{M}_j \subset \mathbb{R}^D$.



The main objective of learning from (samples of) such real-world data:
 seek a most **compact and structured** representation of the data.

Fitting Class Labels via a Deep Network

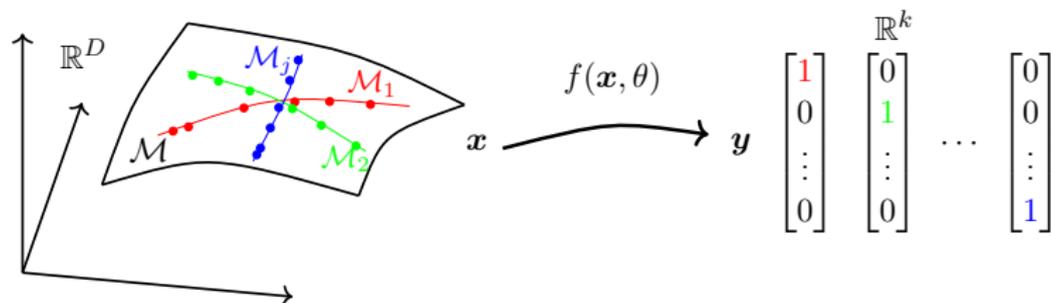


Figure: Black Box DNN for Classification: y is the class label of x represented as a “one-hot” vector in \mathbb{R}^k . To learn a nonlinear mapping $f(\cdot, \theta) : x \mapsto y$, say modeled by a deep network, using cross-entropy (CE) loss.

$$\min_{\theta \in \Theta} \text{CE}(\theta, \mathbf{x}, \mathbf{y}) \doteq -\mathbb{E}[\langle \mathbf{y}, \log[f(\mathbf{x}, \theta)] \rangle] \approx -\frac{1}{m} \sum_{i=1}^m \langle \mathbf{y}_i, \log[f(\mathbf{x}_i, \theta)] \rangle. \quad (3)$$

Prevalence of neural collapse during the terminal phase of deep learning training,
Papayan, Han, and Donoho, 2020.

Fitting Class Labels via a Deep Network

In a supervised setting, using cross-entropy (CE) loss:

$$\min_{\theta \in \Theta} \text{CE}(\theta, \mathbf{x}, \mathbf{y}) \doteq -\mathbb{E}[\langle \mathbf{y}, \log[f(\mathbf{x}, \theta)] \rangle] \approx -\frac{1}{m} \sum_{i=1}^m \langle \mathbf{y}_i, \log[f(\mathbf{x}_i, \theta)] \rangle. \quad (4)$$

Issues (an elephant in the room):

- A large deep neural networks can **fit arbitrary data and labels**.
- Statistical and geometric meaning of internal features **not clear**.
- Task/data-dependent and **not robust nor truly invariant**.

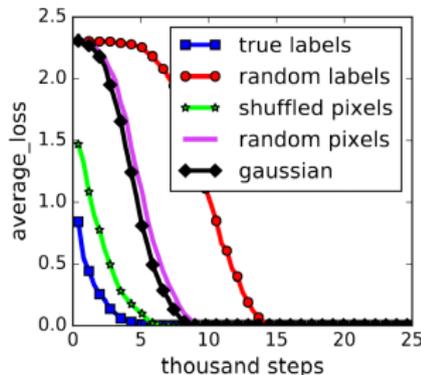


Figure: [Zhang et al, ICLR'17]

What did machines actually “learn” from doing this?

In terms of interpolating, extrapolating, or representing the data?

A Hypothesis: Information Bottleneck

[Tishby & Zaslavsky, 2015]

A feature mapping $f(\mathbf{x}, \theta)$ and a classifier $g(\mathbf{z})$ trained for downstream classification:

$$\mathbf{x} \xrightarrow{f(\mathbf{x}, \theta)} \mathbf{z}(\theta) \xrightarrow{g(\mathbf{z})} \mathbf{y}.$$

The IB Hypothesis: Features learned in a deep network trying to

$$\max_{\theta \in \Theta} \text{IB}(\mathbf{x}, \mathbf{y}, \mathbf{z}(\theta)) \doteq I(\mathbf{z}(\theta), \mathbf{y}) - \beta I(\mathbf{x}, \mathbf{z}(\theta)), \quad \beta > 0, \quad (5)$$

where $I(\mathbf{z}, \mathbf{y}) \doteq H(\mathbf{z}) - H(\mathbf{z}|\mathbf{y})$ and $H(\mathbf{z})$ is the entropy of \mathbf{z} .

- **Minimal** informative features \mathbf{z} that most correlate with the label \mathbf{y}
- Task and label-dependent, consequently sacrificing generalizability, robustness, or transferability

Gap between Theory and Practice (a Bigger Elephant)

For high-dimensional real data,

many statistical and information-theoretic concepts such as entropy, mutual information, K-L divergence, and maximum likelihood:

- curse of **dimensionality** for computation.
- ill-posed for **degenerate** distributions.
- lack guarantees with **finite** (or non-asymptotic) samples.

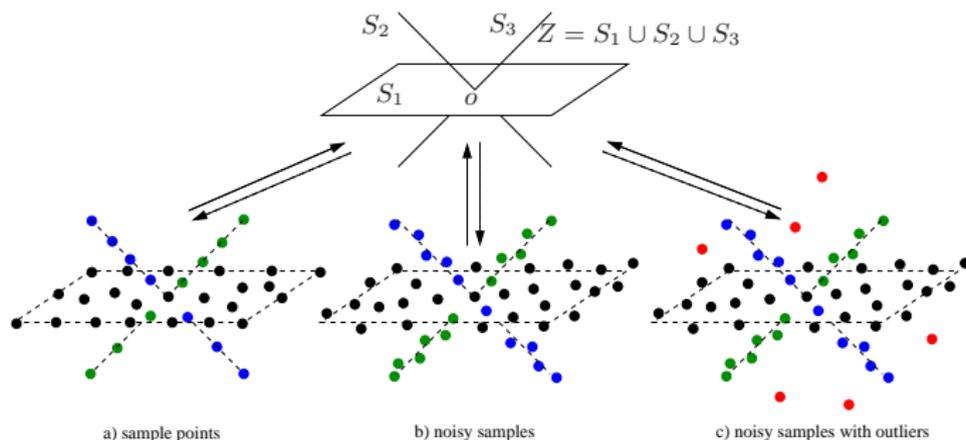
Reality check: principled formulations are replaced with approximate bounds, grossly simplifying assumptions, heuristics, even *ad hoc* tricks and hacks.

How to provide any performance guarantees at all?

A Principled Computational Approach

For high-dim data with mixed **low-dim** structures:

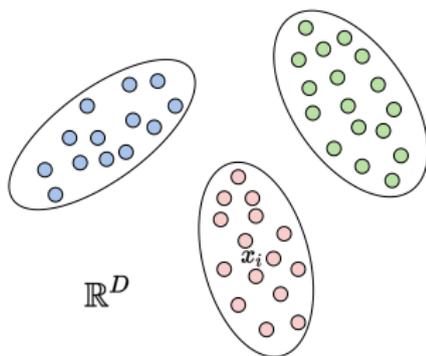
learn to compress, and compress to learn!



Generalized PCA for mixture of subspaces [Vidal, Ma, and Sastry, 2005]

1. Clustering Mixed Data (Interpolation)

Assume data $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m]$
 are i.i.d. samples from a mixture
 of distributions: $p(\mathbf{x}, \theta) = \sum_{j=1}^k \pi_j p_j(\mathbf{x}, \theta)$.



Classic approaches to cluster the data:
 the maximum-likelihood (ML) estimate
 via Expectation Maximization (EM):

$$\max_{\theta, \pi} \mathbb{E} \left[\log \left(\sum_{j=1}^k \pi_j p_j(\mathbf{x}, \theta) \right) \right] \approx \max_{\theta, \pi} \frac{1}{m} \sum_{i=1}^m \log \left(\sum_{j=1}^k \pi_j p_j(\mathbf{x}_i, \theta) \right).$$

Difficulties: ML is not well-defined when distributions are degenerate.

Clustering via Compression

[Yi Ma, Harm Derksen, Wei Hong, and John Wright, TPAMI'07]

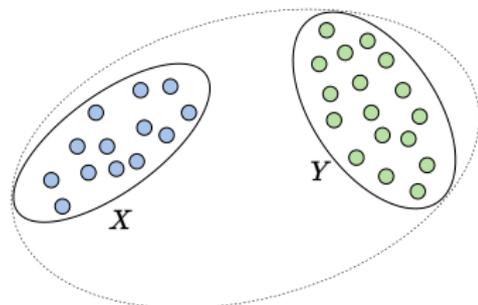
A Fundamental Idea:

Data belong to mixed low-dim structures should be compressible.

Cluster Criterion:

Whether the number of binary bits required to store the data is less (information gain):

$$\#bits(\mathbf{X} \cup \mathbf{Y}) \geq \#bits(\mathbf{X}) + \#bits(\mathbf{Y})?$$



"The whole is greater than the sum of the parts."
– Aristotle, 320 BC

Coding Length Function for Subspace-Like Data

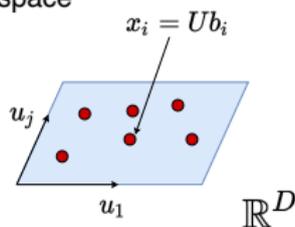
Theorem (Ma, TPAMI'07)

The number of bits needed to encode data $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m] \in \mathbb{R}^{D \times m}$ up to a precision $\|\mathbf{x} - \hat{\mathbf{x}}\|_2 \leq \epsilon$ is bounded by:

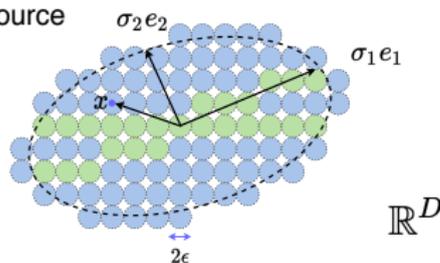
$$L(\mathbf{X}, \epsilon) \doteq \left(\frac{m + D}{2} \right) \log \det \left(\mathbf{I} + \frac{D}{m\epsilon^2} \mathbf{X} \mathbf{X}^\top \right).$$

This can be derived from constructively quantifying SVD of \mathbf{X} or by sphere packing $\text{vol}(\mathbf{X})$ as samples of a noisy Gaussian source.

Linear subspace



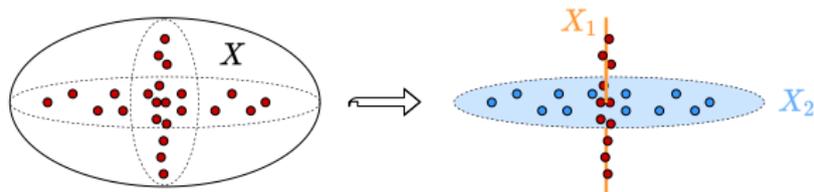
Gaussian source



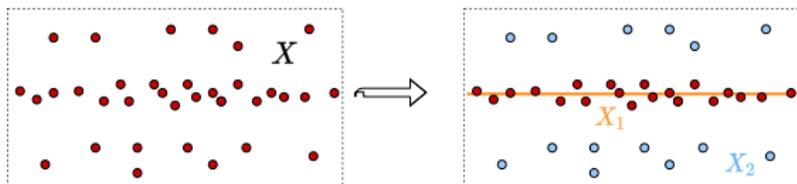
Cluster to Compress

$$L(\mathbf{X}) \geq L^c(\mathbf{X}) \doteq L(\mathbf{X}_1) + L(\mathbf{X}_2) + H(|\mathbf{X}_1|, |\mathbf{X}_2|)?$$

partitioning:



sifting:



A Greedy Algorithm

Seek a partition of the data $\mathbf{X} \rightarrow [\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_k]$ such that

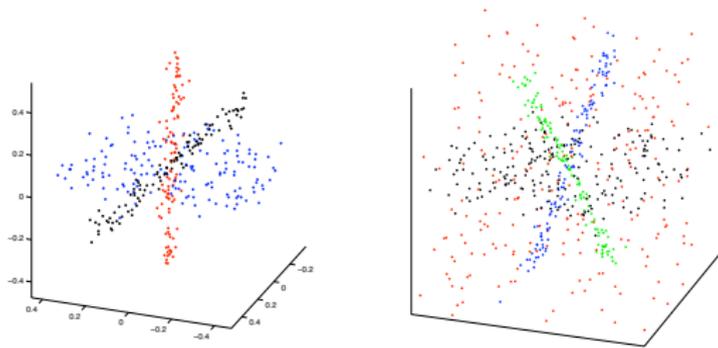
$$\min L^c(\mathbf{X}) \doteq L(\mathbf{X}_1) + \dots + L(\mathbf{X}_k) + H(|\mathbf{X}_1|, \dots, |\mathbf{X}_k|).$$

Optimize with a *bottom-up pair-wise* merging algorithm [Ma, TPAMI'07]:

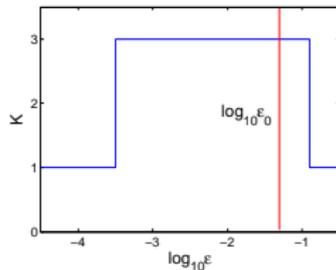
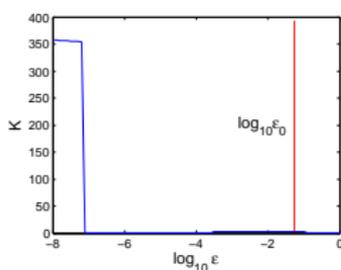
- 1: **input:** the data $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m] \in \mathbb{R}^{D \times m}$ and a distortion $\epsilon^2 > 0$.
- 2: initialize \mathcal{S} as a set of sets with a single datum $\{S = \{\mathbf{x}\} \mid \mathbf{x} \in \mathbf{X}\}$.
- 3: **while** $|\mathcal{S}| > 1$ **do**
- 4: choose distinct sets $S_1, S_2 \in \mathcal{S}$ such that
 $L^c(S_1 \cup S_2) - L^c(S_1, S_2)$ is minimal.
- 5: **if** $L^c(S_1 \cup S_2) - L^c(S_1, S_2) \geq 0$ **then** break;
- 6: **else** $\mathcal{S} := (\mathcal{S} \setminus \{S_1, S_2\}) \cup \{S_1 \cup S_2\}$.
- 7: **end**
- 8: **output:** \mathcal{S}

Surprisingly Good Performance

Empirically, **find global optimum and extremely robust to outliers**



A strikingly sharp **phase transition** w.r.t. quantization ϵ

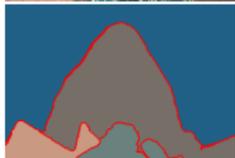


Natural Image Segmentation [Mobahi et.al., IJCV'09]

Compression alone, without any supervision, leads to **state of the art** segmentation on natural images (and many other types of data).



(a) Animals



(b) Buildings



(c) Landscape



(d) People



(e) Water

2. Classify Mixed Data (Extrapolation)

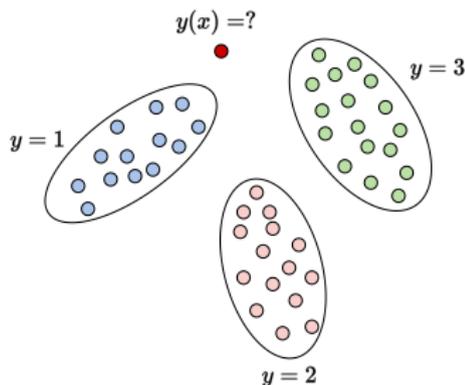
Assume data $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m]$ are i.i.d. samples from a mixture of distributions: $p(\mathbf{x}, \theta) = \sum_{j=1}^k \pi_j p_j(\mathbf{x}, \theta)$.

Classic approach to classify the data is via maximum a posteriori (MAP) classifier:

$$\hat{y}(\mathbf{x}) = \arg \max_j \log p_j(\mathbf{x}, \theta) + \log \pi_j.$$

Difficulties: distributions p_j are hard to estimate and log likelihood is not well-defined when distributions are degenerate.

(probably why SVMs or deep networks prevail instead...)



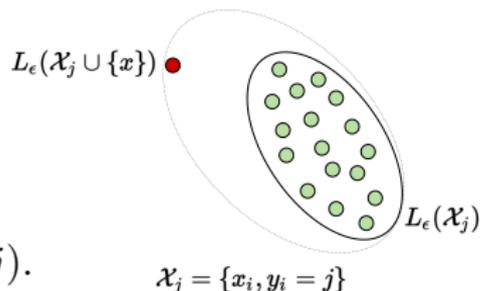
Classify to Compress

[Wright, Tao, Lin, Shum, and Ma, NIPS'07]

A Fundamental Idea:

Count additional #bits needed to encode a query sample \mathbf{x} with data in each class \mathbf{X}_j :

$$\delta L_\epsilon(\mathbf{x}, j) \doteq L_\epsilon(\mathbf{X}_j \cup \{\mathbf{x}\}) - L_\epsilon(\mathbf{X}_j) + L(j).$$



Classification Criterion: Minimum Incremental Coding Length (MICL):

$$\hat{y}(\mathbf{x}) = \arg \min_j \delta L_\epsilon(\mathbf{x}, j).$$

Law of Parsimony: *“Entities should not be multiplied without necessity.”*
–William of Ockham

Asymptotic Property of MICL

Theorem (Wright, NIPS'07)

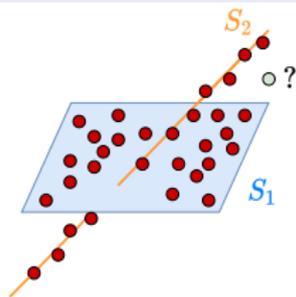
As the number of samples m goes to infinity, the MICL criterion converges at a rate of $O(m^{-1/2})$ to the following criterion:

$$\hat{y}_\epsilon(\mathbf{x}) = \arg \max_j \underbrace{L_G\left(\mathbf{x} \mid \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j + \frac{\epsilon^2}{D} \mathbf{I}\right) + \log \pi_j}_{\text{Regularized MAP}} + \frac{1}{2} D_\epsilon(\boldsymbol{\Sigma}_j),$$

where $D_\epsilon(\boldsymbol{\Sigma}_j) \doteq \text{tr}\left(\boldsymbol{\Sigma}_j(\boldsymbol{\Sigma}_j + \frac{\epsilon^2}{D} \mathbf{I})^{-1}\right)$ is the effective dimension.

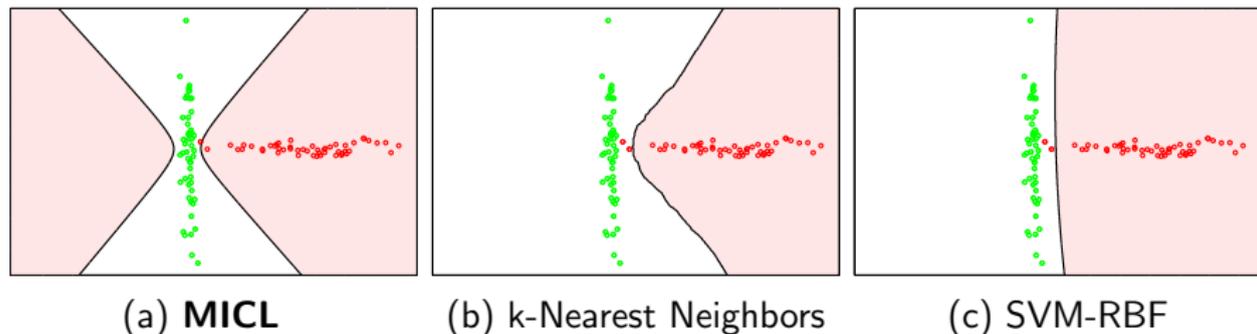
Everything else equal, MICL prefers a class with higher effective dimension.

Err on the side of caution!



Extrapolation of Low-Dim Structure for Classification

Figure: A truly extrapolating (nearest subspace) classifier!



Difficulty in practice: inference computationally costly (non-parametric) and possibly need a kernel (nonlinearity).

Go beyond (non-parametric) data interpolation and extrapolation?

Represent Multi-class Multi-dimensional Data

Given samples

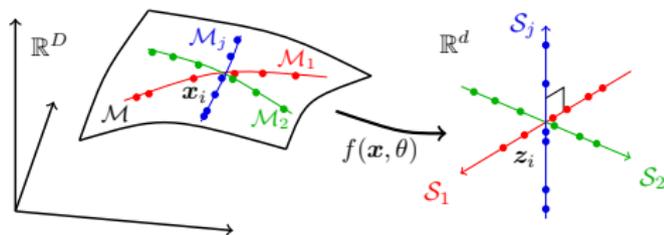
$$\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_m] \subset \cup_{j=1}^k \mathcal{M}_j,$$

seek a good representation

$$\mathbf{Z} = [z_1, \dots, z_m] \subset \mathbb{R}^d$$

through a continuous mapping:

$$f(\mathbf{x}, \theta) : \mathbf{x} \in \mathbb{R}^D \mapsto z \in \mathbb{R}^d.$$



Goals of “re-present” the data:

- **compression**: from high-dimensional samples to compact features.
- **linearization**: from nonlinear structures $\cup_{j=1}^k \mathcal{M}_j$ to linear $\cup_{j=1}^k \mathcal{S}_j$.
- **sparsity**: from separable components \mathcal{M}_j 's to incoherent \mathcal{S}_j 's.
- **self-consistent**: from compact structured \mathbf{Z} back to the data \mathbf{X} .

Seeking a Linear Discriminative Representation (LDR)

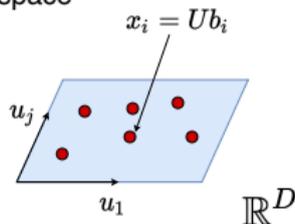
Desiderata: Representation $z = f(x, \theta)$ have the following properties:

- 1 *Within-Class Compressible:* Features of the same class/cluster should be highly compressed in a **low-dimensional** linear subspace.
- 2 *Between-Class Discriminative:* Features of different classes/clusters should be in highly **incoherent** linear subspaces.
- 3 *Maximally Informative:* Dimension (or variance) of features for each class/cluster should be **the same as that of the data**.

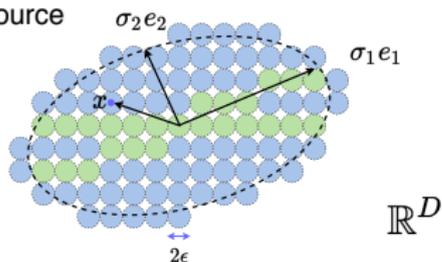
Is there a principled objective for all such properties, together?

Compactness Measure for Linear/Gaussian Representation

Linear subspace



Gaussian source



Theorem (Coding Length, Ma & Derksen TPAMI'07)

The number of bits needed to encode data $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m] \in \mathbb{R}^{D \times m}$ up to a precision $\|\mathbf{x} - \hat{\mathbf{x}}\|_2 \leq \epsilon$ is bounded by:

$$L(\mathbf{X}, \epsilon) \doteq \left(\frac{m + D}{2} \right) \log \det \left(\mathbf{I} + \frac{D}{m\epsilon^2} \mathbf{X} \mathbf{X}^\top \right).$$

This can be derived from constructively quantifying SVD of \mathbf{X} or by sphere packing $\text{vol}(\mathbf{X})$ as samples of a noisy Gaussian source.

Compactness Measure for Linear/Gaussian Representation

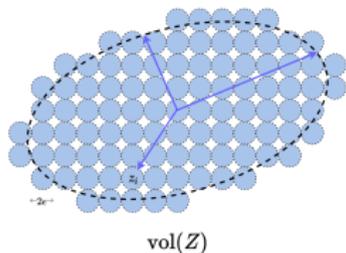
If \mathbf{X} is not (piecewise) linear or Gaussian, consider a **nonlinear** mapping:

$$\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m] \in \mathbb{R}^{D \times m} \xrightarrow{f(\mathbf{x}, \theta)} \mathbf{Z}(\theta) = [\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_m] \in \mathbb{R}^{d \times m}.$$

The average coding length per sample (rate) subject to a distortion ϵ :

$$R(\mathbf{Z}, \epsilon) \doteq \frac{1}{2} \log \det \left(\mathbf{I} + \frac{d}{m\epsilon^2} \mathbf{Z} \mathbf{Z}^\top \right). \quad (6)$$

Rate distortion is an intrinsic measure for the volume of all features.



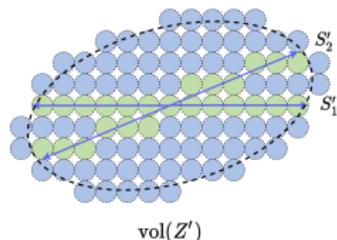
Compactness Measure for Mixed Linear Representations

The features \mathbf{Z} of **multi-class** data

$$\mathbf{X} = \mathbf{X}_1 \cup \mathbf{X}_2 \cup \cdots \cup \mathbf{X}_k \subset \cup_{j=1}^k \mathcal{M}_j.$$

may be partitioned into **multiple** subsets:

$$\mathbf{Z} = \mathbf{Z}_1 \cup \mathbf{Z}_2 \cup \cdots \cup \mathbf{Z}_k \subset \cup_{j=1}^k \mathcal{S}_j.$$



W.r.t. this partition, the **average coding rate** is:

$$R^c(\mathbf{Z}, \epsilon | \mathbf{\Pi}) \doteq \sum_{j=1}^k \frac{\text{tr}(\mathbf{\Pi}_j)}{2m} \log \det \left(\mathbf{I} + \frac{d}{\text{tr}(\mathbf{\Pi}_j)\epsilon^2} \mathbf{Z} \mathbf{\Pi}_j \mathbf{Z}^\top \right), \quad (7)$$

where $\mathbf{\Pi} = \{\mathbf{\Pi}_j \in \mathbb{R}^{m \times m}\}_{j=1}^k$ encode the membership of the m samples in the k classes: the diagonal entry $\mathbf{\Pi}_j(i, i)$ of $\mathbf{\Pi}_j$ is the probability of sample i belonging to subset j . $\Omega \doteq \{\mathbf{\Pi} \mid \sum \mathbf{\Pi}_j = \mathbf{I}, \mathbf{\Pi}_j \geq \mathbf{0}\}$

Measure for Linear Discriminative Representation (LDR)

A fundamental idea: maximize the **difference** between the coding rate of all features and the average rate of features in each of the classes:

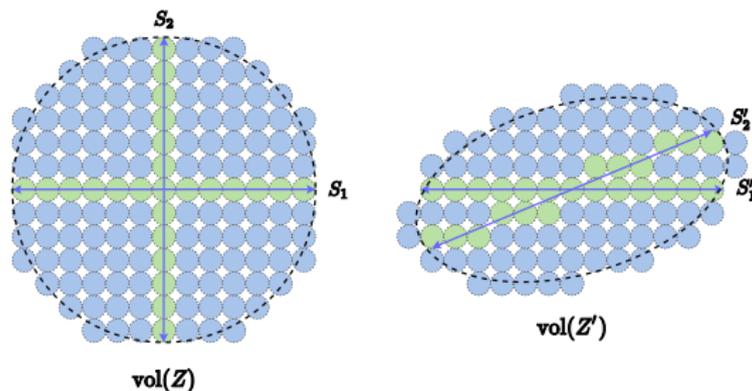
$$\Delta R(\mathbf{Z}, \mathbf{\Pi}, \epsilon) = \underbrace{\frac{1}{2} \log \det \left(\mathbf{I} + \frac{d}{m\epsilon^2} \mathbf{Z}\mathbf{Z}^\top \right)}_R - \underbrace{\sum_{j=1}^k \frac{\text{tr}(\mathbf{\Pi}_j)}{2m} \log \det \left(\mathbf{I} + \frac{d}{\text{tr}(\mathbf{\Pi}_j)\epsilon^2} \mathbf{Z}\mathbf{\Pi}_j\mathbf{Z}^\top \right)}_{R^c}.$$

This difference is called **rate reduction** (measuring **information gain**):

- Large R : **expand** all features \mathbf{Z} as **large** as possible.
- Small R^c : **compress** each class \mathbf{Z}_j as **small** as possible.

Slogan: similarity contracts and dissimilarity contrasts!

Interpretation of MCR²: Sphere Packing and Counting



Example: two subspaces S_1 and S_2 in \mathbb{R}^2 .

- $\log \#(\text{green spheres} + \text{blue spheres}) = \text{rate of span of all samples } R$.
- $\log \#(\text{green spheres}) = \text{rate of the two subspaces } R^c$.
- $\log \#(\text{blue spheres}) = \text{rate reduction } \Delta R$.

Comparison to Contrastive Learning

[Hadsell, Chopra, and LeCun, CVPR'06]

When k is large, a randomly chosen **pair** $(\mathbf{x}_i, \mathbf{x}_j)$ is of high probability belonging to different classes. Minimize the **contrastive loss**:

$$\min -\log \frac{\exp(\langle \mathbf{z}_i, \mathbf{z}'_i \rangle)}{\sum_{j \neq i} \exp(\langle \mathbf{z}_i, \mathbf{z}_j \rangle)}.$$

The learned features of such pairs of samples together with their augmentations \mathbf{Z}_i and \mathbf{Z}_j should have large rate reduction:

$$\max \sum_{ij} \Delta R_{ij} \doteq R(\mathbf{Z}_i \cup \mathbf{Z}_j, \epsilon) - \frac{1}{2}(R(\mathbf{Z}_i, \epsilon) + R(\mathbf{Z}_j, \epsilon)).$$

MCR² contrasts triplets, quadruplets, or any number of sets.

Principle of Maximal Coding Rate Reduction (MCR²)

[Yu, Chan, You, Song, Ma, NeurIPS2020]

Learn a mapping $f(\mathbf{x}, \theta)$ (for a given partition $\mathbf{\Pi}$):

$$\mathbf{X} \xrightarrow{f(\mathbf{x}, \theta)} \mathbf{Z}(\theta) \xrightarrow{\mathbf{\Pi}, \epsilon} \Delta R(\mathbf{Z}(\theta), \mathbf{\Pi}, \epsilon) \quad (8)$$

so as to **Maximize the Coding Rate Reduction (MCR²)**:

$$\begin{aligned} \max_{\theta} \quad & \Delta R(\mathbf{Z}(\theta), \mathbf{\Pi}, \epsilon) = R(\mathbf{Z}(\theta), \epsilon) - R^c(\mathbf{Z}(\theta), \epsilon \mid \mathbf{\Pi}), \\ \text{subject to} \quad & \|\mathbf{Z}_j(\theta)\|_F^2 = m_j, \mathbf{\Pi} \in \Omega. \end{aligned} \quad (9)$$

Since ΔR is *monotonic* in the scale of \mathbf{Z} , one needs to:

normalize the features $z = f(\mathbf{x}, \theta)$ **so as to compare** $\mathbf{Z}(\theta)$ **and** $\mathbf{Z}(\theta')$!

[Batch normalization, Sergey Ioffe and Christian Szegedy, 2015.](#)

[Layer normalization'16, instance normalization'16; group normalization'18...](#)

Theoretical Justification of the MCR² Principle

Theorem (Informal Statement [Yu et.al., NeurIPS2020])

Suppose $\mathbf{Z}^* = \mathbf{Z}_1^* \cup \dots \cup \mathbf{Z}_k^*$ is the optimal solution that maximizes the rate reduction (9). We have:

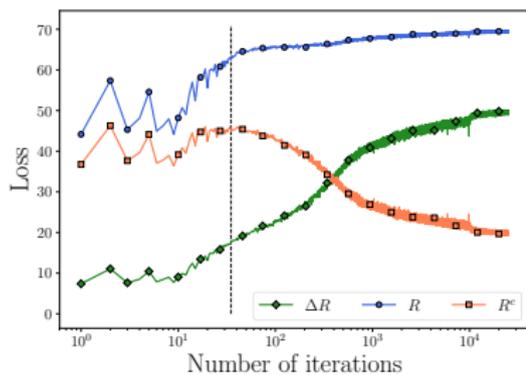
- *Between-class Discriminative: As long as the ambient space is adequately large ($d \geq \sum_{j=1}^k d_j$), the subspaces are all orthogonal to each other, i.e. $(\mathbf{Z}_i^*)^\top \mathbf{Z}_j^* = \mathbf{0}$ for $i \neq j$.*
- *Maximally Informative Representation: As long as the coding precision is adequately high, i.e., $\epsilon^4 < \min_j \left\{ \frac{m_j}{m} \frac{d^2}{d_j^2} \right\}$, each subspace achieves its maximal dimension, i.e. $\text{rank}(\mathbf{Z}_j^*) = d_j$. In addition, the largest $d_j - 1$ singular values of \mathbf{Z}_j^* are equal.*

A new slogan, beyond Aristotle:

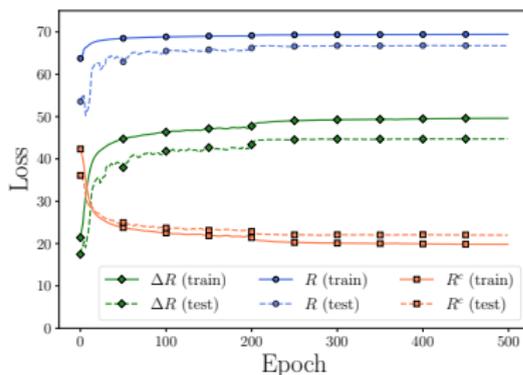
The whole is to be maximally greater than the sum of the parts!

Experiment I: Supervised Deep Learning

Experimental Setup: Train $f(x, \theta)$ as ResNet18 on the CIFAR10 dataset, feature z dimension $d = 128$, precision $\epsilon^2 = 0.5$.



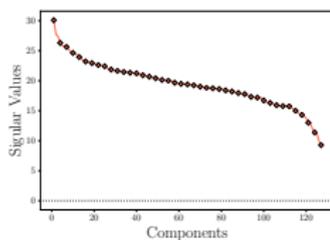
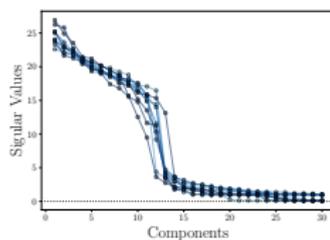
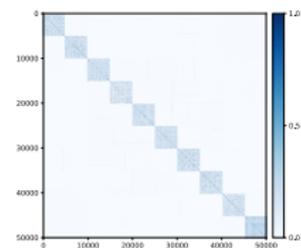
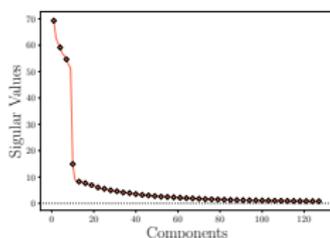
(a)



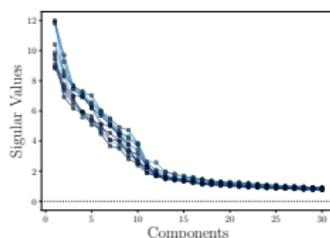
(b)

Figure: (a). Evolution of $R, R^c, \Delta R$ during the training process; (b). Training loss versus testing loss.

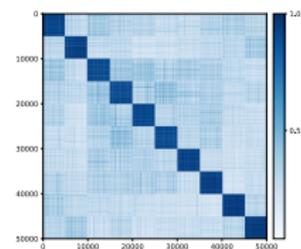
Visualization of Learned Representations \mathbb{Z}

(a) MCR^2 (overall)(b) MCR^2 (PCA of every class)(c) MCR^2 (cosine similarity)

(d) CE (overall)



(e) CE (PCA of every class)

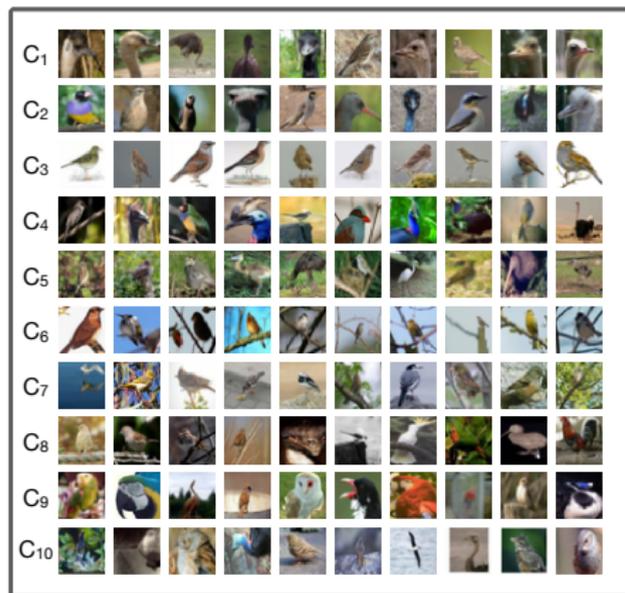


(f) CE (cosine similarity)

Figure: PCA of learned representations from MCR^2 and cross-entropy.

No neural collapse!

Visualization - Samples along Principal Components



(a) Bird



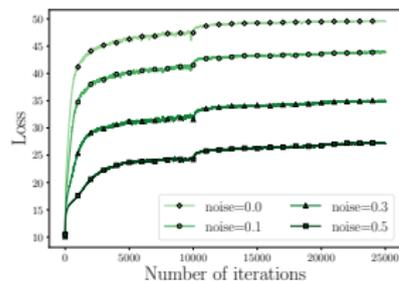
(b) Ship

Figure: Top-10 “principal” images for class - “Bird” and “Ship” in the CIFAR10.

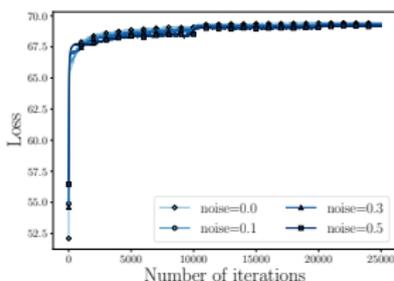
Experiment II: Robustness to Label Noise

	RATIO=0.0	RATIO=0.1	RATIO=0.2	RATIO=0.3	RATIO=0.4	RATIO=0.5
CE TRAINING	0.939	0.909	0.861	0.791	0.724	0.603
MCR ² TRAINING	0.940	0.911	0.897	0.881	0.866	0.843

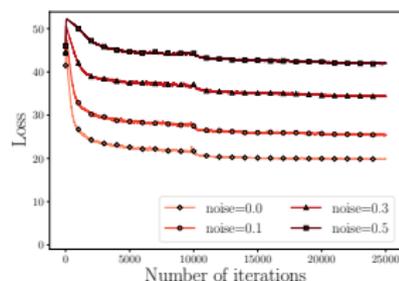
Table 1: Classification results with features learned with labels corrupted at different levels.



(a) $\Delta R(\mathbf{Z}(\theta), \Pi, \epsilon)$



(b) $R(\mathbf{Z}(\theta), \epsilon)$



(c) $R^c(\mathbf{Z}(\theta), \epsilon | \Pi)$

Figure: Evolution of $R, R^c, \Delta R$ of MCR² during training with corrupted labels.

Represent only what can be jointly compressed.

Deep Networks from Optimizing Rate Reduction

$$\mathbf{X} \xrightarrow{f(\mathbf{x}, \theta)} \mathbf{Z}(\theta); \quad \max_{\theta} \Delta R(\mathbf{Z}(\theta), \mathbf{\Pi}, \epsilon).$$

Final features learned by MCR² are more interpretable and robust, **but**:

- The borrowed deep network (e.g. ResNet) is still a “black box”!
- Why is a “deep” architecture necessary, and how wide and deep?
- What are the roles of the “linear and nonlinear” operators?
- Why “multi-channel” convolutions?
- ...

Replace black box networks with entirely “white box” networks?

Projected Gradient Ascent for Rate Reduction

Recall the rate reduction objective:

$$\max_{\mathbf{Z}} \Delta R(\mathbf{Z}) \doteq \underbrace{\frac{1}{2} \log \det \left(\mathbf{I} + \alpha \mathbf{Z} \mathbf{Z}^* \right)}_{R(\mathbf{Z})} - \underbrace{\sum_{j=1}^k \frac{\gamma_j}{2} \log \det \left(\mathbf{I} + \alpha_j \mathbf{Z} \mathbf{\Pi}^j \mathbf{Z}^* \right)}_{R_c(\mathbf{Z}, \mathbf{\Pi})}, \quad (10)$$

where $\alpha = d/(m\epsilon^2)$, $\alpha_j = d/(\text{tr}(\mathbf{\Pi}^j)\epsilon^2)$, $\gamma_j = \text{tr}(\mathbf{\Pi}^j)/m$ for $j = 1, \dots, k$.

Consider directly maximizing ΔR with **projected gradient ascent** (PGA):

$$\mathbf{Z}_{\ell+1} \propto \mathbf{Z}_{\ell} + \eta \cdot \left. \frac{\partial \Delta R}{\partial \mathbf{Z}} \right|_{\mathbf{Z}_{\ell}} \quad \text{subject to} \quad \mathbf{Z}_{\ell+1} \subset \mathbb{S}^{d-1}. \quad (11)$$

Gradients of the Two Terms

The derivatives $\frac{\partial R(\mathbf{Z})}{\partial \mathbf{Z}}$ and $\frac{\partial R_c(\mathbf{Z}, \mathbf{\Pi})}{\partial \mathbf{Z}}$ are:

$$\frac{1}{2} \frac{\partial \log \det(\mathbf{I} + \alpha \mathbf{Z} \mathbf{Z}^*)}{\partial \mathbf{Z}} \Big|_{\mathbf{Z}_\ell} = \underbrace{\alpha (\mathbf{I} + \alpha \mathbf{Z}_\ell \mathbf{Z}_\ell^*)^{-1}}_{\mathbf{E}_\ell \in \mathbb{R}^{d \times d}} \mathbf{Z}_\ell, \quad (12)$$

$$\frac{1}{2} \frac{\partial (\gamma_j \log \det(\mathbf{I} + \alpha_j \mathbf{Z} \mathbf{\Pi}^j \mathbf{Z}^*))}{\partial \mathbf{Z}} \Big|_{\mathbf{Z}_\ell} = \gamma_j \underbrace{\alpha_j (\mathbf{I} + \alpha_j \mathbf{Z}_\ell \mathbf{\Pi}^j \mathbf{Z}_\ell^*)^{-1}}_{\mathbf{C}_\ell^j \in \mathbb{R}^{d \times d}} \mathbf{Z}_\ell \mathbf{\Pi}^j. \quad (13)$$

Hence the gradient $\frac{\partial \Delta R(\mathbf{Z})}{\partial \mathbf{Z}}$ is:

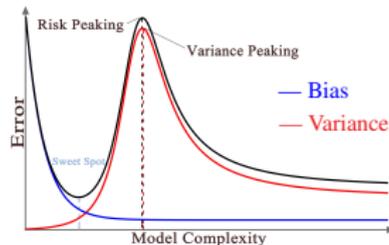
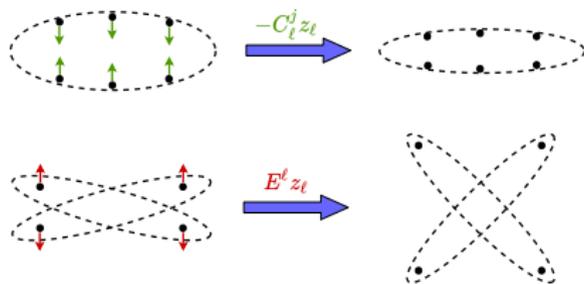
$$\frac{\partial \Delta R}{\partial \mathbf{Z}} \Big|_{\mathbf{Z}_\ell} = \underbrace{\mathbf{E}_\ell}_{\text{Expansion}} \mathbf{Z}_\ell - \sum_{j=1}^k \gamma_j \underbrace{\mathbf{C}_\ell^j}_{\text{Compression}} \mathbf{Z}_\ell \mathbf{\Pi}^j \in \mathbb{R}^{d \times m}. \quad (14)$$

Interpretation of the Linear Operators E and C^j

For any $z_\ell \in \mathbb{R}^d$, we have

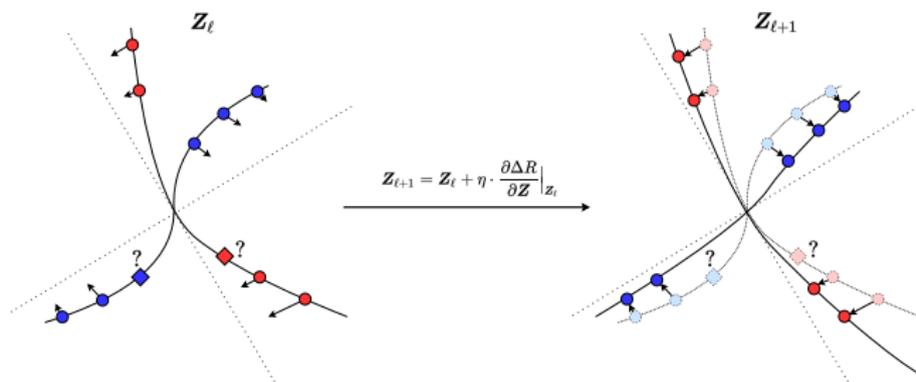
$$E_\ell z_\ell = \alpha(z_\ell - \mathbf{Z}_\ell \mathbf{q}_\ell^*) \quad \text{with} \quad \mathbf{q}_\ell^* \doteq \arg \min_{\mathbf{q}_\ell} \alpha \|z_\ell - \mathbf{Z}_\ell \mathbf{q}_\ell\|_2^2 + \|\mathbf{q}_\ell\|_2^2.$$

$E_\ell z_\ell$ and $C_\ell^j z_\ell$ are the “residuals” of z_ℓ against the subspaces spanned by columns of \mathbf{Z}_ℓ and \mathbf{Z}_ℓ^j , respectively.



Such “auto” ridge regressions **do not overfit** even with redundant random regressors, due to a “double descent” risk [Yang, ICML'20]!

Incremental Deformation via Gradient Flow



Extrapolate the gradient $\frac{\partial \Delta R(\mathbf{Z})}{\partial \mathbf{Z}}$ from training samples \mathbf{Z} to all $z \in \mathbb{R}^d$:

$$\frac{\partial \Delta R}{\partial \mathbf{Z}} \Big|_{z_\ell} = \mathbf{E}_\ell \mathbf{Z}_\ell - \sum_{j=1}^k \gamma_j \mathbf{C}_\ell^j \mathbf{Z}_\ell \underbrace{\mathbf{\Pi}^j}_{\text{known}} \in \mathbb{R}^{d \times m}, \quad (15)$$

$$g(z_\ell, \theta_\ell) \doteq \mathbf{E}_\ell z_\ell - \sum_{j=1}^k \gamma_j \mathbf{C}_\ell^j z_\ell \underbrace{\pi^j(z_\ell)}_{\text{unknown}} \in \mathbb{R}^d. \quad (16)$$

Estimate of the Membership $\pi^j(\mathbf{z}_\ell)$

Estimate the membership $\pi^j(\mathbf{z}_\ell)$ with “softmax” on the residuals $\|\mathbf{C}_\ell^j \mathbf{z}_\ell\|$:

$$\pi^j(\mathbf{z}_\ell) \approx \hat{\pi}^j(\mathbf{z}_\ell) \doteq \frac{\exp(-\lambda \|\mathbf{C}_\ell^j \mathbf{z}_\ell\|)}{\sum_{j=1}^k \exp(-\lambda \|\mathbf{C}_\ell^j \mathbf{z}_\ell\|)} \in [0, 1]. \quad (17)$$

Thus the weighted residuals for contracting:

$$\sigma\left([\mathbf{C}_\ell^1 \mathbf{z}_\ell, \dots, \mathbf{C}_\ell^k \mathbf{z}_\ell]\right) \doteq \sum_{j=1}^k \gamma_j \mathbf{C}_\ell^j \mathbf{z}_\ell \cdot \hat{\pi}^j(\mathbf{z}_\ell) \in \mathbb{R}^d. \quad (18)$$

Many alternatives, e.g. enforcing all features to be in the first quadrant:

$$\sigma(\mathbf{z}_\ell) \approx \mathbf{z}_\ell - \sum_{j=1}^k \text{ReLU}(\mathbf{P}_\ell^j \mathbf{z}_\ell), \quad (19)$$

The ReduNet for Optimizing Rate Reduction

Iterative projected gradient ascent (PGA) :

$$\mathbf{z}_{\ell+1} \propto \mathbf{z}_{\ell} + \eta \cdot \underbrace{\left[\mathbf{E}_{\ell} \mathbf{z}_{\ell} + \sigma \left([C_{\ell}^1 \mathbf{z}_{\ell}, \dots, C_{\ell}^k \mathbf{z}_{\ell}] \right) \right]}_{g(\mathbf{z}_{\ell}, \boldsymbol{\theta}_{\ell})} \quad \text{s.t.} \quad \mathbf{z}_{\ell+1} \in \mathbb{S}^{d-1}, \quad (20)$$

$$f(\mathbf{x}, \boldsymbol{\theta}) = \phi^L \circ \phi^{L-1} \circ \dots \circ \phi^0(\mathbf{x}), \quad \text{with} \quad \phi^{\ell}(\mathbf{z}_{\ell}, \boldsymbol{\theta}_{\ell}) \doteq \mathcal{P}_{\mathbb{S}^{d-1}}[\mathbf{z}_{\ell} + \eta \cdot g(\mathbf{z}_{\ell}, \boldsymbol{\theta}_{\ell})].$$

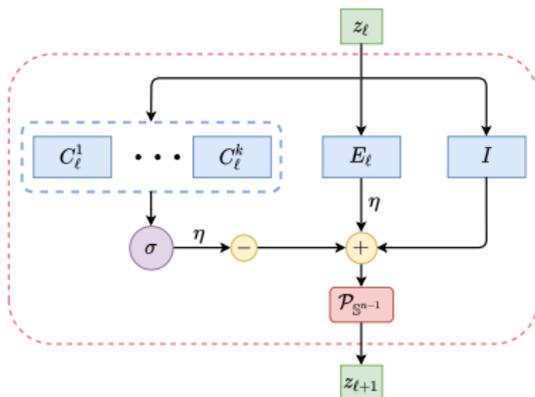


Figure: One layer of the **ReduNet**: one PGA iteration.

The ReduNet versus ResNet or ResNeXt

Iterative projected gradient ascent (PGA):

$$\mathbf{z}_{\ell+1} \propto \mathbf{z}_{\ell} + \eta \cdot \underbrace{\left[\mathbf{E}_{\ell} \mathbf{z}_{\ell} + \sigma \left([C_{\ell}^1 \mathbf{z}_{\ell}, \dots, C_{\ell}^k \mathbf{z}_{\ell}] \right) \right]}_{g(\mathbf{z}_{\ell}, \theta_{\ell})} \quad \text{s.t.} \quad \mathbf{z}_{\ell+1} \in \mathbb{S}^{d-1}, \quad (21)$$

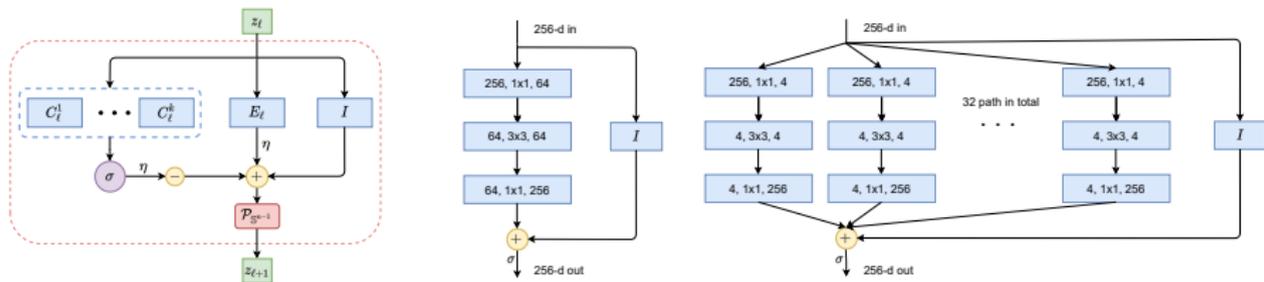


Figure: Left: **ReduNet**. Middle and Right: **ResNet** [He et. al. 2016] and **ResNeXt** [Xie et. al. 2017] (hundreds of layers).

Forward construction instead of back propagation!¹

¹ *The Forward-Forward Algorithm*, G. Hinton, 2022.

The ReduNet versus Mixture of Experts

Approximate iterative projected gradient ascent (PGA) :

$$z_{\ell+1} \propto z_{\ell} + \eta \cdot \underbrace{\left[E_{\ell} z_{\ell} + \sigma([C_{\ell}^1 z_{\ell}, \dots, C_{\ell}^k z_{\ell}]) \right]}_{g(z_{\ell}, \theta_{\ell})} \quad \text{s.t.} \quad z_{\ell+1} \in \mathbb{S}^{d-1}, \quad (22)$$

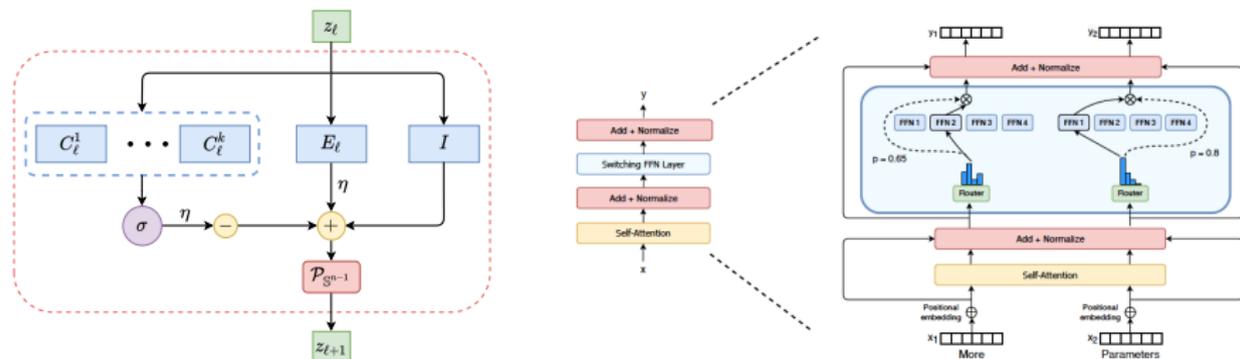


Figure: Left: ReduNet layer. Right: Mixture of Experts [Shazeer et. al. 2017] or Switched Transformer [Fedus et. al. 2021] (1.7 trillion parameters).

Forward construction instead of back propagation!²

² The Forward-Forward Algorithm, G. Hinton, 2022.

ReduNet Features for Mixture of Gaussians

$L = 2000$ -Layers ReduNet: $m = 500, \eta = 0.5, \epsilon = 0.1$.

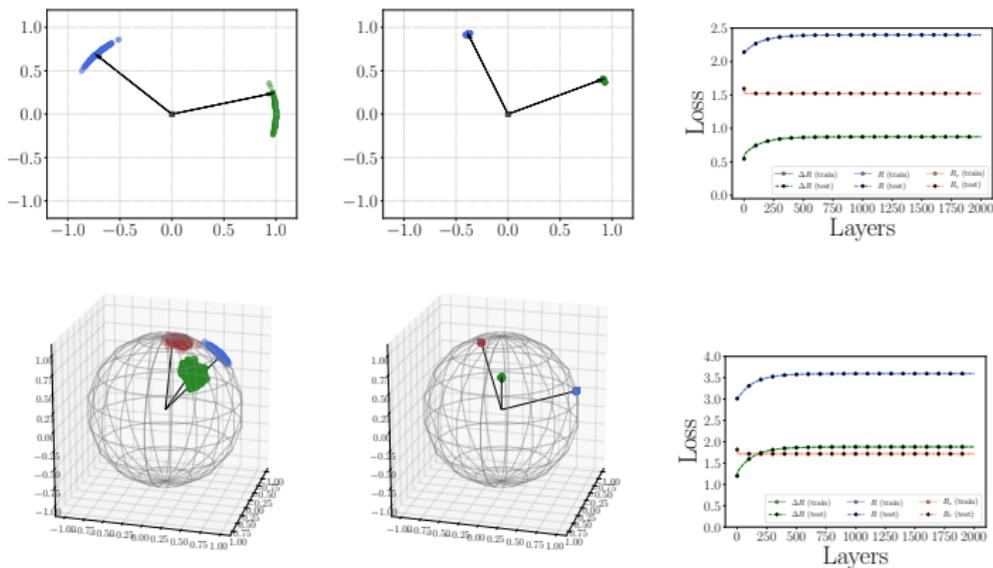


Figure: Left: original samples \mathbf{X} and ReduNet features $\mathbf{Z} = f(\mathbf{Z}, \theta)$ for 2D and 3D Mixture of Gaussians. Right: plots for the progression of values of the rates.

Group Invariant Classification

Feature mapping $f(\mathbf{x}, \boldsymbol{\theta})$ is invariant to a group of transformations:

$$\text{Group Invariance: } f(\mathbf{x} \circ \mathbf{g}, \boldsymbol{\theta}) \sim f(\mathbf{x}, \boldsymbol{\theta}), \quad \forall \mathbf{g} \in \mathbb{G}, \quad (23)$$

where “ \sim ” indicates two features belonging to the same equivalent class.

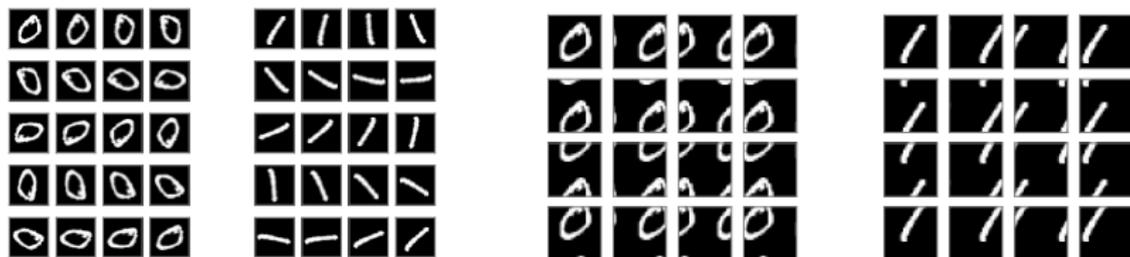


Figure: Left: 1D rotation \mathbb{S}^1 ; Right: 2D cyclic translation \mathcal{T}^2 .

1. Fooling CNNs with simple transformations, Engstrom et.al., 2017.
2. Why do deep convolutional networks generalize so poorly to small image transformations? Azuly & Weiss, 2018.

Group Invariant Classification

Feature mapping $f(x, \theta)$ is invariant to a group of transformations:

$$\text{Group Invariance: } f(x \circ g, \theta) \sim f(x, \theta), \quad \forall g \in \mathbb{G}, \quad (24)$$

where “ \sim ” indicates two features belonging to the same equivalent class.

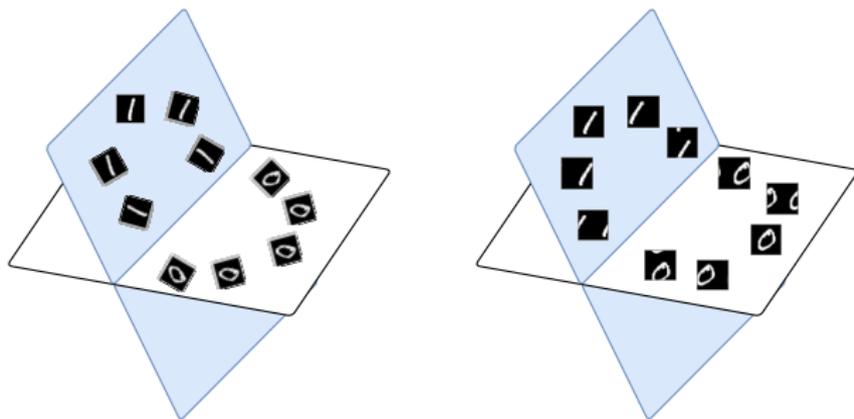


Figure: Embed all equivariant samples to the same subspace.

Circulant Matrix and Convolution

Given a vector $\mathbf{z} = [z_0, z_1, \dots, z_{n-1}]^* \in \mathbb{R}^n$, we may arrange all its circular shifted versions in a circulant matrix form as

$$\text{circ}(\mathbf{z}) \doteq \begin{bmatrix} z_0 & z_{n-1} & \dots & z_2 & z_1 \\ z_1 & z_0 & z_{n-1} & \dots & z_2 \\ \vdots & z_1 & z_0 & \ddots & \vdots \\ z_{n-2} & \vdots & \ddots & \ddots & z_{n-1} \\ z_{n-1} & z_{n-2} & \dots & z_1 & z_0 \end{bmatrix} \in \mathbb{R}^{n \times n}. \quad (25)$$

A circular (or cyclic) convolution:

$$\text{circ}(\mathbf{z}) \cdot \mathbf{x} = \mathbf{z} \circledast \mathbf{x}, \quad \text{where} \quad (\mathbf{z} \circledast \mathbf{x})_i = \sum_{j=0}^{n-1} x_j z_{i+n-j \bmod n}. \quad (26)$$

Convolutions from Cyclic Shift Invariance

Given a set of sample vectors $\mathbf{Z} = [\mathbf{z}^1, \dots, \mathbf{z}^m]$, construct the ReduNet from cyclic-shift augmented families $\mathbf{Z} = [\text{circ}(\mathbf{z}^1), \dots, \text{circ}(\mathbf{z}^m)]$.

Proposition (Convolution Structures of \mathbf{E} and \mathbf{C}^j)

The linear operator in the ReduNet:

$$\mathbf{E} = \alpha \left(\mathbf{I} + \alpha \sum_{i=1}^m \text{circ}(\mathbf{z}^i) \text{circ}(\mathbf{z}^i)^* \right)^{-1}$$

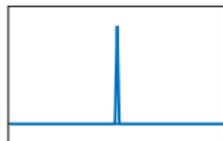
is a circulant matrix and represents a circular convolution:

$$\mathbf{E}\mathbf{z} = \mathbf{e} \circledast \mathbf{z},$$

where \mathbf{e} is the first column vector of \mathbf{E} . Similarly, the operators \mathbf{C}^j associated with subsets \mathbf{Z}^j are also circular convolutions.

Tradeoff between Invariance and Separability

A problem with separability: superposition of shifted “delta” functions can generate any other signals:
 $\text{span}[\text{circ}(\boldsymbol{x})] = \mathbb{R}^n!$



A necessary assumption: \boldsymbol{x} is **sparsely generated** from incoherent dictionaries for different classes:

$$\boldsymbol{x} = [\text{circ}(\mathcal{D}_1), \text{circ}(\mathcal{D}_2), \dots, \text{circ}(\mathcal{D}_k)] \bar{\boldsymbol{z}}.$$

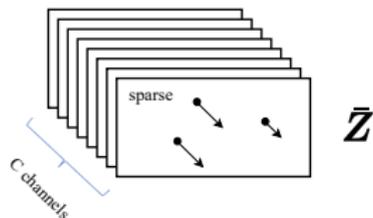


Fish

= Σ 

C convolution kernels

⊗

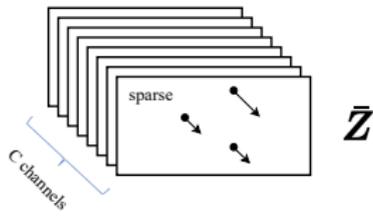


Duck

= Σ 

C convolution kernels

⊗

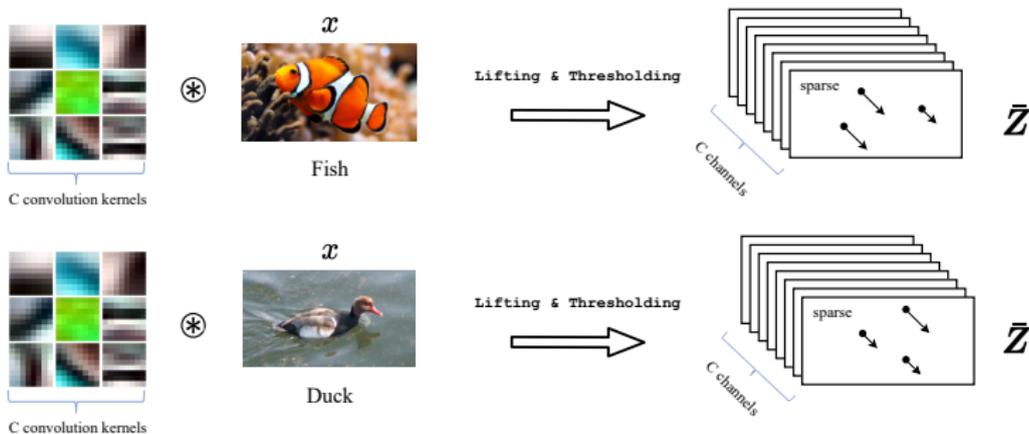


Tradeoff between Invariance and Separability

A basic idea: estimate sparse codes \bar{z} by taking their responses to multiple analysis filters $\mathbf{k}_1, \dots, \mathbf{k}_C \in \mathbb{R}^n$ [Rubinstein & Elad 2014]:

$$\bar{z} = \tau \left[\mathbf{k}_1 \circledast \mathbf{x}, \dots, \mathbf{k}_C \circledast \mathbf{x} \right]^* \in \mathbb{R}^{C \times n}. \quad (27)$$

for some entry-wise “sparsity-promoting” operator $\tau(\cdot)$.



Multi-Channel Convolutions

Given a set of multi-channel sparse codes $\bar{\mathbf{Z}} = [\bar{\mathbf{z}}^1, \dots, \bar{\mathbf{z}}^m]$, construct the ReduNet from their circulant families $\bar{\mathbf{Z}} = [\text{circ}(\bar{\mathbf{z}}^1), \dots, \text{circ}(\bar{\mathbf{z}}^m)]$.

Proposition (Convolution Structures of $\bar{\mathbf{E}}$ and $\bar{\mathbf{C}}^j$)

The linear operator in the ReduNet:

$$\bar{\mathbf{E}} = \alpha \left(\mathbf{I} + \alpha \sum_{i=1}^m \text{circ}(\bar{\mathbf{z}}^i) \text{circ}(\bar{\mathbf{z}}^i)^* \right)^{-1} \in \mathbb{R}^{C_n \times C_n}$$

is a block circulant matrix and represents a multi-channel convolution:

$$\bar{\mathbf{E}}(\bar{\mathbf{z}}) = \bar{\mathbf{e}} \circledast \bar{\mathbf{z}} \in \mathbb{R}^{C_n},$$

where $\bar{\mathbf{e}}$ is the first slice of $\bar{\mathbf{E}}$. Similarly, the operators $\bar{\mathbf{C}}^j$ associated with subsets $\bar{\mathbf{Z}}^j$ are also multi-channel circular convolutions.

Multi-Channel Convolutions

$$\bar{\mathbf{E}}(\bar{\mathbf{z}}) = \bar{\mathbf{e}} \circledast \bar{\mathbf{z}} \in \mathbb{R}^{Cn}, \quad \bar{\mathbf{C}}^j(\bar{\mathbf{z}}) = \bar{\mathbf{c}}^j \circledast \bar{\mathbf{z}} \in \mathbb{R}^{Cn} :$$

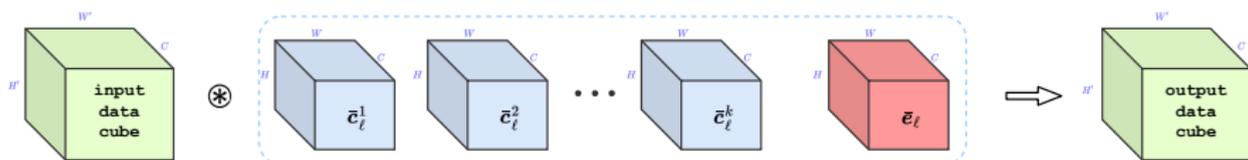


Figure: $\bar{\mathbf{E}}$ and $\bar{\mathbf{C}}^j$ are automatically multi-channel convolutions!

The Convolution ReduNet versus Scattering Network

Iterative projected gradient ascent (PGA) for invariant rate reduction:

$$\bar{z}_{\ell+1} \propto \bar{z}_{\ell} + \eta \cdot \underbrace{\left[\bar{E}_{\ell} \bar{z}_{\ell} + \sigma([\bar{C}_{\ell}^1 \bar{z}_{\ell}, \dots, \bar{C}_{\ell}^k \bar{z}_{\ell}]) \right]}_{g(\bar{z}_{\ell}, \theta_{\ell})}, \quad (28)$$

with each layer being a fixed number of multi-channel convolutions!

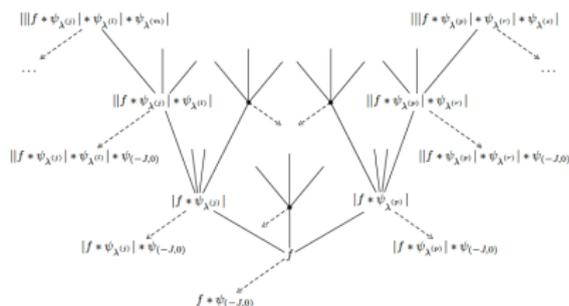
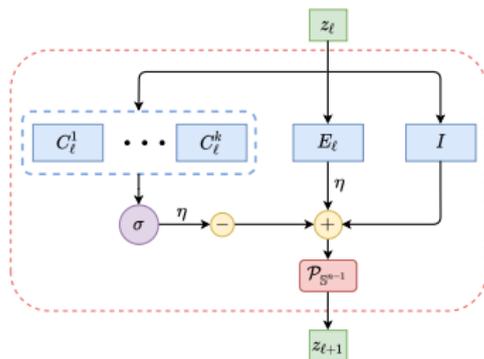


Fig. 2: Scattering network architecture based on wavelet filters and the modulus non-linearity. The elements of the feature vector $\Phi_W(f)$ in (1) are indicated at the tips of the arrows.

Figure: Left: **ReduNet** layer. Right: **Scattering Network** [J. Bruna and S. Mallat, 2013] [T. Wiatowski and H. Blcskei, 2018] (only 2-3 layers).

Fast Computation in Spectral Domain

Fact: all circulant matrices can be simultaneously diagonalized by the *discrete Fourier transform* \mathbf{F} : $\text{circ}(\mathbf{z}) = \mathbf{F}^* \mathbf{D} \mathbf{F}$.

$$\left(\mathbf{I} + \sum_{i=1}^m \text{circ}(\bar{\mathbf{z}}^i) \text{circ}(\bar{\mathbf{z}}^i)^* \right)^{-1} = \left(\mathbf{I} + \begin{bmatrix} \mathbf{F}^* & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \ddots & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{F}^* \end{bmatrix} \begin{bmatrix} \mathbf{D}_{11} & \cdots & \mathbf{D}_{1C} \\ \vdots & \ddots & \vdots \\ \mathbf{D}_{C1} & \cdots & \mathbf{D}_{CC} \end{bmatrix} \begin{bmatrix} \mathbf{F} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \ddots & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{F} \end{bmatrix} \right)^{-1} \in \mathbb{R}^{nC \times nC}$$

where \mathbf{D}_{ij} are all diagonal of size n .

Computing the inverse is $O(C^3 n)$ in the spectral domain, instead of $O(C^3 n^3)$! *Learning convolutional networks for invariant classification is naturally far more efficient in the spectral domain!*

Nature: In visual cortex, neurons encode and transmit information in frequency, hence called “spiking neurons” [Softky & Koch, 1993; Eliasmith & Anderson, 2003].

A “White Box” Deep Convolutional ReduNet by Construction (Spectral Domain)

Require: $\bar{Z} \in \mathbb{R}^{C \times T \times m}$, Π , $\epsilon > 0$, λ , and a learning rate η .

1: Set $\alpha = \frac{C}{m\epsilon^2}$, $\{\alpha_j = \frac{C}{\text{tr}(\Pi^j)\epsilon^2}\}_{j=1}^k$, $\{\gamma_j = \frac{\text{tr}(\Pi^j)}{m}\}_{j=1}^k$.

2: Set $\bar{V}_0 = \{\bar{v}_0^i(p) \in \mathbb{C}^C\}_{p=0, i=1}^{T-1, m} \doteq \text{DFT}(\bar{Z}) \in \mathbb{C}^{C \times T \times m}$.

3: for $\ell = 1, 2, \dots, L$ do

4: for $p = 0, 1, \dots, T - 1$ do

5: Compute $\bar{E}_\ell(p) \in \mathbb{C}^{C \times C}$ and $\{\bar{C}_\ell^j(p) \in \mathbb{C}^{C \times C}\}_{j=1}^k$ as

$$\bar{E}_\ell(p) \doteq \alpha \cdot [\mathbf{I} + \alpha \cdot \bar{V}_{\ell-1}(p) \cdot \bar{V}_{\ell-1}(p)^*]^{-1},$$

$$\bar{C}_\ell^j(p) \doteq \alpha_j \cdot [\mathbf{I} + \alpha_j \cdot \bar{V}_{\ell-1}(p) \cdot \Pi^j \cdot \bar{V}_{\ell-1}(p)^*]^{-1};$$

6: end for

7: for $i = 1, \dots, m$ do

8: for $p = 0, 1, \dots, T - 1$ do

9: Compute $\{\bar{p}_\ell^{ij}(p) \doteq \bar{C}_\ell^j(p) \cdot \bar{v}_\ell^i(p) \in \mathbb{C}^{C \times 1}\}_{j=1}^k$;

10: end for

11: Let $\{\bar{P}_\ell^{ij} = [\bar{p}_\ell^{ij}(0), \dots, \bar{p}_\ell^{ij}(T-1)] \in \mathbb{C}^{C \times T}\}_{j=1}^k$;

12: Compute $\{\hat{\pi}_\ell^{ij} = \frac{\exp(-\lambda \|\bar{P}_\ell^{ij}\|_F)}{\sum_{j=1}^k \exp(-\lambda \|\bar{P}_\ell^{ij}\|_F)}\}_{j=1}^k$;

13: for $p = 0, 1, \dots, T - 1$ do

14: $\bar{v}_\ell^i(p) = \bar{v}_{\ell-1}^i(p) + \eta \left(\bar{E}_\ell(p) \bar{v}_\ell^i(p) - \sum_{j=1}^k \gamma_j \cdot \hat{\pi}_\ell^{ij} \cdot \bar{C}_\ell^j(p) \cdot \bar{v}_\ell^i(p) \right)$;

15: end for

16: $\bar{v}_\ell^i = \bar{v}_\ell^i / \|\bar{v}_\ell^i\|_F$;

17: end for

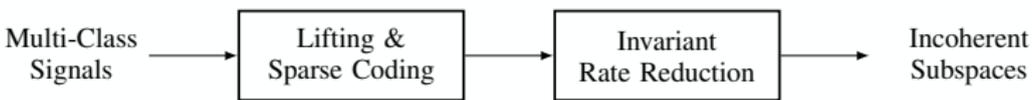
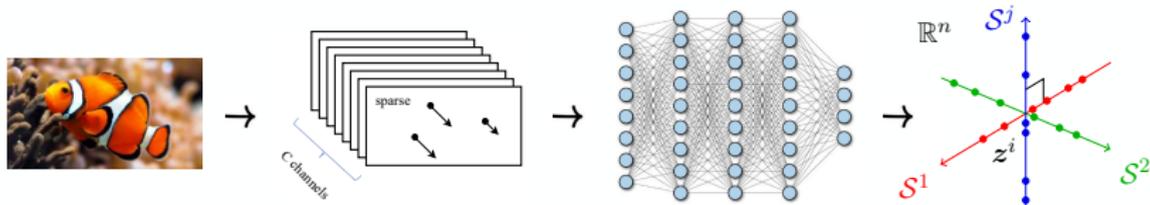
18: Set $\bar{Z}_\ell = \text{IDFT}(\bar{V}_\ell)$ as the feature at the ℓ -th layer;

19: $\frac{1}{2T} \sum_{p=0}^{T-1} \left(\log \det[\mathbf{I} + \alpha \bar{V}_\ell(p) \cdot \bar{V}_\ell(p)^*] - \frac{\text{tr}(\Pi^j)}{m} \log \det[\mathbf{I} + \alpha_j \bar{V}_\ell(p) \cdot \Pi^j \cdot \bar{V}_\ell(p)^*] \right)$;

20: end for

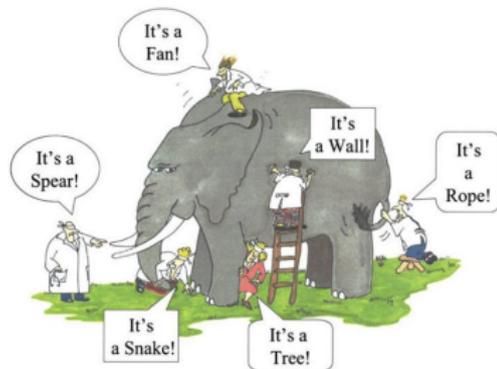
Ensure: features \bar{Z}_L , the learned filters $\{\bar{E}_\ell(p)\}_{\ell, p}$ and $\{\bar{C}_\ell^j(p)\}_{j, \ell, p}$.

Overall Process (the Elephant)



Necessary components:

- **sparse coding** for class separability;
- **deep networks** maximize rate reduction;
- **spectral computing** for shift-invariance;
- convolution, normalization, nonlinearity...



Experiment: 1D Cyclic Shift Invariance of 0 and 1

2000 training samples, 1980 testing, $C = 5$, $L = 3500$ -layers ReduNet.

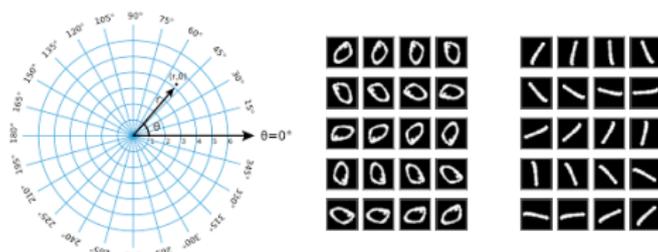


Figure: Left: Multi-channel feature representation of an image in polar coordinates. Right: Example of training/testing samples.

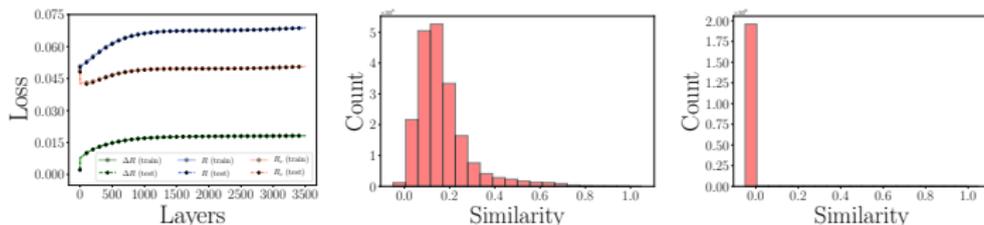


Figure: Left: Rates along the layers; Middle: cross-class cosine similarity among trainings; Right: similarity among testings.

Experiment: 1D Cyclic Shift Invariance of 0 and 1

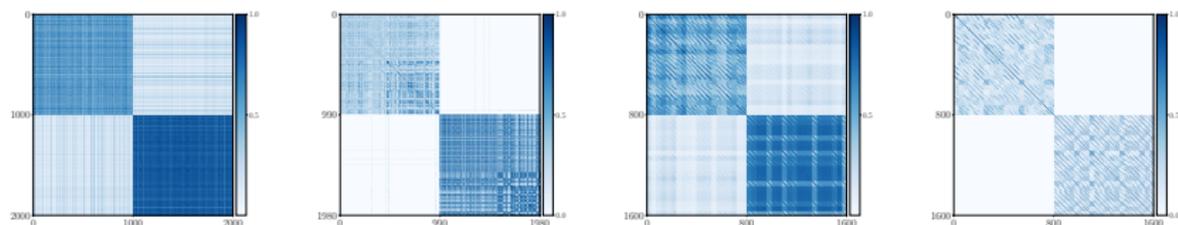


Figure: Left two: heat maps for training and testing. Right two: heat maps for one pair of samples at every possible shift.

Table: Network performance on digits with **all rotations**.

	REDUNET	REDUNET (INVARIANT)
ACC (ORIGINAL TEST DATA)	0.983	0.996
ACC (TEST WITH ALL SHIFTS)	0.707	0.993

1. Fooling CNNs with simple transformations, Engstrom et.al., 2017.
2. Why do deep convolutional networks generalize so poorly to small image transformations? Azulay & Weiss, 2018.

Experiment: 1D Cyclic Shift Invariance of All 10 Digits

100 training samples, 100 testing, $C = 20$, $L = 40$ -layers ReduNet.

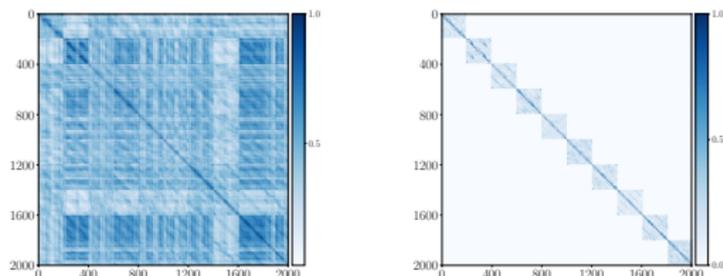


Figure: Heatmaps of cosine similarity among shifted training data $\mathbf{X}_{\text{shift}}$ (left) and learned features $\mathbf{Z}_{\text{shift}}$ (right).

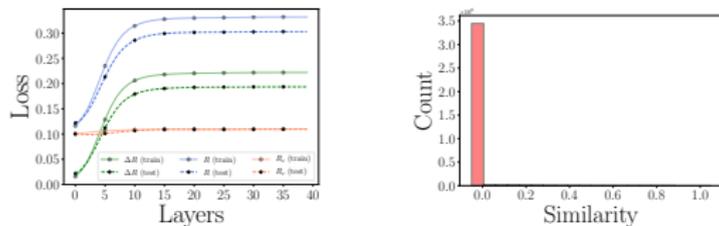


Figure: Left: Rates evolution with iterations; Right: histograms of the cosine similarity (in absolute value) between all pairs of features across different classes.

Experiment: 2D Cyclic Translation Invariance

1000 for training, 500 for testing, $C = 5$, $L = \mathbf{2000}$ -layers ReduNet.

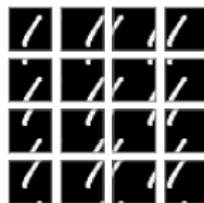
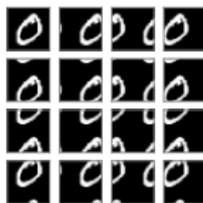
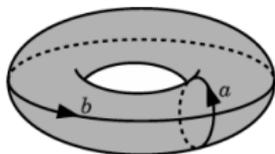
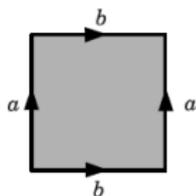


Table: Network performance on digits with **all translations**.

	REDUNET	REDUNET (INVARIANT)
ACC (ORIGINAL TEST DATA)	0.980	0.975
ACC (TEST WITH ALL SHIFTS)	0.540	0.909

1. Fooling CNNs with simple transformations, Engstrom et al., 2017.
2. Why do deep convolutional networks generalize so poorly to small image transformations? Azulay & Weiss, 2018.

Experiment: 2D Cyclic Trans. Invariance of All 10 Digits

100 training samples, 100 testing, $C = 75$, $L = 25$ -layers ReduNet.

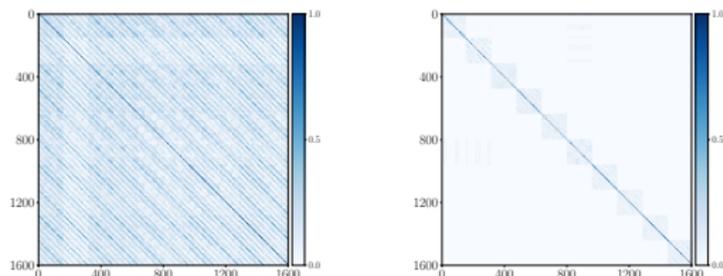


Figure: Heatmaps of cosine similarity among shifted training data $\mathbf{X}_{\text{shift}}$ (left) and learned features $\mathbf{Z}_{\text{shift}}$ (right).

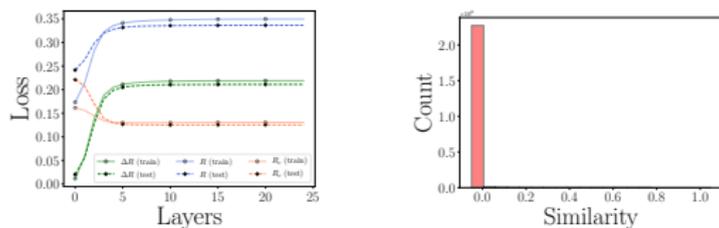


Figure: Left: Rates evolution with iterations; Right: histograms of the cosine similarity (in absolute value) between all pairs of features across different classes.

Experiment: Back Propagation of ReduNet (MNIST)

2D cyclic trans. of 10 digits, 500 training samples, all testing, $C = 16$, $L = 30$ -layers invariant ReduNet.

Initialization	Backpropagation	Test Accuracy
✓	✗	89.8%
✗	✓	93.2%
✓	✓	97.8%

Table: Test accuracy of 2D translation-invariant ReduNet, ReduNet-bp (without initialization), and ReduNet-bp (with initialization) on the MNIST dataset.

- **Backprop:** the ReduNet architecture *can* be fine-tuned by SGD and achieves better standard accuracy after back propagation;
- **Initialization:** using ReduNet for initialization can achieve better performance than the same architecture with random initialization.

Experiment: Back Propagation of ReduNet (CIFAR-10)

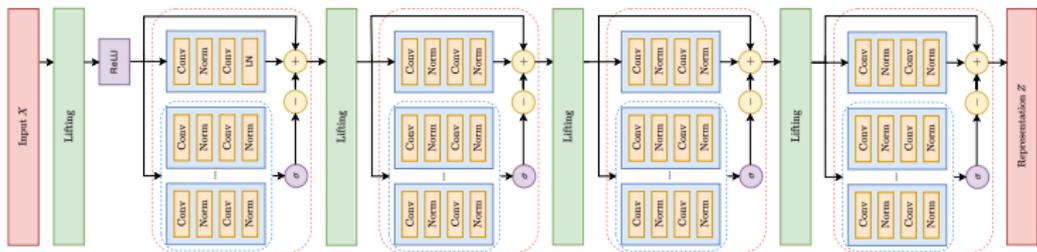


Figure: A ReduNet-inspired architecture

Table: Classification performance of ReduNet-inspired architecture on CIFAR10.

ReLU	TRAIN ACC	TEST ACC
✓	0.9997	0.8327
✗	0.9970	0.6542

Conclusions: Learn to Compress and Compress to Learn!

Principles of Parsimony:

- Clustering via compression: $\min_{\Pi} R^c(\mathbf{X}, \Pi)$
- Classification via compression: $\min_{\pi} \delta R^c(\mathbf{x}, \pi)$
- Representation via maximizing rate reduction: $\max_{\mathbf{Z}} \Delta R(\mathbf{Z}, \Pi)$
- Deep networks via optimizing rate reduction: $\dot{\mathbf{Z}} = \eta \cdot \frac{\partial \Delta R}{\partial \mathbf{Z}}$

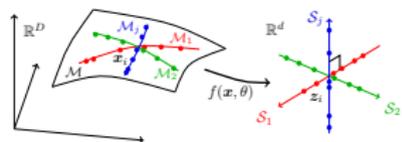
A Unified Framework:

- A principled objective for all settings of learning: **information gain**
- A principled approach to interpret deep networks: **optimization**

“Everything should be made as simple as possible, but not simpler.”
– Albert Einstein

Conclusions: Learn Linear Discriminative Representations

Compared to conventional deep neural networks:



	Conventional DNNs	ReduNets
Objectives	label fitting	information gain
Deep architectures	trial & error	iterative optimization
Layer operators	empirical	projected gradient
Shift invariance	CNNs+augmentation	invariant ReduNets
Initializations	random/pre-design	forward computed
Training/fine-tuning	back prop	forward/back prop
Interpretability	black box	white box
Representations	hidden/latent	incoherent subspaces

Open Problems: Theory

$$\mathbf{MCR}^2: \max_{\mathbf{Z} \subset \mathbb{S}^{d-1}, \mathbf{\Pi} \in \Omega} \Delta R(\mathbf{Z}, \mathbf{\Pi}, \epsilon) = R(\mathbf{Z}, \epsilon) - R^c(\mathbf{Z}, \epsilon | \mathbf{\Pi}).$$

- **Phase transition** phenomenon in clustering via compression?
- Statistical justification for **robustness** of \mathbf{MCR}^2 to label noise?
- **Optimal configurations** for broader conditions and distributions?
- Fundamental **tradeoff** between sparsity and invariance?
- **Jointly optimizing** both representation \mathbf{Z} and clustering $\mathbf{\Pi}$?

$$\mathbf{Joint Dynamics: } \dot{\mathbf{Z}} = \eta \cdot \frac{\partial \Delta R}{\partial \mathbf{Z}}, \quad \dot{\mathbf{\Pi}} = \gamma \cdot \frac{\partial \Delta R}{\partial \mathbf{\Pi}}.$$

Open Problems: Architectures and Algorithms

$$\text{ReduNet: } \bar{\mathbf{z}}_{\ell+1} \propto \bar{\mathbf{z}}_{\ell} + \eta \cdot \left[\bar{\mathbf{e}}_{\ell} \circledast \bar{\mathbf{z}}_{\ell} + \sigma([\bar{\mathbf{c}}_{\ell}^1 \circledast \bar{\mathbf{z}}_{\ell}, \dots, \bar{\mathbf{c}}_{\ell}^k \circledast \bar{\mathbf{z}}_{\ell}]) \right] \in \mathbb{S}^{d-1}.$$

- New architectures from **accelerated** gradient schemes?
- Conditions for channel-wise **separable and short** convolutions?
- Architectures from invariant rate reduction for **other groups**?
- **Transformer**:³

$$\left. \frac{\partial R(\mathbf{Z})}{\partial \mathbf{Z}} \right|_{\mathbf{Z}_{\ell}} = \alpha(\mathbf{I} + \alpha \mathbf{Z}_{\ell} \mathbf{Z}_{\ell}^*)^{-1} \mathbf{Z}_{\ell} \approx \underbrace{\alpha[\mathbf{Z}_{\ell} - \alpha \mathbf{Z}_{\ell}(\mathbf{Z}_{\ell}^* \mathbf{Z}_{\ell})]}_{\text{self-attention head}}.$$

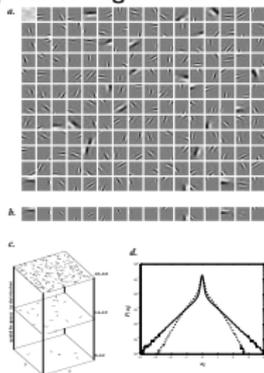
- Algorithmic architectures (or networks) for optimizing $\mathbf{\Pi}$?

³Attention: Self-Expression Is All You Need, Rene Vidal, 2022.

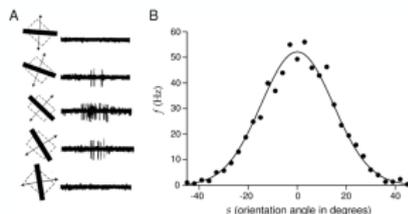
Open Directions: Extensions

- Data with other **dynamical or graphical** structures.
- Better transferability and robustness w.r.t. **low-dim structures**.
- Combine with a **generative model** (a generator or decoder).
- Sparse coding, spectral computing, subspace embedding in **nature**.⁴

sparse coding in visual cortex



Rate coding hypothesis: the signal conveyed by a neuron is in the rate of spiking. Spiking irregularity is largely due to noise and does not convey information.

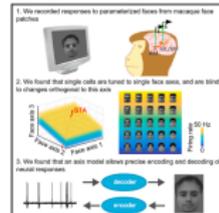


Cell

Article

The Code for Facial Identity in the Primate Brain

Graphical Abstract



Authors

Le Chang, Doris Y. Tsao

Correspondence
lechang@caltech.edu (L.C.),
dortsaoc@caltech.edu (D.Y.T.)

In Brief

Facial identity is encoded via a remarkably simple neural code that relies on the ability of neurons to distinguish facial features along specific axes in face space, disambiguating the long-standing assumption that single face cells encode individual faces.

Highlights

- Facial images can be linearly reconstructed using responses of ~200 face cells.
- Face cells display face tuning along dimensions orthogonal to the axis being coded.
- The axis model is more efficient, robust, and flexible than the exemplar model.
- Face patches ML/MF and AM carry complementary information about faces.

⁴figures from Bruno Olshausen of Neuroscience Dept., UC Berkeley.

From Open-Loop to Closed-Loop Representation

$$\text{MCR}^2 : \mathbf{X} \xrightarrow{f(\mathbf{x}, \theta)} \mathbf{Z}(\theta) : \max_{\theta} \Delta R(\mathbf{Z}(\theta), \mathbf{\Pi}, \epsilon).$$

Features learned are more interpretable, independent, rich, and robust.

Nevertheless:

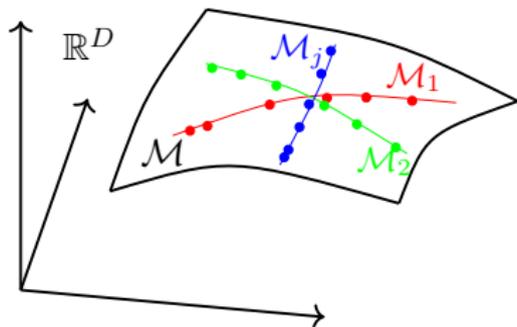
- Anything missing, anything unexpected: $\dim(\mathbf{X}_j) = \dim(\mathbf{Z}_j)$?
- Can we go from the feature \mathbf{Z}_j back to the data \mathbf{X}_j ?
- Is the learned LDR adequate to **generate** real-world (visual) data?

Can we establish an autoencoding between the data and the LDR:

$$\mathbf{X} \xrightarrow{f(\mathbf{x}, \theta)} \mathbf{Z}(\theta) \xrightarrow{g(\mathbf{z}, \eta)} \hat{\mathbf{X}}? \quad (29)$$

Low-dim Autoencoding for High-Dim Data

Assumption: the data \mathbf{X} lie on a low-dimensional submanifold $\mathbf{X} \subset \mathcal{M}$ or multiple ones: $\mathbf{X} \subset \cup_{j=1}^k \mathcal{M}_j$ in a high-dimensional space $\in \mathbb{R}^D$:

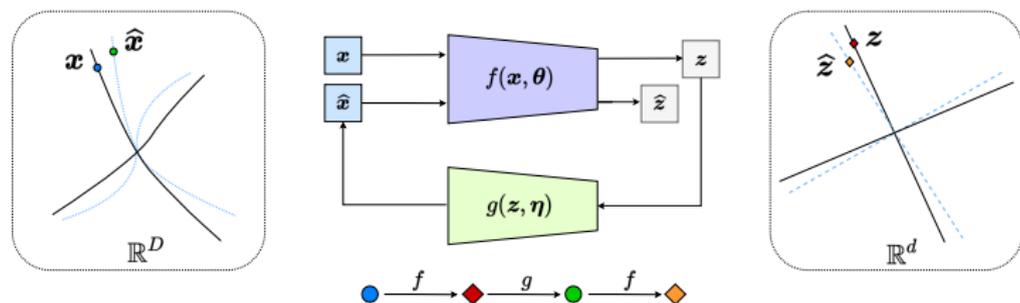


Goal: seeking a low-dim representation \mathbf{Z} in \mathbb{R}^d ($d \ll D$) for the data \mathbf{X} on low-dim submanifolds such that:

$$\mathbf{X} \subset \mathbb{R}^D \xrightarrow{f(\mathbf{x}, \theta)} \mathbf{Z} \subset \mathbb{R}^d \xrightarrow{g(\mathbf{z}, \eta)} \hat{\mathbf{X}} \approx \mathbf{X} \in \mathbb{R}^D. \quad (30)$$

CTRL: Dual Roles of the Encoder and Decoder

f is both an encoder and sensor; and g is both a decoder and controller. They form a **closed-loop feedback control system**:



A closed-loop notion of “**self-consistency**” between Z and \hat{Z} is achieved by a **pursuit-evasion game** between f as a “evader” and g as a “pursuer”:

$$\mathcal{D}(\mathbf{X}, \hat{\mathbf{X}}) \doteq \max_{\theta} \min_{\eta} \sum_{j=1}^k \Delta R \left(\underbrace{f(\mathbf{X}_j, \theta)}_{\mathbf{Z}_j(\theta)}, \underbrace{f(g(f(\mathbf{X}_j, \theta), \eta), \theta)}_{\hat{\mathbf{Z}}_j(\theta, \eta)} \right). \quad (31)$$

Overall CTRL Objective: Self-Consistency & Parsimony

The overall **minimax game** between the encoder f and decoder g :

- f *maximizes* the rate reduction of the features \mathbf{Z} of the data \mathbf{X} ;
- g *minimizes* the rate reduction of the features $\hat{\mathbf{Z}}$ of the decoded $\hat{\mathbf{X}}$.

A minimax program to learn a **multi-class LDR** for data $\mathbf{X} = \cup_{j=1}^k \mathbf{X}_j$:

$$\max_{\theta} \min_{\eta} \underbrace{\Delta R(f(\mathbf{X}, \theta))}_{\text{Expansive encode}} + \underbrace{\Delta R(h(\mathbf{X}, \theta, \eta))}_{\text{Compressive decode}} + \sum_{j=1}^k \underbrace{\Delta R(f(\mathbf{X}_j, \theta), h(\mathbf{X}_j, \theta, \eta))}_{\text{Contrastive \& Contractive}}$$

with $h(\mathbf{x}) \doteq f \circ g \circ f(\mathbf{x})$, or equivalently

$$\max_{\theta} \min_{\eta} \Delta R(\mathbf{Z}(\theta)) + \Delta R(\hat{\mathbf{Z}}(\theta, \eta)) + \sum_{j=1}^k \Delta R(\mathbf{Z}_j(\theta), \hat{\mathbf{Z}}_j(\theta, \eta)).$$

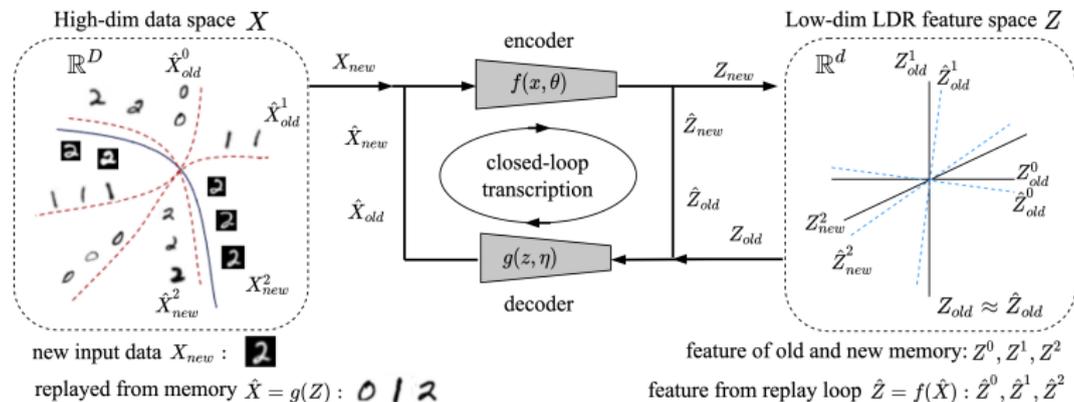
Characteristics of the Overall CTRL Objective

$$\max_{\theta} \min_{\eta} \Delta R(\mathbf{Z}(\theta)) + \Delta R(\hat{\mathbf{Z}}(\theta, \eta)) + \sum_{j=1}^k \Delta R(\mathbf{Z}_j(\theta), \hat{\mathbf{Z}}_j(\theta, \eta)).$$

- **Simplicity:** all terms are **closed-form** rate reduction on features.
- **Self-consistency:** enforced by **closed-loop** encoding and decoding.
- **Structured:** distribution of learned features \mathbf{Z} is **an LDR**.
- **No** need to specify a prior or a surrogate target distribution.
- **No** need of any direct explicit distance between \mathbf{X} and $\hat{\mathbf{X}}$.
- **No** more approximations or bounds for (KL-, JS-, W-) “distances”.
- **No** heuristics or regularizing terms in the objective.

Parsimony and self-consistency are all you need to model \mathbf{X} ?

Incremental Learning via Closed-Loop Transcription

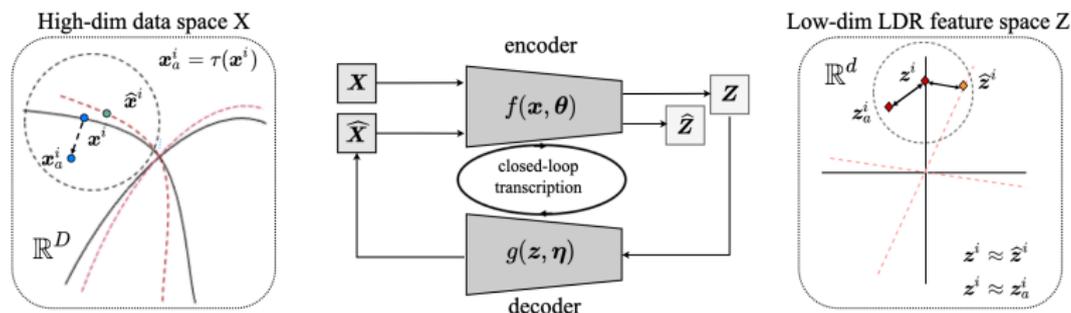


$$\max_{\theta} \min_{\eta} \Delta R(\mathbf{Z}) + \Delta R(\hat{\mathbf{Z}}) + \Delta R(\mathbf{Z}_{new}, \hat{\mathbf{Z}}_{new})$$

subject to $\Delta R(\mathbf{Z}_{old}, \hat{\mathbf{Z}}_{old}) = 0.$ (32)

Incremental Learning of Structured Memory, [arXiv:2202.05411](https://arxiv.org/abs/2202.05411) .
No Catastrophic Forgetting.

Unsupervised Learning via Closed-Loop Transcription



$$\max_{\theta} \min_{\eta} R(\mathbf{Z}) + \Delta R(\mathbf{Z}, \hat{\mathbf{Z}}) \quad (33)$$

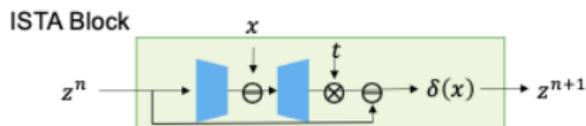
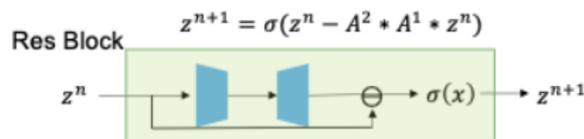
$$\text{subject to } \sum_{i \in N} \Delta R(z_{conv}^i, \hat{z}_{conv}^i) = 0, \text{ and } \sum_{i \in N} \Delta R(z^i, z_a^i) = 0.$$

Unsupervised Learning of Structured Representation,
[arXiv:2210.16782](https://arxiv.org/abs/2210.16782). **No Catastrophic Forgetting.**

Whitebox Network as Sparse Dictionary Learning (SDNet)

Use iterative convolution sparse coding (via ISTA) as layers.

Model size, memory, speed, and accuracy all comparable to ResNet.



$$z^{n+1} = \sigma(z^n - t \cdot A^T \odot (A * z^n - x))$$

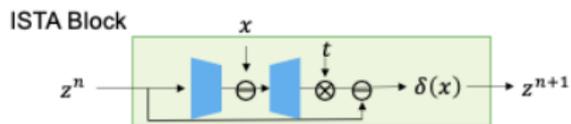
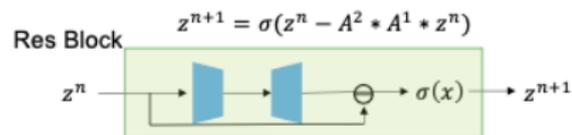
$$z^* = \underset{z}{\operatorname{argmin}} \lambda \|z\|_1 + \frac{1}{2} \|x - A \odot z\|_2^2$$

Dataset	Architecture	Model Size	Top-1 Acc	Memory	Speed
CIFAR-10	ResNet-18 [15]	11.2M	95.54%	1.0 GB	1600 n/s
	ResNet-34 [15]	21.1M	95.57%	2.0 GB	1000 n/s
	MDEQ [21]	11.1M	93.80%	2.0 GB	90 n/s
	SCN [9]	0.7M	94.36%	10.0GB	39 n/s
	SDNet-18 (ours)	11.2M	95.20%	1.2 GB	1500 n/s
	SDNet-34 (ours)	21.1M	95.57%	2.4 GB	900 n/s
CIFAR-100	ResNet-18 [15]	11.2M	77.82%	1.0 GB	1600 n/s
	ResNet-34 [15]	21.1M	78.39%	2.0 GB	1000 n/s
	MDEQ [21]	11.2M	74.12%	2.0 GB	90 n/s
	SCN [9]	0.7M	80.07%	10.0GB	39 n/s
	SDNet-18 (ours)	8.2M	78.31%	1.2 GB	1500 n/s
	SDNet-34 (ours)	10.2M	78.48%	2.4 GB	900 n/s
ImageNet	ResNet-18 [15]	11.7M	68.98%	24.1 GB	2100 n/s
	ResNet-34 [15]	21.5M	72.83%	32.3 GB	1400 n/s
	SCN [9]	9.8M	70.42%	95.1 GB	51 n/s
	SDNet-18 (ours)	11.7M	69.47%	37.6 GB	1800 n/s
	SDNet-34 (ours)	21.5M	72.67%	46.4 GB	1200 n/s

More Stable to Noises and More Robust to Perturbations.

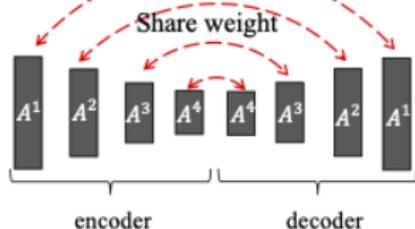
NeurIPS 2022

Transcription via Convolution Sparse Autoencoding



$$z^{n+1} = \sigma(z^n - t \cdot A^T \odot (A * z^n - x))$$

$$z^* = \underset{z}{\operatorname{argmin}} \lambda \|z\|_1 + \frac{1}{2} \|x - A \odot z\|_2^2$$



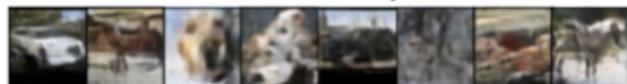
Input



Input+noise

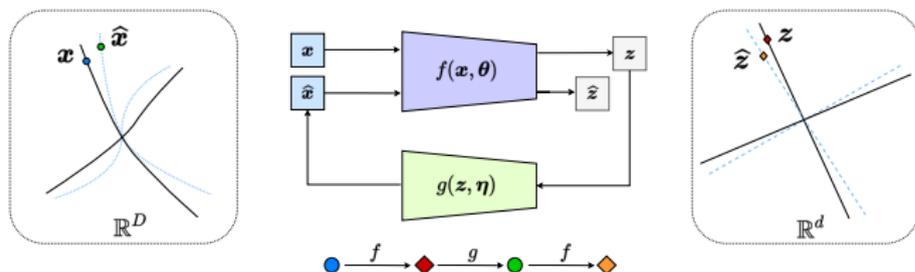


AE w Conv-layer



AE w SD-layer

Conclusions: Closed-Loop Transcription to an LDR



- **a universal learning engine:** transform external real-world data to a compact and structured internal (LDR) representation.
- **parsimony:** optimization of an information theoretical measure of rate reduction via an encoder and decoder.
- **self-consistency:** a pursuit-evasion game between a sensor and generator through a closed-loop feedback system.
- **white box:** learning objectives, network architectures, and learned representations.

Open Mathematical Problems about CTRL

For the closed-loop minimax rate reduction program:

$$\max_{\theta} \min_{\eta} \Delta R(\mathbf{Z}(\theta)) + \Delta R(\hat{\mathbf{Z}}(\theta, \eta)) + \sum_{j=1}^k \Delta R(\mathbf{Z}_j(\theta), \hat{\mathbf{Z}}_j(\theta, \eta)).$$

- **optimality**: characterization of the **equilibrium points**?
- **convergence** of the closed-loop control problem (**infinite-dim**)?
- **linearization** of distribution supports (**plastic manifold learning**)?
- **optimal density** of the distributions (**Brascamp-Lieb inequalities**)?
- **correct model selection** (**no under- or over-fitting**)?
- **guarantees** for approximate **sample-wise auto-encoding**?

Open Directions: Extensions and Connections

- How to **scale up** to hundreds and thousands of classes?
(variational forms for rate reduction, CVPR'22...)
- Internal computational mechanisms for **memory** forming (in Nature)?
(incremental/unsupervised learning without catastrophic forgetting.)
- Better **feedback** for generative quality and discriminative property?
- **Whitebox** architectures for closed-loop transcription (ReduNet like)?
- Closed-loop transcription to **other types of low-dim structures**?
(dynamical, causal, logical, symbolical, graphical, genetic...)

The principles of **parsimony and self-consistency** shall always rule!

References: Learning Transcription via Rate Reduction

- 1 **CTRL**: Closed-Loop Transcription to an LDR via Minimizing Rate Reduction
<https://arxiv.org/abs/2111.06636> (Entropy 2022)
- 2 **Incremental Learning** of Structured Memory via Closed-Loop Transcription
<https://arxiv.org/abs/2202.05411> (under submission)
- 3 **ReduNet**: A Whitebox Deep Network from Rate Reduction (JMLR 2022):
<https://arxiv.org/abs/2105.10446>
- 4 **Representation** via Maximal Coding Rate Reduction (NeurIPS 2020):
<https://arxiv.org/abs/2006.08558>
- 5 **Classification** via Minimal Incremental Coding Length (NIPS 2007):
http://people.eecs.berkeley.edu/~yima/psfile/MICL_SJIS.pdf
- 6 **Clustering** via Lossy Coding and Compression (TPAMI 2007):
<http://people.eecs.berkeley.edu/~yima/psfile/Ma-PAMI07.pdf>

More Supporting Materials

White-box Transformers via Sparse Rate Reduction (Sam Buchanan):

<https://arxiv.org/abs/2306.01129>

A New Textbook:

High-Dim Analysis with Low-Dim Models

<https://book-wright-ma.github.io/>

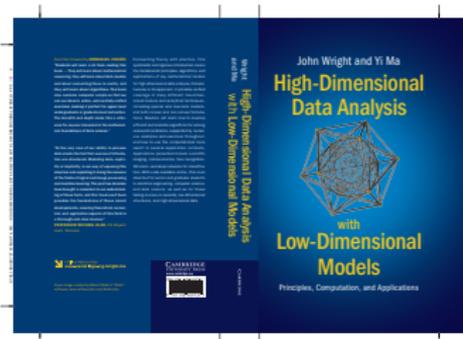
A New Course Berkeley EECS 208:

https://pages.github.berkeley.edu/UCB-EECS208/course_site/

<https://book-wright-ma.github.io/Lecture-Slides/>

A New International Conference:

Conference on Parsimony and Learning: <https://cpal.cc>



Parsimony and self-consistency are governing how to
learn
a **compact and structured** model for real-world data.

Thank you!
Questions, please?

“What I cannot create, I do not understand.”
– Richard Feynman



SIMONS
FOUNDATION