

Deep (Convolution) Networks from First Principles

Yi Ma

EECS Department, UC Berkeley

June 15, 2021

**Yaodong Yu, Kwan Ho Ryan Chan, Chong You,
Chaobing Song, Haozhi Qi, and John Wright**



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

- ① Motivation: Objectives of (Deep) Learning
- ② Prologue: Clustering and Classification via Compression
- ③ Representation via Principle of Maximal Rate Reduction
 - Theoretical justification
 - Experimental results
- ④ Deep Networks from Optimizing Rate Reduction
 - Deep networks as projected gradient ascent
 - Convolution networks from shift invariance
 - Preliminary experiments
- ⑤ Epilogue: Conclusions and Open Problems

High-Dim Data with Mixed Low-Dim Structures

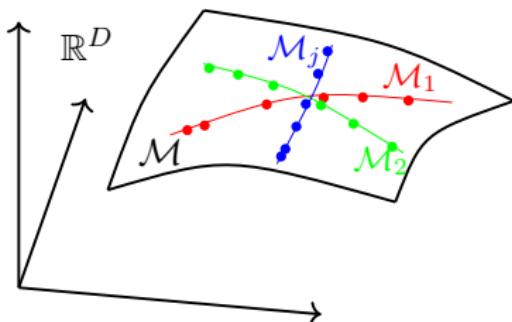


Figure: High-dimensional data $x \in \mathbb{R}^D$ lying on a mixture of low-dimensional submanifolds $\{\mathcal{M}_j\}$.

Three Related Objectives of Learning from Data:

- ① **Interpolation:** Identify which samples belong to the same structure.
- ② **Extrapolation:** Determine to which structure a new sample belongs.
- ③ **Representation:** Find most compact and discriminative representations.

Learning Lumped into a Black Box (Deep Learning)

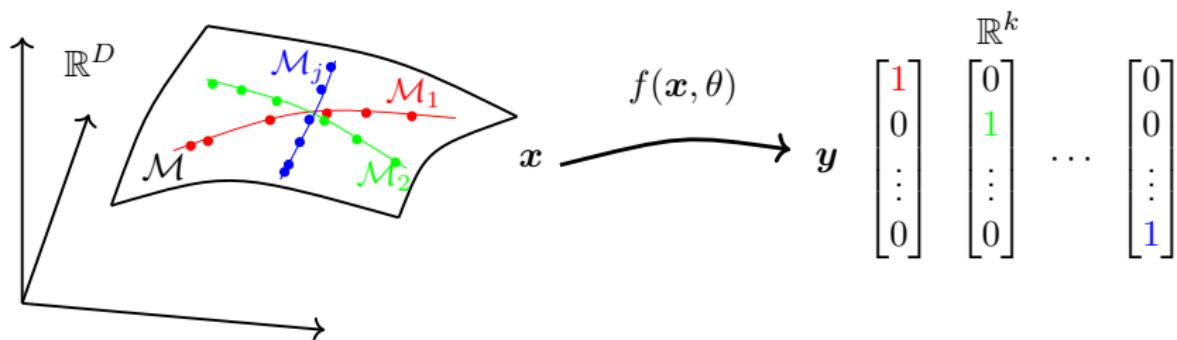


Figure: Black Box Classification: y is the class label of x represented as a “one-hot” vector in \mathbb{R}^k . To learn a nonlinear mapping $f(\cdot, \theta) : x \mapsto y$, say modeled by a deep network.

Fitting Class Labels via a Deep Network

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

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

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**.

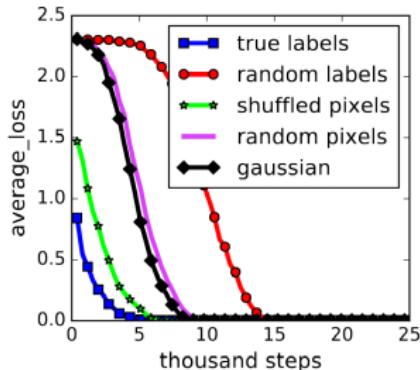


Figure: [Zhang et al, ICLR'17]

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

$$\max_{\theta \in \Theta} \text{IB}(\mathbf{x}, \mathbf{y}, \mathbf{z}(\theta)) \doteq I(\mathbf{z}(\theta), \mathbf{y}) - \beta I(\mathbf{x}, \mathbf{z}(\theta)), \quad \beta > 0, \quad (2)$$

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.

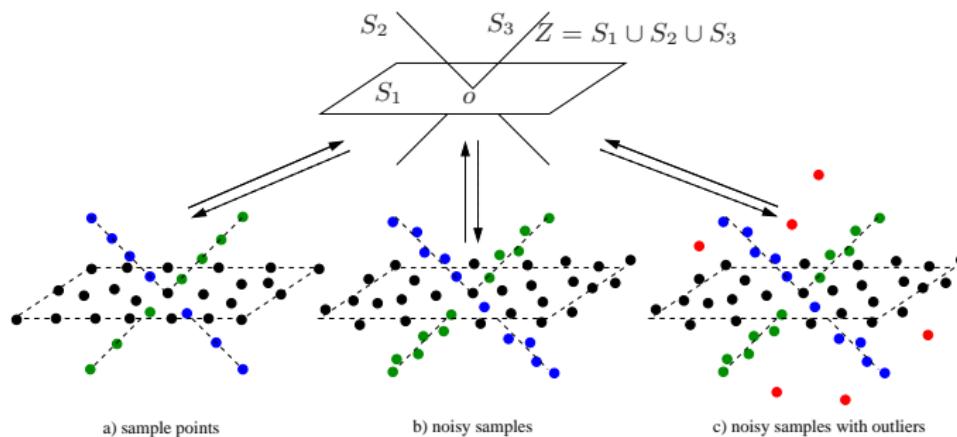
Reality check: principled formulations are replaced with approximate bounds, grossly simplifying assumptions, heuristics, even *ad hoc* tricks and hacks.

How to provide any performance guarantees at all?

A Principled Computational Approach

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

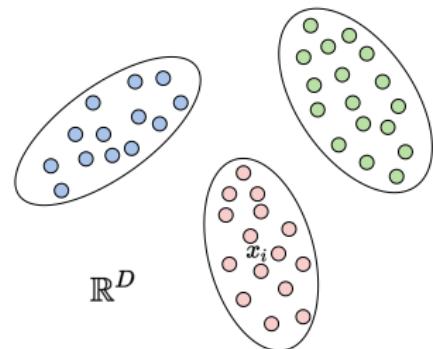
learn to compress, and compress to learn!



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]

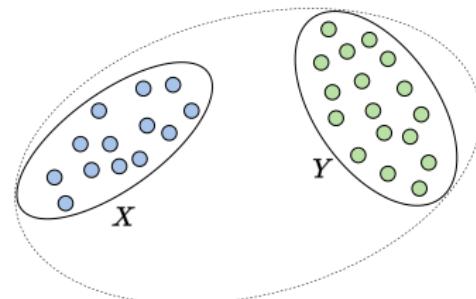
A Fundamental Idea:

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

Cluster Criterion:

Whether the number of binary bits required to store the data:

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



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

Coding Length Function for Subspace-Like Data

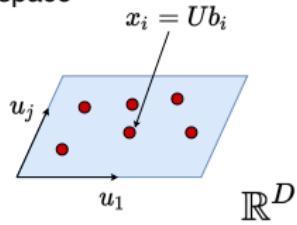
Theorem (Ma, TPAMI'07)

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:

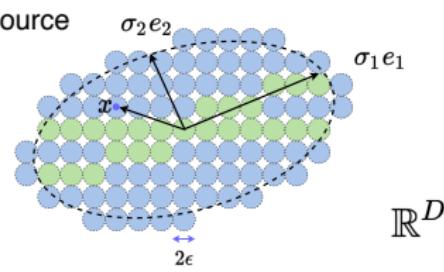
$$L(\mathbf{X}, \epsilon) \doteq \left(\frac{m+D}{2} \right) \log \det \left(\mathbf{I} + \frac{D}{m\epsilon^2} \mathbf{X} \mathbf{X}^\top \right).$$

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



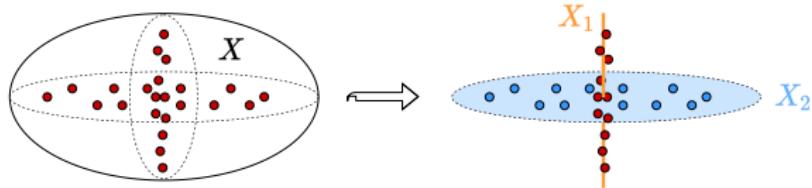
Gaussian source



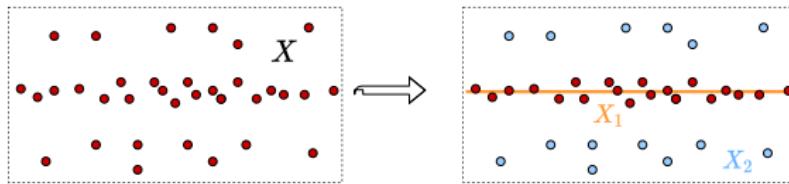
Cluster to Compress

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

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

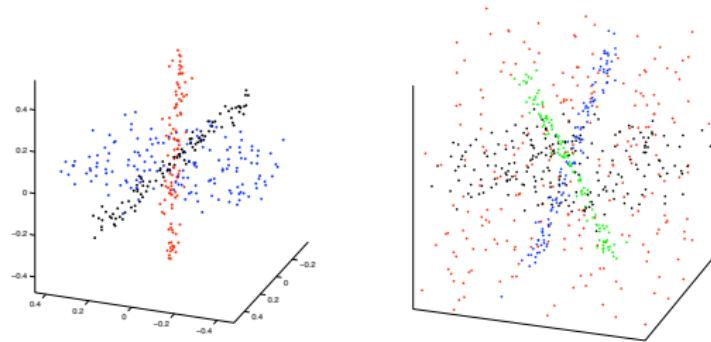
$$\min L^c(\mathbf{X}) \doteq L(\mathbf{X}_1) + \dots + L(\mathbf{X}_k) + H(|\mathbf{X}_1|, \dots, |\mathbf{X}_k|).$$

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

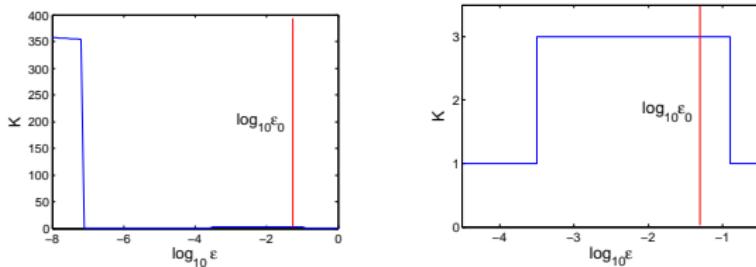
- 1: **input:** the data $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m] \in \mathbb{R}^{D \times m}$ and a distortion $\epsilon^2 > 0$.
- 2: initialize \mathcal{S} as a set of sets with a single datum $\{\mathcal{S} = \{\mathbf{x}\} \mid \mathbf{x} \in \mathbf{X}\}$.
- 3: **while** $|\mathcal{S}| > 1$ **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**

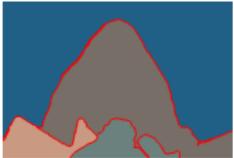


A strikingly sharp **phase transition** w.r.t. quantization ϵ



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

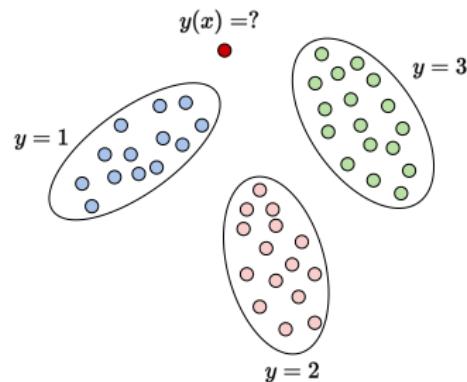
(e) Water

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...)

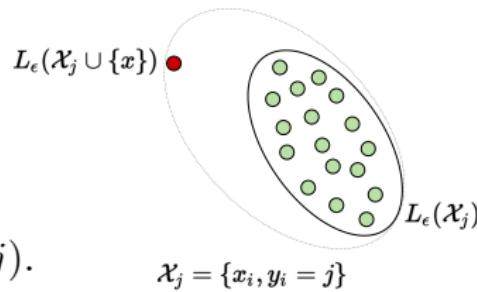
Classify to Compress

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

A Fundamental Idea:

Count additional #bits needed
to encode a query sample 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)

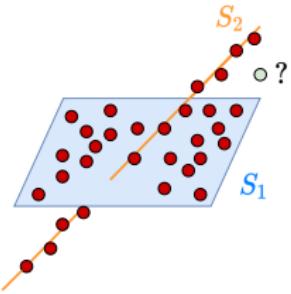
As the number of samples m goes to infinity, the MICL criterion converges at a rate of $O(m^{-1/2})$ to the following criterion:

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

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

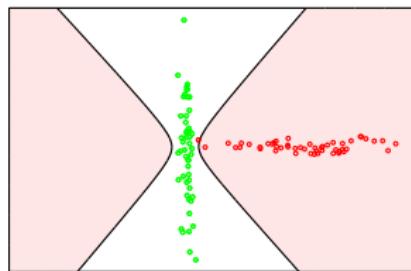
Everything else equal, MICL prefers a class with higher effective dimension.

Err on the side of caution!

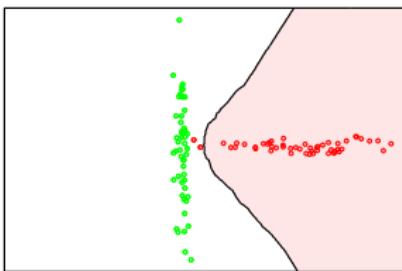


Extrapolation of Low-Dim Structure for Classification

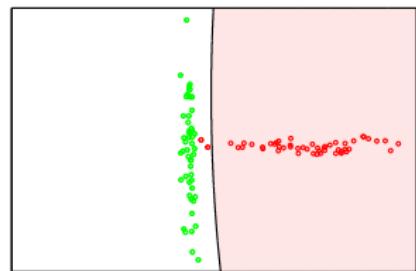
Figure: A truly extrapolating (nearest subspace) classifier!



(a) MICL



(b) k-Nearest Neighbors



(c) SVM-RBF

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

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

3. Represent Mixed Data

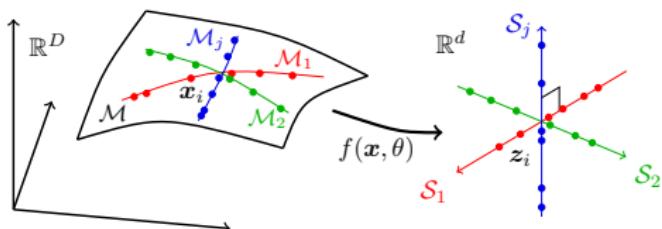
Given

samples $\mathbf{x} \in \mathbb{R}^D$ drawn from
a mixture of k submanifolds

$\mathcal{M} = \{\mathcal{M}_j\}_{j=1}^k$, we

seek a good representation

$\mathbf{z} \in \mathbb{R}^d$ through a continuous
mapping: $f(\mathbf{x}, \theta) : \mathbb{R}^D \rightarrow \mathbb{R}^d$.



Goals of “re-present” the data $f(\mathbf{x}, \theta) : \mathbf{x} \mapsto \mathbf{z}$:

- from non-parametric (samples) to more compact (models).
- from nonlinear structures in \mathbf{x} to linear in \mathbf{z} .
- from separable \mathbf{x} to maximally discriminative \mathbf{z} .

What is a good representation? (Why a deep neural network?)

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 Representation*: Dimension (or variance) of features for each class/cluster should be **as large as possible**.

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

Why not cross entropy? Prevalence of **neural collapse** during the terminal phase of deep learning training, Papyan, Han, and Donoho, 2020.

Measure of Compactness for a Linear Representation

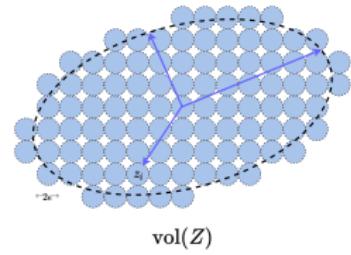
Consider a feature mapping:

$$\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_m] \in \mathbb{R}^{D \times m} \xrightarrow{f(\mathbf{x}, \theta)} \mathbf{Z}(\theta) = [z_1, z_2, \dots, z_m] \in \mathbb{R}^{d \times m}.$$

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

$$R(\mathbf{Z}, \epsilon) \doteq \frac{1}{2} \log \det \left(\mathbf{I} + \frac{d}{m\epsilon^2} \mathbf{Z}\mathbf{Z}^\top \right). \quad (3)$$

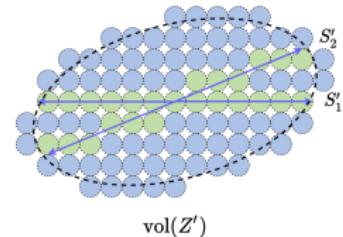
Rate distortion is an intrinsic measure for the volume of all features.



Measure of Compactness for Mixed Representations

The features Z of multi-class data may be partitioned into multiple subsets:

$$Z = Z_1 \cup Z_2 \cup \dots \cup Z_k.$$



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

$$R^c(Z, \epsilon | \Pi) \doteq \sum_{j=1}^k \frac{\text{tr}(\Pi_j)}{2m} \log \det \left(I + \frac{d}{\text{tr}(\Pi_j)\epsilon^2} Z \Pi_j Z^\top \right), \quad (4)$$

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

Learning Representation to Cluster and Classify

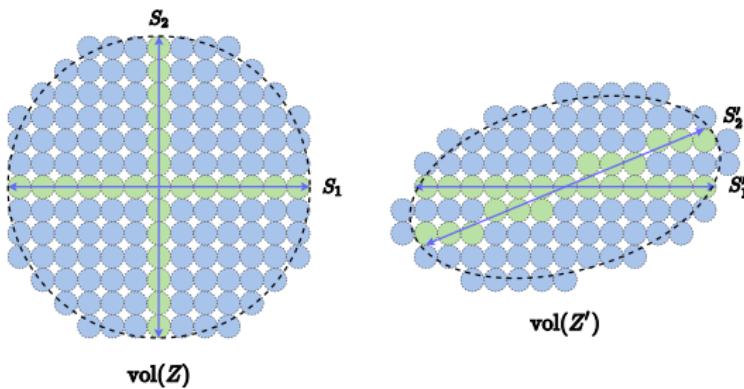
A Fundamental Idea: maximize the **difference** between the coding rate of all features and the average rate of features in the classes:

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

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

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

Maximal Coding Rate Reduction (MCR^2)

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

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

$$\mathbf{X} \xrightarrow{f(\mathbf{x}, \theta)} \mathbf{Z}(\theta) \xrightarrow{\Pi, \epsilon} \Delta R(\mathbf{Z}(\theta), \Pi, \epsilon) \quad (5)$$

so as to **Maximize the Coding Rate Reduction (MCR^2)**:

$$\begin{aligned} \max_{\theta} \quad & \Delta R(\mathbf{Z}(\theta), \Pi, \epsilon) = R(\mathbf{Z}(\theta), \epsilon) - R^c(\mathbf{Z}(\theta), \epsilon \mid \Pi), \\ \text{subject to} \quad & \|\mathbf{Z}_j(\theta)\|_F^2 = m_j, \quad \Pi \in \Omega. \end{aligned} \quad (6)$$

Since ΔR is *monotonic* in the scale of Z , one needs to:

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

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 $Z^* = Z_1^* \cup \dots \cup Z_k^*$ is the optimal solution that maximizes the rate reduction (6). 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. $(Z_i^*)^\top Z_j^* = \mathbf{0}$ 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. $\text{rank}(Z_j^*) = d_j$. In addition, the largest $d_j - 1$ singular values of Z_j^* are equal.

A new slogan, beyond Aristotle:

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

Comparison to Orthogonal Low-Rank Embedding (OLE)

[Lezama, Qiu, Musé, and Sapiro, CVPR'18]

The following loss as a geometric regularizer for the cross entropy (1):

$$\max_{\theta} \text{OLE}(\mathbf{Z}(\theta), \mathbf{\Pi}) \doteq \|\mathbf{Z}(\theta)\|_* - \sum_{j=1}^k \|\mathbf{Z}_j(\theta)\|_*.$$

The nuclear norm $\|\cdot\|_*$ in OLE is convex but non-smooth; the $\log \det(\cdot)$ in MCR² is concave but smooth.

OLE is always negative and reaches the maximum 0 iff the subspaces are orthogonal, regardless of the dimension.

OLE cannot avoid neural collapse whereas MCR² does.

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**:

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

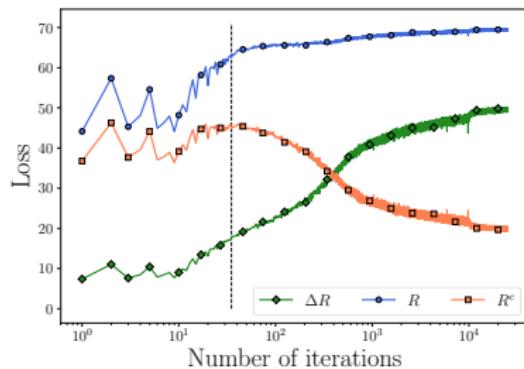
The learned features of such pairs of samples together with their augmentations \mathbf{Z}_i and \mathbf{Z}_j should have large rate reduction:

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

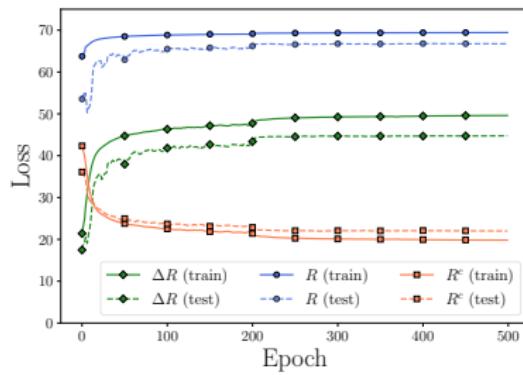
MCR² contrasts triplets, quadruplets, or any number of sets.

Experiment I: Supervised Deep Learning

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



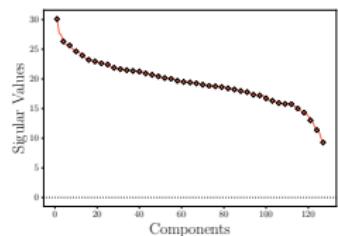
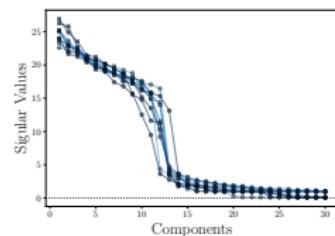
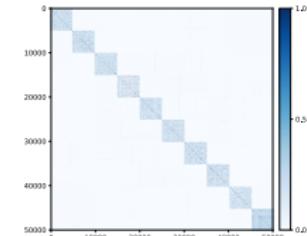
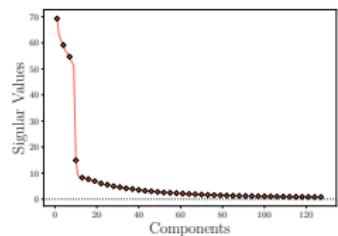
(a)



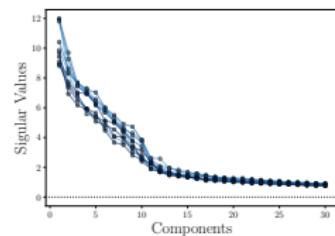
(b)

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

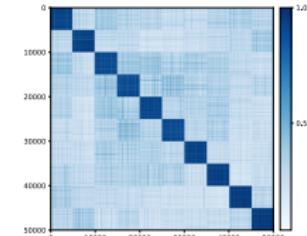
Visualization of Learned Representations Z

(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)

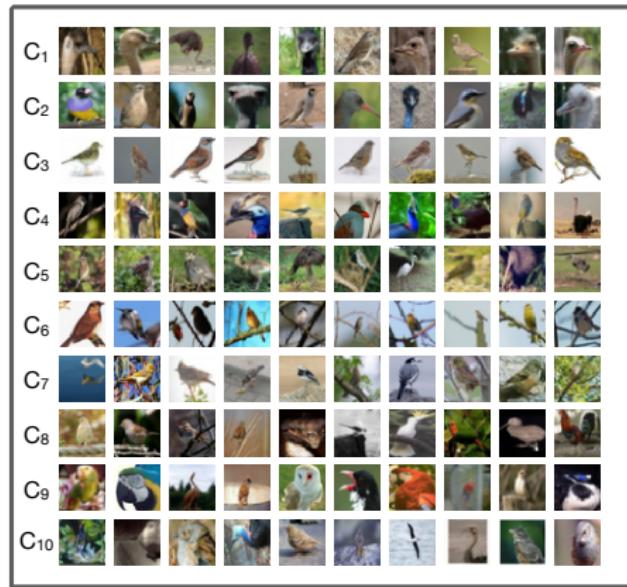


(f) CE (cosine similarity)

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

No neural collapse!

Visualization - Samples along Principal Components



(a) Bird



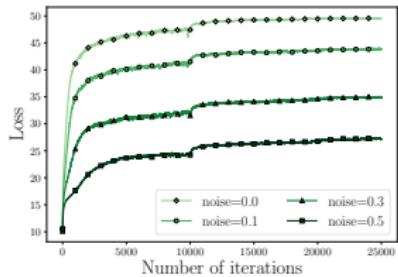
(b) Ship

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

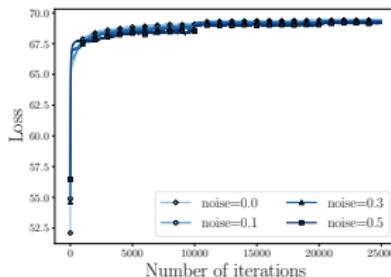
Robustness to Label Noise

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

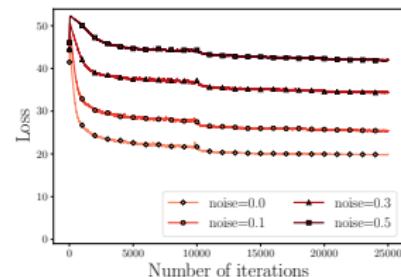
	RATIO=0.1	RATIO=0.2	RATIO=0.3	RATIO=0.4	RATIO=0.5
CE TRAINING	90.91%	86.12%	79.15%	72.45%	60.37%
MCR ² TRAINING	91.16%	89.70%	88.18%	86.66%	84.30%



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



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



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

Figure: Evolution of R , R^c , ΔR of MCR² during training with corrupted labels.

Represent only what can be jointly compressed.

Experiment II: Self-supervised Learning

Without label information, we use the MCR² principle (6) to learn representations Z that are *invariant* to certain class of transformations/augmentations, say \mathcal{T} with a distribution $P_{\mathcal{T}}$.

- ① Given a mini-batch of data $\{\mathbf{x}_j\}_{j=1}^k$, augment each sample \mathbf{x}_j with n augmentations $\{\tau_i(\cdot)\}_{i=1}^n$ randomly drawn from $P_{\mathcal{T}}$.
- ② Label all the augmented samples $\mathbf{X}_j = [\tau_1(\mathbf{x}_j), \dots, \tau_n(\mathbf{x}_j)]$ of \mathbf{x}_j as the j -th class, and Z_j the corresponding learned features.
- ③ Using this self-labeled data, train feature mapping $f(\cdot, \theta)$ via maximizing the MCR² in Eq. (6).

Experimental Setup: Train ResNet18 on CIFAR10 dataset, feature dimension $d = 128$, precision $\epsilon = 0.5$. For every mini-batch, the total number of samples for training is $m = kn$.

Clustering of Real Datasets

Figure: Clustering results on CIFAR10, CIFAR100, and STL10 datasets.

DATASET	METRIC	K-MEANS	JULE	RTM	DEC	DAC	DCCM	MCR ² -CTRL
CIFAR10	NMI	0.087	0.192	0.197	0.257	0.395	0.496	0.630
	ACC	0.229	0.272	0.309	0.301	0.521	0.623	0.684
	ARI	0.049	0.138	0.115	0.161	0.305	0.408	0.508
CIFAR100	NMI	0.084	0.103	-	0.136	0.185	0.285	0.387
	ACC	0.130	0.137	-	0.185	0.237	0.327	0.375
	ARI	0.028	0.033	-	0.050	0.087	0.173	0.178
STL10	NMI	0.124	0.182	-	0.276	0.365	0.376	0.446
	ACC	0.192	0.182	-	0.359	0.470	0.482	0.491
	ARI	0.061	0.164	-	0.186	0.256	0.262	0.290

Set the mini-batch size as $k = 20$, number of augmentations for each sample as $n = 50$ and the precision parameter as $\epsilon^2 = 0.5$. Compare with JULE, RTM [Nina'19], DEC [Xie'16], DAC [Chang'17], and DCCM [Wu'19].

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:

$$\max_{\mathbf{Z}} \Delta R(\mathbf{Z}) \doteq \underbrace{\frac{1}{2} \log \det (\mathbf{I} + \alpha \mathbf{Z} \mathbf{Z}^*)}_{R(\mathbf{Z})} - \underbrace{\sum_{j=1}^k \frac{\gamma_j}{2} \log \det (\mathbf{I} + \alpha_j \mathbf{Z} \boldsymbol{\Pi}^j \mathbf{Z}^*)}_{R_c(\mathbf{Z}, \boldsymbol{\Pi})}, \quad (7)$$

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

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

$$\mathbf{Z}_{\ell+1} \propto \mathbf{Z}_\ell + \eta \cdot \left. \frac{\partial \Delta R}{\partial \mathbf{Z}} \right|_{\mathbf{Z}_\ell} \quad \text{subject to} \quad \mathbf{Z}_{\ell+1} \subset \mathbb{S}^{d-1}. \quad (8)$$

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:

$$\frac{1}{2} \frac{\partial \log \det(\mathbf{I} + \alpha \mathbf{Z} \mathbf{Z}^*)}{\partial \mathbf{Z}} \Big|_{\mathbf{Z}_\ell} = \underbrace{\alpha (\mathbf{I} + \alpha \mathbf{Z}_\ell \mathbf{Z}_\ell^*)^{-1}}_{\mathbf{E}_\ell \in \mathbb{R}^{d \times d}} \mathbf{Z}_\ell, \quad (9)$$

$$\frac{1}{2} \frac{\partial (\gamma_j \log \det(\mathbf{I} + \alpha_j \mathbf{Z} \mathbf{\Pi}^j \mathbf{Z}^*))}{\partial \mathbf{Z}} \Big|_{\mathbf{Z}_\ell} = \gamma_j \underbrace{\alpha_j (\mathbf{I} + \alpha_j \mathbf{Z}_\ell \mathbf{\Pi}^j \mathbf{Z}_\ell^*)^{-1}}_{\mathbf{C}_\ell^j \in \mathbb{R}^{d \times d}} \mathbf{Z}_\ell \mathbf{\Pi}^j. \quad (10)$$

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

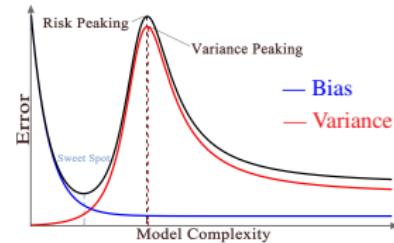
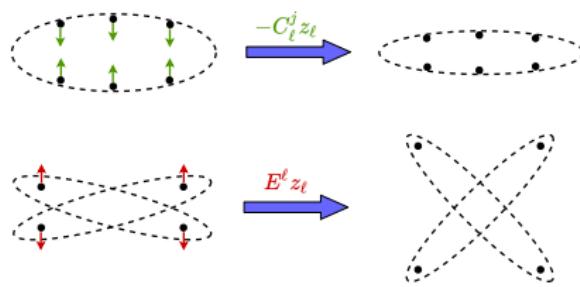
$$\frac{\partial \Delta R}{\partial \mathbf{Z}} \Big|_{\mathbf{Z}_\ell} = \underbrace{\mathbf{E}_\ell}_{\text{Expansion}} \mathbf{Z}_\ell - \sum_{j=1}^k \gamma_j \underbrace{\mathbf{C}_\ell^j}_{\text{Compression}} \mathbf{Z}_\ell \mathbf{\Pi}^j \in \mathbb{R}^{d \times m}. \quad (11)$$

Interpretation of the Linear Operators E and C^j

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

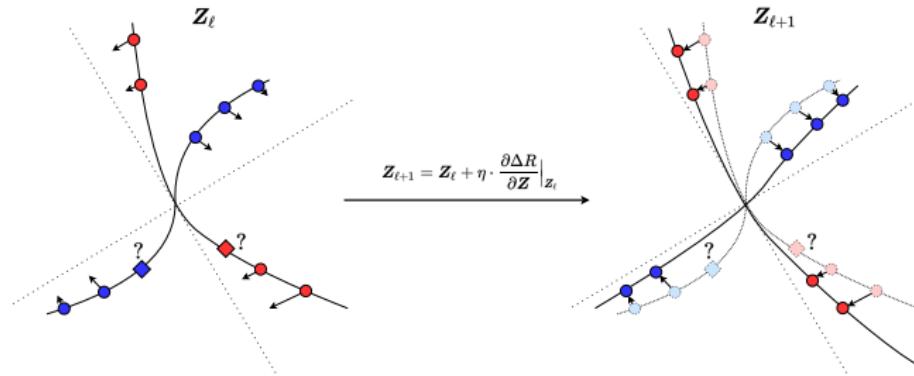
$$E_\ell z_\ell = \alpha(z_\ell - Z_\ell q_\ell^*) \quad \text{with} \quad q_\ell^* \doteq \arg \min_{q_\ell} \alpha \|z_\ell - Z_\ell q_\ell\|_2^2 + \|q_\ell\|_2^2.$$

$E_\ell z_\ell$ and $C_\ell^j z_\ell$ are the “residuals” of z_ℓ against the subspaces spanned by columns of Z_ℓ and 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 $\mathbf{z} \in \mathbb{R}^d$:

$$\frac{\partial \Delta R}{\partial \mathbf{Z}} \Big|_{\mathbf{Z}_\ell} = \mathbf{E}_\ell \mathbf{Z}_\ell - \sum_{j=1}^k \gamma_j \mathbf{C}_\ell^j \mathbf{Z}_\ell \underbrace{\Pi^j}_{\text{known}} \in \mathbb{R}^{d \times m}, \quad (12)$$

$$g(\mathbf{z}_\ell, \boldsymbol{\theta}_\ell) \doteq \mathbf{E}_\ell \mathbf{z}_\ell - \sum_{j=1}^k \gamma_j \mathbf{C}_\ell^j \mathbf{z}_\ell \underbrace{\pi^j(\mathbf{z}_\ell)}_{\text{unknown}} \in \mathbb{R}^d. \quad (13)$$

Estimate of the Membership $\pi^j(z_\ell)$

Estimate the membership $\pi^j(z_\ell)$ with “softmax” on the residuals $\|C_\ell^j z_\ell\|$:

$$\pi^j(z_\ell) \approx \hat{\pi}^j(z_\ell) \doteq \frac{\exp(-\lambda \|C_\ell^j z_\ell\|)}{\sum_{j=1}^k \exp(-\lambda \|C_\ell^j z_\ell\|)} \in [0, 1]. \quad (14)$$

Thus the weighted residuals for contracting:

$$\sigma([C_\ell^1 z_\ell, \dots, C_\ell^k z_\ell]) \doteq \sum_{j=1}^k \gamma_j C_\ell^j z_\ell \cdot \hat{\pi}^j(z_\ell) \in \mathbb{R}^d. \quad (15)$$

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

$$\sigma(z_\ell) \approx z_\ell - \sum_{j=1}^k \text{ReLU}(P_\ell^j z_\ell), \quad (16)$$

The ReduNet for Optimizing Rate Reduction

Iterative projected gradient ascent (PGA) :

$$\mathbf{z}_{\ell+1} \propto \mathbf{z}_\ell + \eta \cdot \underbrace{\left[E_\ell \mathbf{z}_\ell + \sigma([C_\ell^1 \mathbf{z}_\ell, \dots, C_\ell^k \mathbf{z}_\ell]) \right]}_{g(\mathbf{z}_\ell, \boldsymbol{\theta}_\ell)} \quad \text{s.t. } \mathbf{z}_{\ell+1} \in \mathbb{S}^{d-1}, \quad (17)$$

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

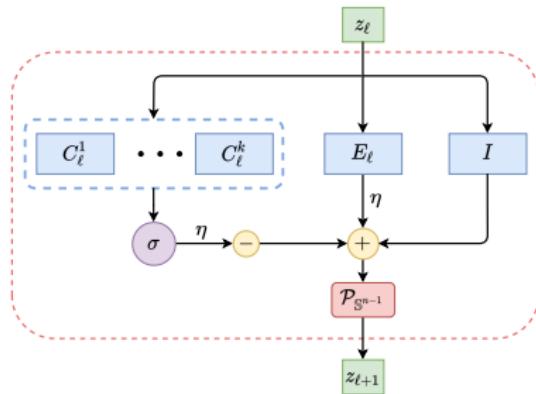


Figure: One layer of the ReduNet: one PGA iteration.

The ReduNet versus ResNet or ResNeXt

Iterative projected gradient ascent (PGA):

$$\mathbf{z}_{\ell+1} \propto \mathbf{z}_\ell + \eta \cdot \underbrace{\left[\mathbf{E}_\ell \mathbf{z}_\ell + \sigma([C_\ell^1 \mathbf{z}_\ell, \dots, C_\ell^k \mathbf{z}_\ell]) \right]}_{g(\mathbf{z}_\ell, \theta_\ell)} \quad \text{s.t. } \mathbf{z}_{\ell+1} \in \mathbb{S}^{d-1}, \quad (18)$$

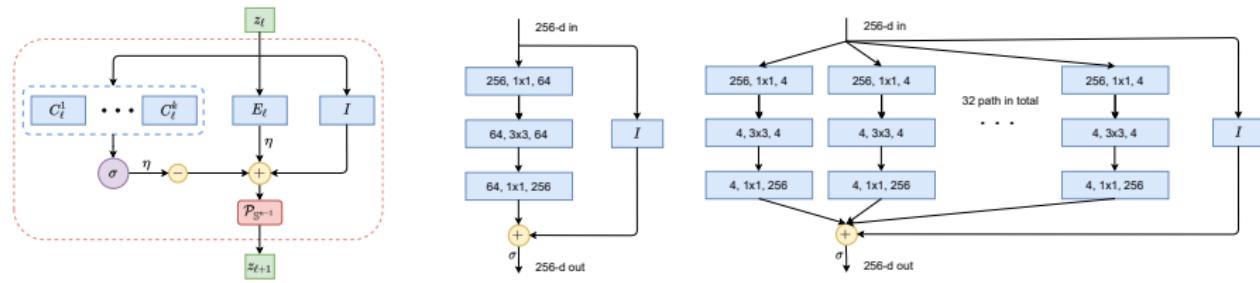


Figure: Left: **ReduNet**. Middle and Right: **ResNet** [He et. al. 2016] and **ResNeXt** [Xie et. al. 2017] (hundreds of layers).

Forward construction instead of back propagation!

The ReduNet versus Mixture of Experts

Approximate iterative projected gradient ascent (PGA) :

$$\mathbf{z}_{\ell+1} \propto \mathbf{z}_\ell + \eta \cdot \underbrace{\left[\mathbf{E}_\ell \mathbf{z}_\ell + \sigma([\mathbf{C}_\ell^1 \mathbf{z}_\ell, \dots, \mathbf{C}_\ell^k \mathbf{z}_\ell]) \right]}_{g(\mathbf{z}_\ell, \theta_\ell)} \quad \text{s.t. } \mathbf{z}_{\ell+1} \in \mathbb{S}^{d-1}, \quad (19)$$

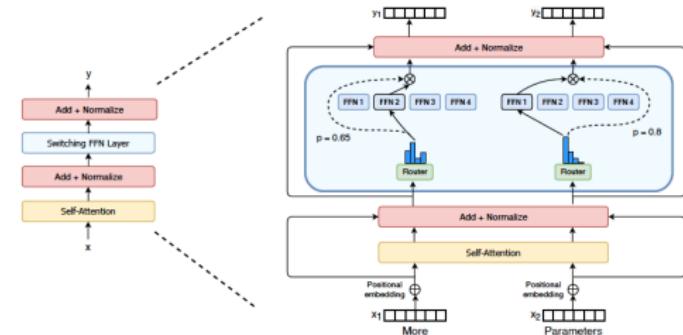
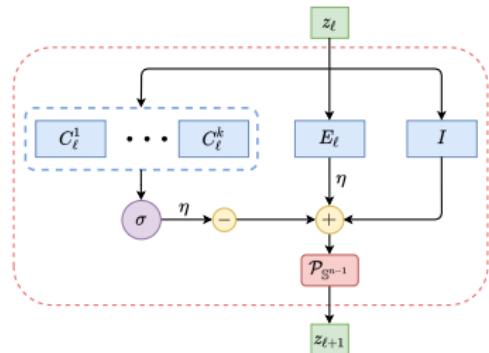


Figure: Left: ReduNet layer. Right: Mixture of Experts [Shazeer et. al. 2017] or Switched Transformer [Fedus et. al. 2021] (1.7 trillion parameters).

Forward construction instead of back propagation!

ReduNet Features for Mixture of Gaussians

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

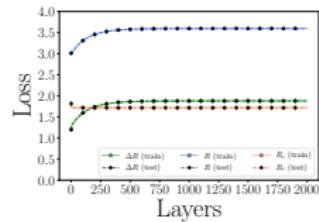
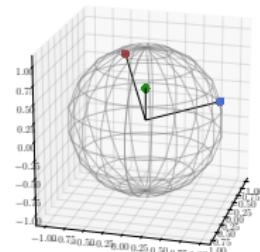
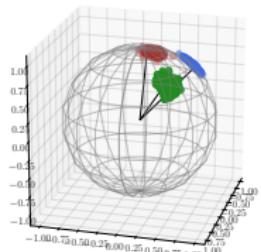
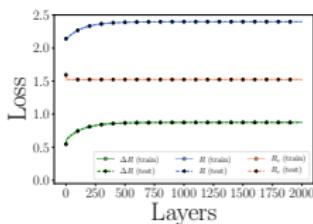
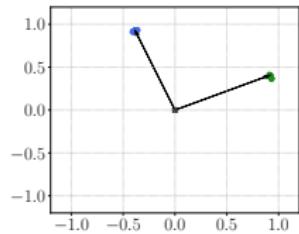
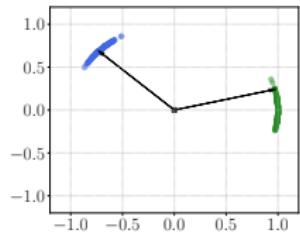


Figure: Left: original samples X and ReduNet features $Z = f(Z, \theta)$ for 2D and 3D Mixture of Gaussians. Right: plots for the progression of values of the rates.

Group Invariant Classification

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

$$\text{Group Invariance: } f(\mathbf{x} \circ \mathbf{g}, \theta) \sim f(\mathbf{x}, \theta), \quad \forall \mathbf{g} \in \mathbb{G}, \quad (20)$$

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

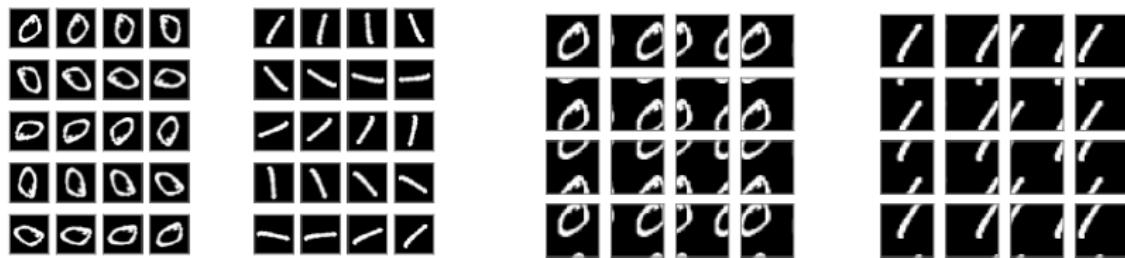


Figure: Left: 1D rotation S^1 ; Right: 2D cyclic translation T^2 .

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

Group Invariant Classification

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

$$\text{Group Invariance: } f(x \circ g, \theta) \sim f(x, \theta), \quad \forall g \in \mathbb{G}, \quad (21)$$

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

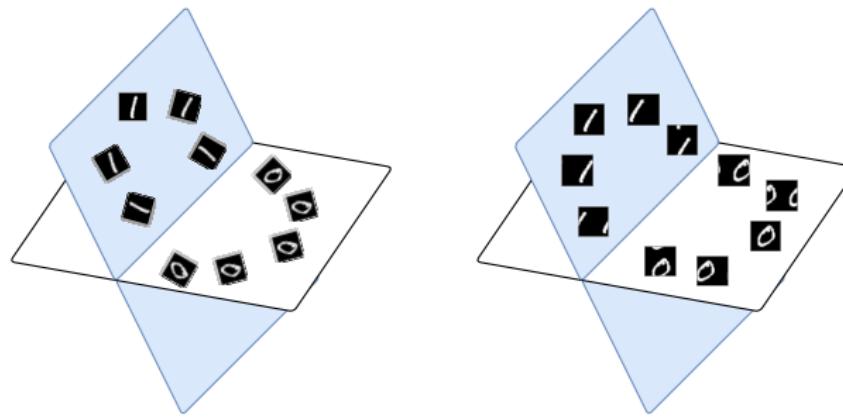


Figure: Embed all equivariant samples to the same subspace.

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

$$\text{circ}(\mathbf{z}) \doteq \begin{bmatrix} z_0 & z_{n-1} & \dots & z_2 & z_1 \\ z_1 & z_0 & z_{n-1} & \cdots & z_2 \\ \vdots & z_1 & z_0 & \ddots & \vdots \\ z_{n-2} & \vdots & \ddots & \ddots & z_{n-1} \\ z_{n-1} & z_{n-2} & \dots & z_1 & z_0 \end{bmatrix} \in \mathbb{R}^{n \times n}. \quad (22)$$

A circular (or cyclic) convolution:

$$\text{circ}(\mathbf{z}) \cdot \mathbf{x} = \mathbf{z} \circledast \mathbf{x}, \quad \text{where} \quad (\mathbf{z} \circledast \mathbf{x})_i = \sum_{j=0}^{n-1} x_j z_{i+n-j \bmod n}. \quad (23)$$

Convolutions from Cyclic Shift Invariance

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

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

The linear operator in the ReduNet:

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

is a circulant matrix and represents a circular convolution:

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

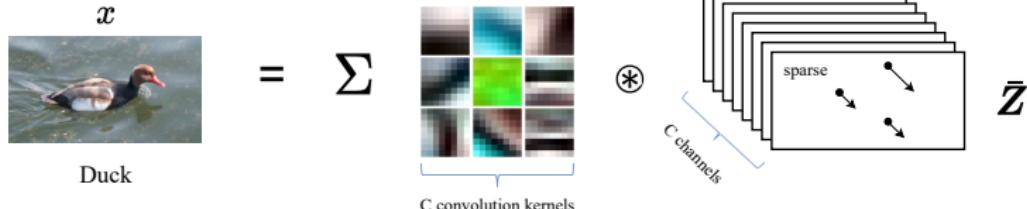
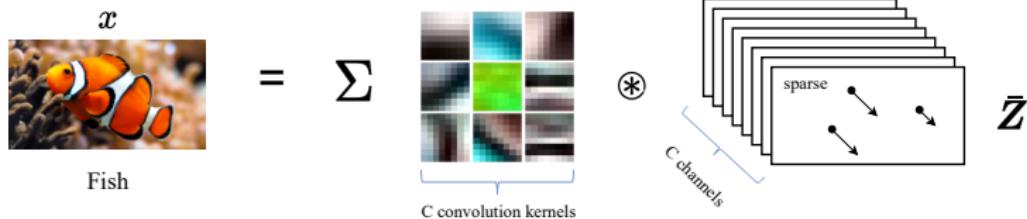
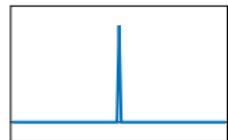
where \mathbf{e} is the first column vector of \mathbf{E} . Similarly, the operators \mathbf{C}^j associated with subsets \mathbf{Z}^j are also circular convolutions.

Tradeoff between Invariance and Separability

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

A necessary assumption: x is **sparsely generated** from incoherent dictionaries for different classes:

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

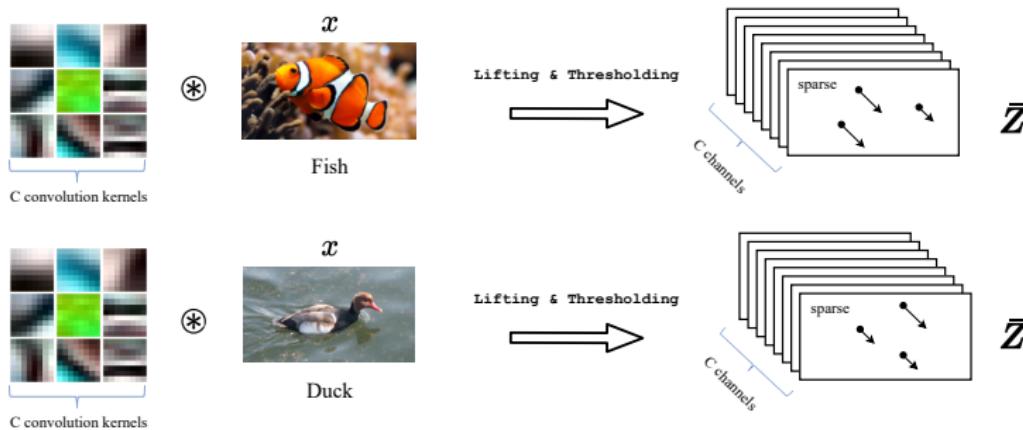


Tradeoff between Invariance and Separability

A basic idea: estimate sparse codes \bar{z} by taking their responses to multiple analysis filters $k_1, \dots, k_C \in \mathbb{R}^n$ [Rubinstein & Elad 2014]:

$$\bar{z} = \tau[k_1 \circledast x, \dots, k_C \circledast x]^* \in \mathbb{R}^{C \times n}. \quad (24)$$

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:

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

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

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

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

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

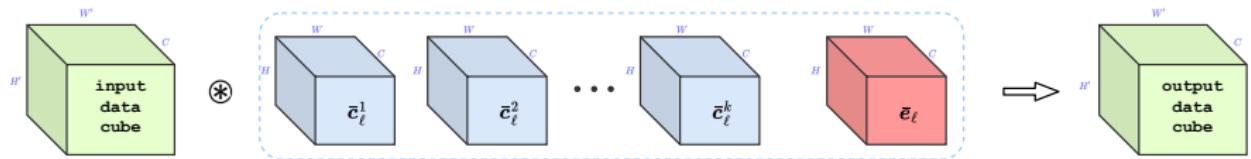


Figure: \bar{E} and \bar{C}^j are automatically multi-channel convolutions!

The Convolution ReduNet versus Scattering Network

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

$$\bar{z}_{\ell+1} \propto \bar{z}_\ell + \eta \cdot \underbrace{\left[\bar{E}_\ell \bar{z}_\ell + \sigma([\bar{C}_\ell^1 \bar{z}_\ell, \dots, \bar{C}_\ell^k \bar{z}_\ell]) \right]}_{g(\bar{z}_\ell, \theta_\ell)}, \quad (25)$$

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

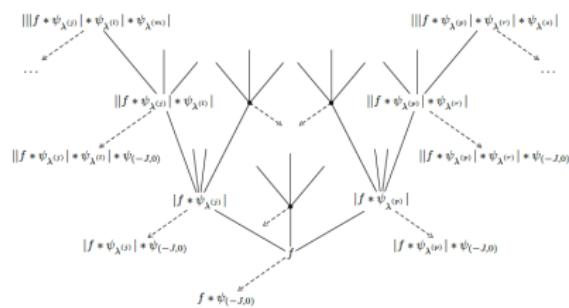
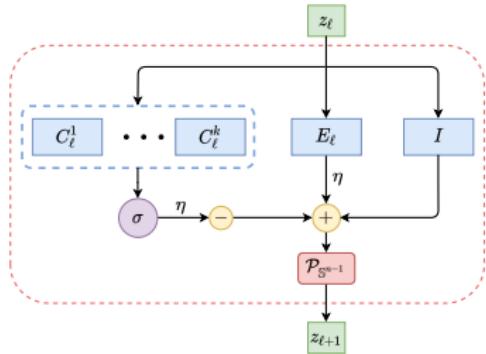


Fig. 2: Scattering network architecture based on wavelet filters and the modulus non-linearity. The elements of the feature vector $\Phi_W(f)$ in (1) are indicated at the tips of the arrows.

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} : $\text{circ}(z) = \mathbf{F}^* \mathbf{D} \mathbf{F}$.

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

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

Computing the inverse is $O(C^3n)$ in the spectral domain, instead of $O(C^3n^3)$! *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)

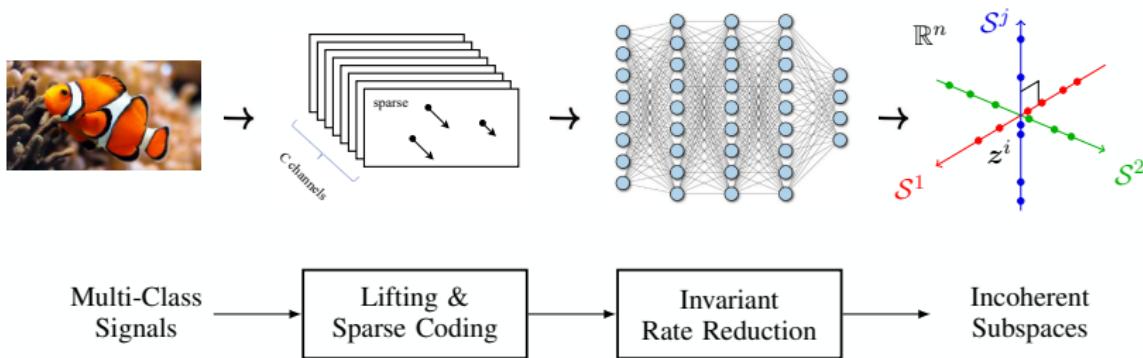
Require: $\bar{\mathbf{Z}} \in \mathbb{C}^{C \times T \times m}$, Π , $\epsilon > 0$, λ , and a learning rate η .

```

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  $\bar{\mathbf{V}}_0 = \{\bar{\mathbf{v}}_0^i(p) \in \mathbb{C}^C\}_{p=0, i=1}^{T-1, m} \doteq \text{DFT}(\bar{\mathbf{Z}}) \in \mathbb{C}^{C \times T \times m}$ .
3: for  $\ell = 1, 2, \dots, L$  do
4:   for  $p = 0, 1, \dots, T - 1$  do
5:     Compute  $\bar{\mathcal{E}}_\ell(p) \in \mathbb{C}^{C \times C}$  and  $\{\bar{\mathcal{C}}_\ell^j(p) \in \mathbb{C}^{C \times C}\}_{j=1}^k$  as
         $\bar{\mathcal{E}}_\ell(p) \doteq \alpha \cdot [\mathbf{I} + \alpha \cdot \bar{\mathbf{V}}_{\ell-1}(p) \cdot \bar{\mathbf{V}}_{\ell-1}(p)^*]^{-1}$ ,
         $\bar{\mathcal{C}}_\ell^j(p) \doteq \alpha_j \cdot [\mathbf{I} + \alpha_j \cdot \bar{\mathbf{V}}_{\ell-1}(p) \cdot \Pi^j \cdot \bar{\mathbf{V}}_{\ell-1}(p)^*]^{-1}$ ;
6:   end for
7:   for  $i = 1, \dots, m$  do
8:     for  $p = 0, 1, \dots, T - 1$  do
9:       Compute  $\{\bar{\mathbf{p}}_\ell^{ij}(p) \doteq \bar{\mathcal{C}}_\ell^j(p) \cdot \bar{\mathbf{v}}_\ell^i(p) \in \mathbb{C}^{C \times 1}\}_{j=1}^k$ ;
10:    end for
11:    Let  $\{\bar{\mathbf{P}}_\ell^{ij} = [\bar{\mathbf{p}}_\ell^{ij}(0), \dots, \bar{\mathbf{p}}_\ell^{ij}(T-1)] \in \mathbb{C}^{C \times T}\}_{j=1}^k$ ;
12:    Compute  $\{\hat{\pi}_\ell^{ij} = \frac{\exp(-\lambda \|\bar{\mathbf{P}}_\ell^{ij}\|_F)}{\sum_{j=1}^k \exp(-\lambda \|\bar{\mathbf{P}}_\ell^{ij}\|_F)}\}_{j=1}^k$ ;
13:    for  $p = 0, 1, \dots, T - 1$  do
14:       $\bar{\mathbf{v}}_\ell^i(p) = \bar{\mathbf{v}}_{\ell-1}^i(p) + \eta \left( \bar{\mathcal{E}}_\ell(p) \bar{\mathbf{v}}_\ell^i(p) - \sum_{j=1}^k \gamma_j \cdot \hat{\pi}_\ell^{ij} \cdot \bar{\mathcal{C}}_\ell^j(p) \cdot \bar{\mathbf{v}}_\ell^i(p) \right)$ ;
15:    end for
16:     $\bar{\mathbf{v}}_\ell^i = \bar{\mathbf{v}}_\ell^i / \|\bar{\mathbf{v}}_\ell^i\|_F$ ;
17:  end for
18:  Set  $\bar{\mathbf{Z}}_\ell = \text{IDFT}(\bar{\mathbf{V}}_\ell)$  as the feature at the  $\ell$ -th layer;
19:   $\frac{1}{2T} \sum_{p=0}^{T-1} \left( \log \det[\mathbf{I} + \alpha \bar{\mathbf{V}}_\ell(p) \cdot \bar{\mathbf{V}}_\ell(p)^*] - \frac{\text{tr}(\Pi^j)}{m} \log \det[\mathbf{I} + \alpha_j \bar{\mathbf{V}}_\ell(p) \cdot \Pi^j \cdot \bar{\mathbf{V}}_\ell(p)^*] \right)$ ;
20: end for
```

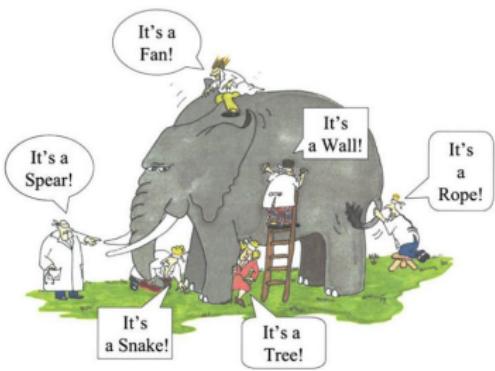
Ensure: features $\bar{\mathbf{Z}}_L$, the learned filters $\{\bar{\mathcal{E}}_\ell(p)\}_{\ell, p}$ and $\{\bar{\mathcal{C}}_\ell^j(p)\}_{j, \ell, p}$.

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.

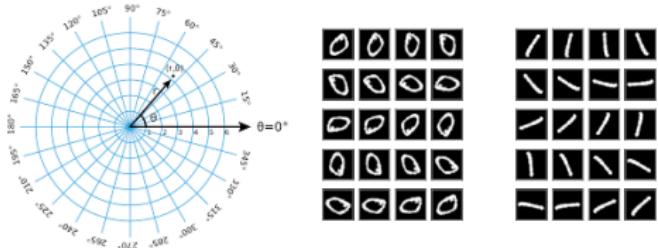


Figure: Left: Multi-channel feature representation of an image in polar coordinates.
Right: Example of training/testing samples.

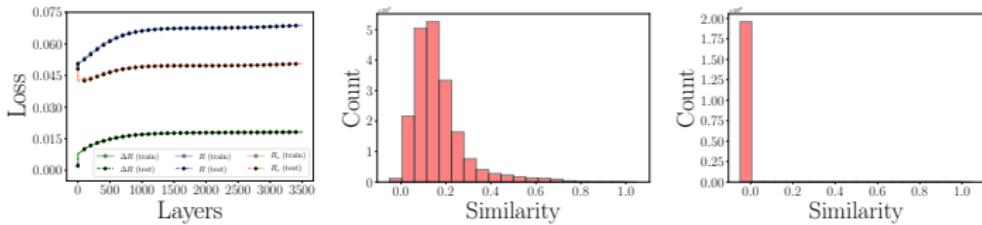


Figure: Left: Rates along the layers; Middle: cross-class cosine similarity among trainings; Right: similarity among testings.

Experiment: 1D Cyclic Shift Invariance of 0 and 1

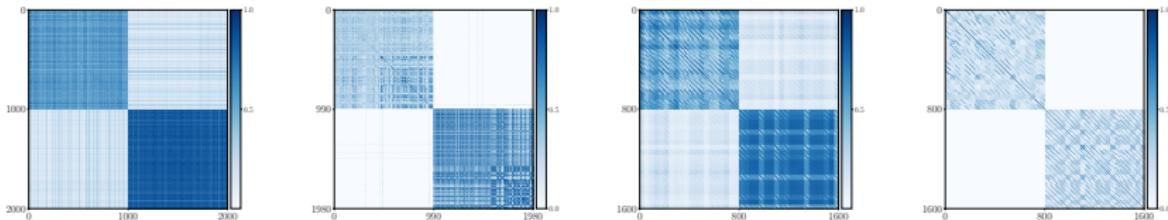


Figure: Left two: heat maps for training and testing. Right two: heat maps for one pair of samples at every possible shift.

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, $C = 20$, $L = \mathbf{40}$ -layers ReduNet.

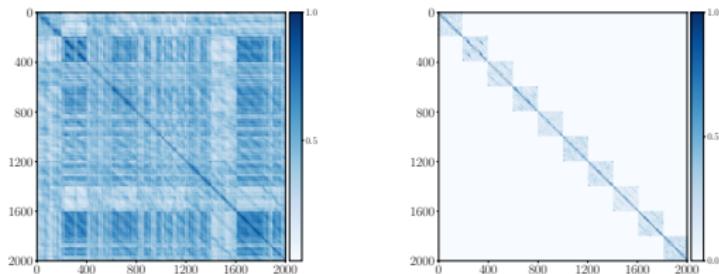


Figure: Heatmaps of cosine similarity among shifted training data X_{shift} (left) and learned features Z_{shift} (right).

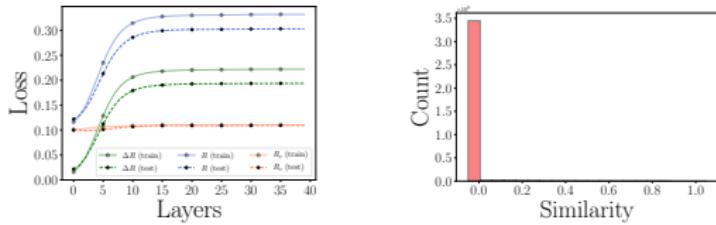


Figure: Left: Rates evolution with iterations; Right: histograms of the cosine similarity (in absolute value) between all pairs of features across different classes.

Experiment: 2D Cyclic Translation Invariance

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

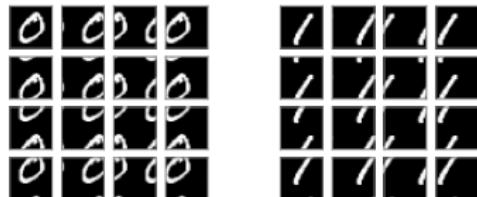
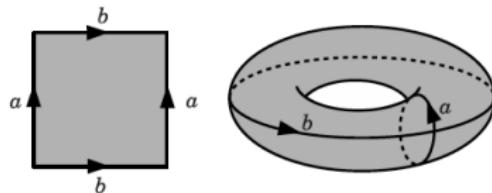


Table: Network performance on digits with **all translations**.

	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, $C = 75$, $L = \mathbf{25}$ -layers ReduNet.

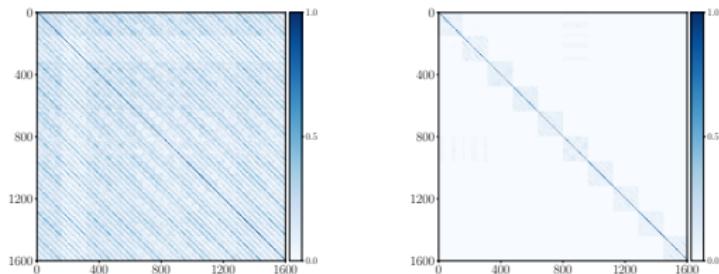


Figure: Heatmaps of cosine similarity among shifted training data X_{shift} (left) and learned features Z_{shift} (right).

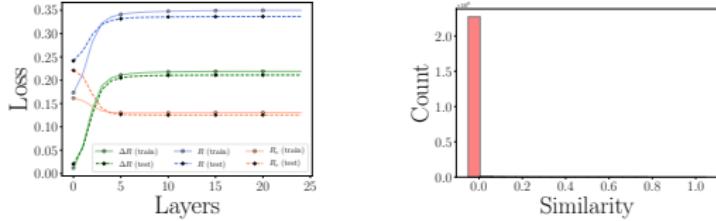


Figure: Left: Rates evolution with iterations; Right: histograms of the cosine similarity (in absolute value) between all pairs of features across different classes.

Experiment: Back Propagation of ReduNet

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.

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: $\dot{\mathbf{Z}} = \eta \cdot \frac{\partial \Delta R}{\partial \mathbf{Z}}$

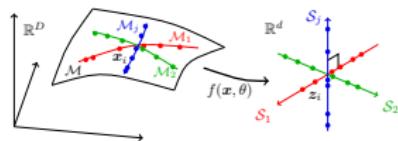
A Unified Framework:

- A principled objective for all settings of learning: **compression**
- 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	Compression ReduNets
Objectives	label fitting	rate reduction
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	unknown	incoherent subspaces

Open Problems: Theory

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

- **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 $\boldsymbol{\Pi}$?

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

Open Problems: Architectures and Algorithms

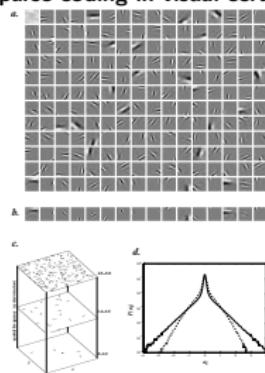
ReduNet: $\bar{z}_{\ell+1} \propto \bar{z}_\ell + \eta \cdot \left[\bar{e}_\ell \circledast \bar{z}_\ell + \sigma([\bar{c}_\ell^1 \circledast \bar{z}_\ell, \dots, \bar{c}_\ell^k \circledast \bar{z}_\ell]) \right] \in \mathbb{S}^{d-1}$.

- New architectures from **accelerated** gradient schemes?
- Conditions for channel-wise **separable and short** convolutions?
- Architectures from invariant rate reduction for **other groups**?
- **Transformer:** $\frac{1}{2} \log \det (\mathbf{I} + \alpha \mathbf{Z} \mathbf{Z}^*) = \frac{1}{2} \log \det \left(\mathbf{I} + \alpha \underbrace{\mathbf{Z}^* \mathbf{Z}}_{\mathbf{z}_i^* \mathbf{U}^* \mathbf{U} \mathbf{z}_j} \right)$?
- Algorithmic architectures (or networks) for optimizing Π ?

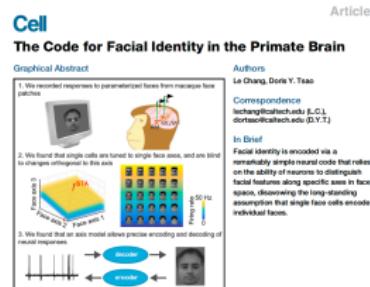
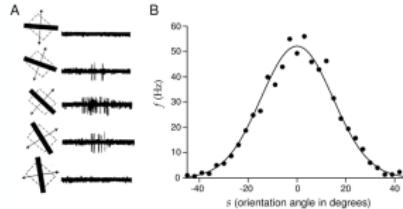
Open Directions: Extensions

- 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**.¹

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.



¹ figures from Bruno Olshausen of Neuroscience Dept., UC Berkeley.

References: Learning via Compression and Rate Reduction

- ① **Clustering** via Lossy Coding and Compression (TPAMI 2007):
<http://people.eecs.berkeley.edu/~yima/psfile/Ma-PAMI07.pdf>
- ② **Classification** via Minimal Incremental Coding Length (NIPS 2007):
http://people.eecs.berkeley.edu/~yima/psfile/MICL_SJIS.pdf
- ③ **Representation** via Maximal Coding Rate Reduction (NeurIPS 2020):
<https://arxiv.org/abs/2006.08558>
- ④ **ReduNet:** A Whitebox Deep Network from Maximizing Rate Reduction:
<https://arxiv.org/abs/2105.10446>
- ⑤ **A New Textbook:** *High-Dim Data Analysis with Low-Dim Models*
<https://book-wright-ma.github.io>

Source Code: Whitebox ReduNet

① Github Link:

<https://github.com/Ma-Lab-Berkeley/ReduNet>

② Google Colab:

https://colab.research.google.com/github/ryanchankh/redunet_demo/blob/master/gaussian3d.ipynb

③ Jupyter Notebook:

https://github.com/ryanchankh/redunet_demo/blob/master/gaussian3d.ipynb

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

Deep (Convolution) Network Architectures are Iterative Optimization Schemes for Compression!

Thank you!
Questions, please?



SIMONS
FOUNDATION