

# *LDMNet*: Low Dimensional Manifold Regularized Neural Networks

Wei Zhu

Duke University

Feb 9, 2018

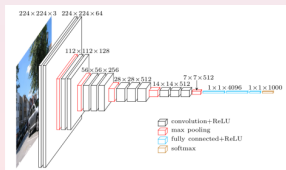
IPAM workshop on deep learning

Joint work with Qiang Qiu, Jiaji Huang, Robert Calderbank, Guillermo Sapiro, and Ingrid Daubechies

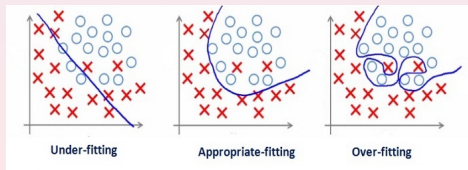
# Neural Networks and Overfitting

Deep neural networks (DNNs) have achieved great success in machine learning research and commercial applications. When large amounts of training data are available, the capacity of DNNs can easily be increased by adding more units or layers to extract more effective high-level features.

However, big networks with millions of parameters can easily overfit even the largest of datasets. As a result, the learned network will have low error rates on the training data, but generalizes poorly onto the test data.



Deep Neural Network



Overfitting

# Regularizations

Many widely-used network regularizations are data-independent.

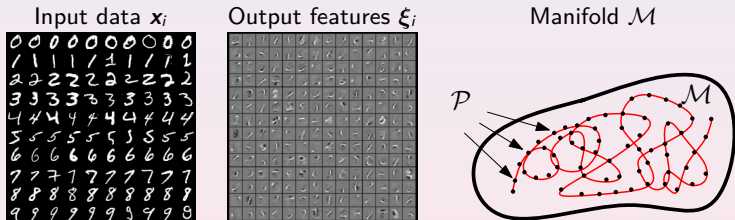
- Weight decay
- Parameter sharing
- DropOut
- DropConnect
- Early stopping
- .....

Most of the data-dependent regularizations are motivated by the empirical observation that data of interest typically lie close to a manifold.

- Tangent distance algorithm
- Tangent prop algorithm
- Manifold tangent classifier
- .....

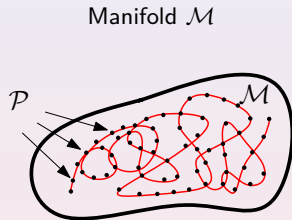
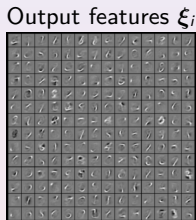
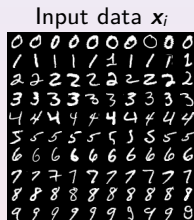
# LDMNet: Low Dimensional Manifold Regularized Neural Networks

Low Dimensional Manifold Regularized Neural Networks (*LDMNet*) incorporates a feature regularization method that focuses on the geometry of both the input data and the output features.



- $\{\mathbf{x}_i\}_{i=1}^N \subset \mathbb{R}^{d_1}$  are the input data.  $\{\xi_i = f_{\theta}(\mathbf{x}_i)\}_{i=1}^N \subset \mathbb{R}^{d_2}$  are the output features.
- $\mathcal{P} = \{(\mathbf{x}_i, \xi_i)\}_{i=1}^N \subset \mathbb{R}^d$  is the collection of the data-feature concatenation  $(\mathbf{x}_i, \xi_i)$ .
- $\mathcal{M} \subset \mathbb{R}^d$  is the underlying manifold, discretely sampled by the point cloud  $\mathcal{P}$ .

## LDMNet: Low Dimensional Manifold Regularized Neural Networks



- The input data  $\{\mathbf{x}_i\}_{i=1}^N$  typically lie close to a collection of low dimensional manifolds, i.e.  $\{\mathbf{x}_i\}_{i=1}^N \subset \mathcal{N} = \cup_{l=1}^L \mathcal{N}_l \subset \mathbb{R}^{d_1}$ .
- The feature extractor,  $f_\theta$ , of a good learning algorithm should be a smooth function over  $\mathcal{N}$ .
- Therefore the concatenation of the input data and output features,  $\mathcal{P} = \{(\mathbf{x}_i, \xi_i)\}_{i=1}^N$ , should sample a collection of low dimensional manifolds  $\mathcal{M} = \cup_{l=1}^L \mathcal{M}_l \subset \mathbb{R}^d$ , where  $d = d_1 + d_2$ , and  $\mathcal{M}_l = \{(\mathbf{x}, f_\theta(\mathbf{x}))\}_{\mathbf{x} \in \mathcal{N}_l}$  is the graph of  $f_\theta$  over  $\mathcal{N}_l$ .

# LDMNet: Low Dimensional Manifold Regularized Neural Networks

We suggest that network overfitting occurs when  $\dim(\mathcal{M}_l)$  is too large after training. Therefore, to reduce overfitting, we explicitly use the dimensions of  $\mathcal{M}_l$  as a regularizer in the following variational form:

$$\begin{aligned} \min_{\theta, \mathcal{M}} \quad & J(\theta) + \frac{\lambda}{|\mathcal{M}|} \int_{\mathcal{M}} \dim(\mathcal{M}(\mathbf{p})) d\mathbf{p} \\ \text{s.t.} \quad & \{(\mathbf{x}_i, f_{\theta}(\mathbf{x}_i))\}_{i=1}^N \subset \mathcal{M}, \end{aligned}$$

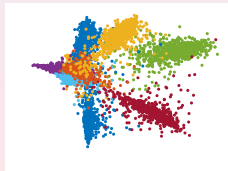
where for any  $\mathbf{p} \in \mathcal{M} = \cup_{l=1}^L \mathcal{M}_l$ ,  $\mathcal{M}(\mathbf{p})$  denotes the manifold  $\mathcal{M}_l$  to which  $\mathbf{p}$  belongs, and  $|\mathcal{M}| = \sum_{l=1}^L |\mathcal{M}_l|$  is the volume of  $\mathcal{M}$ .

# Feature of the MNIST Dataset

10,000 test data



Weight decay

*DropOut**LDMNet*

**Figure:** Test data of MNIST and their features learned by the same network with different regularizers. All networks are trained from the same set of 1,000 images. Data are visualized in two dimensions using PCA, and ten classes are distinguished by different colors.

# Dimension of a Manifold

**Question:** How do we calculate  $\dim(\mathcal{M}(\mathbf{p}))$  from the point cloud

$$\mathcal{P} = \{(\mathbf{x}_i, f_\theta(\mathbf{x}_i))\}_{i=1}^N \subset \mathcal{M}$$

## Theorem

Let  $\mathcal{M}$  be a smooth submanifold isometrically embedded in  $\mathbb{R}^d$ . For any  $\mathbf{p} = (p_i)_{i=1}^d \in \mathcal{M}$ ,

$$\dim(\mathcal{M}) = \sum_{i=1}^d |\nabla_{\mathcal{M}} \alpha_i(\mathbf{p})|^2,$$

where  $\alpha_i(\mathbf{p}) = p_i$  is the coordinate function, and  $\nabla_{\mathcal{M}}$  is the gradient operator on the manifold  $\mathcal{M}$ . More specifically,  $\nabla_{\mathcal{M}} \alpha_i = \sum_{s,t=1}^k \mathbf{g}^{st} \partial_t \alpha_i \partial_s$ , where  $k$  is the intrinsic dimension of  $\mathcal{M}$ , and  $\mathbf{g}^{st}$  is the inverse of the metric tensor.



# Dimension of a Manifold

## Sanity check:

If  $\mathcal{M} = S^1$ , then  $k = \dim(\mathcal{M}) = 1$ ,  $d = \dim(\mathbb{R}^2) = 2$ , and  $\mathbf{x} = \psi(\theta) = (\cos \theta, \sin \theta)^t$  is the coordinate chart.

The metric tensor  $g = g_{11} = \left\langle \frac{\partial \psi}{\partial \theta}, \frac{\partial \psi}{\partial \theta} \right\rangle = 1 = g^{11}$ .

The gradient of  $\alpha_i$ ,  $\nabla_{\mathcal{M}} \alpha_i = g^{11} \partial_1 \alpha_i \partial_1 = \partial_1 \alpha_i \partial_1$  can be viewed as a vector in the ambient space  $\mathbb{R}^2$ :

$$\nabla_{\mathcal{M}}^j \alpha_i = \partial_1 \psi^j \partial_1 \alpha_i$$

Therefore, we have

$$\begin{aligned}\nabla_{\mathcal{M}} \alpha_1 &= \langle \partial_1 \psi^1 \partial_1 \alpha_1, \partial_1 \psi^2 \partial_1 \alpha_1 \rangle = \langle \sin^2 \theta, -\cos \theta \sin \theta \rangle, \\ \nabla_{\mathcal{M}} \alpha_2 &= \langle \partial_1 \psi^1 \partial_1 \alpha_2, \partial_1 \psi^2 \partial_1 \alpha_2 \rangle = \langle -\sin \theta \cos \theta, \cos^2 \theta \rangle.\end{aligned}$$

Hence  $\|\nabla_{\mathcal{M}} \alpha_1\|^2 + \|\nabla_{\mathcal{M}} \alpha_2\|^2 = \sin^2 \theta + \cos^2 \theta = 1$

# LDMNet: Low Dimensional Manifold Regularized Neural Networks

$$\begin{aligned} \min_{\theta, \mathcal{M}} \quad & J(\theta) + \frac{\lambda}{|\mathcal{M}|} \int_{\mathcal{M}} \dim(\mathcal{M}(\mathbf{p})) d\mathbf{p} \\ \text{s.t.} \quad & \{(\mathbf{x}_i, f_{\theta}(\mathbf{x}_i))\}_{i=1}^N \subset \mathcal{M}, \end{aligned} \quad (1)$$

Using Theorem 1, the above variational functional can be reformulated as

$$\begin{aligned} \min_{\theta, \mathcal{M}} \quad & J(\theta) + \frac{\lambda}{|\mathcal{M}|} \sum_{j=1}^d \|\nabla_{\mathcal{M}} \alpha_j\|_{L^2(\mathcal{M})}^2 \\ \text{s.t.} \quad & \{(\mathbf{x}_i, f_{\theta}(\mathbf{x}_i))\}_{i=1}^N \subset \mathcal{M} \end{aligned} \quad (2)$$

where  $\sum_{j=1}^d \|\nabla_{\mathcal{M}} \alpha_j\|_{L^2(\mathcal{M})}^2$  corresponds to the  $L^1$  norm of the local dimension.

How do we solve (2)? Alternate direction of minimization with respect to  $\mathcal{M}$  and  $\theta$ .

# Alternate Direction of Minimization

$$\begin{aligned} \min_{\theta, \mathcal{M}} \quad & J(\theta) + \frac{\lambda}{|\mathcal{M}|} \sum_{j=1}^d \|\nabla_{\mathcal{M}} \alpha_j\|_{L^2(\mathcal{M})}^2 \\ \text{s.t.} \quad & \{(\mathbf{x}_i, f_{\theta}(\mathbf{x}_i))\}_{i=1}^N \subset \mathcal{M} \end{aligned}$$

Given  $(\theta^{(k)}, \mathcal{M}^{(k)})$  at step  $k$  satisfying  $\{(\mathbf{x}_i, f_{\theta^{(k)}}(\mathbf{x}_i))\}_{i=1}^N \subset \mathcal{M}^{(k)}$ , step  $k+1$  consists of the following

- Update  $\theta^{(k+1)}$  and the perturbed coordinate functions  $\alpha^{(k+1)} = (\alpha_1^{(k+1)}, \dots, \alpha_d^{(k+1)})$  as the minimizers of (3) with the fixed  $\mathcal{M}^{(k)}$ :

$$\begin{aligned} \min_{\theta, \alpha} \quad & J(\theta) + \frac{\lambda}{|\mathcal{M}^{(k)}|} \sum_{j=1}^d \|\nabla_{\mathcal{M}^{(k)}} \alpha_j\|_{L^2(\mathcal{M}^{(k)})}^2 \\ \text{s.t.} \quad & \alpha(\mathbf{x}_i, f_{\theta^{(k)}}(\mathbf{x}_i)) = (\mathbf{x}_i, f_{\theta}(\mathbf{x}_i)), \quad \forall i = 1, \dots, N \end{aligned} \tag{3}$$

- Update  $\mathcal{M}^{(k+1)}$ :

$$\mathcal{M}^{(k+1)} = \alpha^{(k+1)}(\mathcal{M}^{(k)}) \tag{4}$$

# Alternate Direction of Minimization

- Update  $\theta^{(k+1)}$  and  $\alpha^{(k+1)} = (\alpha_1^{(k+1)}, \dots, \alpha_d^{(k+1)})$  with the fixed  $\mathcal{M}^{(k)}$ :

$$\begin{aligned} \min_{\theta, \alpha} \quad & J(\theta) + \frac{\lambda}{|\mathcal{M}^{(k)}|} \sum_{j=1}^d \|\nabla_{\mathcal{M}^{(k)}} \alpha_j\|_{L^2(\mathcal{M}^{(k)})}^2 \\ \text{s.t.} \quad & \alpha(\mathbf{x}_i, f_{\theta^{(k)}}(\mathbf{x}_i)) = (\mathbf{x}_i, f_{\theta}(\mathbf{x}_i)), \quad \forall i = 1, \dots, N \end{aligned}$$

- Update  $\mathcal{M}^{(k+1)}$ :

$$\mathcal{M}^{(k+1)} = \alpha^{(k+1)}(\mathcal{M}^{(k)})$$

The manifold update is trivial to implement, and the update of  $\theta$  and  $\alpha$  is an optimization problem with nonlinear constraint, which “can” be solved via the alternating direction method of multipliers (ADMM).

# Alternate Direction of Minimization

$$\begin{aligned} \min_{\theta, \alpha} \quad & J(\theta) + \frac{\lambda}{|\mathcal{M}^{(k)}|} \sum_{j=1}^d \|\nabla_{\mathcal{M}^{(k)}} \alpha_j\|_{L^2(\mathcal{M}^{(k)})}^2 \\ \text{s.t.} \quad & \alpha(\mathbf{x}_i, f_{\theta^{(k)}}(\mathbf{x}_i)) = (f_{\theta^{(k)}}(\mathbf{x}_i) - Z_i^{(k)}), \quad \forall i = 1, \dots, N \end{aligned}$$

Solve the variational problem using ADMM. More specifically,

$$\begin{aligned} \alpha_{\xi}^{(k+1)} = \arg \min_{\alpha_{\xi}} \quad & \sum_{j=d_1+1}^d \|\nabla_{\mathcal{M}^{(k)}} \alpha_j\|_{L^2(\mathcal{M}^{(k)})}^2 \\ & + \frac{\mu |\mathcal{M}^{(k)}|}{2\lambda N} \sum_{i=1}^N \|\alpha_{\xi}(\mathbf{x}_i, f_{\theta^{(k)}}(\mathbf{x}_i)) - (f_{\theta^{(k)}}(\mathbf{x}_i) - Z_i^{(k)})\|_2^2. \\ \theta^{(k+1)} = \arg \min_{\theta} \quad & J(\theta) + \frac{\mu}{2N} \sum_{i=1}^N \|\alpha_{\xi}^{(k+1)}(\mathbf{x}_i, f_{\theta^{(k)}}(\mathbf{x}_i)) - (f_{\theta^{(k)}}(\mathbf{x}_i) - Z_i^{(k)})\|_2^2. \\ Z_i^{(k+1)} = & Z_i^{(k)} + \alpha_{\xi}^{(k+1)}(\mathbf{x}_i, f_{\theta^{(k)}}(\mathbf{x}_i)) - f_{\theta^{(k)}}(\mathbf{x}_i), \end{aligned}$$

where  $\alpha = (\alpha_x, \alpha_{\xi}) = ((\alpha_1, \dots, \alpha_{d_1}), (\alpha_{d_1+1}, \dots, \alpha_d))$ , and  $Z_i$  is the dual variable.

## ADMM

$$\begin{aligned} \alpha_{\xi}^{(k+1)} = \arg \min_{\alpha_{\xi}} & \sum_{j=d_1+1}^d \|\nabla_{\mathcal{M}^{(k)}} \alpha_j\|_{L^2(\mathcal{M}^{(k)})} \\ & + \frac{\mu |\mathcal{M}^{(k)}|}{2\lambda N} \sum_{i=1}^N \|\alpha_{\xi}(\mathbf{x}_i, f_{\theta^{(k)}}(\mathbf{x}_i)) - (f_{\theta^{(k)}}(\mathbf{x}_i) - Z_i^{(k)})\|_2^2. \end{aligned} \quad (5)$$

$$\theta^{(k+1)} = \arg \min_{\theta} J(\theta) + \frac{\mu}{2N} \sum_{i=1}^N \|\alpha_{\xi}^{(k+1)}(\mathbf{x}_i, f_{\theta^{(k)}}(\mathbf{x}_i)) - (f_{\theta}(\mathbf{x}_i) - Z_i^{(k)})\|_2^2. \quad (6)$$

$$Z_i^{(k+1)} = Z_i^{(k)} + \alpha_{\xi}^{(k+1)}(\mathbf{x}_i, f_{\theta^{(k)}}(\mathbf{x}_i)) - f_{\theta^{(k+1)}}(\mathbf{x}_i), \quad (7)$$

Among (5),(6) and (7), (7) is the easiest to implement, (6) can be solved by stochastic gradient descent (SGD) with modified back propagation, and (5) can be solved by the point integral method (PIM) [Z. Shi, J. Sun].

Shi, Z., and Sun, J. (2013). Convergence of the point integral method for the poisson equation on manifolds ii: the dirichlet boundary. arXiv preprint arXiv:1312.4424

# Back Propagation for the $\theta$ Update

$$\begin{aligned}\theta^{(k+1)} &= \arg \min_{\theta} J(\theta) + \frac{\mu}{2N} \sum_{i=1}^N \|\alpha_{\xi}^{(k+1)}(\mathbf{x}_i, f_{\theta^{(k)}}(\mathbf{x}_i)) - (f_{\theta}(\mathbf{x}_i) - Z_i^{(k)})\|_2^2. \\ &= \arg \min_{\theta} \frac{1}{N} \sum_{i=1}^N \ell(f_{\theta}(\mathbf{x}_i), y_i) + \frac{1}{N} \sum_{i=1}^N E_i(\theta),\end{aligned}\quad (8)$$

where  $E_i(\theta) = \frac{\mu}{2} \|\alpha_{\xi}^{(k+1)}(\mathbf{x}_i, f_{\theta^{(k)}}(\mathbf{x}_i)) - (f_{\theta}(\mathbf{x}_i) - Z_i^{(k)})\|_2^2$ . The gradient of the second term with respect to the output layer  $f_{\theta}(\mathbf{x}_i)$  is:

$$\frac{\partial E_i}{\partial f_{\theta}(\mathbf{x}_i)} = \mu \left( f_{\theta}(\mathbf{x}_i) - Z_i^{(k)} - \alpha_{\xi}^{(k+1)}(\mathbf{x}_i, f_{\theta^{(k)}}(\mathbf{x}_i)) \right) \quad (9)$$

This means that we need to only add the extra term (9) to the original gradient, and then use the already known procedure to back-propagate the gradient.

# Point Integral Method for the $\alpha$ Update

$$\alpha_{\xi}^{(k+1)} = \arg \min_{\alpha_{\xi}} \sum_{j=d_1+1}^d \|\nabla_{\mathcal{M}^{(k)}} \alpha_j\|_{L^2(\mathcal{M}^{(k)})} + \frac{\mu |\mathcal{M}^{(k)}|}{2\lambda N} \sum_{i=1}^N \|\alpha_{\xi}(\mathbf{x}_i, f_{\theta^{(k)}}(\mathbf{x}_i)) - (f_{\theta^{(k)}}(\mathbf{x}_i) - Z_i^{(k)})\|_2^2.$$

Note that the objective function above is decoupled with respect to  $j$ , and each  $\alpha_j$  update can be cast into:

$$\min_{u \in H^1(\mathcal{M})} \|\nabla_{\mathcal{M}} u\|_{L^2(\mathcal{M})}^2 + \gamma \sum_{\mathbf{q} \in P} |u(\mathbf{q}) - v(\mathbf{q})|^2, \quad (10)$$

where  $u = \alpha_j$ ,  $\mathcal{M} = \mathcal{M}^{(k)}$ ,  $P = \{\mathbf{p}_i = (\mathbf{x}_i, f_{\theta^{(k)}}(\mathbf{x}_i))\}_{i=1}^N \subset \mathcal{M}$ , and  $\gamma = \mu |\mathcal{M}^{(k)}| / 2\lambda N$ . The Euler-Lagrange equation of (10) is:

$$\begin{cases} -\Delta_{\mathcal{M}} u(\mathbf{p}) + \gamma \sum_{\mathbf{q} \in P} \delta(\mathbf{p} - \mathbf{q})(u(\mathbf{q}) - v(\mathbf{q})) = 0, & \mathbf{p} \in \mathcal{M} \\ \frac{\partial u}{\partial n} = 0, & \mathbf{p} \in \partial \mathcal{M} \end{cases}$$



# Point Integral Method

$$\begin{cases} -\Delta_{\mathcal{M}}u(\mathbf{p}) + \gamma \sum_{\mathbf{q} \in P} \delta(\mathbf{p} - \mathbf{q})(u(\mathbf{q}) - v(\mathbf{q})) = 0, & \mathbf{p} \in \mathcal{M} \\ \frac{\partial u}{\partial n} = 0, & \mathbf{p} \in \partial\mathcal{M} \end{cases}$$

In the point integral method (PIM) [Z. Shi, J. Sun], the key observation is the following integral approximation:

$$\begin{aligned} \int_{\mathcal{M}} \Delta_{\mathcal{M}}u(\mathbf{y}) \bar{R} \left( \frac{\|\mathbf{x} - \mathbf{y}\|^2}{4t} \right) d\mathbf{y} &\approx -\frac{1}{t} \int_{\mathcal{M}} (u(\mathbf{x}) - u(\mathbf{y})) R \left( \frac{\|\mathbf{x} - \mathbf{y}\|^2}{4t} \right) d\mathbf{y} \\ &\quad + 2 \int_{\partial\mathcal{M}} \frac{\partial u}{\partial n}(\mathbf{y}) \bar{R} \left( \frac{\|\mathbf{x} - \mathbf{y}\|^2}{4t} \right) d\tau_{\mathbf{y}}. \end{aligned}$$

The function  $R$  is a positive function defined on  $[0, +\infty)$  with compact support (or fast decay) and

$$\bar{R} = \int_r^\infty R(s) ds.$$

# Local Truncation Error

## Theorem

Let  $\mathcal{M}$  be a smooth manifold and  $u \in C^3(\mathcal{M})$ , then

$$\left\| -\frac{1}{t} \int_{\mathcal{M}} (u(\mathbf{x}) - u(\mathbf{y})) R_t(\mathbf{x}, \mathbf{y}) d\mathbf{y} + 2 \int_{\partial\mathcal{M}} \frac{\partial u}{\partial n}(\mathbf{y}) \bar{R}_t(\mathbf{x}, \mathbf{y}) d\tau_{\mathbf{y}} - \int_{\mathcal{M}} \Delta_{\mathcal{M}} u(\mathbf{y}) \bar{R}_t(\mathbf{x}, \mathbf{y}) d\mathbf{y} \right\|_{L^2(\mathcal{M})} = O(t^{1/4}),$$

where

$$R_t(\mathbf{x}, \mathbf{y}) = \frac{1}{(4\pi t)^{k/2}} R\left(\frac{\|\mathbf{x} - \mathbf{y}\|^2}{4t}\right), \bar{R}_t(\mathbf{x}, \mathbf{y}) = \frac{1}{(4\pi t)^{k/2}} \bar{R}\left(\frac{\|\mathbf{x} - \mathbf{y}\|^2}{4t}\right).$$

# Point Integral Method for the $\alpha$ Update

$$\begin{cases} -\Delta_{\mathcal{M}} u(\mathbf{p}) + \gamma \sum_{\mathbf{q} \in P} \delta(\mathbf{p} - \mathbf{q})(u(\mathbf{q}) - v(\mathbf{q})) = 0, & \mathbf{p} \in \mathcal{M} \\ \frac{\partial u}{\partial n} = 0, & \mathbf{p} \in \partial \mathcal{M} \end{cases}$$

After convolving the above equation with the heat kernel

$R_t(\mathbf{p}, \mathbf{q}) = C_t \exp\left(-\frac{|\mathbf{p}-\mathbf{q}|^2}{4t}\right)$ , we know the solution  $u$  should satisfy

$$-\int_{\mathcal{M}} \Delta_{\mathcal{M}} u(\mathbf{q}) R_t(\mathbf{p}, \mathbf{q}) d\mathbf{q} + \gamma \sum_{\mathbf{q} \in P} R_t(\mathbf{p}, \mathbf{q}) (u(\mathbf{q}) - v(\mathbf{q})) = 0. \quad (11)$$

Combined with Theorem 2 and the Neumann boundary condition, this implies that  $u$  should approximately satisfy

$$\int_{\mathcal{M}} (u(\mathbf{p}) - u(\mathbf{q})) R_t(\mathbf{p}, \mathbf{q}) d\mathbf{q} + \gamma t \sum_{\mathbf{q} \in P} R_t(\mathbf{p}, \mathbf{q}) (u(\mathbf{q}) - v(\mathbf{q})) = 0 \quad (12)$$

# Point Integral Method for the $\alpha$ Update

$$\int_{\mathcal{M}} (u(\mathbf{p}) - u(\mathbf{q})) R_t(\mathbf{p}, \mathbf{q}) d\mathbf{q} + \gamma t \sum_{\mathbf{q} \in P} R_t(\mathbf{p}, \mathbf{q}) (u(\mathbf{q}) - v(\mathbf{q})) = 0$$

Assume that  $P = \{\mathbf{p}_1, \dots, \mathbf{p}_N\}$  samples the manifold  $\mathcal{M}$  uniformly at random, then the above equation can be discretized as

$$\frac{|\mathcal{M}|}{N} \sum_{j=1}^N R_{t,ij} (u_i - u_j) + \gamma t \sum_{j=1}^N R_{t,ij} (u_j - v_j) = 0, \quad (13)$$

where  $u_i = u(\mathbf{p}_i)$ , and  $R_{t,ij} = R_t(\mathbf{p}_i, \mathbf{p}_j)$ . Combining the definition of  $\gamma$  in (10), we can write (13) in the matrix form

$$\left( \mathbf{L} + \frac{\mu}{\tilde{\lambda}} \mathbf{W} \right) \mathbf{u} = \frac{\mu}{\tilde{\lambda}} \mathbf{W} \mathbf{v}, \quad \tilde{\lambda} = 2\lambda/t, \quad (14)$$

where  $\tilde{\lambda}$  can be chosen instead of  $\lambda$  as the hyperparameter to be tuned,  $\mathbf{u} = (u_1, \dots, u_N)^T$ ,  $\mathbf{W}$  is an  $N \times N$  matrix

$$\mathbf{W}_{ij} = R_{t,ij} = \exp \left( -\frac{|\mathbf{p}_i - \mathbf{p}_j|^2}{4t} \right), \quad (15)$$

and  $\mathbf{L}$  is the graph Laplacian of  $\mathbf{W}$ :

# Algorithm for Training *LDMNet*

---

## Algorithm 1 LDMNet Training

---

**Input:** Training data  $\{(\mathbf{x}_i, y_i)\}_{i=1}^N \subset \mathbb{R}^{d_1} \times \mathbb{R}$ , hyperparameters  $\tilde{\lambda}$  and  $\mu$ , and a neural network with the weights  $\theta$  and the output layer  $\xi_i = f_\theta(\mathbf{x}_i) \in \mathbb{R}^{d_2}$ .

**Output:** Trained network weights  $\theta^*$ .

- 1: Randomly initialize the network weights  $\theta^{(0)}$ . The dual variables  $Z_i^{(0)} \in \mathbb{R}^{d_2}$  are initialized to zero.
- 2: **while** not converge **do**
- 3:   1. Compute the matrices  $\mathbf{W}$  and  $\mathbf{L}$  as in (15) and (16) with  $\mathbf{p}_i = (\mathbf{x}_i, f_{\theta^{(k)}}(\mathbf{x}_i))$ .
- 4:   2. Update  $\alpha^{(k+1)}$  in (5): solve the linear systems (14), where
$$\mathbf{u}_i = \alpha_j(\mathbf{p}_i), \quad \mathbf{v}_i = f_{\theta^{(k)}}(\mathbf{x}_i)_j - Z_{i,j}^{(k)}.$$
- 5:   3. Update  $\theta^{(k+1)}$  in (6): run SGD for  $M$  epochs with an extra gradient term (9).
- 6:   4. Update  $Z^{(k+1)}$  in (7).
- 7:   5.  $k \leftarrow k + 1$ .
- 8: **end while**
- 9:  $\theta^* \leftarrow \theta^{(k)}$ .

# Complexity Analysis

The additional computation in Algorithm 1 comes from the update of weight matrices and solving the linear system from PIM once every  $M$  epochs of SGD.

The weight matrix  $\mathbf{W}$  is truncated to only 20 nearest neighbors. To identify those nearest neighbors, we first organize the data points  $\{\mathbf{p}_1, \dots, \mathbf{p}_N\} \subset \mathbb{R}^d$  into a  $k$ -d tree. Nearest neighbors can then be efficiently identified because branches can be eliminated from the search space quickly. Modern algorithms to build a balanced  $k$ -d tree generally at worst converge in  $O(N \log N)$  time, and finding nearest neighbours for one query point in a balanced  $k$ -d tree takes  $O(\log N)$  time on average. Therefore the complexity of the weight update is  $O(N \log N)$ .

Since  $\mathbf{W}$  and  $\mathbf{L}$  are sparse symmetric matrices with a fixed maximum number of non-zero entries in each row, the linear system can be solved efficiently with the preconditioned conjugate gradients method. After restricting the number of matrix multiplications to a maximum of 50, the complexity of the  $\alpha$  update is  $O(N)$ .

## Experimental Setup

We compare the performance of *LDMNet* to widely-used network regularization techniques, weight decay and *DropOut*, using the same underlying network structure.

Unless otherwise stated, all experiments use mini-batch SGD with momentum on batches of 100 images. The momentum parameter is fixed at 0.9. The networks are trained using a fixed learning rate  $r_0$  on the first 200 epochs, and then  $r_0/10$  for another 100 epochs.

For *LDMNet*, the weight matrices and  $\alpha$  are updated once every  $M = 2$  epochs of SGD. For *DropOut*, the corresponding *DropOut* layer is always chosen to have a drop rate of 0.5. All other network hyperparameters are chosen according to the error rate on the validation set.

## MNIST

The MNIST handwritten digit dataset contains approximately 60,000 training images ( $28 \times 28$ ) and 10,000 test images.

Layer	Type	Parameters
1	conv	size: $5 \times 5 \times 1 \times 20$ stride: 1, pad: 0
2	max pool	size: $2 \times 2$ , stride: 2, pad: 0
3	conv	size: $5 \times 5 \times 20 \times 50$ stride: 1, pad: 0
4	max pool	size: $2 \times 2$ , stride: 2, pad: 0
5	conv	size: $4 \times 4 \times 50 \times 500$ stride: 1, pad: 0
6	ReLU ( <i>DropOut</i> )	N/A
7	fully connected	$500 \times 10$
8	softmaxloss	N/A

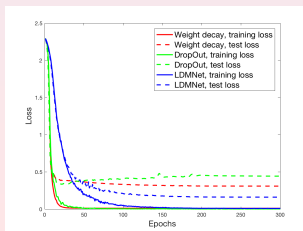
**Table:** Network structure in the MNIST experiments. The outputs of layer 6 are the extracted features, which will be fed into the softmax classifier (layer 7 and 8).



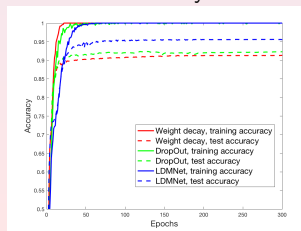
# MNIST

training per class	weight decay	<i>DropOut</i>	<i>LDMNet</i>
50	91.32%	92.31%	<b>95.57%</b>
100	93.38%	94.05%	<b>96.73%</b>
400	97.23%	97.95%	<b>98.41%</b>
700	97.67%	98.07%	<b>98.61%</b>
1000	98.06%	98.71%	<b>98.89%</b>
3000	98.87%	99.21%	<b>99.24%</b>
6000	99.15%	<b>99.41%</b>	99.39%

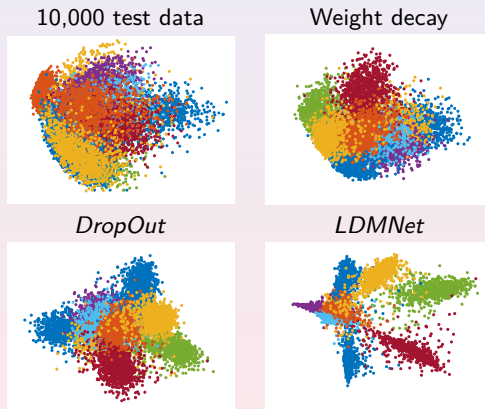
Loss



Accuracy



# MNIST



**Figure:** Test data of MNIST and their features learned by the same network with different regularizers. All networks are trained from the same set of 1,000 images. Data are visualized in two dimensions using PCA, and ten classes are distinguished by different colors.

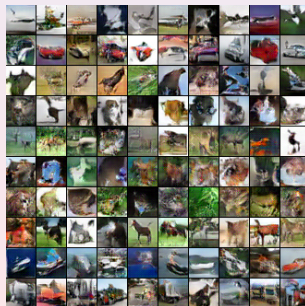
# SVHN and CIFAR-10

SVHN and CIFAR-10 are benchmark RGB image datasets, both of which contain 10 different classes.

SVHN



CIFAR-10



## SVHN and CIFAR 10

Layer	Type	Parameters
1	conv	size: $5 \times 5 \times 3 \times 96$ stride: 1, pad: 2
2	ReLU	N/A
3	max pool	size: $3 \times 3$ , stride: 2, pad: 0
4	conv	size: $5 \times 5 \times 96 \times 128$ stride: 1, pad: 2
5	ReLU	N/A
6	max pool	size: $3 \times 3$ , stride: 2, pad: 0
7	conv	size: $4 \times 4 \times 128 \times 256$ stride: 1, pad: 0
8	ReLU	N/A
9	max pool	size: $3 \times 3$ , stride: 2, pad: 0
10	fully connected	output: 2048
11	ReLU ( <i>DropOut</i> )	N/A
12	fully connected	output: 2048
13	ReLU ( <i>DropOut</i> )	N/A
14	fully connected	$2048 \times 10$
15	softmaxloss	N/A

## SVHN and CIFAR 10

training per class	weight decay	<i>DropOut</i>	<i>LDMNet</i>
50	71.46%	71.94%	<b>74.64%</b>
100	79.05%	79.94%	<b>81.36%</b>
400	87.38%	87.16%	<b>88.03%</b>
700	89.69%	89.83%	<b>90.07%</b>

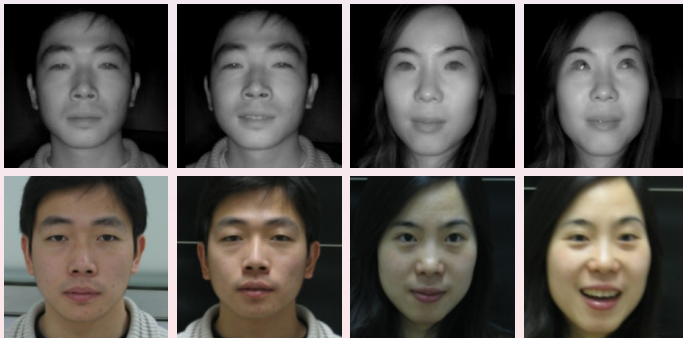
Table: SVHN: testing accuracy for different regularizers

training per class	weight decay	<i>DropOut</i>	<i>LDMNet</i>
50	34.70%	35.94%	<b>41.55%</b>
100	42.45%	43.18%	<b>48.73%</b>
400	56.19%	56.79%	<b>60.08%</b>
700	61.84%	62.59%	<b>65.59%</b>
full data	87.72%		<b>88.21%</b>

Table: CIFAR-10: testing accuracy for different regularizers.

# NIR-VIS Heterogeneous Face Recognition

The objective of the experiment is to match a probe image of a subject captured in the near-infrared spectrum (NIR) to the same subject from a gallery of visible spectrum (VIS) images. The CASIA NIR-VIS 2.0 benchmark dataset is used to evaluate the performance.



**Figure:** Sample images of two subjects from the CASIA NIR-VIS 2.0 dataset after the pre-processing of alignment and cropping. Top: NIR. Bottom: VIS.

# NIR-VIS Heterogeneous Face Recognition

Despite recent breakthroughs for VIS face recognition by training DNNs from millions of VIS images, such approach cannot be simply transferred to NIR-VIS face recognition.

- Unlike VIS face images, we have only limited number of available NIR images.
- The NIR-VIS face matching is a cross-modality comparison.

The authors in [Lezama et al., 2017] introduced a way to transfer the breakthrough in VIS face recognition to the NIR spectrum. Their idea is to use a DNN pre-trained on VIS images as a feature extractor, while making two independent modifications in the input and output of the DNN.

- Modify the input by “hallucinating” a VIS image from the NIR sample.
- Learn an embedding of the DNN features at the output

Lezama, J., Qiu, Q., and Sapiro, G. (2017). Not afraid of the dark: Nir-vis face recognition via cross-spectral hallucination and low-rank embedding. In The IEEE Conference on Computer Vision and Pattern Recognition (CVPR).

# NIR-VIS Heterogeneous Face Recognition

- Modify the input by “hallucinating” a VIS image from the NIR sample.
- Learn an embedding of the DNN features at the output

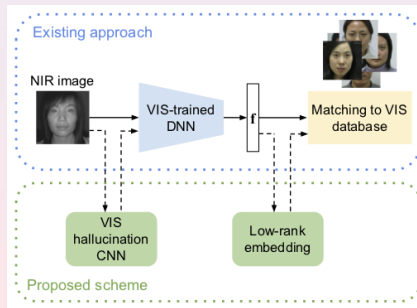


Figure: Proposed procedure in [Lezama et al., 2017]



# NIR-VIS Heterogeneous Face Recognition

We follow the second idea in [Lezama et al., 2017], and learn a nonlinear low dimensional manifold embedding of the output features. we use the VGG-face model as a feature extractor. We then put the 4,096 dimensional features into a two-layer fully connected network to learn a nonlinear embedding using different regularizations.

Layer	Type	Parameters
1	fully connected	output:2000
2	ReLu ( <i>DropOut</i> )	N/A
3	fully connected	output:2000
4	ReLu ( <i>DropOut</i> )	N/A

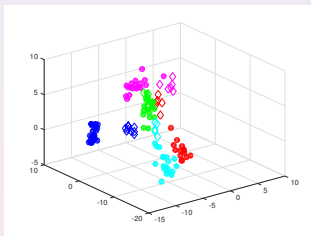
**Table:** Fully connected network for the NIR-VIS nonlinear feature embedding. The outputs of layer 4 are the extracted features.

# NIR-VIS Heterogeneous Face Recognition

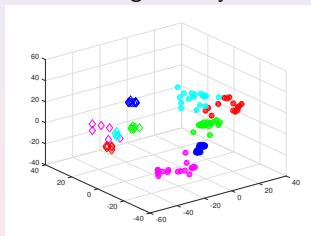
	Accuracy (%)
VGG-face	$74.51 \pm 1.28$
VGG-face + triplet [Lezama et al., 2017]	$75.96 \pm 2.90$
VGG-face + low-rank [Lezama et al., 2017]	$80.69 \pm 1.02$
VGG-face weight Decay	$63.87 \pm 1.33$
VGG-face <i>DropOut</i>	$66.97 \pm 1.31$
VGG-face <i>LDMNet</i>	<b><math>85.02 \pm 0.86</math></b>

# NIR-VIS Heterogeneous Face Recognition

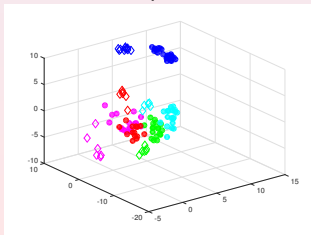
VGG-face



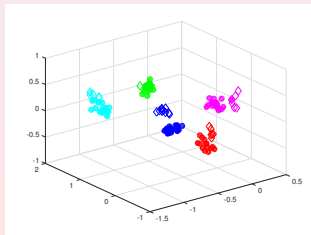
Weight decay



DropOut



LDMNet



# Conclusion

- *LDMNet* is a general network regularization technique that focuses on the geometry of both the input data and the output feature
- *LDMNet* directly minimize the dimension of the manifolds without explicit parametrization.
- *LDMNet* significantly outperforms the widely-used network regularizers.
- Limits: data augmentation,  $O(N \log N)$  complexity.

Thank you