



Northeastern University

Solving Interpretable Kernel Dimension Reduction

Chieh Wu, Jared Miller, Yale Chang, Mario Szaier, Jennifer G. Dy
Dept. of Electrical and Computer Engineering, Northeastern University

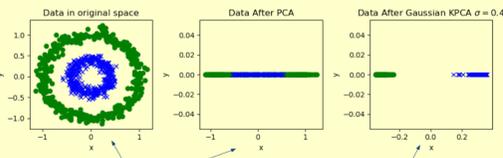
What is IKDR?

Principal Component Analysis (PCA) is the most commonly used Dimension Reduction (DR) technique. It is also an **interpretable** way to reduce the dimension.

We know exactly how the new features relate to the original features.

$$\begin{bmatrix} w_{11} & w_{12} & w_{13} \\ w_{21} & w_{22} & w_{23} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} x_1 w_{11} + x_2 w_{12} + x_3 w_{13} \\ x_1 w_{21} + x_2 w_{22} + x_3 w_{23} \end{bmatrix}$$

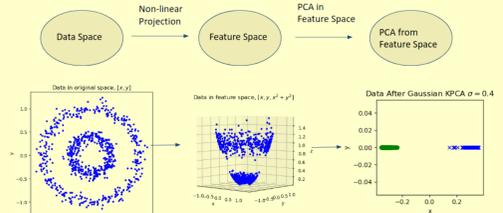
But PCA cannot capture **nonlinear Relationships**.



After Dimension Reduction, you can't tell that Blue and Green are actually separated. Ideally, after Dimension Reduction, samples of the same group should stay close together.

Note: This requires us to also capture non-linear relationships!!!

KPCA captures nonlinear Relationships but not interpretable.



KPCA is very powerful, but

Problem 1: It does not use labels to guide the dimension reduction.

Problem 2: Since KPCA is PCA in the feature space, it's not obvious what they mean.

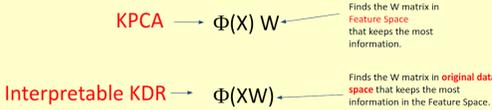
Here is the Gaussian Kernel feature map:

$$\phi(x) = e^{-x^2/2\sigma^2} \left[1, \sqrt{\frac{1}{1!\sigma^2}}x, \sqrt{\frac{1}{2!2\sigma^4}}x^2, \sqrt{\frac{1}{3!3!\sigma^6}}x^3, \dots \right]^T$$

Not too obvious what running PCA on these features mean.

Interpretable Kernel Dimension Reduction (IKDR) solves both problems...

How IKDR produce interpretable results.



HSIC is a general objective for capturing non-linear dependence to achieve IKDR.

$$\max_W HSIC(XW, Y) \text{ s.t. } W^T W = I$$

HSIC(X,Y) measures the non-linear dependence between X and Y in Feature Space. Although this make the solution interpretable, it is very difficult to solve.

In general, many IKDR problems have a common objective.

$$\max_W \sum_{i,j} \Gamma_{i,j} K_{XW_{i,j}} \text{ s.t. } W^T W = I \quad (1)$$

Symmetric Positive Definite Matrix, Kernel Matrix on XW, Grassmann Manifold Constraint.

Where is IKDR used?

Supervised Dimension Reduction for Classification

$$\max_W HSIC(XW, Y) \text{ s.t. } W^T W = I$$

Unsupervised Dimension Reduction for Clustering

$$\max_{W,Y} HSIC(XW, Y) \text{ s.t. } W^T W = I$$

Semi-supervised Dimension Reduction for Clustering Using Multiple Expert Sources

$$\max_{W,Y} \text{Tr}(Y^T \mathcal{L}_W Y) + \mu \text{Tr}(K_{XW} H K_Y H)$$

$$\text{s.t. } \mathcal{L}_W = D^{-\frac{1}{2}} K_{XW} D^{-\frac{1}{2}} W^T W = I, Y^T Y = I$$

Alternative Clustering via Dimension Reduction

$$\max_{W,Y} \text{Tr}(K_{XW} H K_Y H) - \mu \text{Tr}(K_{XW} H K_Y H)$$

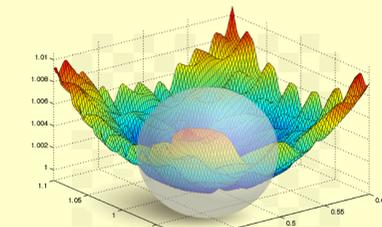
$$\text{s.t. } W^T W = I, Y^T Y = I$$

Publications that used IKDR

- Barshan, Elnaz, et al. "Supervised principal component analysis: Visualization, classification and regression on subspaces and submanifolds." Pattern Recognition 44.7 (2011): 1357-1371.
- Masaeli, Mahdokht, Jennifer G. Dy, and Glenn M. Fung. "From transformation-based dimensionality reduction to feature selection." Proceedings of the 27th International Conference on Machine Learning (ICML-10), 2010.
- Niu, Donglin, Jennifer G. Dy, and Michael I. Jordan. "Multiple non-redundant spectral clustering views." Proceedings of the 27th international conference on machine learning (ICML-10), 2010.
- Niu, Donglin, Jennifer Dy, and Michael I. Jordan. "Dimensionality reduction for spectral clustering." Proceedings of the Fourteenth International Conference on Artificial Intelligence and Statistics, 2011.
- Wu, Chieh, et al. "Iterative spectral method for alternative clustering." International Conference on Artificial Intelligence and Statistics, 2018.
- Chang, Yale, et al. "Clustering with Domain-Specific Usefulness Scores." Proceedings of the 2017 SIAM International Conference on Data Mining. Society for Industrial and Applied Mathematics, 2017.

Why is IKDR Difficult?

Optimizing W is highly non-convex and the solution must intersect the Stiefel Manifold



$$\max_W \sum_{i,j} \Gamma_{i,j} e^{-\frac{W^T(x_i-x_j)(x_i-x_j)^T W}{2\sigma^2}} \text{ s.t. } W^T W = I$$

Existing Solutions

- Dimension Growth
- Optimization Via Stiefel Manifold.
- Optimization Via Grassmann Manifold.
- Stochastic Gradient Descent.

Problems with Existing Solutions

- very slow
- Difficult to implement
- stuck at saddle point
- poor results

Our Solution

The Iterative Spectral Method (ISM)

Our Solution : The Iterative Spectral Method (ISM)

We identified a special family of kernels (The ISM family) with the following properties:

- Each kernel within the family has an associated scaled covariance matrix Φ .
- The most dominant eigenvectors of Φ is the solution to Eq. (1).
- The conic combination of ISM kernels is still in the ISM family.
- The conic combination of Φ s is the associated scaled covariance matrix for the conic combination of kernels.
- If Φ is a function of W, then Φ can be approximated using the 2nd order Taylor series

Formal Definition of the ISM family:

Definition 1. Given $\beta = a(x_i, x_j)^T W W^T b(x_i, x_j)$ with $a(x_i, x_j)$ and $b(x_i, x_j)$ as functions of x_i and x_j , any twice differentiable kernel that can be written in terms of $f(\beta)$ while retaining its symmetric positive semi-definite property is an ISM kernel belonging to the ISM family with an associated Φ matrix defined as

$$\Phi = \frac{1}{2} \sum_{i,j} \Gamma_{i,j} f'(\beta) A_{i,j} \quad (6)$$

where $A_{i,j} = b(x_i, x_j) a(x_i, x_j)^T + a(x_i, x_j) b(x_i, x_j)^T$.

Theorem 3. For any kernel within the ISM family, a Φ independent of W can be approximated with

$$\Phi \approx \text{sign}(\nabla_{\beta} f(0)) \sum_{i,j} \Gamma_{i,j} A_{i,j} \quad (7)$$

The ISM Algorithm:

Algorithm 1 ISM Algorithm

Input : Data X, kernel, Subspace Dimension q

Output : Projected subspace W

Initialization : Initialize Φ_0 using Table 1.

Set W_0 to V_{\max} of Φ_0 .

while $\|\Lambda_i - \Lambda_{i-1}\|_2 / \|\Lambda_i\|_2 < \delta$ do

 Compute Φ using Table 2

 Set W_k to V_{\max} of Φ

end

Examples of Approximations of Φ s

Kernel	Approximation of Φ s
Linear	$\Phi_0 = X^T \Gamma X$
Squared	$\Phi_0 = X^T \mathcal{L}_{\Gamma} X$
Polynomial	$\Phi_0 = X^T \Gamma X$
Gaussian	$\Phi_0 = -X^T \mathcal{L}_{\Gamma} X$
Multiquadratic	$\Phi_0 = X^T \mathcal{L}_{\Gamma} X$

Table 1: Equations for the approximate Φ s for the common kernels.

Examples of Φ s

Kernel	Φ Equations
Linear	$\Phi = X^T \Gamma X$
Squared	$\Phi = X^T \mathcal{L}_{\Gamma} X$
Polynomial	$\Phi = X^T \Psi X, \Psi = \Gamma \odot K_{XW, p-1}$
Gaussian	$\Phi = -X^T \mathcal{L}_{\Psi} X, \Psi = \Gamma \odot K_{XW}$
Multiquadratic	$\Phi = X^T \mathcal{L}_{\Psi} X, \Psi = \Gamma \odot K_{XW}^{(-1)}$

Table 2: Equations for Φ s for the common kernels.

How $K(x, x')$ become $f(\beta)$:

Kernel Name	$f(\beta)$	$a(x_i, x_j)$	$b(x_i, x_j)$
Linear	β	x_i	x_j
Squared	β	$x_i - x_j$	$x_i - x_j$
Polynomial	$(\beta + c)^p$	x_i	x_j
Gaussian	$e^{-\beta/2\sigma^2}$	$x_i - x_j$	$x_i - x_j$
Multiquadratic	$\sqrt{\beta + c^2}$	$x_i - x_j$	$x_i - x_j$

Table 3: Converting common kernels to $f(\beta)$.

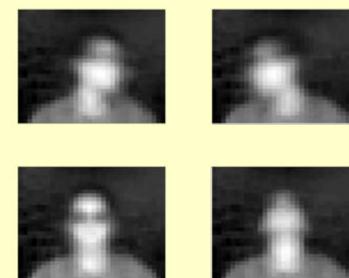
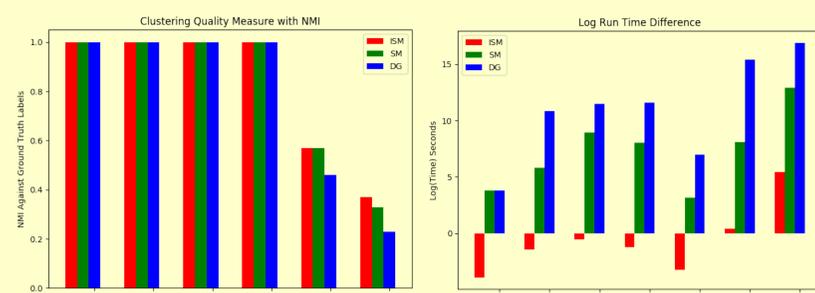
Experimental Results

Supervised	Gaussian				polynomial			
	ISM	DG	SM	GM	ISM	DG	SM	GM
Wine	Time 0.02s ± 0.01s Cost -1311 ± 26 Accuracy 95.0% ± 5%	7.9s ± 2.9s -1201 ± 25 93.2% ± 5.5%	1.7s ± 0.7s -1310 ± 26 95% ± 4.2%	16.8m ± 3.4s -1307 ± 25 95% ± 6%	0.02s ± 0.0s -114608 ± 1752 97.2% ± 3.7%	13.2s ± 6.2s -112440 ± 1719 93.8% ± 3.9%	14.77s ± 0.6s -111339 ± 1652 96.6% ± 3.7%	16.82m ± 3.6s -108892 ± 1590 96.6% ± 2.7%
Cancer	Time 0.08s ± 0.0s Cost -32249 ± 338 Accuracy 97.3% ± 0.3%	4.5m ± 103s -30302 ± 2297 97.3% ± 0.3%	17s ± 12s -31996 ± 499 97.3% ± 0.2%	17.8m ± 80s -30998 ± 560 97.4% ± 0.4%	0.13s ± 0.0s -1894 ± 47 97.4% ± 0.3%	4m ± 1.2m -1882 ± 47 97.3% ± 0.3%	3.3m ± 3s -1737 ± 84 97.4% ± 0.3%	17.5m ± 1.1m -1690 ± 108 97.3% ± 0.3%
Face	Time 0.99s ± 0.1s Cost -3754 ± 31 Accuracy 100% ± 0%	1.92d ± 11h -3431 ± 32 100% ± 0%	10s ± 5s -3749 ± 33 100% ± 0%	22.7m ± 18s -771 ± 28 99.2% ± 0.2%	0.7s ± 0.03s -82407 ± 1670 100% ± 0%	2.1d ± 13.9h -78845 ± 1503 100% ± 0%	5.0m ± 5.7s -37907 ± 15958 99.8% ± 0.2%	21.5m ± 9.8s -3257 ± 517 99.8% ± 0.2%
MINST	Time 13.8s ± 2.3s Cost -639 ± 2.3 Accuracy 99% ± 0%	> 3d N/A N/A	2.5m ± 1.0s -621 ± 5.1 98.5% ± 0.4%	> 3d N/A N/A	12.1s ± 1.4s -639 ± 2 99% ± 0%	> 3d N/A N/A	2.1m ± 3s -620 ± 5.1 99% ± 0%	> 3d N/A N/A
Unsupervised	Time 0.01s Cost -27.4 NMI 0.86	9.9s -25.2 0.86	0.6s -27.3 0.86	16.7m -27.3 0.86	0.02s -1600 0.84	14.4s -1582 0.84	2.9s -1598 0.84	33.5m -1496 0.83
Cancer	Time 0.57s Cost -243 NMI 0.8	4.3m -133 0.79	3.9s -146 0.8	44m -142 0.79	0.5s -15804 0.79	8.0m -14094 0.80	8.8m -15749 0.79	41m -11985 0.80
Face	Time 0.3s Cost -169.3 NMI 0.94	1.3d -167.7 0.95	5.3s -168.9 0.93	55.9m -37 0.89	1.0s -368 0.94	> 3d NA 0.89	22m -348 0.89	1.6d -321 0.89
MINST	Time 1.8h Cost -2105 NMI 0.47	> 3d N/A N/A	1.3d -2001 0.46	> 3d N/A N/A	8.3m -51358 0.32	> 3d N/A 0.32	0.9d -51129 0.32	> 3d N/A N/A

Table 4: Run-time, cost, and objective performance are recorded under supervised/unsupervised objectives. ISM is significantly faster compared to other optimization techniques while achieving lower objective cost.

	Supervised				Unsupervised		
	Linear	Squared	Multiquad	G+P	Linear	Squared	Multiquad
Wine	Time 0.003s ± 0s Accuracy 97.2% ± 2.8%	0.01s ± 0s 96.6% ± 3.7%	0.02s ± 0.01s 97.2% ± 3.7%	0.007s ± 0s 98.3% ± 2.6%	Time 0.02s NMI 0.85	0.04s 0.85	0.06s 0.88
Cancer	Time 0.02s ± 0.002s Accuracy 97.2% ± 0.3%	0.09s ± 0.02s 97.3% ± 0.04%	0.15s ± 0.01s 97.4% ± 0.003%	0.06s ± 0.004s 97.4% ± 0.003%	Time 0.23s NMI 0.80	0.5s 0.79	0.56s 0.84
Face	Time 0.2s ± 0.2s Accuracy 97.3% ± 0.3%	0.3s ± 0.2s 97.1% ± 0.4%	0.3s ± 0.2s 97.3% ± 0.4%	0.5s ± 0.03s 100% ± 0%	Time 0.68s NMI 0.93	0.92s 0.95	3.8s 0.92
MINST	Time 6.4s ± 0.4s Accuracy 99.1% ± 0.1%	17.4s ± 0.4s 99.3% ± 0.2%	10.6m ± 1.9m 99.1% ± 0.1%	17.6s ± 2.5s 99.3% ± 0.2%	Time 3.1m NMI 0.54	4.7m 0.54	52m 0.54

Table 5: Run-time and objective performance are recorded across several kernels within the ISM family. It confirms the usage of Φ or linear combination of Φ in place of kernels.



(a) Identity View (Mean Images)

(b) Pose View (Mean Images)

ISM's Theoretical Foundation

Theorem 1:

Given a full rank Φ with an eigengap as defined by Eq. (80), a fixed point W^* of algorithm 1 satisfies the 2nd order necessary condition using any ISM Kernel.

$$\left(\min_i \bar{\Lambda}_i - \max_j \Lambda_j \right) \geq \mathcal{C}. \quad (80)$$

Theorem 2:

A sequence of subspaces generated by Algorithm 1 contains a converging subsequence.

Theorem 3:

For any kernel within the ISM family, a Φ independent of W can be approximated with

$$\Phi \approx \text{sign}(\nabla_{\beta} f(0)) \sum_{i,j} \Gamma_{i,j} A_{i,j}.$$

Proposition 1:

Any conic combination of ISM kernels is still an ISM kernel.

Corollary 1:

The Φ matrix associated with a conic combination of kernels is the conic combination of Φ s associated with each individual kernel.

Acknowledgments: This work was made possible by the National Science Foundation (NSF IIS-1546428). The PROTECT data is supported by the National Institute of Environmental Health Sciences (P42ES017198).