



Representation Learning by Kernel Autoencoders

Pierre Laforgue, Stephan Cléménçon, Florence d'Alché-Buc

Télécom ParisTech (Chaire Machine Learning for Big Data)

Representation Learning

Autoencoders

Kernel Autoencoders

Experiments

Conclusion & Future Work

Representation Learning

Autoencoders

Kernel Autoencoders

Experiments

Conclusion & Future Work

Representation Learning (RL)

- **A representation:** a collection of features that characterize the observation

Representation Learning (RL)

- **A representation:** a collection of features that characterize the observation
- **Raw data:** redundant, non-relevant, massive
- **Ideal data:** independent, discriminative, informative

Representation Learning (RL)

- **A representation:** a collection of features that characterize the observation
- **Raw data:** redundant, non-relevant, massive
- **Ideal data:** independent, discriminative, informative
- **Feature engineering:** implies domain experts

Representation Learning (RL)

- **A representation:** a collection of features that characterize the observation
- **Raw data:** redundant, non-relevant, massive
- **Ideal data:** independent, discriminative, informative
- **Feature engineering:** implies domain experts
- **Feature/Representation learning:** automate the process



Figure 1: Machine Learning Pipeline

Representation Learning

Autoencoders

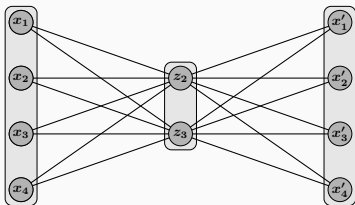
Kernel Autoencoders

Experiments

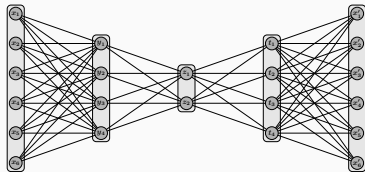
Conclusion & Future Work

Autoencoders (AEs): Principle

- **Idea:** reconstruct the input after having compressed it
- Neural network: symmetric, hour-glass shaped
- *Self-supervised* framework



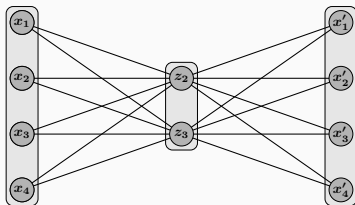
(a) 1 hidden layer AE



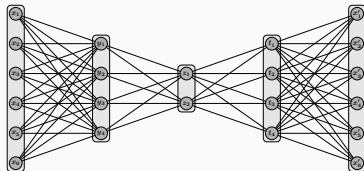
(b) 3 hidden layers AE

Autoencoders: Training

- $z = f_{\mathbf{W}, \mathbf{b}}(x) = \sigma(\mathbf{W}x + \mathbf{b})$ $x' = f_{\mathbf{W}', \mathbf{b}'}(z) = \sigma(\mathbf{W}'z + \mathbf{b}')$
- $\theta^* = \operatorname{argmin}_{\theta} \|x - x'\|^2 = \operatorname{argmin}_{\theta} \|x - f_{\mathbf{W}', \mathbf{b}'} \circ f_{\mathbf{W}, \mathbf{b}}(x)\|^2$
- Encoding $z = \sigma(\mathbf{W}^*x + \mathbf{b}^*)$



(c) 1 hidden layer AE



(d) 3 hidden layers AE

Representation Learning

Autoencoders

Kernel Autoencoders

Experiments

Conclusion & Future Work

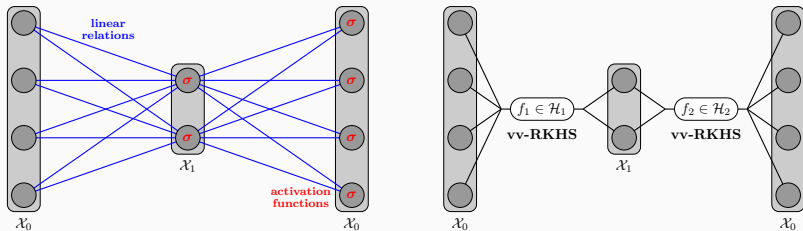


Figure 2: Standard and Kernel 2-layer Autoencoders

Formally

$$\mathbf{AE} : \min_{f_l \in \mathbf{NN}} \frac{1}{n} \sum_{i=1}^n \|x_i - f_L \circ \dots \circ f_1(x_i)\|_{\mathcal{X}_0 = \mathbb{R}^d}^2$$

$$\mathbf{KAE} : \min_{f_l \in \mathbf{vv-RKHS}} \frac{1}{n} \sum_{i=1}^n \|x_i - f_L \circ \dots \circ f_1(x_i)\|_{\mathcal{X}_0 \text{ Hilbert}}^2 + \sum_{l=1}^L \lambda_l \|f_l\|_{\mathcal{H}_l}^2$$

Formally

$$\mathbf{AE} : \min_{f_l \in \text{NN}} \frac{1}{n} \sum_{i=1}^n \|x_i - f_L \circ \dots \circ f_1(x_i)\|_{\mathcal{X}_0 = \mathbb{R}^d}^2$$

$$\mathbf{KAE} : \min_{f_l \in \text{vv-RKHS}} \frac{1}{n} \sum_{i=1}^n \|x_i - f_L \circ \dots \circ f_1(x_i)\|_{\mathcal{X}_0 \text{ Hilbert}}^2 + \sum_{l=1}^L \lambda_l \|f_l\|_{\mathcal{H}_l}^2$$

1. Feasible optimization \rightarrow novel RL algorithm

Formally

$$\mathbf{AE} : \min_{f_l \in \text{NN}} \frac{1}{n} \sum_{i=1}^n \|x_i - f_L \circ \dots \circ f_1(x_i)\|_{\mathcal{X}_0 = \mathbb{R}^d}^2$$

$$\mathbf{KAE} : \min_{f_l \in \text{vv-RKHS}} \frac{1}{n} \sum_{i=1}^n \|x_i - f_L \circ \dots \circ f_1(x_i)\|_{\mathcal{X}_0 \text{ Hilbert}}^2 + \sum_{l=1}^L \lambda_l \|f_l\|_{\mathcal{H}_l}^2$$

1. Feasible optimization \rightarrow novel RL algorithm
2. \mathcal{X}_0 Hilbert non necessarily Euclidean (not only \mathbb{R}^d)

Formally

$$\mathbf{AE} : \min_{f_l \in \mathbf{NN}} \frac{1}{n} \sum_{i=1}^n \|x_i - f_L \circ \dots \circ f_1(x_i)\|_{\mathcal{X}_0 = \mathbb{R}^d}^2$$

$$\mathbf{KAE} : \min_{f_l \in \mathbf{vv-RKHS}} \frac{1}{n} \sum_{i=1}^n \|x_i - f_L \circ \dots \circ f_1(x_i)\|_{\mathcal{X}_0 \text{ Hilbert}}^2 + \sum_{l=1}^L \lambda_l \|f_l\|_{\mathcal{H}_l}^2$$

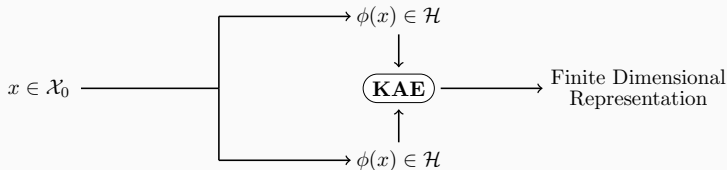
1. Feasible optimization \rightarrow novel RL algorithm
2. \mathcal{X}_0 Hilbert non necessarily Euclidean (not only \mathbb{R}^d)
3. Interesting Hilbert: (kernel) feature space

Autoencoding any data

$$\mathbf{AE} : \min_{f_l \in \mathbf{NN}} \frac{1}{n} \sum_{i=1}^n \|x_i - f_L \circ \dots \circ f_1(x_i)\|_{\mathcal{X}_0 = \mathbb{R}^d}^2$$

$$\mathbf{KAE} : \min_{f_l \in \mathbf{vv-RKHS}} \frac{1}{n} \sum_{i=1}^n \|x_i - f_L \circ \dots \circ f_1(x_i)\|_{\mathcal{X}_0 \text{ Hilbert}}^2 + \sum_{l=1}^L \lambda_l \|f_l\|_{\mathcal{H}_l}^2$$

$$\mathbf{K^2AE} : \min_{f_l \in \mathbf{vv-RKHS}} \frac{1}{n} \sum_{i=1}^n \|\phi(x_i) - f_L \circ \dots \circ f_1(\phi(x_i))\|_{\mathcal{H}}^2 + \sum_{l=1}^L \lambda_l \|f_l\|_{\mathcal{H}_l}^2$$



Representation Learning

Autoencoders

Kernel Autoencoders

Experiments

Conclusion & Future Work

Concentric Circles

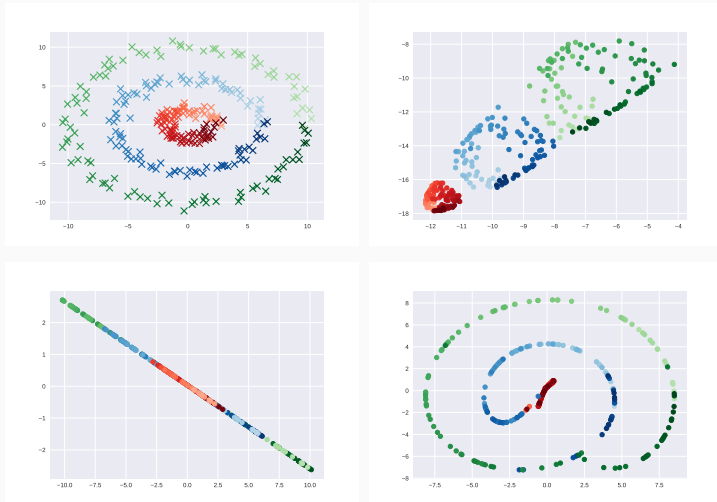


Figure 3: KAE performance on concentric circles

Molecular data

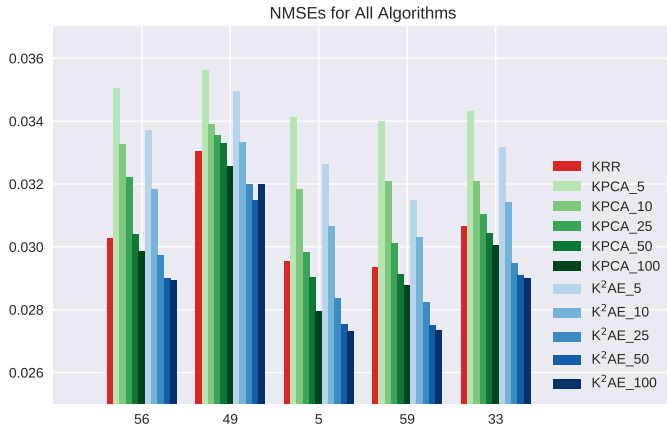


Figure 4: Performance of the different strategies on 5 cancers

Representation Learning

Autoencoders

Kernel Autoencoders

Experiments

Conclusion & Future Work

Conclusion & Future Work

- Flexible tool for Representation Learning
- Advantages from AEs and Kernel Methods
- Extension of standard AEs to any type of data
- Parallel with Kernel PCA
- Combine with a supervised criterion

Preprint available at: <http://arxiv.org/abs/1805.11028>