

Autoencoders

Qian-Yuan TANG
tangqianyuan@gmail.com

July 13 2017 @CSRC, Beijing

Outline

Background: Dimensionality Reduction

- High Dimensional Descriptions vs Low Dimensional Descriptions
- Principle Component Analysis (PCA)
- The Limits of PCA

Autoencoders

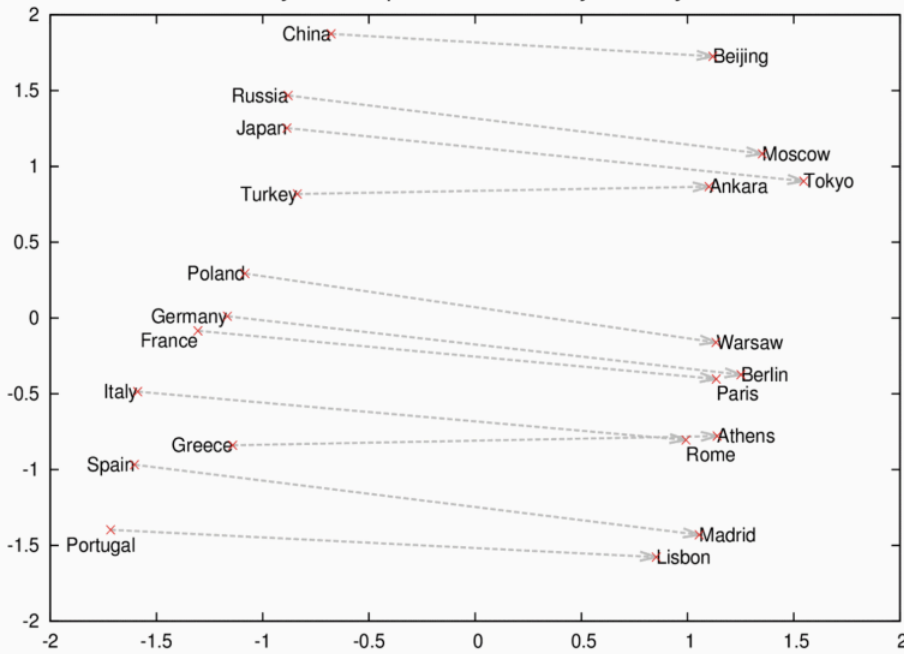
- Basic Introductions
- Stacked Autoencoder and Deep Learning
- Sparse Autoencoder
- Denoising Autoencoder (DAE)
- Contrastive Autoencoder (CAE)
- Applications of Autoencoders

Additional Discussions

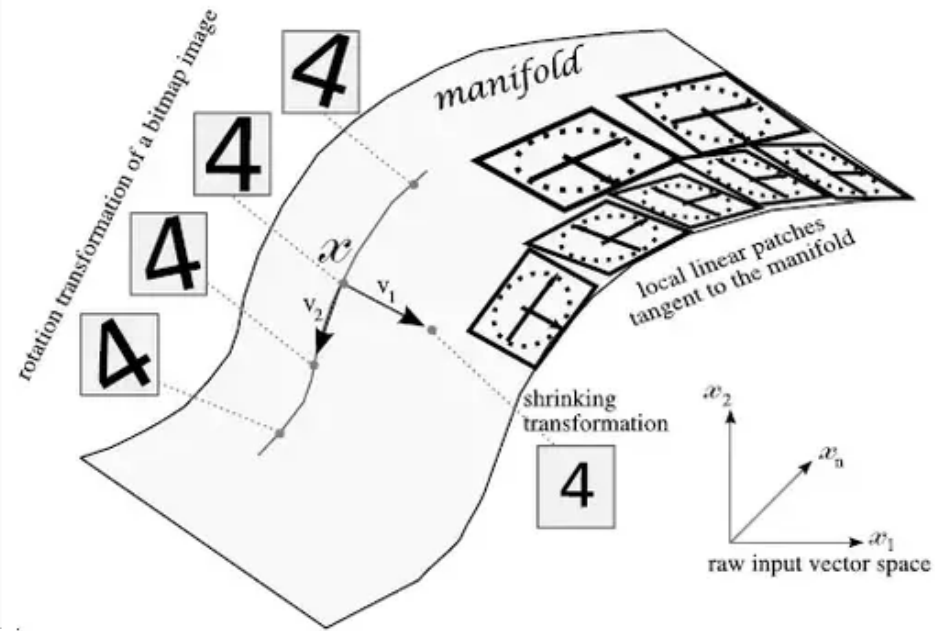
- Information Bottleneck Method
- Autoencoders as generative models (VAE and AAE)

Dimensionality Reduction: Examples (1)

Country and Capital Vectors Projected by PCA

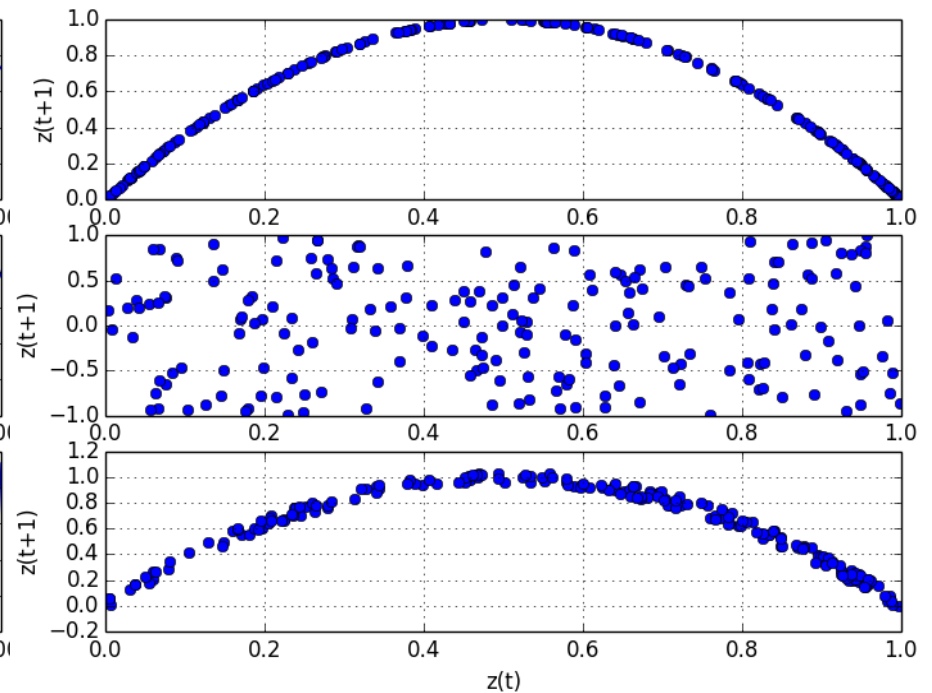
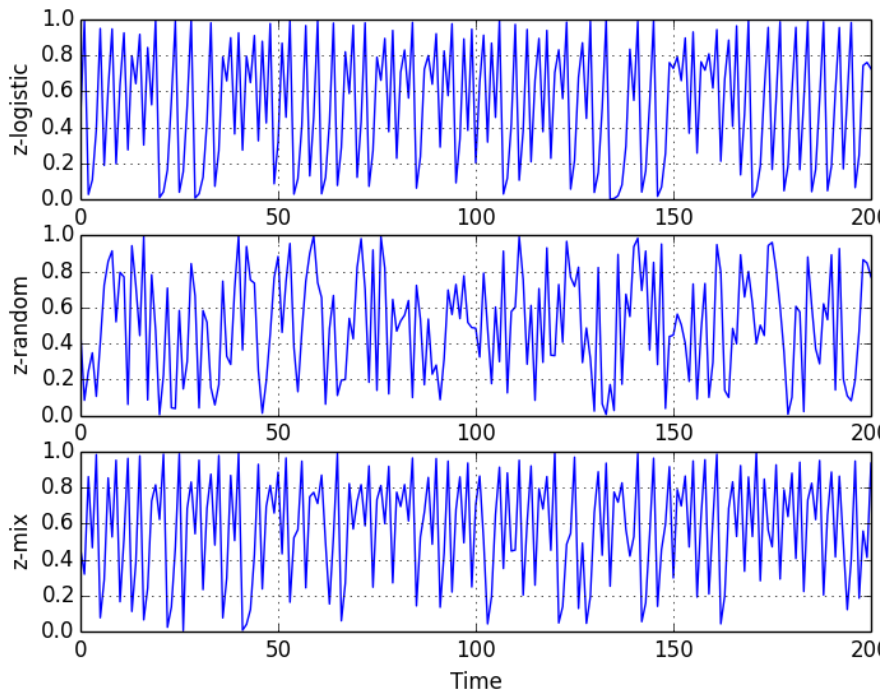


word2vec



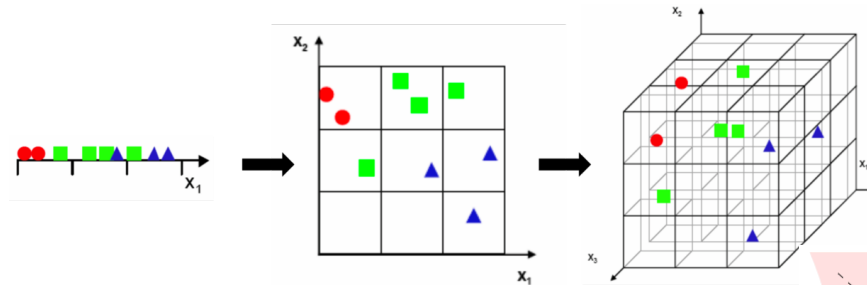
Manifold learning

Dimensionality Reduction: Examples (2)



Dimensionality reduction \rightarrow Revealing the governing equations of a chaotic system

High-dimensional data vs. Low-dimensional data



The curse of the dimensionality!

High-dimensional data

- Detailed descriptions of a system
- Data points become separable (kernel methods)

Low-dimensional data

- Coarse-grained (simplified or compressed) descriptions
- Easy to interpret and visualize
- Noise filtered
- Better generalization (possible applications in transfer learning)
- Reduced time and space complexity

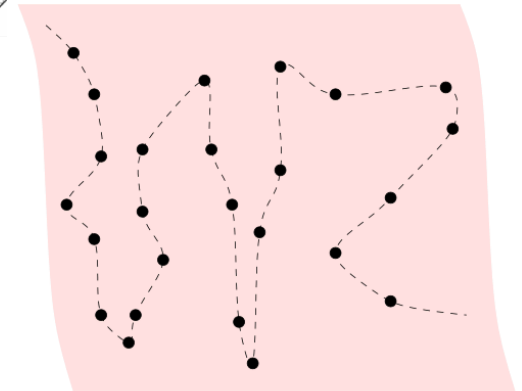


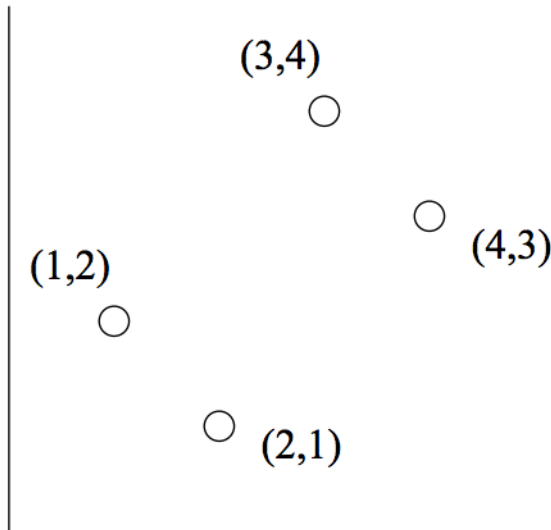
Figure 4.5: Curve or surface?

How to reduce the dimensionality of data intuitively?

- Remove data columns with small variances
- Reducing highly correlated columns

Principle Component Analysis (PCA)

An illustrative example



$$M = \begin{bmatrix} 1 & 2 \\ 2 & 1 \\ 3 & 4 \\ 4 & 3 \end{bmatrix} \quad M^T M = \begin{bmatrix} 1 & 2 & 3 & 4 \\ 2 & 1 & 4 & 3 \end{bmatrix} \begin{bmatrix} 1 & 2 \\ 2 & 1 \\ 3 & 4 \\ 4 & 3 \end{bmatrix} = \begin{bmatrix} 30 & 28 \\ 28 & 30 \end{bmatrix}$$

Eigenvalue eqn. $(30 - \lambda)(30 - \lambda) - 28 \times 28 = 0$

The solution is $\lambda = 58$ and $\lambda = 2$.

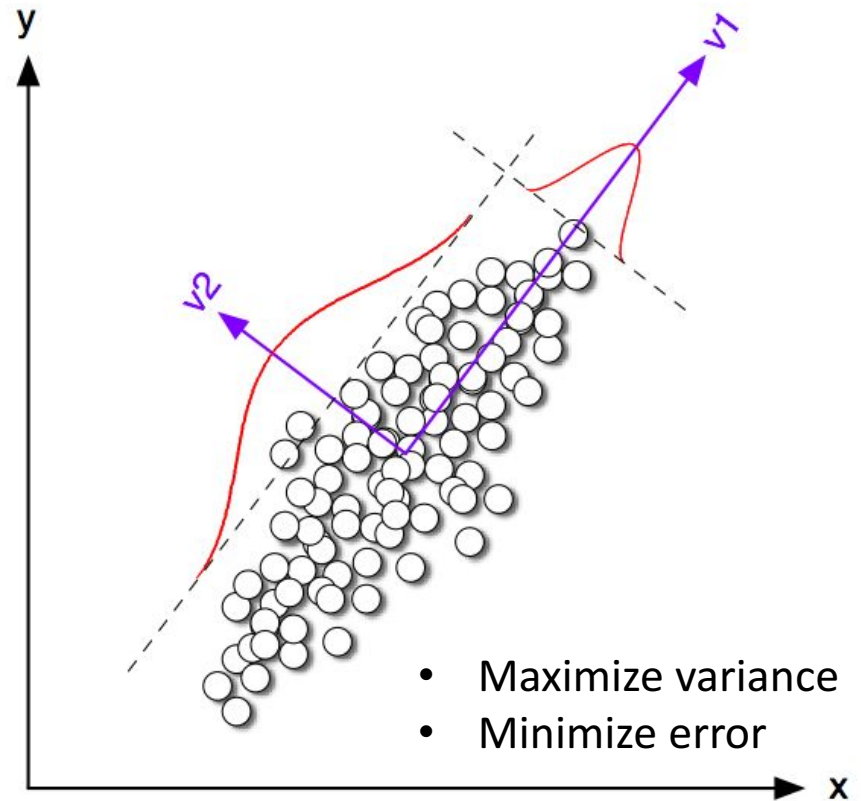
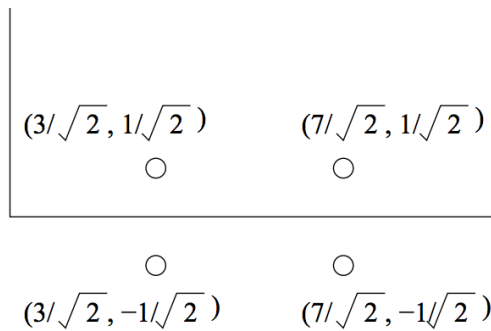
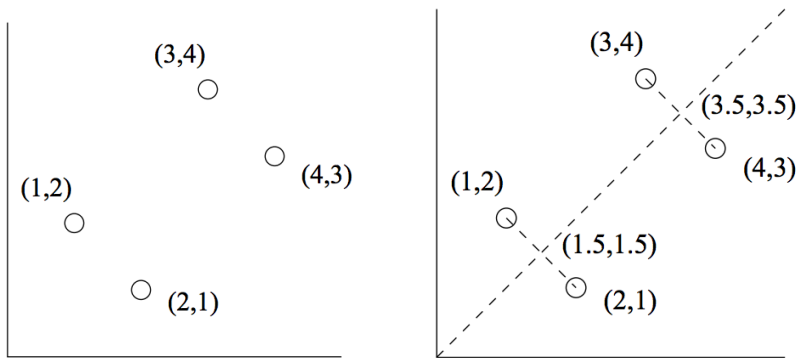
$$\begin{bmatrix} 30 & 28 \\ 28 & 30 \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} = 58 \begin{bmatrix} x \\ y \end{bmatrix} \quad x = y. \quad \begin{bmatrix} 1/\sqrt{2} \\ 1/\sqrt{2} \end{bmatrix}$$

$$\begin{bmatrix} 30 & 28 \\ 28 & 30 \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} = 2 \begin{bmatrix} x \\ y \end{bmatrix} \quad x = -y. \quad \begin{bmatrix} -1/\sqrt{2} \\ 1/\sqrt{2} \end{bmatrix}$$

$$E = \begin{bmatrix} 1/\sqrt{2} & -1/\sqrt{2} \\ 1/\sqrt{2} & 1/\sqrt{2} \end{bmatrix} \quad ME = \begin{bmatrix} 1 & 2 \\ 2 & 1 \\ 3 & 4 \\ 4 & 3 \end{bmatrix} \begin{bmatrix} 1/\sqrt{2} & -1/\sqrt{2} \\ 1/\sqrt{2} & 1/\sqrt{2} \end{bmatrix} = \begin{bmatrix} 3/\sqrt{2} & 1/\sqrt{2} \\ 3/\sqrt{2} & -1/\sqrt{2} \\ 7/\sqrt{2} & 1/\sqrt{2} \\ 7/\sqrt{2} & -1/\sqrt{2} \end{bmatrix}$$

Principle Component Analysis (PCA)

An illustrative example



Variance Maximization

$$E((\mathbf{u} \cdot \mathbf{x})^2) = E((\mathbf{u} \cdot \mathbf{x})(\mathbf{u} \cdot \mathbf{x})^T) = E(\mathbf{u} \cdot \mathbf{x} \cdot \mathbf{x}^T \cdot \mathbf{u}^T)$$

The matrix $\mathbf{C} = \mathbf{x} \cdot \mathbf{x}^T$ contains the correlations (similarities) of the original axes based on how the data values project onto them

So we are looking for w that maximizes $\mathbf{u} \mathbf{C} \mathbf{u}^T$, subject to \mathbf{u} being unit-length

Maximize $\mathbf{u}^T \mathbf{x} \mathbf{x}^T \mathbf{u}$ s.t $\mathbf{u}^T \mathbf{u} = 1$

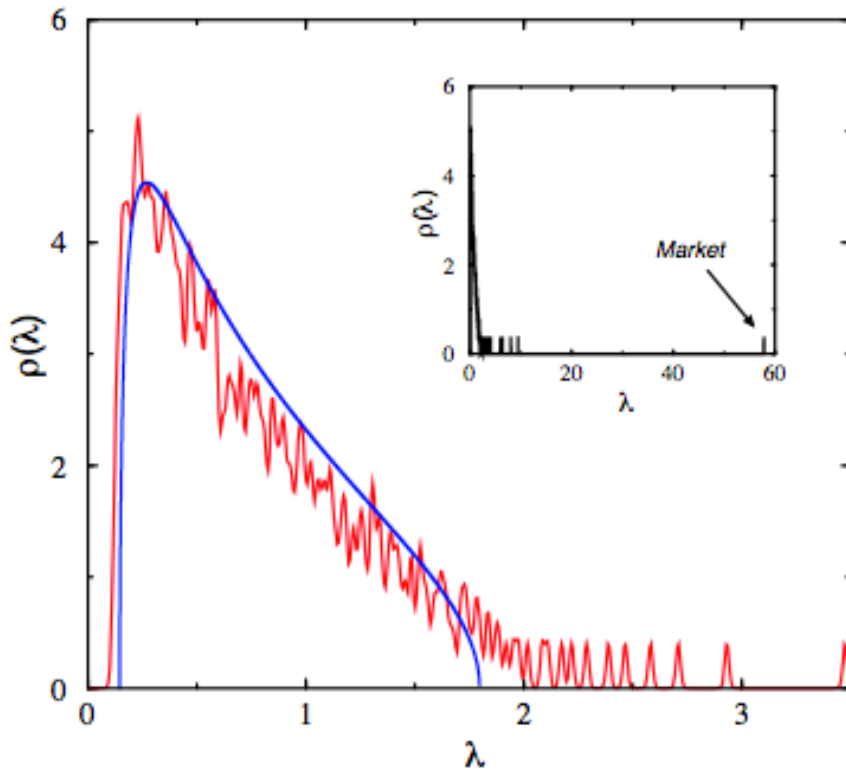
Construct Lagrangian $\mathbf{u}^T \mathbf{x} \mathbf{x}^T \mathbf{u} - \lambda \mathbf{u}^T \mathbf{u}$

Vector of partial derivatives set to zero

$$\mathbf{x} \mathbf{x}^T \mathbf{u} - \lambda \mathbf{u} = (\mathbf{x} \mathbf{x}^T - \lambda \mathbf{I}) \mathbf{u} = 0$$

As $\mathbf{u} \neq \mathbf{0}$ then \mathbf{u} must be an eigenvector of $\mathbf{x} \mathbf{x}^T$ with eigenvalue λ

Eigenvalue Distribution



Eigenvalue distribution of a random matrix

$$\rho(\lambda) = \frac{T}{N} \frac{\sqrt{(\lambda_+ - \lambda)(\lambda - \lambda_-)}}{2\pi\lambda} \quad \text{if } \lambda_+ \leq \lambda \leq \lambda_-$$

$$\lambda_{\pm} = \left[1 \pm \sqrt{\frac{N}{T}} \right]^2.$$

Empirical eigenvalue density for 406 stocks from the S&P 500, and fit using the MP distribution. Note the presence of one large eigenvalue corresponding to the market mode.

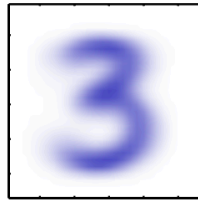
(Picture taken from Financial Applications of Random Matrix Theory: Old Laces and New Pieces, Potters et al.)

Examples of PCA

Original



$M = 1$



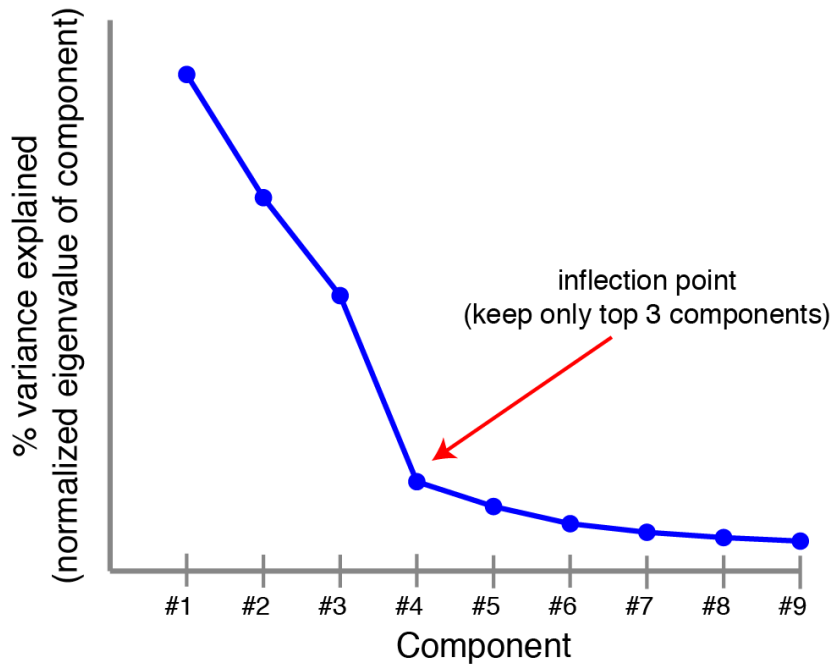
$M = 10$



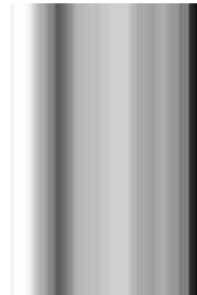
$M = 50$



$M = 250$



PCs # 0



PCs # 10



PCs # 20



PCs # 30



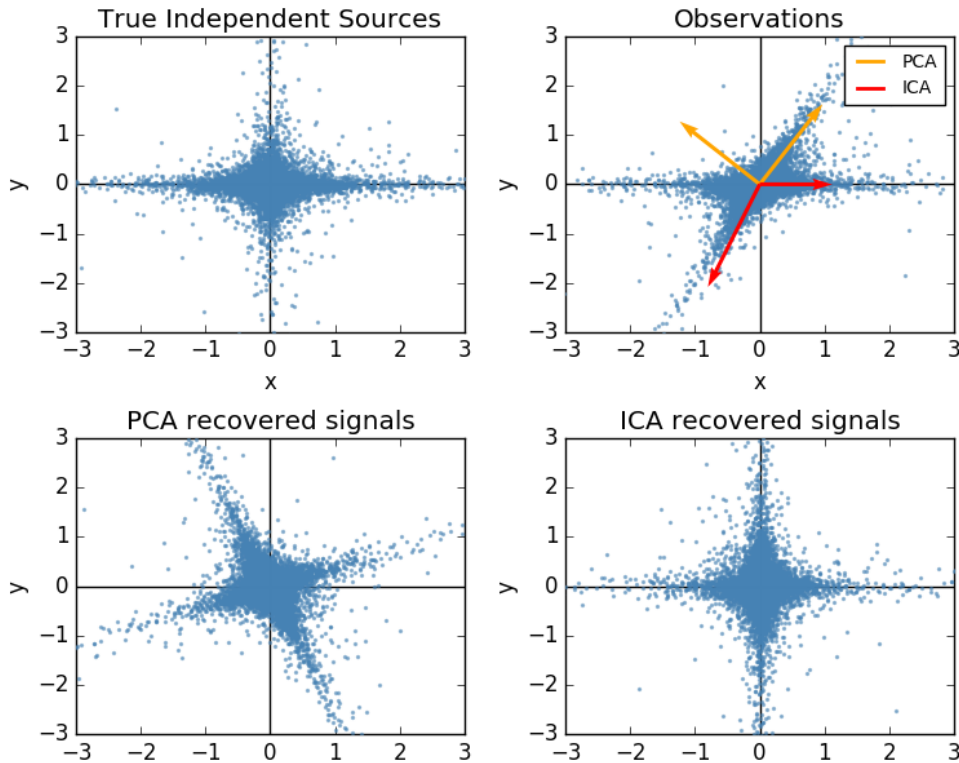
PCs # 40



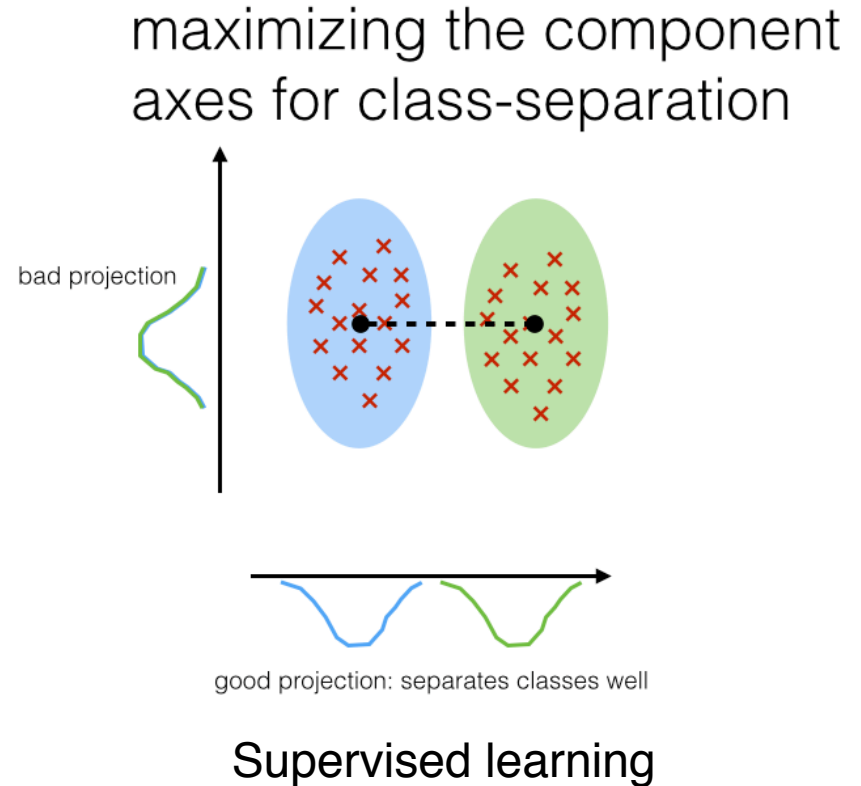
PCs # 50



The limits of PCA (1)

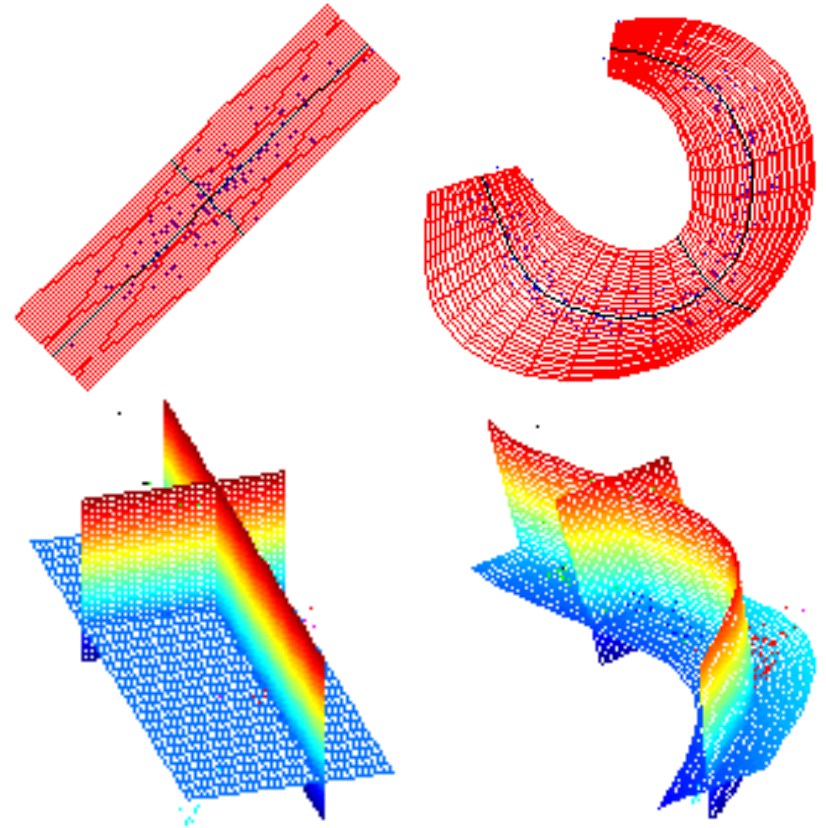
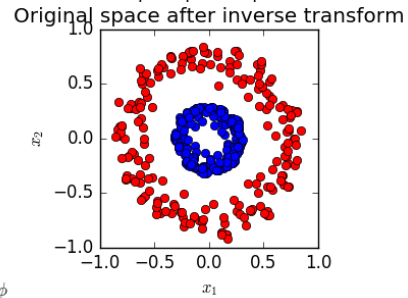
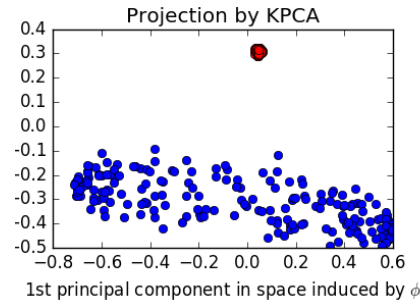
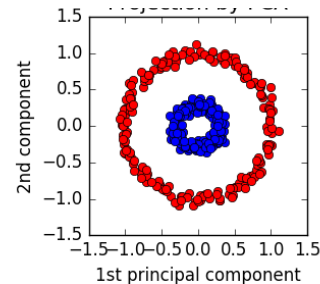
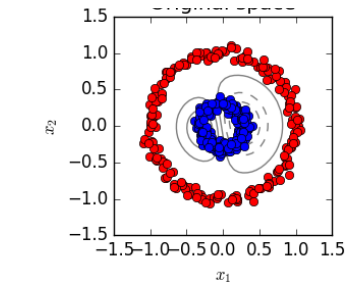
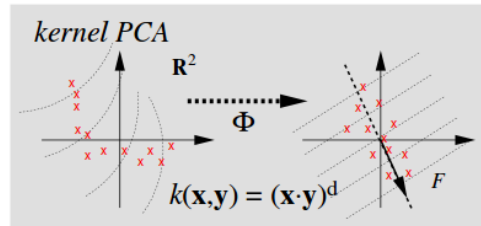
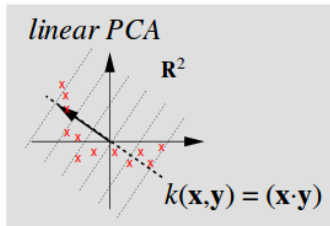


Independent Component Analysis (ICA)



Linear discriminant analysis (LDA)

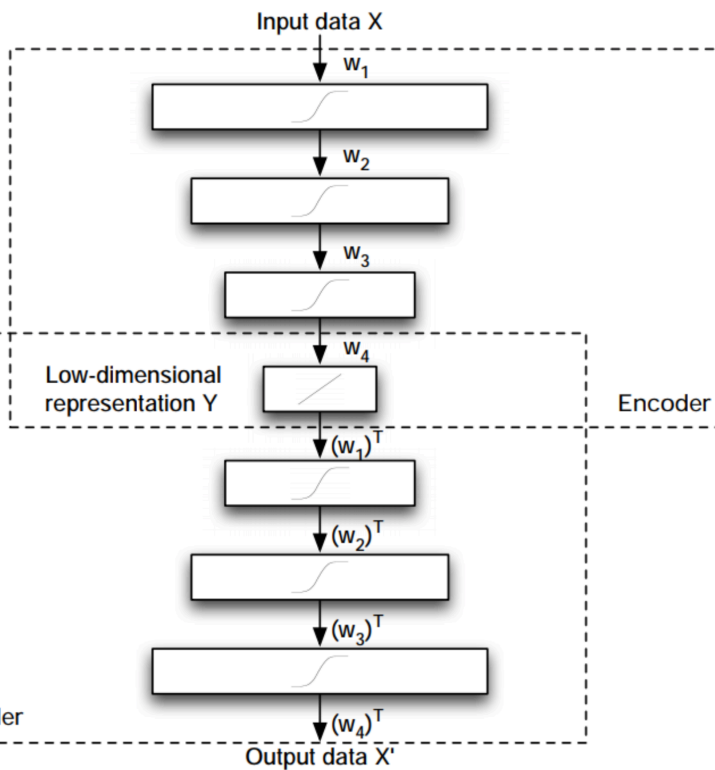
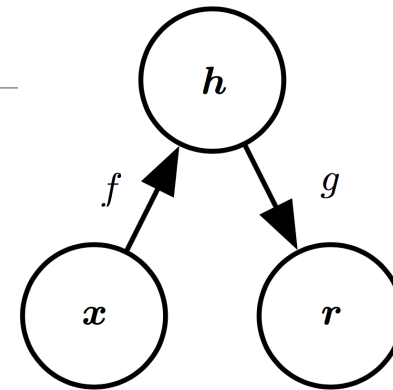
The limits of PCA (2)



Kernel PCA

Nonlinear PCA

Autoencoders



- **Basic setups**

- Input x --- encoder --- decoder --- output r
- Encoder: $h = f(x)$
- Decoder: $r = g(h) = g(f(x))$
- Minimizing loss function: $L(x, g(f(x)))$

- **Discussion**

- When will an autoencoder give the same result as PCA? (Error function)
- How to train a autoencoder? (gradient descent)

Your First Autoencoder

`openExample('nnet/TrainAutoencoderWithSpecifiedOptionsExample')`

%% Train Autoencoder with Specified Options

```
% Load the sample data.
```

```
% Copyright 2015 The MathWorks, Inc.
```

```
X = abalone_dataset;
```

```
%%
```

```
% |X| is an 8-by-4177 matrix defining eight attributes for 4177 different  
% abalone shells: sex (M, F, and I (for infant)), length, diameter, height,  
% whole weight, shucked weight, viscera weight, shell weight. For more  
% information on the dataset, type [help abalone_dataset] in the command  
% line.
```

```
%%
```

```
% Train a sparse autoencoder with hidden size 4, 400 maximum epochs, and  
% linear transfer function for the decoder.
```

```
autoenc = trainAutoencoder(X,4,'MaxEpochs',400,...  
'DecoderTransferFunction','purelin');
```

```
%%
```

```
% Reconstruct the abalone shell ring data using the trained autoencoder.
```

```
XReconstructed = predict(autoenc,X);
```

```
%%
```

```
% Compute the mean squared reconstruction error.
```

```
mseError = mse(X-XReconstructed)
```

The screenshot shows the Neural Network Training (nntraintool) window. The main window title is "Neural Network Training (nntraintool)".

Neural Network: The diagram shows an autoencoder architecture. It starts with an "Input" block with 8 nodes. This feeds into an "Encoder" block with 4 nodes. The encoder consists of a weight matrix "W" and a bias vector "b" added together, followed by a transfer function. The output of the encoder feeds into a "Decoder" block with 8 nodes. The decoder also consists of a weight matrix "W" and a bias vector "b" added together, followed by a transfer function. The output of the decoder is the "Output" block with 8 nodes.

Algorithms:

- Data Division: Training Only (dividetrain)
- Training: Scaled Conjugate Gradient (trainscg)
- Performance: Mean Squared Error with L2 and Sparsity Regularizers (msesparse)
- Calculations: MEX

Progress:

Epoch:	0	400 iterations	400
Time:		0:00:02	
Performance:	0.720	0.0191	0.00
Gradient:	0.488	0.000432	1.00e-06
Validation Checks:	0	0	6

Plots:

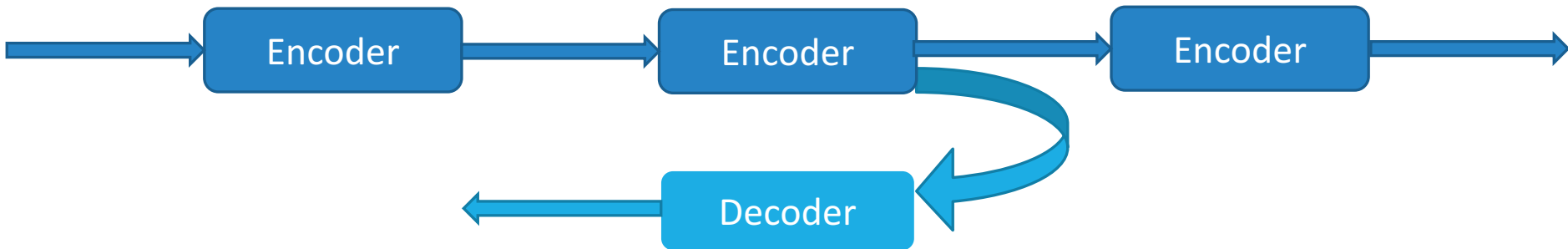
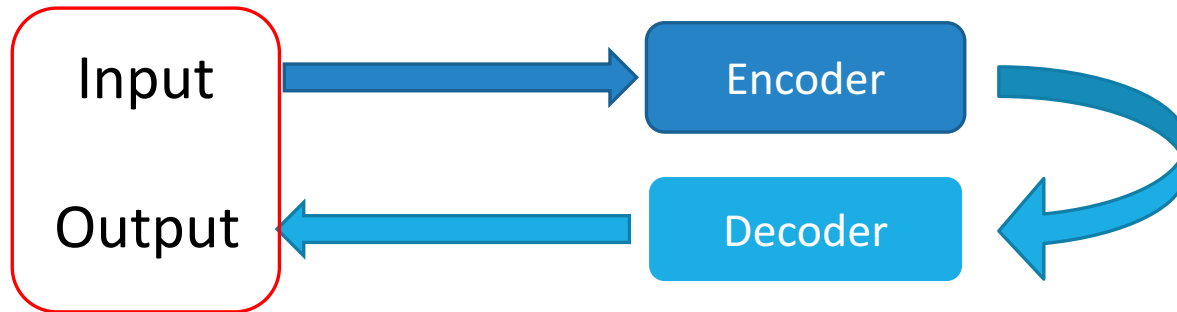
Performance (plotperform)

Plot Interval: 1 epochs

Maximum epoch reached.

Buttons: Stop Training, Cancel

Stacked Autoencoders



Greedy layer-wise training: trains the parameters of each layer individually while freezing parameters for the remainder of the model.

Stacked Autoencoders and Deep Learning

