

AutoEncoders (DRAFT)

Mark A. Austin

University of Maryland

austin@umd.edu

ENCE 688P, Spring Semester 2026

May 25, 2026

AutoEncoder Design and Analysis

Ideal Requirements/Properties:

- **Tied Weights.** Weights in i -th layer of the encoder are equal to the transpose of weights in the i -th layer of the decoder, i.e., $W_i = W_i^T$.
- **Orthogonal Weights.** Weights in the encoder are orthogonal, i.e., $W_i^T \cdot W_i = I$.
- **Uncorrelated Features.** Encoding layer outputs are not correlated.
- **Unit Norm.** The weights have unit norm, i.e.,

$$\sum_{j=1}^p w_{ij}^2 = 1 \text{ for } i = 1 \cdots k. \quad (1)$$

Principal Component Analysis

Interpretation of Principal Components:

- The **first principal component** can be defined as the direction that **maximizes** the **variance** of the projected data.
- The **i-th principal component** is a direction that **maximizes** the **variance** in the projected data and is **orthogonal** to the first $i-1$ principal components.

Applications:

- Exploratory data analysis.
- Dimensionality reduction.

Note: In general, dimensionality reduction loses information. PCA-based reduction procedures tend to minimize information loss.

Principal Component Analysis

Mathematical Procedure: Suppose that our dataset comprises n m -dimensional data points:

$$X = \begin{bmatrix} x_{11} & x_{12} & \cdots & x_{1m} \\ x_{21} & x_{22} & \cdots & x_{2m} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \cdots & x_{nm} \end{bmatrix} \quad (2)$$

Step 1. Compute the mean value for each dimension in the dataset:

$$\bar{X} = [\bar{x}_1 \quad \bar{x}_2 \quad \cdots \quad \bar{x}_m] \quad (3)$$

Principal Component Analysis

Example 1: Define straight line segment + noisy data:

```
def NoisyLineFunction (a,b,x):  
    return a + b*x + 5*(random.random() - 1.0)
```

Generate and Plot Raw Data: (x,y) coordinates:

```
5.00    3.96  
5.25    5.22  
5.50    4.30  
5.75    6.32
```

... data values removed ...

```
19.00   14.44  
19.25   18.02  
19.50   18.62  
19.75   16.12
```


Principal Component Analysis

Compute Principal Components (No components = 2):

```
pca = PCA(n_components=2)
pca.fit(X)
```

Mean values; Eigenvalue and Eigenvectors

```
--- Print mean ...
[12.375      10.5313869]

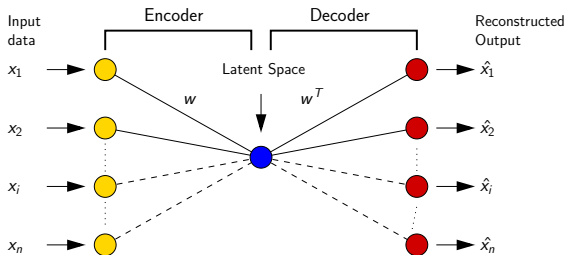
--- Print components (first and second eigenvectors) ...

[ [ 0.75668904  0.65377496]    <-- first eigenvector ...
  [-0.65377496  0.75668904] ]  <-- second eigenvector ...

--- Print variance ...
[32.50694888  1.05218475]
```


Example 1. Simplest Example

Simplest Example. No nonlinear transformation.



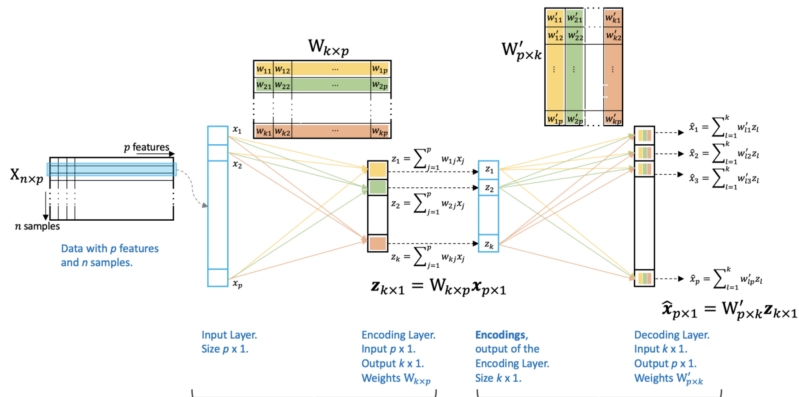
Features:

- A single hidden unit.
- Hidden unit has a **linear activation**.

Use same weight vector w for encoder/decoder. This is PCA.

Example 1. Simplest Example

Schematic for Tied Weights (Autoencoder vs PCA):



Autoencoder \Rightarrow **Encoding**—converting data to encoded features. \Rightarrow **Decoding**—reconstructing data from encoded features.

PCA \Rightarrow **PC transformation**—converting data to PC scores. \Rightarrow **Reconstruction**—reconstructing data from PC scores.

Example 1. Simplest Example

Training Procedure (Cont'd)"

$$\begin{aligned}
 L_2^2 &= \arg \min_W \sum_{i=1}^n \left[x_i^T - x_i^T w^T w \right] \left[x_i - w^T w x_i \right]. \\
 &= \arg \min_W \sum_{i=1}^n \left[x_i^T x_i - 2x_i^T w^T w x_i + x_i^T w^T w w^T w x_i \right].
 \end{aligned}$$

Equation 9 will be minimized when $w^T w$, a $(p \times p)$ matrix, equals I.

Example 1. Simplest Example

Generate Two-Dimensional Linear Dataset:

```
def StraightLineFunction (a,b,x):
    return a + b*x;

....

xvalues = list( np.arange( 5.0, 26.0, 1.00 ) );
yvalues = []
for x in xvalues:
    yvalue = StraightLineFunction( 0.0, 1.0, x)
    yvalues.append( yvalue );
```

Abbreviated Output:

```
[[ 5.  5.]          [22. 22.]
 [ 6.  6.]          [23. 23.]
 [ 7.  7.]          [24. 24.]

.....            [25. 25.]]
                   (21, 2)
```

Example 1. Simplest Example

Statistics of Dataset:

```
--- Data mean: ...
```

```
[15. 15.]
```

```
--- Data covariance: ...
```

```
[[38.5 38.5]
```

```
 [38.5 38.5]]
```

```
--- Eigenvalues: ...
```

```
[77.  0.]
```

```
--- Eigenvectors: ...
```

```
[[ 0.70710678  0.70710678]
```

```
 [-0.70710678  0.70710678]]
```

```
--- Orthogonality of eigenvectors ...
```

```
[[1.  0.]
```

```
 [0.  1.]]
```


Example 1. Simplest Example

Create Train and Test Datasets. Scale Values:

```
--- X_train_scaled data ...
```

```
[ [ 0.90  0.90 ]
  [ 0.15  0.15 ]
  [ 0.80  0.80 ]
  [ 0.40  0.40 ]
  [ 1.00  1.00 ]
  [ 0.75  0.75 ]
  [ 0.70  0.70 ]
  [ 0.60  0.60 ]
  [ 0.00  0.00 ]
  [ 0.05  0.05 ]
  [ 0.85  0.85 ]
  [ 0.45  0.45 ]
  [ 0.55  0.55 ]
  [ 0.50  0.50 ]
  [ 0.30  0.30 ]
  [ 0.95  0.95 ]
  [ 0.10  0.10 ]
  [ 0.65  0.65 ] ]
(18, 2)
```

```
--- X_test_scaled data ...
```

```
[ [ 0.20  0.20 ]
  [ 0.35  0.35 ]
  [ 0.25  0.25 ] ]
(3, 2)
```

```
--- Original test data ...
```

```
[ [ 9.00  9.00 ]
  [ 12.00 12.00 ]
  [ 10.00 10.00 ] ]
(3, 2)
```

Example 1. Simplest Example

Principal Component Analysis: `sklearn.decomposition.PCA()`

```
pca = decomposition.PCA(n_components=2)
pca.fit( X_train_scaled )
```

Summary of PCA Results:

```
--- w_pca (components) ...           --- Fit model, apply transformation
[ [ 0.70710678  0.70710678 ]         [ [-9.4280904e-02 -1.1573652e-18 ]
  [-0.70710678  0.70710678 ] ]     [ 1.1785113e-01 -8.1142219e-19 ] ]

--- w_pca * w_pca^T ...             --- Data Reconstruction ...
[ [ 1.00 -0.00 ]                   [ [ 0.20  0.20 ]
  [-0.00  1.00 ] ]                 [ 0.35  0.35 ]
                                     [ 0.25  0.25 ] ]

--- Explained Variance Ratio:
[ 1.000000e+00 9.966362e-35 ]     --- Reconstruction Loss: 0.0000
```

Example 1. Simplest Example

Autoencoder Architecture: Standard model:

Model: "sequential"

Layer (type)	Output Shape	Param #
dense (Dense)	(None, 1)	3
dense_1 (Dense)	(None, 2)	4

Total params: 7 (28.00 B)

Trainable params: 7 (28.00 B)

Non-trainable params: 0 (0.00 B)

```
encoder = Dense( encoding_dim, activation="linear", input_shape=(input_dim,),  
                use_bias = True)
```

```
decoder = Dense( input_dim, activation="linear", use_bias = True )
```

```
autoencoder = Sequential()
```

```
autoencoder.add(encoder)
```

```
autoencoder.add(decoder)
```

```
autoencoder.compile(metrics=['accuracy'], loss='mean_squared_error',  
                   optimizer='sgd')
```

Example 1. Simplest Example

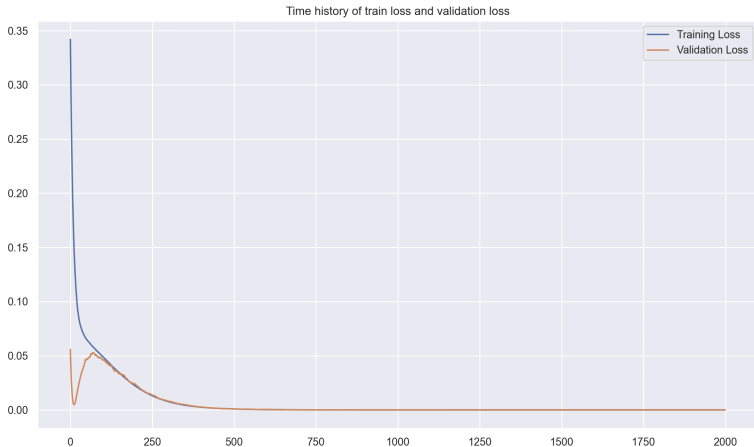
Summary: Linear activation functions.

Two layers:

- **Layer 1:** Input has two nodes, output is one node. Dense connectivity. Two weights + one bias.
- **Layer 2:** Input is one node; output is two nodes. Dense connectivity. Two weights + two bias parameters.

Example 1. Simplest Example

Time History of Learning:



Example 1. Simplest Example

Training Results:

```
--- Encoder weights:
   [ [ 1.79 -0.6 ] ]
```

```
--- Encoder bias:
   [ -0.17 ]
```

```
--- Decoder weights:
   [ [ 0.84 0.84 ] ]
```

```
--- Decoder bias:
   [ 0.14 0.14 ]
```

Testing Results:

```
--- X_test_scaled:  --->
[[0.2  0.2 ]
 [0.35 0.35]
 [0.25 0.25]]
```

```
--- Autoencoder Predictions:
   [ [ 0.20 0.20 ]
     [ 0.35 0.35 ]
     [ 0.25 0.25 ] ]
```

```
--- Original Test Data: --->
[ [ 9.  9.]
  [ 12. 12.]
  [ 10. 10.] ]
```

```
--- Autoencoder Predictions:
   [ [ 9.  9.]
     [ 12. 12.]
     [ 10. 10.] ]
```

Example 1. Simplest Example

Autoencoder Architecture: Tied Weights. No bias parameters.

Model: "sequential_1"

Layer (type)	Output Shape	Param #
dense_2 (Dense)	(None, 1)	3
dense_tied (DenseTied)	(None, 2)	3

Total params: 3 (12.00 B)

Trainable params: 3 (12.00 B)

Non-trainable params: 0 (0.00 B)

Example 1. Simplest Example

Time History of Learning:



Example 1. Simplest Example

Training Results: Tied Weights

```
--- Encoder weights:          --- Encoder bias:
  [ [ 0.305  0.878 ] ]      [ None ]

--- Decoder weights:         --- Decoder bias:
  [ [ 0.305  0.878 ] ]      [ None ]
```

Testing Results: Not great ...

```
--- X_test_scaled:  --->    --- Autoencoder Predictions:
[[0.2  0.2 ]          [ [ 0.11  0.24 ]
 [0.35 0.35]          [ 0.19  0.42 ]
 [0.25 0.25]]         [ 0.13  0.30 ] ]

--- Original Test Data: ---> --- Autoencoder Predictions:
[ [ 9.  9.]           [ [ 4.77  10.86 ]
 [ 12. 12.]           [ 6.36  14.48 ]
 [ 10. 10.] ]         [ 5.30  12.07 ] ]
```

Example 1. Simplest Example

Autoencoder Architecture: Tied Weights. Bias parameters.

Model: "sequential_2"

Layer (type)	Output Shape	Param #
dense_3 (Dense)	(None, 1)	3
dense_tied_1 (DenseTied)	(None, 2)	5

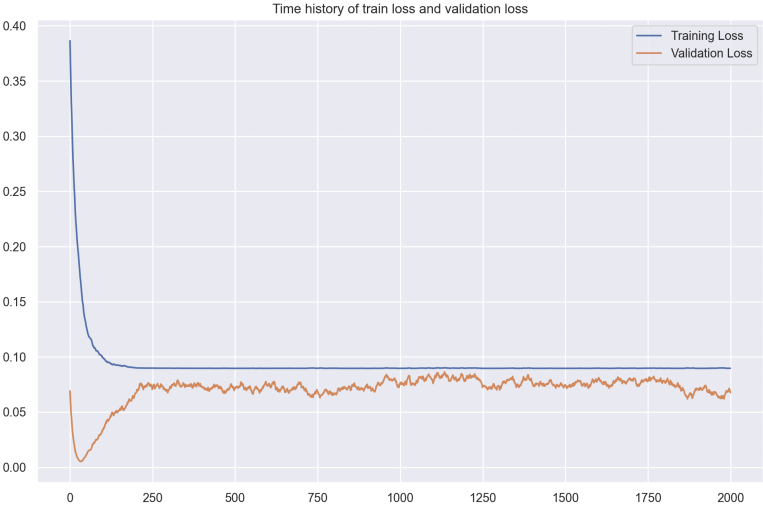
Total params: 5 (20.00 B)

Trainable params: 5 (20.00 B)

Non-trainable params: 0 (0.00 B)

Example 1. Simplest Example

Time History of Learning:



Example 1. Simplest Example

Training Results:

```
--- Encoder weights:          --- Encoder bias:
  [ [ 1.098 -0.847 ] ]       [ 0.00 ]

--- Decoder weights:         --- Decoder bias:
  [ [ 0.84  0.84 ] ]       [ 0.373  0.635 ]
```

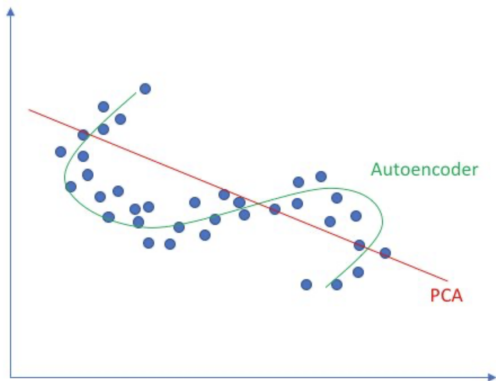
Testing Results: Inferior performance ...

```
--- X_test_scaled:  --->    --- Autoencoder Predictions:
[[0.2  0.2 ]          [ [ 0.43  0.60 ]
 [0.35 0.35]          [ 0.48  0.57 ]
 [0.25 0.25]]         [ 0.45  0.59 ] ]

--- Original Test Data: ---> --- Autoencoder Predictions:
[ [ 9.  9.]           [ [ 3.14  -0.99 ]
 [ 12. 12.]          [ 4.06  -1.52 ]
 [ 10. 10.] ]        [ 3.45  -1.17 ] ]
```

Example 2. Linear vs Nonlinear Dimensionality Reduction

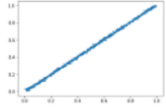
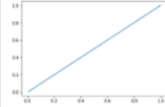
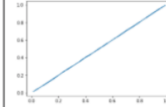
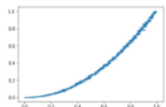
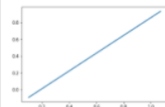
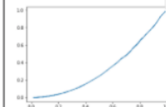
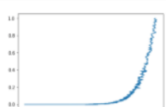


Linear vs Nonlinear Dimensionality Reduction:



Source: Nugroho H. et al., 2020

Example 2. Linear vs Nonlinear Dimensionality Reduction

Linear vs Nonlinear Dimensionality Reduction:

Function	Feature Space	PCA Reconstruction	Auto Encoder Reconstruction
$y=mx+c$			
$y=mx^2+c$			
$y=mx^8+c$			

Example 2. Linear vs Nonlinear Dimensionality Reduction

DL4J: Dataset

Example 2. Linear vs Nonlinear Dimensionality Reduction

DL4J: Neural Network Architecture

Example 3. Anomaly Detection

Example 4. Handwritten Digit Recognition

Problem Statement

- Demonstrate **anomaly detection** on MNIST using simple autoencoder.
- Goal is to **identify digits** that are **unusual**.

MNIST Handwritten Digit Dataset

- MNIST (Modified National Institute of Standards and Technology) is a database of handwritten digits provided by NIST.
- The database contains: 60,000 training images and 10,000 testing images.
- Each image is a scan of a handwritten image 0 through 9.
- Differences among images are due to variations in handwriting style.

Example 4. Handwritten Digit Recognition

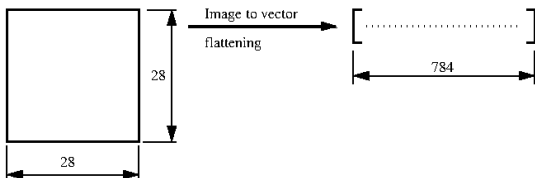
Sample Digits



Example 4. Handwritten Digit Recognition

Solution Procedure

- Images are centered on a 28x28 grid of pixels. Individual pixels take a value 0 through 255.
- Individual 28x28 images \rightarrow 1x784 data vector.



- Create network configuration with 784 inputs/outputs, contracting down to a ten-dimensional embedding vector, i.e., 784 \rightarrow 250 \rightarrow 10 \rightarrow 250 \rightarrow 784.

Example 4. Handwritten Digit Recognition

DL4J: Training Dataset

DL4J: Testing Dataset

Example 4. Handwritten Digit Recognition

DL4J: Network Configuration: 784 -> 250 -> 10 -> 250 -> 784.

```
1 MultiLayerConfiguration conf = new NeuralNetConfiguration.Builder()
2     .seed(12345)
3     .weightInit(WeightInit.XAVIER)
4     .updater(new AdaGrad(0.05))
5     .activation(Activation.RELU)
6     .l2(0.0001)
7     .list()
8     .layer(new DenseLayer.Builder().nIn(784).nOut(250)
9         .build())
10    .layer(new DenseLayer.Builder().nIn(250).nOut(10)
11        .build())
12    .layer(new DenseLayer.Builder().nIn(10).nOut(250)
13        .build())
14    .layer(new OutputLayer.Builder().nIn(250).nOut(784)
15        .activation(Activation.LEAKYRELU)
16        .lossFunction(LossFunctions.LossFunction.MSE)
17        .build())
18    .build();
19
20 MultiLayerNetwork net = new MultiLayerNetwork(conf);
21 net.setListeners(Collections.singletonList(new ScoreIterationListener(10)));
```

Example 4. Handwritten Digit Recognition

DL4J: Summary of Network Model.

```
=====  
LayerName (LayerType)    nIn,nOut    Params    ParamsShape  
=====  
layer0    (DenseLayer)  784,250    196,250    W:{784,250}, b:{1,250}  
layer1    (DenseLayer)  250,10     2,510     W:{250,10}, b:{1,10}  
layer2    (DenseLayer)  10,250     2,750     W:{10,250}, b:{1,250}  
layer3    (OutputLayer) 250,784    196,784    W:{250,784}, b:{1,784}  
-----  
Total Parameters: 398,294  Trainable Parameters: 398,294  
=====
```

Example 4. Handwritten Digit Recognition

Training Procedure:

Compute Reconstruction Error:

Example 4. Handwritten Digit Recognition

Results: Best (left) and Worst (right) Results



[Best, Worst] = [low, high] reconstruction error.

References

- Aggarwal C.C., Neural Networks and Deep Learning, Springer, 2018.
- Nugroho H., Susanty M., Irawan A., Koyimatu M., and Yunita A., Fully Convolutional Variational Autoencoder for Feature Extraction of Fire Detection System, Journal of Computer Science and Engineering, Vol. 13, No. 1, 2020.
- Watt J., Borhani R., Katsaggelos A.K., Machine Learning Refined, Second Edition, Cambridge University Press, 2020.